

# **Quantifying and reducing uncertainty in land use change model projections**

Case studies on the implications of increasing bioenergy demands

## **Quantifying and reducing uncertainty in land use change model projections**

Case studies on the implications of increasing bioenergy demands

This work was carried out within the BE-Basic R&D Program, which was granted a FES subsidy from the Dutch Ministry of Economic affairs, agriculture and innovation (EL&I).

Copyright © Judith A. Verstegen, Faculty of Geosciences, Utrecht University, 2016

ISBN: 978-90-8672-068-2

Print: CPI Koninklijke Wöhrmann B.V., Zutphen, The Netherlands

Cover design: Ivan Soenario (original image: Usina Trapiche S/A)

# **Quantifying and reducing uncertainty in land use change model projections**

Case studies on the implications of increasing bioenergy demands

## **Het kwantificeren en verlagen van onzekerheid in modelprojecties van landgebruiksverandering**

Studies naar de implicaties van de toenemende vraag naar bio-energie  
(met een samenvatting in het Nederlands)

### **Proefschrift**

ter verkrijging van de graad van doctor aan de Universiteit Utrecht  
op gezag van de rector magnificus, prof. dr. G.J. van der Zwaan,  
ingevolge het besluit van het college voor promoties  
in het openbaar te verdedigen op  
vrijdag 19 februari 2016 des middags te 2.30 uur

door

**Judith Anne Verstegen**

geboren op 26 oktober 1985 te Utrecht

Promotoren: Prof. dr. A.P.C. Faaij  
Prof. dr. S.M. de Jong

Copromotoren: Dr. D. Karssenber  
Dr. F. van der Hilst





## Table of contents

Summary.....	9
Samenvatting .....	13
1. Introduction .....	17
2. Spatio-temporal uncertainty in Spatial Decision Support Systems: a case study of changing land availability for bioenergy crops in Mozambique .....	29
3. Identifying a land use change cellular automaton by Bayesian data assimilation .....	55
4. Detecting systemic change in a land use system by Bayesian data assimilation .....	87
5. What can and can't we say about indirect land use change in Brazil using an integrated economic - land use change model?.....	121
6. A spatial optimization approach to find trade-offs between production costs and greenhouse gas emissions: a bioethanol case study .....	155
7. Synthesis .....	183
References .....	197
Dankwoord .....	217
About the author .....	221
List of journal publications.....	223





## Summary

Land use change is a central issue in the sustainability debate, because of its impacts on e.g. climate change, water availability and quality, soil quality and erosion, and biodiversity. Continuing population growth, shifting diets towards higher meat consumption and increasing bioenergy demands call for the exploration of possibilities for sustainable land use change pathways with minimal negative impacts. Spatially explicit demand driven land use change models are tools that support such explorations by projecting the spatial dynamics of a predefined set of land uses over a given period. Different potential future pathways can be assessed with these models by the evaluation of divergent scenarios. Designing land use change models is not straightforward, because the dynamic processes and feedbacks in the land use system are complex and only partially understood. This results in uncertainties in model structure, inputs and parameters, which propagate to the land use change projections and derived impacts. The fact that only a limited number of scenarios can be analysed results in additional uncertainty, i.e. a lack of clarity about the complete set of potential futures, referred to as solution space uncertainty. In land use change impact assessments it is essential to recognize these uncertainties, because management or policy decisions based on erroneous projections can be costly or irreversible, either from an environmental or from an economic point of view. In this thesis, methods are developed to quantify and reduce uncertainty in land use projections. Case studies in this thesis are focused on bioenergy as a driver of land use change, because bioenergy is experiencing a large demand increase and is - being designed as a sustainable alternative for fossil energy - under extra pressure in the sustainability debate to minimize negative impacts.

In this thesis, the solution space uncertainty is quantified by comparing projected land use change impacts with land use change impacts that are minimized based on a selected set of environmental impacts. This is demonstrated for a case study of minimizing production cost and GHG emission for an ethanol supply increase in Goiás, Brazil, for 2030. Results show, for example, that the costs calculated from a land use scenario projection are  $715 \text{ US\$}_{2014} / \text{m}^3$  ethanol, while the minimum attainable production costs are  $656 \text{ US\$}_{2014} / \text{m}^3$  ethanol. This places the scenario in perspective and thereby supports policy makers in the decision making process. The developed methodology has the prospect to identify trade-offs between different impacts and win-win situations for other regions, scales, objectives and commodities.

To quantify error propagation from the model components towards the land use projections the PCRaster Land Use Change model (PLUC) is developed, coupled to a Monte Carlo (MC) analysis scheme. Because the definition of uncertainty in all

model components is an integral part of the modelling framework, the end user can easily evaluate uncertainty for multiple scenarios or case studies. This is relevant as uncertainty in model structure, inputs and parameters often varies between scenarios and especially between case studies. In addition, the embedded coupling allows for land use uncertainty analysis for each model time step, which is important as uncertainty might evolve non-linearly through time.

Methodologies are developed to define uncertainty in each model component and PLUC is applied to dynamically evaluate the effects of these uncertainties on the model output, i.e. the land use projections. Such a full scope error propagation assessment is new in land use change modelling. Generally, we find that: 1) output uncertainty is attribute dependent, e.g. in a projection for Brazil for 2030, the direct land use change (dLUC) area has a coefficient of variation (cv) of only 0.02, while the indirect land use change (iLUC) area has a cv of 0.72, and 2) output uncertainty is scale dependent, decreasing from lower to higher aggregation levels, e.g. the maximum variance in the total potential yield of eucalyptus in a projection for Mozambique for 2030 drops from  $1 \cdot 10^5$  to 680 to 298  $(\text{kg km}^{-2} \text{year}^{-1})^2$  when scaling from cell level to province level to country level, respectively. In general, non-spatial inputs determine output uncertainty at high aggregation levels and spatial inputs determine output uncertainty at low aggregation levels. We also find that the land use system can experience systemic changes, inducing invalidity of the land use model structure. Our method to incorporate such changes into the land use change model increased the projected 95% confidence interval of the land use area per  $25 \times 25 \text{ km}^2$  block by a factor of 2.

To reduce the uncertainty in land use projections, a particle filter is coupled to PLUC. A particle filter is a data assimilation technique that updates prior knowledge about model structure, inputs and parameters by integrating observational data into the model during runtime. This method has the advantage that uncertainty in the observational data, such as classification errors, can be taken into account. The particle filter considerably reduces output uncertainty, e.g. for a case study for São Paulo, Brazil, the 95% confidence intervals of output land use metrics were reduced by at least a factor 3, compared to a run without the particle filter.

Yet, even with reduced uncertainty, the uncertainties in land use projections are large, especially at local scale levels (up to  $100 \times 100 \text{ km}^2$ ) and for long time frames (more than a decade). We conclude that the value of spatially explicit land use change models for answering questions about the future directions the system at local scale levels or for long time frames is limited. For policy implementation strategies this implies that the confidence interval of a land use attribute is often so wide that it is likely to straddle a legislation threshold. Therefore, we deem threshold evaluation for land use indicators, for example iLUC, a very questionable practice. Independent of the implications of these large uncertainties, we at least

created the conditions for policy makers to account for uncertainty in their decisions by presenting the robustness of the land use projections in an understandable way. Communication of uncertainties in land use change models should become common practice, because the users of the land use projections from these models are entitled to the opportunity to grasp the (un)reliability of these projections and related impacts.



## Samenvatting

Landgebruiksverandering is één van de hoofdkwesties in het duurzaamheidsdebat, door haar impacts op bijvoorbeeld klimaatverandering, de beschikbaarheid en kwaliteit van water, de bodemkwaliteit, landdegradatie en biodiversiteit. De aanhoudende populatiegroei, dieetveranderingen richting hogere vleesconsumpties en de groeiende vraag naar bio-energie roepen op tot onderzoek naar de mogelijkheden voor duurzame landgebruikstrajecten met een minimum aan negatieve effecten. Ruimtelijk expliciete, vraaggestuurde landgebruiksmodellen kunnen als instrumentarium dienen voor dergelijk onderzoek doordat ze de dynamieken van een set van landgebruikstypen over een bepaalde tijdsperiode simuleren. Verschillende landgebruikstrajecten kunnen bekeken worden door met zulke modellen uiteenlopende scenario's te analyseren. Het ontwikkelen van landgebruiksmodellen is niet eenvoudig, omdat de dynamische processen en terugkoppelingen in het landgebruikssysteem complex en niet volledig bekend zijn. Dit resulteert in onzekerheden in de modelstructuur, input en parameters, die zich voortplanten naar de landgebruiksprojecties en daarvan afgeleide effecten. Het feit dat slechts een beperkt aantal scenario's geanalyseerd kan worden zorgt voor extra onzekerheid, doordat geen totaalbeeld gevormd kan worden van de volledige set van mogelijke landgebruikstrajecten. We definiëren dit type onzekerheid als 'onzekerheid in de oplossingsruimte'. Bij het analyseren van de effecten van landgebruiksverandering is het belangrijk om onzekerheden te erkennen, omdat management- of beleidsbesluiten gebaseerd op onjuiste landgebruiksprojecties duur of onomkeerbaar kunnen zijn, hetzij vanuit een duurzaamheidsperspectief hetzij vanuit een economisch perspectief. In dit proefschrift zijn methoden ontwikkeld om onzekerheden in landgebruiksprojecties te kwantificeren en te reduceren. Studies in dit proefschrift zijn gericht op bio-energie als drijvende kracht achter landgebruiksverandering, omdat de vraag hiernaar op dit moment sterk toeneemt. Een andere reden is dat bio-energie, bedoeld als duurzaam alternatief voor fossiele energie, onder extra druk staat in het duurzaamheidsdebat om negatieve effecten te verminderen.

In dit proefschrift is de onzekerheid in de oplossingsruimte gekwantificeerd door de geprojecteerde landgebruikseffecten te vergelijken met de landgebruikseffecten die geminimaliseerd zijn op basis van een set van duurzaamheids- en economische criteria. Deze methode is gedemonstreerd op basis van een studie waarin de productiekosten en broeikasgasemissies van een toename in aanbod van ethanol in Goiás, Brazilië, geminimaliseerd worden. De resultaten laten bijvoorbeeld zien dat de productiekosten voor een bepaald scenario  $715 \text{ US}\$_{2014} / \text{m}^3$  ethanol zijn, terwijl de minimum haalbare productiekosten  $656 \text{ US}\$_{2014} / \text{m}^3$  ethanol zijn. Deze informatie plaatst het scenario

in perspectief en ondersteunt daarmee beleidsmakers bij de besluitvorming. De ontwikkelde methode biedt mogelijkheden om afwegingen te maken tussen verschillende effecten en om win-win situaties voor andere regio's, schalen, doelen of goederen te identificeren.

Om de foutenvoortplanting van de verschillende modelcomponenten naar de landgebruiksprojecties te kwantificeren is het PCRaster landgebruiksmodel, PLUC, ontwikkeld, gekoppeld aan een Monte Carlo (MC) analyse schema. Omdat hiermee de definitie van de onzekerheden in alle modelcomponenten een integraal deel van het model is, kan de eindgebruiker eenvoudig de onzekerheid evalueren in meerdere studies of in meerdere scenario's binnen één studie. Dit is relevant aangezien de onzekerheid in modelstructuur, input en parameters vaak varieert tussen verschillende scenario's en met name tussen verschillende studies. Bovendien zorgt de geïntegreerde koppeling ervoor dat de onzekerheid geanalyseerd kan worden in elke tijdstap van de gesimuleerde periode, wat belangrijk is omdat onzekerheid zich vaak niet lineair ontwikkelt over de tijd.

Methodologieën zijn ontwikkeld om de onzekerheid in elke modelcomponent te definiëren en PLUC is toegepast om de dynamische effecten van deze onzekerheid op de modeluitkomsten, de landgebruiksprojecties, te evalueren. Een dergelijke volledige foutenvoortplantingsevaluatie is nieuw binnen de modellering van landgebruiksverandering. Over het algemeen komen we tot twee conclusies. Ten eerste dat de onzekerheid in de modeluitkomsten attribuutafhankelijk is. De variatiecoëfficiënt (cv) van het areaal directe landgebruiksverandering is bijvoorbeeld 0.02 terwijl de cv van het areaal indirect landgebruiksverandering 0.72 is in Brazilië. Ten tweede dat de onzekerheid in de modeluitkomsten schaalafhankelijk is, afnemend van lage naar hoge aggregatieniveaus. De maximale variantie in het totale potentiële gewasopbrengst van eucalyptus in Mozambique loopt bijvoorbeeld terug van  $1 \cdot 10^5$  naar 680 naar 298 ( $\text{kg km}^{-2} \text{jaar}^{-1}$ )<sup>2</sup> wanneer opgeschaald wordt van rastercel ( $1 \text{ km}^2$ ) naar provincie naar landelijk schaalniveau. Over het algemeen wordt de onzekerheid in de modeluitkomsten op hoge aggregatieniveaus bepaald door niet-ruimtelijke input en die op lage aggregatieniveaus door ruimtelijke input. We stellen ook vast dat systeemveranderingen kunnen optreden in het landgebruikssysteem, die ongeldigheid van de modelstructuur veroorzaken. Onze methode om zulke systeemveranderingen in het landgebruiksmodel mee te nemen, resulteerde in een toename met een factor 2 van de breedte van het 95%-betrouwbaarheidsinterval van het landgebruiksareaal per aggregatie blok van  $25 \times 25 \text{ km}^2$ .

Om de onzekerheid in landgebruiksprojecties te verlagen, is er een 'particle filter' gekoppeld aan PLUC. Een particle filter is een data assimilatie techniek die a-priorische waarschijnlijkheden van modelstructuur, input en parameters aanpast

door observatiedata gedurende de modelrun in het model te integreren. Deze methode heeft als voordeel dat de onzekerheid in de observatiedata, bijvoorbeeld als gevolg van classificatiefouten, meegenomen kan worden. Het particle filter reduceert de onzekerheid in de modeluitkomsten aanzienlijk. De breedtes van de 95%-betrouwbaarheidsintervallen van verschillende maatstaven, die zijn afgeleid van landgebruiksverandering, worden bijvoorbeeld verminderd met ten minste een factor 3, vergeleken met de resultaten van een modelrun zonder particle filter.

Desalniettemin zijn, zelfs met gebruik van het particle filter, de onzekerheden in landgebruiksprojecties groot, met name op lokale schaalniveaus (tot 100 x 100 km<sup>2</sup>) en voor lange simulatieperiodes (meer dan een decennium). We concluderen dat de waarde van ruimtelijk expliciete landgebruiksmodellen voor het beantwoorden van vragen over het toekomstige traject van het systeem op lokale schaalniveaus en voor lange tijdsperiodes beperkt is. De implicatie van de gevonden onzekerheden voor beleidsimplementatiestrategieën is dat het betrouwbaarheidsinterval van een maatstaf van landgebruiksverandering vaak zo breed is dat het zeer waarschijnlijk een vastgestelde grenswaarde overspant. Daarom zijn wij van mening dat het vaststellen en evalueren van grenswaarden voor landgebruiksindicatoren, zoals indirecte landgebruiksverandering, een twijfelachtig gebruik is. Onafhankelijk van de implicaties van de gevonden onzekerheden, hebben we in dit proefschrift ten minste de condities gecreëerd voor beleidsmakers om onzekerheden mee te nemen in de besluitvorming, door de robuustheid van landgebruiksprojecties op een begrijpelijke manier te presenteren. De communicatie van onzekerheden in landgebruiksmodellering zou gangbaar moeten worden omdat de gebruikers van deze projecties recht hebben op de mogelijkheid om te bevatten hoe (on)betrouwbaar deze projecties en de daarvan afgeleide effecten zijn.





# 1. Introduction

## 1.1. The significance of land use and land cover change

Land use and land cover change is a central issue in the sustainability debate because of its wide range of environmental impacts (Lambin and Meyfroidt, 2011, Nesheim et al., 2014). Key sustainability issues affected by land use and land cover change are climate change, water availability and quality, soil quality and erosion, and biodiversity. For example, vegetation cover transformation and the conversion of carbon rich lands affect the climate through albedo change and the emission of greenhouse gases (GHGs) (e.g. Stocker et al., 2013). GHG emissions from land use change are difficult to quantify, but are estimated to have been between 10% and 20% of the total annual anthropogenic GHG emissions in the past two decades (Canadell et al., 2007, IPCC, 2014). Furthermore, the expansion of agriculture and accompanying irrigation and fertilizer usage impact water availability and quality (e.g. Scanlon et al., 2007). Piao et al. (2007) reconstruct that land use and land cover changes account for about 50% of the global runoff trend of the last century. Deforestation and forest fragmentation can reduce biodiversity (e.g. Chaplin-Kramer et al., 2015). Sala et al. (2000) estimate that, in the 21<sup>st</sup> century, of all factors influencing biodiversity, land use and land cover changes have the largest global impact, about 1.5 times as large as climate change and more than two times as large as other factors like changes in atmospheric CO<sub>2</sub>, and nitrogen deposition.

Aiming to reach sustainable future pathways with minimal negative impacts, a sound understanding of the land system is required. Land use and land cover change has a wide range of drivers that can be natural or anthropogenic. Changes in *land use*, i.e. a classification of the purpose of exploitation of the land, are by definition anthropogenic. Changes in *land cover*, i.e. a physical classification of the land, can be natural, e.g. ecological succession or natural forest fires (Cihlar and Jansen, 2001), but can also be human induced, e.g. a change from cropland to urban area. In the latter case the land cover change is caused by a land use change, for the given example a change from agricultural use to residential use. In the sustainability debate, the interest in the drivers of land use and land cover change is focused on the anthropogenic drivers, as these are the ones we can try to steer to establish more sustainable pathways. For this reason, the focus in this thesis is on human induced changes and the term 'land use change' is used throughout, instead of 'land use and land cover change'.

The most apparent anthropogenic driver of land use change is the increasing world population (Agarwal et al., 2002, Alexander et al., 2015), as people need space to live, food should be cultivated to feed them, materials like wood and fibres are needed for houses and clothing, etc. Another driver of land use change is the

ongoing dietary transition towards an increased consumption of meat, influenced by a rise in GDP per capita. The production of animal products accounts for 65% (439 Mha) of the land use change over the past 50 years (Alexander et al., 2015). More recently, awareness of the exhaustibility and environmental impacts of fossil energy have, among other causes, resulted in an increasing land requirement for bioenergy crop production (Popp et al., 2014). Although the current area used for bioenergy crops is relatively small, 1.8% (81 Mha) of the global agricultural land, it had a share of 36% (3.2 Mha per year) in the net agricultural expansion area since 1994 (Alexander et al., 2015). Globalization increases the average distance between the locations of production and consumption (Lambin and Meyfroidt, 2011), which complicates tracking the connection between these land use change drivers and their impacts.

Ongoing changes in global population, diets, and bioenergy demands are expected towards the future. The world population is expected to continue to grow with about 1% per year (World Bank, 2015), the developing countries are likely to steadily shift to levels of meat consumption typical of western diets (Alexandratos and Bruinsma, 2012), and the global biofuel demand is projected to grow from its current 3.0% of the road transport fuel demand to 3.5% by 2020 (OECD/IEA, 2014). Given the wide range of impacts of land use change and the awareness that these can, at least partly, be steered by policies, the implications of the future demands for land should be assessed and the possibilities for sustainable pathways should be explored to advise policy makers. The wide range of impacts play at different time scales and thereby require different time horizons for land use projections, e.g., for a government period (~4 years), a policy period (10-25 years), or a climate change period (50-100 years).

In the impact assessments it is essential to know how solid the land use change projections are, because decisions based on erroneous projections can be costly, either from an environmental or from an economic point of view, due to the frequent irreversibility of these decisions. In other words, there is a need to recognize the uncertainty in these projections. Knowing the source of the uncertainty also offers means to reduce the uncertainty, for example by intensifying the data collection regarding this source. The aim of this thesis is to develop methods to quantify and reduce uncertainty in land use projections. In the following, background information is provided on land use change models used to generate land use change projections (section 1.2), the components that such models consist of (section 1.3), the importance and difficulties of uncertainty quantification (section 1.4), and the case study this thesis is focused on (section 1.5). After that, the research problem and research questions are defined (section 1.6).

## 1.2. Land use change models

Land use change can be assessed in different ways. Agricultural statistics (census) databases, such as FAOSTAT (FAO, 2015), help to construct past trends of the areas of the different land use types, and remote sensing can be used to spatially evaluate land cover changes (see e.g. Allen et al., 2013 for a review). Since future trends depend on the aforementioned drivers, they are often not simply an extrapolation of past trends. Because of this, simulation models, which incorporate a theory on how the land use system relates to the drivers, are used to project future land use trends (Veldkamp and Lambin, 2001). Land use change impacts can vary widely over space, for example, land use change related GHG emissions depend on initial land use, soil type and climate conditions (e.g. van der Hilst et al., 2014). Therefore, spatially explicit land use change models are suitable to study the expected impacts.

Spatially explicit land use change models start with an initial land use situation for a given case study area and use an inferred transition function, representing the processes of change, to simulate the expansion and contraction of a predefined set of land use types over a given period. Land use change models help to improve our understanding of the land system by establishing cause-effect relations and testing these on historic data. In this way they help to identify the drivers of land use change and their relative importance. In addition, the models can be used to explore future land use pathways for different scenarios (Lambin et al., 2000, Magliocca et al., 2015), to answer questions of the 'what-if'-type. For example, what will be the response of the land use system if a zoning policy is introduced? When a simulation model is used together with a visualization tool targeting to evaluate such what-if questions to support decision making, it is sometimes called a Spatial Decision Support System (SDSS) (Geertman and Stillwell, 2004).

Various land use change modelling approaches exist. Agent based models (ABMs) try to capture the complexity of human decision making by defining the actors of change as autonomous entities (Parker et al., 2003, Matthews et al., 2007). What a single agent represents depends on the scale of the study, and ranges from an individual person (farmer, citizen) to a community of hundreds of people. In general, ABMs apply to local studies, mainly because for a large area both the acquisition of sufficient data to parameterize the model at the individual level for all relevant processes, and the computation time become problematic. In addition, the poor level of detail (coarse resolution) and the aggregation of several actors into single agents, inherent at large scales, results in a more homogenous set of agent types. This reduces the added value of the agent based approach (Evans and Kelley, 2004) and casts doubts about the validity of behavioural assumptions at this scale level (Verburg et al., 2004).

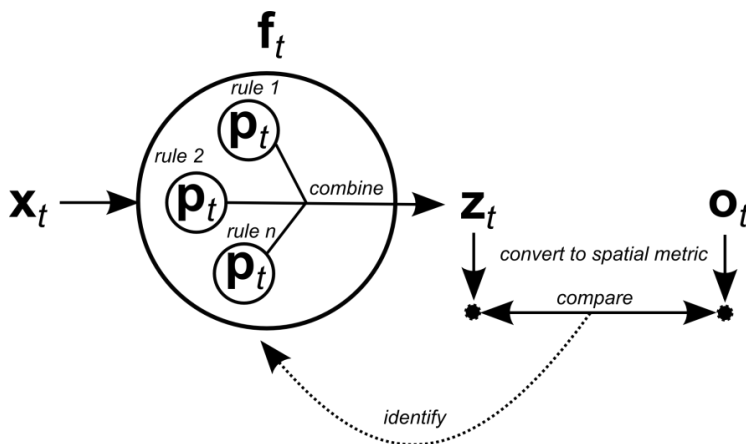
In a cellular automaton (CA) modelling approach land use change is simulated through local interactions in a space that consists of a set of discrete cells, a raster (Wolfram, 1984). Each cell has a state, representing the land use type (for example forest or cropland) present in that cell at that point in time. The advantage of a CA is that it is composed of relatively simple rules that can lead to complex spatial patterns, as observed in land use change (Santé et al., 2010). Cellular automata can be divided into two groups: pure CAs and constrained CAs. Pure CAs, e.g. SLEUTH (Chaudhuri and Clarke, 2013), determine the total area of land use change endogenously, by a bottom-up approach, based on historic data. In constrained CAs the total amount of change for the study area is determined exogenously, as an input or model boundary condition (White and Engelen, 2000). A model with an exogenous demand is considered more suitable for the simulation of human-induced land use changes, which are often steered by drivers outside of the land use system, demanding for the activity taking place or the commodity being produced on the land (Verburg et al., 2004). Models that constrain the amount of change are also called demand driven land use change models. Well-known examples of demand driven models are LandSHIFT (Schaldach et al., 2011) and models of the CLUE family (e.g. Verburg et al., 2002, Verburg and Overmars, 2009). Having the demand as an exogenous variable also allows for what-if questions regarding demand changes, e.g., what will be the response of the land use system if developing countries switch to a more meat-based diet like the developed countries? For this reason, demand driven, or constrained, CAs are more informative than pure CAs in the sustainability debate.

### **1.3. Model components**

All simulation models consist of different model components. The inputs of the model form the first component. These are in a demand driven land use change model the demands for all land use types over a particular period of time, and a number of spatial attributes that influence the locations of land use change, for instance potential crop yield. The set of inputs for time  $t$  is denoted as  $\mathbf{x}_t$  (Figure 1.1). The second component is the system state at time  $t$ , denoted as  $\mathbf{z}_t$  (Figure 1.1), for instance land use. The system state of time  $t - 1$  (the previous time step), is also one of the inputs to the calculation of the system state at time  $t$ , because, in almost all systems and certainly in all land use systems, the current system state depends on the past system state. The core of a simulation model is the set of transition rules that represent the processes that lead to the change in the system state over time. In the case of a land use change model a transition rule is a function calculating the suitability of each location for a particular land use type, with respect to a spatial attribute that influences the allocation of that land use

type (an input). The third component, encompassing the selection and combination technique of the transition rules, is called the model structure  $f_t$  (Figure 1.1). The transition rules can contain parameters  $p_t$  defining the characteristics of the process represented by the transition rule. A land use change model can have, for example, a parameter that specifies whether the relation between the attribute values and the suitability values is linear or exponential. The parameters are the fourth model component.

If a land use model is to be used to explore future land use pathways, it is essential to find the model structure and parameter values that result in an optimal model representation of the studied land use system. Identification of the model structure and parameter values can be accomplished through comparison of the modelled system state and observations of the real system  $o_t$  for the same time step  $t$  (Figure 1.1) and subsequently selecting the parameter values and model structure that minimize the difference between modelled and observed land use. Although not explicitly part of the simulation model, the observations of the real system can be considered the fifth and final model component, because they do contribute to the model formulation. The observations can be land use maps, but also aggregated measures or spatial metrics of configuration (a measure of spatial patterns) or composition (a measure of the proportions of different classes) (Csillag and Boots, 2005).



**Figure 1.1: Conceptual model of a simulation model:  $f_t$  represents the processes of change in the system state over time, i.e. the set of transition rules and the way to combine them,  $x_t$  represents all inputs, and  $p_t$  contains the parameters. Identification of the parameters and model structure is based on a comparison between the simulated system state  $z_t$ , or a derived spatial metric, with the observed system state  $o_t$ , or a derived spatial metric.**

## 1.4. Uncertainty

Constructing a land use change model is not straightforward. The dynamic processes and feedback loops in the land use system are complex and only partially understood (Manson, 2007, Verburg et al., 2013). Therefore, the model structure is by definition a strong simplification of this system. In addition, the model inputs and parameter values are uncertain and the observations of the real system contain errors. For example, land use maps are usually created from remotely sensed data, through classification, which is prone to errors (e.g. Burnicki, 2011). Also, land use areas from agricultural statistics tend to be rough estimates and can have large gaps (missing data values) that have to be interpolated (Lokupitiya et al., 2007). Furthermore, there can be discrepancies between the data and the desired model setup in e.g. land use classes or spatial resolution, and bridging these discrepancies can add additional errors (Schroeder, 2007). These errors propagate through the model, which generates uncertainty in the modelled system state, i.e. in the land use projections. The land use system behaves non-linearly through time, and therefore the uncertainty of the projections is likely to be non-linear too (Manson, 2007). Thus, propagating uncertainty should be calculated iteratively.

Yet, the uncertainty resulting from errors propagating through the different model components is not the only uncertainty associated with the land use projections. Another factor, that could also be termed an uncertainty, is related to the practice of scenario analysis. Namely, the analysis of chosen scenarios merely provides information on land use changes and the related impacts of a specific set of future story lines. It gives no information on how these different scenario outputs relate to the complete set of potential futures (Seppelt et al., 2013). For example, in an assessment of the environmental impacts through land use change of a certain increase in demand for a crop, one could use a land use change model to evaluate two scenarios. Presumably, one scenario will result in lower environmental impacts than the other one, but this approach does not ascertain that these are the lowest obtainable impacts. A calculation of how far off these 'low' environmental impacts are from the lowest obtainable impacts would make it possible to better place the scenario results into perspective, even if the lowest possible environmental impacts are not reachable given the available policy instruments. In this thesis, this uncertainty about the complete set of potential futures is defined as 'solution-space uncertainty'.

The uncertainty in land use change model projections is rarely communicated to the end users, which is problematic given that a land use projection that is erroneous or has not been put into perspective can result in wrong conclusions (Pontius Jr. and Spencer, 2005, Moulds et al., 2015). One reason that uncertainty is not communicated is that the land use change models currently in use do not have accompanying tools for uncertainty quantification. To our knowledge, only the

pure CA SLEUTH includes tools to allow probabilistic estimates of land use change projections (e.g. Clarke et al., 1997). Most land use change models are implemented in compiled programming languages of which the source code is not made available (Rosa et al., 2014), making it difficult for modellers to add iterative uncertainty analyses tools to existing models.

Another reason that uncertainty is not communicated, is that end users of land use change models, or SDSSs in general, often fail to appreciate uncertainty (Foody, 2003). They consider the fact that uncertainty can play a role in decision making a reason not to deal with it, regarding it as a risk or trouble, while quantified uncertainty is actually useful additional information about the robustness and trustworthiness of the projections and related impacts (Aerts et al., 2003b). In addition, information about the uncertainty of the different model components signifies which model components have the prospects for and most impact on uncertainty reduction, i.e. making model projections more trustworthy. We, as scientists, should not take the end users' anxiety towards uncertainty as an excuse not to calculate and communicate the uncertainty in land use change projections. It is our responsibility to clarify that this uncertainty does not undermine science, but is an intrinsic part of our (mis)understanding of the modelled processes and (un)comprehensiveness of the data for the given case study, and that this information can be used to the end user's benefit (Ivanovic and Freer, 2009). Therefore, we claim, alongside others (Ivanovic and Freer, 2009, e.g. Uusitalo et al., 2015, O'Hagan, 2012, Bastin et al., 2013), that uncertainty in simulation model outputs should be quantified and communicated to the end users.

Although other approaches are available for quantifying uncertainty in simulation models (Uusitalo et al., 2015), e.g. fuzzy logic (Nguyen et al., 2007), stochastic modelling is chosen in this thesis, because of its ease of implementation and strong roots in mathematics (Aerts et al., 2003b). Stochastic modelling involves defining probability distributions for each variable that is considered to be uncertain and using random sampling to draw values from these probability distributions. The uncertain variables defined by probability distributions in the stochastic model can be any of the model components, inputs, parameters, model structure or observations, except for the system state itself, as that is the component we want to calculate the uncertainty for. The uncertainty is propagated towards the output of the simulation model, in our case the land use projections. Monte Carlo (MC) simulation is a solution scheme that involves running the stochastic model a large number of times, each time with a random value from each of the probability distributions of the stochastic variables. This results in different land use projections for these model realizations. Over the different projections summary statistics or spatial metrics can be calculated, quantifying the range of the total

ensemble and thereby the uncertainty per time step. This uncertainty is visualized in the SDSS for the end users.

Once output uncertainty is known, methods can be applied to reduce this uncertainty. As there are currently no stochastic demand driven land use change models, there are also no methods for uncertainty reduction. There are, however, methods to find the transition rules, or model structure. Regression on a land use map is the most commonly applied model structure identification method (Verburg et al., 2002, Verburg et al., 1999, Aguiar et al., 2007, Diogo et al., 2014). In recent years, methods originating from artificial intelligence have become more prevalent, like neural networks (Dai et al., 2005, Li and Yeh, 2002) and swarm intelligence algorithms (Liu et al., 2008, Feng et al., 2011).

There are two problems in the current model structure identification methods. Firstly, the uncertainty in the observational data, from e.g. classification of remote sensing images and discretization, of the system that the model should reproduce is not taken into account. Ignoring uncertainty in the empirical data may lead to an underestimation of model output uncertainty. Secondly, the model performance should be optimized based on the modelling aim, on the attribute one wants to project. For example, if the effect of future land use change on animal passageways is studied, connectivity of patches is an important characteristic of the land use projections, but if GHG emissions of land use change are studied it is more important in which climate zone the change occurs and which previous land use type is replaced, among other things. Other studies often fail to indicate which target is used to calibrate on (e.g. Aguiar et al., 2007, van Vliet et al., 2012), or calibrate on a cell by cell comparison statistic, e.g. Kappa (e.g. van Delden et al., 2010, Mancosu et al., 2015) or revised Kappa (e.g. van Vliet et al., 2011, Blečić et al., 2015). We do not consider a cell-by-cell comparison a good practice as the aim of land use change models is usually not, and should not be, to simulate precisely the land use of each single cell in each year (Parker et al., 2008). More realistic is to try to capture the spatio-temporal patterns relevant to the modelling aim. Therefore, there is a need for a coupled land use change modelling and calibration framework, to allow recalibration of the model when an end user needs it for a new purpose or case study.

## **1.5. Bioenergy**

There is a myriad of end users of land use change models, connected to the wide range of drivers and environmental impacts of land use change. Extent-wise, the most important land conversion is natural vegetation to agricultural land, i.e. cropland and pastoral land (Lambin and Meyfroidt, 2011). An important question in



the sustainability debate is to which extent this conversion is currently induced by a demand for bioenergy and how this will change in the future, as e.g. the global biofuel demand is projected to grow from its current 3.0% of the road transport fuel demand to 3.5% by 2020 (OECD/IEA, 2014)

In this thesis, the developed methods to quantify and reduce uncertainty for spatially explicit land use change models are illustrated with case studies about the implications of increasing demands for bioenergy crops. Bioenergy crops are, compared to food, feed and fibre crops, under extra pressure in the sustainability debate, because they are targeted to be a sustainable alternative for fossil energy. While bioenergy crops are more sustainable than fossil sources in the sense that they are renewable and that they offer energy security for countries without (or with small or expensive to extract) oil and gas reserves, there are other aspects in which bioenergy crops are not by definition sustainable. One of the most discussed aspects herein is the total amount of GHG emissions from bioenergy production when land use change effects are taken into account, about which a heated debate started seven years ago with two papers, from Searchinger et al. (2008) and Fargione et al. (2008). GHG emission estimates of bioenergy production vary widely because there is no consensus about the carbon accounting approach (O'Brien et al., 2015), but also because carbon stocks vary widely over space, so the carbon stock change is dependent on where exactly the bioenergy crop will be cultivated (e.g. van der Hilst et al., 2014). For this latter problem, the uncertainty assessment of land use change projections would be of help. That is, if the range of possible land use change patterns is known, this range can be used to assess the range in carbon stock changes. This ranges given more information than a single value, of which the accuracy is unknown.

According to the IPCC (2011), the global technical potential for bioenergy from biomass for 2050 is 500 EJ per year, of which 70 EJ per year could be supplied from marginal or degraded land (500 Mha) and 120 EJ per year from good quality agricultural and pastoral land (300 Mha). However, the future demand for and the potential supply of bioenergy depend on, among other things, the implementation and enforcement of energy and sustainability policies, feedstock choice, improvements in agricultural management, and prices of fossil resources (e.g. Dornburg et al., 2010). Because of this, spatially explicit demand driven land use change models are particularly applicable to bioenergy case studies, since they allow for the assessment of scenarios with contrasting assumptions for these key variables. Spatially explicit land use change models have already been used to answer bioenergy related questions (van der Hilst et al., 2014, e.g. Lapola et al., 2010, Hellmann and Verburg, 2011). Recently, models have been especially of use in the projection of indirect land use change (iLUC), i.e. the cascading effect of a change in land use outside the bioenergy feedstock cultivation area, induced by a

change in use or production quantity of that feedstock (see Wicke et al., 2012 for a review). But only one of these model assessments of iLUC was done using a spatially explicit land use change model (Lapola et al., 2010). More importantly, although there is a huge divergence of estimated GHG effects of iLUC between the different studies, from 200% below up to 1700% above the carbon footprint of fossil fuels (Finkbeiner, 2014), evaluation of the uncertainties in the projections of one model, resulting from errors propagating through the different model components within this model, has not been performed. A few studies have assessed the solution space of particular bioenergy case studies, but all of them are spatially aggregated (e.g. Lautenbach et al., 2013, Akgul et al., 2012) or use very small case study areas, which limits the general applicability of the conclusions (Chikumbo et al., 2015). Also, none of these assesses the position of one or more scenario projections within this solution space.

## **1.6. Problem definition, research questions and thesis outline**

From the above, it can be concluded that it is essential to be able to spatially project land use changes through time. Yet, the spatially explicit land use change models that are developed to do so, contain uncertainties in their different model components (inputs, model structure, parameters and empirical observations) and in the position of their scenario projections in the total solution space. Commonly used land use change models do not quantify these uncertainties. This leaves the end user unaware of the accuracy of the modelled land use projections and uninformed about the sources of uncertainty and possibilities for uncertainty reduction. The aim of this thesis is to develop methods to quantify and reduce uncertainty for spatially explicit land use change models. Case studies in this thesis are focused on bioenergy as a driver of land use change, because the demand for bioenergy is expected to increase in the near future and the impacts of this increase are an important issue in the sustainability debate. The following research questions are assessed in the chapters of this thesis:

### *1. How can uncertainty in land use change projections be quantified?*

In Chapter 2, Monte Carlo simulation is applied to quantify the uncertainty in land use change projections. Hereto uncertainties in all model components have to be estimated: inputs (Chapter 2, 5), parameters (Chapter 2, 3, 4), model structure (Chapters 3, 4, 5) and observations (Chapters 3, 4, 5). The methods to estimate these uncertainties and the effects of these uncertainties on different model outputs at various scales are discussed. Chapter 2 focuses on the uncertainty in potential bioenergy crop yield and on uncertainty visualization techniques for

SDSSs. Chapter 3 quantifies uncertainties in different spatial landscape metrics that are of relevance for land use change impacts, such as the number of interconnected patches of a land use type, which is of relevance to e.g. biodiversity. Chapter 4 demonstrates the effect of a non-stationary model structure on the total area of a bioenergy crop in different regions. Chapter 5 focuses on the uncertainty in land use change at different scale levels. In Chapter 6, a scenario projection of land use is compared with the 'optimal' land use configurations, given the objectives to minimize biofuel production costs and GHG emissions. This provides information on the solution-space uncertainty.

### *2. How can uncertainty in land use change projections be reduced?*

Monte Carlo simulation is also used in Chapter 3 and 5, but now coupled to a particle filter, a data assimilation technique, to reduce the uncertainty in the model with observations of land use change over time. Errors in observations used to reduce the uncertainty are taken into account in both chapters.

### *3. What are the contributions of different model components to the uncertainty in land use change projections?*

All chapters take into account uncertainty in different model components: inputs (Chapter 2, 5), parameters (Chapter 2, 3, 4), model structure (Chapters 3, 4, 5) and observations (Chapters 3, 4, 5). We aim to evaluate how large the influence of each of these components on the uncertainty in projections is. This is evaluated for various model output attributes at different spatial scale levels. This information indicates which components have the highest need for improvement. Future data collection and model development efforts can be tailored to this.

### *4. What are the implications of the land use change projection uncertainties for bioenergy implementation strategies?*

The impacts of the findings of this thesis for bioenergy implementation strategies, like policy formulations, are discussed in all chapters for different bioenergy case studies. Chapter 2 explores the potential for sugar cane and eucalyptus for bioenergy in Mozambique up to 2030. Chapter 3 and 4 focus on sugar cane expansion in the state of São Paulo, in Brazil. Chapter 5 assesses the whole of Brazil, taking into account the expansion of sugar cane but the iLUC effects. Finally, Chapter 6 is focused on the allocation of sugar cane fields and ethanol mills in the state Goiás, one of the new expansion regions of sugar cane in Brazil.

Within the scope of this thesis the PCRaster Land Use Change (PLUC) model was developed. PLUC is a demand driven CA, similar to other land use change models such as CLUE (Verburg and Overmars, 2009, Verburg et al., 1999) and LandSHIFT (Schaldach et al., 2011). It is built within the PCRaster Python framework (Karssenberget al., 2010) that facilitates integration of spatio-temporal modelling and uncertainty analysis. Hereto, it uses the PCRaster Python library (Karssenberget al., 2007) that contains spatio-temporal functions for rasters. PLUC can be tailored to any case study and any region as the user can define the number of land use types to be modelled and transition functions valid for these land use types.

## 2. Spatio-temporal uncertainty in Spatial Decision Support Systems: a case study of changing land availability for bioenergy crops in Mozambique

**Judith A. Versteegen, Derek Karssenbergh, Floor van der Hilst, André P.C. Faaij (2012), Computers, Environment and Urban Systems 36, 30-42.**

**Abstract** - Spatial Decision Support Systems (SDSSs) often include models that can be used to assess the impact of possible decisions. These models usually simulate complex spatio-temporal phenomena, with input variables and parameters that are often hard to measure. The resulting model uncertainty is, however, rarely communicated to the user, so that current SDSSs yield clear, but therefore sometimes deceptively precise outputs. Inclusion of uncertainty in SDSSs requires modelling methods to calculate uncertainty and tools to visualize indicators of uncertainty that can be understood by its users, having mostly limited knowledge of spatial statistics. This research makes an important step towards a solution of this issue. It illustrates the construction of the PCRaster Land Use Change model (PLUC) that integrates simulation, uncertainty analysis and visualization. It uses the PCRaster Python framework, which comprises both a spatio-temporal modelling framework and a Monte Carlo analysis framework that together produce stochastic maps, which can be visualized with the Aguila software, included in the PCRaster Python distribution package. This is illustrated by a case study for Mozambique in which it is evaluated where bioenergy crops can be cultivated without endangering nature areas and food production now and in the near future, when population and food intake per capita will increase and thus arable land and pasture areas are likely to expand. It is shown how the uncertainty of the input variables and model parameters effects the model outcomes. Evaluation of spatio-temporal uncertainty patterns has provided new insights in the modelled land use system about, e.g., the shape of concentric rings around cities. In addition, the visualization modes give uncertainty information in an comprehensible way for users without specialist knowledge of statistics, for example by means of confidence intervals for potential bioenergy crop yields. The coupling of spatio-temporal uncertainty analysis to the simulation model is considered a major step forward in the exposure of uncertainty in SDSSs.

## 2.1. Introduction

Spatial Decision Support Systems (SDSSs) are interactive, computer-based systems that include simulation models and visualization tools designed to assess the impact of possible decisions (Geertman and Stillwell, 2004). With SDSSs, planners can investigate the effects of different scenarios, and explore intervention possibilities by adjusting model inputs in a user interface accessible to users that do not have expert knowledge of modelling theory and technology. The models in SDSSs usually simulate change over time of a spatial phenomenon or, more likely, a number of spatial phenomena that interact with each other. These dynamic processes and interactions tend to be complex and are rarely fully understood (Manson, 2007). As simulation models are simplifications of open, complex systems, model output errors are inherent as a result of the debatable choice and conceptualization of relevant sub-processes, uncertainty in model parameters and input variables, and the discretization of information. These errors propagate through the model because the state of the modelled system at a certain moment in time is a function of its state in the past. This generates uncertainty in model outputs.

Whereas scientists are familiar with the concept of uncertainty and methods to quantify it (Brown and Heuvelink, 2007, Chen et al., 2011, Goodchild, 2004, Heuvelink, 1998), SDSS users tend to seek certainty and deterministic solutions (Bradshaw and Borchers, 2000). In other words, they demand practical models with clear and unambiguous results to facilitate decision making. The result of this desire for clarity and simplicity is that SDSSs tend to underestimate, if not ignore, uncertainty (Foody, 2003). They thus yield clear, but therefore sometimes deceptively precise outputs.

Decisions based on misinterpreted or erroneous model output can be costly due to the irreversibility of such decisions. Uncertainty thus needs to be communicated clearly. Therefore we claim, together with others (Foody, 2003, e.g. Aerts et al., 2003b, Ivanovic and Freer, 2009, Ma et al., 2007, Oreskes et al., 1994), that instead of obscuring uncertainty users should be made more aware of uncertainty in SDSSs. Difficulties in including and communicating uncertainty in SDSSs are: 1) the concepts and measures of uncertainty can be somewhat difficult to grasp for users without specialist knowledge of statistics, 2) uncertainty is input-dependent, 3) uncertainty varies over space, time and aggregation level, and 4) software packages that can integrate spatio-temporal modelling, uncertainty analysis and visualization are rare.

The first difficulty arises from the fact that it is unlikely that the average user of an SDSS has skills in handling uncertainty at a comparable level as the modellers themselves (Foody, 2003). Therefore, modellers should aid their end users by

providing intuitive and insightful indicators of uncertainty and straightforward tools to visualize these (Aerts et al., 2003b).

The second problem is that the model uncertainty cannot be calculated on forehand by the modeller, as output uncertainty depends on model inputs and parameters. Currently, uncertainty in simulation models is sometimes assessed by providing a (static) map of output uncertainty of the final time step of a standard run of the simulation model (Brown et al., 2005, Chang et al., 2008, Eckhardt et al., 2003). Clearly, this does not suffice for an SDSS, in which investigating the effects different model settings is the main goal, so that inputs and parameters are altered frequently. For that reason, uncertainty analysis should be automatically calculated by the model itself, every time it is run with different settings.

The third difficulty is that uncertainty varies over space and time, because complex systems behave non-linearly. This makes that providing an uncertainty map of the final time step only is not sufficient when it comes to complex, non-linear systems (Ligmann-Zielinska and Sun, 2010). So, an uncertainty map is needed for each time step, in order to allow *iterative uncertainty analysis* (Manson, 2007), as for example demonstrated by Gorsevski et al. (2006). In addition, it has been shown by others (Pontius Jr. and Spencer, 2005, Hiemstra et al., 2012, Kok et al., 2001) that uncertainty is highly dependent on the level of spatial aggregation. Usually, uncertainty becomes lower at a coarser scale, because rearrangements in the landscape at locations in close proximity cancel each other out when they are aggregated. This is highly relevant in SDSSs, as their end users operate at different managerial levels (e.g., town, district, province, country), related to spatial scale levels (local, regional, national). To be able to aid these different end users, it should be possible in an SDSS to assess uncertainty at different levels of aggregation. Ideally, inclusion of such methods should be possible without too much additional work on the side of the model developer.

The final problem is that most software packages are either dedicated to model development, e.g., Stella (2010) and NetLogo (2010), or to uncertainty analysis (for a package overview see Goovaerts, 2010), or to visualization, e.g., ArcGIS (ESRI, 2011). Using such packages to construct an uncertainty-inclusive SDSS would require a complex coupling mechanism. Also, existing visualization tools do not explicitly support methods to visualize uncertain spatio-temporal data. A possible solution to this problem is the PCRaster model construction framework (Karszenberg et al., 2010, PCRaster, 2010), which offers a combined interface for spatio-temporal modelling and uncertainty analysis and includes a visualization tool for stochastic data in its distribution package.

The objective of this paper is to construct an SDSS that integrates simulation, iterative uncertainty analysis, and visualization to facilitate end users at different managerial levels to take uncertainty into account in decision making. This is

illustrated by a case study of bioenergy-crop potentials in Mozambique. Although some studies have been conducted to assess the area of potentially available land for bioenergy crops in Mozambique (Batidzirai et al., 2006, Watson, 2011), none of these studies was carried out in both a temporally dynamic and spatially explicit way. The PCRaster Land Use Change model (PLUC) is developed to evaluate where bioenergy crops can be cultivated without entering into competition with other important land uses from an economic or sustainability point of view, now and in the near future when population and food intake per capita and thus arable land is likely to increase. We show that PLUC allows stochastic model inputs and produces interactive visualizations of forecast uncertainty, in space and time, at a range of spatial aggregation levels. These visualizations can be used by decision makers to evaluate possible locations and potential yields for bioenergy crops.

The next section of this paper describes the concepts and methods of the PCRaster Python framework, outlines how the framework is applied to construct the land use change model for Mozambique, explains the error models of the different stochastic inputs, and illustrates the mode of implementation. The results section shows different visualization modes of uncertainty indicators and their potential usage, drawing on the outputs of PLUC for the Mozambique case study. The final section discusses the advantages and shortcomings of the uncertainty-inclusive simulation model.

## **2.2. Methodology**

### **2.2.1. Software framework**

Although other approaches exist to include uncertainty in a model, such as fuzzy logic (e.g. Nguyen et al., 2007, Robinson, 2003), stochastic modelling has the advantage that it has a strong root in mathematics. In stochastic modelling, a model input is defined by a probability distribution of all possible values. This uncertainty is propagated through the model using a numerical solution scheme. Monte Carlo simulation is such a solution scheme, which is attractive because of its general applicability and ease of implementation (Aerts et al., 2003b). It involves running the model a large number of times, each time drawing a realization from the input probability distribution(s). For spatial models this results in different spatial patterns for the different model realizations, i.e. model runs or samples.

The PCRaster model construction framework (Karssenberget al., 2010, PCRaster, 2010) facilitates this integration of spatio-temporal modelling and uncertainty analysis through the PCRaster Python library (Karssenberget al., 2007). This library provides a large set of spatio-temporal functions on raster maps, embedded in the Python language (Python software foundation, 2014). Both a spatio-temporal



modelling framework and a Monte Carlo analysis framework are present as a Python class. These classes include methods to write the simulation results and uncertainty indicators to disk as maps, which can be visualized with the Aguila software (Pebesma et al., 2007), included in the PCRaster Python distribution package. To allow construction of a spatio-temporal model that permits stochastic inputs and assessment of the resulting uncertainty, three main methods are provided by the framework that together form the schedule in Table 2.1 (Karsenberg and de Jong, 2006). Firstly, there is a loop for evaluation of the spatio-temporal process itself (line 2). Herein, the modeller can program the equations that represent the change of the system state over a time step (line 3). The framework provides for this purpose a number of functions particularly designed for spatial and stochastic operations. Secondly, a loop over this spatio-temporal model is performed to generate the Monte Carlo samples (line 1). Finally, summary statistics are computed over all Monte Carlo samples (line 4), representing the uncertainty in the model output within a time step or over the whole simulated period of time.

The individual Monte Carlo sample results and summary statistics are written to disk with a function from the PCRaster Python library. This function uses rules for file names defined by the modelling framework, so that they can directly be visualized with the Aguila software that recognizes these name conventions. Temporal deterministic and stochastic data can be viewed by animation or toggling through time, of which the last technique was considered helpful by decision makers (experts as well as novices) for visualizing uncertain data in a study of Aerts et al. (2003a). In addition, stochastic outputs are visualized as maps or plots of mean, standard deviation, confidence interval, exceedance probability, or cumulative probability distribution of a variable. These visualizations are interactive, which allows users to explore the data (Karsenberg et al., 2010, Pebesma et al., 2007).

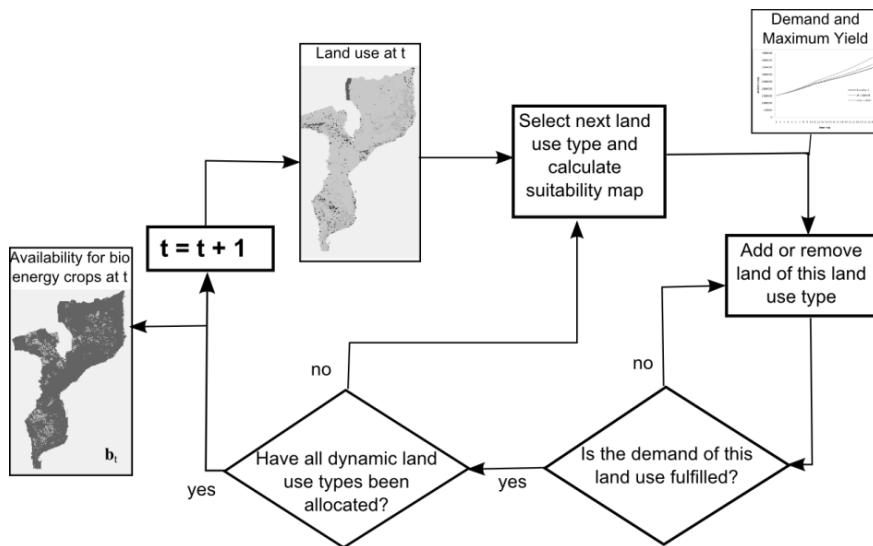
**Table 2.1: Modelling schedule of the PCRaster Python framework**

```
1 for each MC sample:
2   for each time step:
3     solve system state equation
4   compute summary statistics
```

### 2.2.2. Land use change model

The potential to employ the PCRaster Python framework for an SDSS including uncertainty is illustrated by the construction of a land use change model, meant to aid in evaluating where bioenergy crop plantations can be allocated. Until now, bioenergy potentials were mostly assessed in a spatially aggregated (e.g. Hoogwijk et al., 2005) or temporally static way (e.g. van der Hilst et al., 2010). So, in this study a both spatially explicit and temporally dynamic model is created for Mozambique. The case is relevant, because Mozambique is considered promising for bioenergy crop production by its vast amounts of available land (Smeets et al., 2007), favourable environmental conditions for cultivation (Batidzirai et al., 2006), and relatively low productivity of current agriculture (Arndt et al., 2010), which offers potential for improvement. However, over the past decade the area of forests and woodlands has decreased substantially due to an increase in cultivated areas (Jansen et al., 2008). Now as well as in the future, cultivation of bioenergy crops should not add to that effect and not endanger other important land uses, either from an economic, e.g., food crops and livestock, or from a sustainability point of view, e.g., conservation areas (Haberl et al., 2010). The population is expected to increase and its diet is expected to change as well, which induces further shifts in land use. The direction and extent of this shift depends on the agricultural and livestock productivity, which is expected to improve. The trends herein are derived from literature (e.g. FAO, 2003, INE, 2003), fieldwork, and meetings with national and local authorities. More background on the modelled processes and used data is provided in the twin-publication of this paper (van der Hilst et al., 2012). This section focuses on the set-up of the model.

The main model component is the state transition function representing the change in spatially distributed land use over a time step. Many models of land use dynamics have been constructed before (see for overviews Agarwal et al., 2002, Parker et al., 2003, Verburg et al., 2004). Some focus on one specific land use conversion, such as urbanization (Batty, 2005, Ligtenberg et al., 2009), but more often several land uses in the area compete for new locations (Verburg and Overmars, 2009, Lei et al., 2005). We adopt the latter approach, in which we focus on active change in agricultural land use types and forest as the future distribution of land reserved for these land uses is the main issue in the view of potential locations for bioenergy crop plantations. By *active change*, we mean that expansion or contraction of the total area of this land use is explicitly steered by certain drivers. Other land use types on the land use map can change passively, by expansion or contraction of an active land use type.



**Figure 2.1: Conceptual model of land use change.**

The land use change is steered by two factors: 1) the demand of the population for food, non-food crops (e.g., cotton and tobacco) and wood, and 2) the growth rate of yield, defined by agricultural and livestock productivity. The conceptual model (Figure 2.1) includes three loops. The first one loops over time and coincides with line 2 in Table 2.1. The other two belong to the actual state transition function (line 3 in Table 2.1); one loops over all active land use types and the other checks iteratively whether this land use type should expand, contract or has met its demand. The actual location of the expansion or contraction of the land use types is determined by suitability factors, like distance to cities and transport networks, current land use in the neighbourhood, and location-specific yield due to characteristics of the soil and climate. Areas occupied by other economically important land uses, physically constrained areas, and protected land uses are excluded and the remaining land is potentially available for bioenergy crops. But the bioenergy crops are not included as a land use type, which means they are not allocated. The model is explained in more detail in the following.

An important step in model implementation is the choice of the support to represent processes (Hengl, 2006, Pan et al., 2010, Bierkens et al., 2000). The support refers to the size in the spatial domain (i.e. spatial discretisation), and the temporal domain (i.e. time step duration) over which processes are considered homogeneous. Model inputs and parameters need to be representative for the support used (c.f. Bierkens et al., 2000). The following criteria were used to select a suitable support. 1) The scale at which the model output is required in order to answer the end user's questions (Evans and Kelley, 2004), also called *policy scale* (Bierkens et al., 2000). In our case, policy makers at different spatial scale levels are

involved, ranging from national to local scale. Information is needed on changes in land use over the coming decades. 2) The *process scale*, i.e. the scale of natural variability of the studied process (Blöschl, 1999). We do not aim to study land use changes at parcel level; the smallest transformation of interest is at the level of a small community of farmers. 3) *Observational data availability*. If the input data does not match the resolution of the model, upscaling or downscaling could be applied, but one should be careful to do so, as spatio-temporal probability distributions alter with changing resolution (Bierkens et al., 2000). Thus it is preferable to choose a support that matches the support of available observational data. 4) *Processing power*, as calculation time increases exponentially with the number of cells (Hengl, 2006). As PLUC is designed as an SDSS, it should be usable on a desktop pc. The Monte Carlo simulations should be performable in a reasonable amount of time and should result in a manageable amount of data. As the system studied constitutes of a number of subsystems that operate at different scales, while data availability differs between subsystems, we implement our models at three levels of spatial support: 1) country level, for data that is only available for the country as a whole, like economic trends, 2) land use type level, for crop-dependent but location-independent variables, like maximum possible product yield, and 3) cell level, for location-specific information, like population density. These three levels are implemented using gridded maps with a cell size of 1 km<sup>2</sup>. This means that technically, all information is discretized into cells of 1 km<sup>2</sup>, but in fact the three abovementioned levels of spatial support are used in input data and process descriptions. The model uses a time step ( $\Delta t$ ) of one year, with time step  $t = 1, 2, \dots, T$ . Most equations are evaluated separately for each of the  $N$  land use types, with  $n = 1, 2, \dots, N$ . The model is run for 25 years and the following dynamic land use types are defined: cropland, mosaic cropland-pasture (grazed grassland), mosaic cropland-grassland (not grazed), pasture, forest. However, any cell size, model period, and number and type of land uses could be defined by the user.

The demand  $d_{n,t}$  (kg · year<sup>-1</sup>) for products from land use type  $n$  at time step  $t$  is:

$$d_{n,t} = p_t * i_{n,t} * r_t \quad , \text{ for each } n \text{ in each } t \quad 2.1$$

In Equation 2.1,  $p_t$  denotes the number of inhabitants in the country at  $t$ , intake  $i_{n,t}$  (kg · caput<sup>-1</sup> · year<sup>-1</sup>) specifies the demand per capita of products from land use type  $n$  at  $t$ , and the self-sufficiency ratio  $r_t(-)$  is the extent to which the food demands are met by domestic supply at  $t$ .

Potential yield  $\mathbf{p}_{n,t}$  ( $\text{kg} \cdot \text{km}^{-2} \cdot \text{year}^{-1}$ ) is the yield of products from land use type  $n$  at  $t$  if the cell would be occupied by that land use type:

$$\mathbf{p}_{n,t} = m_{n,t} * \mathbf{f}_n, \text{ for each } n \text{ in each } t \quad 2.2$$

In Equation 2.2,  $m_{n,t}$  is the maximum possible product yield ( $\text{kg} \cdot \text{km}^{-2} \cdot \text{year}^{-1}$ ) of products from land use type  $n$  at  $t$ , which can increase through time due to technological improvements in the agricultural or livestock sector, i.e. increased productivity (FAO, 2003). The location-specific variable  $\mathbf{f}_n \in [0,1]$  is the actual fraction of this yield that can be reached in a cell, depending on factors like soil type, climate, and water availability. Note that a bold font indicates that a variable is a spatial field, i.e. information at cell level.

The total product yield of a land use type is calculated using the land use map at  $t$ . First, a spatial field of the current yield  $\mathbf{c}_{n,t}$  ( $\text{kg} \cdot \text{km}^{-2} \cdot \text{year}^{-1}$ ) of land use type  $n$  at  $t$  is constructed that contains the value of  $\mathbf{p}_{n,t}$  for cells that are currently occupied by type  $n$  at  $t$  and zero for cells occupied by other land use types. The total product yield  $y_{n,t}$  (kg) of this land use type is therefore:

$$y_{n,t} = \sum(\mathbf{c}_{n,t} * a), \text{ for each } n \text{ in each } t \quad 2.3$$

In Equation 2.3, the summation sign indicates summation over the whole spatial field and  $a$  is the cell size in  $\text{km}^2$  ( $1 \text{ km}^2$  in this study).

To determine where a certain land use expands or contracts, every land use type is assigned a number of suitability factors. In total, nine suitability factors have been implemented in PLUC. The number and kind of suitability factors differ per land use type.

Table 2.2 shows which suitability factors were implemented for every dynamic land use type in the case study of Mozambique and which weights were assigned to them (see Equation 2.4 below). For a detailed explanation of why these factors and weights were chosen, the reader is referred to van der Hilst et al. (2012).

**Table 2.2: Weights per suitability factor  $i$  per land use type  $n$ . The suitability factors concern: spatial autocorrelation (1), distance to roads (2), water (3), and cities (4), yield (5), population density (6), livestock density (7), distance to plot edge (8), and conversion elasticity (9).**

land use type ( $n$ )	suitability factor ( $i$ )								
	1	2	3	4	5	6	7	8	9
cropland	0.20	0.10	0.10	0.10	0.20	0.20	0	0	0.10
cropland-grassland	0.20	0.10	0.10	0.10	0.20	0.20	0	0	0.10
cropland-pasture	0.20	0.10	0.10	0.10	0.15	0.15	0.10	0	0.10
pasture	0.30	0.05	0.15	0.05	0.1	0.05	0.20	0	0.10
forest	0.25	0	0	0.20	0.05	0.30	0	0.20	0

The spatial autocorrelation suitability factor (1) assumes that land uses attract land uses of the same or a related type (e.g., cropland is related to mosaic cropland-pasture and mosaic cropland-grassland), i.e. that related land uses tend to cluster. It determines the area of land use type  $n$  or types that are related to  $n$  in the neighbourhood of a cell with land use type  $n$ . The size of the neighbourhood is determined by the window length  $l$  (m), i.e. the length of the square window around the centre cell, and the area occupied by the same or related land use types is calculated relative to the size of the window. This suitability factor thus introduces positive feedback loop into the model: if a land use type is allocated at a certain location, the suitability value in the neighbourhood of that location increases for this land use type, so that more of this land use might be allocated in the area in the next time step.

The suitability factors for distance to roads (2), water (3), cities (4) and edge of the plot (8) determine suitability based on the shortest Euclidean distance to the object under consideration, i.e. roads, water, cities, or edges of plots, i.e. spatially connected areas of a uniform land use type. For crops and pasture typically a location close to roads and cities is preferred in order to minimize transport costs, and close to water for irrigation. Wood is preferably harvested at the edge of the forest, because this makes harvesting easier. The distance suitability functions (2, 3, 4 and 8) all have two input parameters: relation type between distance and suitability (linear, exponential or inversely proportional) and range  $r_n$  (m). The range indicates the maximum distance of effect of the feature, e.g., the maximum distance of effect of a road on the land use type pasture is set at 5 km, based on fieldwork, expert knowledge and literature (e.g. Jansen et al., 2008). Cells at a distance of more than the range have a suitability value of zero for this suitability factor.

The yield fraction (5), population density (6) and livestock density (7) suitability factors relate to the fraction of the maximum that can be found in a cell (see Equation 2.2), where the maximum refers to maximum yield ( $m_{n,t}$ ), maximum population density and maximum livestock density in factor 5, 6, and 7, respectively. Per land use type the direction (increasing or decreasing) and relation type between fraction and suitability (linear, exponential or inversely proportional) can be indicated, like in the distance suitability function, e.g., cropland is preferably located on cells with a high yield, but wood might preferably be harvested from cells with a low yield, i.e. biomass, because harvesting is easier in sparse forest.

The current land use (9) suitability factor indicates the compliance of a certain land use type to be transformed into the land use type that implements the suitability factor. This factor is sometimes referred to as conversion elasticity (Verburg and Overmars, 2009), e.g., it is more preferable for pasture to be placed on a cell that is currently defined as 'abandoned' than on a cell that is currently defined as

'cropland', as the second case involves a greater loss of economic value. Note that most spatial fields resulting from the suitability factors ( $\mathbf{u}_{i,n,t}$ , see Equation 2.4 below) remain the same over time, as the location of the features it relates to, e.g., roads, does not change over time in the model. This is the case for the factors 2, 3, 4, 5, 6, and 7. However, the spatial fields of suitability factors related to dynamic land use types (1, 8, and 9) do change over time and thus establish feedback loops in the land use system.

For every land use type a total suitability map  $\mathbf{s}_{n,t} \in [0,1]$ , indicating the aggregated appropriateness of a given location for land use  $n$  at time step  $t$ , is computed from its suitability factors:

$$\mathbf{s}_{n,t} = \sum_{i=1}^9 (w_{i,n} * \mathbf{u}_{i,n,t}) \quad , \text{ for each } n \text{ in each } t \quad 2.4$$

$$\text{with } \sum_{i=1}^9 (w_{i,n}) = \mathbf{1}$$

In Equation 2.4,  $\mathbf{u}_{i,n,t} \in [0,1]$  is the spatial field resulting from suitability factor  $i$  for land use type  $n$  at time step  $t$ , and  $w_{i,n} \in [0,1]$  is the weight of suitability factor  $i$  for land use type  $n$  (see Table 2.2).

Now, all information is available to allocate the land uses. The land use types have a certain hierarchy determined mainly by their economic importance. In Africa, the current trend is that agricultural intensification takes place on land that is now used as mosaic cropland-pasture or mosaic cropland-grassland. This extensive land use becomes less common and is moved to less fertile grounds, while areas even less fertile, like mountain areas and forests, come into use for grazing of livestock (Lambin et al., 2001). Therefore, we assume the following order of allocation for the dynamic land uses: cropland, mosaic cropland-pasture, mosaic cropland-grassland, pasture, forest. The allocation schema can be written in pseudo-code as in Table 2.3. When the land use type expands, it allocates new cells of this type at locations with the highest suitability, i.e.  $\max(\mathbf{s}_{n,t})$ , and when it contracts it removes cells of this type with the lowest suitability, i.e.  $\min(\mathbf{s}_{n,t})$ . Cells are converted to or removed from this land use until the total yield  $y_{n,t}$  equals the total demand  $d_{n,t}$ . In Table 2.3 the number 99 in line 7 refers to the land use category 'abandoned'. For land use type forest the class 'abandoned' is named 'deforested' in order to be able to distinguish cleared forest from deserted agricultural land on the resulting land use map.

**Table 2.2: Pseudo-code for land use allocation procedure. Allocation of each land use type  $n$  in each time step  $t$  proceeds until yield fulfills demand.**

```

1  if  $d_{n,t} > y_{n,t}$ :
2      while  $d_{n,t} > y_{n,t}$ :
3          convert cell with  $\max(\mathbf{s}_{n,t})$  to  $n$ 
4          update  $y_{n,t}$ 
5  else if  $d_{n,t} < y_{n,t}$ :
6      while  $d_{n,t} < y_{n,t}$ :
7          convert cell with  $\min(\mathbf{s}_{n,t})$  to 99
8          update  $y_{n,t}$ 
9  else:
10     do nothing

```

When allocation of one land use type is finished, allocation of the next type is performed, with the restriction that it cannot convert cells with a land use type that has already been allocated in that time step. Deforested areas become forest again when they are left fallow for 10 years. The regenerated forest can be harvested once more to fulfil the wood demand.

At the end of each time step, when the land use map has changed according to the demands of the different land use types, it is determined which cells are potentially available for bioenergy crops. This is done by excluding all areas occupied by crops, pasture, steep slopes (calculated from the digital elevation model), roads, water, cities, forest concession areas, community areas, and nature reservation areas. This results in a Boolean map ( $\mathbf{b}_t$ ) where cells are available (True) or unavailable (False) for bioenergy crops at time step  $t$ . Although the bioenergy crops are not allocated on the land use map, they do have their own maximum possible product yield  $\mathbf{m}_t$  and yield fraction map  $\mathbf{f}$ , so the bioenergy crop yield per location and in total for the available area  $\mathbf{b}_t$  can be calculated using Equations 2.2 and 2.3. These results can be used to assess the available area and yield on provincial and national level, in order to take the influence of spatial aggregation into account.

### 2.2.3. Error models

The various projections of population growth, diet change and technological improvements in the agricultural sector differ significantly (Arndt et al., 2010, FAO, 2003, UNDP, 2008), so applying one of these datasets deterministically to quantify drivers presumably ignores a large input error. Also, a number of model parameters are uncertain as they can only be estimated by expert knowledge, because extensive model calibration datasets are currently not available. PLUC takes each of these input errors into account in calculating the forecast



uncertainty. Model drivers and parameters that are uncertain are defined here as lumped or spatially distributed stochastic variables. The variables can be divided into three groups according to their data type: single value, spatial field and time series (Table 2.4).

Two stochastic variables represent a single value: the window length  $l$  used in suitability factor 1 and the range  $r_n$ , used in suitability factors 2, 3, 4, and 8. For window length  $l$  a normal error model is used:

$$l = \mu_l + Z_l * \sigma_l \tag{2.5}$$

with  $Z_l \sim N(0,1)$

In Equation 2.5,  $\mu_l$  is the mean of  $l$ , i.e. the value that would be used in a deterministic run. In our case study a value of 3 km is used for  $\mu_l$ , as this means that only the direct neighbours are taken into account in our 1 x 1 km raster. In Equation 2.5,  $\sigma_l$  is the standard deviation, which is set to 1 km. It should be noted that  $l$  can attain values such that the window cuts through cells. This is not a problem as the suitability factor calculates the area in the window occupied by attracting land use types, not the number of cells.

For the range  $r_n$  an error model with a uniform distribution is used:

$$r_n = \sqrt{a} + Z_r * 2 * \mu_{r,n} \tag{2.6}$$

with  $Z_r \sim U(0, 1)$

In Equation 2.6,  $a$  is the cell size in  $\text{km}^2$ . This means that realizations of the range vary between the cell length and twice the mean value that is used in a deterministic run. This lower limit is used, because land use cannot be allocated on the feature under consideration, e.g., a road.

**Table 2.4: Stochastic variables of the land use change model and their data type, error model and standard deviation ( $\sigma$ ), if applicable. For explanation of error models, see main text.**

stochastic variable	data type	error model	$\sigma$
window length ( $l$ )	single value	normal	3 km
range ( $r_n$ )	single value	uniform	-
elevation ( $h$ )	spatial field	normal	1 m
yield fraction ( $f_n$ )	spatial field	relative normal	0.2
population density	spatial field	relative normal	0.1
livestock density	spatial field	relative normal	0.1
maximum yield ( $m_{n,t}$ )	time series	relative normal	0.1
demand ( $d_{n,t}$ )	time series	uniform	-

Four stochastic input variables represent a spatial field, i.e. a raster map of values. For surface elevation  $\mathbf{h}$  a normal error is used:

$$\begin{aligned} \mathbf{h} &= \boldsymbol{\mu}_h + \mathbf{Z}_h * \sigma_h \\ \text{with } \mathbf{Z}_h &\sim \mathbf{N}(\mathbf{0}, \mathbf{1}) \end{aligned} \quad 2.7$$

In Equation 2.7,  $\boldsymbol{\mu}_h$  is the original elevation map and  $\sigma_h$  is the standard deviation, for which a value of 1 m is used. Note that the normal error  $\mathbf{Z}_h$  is a spatial field, which means that a separate value is drawn for each cell. The other three stochastic spatial fields are the yield fraction, population density and livestock density. For all three a relative normal error model is used. This means that the error (in this case a normal error) is higher for higher mean values. For example, the yield fraction  $\mathbf{f}_n$  is defined as:

$$\begin{aligned} \mathbf{f}_n &= \boldsymbol{\mu}_f + \mathbf{Z}_f * \sigma_f * \boldsymbol{\mu}_f \\ \text{with } \mathbf{Z}_f &\sim \mathbf{N}(\mathbf{0}, \mathbf{1}) \end{aligned} \quad 2.8$$

In Equation 2.8,  $\boldsymbol{\mu}_f$  is the original yield fraction map and  $\sigma_f$  is the standard deviation, for which a value of 0.2 is used. The population density and livestock density both have a standard deviation of 0.1. The standard deviation of the yield fraction is higher, because we found several different spatial data sets of the yield that were distinctively different, which indicates a large input error.

Finally, two stochastic time series are used. The first is the maximum yield  $m_{n,t}$ , to which a relative normal error model is assigned, similar to the one explained in Equation 2.8:

$$\begin{aligned} m_{n,t} &= \mu_{m,n,t} + Z_m * \sigma_m * \mu_{m,n,t} \\ \text{with } Z_m &\sim N(0,1) \end{aligned} \quad 2.9$$

In Equation 2.9,  $\mu_{m,n,t}$  is the mean of the of the maximum yield of land use type  $n$  at time step  $t$ , i.e. the expected value obtained from observation data or expert knowledge, and  $\sigma_m$  is the standard deviation, for which a value of 0.1 is used. The second time series is for the demand  $d_{n,t}$ . It uses an error model based on a uniform distribution between the upper and lower limit of the attribute:

$$\begin{aligned} d_{n,t} &= l_t + Z_d * (u_{n,t} - l_{n,t}) \\ \text{with } Z_d &\sim U(0,1) \end{aligned} \quad 2.10$$

In Equation 2.10,  $l_{n,t}$  and  $u_{n,t}$  are the lower limit and upper limit of the demand of land use type  $n$  at time step  $t$ . The limits of the demand  $d_{n,t}$  are determined by upper and lower limits of population  $p_t$ , intake  $i_{n,t}$  and self-sufficiency ratio  $r_t$  (Equation 2.1) predicted by the FAO (2003) and expert knowledge. Note that, although the resulting variables from Equation 2.9 and 2.10,  $m_{n,t}$  and  $d_{n,t}$ , change

over time, the stochastic variables  $Z_d$  and  $Z_m$  are drawn once, at the start of the simulation. These variables are used on all land use types, because they simulate the effect of the rate of increase in population, which cannot be different for the different land use types in the same model run.

#### 2.2.4. Implementation

The schedule of PLUC is given in Table 2.5. The PCRaster Python framework consists of four main methods that represent the scheme in Table 2.1: the `premcloop` (line 6 in Table 2.5) is evaluated only once, the `initial` (line 21) is evaluated once for each realization, the `dynamic` (line 33) is evaluated once for each time step in each realization, and the `postmcloop` (line 54) calculates the descriptive statistics over the Monte Carlo samples. Some variables have the prefix `self`, as they are defined as member variables to allow usage over the four different methods.

When the `LandUseChangeModel` is initiated it calls PCRaster Python's `DynamicModel` (line 3) and `MonteCarloModel` (line 4). Among other things, these allow retrieving time steps (line 64) and Monte Carlo samples (line 65). Next, the `premcloop` (line 6) imports input maps in lines 7-16. A separate file, defined by the model builder, `parameters.py`, defines all non-spatial inputs. This is done to prevent that the end user has to make changes in the main model scheme. All its variables and parameters are imported in lines 17-20.

The `initial` method (line 21) is used to define initial or temporally constant stochastic variables for which a realization is drawn for each Monte Carlo sample. An example of an initial variable is the environment, i.e. the land use map, initiated in line 22. Examples of temporally constant stochastic variables are  $Z_d$  and  $Z_m$  used in Equation 2.10 and 2.9 for the calculation of demand  $d_{n,t}$  and maximum yield  $m_{n,t}$ . For  $Z_d$  a value is drawn from a uniform distribution between 0 and 1 with the PCRaster Python `mapuniform()` function (lines 23). For the  $Z_m$  a value from a normal distribution, from the PCRaster Python `mapnormal()` function, is multiplied by the standard deviation defined by the user (line 24).

**Table 2.3. Main scheme of land use change model. Three dots and a discontinuity in the line numbering indicate omitted sections.**

```

1     class LandUseChangeModel(DynamicModel, MonteCarloModel):
2         def __init__(self):
3             DynamicModel.__init__(self)
4             MonteCarloModel.__init__(self)
5             setclone('landuse')

6     def premcloop(self):
7         self.initialEnvironment = self.readmap('landuse')
...
17        self.landUseList = parameters.getLandUseList()
...

21    def initial(self):
22        self.environment = self.initialEnvironment
23        self.demandStoch = mapuniform()
24        self.maxYieldStoch = mapnormal() * self.sdYield
25        self.landUse = LandUse(self.landUseList, self.environment)
26        self.landUse.drawRealizationsParams(self.stochParams)
27        self.landUse.createLandUseTypeObjects(self.relatedTypeDict
        , self.suitFactorDict, self.weightDict, self.varDict)
...

33    def dynamic(self):
34        demandUp = timeinputscalar('deUp.tss', self.environment)
35        demandLow = timeinputscalar('deLow.tss', self.environment)
36        demandDiff = (demandUp - demandLow)
37        demand = demandDiff * self.demandStoch + demandLow
...
42        self.landUse.calculateSuitabilityMaps()
43        self.landUse.allocate(maxYield, demand)
44        self.landUse.growForest()
45        self.environment = self.landUse.getEnvironment()
46        self.report(self.environment, 'landUse')
47        eu,euPr,euTo = self.landUse.getBioPotential(self.bioNoGo,\
        self.provinces)
...

56    def postmcloop(self):
57        name = ['eu', 'euPr', 'euTo']
58        mcaveragevariance(name, self.sampleNumbers(), \
        self.timeSteps())
59        name = ['eY', 'eYPr', 'eYTo']
60        percent = [0.05,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,0.95]
61        mcpercentiles(name, percent, self.sampleNumbers(), \
        self.timeSteps())

62    nrOfTimeSteps = parameters.getNrTimesteps()
63    nrOfSamples = parameters.getNrSamples()
64    myModel = LandUseChangeModel()
65    dynamicModel = DynamicFramework(myModel, nrOfTimeSteps)
66    mcModel = MonteCarloFramework(dynamicModel, nrOfSamples)
67    mcModel.run()

```

Next, the class `LandUse`, defined by the model builder, is instantiated (line 25). It is used to keep track of the changing land use map. This class has a method to calculate all other realizations for the stochastic variables in Table 2.4 (line 26). This `LandUse` class also instantiates  $N$  objects (i.e. one for each land use type) of the class `LandUseType` (line 27) that handle the land use type specific tasks, like computing suitability maps and allocating land (Equations 2.2, 2.3, 2.4 and Table 2.3). These methods are implemented with functions from the PCRaster Python library, including point, neighbourhood, and global operations (Burrough and McDonnell, 1998). In lines 28 to 32 other initial actions are taken, like computing a map of the distance to roads, needed by the `LandUseType` class for calculation of the suitability maps.

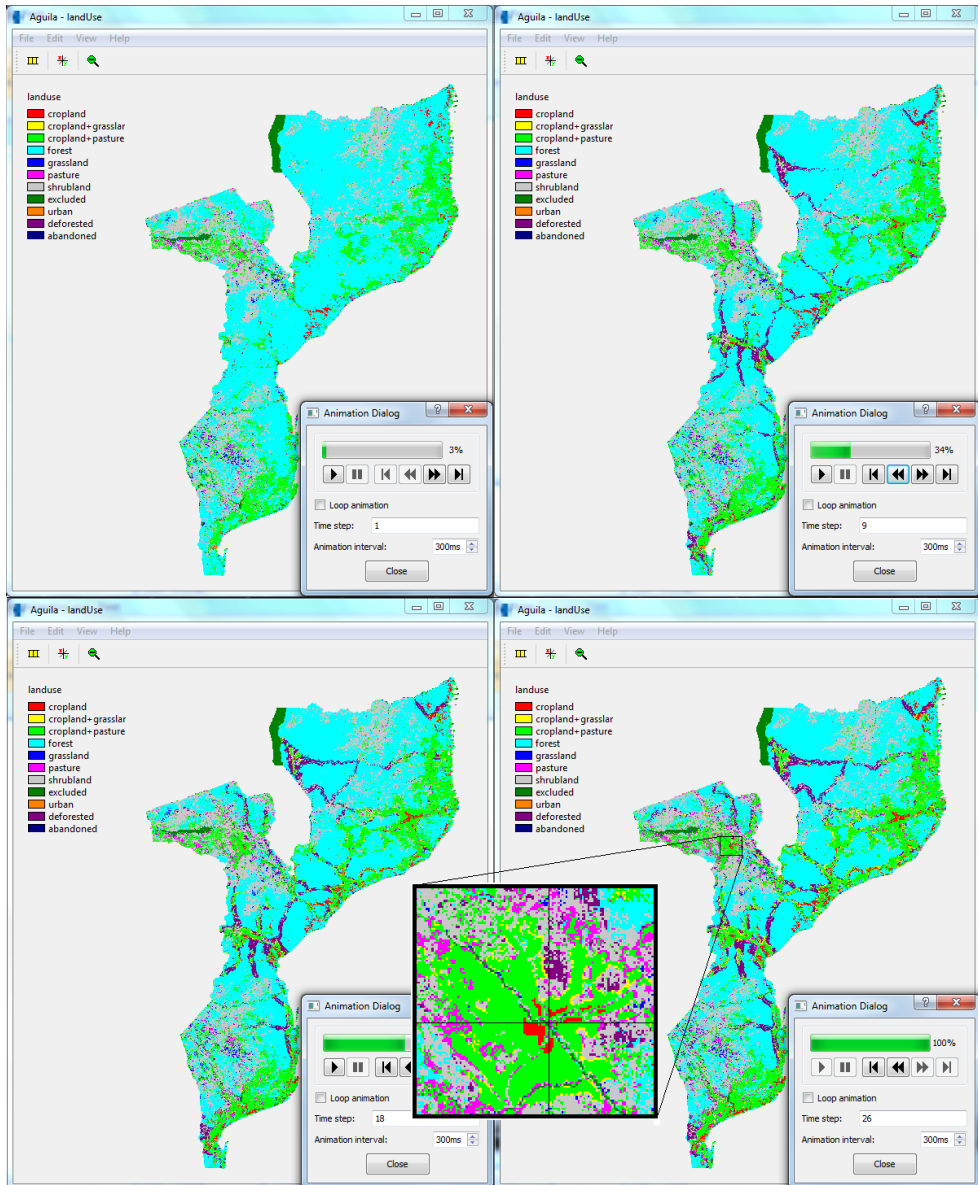
In the `dynamic` method (line 33), the temporal components of the model are evaluated. Demand is defined as a stochastic input variable by two time series per land use type (lines 34 and 35), as explained in the error model section. The variable `demandStoch`, drawn in line 23, is used as the position between the upper and lower bound,  $l_{n,t}$  and  $u_{n,t}$ , to calculate the realization for  $d_{n,t}$  (Equation 2.10) (lines 37 and 38). In lines 38-41 the maximum yield  $m_{n,t}$  is determined with the variable `maxYieldStoch` from line 24 (see Equation 2.9) in a similar way. Next, the suitability maps are calculated (line 42), and the allocation procedure, explained in Table 2.3, is called (line 43 in Table 2.5). The total land use map is updated (line 45) and saved to disk with the PCRaster Python function `report()` (line 46) that creates file extensions recognizable for the Aguila software. Next, it is determined which space is left over for bioenergy crops. In lines 47 and 48 it is determined where the bioenergy crop eucalyptus could be allocated and what its potential yield is on a cell-basis (`eu`), aggregated per province (`euPr`), and in total for the whole country (`euTo`).

In the `postmclloop`, the PCRaster Python function `mcaveragevariance` (line 58) calculates mean, variance and standard error of the files defined in line 57, for all time steps in all samples. The function `mcpercentiles` (line 61) computes the percentiles specified in the list in line 60 for a new list of files (line 59). All these estimators of uncertainty are automatically saved to disk. Note that these last two methods from the PCRaster Python library are the part of the model that calculate uncertainty. So, addition of just these two methods can turn any model into an uncertainty-inclusive model, given that it employs stochastic variables.

## 2.3. Results

The model was run in deterministic mode and in Monte Carlo mode using 500 samples to project land use change in Mozambique from 2005-2030. The Monte

Carlo run implements the error models defined in Table 2.4, while the deterministic run takes the mean of each of these variables. This section focuses on added value of the uncertainty analysis, an in-depth discussion of the simulated land use patterns and potential bioenergy crop areas is provided by van der Hilst et al. (2012).

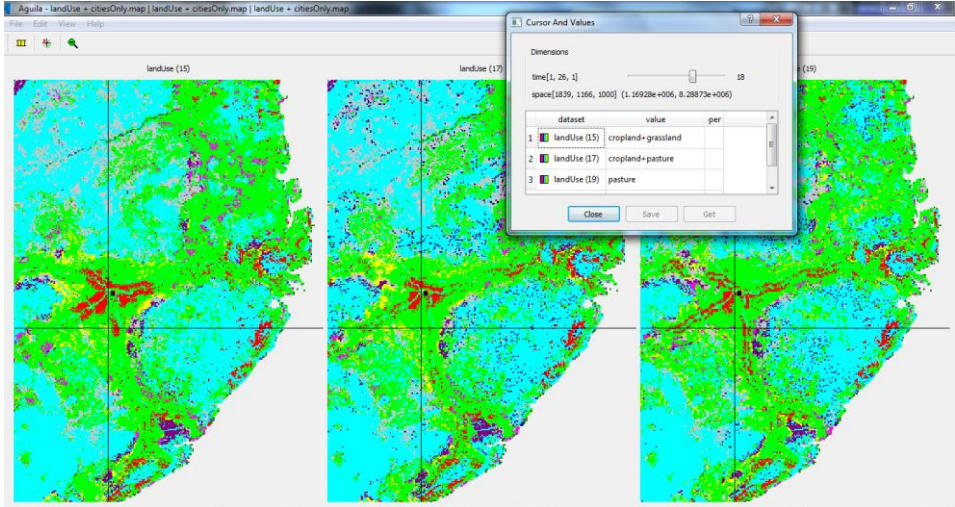


**Figure 2.2:** Screenshot of AguilA resulting from a deterministic model run showing land use in Mozambique in 2005 (time step 1, top left), 2013 (time step 9, top right), 2022 (time step 18, bottom left) and 2030 (time step 26, bottom right) with a close up around the city Tete.

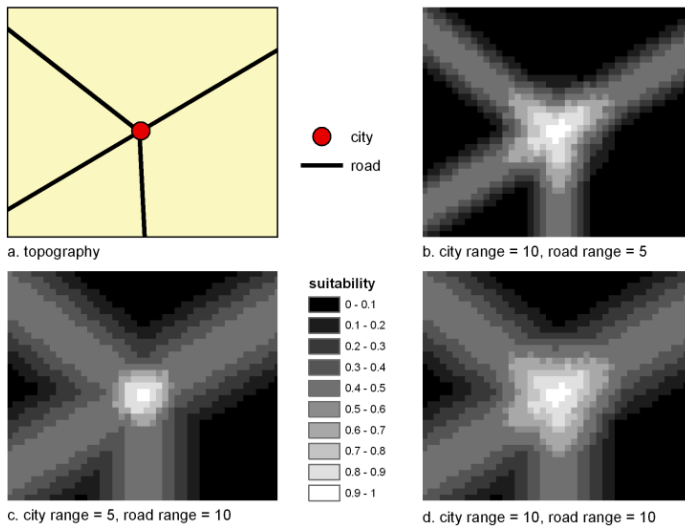
**Table 2.6: Simulated areas (km<sup>2</sup>) for the dynamic land use types in 2005, 2013, 2022 and 2030.**

year	cropland	mosaic cropland-pasture	mosaic cropland-grassland	pasture	forest
2005	9382	123737	2	7287	466205
2013	11639	135188	5230	6835	418654
2022	14137	146212	13809	6179	391540
2030	16483	155297	23405	5618	368805

Figure 2.2 shows how some land use maps resulting from a deterministic run of the model can be spatially and temporally explored using the Aguila visualization tool. Table 2.6 quantitatively summarizes the areas occupied by the dynamic land use types. Table 2.6 and Figure 2.2 give an impression of the expansions and contractions of the different land use types. The spider-web-like pattern of deforestation (dark purple) is a result of wood harvesting near roads. Roads namely form both the edge the forest plot, so that harvesting is easier, and a transportation possibility. The bands of deforestation become broader over time, but stabilize more or less in the last two time steps shown, as forest then starts to regrow at cells that were emptied at the beginning of the simulation. Cropland (red) expands, mainly around cities at the cost of (mosaic) cropland-pasture (light green). This can be seen as agricultural intensification and specialization. The land use type cropland-pasture that represents the extensive self subsistence farming practices including extensive cultivation of crops and grazing of livestock on the same plot, is relocated to areas that have been cleared by the harvest of wood. Areas of (mosaic) cropland-grassland (yellow) expand as well, but more at the outer border (away from cities) of existing cropland-grassland and cropland-pasture plots, because they have less economic value and are thus allocated further away from population centres (see close up in the lower right panel in Figure 2.2). Pasture intensification, i.e. conversion from cropland-pasture to 'pure' pasture (light purple), takes place primarily in the North-East of Mozambique, where the largest concentrations of livestock, in this case goats, are present. For pasture it is less essential than for crops that they are located close to a market place, as animals are self-transporting and are taken to the market less frequently (von Thünen, 1966), so they are located even further away from the city than cropland-grassland. In this way concentric rings of land use evolve, as predicted by e.g., William Alonso's Bid Rent Theory (Alonso, 1964) and the Von Thünen model (von Thünen, 1966). A lot of information can be derived from this deterministic output, but it gives no information about the certainty of the observations, i.e. how general and how certain are the observed patterns?



**Figure 2.3:** Three different realizations of land use in 2022 (time step 18) zoomed in around the city Nampula, indicated in black, and the cursor window (top right) showing the land use type of the selected cell (cross in map views). The legend is the same as the one given in Figure 2.2.



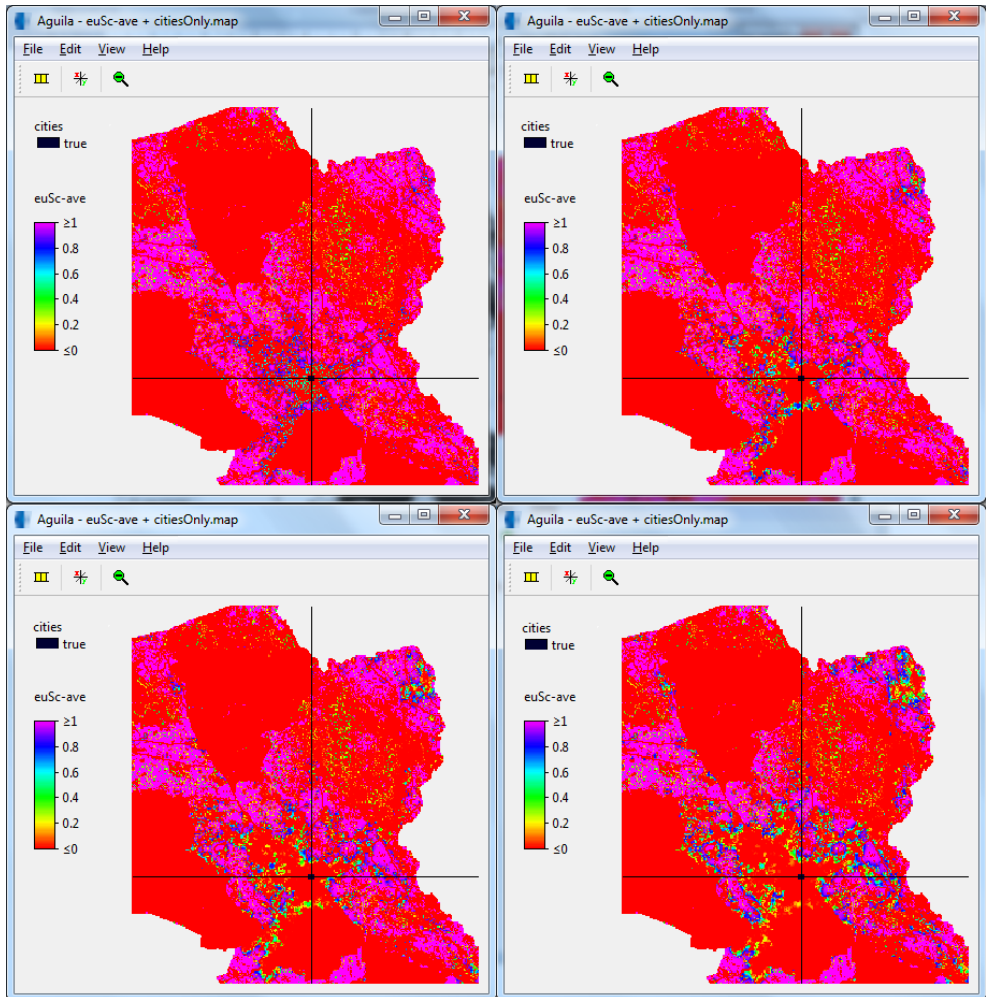
**Figure 2.4:** Artificial dataset showing the effect of the range parameters of city and road on the total suitability and consequently on the shape of the concentric 'rings'.

To study this, AguilA allows visualization of different Monte Carlo runs in linked views, which make comparison easy. Figure 2.3 shows the variation within three out of the 500 generated samples by providing a close up of land use in 2022 in the area around the city Nampula. It can be seen that the overall pattern of concentric rings is the same for all three realizations, but some differences are visible (see, e.g., the values of selected cell in the cursor window). One is that it is evident that



the centre and right image have some abandoned (dark blue) cells, while the left one has none. This can be explained by the fact that the values for demand ( $d_{n,t}$ ) and maximum yield ( $m_{n,t}$ ) are sampled separately. Cells that are abandoned on the centre and right image are most often classified as cropland-pasture (light green) in the left one, so we take this land use type as an example for the explanation. It is given as a model input that both the demand and maximum yield of cropland-pasture increase over the modelled period. The growth rate however, differs per realization, depending on the values of  $Z_d$  and  $Z_m$ . As a result, three different situations are possible. The first situation applies for the centre and right image: the demand for cropland pasture has increased, but agricultural and livestock productivity have increased so much that in the area currently occupied by cropland-pasture too much yield is generated. As a result the land use type cropland-pasture contracts in order to balance yield with demand. The second situation is that the agricultural and livestock productivity have not increased enough to counteract the effect of increasing demand, so the land use type expands. The last situation is an equilibrium, in which the increases in demand and maximum yield are in balance, so that no expansion or contraction is necessary. It cannot be derived directly from the left map in Figure 2.3 which of the last two situations has occurred there, but it can be assumed that it is the second situation, expansion, as an exact equilibrium situation is very improbable. The described process is complicated by the fact that land of a certain type can be taken by another type during the simulation, so the current yield has to be updated constantly to check which situation is at hand. This stresses the fact that model output cannot be directly related to the input uncertainties due to the numerous non-linearities in the model.

Another observation that can be made in Figure 2.3 is the shape of the concentric 'rings' around Nampula. If we focus on cropland (red), the centre image shows circular clustering around the city, while in the right image it has a more star-like shape, with lumps around the four roads connecting to Nampula. The left image has a shape somewhat in between. This difference is an effect of two suitability factors: distance to roads (suitability factor 2) and distance to cities (suitability factor 4). For cropland both factors are used (see Table 2.2), with the same weight  $w_{i,n}$ , which means they have equal effects on the total suitability for cropland. For both suitability factors a separate realization is made for the stochastic parameter  $Z_r$ , which determines the range  $r_n$ , i.e. maximum distance of effect. Figure 2.4 illustrates how these two range parameters effect the sum of the suitability factors 2 and 4. The figure shows that a smaller value for the range of roads results in more clustering around roads and consequently star-shaped concentric 'rings' (Figure 2.4b), while a smaller value for the range of cities results in more clustering around cities and consequently circular concentric 'rings' (Figure 2.4c), Similar values for the two ranges results in a shape in between (Figure 2.4d).

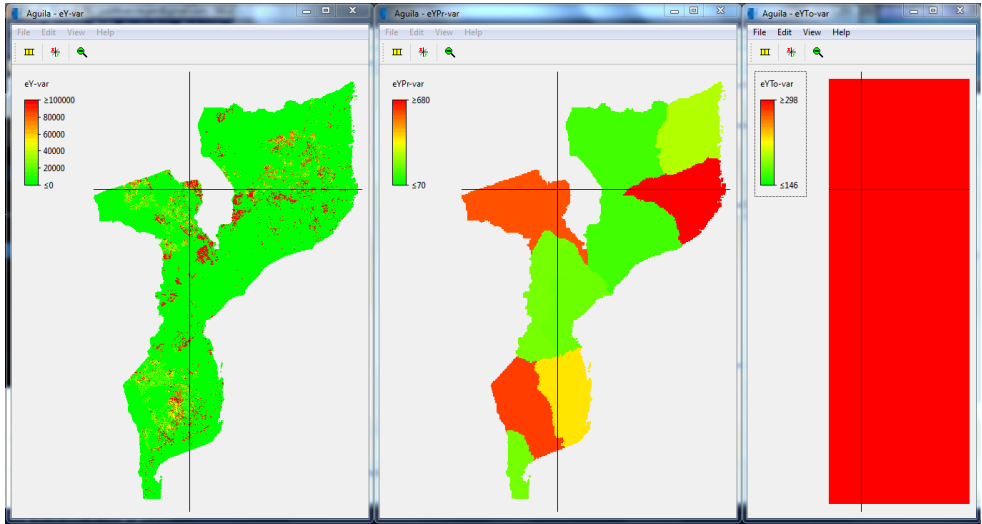


**Figure 2.5: Probability that a cell is available for bioenergy (eucalyptus) in 2005 (time step 1, top left), 2013 (time step 9, top right), 2022 (time step 18, bottom left) and 2030 (time step 26, bottom right) zoomed in around the city Tete indicated in black.**

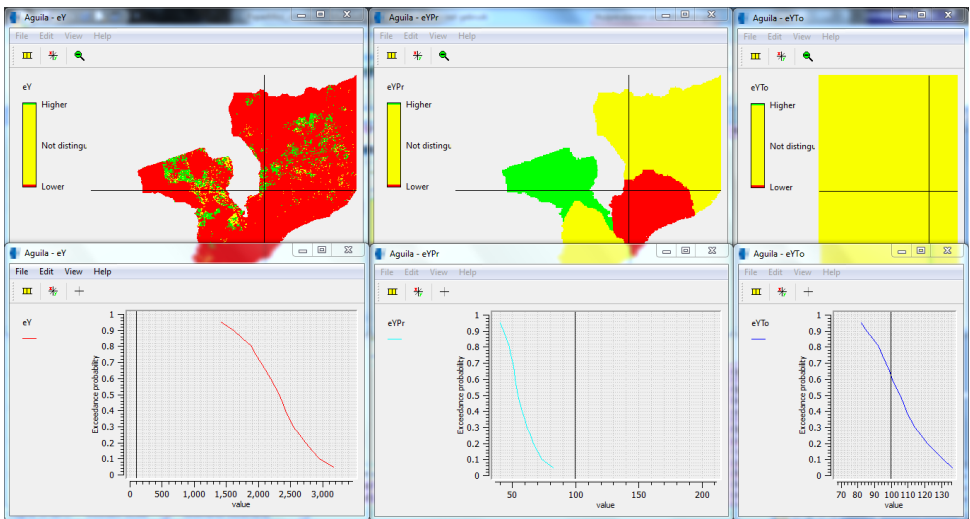
The total area covered by the agricultural land use types and their spatial distribution determine what land is available for bioenergy crops. For this purpose, it is not very convenient to look at all realizations separately. Therefore, a result from the summary statistics over all 500 Monte Carlo samples is used (see line 58 in Table 2.5). Figure 2.5 shows the probability that the bioenergy crop eucalyptus can be cultivated in a cell at a certain point in time without interfering with other important land uses. A value 1 indicates that a cell is certainly available, i.e. it was available in all realizations, a value 0 that it is certainly unavailable, i.e. it was available in none of the realizations, and any value in between indicates uncertainty in availability. In 2005 some land is available for eucalyptus around

Tete, but in 2013 a ring has formed around the city that is certainly unavailable. This can be explained by the formation of concentric rings of agriculture around the city that may not be disturbed by the cultivation of eucalyptus. The zone around Tete that is unavailable for eucalyptus becomes larger over time. It can be seen that the edges of the ring have a value somewhere between 0 and 1. This is because the size and shape of the concentric rings differ between the samples, so that cells at the edges of the concentric ring are in some samples occupied by agricultural land use types in the considered year and in some samples not. In the first case they are not available for eucalyptus and the second they are. This uncertainty ring 'moves' away from the city and roads through time, as the concentric rings around the city grow.

The planning of bioenergy crop plantations is not only of interest to local decision makers, but also of concern at higher managerial levels, e.g., province and country level. At these levels, the possible yield of eucalyptus per province and for the whole country at a certain point in the future and the uncertainty in these predictions are relevant. Figure 2.6 shows the variance in potential eucalyptus yield at three different aggregation levels. The yield is calculated per  $\text{km}^2$  at all three levels for comparability reasons. The variance in yield is a combined result of cell availability (Figure 2.5), yield fraction and maximum yield of eucalyptus as explained in the following. Most cells in Figure 2.6 have a variance of zero, because they have an availability probability of zero, or because the soil is infertile (yield fraction of zero). Some cells have a variance slightly above zero; these cells have an availability probability of one, but their yield differs because of the stochastic parameters in yield fraction ( $Z_f$ ) and maximum yield ( $Z_m$ ). Finally, there are some cells with very high variances; these cells are sometimes unavailable, and then have a yield of zero, and sometime available, and then have a yield dependent on the stochastic parameters yield fraction and maximum yield. This generates very large variances. It is evident that the maximum variance (maximum value of the value scale bar shown on the left side of the three maps) becomes much lower when scaling from cell level ( $1 \cdot 10^5 \text{ (kg} \cdot \text{km}^{-2} \cdot \text{year}^{-1})^2$ ) to province level ( $680 \text{ (kg} \cdot \text{km}^{-2} \cdot \text{year}^{-1})^2$ ) to country level ( $298 \text{ (kg} \cdot \text{km}^{-2} \cdot \text{year}^{-1})^2$ ). This is because local differences between samples are levelled out at higher aggregation levels.



**Figure 2.6: Variance in yield ((10-5 kg-km-2-year-1)<sup>2</sup>) per cell (left), per province (middle) and for the country as a whole (right) in 2030.**



**Figure 2.7: Top: 95% confidence interval for a yield above a threshold of 100 kg-km-2-year-1 in 2030. Higher = certainly above, Lower = certainly below, Not distinguishable = confidence interval straddles the threshold. Bottom: probability of exceedance (y-axis) for different yield values (x-axis). The vertical black line indicates the threshold value. Left = per cell, centre = per province, right = for the whole country.**

Other uncertainty information that can be used are the calculated percentiles (see line 61 in Table 2.5). The lower part of Figure 2.7 shows how Aguila visualizes these percentiles as exceedance probabilities. The s-shape of the curves indicates that eucalyptus yield is normally distributed at all aggregation levels. The width of the curves, in terms of the range of values they cover, decreases with increasing

aggregation level. This is again an indication that uncertainty decreases at higher aggregation levels. By using the option in Aguila to show confidence intervals of exceedance probabilities, the percentiles can be used to check, for example, whether a eucalyptus yield of more than  $100 \text{ kg} \cdot \text{km}^{-2} \cdot \text{year}^{-1}$  can be achieved at a certain point in time. The upper part of Figure 2.7 is the result of this query. The maps in Figure 2.7 show that on a country level (yellow) no definite answer can be given, but at a province level the province Tete (green) can definitely fulfil this condition.

## 2.4. Discussion and conclusions

We have shown how a modeller can construct a Spatial Decision Support System (SDSS) that integrates simulation, uncertainty analysis and visualization. This is considered useful in the light that current SDSSs tend to ignore uncertainty (Foody, 2003, Ivanovic and Freer, 2009). The advantages of the constructed PCRaster Land Use Change model (PLUC) is that the uncertainty analysis is coupled to the model, so that the output uncertainty indicators adapt automatically to changes in inputs or parameters, and that the analysis is iterative (Manson, 2007), i.e. evaluated at each time step, which is important in non-linear models. We claim that the output maps and graphs of uncertainty distribution in space and time provide an intuitive way for end users to take uncertainty into account in their decisions. The mode of visualization for uncertain spatio-temporal data was considered suitable by end users without specialist knowledge of statistics in study by Aerts et al. (2003a).

Hiemstra and Karssenbergh (2012) argue that Monte Carlo results can give a non-expert user the impression that it is hard to make any decision at all, because of the large number of cells that are reported as uncertain. Although it also good to stress the limited extent to which models can answer certain questions about complex systems (Manson, 2007), SDSSs usually *do* have merits for their end users. We have shown that by providing a means to display for example confidence intervals, with easily understandable qualitative categories *lower*, *higher*, and *not distinguishable* instead of difficult to interpret continuous measures such as full probability distributions, uncertainty information can be used to visualize locations where decisions can be made given a predefined confidence level. A disadvantage is that the error models of the input variables and parameters need to be specified, which is not a straightforward task, especially for users inexperienced in statistics.

Another advantage of the stochastic land use change model is that it provides insights that could have been missed in a deterministic model. For example, where the probability that bioenergy crops could be cultivated is high or low. Running the model at a finer scale on sites with a high potential bioenergy-crop yield, and

investigating the effect of actual allocation is the next stage of our research. A different example of the added value of stochastic modelling is the observation of the difference in the shape of concentric rings of land use around cities, which has provided an insight in the combined influence of the range parameters of cities and roads on the resulting land use patterns. This indicates the huge effect that such parameters can have on model outcomes and thereby emphasizes the caution that should be taken in setting such parameters deterministically.

Although the error propagation modelling provides information on the uncertainties in model outputs, it is not a means to evaluate the quality of the internal structure of the model, or to validate the model by comparing model outputs and independent observational data. Evaluation of the internal model structure would be possible in principle by using detailed observational data on certain sub-processes in the model. However, this data is currently not available for Mozambique. Also, validation is difficult because land use change in Mozambique in the past decades is characterized by civil war, independence and natural disasters, and therefore we do not expect steady continuation of past trends.

A disadvantage for end users of the usage of the Monte Carlo method in an SDSS, also concluded from other studies (Aerts et al., 2003b, Ligmann-Zielinska and Sun, 2010), is computation time. End users do not always have the time to wait at length for their model output. On a desktop pc it took about two days to run the 500 samples. Especially calculation of percentiles is computationally demanding. Another drawback is that PLUC now takes into account data uncertainties and model parameter uncertainties, but not model structure uncertainties, e.g., about the selected error models (Brown and Heuvelink, 2007), model rules, and raster resolution, which are debatable as well. This could be solved by creating a probability distribution over different plausible rules or over the range between coarsest legible resolution and finest legible resolution (Hengl, 2006), so that these can also be sampled in the Monte Carlo simulation. The latter is complex, as spatio-temporal probability distributions alter with changing resolution (Bierkens et al., 2000). Also, having more stochastic parameters complicates model parameterization and raises the number of samples required and thus increases computation time even more. Nevertheless, it is important that uncertainty in simulation models in SDSSs, which grow ever more complex, is somehow evaluated and communicated. This paper shows that this can be accomplished almost without any additional work on the modellers side, which is a major step forward in the exposure of uncertainty in SDSSs.

### **3. Identifying a land use change cellular automaton by Bayesian data assimilation**

**Judith A. Versteegen, Derek Karssenberg, Floor van der Hilst, André P.C. Faaij (2014), *Environmental Modelling & Software*, 53, 121-136.**

**Abstract** - We present a Bayesian method that simultaneously identifies the model structure and calibrates the parameters of a cellular automaton (CA). The method entails sequential assimilation of observations, using a particle filter. It employs prior knowledge of experts to define which processes might be important in the system, and uses empirical information from observations to identify which ones really are and how these processes should be parameterized. In a case study for the São Paulo state in Brazil, we identify a land use change CA simulating sugarcane cropland expansion from 2003 to 2016. Eight annual observation maps of sugar cane cultivation are used, split over space and time for calibration and validation. It is shown that the identified CA can properly reproduce the observations, and has a minimum reduction factor of 3 in root mean square error compared to a Monte Carlo simulation without particle filter. In the part of the study area where no observational data are assimilated (validation area), there is little reduction in model performance compared to the part with observational data. So, incomplete datasets, regional land survey data, or clouded remote sensing images can still provide useful information for this particle filter method, which is an advantage because good quality land use maps are rare. Another advantage is that in our approach the output uncertainty encompasses errors from expert knowledge, model structure, parameters and observation (calibration) data. This can, in our opinion, be very useful for example to determine up to what future period the results are a secure basis for decisions and policy making.

### 3.1. Introduction

A Cellular Automaton (CA) represents spatio-temporal change as local interactions of different entities and processes in a raster environment (Santé et al., 2010). The fact that a CA consists of relatively simple rules that can lead to complex patterns, makes it suitable to study complex system behaviour, which is currently considered important in environmental systems research (Manson, 2007, Page, 2011, Johnson, 2010, Grimm and Railsback, 2012). Therefore, cellular automata are applied in many environmental modelling domains, like fire propagation (Berjak and Hearne, 2002), vegetation spreading (Kéfi et al., 2007), and urban or land use change modelling (Verburg et al., 2004, Batty, 2005, Lauf et al., 2012). In CA development, one can distinguish between model structure identification, i.e. finding the set of processes to be represented in the model, conceptualized into the set of transition rules, and model calibration, i.e. finding the correct parameterization of these processes. In urban and land use change modelling, finding the set of transition rules is problematic (Santé et al., 2010, Straatman et al., 2004), which possibly poses limitations on the reliability and therefore the usability of these models.

Transition rule derivation can be done in a number of ways. 1) From fundamental, e.g., physical or chemical, laws (e.g., Collin et al., 2011). This is difficult in land use change modelling, as most fundamental laws in this field do not provide a quantitative process description. Yet, some have successfully applied physical laws to simulate land use expansion, mainly aimed at cities (Batty, 2012, Bettencourt, 2013). 2) By experts, who have experience-based knowledge of the study area. This is widely done in land use change modelling (e.g., van der Hilst et al., 2012, Yu et al., 2011), but it is somewhat subjective. 3) From empirical data. It is recognized that this is challenging in land use change modelling (Straatman et al., 2004, Hansen, 2012), but it is still important to continue exploring this option, because there is a need to find a more evidence-based approach to set up a land use change model.

One can combine the benefits of expert knowledge and empirical data by using a method for transition rule derivation in which the prior knowledge (definition of potential model structures) is defined by experts, and the posterior knowledge (identification of the best structure) is attained by empirical data. Our objective is to devise such a method, which we believe should fulfil two requirements. The first requirement is that the method should be able to quantify uncertainty (Aerts et al., 2003b, Rasmussen and Hamilton, 2012), i.e. it should not only be able to select the best model structure from all potential model structures defined by the prior knowledge, but it should give the likelihood of each individual structure being correct. In this way, a stochastic CA is obtained, which combines all potential model structures and parameters in an optimal way. The most important advantage of this is that confidence intervals of the modelled land use projections



can be defined, such that policy makers can decide up to what point in time the projections are reliable enough to be a foundation for their policies. The second requirement is that, herein, one should not only take into account uncertainty in the prior information, but also in the empirical data, the observations of land use, used to update the priors (Fang et al., 2006). Ignoring uncertainty in the empirical data may lead to an underestimation of model output uncertainty.

The combined requirements of prior knowledge, observation uncertainty, and posterior knowledge with output uncertainty lead towards Bayesian methods, which start out with prior knowledge, and then assemble model uncertainty and observation uncertainty to end up with posterior knowledge including uncertainty information. Therefore, we show a method for model structure identification and calibration using the particle filter, a sequential Bayesian estimation, or data assimilation, technique (van Leeuwen, 2009). Data assimilation techniques update the prior knowledge during model runtime at time steps when observations are available. We will use this property to sequentially update both the model rules and their parameters. Data assimilation techniques are increasingly being used to calibrate spatio-temporal models in a wide range of different fields in the environmental sciences, such as oceanography (van Leeuwen, 2003), hydrology (Salamon and Feyen, 2009), and atmospheric transport (Hiemstra et al., 2012), but have, to our knowledge, not yet been applied for model structure identification. Recently, their potential has been recognized in the land use change field (van der Kwast et al., 2011, Zhang et al., 2011).

The approach that is most often used in land use change modelling to define the model structure is regression on a land use map (Verburg et al., 2002, Verburg et al., 1999, Aguiar et al., 2007, Diogo et al., 2014). This method mostly results in only one deterministic model structure, without uncertainty in either the observations used to construct the regression model or in the model itself, and therefore does not meet our requirements. In the last decade, model rule identification methods originating from artificial intelligence have become popular, like neural networks (Dai et al., 2005, Li and Yeh, 2002), and swarm intelligence algorithms (Liu et al., 2008, Feng et al., 2011). These, however, do not take into account observation uncertainty, the second requirement. Moreover, they result in black-box models (Li and Yeh, 2002), i.e. they do not provide explicit posterior knowledge. Bayesian land use model structure identification has been performed before by Kocabas and Dragicevic (2007). They apply a Bayesian network and an influence diagram. However, they do not include observation uncertainty.

In this study, we evaluate the performance of the particle filter method for model structure identification and calibration of a land use change CA. Furthermore, we assess the effect of the amount of observational data assimilated, because time series of good quality land use maps are often absent (Straatman et al., 2004). We

also consider the effect of a pre-set (expert-based) model structure, to represent the situation of a model structure identification determined beforehand, which is now common practice in land use change modelling. In all approaches we provide confidence intervals with the land use projections, useful as a decision criterion for policy makers.

The assessments are carried out on a case study of the expansion of sugar cane fields in the São Paulo state in Brazil, using an adapted form of the PCRaster Land Use Change model (PLUC) (Chapter 2). As the sugar cane is partly used to produce ethanol, this case study is relevant in view of the current debate on the sustainability of bioenergy from dedicated crops when land use change is taken into account (Lapola et al., 2010, Hellmann and Verburg, 2011). São Paulo is especially interesting because it has a long history in ethanol production (Walter et al., 2011) and very good observational data availability (Rudorff et al., 2010).

The next section provides a definition of the problem of transition rule identification in a CA, a brief explanation of data assimilation, a description of the case study, an outline of the prior information about the land use change model structure and parameters, details of the performance measures used, a description of the observational data, and a scenario sketch. This is followed by a combined results and discussion section, and a conclusion section.

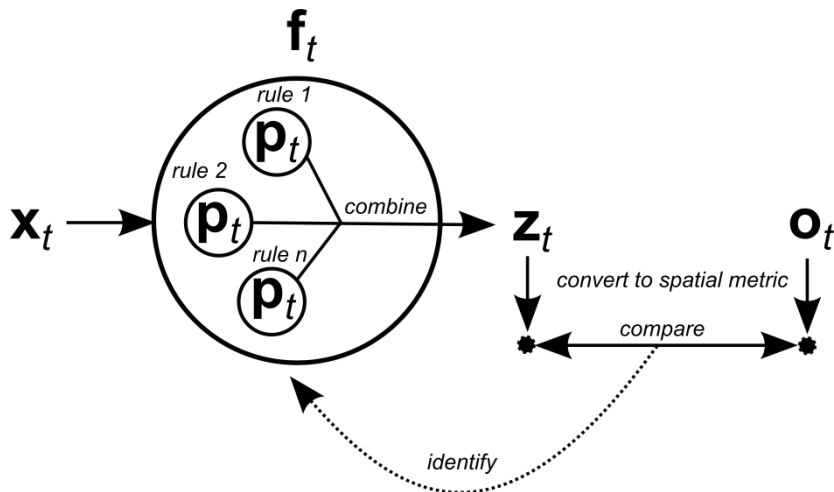
## **3.2. Methods**

### ***3.2.1. Model structure and parameter identification in a land use change cellular automaton***

A cellular automaton (CA) consists of a set of transition rules representing the processes that lead to change in the system state over time and rules to combine these transition rules (Figure 3.1). In the case of a land use change CA a transition rule is a function calculating the suitability of each location (cell) for a particular land use type, with respect to a spatial attribute that influences the allocation of that land use type, for instance the slope or the distance to roads. So, a land use change CA contains for each land use type a set of transition rules. The transition rules contain parameters defining the characteristics of the process represented by the transition rule, for example an exponent in an exponential relationship between the suitability value and slope. The transition rules need to be selected and combined such that they represent the key processes that steer the spatial allocation of land use change. This can be accomplished by selecting from a set of candidate transition rules. This could be done either in a Boolean fashion, by switching transition rules on or off, or in a continuous fashion, by weighting each

transition rule. We refer to this selection of transition rules as model structure identification.

In modelling, it is essential to find the model structure and parameter values that result in an optimal model representation of the studied land use system. Identification of the model structure and parameter values can be accomplished through comparison of the modelled system, with certain transition rules and parameter values, and observations of the real system (Figure 3.1, right side), subsequently selecting the parameter values and model structure that minimize the difference between modelled and observed land use. Parameter identification, or calibration, has become common practice in land use change CA modelling, although the applied method differs per study (Santé et al., 2010). But methods to identify the transition rules, or model structure, are generally lacking (Straatman et al., 2004). Here, we propose a technique to simultaneously identify the parameters and the model structure using observational data.



**Figure 3.1: Conceptual model of a general CA:  $f_t$  represents the processes of change in the system state over time, i.e. the set of transition rules and the way to combine them,  $x_t$  represents all inputs, usually spatial attributes, and  $p_t$  contains the parameters. Model calibration refers to identifying  $p_t$ , model structure identification refers to selecting the transition rules  $f_t$ . Identification of the parameters and model structure is based on a comparison between the land use map  $z_t$ , or a derived spatial metric, with the observed land use map  $o_t$ , or a derived spatial metric. The parameter values and model structure with the smallest difference between  $z_t$  and  $o_t$  are considered optimal.**

To summarize, a general CA, with the system state variable(s)  $\mathbf{z}_t$  and initial state(s)  $\mathbf{z}_0$ , can be defined as:

$$\mathbf{z}_t = \mathbf{f}_t(\mathbf{z}_{t-1}, \mathbf{x}_t, \mathbf{p}_t), \text{ for each } t = 1, 2, \dots, T \quad 3.1$$

In Equation 3.1,  $\mathbf{f}_t$  is the set of transition rules at time step  $t$ , representing the processes that lead to change in the system state over time. The vector  $\mathbf{x}_t$  represents all inputs, usually spatial attributes, and boundary conditions and  $\mathbf{p}_t$  contains the parameters. In a stochastic model, the uncertain parts of the system are described stochastically. So, we have  $p(\mathbf{f}_t)$ , a probability distribution of possible transition rules,  $p(\mathbf{z}_{t-1})$  the probability distribution of the previous system states,  $p(\mathbf{x}_t)$ , the probability distribution of inputs and boundary conditions, and  $p(\mathbf{p}_t)$ , the probability distribution of the parameters. In the case that no observational data are used, these distributions together determine the shape of the resulting probability distribution of the state variable, referred to as  $p(\mathbf{z}_t)$ . Yet, our aim is to use observations to simultaneously identify  $p(\mathbf{f}_t)$  and  $p(\mathbf{p}_t)$  in such a way that the model output matches the observations as closely as possible.

### 3.2.2. General particle filter framework

If we want to use the information comprised in system observations  $\mathbf{o}_t$  to select the transition rules (model structure identification) and parameterizations (calibration), that perform well and to incorporate this knowledge into the model,  $p(\mathbf{z}_t)$  should be updated (the ‘identify’ step in Figure 3.1). Bayes’ rule updates a probability distribution of a variable, when evidence, i.e. an observation, of this variable arrives. So, for the time steps at which observational data are available the following equation is evaluated.

$$p(\mathbf{z}_t | \mathbf{o}_t) = \frac{p(\mathbf{o}_t | \mathbf{z}_t) \cdot p(\mathbf{z}_t)}{p(\mathbf{o}_t)}, \text{ for each } t \quad 3.2$$

In Equation 3.2,  $p(\mathbf{o}_t)$  is the probability distribution of the observations, i.e. the measurement data and their uncertainty. Thereby, model structure identification using Bayes’ rule fulfils our second requirement of taking into account observation uncertainty.  $p(\mathbf{o}_t | \mathbf{z}_t)$  is the joint probability density of the observations at  $t$  given the model state, which can be seen as the likelihood that the observations occur given the model. The posterior probability  $p(\mathbf{z}_t | \mathbf{o}_t)$  is the probability distribution of the state variable  $p(\mathbf{z}_t)$  adjusted to the observations. Hence, Bayes’ rule quantifies output (posterior) uncertainty given observation and input (prior) uncertainty, thereby satisfying our first requirement.

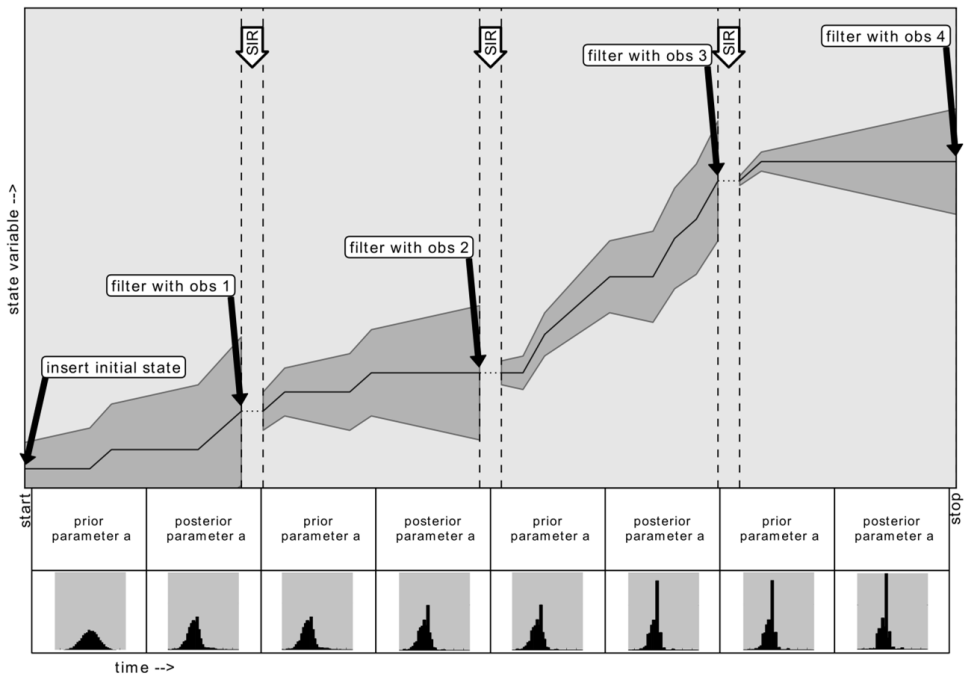
Numerically, Equation 3.1 is often solved using Monte Carlo analysis, which represents probability distributions by a number of realizations,  $N$ , of the model.

Several sequential data assimilation techniques are available to solve Equation 3.2. Most data assimilation techniques are based on filtering theory (Jazwinski, 1970): they filter the Monte Carlo realizations sequentially over time. The most well-known filter technique is probably the Ensemble Kalman filter, first introduced by Evensen (1994). This filter is, however, not guaranteed to work with non-Gaussian distributions and non-linear systems and is thus not suitable for identifying transition rules in complex systems (Pasetto et al., 2012). More importantly, the standard version of the Ensemble Kalman filter allows updates only for the variables for which observations are available, which makes updating model structure and parameters in a land use change CA impossible, as these are not observable in the real world. Versions of the Ensemble Kalman filter that do allow this require very strict premises that do not hold for cellular automata. Therefore, we have selected another data assimilation technique that allows updating all, also non-observed, model variables and can handle non-Gaussianity and non-linearity: the sequential importance resampling (SIR) particle filter (van Leeuwen, 2009), hereafter simply referred to as the particle filter. Here, we only provide a short description of the particle filter. For a more extensive introduction into the particle filter, see e.g., Arampulam et al. (2002), Bengtsson et al. (2008), and van Leeuwen (2009). At each time step for which observational data are available the particle filter uses Bayes' rule (Equation 3.2) to assess the probability that a certain Monte Carlo realization, here called particle, and the observed data can be considered equal (Hartig et al., 2011). Herein, the following steps are taken (Figure 3.2):

1.  $N$  realizations are drawn from the initial probability distributions of model structures  $\mathbf{f}_t$ , inputs  $\mathbf{x}_t$ , and parameters  $\mathbf{p}_t$  (Equation 3.1), resulting in a total number of  $N$  particles.
2. For all  $N$  particles the land use change model is run up to the next filter moment, i.e. the next moment for which observational data are available.
3. The posterior probability that the modelled state at that moment is correct given the observations with their uncertainty, is calculated for each of the particles. A posterior probability of one indicates a perfect match and a posterior probability of zero a complete mismatch.
4. Now, the sequential importance resampling (SIR) is performed:  $N$  particles are drawn to be progressed to the next observation moment with probabilities that are proportional to the posterior probabilities calculated in step 3. This procedure causes particles with a high posterior probability to be copied often (drawn several times) and particles with a low posterior probability to be removed (never drawn).
5. Steps 2 to 4 are repeated until all filter moments are completed and the model has reached the final time step. This means that at each filter moment the initial distributions obtained in step 1 are narrowed, i.e. the

number of unique particles diminishes over time (see e.g., the histograms of parameter  $a$  in Figure 3.2).

Note that, whenever a particle is copied, all model components within it are copied. Within  $p(\mathbf{z}_t|\mathbf{o}_t)$  the transition rules  $p(\mathbf{f}_t)$  and parameters  $p(\mathbf{p}_t)$  are thereby updated as well. So, after assimilation of all observations, i.e. after the last filter moment, the best model structure is identified and the CA is calibrated, i.e. the posterior probability distributions of  $\mathbf{f}_t$  and  $\mathbf{p}_t$  are obtained.



**Figure 3.2: Functioning of the particle filter. ‘Obs 1’ means observations at filter moment 1, the solid dark grey line indicates the median system state, grey areas represent the confidence interval. Histograms underneath the plots illustrate the effect of the filter moments on a parameter, referred to as parameter  $a$ . The prior distribution of parameter  $a$  at filter moment  $t$  is always equal to the posterior of parameter  $a$  at filter moment  $t-1$ . The effect on the transition rules is the same (not shown).**

For conducting these steps, the PCRaster Python framework is used, which is freely available via <http://pcraster.geo.uu.nl> (Karszenberg et al., 2010). Steps 1 and 2 involve running the land use change model in Monte Carlo mode, as explained in Chapter 2. Step 3 is achieved by solving Bayes' theorem for each particle:

$$p(\mathbf{z}_t^i | \mathbf{o}_t) = \frac{p(\mathbf{o}_t | \mathbf{z}_t^i) \cdot p(\mathbf{z}_t^i)}{\sum_{j=1}^N p(\mathbf{o}_t | \mathbf{z}_t^j) \cdot p(\mathbf{z}_t^j)}, \text{ for each } i = 1, 2, \dots, N \quad 3.3$$

In Equation 3.3,  $p(\mathbf{z}_t^i)$  is the prior probability of model realization  $i$ , which is always equal to  $1/N$  because the same number of particles is drawn at each filter moment (step 4). If the observations are not of the state variable, but of a derived spatial metric, like relative proportions of land use in a subarea as it is often found in census data, the model state  $\mathbf{z}_t^i$  has to be converted to that measure before filtering.

In Equation 3.3,  $p(\mathbf{z}_t^i | \mathbf{o}_t)$  is the posterior probability of particle  $i$  and  $p(\mathbf{o}_t | \mathbf{z}_t^i)$  is the probability of the observations given particle  $i$ . Under the assumption that the observation error has a Gaussian distribution, the latter can be calculated as (van Leeuwen, 2009):

$$p(\mathbf{o}_t | \mathbf{z}_t^i) = e^{-1/2[\mathbf{o}_t - \mathbf{z}_t^i]^T \mathbf{R}_t^{-1} [\mathbf{o}_t - \mathbf{z}_t^i]}, \text{ for each } t \quad 3.4$$

In Equation 3.4,  $\mathbf{R}_t$  is the covariance matrix of the observation error and  $T$  indicates matrix transposition. Going through steps 1-5, the procedure 'filters' the ensemble of particles because many particles do not match the observations, receive low weights, and are thus not drawn and not progressed to the next observation moment. So, although the number of particles remains the same, due to the resampling in step 4, the variation in the particles in terms of their uniqueness in the transition rules and parameters diminishes. This means that the initial probability distributions of these model components are narrowed. Hence, the particle filter has identified which transition rules are most likely to be valid (model structure), and in what ranges the parameters are most likely to fall. The model has thereby been calibrated.

### 3.2.3. Identifying transition rules of a land use change CA

#### Case study

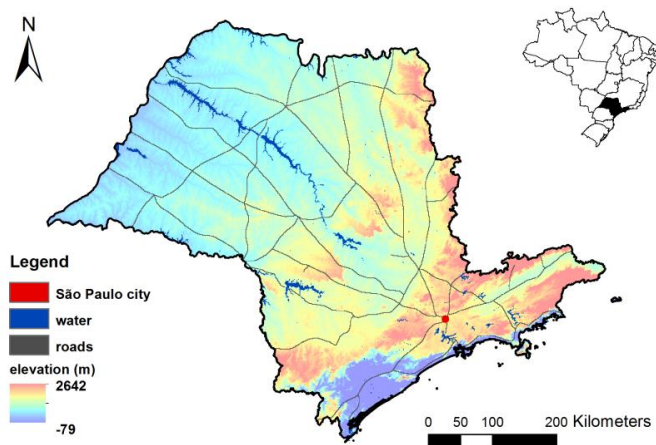
A case study is defined to test the usability of the particle filter for model structure identification and calibration of a land use change CA. An important current debate in the land use change domain is whether bioenergy from dedicated crops is still sustainable when land use change (direct and indirect) is taken into account, in view of, e.g., carbon emissions (Searchinger et al., 2008, Fargione et al., 2008,

Lapola et al., 2010), rising food prices (von Braun, 2008), and biodiversity (Hellmann and Verburg, 2010). For all these aspects it is important to know where bioenergy crops have expanded in the past and are likely to expand in the future. Such projections can be made with a land use change CA.

A key player in the bioenergy market is Brazil, mainly with the production of ethanol from sugar cane. Within Brazil, the state of São Paulo (Figure 3.3) has the longest history as well as the largest share in sugar cane production (about 60% of the national production in recent years). In addition it still experiences a significant production growth (Walter et al., 2011, Sparovek et al., 2009), especially since the introduction of the flex-fuel car in 2003 (Macedo and Seabra, 2008). The actuality of the debate, together with availability of an annual spatial dataset of sugarcane cropland distribution from the Canasat project of the National Institute for Space Research in Brazil (INPE) (Rudorff et al., 2010) as observational data, makes sugar cane cropland expansion in the São Paulo state a suitable case for testing the merits of the particle filter for identifying a CA.

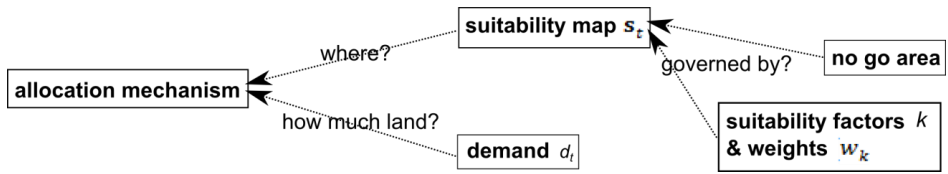
### *Transition rules and parameters: prior information*

For simulating sugar cane cropland expansion an adapted form of the PCRaster Land Use Change model (PLUC) (Chapter 2) is used. The land use change transition function,  $f_t$  in Equation 3.1, is regulated by the spatial allocation mechanism, which relies on the land demand and the suitability map (Figure 3.4).



**Figure 3.3: Study area**

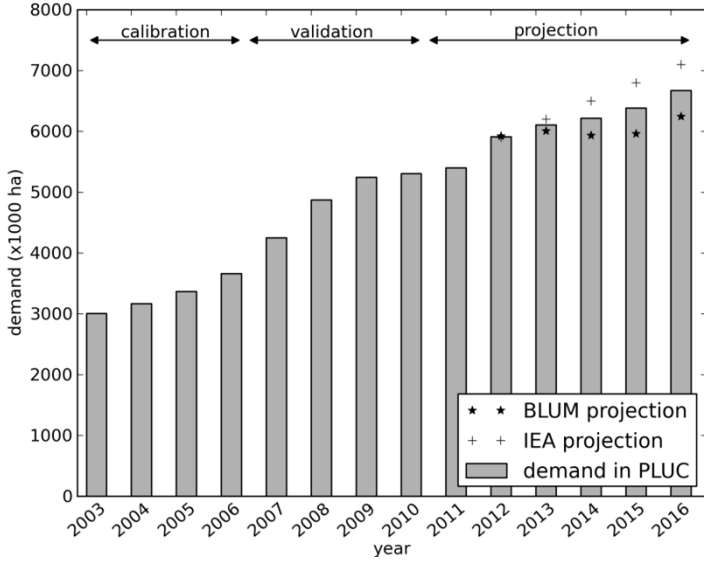




**Figure 3.4: Schematic representation of the land use change transition function.**

Different allocation mechanisms are possible, involving various degrees of competition between land use types. For instance, a model can fulfil the allocation of all land use types one by one (first the complete land claim of one land use type and then that of the next land use type), or assign to each cell the land use type with the highest suitability in that cell, or use a stochastic variant of the latter method. We use, however, a fixed, deterministic allocation mechanism, meaning that the cell with highest suitability is allocated first. This can be justified by the fact that there is little or no difference between different allocation mechanisms when considering only a single land use type.

The amount of land that is allocated or removed is steered by the demand ( $d_t$ ) for products associated with the land use types, in this case sugar cane. In the current study, demand is expressed in hectares of cultivated land. All maps are resampled to a one-kilometre resolution and projected to the Albers Equal Area projection to preserve correctness of area, to ensure that the correct number of hectares is allocated. During the calibration and validation phase, 2003 to 2010, the demand is known from the observational data; it is simply the total area of sugar cane cultivation on the Canasat map per year (Figure 3.5). For 2011, the Canasat map is not available to us, but the total area is, so the demand is known. Two data sources are used to construct the demand between 2012 and 2016. From the Brazilian Land Use Model (BLUM) (ICONE, 2012, Nassar et al., 2008), an economic partial equilibrium model, preliminary results of the future development of harvested area of sugar cane in São Paulo are used. For sugar cane the harvested area is always smaller than the cultivated area, because sugar cane is a semi-perennial crop: after about six to eight harvests the cycle is interrupted and the area is renovated for a year. The harvested area from BLUM is converted to cultivated area by adding the average fraction of sugar cane fields under renewal, which is derived from Canasat data. From the Brazilian agricultural economics institute, IEA, a study is used that estimates the cultivated area of sugar cane in São Paulo up to 2016 (Torquato, 2006). As we have equal trust in both sources, the demand in the land use change CA from 2012 to 2016 is the mean of the two time series created from these sources (Figure 3.5). It would be better to take into account the input uncertainty coming from the inconsistency between the two time series, but in this study we want to show how the uncertainty in the model structure and parameters propagates, without interference with the demand uncertainty.



**Figure 3.5: Demand for sugar cane area from 2003 to 2016. Demand in the calibration and validation phases comes from the Canasat maps. Demand in the projection phase is the mean of two projection time series from BLUM and IEA.**

The preferred location of the expansion of sugar cane is regulated by the total suitability map, an assembly of the no-go areas, and all suitability factors (Figure 3.4). The no-go map is derived from the sugar cane zoning for the São Paulo state (Padua Junior et al., 2012). These cells are masked from the total suitability map and therefore not available for land use change. The total suitability map  $\mathbf{s}_t \in [0,1]$  for sugar cane at time step  $t$  is:

$$\mathbf{s}_t = \sum_{k=1}^K (w_k \cdot \mathbf{u}_{k,t}), \text{ for each } t$$

$$\text{with } \sum_{k=1}^K (w_k) = 1 \quad 3.5$$

$$\text{and } \mathbf{u}_{k,t} = h(\mathbf{x}_{k,t}, \mathbf{p}_{k,t})$$

In Equation 3.5,  $k$  is the suitability factor, with  $k = 1, 2, \dots, K$  and  $w_k \in [0,1]$  is the spatially and temporally uniform weight of factor  $k$ . Furthermore,  $\mathbf{u}_{k,t} \in [0,1]$  is the suitability map for suitability factor  $k$ . The function  $h()$  uses the spatial attribute  $\mathbf{x}_{k,t}$  and parameter  $\mathbf{p}_{k,t}$  to create the proxy for land use change, and then normalizes it, i.e. linearly transforms it to a scale between 0 and 1, to obtain  $\mathbf{u}_{k,t}$ . The transformation is linear, because the actual shape of the relation (linear, convex, concave) between  $\mathbf{u}_{k,t}$  and  $\mathbf{x}_{k,t}$  is determined by the parameters  $\mathbf{p}_{k,t}$  within  $\mathbf{u}_{k,t}$ , discussed later in this section per suitability factor  $k$ .

The model structure of the CA,  $p(\mathbf{f}_t)$  (Equation 3.1), is formed by the weights  $w_k$ , because they determine if a certain process, or suitability factor, is of influence on

the total suitability map, and how large this influence is. The prior distribution of the weights is constructed using the following procedure, making sure that every weight covers the complete range  $[0,1]$  and that the constraint  $\sum_{k=1}^K(w_k) = 1$  is preserved. The weight of the first (randomly chosen) suitability factor  $k$  is drawn uniformly between zero and one. The next weight is drawn between zero and one minus the sum of the previous weights. The last weight is one minus the sum of all others. In this way, all weights have the same prior probability distribution, and all weights can become high (close to one) and are also often set to (close to) zero, i.e. the process is (almost) switched off. By assimilating observations, some weights can converge to zero over time, indicating that the processes, which the associated suitability factors embody, are irrelevant in the observed system. Other factors will prove to be important. So, by determining which suitability factors are relevant, we identify the CA model structure.

The ‘candidate’ suitability factors and a short explanation of the processes they represent are listed in Table 3.1. They are referred to as candidate suitability factors, because over the course of the calibration they either prove to be relevant, by obtaining a weight above zero, or to be irrelevant, by receiving a weight of zero. The candidate factors are derived from informal discussions with experts and literature review (Lapola et al., 2010, Walter et al., 2011, Rudorff et al., 2010, Macedo and Seabra, 2008, Sparovek et al., 2007, Sparovek et al., 2012, de Souza Soler and Verburg, 2010, Aguiar et al., 2011). Sugar cane in the neighbourhood ( $k = 1$ ) is expected to be important because larger plantations usually require less money per hectare as equipment and infrastructure can be shared (economies of scale). Also, a group of existing sugar cane fields usually already has a mill, in which the sugar cane is crushed, in the vicinity, so the sugar cane from new fields in the neighbourhood could go to the same mill. Travel time, and thereby transportation costs, to São Paulo city ( $k = 2$ ) could be of influence because the ethanol is distributed through there, so São Paulo city is the main market. It is assumed that transportation occurs by truck only (Macedo and Seabra, 2008). Potential yield ( $k = 3$ ) is important for the potential profits per hectare. We use a potential yield map created from physical landscape properties and climate data by the IIASA (Tóth et al., 2012). Slope ( $k = 4$ ) is critical, because the São Paulo state tries to eliminate pre-harvest burning, with its negative impacts on human health and on the environment due to the emission of pollutant gases (Aguiar et al., 2011). Pre-harvest burning can be banned when manual harvesting is replaced by mechanical harvesting, which does not require burning. However, the harvest machines cannot operate on sloping ground. The state law that promotes sustainable production practices for sugarcane in São Paulo State therefore induces new sugar cane agrarians to avoid slopes above 12%. Random noise ( $k = 5$ ) is used as a suitability factor to include local allocation choices that cannot be captured by a specific process or attribute at this model scale.

**Table 3.1: Candidate suitability factors for sugar cane in São Paulo**

k	Suitability factor	Process represented
1	Sugar cane in neighbourhood	Economies of scale
2	Travel time to São Paulo*	Transportation costs to the main market
3	Potential yield	Profits
4	Slope	Mechanization potential
5	Random noise	Local allocation choices (unexplained)

\* Calculated as distance divided by speed. Speeds on different road types are taken from de Souza Soler & Verburg (2010).

The tuning of the parameters of the model, or calibration, relates to finding the posterior distribution of all parameters,  $p(\mathbf{p}_t)$ , within the suitability factors. In total, there are five parameters to calibrate:  $\mathbf{p}_t = [f, l, a_2, a_3, a_4]$ , which are explained in the following. Suitability factor 1, the neighbourhood effect, is defined as:

$$\mathbf{u}_{1,t} = \text{norm}(-\mathbf{x}_{1,t}^2 + 2 \cdot f \cdot \left(\frac{l}{c}\right)^2 \cdot \mathbf{x}_{1,t}), \text{ for each } t \quad 3.6$$

In Equation 3.6, the number of neighbours being sugar cane in the neighbourhood window in a certain time step is  $\mathbf{x}_{1,t}$ . Parameter  $c$  is the cell length, in this study fixed at 1000 m, and parameter  $l$  (m) is the window length of the window that determines whether or not a cell belongs to the neighbourhood. And  $f$  is the ‘preferred’ fraction of neighbours being sugar cane of the total number of neighbours,  $(l/c)^2$ , within the window. The function  $\text{norm}()$  normalizes its contents. The prior distribution of  $l$  is lognormal,  $l = e^{Z_l}$ , with  $Z_l \sim N(8.5, 0.7)$ , which results in a median of around 5000 m and a mean of 3000 m. For these values, the window,  $(l/c)^2$ , is the extended and direct Moore neighbourhood, two of the most commonly used neighbourhood types. The prior distribution of  $f$  is uniform,  $f = Z_f$ , with  $Z_f \sim U(0, 1)$ . If, for example,  $f$  equals 0.5, the highest suitability ( $\mathbf{u}_{1,t} = 1$ ) occurs where half of the neighbours in the window is sugar cane. A reason why it could be the case that the neighbourhood shows ‘gaps’ without sugar cane is that a law, the Forest Code, requires that a certain portion of private farmland is set aside for natural vegetation preservation. In São Paulo, being outside the Legal Amazon Region, this is 20% (Sparovek et al., 2012). Another reason why more sugar cane in the neighbourhood is not always better is that the mill in which the sugar cane is crushed, has a maximum crushing capacity per season. In São Paulo the average maximum capacity is  $1.9 \cdot 10^6$  tonnes (Walter et al., 2011). When this capacity is reached, it is not necessarily an economic advantage anymore to create new sugar cane fields close to existing ones, as a new mill has to be built anyway.

Suitability factors 2, 3 and 4 from Table 3.1 are calculated as:

$$\mathbf{u}_k = \text{norm}(1 - \mathbf{x}_k^{a_k}), \text{ for } k = 2, 4 \quad 3.7$$

$$\mathbf{u}_k = \text{norm}(\mathbf{x}_k^{a_k}), \text{ for } k = 3$$

In Equation 3.7,  $\mathbf{x}_k$  is the attribute of suitability factor  $k$ . In the case of attributes travel time to São Paulo ( $k = 2$ ) and slope ( $k = 4$ ), lower attribute values lead to a higher suitability, while for potential yield ( $k = 3$ ) higher attribute values lead to a higher suitability. The parameter  $a_k$  determines the shape of the suitability function. A value of one results in a linear function, meaning that suitability increases or decreases linearly with the increase in the value of the attribute. For  $0 < a_k < 1$ , the shape of  $\mathbf{u}_k$  is concave, and for  $a_k > 1$ , the shape is convex. The prior distribution of  $a_k$  is lognormal,  $a_k = e^{Z_{a_k}}$ , with  $Z_{a_k} \sim N(0, 1.8)$ , which results in a distribution of  $a_k$  with a median of one, i.e. a linear relation between  $\mathbf{u}_k$  and  $\mathbf{x}_k$ . The uncertainty of the attributes,  $\mathbf{x}_k$ , is based on information provided with the datasets (Tóth et al., 2012, Jarvis et al., 2008). Only the locations of roads and São Paulo city, used for suitability factor 2, are assumed to be known, and therefore used deterministically. Note that these suitability factors and their uncertainty remain static over time. In reality they change, e.g., new roads can be build and potential yield can decline due to land degradation, but these processes are not taken into account due to a lack of data.

Suitability factor, 5, has no parameters. Its suitability is equal to its attribute, which is a uniformly distributed random spatial field, varying between zero and one, drawn separately in each time step, so  $\mathbf{u}_{5,t} = \mathbf{x}_{5,t}$ .

### *Spatial metrics*

The purpose of land use change models is usually not, and should not be, to simulate precisely the land use of each single cell in each year (Parker et al., 2008). More realistic is to try to capture certain spatio-temporal patterns. Therefore aggregated measures or spatial metrics are often more useful for calibration than location-based methods (Pijanowski et al., 2006). To correctly identify the system dynamics, it is essential to evaluate multiple system characteristics. So, multiple spatial metrics should be assessed in the model structure identification and calibration, observed at different system levels (Grimm and Railsback, 2012). Three spatial metrics were selected based on their complementarity (global vs. regional, configuration vs. composition (Csillag and Boots, 2005)). 1) The fraction of sugar cane in 150 x 150 km blocks. This metric ensures that the demand for sugar cane in the São Paulo state is distributed in correct proportions over regional areas. The metric is regional and based on composition. 2) The total number of

interconnected sugar cane patches. This signifies whether all sugar cane cropland is connected into one patch or distributed over many patches. 3) The landscape shape index,  $q_t$ , calculated as (Pijanowski et al., 2002):

$$q_t = e_t / \min(e_t) \quad 3.8$$

Herein,  $e_t$  (m) is the total length of the edge of the sugar cane patches in time step  $t$ ,  $\min(e_t)$  (m) is the minimum total length of edge for a maximally aggregated sugar cane patch, attained if all sugar cane is grouped into one square patch. So, it indicates the shape of the patches, e.g., very compact or more crooked. The last two metrics are global and based on configuration, so that a balance between global and regional, and composition and configuration is obtained in the model calibration and validation.

Validation is done using the same spatial metrics used for calibration, as the performance criteria should match the model purpose and thus the calibration criteria (Rykiel Jr., 1996). The root mean square error (RMSE) and 95% confidence intervals are used to quantitatively compare the spatial metrics obtained from the model projections and from the observational data of the same period.

### *Observational data*

The observational data are eight annual maps of sugar cane occurrence, classified from Landsat images by INPE for the Canasat project (Rudorff et al., 2010), with a resolution of 30 m and a temporal extent from 2003 to 2010. The data are resampled to a 1 km resolution.

In order to solve Equation 3.4 (particle filter), not only the mean observations,  $\mathbf{o}_t$ , i.e. the observed spatial metrics, have to be known, but also their error covariance,  $\mathbf{R}_t$ . As the error of the Canasat maps is not known, at least not for each cell, the error is determined by generating possible realizations of these maps. Hereto we use a stochastic simulation procedure that is widely applied to predict the uncertainty of an attribute at unknown locations, given a set of known locations (Pebesma and Wesseling, 1998). We use it to assess the uncertainty of an attribute at locations for which we already know the attribute value. Herein the following actions are taken for each time step in the calibration period:

1. Experimental semi-variances (Cressie, 1993) are calculated and plotted. This is done based on a Boolean map, in which sugar cane is one and no sugar cane is zero.
2. A semi-variogram model (Cressie, 1993) is fitted using the software gstat (Pebesma, 2004). An exponential model was chosen, because it yielded the best fit and because this model is usually a good choice when several

patterns interfere (Burrough and McDonnell, 1998), which can be expected for land use patterns that are often governed by many drivers of location.

3. Cross-validation (Burrough and McDonnell, 1998) is performed to check the semi-variogram model.
4. Gaussian simulation (Pebesma and Wesseling, 1998) is used to create 100 potential spatial fields of scalar values, in which values close to zero indicate that there is probably no sugar cane and values close to one that there probably is.
5. A threshold is applied to these 100 fields to turn them into Boolean maps again. To include not only configuration errors (sugar cane located in the wrong cell), but also composition errors (the total area of sugar cane in the map is wrong), this threshold is a normally distributed stochastic parameter with a mean of 0.5 and a standard deviation of 0.1. As a result, some of the realizations will have a larger sugar cane cropland coverage than others.
6. From each of the 100 realizations, the three spatial metrics are derived.
7. The covariance matrix  $\mathbf{R}_t$  is calculated from these spatial metrics.

## Scenarios

To summarize, the land use change model (Figure 3.4) is run for the case study. During run time Equation 3.3 is applied at each filter moment to find  $p(\mathbf{z}_t^i | \mathbf{o}_t)$ , the posterior probability of land use change model particle  $i$ . The system state  $\mathbf{z}_t^i$  and the observations  $\mathbf{o}_t$  are compared in the form of the three spatial metrics (section 2.3.3), derived from respectively the model output and the Canasat maps. Because  $p(\mathbf{z}_t^i | \mathbf{o}_t)$  contains, besides the system state, also the transition rules  $p(\mathbf{f}_t)$ , where  $\mathbf{f}_t \leftarrow w_k$  (Equation 3.5) and the parameters  $p(\mathbf{p}_t)$ , where  $\mathbf{p}_t = [f, l, a_2, a_3, a_4]$  (Equations 3.6 and 3.7), these are updated as well when the ensemble of model runs is updated using sequential importance resampling.

**Table 3.2: Calibration and validation scenarios**

#	purpose	model structure (weights)	parameters	filter in t =	# of blocks
1	Reference case	stochastic	stochastic	-	-
2	Filtering	stochastic	stochastic	2004, 2005	5
3	Less observational data	stochastic	stochastic	2004, 2005	2
4	Model rules preliminarily set	deterministic	stochastic	2004, 2005	5

As one should not use the same set of data for calibration and validation, we use a split-sample approach over space and time in all scenarios. Two years (2004 - 2005) are used for calibration and five (2006-2010) for validation. More years are used for validation than for calibration to be able to study to what extent the performance decreases over time. In addition, the data are split over space: half of the ten 150 x 150 km blocks is used for calibration (from this point onwards called calibration blocks), and the other half is not (from this point onwards called validation blocks). In this way, the performance of the calibration blocks and the validation blocks can be evaluated separately, to assess to what extent the model can be used for areas where no calibration data are available. Then, we use the calibrated and validated model for land use change projections up to 2016.

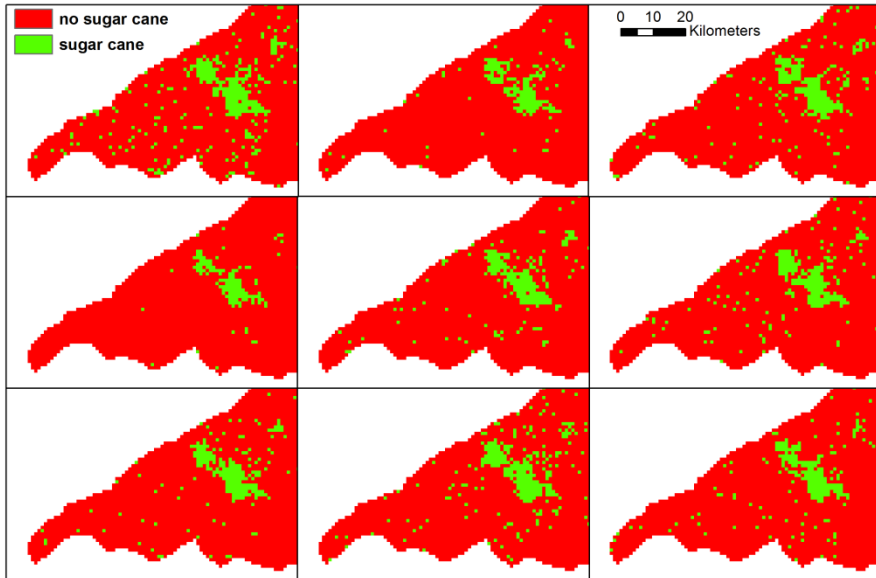
Four calibration scenarios were designed (Table 3.2). The first is the reference case, i.e. what would happen without filtering, thus only Monte Carlo simulation. The second scenario is meant to show the effect of the particle filter and to evaluate the two different sets of five blocks (sample split over space), as explained above. The third scenario uses fewer blocks, i.e. a smaller spatial coverage, to test the influence of observational data availability, as time series of good quality observations are often not obtainable (Straatman et al., 2004). In reality, this could represent a situation in which one has incomplete data for calibration, e.g., remote sensing images partly covered by clouds. The fourth scenario is meant to discuss the matter of model structure identification. In many calibration efforts the model rules are preliminarily set and only the parameters are calibrated. This scenario represents that situation in order to evaluate the difference between pre-set (this scenario) and stochastic model rules (scenario 2). All scenarios are run using  $N = 25000$  particles.

### **3.3. Results and Discussion**

#### **3.3.1. Realizations of observations**

For the filter moments, realizations of the observations were created, using Gaussian simulation as explained in section 2.3.4. These realizations represent potential instances of sugar cane cropland maps and are used to determine the covariance matrix,  $\mathbf{R}_t$ . Figure 3.6 shows a close up of nine of the hundred realizations for 2005. On the one hand one can identify the configuration errors (sugar cane cropland located in the wrong cell). For example, variations in the shape of the large patch in the middle of the close up. On the other hand, the effects of composition errors (the total area of sugar cane cropland in the map is wrong) are visible. The upper left realization, for instance, shows a larger total sugar cane cropland area than the one below it.





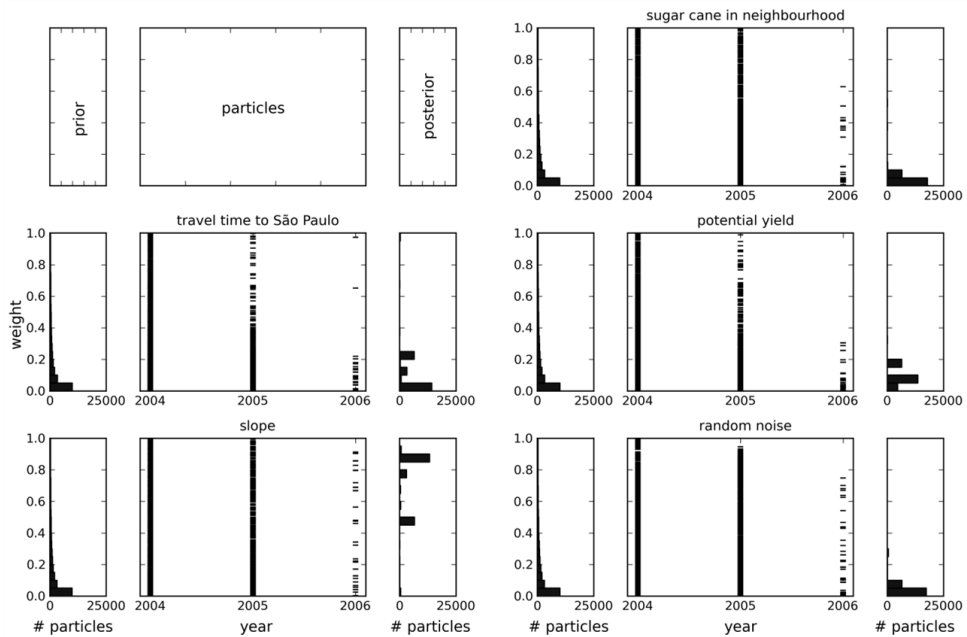
**Figure 3.6: Close up of the most Western wedge of São Paulo of nine realizations of the observations from 2005.**

### 3.3.2. Model structure identification

The evolution of the model structure  $p(\mathbf{f}_t)$  is illustrated by the evolution of the weights  $w_k$  of the five candidate suitability factors over time for scenario 2 (Figure 3.7). All weights have the same prior lognormally shaped distribution between zero and one. Over time, some particles are filtered out and others are copied. After the first filter moment (2004), of which the results can be seen in 2005 (Figure 3.7), the effects are small. But, for example, for the random noise suitability factor, the weight distribution is narrowed: a land use change model structure in which random noise is the only relevant suitability factor is rejected. Also, for potential yield and travel time to São Paulo high weights have become less prevalent in the ensemble. In the second filter moment (2005), the distributions converge further, e.g., the weight of slope converges towards high values; its posterior distribution has two peaks around 0.9 and 0.5. This means that slope is a very important factor in the allocation of sugar cane cropland. Slope was expected to be important, because the São Paulo state tries to eliminate pre-harvest burning, as explained in section 0. We had, however, not foreseen that slope is so much more important than the other factors. Travel time to São Paulo city, for example, was expected to be more important. The results that it is not (it has a weight distribution between 0.0 and 0.2), could be because the whole study area is relatively close to São Paulo city. Another explanation can be the fact that the sugar cane is not transported to the city directly. First the sugar cane is transported to sugar cane mills, where it is

converted to ethanol. The ethanol in its turn, is usually transported to São Paulo city. Ethanol, though, is a much higher value and higher energy density product, for which transportation costs, and therefore travel times, are less important. We have considered to include travel time to the sugar cane mills as a suitability factor, but the problem with this is that new mills emerge quite often, and since we do not know where, forecasting becomes problematic with this model structure. An additional reason for relatively low importance of travel time could be that an ethanol pipeline is present running through the high-density sugar cane cropland area in the middle north part of São Paulo state to São Paulo city. More pipelines are planned in the future. Macedo and Saebra (2008) expect that by 2020 20% of the ethanol in Brazil will be transported through pipelines.

Some weights converge to values close to zero, as is the case for sugar cane in the neighbourhood and random noise, which means the processes are less relevant. The fact that random noise obtains a low weight in the posterior is a good sign. It indicates that unexplained local allocation choices have little influence on the regional and global sugar cane pattern. However, it does not converge to zero entirely, so it is not completely irrelevant.



**Figure 3.7:** Evolution of the weights  $w_k$  of the candidate suitability factors over time for scenario 2, with filtering in 2004 and 2005. For each suitability factor the black horizontal lines in the centre panel are a random selection of 10% (for visualization purposes) of the particles, the bars on the left represent the prior distribution, and the bars on the right represent the posterior distribution of the weight (see also diagram structure on the top left).

### 3.3.3. Calibration

At the same time the parameters within the suitability factors have been calibrated (Figure 3.8). The prior of the logarithm of window length  $l$ , of the sugar cane in neighbourhood suitability factor, is normally distributed and the posterior has converged to values around 8. This means that the best neighbourhood is the direct Moore neighbourhood, as  $e^8 \approx 3000$  m, which is 3 cell lengths. The prior distribution of the preferred fraction of neighbours being sugar cane  $f$  is uniformly distributed between zero and one. The posterior distribution has two peaks. The largest lies around 0.5, indicating that the preferred number of neighbours being sugar cane, is half of the total number of neighbours in the neighbourhood. Because the posterior distribution of window length  $l$  has a value of about 3 cells, this means that the highest suitability for the neighbourhood function is reached when four out of the eight direct neighbours are sugar cane. This could imply that farmers set aside part of their farmland to meet the terms of the Forest Code. But to verify this, a further study on the land use of the non-sugar-cane cells in the window should be done. When these cells are indeed natural surroundings, it could be the case, but when this land is used otherwise, e.g., for other crops or grazing of livestock, a different reason exists for the scattered pattern. At the moment we are unable to examine this, because we have no map of the other land uses in the São Paulo state that is sufficiently detailed in space, time and attributes.

The lower three panels in Figure 3.8 represent the logarithm of the parameter  $a_k$ , which determines the shape of the suitability functions of travel time to São Paulo ( $k = 2$ ), potential yield ( $k = 3$ ) and slope ( $k = 4$ ). The prior of the natural log of all three parameters is normally distributed, with a median of zero, so that the median of  $a_k = e^0 = 1$ , which results in a linear relationship between the attribute value and the suitability. The posterior distribution of  $a_2$  for travel time to São Paulo has its peaks above zero. Applying Equation 3.7, we see that the resulting shape of  $\mathbf{u}_2$  plotted against  $\mathbf{x}_2$  is convex (Figure 3.9). This means that up to a travel time of about ten hours from São Paulo city the suitability is high and almost constant, and further away it quickly drops to zero. On the map this shift arises far away from main roads in the North-western and South-western parts of the study area only. From that point onwards transportation costs possibly become too high<sup>1</sup>.

---

<sup>1</sup> If costs versus revenues are truly the reason for the position of the tipping point of  $\mathbf{a}_2$ , it should be noted that this point is very sensitive to e.g., changes in fuel prices. The model could be improved by calculating transportation costs instead of travel time, so that these economic variables are reflected in the suitability factor. However, obtaining data to accurately do that is time-consuming, as e.g., regulations, taxes and fuel prices can vary widely per administrative unit in Brazil.

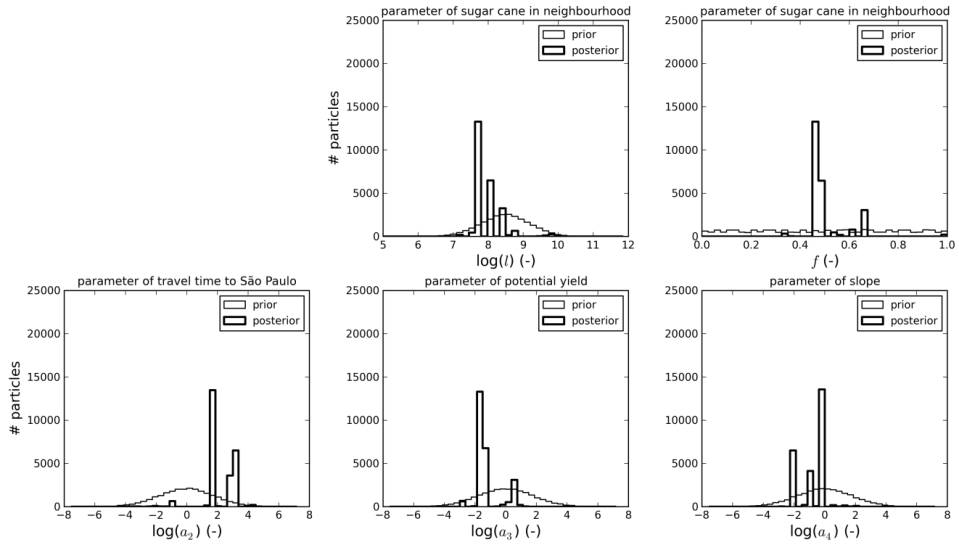
The posterior distribution of  $a_3$  for potential yield has its largest peak around -2. Applying Equation 3.7, we see that  $\mathbf{u}_3$  plotted against  $\mathbf{x}_3$  has a convex shape (Figure 3.9). The curve shows that even soils with relatively low yield, in the region of 15% of the maximum attainable yield, are suitable for sugar cane cultivation. This is in line with information provided to us by experts from CTBE (Brazilian Bioethanol Science and Technology Laboratory), who stated that in the São Paulo state all soils are good enough to cultivate sugar cane and the soil must be prepared anyway, so the precise quality is not that important. The climate is also uniformly good; in the entire state sugar cane can be cultivated without irrigation.

The posterior distribution of slope has its largest peak at approximately zero. Therefore, the median of  $\mathbf{u}_4$  plotted against  $\mathbf{x}_4$  has a linear shape (Figure 3.9), meaning that the suitability decreases linearly with the increase of slope. It should be noted that slopes higher than 12% are included in the no-go area. Therefore, in the model sugar cane cannot be allocated on high slopes anyway, so the part of the graph where  $x_4 > 12\%$  has no effect in the CA.

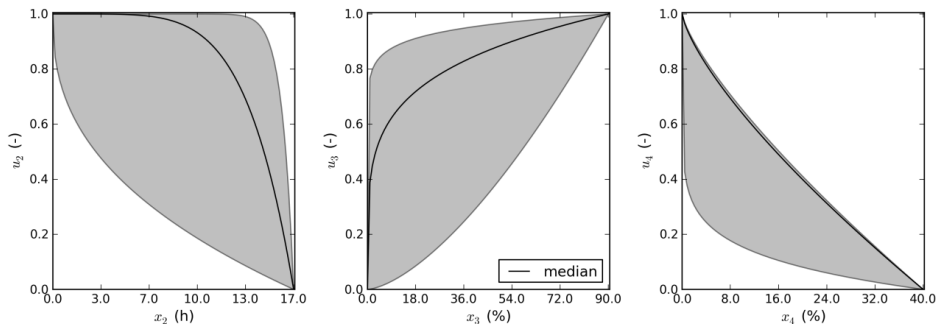
In scenario 3, in which less observational data are assimilated, the parameters converge to similar values, but the distributions remain much broader. This is because less information is available on whether a certain parameterization performs well. In scenario 4, parameter distributions are narrower again. They have converged to different values than scenario 2 and 3, to correct for the fixed model structure, which is different from the optimal model structure in scenario 2 (Figure 3.7). For example, window length  $l$ , obtains a median of about 13 km. Apparently, the fact that sugar cane is the neighbourhood was given a too high weight, can be partly compensated by calculating the number of neighbours in a larger window.

### 3.3.1. Validation

Figure 3.10 compares the root mean square error (RMSE) of the validation time steps 4 to 8 (2006-2010) of all scenarios for all three spatial metrics for the calibration and the validation blocks separately. The RMSE is a frequently used measure of the differences between a modelled and an observed variable. As the measure is scale dependent (Hyndman and Koehler, 2006), it has a relative meaning only, so we merely use it to compare between scenarios. For absolute comparison, the modelled spatial metrics and the observed spatial metrics are compared for scenarios 1, 2, and 4 in Figure 3.11.



**Figure 3.8: Prior (thin line) and posterior (thick line) distributions of all parameters for scenario 2.**



**Figure 3.9: Suitability distributions  $u_k$  for the attributes  $x_k$ : travel time to São Paulo ( $k = 2$ ), potential yield ( $k = 3$ ), and slope ( $k = 4$ ), given the median (solid black line) and the 95% confidence interval (grey area) of the posterior distribution of  $a_k$ .**

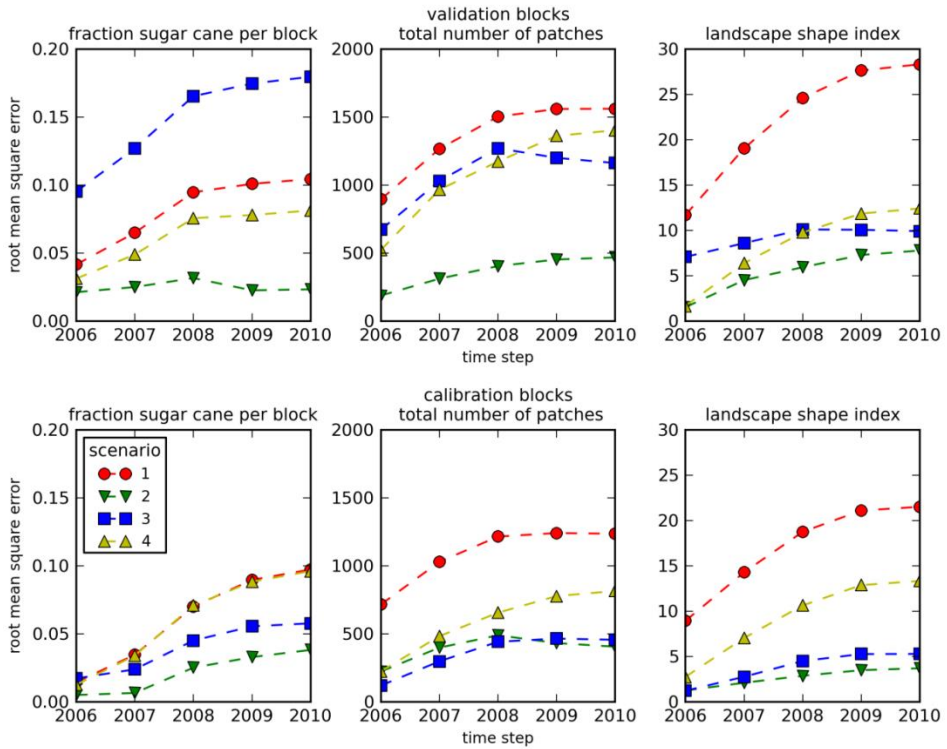
In general, it can be observed that in the two scenarios in which the particle filter was applied, the width of the 95% confidence interval is significantly smaller, and the relative difference in width between the scenarios with and the scenario without filtering increases over time. For example, the width of the confidence interval in scenario 2 is less than 10% of the width in scenario 1 in 2016 (Figure 3.11). The scenarios in which the particle filter was applied have a lower RMSE (Figure 3.10). This implies that the particle filter reduces the uncertainty in the ensemble of model runs in a way that brings the output values closer to the observed values. For the calibration blocks, the general trends are similar to the

validation blocks. Overall, the performance is a bit better in calibration blocks than in the validation blocks, as expected.

Scenario 2 outperforms the other scenarios regarding all metrics throughout the modelled period, except for two years (Figure 3.10). Concerning the fraction sugar cane per block, the median of the model output and the observations differ less than 1% up to 2006, and for three out of the five blocks even all the way up to 2010 (Figure 3.11). In scenario 1, the divergence of the two medians starts already in 2004 and concerns four instead of two blocks. This shows that the particle filter has managed to resample the ensemble of particles in such a way that the sugar cane expansion behaviour is corrected in the two other blocks, however not all the way up to 2010. The median of the model output and the observations for landscape shape index differ less than 1% up to 2006 for both scenario 1 and 2. After that, the modelled landscape shape index is too low (maximum of 10%), and the observations do not fall within the 95% confidence interval in scenario 2. So, the particle filter does narrow the confidence interval and reduce the RSME over the complete validation period, but further away from the filter moments the ensemble is not successful anymore in projecting either the fraction of sugar cane for two out of the five blocks or the landscape shape index. This can be a result of either incorrect model identification, or non-stationarity in the land use system itself. The implication of this is that, if one wants to use the identified CA to assess an impact of land use change that uses these blocks or the landscape shape index as a criterion, the reliability of this impact analysis is low. Yet, the fact that there is little difference in performance between the calibration and validation blocks for scenario 2 (Figure 3.10) indicates that the identified model structure and parameters perform equally well in the part of the study area for which no observations were assimilated. With information from half of the study area, the particle filter is able to identify model structure and parameters resulting in the same performance in the other half of the area.

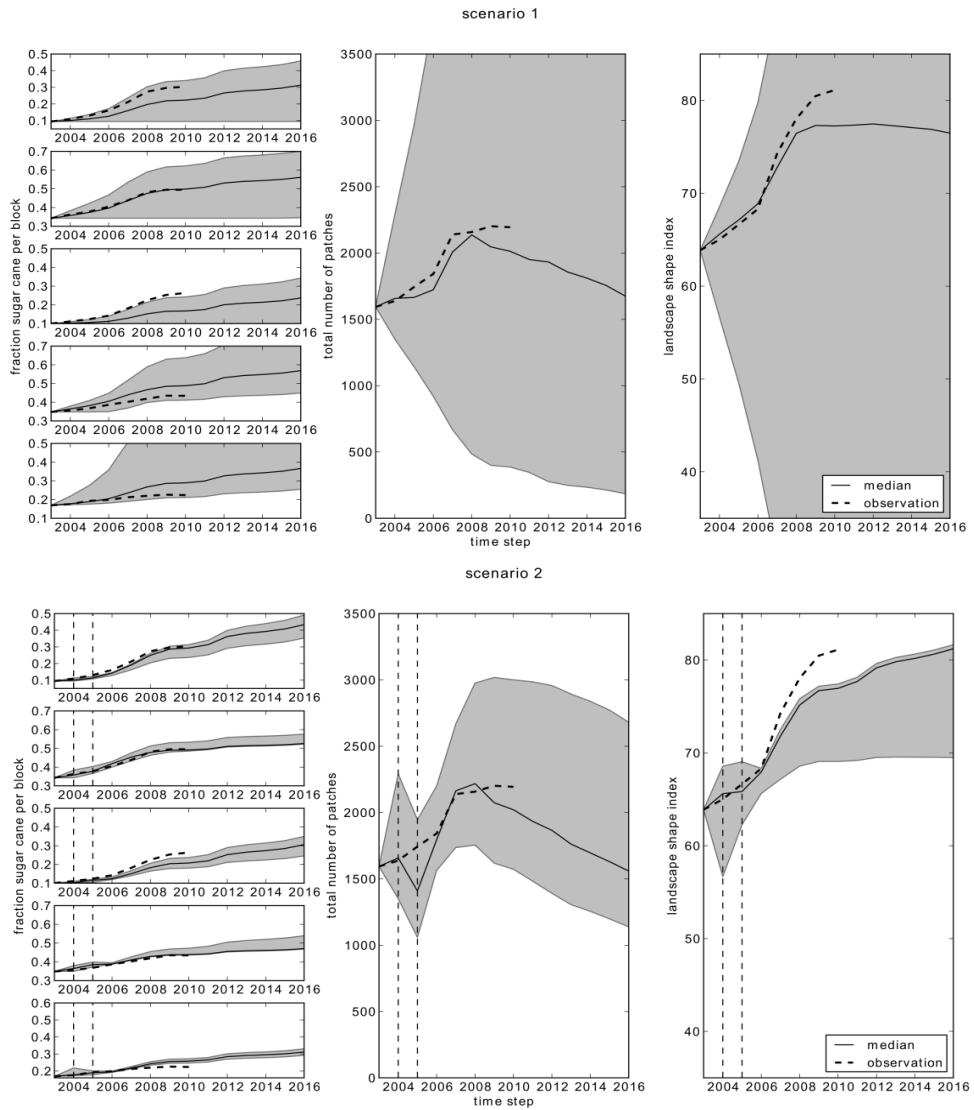
The effect of less observational data (scenario 3) is largest for the fraction of sugar cane per block (Figure 3.10), where it performs even worse than the reference scenario (scenario 1). The RMSE is more than four times as large as for scenario 2. For the calibration blocks scenario 3 seems to perform better, but is calculated only for the two calibration blocks instead of five as is the case for the other scenarios. Because of this it is problematic to compare the RMSE of scenario 3 with the other ones for the calibration blocks.

For scenario 4 (pre-set model structure but calibrated parameters), the decrease in performance is especially large for the number of patches (Figure 3.10). Figure 3.11 shows that this scenario is completely unable to capture the trend in this metric. It predicts continuous decrease from the start, while there should be an increase.



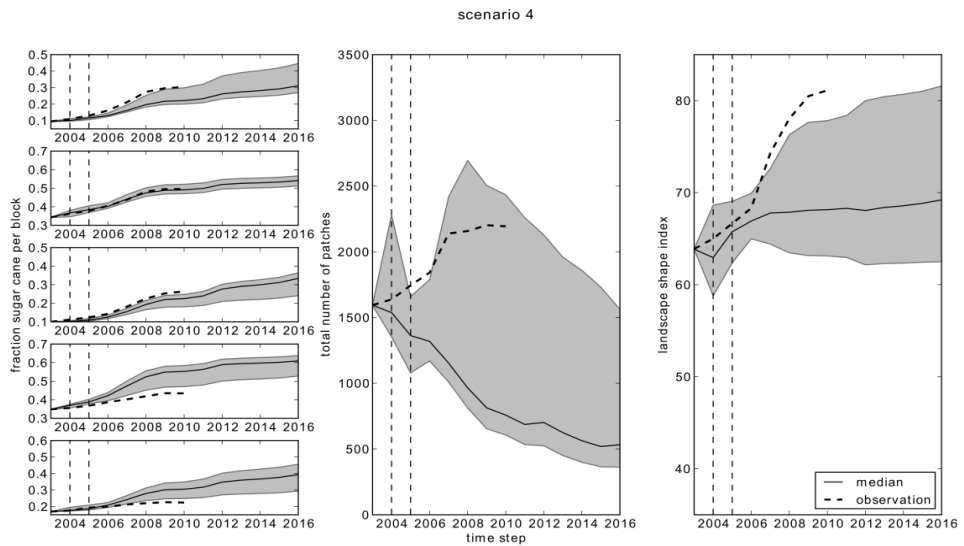
**Figure 3.10: Root mean square error for all spatial metrics for all scenarios for the validation blocks (top) and calibration blocks (bottom) for scenario 1 (reference case), 2 (particle filter), scenario 3 (less observational data), and 4 (preliminarily set model rules).**

To further compare the predictive power of the reference and the particle filter scenario, Figure 3.12 shows the maps of the probability of sugar cane cropland coverage in 2007 and 2016 for both scenarios the observational data in 2007. In 2007, two filter moments have passed in scenario 2, so its land use map differs a lot from the one of the reference scenario. In scenario 1, sugar cane cropland has expanded along edges of existing patches. The area in the Southern part of São Paulo state has probabilities of zero because this is defined as the no-go area for sugar cane (Padua Junior et al., 2012). Furthermore, almost all cells have a, although low, probability to be occupied by sugar cane in 2007. In the output of scenario 2, the uncertainty of where sugar cane will be located has been greatly reduced; most cells are either red (probability of zero) or green (probability of one). The majority of the expansion has taken place in the North-West. Both the location and the configuration (scattered) of the expansion in this scenario result resemble the observations much better than scenario 1.

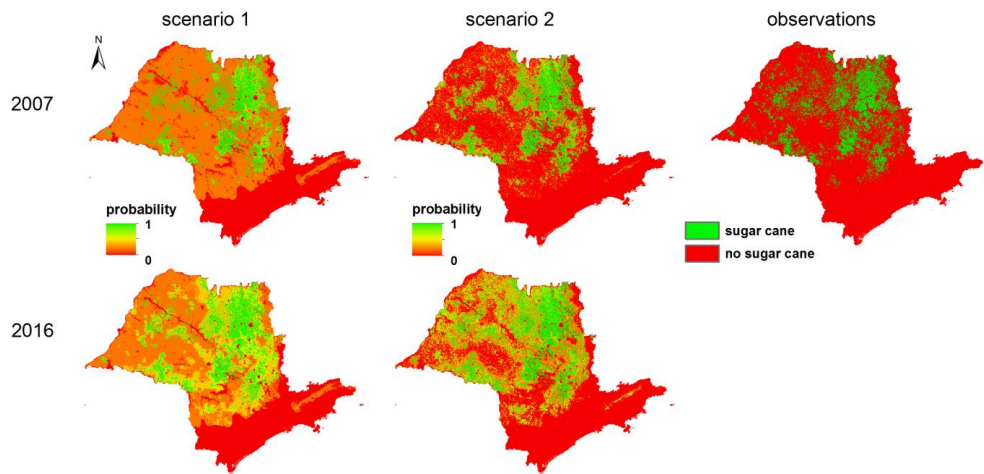


**Figure 3.11 (continues on next page, caption provided there).**





**Figure 3.11 (continued):** Comparison of modelled median (solid line), 95% confidence interval (grey area) and observations (dashed line) of the spatial metrics for the five validation blocks for scenario 1 (reference case), 2 (particle filter) and 4 (preliminarily set model rules). The filter moments are indicated with vertical dashed lines.



**Figure 3.12:** The probability of sugar cane cultivation for scenario 1 (left) and scenario 2 (centre), and the observations (right), for 2007 and 2016.

For 2016 the same differences between scenario 1 and 2 appear; scenario 1 is much more uncertain, although the uncertainty in scenario 2 has also increased, as the time period since the last filter moment is long. Scenario 2 indicates the highest probability of expansion mainly in the western part of São Paulo state.

Obviously, there are no observations for 2016, so validity cannot be checked. Yet, the sugar cane cropland expansion sites up to 2020, given by Lapola et al. (2010, p. 2) resemble the locations with a high probability in our results closely.

To ensure that the model structure and parameters obtained with the particle filter are not dependent on partition of the blocks for calibration and validation, scenario 2 was repeated with a different, again randomly drawn, block division. In Figure 3.13, Figure 3.14, and Figure 3.15 in the Appendix A the equivalents of respectively Figure 3.7, Figure 3.8, and Figure 3.11 are shown for the new sets of calibration and validation blocks. The results are deemed comparable, so we reject the possibility that the results shown in this section are a product of the partition of the area into calibration and validation blocks.

### **3.4. Conclusion**

The method used here simultaneously identifies the model structure and calibrates the parameters of a land use change cellular automaton (CA) by sequential assimilation of observations, using a particle filter. The method uses the (subjective) knowledge of experts to define which processes or drivers might be important in the system, and applies the (objective) information from observations to adjust the model structure and calibrate the parameters. With the candidate suitability factors chosen in this study and the observational data used, the particle filter identified a probability distribution of model structures and associated parameters that could forecast the chosen spatial metrics fairly good. Nevertheless, performance clearly decreased over time, further away from the filter moments (Figure 3.10 and Figure 3.11). So, using this calibration technique, information about the land use system is gained and short-term land use projections are clearly improved, but projections of more than a few years ahead are not very reliable. It remains, however, a question whether this is a deficit of the calibration technique, or a result of non-stationarity in the land use system itself. The assessment of the persistence of land use change drivers and possible changes in their relative importance over time, is the next stage of our research.

In areas where no observational data were assimilated (validation blocks), the model performed about as well as in the areas for which observations were assimilated (calibration blocks) (Figure 3.10). So, spatially incomplete datasets, regional land survey data, or clouded remote sensing images can still provide valuable information for this CA identification. Also, data gaps in time are not a problem, as the sequence of filter moments does not need to be continuous, i.e. there may be gaps in the time series of observations that is assimilated. These two characteristics of the used method are a considerable advantage given the fact

that time series of good quality land use maps are rare (Straatman et al., 2004). Only when the area of observational data availability became significantly smaller, the model performed a lot worse in the unknown areas (validation blocks). Besides the improvements on the land use change model itself, the technique also improves land use system knowledge, because the result of the model structure identification provides information on the relative importance of the drivers. For example, in our case study slope turned out to be much more important for sugar cane cropland allocation than expected in advance.

It is shown that the particle filter method can be used not only to calibrate the parameters inside the transition rules, but also to find the relative importance of the transition rules, the model structure. The importance of taking into account the identification of the model structure is shown by running a scenario with a pre-defined model structure. As expected, it performed worse compared to the scenario in which the model structure was identified by the particle filter. However, the choice for this specific, pre-defined structure was of course arbitrary, so this result should be interpreted with caution. A 'real' expert might have chosen a completely different model structure, possibly performing better.

Before carrying out the method presented in this paper, one should also contemplate the aim of the CA modelling effort. The calibration target, i.e. the spatial pattern that the model should be able to reproduce, should match the modelling aim. If, for instance, the effect of future land use change on animal passageways is studied, connectivity of patches is probably an important characteristic. Hence, one or more measures of connectivity should be used as a calibration and validation target. In this study, three spatial metrics were used. It turned to be complicated to correctly reproduce the landscape shape index in the longer run. This should be taken into account when studying impacts that rely on these properties.

Although it is an advantage that the search space, i.e. the prior distributions of parameters and model structure, can be defined by experts, it should be kept in mind that this prior information has a large effect on the outcomes. The selection of candidate suitability factors and the prior distribution of their weights should be performed carefully. The potential solution to incorporate a huge number of candidate suitability factors to make sure that all possible drivers are considered, is not feasible, because the addition of one parameter makes the number of required particles increase exponentially (Bengtsson et al., 2008). Many CAs have such a large number of parameters that computation time and disk space become severe constraints. Three possible solutions, that can also be combined, are: 1) to fix parameters having little influence on outcomes or having evident values and to calibrate only the remaining ones, 2) to apply a more advanced particle filter scheme giving similar results with a lower number of particles thus reducing the

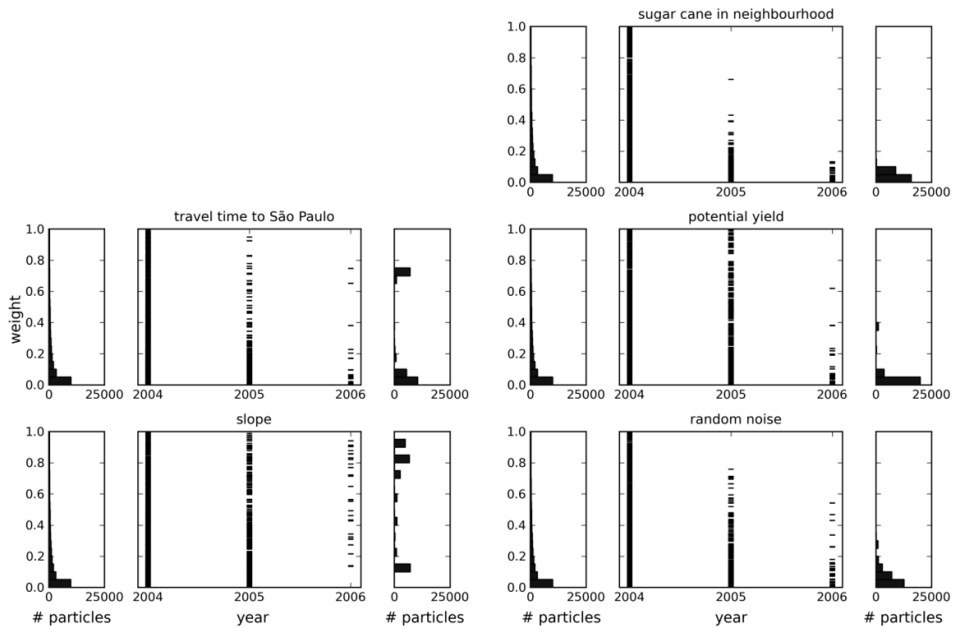
required run time and disk space (Spiller et al., 2008, Jeremiah et al., 2012), and 3) to use super computers or cluster machines, to allow for a larger number of particles in the data assimilation scheme, thereby enabling an increase in the number of parameters that can be calibrated.

An advantage of the method shown is that model uncertainty and observation uncertainty are taken into account. It must be acknowledged that the uncertainty in the observational data is constructed by making realizations using Gaussian simulation. So, the assumed observation uncertainty can be incorrect, but at least observation uncertainty is not ignored, as it is in many other predictive studies (Ivanovic and Freer, 2009). This means that our output uncertainty encompasses errors from model structure, parameters and observation (calibration) data. Such a full scope error propagation assessment is, to our knowledge, new in land use change CA modelling. In our opinion, it can be very useful, for example, to determine for which future time frame the results are reliable enough to base decisions and policies on. A practical application of how uncertainty information can be used in decision making is given by, e.g., Aerts et al. (2003b) and Chapter 2.

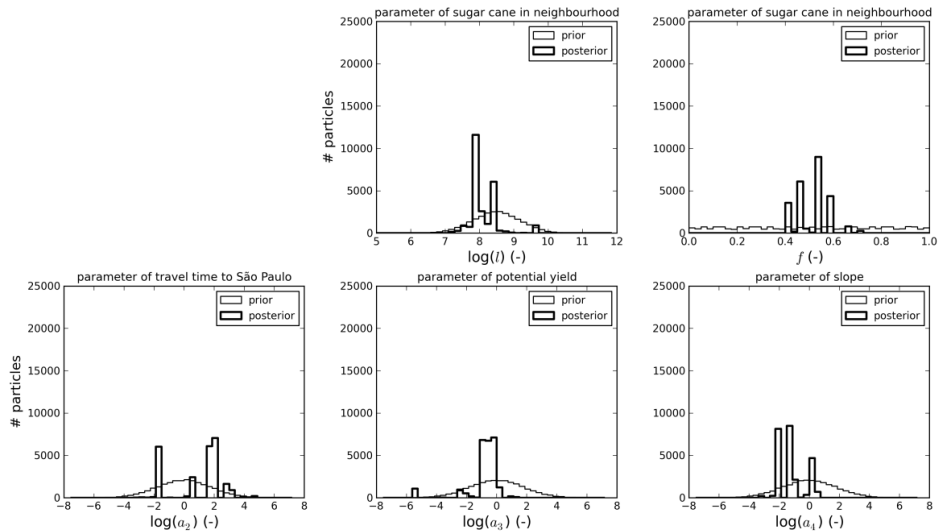
### **3.5. Acknowledgements**

This work was carried out within the BE-Basic R&D Program, which was granted a FES subsidy from the Dutch Ministry of Economic affairs, agriculture and innovation (EL&I). We are grateful to the National Institute for Space Research in Brazil (INPE) for providing the Canasat maps that were used as observational data.

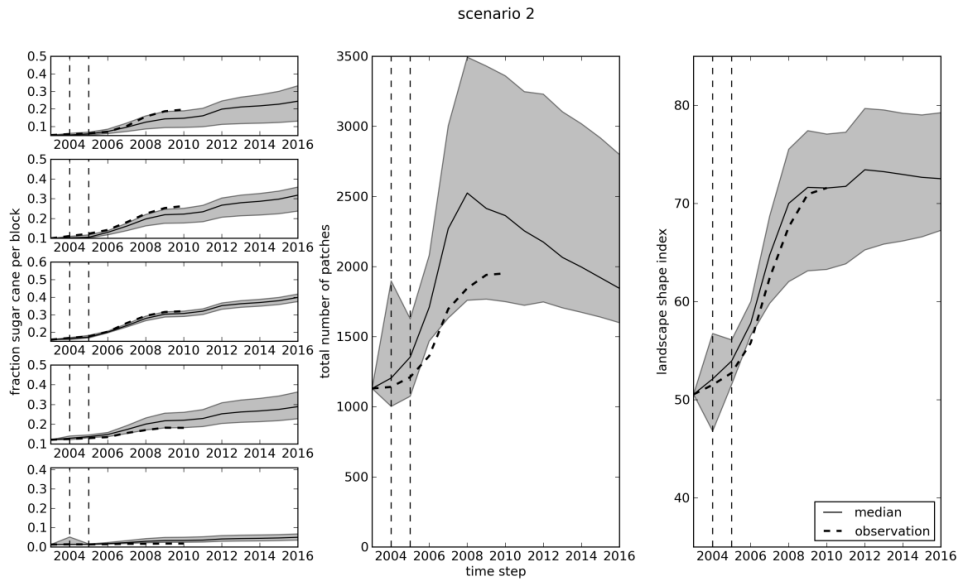
### 3.6. Appendix A



**Figure 3.13: Evolution of the weights  $w_k$  of the candidate suitability factors over time for the new block division, with filtering in 2004 and 2005. For each suitability factor the black horizontal lines in the centre panel are a random selection of 10% of the particles, the bars on the left represent the prior distribution, and the bars on the right represent the posterior distribution of the weight.**



**Figure 3.14: Prior (thin line) and posterior (thick line) distributions of all parameters for the new block division.**



**Figure 3.15: Comparison of modelled medians (solid lines), 95% confidence intervals (grey areas) and observations (dashed lines) of the spatial metrics for the validation blocks for scenario 2 (particle filter). The filter time steps are indicated with vertical dashed lines.**

#### **4. Detecting systemic change in a land use system by Bayesian data assimilation**

**Judith A. Versteegen, Derek Karssenbergh, Floor van der Hilst, André P.C. Faaij (2016), *Environmental Modelling & Software*, 75, 424-438.**

**Abstract** - A spatially explicit land use change model is typically based on the assumption that the relationship between land use change and its explanatory processes is stationary. This means that model structure and parameterization are usually kept constant over the model runtime, ignoring potential systemic changes in this relationship resulting from societal changes. We have developed a methodology to test for systemic changes and demonstrate it by assessing whether or not a land use change model with a constant model structure is an adequate representation of the land use system given a time series of observations of past land use. This was done by assimilating observations of real land use into a land use change model, using a Bayesian data assimilation technique, the particle filter. The particle filter was used to update the prior knowledge about the model structure, i.e. the selection and relative importance of the explanatory processes for land use change allocation, and about the parameters. For each point in time for which observations were available the optimal model structure and parameterization were determined. In a case study of sugar cane expansion in Brazil, it was found that the assumption of a constant model structure was not fully adequate, indicating systemic change in the modelling period (2003-2012). The systemic change appeared to be indirect: a factor has an effect on the demand for sugar cane, an input variable, in such a way that the transition rules and parameters have to change as well. Although an inventory was made of societal changes in the study area during the studied period, none of them could be directly related to the onset of the observed systemic change in the land use system. Our method which allows for systemic changes in the model structure resulted in an average increase in the 95% confidence interval of the projected sugar cane fractions of a factor of two compared to the assumption of a stationary system. This shows the importance of taking into account systemic changes in projections of land use change in order not to underestimate the uncertainty of future projections.

## 4.1. Introduction

Land use change (LUC) is the result of complex interactions between socio-economic and environmental processes (Verburg, 2006, Brown et al., 2008). To simulate potential development pathways in the land use system, scenario storylines are used in combination with land use change models. Various modelling approaches have been designed for this. Several such approaches are founded on the conceptual distinction between 1) the quantity of change per land use type, also called demand, and 2) the spatial allocation of this change (Pontius Jr. and Neeti, 2010). The quantity of change can be seen as a model input, because it is dictated by the scenario storyline, and the spatial allocation of change is defined by the model structure and parameters. The model structure consists of a set of suitability factors that serve as proxies for the land use system being socio-economic and environmental processes that regulate the location of change (Schaldach et al., 2011, e.g., Verburg et al., 2002, van der Hilst et al., 2012), such as topography, accessibility, and potential revenues, and the way these factors interrelate.

Although there are some exceptions (e.g., Clarke et al., 1997, Carlson et al., 2012), the selection, relative importance and parameterization of the suitability factors, i.e. the model structure and its parameters, are in current applications usually kept constant over model runtime. A crucial assumption, implicit in this method, is that the relationship between LUC and its explanatory processes is stationary (Manson, 2007). This assumption ignores potential systemic changes in this relationship resulting from societal changes including technological, political or economic developments. A systemic change is a fundamental change in system structure. Because the notion of 'fundamental' is subjective, we recognize systemic change in the context of models by: "a system state change that cannot be simulated using a constant model structure and/or parameterization". This definition is further explained in section 2.1. Our aim is to develop a general methodology, applicable to any type of model, to test for systemic change given this definition. We demonstrate this methodology in a case study with a land use change model and try to answer the following questions: 1) Is the assumption of a spatially explicit LUC model with a constant model structure and parameters, as generally used in the land use change community, an adequate representation of the land use system, or do observations of past land use over time indicate systemic changes? 2) If systemic changes occur, can these be related to known societal changes? 3) How does the inclusion of systemic changes in the model affect model projection uncertainty?

Evaluation of the stationarity of the relationship between land use and a set of spatial attributes has been done by others (Aspinall, 2004, Bakker et al., 2011, Bakker and Veldkamp, 2012). These studies use logistic regression, separately from



the land use change model. Therefore, they do not gain information on how to implement the (either changing or constant) spatial attributes into the model, in other words, how to turn these attributes into suitability factors, which restricts their value for the challenge of modelling systemic change. In addition, they often do not take into account uncertainty in the model and/or observational data and compare only two points in time, with the exception of Bakker et al. (2011), who compare three points in time.

To overcome these restrictions, we assimilate a time series of observations of real land use into a spatially explicit LUC model to find the best model configuration for different points in time. A similar approach has been demonstrated in the field of hydrology by Merz et al. (2011). Here, we use the particle filter (van Leeuwen, 2009), which is a Bayesian estimation or data assimilation technique. A particle filter updates the prior knowledge about the model structure and parameters during model runtime at points in time for which observations are available. In this way we assess the land use system structure, or model structure, as a whole, instead of only its components independently, in a fashion similar to our previous study (Chapter 3). Unlike in this previous study, we apply the particle filter here separately for each year for which a land use map is available. By following this approach, optimal model structures and parameterizations, can be obtained for these different points in time. This allows us to create a time series of the evolution of the model structure and parameters. Two stationarity tests, a distribution comparison test and the Runs test (Wald and Wolfowitz, 1940), are used to assess deviations in this time series and to check whether these deviations can be attributed to randomness or not. If not, this indicates systemic change. An important advantage of the particle filter compared to, for example, logistic regression is that it provides posterior knowledge including uncertainty, which enables providing confidence intervals for the identified model structures and the associated land use change projections.

We have set up a spatially explicit land use change model simulating sugar cane cropland expansion in the state of São Paulo in Brazil for the period 2003 to 2012, for which a time series of sugar cane occurrence maps of high quality (Rudorff et al., 2010, Adami et al., 2012a) is available as observational data. This case is suitable for testing our approach, because there are a number of societal changes in Brazil in the studied period that might have caused systemic change in the sugar cane expansion patterns, e.g., the economic crisis in 2008, and the adaptation of the Forest Code in 2012. We consider four suitability factors that could potentially be of importance in the spatial allocation of new sugar cane fields: sugar cane in the neighbourhood, distance to the sugar cane processing mills, potential yield, and slope. First, we use a synthetic dataset generated by the model for the years 2004-2012 and demonstrate that the particle filter can reproduce the model

structure and parameters which were applied to generate this synthetic dataset. Second, we use the real observations as observational data, to find the best fitting model structure and parameters for the real system for each of the years. Significant variations in optimal model structure and parameters between consecutive years would indicate that the changes in these years cannot be simulated using the stationarity assumption, and therefore signify systemic changes. In our study we try to relate the changes in model structure over time to the societal changes identified beforehand.

Next, sugar cane expansion is projected for the years 2013-2022. For this projection phase, the model is run with model structure and parameters varying over time. The trend in this variation, if any, depends on the connection between societal changes and the variation in model structure and parameters found for the time span between 2004 and 2012. This run is compared to a run with a classical model having a constant model structure and parameters. Differences in system state behaviour and uncertainty are evaluated.

The next section explores the concept of systemic change in the context of models, and provides explanations of the land use change model, the particle filter technique, and the stationarity analysis. Section 3 describes the case study, mentions potential causes of systemic changes in the case study area, and delineates the different model runs. Section 4, and 5 and 6 are the results, discussion, and conclusion sections.

## 4.2. Methods

### 4.2.1. Systemic change

In the introduction, systemic change was defined as a change in the system indicated by a system state change that cannot be simulated using a constant model structure and/or parameterization. The system state as a function of the model structure can be described as:

$$\mathbf{y}_t = \mathbf{f}(\mathbf{y}_{t-1}, \mathbf{x}_t, \mathbf{p}), \text{ for each } t = 1, 2, \dots, T \quad 4.1$$

In Equation 4.1,  $\mathbf{y}_t$  is the system state at the time step  $t$ ,  $\mathbf{f}$  is the set of transition rules, representing the processes that lead to change in the system state over time, and the way they are implemented and combined, i.e. the model structure. The vector  $\mathbf{x}_t$  represents all inputs, both spatial and non-spatial, and  $\mathbf{p}$  contains the parameters of the transition rules  $\mathbf{f}$ . In the case of a spatially explicit land use change model the spatial inputs are the input maps to calculate the suitability factors and the non-spatial input is the demand.

Systemic change means that a certain  $\mathbf{f}$  and/or  $\mathbf{p}$ , which was at previous time steps able to give an accurate representation of the change in  $\mathbf{y}_t$ , has to be altered at a certain point in time to remain able to correctly simulate  $\mathbf{y}_t$ . This coincides with the special issue editors' description of systemic change as "entities no longer interrelating in a particular way" or "changes in the set of exogenous variables to which the system is sensitive".

The systemic change visible in  $\mathbf{f}$  and/or  $\mathbf{p}$  can be either direct or indirect. Direct means that an action, e.g. a policy change, directly affects  $\mathbf{f}$  and/or  $\mathbf{p}$ . Indirect denotes that the action has an effect on the inputs  $\mathbf{x}_t$ , in such a way that the transition rules and parameters have to change as well (Filatova and Polhill, 2012). In other words, the behaviour of the inputs over time suddenly changes, beyond the function domain of  $\mathbf{f}$ , with the result that  $\mathbf{f}$  becomes invalid.

#### 4.2.2. Land use change model

The land use change model applied in our study is a branch of the PCRaster Land Use Change model (PLUC) (Chapter 2), a spatially explicit LUC model. It is, like many other land use change models (Pontius Jr. and Neeti, 2010), grounded on the conceptual distinction between 1) the quantity of change per land use type, and 2) the spatial allocation of this change. The total quantity of land required per land use type, also called the demand  $\mathbf{x}_{d,n,t}$ , is an input (present in  $\mathbf{x}_t$  in Equation 4.1), defined by historical data in historical runs, the identification phase, and by the scenario storyline in runs for future land use change, the projection phase. The change in demand, which can denote expansion as well as contraction, is allocated using the total suitability map  $\mathbf{s}_{n,t}$ , a weighted sum of the suitability factors for that land use type:

$$\begin{aligned} \mathbf{s}_{n,t} &= \sum_{k=1}^{K_n} (w_{n,k} \cdot \mathbf{u}_{n,k,t}), \text{ for each } n \text{ in each } t \\ \text{with } \sum_{k=1}^{K_n} (w_{n,k}) &= 1 \\ \text{and } \mathbf{u}_{n,k,t} &= h(\mathbf{x}_{n,k,t}, \mathbf{p}_{n,k}) \end{aligned} \tag{4.2}$$

In Equation 4.2,  $t$  is the time step in years, with  $t = 1, 2, \dots, T$ ;  $n$  is the land use type, with  $n = 1, 2, \dots, N$ ; and  $k$  is the suitability factor, with  $k = 1, 2, \dots, K_n$ . Furthermore,  $\mathbf{u}_{n,k,t} \in [0,1]$  is the suitability map for suitability factor  $k$ ; and  $w_{n,k} \in [0,1]$  is the weight of factor  $k$ , which denotes the importance of this specific proxy in the total suitability map  $\mathbf{s}_{n,t}$ . The suitability factors and their weights together establish the model structure of the LUC model, and are, just like the parameters  $\mathbf{p}_{n,k}$ , temporally and spatially constant as in most land use change models. The function  $h()$  uses the spatial attribute  $\mathbf{x}_{n,k,t}$  and parameter(s)  $\mathbf{p}_{k,t}$  to create the proxy for land use change, and normalizes it, i.e. linearly transforms it to a value between 0

and 1, to obtain the suitability map for suitability factor  $k$ ,  $\mathbf{u}_{n,k,t}$ . The transformation is linear, because the actual shape of the relation (linear, convex, concave) between  $\mathbf{u}_{n,k,t}$  and  $\mathbf{x}_{n,k,t}$  is determined by the parameters  $\mathbf{p}_{n,k}$  within  $\mathbf{u}_{n,k,t}$ , discussed later. If required, areas where expansion is not allowed (no-go areas) can be masked out in the total suitability map, so that no change can occur in these cells.

Two types of suitability factors exist in PLUC: attraction or repulsion factors and feedback effects. Attraction/repulsion factors represent the attracting or repelling effect of a spatial attribute on a land use type. They are defined as:

$$\mathbf{u}_{n,k,t} = \mathit{norm}(c_{n,k} \cdot \mathbf{x}_{n,k,t}^{a_{n,k}}), \quad \text{for } k \in \text{attraction/repulsion} \quad 4.3$$

In Equation 4.3,  $\mathbf{x}_{n,k,t}$  is the attribute of suitability factor  $k$ , e.g., potential yield. The parameter  $a_{n,k}$  determines the shape of the suitability function. A value of 1 results in a linear function, meaning that suitability increases or decreases linearly with the increase in the value of the attribute. For  $0 < a_k < 1$ , the shape of  $\mathbf{u}_{n,k,t}$  is concave, and for  $a_{n,k} > 1$ , the shape is convex. Whether a certain attribute attracts (higher attribute values lead to a higher suitability) or repels (lower attribute values lead to a higher suitability) land use type  $n$ , is determined by the constant  $c_{n,k}$ , having a value of 1 in case of attraction and -1 in case of repulsion. The function  $\mathit{norm}()$  normalizes its contents, so that  $\mathbf{u}_{n,k,t} \in [0,1]$ .

Feedback effects characterize the positive effect of the presence of a land use type on the allocation of another land use type in a pre-defined neighbourhood. They represent temporal feedback in the system because the land use updated in the previous time step,  $\mathbf{y}_{t-1}$ , is used as an input, and thereby generate non-linear system behaviour. They are calculated as:

$$\mathbf{u}_{n,k,t} = \mathit{norm}(\mathbf{x}_{n,k,t-1}^2 + 2 \cdot f_{n,k} \cdot l_{n,k,w}^2 \cdot \mathbf{x}_{n,k,t-1}), \quad \text{for } k \in \text{feedback} \quad 4.4$$

In Equation 4.4,  $\mathbf{x}_{n,k,t-1}$  is the number of neighbours of the land use type of interest (usually land use type  $n$  itself) in the neighbourhood window in time step  $t - 1$ , derived from the system state in the previous time step  $\mathbf{y}_{t-1}$ . Parameter  $l_{n,k,w}$  is the window length (in number of cells) of the window that determines whether or not a cell belongs to the neighbourhood. Furthermore,  $f_{n,k}$  is the ‘preferred’ fraction of neighbours of the land use type of interest appearing in the total number of neighbours, that is  $l_{n,k,w}^2$ , within the window. Equation 4.4 creates a parabolic shape of  $\mathbf{u}_{n,k,t}$  against  $\mathbf{x}_{n,k,t-1}$ . The rationale for this, is that e.g., financial and policy related principles can lead to a specific optimal number of neighbours. More as well as less neighbours than this optimum reduces the suitability  $\mathbf{u}_{n,k,t}$ . For example, for farmers it can be advantageous when some of their neighbours cultivate the same crop, for they can share machinery, but it can

also be disadvantageous, as the land price increases due to this favourable situation under a growing demand.

Allocation of the demand  $\mathbf{x}_{d,n,t}$  occurs at each  $t$  by sorting  $\mathbf{s}_{n,t}$  and allocating cells to land use type  $n$  until  $\mathbf{x}_{d,n,t}$  has been fulfilled. This mechanism (Equation 4.1 plus the allocation) is  $\mathbf{f}$  in Equation 4.1. Systemic change concerning  $\mathbf{f}$  can happen through a required change in  $\mathbf{w}_{n,k}$ . The land use map resulting from applying  $\mathbf{f}$  is  $\mathbf{y}_t$  in Equation 4.1. The inputs  $\mathbf{x}_t$  in Equation 4.1 are  $\mathbf{x}_{d,n,t}$  and  $\mathbf{x}_{n,k,t}$  or  $\mathbf{x}_{n,k,t-1}$ , depending on  $k$ . The parameters  $\mathbf{p}$  are  $a_{n,k}$ ,  $l_{n,k,w}$ , and  $f_{n,k}$ .

We do not use a single, deterministic model run to simulate land use change, but an ensemble of runs, a Monte Carlo simulation (Aerts et al., 2003b, Chapter 2). For land use type  $n$ , in each cell, in each year, this ensemble represents the probability distribution of the occurrence of  $n$ . The ensemble is created by sampling from the prior probability distributions of the weights ( $\mathbf{w}_{n,k}$ ) and parameters ( $\mathbf{p}_{n,k}$ ) of the suitability factors (Equations 4.2 to 4.4), and running the land use change model for each of the ensemble members  $i$ , with  $i = 1, 2, \dots, I$ .

### 4.2.3. Particle filter

The ensemble of runs represents the range of possible model outcomes in each year given the uncertainty in model structure and parameters. The sequential importance resampling (SIR) particle filter (van Leeuwen, 2009) is a Bayesian estimation technique that uses observations to reduce the uncertainty in the ensemble, in our case to identify the optimal model structure and parameters. At a time step when observations are available, i.e. a filter moment, the particle filter solves Bayes' theorem for each ensemble member  $i$ , also called particle:

$$p(\mathbf{z}_t^i | \mathbf{o}_t) = \frac{p(\mathbf{o}_t | \mathbf{z}_t^i) \cdot p(\mathbf{z}_t^i)}{p(\mathbf{o}_t)} = \frac{p(\mathbf{o}_t | \mathbf{z}_t^i) \cdot p(\mathbf{z}_t^i)}{\sum_{j=1}^N p(\mathbf{o}_t | \mathbf{z}_t^j) \cdot p(\mathbf{z}_t^j)}, \text{ for each } i = 1, 2, \dots, N \quad 4.5$$

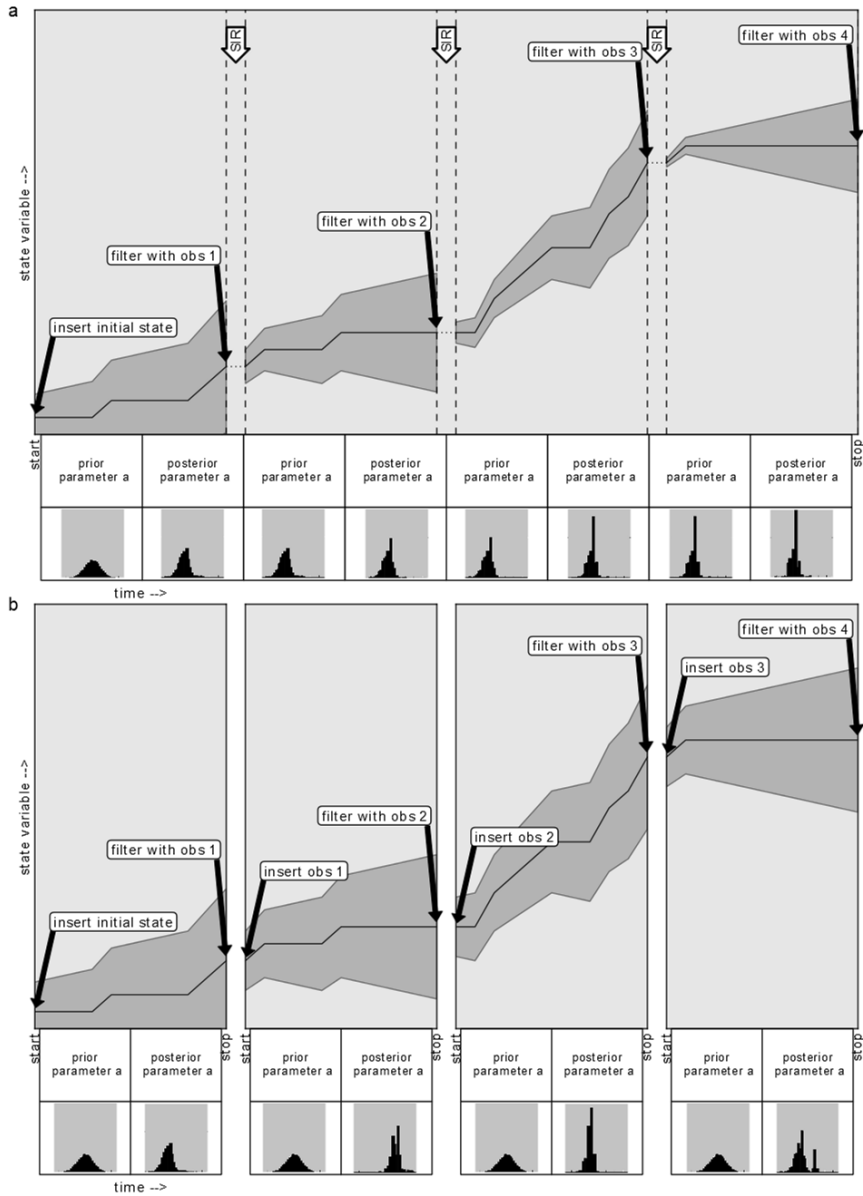
In Equation 4.5,  $p(\mathbf{z}_t^i | \mathbf{o}_t)$  is the posterior probability of the model state  $\mathbf{z}_t^i$  of ensemble member  $i$ . The model state consist of the system states as well as the transition rules, inputs and parameters, i.e.  $\mathbf{z}_t = (\mathbf{y}_t, \mathbf{f}, \mathbf{x}_t, \mathbf{p})$  (see Equation 4.1). Thus, when  $p(\mathbf{z}_t^i | \mathbf{o}_t)$  is updated, not only the system state  $p(\mathbf{y}_t)$  is updated, but the probability distributions of the weights,  $p(\mathbf{w}_{n,k})$ , and of the parameters,  $p(\mathbf{p}_{n,k})$  are updated as well, because they are enclosed in the same ensemble member. Furthermore in Equation 4.5,  $p(\mathbf{z}_t^i)$  is the prior probability of ensemble member  $i$ , and  $p(\mathbf{o}_t)$  is the probability distribution of the observations, i.e. the measurement data and their uncertainty. If the observations are not of the state variable, but of a derived measure, the modelled system state  $\mathbf{y}_t^i$  has to be converted to that measure before filtering. The prior probability of ensemble member  $i$  is always

equal to  $1/l$  due to the sequential importance resampling (SIR) strategy (van Leeuwen, 2009). Namely, SIR samples a new set of  $l$  ensemble members after each filter moment, where the probability that an ensemble member is resampled equals the posterior probability of ensemble member  $i$  in that filter moment. Finally,  $p(\mathbf{o}_t | \mathbf{z}_t^i)$  is the probability of the observations given ensemble member  $i$ . Under the assumption that the observation error has a Gaussian distribution, the latter can be calculated as (van Leeuwen, 2009):

$$p(\mathbf{o}_t | \mathbf{z}_t^i) = e^{-1/2[\mathbf{o}_t - \mathbf{H}(\mathbf{z}_t^i)]^T \mathbf{R}_t^{-1} [\mathbf{o}_t - \mathbf{H}(\mathbf{z}_t^i)]}, \text{ for each } t \quad 4.6$$

In Equation 4.6,  $\mathbf{H}$  is the measurement operator that transforms the model state to the observation, i.e. it selects the modelled system state  $\mathbf{y}_t^i$  from  $\mathbf{z}_t^i$  and, if necessary, converts it to the same support as the observations of the system state.  $\mathbf{R}_t$  is the covariance matrix of the observation error and  $T$  indicates matrix transposition. The diagonal elements of  $\mathbf{R}_t$  represent variance of the observation error,  $\sigma_{o,t}^2$ . The off-diagonal elements of  $\mathbf{R}_t$  are relevant only when observation errors are correlated over space and/or time, otherwise they are zero.

We apply Equations 4.5 and 4.6 in two different ways (Figure 4.1). The traditional way to use the particle filter is sequentially (see Figure 4.1a) in order to, in the end, obtain  $p(\mathbf{z}_T | \mathbf{o}_T)$ , the posterior of the model state at the final time step  $T$ . We have used this method before to identify model structure (Chapter 3). This approach assumes a constant model structure and constant parameters: observations of the system state at a certain point in time are used as *additional* information about the best-fit model structure and parameters. So, although the number of ensemble members remains the same, due to the SIR, the variation in the ensemble members in terms of their uniqueness in parameter values will typically diminish over time. This means that the probability distributions of these parameters are gradually narrowed. This, however, also means that the approach does not work when parameters 'in reality' change over time, because there is no stationary parameter value the model state can converge towards. Therefore, this approach is suitable only under the assumption of model structure and parameter stationarity, i.e. no systemic changes.



**Figure 4.1: Functioning of the particle filter, (a) in the traditional approach with sequential importance resampling (SIR), and (b) in the approach used to assess the presence of systemic changes. ‘Obs 1’ means observations at filter moment 1, the solid dark grey line indicates the median system state, grey areas represent the confidence interval. Histograms underneath the plots illustrate the effect of the filter moments on a general parameter  $a$ . In panel a, the prior of parameter  $a$  at filter moment  $t$  is always equal to the posterior of parameter  $a$  at filter moment  $t-1$ . In panel b this relation is not present; all priors are the same and not dependent on any of the posteriors.**

Because in this paper we want to validate that assumption, we also apply Equations 4.5 and 4.6 in an atypical way (Figure 4.1b). Observations of the system state at a certain point in time are used as *distinctive* information about the best-fit model structure and parameters at that point in time. Equations 4.5 and 4.6 are not applied sequentially, but separately at each filter moment. So, the model is initiated with the system state provided by  $\mathbf{o}_{t-x}$  and is run up to the next filter moment, obtaining  $p(\mathbf{z}_{t-x,t}|\mathbf{o}_t)$ , where  $x$  is the number of time steps between model start and filter moment. Next, the model is initiated with the system state provided by  $\mathbf{o}_t$ , with all parameters set to their initial prior probability distribution again, and the model is run up to the next filter moment. Using this approach, we find a distinctive  $p(\mathbf{z}_{t-x,t}|\mathbf{o}_t)$  at every filter moment, not limited by previous observations, valid for period  $t - x$  to  $t$  (one subplot in Figure 4.1b). Hence, we can explore whether the optimal model structure and/or parameters vary significantly over time. If one of them does, the various observed system states cannot be simulated using a constant model structure and/or parameters, reflecting systemic change.

#### 4.2.4. Stationarity analysis

Of course it can be viewed in a qualitative way whether the posterior probability distributions of the weights,  $p(\mathbf{w}_{n,k})$ , and of the parameters,  $p(\mathbf{p}_{n,k})$  vary notably over time, but it is more objective to use a quantitative test. Two general approaches exist to quantitatively test for stationarity: parametric and non-parametric tests (Grazzini, 2012). Parametric tests rely on the assumption that the distribution of the variable being sampled is known. Because this is often not the case for models of complex systems, like land use change models, a non-parametric test is more applicable (Grazzini, 2012, Lanzante, 1996).

First, a non-parametric distribution comparison test is applied to check to what extent the distribution of a parameter at time  $t$  differs from the average distribution over all other time steps  $T^*$ . In this test we first calculate a random variable  $\mathbf{d}_{n,k,t}$ , which represents the difference between the posterior distribution of a parameter at a particular time step and the mean of the parameter over all time steps. This is done in an approach similar to bootstrapping (Efron and Tibshirani, 2003), by subtracting a value randomly taken from the posterior distribution of a parameter at time step  $t$  from a value randomly taken from the average distribution of this parameter over  $T^*$ . This is done  $l$  times, i.e. as many times as we have Monte Carlo samples. If the posterior distribution of a parameter at time  $t$  and the average distribution of this parameter over  $T^*$  are the same, the resulting distribution  $p(\mathbf{d}_{n,k,t})$  is centred on zero. So, under the null hypothesis of stationarity ( $H_0$ ),  $p(\mathbf{d}_{n,k,t})$  has an expected value of zero. To test this hypothesis,



taking into account the full distribution of the parameter, we infer whether or not zero falls within the confidence interval of  $p(\mathbf{d}_{n,k,t})$ :

$$H_0: Q_{\alpha/2}(p(\mathbf{d}_{n,k,t})) < 0 < Q_{1-(\alpha/2)}(p(\mathbf{p}_{n,k,t})) , \text{ for each } t \quad 4.7$$

In Equation 4.7,  $Q_{\alpha/2}(p(\mathbf{d}_{n,k,t}))$  is the  $\alpha/2$  percentile of  $p(\mathbf{d}_{n,k,t})$ , e.g. the 2.5<sup>th</sup> percentile when verifying if  $\alpha < 0.05$ . If the  $H_0$  is rejected, there is systemic change in time  $t$ . The same is done for the weights.

The test above indicates the probability of systemic change in a certain parameter at a certain time step, but does not tell anything about the temporal correlation of potential deviations. To investigate this temporal correlation we use the non-parametric Wald-Wolfowitz test, also called Runs test (Wald and Wolfowitz, 1940). This test was successfully applied before to test for stationarity in a complex system model by Grazzini (2012). Given a time series and a function that aims to explain the trend in these time series, the values in the time series should be randomly distributed above and below the function, uncorrelated over time, if the function gives an adequate description of the time series. This is true regardless of the shape of the error distribution in the time series. Values above the function are labelled as + and values below the function are labelled as -. In the Runs test, a run is defined as a sequence of identical instances, i.e. either pluses or minuses. For example, the series +, +, +, -, +, -, -, +, + contains five runs, namely one run of three pluses, followed by a run of one minus, etc. In a time series containing a number of  $N_{+-}$  values, with random temporally uncorrelated errors, the expected, or mean, number of runs,  $\mu_r$ , is (Wald and Wolfowitz, 1940):

$$\mu_r = \frac{2N_+N_-}{N_{+-}} + 1 \quad 4.8$$

And the variance of  $\mu_r$  is:

$$\sigma_r^2 = \frac{2N_+N_-(2N_+N_- - N_{+-})}{N_{+-}^2(N_{+-} - 1)} \quad 4.9$$

In Equation 4.8 and 4.9,  $N_+$  is the number of plus instances,  $N_-$  is the number of minus instances. Using this probability distribution of the number of runs, the Runs test checks whether the null hypothesis of randomness can be rejected or not. In other words, it checks whether the distribution of values in the time series above and below the function can be considered random, uncorrelated over time, given a certain level of significance  $\alpha$ . If it is not, the function is not an appropriate representation of the time series. When one uses a constant value for this function, one can test whether this constancy, of e.g. a parameter, is a correct representation of this parameter in the studied system. If it is not, there is a systemic change in our definition (section 4.2.1).

### 4.3. Case study

#### 4.3.1. Background

A case study on sugar cane cropland expansion in São Paulo state is set up to examine the presence of systemic changes in a land use system. Expansion of the sugar cane area is related to an increasing demand for both sugar and ethanol. Brazilian ethanol production from sugar cane is seen as one of the most efficient biofuel technologies currently available, and it is profitable without subsidies, so in the future the expansion is expected to continue (Walter et al., 2011, Sparovek et al., 2009, Cerqueira Leite et al., 2009). Sugar cane cultivation in Brazil is concentrated in the South-Central region (Rudorff et al., 2010). This region, which includes São Paulo state, has been annually mapped since 2003 in the Canasat project (Rudorff et al., 2010) with an overall thematic accuracy of 98% (Adami et al., 2012a), providing a reliable time series of observational data for the particle filter.

Within the period between 2003 and 2012, a number of societal changes in Brazil might have caused systemic change in the sugar cane expansion. The review provided here is a reflection of societal changes that were deemed of importance by the authors after consultation of literature and experts, and should therefore not be considered exhaustive. We focus on developments that potentially change the spatial allocation of sugar cane, not the quantity of change (demand) as in PLUC the demand is exogenous (section 4.2.2). We provide hypotheses of how they might change the model structure or parameters when this is not immediately clear.

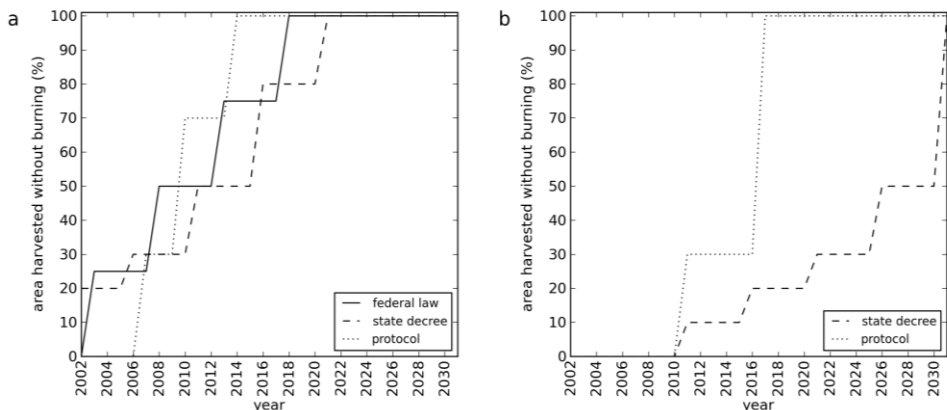


Figure 4.2: Schemes specifying the required area without pre-harvest burning as a proportion of the total sugar cane area (a) for slopes  $\leq 12\%$  and (b) for slopes  $> 12\%$ .

Some policy changes in the last decade may have had, and are in fact established to have, an effect on the land use system. According to Sparovek et al. (2010) “Land-use modellers exploring the Brazilian case generally pay little attention to the influence of legal aspects, i.e., how Brazilian regulations influence agriculture, including the size and spatial distribution of the expansion potential”. Therefore we discuss these aspects in detail, and explore whether they can be traced as systemic changes in our case study. Firstly, sugar cane straw is currently often burned before harvesting the sugar cane to improve the safety of the ‘cane cutters’ and to increase the yield. The Brazilian government tries to eliminate pre-harvest burning because it has negative effects on human health and on the environment due to the emission of pollutant gases (Aguiar et al., 2011). Replacing manual harvesting by mechanical harvesting can eliminate pre-harvest burning, because mechanical harvesting does not require burning. However, the harvest machines cannot operate on sloping ground; 12% is considered the maximum slope for mechanical harvesting (Macedo, 2007). Therefore, schemes are established, specifying the maximum area on which pre-harvest burning can be practiced as a proportion of the total sugar cane area, per slope category per year (Figure 4.2). With the most recent of the three schemes, the Green Ethanol protocol, compliance is not obligatory, but can be advantageous for the producers because the protocol resembles importer’s preferences and offers a first step towards certification. By 2008, 145 out of the 177 ethanol plants in São Paulo complied to the protocol, representing 89% of total cane crushing (Lucon and Goldemberg, 2010). However, a study by Aguiar et al. (2011) for the years 2009 and 2010 shows that in some areas the harvesting system was shifted in the wrong direction in these years, i.e. from green harvest to harvest with burning.

The second important policy change is the initiation of sugarcane agro-environmental zoning (AEZ) in São Paulo in 2008 (Lucon and Goldemberg, 2010). Using eight physical indicators (climate, surface water, slope, ground water, biodiversity protection areas, biodiversity connectivity, and integral protection units) a map with four categories is created: suitable, moderately restricted, highly restricted, and unsuitable for sugar cane cultivation (Padua Junior et al., 2012). The initiative of the state zoning has led to the launch of the federal Sugarcane Agro-ecological Zoning (ZAE Cana) in 2009 (Lucon and Goldemberg, 2010). It uses similar physical indicators, and in addition aims to protect the Amazon and Pantanal biomes and Upper Paraguay River Basin. The initiation of these zonings might increase the importance of physical suitability factors in the model structure.

A third policy change that might have had an effect on the spatial allocation of sugar cane is the adaptation of the Forest Act, or Forest Code, in 2012. The Forest Act, established in 1965, is the main legal framework in Brazil for natural vegetation (not only forest) conservation. Among other things, it specified the

fraction of the farmland that should be set aside for biodiversity conservation, meaning that natural vegetation should be kept in place. This fraction was 35% inside the Legal Amazon, and 20% outside. A low compliance in the past and several amendments of the Forest Code between 1965 and 2012 that allowed farmers to sometimes preserve lower fractions of farmland, had put a large part of the Brazilian farmers in an illegitimate situation (Sparovek et al., 2012). As this illegality was a national and international (certification) problem and total compliance with the prescribed fractions through vegetation restoration would be too costly, a revision of the Forest Act was accepted in 2012. The exact rules in the new Forest Act remain vague up to this point in time. It is expected that it maintains the preservation requirements for future expansion, but legalizes the farmers' situation for those who deforested illegally before 2008 (Costa and Gray, 2011). The story goes that, as the process of the revision started already in 2009, farmers anticipated on the new Forest Code since then by accelerated illegal forest cutting, hoping that amnesty would be granted, but, as far as we know, this is not scientifically proven. In the model the new Forest Code might affect the neighbourhood suitability factor or the parameters herein that determine the 'preferred' fraction of neighbours,  $f_{n,k}$  (Equation 4.4), as a percentage of the land cannot be used for cultivation.

Economically, many developments have taken place affecting the sugar cane sector, and it goes beyond the scope of this paper to discuss them in detail. Obviously, the crisis in 2008 may have affected the sugar cane demand (model input) but also the spatial distribution. The crisis led to a discontinuation of investments, forcing farmers to produce at older, less productive sites (Gómez Jr., 2013), and to postpone modernization of agricultural machinery (Aguiar et al., 2011). The latter might cause farmers to care less about the slope of the field, since they cannot afford machinery to harvest mechanically anyway.

Finally, a shift not in the human but in the environmental system that could affect the allocation of sugar cane is that Brazil has experienced some bad harvests between 2009 and 2011. Aguiar et al. (2010) report that from season 2006/2007 to 2008/2009 the area of sugar cane left unharvested has gone up from 3.0 to 4.1 to 11.6%. They believe that this was related to unfavourable harvest weather conditions as well as delays in constructing planned mills, with the result that the mills were not operational in time. In 2009/2010, the unharvested area rose to 18.1%, thereafter decreasing to 6.9% and 1.0%, reaching 4.1% in the season 2012/2013 (Aguiar et al., 2011, Aguiar, Personal communication, July 17th 2014).

### 4.3.2. Model setup

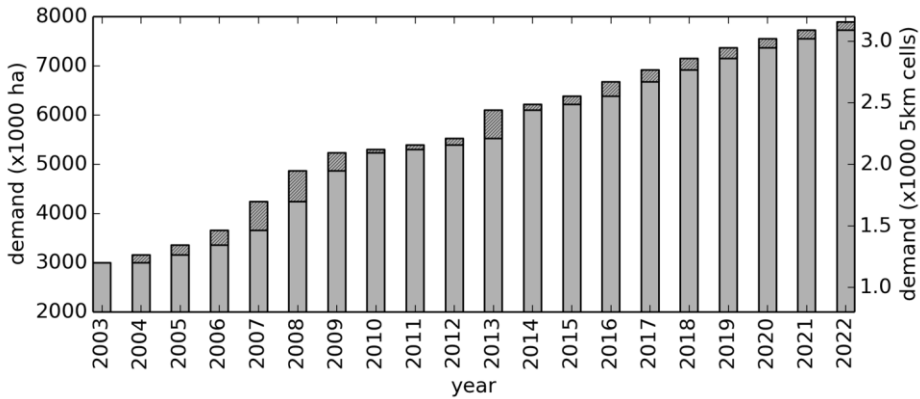
The land use change model described in section 4.2.2, is applied to the São Paulo case study as follows (Table 4.1). Sugar cane ( $n = 1$ ) is the only 'active' land use type, i.e. a land use type responding to a demand, thus  $N = 1$ . The suitability map for sugar cane expansion  $\mathbf{s}_{1,t}$  (Equation 4.2) is calculated using four suitability factors ( $K_1 = 4$ ), derived from discussions with experts and literature review (Lapola et al., 2010, Walter et al., 2011, Rudorff et al., 2010, Macedo and Seabra, 2008, Sparovek et al., 2007, Sparovek et al., 2012, de Souza Soler and Verburg, 2010, Aguiar et al., 2011, Adami et al., 2012a) (Table 4.1). A similar model setup has been calibrated before with the particle filter, resulting in a minimum reduction factor of 3 in the root mean square error of three spatial metrics compared to the reference model (Chapter 3). Sugar cane in the neighbourhood ( $k = 1$ ) is expected to be important because larger plantations require less investment costs per hectare as equipment and infrastructure can be shared. The distance from the field to the sugar cane mill ( $k = 2$ ) determines the transportation costs of sugar cane to the processing unit. The distance to the mill is expected to be more eminent than the distance from the mill to the distribution centre, because the end product (ethanol or sugar) has a higher energy density than the sugar cane and thus lower transport costs per energy unit. Potential yield ( $k = 3$ ), an indicator linking agro-climate conditions to crop requirements, is important for the potential revenues per hectare. Slope ( $k = 4$ ) defines the potential for sugar cane harvest mechanization (see section 4.3.1). Table 4.1 gives more details about the parameterization of the processes described above.

The total area of São Paulo state is about 250000 km<sup>2</sup> and a resolution of 5 km is used. We purposefully model at a resolution larger than the average farm size in Brazil. The land tenure system in Brazil is complex and includes farms managed by the mill, farms held by the mill and leased to a farmer, farms held by the farmer with a contract obligation to deliver to a certain mill for a fixed price, farms held by a farmer having no such contract, and other constructions (see e.g., Sparovek et al., 2007). Spatial data on this with a complete coverage is not available so modelling at farm level would have limited value as the inputs of such a model are highly uncertain. Apart from this, a much finer resolution would increase model run time which hampers the application of the Monte Carlo based particle filter. By aggregating to a cell size of several farms we hope to average out the effects related to the land tenure situation of the individual farmer and focus more on general sugar cane expansion trends, resulting in a model with a relatively short run time.

**Table 4.1: Suitability factors for sugar cane in São Paulo, including type of suitability factor used (attraction/repulsion or neighbourhood effect), the probability distributions of the parameters ( $p_{1,k,t}$ ), and the data sources. The probability distributions of the stochastic variables, e.g.,  $Z_{w_{1,k}}$ , represent prior distributions; during filtering they change.**

$k$	1	2	3	4
<b>suitability factor</b>	sugar cane in neighbourhood	distance to sugar cane mills	potential yield	slope
<b>process to represent</b>	economies of scale	transportation costs to processing units	profits	mechanization potential
<b>weights (<math>w_{1,k}</math>) (Equation 4.2)</b>	$w_{1,1} = \frac{Z_{w_{1,1}}}{\sum_{k=1}^4 (Z_{w_{1,k}})}$ ,	$w_{1,2} = \frac{Z_{w_{1,2}}}{\sum_{k=1}^4 (Z_{w_{1,k}})}$ ,	$w_{1,3} = \frac{Z_{w_{1,3}}}{\sum_{k=1}^4 (Z_{w_{1,k}})}$ ,	$w_{1,4} = \frac{Z_{w_{1,4}}}{\sum_{k=1}^4 (Z_{w_{1,k}})}$ ,
	with $Z_{w_{1,k}} \sim U(0,1)$			
<b>type of suitability factor</b>	neighbourhood effect (Equation 4.4)	attraction/repulsion (Equation 4.3)	attraction/repulsion (Equation 4.3)	attraction/repulsion (Equation 4.3)
<b>parameters (<math>p_{1,k,t}</math>) (Equations 4.3 and 4.4)</b>	$l_{1,1,w} = e^{z_l}$ , $Z_{l_w} \sim N(9.6,0.7)$ $f_{1,1} = Z_{f_{1,1}}$ , $Z_{f_{1,1}} \sim U(0,1)$	$a_{1,2} = e^{z_{a_{1,2}}}$ , $Z_{a_{1,2}} \sim N(0,1.8)$ $c_{1,2} = -1$	$a_{1,3} = e^{z_{a_{1,3}}}$ , $Z_{a_{1,3}} \sim N(0,1.8)$ $c_{1,3} = 1$	$a_{1,4} = e^{z_{a_{1,4}}}$ , $Z_{a_{1,4}} \sim N(0,1.8)$ $c_{1,4} = -1$
<b>original map attribute for <math>x_{1,k,t}</math></b>	sugar cane	location of mills	potential yield	digital elevation model
<b>map source</b>	Rudorff et al., 2010	Picoli, 2013	Tóth et al., 2012	Farr et al., 2007

The model is run in two phases: an identification phase, in which we identify the relative importance of the suitability factors and their parameterization (Table 4.1) and try to recognize systemic changes; and a projection phase, in which we use this information to propagate the land use change (further explained in section 4.3.4). In total  $T$  is 20 time steps, with  $t = 1$  representing the year 2003. The initial land use map (a Boolean map: 1 = sugar cane, 0 = no sugar cane) is the 2003 Canasat map (Rudorff et al., 2010), which has a resolution of 30 m, resampled to the model resolution (5000 m). The resampling is done in such a way that the total sugar cane area in the resampled map matches the total sugar cane area in the original map, i.e. the demanded area for sugar cane is harmonized because this is an important model input. Note that the class 'no sugar cane' is passive: it has no demand and can only change through conversion by the active land use type 'sugar cane'.



**Figure 4.3: Demand per year for the identification phase (2003-2012) and the projection phase (2012-2022). The demand increase with respect to the previous year (area for the model to allocate) is hatched to enhance visibility of variation over time in this demand increase.**

So, the input variable ‘demand’ is in the identification phase simply the total area of sugar cane found in the Canasat maps (Figure 4.3). In the projection phase two data sources are used to construct the demand. The first is the Brazilian Land Use Model (BLUM) (ICONE, 2012, Nassar et al., 2008), an economic partial equilibrium model, and the second is the Brazilian agricultural economics institute, IEA (Torquato, 2006). As we have equal trust in both sources, the demand in PLUC from 2013 to 2022 is the mean of the two time series created from these sources (Figure 4.3).

#### 4.3.3. Particle filter setup

The data assimilation framework in the PCRaster Python framework (Karszenberg et al., 2010) is used for the particle filtering. The data used to create the observational data are nine annual maps of sugar cane occurrence (Rudorff et al., 2010), from 2004 to 2012 (the data of 2003 is used as the initial system state). We, together with others (e.g., Parker et al., 2008), believe that the purpose of land use change models is not, and should not be, to simulate precisely the land use of each single cell in each year. For this reason, we do not use the sugar cane map directly as  $\mathbf{o}_t$  (Equation 4.5), but calculate the fraction of sugar cane in 25 x 25 km blocks and take that as  $\mathbf{o}_t$ . In total the study area consists of 473 of such blocks, making the length of the array  $\mathbf{o}_t$  473. Obviously, we also convert the model output to that measure (fraction of sugar cane in 25 x 25 km blocks) before filtering.

The observational data has two error sources: 1) errors in the classification of the remote sensing image and 2) errors from the upscaling to a larger cell size. We

assume that there is no spatial or temporal correlation in the errors of the observational data, so only the diagonal elements of  $\mathbf{R}_t$  need to be defined, i.e. the variances of the observation error,  $\sigma_{o,t}^2$ .

A study by Adami et al. (2012b) shows that the user's accuracy (the probability that a cell classified as a certain class is actually that class (Lillesand et al., 2003)) of the Canasat data is 0.97 for the sugar cane class and 0.98 for the no-sugar cane class. To obtain the standard deviation,  $\sigma_{o,t,u}$  for the 25 x 25 km blocks belonging to these user's accuracies, we simulate for every potential fill of a block (0-100% sugar cane)  $1 \cdot 10^5$  events where sugar cane cells have a probability of 0.03 to become no-sugar cane, and no-sugar cane cells have a probability of 0.02 to become sugar cane. The results indicate that  $\sigma_{o,t,u}$  is linearly related to the fraction of sugar cane per block, as:

$$\sigma_{o,t,u}^2 = (2.8 \cdot 10^{-2} + 1 \cdot 10^{-2} \cdot \mathbf{o}_t)^2, \text{ for each } t \quad 4.10$$

The upscaling error arises from the fact that in the modelled land use in a cell is either sugar cane or no sugar cane (Boolean), while in the data the fraction of sugar cane per cell is given, leading to a difference between model output and observations. This error has a maximum of 100% (a cell has in reality a sugar cane fraction of 0.5, but the model output is 0 or 1) and a minimum of 0%. The distribution of this error is difficult to estimate, because it depends on the observed values. For this reason we assume normality of this error (required to fulfil the conditions of Equation 4.6) and apply its maximum possible size of  $0.5 \cdot \mathbf{o}_t$  to all observations.

Combining these two error sources, the total variance of the observations is:

$$\sigma_{o,t}^2 = (2.8 \cdot 10^{-2} + 5.1 \cdot 10^{-1} \cdot \mathbf{o}_t)^2, \text{ for each } t \quad 4.11$$

We are aware of the fact that we have made a strong assumption about the shape and magnitude of the variance, but we want to stress that the stationarity test used to detect systemic change (see sections 2.4 and 3.4) employs only the mean of the posterior and not the full distribution. This diminishes the effect of this assumption on the conclusions about systemic change.

#### 4.3.4. Stationarity analysis setup

For the distribution comparison test (Equation 4.7), instead of assuming an  $\alpha$  value and reporting whether or not  $H_0$  is rejected, we calculate the  $\alpha$  belonging to the tipping point between rejection and no rejection, which is more transparent. This  $\alpha$  value gives the probability that  $H_0$  is unjustly rejected, i.e. the probability that we assume systemic change in this parameter at  $t$  while in fact there is stationarity. Hence, low  $\alpha$  values indicate a high probability of systemic change. The Runs test is



applied as follows. For the posterior distributions of the weights ( $w_{n,k}$ ) and parameters ( $\mathbf{p}_{n,k}$ ) of the suitability factors (Equations 4.2 to 4.4), obtained separately for each observation time frame ( $p(\mathbf{z}_{t-x,t}|\mathbf{o}_t)$ , Figure 4.1b), the overall mean is obtained. This overall mean is the function aiming to explain the trend in time series of posteriors. Next, the mean per posterior distribution, i.e. per observation time, is obtained, and assigned a + if it is above the overall mean and a – if it is below. On this sequence the Runs test is applied, and the p-value is reported, indicating the probability that the pattern found in the deviations from the mean is random. If the null hypothesis of randomness is rejected,  $\mathbf{f}$  and/or  $\mathbf{p}$  in Equation 4.1 cannot be considered stationary, so a systemic change is present. Note that the Runs test checks only *if* an average value in the time series is above or below the mean and not *how much* it is above or below. Therefore, the detection of systemic change should be based on the combined results of visual inspection of the means, the distribution comparison test and the Runs test.

#### 4.3.5. Scenarios

Three scenarios are run (Table 4.2). By the word ‘scenario’ we do not mean a scenario storyline, i.e. a potential development pathway of the land use system, but a model setup designed to investigate a specific property of the used method or studied system. All scenarios are run using an ensemble of 5000 members.

In the first scenario, the ability of the particle filter to detect the correct weights ( $w_{n,k}$ ) and parameters ( $\mathbf{p}_{n,k,t}$ ) is tested using a synthetic dataset. The synthetic dataset is created by running the model deterministically with the demand equal to the demand in the Canasat data (Rudorff et al., 2010) (Figure 4.3) and the settings specified in Table 4.3. These settings do not change over time. Next, the model is run stochastically and the particle filter is applied separately for each year (the method shown in Figure 4.1b) from 2004 to 2012 using the synthetic data as observations  $\mathbf{o}_t$ . If the method is working correctly, the distributions of the weights ( $w_{n,k}$ ) and parameters ( $\mathbf{p}_{n,k}$ ) should converge to the values in Table 4.3 in all years, i.e. the particle filter should trace back settings that were applied to generate the synthetic dataset.

**Table 4.2: Scenario setup regarding observational data (synthetic or Canasat (Rudorff et al., 2010)), filter method, and projection method.**

scenario	observational data	filter (2004 - 2012)	projection (2013 - 2022)
1	synthetic	method Figure 4.1b	-
2	Canasat	method Figure 4.1b	using a trend or random posterior from all $p(\mathbf{z}_{t-x,t} \mathbf{o}_t)$ depending on the results
3	Canasat	method Figure 4.1a	using posteriors from 2012

**Table 4.3: Model settings for the synthetic dataset: the weights ( $w_{1,k}$ ) and parameters ( $\mathbf{p}_{1,k}$ ),  $l_{1,k,w}$ , neighborhood window length,  $f_{1,k}$ , neighborhood fill, and  $a_{1,k}$ , suitability function shape parameter, for  $k = 1, 2, 3, 4$ .**

$k$	1, sugar cane in neighbourhood	2, distance to sugar cane mills	3, potential yield	4, slope
weights ( $w_{1,k}$ )	$w_{1,1} = 0.25$	$w_{1,2} = 0.25$	$w_{1,3} = 0.25$	$w_{1,4} = 0.25$
parameters ( $\mathbf{p}_{1,k}$ )	$l_{1,1,w} = 15000$ m $f_{1,1} = 0.5$	$a_{1,2} = 1$	$a_{1,3} = 1$	$a_{1,4} = 1$

In the second scenario, we run a new ensemble, applying the particle filter separately for each year, but now with the Canasat observational data of sugar cane distribution. Potentially, significantly different posterior distributions of the weights and parameters are obtained in each year. If any kind of trend or connection to the societal changes can be detected in distributions, this trend is prolonged in the projection phase (2013-2022). If no trend is apparent, and no connection to societal changes can be found, we assume that for each time step in the projection phase any of the systems found in the identification phase can be valid. So in each projection year (2013-2022) a posterior model state  $p(\mathbf{z}_{t-x,t}|\mathbf{o}_t)$ , containing probability distributions of the weights,  $p(w_{n,k})$ , and of the parameters,  $p(\mathbf{p}_{n,k})$ , is drawn randomly from all posteriors model states of the identification years (2004-2012). The projection phase of scenario 2 is run five times to cover the uncertainty arising from the diverse sequences of posteriors that are drawn.

In the third scenario, the traditional particle filter method with Sequential Importance Resampling (Figure 4.1a) is applied, again with the Canasat data. During the projection phase, the posterior distribution of the final year is used,  $p(\mathbf{z}_T|\mathbf{o}_T)$ , because this posterior contains information from the whole identification phase.

## 4.4. Results

### 4.4.1. Identification with synthetic data

For scenario 1, the means of the posterior distributions of the weights,  $w_{1,k}$ , converge to values around 0.25 (Figure 4.4), as expected since these were the values used to generate the synthetic dataset (Table 4.3). The weight of distance to sugar cane mills,  $w_{1,2}$ , is on average 0.03 too low, 0.22, and the weight of potential yield,  $w_{1,3}$ , is on average 0.03 too high, 0.28. The other two are on average exactly 0.25. No significant trends are visible over time, which is confirmed by the distribution comparison test (average  $\alpha$  values all very high,  $> 0.85$ ) and the Runs test (Table 4.4). The complete posterior distributions of the weights and the posterior distributions of the parameters are given in Appendix A. The distributions of parameters  $a_{1,2}$ ,  $a_{1,3}$ , and  $a_{1,4}$  in all years, and the distributions of the parameters  $l_{1,1,w}$  and  $f_{1,1}$  in all years but 2006 to 2009 remain broad, indicating that the parameters are difficult to identify. The Runs test concludes non-stationarity for four out of the five parameters, but the distribution comparison test does not for any  $\alpha$  below 0.2.

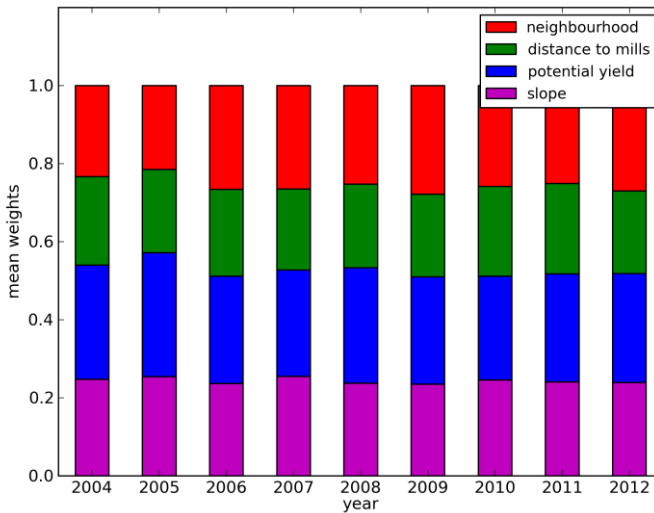
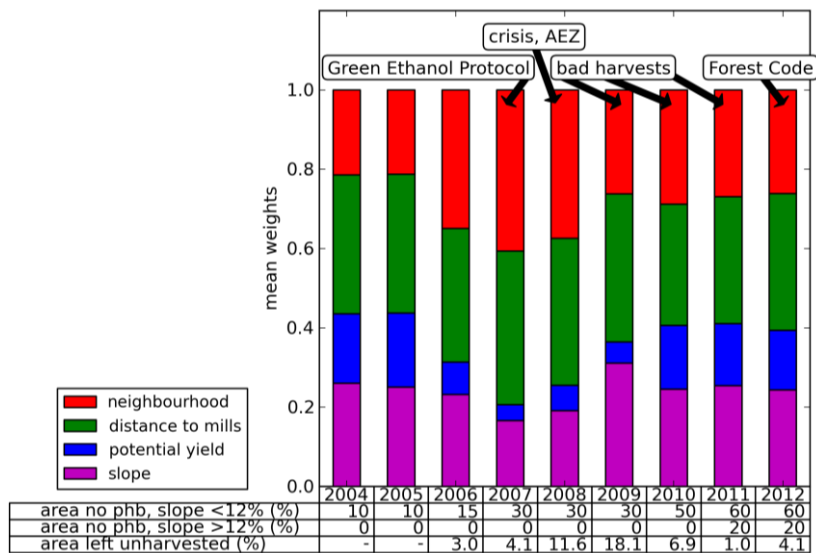


Figure 4.4: Mean of the posterior distributions of the weights of the suitability factors,  $w_{1,k}$ , obtained with synthetic observational data (scenario 1).

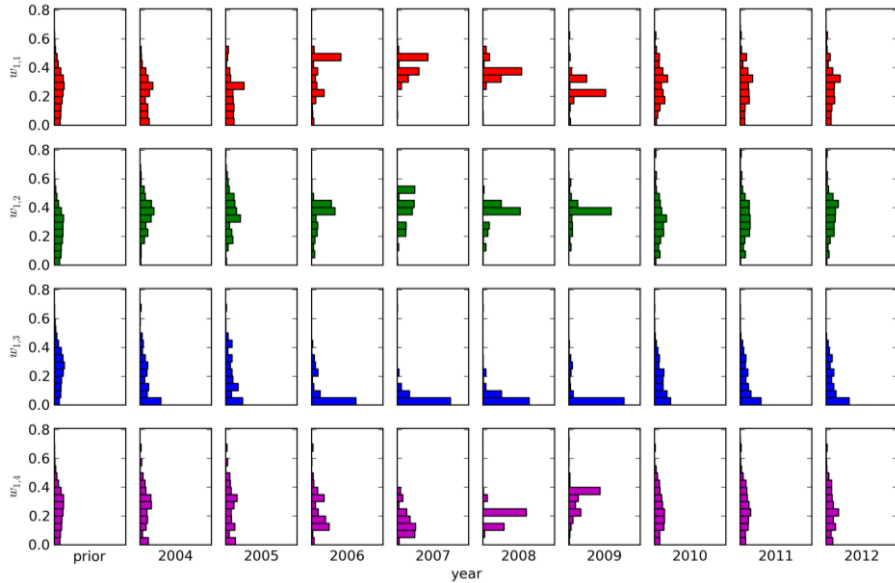
#### 4.4.2. Identification with Canasat data

For scenario 2, the means of the posterior distributions of the weights,  $w_{1,k}$ , (Figure 4.5) imply that all four selected suitability factors are relevant in the land use system, because none of them receives a zero weight. On average, distance to the mills appears to be the most important suitability factor in determining where sugar cane expands. Next most important are sugar cane in the neighbourhood and slope, switching over time between second and third most important. Least important, but still relevant with an average weight of 0.1, is potential yield.

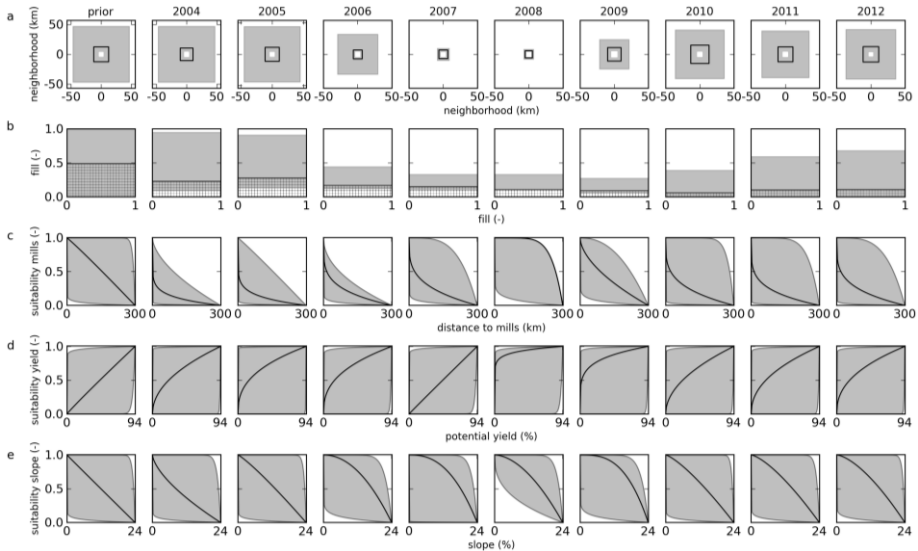
The mean weight of distance to sugar cane mills appears to be stationary (Figure 4.5). This is confirmed by the Runs test, using, for example, a 5% significance level (Table 4.4). The distribution comparison test confirms this; it has the highest average  $\alpha$  value and also the most constant  $\alpha$  value over time, indicating the highest probability of stationarity. The mean weights of the other factors clearly change over time (Figure 4.5). In the period 2006 to 2008, the mean weight of sugar cane in the neighbourhood is 54% higher than in the other years, and the weight of potential yield is 24% lower that persists one year longer. This non-stationarity, indicating systemic change, is confirmed for all three factors by the distribution comparison test, with low  $\alpha$  values in especially 2007 and 2008, and the Runs test, and is strongest for potential yield (Table 4.4).



**Figure 4.5: Mean of the posterior distributions of the weights of the suitability factors,  $w_{1,k}$ , obtained with Canasat observational data (Rudorff et al., 2010) (scenario 2). Occurrences of societal changes, discussed in section 3.1, are indicated above the bar graph. Minimum percentage of sugar cane area that should be harvested without pre-harvest burning (phb) per year (average of the state and Green Ethanol Protocol requirements, Figure 4.2) is indicated below the bar graph, together with the percentage of sugar cane area left unharvested.**



**Figure 4.6: Posterior distributions of the weights of the suitability factors sugar cane in neighbourhood ( $w_{1,1}$ ), distance to mills ( $w_{1,2}$ ), potential yield ( $w_{1,3}$ ), and slope ( $w_{1,4}$ ), obtained with Canasat observational data (Rudorff et al., 2010) (scenario 2).**

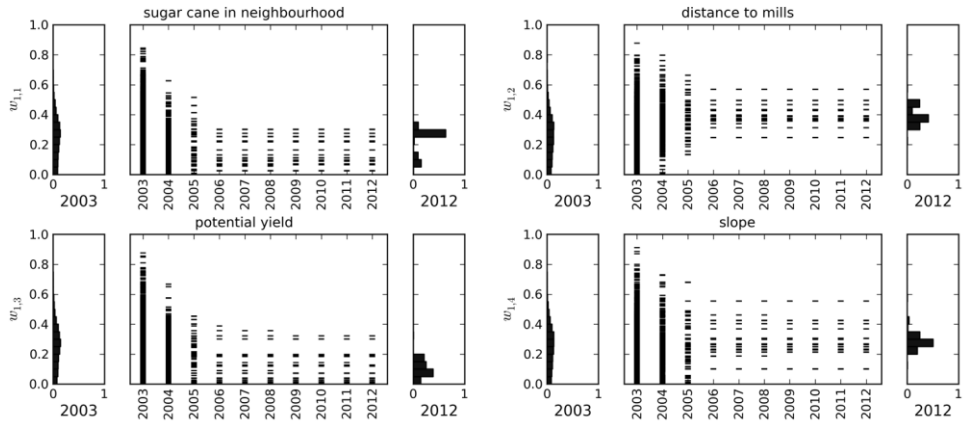


**Figure 4.7: Representation of the posterior distributions of the parameters of the suitability factors, obtained with Canasat observational data (Rudorff et al., 2010) (scenario 2): a, window that determines whether or not a cell belongs to the neighbourhood ( $I_{1,1w}^2$ ), b, 'preferred' fraction of sugar cane neighbours (hatched area) within the window ( $f_{1,1}$ ), c, suitability for distance to mills  $u_{1,2,t}$  plotted against distance to mills  $x_{1,2,t}$ , d, suitability for potential yield  $u_{1,3,t}$  plotted against potential yield  $x_{1,3,t}$ , and e, suitability for slope  $u_{1,4,t}$  plotted against slope  $x_{1,4,t}$ . Black lines represent the median of the parameter value, grey areas are 95% confidence intervals.**

In the parameters,  $\mathbf{p}_{n,k,t}$ , the systemic change is visible as well (Figure 4.7), although the uncertainty in the parameter values is mostly high. From 2006 to 2008 the mean window size of the neighbourhood,  $l_{1,1,w}$ , is about 50% smaller than in the other years. The  $\alpha$  values are very low, 0.02 and 0.01 in 2007 and 2008 (Table 4.4), indicating systemic change. The neighbourhood fill,  $f_{1,1}$ , on the contrary, is more stable over time, around a value of 0.1, except in the first two years, when it is about 0.2. For the parameters  $a_{1,k}$ , for  $k = 1, 2, 3$ , the stability over time is more difficult to observe because of the large uncertainty. This can also be concluded from the distribution comparison test, which gives few low  $\alpha$  values, except for  $a_{1,2}$  in 2004 to 2006. For all weights (Figure 4.6) and the parameters  $l_{1,1,w}$  and  $f_{1,1}$  (Figure 4.7) the posteriors in the period 2006 to 2009 are narrower than in the other years.

**Table 4.4: Results of the two stationarity tests for the weights (weight of neighborhood  $w_{1,1}$ , weight of distance to mills  $w_{1,2}$ , weight of potential yield  $w_{1,3}$ , and weight of slope  $w_{1,4}$ ) and the parameters (window length  $l_{1,1,w}$ , neighborhood fill  $f_{1,1}$ , and the shape parameters for the attraction/repulsion suitability factors  $a_{1,k}$  for  $k = 1, 2, 3$ ) for scenario 1 (with synthetic data that is supposed to be stationary) and scenario 2 (with the Canasat data (Rudorff et al., 2010)). For the distribution comparison test, the  $\alpha$  (Equation 4.7) is given belonging to the tipping point between rejection and no rejection. This  $\alpha$  value gives the probability that the null hypothesis of stationarity is unjustly rejected. The p-value of the Runs test gives the probability that the pattern found in the deviations from the mean is random. Hence, for both tests: low values (red colors) indicate a high probability of systemic change and high values (green colors) indicate a high probability of stationarity.**

	year	2004	2005	2006	2007	2008	2009	2010	2011	2012			
variable	scenario	probability ( $\alpha$ in Equation 4.7) of unjustly rejecting stationarity										average $\alpha$	p-value Runs test
$w_{1,1}$	1	0.82	0.65	0.93	0.93	0.91	0.79	0.98	0.96	0.9	0.87	0.66	
	2	0.54	0.48	0.68	0.07	0.11	0.64	0.97	0.85	0.83	0.57	0.05	
$w_{1,2}$	1	0.99	0.9	0.98	0.92	0.98	0.93	0.94	0.95	0.89	0.94	0.66	
	2	0.96	0.97	0.89	0.67	0.6	0.62	0.69	0.79	0.99	0.8	0.15	
$w_{1,3}$	1	0.97	0.82	0.91	0.84	0.91	0.87	0.87	0.95	0.97	0.9	0.21	
	2	0.86	0.74	0.49	0.16	0.25	0.31	0.85	0.95	1	0.62	0.04	
$w_{1,4}$	1	0.99	0.94	0.97	0.88	0.98	0.96	1	0.98	0.99	0.97	0.66	
	2	0.86	0.95	0.86	0.41	0.41	0.45	0.97	0.96	0.93	0.76	0.05	
$l_{1,1,w}$	1	1	0.92	0.82	0.17	0.28	0.38	0.75	0.93	0.93	0.69	0.04	
	2	1	0.86	0.31	0.02	0.01	0.61	0.79	0.97	0.92	0.61	0.04	
$f_{1,1}$	1	0.78	0.99	0.76	0.32	0.62	0.33	0.71	0.8	0.87	0.69	0.04	
	2	0.44	0.16	0.87	0.88	0.34	0.21	0.22	0.53	0.59	0.47	0.01	
$a_{1,2}$	1	0.58	0.7	0.57	0.47	0.77	0.43	0.56	0.57	0.7	0.59	0.04	
	2	0.05	0.08	0.04	0.94	0.3	0.83	0.36	0.23	0.24	0.34	0.59	
$a_{1,3}$	1	0.86	0.92	0.52	0.23	0.43	0.26	0.71	0.76	0.62	0.59	0.04	
	2	0.54	0.57	0.33	0.46	0.53	0.19	0.43	0.43	0.37	0.43	0.01	
$a_{1,4}$	1	0.63	0.57	0.54	0.44	0.51	0.6	0.57	0.6	0.61	0.56	0.15	
	2	0.41	0.54	0.49	0.7	0.85	0.88	0.58	0.64	0.65	0.64	0.66	



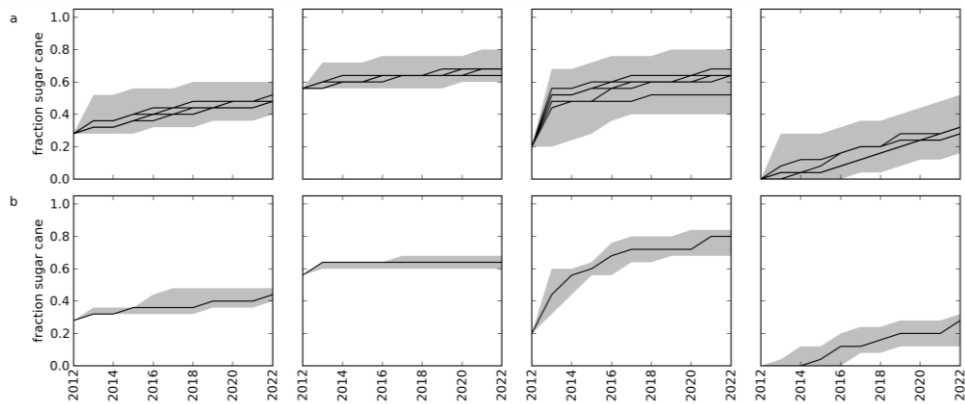
**Figure 4.8: Evolution of the weights for the four suitability factors,  $w_{1,k}$ , using the traditional particle filter approach and the Canasat data (Rudorff et al., 2010) (scenario 3). In the main panels the black horizontal lines represent the ensemble members. The smaller panels on the left and right give the full prior (2003) and posterior (2012) distributions.**

The start of the systemic change, 2006, is a year with no identified societal changes (Figure 4.5). The ‘recovery’ period of the system, 2009 to 2010, coincides with the years of bad harvests, and a percentage of sugar cane fields left unharvested more than twice as high as in previous years (Aguar et al., 2010). In 2010, the average area on which pre-harvest burning is forbidden increases from 30 to 50% for slopes below 12%. The potential connection between the societal changes and the observed systemic change is considered in the discussion section (4.5.2).

Filtering with Canasat data (Rudorff et al., 2010) using the traditional particle filter method (scenario 3) (Figure 4.8) yields the same order of importance of the suitability factors (distance to mills, neighbourhood and slope, potential yield) as found in scenario 2 (Figure 4.5). However, the posterior distributions of the weights in scenario 3 are narrower than the posterior distributions (per year) in scenario 2 (Figure 4.6). And, obviously, variation in the weights over time cannot be detected in scenario 3, as only one posterior distribution is obtained per weight (Figure 4.8).

#### 4.4.3. Projection

Because no trend is detected in the posteriors of the weights and parameters in the identification phase that could be extrapolated, in each projection year (2013-2022) a posterior is drawn randomly from the posteriors of the identification years (2004-2012). The projection phase of scenario 2 is run five times to cover the uncertainty arising from the diverse potential sequences of drawn posteriors. Scenario 3 is run once, because it always yields the same result.



**Figure 4.9: Projection of fraction of sugar cane in 4 random 25 x 25 km blocks out of the total of 473 blocks in the study area for, a, scenario 2, and, b, scenario 3. Scenario 2 is run five times, to show how the results differ when other posteriors for the weights and parameters are used. Black lines represent the median of the block value, grey areas are 95% confidence intervals (for scenario 2 calculated over all values of the five runs together).**

The projection of the fraction of sugar cane per 25 x 25 km block shows little difference between scenario 2 and 3 in the median expansion trend of the selected blocks (Figure 4.9); the lines in the upper and lower panel have similar courses and end up at similar values in 2022. However, the 95% confidence interval in this trend in scenario 2 is on average twice as large as in scenario 3.

## 4.5. Discussion

### 4.5.1. Identification with synthetic data

From the fact that in scenario 1 the weights of the suitability factors converge to approximately the correct mean value of 0.25 (maximum error is 0.03), we conclude that the particle filter is successful in inferring the weights. For the variation in this mean over time, the Runs test gives high p-values (0.21-0.66), meaning that there is no reason to reject the null hypothesis that this variation is caused by randomness: the weights are almost certainly stationary, as expected.

The relatively weak convergence of the parameter values in most years indicates that the parameters perform equally well (or badly) over their complete prior distribution. So we conclude that the sugar cane distribution data does not contain sufficient information for inferring the parameters using the particle filter at this resolution. This could be related to the spatial averaging of the model results and observations to the 25 x 25 km blocks in the particle filter. The parameters should be stationary over time, but the Runs test denotes non-stationarity. However, this result is unreliable because of the large uncertainty in the parameters. The



distribution comparison test, that does take into account the full posterior distributions, does not indicate systemic change, supporting the conclusion that the low identifiability of the parameter makes the Runs test unsuitable. In future research other types of data may be used for inferring the parameters, or the parameters may be fixed and only the weights calibrated. So, in this case study the test for systemic change in model structure and/or parameterization should be focused on the model structure.

#### *4.5.2. Identification with Canasat data*

In scenario 2, non-stationarity is observed for three out of the four weights of the suitability factors (Table 4.4), indicating a period of systemic change. It concerns the weights of sugar cane in the neighbourhood, potential yield, and slope, in the period 2006 to 2008/2009. The weight of the neighbourhood suitability factor becomes higher, while the weights of slope and potential yield become lower (Figure 4.5). This implies that in determining where to create a new sugar cane field, existing fields in the neighbourhood become more important, and the slope of the land and potential yield, and therefore the mechanization potential and expected revenues, become less important. In addition, it appears that the neighbourhood window that is used to look for existing sugar cane fields in the vicinity becomes smaller (Figure 4.7), i.e. expansion of sugar cane occurs closer to existing fields. However, in the previous section it was concluded that we should be careful in interpreting the results for the parameters, because the scenario with synthetic data implied that they cannot be fully trusted.

The systemic change is gradual and reaches its maximum in 2007 (Table 4.4). Looking at the societal changes in the studied period that we considered of possible influence (section 4.3.1), there are two policy changes: the adoption of one of the schemes for the phasing out of pre-harvest burning (Aguiar et al., 2011, Gallardo and Bond, 2011) and the implementation of the agro-environmental zoning (Lucon and Goldemberg, 2010). Yet, these two policies are expected to increase the importance of the suitability factor slope, while in the identified systemic change period the weight decreases. The economic crisis in 2008 is also an unlikely candidate for the cause of the change, because the systemic change clearly starts already in 2006, when the crisis was not yet foreseen. In conclusion, given our shortlist of potential causes for systemic changes, no cause can be found for the onset of systemic change. Nevertheless, the low weight of the suitability factor potential yield in 2008 and 2009 confirms the conclusion of journalists that the crisis forced, by a discontinuation of investments, sugar cane production at older, less productive sites (Gómez Jr., 2013). A potential explanation of the recovery of the system in 2009, to its functioning before 2006, are the bad harvests

from 2009 to 2011 and the consequential larger share of fields left unharvested. The changed system in which the importance of potential yield was low, was possibly not maintainable anymore, because bad weather usually has a relatively large influence on already low yielding soils.

Remarkable is the relation between the demand increase (Figure 4.3) and the systemic change. The years 2006-2009 have a demand increase above average, the exact same period as the systemic change, with maxima in 2007 and 2008. This implies an indirect systemic change: an action has an effect on the input 'demand', in such a way that the transition rules and/or parameters have to change as well (Filatova and Polhill, 2012). A possible reasoning behind the connection between fast demand increase and an increase in the weight of the neighbourhood suitability factor is that a sudden upsurge in demand increase is difficult to predict for farmers. New farmers, searching new locations, may not have time to respond to the upsurge, but existing farmers, already having the machines and infrastructure, can expand their existing fields in response to the upsurge. This results in the fact that expansion is more guided by existing sugar cane cultivation than by optimal conditions (slope and potential yield) for new fields. The existing fields are already close to mills, so this suitability factor remains of equal importance. However, other explanations might be possible as well.

The fact that most of the societal changes cannot be traced in our results does not mean that they have no effect at all on the sugar cane expansion system in São Paulo; it only means that they have no effect on the sugar cane expansion system in São Paulo given our model setup, observational data and resolution. Optimal model structure, and consequently stationarity, is thus different at different resolutions, as also noted by Pontius Jr. and Spencer (Pontius Jr. and Spencer, 2005). For example, one would expect that the adoption of a new scheme for the phasing out of pre-harvest burning (Figure 4.2) results in an increase of the weight of the suitability factor slope, but this was not observed in this study. At a different resolution, considering a longer time period (there might be a time lag), considering different suitability factors, or using a different implementation of the currently used suitability factors the effects can possibly be observed. In this case we studied the systemic changes using block averages of sugar cane coverage as model outputs and observations. If systemic changes are studied based on spatial patterns, like number of patches, or landscape shape index (Pijanowski et al., 2002) the conclusion can be different. However, the advantage of using areal averages of land use is that this measure cannot be derived only from land use maps, but also from agricultural statistics databases. Where time series of land use maps of sufficient length are not always available, time series of agricultural statistics usually are, showing the large applicability of the shown approach. But, to be able to draw a conclusion on the impact of societal changes, or more specific the

effectiveness of policies, the methodology should be applied with various model setups, with different spatial patterns in the observational data, at various resolutions. This was, however, not the aim of our study. Also, the single land use type is an oversimplification. Applying the particle filter on a land use model with multiple land use types, using an agricultural statistics database as observations, is the next stage of our research.

The lower standard deviations in the period 2006 to 2009 compared to the other years can be explained by the information content of the observations. In the given period the demand increase is relatively large, as mentioned before. This implies allocation of a relatively high number of sugar cane cells in those years. With this greater amount of change, the particle filter can better detect the optimal relative importance and parameterization of the suitability factors, so the convergence of the probability distributions will be stronger. As a reference: when there is no demand increase or decrease at all, the particle filter can never identify the optimal model structure, because the observations contain no information (no change). It is important to note here that the information content of the observations is not the reason for the detected systemic change, because in scenario 1 (synthetic data) the same demand time series was applied, and the model structure was stationary.

#### *4.5.3. Projection*

The 95% confidence interval for the projected fraction of sugar cane per block is twice as large for scenario 2 compared to scenario 3 (Figure 4.9), indicating that the use of a different posterior in each year results in a higher uncertainty regarding the dynamics of fraction of sugar cane in a block. Still, caution should be taken in generalizing this quantitatively. If it is true that the systemic change in the identification period is related to changes in demand increase, the model structure used in the projection period should depend on projected demand. Nonetheless, if one assumes that different system structures that have existed in the past are valid, in any order, in the future, uncertainty in the projection of land use change becomes considerably higher. Although it was not analysed in this study, it is even possible that the uncertainty arising from the potential systemic changes in the future is so large that variation in the results of different storylines completely disappears. This is something that should be kept in mind when conducting land use change projections, especially over long time intervals.

Instead of representing changes in the model structure by a random approach, as is done here, it would be preferable to extend the model by including processes representing the systemic change itself. This would enable better forecasting of future changes as variation in the model structure becomes a function of the state

or inputs of the modelled system itself. For example, one can connect the land use change model to a transportation model to account for changes in accessibility (e.g., Aljoufie et al., 2013), or to an erosion model to represent changing landscape features (e.g., Claessens et al., 2009). However, in our study the causes for the systemic change were not clearly identified, so dealing with systemic change by including the societal changes causing them was not achievable. We foresee that this will also be unachievable in many other land use change modelling studies, as often the knowledge of the system and the data availability are insufficient to fully understand and model the systemic changes.

#### **4.6. Conclusion**

Our first aim in this paper was to develop a general methodology, applicable to any type of model, to test for systemic change. In the methodology observations of the real system are assimilated into the model, using a particle filter (van Leeuwen, 2009). The particle filter was used to update the prior knowledge about the model structure, in the case of our land use change model the selection and relative importance of suitability factors, and parameters during model runtime at years for which observations of real land use were available (2004-2012) (see also Chapter 3). Using the particle filter separately for each point in time for which a land use map was available, we have obtained optimal model structures for these different points in time.

One limitation of our methodology is the strong assumption about the uncertainty in the observations. Also, the land use change model that was used to test the methodology was relatively simple, with only one active land use type. Another problem is that the two statistical tests used to provide evidence for the systemic changes, did not always give high significance levels. Therefore, we hope that this study serves as an eye opener to the potential presence of systemic change, in land use systems as well as in other modelling domains, and as a first step towards a sound methodology to test for systemic changes.

Given these limitations, we still believe that some conclusions can be drawn about systemic change in our case study of sugar cane expansion in the São Paulo state in Brazil. Here, the assumption of a constant model structure was not an adequate representation of the land use system given a time series of observations of past land use. A visual inspection and an analysis of the quantity of variation in the distinctive posterior distributions of the suitability factor weights and parameters, as well as the outcome of two statistical tests on these distributions have provided a strong indication of non-stationarity in the model structure and parameters, i.e. systemic change.

The systemic change appeared to be indirect: something has an effect on the input demand for sugar cane, in such a way that the transition rules and parameters have to change as well (Filatova and Polhill, 2012). But, although an inventory was made of societal changes in the study area during the studied period, none of these could be related to the onset of the observed systemic change in the land use system in 2006. The recovery of the system, in 2008 or 2009, might be related to a few years of bad harvests, forcing farmers to focus more on potential yield when selecting a new field.

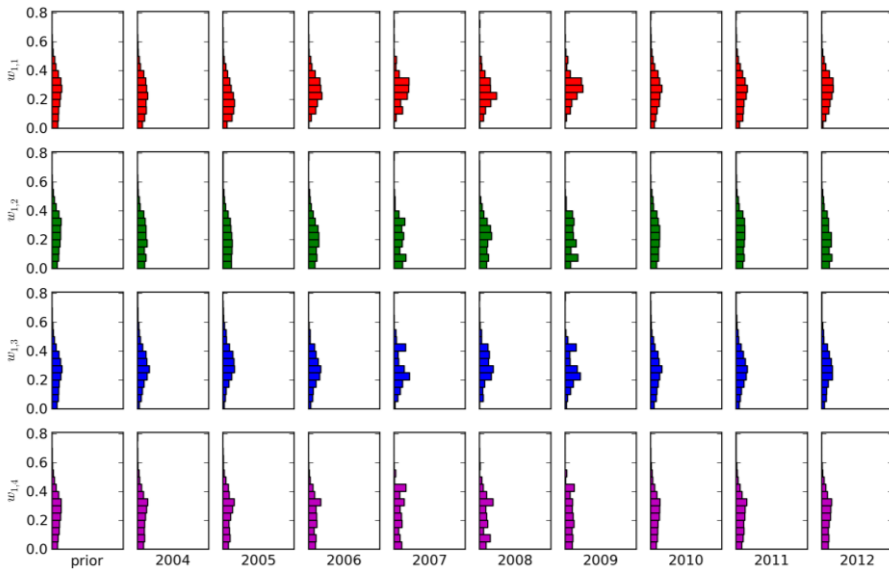
Because no clear reason was detected for the model structure and parameter changes in the identification period, we assumed that a future land use system could be any of the land use systems found in the identification period. Applying this resulted in an increase of the 95% confidence interval of the projected fraction of sugar cane by a factor of two compared to the assumption that the future land use system is a combination of all land use systems found in the identification period, in a stationary way.

In view of the above, we recommend land use change modellers to check, if permitted by data availability, whether or not the system was stationary in the past and if potential causes can be found for detected non-stationarity. The methodology proposed in this paper can be used for such an analysis although it certainly needs further evaluation given the limitations of this study described above. Non-stationarity in land use change projections is challenging to model, because it is difficult to determine when the system will change and how. We cannot expect land use change modellers to incorporate systemic changes in their models. Nonetheless, we believe that they should be more aware, and communicate more clearly, that what they try to project is at the limits, and perhaps beyond the limits, of what is still projectable, because systemic changes seem to occur in reality.

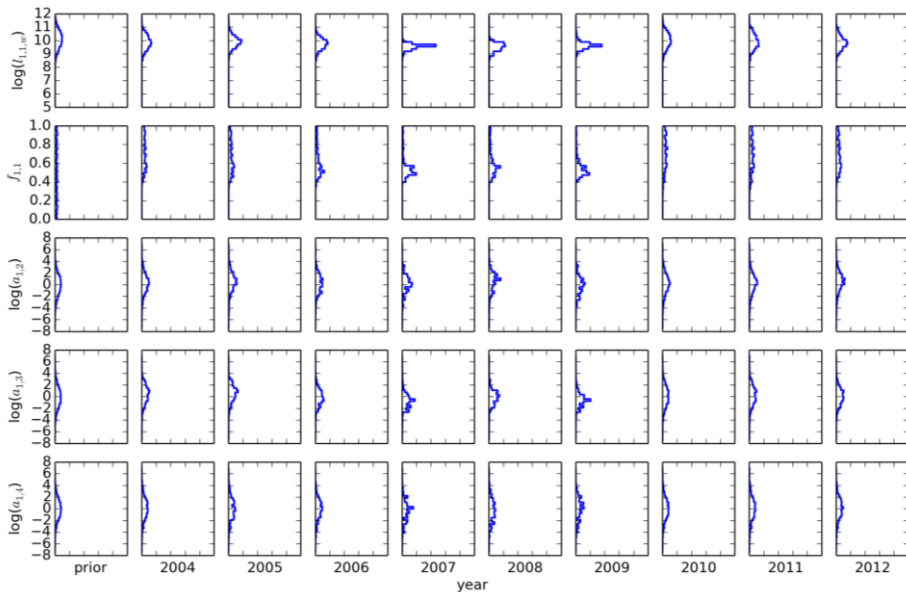
#### **4.7. Acknowledgements**

This work was carried out within the BE-Basic R&D Program, which was granted a FES subsidy from the Dutch Ministry of Economic affairs, agriculture and innovation (EL&I). We thank the Brazilian National Institute for Space Research (INPE), and especially Bernardo Rudorff, for providing the Canasat maps that were used as observational data. Six anonymous reviewers and the special issue editors are thanked for their contributions.

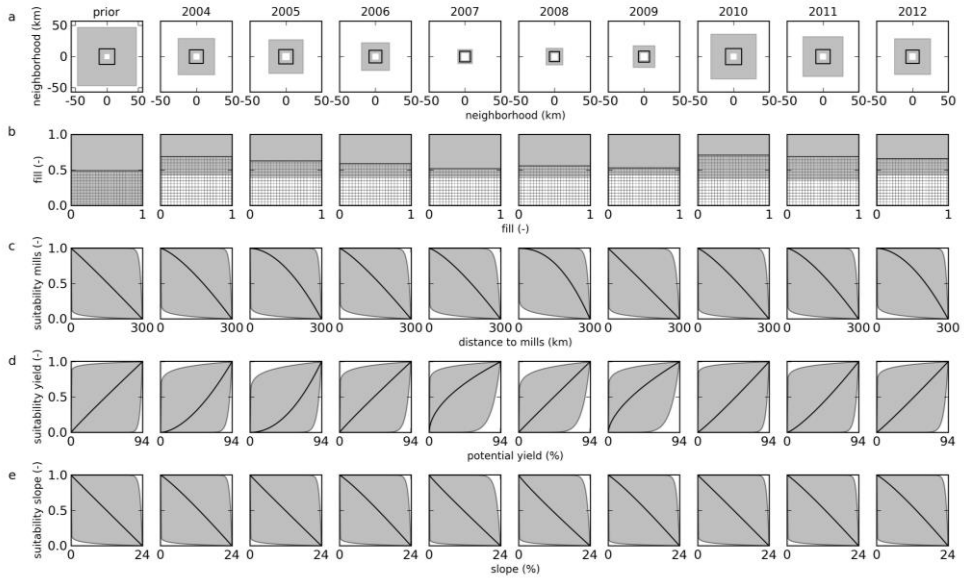
## 4.8. Appendix A



**Figure 4.10:** Posterior distributions of the weights of the suitability factors sugar cane in neighbourhood ( $w_{1,1}$ ), distance to mills ( $w_{1,2}$ ), potential yield ( $w_{1,3}$ ), and slope ( $w_{1,4}$ ), obtained with synthetic observational data (scenario 1).



**Figure 4.11:** Posterior distributions of the parameters of the suitability factors, the length of the neighbourhood window ( $l_{1,1,w}$  (m)), the ‘preferred’ fraction of sugar cane neighbors within the window ( $f_{1,1}$ ), and  $a_{1,k}$ , the suitability function shape parameter, for  $k = 2, 3, 4$ , obtained with synthetic observational data (scenario 1).



**Figure 4.12: Representation of the posterior distributions of the parameters of the suitability factors, obtained with synthetic observational data (scenario 1): a, window that determines whether or not a cell belongs to the neighbourhood ( $I_{n,k,w}^2$ ), b, 'preferred' fraction of sugar cane neighbours (hatched area) within the window ( $f_{1,1}$ ), c, suitability for distance to mills  $u_{1,2,t}$  plotted against distance to mills  $x_{1,2,t}$ , d, suitability for potential yield  $u_{1,3,t}$  plotted against potential yield  $x_{1,3,t}$ , and e, suitability for slope  $u_{1,4,t}$  plotted against slope  $x_{1,4,t}$ . Black lines represent the median of the parameter value, grey areas are 95% confidence intervals.**





## 5. What can and can't we say about indirect land use change in Brazil using an integrated economic - land use change model?

**Judith A. Versteegen, Floor van der Hilst, Geert Woltjer, Derek Karssenber, Steven M. de Jong, André P.C. Faaij (2015), Global Change Biology Bioenergy, early view.**

**Abstract** - It is commonly recognized that large uncertainties exist in modelled biofuel induced indirect land use change, but until now, spatially explicit quantification of such uncertainties by means of error propagation modelling has never been performed. In this paper, we demonstrate a general methodology to stochastically calculate direct and indirect land use change (dLUC and iLUC) caused by an increasing demand for biofuels, with an integrated economic – land use change model. We use the global Computable General Equilibrium model MAGNET, connected to the spatially explicit land use change model PLUC. We quantify important uncertainties in the modelling chain. Next, dLUC and iLUC projections for Brazil up to 2030 at different spatial scales and the uncertainty herein are assessed. Our results show that cell ( $5 \times 5 \text{ km}^2$ ) based probabilities of dLUC range from 0 to 0.77, and of iLUC from 0 to 0.43, indicating that it is difficult to project exactly where dLUC and iLUC will occur, with more difficulties for iLUC than for dLUC. At country level, dLUC area can be projected with high certainty, having a coefficient of variation (cv) of only 0.02, while iLUC area is still uncertain, having a cv of 0.72. The latter means that, considering the 95% confidence interval, the iLUC area in Brazil might be 2.4 times as high or as low as the projected mean. Because this confidence interval is so wide that it is likely to straddle any legislation threshold, our opinion is that threshold evaluation for iLUC indicators should not be implemented in legislation. For future studies we emphasize the need for provision of quantitative uncertainty estimates together with the calculated LUC indicators, to allow users to evaluate the reliability of these indicators and the effects of their uncertainty on the impacts of land use change, like greenhouse gas emissions.

## 5.1. Introduction

Governments throughout the world have set mandatory biofuel targets for the transport sector, aiming at mitigating climate change, improving energy security, and stimulating rural development (Sorda et al., 2010). Currently, one of the central problems in the biofuel arena is the premise of biofuel induced land use change (IPCC, 2011, Finkbeiner, 2014, Warner et al., 2014, Creutzig et al., 2012). These land use changes can have negative impacts like carbon stock loss, rising food prices, loss of biodiversity, and water scarcity, reducing the eligibility of the feedstock as a sustainable source for biofuels. An increased demand for biofuel feedstocks can lead to direct land use change (dLUC): land use is changed from some previous use to the biofuel feedstock. This, in turn, can lead to indirect land use change (iLUC): a change of land use outside the biofuel feedstock cultivation area, induced by a change in use or production quantity of that biofuel feedstock. This can happen either when the agricultural land use type converted to the biofuel feedstock is displaced to elsewhere, in order to continue to meet the demand for its agricultural products, or when the direct conversion triggers a change in the price of agricultural products, causing land to be taken into (or out of) production elsewhere (Wicke et al., 2012). The question to be tackled is to what extent the global increase in demand for biofuels (Broch et al., 2013) leads to dLUC and iLUC and how the negative effects can be minimized.

Direct land use changes are unambiguously visible in both historical data and spatial land use change model results. DLUC takes place wherever a bioenergy crop field appears and consequently displaces the previous land use. On the contrary, iLUC cannot be directly observed (Finkbeiner, 2014), because if e.g. pasture displaces forest in the presence of an expansion of bioenergy cropland over pasture, this does not necessarily mean that the pasture displacement is caused by the expansion of bioenergy cropland. The pasture might have caused deforestation for a reason unrelated to bioenergy. In other words, the indirect effects of a particular demand increase cannot be identified from historical data because the effects are intertwined with a wide range of processes from which the effects are also present in these data (Overmars et al., 2011, O'Hare et al., 2011). Separate identification is only possible by comparing all land use changes with and without the demand increase for bioenergy, which can be done using a simulation model (Creutzig et al., 2015).

The processes governing dLUC and iLUC range from global to local scale. For example, the impact of the biofuel targets on demands for feedstocks in different parts of the world is a global market issue. On the other hand, at which location the land use changes and which previous land use is replaced, is primarily steered by local factors, such as accessibility and biophysical conditions (Meyfroidt et al., 2013). Likewise, the impacts of the land use change are highly location-dependent

(e.g. van der Hilst et al., 2014). Therefore, a sound approach to model iLUC is by using a global economic model coupled to a spatially explicit land use change (LUC) model to take both the global and local scale level into account, as for example demonstrated by Lapola *et al.* (2010).

It is commonly recognized that there is a large uncertainty in modelled iLUC (Wicke et al., 2012, Finkbeiner, 2014, Creutzig et al., 2015, Mathews and Tan, 2009, Malins, 2013). The uncertainties arise from model structure uncertainty (Refsgaard et al., 2006, Chapter 4), from data (inputs, calibration dataset and initial system state) (Dendoncker et al., 2008), and from model coupling (Ray et al., 2012). For iLUC in particular, uncertainty in reported values also stems from the fact that the assumptions, the employed models and the validity of these models are often not clearly communicated (Mathews and Tan, 2009). Information quantifying uncertainty in iLUC is critical to evaluate whether or not iLUC indicators are reliable enough to be included in legislation, to identify which parts of the modelling chain have the highest priority for improvement, i.e. cause most uncertainty, and to assess how this uncertainty propagates to the impacts of iLUC, like greenhouse gas (GHG) emissions (e.g. Plevin et al., 2015). Uncertainty information can be obtained by 1) being explicit about the applied models, the processes included in these models and the parameter settings used, as well as the uncertainty in the various model components and the performance of these models (Broch et al., 2013, Mathews and Tan, 2009), and 2) assessment of the magnitude of the output uncertainty by e.g. doing Monte Carlo analyses of iLUC (Wicke et al., 2012, Warner et al., 2014, Plevin et al., 2015, Nelson et al., 2014, Wicke et al., 2015). Uncertainty should be assessed at different spatial scales because different types of impacts play a role at different scales and it is known that uncertainty is highly scale-dependent (Pontius Jr. and Spencer, 2005, e.g. Chapter 2). Yet, such information is currently scarcely reported for iLUC; a status we aim to improve with this paper.

We have set up a model study with the global Computable General Equilibrium (CGE) model MAGNET (e.g. Woltjer and Kuiper, 2014, Kavallari et al., 2014), integrated with the spatially explicit land use change model PLUC (e.g. Chapter 2). With this integrated model we project land use change caused by an increasing demand for biofuels up to 2030 for Brazil, one of the main bioethanol producers in the world. Since Brazil holds the world's major potential for agricultural expansion (Alexandratos and Bruinsma, 2012), production and export of bioethanol are likely to increase in the future (Walter et al., 2014, IEA, 2013, OECD/Food and Agriculture Organization of the United Nations, 2014ood and Agriculture Organization of the United Nations 2014). Yet, the country also maintains the largest area of natural remnants, with high carbon stocks and high levels of biodiversity, stressing the need to assess potential negative impacts. For this case study we seek to answer the following research questions: 1) What are the dLUC and iLUC projections for

Brazil up to 2030 at different spatial scales and what is the uncertainty herein? 2) What are the sources of uncertainty for each step in the model chain and how do these uncertainties influence dLUC and iLUC projections? 3) What is the contribution of the economic and land use change model to the uncertainty in dLUC and iLUC at the different spatial scales?

The next section introduces the Brazilian case study, presents the LUC model, the CGE model, and the way they are coupled, describes the calibration method, defines the projection scenario for the increased demand for biofuels, and explains how iLUC is derived from the results. Section three illustrates the results for the three research questions. The final section discusses these results in light of the research questions and gives suggestions for further research.

## **5.2. Materials and methods**

### **5.2.1. Overview**

The projection of dLUC and iLUC in Brazil caused by an increasing demand for biofuels and the uncertainty herein is performed using MAGNET (Woltjer and Kuiper, 2014), a global Computable General Equilibrium (CGE) model, connected to the land use change model PLUC (e.g. Chapter 2), tailored to Brazil (Figure 5.1). For 2006 an initial land use map is created by combining tabular area data per land use type and land use maps with satellite data. This map is used as the initial system state for PLUC. Next, PLUC is calibrated from 2007 until 2012 based on trends per land use type from agricultural statistics databases. To project the dLUC and iLUC effects of the biofuel mandates, we define both a 'biofuel scenario' that includes these mandates and a 'reference scenario' that does not include them. For both scenarios, MAGNET determines the supply and demand of all commodities in all world regions up to 2030 and, related to that, the area they occupy. This 2013 – 2030 time series of land area demands per land use type for Brazil is then input for the spatially explicit land use change projection up to 2030 by PLUC. The PLUC outputs are a time series of land use maps. By comparison of the maps of the two scenarios dLUC and iLUC are assessed.

In the model chain, uncertainties in the inputs, calibration dataset, initial system state and model structure are quantified, part of which propagates through the model coupling (Figure 5.1). To quantify uncertainty in MAGNET, it is run with two different parameter sets, resulting in an upper and a lower demand limit. PLUC, including the generation of the initial land use map, the calibration, and the demand coming from MAGNET, is used stochastically by running it in Monte Carlo mode (Figure 5.1).

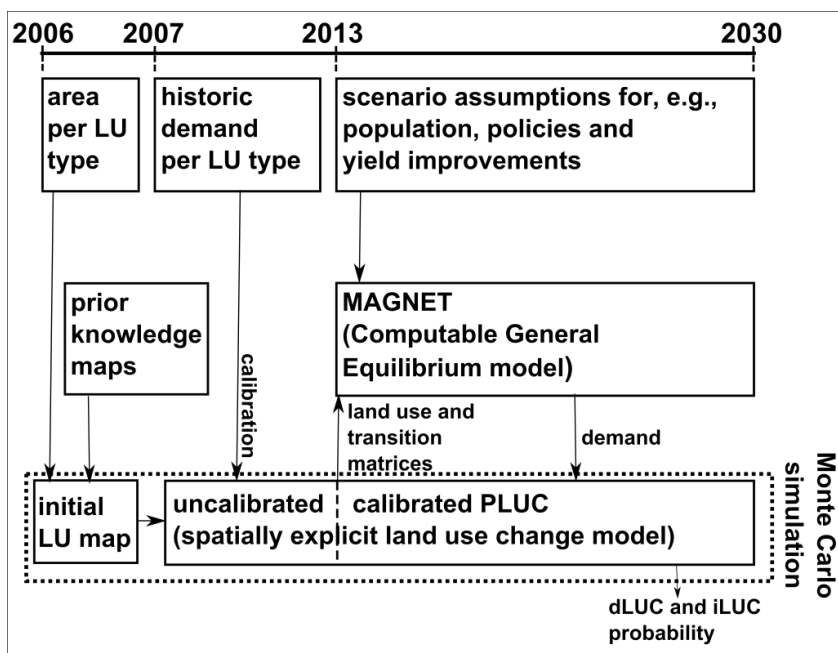


Figure 5.1: Overview of the modelling chain and model run-time frame to simulate the probability of dLUC and iLUC in Brazil up to 2030.

### 5.2.2. Case study

Brazil has been producing bioethanol from sugar cane since the beginning of the 20<sup>th</sup> century and has been exporting the ethanol since 1989 (Andrade de Sá et al., 2013). Sugar cane currently occupies the third largest area of all crops in Brazil, topped only by soy and maize (although a large quantity of the maize is cultivated as second crop) (IBGE, 2013b). The main sugar cane production areas are the Central South region and the Northeast region. Recent expansion has mainly taken place in the Central South region: in the past decade the total area dedicated to sugar cane cultivation has more than doubled in that region (Rudorff et al., 2010). It is expected that future expansion will also predominantly occur in the Central South region (Lapola et al., 2010, Nassar et al., 2008). According to Adami *et al.* (2012b) over 99% of all sugar cane expansion in the last decade has taken place over existing agricultural land, signifying that the direct effect of increasing ethanol demand on deforestation is negligible. However, deforestation can still take place through iLUC, which is also shown by others (e.g. Lapola et al., 2010, de Souza Ferreira Filho and Horridge, 2014).

### 5.2.3. *Initial land use map and land use change model*

We distinguish eleven different land use types  $n$ , where  $n = 1, 2, \dots, 11$ : urban, water, natural forest, rangeland, planted forest, crops (excluding sugar cane), grass and shrubs, sugar cane, planted pasture, bare soil and abandoned agricultural land. Planted pasture and natural pasture (rangeland) are modelled separately because the extensively managed, naturally vegetated rangelands have a stocking rate of about 70% lower than the intensively managed planted pastures (IBGE, 2006, Aguiar and d'Athayde, 2014). Cropland includes both annual and permanent crops. Sugar cane is modelled as a separate land use type to be able to evaluate where sugar cane expands in reaction to the increased ethanol demand and which other land uses it replaces.

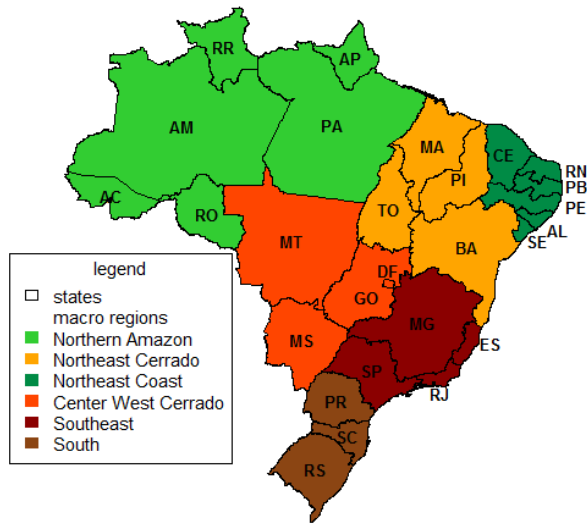
PLUC (PCRaster Land Use Change model) (van der Hilst et al., 2014, Diogo et al., 2014, van der Hilst et al., 2012, Chapter 2) is founded on the separation between the quantity of change per land use type, and the spatial allocation of this change, like many other land use change models (Pontius Jr. and Neeti, 2010). The quantity of land demanded per land use type  $n$  is called 'demand'  $d_{n,t}$ , in which  $t$  is the time step in years, with  $t = 1, 2, \dots, T$ . The total area per land use type in the demand time series (tabular area data from agricultural statistics) for the initial year of the simulation should match the total area per land use type in the initial land use map, i.e. the initial system state of the model. If the time series and initial map are coming from different sources, which is likely, a perfect match is obviously never going to be the case. You and Wood (2005) provide a deterministic method to create a land use map that matches the time series, by spatially disaggregating land use areas per administrative region from the time series into raster cells within that region, using prior knowledge maps. We apply this procedure, using municipalities as administrative regions (5566 in total for Brazil), to create an initial land use map for Brazil with a cell size of  $5 \times 5 \text{ km}^2$  for the year 2006. This year is chosen because it was the year in the recent past (to have a calibration period) with the best data availability for both the tabular and prior knowledge map data. Compared to You and Wood (2005) we do a few things differently, most importantly adding a method to make a stochastic map instead of a deterministic map, in order to include uncertainty arising from errors in the initial land use map into the model chain, as explained in Appendix A.

**Table 5.1: Suitability factors,  $k$ , of the active land use types,  $n$ , for the Brazilian case study.**

$n$	land use type	$k$	process represented	suitability factor
4	rangeland	1 2 3	economies of scale transportation costs potential profits per hectare	$n$ in the neighbourhood distance to roads potential yield of $n$
5	planted forest	1 2 3	economies of scale transportation costs potential profits per hectare	$n$ in the neighbourhood distance to roads potential yield of $n$
6	crops	1 2 3 4 5	economies of scale transportation costs potential profits per hectare costs to make the land cultivatable double cropping potential	$n$ in the neighbourhood travel time to hubs for $n$ potential yield of $n$ conversion elasticity growing season length
8	sugar cane	1 2 3 4	economies of scale transportation costs potential profits per hectare costs to make the land cultivatable	$n$ in the neighbourhood travel time to hubs for $n$ potential yield of $n$ conversion elasticity
9	planted pasture	1 2 3	economies of scale transportation costs potential profits per hectare	$n$ in the neighbourhood distance to hubs for $n$ potential yield of $n$

Out of the eleven land use types considered in PLUC, five are assumed to respond to changes in the economy by expanding or contracting: rangeland, planted forest, crops, sugar cane and planted pasture. These *active* land use types are demand-driven (Table 5.1). The other six land use types do not have demands. They are either *passive*, meaning that they can contract or expand due to the dynamics of the active land use types, or *static*, meaning that they cannot change and are thus fixed on the map. Passive land use types are natural forest, grass and shrubs, bare soil, and abandoned agricultural land. Abandoned land originates when an active land use type contracts; it is not present in the initial land use map. Static land use types are urban and water.

The demands for the five dynamic land use types over time in Brazil have been subdivided into six regions (Figure 5.2), corresponding to the macro regions defined by the Brazilian Institute of Geography and Statistics (IBGE). We added one region by splitting the Northeastern macro region into two regions, as suggested by Nassar *et al.* (2010), because the Northeast coast differs significantly from the Northeast Cerrado (savannah) in terms of agricultural production.



**Figure 5.2: The six macro regions in Brazil (six different colours) used as demand input units in PLUC and the 27 states (black lines), or in fact 26 states and one federal district, used as calibration units. The state name abbreviations are: AC = Acre, AL = Alagoas, AM = Amazonas, AP = Amapá, BA = Bahia, CE = Ceará, DF = Distrito Federal, ES = Espírito Santo, GO = Goiás, MA = Maranhão, MG = Minas Gerais, MS = Mato Grosso do Sul, MT = Mato Grosso, PA = Pará, PB = Paraíba, PI = Piauí, PR = Paraná, RJ = Rio de Janeiro, RN = Rio Grande do Norte, RO = Rondônia, RR = Roraima, RS = Rio Grande do Sul, SC = Santa Catarina, SE = Sergipe, SP = São Paulo, TO = Tocantins.**

In PLUC, the spatial allocation is regulated by spatial attributes that serve as proxies for important drivers of location, i.e. processes that determine where a land use type expands or contracts. These are called suitability factors  $k$ , with  $k = 1, 2, \dots, K_n$  (each active land use type  $n$  can have a different number of suitability factors). For each  $n$  defined as active, a weighted sum of these suitability factors forms the total suitability map. In one model time step, representing one year, the demands of the active land use types are allocated sequentially for each macro region, as follows. For the first active land use type  $n$  the total suitability map is sorted, and cells are allocated to  $n$ , starting with the cell with the highest suitability value that is not yet of type  $n$ , until  $d_{n,t}$  is fulfilled. Next, the same is done for the second land use type in the sequence, with the exception that cells occupied by the first land use type cannot be changed. This procedure continues until the demands of all active land use types in all macro regions have been allocated (see also Appendix B).

The suitability factors for the Brazilian case study are given in Table 5.1. To represent economies of scale ( $k = 1$ ), the number of neighbours of the same land use type is counted in a square window of 5 by 5 cells ( $25 \times 25 \text{ km}^2$ ). For transportation costs ( $k = 2$ ), the travel time to hubs is used as a proxy. This is the time it takes to transport the products originating from the land use type to the nearest production facility. For planted forest we have no data about the location



of hubs (e.g., saw mills) and for rangelands we believe that the livestock hubs are of lower importance, because livestock from rangeland is often 'finished' elsewhere before being slaughtered. Therefore, for these two land uses we apply distance to roads as the proxy for transportation costs. Potential profits per hectare ( $k = 3$ ) are represented by potential yield maps, using IIASA's GAEZ data (Tóth et al., 2012). Since, to our knowledge, no potential yield map exists for woody biomass, we use IIASA's map of the length of the growing season as a proxy for the potential yield of planted forest. The costs to make the land cultivatable ( $k = 4$ ) are estimated using a conversion elasticity, i.e. a fraction indicating the ease with which a certain land use type can be transformed into the land use type that implements the suitability factor, especially relevant for crops. Double cropping potential ( $k = 5$ ) is an important suitability factor in Brazil, indicated by the rapid increase in double or even triple cropped area over the last decade (Conab, 2014, Galford et al., 2008). We do not have a map of double cropping potential, so we use the growing season length as a proxy, which is supported by an analysis of the relation between these two by Arvor *et al.* (2014).

The no-go map, i.e. areas where expansion is not allowed, is an overlay of military areas, areas of indigenous people, and federal and state conservation units (Gurgel et al., 2009). Conservation policies or initiatives which have historically not been well enforced, such as the Forest Act (Sparovek et al., 2012), the soy moratorium (Rudorff et al., 2011), and the sugar cane zoning (Padua Junior et al., 2012), are not taken into account in this simulation. We are preparing another paper, in which we include more scenarios with, among other things, stricter nature conservation rules (van der Hilst et al., in prep.).

We use a Monte Carlo simulation with 5000 realizations. The weights of the suitability factors and the order of allocation are modelled stochastically. Their prior probability distributions are uninformed (see Methods S2).

#### **5.2.4. Calibration**

The aim of the calibration phase, 2007 to 2012, is to narrow the probability distributions of all stochastic elements: the order of the land use types and all weights of the suitability factors (Table 5.1). The model calibration is performed using a Bayesian data assimilation technique, the sequential importance resampling (SIR) particle filter (van Leeuwen, 2009). In short, the SIR particle filter compares the land use system simulated by PLUC and observations of the land use system from the real world, taking into account uncertainty in these observations. Next, it updates the Monte Carlo ensemble in such a way that well-performing realizations are progressed and poorly-performing realizations are discarded. An

extensive explanation of this model structure identification and calibration method for a case study in the São Paulo state is provided in Chapter 3.

For calibration, a time series of land use/cover data is required as observational data. For Brazil we use time series of areal data per land use type per state (Figure 5.2). These time series are derived using info from agricultural statistics databases (IBGE, 2013b, see van der Hilst et al., in prep., IBGE, 2013a, ABRAF, 2013). These observational data are not error free. Between the yearly (IBGE, 2013b (for crops), IBGE, 2013a (for livestock)) and the 10-yearly (IBGE, 2006) census data sources, areas and area increases differ from zero up to more than 100%. As one cannot calculate a standard deviation based on two values, we make an educated guess of the average error based on these data sources. Under the assumption that the observational errors are uncorrelated over space and time, we assign an observation error to the observed increase in area with a standard deviation of 20% of the observed increase in that time step.

After calibration, a land use matrix, summarizing the total areas per land use type in 2012, is computed per macro region, representing the initial system state for MAGNET (Figure 5.1). In addition, a land transition matrix is calculated per macro region, to be used for the calibration of MAGNET. These six land transition matrices show the average area of conversion from every land use type to every other land use type derived from PLUC over the whole calibration period.

As a measure of model performance, we calculate root mean squared error (RMSE), the root of the summed squared differences between the median of the modelled area and observed area over all states. We determine the reduction in RMSE (%) for the results of the calibrated model, i.e. with uncertainty reduced by the SIR particle filter, compared to the non-calibrated model. To evaluate the effect of calibration, we apply a split-sample approach: PLUC is calibrated using data from 2007 to 2009 and the model reduction in RMSE is evaluated from 2010 to 2012. This split-sample approach is used only to evaluate the effect of calibration. The model parameters we use for the projection, integrated with the MAGNET model, are calibrated based on all available observational data (2007-2012).

### 5.2.5. *Economic (CGE) model*

The growing demand for food, feed, fibre and bioenergy requires an increased agricultural output. This can be reached by raising inputs like fertilizers, machinery and labour (bound by technological limitations), i.e. expansion at the *intensive margin*, or by converting new land to agriculture, i.e. expansion at the *extensive margin* (Hertel, 2011), which can result in iLUC. At what ratio both alternatives are applied in face of a growing demand depends on e.g., land availability, prices and

policies that vary worldwide. To evaluate how demand grows over time and to assess to what extent this demand is fulfilled by expansion at the intensive and extensive margins, we use a global Computable General Equilibrium (CGE) model (Rose, 1995). Key parameters in CGE models are the elasticities, simulating behavioural responses, for example the response of the demand for a commodity to a change in price or the response of consumption to a change in GDP per capita.

The CGE model used is MAGNET (Modular Applied GeNeral Equilibrium Toolbox) (see for an extensive explanation Woltjer and Kuiper, 2014). This is the modularized and improved version of LEITAP (e.g. Banse et al., 2011, Hoefnagels et al., 2013). MAGNET uses the GTAP database version 8 (Narayanan et al., 2012), in an extended and adaptable form. For this case study, we use the database with 42 sectors (including various ethanol sectors that take into account co- and by-products like molasses and electricity, and a difference between planted pasture and rangeland), 45 commodities and 15 regions, of which Brazil is one. Brazil has been subdivided into six regions, matching the input macro regions for PLUC (Figure 5.2). These six macro regions are a subdivision in MAGNET in terms of agricultural production and land area only; for international trade Brazil is considered as one region. Total land availability per macro region is calculated from the no-go map.

In order to model land cover change, a regional land transition approach has been developed that is inspired on the work of de Souza Ferreira Filho and Horridge (2014) and further developed by Woltjer (2013). Herein, the area of land that is changed from one particular land use type  $n$  to another one  $m$ , depends on the land transition elasticity  $e_{n,m}$ . Using expert knowledge and trial and error, we test for all combinations of  $n$  and  $m$  for what values of  $e_{n,m}$  MAGNET can best reproduce the 2012 system state given by the land use matrix from PLUC, and the transitions given by the land transition matrix.

To assess the uncertainty related to the key parameters in the economic model, two runs are performed, one with considerably higher (200%) and one with considerably lower (25%) land transition elasticities  $e_{n,m}$  than the values found by the procedure above. This results in two demand time series per land use type, one for the upper land transition elasticities,  $d_{u,n,t}$ , and one for the lower land transition elasticities,  $d_{l,n,t}$ , where all potential lines between these time series are assumed to have equal likelihood:

$$d_{n,t} = d_{l,n,t} + Z_d \cdot (d_{u,n,t} - d_{l,n,t}), \quad \text{with } Z_d \sim U(0,1), \quad 5.1$$

for each active  $n$  in each  $t$

Equation 5.1 shows that the demand input of PLUC  $d_{n,t}$  in the projection phase has an error model based on a uniform distribution between  $d_{u,n,t}$  and  $d_{l,n,t}$ .

### 5.2.6. *Projection*

In the projection from 2013 to 2030 the socio-economic developments are based on the Shared Socioeconomic Pathways (SSPs) (O'Neill et al., 2014). The SSPs quantify global drivers of the energy-economy-land use system such as demographics and economic development. In these pathways projections are included on population and GDP growth. We use SSP2, the Middle of the Road pathway with some additional assumptions on for example the agricultural intensification over time (see van der Hilst et al., in prep.).

Using SSP2 and these assumptions, MAGNET is run up to 2030, providing total land areas occupied by all land use types for all world regions and the six macro regions in Brazil for the years 2013, 2015, 2020, 2025 and 2030. Yearly demand time series for the six macro regions to serve as an input for PLUC are obtained by a linear interpolation between these years and an aggregation of the areas of all individual crops, except sugar cane, into the single class cropland.

To evaluate the future dLUC and iLUC effects caused by current and planned ethanol mandates worldwide, we define both a 'biofuel scenario' including these mandates and a 'reference scenario' excluding them. This does not mean that there is no increase in the demand for sugar cane in the reference scenario, only that there is no (additional) increase originating from the increased ethanol demand. All other inputs and parameters of both models are kept the same as in the 'biofuel scenario'.

### 5.2.7. *Direct land use change (dLUC) and indirect land use change (iLUC)*

Normally, direct land use change can be assessed using one scenario, as the difference between current and projected land use. In our case, however, we want to assess dLUC from sugar cane caused by the biofuel mandates, i.e. only sugar cane expansion for ethanol. Therefore we want to exclude sugar cane expansion that is a result of an increased demand for sugar over time. Hence, both dLUC and iLUC originating from the mandates are assessed through the difference between the reference and the biofuel scenario (Table 5.2) in 2030. A grid cell that is sugar cane in the biofuel scenario, and something else in the reference scenario, is considered dLUC, i.e. sugar cane expansion resulting from the biofuel mandates. A grid cell that is nature in the reference scenario and agricultural land but not sugar cane in the biofuel scenario, is considered iLUC. The opposite effects exist as well. A grid cell that is sugar cane in the reference scenario and something else in the biofuel scenario is negative dLUC (neg\_dLUC), and a grid cell that is agriculture in the reference scenario and nature or abandoned land in the biofuel scenario is negative iLUC (neg\_iLUC).

Especially for iLUC this opposite effect might appear in the real world. If, for example, an area of 10000 ha of wheat fields is present, and 80% of this area is taken over by sugar cane for ethanol, then the remaining 20% of wheat land might be abandoned because the advantages of economies of scale have disappeared. The 8000 ha of displaced wheat land and the 2000 ha of wheat land now grown elsewhere, make 10000 ha of iLUC. In our methodology we count the abandoned land as -2000 ha of iLUC (and therefore we call it neg\_iLUC (Table 5.2)), coming to a total of 8000 ha iLUC, which was indeed the area of land shifted by sugar cane.

**Table 5.2: Classification of differences in land use between the reference and the biofuel scenario that are considered undesirable effects of increasing ethanol demand (dLUC and iLUC, dark grey), and the opposite effects (neg\_dLUC and neg\_iLUC, light grey). The class 'other agriculture' includes rangeland, planted forest, crops, and planted pasture. The class 'nature' includes natural forest, grass and shrubs, bare soil, and abandoned agricultural land; thereby assuming that land will eventually become nature when left abandoned. Zero stands for no difference, i.e. neither (neg\_)dLUC nor (neg\_)iLUC.**

		biofuel scenario		
		sugar cane	other agriculture	nature
reference scenario	sugar cane	0	neg_dLUC	neg_dLUC
	other agriculture	dLUC	0	neg_iLUC
	nature	dLUC	iLUC	0

To compare outcomes at different spatial scales, we focus our analysis on local, regional and national level, calculated from output of PLUC. At the regional level we use  $250 \times 250 \text{ km}^2$  blocks. We do not use administrative levels, like states, because these differ in size and are thus problematic to compare. The coefficient of variation (cv) (standard deviation of dLUC or iLUC area over all Monte Carlo realizations divided by the mean of dLUC or iLUC area over all Monte Carlo realizations) is used as the measure of uncertainty. Since this measure of uncertainty is standardized by the mean, the cv is comparable between dLUC and iLUC and between regions with different magnitudes of dLUC or iLUC. As the local level we use probabilities of dLUC and iLUC in single cells ( $5 \times 5 \text{ km}^2$ ).

### **5.2.8. Contribution of the two models to total output uncertainty**

We compare the contribution of the two models to the total output uncertainty, by running the projection until 2030 three times, all three with 5000 realizations. One Monte Carlo run is with both models stochastic (the default run used in all analysis described above). One run is with only PLUC stochastic (including the uncertainty in the initial land use map and calibration time series). In this run the demand  $d_{n,t}$  is fixed at the mean between the upper and lower time series, by setting  $Z_d$  (Equation 5.1) to 0.5 for all Monte Carlo realizations to exclude uncertainty from MAGNET. The uncertainty in the output of this run is thus caused by uncertainty in PLUC only. The final run is with only MAGNET stochastic. In PLUC the weights, the order of allocation and the land use map for 2012 are fixed by taking the medians hereof from the calibrated model, to exclude uncertainty from PLUC. This run results in information about output uncertainty caused by MAGNET. For the three runs we compare the mean and the coefficient of variation in dLUC and iLUC area at the different spatial scales.

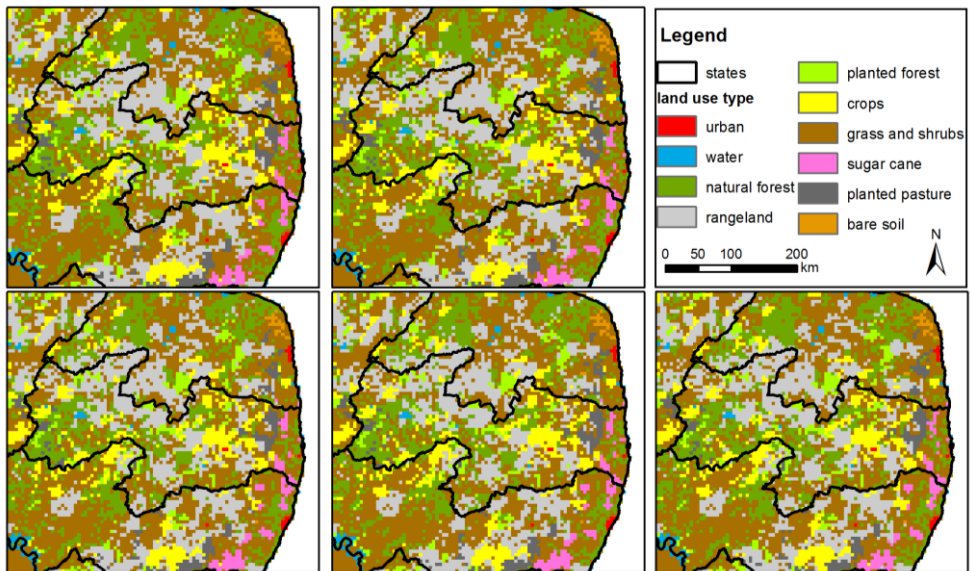
## **5.3. Results**

### **5.3.1. Sources of uncertainty for each step in the model chain and their influence on dLUC and iLUC projections**

#### ***Initial land use map***

For the initial year, 2006, a land use map was created for each Monte Carlo realization to serve as the initial system state (Figure 5.3). The total area per land use type per macro region is the same for all realizations and also the locations of individual patches within the macro region are the same, but the shape of these patches differs slightly, see for example the patch of sugar cane at the bottom of the map view in Figure 5.3. The patches of land use types that were assumed to be

known precisely, being urban, water and bare soil (see Appendix A) always have the same shape, see for example the shape of the city Natal, in the northeast of the map. We can conclude that the uncertainty in the initial land use map is very local, important only at cell level. The effect for projected iLUC will mainly be that when sugar cane expands in a certain grid cell, uncertainty in the initial land use map makes that in some realizations it expands over agricultural land, which may result in iLUC through displacement (depending on the demand trend for the displaced agricultural land use type), and in other realizations over nature, not resulting in iLUC, because there is no displacement effect.



**Figure 5.3: Five out of the 5000 realizations of the initial land use map (year 2006) zoomed in to the state Paraíba, in the Northeast Coast region of Brazil (see Figure 5.2).**

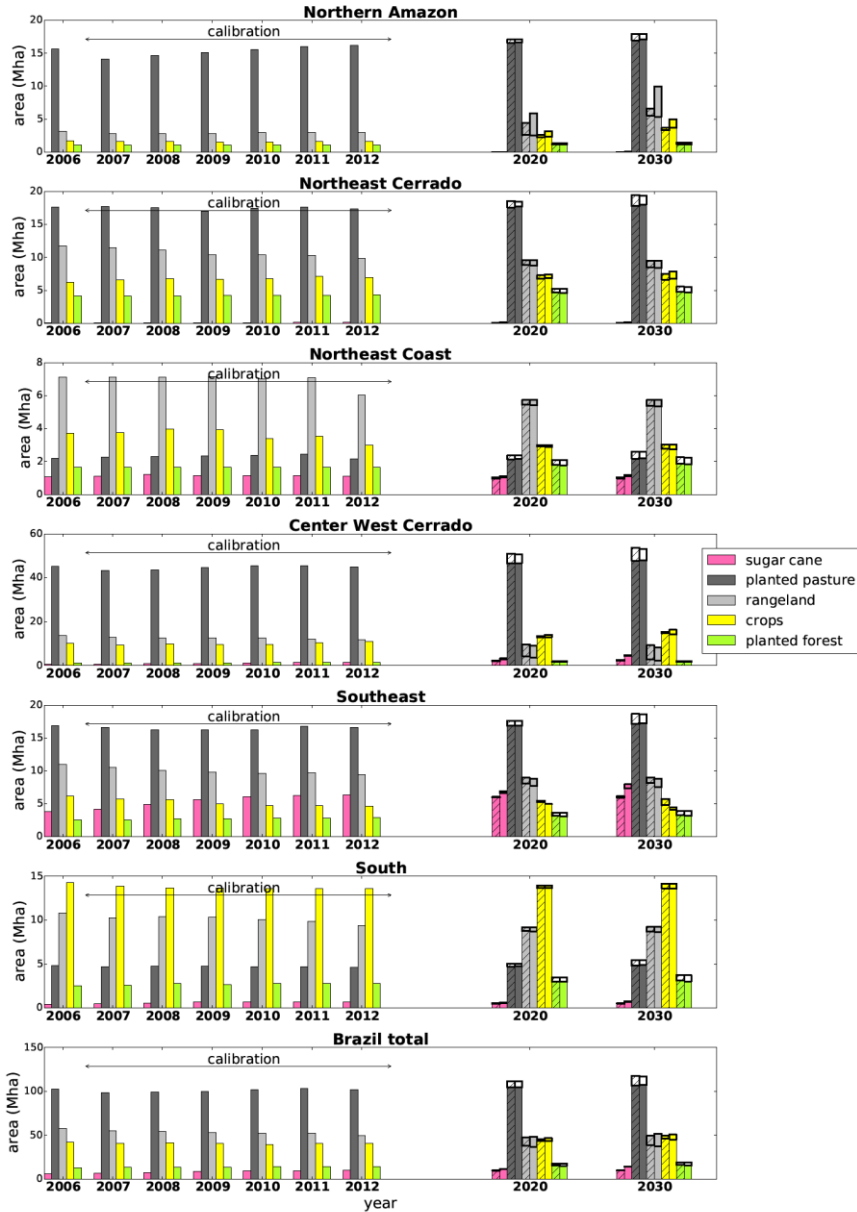


Figure 5.4: Demand for the five dynamic land use types in the six macro regions and Brazil total for the initial year, for the calibration period (using data from IBGE (2013b, 2013a) and ABRAF (2013)), and for two out of the five output years in the projection period (output from the MAGNET model). The ranges of the y-axes differ between macro regions to improve the visibility of trends. In the projection period, the hatched bar is the reference scenario and the non-hatched bar is the biofuel scenario. The thick box on top of the bars indicates the uncertainty in the output, i.e. the difference between  $d_{u,n,t}$  (elasticities set to 200%) and  $d_{l,n,t}$  (elasticities set to 25%). In the case of a filled box  $d_{l,n,t}$  is higher than  $d_{u,n,t}$ , and in case of an unfilled box  $d_{l,n,t}$  is lower.



### *Land use change model, calibration*

The input demand time series that were constructed from agricultural statistics for calibration are shown in Figure 5.4 (2007-2012, indicated by an arrow). After calibration using this demand and the observations, out of the 120 possible sequences (see Appendix B) for the order of allocation of the land use types, 72 obtain a posterior probability of zero, i.e. they are not present anymore in the ensemble. So, 48 unique sequences remain, with posterior probabilities ranging between 0.002 and 0.19. The land use sequence with the highest posterior probability is planted pasture – planted forest – sugar cane – rangeland – crops. An analysis of all other sequences and their posterior probabilities reveals that there is a dichotomy in this most common sequence. Planted pasture, planted forest and sugar cane usually (in about 80% of the realizations) come in the first part of the sequence, and rangeland and crops in the last part, but the order among them fluctuates.

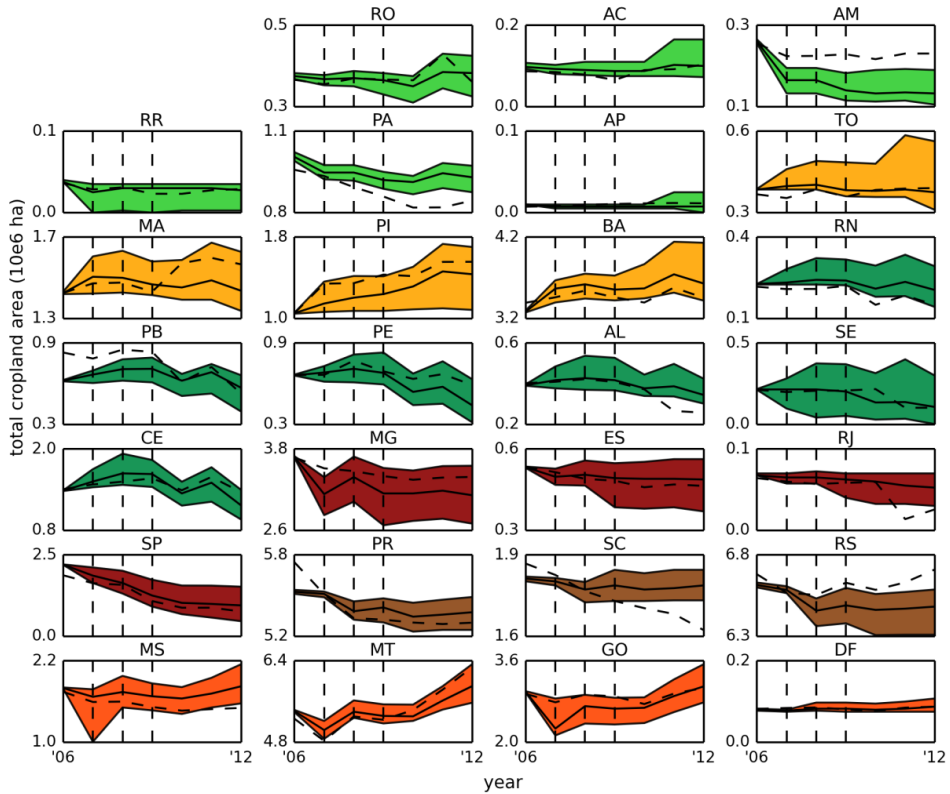
Sugar cane can only displace land use types coming after it in the sequence. So, in 80% of the realizations it predominantly replaces crops and rangeland. An important consequence of this calibration result with regard to iLUC is that the iLUC within a macro region will originate mainly from the displacement of crops and rangeland, which is in line with the findings of Lapola *et al.* (2010).

The weights of the suitability factors have been calibrated as well (Table 5.3). In general, suitability factor  $k = 1$ , representing  $n$  in the neighbourhood, obtains a high weight. This means that a land use type is likely to expand in regions in which it is already cultivated. This factor has the highest median posterior weight for rangeland, sugar cane and planted pasture. Accordingly, dLUC will take place close to existing sugar cane patches. For cropland the double cropping potential ( $k = 5$ ) is the most important suitability factor for expansion. Galford *et al.* (2008) have found in a case study in Matto Grosso (an important expansion region, see Figure 5.2) by means of remote sensing that newly established cropland is usually single cropped, but is converted to double cropping after two to three years. The high weight for the double cropping potential factor indicates that this potential already plays a role at the establishment of the cropland, while the actual implementation of double cropping takes place a few years later. As a consequence, the location of iLUC in the case of displaced cropland is likely to be a location with a high double cropping potential. In conclusion, the calibrated land use change model mainly influences the location of dLUC and iLUC within the macro region, i.e. distribution between and also within states in a macro region.

**Table 5.3: The mean, first quartile and third quartile of the weights of the suitability factors  $k$  of all active land use types  $n$  resulting from the calibration.**

$n$	name	$k$	suitability factors	1 <sup>st</sup> quartile	median	3 <sup>rd</sup> quartile
4	rangeland	1	$n$ in the neighbourhood	0.35	0.46	0.60
		2	distance to roads	0.24	0.35	0.45
		3	potential yield of $n$	0.03	0.19	0.28
5	planted forest	1	$n$ in the neighbourhood	0.19	0.29	0.36
		2	distance to roads	0.27	0.34	0.37
		3	potential yield of $n$	0.32	0.37	0.51
6	crops	1	$n$ in the neighbourhood	0.11	0.22	0.36
		2	travel time to hubs for $n$	0.05	0.14	0.22
		3	potential yield of $n$	0.05	0.11	0.20
		4	conversion elasticity	0.12	0.23	0.33
		5	growing season length	0.21	0.30	0.36
8	sugar cane	1	$n$ in the neighbourhood	0.23	0.29	0.36
		2	travel time to hubs for $n$	0.21	0.28	0.33
		3	potential yield of $n$	0.15	0.22	0.26
		4	conversion elasticity	0.17	0.21	0.24
9	planted pasture	1	$n$ in the neighbourhood	0.40	0.53	0.66
		2	distance to hubs for $n$	0.33	0.45	0.56
		3	potential yield of $n$	0.00	0.02	0.03

The results above were based on calibration over 2007-2012. To show the effect of calibration, we have applied a split-sample approach, with calibration only from 2007 to 2009, to allow a comparison with observational data in the validation period from 2010 to 2012. We compare the modelled against observed area of cropland (Figure 5.5), because this land use type gives the most information on model performance since it both expands and contracts in the calibration period. For most states without a clear break in their trend, for example Roraima (RR), Ceará (CE), and Mato Grosso (MT), the modelled median remains good in the validation period. However, the areas of states that do show a trend break, e.g. Maranhão (MA) and Rio de Janeiro (RJ), are poorly simulated, although at least for Maranhão the observed cropland area falls within the 95% confidence interval of the modelled area.



**Figure 5.5: Modelled and observed areas per state, for the land use type cropland as an example. Vertical dashed lines: calibration years, dashed lines: observed area, solid lines: modelled median area, colored planes: 95% confidence interval of the modelled area, where the color corresponds to the color of the macro region in Figure 5.2, to be able to quickly see which states belong to the same macro region. The state name abbreviations are: AC = Acre, AL = Alagoas, AM = Amazonas, AP = Amapá, BA = Bahia, CE = Ceará, DF = Distrito Federal, ES = Espírito Santo, GO = Goiás, MA = Maranhão, MG = Minas Gerais, MS = Mato Grosso do Sul, MT = Mato Grosso, PA = Pará, PB = Paraíba, PI = Piauí, PR = Paraná, RJ = Rio de Janeiro, RN = Rio Grande do Norte, RO = Rondônia, RR = Roraima, RS = Rio Grande do Sul, SC = Santa Catarina, SE = Sergipe, SP = São Paulo, TO = Tocantins.**

**Table 5.4: Reduction in root mean square error (%) of the median area of land use type  $n$  for the calibrated model compared to the reference case (Monte Carlo run without particle filter), summed over all states, given per year for the validation period (2010 – 2012).**

$n$	land use type	year		
		2010	2011	2012
4	rangeland	-1	-2	-4
5	planted forest	-6	0	2
6	crops	53	50	42
8	sugar cane	26	21	24
9	planted pasture	33	35	25

To summarize the effect of calibration for all land use types, we compare the root mean square error (RMSE) in area summed over all states of the calibrated and non-calibrated model (Table 5.4). For crops, sugar cane and planted pasture a considerable RMSE reduction is achieved. The highest reduction is achieved for crops, with a maximum of 53% in 2010. For sugar cane the average reduction is 24%. A significant reduction for sugar cane is important, since it is the land use type of main interest. Being able to correctly project the location of sugar cane expansion, connotes correct modelling of dLUC, which is the first step in also correctly projecting iLUC since the two are chained. For rangeland and planted forest the calibration does not bring the modelled median area per state closer to the observed area. The modelled median even becomes worse, although not significantly, only a few percent. The reason why PLUC cannot find weights for the suitability factors that result in a correct projection, is probably the poor data availability for these two land use types. For example, for the initial land use map, no good prior knowledge maps were available (see Appendix A), and for the suitability factors we have no information about the locations of the hubs for these land use types.

### *Economic model, projection*

Demands are projected by MAGNET per land use type for 2013, 2015, 2020, 2025 and 2030. To illustrate the trend, the demands for 2020 and 2030 for the reference and the biofuel scenario and the uncertainty herein are shown in Figure 5.4. An interesting result is that the uncertainty within a scenario is often higher than the difference between the scenarios. This indicates that it can be problematic to draw conclusions about the effect of, for example, a policy by means of comparing scenarios from the CGE model. If the land transition elasticities are uncorrelated between the two scenarios, the large uncertainty makes that the policy effects might be negative as well as positive. Yet, we believe that although the elasticities are uncertain, they are correlated between the two scenarios, as these scenarios represent the same system, as long as the difference between scenarios is not too large. Others doubt this; a discussion that is known in economic modelling as the Lucas critique. Lucas (1976) argues in his work that the parameters in economic models are not policy-invariant, and that they would therefore change when a policy is implemented. This discussion is interesting, but goes beyond the scope of this paper. Nevertheless, we should be aware, that if Lucas is correct, the uncertainties in dLUC and iLUC shown in the next sections might be significantly higher.

In the reference scenario sugar cane mainly expands in the Center West Cerrado and the Southeast (together called the Central South). The extra demand for sugar cane for ethanol from the mandates (biofuel scenario) also mainly ends up in these

two regions. In the biofuel scenario, the total area of sugar cane in the Center West Cerrado almost triples by 2030 compared to 2012.

The difference between the reference scenario and the biofuel scenario for the other land use types within Brazil is the largest in the Southeast (Figure 5.4). In this macro region, the areas of crops and rangeland are significantly smaller in the biofuel scenario than in the reference scenario. As the productivity of all land use types are roughly the same in these two scenarios, this decrease in area means that MAGNET assumes that these areas of crops and rangeland are displaced by sugar cane. The displaced land uses are shifted to the Northeast Cerrado and the Northern Amazon: here crops and rangeland occupy a larger area in the biofuel scenario than in the reference scenario (Figure 5.4, difference between the hatched and plain bars).

### *Conceptual differences between the two models*

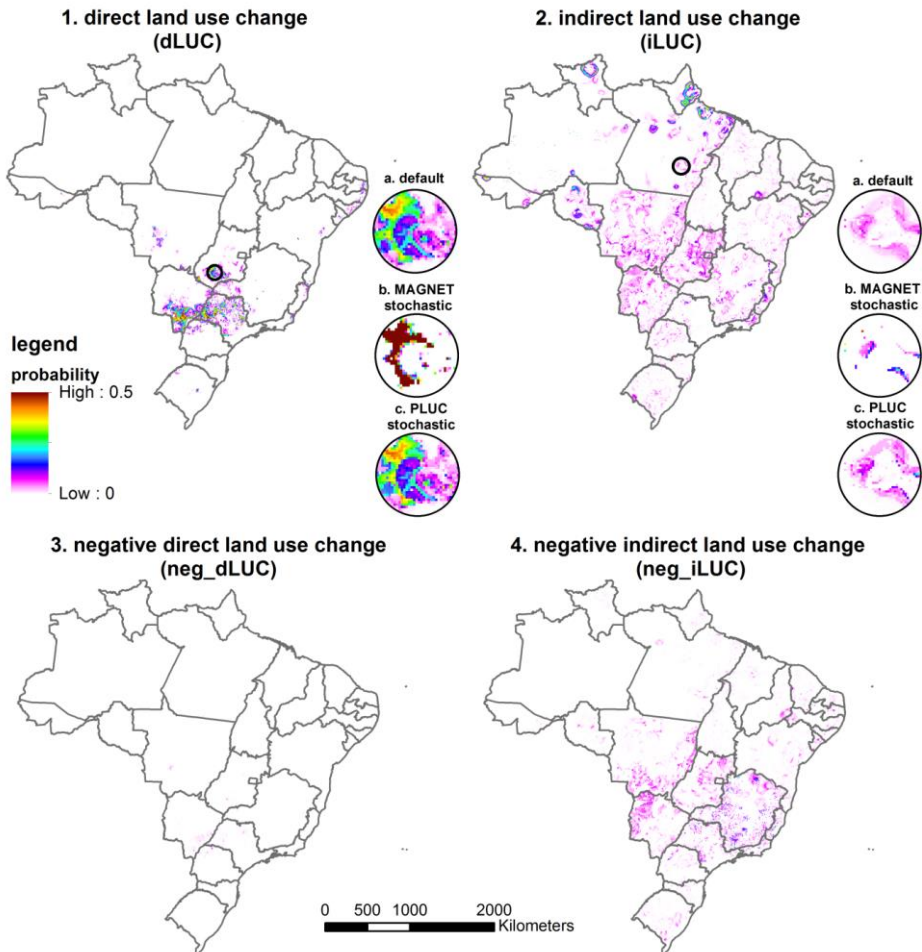
Despite the 'shared' conversion matrix between MAGNET and PLUC and despite the fact that MAGNET provides the demand as an input for PLUC, the conversion dynamics between the two models differ because of conceptual differences between the models. A result of this is that the area of iLUC for the whole of Brazil calculated from MAGNET differs from the area calculated by PLUC (further discussed later on), although ideally these two would be the same. This problem does not occur for dLUC, only for iLUC, and in the following we explain why.

The origin of the problem is that in PLUC sugar cane expands in the projection period, besides over cropland and rangeland, also often over planted pasture; a displacement also observed in other studies (e.g. Rudorff et al., 2010, Adami et al., 2012b). In MAGNET, however, the area of conversion from planted pasture to sugar cane is negligible. This displacement of planted pasture in PLUC, not present in MAGNET, has two effects. One is that in PLUC the area of planted pasture is decreased in a region, such that, in that same year, planted pasture should expand (in addition to the expansion caused by a potential increase in demand already given by MAGNET) elsewhere in that region in order to make up for the lost acreage. This causes iLUC in the region, not anticipated by MAGNET. Another effect is that the areas of rangeland and/or crops in PLUC are larger than dictated by the demand in MAGNET for that year, so that these land uses will contract, resulting in abandoned land. This causes negative iLUC, which by definition never occurs in MAGNET. It can be debated which of the two models, if any, is correct. But the most important implication for our study is that uncertainty in iLUC projections does not only stem from uncertainties of parameters and model structure within one model component, but also from the dissimilarity in model concepts between the two models within the integrated model chain.

### ***5.3.2. DLUC and iLUC projections for Brazil up to 2030 at different spatial scales and the uncertainty herein***

The direct land use change as a result of an increased ethanol production from 2013 to 2030 mainly takes place in the Central South region. The highest cell (5 x 5 km<sup>2</sup>) based probabilities, up to 0.77, exist in the South of Mato Grosso do Sul and the West of São Paulo (Figure 5.6, frame 1). The highest probabilities of indirect land use change, with a maximum of 0.43, occur in the Amazonian states Rondônia, Amapá, and Roraima (Figure 5.6, frame 2). Probabilities in these three small states are high because they are the only places in the Northern Amazon where any agricultural land use type can expand, as the rest of the Northern Amazon has very few roads, almost no existing agriculture and thus few hubs, and many protected areas. Implementation of new roads in the Amazon could change the spatial distribution drastically, but this is not included into the model due to limited spatial planning data availability. In the other macro regions, there are more options for expansion, and there is more variation in the suitability maps (best locations for expansion) between the different land use types and between the individual Monte Carlo realizations, i.e. more uncertainty. In these other macro regions, iLUC locations with high probabilities are the frontier of the sugar cane expansion area (Goiás, Mato Grosso do Sul, and Mato Grosso) as well as the 'arc of deforestation', the transition area from cultivated land to mainly natural vegetation (Mato Grosso, Pará and Rondônia).

As expected, there are only very few cells experiencing negative dLUC, and with negligibly low probabilities, with a maximum of 0.07 (Figure 5.6, frame 3). Conversely, negative iLUC (land abandonment in the biofuel scenario and not in the reference scenario, Figure 5.6, frame 4) does appear, with probabilities up to 0.48, mainly in Espírito Santo, Minas Gerais and the Pantanal, which is the wetland area in the West of Mato Grosso do Sul and the South of Mato Grosso. These are areas where the suitability for most agriculture is low, resulting in land abandonment when the demand in the biofuel scenario is lower than in the reference scenario (see also the discussion in the previous section). With lower probabilities, up to 0.1, this effect also occurs in the rest of the Central South.

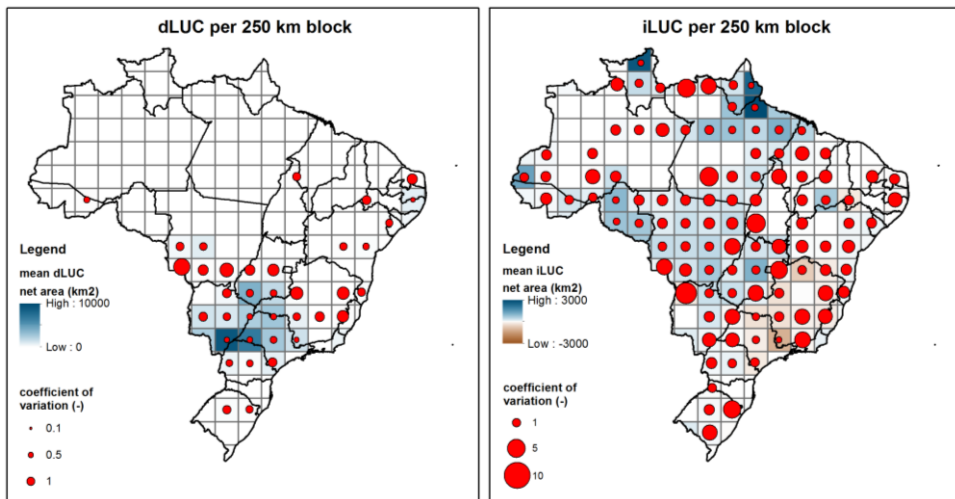


**Figure 5.6: Probability of 1) dLUC, 2) iLUC, 3) negative dLUC, and 4) negative iLUC per grid cell. Probabilities are shown at a common scale from 0 to 0.5. Only for dLUC a few cells with higher probabilities exist, up to 0.77, but stretching the scale up to 0.77 reduces discernibility between different cells with low probability in all four maps. For dLUC and iLUC map detail frames at a location in the expansion area are provided, showing probabilities for the three runs: a) the default with both models stochastic, b) with only PLUC stochastic, and c) with only MAGNET stochastic.**

In the following we add up dLUC and negative dLUC, and iLUC and negative iLUC, obtaining net dLUC and net iLUC. Scaling up to 250 x 250 km<sup>2</sup> blocks (Figure 5.7), we can calculate the coefficient of variation (cv), indicating relative uncertainty in dLUC and iLUC. Clearly, the uncertainty in iLUC is generally larger than in dLUC. The median cv over all selected blocks for dLUC is 0.91, while for iLUC it is 1.61. This is caused by the fact that dLUC is affected by the dynamics of sugar cane only, while

iLUC is an effect of the interplay of all land use types, thereby being subjected to the uncertainties in all weights of all suitability factors (Table 5.3) and the order of allocation. The maximum cv value of dLUC is 4. The maxima occur at the expansion frontier of sugar cane, through Mato Grosso, Goiás and Minas Gerais. The maximum cv value of iLUC is 6. This means that, when considering the 95% confidence interval, mean iLUC values might be as much as 13 times as high or as low. In a nutshell, the uncertainty in these blocks is so high that we can say practically nothing about expected iLUC there, except when mean iLUC is (very close to) zero (13 times zero is still zero). Coefficients of variation in the arc of deforestation generally range from 1 to 3, which is a bit better, but still very uncertain.

Comparing Figure 5.6 and Figure 5.7 it becomes apparent that blocks of maxima in cv of iLUC area correspond to regions where both iLUC and negative iLUC might appear. When some Monte Carlo realizations have negative iLUC values and others positive iLUC values, the standard deviation is large, and correspondingly the cv. In these blocks the net iLUC effect might be positive as well as negative, so that the impacts on for example biodiversity, might be negative as well as positive respectively.



**Figure 5.7: Mean net area (km<sup>2</sup>) (colour of the block) and the coefficient of variation (cv) (-) (size of the red circle) of dLUC (left) and iLUC (right), per 250 x 250 km<sup>2</sup> block. For the display of the cv, blocks smaller than 31250 km<sup>2</sup> (half of a 250 x 250 km<sup>2</sup> block, occurring at the map edges) are filtered out, as the cv is heavily influenced by the support size of the block. Also blocks with mean dLUC or iLUC smaller than 25 km<sup>2</sup> (one cell) are filtered out, because when the mean goes to zero, the cv becomes infinite.**



When looking at the cv for dLUC and iLUC area for Brazil as a whole (Table 5.5, both models stochastic), the values are many times smaller. The total amount for the whole of Brazil can be determined about nine times as precise for dLUC and about two times as precise for iLUC compared to the median of the 250x250 km<sup>2</sup> blocks, because small scale errors balance each other out when aggregating.

### *5.3.3. Contribution of the economic and land use change model to the uncertainty in dLUC and iLUC at different spatial scales*

The cv of the dLUC area at national level is for 100% caused by MAGNET (Table 5.5), which is logical, as MAGNET determines the total demand for sugar cane, and PLUC only allocates it within the macro regions. For iLUC area this is not the case. The cv value of iLUC area for the run with only MAGNET stochastic is about sixteen times higher than cv value of iLUC area for the run with only PLUC stochastic, so about 93% of uncertainty in iLUC stems from MAGNET. Yet, the exact contribution of both models cannot be determined, because errors from the two models partly compensate each other (the cv values of the two runs do not add up to the cv value of 0.72 found in the default run). The reason that uncertainty in iLUC at national level is not fully determined by MAGNET is that in PLUC iLUC can occur within a macro region that is additional to the iLUC between macro regions from MAGNET.

At the grid cell level, many cells have some probability of experiencing dLUC or iLUC when both models are stochastic (Figure 5.6, detail frame 1a.). With only PLUC stochastic and MAGNET deterministic, there is in general not much difference with the results of the default run (Figure 5.6, panel 1c.), although somewhat fewer cells have a probability above zero on dLUC and iLUC. This indicates that only a small part of the uncertainty at cell level is caused by MAGNET. With only MAGNET stochastic and PLUC deterministic, compared to the default run less than half of the cells have a probability above zero on dLUC and iLUC, and the ones that have, have a relatively high probability, indicating much lower uncertainty (Figure 5.6, panel 1b.). The uncertainty now mainly exist at the edges of the expansion patches, caused by the variation in demand from MAGNET. For iLUC (Figure 5.6, panels 2a., 2b., 2c.) the same reasoning applies. In conclusion, uncertainty at grid cell level is mainly caused by uncertainty in PLUC, for both dLUC and iLUC.

**Table 5.5: Total area (Mha), standard deviations (sd) (Mha), and coefficients of variation (cv) (-) of dLUC and iLUC for Brazil for three different runs: 1) both the land use change model and the economic model stochastic, 2) land use change model stochastic and the economic model deterministic, 3) the land use change model deterministic and the economic model stochastic**

run	mean dLUC	sd dLUC	cv dLUC	mean iLUC	sd iLUC	cv iLUC
both stochastic	4.20	0.10	0.02	3.13	2.25	0.72
PLUC stochastic	4.21	0.00	0.00	3.15	0.61	0.06
MAGNET stochastic	4.20	0.10	0.02	2.62	2.19	0.84

## 5.4. Discussion

In this paper we have demonstrated a general methodology to calculate direct and indirect land use change (dLUC and iLUC) stochastically with an integrated economic – land use change model, taking into account important uncertainties in all components of the modelling chain. The proficiencies of this methodology were shown for a case study of land use change in Brazil up to 2030, steered by current and planned ethanol mandates worldwide. Here, we shortly discuss the answers to our three research questions and give recommendations for further studies.

### 5.4.1. What are the dLUC and iLUC projections for Brazil up to 2030 at different spatial scales and what is the uncertainty herein?

Cell (5 x 5 km<sup>2</sup>) based probabilities of dLUC range from 0 to 0.77, and of iLUC from 0 to 0.43. Thus, given our scenario assumptions, there is no single cell in Brazil for which it can be said with certainty that dLUC or iLUC will take place up to 2030. So, it is difficult to project exactly where dLUC and iLUC will occur, but it is certain that it will occur (there are no Monte Carlo realization without dLUC or iLUC effects). Yet, overall locations of iLUC are in line with the locations projected by Lapola *et al.* (2010). For dLUC our study shows some locations with high probabilities in Mato Grosso do Sul and Goiás, where Lapola *et al.* (2010) do not project dLUC. As there are the ‘new’ expansion areas, this inconsistency is likely caused by the fact that their projection is for 2020 while we project up to 2030. Also, in our projections there are many cells for which it can be concluded with certainty that no dLUC or iLUC will take place in 2030, which is surely relevant information. In 250 x 250 km<sup>2</sup> blocks, the coefficient of variation (cv) ranges from 0 to 4 for dLUC and from 0 to 6 for iLUC. Large cv values for dLUC occur at the frontier of sugar cane expansion. High cv values for iLUC occur where both iLUC and the opposite effect (agriculture in the reference scenario is abandoned land in the biofuel scenario), introduced in this paper, might take place.

The uncertainty in iLUC area and location is generally higher than in dLUC, because iLUC is caused by the interplay of various land use types that each have their

uncertain model parameters, while dLUC is mainly affected by the parameters for sugar cane. Uncertainty in dLUC and iLUC is lower at higher aggregation levels. For iLUC the decrease in uncertainty by aggregation is smaller. At country level, the cv for iLUC is 36 times higher than for dLUC in our case study. At this level, dLUC can be projected with high certainty, having a cv of only 0.02, while iLUC is still uncertain, having a cv of 0.72. Thus, to answer the question posed in the title, what we can and cannot say about iLUC: we can merely say things about iLUC with high uncertainties. Estimated iLUC areas, even at country level, might as well be 2.4 times as high or as low, given the 95% confidence interval.

#### ***5.4.2. What are the sources of uncertainty for each step in the model chain and how do these uncertainties influence dLUC and iLUC projections?***

Uncertain components in the land use change model are 1) the initial land use map, causing uncertainty at cell level, 2) the order of allocation of the land uses, causing uncertainty in especially iLUC, and 3) the selection and weights of the suitability factors for allocation of the land use types, causing mainly uncertainty at intermediate aggregation levels, like states. The reduction in root mean square error in the modelled median land use areas per state by model calibration, compared to a non-calibrated model is on average 20%. Poor-performing land use types are rangeland and planted forest, probably due to poor data availability for the drivers of location of land use change.

For the economic model we have assessed only the effect of one of the most critical parameters: the land transition elasticities that simulate the likelihood of particular land transitions. The uncertainty caused by varying these elasticities mainly plays a role at national level. At the cell level, it only causes uncertainty at the edge of the patch of expansion.

The final aspect generating uncertainty in the output is the difference in the conceptual model between the economic and the land use change model concerning the reclaiming of abandoned land. This conceptual difference affects the total amount of iLUC and the opposite iLUC effect.

#### ***5.4.3. What is the contribution of the economic and land use change model to the uncertainty in dLUC and iLUC at the different spatial scales?***

At the cell level, uncertainty is primarily determined by the land use change model. Going to higher aggregation levels the influence of the uncertainty in the Computable General Equilibrium (CGE) model on output uncertainty increases. At national level, the cv of dLUC is caused by the CGE model for 100%. The contribution of the economic model to the cv of iLUC at this level is about 93%,

although this cannot be determined precisely, because errors from the two models partly compensate each other.

#### ***5.4.4. Implications and recommendations***

From the above we can conclude that projected iLUC areas and locations are highly uncertain. Based on the case study, our opinion is that threshold evaluation for iLUC indicators should not be implemented in legislation. Thresholds (cf. Malins, 2013) have no use when the model, used to check whether an indicator for a specific case is above or below this threshold, gives an output confidence interval that straddles the threshold. This is likely to happen considering the high uncertainties found in our study. As most iLUC (or LUC) indicators in legislation are provided in terms of greenhouse gas (GHG) emissions generated, the impacts of the uncertainty in dLUC and iLUC projections on GHG emissions should be assessed in order to underpin our conclusion. Error propagation assessment for other impacts, like biodiversity and water availability, is also desirable. Our opposition to thresholds for iLUC factors in legislation does not mean we favour negligence of biofuel induced land use change. We propose, in line with e.g., Finkbeiner (2014) and Mathews and Tan (2009), a change of focus from quantifying iLUC to taking proactive measures to mitigate iLUC, even though the effectiveness of these measures might be difficult to quantify.

Our quantification of the sources of uncertainty allows identification of the parts of the modelling chain having the highest priority for improvement. If one wants better estimates of dLUC and iLUC at cell level for a given case study, for example to be able to better quantify local GHG emissions caused by the biofuel targets, one should focus on improving the land use change model. Spatially explicit input data could be improved, especially for land use types that are problematic to derive from remote sensing: rangeland (problematic to distinguish from natural savannah) and planted forest (problematic to distinguish from natural forest). And data that are now only included at an aggregate level, like land management and yield level, could be included spatially to account for spatial variation. Also, better information on data accuracy would be helpful. Due to the lack of accuracy information we had to make strong assumptions on the errors in the maps used to create the initial land use map and the observational data used for calibration.

If one wants better estimates of dLUC and iLUC at country level, one should focus on improving the economic model. Our current estimates of uncertainty in the CGE model might be underestimated, because we have evaluated the uncertainty from the land transition elasticities only, while land use changes might be sensitive to other parameters as well (Kavallari et al., 2014). Yet, making other parameters

stochastic could also reduce uncertainty, when they cancel out each other's errors. It would be good if making parameters stochastic and running Monte Carlo simulations would become common practice in economic modelling. Another option for better country level estimates might be the usage of a whole different type of model or tool, obviously also stochastic, although our current study gives cannot ascertain whether and to what extent that could reduce uncertainty.

One thing that could improve iLUC estimates at all spatial scales is a better match between the economic and land use change model. The best solution would be to link the economic and the land use model with a hard link that includes a feedback, as also suggested by Wicke *et al.* (2015). However, there is an inherent risk that this feedback loop is infinite, meaning that the land use dynamics cannot be resolved, and there are many technical obstacles that complicate hard linking.

Yet, even if improved models, or improved model connections are used, in all cases we strongly advise to provide quantitative uncertainty estimates together with the calculated dLUC, iLUC or LUC indicators so that users of these indicators can evaluate the reliability of the indicators.

## **5.5. Acknowledgements**

This work was carried out within the BE-Basic R&D Program, which was granted a FES subsidy from the Dutch Ministry of Economic affairs, agriculture and innovation (EL&I). We thank two anonymous reviewers for their valuable comments.

## **5.6. Appendix A: Initial land use map**

The demand time series (tabular area data) for the initial year of simulation, 2006, should match the total area per land use type in the initial land use map. If they do not, as in our case, three options exist to harmonize the two: 1) adapt the time series to the map, 2) adapt the map to the time series, and 3) do something intermediate. The first option is relatively easy. One can derive relative increase and/or decreases in area from the time series and apply this trend to the total area in the initial land use map. But if, for the case study under consideration, one has a higher confidence in the area values of the time series than in the map, or if there is no single map containing all land use types one wants to model, the second or third option is preferable. Because for Brazil no recent, coherent land use map is available for the whole country, we apply the second method, adapting the map to the time series.

You and Wood (2005) provide a deterministic method to create a land use map that matches the time series, by spatially disaggregating land use areas per administrative region from the time series into raster cells within that region. They base the disaggregation on prior knowledge about the likely distribution of the land use types, originating from one or several land use maps and potentially other sources like maps of potential yield for crops, farm size or cropping intensity.

We apply the procedure of You and Wood (2005) to create an initial land use map for Brazil on a raster with a cell size of 1 km<sup>2</sup>. Differently from You and Wood, in our method the land uses are disaggregated sequentially instead of simultaneously, because our confidence in the prior knowledge maps (Table 5.6) varies for the different land use types. In addition to the prior knowledge maps given in Table 5.6, a map of natural remnants (CSR/Ibama et al., 2013) is used as prior knowledge about where it is *unlikely* that agricultural (unnatural) land uses, i.e. planted forests, crops, sugar cane and planted pasture, are present. As administrative regions we use municipalities, the smallest units in the IBGE Census data (IBGE, 2006), 5566 in total for Brazil. The result of the disaggregation is a map  $\mathbf{a}_n$  for every land use type  $n$ , giving per cell the fraction that is occupied by that land use type  $n$ . Cells for which the sum of the fractions of all land use types is smaller than one ( $\sum_{n=1}^{10}(\mathbf{a}_n) < 1$ ) can be 'filled' up to one with natural land uses. In our case study these cells are filled with natural forest<sup>2</sup> and grassland, relative to the shares in the GlobCover class (Arino et al., 2008) of that cell.

Some land use change models work with fractions of land use types per cell, e.g., IMAGE (Bouwman et al., 2006) and CLUE (Verburg et al., 1999), but PLUC (Chapter 2) can handle only a single land use type per cell, similar to CLUE-S (Verburg et al., 2002). Therefore, the collection of fraction maps has to be transformed into one nominal land use map. For every land use type  $n$ ,  $\mathbf{a}_n$  is first converted to a Boolean map, keeping the total area per land use type the same as in the tabular data. This can be done by sorting the map  $\mathbf{a}_n$  for each administrative region separately, and assigning cells to land use type  $n$  (set Boolean true) starting with the highest value in  $\mathbf{a}_n$  until the total area of these cells equals the area for  $n$  in the tabular data in that administrative region. Again, this should be sequentially for all land use types, masking out cells that are already occupied by previously allocated land use types. In the end, the Boolean maps for all  $n$  are combined to create a single nominal land use map. This is performed using a cell size of 5 x 5 km<sup>2</sup>.

---

<sup>2</sup> Note that this produces a map with a larger area of natural forest than indicated by the tabular data. The tabular data for forest in Brazil from the IBGE Census, however, report natural forest areas in the Amazon that we believed were questionably low, justifying our approach.

**Table 5.6: Sources of the tabular area per land use type per municipality, sources of the prior knowledge maps, and the standard deviation,  $\sigma_{a,n}$ , used to create an ensemble of initial land use maps (Equation 5.2).**

$n$	name	source(s) of tabular area per LU type in 2006	source(s) of prior knowledge map(s)	$\sigma_{a,n}$
1	urban	GlobCover (Arino et al., 2008)	GlobCover (Arino et al., 2008)	0
2	water	GlobCover (Arino et al., 2008)	GlobCover (Arino et al., 2008)	0
3	natural forest	IBGE Census (IBGE, 2006) <sup>1</sup>	GlobCover (Arino et al., 2008)	0.1
4	rangeland	IBGE Census (IBGE, 2006) <sup>1</sup>	GlobCover (Arino et al., 2008) and Global Pasture data (Ramankutty and Foley, 1999) <sup>2</sup>	0.1
5	planted forest	IBGE Census (IBGE, 2006) <sup>1</sup>	GlobCover (Arino et al., 2008) <sup>3</sup>	0.1
6	crops	IBGE Census (IBGE, 2006) <sup>1</sup> and Conab (Conab, 2014) <sup>4</sup>	GlobCover (Arino et al., 2008)	0.1
7	grass and shrubs	-	-	0.1
8	sugar cane	Canasat (Rudorff et al., 2010) <sup>5</sup> , IBGE Census (IBGE, 2006) <sup>1</sup> and Conab (Conab, 2014) <sup>4</sup>	Canasat (Rudorff et al., 2010) <sup>5</sup> and M3 crop data (Monfreda et al., 2008)	0.1
9	planted pasture	IBGE Census (IBGE, 2006) <sup>1</sup>	PROBIO (MMA, 2008)	0.1
10	bare soil	GlobCover (Arino et al., 2008)	GlobCover (Arino et al., 2008)	0

<sup>1</sup> NoData in the IBGE data is presumed to be zero, following IBGE's own approach for scaling up from municipality to e.g., state level.

<sup>2</sup> Because the Global pasture data is coarse, the correlation between the presence of pasture and the GlobCover classes was assessed. Next, each cell was assigned the correlation value based on its GlobCover class, and this map (thus having the fine resolution of GlobCover) was used as prior knowledge map.

<sup>3</sup> Locations of planted forest in GlobCover are distinguished from locations of natural forest based on the map of natural remnants (CSR/Ibama et al., 2013). Planted forest is allocated outside of the remnants and natural forest inside, furthermore based on the same prior knowledge map.

<sup>4</sup> The Conab data is used to calculate for each crop the percentage that is harvested as second or third crop. The Conab data is per state, so for all municipalities in that state the IBGE area is reduced by that percentage, because second and third crops do not require land (they are cultivated on the same land as the first harvest).

<sup>5</sup> The Canasat map covers only the Central South region in Brazil, which includes the states Goiás, Minas Gerais, Mato Grosso, Mato Grosso do Sul, Paraná and São Paulo from 2003 onwards and Espírito Santo and Rio de Janeiro from 2010 onwards. Therefore, we use IBGE Census data in combination with Conab data (see 4) as sources of tabular area data per LU type in 2006 and the M3 crop data as prior knowledge map for the other states.

Uncertainty is taken into account by adding random noise to  $\mathbf{a}_n$ . If this is done for each Monte Carlo realization, every realization gets a slightly different initial land use map, varying especially at locations at which the value of  $\mathbf{a}_n$  is close the point at which the total area for  $n$  is met. Under the assumption that the error of  $\mathbf{a}_n$  is linearly related to  $\mathbf{a}_n$ , the error model is defined as:

$$\mathbf{A}_n = \mathbf{a}_n + Z_a \quad 5.2$$

with  $Z_a \sim N(0, \sigma_{\mathbf{a},n} * \mathbf{a}_n)$

In Equation 5.2,  $\mathbf{A}_n$  is the stochastic fraction map for land use type  $n$ , and  $\sigma_{\mathbf{a},n}$  is the standard deviation. For all land use types  $\sigma_{\mathbf{a},n}$  is set at 0.1, except for the land use types which are the easiest to classify in the remote sensing image and for which the area as well as the prior knowledge are therefore derived from GlobCover; for these  $\sigma_{\mathbf{a},n}$  is set to 0 (Table 5.6).

## 5.7. Appendix B: Land use change model

The total suitability map  $\mathbf{s}_{n,t}$  of the active land use types is defined as:

$$\mathbf{s}_{n,t} = \sum_{k=1}^{K_n} (w_{n,k} \cdot \mathbf{u}_{n,k,t}), \text{ for each active } n \text{ in each } t \quad 5.3$$

with  $\sum_{k=1}^{K_n} (w_{n,k}) = 1$

In Equation 5.3,  $k$  is the suitability factor, with  $k = 1, 2, \dots, K_n$  (each active land use type  $n$  can have a different number of suitability factors). Furthermore,  $\mathbf{u}_{n,k,t} \in [0,1]$  is the normalized suitability map for land use type  $n$  for suitability factor  $k$  at time  $t$ ; and  $w_{n,k} \in [0,1]$  is the weight of factor  $k$ , denoting the importance of the drivers of location in the total suitability map  $\mathbf{s}_{n,t}$ . The weights are part of the model structure (Chapter 4). Another structural element is the order in which the land uses allocate their demands. These two elements are modelled stochastically, i.e. defined by a probability distribution of all possible values, to take into account model structure uncertainty. For the stochastic model we use a Monte Carlo simulation with 5000 realizations. The prior probability distribution of the weights of the suitability factors in this ensemble is defined as:

$$w_{n,k} = \frac{Z_{w_{n,k}}}{\sum_{k=1}^{K_n} (Z_{w_{n,k}})}, \quad \text{with } Z_{w_{n,k}} \sim U(0,1), \quad 5.4$$

for each active  $n$  for each  $k$

In Equation 5.4,  $U(0,1)$  denotes a uniform distribution between zero and one. Equation 5.4 ensures that the sum of the stochastic weights is 1. The order of the active land use types is also randomized, meaning that the order in which the land use types are allocated also differs per Monte Carlo realization. At the start of the



simulation, for each Monte Carlo realization, a weight is drawn for each suitability factor for each active land use type and a sequence is drawn. Given the fact that we model five dynamic land use types, the prior probability distribution of this sequence consists of  $5! = 120$  possible sequences, which are given equal prior probabilities. The posterior probabilities of the weights and allocation sequence is determined by the calibration (section 5.2.4).



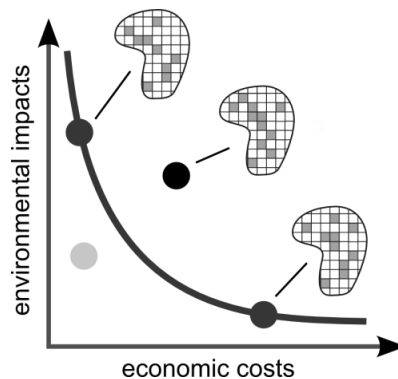
## 6. A spatial optimization approach to find trade-offs between production costs and greenhouse gas emissions: a bioethanol case study

Judith A. Versteegen, Jan Gerrit Geurt Jonker, Derek Karssenbergh, Floor van der Hilst, Oliver Schmitz, Steven M. de Jong, and André P.C. Faaij (in prep.).

**Abstract** - Ideally, land is used for different food, feed, fibre and bioenergy commodities in a way that economic benefits coincide with positive environmental impacts. Regrettably, often trade-offs exist between economic and environmental objectives. A Pareto frontier can quantitatively express these trade-offs. In this paper we demonstrate how a Pareto frontier can be constructed for a supply increase of a given commodity and how information derived from it could aid to formulate management or policy recommendations. This is illustrated by applying a spatially explicit optimization approach to a case study in which we aim to minimize production costs and GHG emissions of bioethanol for the state Goiás, Brazil, in 2030, for different carbon prices. We find that the Pareto frontier ranges between minimum costs of 656 US\$<sub>2014</sub> / m<sup>3</sup> ethanol and minimum GHG emissions of -399 kg CO<sub>2</sub>-eq / m<sup>3</sup> ethanol, i.e. carbon sequestration. At a moderate increase in carbon price from 0 to 10 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq, the frontier shows a steep decrease in GHG emissions, a reduction of more than 50%, while production costs practically remain at their minimum, signifying an excellent opportunity for Brazilian bioethanol emission reduction policies. The developed methodology has the prospect to identify trade-offs and win-win situations for other regions, scales, objectives and commodities.

## 6.1. Introduction

The intensifying pressure on the scarcely available land, by e.g., food, feed, fibre and bioenergy production requirements, leads to a growing debate about the potential negative environmental impacts of land use changes (Lambin and Meyfroidt, 2011). In the most sustainable situation, land is used for the different commodities in such a way that economic benefits are combined with positive environmental impacts. But in reality, trade-offs exist between economic and environmental objectives. Such trade-offs can be expressed quantitatively in a Pareto frontier, which shows all (in this case land use) alternatives for which it is impossible to improve one objective, without impairing another (Figure 6.1). In other words, it gives all 'optimal alternatives' given the different objectives (Seppelt et al., 2013). A Pareto frontier can help to structure the aforementioned debate, by showing the actors involved how much will be lost in one objective for a gain in another objective. Our aim is to demonstrate how a Pareto frontier between economic and environmental objectives can be constructed for a supply increase of a given commodity and how information derived from it could aid to formulate management or policy recommendations.



**Figure 6.1: The concept of a Pareto frontier (solid line) between an economic and an environmental objective with two optimal alternatives (dark grey dots) on the frontier and the spatial configurations belonging to these alternatives. In the spatial configuration the grey cells are cultivated for the commodity under consideration and the white ones are not. The black dot, accompanied by its the spatial configuration, is a suboptimal alternative. The light grey dot is an infeasible point and consequently has no spatial configuration belonging to it.**

Many constituents of economic and environmental objectives vary widely over space. For example, land use change related greenhouse gas (GHG) emissions depend on previous land use, soil type and climate conditions (e.g. van der Hilst et al., 2014). Also, the spatial configuration of the existing sourcing locations of the commodity is of relevance. Therefore, the spatial dimension should be taken into account. The few previous studies of Pareto frontier construction that included spatial data either used regionally aggregated data (e.g. Lautenbach et al., 2013,

Akgul et al., 2012), or analysed very small areas, making it difficult to draw generalizable conclusions relevant for regional or national policy making (e.g. Chikumbo et al., 2015, study area: 1500 ha).

It should be realized that a Pareto frontier shows the economic and environmental impacts of the *optimal* land use alternatives. In reality however, the land use system behaves suboptimal (Lambin, 2004). As an indication of the effort required to reach a point on the Pareto frontier, an assessment of economic and environmental impacts of the commodity under consideration given scenario projections with the current, suboptimal, land use system (e.g., the black dot in Figure 6.1) is decisive information (Seppelt et al., 2013).

Among the most debated commodities at the moment are biofuels. Scenario projections have been performed to study future economic benefits (e.g. Jonker et al., 2015, van der Hilst and Faaij, 2012) and potential environmental impacts (e.g. Warner et al., 2013, Fargione et al., 2010, Chapter 5, Gibbs et al., 2008) of biofuels. Such scenario projections give insights in the future sketched by the scenario, but do not indicate if this solution is optimal or suboptimal and do not provide information about potential other futures and the associated trade-offs between economic and environmental impacts. Therefore, we address a case study of simultaneously minimizing production costs and GHG emissions of a 2030 bioethanol supply for the state Goiás, in Brazil. The production of bioethanol from sugar cane has two important location variables, the locations where the sugar cane is cultivated and the locations of the mills that process the sugar cane into bioethanol (de Meyer et al., 2014)<sup>3</sup>. Consequently, the control variables in our optimization are the location of sugar cane fields and the location and scale (processing capacity) of the mills. In addition we use a land use change model to obtain a projection with current (suboptimal) trends. Our main research question is: What is the Pareto frontier between the costs and GHG emissions of bioethanol production from sugar cane in Goiás, and what are the spatial patterns of sugar cane fields and processing mills belonging to different points on this Pareto frontier? Related to the above mentioned issues, we have two sub questions: a) How do these spatial patterns at different points on the Pareto frontier differ and what drives these differences? and b) How do the optimal costs and GHG emissions compare to the costs and GHG emissions projected for Goiás for 2030 using current land use change trends?

---

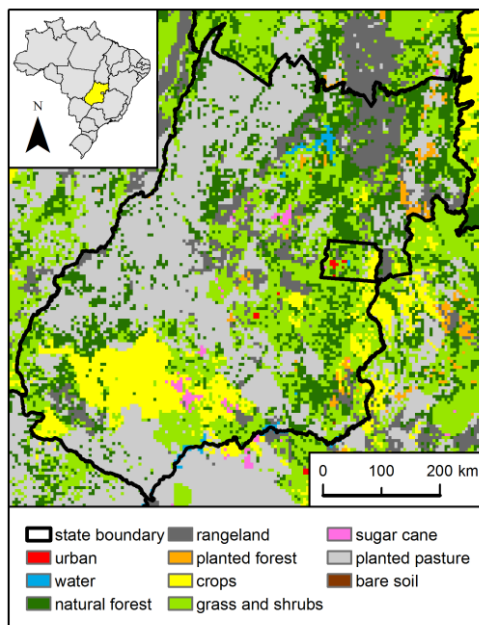
<sup>3</sup> In this study we do not consider the distribution of the ethanol to the customers.

## 6.2. Methods

### 6.2.1. Case study area

Brazil started producing bioethanol from sugar cane at the beginning of the 20<sup>th</sup> century, motivated by a gasoline import burden and a sugar production surplus (Walter et al., 2014). In the harvest season 2013/2014 ethanol production had reached 27.5 million m<sup>3</sup> (UNICA, 2015). The second largest proportion of this, 3.9 million m<sup>3</sup> (UNICA, 2015), was produced in the state Goiás, recently experiencing a fast growth of sugar cane area (Adami et al., 2012b). Goiás (Figure 6.2) is chosen as a case study area, because also a large share of the future sugar cane expansion is expected to occur here (Lapola et al., 2010, Chapter 5), making the evaluation of trade-offs between production costs and GHG emissions of bioethanol highly relevant.

Goiás is 340 000 km<sup>2</sup>, roughly the size of Germany. It is principally a large plateau with a tropical climate. It is characterized by a vast area of planted pastures, especially in the West (Figure 6.2). Furthermore, there is a sizeable patch of cropland in the Southwest, with mainly soy, corn and sugar cane (IBGE, 2013b). Most forest areas are found in the centre North part, which has a more mountainous character. We model Goiás at a 5 x 5 km<sup>2</sup> resolution (Figure 6.2).



**Figure 6.2: Land use map for Goiás, Brazil for 2006 at a 5 x 5 km<sup>2</sup> resolution (from land use map of Brazil in Chapter 5). The upper left frame shows all states in Brazil (black lines) with Goiás indicated in yellow.**

### 6.2.2. Objective function

The production of bioethanol from sugar cane can be divided into four steps:

1. Acquisition and preparation of land for the sugar cane plantation;
2. Sugar cane cultivation and harvest;
3. Transportation of the harvested cane to the mill;
4. Processing of the sugar cane into bioethanol.

All four steps involve both a production cost and a GHG emission component. Total production costs,  $c$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol), are therefore:

$$c = c_l + c_c + c_t + c_p \quad 6.1$$

In Equation 6.1,  $c_l$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) are the land costs, mainly depending on potential sugar cane yield and the costs to convert the initial land use to sugar cane,  $c_c$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) are the cultivation costs, partly yield dependent (e.g. fertilizer costs) and partly area-dependent (e.g. machinery costs),  $c_t$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) are the transport costs, depending on the distances between the fields and the processing mills, and  $c_p$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) are the processing costs, depending on the scale of the mill (processing capacity). Correspondingly, total emissions  $e$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol) are:

$$e = e_l + e_c + e_t + e_p \quad 6.2$$

In Equation 6.2,  $e_l$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol) are the land emissions, depending mainly on the carbon stocks of the replaced land use type,  $e_c$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol) are the cultivation emissions, partly yield dependent (e.g. fertilizer emissions) and partly area-dependent (e.g. machinery emissions),  $e_t$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol) are the transport emissions, depending on the distances between the fields and the processing mills, and  $e_p$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol) are the processing emissions, a fixed amount per m<sup>3</sup> ethanol produced. Note that we calculate costs and emissions at the 'factory gate', meaning that the revenues from selling the ethanol and the avoided emissions from the replacement of fossil fuel are not included. The methods to calculate the cost and emission components are based on two papers by Jonker et al. (2015, in prep.). The details of these calculations are provided in Appendix A.

One way to stimulate reduction of GHG emissions for agricultural products is to charge the producer for these emissions using a carbon price (Smith et al., 2008, Chen et al., 2012). When a carbon pricing system is established the two objectives of minimizing production costs and GHG emissions can be combined into a single objective:

$$x = c + e \cdot p \tag{6.3}$$

In Equation 6.3,  $x$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) are the aggregate costs (production costs plus GHG costs) that we aim to minimize and  $p$  (US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq) is the carbon price. The Pareto frontier between the production costs ( $c$ ) and GHG emissions ( $e$ ) of bioethanol is found by minimizing the aggregate costs ( $x$ ) for different carbon prices ( $p$ ).

### 6.2.3. Control variables and optimization model

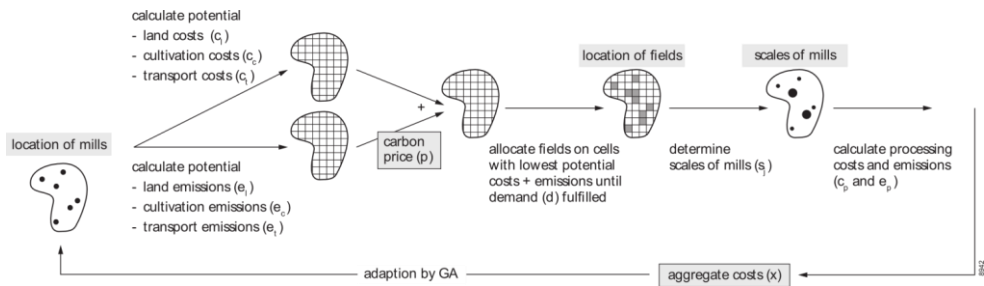
For a given bioethanol supply increase  $d$  (m<sup>3</sup> ethanol) and carbon price  $p$ , the aggregate costs  $x$  are controlled by 1) the locations of the sugar cane fields, 2) the number of processing mills and their scale (processing capacity), and 3) the locations of these mills (de Meyer et al., 2014). The relation between the aggregate costs and each of these three control variables is non-linear and the set of potential solutions is too large for exhaustive search. Metaheuristics are designed to find near-optimal solutions in an acceptable computation time for such cases, using some kind of ‘smart’ search through the solution space (Blum and Roli, 2003). We use a genetic algorithm (GA), because this metaheuristic has proved to generate good results for optimization problems similar to ours (e.g., Li and Yeh, 2005, Stewart et al., 2004). A GA mimics the process of natural selection in a population of solutions, called individuals, similar to the set of samples in a Monte Carlo approach (Bennett et al., 1998). This population evolves through a number of rounds, called generations, towards better solutions. The best performing individual of the final evolved population is the optimal solution. See Appendix 3 for a more elaborate explanation of the GA.

It is not feasible to simultaneously optimize the locations of all fields, the number of mills, and all locations and scales of these mills. Such a large number of variables for a large case study area becomes unworkable for a GA (Haupt and Haupt, 2004). Therefore, take a sequential approach. We fix the number of mills and let the variables that the GA determines be the coordinates of these mills (Figure 6.3). For each individual in each generation, the following method is applied to calculate the objective value  $x$ . Once the locations of the mills are known, sugar cane fields are



placed at the locations with the lowest  $(c_l + c_c + c_t) + p \cdot (e_l + e_c + e_t)$ <sup>4</sup> until the supply is fulfilled. The total area of sugar cane fields can differ per individual in the GA, i.e. per solution, depending on whether sugar cane is allocated on high yielding (small area) or low yielding (large area) locations. When the fields are allocated, it is known which field delivers sugar cane to which mill (lowest  $c_t + p \cdot e_t$ ) and thus the scale of each mill (tonne cane / year) can be calculated. Although the number of placed mills is fixed, this control variable is still optimized. If no fields are assigned to a mill, the production costs and GHG emissions of this mill are not counted. Therefore, in reality, only the maximum number of mills is fixed, not the actual number of mills (see also Appendix A). It is ascertained that the scale of each mill cannot exceed the maximum attainable scale  $S$  (tonne cane / year). Now all control variables are known, so the objective value  $x$  can be calculated. Next, the GA adapts the locations of the mills of a selected proportion of the individuals (see Appendix 3) and the process is repeated (Figure 6.3).

The optimization is implemented in the Python programming language (Python software foundation, 2014), using the AMORI software (AMORI, 2009) for the GA and the PCRaster Python framework (Karssenberget al., 2010) for the calculation of the phenotypes and the objective values of the individuals.



**Figure 6.3: Conceptual model of control variables and calculation of the objective value.**

#### 6.2.4. Scenario projection using current trends

In reality, land use systems behave suboptimal with respect to production costs and GHG emissions (Lambin, 2004). To realize how far the real situation is from the optimal situation, the optimal production costs and GHG emissions of bioethanol, calculated as described above, are compared to the production costs and GHG emissions that would be obtained in 2030 given a scenario projection of the current land use system. Hereto, we use an integrated economic - land use change

<sup>4</sup> Processing costs cannot be calculated yet, because they depend on the scale of the mill, which can only be determined once the fields are allocated.

model, allocating sugar cane fields (and other land uses) in Brazil with allocation rules calibrated on historic data (Chapter 5). The results of these scenario projections for different scenarios are described by van der Hilst et al. (in prep.). For this study, we use the results for Goiás (same supply increase as in the optimization, see next section) from the Business-As-Usual (BAU) scenario that includes an increased ethanol demand due to current and planned ethanol mandates worldwide.

The land use change model does not project the mills. In another study (Jonker et al., in prep.), we have optimized the number of processing mills, their scales and locations based on costs, for the projected locations of the sugar cane fields in 2030, using a mixed integer linear programming model. We use the output of this study as our projection of the mills. Together we call the locations of the sugar cane fields, number of processing mills and their scale locations of these mills generated in this way the ‘scenario projection’. Costs and GHG emissions of this projection are obtained using Equations 6.1 and 6.2. Note that part of this configuration is still optimized, although this optimization was based on production costs only, and the sugar cane field locations were not a control variable.

### **6.2.5. Data and boundary conditions**

The total increase in supply for bioethanol for Goiás for 2030  $d$  is 10.2 million  $\text{m}^3$  ethanol, derived from the total production of sugar cane in 2030 projected by the land use change model (van der Hilst et al., in prep.) minus the total production in the initial land use map (Chapter 5), and an assumed conversion efficiency  $\eta$  of  $0.09 \text{ m}^3$  ethanol / tonne cane (Jonker et al., 2015). We derive the supply from this projection to ensure that the points on the Pareto frontier and the scenario projection are equal in terms of total ethanol production. Other cost parameter values are derived from Jonker et al. (2015). The most important limitation of these values is that they were collected for the state São Paulo, while our current case study is for the state Goiás. It is possible that the values differ for Goiás in reality, for example higher or lower labour costs or different machinery usage. The same goes for cultivation and processing emission parameter values; land emission and transport emission parameter values are general, i.e. not state specific (see also Jonker et al., in prep.). All data values and sources are provided in Appendix B.

It is assumed that sugar cane present in the initial land use map goes to existing mills; this is not remodelled. New sugar cane fields cannot be allocated on raster cells that are urban, water or sugar cane in the initial land use map (as this would generate no additional ethanol compared to the initial situation). The number of mills  $M$  placed is 30 and each mill has a maximum scale of 5.5 million tonne cane / year (Jonker et al., in prep.). The GA is run with a population of 1000 individuals for

five different carbon prices of 0, 10, 100, 200 and 400 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq<sup>5</sup>. It is also run once optimizing on emissions only, to get the minimum attainable emissions (minimum attainable costs are reached at a carbon price of 0 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq). The GA settings are determined by performance tests as shown in Appendix 3.

### 6.3. Results and Discussion

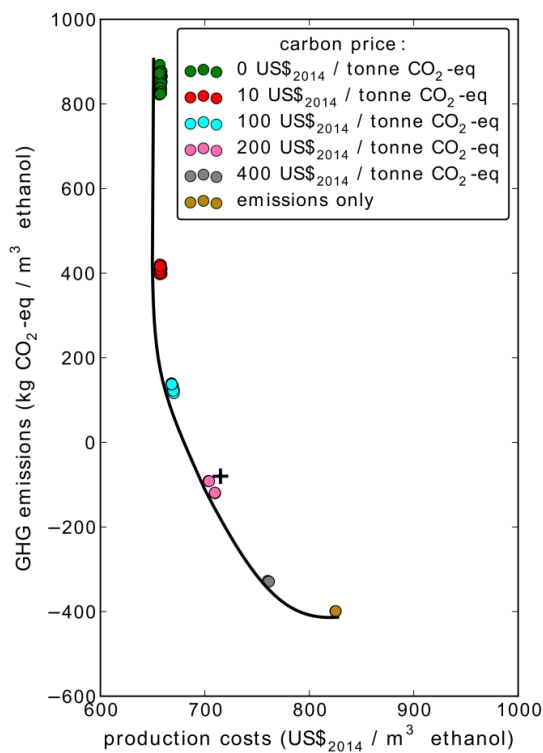
#### 6.3.1. Pareto frontier and spatial patterns belonging to different points on this frontier

The minimum attainable emissions for a production of 10.2 million m<sup>3</sup> ethanol in Goiás in 2030 are  $-399 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol, i.e. carbon sequestration of  $399 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol, and the minimum attainable production costs (excluding GHG costs) are 656 US\$<sub>2014</sub> / m<sup>3</sup> ethanol (Figure 6.4). These minimum costs are similar to the 520 US\$<sub>2010</sub> / m<sup>3</sup> (approximately 650 US\$<sub>2014</sub> / m<sup>3</sup> ethanol) calculated by Jonker et al. (2015). The Pareto frontier shows the trade-offs between production costs and GHG emissions between these two extremes. Interesting from a policy perspective is that, at a carbon price of only 10 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq, which is roughly the current carbon price in Europe (EEX, 2015), GHG emissions are reduced by more than 50%, compared to a zero carbon price in our simulations, with a negligible increase in production costs (Figure 6.4). Also, emission savings required by the Renewable Energy Directive (RED) (European Parliament and Council of the European Union, 2009) are reached at this carbon price<sup>6</sup> (Figure 6.5b).

---

<sup>5</sup> No carbon pricing system is currently installed in Brazil and it is unknown if and when it will be installed (Dahan et al., 2015).

<sup>6</sup> Note that RED also requires inclusion of GHG emissions from the distribution of ethanol (transport from the mills to the customers). This is not included in our study but we expect its contribution to be small because the emissions of sugar cane transport are already small (Figure 6.5b) and the ethanol has a much higher energy density than sugar cane, and therefore lower GHG emissions per m<sup>3</sup> ethanol (e.g. Hamelinck et al., 2005a).



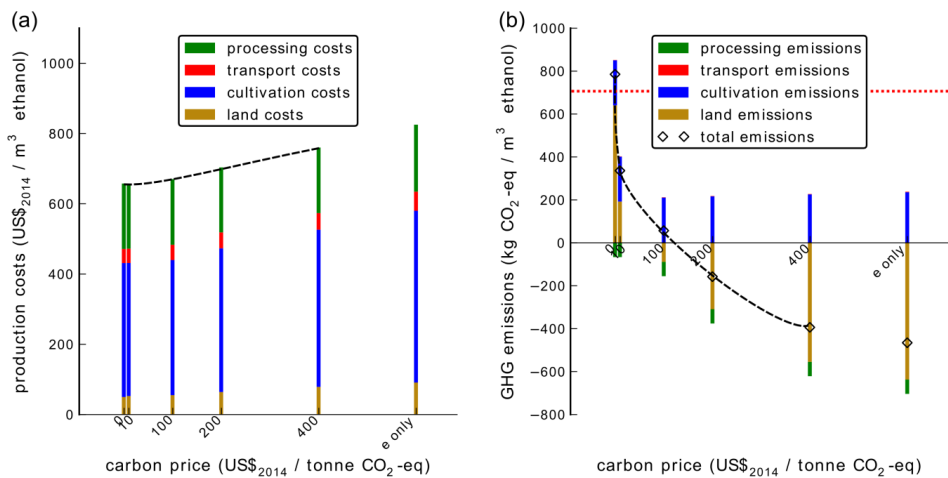
**Figure 6.4: Pareto frontier (black line) between production costs (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) and GHG emissions (kg CO<sub>2</sub>-eq / m<sup>3</sup> ethanol) for a supply increase of 10.2 million m<sup>3</sup> ethanol in Goiás in 2030. For each carbon price (0, 10, 100, 200 and 400 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq), and the optimization on emissions only, the best 10% of the final GA population is plotted, where each individual is indicated by a coloured circle. The black plus sign indicates the production costs and GHG emissions for the scenario projection using current trends.**

As expected, GHG emissions decrease with increasing carbon price, while production costs increase with increasing carbon price (Figure 6.4), mainly through an increase in cultivation costs and also a slight increase in land costs (Figure 6.5a). Cultivation costs and land costs depend on the yield of the sugar cane field (section 6.2.2 and Appendix A). Consequently, at low carbon prices, when the aggregate costs are mainly driven by production costs, all sugar cane fields and mills are concentrated in high-yielding areas (Figure 6.6, e.g. location 1 and 3), and the scale of the mill depends on the size of the high-yield area. The GHG emissions of ethanol production are dominated by the emissions of land use change. Therefore, reductions in GHG emissions along the Pareto frontier are reached through a reduction in land emissions (Figure 6.5b). Land emissions mainly depend on the type of land use replaced by sugar cane and carbon is sequestered when this is cropland. Hence, at high carbon prices, sugar cane fields are placed at locations that used to be cropland (Figure 6.6, e.g. location 2). However, not only croplands

are replaced. The yield remains an important driver of the sugar cane field pattern at high carbon prices (Figure 6.6, e.g. location 1), because land emissions and cultivation emissions are yield dependent too (section 6.2.2 and Appendix A). Therefore, low-yielding cropland is never preferred over high-yielding land of another land use type.

The spatial patterns of sugar cane fields and mills belonging to different carbon prices provide interesting information for management strategies. Win-win configurations that are optimal at all points along the Pareto frontier (Figure 6.6, e.g. location 1), are good locations for investment, because they are robust, i.e. independent of the future carbon pricing policy.

From this analysis, two general conclusions can be derived. Firstly, yield is the main determinant of ethanol production costs and also an important driver of ethanol GHG emissions. The Pareto frontier would change drastically if the yield would change, either spatially (pattern of Figure 6.6e changes), or in the absolute sense (absolute yield in tonne cane / ha belonging to the value of 1 in Figure 6.6e changes). Therefore, improvements in management practices and the introduction of higher yielding sugar cane varieties can positively impact costs as well as emissions. This conclusion also gives rise to the need to examine the effect of climate change, because through a change in yields this can have major impacts on costs and emissions in a positive as well as in a negative direction (e.g. Holzkämper et al., 2015).



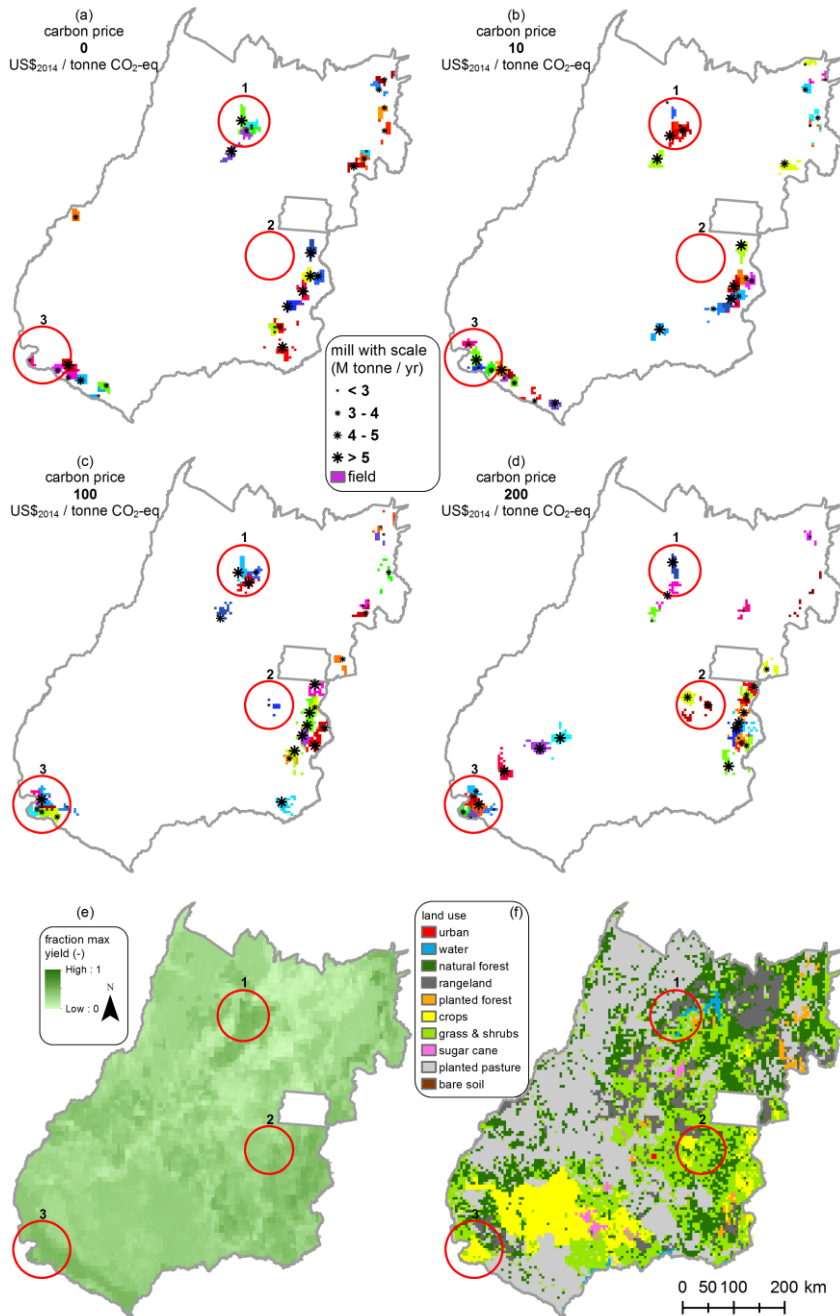
**Figure 6.5: (a) Production cost breakdown and (b) GHG emission breakdown for all carbon prices. ‘E only’ are the results of the optimization on emissions only. The black dashed lines are the interpolation between carbon price and production costs or GHG emissions and the red dotted line indicates the maximum GHG emissions allowed for biofuels produced in installations in which production started on or after 1 January 2017 according to the Renewable Energy Directive (RED) (European Parliament and Council of the European Union, 2009).**

Secondly, the shape of the Pareto frontier is subjected to the magnitude of the ethanol supply increase, as indicated for costs by others through cost-supply curves (e.g. van der Hilst and Faaij, 2012). For the supply analysed in this study, it was possible to allocate all sugar cane fields on high-yielding locations, but with an increase in supply lower yielding locations will have to be used and consequently average costs per m<sup>3</sup> ethanol will increase. Along the same line of reasoning GHG emissions will increase with fewer possibilities to pick field locations with low carbon stocks. A limitation in our methodology is that it does not take into account indirect land use land change (iLUC). ILUC is the cascading effect of a land use change: for instance, when a sugar cane is allocated on land previously used for crop cultivation, the crop production has to be moved to elsewhere, in order to sustain the demand for this crop (Wicke et al., 2012). ILUC is likely to cause GHG emissions, additional to the direct emissions. Since they are (indirectly) caused by the land allocation of sugar cane, it could be argued that these emissions should be included in the commodity's Pareto frontier.

Incorporating iLUC requires determination of where the displaced land use type reappears. This entails the use of either a spatial land use change model determining the dynamics of all displaceable land use types (e.g. Chapter 5) or a combined optimization model for ethanol production and all other commodities (e.g. Lautenbach et al., 2013). Both options are data intensive and therefore require massive run times. And modelling Goiás only is not sufficient, because the reallocation can take place anywhere in the world. Yet, it is obvious that, if iLUC could be included in the optimization, 1) it is very likely that higher GHG emissions are obtained and 2) placing sugar cane on croplands becomes much less favourable, because the crop production has to be moved to elsewhere and eventually some natural vegetation will probably be converted<sup>7</sup>. This means that the position and shape of the Pareto frontier changes as well as the spatial patterns of sugar cane fields and mills belonging to the points on the frontier.

---

<sup>7</sup> The area of natural vegetation eventually converted is often not equal to the area of sugar cane allocated, for example because the location where the agricultural land has been moved, has a higher or lower productivity or different management practices. Or because of price effects, which can only be determined if an economic equilibrium model is used. These things are additional complications when trying to include iLUC in the optimization.



**Figure 6.6: Spatial patterns of sugar cane fields and mills belonging to the optimal individuals for different carbon prices of (a) 0; (b) 10; (c) 100; and (d) 200 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq. Fields with the same colour belong to the same mill. At the bottom the input maps (e) yield fraction (Tóth et al., 2012) and (f) initial land use (2006) (Chapter 5) are displayed. The red circles with numbers are locations referred to in the main text.**

### 6.3.2. *Costs and emissions optimized versus scenario projection*

The scenario projection using current trends generates production costs of 715 US\$<sub>2014</sub> / m<sup>3</sup> ethanol and GHG emissions of  $-80 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol (Figure 6.4). As expected, current trends result in somewhat higher production costs than the minimum attainable, 59 US\$<sub>2014</sub> / m<sup>3</sup> ethanol higher. Yet, the scenario projection is surprisingly close to the Pareto frontier, near our optimal results for a carbon price of 200 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq (Figure 6.4), thus with GHG emissions lower than we expected. We supposed that, with no carbon pricing system installed, ethanol producers would not be too concerned about GHG emissions. Still, GHG emissions are  $865 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol lower than what our optimization comes to at zero carbon price, mainly because about 80% of the sugar cane field expansion is projected to occur on croplands (Jonker et al., in prep.).

We believe that two factors contribute to these relatively low GHG emissions of the scenario projection. Firstly, not only does the conversion of croplands result in low GHG emissions, croplands have two other advantages why they are often selected to be converted to sugar cane by the land use change model, used to create the scenario projection (Chapter 5): 1) they are easier, i.e. cheaper, to convert compared to non-agricultural land, and 2) sugar cane is likely to be allocated in the neighbourhood of existing sugar cane fields and mills, because this creates economies of scale, and in Goiás a large share of the existing sugar cane fields happens to be bordered by croplands (Figure 6.2).

The second reason is that, although no carbon pricing system is currently installed in Brazil, current trends can take into account GHG emissions through other sustainability regulations. The RED (Figure 6.5b) is not likely to have an effect, because it is a European regulation and almost all Brazilian ethanol is currently used domestically. But there are some Brazilian sustainability regulations. Examples are the federal ecological zoning for sugar cane and certification of sustainably produced ethanol (Lucon and Goldemberg, 2010). Although such regulations were not explicitly included in the land use change model, they might have been implicitly captured in the model structure through calibration on historic data.

Note that our projection assumes that current trends are continued towards 2030, while in reality a land use system can change quite abruptly (Chapter 4). Therefore, the low GHG emissions of the scenario projection are uncertain and by no means undermine the potential benefits of a carbon pricing system in Brazil.



## 6.4. Conclusion

This study has spatially assessed trade-offs between production costs (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) and GHG emissions (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol) of a 2030 bioethanol supply increase of 10.2 million m<sup>3</sup> ethanol in the state Goiás, Brazil, by optimizing the location of sugar cane fields and the location and size of processing mills. The Pareto frontier between the production costs and GHG emissions of bioethanol has been obtained by carrying out this optimization for different carbon prices (US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq). The minimum attainable GHG emissions are  $-399 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol, i.e. carbon sequestration, and the minimum attainable production costs (excluding GHG costs) are 656 US\$<sub>2014</sub> / m<sup>3</sup> ethanol. The Pareto frontier ranges between those two extremes and shows a steep decrease in GHG emissions while production costs practically remain at their minimum at a relatively small increase in carbon price from 0 to 10 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq. In-between these prices the emission savings required by the Renewable Energy Directive (RED) (European Parliament and Council of the European Union, 2009) are reached, which is promising for potential sustainability policies.

The main spatial determinant of production costs is the potential yield and the main determinants of the GHG emissions are the replaced land use type, where current cropland has a high sequestration potential, and also potential yield. This information can be used to restrict zoning for sugar cane fields to current cropland fields while monitoring that, for the given supply, the area of high potential yield in these regions is large enough to reach a competitive production price per m<sup>3</sup> ethanol.

Production costs and GHG emissions projected for Goiás for 2030 using current land use change trends are 715 US\$<sub>2014</sub> / m<sup>3</sup> ethanol and  $-80 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol respectively. We expected suboptimal results with high emissions because no carbon pricing system is currently installed in Brazil, but these values are relatively close to the Pareto frontier, near our optimal results for a carbon price of 200 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq.

Our study illustrates how much more information can be derived from the construction of a Pareto frontier between economic and environmental objectives, compared to a scenario projection, the currently most frequently used approach for the assessment of potential futures. A Pareto frontier shows the trade-offs between the two objectives, interesting for producers, and, in our case study, indicates at which carbon price a particular GHG emission reduction is reached, information of interest for policy making. The spatial configurations of the optimized commodity give, in contrast to scenario projections, information about the robustness or uncertainty of promising locations under different (in our case carbon pricing) policies, thereby facilitating management decisions. The developed

methodology is general and can easily be applied to other regions, scales, objectives and commodities.

## 6.5. Acknowledgements

This work was carried out within the BE-Basic R&D Program, which was granted a FES subsidy from the Dutch Ministry of Economic affairs, agriculture and innovation (EL&I). We thank the MSc students Luis Manuel de Vries and Sanne Hettinga for the implementation of the genetic algorithm and the explorative study on optimizing the locations of sugar cane mills respectively.

## 6.6. Appendix A: Equations of the cost and emission components

### 6.6.1. General

This section provides all equations used for the cost and emission components of the allocation model. In the equations the following notations are used throughout. *Italic variables* are non-spatial and **bold variables** are spatial, i.e. a map of values at raster cells instead of a single value. The subscript  $i$  for a spatial variable indicates the selection of a cell in the map on which sugar cane is allocated, for  $i = 1, 2, \dots, I$ . The value for  $I$ , the total number of cells with sugar cane, varies per model run (per individual in the genetic algorithm (GA) population), depending on whether sugar cane is allocated on high yielding (low  $I$ ) or low yielding (high  $I$ ) cells.

Many of the equations contain the yield of sugar cane. Sugar cane is a semi-perennial crop. This means that after planting, it can be harvested for some consecutive years. In Brazil, a 6-year cycle is most common: the first harvest takes place 12 or 18 months after planting and the subsequent harvests once every year for four years (Macedo et al., 2008). During these four years, the yield gradually decreases. In the next year, the field is renewed. The yield we use, is the average yield over this 6-year cycle. The exception are the cultivation costs, where we do specifically account for the changing yield over the cycle, because here it influences the total costs. The yield of sugar cane in a cell,  $\mathbf{y}_i$  (tonne cane), is constructed from a map showing the relative yield distribution over space and an average (over the cycle) maximum attainable yield:

$$\mathbf{y}_i = \mathbf{f}_i \cdot m \cdot a \quad \text{for each } i \quad 6.4$$

In Equation 6.4,  $m$  (tonne cane / ha) is the maximum attainable yield, in our case for 2030. Furthermore,  $\mathbf{f}_i \in [0, 1]$  (-) is the fraction of the maximum yield that can

be obtained. And  $a$  (ha) is the cell area, which is constant over space since we use an Albers Equal Area map projection.

Many of the cost and emission components are expressed per tonne cane. Because we are interested in the costs of the end product, ethanol, all components are converted to costs and emissions per m<sup>3</sup> ethanol, using the conversion efficiency  $\eta$  (m<sup>3</sup> ethanol / tonne cane). We assume a single conversion efficiency for all mills allocated.

### 6.6.2. Costs

All cost equations are derived from the equations by Jonker et al. (2015). For more detailed information we refer to that study. The land costs,  $c_l$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol), consist of two parts: the costs to buy the land and the land conversion costs to make the land cultivatable for sugar cane. Both parts contain area dependent costs and yield dependent costs:

$$c_l = \sum_{i=1}^I \frac{\alpha \cdot a \cdot (a_l \cdot (\mathbf{b}_i + \mathbf{l}_i) + b_l \cdot \mathbf{y}_i \cdot (\mathbf{b}_i + \mathbf{l}_i))}{\mathbf{y}_i \cdot \eta} \quad 6.5$$

In Equation 6.5,  $\mathbf{b}_i$  (US\$<sub>2014</sub> / ha) are the costs to buy the land and  $\mathbf{l}_i$  (US\$<sub>2014</sub> / ha) are the land conversion costs. Both vary over space based on region and/or current land use type. The factors  $a_l$  (-) and  $b_l$  (ha / tonne) distinguish between the area dependent costs and yield dependent parts. The factor  $\alpha$  (-) is the annuity factor that transforms the total costs to yearly costs:  $\alpha = r / (1 - (1 + r)^{-L})$ . Herein,  $r$  (-) is the discount rate, i.e. the time value of money according to the theory of time preference, and  $L$  (years) is the lifetime or amortization period (Blok, 2006, p. 195, Equation 11.2b). Again,  $a$  is the cell area.

The cultivation costs,  $c_c$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol), contain many factors, like fertilizer application and machinery use. A simplified equation, summing all cost components in area-dependent and yield dependent cultivation costs, is derived from the equation given by Jonker et al. (2015). The initial investment costs are annualized, but this cannot simply be done using  $\alpha$  as in Equation 6.5, because the yearly costs vary over the 6-year sugar cane cultivation cycle:

$$c_c = \frac{1}{\eta} \sum_{i=1}^I \left( \frac{\sum_{t=1}^6 (a_{c,t} \cdot a) / \sum_{t=1}^6 (1 + r^t)}{\sum_{t=1}^6 (\mathbf{y}_i \cdot \mathbf{y}_t) / \sum_{t=1}^6 (1 + r^t)} + \frac{\sum_{t=1}^6 (b_{c,t}) / \sum_{t=1}^6 (1 + r^t)}{\sum_{t=1}^6 \mathbf{y}_t / \sum_{t=1}^6 (1 + r^t)} \right) \quad 6.6$$

In Equation 6.6,  $\mathbf{y}_t$  (-) is the yield factor for year  $t = 1, 2, 3, 4, 5, 6$ , decreasing over the 6-year sugar cane cultivation cycle,  $a_{c,t}$  (US\$<sub>2014</sub> / ha) are the area-dependent cultivation costs at year  $t$ , for example the costs of machinery,  $b_{c,t}$  (US\$<sub>2014</sub> / tonne

cane) are the yield dependent cultivation costs at year  $t$ , for example for fertilizers. Again,  $r$  (-) is the discount rate.

The costs of transporting the sugar cane to the mill,  $c_t$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol), are calculated over the road network. It is assumed that the truck will take the fastest route. The fastest route is determined by applying a least cost path algorithm on the road map with speed differing per road type and off-road, where the 'costs' per kilometre are one divided by the speed (higher speed are lower 'costs' because it is faster). The 'costs' themselves are not used, only the route, to compute the average speed and diesel use along the routes and total distance of the routes per sugar cane cell:

$$c_t = \sum_{i=1}^I \frac{(a_t/v_i + b_i) \cdot d_i + o_t}{\eta} \quad 6.7$$

In Equation 6.7,  $v_i$  (km / hour) is the average truck speed, and  $d_i$  (km) is the total distance of the fastest route to the nearest (time-wise) mill. Furthermore in Equation 6.7,  $a_t$  (US\$<sub>2014</sub> / tonne-hour) are the annual costs of the truck,  $b_i$  (US\$<sub>2014</sub> / tonne-km) are the diesel costs per tonne cane that differ per field depending on the road types in the route to the mill, and  $o_t$  (US\$<sub>2014</sub> / tonne cane) are the costs of loading and unloading the truck. When, during the allocation process of the fields, one of the mills has reached the maximum capacity, the spatial variables are updated for all cells that had this 'full' mill as their closest mill, according to the fastest route to the next nearest mill.

The costs of processing the sugar cane (converting it to ethanol),  $c_p$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol), include capital depreciation, operational costs and revenues from electricity generation. In contrast to the other cost components, the processing costs are calculated per mill instead of per field, because the costs depend on the scale of the mill. Therefore, the total processing costs are obtained by summing over all *active* mills,  $j = 1, 2, \dots, J$ . The notion 'active' indicates all mills to which fields are assigned. Mills to which no fields are assigned, do have a location in theory (they have a x and y coordinate in the GA), but do not contribute to the total costs of the individual and are thus excluded from the analysis:

$$c_p = \sum_{j=1}^J \left( \frac{(\alpha \cdot (a_p \cdot s_j) + b_p)}{s_j \cdot q \cdot \eta} + c_o - g_e \cdot r_e \right) \quad 6.8$$

In Equation 6.8,  $s_j$  (tonne cane / hour) is the scale of the mill, indicating the sugar cane processing capacity, and  $s_j \cdot q \cdot \eta$  (m<sup>3</sup> ethanol) is the annual output of the mill, in which  $q$  (hours) is the number of hours per year the mill is in running. In this study,  $q$  is assumed to be the same for all mills. Furthermore,  $a_p$  (US\$<sub>2014</sub>-hour / tonne) is a cost factor that decreases with the scale of the mill, representing the

advantages of economies of scale, while  $b_p$  (US\$<sub>2014</sub>) is a fixed cost factor. Moreover,  $c_o$  (US\$<sub>2014</sub> / m<sup>3</sup> ethanol) are the fixed operation costs, and  $g_e$  (kWh / m<sup>3</sup> ethanol) is the electricity surplus. This electricity is generated from bagasse, a fibrous product left over after the sugary juice is extracted from the sugar cane. The electricity surplus, the part of the generated electricity the mill does not need for the sugar cane processing, can be sold to the grid;  $r_e$  (US\$<sub>2014</sub> / kWh) are the revenues obtained for this surplus. Again, as with land costs,  $\alpha$  (-) is the annuity factor that transforms the total costs to yearly costs.

### 6.6.3. Emissions

The land emissions,  $e_l$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol), from carbon stock changes are calculated using the IPCC approach (IPCC, 2006). This approach involves five carbon pools: above ground biomass, below ground biomass, dead wood, litter, and soil organic carbon (SOC). In line with the Tier 1 approach of the IPCC, an equilibrium is assumed in the dead wood and litter stocks, i.e. they are considered not to change:

$$e_l = \frac{44}{12} \sum_{i=1}^I \frac{(s_{i,2006} - s_{i,c} + b_{i,2006} - b_{i,c}) \cdot a}{h \cdot y_i \cdot \eta} \quad 6.9$$

In Equation 6.9,  $s_{i,2006}$  (tonne C / ha) is the mineral soil organic carbon for 2006, the reference year. The mineral SOC is calculated given the soil type, climate, land use and management (IPCC, 2006). Next,  $s_{i,c}$  (tonne C / ha) is the mineral soil organic carbon when all cells  $i = 1, 2, \dots, I$  are converted to sugar cane. This means that land use and management are changed, while soil type and climate remain the same. Organic SOC is not considered in Equation 6.9, because our study area does not contain organic soils. Furthermore,  $b_{i,2006}$  (tonne C / ha) is the above and below ground biomass for 2006, based on land use and productivity. In the same fashion as with  $s_{i,c}$ ,  $b_{i,c}$  (tonne C / ha) is the above and below ground biomass in all cells  $i = 1, 2, \dots, I$  where sugar cane is cultivated. The total carbon stock changes are divided over a time horizon  $h$  (years). The factor 44/12 is to convert from C to CO<sub>2</sub>-eq.

The cultivation emissions,  $e_c$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol), originate from the use of diesel (for machinery), fertilizers, agrochemicals and other chemicals:

$$e_c = \sum_{i=1}^I \frac{l_c \cdot a + y_i \cdot k_c}{y_i \cdot \eta} \quad 6.10$$

In Equation 6.10,  $k_c$  (tonne CO<sub>2</sub>-eq / ha) are the diesel emissions from the machinery, and  $l_c$  (tonne CO<sub>2</sub>-eq / tonne cane) are the yield-dependent emissions,

including primarily fertilizer emissions. These are mainly N<sub>2</sub>O emissions, converted to CO<sub>2</sub>-eq.

The transport emissions,  $e_t$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol), are diesel emissions from the trucks transporting the sugar cane to the mill:

$$e_t = \sum_{i=1}^I \frac{\mathbf{k}_i \cdot \mathbf{d}_i}{l_t \cdot \eta} \quad 6.11$$

In Equation 6.11,  $\mathbf{k}_i$  (tonne CO<sub>2</sub>-eq / km) are the diesel emissions per tonne cane that differ per field depending on the road types in the route to the mill, including a factor to correct for the empty return of the truck, and  $l_t$  (tonne) is the load of a full sugar cane truck.

The processing emissions,  $e_p$  (tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol), are assumed not to differ per mill, in contrast to the processing costs; it is a fixed emission per tonne cane:

$$e_p = \frac{k_p - l_p}{\eta} \quad 6.12$$

In Equation 6.12,  $k_p$  (tonne CO<sub>2</sub>-eq / tonne cane) includes all processing emissions and  $l_p$  (tonne CO<sub>2</sub>-eq / tonne cane) are the emissions avoided by electricity production.

## 6.7. Appendix B: Input data

This section describes the data used for the Goiás case study within the equations given in the previous section. The values of non-spatial variables are given in Table 6.1 for costs and in Table 6.2 for emissions. The data sources for all maps are given in Table 6.3.

One of the most important variables in the cost calculations is the maximum attainable yield,  $m$ . The value of  $m$  is determined for 2012 by finding the  $m$  for which  $\sum_{i=1}^I \mathbf{y}_i = q$ , where  $i = 1, 2, \dots, I$  are in this case all cells that are projected to be sugar cane for 2012 by a combination of data from the Canasat project (Rudorff et al., 2010) and a model projection from the PCRaster Land Use Change model PLUC (van der Hilst et al., in prep.). Furthermore,  $q$  is the total sugar cane production reported by the Brazilian Sugarcane Industry Association UNICA (UNICA, 2015). The value of  $m$  for 2030 (Table 6.1) is found by applying a yield trend over time from Jonker et al. (2015) to the 2012 value.

Regarding land emissions we assume, in line with the IPCC method (IPCC, 2006), that the above and below ground biomass of cropland is zero, because the crops

are fully harvested each year. For planted pasture, it is assumed that all above ground biomass is eaten by the livestock each year, so the biomass stock of pasture is only its below ground biomass. Rangelands have a stocking rate of about 70% lower than pastures (Aguar and d'Athayde, 2014), so we assume that only 30% of the above ground biomass is eaten each year, 70% remains in stock. Along similar lines, we assume that the above ground sugar cane is harvested each year and that the roots remain intact. The biomass stock of planted forest is also its below ground biomass only, as all carbon in the above ground stock is eventually harvested.

**Table 6.1 (continues on next page): Non-spatial data for the sugar cane production costs. All values are for the year 2030 and expressed in US\$<sub>2014</sub>; if values in the source were in another monetary unit, they are converted using the IGP-DI index (Banco Central do Brasil, 2015).**

cost component	variable	unit	symbol	value	source		
general	maximum yield	tonne/ha	$m$	212	Rudorff et al., 2010, see explanation in main text, van der Hilst et al., in prep., UNICA, 2015		
	cell area	ha	$a$	2500	-		
	conversion efficiency	m <sup>3</sup> ethanol / tonne cane	$\eta$	0.09	Jonker et al., 2015		
land	annuity factor	-	$\alpha$	0.13*	Jonker et al., 2015		
	correction factor	-	$a_l$	0.33	FNP Informa economics, 2012		
	correction factor	ha / tonne	$b_l$	$6.67 \cdot 10^{-3}$	FNP Informa economics, 2012		
cultivation	yield factor in year $t$	-	$y_t$	$t$	$y_t$	Macedo et al., 2004	
				1	0		
				2	1.29		
				3	1.09		
				4	0.95		
				5	0.87		
				6	0.81		
	area-dependent cultivation costs in	US\$ <sub>2014</sub> / ha		$a_{c,t}$	$t$	$y_t$	Jonker et al., 2015
					1	2577	
					2	1124	
3					1124		

	year $t$			4	1124	
				5	1124	
				6	1050	
	yield dependent cultivation costs in year $t$	US\$ <sub>2014</sub> / tonne cane	$b_{c,t}$	$t$	$y_t$	Jonker et al., 2015
				1	8.3	
				2	15.0	
				3	15.9	
				4	16.3	
				5	16.9	
				6	17.2	
	discount rate	-	$r$	0.12		Jonker et al., 2015
<b>transport</b>	capital depreciation of the truck	US\$ <sub>2014</sub> / tonne-hour	$a_t$	2.68		Jonker et al., 2015
	truck loading and unloading	US\$ <sub>2014</sub> / tonne cane	$o_t$	2.00		Jonker et al., 2015
<b>processing</b>	annuity factor	-	$\alpha$	0.13 <sup>*</sup>		Jonker et al., 2015
	period per year the mills runs	hours	$q$	170 · 24		Dias et al., 2011
	cost factor decreasing with scale	US\$ <sub>2014</sub> -hour / tonne	$a_p$	75.57 <sup>**</sup> 59.09 <sup>***</sup>		Jonker et al., 2015
	fixed cost factor	US\$ <sub>2014</sub>	$b_p$	40 · 10 <sup>6**</sup> 100 · 10 <sup>6***</sup>		Jonker et al., 2015
	operation costs of the mill	US\$ <sub>2014</sub> / m <sup>3</sup> ethanol	$c_o$	98.67		Jonker et al., 2015
	electricity surplus	kWh / m <sup>3</sup> ethanol	$g_e$	906.67		Jonker et al., 2015
	revenues from electricity	US\$ <sub>2014</sub> / kWh	$r_e$	0.07		Jonker et al., 2015

\* Calculated for an amortization period  $L$  of 20 years with a 12% interest rate  $r$  (Blok, 2006, p. 195, Equation 11.2b)

\*\* For mills with a scale smaller than 1000 tonne / hour

\*\*\* For mills with a scale equal to or larger than 1000 tonne / hour



**Table 6.2: Non-spatial data for the sugar cane production emissions. All values are for the year 2030.**

emission component	variable	unit	symbol	value	source
land	time horizon	years	$h$	20	European Parliament and Council of the European Union, 2009, IPCC, 2006
cultivation	yield-dependent emissions	tonne CO <sub>2</sub> -eq / tonne cane	$k_c$	$15.22 \cdot 10^{-3}$	quantities: Jonker et al., 2015, emission per component: Macedo et al., 2008
	area-dependent emissions	tonne CO <sub>2</sub> -eq / ha	$l_c$	$365.91 \cdot 10^{-3}$	quantities: Jonker et al., 2015, emission per component: Macedo et al., 2008, Seabra et al., 2011
transport	truck load	tonne	$l_t$	30	Jonker et al., 2015, CTBE, 2012
processing	processing emissions	tonne CO <sub>2</sub> -eq / tonne cane	$k_p$	$4.45 \cdot 10^{-3}$	Jonker et al., 2015, Seabra et al., 2010
	emissions avoided by electricity production	tonne CO <sub>2</sub> -eq / tonne cane	$l_p$	$10.44 \cdot 10^{-3}$	energy mix of Brazil (excluding bagasse) and related emissions: IEA, 2013, surplus quantity: Jonker et al., 2015

**Table 6.3 (continues on next page): Sources of all spatial data for sugar cane production costs and emissions.**

component	variable	unit	symbol	source
general	yield fraction	tonne / ha	$f_i$	Tóth et al., 2012
land costs	land tenure costs	US\$ <sub>2014</sub> / ha	$b_i$	FNP Informa economics, 2012*
	land conversion costs	US\$ <sub>2014</sub> / ha	$l_i$	FNP Informa economics, 2012**
transport costs	speed	km / hour	$v_i$	speeds on different road types adapted from (to correct for trucks going slower): de Souza Soler and Verburg, 2010
	diesel costs	US\$ <sub>2014</sub> / tonne-km	$b_i$	adapted from for different speeds: Jonker et al., 2015
	distance	km	$d_i$	calculated over road network from: UFG, 2015
conversion costs	scale	tonne cane / hour	$s_j$	determined by the model, max scale for 2030 set at 1348 tonne cane / hour 5.5 Mtonne / year) (MME, 2013)
land emissions	mineral soil organic carbon in 2006	tonne C / ha	$s_{i,2006}$	2006 land use map: Chapter 5, values for carbon dependent on land use type, soil, climate and management level: IPCC, 2006, soil map: Batjes, 2010, climate map: Hijmans et al., 2005, Bernoux et al., 2006
	mineral soil organic carbon when sugar cane is cultivated	tonne C / ha	$s_{i,c}$	values for carbon dependent on land use type, soil, climate and management level: IPCC, 2006, soil map: Batjes, 2010, climate map: Hijmans et al., 2005
	total (above + below ground) biomass stock for 2006	tonne C / ha	$b_{i,2006}$	above ground biomass: maximum yield assumptions by the authors together with yield fraction map by Tóth et al., 2012, 2006 land use map: Chapter 5, ratio below to above ground biomass: IPCC, 2006, Jangpromma et al., 2012, de Miranda et al., 2014, Epron et al., 2013

	total biomass stock when sugar cane is cultivated	tonne C / ha	$\mathbf{b}_{i,c}$	above ground biomass: maximum yield assumptions by the authors together with yield fraction map by Tóth et al., 2012, 2006 land use map: Chapter 5, ratio below to above ground biomass: IPCC, 2006, Jangpromma et al., 2012
<b>transport emissions</b>	diesel emissions	tonne CO <sub>2</sub> -eq / km	$\mathbf{k}_i$	quantities: Jonker et al., 2015, Macedo et al., 2008, Hamelinck et al., 2005b

\* The FNP (FNP Informa economics, 2012) specifies land value per micro region in Goiás per land use type (distinction between natural vegetation, pasture and cropland). A land value map was made using the map of micro regions in Brazil and the land use map of 2006 (Chapter 5). In addition the FNP indicates a higher land value around the cities of Rio Verde and Santa Helena de Goiás. This higher land value was assigned to all grid cells within a buffer of 50 km around these cities.

\*\* The FNP (FNP Informa economics, 2012) specifies land conversion costs separately for nature and agriculture. A conversion cost map was made by linking these values to the land use map of 2006 (Chapter 5).

## 6.8. Appendix 3: Genetic algorithm

A GA searches the solution space by mimicking evolutionary processes. It starts with a population of  $N$  candidate solutions, also called individuals (Figure 6.7). Each individual has a genotype, consisting of a bit-string of genes representing the control variables of the problem, and a phenotype, the ‘appearance’ resulting from the genotype (Bennett et al., 1998). The fitness of each individual in the population is calculated by evaluating this phenotype against the objective(s). The best-performing individuals (a predetermined fraction of the population) are selected to ‘reproduce’. This is done by crossover, also called recombination, and mutation (Blum and Roli, 2003). Crossover is the process of taking genes from two parents and combining them into a new genotype. Mutation alters a bit in one or more randomly selected genes. The new generation, the parents and the children together, generally has a higher fitness than the previous generation. The GA is configured to terminate when the optimum has been found.

The settings of the parameters of our GA are tuned by systematic variation and monitoring the effect on the variance in the population and on the objective value of the best individual of the final population. First the fraction of the population to reproduce and the mutation rate are optimized for a population of 100 individuals.

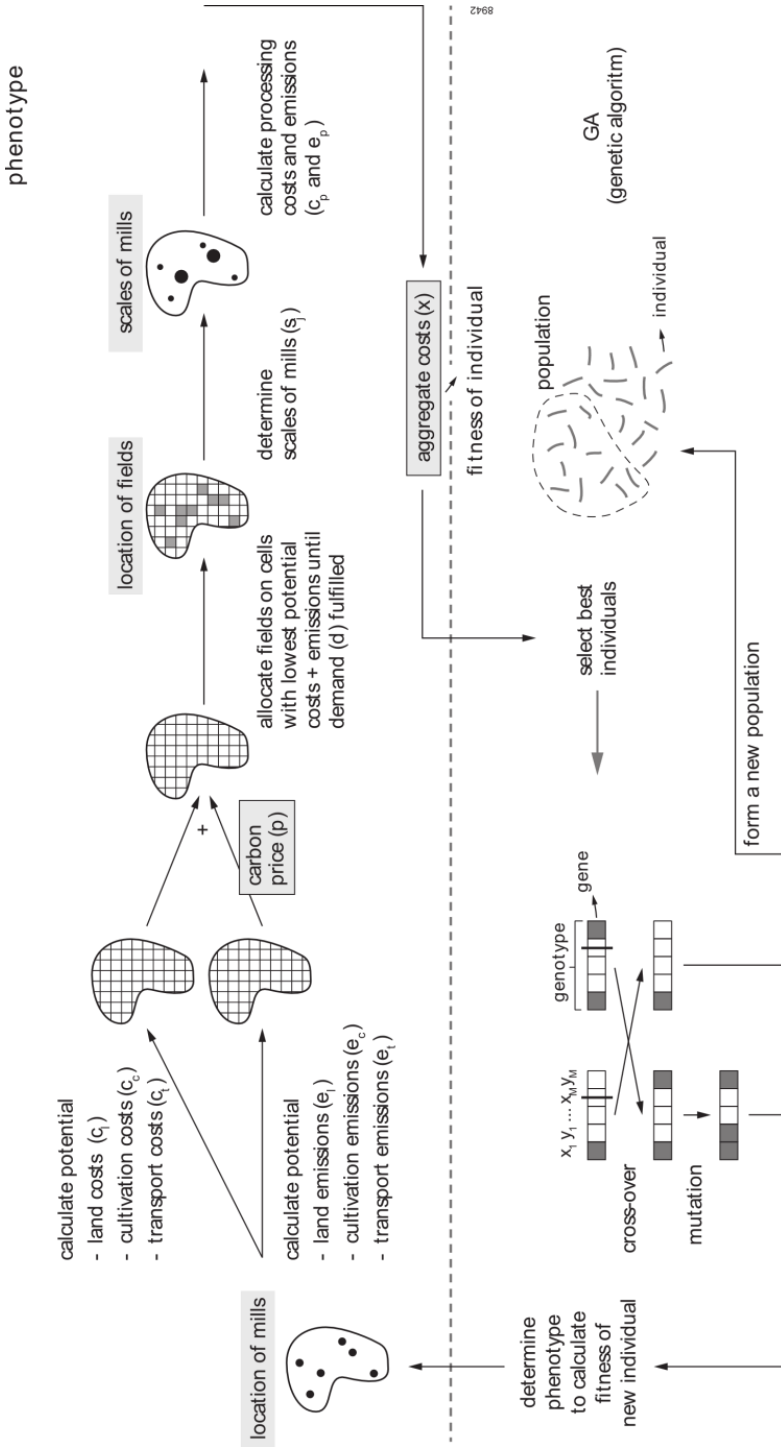
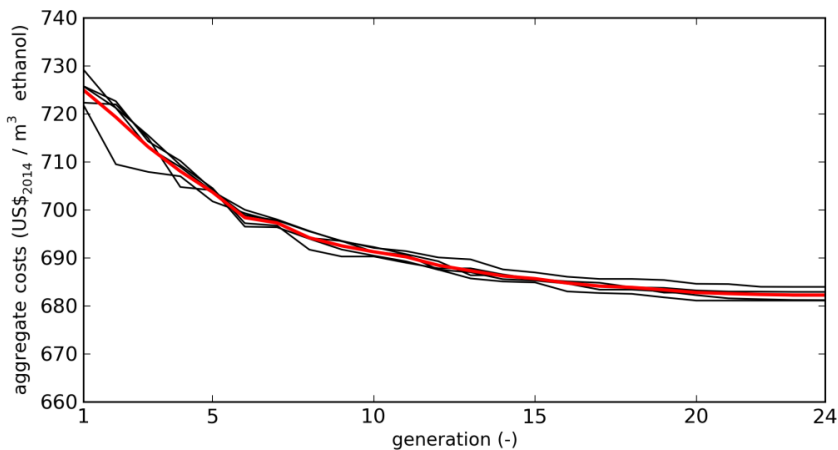


Figure 6.1: Conceptual model of the optimization with the genetic algorithm (below the dashed line) and the control variables and calculation of the objective values (above the dashed line).

A fraction of the population to reproduce of 0.2 means that the best 20% of the population is progressed to the next generation and this 20% creates the new 80% of the population by cross-over. When the fraction of the population to reproduce is too low, the objective value stabilizes too early, before the optimum is reached, because too little variation remains in the population. If the fraction is too high, individuals with a relatively low fitness reproduce, thereby not improving the fitness of the next generation.

The mutation rate is the fraction of the total population that will be mutated. If the mutation rate is too low, the GA can become stuck at a local optimum, while if it is too high, the genotypes of the individuals with a high fitness change too much and there is no convergence towards the optimum objective value (e.g. Bennett et al., 1998). For our optimization problem the fastest conversion towards the lowest minimum reached was with a population fraction to reproduce of 0.1 and a mutation rate of 0.3. During cross-over, individuals are split at two locations in the bit-string. The maximum number of bits to mutate in a single individual is two.

Next, we increased the population size and number of generations with these settings until no improvement in the objective function was reached anymore. This was at a population of 1000 (no improvement anymore for 10000) for 24 generations (no improvement anymore for 25). One run takes 28 hours on a Linux server with 16 GB RAM and 24 cores with 2 GHz Intel Xeon processors. Running the GA five times with these parameter settings for a single carbon price of 100 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq proved that at generation 24 the objective value of the best final individual is always stable for some consecutive generations, and has a maximum variation of 2.8 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq between the different runs (Figure 6.8).



**Figure 6.8:** Development of the objective value  $x$  over the generations of the GA for five different runs (black lines) at a carbon price  $p$  of 100 US\$<sub>2014</sub> / tonne CO<sub>2</sub>-eq. The red line indicates the mean over the five runs.



## **7. Synthesis**

### **7.1. Context**

Land use change caused by human drivers can have large impacts on, for instance, climate change, water availability and quality, soil quality and erosion, and biodiversity (Lambin and Meyfroidt, 2011, Nesheim et al., 2014). Assessing the future land use change effects of human drivers requires land use change models. Designing such models is not straightforward, because the dynamic processes and feedback loops in the land use system are complex and only partially understood (Manson, 2007, Verburg et al., 2013). This results in uncertainties in model structure, inputs and parameters, which propagate to the land use change projections. Another type of uncertainty in these projections, the solution space uncertainty, originates from the fact that only a limited number of scenarios can be analysed, causing a lack of clarity about the complete set of potential futures. Current land use change models do not quantify and communicate these uncertainties and thereby may create a false sense of certainty. The aim of this thesis was to develop methods to quantify and reduce uncertainty in land use projections. The focus herein was on bioenergy case studies because bioenergy is experiencing a large demand increase and the land use change impacts of this demand increase are crucial for its sustainability and therefore need to be examined. Four research questions were defined in Chapter 1. They are answered and discussed in sections 7.2 to 7.5. Section 7.6 and 7.7 present the future perspectives and recommendations for land use change modelling in general and for bioenergy case studies specifically.

### **7.2. How can uncertainty in land use change projections be quantified?**

Uncertainties in land use change projections of two types are evaluated in this thesis: uncertainties from errors propagating through the model to the outputs and solution space uncertainty (Table 7.1). To quantify uncertainty of error propagation, Monte Carlo (MC) analysis is used. The PCRaster Land Use Change model (PLUC) is developed, a stochastic land use change cellular automaton (CA) coupled to a MC scheme using the PCRaster Python framework (Karssenberget al., 2010). This MC analysis scheme also calculates summary statistics or spatial metrics, like the probability of occurrence of a particular land use type in a cell (Chapters 2, 3, and 5), the total number of patches (Chapter 3), or average area of a land use type at a higher aggregation level (Chapters 3, 4, and 5) per year.

**Table 7.1: Overview of the variables in the model components of which uncertainty is taken into account per chapter.**

type of uncertainty in projections	model component	variable	Chapter				
			2	3	4	5	6
propagating errors	non-spatial inputs	demand	x			x	
		maximum yield	x				
	spatial inputs	spatial attribute suitability factors	x				
		potential yield fraction	x				
		initial land use map				x	
	model structure	weights of suitability factors		x	x	x	
		order of allocation				x	
		systemic change			x		
	parameters	parameters of suitability factors	x	x	x		
	observations	land use composition		x	x	x	
land use configuration			x				
solution space	-	-				x	

The embedded coupling of the MC scheme with the land use change model allows for two activities that are important in uncertainty modelling. Firstly, because the definition of uncertainty in all model components is an integral part of the modelling framework, the end user can easily evaluate uncertainty for multiple scenarios or case studies. This is necessary as uncertainty in inputs, parameters, model structure or observations often varies between scenarios and especially between case studies. Secondly, the embedded coupling allows for iterative output uncertainty analysis, i.e. for each model time step, because the modelling framework comes with calculation and visualisation routines for spatio-temporal stochastic variables. This is important for complex systems, of which the land use system is an example, because they behave non-linearly through time (Manson, 2007). Therefore, an uncertainty map of the final time step is not always sufficient information for decision making, since the ranges of output uncertainty are likely to be non-linear through time too. PLUC is applied in Chapters 2 to 5 to iteratively evaluate the effects of uncertainties in all different model components: inputs, parameters, model structure and observations used to reduce the uncertainty (Table 7.1). Such a full scope error propagation assessment is new in land use change modelling.

To be able to evaluate the error propagation effects of uncertainties in different model components, the uncertainties in these components have to be estimated. To estimate uncertainty in the non-spatial inputs two methods have been used. In the first method (Chapter 2), the probability distribution is defined using a relative normal error model. This means that the error is higher for higher mean (expected)



values. The standard deviation herein is estimated based on predicted upper and lower limits for driving factors like population numbers. In the second method (Chapter 5), a general computable equilibrium (CGE) model is applied to estimate the demands for the different land use types. The probability distributions of these demands are estimated by varying the land transition elasticities, the most uncertain parameters in the CGE model. For the uncertainty in the spatial inputs, like the digital elevation model to calculate the slope, used as a suitability factor, normal and relative normal error models are applied on each cell value (Chapter 2). To represent uncertainty in the initial land use map, a method is developed to generate a separate realization of this initial land use map for each MC run based on presumed errors in the agricultural statistics that were used in creation of the map (Chapter 5).

To quantify uncertainty in the model structure, methods were developed to vary the order of allocation of the dynamic land use types (Chapter 5) and the weights of the suitability factors of these land use types (Chapter 3, 4, and 5). A problem herein is that one cannot use a simple error model, such as the normal or relative normal distribution, to generate a realization for a single transition rule (e.g. a weight or the position of a land use type in the order of allocation), because the overall model structure has to be coherent. For the given variables this means that two land use types cannot be in the same position in the order of allocation and the weights of all suitability factor together should sum to a value of one. Therefore, specific error models were developed in Chapters 3 to 5 to ensure a coherent model structure. An additional source of model structure uncertainty is that the model structure can become invalid through time because of a systemic change. This denotes that a certain model structure, which was at previous time steps able to give an accurate representation of the land use change system, has to be altered at a certain point in time to remain able to simulate this system. In Chapter 4, we have developed a method to quantify uncertainty stemming from this non-stationarity of a land use system. We assessed which values for weights and parameters were valid for different points in time and how the variation in these values (systemic change) affects projection uncertainty. A visual inspection and an analysis of the quantity of this variation, as well as the outcome of two statistical tests have provided a strong indication of non-stationarity for three out of the four weights of the suitability factors, indicating a period of systemic change.

For the model parameters there is not prior indication that they are correlated. Hence, simple error models are used such as the normal and uniform distribution to represent uncertainty in the parameters (Chapter 2 to 4). Herein, the standard deviation of the normally distributed parameters is estimated based on expert knowledge.

Uncertainty in the observational data, which are used to reduce projection uncertainty (see next section), is estimated in various ways throughout this thesis. In Chapter 3, we use a stochastic simulation method to simulate the observation variances and covariances. In Chapter 4, it is assumed that there is no spatial or temporal correlation in the errors of the observational data because the data are used at a higher aggregation level, and therefore only the variance has to be estimated, which is done from 1) errors in the classification of the remote sensing image, using the reported classification accuracy (Adami et al., 2012b), and 2) errors from the upscaling to a larger cell size. In Chapter 5, we estimate the observation error by comparing two data sources.

After running the MC analysis, uncertainty in the projected land use change for the required point in time is transformed to the attribute and scale of interest for the end user. We have shown that output uncertainty is highly scale dependent. In Chapter 2, the maximum variance in the total potential yield of eucalyptus in Mozambique drops from  $1 \cdot 10^5$  to 680 to 298 ( $\text{kg km}^{-2} \text{year}^{-1}$ )<sup>2</sup> when scaling from cell level to province level to country level, respectively. Chapter 5 illustrates that the maximum coefficient of variation (cv) in indirect land use change area decreases from 6 to 0.72, when going from regional (250 x 250 km<sup>2</sup> blocks) to national (Brazil) level. This is because local differences between realizations are levelled out at higher aggregation levels.

Uncertainty in the outputs does not only depend on the scale level, but also on the evaluated attribute. For Brazil as a whole, the direct land use change (dLUC) area has a coefficient of variation (cv) of only 0.02, while the indirect land use change (iLUC) area has a cv of 0.72 (Chapter 5). The uncertainty in iLUC area and location is generally higher than in dLUC area, because iLUC is caused by the interplay of various land use types that each have their uncertain demand and model structure, while dLUC is mainly affected by the demand and weights of the suitability factors for sugar cane. It is relevant for the end user to realize that some output variables respond differently to alterations in parameters and parameter uncertainties than other output variables.

Besides the propagated errors in the land use projections, solution space uncertainty exists in these projections. This solution space uncertainty can be quantified by optimization of certain objectives and Pareto frontier construction (Seppelt et al., 2013). A Pareto frontier provides all (in this case land use) alternatives for which it is impossible to improve one objective, without impairing another. Scenarios can be used to quantify impacts, but an optimization study in combination with a Pareto frontier visualisation demonstrates the lowest attainable values for all impacts and the trade-offs between different impacts (Chapter 6). We plea for a more extensive use of optimization in land use change studies that focus on impacts.

Thus, error propagation and solution space uncertainty effects on land use change projections have been quantified. Assessing the error propagation effects of uncertainties in different model components through MC analysis requires estimation of the uncertainties in these components. Although methods have been given for doing this, often strong assumptions had to be made about data accuracies. Data products about land use often include very little or no information on accuracy, or provide an overall value only, i.e. no accuracy information at cell level. Therefore we request the land use data suppliers to provide, preferably cell-based, accuracy information on their data products, to allow us to improve on this point. Another problem is that the model structure is a strong abstraction of the processes that can be observed. For example, the weights and the order of allocation of the land use types are variables for which values cannot be observed in reality. A model setup representing processes at a more detailed level could help to overcome this. Section 7.6 further elaborates on this potential solution.

### **7.3. How can uncertainty in land use change projections be reduced?**

Current problems in land use change modelling are that there are no methods for uncertainty reduction, that uncertainty in the observational data is not taken into account and that calibration is often not targeted on the attribute of interest related to the modelling aim. These problems are solved through the coupling of a particle filter to the PLUC model. A particle filter is a data assimilation technique (van Leeuwen, 2009) that updates prior knowledge about model structure, parameters and inputs during model runtime. The prior knowledge is represented by the probability distributions of all uncertain inputs, parameters and model structure, represented in the MC ensemble. The update occurs at time steps when observations of state variables or derived spatial metrics are available. It thereby takes into account the uncertainty in these observations. This method has the advantage that subjective knowledge of experts can be used to define uncertainty, in the form of probability distributions, in model structure, parameters and inputs, but then objective knowledge of observations is used to adjust these probability distributions. An additional advantage is that the posterior (calibrated) distribution of the model structure and parameters discloses information about how the land use system functions.

Obviously, the main advantage of the particle filter is that it can reduce the uncertainty in land use change projections and derived summary statistics or metrics. In Chapter 3, the 95% confidence intervals of three spatial metrics (landscape shape index, sugar cane fraction in 150 x 150 km<sup>2</sup> blocks, and number of sugar cane patches) were reduced by at least a factor 3 by the particle filter, compared to an MC run without the particle filter. In Chapter 5, the root mean

squared error (RMSE) in the modelled median land use areas per state were reduced by 20% on average, compared to an MC run without the particle filter. Here the average uncertainty reduction is relatively low compared to Chapter 3, because the RMSE of the land use types rangeland and planted forest is not improved at all, probably due to poor data availability for the drivers of location of land-use change and because the observational data of these land use types are problematic to derive from remote sensing. The performance improvement by the particle filter holds for most spatial metrics when observational data are incomplete, in space or in time. This was demonstrated in Chapter 3 by using only half of the observational data. The decrease in RMSE of the landscape shape index and the number of patches compared to a run without particle filter was on average still ~50% five years after the calibration period, compared to ~70% when using all data. This is a considerable advantage given the fact that time series of good quality land use maps are rare, and thus missing data in a calibration time series is common (Straatman et al., 2004).

A challenge in applying the particle filter is the estimation of the prior probability distributions of inputs, parameters, model structures and observational data. Although it is an advantage that these distributions can be defined by experts, it should be kept in mind that this prior information has an effect on the outcomes. This effect has not been assessed quantitatively, but we stress that candidate suitability factors and the prior distribution of their weights should be selected sensibly. The potential solution to use a larger number of candidate suitability factors to ensure that all possible drivers are considered is not feasible, because the addition of only one parameter causes an exponential increase in the number of required particles (Bengtsson et al., 2008). Three solutions can be considered: 1) make a selection of parameters stochastic, preferably based on sensitivity analyses that indicate which inputs and parameters have most influence on the output, while fixing the others at a single deterministic value. Reducing the number of stochastic inputs and parameters, reduces the number of MC realizations required. 2) Use super computers or cluster machines, to allow a rapid calculation of all MC realizations. For example, on our Linux server with 24 cores the 500 sample Mozambique run takes two hours instead of almost two days. 3) Use of a more advanced particle filter scheme giving similar results with a lower number of particles and thus shorter run time (Spiller et al., 2008, Jeremiah et al., 2012).

While the use of the particle filter reduces uncertainty, land use projections of more than a few years ahead are not very reliable, as the 95% confidence interval of all spatial metrics quickly increase in width over the projection period in Chapter 3. For example, if one wants to be able to predict the number of patches with a maximum 95% confidence interval width of 1000 patches, which is wide considering that it is about 50% of the mean, this is possible only for a projection

period of three years ahead. This is very problematic given that 1) policy analysis and implementation typically take longer than that, and 2) 'three years ahead' often means that the current point in time is not even reached, because the initial land use map generally dates from a few years ago. For example, in Chapter 5, the most recent initial land use map we could construct for Brazil was for 2006. The calibration period was 2006 to 2012, so three years ahead from that means a projection for 2015, the same year that the chapter was published as an article and thus not really a future projection.

We conclude that a particle filter has the ability to reduce the uncertainty in land use projections, but not enough to ensure reliable projections for time frames longer than a decade, even for spatially aggregated attributes such as the total number of patches in the study area.

#### **7.4. What are the contributions of different model components to the uncertainty in land use change projections?**

PLUC is applied in Chapters 2 to 5 to evaluate the effects of uncertainties in different model components on different attributes related to land use or land use change (Table 7.1). It is difficult to make statements about which model component contributes most to the land use projection uncertainty, because, as discussed in section 7.2, uncertainty depends on which attribute is evaluated. Thereby, the contribution to uncertainty of the components also depends on which attribute is evaluated. However, we did quantify the contribution to uncertainty of different components at different spatial scales. In general, non-spatial inputs determine projection uncertainty at high aggregation levels (coarse scale) and spatial inputs determine projection uncertainty at low aggregation levels (fine scale).

For example, in Chapter 3, at high aggregation levels (national, i.e. the whole of Mozambique, or provincial), the variance in total potential eucalyptus yield is mainly determined by the interplay of demand and maximum yield, because local uncertainties, stemming from variations within the maps, are levelled out, as also explained in section 7.2. At the cell ( $1 \times 1 \text{ km}^2$ ) level, the local differences do matter. Most cells have a variance of zero, because their potential yield fraction is zero or because they have a bioenergy availability probability of zero, i.e. they are currently or in the future used for the production of food, feed or fibres and should therefore not be used for bioenergy crop cultivation. Some cells have a variance slightly above zero; these cells have an availability probability of one, but their yield differs because of the stochastic parameters in yield fraction and maximum yield. Finally, there are some cells with very high variances; these cells are in some MC

runs unavailable, and then have a yield of zero, and in others available, and then have a yield dependent on the stochastic parameters yield fraction and maximum yield. This generates very large variances, up to  $1 \cdot 10^5$  ( $\text{kg km}^{-2} \text{ year}^{-1}$ )<sup>2</sup>.

The fact that the contribution of different components to output uncertainty depends on the evaluated attribute is shown in Chapter 5. At national level, the cv of the dLUC area is caused by the uncertainties in the CGE model for 100%. The contribution of the uncertainties in the economic model to the cv of iLUC area at this level is about 93%, although this cannot be determined precisely, because errors from the two models partly compensate each other.

Another factor causing uncertainties in the outputs is the uncertainty from non-stationarity in model structure and parameters (Chapter 4). The systemic change increased the width of the projected 95% confidence interval of the sugar cane area per  $25 \times 25 \text{ km}^2$  block by a factor 2, compared to a run without systemic change. The systemic change appeared to be indirect: something has an effect on the input demand for sugar cane, in such a way that the transition rules and parameters have to change as well (Filatova and Polhill, 2012). But, although an inventory was made of societal changes in the study area during the studied period, none of these could be related to the onset of the observed systemic change in the land use system in 2006. We do realize that our method to detect systemic changes has drawbacks. Therefore, we suggest it should be further developed, but also stress that for any case study it should always be questioned whether or not continuation of past trends is expected.

As for the other components, such as observational data and parameters, it is difficult to draw general conclusions about their contribution to output uncertainty because we did not systematically turn uncertainty in these components on and off to test the sensitivity of the land use change projections to these uncertainties. We suggest this as a line for further research, although it still has to be proven whether general conclusions, valid throughout scenarios and case study areas, can be derived from such experiments.

To conclude, it is shown that quantification of uncertainties of the different model components gives insights in the contribution of these components to different output variable at different scales. The assessment of the propagation of uncertainties enables the identification of the components with the highest priority for improvement given the aim of the end user.

## **7.5. What are the implications of the land use change projection uncertainties for bioenergy implementation strategies?**

The uncertainty in land use change model projections is rarely communicated to the end users, which is problematic given that a land use projection that is erroneous or has not been put into perspective can result in wrong conclusions (Pontius Jr. and Spencer, 2005, Moulds et al., 2015). In this thesis, uncertainty was not only quantified and reduced, but also visualized in a Spatial Decision Support System (SDSS) (Chapter 2). We deem that, independent of the implications of uncertainty, we at least created the conditions for end users to take into account uncertainty in their decisions by presenting uncertainty in the land use projections in an understandable way. Next, the implications for bioenergy implementation strategies of the following aspects are discussed: 1) projection uncertainty, 2) land use system non-stationarity, 3) solution space uncertainty, and 4) input data.

The most important implication of the presented uncertainties in land use change projections for the bioenergy community is the achieved comprehension that the direct effects, and even more the indirect effects of bioenergy crop expansion, are difficult to project. In Chapter 5, cell-based ( $5 \times 5 \text{ km}^2$ ) probabilities of dLUC range from 0 to 0.77, and of iLUC from 0 to 0.43. In  $250 \times 250 \text{ km}^2$  blocks, the maximum cv is 4 for dLUC and 6 for iLUC. For Brazil as a whole, the dLUC area has a cv of only 0.02, while the iLUC area has a cv of 0.72 (Chapter 5). The latter means that, considering the width of the 95% confidence interval, the iLUC area in Brazil might be 2.4 times as high or as low as the projected mean. Because this confidence interval is so wide, it is likely it will straddle any selected legislation threshold. Thus, threshold evaluation for iLUC indicators is a very questionable practice. Therefore, we propose, in line with others (e.g. Finkbeiner, 2014, Mathews and Tan, 2009), a change of focus from quantifying iLUC to taking proactive measures to mitigate iLUC, even though we know that the effectiveness of these measures is difficult to quantify.

One potential option to reduce uncertainty in a linked CGE - land use change model would be to link the models with a hard link that includes a feedback, as also suggested by Wicke et al. (2012). However, there is an inherent risk that this feedback loop is infinite, meaning that the land use dynamics cannot be resolved. In addition, there are many technical obstacles that complicate hard linking, like dissimilar programming languages or dissimilar spatial and temporal discretizations between the two models (Schmitz et al., 2014). And also, hard linking might partly solve the uncertainties stemming from the conceptual differences between the CGE and the land use change model, but it does not solve the structural uncertainties within the CGE and land use change model. Besides the coupling of a CGE model to the land use change model, other multi-model approaches could be

beneficial, not only through coupling but also through collaboration between modellers in harmonizing input data and parameters (Wicke et al., 2015).

At the aggregation level of Brazil, dLUC can be to a certain extent be projected. It has a cv of 0.02 for Brazil. Note however, that at the aggregation level of the whole of Brazil the CGE model directly projects dLUC; no spatially explicit land use change model is required at this level. At the level of the 250 x 250 km<sup>2</sup> blocks, the maximum cv of 4 for dLUC is still very high, but there are also areas about which something can be said with a cv of about 0.5. Land use change projections at this level are functional in our opinion to select areas of interest for further impact analyses, for example using optimization approaches like in Chapter 6 (see next section).

A land use system can be non-stationary, i.e. a systemic change can occur, as illustrated by Chapter 4. There was an indication that the observed systemic change was indirect: something has an effect on the input demand for a particular land use type, in such a way that the transition rules and parameters have to change as well (Filatova and Polhill, 2012). In other words, the demand trend over time suddenly changes, beyond the function domain of the transition rules, with the result that the transition rules become invalid. Although it could not be tested, it can be expected that the introduction of a new bioenergy crop in an area and the implementation of new policies related to that have such an effect. If the bioenergy crop is a new crop in that area, the demand was zero before, so the trend changes by definition, and if the bioenergy crop was already in cultivation in the area for food, feed or fibres, the demand is expected to experience a sudden upsurge. These changes can commence systemic change, which is not taken into account in the land use change model projections currently used. Therefore, the current land use change models might not be very suitable for such ex-ante evaluations of new technologies/crops. However, this should be tested on a case study, and preferably on several case studies, with data on the introduction period of the bioenergy crop.

The solution space uncertainty quantification shown in this thesis has constructive implications for bioenergy implementation strategies, especially because bioenergy certification schemes require quantification of different environmental and socio-economic impacts (e.g. Roundtable on Sustainable Biofuels, 2010). In Chapter 6, we demonstrate how a Pareto frontier between production cost and GHG emission objectives can be constructed for an ethanol supply increase and how information derived from it could aid to formulate management or policy recommendations. The scenario projection for Goiás for 2030 using current land use change trends results in production costs of 715 US\$<sub>2014</sub> / m<sup>3</sup> ethanol and GHG emissions of  $-80 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol. The Pareto frontier shows that the minimum attainable production costs are 656 US\$<sub>2014</sub> / m<sup>3</sup> ethanol and the minimum



attainable GHG emissions are  $-399 \cdot 10^{-3}$  tonne CO<sub>2</sub>-eq / m<sup>3</sup> ethanol. This puts the scenario in the scope of what is possible and thereby supports the decision making process. The developed methodology has the prospect to identify trade-offs and win-win situations for other regions, scales, objectives and commodities.

Final implications are related to the input data. The large uncertainty in two important land use change model inputs considerably decreases the value of land use change models for the evaluation of bioenergy related questions: potential yield maps and demands. In all chapters in this thesis, potential yield maps for different crops and for pasture are used from the IIASA (Tóth et al., 2012). In Chapter 2, the demand for land is given in tonnes of products, and the potential yield maps are used to rate the total mass of the product that can be obtained from a cell. As such, the potential yield maps have a large influence on the results. For Brazil, however, the consulted experts had doubts about the validity of these yield maps. Therefore, in Chapter 3 to 5, the demand is given in areal units, and the potential yield maps are used as a suitability factor only. In Chapter 6, the potential yield for sugar cane is one of the main determinants of the calculated production price and GHG emissions of bioethanol. It is helpful that global potential yield maps are freely available, but more trustworthy potential yield maps, including cell based projections of yield developments, and, perhaps more importantly, cell based information on the accuracy of these maps, would be a huge asset for future land use change projection efforts.

Also related to input data, in Chapter 2, the uncertainty in the demand input, constructed through an extrapolation of population and diet, is one of the main factors determining output uncertainties (see previous paragraphs). So, even though the applied land use change model has its intrinsic uncertainties too, the main unknown factor in determining which part of the land will be available is this demand. In Chapter 5, it was tried to solve this issue, by coupling a CGE model to the land use change model. However, two new issues arise here. Firstly, the economic model turns out to generate large uncertainties in projected land areas too (ranges of ~10-20%). And secondly, the dissimilarity in model concepts between the two models within the integrated model chain augments output uncertainty.

It is positive that spatially explicit analyses are starting to become prevalent in bioenergy impact assessment. The uncertainty quantification demonstrated in this thesis can help to guide users what such analysis can and cannot be used for. The next section elaborates on potential future directions for land use change projections.

## **7.6. Future perspectives for spatially explicit land use change projections**

Given the very wide uncertainty ranges of the projected attributes evaluated in this thesis, we conclude that the value of spatially explicit land use change models for answering questions about the future directions of the system at local scale levels (up to 100 x 100 km<sup>2</sup>) and for long time frames (more than a decade) is limited. We add the remark that these conclusions are based on model studies with PLUC only. Although we deem that the conclusions are applicable to demand-driven land use change models in general, this cannot be proven. We challenge other modellers to make their models stochastic and prove us wrong in our scepticism about the use of land use change models at local scales and for long time frames. Yet, there are conditions under which, or aims for which spatially explicit land use change models can be valuable decision support tools in our opinion.

At higher aggregation levels uncertainties can be reasonable. Conclusions about regions, e.g. state or province level, can be drawn, but not for all attributes of interest and not for all time frames. Therefore, even when results are aggregated, it is important that uncertainties are quantified. Our research facilitates modellers herein. An interesting application of regional assessments is to identify a region of interest, for example a region where a bioenergy crop can be expected to expand with reasonable confidence. For this region of interest an optimization and Pareto frontier approach can be applied to see how bioenergy cropland expansion can be regulated to obtain minimal negative impacts.

Uncertainties of differences between two scenarios are often lower than the uncertainties within one scenario (Chapter 5). This is valid under the assumption that, although e.g. a parameter value might be unknown, it will probably be the same value for different scenarios. Therefore, land use change models are more valuable for the relative comparison of impacts between two (or more) scenarios than for the assessment of the absolute impacts of a single scenario. Problematic herein is that the uncertainty between two scenarios becomes smaller when the scenarios are more divergent in terms of demands, but the likelihood that the assumption about the correspondence of parameter values and model structures holds for these scenarios becomes lower. This is because (indirect) systemic change is more likely to occur (Chapter 4).

Under the condition that the demand can be accurately estimated, communication of information via visualization is one of the advantages of land use change projections. When society is informed that, for example, 35 billion litres of ethanol will be produced in Brazil by 2020 (example from Lapola et al., 2010), only very few people will be able to grasp this number, even when converted to areal units, e.g. 90 000 km<sup>2</sup>. In the media, such numbers are usually converted to units of

something people are familiar with, in the Netherlands usually the number of soccer fields. Even for a soccer player like myself, conceiving the extent of anything above about 10 soccer fields is problematic, which is certainly at stake in the example of 90 000 km<sup>2</sup>. The meaning of such a number is easier to grasp when the extent is shown on a map. For this purpose, policy makers and other end users might benefit from the visualization of land use change shown by land use change models. Yet, hereto the sophistication level that land use change models currently have is not necessary, because it is of lesser relevance where exactly the land use is allocated. Our suggestion is to create a 'light' land use change model with fewer parameters as a web-based simulation tool. Such a web-based tool has, additional to the benefit of visualisation, numerous advantages like the ability to use the model at all times at all places (with internet access), easy access for everyone, provision of computation power by the server instead of the user's computer, and easy maintenance by the modeller for all users at once (Byrne et al., 2010).

Land use change models might be more valuable for local assessments if they become more process based. For example, all suitability maps of the suitability values could be expressed in terms of costs (e.g. Koomen et al., 2015), instead of as normalized values between zero and one. The factor 'distance to roads' can be transformed to transportation costs, the factor 'potential yield' can be transformed to cultivation costs minus revenues, et cetera, similar to the approach used in Chapter 6 to calculate production costs of ethanol. This approach has two advantages. Firstly, all suitability maps can be summed without using the intangible and immeasurable 'weights'. Secondly, there is a chance that systemic change, present in more conceptual models, can be avoided, because the modelled processes have been described more realistically. This reduces the uncertainty from the wrongfully stationary model structure. Two disadvantages are that some suitability factors are difficult to express in terms of money and that the new model will have more parameters, including ones that are problematic to project into the future, e.g. petrol prices that are included in transportation costs.

## **7.7. Recommendations**

- Better-quality data are required as the inputs and observational data of land use change models, e.g. time series from classified remote sensing images, agricultural statistics databases and potential yield maps, preferably with a global coverage. In addition, more information is required on uncertainty in these data products. It would be helpful for land use change modellers if uncertainty in these sources was quantified and distributed alongside the products.

- Uncertainty in spatially explicit land use change model projections should be quantified, reduced by observations, and communicated to the end users, even when the accuracies of the model components are difficult to estimate; concealing uncertainty is not the solution. Users of land use change models indicators should have the opportunity to grasp the (un)reliability of these models.
- Spatial aggregation of land use projections is recommended because local uncertainties are too high. Depending on the case study and the attribute of interest, projections at higher aggregation levels might be reliable enough to support decisions, as uncertainty is highly scale-dependent with lower uncertainties at higher spatial scale levels.
- If one really wants to use spatially explicit land use change models to project future land use changes locally, other approaches should be examined to try to improve the predicative value of the land use projections. Suggestions for improvements are collaborations with modellers from other domains, hard links with models from these domains, and more process based model rules in the land use change model.
- The creation of a 'light', web-based version of a spatially explicit land use change model, including visualization tools, is recommended. A quick, uncertainty-inclusive, accessible and understandable model for a broad public, would be constructive for decision making throughout the world and throughout decision making levels.
- More research is required regarding systemic change in land use systems. An especially interesting question is whether the land use system remains stationary when a new crop with new technologies, e.g. a bioenergy crop, is introduced.
- The usage of optimization and Pareto frontiers should be expanded in land use change modelling, because scenario studies alone do not describe the complete solution space of potential impacts with the result that win-win land configurations might be overlooked.

## References

- ABRAF, 2013. Yearbook Statistical ABRAF 2013, base year 2012, Brazilian Association of Forest Plantation Producers (ABRAF), Brasilia, Brazil.
- Adami, M., Mello, M.P., Aguiar, D.A., Rudorff, B.F.T., de Souza, A.F., 2012a. A web platform development to perform thematic accuracy assessment of sugarcane mapping in South-Central Brazil. *Remote Sensing* 4(10), 3201-3214.
- Adami, M., Rudorff, B.F.T., Freitas, R., Aguiar, D.A., Sugawara, L.M., Mello, M.P., 2012b. Remote Sensing Time Series to Evaluate Direct Land Use Change of Recent Expanded Sugarcane Crop in Brazil. *Sustainability* 4(4), 574-585.
- Aerts, J.C.J.H., Clarke, K.C., Keuper, A.D., 2003a. Testing popular visualization techniques for representing model uncertainty. *Cartography and Geographic Information Science* 30(3), 249-261.
- Aerts, J.C.J.H., Goodchild, M.F., Heuvelink, G.B.M., 2003b. Accounting for Spatial Uncertainty in Optimization with Spatial Decision Support Systems. *Transactions in GIS* 7(2), 211-230.
- Agarwal, C., Green, G.M., Grove, J.M., Evans, T.P., Schweik, C.M., 2002. A Review and Assessment of Land-Use Change Models: Dynamics of Space, Time, and Human Choice, US Department of Agriculture, Forest Service, 1.
- Aguiar, A.P.D., Câmara, G., Escada, M.I.S., 2007. Spatial statistical analysis of land-use determinants in the Brazilian Amazonia: Exploring intra-regional heterogeneity. *Ecological Modelling* 209, 169-188.
- Aguiar, D.A., Personal communication, July 17th 2014.
- Aguiar, D.A., da Silva, W.F., Rudorff, B.F.T., Adami, M., 2010. Canasat project: monitoring the sugarcane harvest type in the state of São Paulo, Brazil, in Wagner, W., Székely, B. (Eds.), *Proceedings of the International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences (ISPRS) TC VII Symposium – 100 Years. ISPRS, Vienna, Austria, 10-15.*
- Aguiar, D.A., Rudorff, B.F.T., Silva, W.F., Adami, M., Mello, M.P., 2011. Remote Sensing Images in Support of Environmental Protocol: Monitoring the Sugarcane Harvest in São Paulo State, Brazil. *Remote Sensing* 3(12), 2682-2703.
- Aguiar, G.A.M., d'Athayde, H.P., 2014. Information about cattle production and agriculture in Brazil, SCOT Consultoria, Bebedouro, São Paulo, Brazil.
- Akgul, O., Shah, N., Papageorgiou, L.G., 2012. An optimisation framework for a hybrid first/second generation bioethanol supply chain. *Computers and Chemical Engineering* 42, 101-114.
- Alexander, P., Rounsevell, M.D.A., Dislich, C., Dodson, J.R., Engström, K., Moran, D., 2015. Drivers for global agricultural land use change: The nexus of diet, population, yield and bioenergy. *Global Environmental Change* 35, 138-147.
- Alexandratos, N., Bruinsma, J., 2012. World Agriculture Towards 2030/2050: The 2012 Revision. ESA Working paper No. 12-03, Food and Agriculture Organization of the United Nations, Agricultural Development Economics Division, Rome, Italy.

- Aljoufie, M., Zuidgeest, M., Brussel, M., van Vliet, J., van Maarseveen, M., 2013. A cellular automata-based land use and transport interaction model applied to Jeddah, Saudi Arabia. *Landscape and Urban Planning* 112(1), 89-99.
- Allen, T.R., Wang, Y., Crawford, T.W., 2013. Remote Sensing of Land Cover Dynamics. *Treatise on Geomorphology* 3, 80-102.
- Alonso, W., 1964. *Location and land use: Toward a general theory of land rent*. Harvard University Press, London, United Kingdom.
- AMORI, 2009. Automatic Model Optimization Reference Implementation, Available online at <http://sourceforge.net/projects/amori>.
- Andrade de Sá, S., Palmer, C., di Falco, S., 2013. Dynamics of indirect land-use change: Empirical evidence from Brazil. *Journal of Environmental Economics and Management* 65(3), 377-393.
- Arino, O., Bicheron, P., Achard, F., Latham, J., Witt, R., Weber, J.-., 2008. GlobCover: The most detailed portrait of Earth. *European Space Agency Bulletin* 2008(136), 24-31.
- Arndt, C., Benfica, R., Tarp, F., Thurlow, J., Uaiene, R., 2010. Biofuels, poverty, and growth: a computable general equilibrium analysis of Mozambique. *Environment and Development Economics* 15, 81-105.
- Arulampalam, M.S., Maskell, S., Gordon, N., Clapp, T., 2002. A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *IEEE Transactions on Signal Processing* 50(2), 174-188.
- Arvor, D., Dubreuil, V., Ronchail, J., Simões, M., Funatsu, B.M., 2014. Spatial patterns of rainfall regimes related to levels of double cropping agriculture systems in Mato Grosso (Brazil). *International Journal of Climatology* 34(8), 2622-2633.
- Aspinall, R., 2004. Modelling land use change with generalized linear models - A multi-model analysis of change between 1860 and 2000 in Gallatin Valley, Montana. *Journal of environmental management* 72, 91-103.
- Bakker, M.M., Hatna, E., Kuhlman, T., Mùcher, C.A., 2011. Changing environmental characteristics of European cropland. *Agricultural Systems* 104(7), 522-532.
- Bakker, M.M., Veldkamp, A., 2012. Changing relationships between land use and environmental characteristics and their consequences for spatially explicit land-use change prediction. *Journal of Land Use Science* 7(4), 407-424.
- Banco Central do Brasil, 2015. *Price Indices in Brazil; with information up to March 2015*, Banco Central do Brasil, Economic Policy Board, Investor Relations and Special Studies Department, Brazil.
- Banse, M., van Meijl, H., Tabeau, A., Woltjer, G., Hellmann, F., Verburg, P.H., 2011. Impact of EU biofuel policies on world agricultural production and land use. *Biomass and Bioenergy* 35(6), 2385-2390.
- Bastin, L., Cornford, D., Jones, R., Heuvelink, G.B.M., Pebesma, E., Stasch, C., Nativi, S., Mazzetti, P., Williams, M., 2013. Managing uncertainty in integrated environmental modelling: The UncertWeb framework. *Environmental Modelling and Software* 39, 116-134.
- Batidzirai, B., Faaij, A.P.C., Smeets, E.M.W., 2006. Biomass and bioenergy supply from Mozambique. *Energy for Sustainable Development* 10(1), 54-81.

- Batjes, N.H., 2010. IPCC default soil classes derived from the Harmonized World Soil Data Base (Ver. 1.1), Carbon Benefits Project (CBP) and ISRIC - World Soil Information, Wageningen, The Netherlands (with dataset).
- Batty, M., 2012. Building a science of cities. *Cities* 29, S9-S16.
- Batty, M., 2005. Agents, cells, and cities: New representational models for simulating multiscale urban dynamics. *Environment and Planning A* 37(8), 1373-1394.
- Bengtsson, T., Bickel, P., Li, B., 2008. Curse-of-dimensionality revisited: Collapse of the particle filter in very large scale systems. *Probability and Statistics: Essays in Honor of David A.Freedman* 2, 316-334.
- Bennett, D.A., Armstrong, M.P., Wade, G.A., 1998. Exploring the solution space of semi-structured geographical problems using genetic algorithms. *Transactions in GIS* 3(1), 51-71.
- Berjak, S.G., Hearne, J.W., 2002. An improved cellular automaton model for simulating fire in a spatially heterogeneous Savanna system. *Ecological Modelling* 148(2), 133-151.
- Bernoux, M., Cerri, C.C., Cerri, C.E.P., Siqueira Neto, M., Metay, A., Perrin, A.-., Scopel, E., Razafimbelo, T., Blavet, D., Piccolo, M.D.C., Pavei, M., Milne, E., 2006. Cropping systems, carbon sequestration and erosion in Brazil, a review. *Agronomy for Sustainable Development* 26(1), 1-8.
- Bettencourt, L.M.A., 2013. The origins of scaling in cities. *Science* 340(6139), 1438-1441.
- Bierkens, M.F.P., Finke, P.A., de Willigen, P., 2000. Upscaling and downscaling methods for environmental research. Kluwer Academic Publishers, Dordrecht, The Netherlands.
- Blecic, I., Cecchini, A., Trunfio, G.A., 2015. How much past to see the future: a computational study in calibrating urban cellular automata. *International Journal of Geographical Information Science* 29(3), 349-374.
- Blok, K., 2006. Introduction to energy analysis. Techne Press, Amsterdam, The Netherlands.
- Blöschl, G., 1999. Scaling issues in snow hydrology. *Hydrological Processes* 13, 2149-2175.
- Blum, C., Roli, A., 2003. Metaheuristics in Combinatorial Optimization: Overview and Conceptual Comparison. *ACM Computing Surveys* 35(3), 268-308.
- Bouwman, A.F., Kram, T., Klein Goldewijk, K., 2006. Integrated modelling of global environmental change - An overview of IMAGE 2.4, Netherlands Environmental Assessment Agency (MNP), Bilthoven, The Netherlands.
- Bradshaw, G.A., Borchers, J.G., 2000. Uncertainty as information: narrowing the science-policy gap. *Conservation Ecology* 4(1), online version.
- Broch, A., Hoekman, S.K., Unnasch, S., 2013. A review of variability in indirect land use change assessment and modeling in biofuel policy. *Environmental Science and Policy* 29, 147-157.
- Brown, D.G., Page, S., Riolo, R., Zellner, M., Rand, W., 2005. Path dependence and the validation of agent-based spatial models of land use. *International Journal of Geographical Information Science* 19(2), 153-174.
- Brown, D.G., Robinson, D.T., An, L., Nassauer, J.I., Zellner, M., Rand, W., Riolo, R., Page, S.E., Low, B., Wang, Z., 2008. Exurbia from the bottom-up: Confronting empirical challenges to characterizing a complex system. *Geoforum* 39(2), 805-818.

- Brown, J.D., Heuvelink, G.B.M., 2007. The Data Uncertainty Engine (DUE): A software tool for assessing and simulating uncertain environmental variables. *Computers and Geosciences* 33(2), 172-190.
- Burnicki, A.C., 2011. Modeling the probability of misclassification in a map of land cover change. *Photogrammetric Engineering and Remote Sensing* 77(1), 39-50.
- Burrough, P.A., McDonnell, R.A., 1998. *Principles of geographical information systems*. Oxford University Press, Oxford, United Kingdom.
- Byrne, J., Heavey, C., Byrne, P.J., 2010. A review of Web-based simulation and supporting tools. *Simulation Modelling Practice and Theory* 18(3), 253-276.
- Canadell, J.G., Le Quééré, C., Raupach, M.R., Field, C.B., Buitenhuis, E.T., Ciais, P., Conway, T.J., Gillett, N.P., Houghton, R.A., Marland, G., 2007. Contributions to accelerating atmospheric CO<sub>2</sub> growth from economic activity, carbon intensity, and efficiency of natural sinks. *Proceedings of the National Academy of Sciences of the United States of America* 104(47), 18866-18870.
- Carlson, K.M., Curran, L.M., Ratnasari, D., Pittman, A.M., Soares-Filho, B.S., Asner, G.P., Trigg, S.N., Gaveau, D.A., Lawrence, D., Rodrigues, H.O., 2012. Committed carbon emissions, deforestation, and community land conversion from oil palm plantation expansion in West Kalimantan, Indonesia. *Proceedings of the National Academy of Sciences of the United States of America* 109(19), 7559-7564.
- Cerqueira Leite, R.C.d., Verde Leal, M.R.L., Barbosa Cortez, L.A., Griffin, W.M., Gaya Scandiffio, M.I., 2009. Can Brazil replace 5% of the 2025 gasoline world demand with ethanol? *Energy* 34(5), 655-661.
- Chang, N.B., Parvathinathan, G., Breeden, J.B., 2008. Combining GIS with fuzzy multicriteria decision-making for landfill siting in a fast-growing urban region. *Journal of Environmental Management* 87(1), 139-153.
- Chaplin-Kramer, R., Sharp, R.P., Mandle, L., Sim, S., Johnson, J., Butnar, I., Milà l Canals, L., Eichelberger, B.A., Ramler, I., Mueller, C., McLachlan, N., Yousefi, A., King, H., Kareiva, P.M., 2015. Spatial patterns of agricultural expansion determine impacts on biodiversity and carbon storage. *Proceedings of the National Academy of Sciences of the United States of America* 112(24), 7402-7407.
- Chaudhuri, G., Clarke, K.C., 2013. The SLEUTH Land Use Change Model: A Review. *Environmental Resources Research* 1(1), 88-105.
- Chen, H., Wood, M.D., Linstead, C., Maltby, E., 2011. Uncertainty analysis in a GIS-based multi-criteria analysis tool for river catchment management. *Environmental Modelling and Software* 26(4), 395-405.
- Chen, X., Khanna, M., Yeh, S., 2012. Stimulating learning-by-doing in advanced biofuels: Effectiveness of alternative policies. *Environmental Research Letters* 7(4), Article number 045907.
- Chikumbo, O., Goodman, E., Deb, K., 2015. Triple Bottomline Many-Objective-Based Decision Making for a Land Use Management Problem. *Journal of Multi-Criteria Decision Analysis* 22, 133-159.
- Cihlar, J., Jansen, L.J.M., 2001. From land cover to land use: A methodology for efficient land use mapping over large areas. *Professional Geographer* 53(2), 275-289.



- Claessens, L., Schoorl, J.M., Verburg, P.H., Geraedts, L., Veldkamp, A., 2009. Modelling interactions and feedback mechanisms between land use change and landscape processes. *Agriculture, Ecosystems and Environment* 129, 157-170.
- Clarke, K.C., Hoppen, S., Gaydos, L.J., 1997. A self-modifying cellular automaton model of historical urbanization in the San Francisco Bay area. *Environment and Planning B: Planning and Design* 24(2), 247-261.
- Collin, A., Bernardin, D., Sero-Guillaume, O., 2011. A physical-based cellular automaton model for forest-fire propagation. *Combustion Science and Technology* 183(4), 347-369.
- Conab, 2014. Séries Históricas Relativas às Safras 1976/77 a 2013/14 de Área Plantada, Produtividade e Produção.
- Costa, P.A.L., Gray, K., 2011. Brazilian Senate Approves New Forest Code - New law guarantees preservation of 80% of the Amazon Rainforest.
- Cressie, N.A.C., 1993. *Statistics for spatial data*. Wiley, New York, USA.
- Creutzig, F., Popp, A., Plevin, R., Luderer, G., Minx, J., Edenhofer, O., 2012. Reconciling top-down and bottom-up modelling on future bioenergy deployment. *Nature Climate Change* 2(5), 320-327.
- Creutzig, F., Ravindranath, N.H., Berndes, G., Bolwig, S., Bright, R., Cherubini, F., Chum, H.L., Corbera, E., Delucchi, M., Faaij, A.P.C., Fargione, J., Haberl, H., Heath, G., Lucon, O., Plevin, R., Popp, A., Robledo-Abad, C., Rose, S., Smith, P., Stromman, A., Suh, S., Masera, O., 2015. Bioenergy and climate change mitigation: An assessment. *Global Change Biology Bioenergy* 7(5), 916-944.
- Csillag, F., Boots, B., 2005. Toward Comparing Maps as Spatial Processes, in Fisher, P. (Ed.), *Developments in Spatial Data Handling; 11th International Symposium on Spatial Data Handling*. Springer Berlin Heidelberg, 641-652.
- CSR/Ibama, MMA, SBF, 2013. Projeto de Monitoramento do Desmatamento dos Biomas Brasileiros por Satélite - PMDBBS.
- CTBE, 2012. CanaSoft (version 0.15). Technological Assessment Program - Brazilian Bioethanol Science and Technology Laboratory (CTBE).
- Dahan, L., Rittenhouse, K., Francis, D., Sopher, P., Schwartz, J., De Clara, S., 2015. *The World's Carbon Markets: A Case Study Guide to Emissions Trading - Brazil: an emissions trading case study*, International Emissions Trading Association (IETA), Paris, France.
- Dai, E., Wu, S., Shi, W., Cheung, C.-., Shaker, A., 2005. Modeling Change-Pattern-Value dynamics on land use: An integrated GIS and Artificial Neural Networks approach. *Environmental Management* 36(4), 576-591.
- de Meyer, A., Cattrysse, D., Rasinmäki, J., Van Orshoven, J., 2014. Methods to optimise the design and management of biomass-for-bioenergy supply chains: A review. *Renewable and Sustainable Energy Reviews* 31, 657-670.
- de Miranda, S.D.C., Bustamante, M., Palace, M., Hagen, S., Keller, M., Ferreira, L.G., 2014. Regional variations in biomass distribution in Brazilian Savanna Woodland. *Biotropica* 46(2), 125-138.
- de Souza Ferreira Filho, J.B., Horridge, M., 2014. Ethanol expansion and indirect land use change in Brazil. *Land Use Policy* 36, 595-604.

- de Souza Soler, L., Verburg, P.H., 2010. Combining remote sensing and household level data for regional scale analysis of land cover change in the Brazilian Amazon. *Regional Environmental Change* 10(4), 371-386.
- Dendoncker, N., Schmit, C., Rounsevell, M., 2008. Exploring spatial data uncertainties in land-use change scenarios. *International Journal of Geographical Information Science* 22(9), 1013-1030.
- Dias, M.O.S., Cunha, M.P., Jesus, C.D.F., Rocha, G.J.M., Pradella, J.G.C., Rossell, C.E.V., Maciel Filho, R., Bonomi, A., 2011. Second generation ethanol in Brazil: Can it compete with electricity production? *Bioresource technology* 102(19), 8964-8971.
- Diogo, V., van der Hilst, F., van Eijck, J., Faaij, A., Versteegen, J.A., Hilbert, J., Carballo, S., Volante, J., 2014. Combining empirical and theory-based land use modelling approaches to assess future availability of land and economic potential for sustainable biofuel production: Argentina as a case study. *Renewable & Sustainable Energy Reviews* 34, 208-224.
- Dornburg, V., van Vuuren, D.P., van de Ven, G.W.J., Langeveld, H., Meeusen, M., Banse, M., Van Oorschot, M., Ros, J., van den Born, G.J., Aiking, H., Londo, M., Mozaffarian, H., Verweij, P., Lysen, E., Faaij, A.P.C., 2010. Bioenergy revisited: Key factors in global potentials of bioenergy. *Energy and Environmental Science* 3(3), 258-267.
- Eckhardt, K., Breuer, L., Frede, H., 2003. Parameter uncertainty and the significance of simulated land use change effects. *Journal of Hydrology* 273, 164-176.
- EEX, 2015. European Emission Allowances. 2015.
- Efron, B., Tibshirani, R.J., 2003. *An introduction to the bootstrap*. CRC Press LLC, Boca Raton, Florida.
- Epron, D., Nouvellon, Y., Mareschal, L., Moreira, R.M.E., Koutika, L., Geneste, B., Delgado-Rojas, J.S., Laclau, J., Sola, G., Gonçalves, J.L.D.M., Bouillet, J., 2013. Partitioning of net primary production in Eucalyptus and Acacia stands and in mixed-species plantations: Two case-studies in contrasting tropical environments. *Forest Ecology and Management* 301, 102-111.
- ESRI, 2011. ArcGIS internet site, Available online at: <http://www.esri.com>. 2011.
- European Parliament and Council of the European Union, 2009. DIRECTIVE 2009/28/EC OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL of 23 April 2009 on the promotion of the use of energy from renewable sources and amending and subsequently repealing Directives 2001/77/EC and 2003/30/EC.
- Evans, T.P., Kelley, H., 2004. Multi-scale analysis of a household level agent-based model of landcover change. *Journal of Environmental Management* 72, 57-72.
- Evensen, G., 1994. Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics. *Journal of Geophysical Research* 99(C5), 10,143-10,162.
- Fang, S., Gertner, G., Wang, G., Anderson, A., 2006. The impact of misclassification in land use maps in the prediction of landscape dynamics. *Landscape Ecology* 21(2), 233-242.
- FAO, 2015. FAOSTAT, Food and Agriculture Organisation of the United Nations, Statistics Division, Available online at: <http://faostat3.fao.org/home/E>.
- FAO, 2003. *World agriculture towards 2015/2030: an FAO perspective*, in Bruinsma, J. (Ed.). Food and Agriculture Organisation, Rome, Italy.

- Fargione, J., Hill, J., Tilman, D., Polasky, S., Hawthorne, P., 2008. Land clearing and the biofuel carbon debt. *Science* 319(5867), 1235-1238.
- Fargione, J.E., Plevin, R.J., Hill, J.D., 2010. The ecological impact of biofuels. *Annual Review of Ecology, Evolution, and Systematics* 41, 351-377.
- Farr, T.G., Rosen, P.A., Caro, E., Crippen, R., Duren, R., Hensley, S., Kobrick, M., Paller, M., Rodriguez, E., Roth, L., Seal, D., Shaffer, S., Shimada, J., Umland, J., Werner, M., Oskin, M., Burbank, D., Alsdorf, D., 2007. The Shuttle Radar Topography Mission. *Reviews of Geophysics* 45, Article number RG2004.
- Feng, Y., Liu, Y., Tong, X., Liu, M., Deng, S., 2011. Modeling dynamic urban growth using cellular automata and particle swarm optimization rules. *Landscape and Urban Planning* 102(3), 188-196.
- Filatova, T., Polhill, J.G., 2012. Shocks in coupled socio-ecological systems: what are they and how can we model them? R. Seppelt, A. A. Voinov, S. Lange and D. Bankamp (Eds.), International Environmental Modelling and Software Society (iEMSs), 2012 International Congress on Environmental Modelling and Software, Managing Resources of a Limited Planet: Pathways and Visions under Uncertainty, Sixth Biennial Meeting, Leipzig, Germany, Leipzig, Germany, July 1-5, 2012, 2619-2630.
- Finkbeiner, M., 2014. Indirect land use change – Help beyond the hype? *Biomass and Bioenergy* 62, 218-221.
- FNP Informa economics, 2012. *Agriannual 2012*, São Paulo, Brazil.
- Foody, G.M., 2003. Uncertainty, knowledge discovery and data mining in GIS. *Progress in Physical Geography* 27(1), 113-121.
- Galford, G.L., Mustard, J.F., Melillo, J.M., Gendrin, A., Cerri, C.C., Cerri, C.E.P., 2008. Wavelet analysis of MODIS time series to detect expansion and intensification of row-crop agriculture in Brazil. *Remote Sensing of Environment* 112(2), 576-587.
- Gallardo, A.L.C.F., Bond, A., 2011. Capturing the implications of land use change in Brazil through environmental assessment: Time for a strategic approach? *Environmental Impact Assessment Review* 31(3), 261-270.
- Geertman, S., Stillwell, J., 2004. Planning support systems: An inventory of current practice. *Computers, Environment and Urban Systems* 28(4), 291-310.
- Gibbs, H.K., Johnston, M., Foley, J.A., Holloway, T., Monfreda, C., Ramankutty, N., Zaks, D., 2008. Carbon payback times for crop-based biofuel expansion in the tropics: The effects of changing yield and technology. *Environmental Research Letters* 3(3), Article number 034001.
- Gómez Jr., C., 2013. Rise of ethanol in Brazil? *SciTechDaily* .
- Goodchild, M.F., 2004. A general framework for error analysis in measurement-based GIS. *Journal of Geographical Systems* 6(4), 323-324.
- Goovaerts, P., 2010. Geostatistical Software, in Fischer, M.M., Getis, A. (Eds.), *Handbook of Applied Spatial Analysis*. Springer, Heidelberg, 129-138.
- Gorsevski, P.V., Gessler, P.E., Boll, J., Elliot, W.J., Foltz, R.B., 2006. Spatially and temporally distributed modeling of landslide susceptibility. *Geomorphology* 80, 178-198.
- Grazzini, J., 2012. Analysis of the Emergent Properties: Stationarity and Ergodicity. *Journal of Artificial Societies and Social Simulation* 15, March 12 2013.

- Grimm, V., Railsback, S.F., 2012. Pattern-oriented modelling: A 'multi-scope' for predictive systems ecology. *Philosophical Transactions of the Royal Society B: Biological Sciences* 367(1586), 298-310.
- Gurgel, H.C., Hargrave, J., França, F., Holmes, R.M., Ricarte, F.M., Dias, B.F.S., Rodrigues, C.G.O., de Brito, M.C.W., 2009. Unidades de conservação e o falso dilema entre conservação e desenvolvimento. *Boletim regional, urbano e ambiental* 3, 109-119.
- Haberl, H., Beringer, T., Bhattacharya, S.C., Erb, K.H., Hoogwijk, M., 2010. The global technical potential of bio-energy in 2050 considering sustainability constraints. *Current Opinion in Environmental Sustainability* 2, 394-403.
- Hamelinck, C.N., Suurs, R.A.A., Faaij, A.P.C., 2005a. International bioenergy transport costs and energy balance. *Biomass and Bioenergy* 29(2), 114-134.
- Hamelinck, C.N., Van Hooijdonk, G., Faaij, A.P.C., 2005b. Ethanol from lignocellulosic biomass: Techno-economic performance in short-, middle- and long-term. *Biomass and Bioenergy* 28(4), 384-410.
- Hansen, H.S., 2012. Empirically derived neighbourhood rules for urban land-use modelling. *Environment and Planning B: Planning and Design* 39(2), 213-228.
- Hartig, F., Calabrese, J.M., Reineking, B., Wiegand, T., Huth, A., 2011. Statistical inference for stochastic simulation models - theory and application. *Ecology Letters* 14(8), 816-827.
- Haupt, R.L., Haupt, S.E., 2004. The Binary Genetic Algorithm, in *Practical Genetic Algorithms*, 2nd ed. Wiley, Hoboken, New Jersey, 27-66.
- Hellmann, F., Verburg, P.H., 2011. Spatially explicit modelling of biofuel crops in Europe. *Biomass and Bioenergy* 35(6), 2411-2424.
- Hellmann, F., Verburg, P.H., 2010. Impact assessment of the European biofuel directive on land use and biodiversity. *Journal of Environmental Management* 91(6), 1389-1396.
- Hengl, T., 2006. Finding the right pixel size. *Computers and Geosciences* 32(9), 1283-1298.
- Hertel, T.W., 2011. The global supply and demand for agricultural land in 2050: A perfect storm in the making? *American Journal of Agricultural Economics* 93(2), 259-275.
- Heuvelink, G.B.M., 1998. Error propagation in environmental modelling with GIS. Taylor & Francis, London, United Kingdom.
- Hiemstra, P.H., Karssenbergh, D., van Dijk, A., de Jong, S.M., 2012. Using the particle filter for nuclear decision support. *Environmental Modelling and Software* 37, 78-89.
- Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., Jarvis, A., 2005. Very high resolution interpolated climate surfaces for global land areas. *International Journal of Climatology* 25(15), 1965-1978.
- Hoefnagels, R., Banse, M., Dornburg, V., Faaij, A., 2013. Macro-economic impact of large-scale deployment of biomass resources for energy and materials on a national level-A combined approach for the Netherlands. *Energy Policy* 59, 727-744.
- Holzämper, A., Klein, T., Seppelt, R., Fuhrer, J., 2015. Assessing the propagation of uncertainties in multi-objective optimization for agro-ecosystem adaptation to climate change. *Environmental Modelling & Software* 66, 27-35.
- Hoogwijk, M., Faaij, A.P.C., Eickhout, B., de Vries, Turkenburg, W.C., 2005. Potential of biomass energy out to 2100, for four IPCC SRES land-use scenarios. *Biomass and Bioenergy* 29(4), 225-257.

- Hyndman, R.J., Koehler, A.B., 2006. Another look at measures of forecast accuracy. *International Journal of Forecasting* 22(4), 679-688.
- IBGE, 2013a. Pesquisa Pecuária Municipal 2006-2012, Available online at: <http://www.sidra.ibge.gov.br/>.
- IBGE, 2013b. Produção Agrícola Municipal 2006-2012, Available online at: <http://www.sidra.ibge.gov.br/>.
- IBGE, 2006. Censo Agropecuário 2006, Available online at: <http://www.sidra.ibge.gov.br/>.
- ICONE, 2012. ICONE - Institute for International Trade Negotiations, Brief Description for the Brazilian Land Use Model - BLUM.
- IEA, 2013. World Energy Outlook 2013, International Energy Agency, Paris, France.
- INE, 2003. Censo Agro-Pecuário 1999-2000, resultados Temáticos, Instituto Nacional de Estatística, Maputo.
- IPCC, 2014. Climate Change 2014: Synthesis Report. Contribution of Working Groups I, II and III to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change, in Pachauri, R.K., Meyer, L.A. (Eds.). IPCC, Geneva, Switzerland.
- IPCC, 2011. IPCC Special Report on Renewable Energy Sources and Climate Change Mitigation, in Edenhofer, O., Pichs-Madruga, R., Sokona, Y., Seyboth, K., Matschoss, P., Kadner, S., Zwickel, T., Eickemeier, P., Hansen, G., Schlömer, S., von Stechow, C. (Eds.). Cambridge University Press, Cambridge, United Kingdom and New York, USA.
- IPCC, 2006. 2006 IPCC Guidelines for National Greenhouse Gas Inventories, in Eggleston, H.S., Buendia, L., Miwa, K., Ngara, T., Tanabe, K. (Eds.). The Intergovernmental Panel on Climate Change (IPCC), Hayama, Kanagawa, Japan.
- Ivanovic, R.F., Freer, J.E., 2009. Science versus politics: truth and uncertainty in predictive modelling. *Hydrological Processes* 23, 2549-2554.
- Jangpromma, N., Thammasirirak, S., Jailsil, P., Songsri, P., 2012. Effects of drought and recovery from drought stress on above ground and root growth, and water use efficiency in sugarcane (*Saccharum officinarum* L.). *Australian Journal of Crop Science* 6(8), 1298-1304.
- Jansen, L.J.M., Bagnoli, M., Focacci, M., 2008. Analysis of land-cover/use change dynamics in Manica Province in Mozambique in a period of transition (1990-2004). *Forest Ecology and Management* 254(2), 308-326.
- Jarvis, A., Reuter, H.I., Nelson, A., Guevara, E., 2008. Hole-filled seamless SRTM data V4, International Centre for Tropical Agriculture (CIAT), available from <http://srtm.csi.cgiar.org>.
- Jazwinski, A.H., 1970. Stochastic processes and filtering theory. Academic Press, New York.
- Jeremiah, E., Sisson, S.A., Sharma, A., Marshall, L., 2012. Efficient hydrological model parameter optimization with Sequential Monte Carlo sampling. *Environmental Modelling and Software* 38, 283-295.
- Johnson, B.R., 2010. Eliminating the mystery from the concept of emergence. *Biology and Philosophy* 25(5), 843-849.
- Jonker, J.G.G., van der Hilst, F., Junginger, H.M., Cavalett, O., Chagas, M.F., Faaij, A.P.C., 2015. Outlook for ethanol production costs in Brazil up to 2030, for different biomass crops and industrial technologies. *Applied Energy* 147, 593-610.

- Jonker, J.G.G., van der Hilst, F., Junginger, H.M., Versteegen, J.A., Lin, T., Rodríguez, L.F., Faaij, A.P.C., in prep. Supply chain optimization of sugarcane first generation and eucalyptus second generation ethanol production in Brazil.
- Karszenberg, D., de Jong, K., 2006. Towards improved solution schemes for Monte Carlo simulation in environmental modeling languages., in Oosterom, P.J.M. (Ed.), *Geo-Information and Computational Geometry*. NGC (Nederlandse Commissie voor Geodesie, in English: Netherlands Geodetic Commission), Delft.
- Karszenberg, D., de Jong, K., van der Kwast, J., 2007. Modelling landscape dynamics with Python. *International Journal of Geographical Information Science* 21(5), 483-495.
- Karszenberg, D., Schmitz, O., Salamon, P., de Jong, K., Bierkens, M.F.P., 2010. A software framework for construction of process-based stochastic spatio-temporal models and data assimilation. *Environmental Modelling and Software* 25, 489-502.
- Kavallari, A., Smeets, E., Tabeau, A., 2014. Land use changes from EU biofuel use: A sensitivity analysis. *Operational Research* 14(2), 261-281.
- Kéfi, S., Rietkerk, M., Alados, C.L., Pueyo, Y., Papanastasis, V.P., ElAich, A., De Ruiter, P.C., 2007. Spatial vegetation patterns and imminent desertification in Mediterranean arid ecosystems. *Nature* 449(7159), 213-217.
- Kocabas, V., Dragičević, S., 2007. Enhancing a GIS cellular automata model of land use change: Bayesian networks, influence diagrams and causality. *Transactions in GIS* 11(5), 681-702.
- Kok, K., Farrow, A., Veldkamp, A., Verburg, P.H., 2001. A method and application of multi-scale validation in spatial land use models. *Agriculture, Ecosystems and Environment* 85, 223-238.
- Koomen, E., Diogo, V., Dekkers, J., Rietveld, P., 2015. A utility-based suitability framework for integrated local-scale land-use modelling. *Computers, Environment and Urban Systems* 50, 1-14.
- Lambin, E.F., 2004. Modelling Land-Use Change, in Wainwright, J., Mulligan, M. (Eds.), *Environmental Modelling - Finding Simplicity in Complexity*, 1st ed. John Wiley & Sons, Chichester, West Sussex, England, 245-254.
- Lambin, E.F., Meyfroidt, P., 2011. Global land use change, economic globalization, and the looming land scarcity. *Proceedings of the National Academy of Sciences of the United States of America* 108(9), 3465-3472.
- Lambin, E.F., Rounsevell, M.D.A., Geist, H.J., 2000. Are agricultural land-use models able to predict changes in land-use intensity? *Agriculture, Ecosystems and Environment* 82, 321-331.
- Lambin, E.F., Turner, B.L., Geist, H.J., Agbola, S.B., Angelsen, A., Bruce, J.W., Coomes, O.T., Dirzo, R., Fischer, G., Folke, C., George, P.S., Homewood, K., Imbernon, J., Leemans, R., Li, X., Moran, E.F., Mortimore, M., Ramakrishnan, P.S., Richards, J.F., Skånes, H., Steffen, W., Stone, G.D., Svedin, U., Veldkamp, T.A., Vogel, C., Xu, J., 2001. The causes of land-use and land-cover change: Moving beyond the myths. *Global Environmental Change* 11(4), 261-269.
- Lanzante, J.R., 1996. Resistant, robust and non-parametric techniques for the analysis of climate data: Theory and examples, including applications to historical radiosonde station data. *International Journal of Climatology* 16(11), 1197-1226.

- Lapola, D.M., Schaldach, R., Alcamo, J., Bondeau, A., Koch, J., Koelking, C., Priess, J.A., 2010. Indirect land-use changes can overcome carbon savings from biofuels in Brazil. *Proceedings of the National Academy of Sciences of the United States of America* 107(8), 3388-3393.
- Lauf, S., Haase, D., Hostert, P., Lakes, T., Kleinschmit, B., 2012. Uncovering land-use dynamics driven by human decision-making - A combined model approach using cellular automata and system dynamics. *Environmental Modelling and Software* 27-28, 71-82.
- Lautenbach, S., Volk, M., Strauch, M., Whittaker, G., Seppelt, R., 2013. Optimization-based trade-off analysis of biodiesel crop production for managing an agricultural catchment. *Environmental Modelling and Software* 48, 98-112.
- Lei, Z., Pijanowski, B.C., Alexandridis, K.T., Olson, J., 2005. Distributed modeling architecture of a multi-agent-based behavioral economic landscape (MABEL) model. *Simulation* 81(7), 503-515.
- Li, X., Yeh, A.G.-, 2005. Integration of genetic algorithms and GIS for optimal location search. *International Journal of Geographical Information Science* 19(5), 581-601.
- Li, X., Yeh, A.G.-, 2002. Neural-network-based cellular automata for simulating multiple land use changes using GIS. *International Journal of Geographical Information Science* 16(4), 323-343.
- Ligmann-Zielinska, A., Sun, L., 2010. Applying time-dependent variance-based global sensitivity analysis to represent the dynamics of an agent-based model of land use change. *International Journal of Geographical Information Science* 24(12), 1829-1850.
- Ligtenberg, A., Beulens, A., Kettenis, D., Bregt, A.K., Wachowicz, M., 2009. Simulating knowledge sharing in spatial planning: an agent-based approach. *Environment and Planning B: Planning and Design* 36(4), 644-663.
- Lillesand, T.M., Kiefer, R.W., Chipman, J.W., 2003. *Remote sensing and image interpretation*, 5th ed. Wiley.
- Liu, X., Li, X., Liu, L., He, J., Ai, B., 2008. A bottom-up approach to discover transition rules of cellular automata using ant intelligence. *International Journal of Geographical Information Science* 22, 1247-1269.
- Lokupitiya, E., Breidt, F.J., Lokupitiya, R., Williams, S., Paustian, K., 2007. Deriving Comprehensive County-Level Crop Yield and Area Data for U.S. Cropland. *Agronomy Journal* 99(3), 673-681.
- Lucas Jr., R.E., 1976. *Econometric policy evaluation: A critique*. American Elsevier (Ed.), Carnegie-Rochester Confer. Series on Public Policy, New York, 1, 19-46.
- Lucon, O., Goldemberg, J., 2010. São Paulo - The "other" Brazil: Different pathways on climate change for state and federal governments. *Journal of Environment and Development* 19(3), 335-357.
- Ma, L., Arentze, T., Borgers, A., Timmermans, H., 2007. Modelling land-use decisions under conditions of uncertainty. *Computers, Environment and Urban Systems* 31(4), 461-476.
- Macedo, I.C., 2007. *Sugar Cane's Energy - Twelve studies on Brazilian sugar cane agribusiness and its sustainability*.
- Macedo, I.C., Leal, M.R.L.V., Silva, J.E.A.R., 2004. *Assessment of greenhouse gas emissions in the production and use of fuel ethanol in Brazil*, Government of the State of São Paulo, São Paulo, Brazil.

- Macedo, I.C., Seabra, J.E.A., 2008. Chapter 4: Mitigation of GHG emissions using sugarcane bioethanol, in Zuurbier, P., van de Vooren, J. (Eds.), *Sugarcane Ethanol: Contributions to Climate Change Mitigation and the Environment*. Wageningen Academic Publishers, Wageningen, The Netherlands, 95-110.
- Macedo, I.C., Seabra, J.E.A., Silva, J.E.A.R., 2008. Green house gases emissions in the production and use of ethanol from sugarcane in Brazil: The 2005/2006 averages and a prediction for 2020. *Biomass and Bioenergy* 32(7), 582-595.
- Magliocca, N.R., van Vliet, J., Brown, C., Evans, T.P., Houet, T., Messerli, P., Messina, J.P., Nicholas, K.A., Ornetsmüller, C., Sagebiel, J., Schweizer, V., Verburg, P.H., Yu, Q., 2015. From meta-studies to modeling: Using synthesis knowledge to build broadly applicable process-based land change models. *Environmental Modelling & Software* 72, 10-20.
- Malins, C., 2013. A model-based quantitative assessment of the carbon benefits of introducing iLUC factors in the European Renewable Energy Directive. *Global Change Biology Bioenergy* 5(6), 639-651.
- Mancosu, E., Gago-Silva, A., Barbosa, A., de Bono, A., Ivanov, E., Lehmann, A., Fons, J., 2015. Future land-use change scenarios for the Black Sea catchment. *Environmental Science and Policy* 46, 26-36.
- Manson, S.M., 2007. Challenges in evaluating models of geographic complexity. *Environment and Planning B: Planning and Design* 34, 245-260.
- Mathews, J.A., Tan, H., 2009. Biofuels and indirect land use change effects: The debate continues. *Biofuels, Bioproducts and Biorefining* 3(3), 305-317.
- Matthews, R.B., Gilbert, N.G., Roach, A., Polhill, J.G., Gotts, N.M., 2007. Agent-based land-use models: A review of applications. *Landscape Ecology* 22(10), 1447-1459.
- Merz, R., Parajka, J., Blöschl, G., 2011. Time stability of catchment model parameters: Implications for climate impact analyses. *Water Resources Research* 47(2).
- Meyfroidt, P., Lambin, E.F., Erb, K.-., Hertel, T.W., 2013. Globalization of land use: Distant drivers of land change and geographic displacement of land use. *Current Opinion in Environmental Sustainability* 5(5), 438-444.
- MMA, 2008. Cobertura Vegetal do Bioma Amazônia, Cobertura Vegetal do Bioma Caatinga, Cobertura Vegetal do Bioma Cerrado, Cobertura Vegetal do Bioma Mata Atlântica, Cobertura Vegetal do Bioma Pampa, Cobertura Vegetal do Bioma Pantanal.
- MME, 2013. Plano decenal de expansão de energia 2021, MME, Brasília, Brazil.
- Monfreda, C., Ramankutty, N., Foley, J.A., 2008. Farming the planet: 2. Geographic distribution of crop areas, yields, physiological types, and net primary production in the year 2000. *Global Biogeochemical Cycles* 22(1).
- Moulds, S., Buytaert, W., Mijic, A., 2015. An open and extensible framework for spatially explicit land use change modelling: The lulcc R package. *Geoscientific Model Development* 8(10), 3215-3229.
- Narayanan, G., Badri, A.A., McDougall, R., 2012. *Global Trade, Assistance, and Production: The GTAP 8 Data Base*.
- Nassar, A.M., Antoniazzi, L.B., Moreira, M.R., Chiodi, L., Harfuch, L., 2010. *An Allocation Methodology to Assess GHG Emissions Associated with Land Use Change*, Institute for International Trade Negotiations.
- Nassar, A.M., Rudorff, B.F.T., Antoniazzi, L.B., Aguiar, D.A., Bacchi, M.R.P., Adami, M., 2008. Chapter 3: Prospects of the sugarcane expansion in Brazil: impacts on direct and



- indirect land use changes, in Zuurbier, P., van de Vooren, J. (Eds.), *Sugarcane Ethanol: Contributions to Climate Change Mitigation and the Environment*. Wageningen Academic Publishers, Wageningen, 63-112.
- Nelson, G.C., van der Mensbrugge, D., Ahammad, H., Blanc, E., Calvin, K., Hasegawa, T., Havlik, P., Heyhoe, E., Kyle, P., Lotze-Campen, H., von Lampe, M., Mason d'Croz, D., van Meijl, H., Müller, C., Reilly, J., Robertson, R., Sands, R.D., Schmitz, C., Tabeau, A., Takahashi, K., Valin, H., Willenbockel, D., 2014. Agriculture and climate change in global scenarios: Why don't the models agree. *Agricultural Economics (United Kingdom)* 45(1), 85-101.
- Nesheim, I., Reidsma, P., Bezlepkina, I., Verburg, R., Abdeladhim, M.A., Bursztyn, M., Chen, L., Cissé, Y., Feng, S., Gicheru, P., Hannes, J.K., Novira, N., Purushothaman, S., Rodrigues-Filho, S., Sghaier, M., 2014. Causal chains, policy trade offs and sustainability: Analysing land (mis)use in seven countries in the South. *Land Use Policy* 37, 60-70.
- NetLogo, 2010. NetLogo software, available online at: [Http://ccl.northwestern.edu/netlogo/](http://ccl.northwestern.edu/netlogo/).
- Nguyen, T.G., de Kok, J.L., Titus, M.J., 2007. A new approach to testing an integrated water systems model using qualitative scenarios. *Environmental Modelling and Software* 22(11), 1557-1571.
- O'Brien, M., Schütz, H., Bringezu, S., 2015. The land footprint of the EU bioeconomy: Monitoring tools, gaps and needs. *Land Use Policy* 47, 235-246.
- OECD/Food and Agriculture Organization of the United Nations, 2014. *OECD-FAO Agricultural Outlook 2014*, OECD Publishing, Paris, France.
- OECD/IEA, 2014. *Renewable Energy 2014. Medium-Term Market Report. Market Analysis and Forecasts to 2020*, International Energy Agency, Paris, France.
- O'Hagan, A., 2012. Probabilistic uncertainty specification: Overview, elaboration techniques and their application to a mechanistic model of carbon flux. *Environmental Modelling and Software* 36, 35-48.
- O'Hare, M., Delucchi, M., Edwards, R., Fritsche, U., Gibbs, H., Hertel, T., Hill, J., Kammen, D., Laborde, D., Marelli, L., Mulligan, D., Plevin, R., Tyner, W., 2011. Comment on "Indirect land use change for biofuels: Testing predictions and improving analytical methodologies" by Kim and Dale: Statistical reliability and the definition of the indirect land use change (iLUC) issue. *Biomass and Bioenergy* 35(10), 4485-4487.
- O'Neill, B.C., Krieglner, E., Riahi, K., Ebi, K.L., Hallegatte, S., Carter, T.R., Mathur, R., van Vuuren, D.P., 2014. A new scenario framework for climate change research: The concept of shared socioeconomic pathways. *Climatic Change* 122(3), 387-400.
- Oreskes, N., Shrader-Frechette, K., Belitz, K., 1994. Verification, validation, and confirmation of numerical models in the earth sciences. *Science* 263(5147), 641-646.
- Overmars, K.P., Stehfest, E., Ros, J.P.M., Prins, A.G., 2011. Indirect land use change emissions related to EU biofuel consumption: An analysis based on historical data. *Environmental Science and Policy* 14(3), 248-257.
- Padua Junior, A.L., Costa Pasini, A.C., Comatsu, C.E., Casarin, D.C.P., Michelino, G.G., von Glhen, H.C., da Silva, I.X., de Moreas, J.F.L., de Carvalho, J.P., Sandoval, M., Valeriano, M., Araujo, N., Brunini, O., Vedovelo, R., Viegas, R., Campaign, R.C.F., Adami, S.F., 2012.

Agro-environmental Zoning - Green Ethanol - Environmental System for São Paulo - Government of São Paulo, Available online at: <http://www.ambiente.sp.gov.br/etanolverde/zonamento-agroambiental/>.

- Page, S.E., 2011. Diversity and complexity. Princeton University Press, Princeton, USA.
- Pan, Y., Roth, A., Yu, Z., Doluschitz, R., 2010. The impact of variation in scale on the behavior of a cellular automata used for land use change modeling. *Computers, Environment and Urban Systems* 34(5), 400-408.
- Parker, D.C., Hessel, A., Davis, S.C., 2008. Complexity, land-use modeling, and the human dimension: Fundamental challenges for mapping unknown outcome spaces. *Geoforum* 39(2), 789-804.
- Parker, D.C., Manson, S.M., Janssen, M.A., Hoffmann, M.J., Deadman, P., 2003. Multi-Agent Systems for the Simulation of Land-Use and Land-Cover Change: A Review. *Annals of the Association of American Geographers* 93(2), 314-337.
- Pasetto, D., Camporese, M., Putti, M., 2012. Ensemble Kalman filter versus particle filter for a physically-based coupled surface-subsurface model. *Advances in Water Resources* 47, 1-13.
- PCRaster, 2010. PCRaster software, available online at: <Http://pcraster.geo.uu.nl>.
- Pebesma, E.J., 2004. Multivariable geostatistics in S: the gstat package. *Computers & Geosciences* 30(7), 683-691.
- Pebesma, E.J., de Jong, K., Briggs, D., 2007. Interactive visualization of uncertain spatial and spatio-temporal data under different scenarios: An air quality example. *International Journal of Geographical Information Science* 21(5), 515-527.
- Pebesma, E.J., Wesseling, C.G., 1998. Gstat: A program for geostatistical modelling, prediction and simulation. *Computers and Geosciences* 24(1), 17-31.
- Piao, S., Friedlingstein, P., Ciais, P., De Noblet-Ducoudré, N., Labat, D., Zaehle, S., 2007. Changes in climate and land use have a larger direct impact than rising CO<sub>2</sub> on global river runoff trends. *Proceedings of the National Academy of Sciences of the United States of America* 104(39), 15242-15247.
- Picoli, M., 2013. Brazilian Sugarcane Mills Map - 2013.
- Pijanowski, B.C., Alexandridis, K.T., Müller, D., 2006. Modelling urbanization patterns in two diverse regions of the world. *Journal of Land Use Science* 1, 83-108.
- Pijanowski, B.C., Brown, D.G., Shellito, B.A., Manik, G.A., 2002. Using neural networks and GIS to forecast land use changes: A Land Transformation Model. *Computers, Environment and Urban Systems* 26(6), 553-575.
- Plevin, R.J., Beckman, J., Golub, A.A., Witcover, J., O'Hare, M., 2015. Carbon accounting and economic model uncertainty of emissions from biofuels-induced land use change. *Environmental Science and Technology* 49(5), 2656-2664.
- Pontius Jr., R.G., Neeti, N., 2010. Uncertainty in the difference between maps of future land change scenarios. *Sustainability Science* 5(1), 39-50.
- Pontius Jr., R.G., Spencer, J., 2005. Uncertainty in extrapolations of predictive land-change models. *Environment and Planning B: Planning and Design* 32(2), 211-230.
- Popp, J., Lakner, Z., Harangi-Rákos, M., Fári, M., 2014. The effect of bioenergy expansion: Food, energy, and environment. *Renewable and Sustainable Energy Reviews* 32, 559-578.

- Python software foundation, 2014. Python language reference, version 2.7, Available online at: [Http://www.python.org](http://www.python.org).
- Ramankutty, N., Foley, J.A., 1999. Estimating historical changes in global land cover: Croplands from 1700 to 1992. *Global Biogeochemical Cycles* 13(4), 997-1027.
- Rasmussen, R., Hamilton, G., 2012. An approximate bayesian computation approach for estimating parameters of complex environmental processes in a cellular automata. *Environmental Modelling and Software* 29(1), 1-10.
- Ray, D.K., Pijanowski, B.C., Kendall, A.D., Hyndman, D.W., 2012. Coupling land use and groundwater models to map land use legacies: Assessment of model uncertainties relevant to land use planning. *Applied Geography* 34(3), 356-370.
- Refsgaard, J.C., van der Sluijs, J.P., Brown, J., van der Keur, P., 2006. A framework for dealing with uncertainty due to model structure error. *Advances in Water Resources* 29(11), 1586-1597.
- Robinson, V.B., 2003. A perspective on the fundamentals of fuzzy sets and their use in geographic information systems. *Transactions in GIS* 7(1), 3-30.
- Rosa, I.M.D., Ahmed, S.E., Ewers, R.M., 2014. The transparency, reliability and utility of tropical rainforest land-use and land-cover change models. *Global Change Biology* 20(6), 1707-1722.
- Rose, A., 1995. Input-output economics and computable general equilibrium models. *Structural Change and Economic Dynamics* 6(3), 295-304.
- Roundtable on Sustainable Biofuels, 2010. RSB Principles & Criteria for Sustainable Biofuel Production.
- Rudorff, B.F.T., Adami, M., Aguiar, D.A., Moreira, M.A., Mello, M.P., Fabiani, L., Amaral, D.F., Pires, B.M., 2011. The soy moratorium in the Amazon biome monitored by remote sensing images. *Remote Sensing* 3(1), 185-202.
- Rudorff, B.F.T., Aguiar, D.A., Silva, W.F., Sugawara, L.M., Adami, M., Moreira, M.A., 2010. Studies on the Rapid Expansion of Sugarcane for Ethanol Production in São Paulo State (Brazil) Using Landsat Data. *Remote Sensing* 2(4), 1057-1076.
- Rykiel Jr., E.J., 1996. Testing ecological models: The meaning of validation. *Ecological Modelling* 90(3), 229-244.
- Sala, O.E., Chapin III, F.S., Armesto, J.J., Berlow, E., Bloomfield, J., Dirzo, R., Huber-Sanwald, E., Huenneke, L.F., Jackson, R.B., Kinzig, A., Leemans, R., Lodge, D.M., Mooney, H.A., Oesterheld, M., Poff, N.L., Sykes, M.T., Walker, B.H., Walker, M., Wall, D.H., 2000. Global biodiversity scenarios for the year 2100. *Science* 287(5459), 1770-1774.
- Salamon, P., Feyen, L., 2009. Assessing parameter, precipitation, and predictive uncertainty in a distributed hydrological model using sequential data assimilation with the particle filter. *Journal of Hydrology* 376, 428-442.
- Santé, I., García, A.M., Miranda, D., Crecente, R., 2010. Cellular automata models for the simulation of real-world urban processes: A review and analysis. *Landscape and Urban Planning* 96(2), 108-122.
- Scanlon, B.R., Jolly, I., Sophocleous, M., Zhang, L., 2007. Global impacts of conversions from natural to agricultural ecosystems on water resources: Quantity versus quality. *Water Resources Research* 43(3).

- Schaldach, R., Alcamo, J., Koch, J., Kölling, C., Lapola, D.M., Schüngel, J., Priess, J.A., 2011. An integrated approach to modelling land-use change on continental and global scales. *Environmental Modelling and Software* 26(8), 1041-1051.
- Schmitz, O., Salvadore, E., Poelmans, L., van der Kwast, J., Karssenber, D., 2014. A framework to resolve spatio-temporal misalignment in component-based modelling. *Journal of Hydroinformatics* 16(4), 850-871.
- Schroeder, J.P., 2007. Target-Density Weighting Interpolation and Uncertainty Evaluation for Temporal Analysis of Census Data. *Geographical Analysis* 39(3), 311-335.
- Seabra, J.E.A., Macedo, I.C., Chum, H.L., Faroni, C.E., Sarto, C.A., 2011. Life cycle assessment of Brazilian sugarcane products: GHG emissions and energy use. *Biofuels, Bioproducts and Biorefining* 5(5), 519-532.
- Seabra, J.E.A., Tao, L., Chum, H.L., Macedo, I.C., 2010. A techno-economic evaluation of the effects of centralized cellulosic ethanol and co-products refinery options with sugarcane mill clustering. *Biomass and Bioenergy* 34(8), 1065-1078.
- Searchinger, T.D., Heimlich, R., Houghton, R.A., Dong, F., Elobeid, A., Fabiosa, J., Tokgoz, S., Hayes, D., Yu, T.-., 2008. Use of U.S. croplands for biofuels increases greenhouse gases through emissions from land-use change. *Science* 319(5867), 1238-1240.
- Seppelt, R., Lautenbach, S., Volk, M., 2013. Identifying trade-offs between ecosystem services, land use, and biodiversity: A plea for combining scenario analysis and optimization on different spatial scales. *Current Opinion in Environmental Sustainability* 5(5), 458-463.
- Smeets, E.M.W., Faaij, A.P.C., Lewandowski, I.M., Turkenburg, W.C., 2007. A bottom-up assessment and review of global bio-energy potentials to 2050. *Progress in Energy and Combustion Science* 33(1), 56-106.
- Smith, P., Martino, D., Cai, Z., Gwary, D., Janzen, H., Kumar, P., McCarl, B., Ogle, S., O'Mara, F., Rice, C., Scholes, B., Sirotenko, O., Howden, M., McAllister, T., Pan, G., Romanenkov, V., Schneider, U., Towprayoon, S., Wattenbach, M., Smith, J., 2008. Greenhouse gas mitigation in agriculture. *Philosophical Transactions of the Royal Society B: Biological Sciences* 363(1492), 789-813.
- Sorda, G., Banse, M., Kemfert, C., 2010. An overview of biofuel policies across the world. *Energy Policy* 38(11), 6977-6988.
- Sparovek, G., Barretto, A.G.d.O.P., Berndes, G., Martins, S., Maule, R., 2009. Environmental, land-use and economic implications of Brazilian sugarcane expansion 1996-2006. *Mitigation and Adaptation Strategies for Global Change* 14(3), 285-298.
- Sparovek, G., Berndes, G., Barretto, A.G.d.O.P., Klug, I.L.F., 2012. The revision of the Brazilian Forest Act: increased deforestation or a historic step towards balancing agricultural development and nature conservation? *Environmental Science & Policy* 16, 65-72.
- Sparovek, G., Berndes, G., Egeskog, A., de Freitas, F.L.M., Gustafsson, S., Hansson, J., 2007. Sugarcane ethanol production in Brazil: An expansion model sensitive to socioeconomic and environmental concerns. *Biofuels, Bioproducts and Biorefining* 1(4), 270-282.
- Sparovek, G., Berndes, G., Klug, I.L.F., Barretto, A.G.d.O.P., 2010. Brazilian agriculture and environmental legislation: Status and future challenges. *Environmental Science and Technology* 44(16), 6046-6053.

- Spiller, E.T., Budhiraja, A., Ide, K., Jones, C.K.R.T., 2008. Modified particle filter methods for assimilating Lagrangian data into a point-vortex model. *Physica D: Nonlinear Phenomena* 237, 1498-1506.
- Stella, 2010. Stella software, available online at: [Http://www.iseesystems.com](http://www.iseesystems.com).
- Stewart, T.J., Janssen, R., van Herwijnen, M., 2004. A genetic algorithm approach to multiobjective land use planning. *Computers and Operations Research* 31(14), 2293-2313.
- Stocker, B.D., Roth, R., Joos, F., Spahni, R., Steinacher, M., Zaehle, S., Bouwman, L., Xu-Ri, Prentice, I.C., 2013. Multiple greenhouse-gas feedbacks from the land biosphere under future climate change scenarios. *Nature Climate Change* 3(7), 666-672.
- Straatman, B., White, R., Engelen, G., 2004. Towards an automatic calibration procedure for constrained cellular automata. *Computers, Environment and Urban Systems* 28, 149-170.
- Torquato, S.A., 2006. Cana-de-açúcar para indústria: o quanto vai precisar crescer. *Análises e Indicadores do Agronegócio* 1(10), online version.
- Tóth, G., Kozłowski, B., Prieler, S., Wiberg, D., 2012. Global Agro-ecological Zones (GAEZ v3.0), IIASA and FAO, IIASA, Laxenburg, Austria and FAO, Rome, Italy.
- UFG, 2015. Laboratório de Processamento de Imagens e Geoprocessamento (LAPIG) do Instituto de Estudos Sócio-Ambientais (IESA) da Universidade Federal de Goiás (UFG), Available online at: <http://www.lapig.iesa.ufg.br/lapig/>. 2013.
- UNDP, 2008. World Population Prospects: The 2008 Revision.
- UNICA, 2015. UNICADATA, cane harvest reports, production history, prices & quotes, exports & imports, fuel consumption, auto sales & fleet size, Available online at: [www.unicadata.com.br](http://www.unicadata.com.br). 2015.
- Uusitalo, L., Lehtikoinen, A., Helle, I., Myrberg, K., 2015. An overview of methods to evaluate uncertainty of deterministic models in decision support. *Environmental Modelling & Software* 63, 24-31.
- van Delden, H., Stuczynski, T., Ciaian, P., Paracchini, M.L., Hurkens, J., Lopatka, A., Shi, Y.-., Prieto, O.G., Calvo, S., van Vliet, J., Vanhout, R., 2010. Integrated assessment of agricultural policies with dynamic land use change modelling. *Ecological Modelling* 221(18), 2153-2166.
- van der Hilst, F., Dornburg, V., Sanders, J.P.M., Elbersen, B., Graves, A., Turkenburg, W.C., Elbersen, H.W., van Dam, J.M.C., Faaij, A.P.C., 2010. Potential, spatial distribution and economic performance of regional biomass chains: The North of the Netherlands as example. *Agricultural Systems* 103(7), 403-417.
- van der Hilst, F., Faaij, A.P.C., 2012. Spatiotemporal cost-supply curves for bioenergy production in Mozambique. *Biofuels, Bioproducts and Biorefining* 6(4), 405-430.
- van der Hilst, F., Verstegen, J.A., Karssenbergh, D., Faaij, A.P.C., 2012. Spatio-temporal land use modelling to assess land availability for energy crops - illustrated for Mozambique. *Global Change Biology Bioenergy* 4(6), 859-874.
- van der Hilst, F., Verstegen, J.A., Woltjer, G., Smeets, E.M.W., Faaij, A.P.C., in prep. Quantification and allocation of land use change resulting from the expansion of biofuel production.
- van der Hilst, F., Verstegen, J.A., Zheliezna, T., Drozdova, O., Faaij, A.P., 2014. Integrated spatiotemporal modelling of bioenergy production potentials, agricultural land use,

- and related GHG balances; demonstrated for Ukraine. *Biofuels, Bioproducts and Biorefining* 8(3), 391-411.
- van der Kwast, J., Canters, F., Karssenbergh, D., Engelen, G., Van De Voorde, T., Uljee, I., De Jong, K., 2011. Remote sensing data assimilation in modeling urban dynamics: Objectives and methodology. A. Stein, E. J. Pebesma and G. B. M. Heuvelink (Eds.), *Procedia Environmental Sciences*, Enschede, 23-25 March 2011, 7: *Spatial Statistics 2011: Mapping Global Change*, 140-145.
- van Leeuwen, P.J., 2009. Particle filtering in geophysical systems. *Monthly Weather Review* 137(12), 4089-4114.
- van Leeuwen, P.J., 2003. A variance-minimizing filter for large-scale applications. *Monthly Weather Review* 131(9), 2071-2084.
- van Vliet, J., Bregt, A.K., Hagen-Zanker, A., 2011. Revisiting Kappa to account for change in the accuracy assessment of land-use change models. *Ecological Modelling* 222(8), 1367-1375.
- van Vliet, J., Hurkens, J., White, R., van Delden, H., 2012. An activity-based cellular automaton model to simulate land-use dynamics. *Environment and Planning B: Planning and Design* 39(2), 198-212.
- Veldkamp, A., Lambin, E.F., 2001. Editorial: Predicting land-use change. *Agriculture, Ecosystems and Environment* 85, 1-6.
- Verburg, P.H., 2006. Simulating feedbacks in land use and land cover change models. *Landscape Ecology* 21(8), 1171-1183.
- Verburg, P.H., De Koning, G.H.J., Kok, K., Veldkamp, A., Bouma, J., 1999. A spatial explicit allocation procedure for modelling the pattern of land use change based upon actual land use. *Ecological Modelling* 116(1), 45-61.
- Verburg, P.H., Mertz, O., Erb, K.-., Haberl, H., Wu, W., 2013. Land system change and food security: Towards multi-scale land system solutions. *Current Opinion in Environmental Sustainability* 5(5), 494-502.
- Verburg, P.H., Overmars, K.P., 2009. Combining top-down and bottom-up dynamics in land use modeling: Exploring the future of abandoned farmlands in Europe with the Dyna-CLUE model. *Landscape Ecology* 24(9), 1167-1181.
- Verburg, P.H., Schot, P.P., Dijst, M.J., Veldkamp, A., 2004. Land use change modelling: Current practice and research priorities. *GeoJournal* 61(4), 309-324.
- Verburg, P.H., Soepboer, W., Veldkamp, A., Limpiada, R., Espaldon, V., Mastura, S.S.A., 2002. Modeling the spatial dynamics of regional land use: The CLUE-S model. *Environmental management* 30(3), 391-405.
- von Braun, J., 2008. Rising food prices: What should be done? *EuroChoices* 7(2), 30-35.
- von Thünen, J.H., 1966. *Der isolierte staat in beziehung auf landwirtschaft un nationalökonomie*. Fischer, Stuttgart.
- Wald, A., Wolfowitz, J., 1940. On a test whether two samples are from the same population. *The Annals of Mathematical Statistics* 11(2), 147-162.
- Walter, A., Dolzan, P., Quilodrán, O., De Oliveira, J.G., Da Silva, C., Piacente, F., Segerstedt, A., 2011. Sustainability assessment of bio-ethanol production in Brazil considering land use change, GHG emissions and socio-economic aspects. *Energy Policy* 39(10), 5703-5716.

- Walter, A., Galdos, M.V., Scarpore, F.V., Leal, M.R.L.V., Seabra, J.E.A., da Cunha, M.P., Picoli, M.C.A., de Oliveira, C.O.F., 2014. Brazilian sugarcane ethanol: Developments so far and challenges for the future. *Wiley Interdisciplinary Reviews: Energy and Environment* 3(1), 70-92.
- Warner, E., Inman, D., Kunstman, B., Bush, B., Vimmerstedt, L., Peterson, S., Macknick, J., Zhang, Y., 2013. Modeling biofuel expansion effects on land use change dynamics. *Environmental Research Letters* 8(1).
- Warner, E., Zhang, Y., Inman, D., Heath, G., 2014. Challenges in the estimation of greenhouse gas emissions from biofuel-induced global land-use change. *Biofuels, Bioproducts and Biorefining* 8(1), 114-125.
- Watson, H.K., 2011. Potential to expand sustainable bioenergy from sugarcane in southern Africa. *Energy Policy* 39(10), 5746-5750.
- White, R., Engelen, G., 2000. High-resolution integrated modelling of the spatial dynamics of urban and regional systems. *Computers, Environment and Urban Systems* 24(5), 383-400.
- Wicke, B., van der Hilst, F., Daioglou, V., Banse, M., Beringer, T., Gerssen-Gondelach, S., Heijnen, S., Karssenberg, D., Laborde, D., Lippe, M., van Meijl, H., Nassar, A., Powell, J., Prins, A.G., Rose, S.N.K., Smeets, E.M.W., Stehfest, E., Tyner, W.E., Versteegen, J.A., Valin, H., van Vuuren, D.P., Yeh, S., Faaij, A.P.C., 2015. Model collaboration for the improved assessment of biomass supply, demand, and impacts. *Global Change Biology Bioenergy* 7, 422-437.
- Wicke, B., Verweij, P., Van Meijl, H., Van Vuuren, D.P., Faaij, A.P.C., 2012. Indirect land use change: Review of existing models and strategies for mitigation. *Biofuels* 3(1), 87-100.
- Wolfram, S., 1984. Cellular automata as models of complexity. *Nature* 311(5985), 419-424.
- Woltjer, G.B., 2013. Forestry in MAGNET: a new approach for land use and forestry modelling, Statutory Research Tasks Unit for Nature & the Environment (WOT Natuur & Milieu), Wageningen, The Netherlands.
- Woltjer, G.B., Kuiper, M.H., 2014. The MAGNET Model: Module description, LEI Wageningen UR (University & Research centre), Wageningen, The Netherlands.
- World Bank, 2015. World DataBank, Health Nutrition and Population Statistics: Population estimates and projections, Available online at: <http://databank.worldbank.org/data/reports.aspx?source=health-nutrition-and-population-statistics>. 2015.
- You, L., Wood, S., 2005. Assessing the spatial distribution of crop areas using a cross-entropy method. *International Journal of Applied Earth Observation and Geoinformation* 7(4), 310-323.
- Yu, J., Chen, Y., Wu, J., Khan, S., 2011. Cellular automata-based spatial multi-criteria land suitability simulation for irrigated agriculture. *International Journal of Geographical Information Science* 25(1), 131-148.
- Zhang, Y., Li, X., Liu, X., Qiao, J., 2011. The CA model based on data assimilation. *Yaogan Xuebao - Journal of Remote Sensing* 15(3), 475-482.





## Dankwoord

“He hallo! / Kijk mij eens het gras zien groeien / En iedereen maar denken da'k niks doe.”<sup>8</sup> Veel mensen in mijn omgeving hebben me de laatste jaren gevraagd wanneer ik nou klaar ben met studeren en eindelijk ga werken. Zelf had ik niet het idee dat ik niks deed, een antwoord dat nooit kon overtuigen. Nu kan ik dan toch zeggen dat ik klaar ben. Maar ik heb het gras niet in m'n eentje gegroeid. Degenen die daaraan de hebben bijgedragen wil ik hier bedanken.

“Maar niets minder waar / Want al ben ik dan niet klaar / Ik ben de wedstrijd met de Beatles en de buren moe.” Derek, zonder jou had ik de wedstrijd niet uitgespeeld. Gelukkig vroeg jij me heel nuchter wat dan mijn eisen waren om door te kunnen gaan. Met je grote conceptuele denkvermogen, integriteit en precisie ben je de beste begeleider die ik me had kunnen wensen. Je gedrevenheid en humor maken dat ik graag met je samenwerk en hoop dat we dat in de toekomst nog eens kunnen doen. Floor, ik vind het bewonderenswaardig hoe jij je door de, door Derek en mij opgezette en vaak ‘ver van jouw bed’, teksten heen werkte en daar zeer zinvolle feedback op kon geven. Aan de ene kant leerde je me wat voor soorten gras je allemaal kan zien groeien in de bio-energie wereld (vingergras, suikerriet, ...), en aan de andere kant liet je me helemaal mijn eigen weg gaan. Super combi! En onze reisjes, of het nou naar den Haag was of naar de VS, leverden altijd interessante situaties op... André, jij had vaak minder tijd voor die technische teksten en schreef er ook schitterend eerlijk bij ‘niet gelezen, te wiskundig’. Voorts, ik weet dat alle PhD studenten dit over je zeggen, maar het is onontkoombaar: je enthousiasme is aanstekelijk en onbeschrijfelijk waardevol in een promotietraject. Steven, jij kwam bij mijn begeleidingsteam toen André weg ging uit Utrecht. Je hebt inhoudelijk dus minder bijgedragen aan het proefschrift, maar het gaf me gemoedsrust om een achterevoerder te hebben. En ook in andere rollen was je waardevol: als stok achter de deur voor mijn BKO, en als de soigneur die mij op veldwerk ontbijt kwam brengen toen ik een kapot was omdat ik dacht dat het slim was in één ruk door van Dallas via Amsterdam naar Savournon te reizen.

De werksfeer wordt naast het begeleidingsteam ook bepaald door collega's. Allereerst dank ik iedereen van het Copernicus instituut. Vooral in het laatste jaar, toen jullie mij gevonden hadden als 'GISser', heb ik genoten van de samenwerking. Ik dank ook het gehele, voormalige en huidige, BE-Basic team. Speciale bedankjes zijn er voor mijn kamergenoten. Gert-Jan, zonder onze gezamenlijke evaluatie van het voetbalweekend kon mijn werkweek niet beginnen. En Anne Sjoerd, wetend

---

<sup>8</sup> Dit citaat komt uit 'De Beatles en de buren' van Acda en de Munnik, evenals alle volgende citaten in dit dankwoord (samen het hele nummer), tenzij anders aangegeven.

dat er altijd iemand met een brede lach achter m'n rug zat, kon ik nooit lang in een tegenslag blijven hangen. In het bijzonder dank ik de collega's die goede vrienden geworden zijn. Boudewijn, onze gedachtewisselingen over wetenschap en onderwijs waren inspirerend, evenals de whisky. Karina, your openness and fervour are marvellous; I promise to come visit you in Brazil.

Ook bedank ik alle collega's bij fysieke geografie. Mijn kamergenoten, Edwin, Niko, en later Jannis, jullie waren een welkome afleiding van het werk met grappige of goede gesprekken, een goed team voor het oplossen van technische problemen en een goed testpanel voor mijn baksels. De eerstejaars veldwerken aardwetenschappen die ik mocht begeleiden behoorden tot de mooiste perioden in m'n PhD. Martin, bedankt voor de altijd weer feilloze organisatie. En alle docenten en studenten, dank voor jullie gezelligheid en inzet. Ik heb ook erg genoten van het vaste clubje waarmee we Champions League wedstrijden keken (en hopelijk kijken) in O'Connells: Edwin, Jannis, Kim, Koko, en sporadisch wat anderen. Ook dank ik het PCRaster team, Kor en Oliver, voor de softwarematige hulp. En Oliver, dankjewel dat ik altijd bij je terecht kon met vragen over wat dan ook, en steevast een bruikbaar, innemend en humoristisch antwoord kreeg.

Geert, dankjewel voor onze prettige samenwerking bij de koppeling van MAGNET en PLUC. Everyone at CTBE, thank you for the help during my stay in Brazil. Everyone at INPE in Brazil, and especially Bernardo, many thanks for the Canasat data, an invaluable part of many of my analyses.

"He hallo! / Kijk mij eens de wolken breken / Liggend op mijn rug hier in het gras." Het TNO voetbal: wat voor week ik ook had, ik kon er al die jaren (ook al voor m'n PhD) op vertrouwen dat er met lunchtijd op donderdag een leuke pot voetbal gespeeld werd. Bij naam bedank ik de harde kern van de laatste jaren, Geert, Mart, Enrico, Bas, Simon, Jens, Sjef, Frank (2x), Jonathan, Jan, en Olwijn, maar uiteraard ook dank aan iedereen die af en toe meedeed en sorry aan degenen die ik vergeten ben. Verder qua voetbal, Sporting '70 vrouwen 3 zaterdag en vrouwen 1 zondag, bedankt voor de weekendafleiding!

"Verslagen maar preciezer / De gelukkige verliezer / Van wat een wedstrijd met de Beatles en de burens was." Sanne en Michiel, het was jammer dat ik jullie kwijt raakte als collega's en kamergenoten maar ik heb er twee fantastische vrienden voor teruggekregen. Michiel, altijd fascinerend hoe je de grenzen van mijn denkbeelden ter discussie stelt; blijven doen. Sanne, van jouw onbevangenheid kan ik veel leren, je maakt me altijd vrolijk. Dankjewel dat je mijn paranymf wilt zijn.

"Alles moet aparter dan apart / Alles moet unieker dan uniek / Alles moet bijzonder / En dat alles maakt dat alles weer moet." Als er teveel moest, waren daar altijd mijn beste vrienden om op terug te vallen. Marijn, dankjewel voor je goede raad en de wandelingen in het weekend. Jules, waar zou ik zijn zonder jouw nuchtere

bespiegelingen. Christine, ik zie je niet vaak, maar onze vakanties samen zijn altijd intens en verrassend. En Siebe, jij verblijdt me vaak met leuke ideeën voor uitjes en activiteiten, waardoor we nu rock 'n roll dansen samen. Top dat je er altijd voor me bent er en ook nu als mijn paranimf.

“Alles moet aparter dan apart / Alles moet unieker dan uniek / Alles moet bijzonderder dan / Alles wat allang bijzonder was.” Mama, dankjewel voor de stevige basis die je me gegeven hebt. Je heb me ook geleerd om strak te plannen en vast te houden aan die planning, de eigenschap die me tot de selecte groep deed behoren van promovendi die binnen de tijd hun proefschrift afkregen. Papa, jij hebt me al jong wetenschappelijke bezieling gegeven. Bijvoorbeeld, wanneer ik als kleuter met jou de fiets zat en er een brandweerwagen langskwam, legde jij me geduldig het Dopplereffect uit. Als eerbetoon aan jou heb ik altijd mijn tweede naam voluit op mijn artikelen staan. Jij zult begrijpen waarom. Mensen kijken me vaak meewarig aan als ik zeg twee docenten als ouders te hebben, maar ik heb er een groot enthousiasme voor doceren aan overgehouden. Mama, bedankt dat je me bijbaantjes bezorgd hebt in wiskundebijles en examensurveillance. En papa bedankt dat je mij als tiener stiekem jouw tentamens voor de HU liet nakijken... Sandra, lief zusje, bedankt voor alles wat je met me deelt. Brent, blij dat jij er bent. Kristian en Rutger, de handigste broertjes ter wereld, ik zie jullie niet zo vaak; dat komt doordat ik een Honda rijd ;-).

“Wat de Beatles en de burens was.” Ewout, jij was er voor me op het kritische moment in m'n PhD. Jammer dat 'ik moet je kunnen zien!' achteraf misschien niet alleen een grap was over het inparkeren van een caravan. Enrico, dank je voor je steun, vertrouwen en gezelligheid in het laatste jaar, en voor Indonesië en pizza op vrijdag. “La la la / la la la / la ... / En voor de rest moet je bij de Beatles en de burens zijn.”

Toekomst, de Kans op iets Moois: “Une valse à mille temps / Une valse à mille temps / Une valse à mille temps / Offre seule aux amants / Trois cent trente-trois fois le temps / De bâtir un roman.”<sup>9</sup>

---

<sup>9</sup> Dit citaat komt uit 'La Valse à mille temps' van Jacques Brel.



## About the author

Judith Versteegen was born in Utrecht, The Netherlands, on October 26<sup>th</sup> 1985. She obtained a Bachelor of Science (BSc) in Earth Sciences at Utrecht University in 2007. After that, she worked for a year at the Netherlands Organisation for Applied Scientific Research, TNO, where she continued to work during her Master of Science (MSc) in Geo-Information Science at Wageningen University. She finished this Masters cum laude in 2010 with a thesis on agent based modelling of the spatial planning process of urban expansion. After this, she was solicited to build a land use change model for the PhD project of Floor van der Hilst, supervised by Prof. dr. André Faaij, at the Copernicus Institute of Sustainable Development at Utrecht University. The aim of this modelling exercise was to assess future bioenergy potentials in Mozambique. The pleasant cooperation and valuable research outcomes resulted in an offer to Judith to continue this line of research as a PhD candidate in the BE-Basic program, a public-private partnership developing industrial bio-based solutions for a sustainable society. Committed to pursue a career in academia and enjoying the modelling work she eagerly accepted this offer.



During her PhD, Judith spent some time in Brazil for working visits and she presented her work at international conferences in the United States, Portugal, Germany, Austria and the Netherlands. Within the Copernicus Institute, Judith readily served as an informal 'helpdesk' for GIS-related questions. She also contributed to the educational programmes of Earth Sciences and Geographical Information Management and Applications (GIMA) by supervising MSc students during their theses, supervising BSc students during fieldwork, giving lectures and computer practicals and preparing course materials for a variety of BSc and MSc courses. Judith's application for the basic teaching qualification for higher education (BKO) is currently pending and is expected to be granted early 2016.

At the end of her PhD, Judith, enjoying to help out her colleagues with GIS-challenges, became appealed by GIS-service oriented governmental organizations and companies. In January 2016, she has started working for Ordina, an information technology (IT) service company, as a GEO-IT software engineer.



## List of journal publications

Verstegen, J.A., Jonker, J.G.G., Karssenber, D., van der Hilst, F., Schmitz, O., de Jong, S.M., Faaij, A.P.C. (in prep.). A spatial optimization approach to find trade-offs between production costs and greenhouse gas emissions: a bioethanol case study.

Goh, C.S., Junginger, M., Verstegen, J.A., Faaij, A.P.C., Wicke B. (in prep.). Linking carbon stock change from land-use change to consumption of agricultural products: A structural review of approaches, methodologies and key functions.

Verstegen, J.A., Karssenber, D., van der Hilst, F. & Faaij, A.P.C. (2016). Detecting systemic change in a land use system by Bayesian data assimilation. *Environmental Modelling & Software*, 75, 424-438.

Wicke, B., van der Hilst, F., Daioglou, V., Banse, M., Beringer, T., Gerssen-Gondelach, S., Heijnen, S., Karssenber, D., Laborde, D., Lippe, M., van Meijl, H., Nassar, A., Powell, J., Prins, A.G., Rose, S.N.K., Smeets, E.M.W., Stehfest, E., Tyner, W.E., Verstegen, J.A., Valin, H., van Vuuren, D.P., Yeh, S., Faaij, A.P.C., (2015). Model collaboration for the improved assessment of biomass supply, demand, and impacts. *Global Change Biology Bioenergy*, 7(3), 422-437.

Verstegen, J.A., van der Hilst, F., Woltjer, G., Karssenber, D., de Jong, S.M. & Faaij, A.P.C. (2015). What can and can't we say about indirect land-use change in Brazil using an integrated economic – land-use change model? *Global Change Biology Bioenergy*, early view.

Verstegen, J.A., Karssenber, D., Hilst, F. van der & Faaij, A.P.C. (2014). Identifying a land use change cellular automaton by Bayesian data assimilation. *Environmental Modelling & Software*, 53, 121-136.

Diogo, V., van der Hilst, F., van Eijck, J., Faaij, A., Verstegen, J.A., Hilbert, J., Carballo, S., Volante, J. (2014). Combining empirical and theory-based land use modelling approaches to assess future availability of land and economic potential for sustainable biofuel production: Argentina as a case study. *Renewable & Sustainable Energy Reviews*, 34, 208-224.

van der Hilst, F., Verstegen, J.A., Zheliezna, T., Drozdova, O., Faaij, A.P. (2014). Integrated spatiotemporal modelling of bioenergy production potentials, agricultural land use, and related GHG balances; demonstrated for Ukraine. *Biofuels, Bioproducts and Biorefining*, 8(3), 291-411.

Verstegen, J.A., Karssenber, D., van der Hilst, F., & Faaij, A.P.C. (2012). Spatio-temporal uncertainty in Spatial Decision Support Systems: A case study of changing land availability for bioenergy crops in Mozambique, *Computers, Environment and Urban Systems* 36, 30-42.

van der Hilst, F., Verstegen, J.A., Karssenber, D., Faaij, A.P.C. (2012). Spatio-temporal land use modelling to assess land availability for energy crops - illustrated for Mozambique. *Global Change Biology Bioenergy* 4(6), 859-874.