

NEW CONSTRUCTIONS OF SATISFACTION CLASSES

ALI ENAYAT AND ALBERT VISSER

ABSTRACT. We use model-theoretic ideas to present a perspicuous and versatile method of constructing full satisfaction classes on models of Peano arithmetic. We also comment on the ramifications of our work on issues related to conservativity and interpretability.

1. Introduction

In our forthcoming paper [2] we explore satisfaction classes over a wide variety of ‘base theories’ ranging from weak fragments of arithmetic to systems of ZF set theory and beyond. This note provides a synopsis of some of this work in the context of the most popular base theory adopted in investigations of axiomatic theories of truth, namely PA (Peano Arithmetic).

The notion of a satisfaction class was first introduced and investigated by Krajewski in his 1976 paper [11]. Two noteworthy accomplishments of [11] are the following results:

- (1) If a countable model of a ‘base theory’ (such as PA) carries at least one full satisfaction class, then it carries continuum-many full satisfaction classes.
- (2) Every model of ZF has an elementary extension that carries a full satisfaction class.

The question whether the analogue of (2) holds for PA remained open until the appearance of the joint work [10] of Kotlarski, Krajewski, and Lachlan in 1981, in which the rather exotic proof-theoretic technology of ‘ \mathcal{M} -logic’ (an infinitary logical system based on a nonstandard model \mathcal{M}), was invented to construct ‘truth classes’ over countable recursively saturated models of PA.¹ This model-theoretic result can be used to show that the analogue of (2) does indeed hold for PA, which in turn can be used to show that PA^{FT}

Date: November 24, 2012.

This research was partially supported by a grant from the Descartes Center of Utrecht University, which supported the first author’s visit to Utrecht to work closely with the second author.

¹As explained in Section 4, a truth class is essentially a well-behaved kind of satisfaction class. The \mathcal{M} -logic methodology was further elaborated to establish refined constructions of full truth classes by Smith [15,16,17], Kaye [9], and Engström [3].

is conservative over PA , where $\text{PA}^{\text{FT}} = \text{PA} + \text{“T is a full truth class”}$. The conservativity of PA^{FT} over PA has attracted considerable philosophical attention, especially in relation to the grand debate concerning deflationism.²

In this paper we present a perspicuous method for the construction of full satisfaction classes that is dominantly based on *model-theoretic techniques* (e.g., expanding the language, compactness, and elementary chains). As we shall see, our construction method is quite versatile and can be used to construct many (if not all) of the results that have hitherto been only possible to establish with the use of \mathcal{M} -logic machinery. Furthermore, the method can also be employed to build new types of full satisfaction classes (see Section 6).

We present the necessary preliminaries in Section 2, and then in Section 3 we concentrate on the basic form of our new construction of full satisfaction classes, where it is used to show that every model of PA has an elementary extension that carries a full satisfaction class. The versatility of the methodology of Section 3 is illustrated in Section 4, in which an appropriate modification of the method is used to construct truth classes for models of PA . As explained in Section 5, certain arithmetizations of our construction can also be employed to establish that (1) PA^{FT} is interpretable in PA ; and (2) the conservativity of PA^{FT} over PA can be verified in PRA (Primitive Recursive Arithmetic). Finally, in Section 6 we briefly describe further applications of the methods introduced in this paper.

Acknowledgments. We are grateful to the editors of this volume for their interest in our work. Thanks also to Volker Halbach, Fredrik Engström, and James Schmerl for helpful feedback on preliminary drafts of this paper. We are particularly indebted to Schmerl for catching an inaccuracy in an earlier formulation of Lemma 3.1, and for his suggestion to distill the results of our paper [2] in this form for wider dissemination.

2. Preliminaries

2.1. Definition. Throughout this paper PA refers to Peano arithmetic formulated in a *relational language* \mathcal{L}_{PA} using the logical constants $\{\neg, \vee, \exists, =\}$. Note that in this formulation PA has no constant symbols; the arithmetical operations of addition and multiplication are construed as ternary relations; and conjunction, universal quantification, and other logical constants are taken as defined notions in the usual way.

²A recent noteworthy paper in this connection is McGee’s [13].

It is well known that PA has more than sufficient expressive machinery to handle syntactic notions. The following list of \mathcal{L}_{PA} -formulae will be useful here.³

- $\text{Form}(x)$ is the formula expressing “ x is the code of an \mathcal{L}_{PA} -formula using variables $\{v_i : i \in \mathbb{N}\}$, and the non-logical symbols available in \mathcal{L}_{PA} ”.
- $\text{Asn}(x)$ is the formula expressing “ x is the code of an assignment”, where an assignment here simply refers to a function whose domain consists of a (finite) set of variables. We use α and its variants (α' , α_0 , etc.) to range over assignments.
- $y \in \text{FV}(x)$ is the formula expressing “ $\text{Form}(x)$ and y is a free variable of x ”.
- $y \in \text{Dom}(\alpha)$ is the formula expressing “the domain of α includes y ”.
- $\text{Asn}(\alpha, x)$ is the following formula expressing “ α is an assignment for x ”:

$$(\text{Form}(x) \wedge \text{Asn}(\alpha) \wedge \forall y(y \in \text{Dom}(\alpha) \leftrightarrow y \in \text{FV}(x))).$$

- $x \triangleleft y$ is the formula expressing “ x is the code of an immediate subformula of the \mathcal{L}_{PA} -formula coded by y ”, i.e., $x \triangleleft y$ abbreviates the conjunction of $\text{Form}(y)$ and the following disjunction:

$$(y = \neg x) \vee \exists z((y = x \vee z) \vee (y = z \vee x)) \vee \exists i(y = \exists v_i x).^4$$

The theory PA^{FS} (read as “PA with full satisfaction”) is formulated in an *expansion* of the language \mathcal{L}_{PA} by adding a new *binary* predicate $\text{S}(x, y)$. The binary/unary distinction is of course not an essential one since PA has access to a definable pairing function. However, the binary/unary distinction *at the conceptual level* marks the key difference between satisfaction classes and truth classes (the latter are discussed in Section 4). PA^{FS} is defined below with the help of a collection of sentences $\text{Tarski}(\text{S}, \text{F})$.

When reading the definition below it is helpful to bear in mind that $\text{Tarski}(\text{S}, \text{F})$ expresses:

F is a subset of Form that is closed under immediate subformulae; each member of S is an ordered pair of the form (x, α) , where $x \in \text{F}$ and α is an assignment for x ; and S satisfies Tarski’s compositional clauses.

³All of the formulae in the list can be arranged to be Σ_1 -formulae in the sense of Definition 2.4.

⁴Technically speaking, this formula should be written so as to distinguish the logical operations of the meta-language with those of the object-language. For example, using Feferman’s commonly used ‘dot-convention’, one would write:

$$(y = \neg \dot{x}) \vee \exists z \left((y = x \dot{\vee} z) \vee (y = z \dot{\vee} x) \right) \vee \exists i (y = \exists v_i \dot{x}).$$

However, since the difference between the two kinds of operation will be always clear from the context, we have opted for the lighter notation.

2.2. Definition. $\text{PA}^{\text{FS}} := \text{PA} \cup \text{Tarski}(\text{S}, \text{Form})$, where $\text{Tarski}(\text{S}, \text{F})$ is the conjunction of the universal generalizations of the formulae $\text{tarski}_0(\text{S}, \text{F})$ through $\text{tarski}_4(\text{S}, \text{F})$ described below, all of which are formulated in $\mathcal{L}_{\text{PA}} \cup \{\text{F}(\cdot), \text{S}(\cdot, \cdot)\}$, where S and F do not appear in \mathcal{L}_{PA} .

In the following formulae R ranges over the *relations in* \mathcal{L}_{PA} ; t, t_0, t_1, \dots are *metavariables*, e.g, we write $R(t_0, \dots, t_{n-1})$ instead of $R(v_{i_0}, \dots, v_{i_{n-1}})$; and $\alpha' \supseteq \alpha$ abbreviates

$$(\text{Dom}(\alpha') \supseteq \text{Dom}(\alpha)) \wedge \forall t \in \text{Dom}(\alpha) \alpha(t) = \alpha'(t).$$

- $\text{tarski}_0(\text{S}, \text{F}) := (\text{F}(x) \rightarrow \text{Form}(x)) \wedge (\text{S}(x, \alpha) \rightarrow (\text{F}(x) \wedge \text{Asn}(\alpha, x))) \wedge (y \triangleleft x \wedge \text{F}(x) \rightarrow \text{F}(y)).$
- $\text{tarski}_{1,R}(\text{S}, \text{F}) := \left(\text{F}(x) \wedge (x = \ulcorner R(t_0, \dots, t_{n-1}) \urcorner) \wedge \text{Asn}(\alpha, x) \wedge \bigwedge_{i < n} \alpha(t_i) = a_i \right) \rightarrow (\text{S}(x, \alpha) \leftrightarrow R(a_0, \dots, a_{n-1})).$
- $\text{tarski}_2(\text{S}, \text{F}) := (\text{F}(x) \wedge (x = \neg y) \wedge \text{Asn}(\alpha, x)) \rightarrow (\text{S}(x, \alpha) \leftrightarrow \neg \text{S}(y, \alpha)).$
- $\text{tarski}_3(\text{S}, \text{F}) := (\text{F}(x) \wedge (x = y_1 \vee y_2) \wedge \text{Asn}(\alpha, x)) \rightarrow (\text{S}(x, \alpha) \leftrightarrow (\text{S}(y_1, \alpha \upharpoonright \text{FV}(y_1)) \vee \text{S}(y_2, \alpha \upharpoonright \text{FV}(y_2)))).$
- $\text{tarski}_4(\text{S}, \text{F}) := (\text{F}(x) \wedge (x = \exists t y) \wedge \text{Asn}(\alpha, x)) \rightarrow (\text{S}(x, \alpha) \leftrightarrow \exists \alpha' \supseteq \alpha \text{S}(y, \alpha')).$

2.3. Definition. Suppose $\mathcal{M} \models \mathcal{L}_{\text{PA}}$, $F \subseteq M$, and S is a binary relation on M .⁵

(a) S is an *F-satisfaction class* if $(\mathcal{M}, S, F) \models \text{Tarski}(\text{S}, \text{F})$.⁶

(b) Let $\omega_{\mathcal{M}}$ be the well-founded initial segment of \mathcal{M} that is isomorphic to the ordinal ω . We say that F is the set of *standard* \mathcal{L}_{PA} -formulae of \mathcal{M} if

$$F = \text{Form}^{\mathcal{M}} \cap \omega_{\mathcal{M}}.$$

In this case there is a unique *F-satisfaction class* on \mathcal{M} , known as the *Tarskian satisfaction class* on \mathcal{M} .

(c) S is a *full satisfaction class* on \mathcal{M} if S is an *F-satisfaction class* for $F := \text{Form}^{\mathcal{M}}$. This is equivalent to $(\mathcal{M}, S) \models \text{PA}^{\text{FS}}$.

2.4. Definition. $\Sigma_0 = \Pi_0$ = the collection of \mathcal{L}_{PA} -formulae all of whose quantifiers are of the form $\exists x < y \varphi$ or $\forall x < y \varphi$; Σ_{n+1} consists of formulae of the form $\exists x_0 \dots \exists x_{k-1} \varphi$, where $\varphi \in \Pi_n$; and Π_{n+1} consists of formulae of the form $\forall x_0 \dots \forall x_{k-1} \varphi$, where $\varphi \in \Sigma_n$. Here k ranges over ω , with

⁵Throughout the paper we use the convention of using M, M_0, N , etc. to denote the universes of discourse of structures $\mathcal{M}, \mathcal{M}_0, \mathcal{N}$, etc.

⁶Note that the closure of F under direct subformulae does not guarantee that F should also contain ‘infinitely deep’ subformulae of a nonstandard formula in F .

the understanding that $k = 0$ corresponds to an empty block of quantifiers; this convention leads to the pleasant consequence that $\Sigma_n \subseteq \Sigma_{n+1}$ and $\Pi_n \subseteq \Pi_{n+1}$ for all n .

2.5. Theorem. (Mostowski [9], [6]) *For each nonzero $n < \omega$ there is a binary Σ_n -formula $\text{Sat}_n(x, y)$ such that*

$$\text{PA} \vdash \text{Tarski}(\text{Sat}_n, \Sigma_n),$$

where Σ_n is (the arithmetization of) the set of codes of formulae in Σ_n .

3. The Basic Construction

In this section we explain the basic methodology of building satisfaction classes using tools from model theory. The following lemma lies at the heart of the main result of this section.

3.1. Lemma. *Let $\mathcal{N}_0 \models \text{PA}$, $F_1 := \text{Form}^{\mathcal{N}_0}$, $F_0 \subseteq F_1$, and suppose S_0 is an F_0 -satisfaction class. Then there is an elementary extension \mathcal{N}_1 of \mathcal{N}_0 that carries an F_1 -satisfaction class $S_1 \supseteq S_0$ and $(c, \alpha) \in S_0$ whenever $c \in F_0$, $\alpha \in N_0$, and $(c, \alpha) \in S_1$.*

Proof: Let $\mathcal{L}_{\text{PA}}^+(\mathcal{N}_0)$ be the language obtained by enriching \mathcal{L}_{PA} with constant symbols for each member of N_0 , and new *unary* predicates U_c for each $c \in \text{Form}^{\mathcal{N}_0}$. It helps to have in mind that the intended interpretation of U_c is $\{\alpha \in A_c : S_1(c, \alpha)\}$, where $A_c := \{\alpha : \mathcal{N}_1 \models \text{Asn}(\alpha, c)\}$.

We first wish to describe a new set of axioms

$$\Theta := \{\theta_c : c \in F_1\}$$

formulated in $\mathcal{L}_{\text{PA}}^+(\mathcal{N}_0)$, where θ_c stipulates ‘local Tarskian behavior’ for U_c .

If $R \in \mathcal{L}_{\text{PA}}$ and $\mathcal{N}_0 \models c = \ulcorner R(t_0, \dots, t_{n-1}) \urcorner$, then

$$\theta_c := \forall \alpha (U_c(\alpha) \leftrightarrow \text{Asn}(\alpha, c) \wedge R(\alpha(t_0), \dots, \alpha(t_{n-1}))).$$

If $\mathcal{N}_0 \models c = \neg d$, then

$$\theta_c := \forall \alpha (U_c(\alpha) \leftrightarrow \text{Asn}(\alpha, c) \wedge \neg U_d(\alpha)).$$

If $\mathcal{N}_0 \models c = d_1 \vee d_2$, then

$$\theta_c := \forall \alpha (U_c(\alpha) \leftrightarrow \text{Asn}(\alpha, c) \wedge (U_{d_1}(\alpha \upharpoonright \text{FV}(d_1)) \vee U_{d_2}(\alpha \upharpoonright \text{FV}(d_2)))).$$

If $\mathcal{N}_0 \models c = \exists v_a b$, then

$$\theta_c := \forall \alpha (U_c(\alpha) \leftrightarrow \text{Asn}(\alpha, c) \wedge \exists \alpha' \supseteq \alpha U_b(\alpha') \wedge \text{Asn}(\alpha', b)).$$

Let

$$\Gamma := \{U_c(\alpha) : c \in F_0 \text{ and } (c, \alpha) \in S_0\} \cup \{\neg U_c(\alpha) : c \in F_0 \text{ and } (c, \alpha) \notin S_0\},$$

and let

$$\text{Th}^+(\mathcal{N}_0) := \text{Th}(\mathcal{N}_0, a)_{a \in N_0} \cup \Theta \cup \Gamma.$$

We now proceed to show that $\text{Th}^+(\mathcal{N}_0)$ is consistent by demonstrating that each finite subset of $\text{Th}^+(\mathcal{N}_0)$ is interpretable in (\mathcal{N}_0, S_0) . To this end, suppose T_0 is a finite subset of $\text{Th}^+(\mathcal{N}_0)$ and let C consist of the collection of $c \in F_0$ such that U_c appears in T_0 . If $C = \emptyset$, T_0 is readily seen to be consistent, so we shall assume that $C \neq \emptyset$ for the rest of the argument.

Our goal is to construct subsets $\{U_c : c \in C\}$ of N_0 such that the following two conditions hold when U_c is interpreted by U_c :

- (1) $(\mathcal{N}_0, U_c)_{c \in C} \models \{\theta_c : c \in C\}$, and
- (2) For $c \in C \cap F_0$, $U_c = \{\alpha \in N_0 : (c, \alpha) \in S_0\}$.

We shall construct $\{U_c : c \in C\}$ *in stages*, beginning with the simplest formulae in C , and working our way up using Tarski rules for more complex ones. Recall that $c \triangleleft d$ expresses “ c is a direct subformula of d ”. Define \triangleleft^* on C by:

$$c \triangleleft^* d \text{ iff } (c \triangleleft d)^{\mathcal{N}_0} \text{ and } \theta_d \in T_0 \cap \Theta.$$

Note that whenever $c \triangleleft^* d$, then for all $c' \triangleleft d$ we have $c' \in C$ and $c' \triangleleft^* d$. The finiteness of C implies that (C, \triangleleft^*) is well-founded, which in turn helps us define a useful measure of complexity for $c \in C$ using the following recursive definition:

$$\text{rank}_C(c) := \sup\{\text{rank}_C(d) + 1 : d \in C \text{ and } d \triangleleft^* c\}.$$

Note that for $c \in C$, $\text{rank}_C(c) = 0$ precisely when $\theta_c \notin T_0 \cap \Theta$. Next, let

$$C_i := \{c \in C : \text{rank}_C(c) \leq i\}.$$

Observe that $C_0 \neq \emptyset$ (since C is finite and nonempty), and that if $c \in C_{i+1}$, then the codes of all immediate subformulae of the formula coded by c are in C_i . This observation ensures that the following recursive clauses yield a well-defined U_c for each $c \in C$.

- If $c \in C_0$ then $U_c := \begin{cases} \{\alpha : (c, \alpha) \in S_0\}, & \text{if } c \in F_0; \\ U_c := \emptyset, & \text{if } c \notin F_0. \end{cases}$

- If $c \in C_{i+1} \setminus C_i$ and $c = \neg d$, then

$$U_c := \{\alpha \in A_c : \alpha \notin U_d\}.$$

- If $c \in C_{i+1} \setminus C_i$ and $c = a \vee b$, then

$$U_c := \{\alpha \in A_c : \alpha \upharpoonright \text{FV}(a) \in U_a \text{ or } \alpha \upharpoonright \text{FV}(b) \in U_b\}.$$

- If $c \in C_{i+1} \setminus C_i$ and $c = \exists v_a b$, then

$$U_c := \{\alpha \in A_c : \exists \alpha' \in N (\alpha \subseteq \alpha' \text{ and } \alpha' \in U_b)\}.$$

Note that in the first item above, the choice of $U_c := \emptyset$ when $c \in C_0$ and $c \notin F_0$ is completely arbitrary.⁷ Also, in the third item above where $c = a \vee b$, both a and b will be in C_i , thanks to the properties of \triangleleft^* .

It is routine to verify, using induction on $\text{rank}_C(c)$, that (1) and (2) hold for $(\mathcal{N}_0, U_c)_{c \in C}$. More specifically, if $\text{rank}_C(c) = 0$, then $\theta_c \notin T_0 \cap \Theta$, so (1) is vacuously satisfied, and (2) is satisfied by design. On the other hand, when $\text{rank}_C(c) > 0$ then (1) is satisfied since U_c is defined so as to comply with Tarski conditions; and (2) is satisfied since S_0 is an F_0 -satisfaction class. This concludes the proof of the consistency of arbitrary finite subsets T_0 of $\text{Th}^+(\mathcal{N}_0)$, which in turn shows that $\text{Th}^+(\mathcal{N}_0)$ has a model, i.e., some elementary extension \mathcal{M}_1 of \mathcal{M}_0 has an expansion \mathcal{N}_1^+ of the form

$$\mathcal{N}_1^+ := (\mathcal{N}_1, U_c)_{c \in F_1}$$

with the property that $\mathcal{N}_1^+ \models \text{Th}^+(\mathcal{N}_0)$. Let S_1 be the binary relation on N_1 defined via

$$S_1(c, \alpha) \Leftrightarrow \alpha \in U_c.$$

It is evident that S_1 is an F_1 -satisfaction class, $S_1 \supseteq S_0$ and $(c, \alpha) \in S_0$ whenever $c \in F_0$, $\alpha \in N_0$, and $(c, \alpha) \in S_1$. \square

3.2. Theorem. *Let \mathcal{M}_0 be a model of PA of any cardinality.*

- If S_0 is an F_0 -satisfaction class on \mathcal{M}_0 , then there is an elementary extension \mathcal{M} of \mathcal{M}_0 that carries a full satisfaction class that extends S_0 .*
- There is an elementary extension \mathcal{M} of \mathcal{M}_0 that carries a full satisfaction class.*

Proof: Note that (b) is an immediate consequence of (a) since we may choose F_0 to be the set of *atomic* \mathcal{M}_0 -formulae and S_0 to be the obvious satisfaction predicate for F_0 . To establish (a), we note that by Lemma 3.1 there is an elementary extension \mathcal{M}_1 of \mathcal{M}_0 that carries an F_1 -satisfaction class, where $F_1 := \text{Form}^{\mathcal{M}_0}$. Lemma 3.1 allows this argument to be carried out ω -times to yield two sequences $\langle \mathcal{M}_i : i \in \omega \rangle$ and $\langle S_i : i \in \omega \rangle$ that satisfy the following properties for each $i \in \omega$:

- $\mathcal{M}_i \prec \mathcal{M}_{i+1}$;
- S_{i+1} is an F_{i+1} -satisfaction class on \mathcal{M}_{i+1} with $F_{i+1} := \text{Form}^{\mathcal{M}_i}$; and
- $S_i = S_{i+1} \cap \{(c, \alpha) : c \in F_i, \mathcal{M}_i \models \text{Asn}(\alpha, c)\}$.

Let $\mathcal{M} := \bigcup_{i \in \omega} \mathcal{M}_i$, and $S := \bigcup_{i \in \omega} S_i$. Tarski's elementary chain theorem and

- together imply that \mathcal{M} elementarily extends \mathcal{M}_0 . It is easy to see, using (2) and (3), that S is a full satisfaction class on \mathcal{M} . \square

⁷As shown in [2] this feature can be exploited to construct 'pathological' satisfaction classes, such as the one mentioned at the end of Section 6 of this paper.

Theorem 3.2, when coupled with the completeness theorem of first order logic, immediately yields the following conservativity result.

3.3. Corollary. PA^{FS} is a conservative extension of PA.

Proof: Suppose not. Then for some arithmetical sentence φ we have:

- (1) $\text{PA}^{\text{FS}} \vdash \varphi$, and
- (2) $\text{PA} \not\vdash \varphi$.

Since (2) implies that $\text{PA} \cup \{\neg\varphi\}$ is consistent, by the completeness theorem for first order logic, there is a model $\mathcal{M}_0 \models \text{PA} \cup \{\neg\varphi\}$. On the other hand, by part (b) of Theorem 3.2 there is an elementary extension \mathcal{M}_1 of \mathcal{M}_0 that carries a full satisfaction class, and therefore by (1) $\mathcal{M}_1 \models \varphi$. This contradicts the fact that \mathcal{M}_1 elementarily extends \mathcal{M}_0 . \square

3.4. Corollary. Every resplendent model of PA carries a full satisfaction class. In particular, every countable recursively saturated model of PA carries a full satisfaction class.

Proof: The first claim directly follows from the definition of a resplendent model. The second claim follows from the first claim, when coupled with the key result that countable recursively saturated models are resplendent (see [9, Section 15.2] for more detail). \square

4. Truth Classes

With the exception of Krajewski's original paper [11], what we refer to as a 'truth class' here has been dubbed 'satisfaction class' in the model-theoretic literature. More specifically, Krajewski [11] employed the framework of satisfaction classes over base theories formulated in relational languages as in this paper, however, the later series of papers [10], [16], and [17] all used the framework of truth classes over Peano arithmetic formulated in a relational language, augmented with 'domain constants'. Later, Kaye [9] developed the theory of satisfaction classes over models of PA in languages incorporating function symbols; his work was extended by Engström [3] to truth classes over models of PA in functional languages.

As explained in this section, there is a simple canonical correspondence between truth classes over models of PA (in a relational language) and certain types of satisfaction classes, here referred to as 'extensional'. The main aim of this section is to demonstrate that the method of building satisfaction classes in the previous section can be conveniently modified so as to yield full *extensional* satisfaction classes (and thereby: full truth classes) over appropriate models of PA.

Within PA one can easily define an injective function c that yields the code for a constant symbol \bar{x} for each member x of the domain. This enables PA to internally represent the language $\mathcal{L}_{\text{PA}}^+ = \mathcal{L}_{\text{PA}} + \text{'domain constants'}$. We can then add a unary predicate $\text{T}(x)$ denoting a *truth class* (instead of a binary

predicate $S(x, y)$ for a satisfaction class) to \mathcal{L}_{PA} , whose intended interpretation is “ x is the code of a true *sentence* σ ”, where σ is an arithmetical sentence formulated in a language $\mathcal{L}_{\text{PA}}^+$. We will make this more precise in the following definition.

4.1. Definition. $\text{PA}^{\text{FT}} := \text{PA} \cup \text{Tarski}(\text{T})$, where $\text{Tarski}(\text{T})$ is the conjunction of the universal generalizations of $\text{tarski}_0(\text{T})$ through $\text{tarski}_4(\text{T})$, all formulated in the language $\mathcal{L}_{\text{PA}} \cup \{\text{T}(\cdot)\}$, as described below.⁸ In what follows $\text{Sent}(x)$ is the \mathcal{L}_{PA} -formula that expresses “ x is a formula of $\mathcal{L}_{\text{PA}}^+$ with no free variables”, and R ranges over relations symbols in \mathcal{L}_{PA} .

- $\text{tarski}_0(\text{T}) := (\text{T}(x) \rightarrow \text{Sent}(x))$.
- $\text{tarski}_{1,R}(\text{T}) := (\ulcorner R(\bar{t}_0, \dots, \bar{t}_{n-1}) \urcorner = x) \rightarrow (R(t_0, \dots, t_{n-1}) \leftrightarrow \text{T}(x))$.
- $\text{tarski}_2(\text{T}) := (x = \neg y) \rightarrow (\text{T}(x) \leftrightarrow \neg \text{T}(y_1))$.
- $\text{tarski}_3(\text{T}) := (x = y_1 \vee y_2) \rightarrow (\text{T}(x) \leftrightarrow (\text{T}(y_1) \vee \text{T}(y_2)))$.
- $\text{tarski}_4(\text{T}) := (x = \exists v_i \varphi) \rightarrow (\text{T}(x) \leftrightarrow \exists z \text{T}(\varphi(\bar{z})))$.

T is a *full truth class* on \mathcal{M} if $(\mathcal{M}, T) \models \text{PA}^{\text{FT}}$.

4.2. Definition. A *substitution* for a formula ψ of \mathcal{L}_{PA} is a function

$$\sigma : \text{FV}(\psi) \rightarrow \text{Var}$$

such that σ respects substitutability in the ‘usual way’, i.e., if x is a free variable of ψ , then x is not in the scope of any quantifier that binds $\sigma(x)$. Given ψ and σ as above, let $\psi * \sigma$ be the formula obtained from ψ by applying the substitution σ , and A be the set of pairs (φ, α) such that α is an assignment for the formula φ . This allows us to define a key equivalence relation \sim on A by decreeing that $(\varphi_0, \alpha_0) \sim (\varphi_1, \alpha_1)$ iff there is some $(\psi, \beta) \in A$, and there are substitutions σ_0 and σ_1 for ψ , with

$$\varphi_i = \psi * \sigma_i \text{ and } \beta = \alpha_i \circ \sigma_i, \text{ for } i = 0, 1.$$

In the above, $\alpha_i \circ \sigma_i$ is the composition of α_i and σ_i . It is important to bear in mind that, intuitively speaking, $(\varphi_0, \alpha_0) \sim (\varphi_1, \alpha_1)$ means that φ_0 and φ_1 are the same except for their free variables, and for all variables x and y , if x occurs freely in the same position in φ_0 as y does in φ_1 , then $\alpha_0(x) = \alpha_1(y)$.

- An F -satisfaction class S is *extensional* if for all φ_0 and φ_1 in F , $\mathcal{M} \models (\varphi_0, \alpha_0) \sim (\varphi_1, \alpha_1)$ implies $(\varphi_0, \alpha_0) \in S$ iff $(\varphi_1, \alpha_1) \in S$.⁹

The following proposition describes a canonical correspondence between extensional satisfaction classes and truth classes. The routine but laborious proof is left to the reader.

⁸ PA^{FT} is the relational analogue of the theory of CT^\dagger in Halbach’s monograph [8]. The base theory of CT^\dagger is PA formulated in a *functional* language. The conservativity of CT^\dagger over the functional language version of PA can also be established using the techniques of this paper (see Section 6).

⁹Note that an extensional satisfaction predicate need not be closed under re-naming of bound variables.

- In what follows \mathbf{c} is the \mathcal{M} -definable injection $m \mapsto_{\mathbf{c}} \bar{m}$ that designates a constant symbol \bar{m} for each $m \in M$, and $\varphi(\mathbf{c} \circ \alpha)$ is the sentence in the language $\mathcal{L}_{\text{PA}}^+$ obtained by replacing each occurrence of a free variable x of φ with the constant symbol \bar{m} , where $\alpha(x) = m$.

4.3. Proposition. *Suppose $\mathcal{M} \models \text{PA}$, T is a full truth class on \mathcal{M} , and S is an extensional full satisfaction class on \mathcal{M} .*

- (a) $\mathcal{S}(T)$ is an extensional satisfaction class on \mathcal{M} , where $\mathcal{S}(T)$ is defined as the collection of ordered pairs (φ, α) such that $\varphi(\mathbf{c} \circ \alpha) \in T$.
- (b) $\mathcal{T}(S)$ is a truth class on \mathcal{M} , where $\mathcal{T}(S)$ is defined as the collection of $\varphi \in \mathcal{L}_{\text{PA}}^+$ such that for some $\psi \in \mathcal{L}_{\text{B}}^+$ and some assignment α for ψ , $\varphi = \psi(\mathbf{c} \circ \alpha)$ and $(\psi, \alpha) \in S$.
- (c) $\mathcal{S}(\mathcal{T}(S)) = S$, and $\mathcal{T}(\mathcal{S}(T)) = T$.

Before describing the construction of extensional satisfaction classes we need the preliminaries presented in Definition 4.4 and Lemma 4.5.

4.4. Definition.

- (a) Given formulae φ_0 and φ_1 of \mathcal{L}_{PA} , we write $\varphi_0 \approx \varphi_1$ if there is a formula ψ , and substitutions σ_0 and σ_1 for ψ such that $\varphi_i \approx \psi * \sigma_i$ for $i = 0, 1$.
- (b) Given $c \in \text{Form}^{\mathcal{M}}$, let $\mathbf{TC}_{\mathcal{M}}(c)$ be the *externally defined* transitive closure of c with respect to the direct subformula relation, i.e.,

$$\mathbf{TC}_{\mathcal{M}}(c) := \bigcup_{n < \omega} \mathbf{TC}_{\mathcal{M}}(c, n),$$

where $\mathbf{TC}_{\mathcal{M}}(c, 0) := \{c\}$ and

$$\mathbf{TC}_{\mathcal{M}}(c, n+1) := \{x \in M : x \triangleleft^{\mathcal{M}} d \text{ for some } d \in \mathbf{TC}_{\mathcal{M}}(c, n)\}.$$

The following lemma presents salient features of the two equivalence relations \sim and \approx .

4.5. Lemma. *Let \sim be as in Definition 4.2; and $\mathbf{TC}_{\mathcal{M}}(c)$ and \approx be as in Definition 4.4.*

- (i) *If $d \in \mathbf{TC}_{\mathcal{M}}(c)$ and $d \neq c$, then $\neg(c \approx d)$.*
- (ii) *\approx preserves the principal connectives, i.e., it relates negations to negations, disjunctions to disjunctions, and existential formulae to existential formulae with the same bound variable. Moreover, if $\neg c \approx \neg d$, then $c \approx d$; if $c \vee d \approx c' \vee d'$, then $c \approx c'$ and $d \approx d'$; and if $\exists t c \approx \exists t' c'$, then $t = t'$ and $c \approx c'$.*
- (iii) *If $(\varphi_0, \alpha_0) \sim (\varphi_1, \alpha_1)$, then $\varphi_0 \approx \varphi_1$.*
- (iv) *If $(\neg\varphi_0, \alpha_0) \sim (\neg\varphi_1, \alpha_0)$, then $(\varphi_0, \alpha_0) \sim (\varphi_1, \alpha_1)$.*
- (v) *If $(\varphi_0 \vee \varphi_1, \alpha) \sim (\varphi'_0 \vee \varphi'_1, \alpha')$, then $(\varphi_0, \alpha \upharpoonright \text{FV}(\varphi_0)) \sim (\varphi'_0, \alpha' \upharpoonright \text{FV}(\varphi'_0))$ and $(\varphi_1, \alpha \upharpoonright \text{FV}(\varphi_1)) \sim (\varphi'_1, \alpha' \upharpoonright \text{FV}(\varphi'_1))$.*

(vi) If $\varphi = \exists t \psi$, and $\varphi' = \exists t' \psi'$, and $(\varphi, \alpha) \sim (\varphi', \alpha')$, then $t = t'$ and for some e

$$(\varphi, \alpha[t : e]) \sim (\varphi', \alpha'[t' : e]).^{10}$$

The next Lemma presents a variant of Lemma 3.1 that is our main tool for constructing extensional satisfaction classes.

4.6. Lemma. *Let $\mathcal{N}_0 \models \text{PA}$, $F_1 := \text{Form}^{\mathcal{N}_0}$, $F_0 \subseteq F_1$, and suppose S_0 is an extensional F_0 -satisfaction class. Then there is an elementary extension \mathcal{N}_1 of \mathcal{N}_0 that carries an extensional F_1 -satisfaction class $S_1 \supseteq S_0$ and $(c, \alpha) \in S_0$ whenever $c \in F_0$, $\alpha \in N_0$, and $(c, \alpha) \in S_1$.*

Proof. Let Θ and Γ be as in the proof of Lemma 3.1, and let

$$\text{Th}^+(\mathcal{N}_0) := \text{Th}(\mathcal{N}_0, a)_{a \in N_0} \cup \Theta \cup \Gamma \cup \Delta,$$

where $\Delta := \{\delta_{cc'} : c, c' \in F_1\}$, and

$$\delta_{cc'} := \forall \alpha \forall \alpha' ((c, \alpha) \sim (c', \alpha') \rightarrow (\mathbf{U}_c(\alpha) \leftrightarrow \mathbf{U}_{c'}(\alpha'))).$$

The proof of the lemma would be complete once we verify that $\text{Th}^+(\mathcal{N}_0)$ has a model. To this end, we shall demonstrate that every finite subset T_0 of $\text{Th}^+(\mathcal{N}_0)$ is interpretable in \mathcal{N} . Let C be the collection of $c \in F_1$ such that c appears in T_0 . Also, let \triangleleft^* and $\text{rank}_C(c)$ be precisely as in the proof of Lemma 3.1.

We can extend C to another finite set \overline{C} so that it satisfies a certain closure property, namely: whenever we have $c \approx c'$ and $d \triangleleft^* c$, where c, c' and d are all in \overline{C} , then there is some $d' \in \overline{C}$ such that $d' \triangleleft^* c'$ with $d \approx d'$. This can be done simply by adding any missing direct subformulae d' by an appropriate recursion.¹¹ By replacing C by \overline{C} we may therefore additionally assume:

(#) If c and c' are both in C with $c \approx c'$, then $\text{rank}_C(c) = \text{rank}_C(c')$.

As in the proof of Lemma 3.1, we then recursively construct $\{U_c : c \in C\}$ such that:

- (1) $(\mathcal{N}_0, U_c)_{c \in C} \models \{\theta_c : c \in C\}$ and
- (2) For $c \in C \cap F_0$, $U_c = \{\alpha \in N_0 : (c, \alpha) \in S_0\}$.

¹⁰Here $\alpha[t : e]$ is the assignment obtained by redefining the value of α at the variable t to be e if $t \in \text{Dom}(\alpha)$; note that $\alpha[t : e] = \alpha$ if $t \notin \text{Dom}(\alpha)$.

¹¹More specifically, first define \triangleleft° on C by $d' \triangleleft^\circ c$ iff $d' \triangleleft^* c' \approx c$, for some $c' \in C$. Since \triangleleft° is cycle-free, C is well-founded, and therefore lends itself to a ranking function $\text{rank}_C^\circ(c)$. Let $n = \max\{\text{rank}_C^\circ(c) : c \in C\}$, and for $0 \leq i \leq n$ define $D_i := \{c \in C : \text{rank}_C^\circ(c) = i\}$. Next use a 'backward' recursion to define E_n, E_{n-1}, \dots, E_0 via:

- $E_n := D_n$;
- $E_{n-(i+1)} := D_{n-(i+1)} \cup \{d : d \triangleleft^{\mathcal{N}_0} c \text{ for some } c \in E_{n-i}\}$.

Finally, let $\overline{C} := E_n \cup \dots \cup E_0$. It is easy to see that \overline{C} is finite, extends C , and has the desired closure property.

It remains to show:

$$(3) (\mathcal{N}_0, U_c)_{c \in C} \models \{\delta_{cc'} : c, c' \in C\}.$$

We establish (3) by using induction on $\text{rank}_C(c)$ to show that $\forall c \in C P(c)$, where $P(c)$ abbreviates:

$$\forall c' \in C (\mathcal{N}_0, U_c)_{c \in C} \models \forall \alpha \forall \alpha' ((c, \alpha) \sim (c', \alpha') \rightarrow (U_c(\alpha) \leftrightarrow U_{c'}(\alpha'))).$$

If $\text{rank}_C(c) = 0$ and $(c, \alpha) \sim (c', \alpha')$, then by part (iii) of Lemma 4.5 we have $c \approx c'$, which in turn by (#) assures us that $\text{rank}_C(c') = 0$. This makes it clear that $P(c)$ holds when $\text{rank}_C(c) = 0$ since S_0 is assumed to be an *extensional* F_0 -satisfaction class.

To verify the inductive step, suppose:

$$(4) P(x) \text{ holds for all } x \in C \text{ with } \text{rank}_C(x) = i.$$

Let $c \in C$ with $\text{rank}_C(c) = i+1$, and suppose $(c, \alpha) \sim (c', \alpha')$, where $c = \exists t d$. Then $c' = \exists t' d'$, and $d \approx d'$ by part (ii) of Lemma 4.5. Observe that thanks to (#) we have:

$$(5) \text{rank}_C(c') = i + 1 \text{ and } \text{rank}_C(d) = \text{rank}_C(d') = i.$$

Now if $\alpha \in U_c$, then $\alpha[t : e] \in U_d$ for some e by (1), and therefore by part (vi) of Lemma 4.5, we obtain:

$$(6) (d, \alpha[t : e]) \sim (d', \alpha'[t' : e]).$$

Using (4), (5), and (6) we may now conclude that $\alpha'[t : e] \in U_{d'}$, which by (1) yields $\alpha' \in U_{d'}$, thus completing the verification of the quantificational case (by symmetry). A similar reasoning can be carried out for propositional cases. This concludes the proof of consistency of T_0 .

The rest is precisely as before: the consistency of $\text{Th}^+(\mathcal{N}_0)$ implies that there is an elementary extension \mathcal{N}_1 of \mathcal{N}_0 that has an expansion $\mathcal{N}_1^+ := (\mathcal{N}_1, S) \models \text{Th}^+(\mathcal{N}_0)$, and the binary relation S_1 on \mathcal{N}_1 defined via

$$S_1(c, \alpha) \Leftrightarrow \alpha \in U_c$$

has the property that S_1 is an extensional F_1 -satisfaction, $S_1 \supseteq S_0$, and $(c, \alpha) \in S_0$ whenever $c \in F_0$, $\alpha \in N_0$, and $(c, \alpha) \in S_1$. \square

4.7. Theorem. *Let $\mathcal{M}_0 \models \text{PA}$. There is an elementary extension \mathcal{M} of \mathcal{M}_0 that carries a full extensional satisfaction class.*

Proof: Since the satisfaction class S_0 on the collection F_0 of atomic formulae of \mathcal{M}_0 is extensional, we may use Lemma 4.6 instead of Lemma 3.1 in order to carry out the elementary chain argument of Theorem 3.2. \square

By coupling Theorem 4.7 with part (b) of Proposition 4.3 we obtain:

4.8. Corollary. *Every model of PA has an elementary extension that carries a full truth class.*

Finally, the line of reasoning employed in the proof of Corollary 3.3 shows, using Corollary 4.8, that:

4.9. Corollary. PA^{FT} is a conservative extension of PA.

5. Arithmetization, Interpretability, and Conservativity

Here we briefly discuss the *arithmetization* of the constructions of the previous two sections, with an eye towards issues connected with interpretability and conservativity. As explained in [2, Section 4] the compactness and elementary chain argument employed in the proofs of Theorems 3.2 and 4.7 can be implemented in the fragment $\text{I}\Sigma_2$ of PA with the help of the ‘Low Basis Theorem’ of Recursion Theory. Coupled with Orey’s Compactness Theorem, this can be used to establish the following:

5.1. Theorem. [2] PA^{FT} is interpretable in PA.¹²

On the other hand, the technology of LL_1 -sets¹³ of [6, Theorem 4.2.7.1, p. 104] can be used to show that the proofs of both theorems 3.2 and 4.7 can even be implemented in the fragment $\text{I}\Sigma_1$ of PA. In light of the fact that the statement “ PA^{FT} is conservative over PA” is a Π_2 -statement, and $\text{I}\Sigma_1$ is well known¹⁴ to be Π_2 -conservative over PRA, we obtain the following:

5.2. Theorem. [2] The conservativity of PA^{FT} over PA can be verified in PRA.¹⁵

5.3. Remark. The verification of the conservativity of PA^{FT} over PA within PRA was first claimed by Halbach in [7], using cut-elimination.¹⁶ Later, Fischer [4] gave a proof, based on the cut-elimination argument in [7], to show that PA^{FT} is interpretable in PA. Unfortunately, a gap was discovered recently (by Fujimoto) in the cut-elimination argument in [7], which in turn impaired Fischer’s interpretability claim. Happily, Leigh [12] has succeeded in developing a proof-theoretic demonstration of the conservativity of PA^{FT} over PA that is implementable in PRA. Moreover, [12, Theorem 1] can be

¹²Indeed B^{FS} turns out to be interpretable in B for all base theories B that have access to the full scheme of induction over their ambient ‘numbers’. In particular, ACA^{FS} is interpretable in ACA. On the other hand, as shown in [2, Section 8], ACA_0^{FS} is *not* interpretable in ACA_0 (more generally, B^{FS} is shown to be *not* interpretable in B, if B is finitely axiomatizable).

¹³ LL_1 -sets are a special type of ‘low’ sets.

¹⁴This classical result was independently established by Mints, Parsons, and Takeuti, using proof-theoretic methods. The work of Paris and Kirby (described in [14, IX.3]), and more recently Avigad [1] has also provided model-theoretic demonstrations of this conservativity result.

¹⁵Indeed, by using the technique of Friedman [5], this conservativity result is already verifiable in the fragment SEFA (Superexponential Arithmetic) of PRA.

¹⁶Halbach’s base theory in his work is the usual version of PA that is formulated in a functional language.

used to verify the interpretability of PA^{FT} over PA , by using Fischer's strategy in [4].¹⁷ Therefore, Theorems 5.1 and 5.2 can be arrived at via two completely different routes.

6. Further Results

In Section 4 we saw that the core methodology of Section 3 can be fine-tuned to build full extensional satisfaction classes. Indeed, as shown in [2] one can strengthen Theorem 4.7 by imposing further desirable conditions on the satisfaction class S . For example, every model \mathcal{M}_0 of PA has an elementary extension \mathcal{M} that carries a full extensional satisfaction class S that satisfies all of the following additional properties:

- (1) $\text{Sat}_n^{\mathcal{M}} \subseteq S$ for all $n \in \omega$ (see Theorem 2.5 for Sat_n).
- (2) If $c \in \text{Form}^{\mathcal{M}}$ and $\mathcal{M} \models$ “ c is an axiom of PA ”, then S deems c to be ‘true’.¹⁸
- (3) If c and c' are \mathcal{M} -formulae such that $\mathcal{M} \models$ “ c' is an alphabetic variant¹⁹ of c ”, then $(c, \alpha) \in S$ iff $(c', \alpha) \in S$.

Furthermore, the third condition above can be strengthened by accomodating a combination of extensional equivalence and alphabetic equivalence, thereby yielding truth classes that are closed under alphabetic equivalence. We have also shown that a small dose of condition (1) can be used to build full truth classes over models of arithmetical theories formulated in *functional* languages (this result will appear in the projected sequel to [2]).

One can also use the method of Section 3 to build bizarre satisfaction classes. For example, as shown in [2, Section 5], every model \mathcal{M}_0 of PA has an elementary extension \mathcal{M} that has a full satisfaction class S that exhibits the following pathology:

$$\{a \in M : (\sigma_a, \alpha_{\text{Null}}) \in S\} = \omega_{\mathcal{M}},$$

where $\omega_{\mathcal{M}}$ is the well-founded initial segment of \mathcal{M} that is isomorphic to ω , and σ_a is defined for all $a \in M$ by a recursion within \mathcal{M} via the following clauses:

- $\sigma_0 := \exists v_0 (v_0 = v_0)$ (or $\sigma_0 =$ any other logically valid sentence);
- $\sigma_{n+1} := (\sigma_n \vee \sigma_n)$.

¹⁷We are grateful to Graham Leigh for his kind permission to quote his unpublished work here.

¹⁸As remarked in the last sentence of [10], this condition can also be arranged using the machinery of \mathcal{M} -logic. Note that ‘axioms of PA ’ in the sense used here do not include the logical axioms.

¹⁹ c' is an alphabetic variant of c if c' is obtainable from c by the usual rules of re-naming the bound variables of c .

REFERENCES

- [1] J. Avigad, *Saturated models of universal theories*, **Annals of Pure and Applied Logic** vol. 118 (2002), pp. 219-234.
- [2] A. Enayat and A. Visser, *Full satisfaction classes in a general setting* (Part I), to appear in **Logic Group Preprint Series**, available at [http://www.phil.uu.nl/preprints/lgps/Utrecht Preprint Series](http://www.phil.uu.nl/preprints/lgps/Utrecht%20Preprint%20Series).
- [3] F. Engström, *Satisfaction classes in nonstandard models of first-order arithmetic*, **arXiv.org.math**, available at: <http://arxiv4.library.cornell.edu/abs/math/0209408v1>
- [4] M. Fischer, *Minimal truth and interpretability*, **Review of Symbolic Logic**, vol. 2 (2009), pp. 799–815.
- [5] H. Friedman, *Finitist proofs of conservation*, **FOM Archives**, available at <http://cs.nyu.edu/pipermail/fom/1999-September/003405.html>.
- [6] P. Hájek and P. Pudlák, **Metamathematics of First-Order Arithmetic**, Springer, 1993.
- [7] V. Halbach, *Conservative theories of classical truth*, **Studia Logica**, vol. 62 (1999), pp. 353-370.
- [8] -----, **Axiomatic Theories of Truth**, Cambridge University Press, Cambridge, 2011.
- [9] R. Kaye, **Models of Peano Arithmetic**, Oxford Logic Guides, Oxford University Press, Oxford, 1991.
- [10] H. Kotlarski, S. Krajewski, and A.H. Lachlan, *Construction of satisfaction classes for nonstandard models*, **Canadian Mathematical Bulletin**. vol. 24 (1981), pp. 283–293.
- [11] S. Krajewski, *Nonstandard satisfaction classes*, in **Set Theory and Hierarchy Theory: A Memorial Tribute to Andrzej Mostowski** (ed. W. Marek et al.) Lecture Notes in Mathematics, vol. 537, Springer-Verlag, Berlin, 1976, pp. 121-144.
- [12] G. E. Leigh, *Deflating truth*, manuscript (November 2012).
- [13] V. McGee, *In praise of the free lunch: why disquotationalists should embrace compositional semantics*, in **Self-reference**, CSLI Lecture Notes, 178, CSLI Publ., Stanford, pp. 95-120.
- [14] S. Simpson, **Subsystems of second order arithmetic**, Perspectives in Mathematical Logic, Springer-Verlag, Berlin, 1999.
- [15] S. T. Smith, **Non-standard Syntax and Semantics and Full Satisfaction Classes**, Ph.D. thesis, Yale University, New Haven, Connecticut, 1984.
- [16] -----, *Nonstandard characterizations of recursive saturation and resplendency*, **Journal of Symbolic Logic**, vol. 52 (1987), pp. 842-863.
- [17] -----, *Nonstandard definability*, **Annals of Pure and Applied Logic**, vol. 42 (1989), pp. 21–43.

ALI ENAYAT

DEPARTMENT OF MATHEMATICS AND STATISTICS, AMERICAN UNIVERSITY, 4400 MASS. AVE. NW, WASHINGTON, DC 20016-8050, USA.
enayat@american.edu

ALBERT VISSER

DEPARTMENT OF PHILOSOPHY, BESTUURSGEBOUW, HEIDELBERGLAAN 6, 584 CS UTRECHT, THE NETHERLANDS.
albert.visser@phil.uu.nl