

Active learning for screening prioritization in systematic reviews

A simulation study

Gerbrich Ferdinands Raoul Schram Jonathan de Bruin Ayoub Bagheri
Daniel Oberski Lars Tummers Rens van de Schoot

10 August, 2020

1 **Abstract**

2 **Background** Conducting a systematic review requires great screening effort. Various tools have been
3 proposed to speed up the process of screening thousands of titles and abstracts by engaging in active learning.
4 In such tools, the reviewer interacts with machine learning software to identify relevant publications as
5 early as possible. To gain a comprehensive understanding of active learning models for reducing workload
6 in systematic reviews, the current study provides an methodical overview of such models. Active learning
7 models were evaluated across four different classification techniques (naive Bayes, logistic regression, support
8 vector machines, and random forest) and two different feature extraction strategies (TF-IDF and doc2vec).
9 Moreover, models were evaluated across six systematic review datasets from various research areas to assess
10 generalizability of active learning models across different research contexts.

11 **Methods** Performance of the models were assessed by conducting simulations on six systematic review
12 datasets. We defined desirable model performance as maximizing recall while minimizing the number of
13 publications needed to screen. Model performance was evaluated by recall curves, WSS@95, RRF@10, and
14 ATD.

15 **Results** Within all datasets, the model performance exceeded screening at random order to a great degree.
16 The models reduced the number of publications needed to screen by 91.7% to 63.9%.

17 **Conclusions** Active learning models for screening prioritization show great potential in reducing the
18 workload in systematic reviews. Overall, the Naive Bayes + TF-IDF model performed the best.

19 **Systematic Review registrations** Not applicable.

20 **Keywords:** systematic reviews, active learning, screening prioritization, researcher-in-the-loop, title-and-
21 abstract screening, automation, text mining.

22 Background

23 Systematic reviews are top of the bill in research. A systematic review brings together all available studies
24 relevant to answer a specific research question [1]. Systematic reviews inform practice and policy [2] and
25 are key in developing clinical guidelines [3]. However, systematic reviews are costly because to identify
26 publications relevant to answering the research question, they among else involve the manual screening of
27 thousands of titles and abstracts.

28 Conducting a systematic review typically requires over a year of work by a team of researchers [4]. Nevertheless,
29 systematic reviewers are often bound to a limited budget and timeframe. Currently, the demand for systematic
30 reviews exceeds the available time and resources by far [5]. Especially when answering the research question
31 at hand is urgent, it is extremely challenging to provide a review that is both timely and comprehensive.

32 To ensure a timely review, reducing the workload in systematic reviews is essential. With advances in machine
33 learning (ML), there has been wide interest in tools to reduce the workload in systematic reviews [6]. Various
34 ML models have been proposed, aiming to predict whether a given publication is relevant or irrelevant to the
35 systematic review. Previous findings suggest that such models potentially reduce the workload with 30-70%
36 at the cost of losing 5% of relevant publications, in else, a 95% recall [7].

37 A well-established approach to increase the efficiency of title and abstract screening is screening prioritization
38 [8, 9]. In screening prioritization, the ML model presents the reviewer with the publications that are most
39 likely to be relevant first, thereby expediting the process of finding all of the relevant publications. Such
40 an approach allows for substantial time-savings in the screening process as the reviewer can decide to stop
41 screening after a sufficient number of relevant publications have been found [10]. Moreover, the early retrieval
42 of relevant publications facilitates a faster transition of those publications to the next steps in the review
43 process [8].

44 Recent studies have demonstrated the effectiveness of screening prioritization by means of active learning
45 models [10, 11, 12, 13, 14, 15, 16]. With active learning, the ML model can iteratively improve its predictions
46 on unlabeled data by allowing the model to select the records from which it wants to learn [17]. The model
47 proposes these records to a human annotator who provides the records with labels, which the model then
48 uses to update its predictions. The general assumption is that by letting the model select which records are
49 labeled, the model can achieve higher accuracy more quickly while requiring the human annotator to label as
50 few records as possible [18]. Active learning has proven to be an efficient strategy in large unlabeled datasets
51 where labels are expensive to obtain [18]. This makes the screening phase in systematic reviewing an ideal
52 candidate solution for such models, because typically labeling a large number of publications is very costly.

53 When active learning is applied in the screening phase, the reviewer screens publications that are suggested
54 by an active learning model. Subsequently, the active learning model learns from the reviewers' decision
55 ('relevant', 'irrelevant') and uses this knowledge to update its predictions and to select the next publication
56 to be screened by the reviewer.

57 The application of active learning models in systematic reviews has been extensively studied [10, 11, 12, 15, 16].
58 While previous studies have evaluated active learning models in many forms and shapes [10, 11, 12, 13, 14, 15,
59 19, 20, 21], ready-to-use software tools implementing such models (Abstrackr [22], Colandr [23], FASTREAD
60 [11], Rayyan [24], and RobotAnalyst [25]) currently use the same classification technique to predict relevance
61 of publications, namely support vector machines (SVM). It was found [26, 27] that different classification
62 techniques can serve different needs in the retrieval of relevant publications, for example the desired balance
63 between recall and precision. Therefore, it is essential to evaluate different classification techniques in the
64 context of active learning models. The current study investigates active learning models adopting four
65 classification techniques: naive Bayes (NB), logistic regression (LR), SVM, and random forest (RF). These
66 are widely adopted techniques in text classification [28] and are fit for software tools to be used in scientific
67 practice due to their relatively short computation time.

68 Another component that influences model performance is how the textual content of titles and abstracts is
69 represented in a model, called the feature extraction strategy [19, 20, 29]. One of the more sophisticated
70 feature extraction strategies is doc2vec (D2V), also known as paragraph vectors [30]. D2V learns continuous
71 distributed vector representations for pieces of text. In distributed text representations, words are assumed
72 to appear in the same context when they are similar in terms of a latent space, the "embedding". A word
73 embedding is simply a vector of scores estimated from a corpus for each word; D2V is an extension of this idea
74 to document embeddings. Embeddings can sometimes outperform simpler feature extraction strategies such
75 as term frequency-inverse document frequency (TF-IDF). They can be trained on large corpora to capture
76 wider semantics and subsequently applied in a specific systematic reviewing application [30]. Therefore, it is
77 interesting to compare models adopting D2V to models adopting TF-IDF.

78 Lastly, previous studies have mainly focussed on reviews from a single scientific field, like medicine [15, 16] or
79 computer science [10, 11]. To draw conclusions about the general effectiveness of active learning models, it
80 is essential to evaluate models on reviews from varying research contexts [7, 31]. To our knowledge, Miwa
81 et al [12] were the only researchers to make a direct comparison between systematic reviews from different
82 research areas, such as the social and the medical sciences. They found that the application of active learning
83 to systematic reviews was more difficult on a systematic review from the social sciences due to the different
84 nature of the vocabularies used. Thus, it is of interest to evaluate model performance across different research

85 contexts, namely social science, medical science and computer science.

86 Taken together, for a more comprehensive understanding of active learning models in the context of systematic
87 reviewing, a methodical evaluation of such models is required. The current study aims to address this issue
88 by answering the following research questions:

89 **RQ1** What is the performance of active learning models across four classification techniques?

90 **RQ2** What is the performance of active learning models across two feature extraction strategies?

91 **RQ3** Does the performance of active learning models differ across six systematic reviews from four research
92 areas?

93 The purpose of this paper is to show the usefulness of active learning models for reducing workload in title
94 and abstract screening in systematic reviews. We adopt four different classification techniques (NB, LR, SVM,
95 and RF) and two different feature extraction strategies (TF-IDF and D2V) for the purpose of maximizing
96 the number of identified relevant publications, while minimizing the number of publications needed to screen.
97 Models were assessed by conducting a simulation on six systematic review datasets. To assess generalizability
98 of the models across research contexts, datasets containing previous systematic reviews were collected from
99 the fields of medical science [32, 33, 34], computer science [11], and social science [35, 36]. The models,
100 datasets and simulations are implemented in a pipeline of active learning for screening prioritization, called
101 **ASReview** [37]. **ASReview** is a generic open source tool, encouraging fellow researchers to replicate findings
102 from previous studies. To facilitate usability and acceptability of ML-assisted title and abstract screening in
103 the field of systematic review our scripts and data used are openly available.

104 **Methods**

105 **Technical details**

106 What follows is a more detailed account of the active learning models to clarify the choices made in the
107 design of the current study.

108 **Task description**

109 The screening process of a systematic review starts with all publications obtained in the search. The task is
110 to identify which of these publications are relevant, by screening their titles and abstracts. In *active learning*
111 for screening prioritization, the screening process proceeds as follows:

- 112 • Start with the set of all unlabeled records (titles and abstracts)

- 113 • The reviewer provides a label for a few, e.g. 5-10 records, creating a set of labeled records. The label
114 can be either *relevant* or *irrelevant*.
- 115 • The active learning cycle starts:
 - 116 1. A classifier is trained on the labeled records
 - 117 2. The classifier predicts relevancy scores for all unlabeled records
 - 118 3. Based on the predictions by the classifier, the model selects the record with the highest relevancy
119 score
 - 120 4. The model requests the reviewer to screen this record
 - 121 5. The reviewer screens the record and provides a label, *relevant* or *irrelevant*.
 - 122 6. The newly labeled record is moved to the training data
 - 123 7. Back to step 1
 - 124 8. Repeat this cycle until the reviewer decides to stop [10] or until all records have been labeled

125 In this active learning cycle, the model incrementally improves its predictions on the remaining unlabeled
126 title and abstracts. Relevant titles and abstracts are identified as early in the process as possible. A more
127 technical description of the active learning cycle can be found in Additional file 1.

128 This case is an example of pool-based active learning, as the next record to be queried is selected by predicting
129 relevancy for all records in a fixed pool [17]. Another form of active learning is stream-based active learning,
130 in which the data is regarded as a stream instead of a fixed pool, in which the model selects one record at
131 a time and then decides whether or not to query this record. This approach of active learning is preferred
132 when it is expensive or impossible to exhaustively search the data for selecting the next query. A possible
133 application of stream-based active learning is living systematic reviews, as the review is continually updated
134 as new evidence becomes available. For an example see the study by Wynants et al. [38].

135 **Class imbalance problem**

136 Typically, only a fraction of the publications belong to the relevant class (2.94%, [4]). To some extent, this
137 fraction is under the control of the researcher through the search criteria: if the researcher narrows the search
138 query, it will generally result in a higher proportion of relevant publications. However, in most applications
139 this practice would yield an unacceptable number of false negatives (erroneously excluded papers) in the
140 querying phase of the review process. For this reason, the querying phase in most practical applications

141 would yield a very low percentage of relevant publications. Because there are generally far fewer examples of
142 relevant than irrelevant publications to train on, the class imbalance causes the classifier to miss relevant
143 publications [7]. Moreover, classifiers can achieve high accuracy but still fail to identify any of the relevant
144 publications [15].

145 Previous studies have addressed the class imbalance problem by rebalancing the training data in various ways
146 [7]. To decrease the class imbalance in the training data, we rebalance the training set by a technique we
147 propose to call “dynamic resampling” (DR). DR undersamples the number of irrelevant publications in the
148 training data, whereas the number of relevant publications are oversampled such that the size of the training
149 data remains the same. The ratio between relevant and irrelevant publications in the rebalanced training
150 data is not fixed, but dynamically updated and depends on the number of publications in the available
151 training data, the total number of publications in the dataset, and the ratio between relevant and irrelevant
152 publications in the available training data. Additional file 2 provides a detailed script to perform DR.

153 **Classification**

154 To make relevancy predictions on the unlabeled publications, a classifier is trained on features from the
155 training data. The performance of the following four classifiers is explored:

- 156 • Support vector machines (SVM) - SVMs separate the data into classes by finding a multidimensional
157 hyperplane [39, 40].
- 158 • L2-regularized logistic regression (LR) - models the probabilities describing the possible outcomes
159 by a logistic function. The classifier uses regularization, shrinking coefficients of features with small
160 contributions to the solution towards zero.
- 161 • Naive Bayes (NB) is a supervised learning algorithm often used in text classification. Based on Bayes’
162 theorem, with the ‘naive’ assumption that all features are independent given the class value [41].
- 163 • Random forests (RF) is a supervised learning algorithm where a large number of decision trees are fit
164 on samples obtained from the original data by sampling both rows (bootstrapped samples) and columns
165 (feature samples). In prediction mode, each tree casts a vote on the class, and the final prediction is the
166 class that received the most votes [42].

167 **Feature extraction**

168 To predict relevance of a given publication, the classifier uses information from the publications in the dataset.
169 Examples of information are titles and abstracts. However, a model cannot make predictions from the titles

170 and abstracts as they are; their textual content needs to be represented numerically as feature vectors. This
171 process of numerically representing textual content is referred to as ‘feature extraction’.

172 TF-IDF is a specific way of assigning scores to the cells of the “document-term matrix” used in all bag-of-words
173 representations. That is, the rows of the document-term matrix represent the documents (titles and abstracts)
174 and the columns represent all words in the dictionary. Instead of simply counting the number of times each
175 word occurred in the given document, TF-IDF assigns a score to a word relative to the number of documents
176 the word occurs. The idea behind weighting words by their rarity is that surprising word choices should
177 subsequently make for more discriminative features [43]. A disadvantage of TF-IDF and other bag-of-words
178 methods is that they do not take the ordering of words into account, thereby ignoring syntax. However, in
179 practice, TF-IDF is often found to be a strong baseline [44].

180 In recent years, a range of modern methods have been developed that often outperform bag-of-words
181 approaches. Here, we consider doc2vec, an extension of the classic word2vec embedding [30]. In word
182 embedding models, whether a word did or did not happen to appear in a specific context is predicted by
183 its similarity to that context in a latent space - the “embedding”. The context is usually a sliding window
184 across training sentences. For example, if the window “child ate cookies” occurs in the training data, this
185 might be compared with a random ‘negative’ window that did not occur, such as “child lovely cookies”. The
186 tokens “child” and “cookies” are then assigned scores (vectors) that give a higher inner product with the
187 “child” vector, and a smaller product with “lovely”. The word vectors of “ate” and “lovely” are similarly
188 updated. Typically the embedding dimension is a few hundred, i.e. each word vector contains some two
189 hundred scores. Note that if “cookies” previously co-occurred frequently with “spinach”, then the above
190 also indirectly makes “ate” more similar to “spinach”, even if these two words have not been observed yet
191 in the same context. Thus, the distributed representation learns something of the meaning of these words
192 through their occurrence in similar contexts. D2V performs such a procedure while including a paragraph
193 identifier, allowing for paragraph embeddings - or, in our case, embeddings for titles and abstracts. In short,
194 D2V converts each abstract into a vector of a few hundred scores, which can be used to predict relevancy.

195 **Query strategy**

196 The active learning model can adopt different strategies in selecting the next publication to be screened by
197 the reviewer. A strategy mentioned before is selecting the publication with the highest probability of being
198 relevant. In the active learning literature this is referred to as certainty-based active learning [17]. Another
199 well-known strategy is uncertainty-based active learning, where the instances that are presented next are
200 those instances on which the model’s classifications are the least certain, i.e. close to 0.5 probability [17].

201 Further strategies include selecting the next instance to optimize for various criteria, including: model fit
202 (MLI), model change (MMC), parameter estimate accuracy (EVR), and expected (EER) or worst-case (MER)
203 prediction accuracy [45]. Although uncertainty sampling is not explicitly motivated by the optimization of
204 any particular criterion, intuitively it can be seen as attempting to improve the model’s accuracy by reducing
205 uncertainty about its parameter estimates.

206 Simulation-based comparisons of these methods across different domains have yielded an ambiguous picture
207 of their relative strengths [12, 45]. What has become clear from such studies is that the features of the task
208 at hand determine the effectiveness of active learning strategies (“no free active lunch”). For example, if
209 a linear classifier is used for a task that also happens to have a Bayes optimal linear decision boundary, a
210 model-based approach such as Fisher information reduction can be expected to perform well, whereas the
211 same technique can be disastrous when the model is misspecified - a fact that cannot be known in advance.
212 Furthermore, the criteria mentioned above differ from the task of title and abstract screening in systematic
213 reviews: here, the aim is not to obtain an accurate model, but rather to end up with a list of records belonging
214 to the relevant class [46]. This is the criterion corresponding intuitively to certainty-based sampling. For this
215 reason, we choose to focus on certainty-based sampling strategies as the baseline strategy for active learning
216 in systematic reviewing. However, different strategies may outperform our baseline in specific applications.

217 **Simulation study**

218 This section describes the simulation study that was carried out to answer the research questions.

219 **Set-up**

220 To address RQ1, four models were investigated combining each classifier with TF-IDF feature extraction:

- 221 1. SVM + TF-IDF
- 222 2. NB + TF-IDF
- 223 3. RF + TF-IDF
- 224 4. LR + TF-IDF

225 To address RQ2, the classifiers were combined with D2V feature extraction, leading to the following three
226 combinations:

- 227 5. SVM + D2V
- 228 6. RF + D2V
- 229 7. LR + D2V

230 The combination NB + D2V could not be tested because the multinomial naive Bayes classifier¹ requires
231 a feature matrix with positive values, whereas the D2V feature extraction approach² produces a feature
232 matrix that can contain negative values. The performance of the seven models was evaluated by simulating
233 every model on six systematic review datasets, addressing RQ3. Hence, 42 simulations were carried out,
234 representing all model-dataset combinations.

235 Instead of having a human reviewer label publications manually, the screening process was simulated by
236 retrieving the labels in the data. Each simulation started with an initial training set of one relevant and one
237 irrelevant publication to represent a challenging scenario where the reviewer has very little prior knowledge
238 on the publications in the data. The model was retrained each time after a publication had been labeled. A
239 simulation ended after all publications in the dataset had been labeled. To account for sampling variance,
240 every simulation was repeated 15 times. To account for bias due to the content of the initial publications,
241 the initial training set was randomly sampled from the dataset for each of the 15 trials. Although varying
242 over trials, the 15 initial training sets were kept constant for each dataset to allow for a direct comparison of
243 models within datasets. A seed value was set to ensure reproducibility. The simulation study was carried out
244 using the ASReview simulation extension [47]. For each simulation, hyperparameters were optimized through
245 a Tree of Parzen Estimators (TPE) algorithm [48] to arrive at maximum model performance.

246 Simulations were carried out in ASReview version 0.9.3 [47]. Analyses were carried out using R version 3.6.1
247 [49]. The simulations were carried out on Cartesius, the Dutch national supercomputer.

248 Datasets

249 The models were simulated on a convenience sample of six systematic review datasets. The data selection
250 process was driven by two factors. Firstly, datasets are collected from various research areas to assess
251 generalizability of the models across research contexts (RQ3). Secondly, all original data files have to be
252 openly published with a CC-BY license. Datasets are available through ASReview’s systematic review
253 datasets GitHub³.

254 The Wilson dataset [50] - from the field of medicine - is from a review on the effectiveness and safety of
255 treatments of Wilson Disease, a rare genetic disorder of copper metabolism [33]. From the same field, the
256 ACE dataset contains publications on the efficacy of Angiotensin-converting enzyme (ACE) inhibitors, a
257 treatment drug for heart disease [32]. Additionally, the Virus dataset is from a systematic review on studies

¹https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html#sklearn.naive_bayes.MultinomialNB

²<https://radimrehurek.com/gensim/models/doc2vec.html>

³<https://github.com/asreview/systematic-review-datasets>

258 that performed viral Metagenomic Next-Generation Sequencing (mNGS) in farm animals [34]. From the field
259 of computer science, the Software dataset contains publications from a review on fault prediction in software
260 engineering [51]. The Nudging dataset [52] belongs to a systematic review on nudging healthcare professionals
261 [35], stemming from the social sciences. From the same research area, the PTSD dataset contains publications
262 on studies applying latent trajectory analyses on posttraumatic stress after exposure to traumatic events
263 [36]. Of these six datasets, ACE and Software have been used for model simulations in previous studies on
264 ML-aided title and abstract screening [11, 32].

265 Data were preprocessed from their original source into a dataset, containing title and abstract of the
266 publications obtained in the initial search. Duplicates and publications with missing abstracts were removed
267 from the data. Datasets were labeled to indicate which candidate publications were included in the systematic
268 review, thereby denoting relevant publications. All datasets consisted of thousands of candidate publications,
269 of which only a fraction was deemed relevant to the systematic review. For the Virus and the Nudging
270 dataset, this proportion was about 5 percent. For the remaining six datasets, the proportions of relevant
271 publications were centered around 1-2 percent. (Table 1).

272 **Evaluating performance**

273 Model performance was assessed by three different measures, Work Saved over Sampling (WSS), Relevant
274 References Found (RRF), and Average Time to Discovery (ATD). WSS indicates the reduction in publications
275 needed to be screened, at a given level of recall [32]. Typically measured at a recall level of 95%, WSS@95
276 yields an estimate of the amount of work that can be saved at the cost of failing to identify 5% of relevant
277 publications. In the current study, WSS is computed at 95% recall. RRF@10 represents the proportion of
278 relevant publications that are found after screening 10% of all publications.

279 Both RRF and WSS are sensitive to the position of the cutoff value and the distribution of the data.
280 Moreover, WSS makes assumptions about the acceptable recall level whereas this level might depend on the
281 research question at hand [7]. Therefore, we introduce the ATD, the average fraction of non-reviewed relevant
282 publications during the review (except the relevant publications in the initial training set). The ATD is an
283 indicator of performance throughout the entire screening process instead of performance at some arbitrary
284 cutoff value. The ATD is computed by taking the average of the Time to Discovery (TD) of all relevant
285 publications. The TD for a given relevant publication i is computed as the fraction of publications needed to
286 screen to detect i . Additional file 3 provides a detailed script to compute the ATD.

287 Furthermore, model performance was visualized by plotting recall curves. Plotting recall as a function of
288 the proportion of screened publications offers insight in model performance throughout the entire screening

289 process [11, 13]. The curves give information in two directions. On the one hand they display the number of
290 publications that need to be screened to achieve a certain level of recall, but on the other hand they present
291 how many relevant publications are identified after screening a certain proportion of all publications (RRF).
292 For each simulation, the RRF@10, WSS@95, and ATD are reported as means over 15 trials. To indicate the
293 spread of performance within simulations, the means are accompanied by an estimated standard deviation
294 \hat{s} . To compare the overall performance across datasets, median performance is reported for every dataset,
295 accompanied by the Median Absolute Deviation (MAD), indicating variability between models within a
296 certain dataset. Recall curves are plotted for each simulation, representing the average recall over 15 trials \pm
297 the standard error of the mean.

298 **Results**

299 This section proceeds as follows: Firstly, as an example the results of the Nudging dataset are discussed in
300 detail to provide a basis for answering the research questions. Secondly, the results are presented for each
301 research question over all datasets.

302 **Evaluation on the Nudging dataset**

303 Figure 1a shows the recall curves for all simulations on the Nudging dataset. As described in the previous
304 section, these curves plot recall as a function of the proportion of publications screened. The curves represent
305 the average recall over 15 trials \pm the standard error of the mean in the direction of the y-axis. The x-axis
306 is cut off at 40% since at this point in screening all models had already reached 95% recall. The dashed
307 horizontal lines indicate the RRF@10 values, the dashed vertical lines the WSS@95 values. The dashed black
308 diagonal line corresponds to the expected recall curve when publications are screened in a random order.

309 The recall curves were used to examine model performance throughout the entire screening process and
310 to make a visual comparison between models within datasets. For example in Figure 1a, after screening
311 about 30% of the publications all models had already found 95% of the relevant publications. Moreover,
312 after screening 5% the green curve - representing the RF + TF-IDF model - splits away from the others
313 and remains to be the lowest of all curves until about 30% of publications have been screened. Hence, from
314 screening 5 to 30 percent of publications, the RF + TF-IDF model was the slowest in finding the relevant
315 publications. The ordering of the remaining recall curves changes throughout the screening process, but
316 maintain relatively similar performance at face value.

317 Figure 1b shows a subset of the recall curves in Figure 1a, namely the curves of the first four models to

318 allow for a visual comparison across classification techniques adopting the TF-IDF feature extraction strategy.
319 Figure 1c shows recall curves for the remaining three models to compare the models using D2V feature
320 extraction. Figures 1d to 1f compare recall curves for models adopting the TF-IDF feature extraction strategy
321 to recall curves for their D2V-using counterparts.

322 It can be seen from Table 2 that in terms of ATD, the best performing models on the Nudging dataset were
323 SVM + D2V and LR + D2V, both with an ATD of 8.8%. This indicates that the average proportion of
324 publications needed to screen to find a relevant publication was 8.8% for both models. In the SVM + D2V
325 model, the standard deviation was 0.33, whereas for the LR + D2V model $\hat{s} = 0.47$. This indicates that for
326 the SVM + D2V model, the ATD values of individual trials were closer to the overall mean compared to
327 the LR + D2V model, meaning that the SVM + D2V model performed more stable across different initial
328 training datasets. Median ATD for this dataset was 9.5% with an MAD of 1.05, indicating that for half of
329 the models, the ATD was within 1.05 percentage point distance from the median ATD.

330 As Table 3 shows, the highest WSS@95 value on the Nudging dataset was achieved by the NB + TF-IDF
331 model with a mean of 71.7%, meaning that this model reduced the number of publications needed to screen
332 by 71.7% at the cost of losing 5% of relevant publications. The estimated standard deviation of 1.37 indicates
333 that in terms of WSS@95, this model performed the most stable across trials. The model with the lowest
334 WSS@95 value was RF + TF-IDF ($\bar{x} = 64.9\%$, $\hat{s} = 2.50$). Median WSS@95 of these models was 66.9%, with
335 a MAD of 3.05, indicating that of all datasets, the WSS@95 values of the models simulated on the Nudging
336 dataset varied the most within the Nudging dataset.

337 As can be seen from the data in Table 4, LR + D2V was the best performing model in terms of RRF@10,
338 with a mean of 67.5% indicating that after screening 10% of publications, on average 67.5% of all relevant
339 publications had been identified, with a standard deviation of 2.59. The worst performing model was RF +
340 TF-IDF ($\bar{x} = 53.6\%$, $\hat{s} = 2.71$). Median performance was 62.6%, with an MAD of 3.89 indicating again that
341 of all datasets, the RRF@10 values were most dispersed for models simulated on the Nudging dataset.

342 Overall evaluation

343 Recall curves for the simulations on the five remaining datasets are presented in Figure 2. For the sake of
344 conciseness, recall curves are only plotted once per dataset, like in Figure 1a for the Nudging dataset. Please
345 refer to Additional file 4 for figures presenting subsets of recall curves for the remaining datasets, like in
346 Figure 1b-f.

347 First of all, as the recall curves exceed the expected recall at screening at random order by far for all datasets,

348 the models were able to detect the relevant publications much faster compared to when screening publications
349 at random order. Even the worst results outperform this reference condition. Across simulations, the ATD
350 was at maximum 11.8% (in the Nudging dataset), the WSS@95 at least 63.9% (in the Virus dataset), and
351 the lowest RRF@10 was 53.6% (in the Nudging dataset). Interestingly, all these values were achieved by the
352 RF + TF-IDF model.

353 Similar to the simulations on the Nudging dataset (Figure 1a), the ordering of recall curves changes throughout
354 the screening process, indicating that some models perform better at the start of the screening phase whereas
355 others models take the lead later on. Moreover, the ordering of models in the Nudging dataset (Figure 1a) is
356 not replicated in the remaining five datasets (Figure 2).

357 **RQ1 - Comparison across classification techniques**

358 The first research question was aimed at evaluating the four models adopting either the NB, SVM, LR or
359 RF classification technique combined with TF-IDF feature extraction. When comparing ATD-values of the
360 models (Table 2), the NB + TF-IDF model ranked first in the ACE, Virus, and Wilson dataset, shared first
361 in the PTSD and Software dataset, and second in the Nudging dataset in which the SVM + D2V and LR +
362 D2V models achieved the lowest ATD value. The RF + TF-IDF ranked last in all of the datasets except for
363 the ACE and the Wilson dataset, in which the RF + D2V model achieved the highest ATD-value.

364 Additionally, in terms of WSS@95 (Table 3) the ranking of models was strikingly similar across all datasets.
365 In the Nudging, ACE, and Virus dataset, the highest WSS@95 value was always achieved by the NB +
366 TF-IDF model, followed by LR + TF-IDF, SVM + TF-IDF, and RF + TF-IDF. In the PTSD and the
367 Software dataset this ranking applied as well, except that two models showed the same WSS@95 value. The
368 ordering of the models for the Wilson dataset was NB + TF-IDF, RF + TF-IDF, LR + TF-IDF and SVM +
369 TF-IDF.

370 Moreover, in terms of RRF@10 (Table 4) the NB + TF-IDF model achieved the highest RRF@10 value in the
371 ACE and Virus dataset. Within the PTSD dataset, LR + TF-IDF was the best performing model, for the
372 Software and Wilson dataset this was SVM + D2V, and for the Nudging dataset LR + D2V performed best.

373 Taken together, these results show that while all four models perform quite well, the NB + TF-IDF model
374 demonstrates high performance on all measures across all datasets, whereas the RF + TF-IDF model never
375 performed best on any of the measures across all datasets.

376 **RQ2 - Comparison across feature extraction techniques**

377 This section is concerned with the question of how models using different feature extraction strategies relate
378 to each other. The recall curves for the Nudging dataset (Figure 1d-f) show a clear trend of the models
379 adopting D2V feature extraction outperforming their TF-IDF counterparts. This trend also shows from
380 the WSS@95 and RRF@10 values indicated by the vertical and horizontal lines in the figure. Likewise, the
381 ATD values (Table 2) indicate that for the models adopting a particular classification technique, the models
382 adopting D2V feature extraction always achieved a lower ATD-value than the model adopting TF-IDF feature
383 extraction.

384 In contrast, this pattern of models adopting D2V outperforming their TF-IDF counterparts in the Nudging
385 dataset is not replicated across other datasets. Whether evaluated in terms of recall curves, WSS@95,
386 RRF@10, or ATD, the findings were mixed. Neither one of the feature extraction strategies showed superior
387 performance within certain datasets nor within certain classification techniques.

388 **RQ3 - Comparison across research contexts**

389 First of all, models showed much higher recall curves for some datasets than for others. While performance
390 of the PTSD (Figure 2a) and Software datasets (Figure 2b) was quite high, performance was much lower
391 across models for the Nudging (Figure 1a) and Virus (Figure 2d) datasets. The models simulated on the
392 PTSD and Software datasets also demonstrated high performances in terms of the median ATD, WSS@95,
393 and RRF@10 values for these models (Table 2, 3, and 4).

394 Secondly, variability of between-model performance differed across datasets. For the PTSD (Figure 2a),
395 Software (Figure 2b), and the Virus (Figure 2d) datasets, recall curves form a tight group meaning that within
396 these datasets, the models performed similarly. In contrast, for the Nudging (Figure 1a), ACE (Figure 2c),
397 and Wilson (Figure 2e) dataset, the recall curves are much further apart, indicating that model performance
398 was more dependent on the adopted classification technique and feature extraction strategy. The MAD values
399 of the ATD, WSS@95 and RRF@10 confirm that model performance is less spread out within the PTSD,
400 Software, and Virus datasets than within the Nudging, ACE, and Wilson datasets. Moreover, the curves for
401 the ACE (Figure 2c) and Wilson (Figure 2e) datasets show a larger standard error of the mean compared the
402 other datasets.

403 Taken together, although model performance is very data-dependent, there does not seem to be a distinction
404 in performance between the datasets from the biomedical sciences (ACE, Virus, and Wilson) and datasets
405 from other fields (Nudging, PTSD, and Software).

406 **Discussion**

407 The current study evaluates the performance of active learning models for the purpose of identifying relevant
408 publications in systematic review datasets. It has been one of the first attempts to examine different
409 classification strategies and feature extraction strategies in active learning models for systematic reviews.
410 Moreover, this study has provided a deeper insight into the performance of active learning models across
411 research contexts.

412 **Active learning-based screening prioritization**

413 All models were able to detect 95% of the relevant publications after screening less than 40% of the total
414 number of publications, indicating that active learning models can save more than half of the workload in
415 the screening process. In a previous study, the ACE dataset was used to simulate a model that did not use
416 active learning, finding a WSS@95 value of 56.61% [32], whereas the models in the current study achieved
417 far superior WSS@95 values varying from 68.6% to 82.9% in this dataset. In another study [11] that did
418 use active learning, the Software dataset was used for simulation and a WSS@95 value of 91% was reached,
419 strikingly similar to the values found in the current study which ranged from 90.5% to 92.3%.

420 **Classification techniques**

421 The first research question in this study sought to evaluate models adopting different classification techniques.
422 The most important finding to emerge from these evaluations was that the NB + TF-IDF model consistently
423 performed as one of the best models. Our results suggest that while SVM performed fairly well, the LR and
424 NB classification techniques are good if not superior alternatives to this default classifier in software tools.
425 Note that LR and NB were always good methods for text classification tasks [53].

426 **Feature extraction strategy**

427 The overall results on models adopting D2V versus TF-IDF feature extraction strategy remain inconclusive.
428 According to our findings, models adopting D2V do not outperform models adopting the well-established
429 TF-IDF feature extraction strategy. Given these results, preference goes out to the TF-IDF feature extraction
430 technique as this relatively simple technique will lead to a model that is easier to interpret. Another advantage
431 of this technique is its short computation time.

432 **Research contexts**

433 Difficulty of applying active learning is not confined to any particular research area. The suggestion that
434 active learning is more difficult for datasets from the social sciences compared to data from the medical
435 sciences [12] does not seem to be the case. A possible explanation for this is that this difficulty has to be
436 attributed to factors more directly related to the systematic review at hand, such as the proportion of relevant
437 publications or the complexity of inclusion criteria used to identify relevant publications [16, 54]. Although
438 the current study did not investigate the inclusion criteria of systematic reviews, the datasets on which the
439 active learning models performed worst, Nudging and Virus, were interestingly also the datasets with the
440 highest proportion of relevant publications, 5.4% and 5.0%, respectively.

441 **Limitations and future research**

442 When applied to systematic reviews, the success of active learning models stands or falls with the gener-
443 alizability of model performance across unseen datasets. In our study, is important to bear in mind that
444 model hyperparameters were optimized for each model-dataset combination. Thus, the observed results
445 reflect the maximum model performance for each presented datasets. The question remains whether model
446 performance generalizes to datasets for which the hyperparameters are not optimized. Further research
447 should be undertaken to determine the sensitivity of model performance to the hyperparameter values.

448 Additionally, while the sample of datasets in the current study is diverse compared to previous studies, the
449 sample size ($n=6$) does not allow for investigating how model performance relates to characteristics of the
450 data, such as the proportion of relevant publications. To build more confidence in active learning models for
451 screening publications, it is essential to identify how data characteristics affect model performance. Such a
452 study requires more data on systematic reviews. Thus, a more thorough study depends on researchers to
453 openly publish their systematic review datasets.

454 Moreover, the runtime of simulations varied widely across models, indicating that some models take longer
455 to retrain after a publication has been labeled than other models. This has important implications for the
456 practical application of such models, as an efficient model should be able to keep up with the decision-making
457 speed of the reviewer. Further studies should take into account the retraining time of models.

458 **Conclusions**

459 Overall, the findings confirm the great potential of active learning models to reduce the workload for systematic
460 reviews. The results shed new light on the performance of different classification techniques, indicating that

461 the NB classification technique is superior to the widely used SVM. As model performance differs vastly
462 across datasets, this study raises the question which factors cause models to yield more workload savings for
463 some systematic review datasets than for others. In order to facilitate the applicability of active learning
464 models in systematic review practice, it is essential to identify how dataset characteristics relate to model
465 performance.

466 **Declarations**

467 **List of abbreviations**

468 ATD - Average Time to Discovery

469 D2V - doc2vec

470 LR - Logistic regression

471 MAD - Median Absolute Deviation

472 ML - Machine Learning

473 NB - Naive Bayes

474 PTSD - Post Traumatic Stress Disorder

475 RF - Random forest

476 RRF - Relevant References Found

477 SD - Standard Deviation

478 SEM - Standard Error of the Mean

479 SVM - Support vector machines

480 TF-IDF - Term Frequency - Inverse Document Frequency

481 TPE - Tree of Parzen Estimators

482 TD - Time to Discovery

483 WSS - Work Saved over Sampling

484 **Ethics approval and consent to participate**

485 This study has been approved by the Ethics Committee of the Faculty of Social and Behavioural Sciences of
486 Utrecht University, filed as an amendment under study 20-104.

487 **Consent for publication**

488 Not applicable.

489 **Availability of data and materials**

490 All data and materials are stored in the GitHub repository for this paper, [https://github.com/asreview/paper-](https://github.com/asreview/paper-evaluating-models-across-research-areas)
491 [evaluating-models-across-research-areas](https://github.com/asreview/paper-evaluating-models-across-research-areas). This repository contains all systematic review datasets used during
492 this study and their preprocessing scripts, scripts for the hyperparameter optimization, the simulations, the
493 processing and analysis of the results of the simulations, and for the figures and tables in this paper. The raw
494 output files of the simulation study are stored on the Open Science Framework, <https://osf.io/7mr2g/> and
495 <https://osf.io/ag2xp/>.

496 **Competing interests**

497 The authors declare that they have no competing interests.

498 **Funding**

499 This project was funded by the Innovation Fund for IT in Research Projects, Utrecht University, The
500 Netherlands. Access to the Cartesius supercomputer was granted by SURFsara (ID EINF-156). Both the
501 Innovation Fund and SURFsara had no role whatsoever in the design of the current study, nor in the data
502 collection, analysis and interpretation, nor in writing the manuscript.

503 **Author's contributions**

504 RvdS, RS, JdB and GF designed the study. RS developed the DR balance strategy and ATD metric, and
505 wrote the programs required for hyperparameter optimization and cloud computation. RS, JdB, and DO
506 designed the architecture required for the simulation study. GF extracted and analyzed the data and drafted
507 the manuscript. RvdS, AB, RS, DO, LT, and JdB assisted with writing the paper. LT, DO, AB, and RvdS
508 provided domain knowledge. All authors read and approved the final manuscript.

509 **Acknowledgements**

510 We are grateful for all researchers who have made great efforts to openly publish the data on their systematic
 511 reviews, special thanks go out to Rosanna Nagtegaal.

Figures

Figure 1: Recall curves of different models for the Nudging dataset, indicating how fast the model finds relevant publications during the process of screening publications. Figure a displays curves for all seven models at once. Figures b to f display curves for several subsets of those models to allow for a more detailed inspection of model performance.

Figure 2: Recall curves of all seven models for (a) the PTSD, (b) Software, (c) ACE, (d) Virus, and (e) Wilson dataset.

Tables

Table 1: Statistics on the datasets obtained from six original systematic reviews.

Dataset	Candidate publications	Relevant publications	Proportion relevant (%)
Nudging	1,847	100	5.4
PTSD	5,031	38	0.8
Software	8,896	104	1.2
ACE	2,235	41	1.8
Virus	2,304	114	5.0
Wilson	2,333	23	1.0

Table 2: ATD values ($\bar{x}(\hat{s})$) for all model-dataset combinations. For every dataset, the best results are in bold. Median (MAD) is given for all datasets.

	Nudging	PTSD	Software	ACE	Virus	Wilson
SVM + TF-IDF	10.1 (0.18)	2.1 (0.13)	1.9 (0.04)	7.1 (1.15)	8.5 (0.17)	4.0 (0.32)
NB + TF-IDF	9.3 (0.29)	1.7 (0.11)	1.4 (0.03)	4.9 (0.51)	8.2 (0.22)	3.9 (0.35)
RF + TF-IDF	11.7 (0.44)	3.3 (0.26)	2.0 (0.09)	6.8 (0.74)	10.5 (0.42)	5.6 (1.15)
LR + TF-IDF	9.5 (0.19)	1.7 (0.10)	1.4 (0.01)	5.9 (1.17)	8.3 (0.24)	4.3 (0.32)
SVM + D2V	8.8 (0.33)	2.1 (0.15)	1.4 (0.05)	6.1 (0.33)	8.4 (0.21)	4.5 (0.30)
RF + D2V	10.3 (0.87)	3.0 (0.33)	1.6 (0.09)	7.2 (1.26)	9.2 (0.43)	7.2 (1.49)
LR + D2V	8.8 (0.47)	1.9 (0.16)	1.4 (0.04)	5.4 (0.18)	8.3 (0.40)	4.7 (0.30)
median (MAD)	9.5 (1.05)	2.1 (0.48)	1.4 (0.12)	6.1 (1.11)	8.4 (0.18)	4.5 (0.64)

Table 3: WSS@95 values ($\bar{x}(\hat{s})$) for all model-dataset combinations. For every dataset, the best results are in bold. Median (MAD) is given for all datasets.

	Nudging	PTSD	Software	ACE	Virus	Wilson
SVM + TF-IDF	66.2 (2.90)	91.0 (0.41)	92.0 (0.10)	75.8 (1.95)	69.7 (0.81)	79.9 (2.09)
NB + TF-IDF	71.7 (1.37)	91.7 (0.27)	92.3 (0.08)	82.9 (0.99)	71.2 (0.62)	83.4 (0.89)
RF + TF-IDF	64.9 (2.50)	84.5 (3.38)	90.5 (0.34)	71.3 (4.03)	63.9 (3.54)	81.6 (3.35)
LR + TF-IDF	66.9 (4.01)	91.7 (0.18)	92.0 (0.10)	81.1 (1.31)	70.3 (0.65)	80.5 (0.65)
SVM + D2V	70.9 (1.68)	90.6 (0.73)	92.0 (0.21)	78.3 (1.92)	70.7 (1.76)	82.7 (1.44)
RF + D2V	66.3 (3.25)	88.2 (3.23)	91.0 (0.55)	68.6 (7.11)	67.2 (3.44)	77.9 (3.43)
LR + D2V	71.6 (1.66)	90.1 (0.63)	91.7 (0.13)	77.4 (1.03)	70.4 (1.34)	84.0 (0.77)
median (MAD)	66.9 (3.05)	90.6 (1.53)	92.0 (0.47)	77.4 (5.51)	70.3 (0.90)	81.6 (2.48)

Table 4: RRF@10 values (\bar{x} , (\hat{s})) for all model-dataset combinations. For every dataset, the best results are in bold. Median (MAD) is given for all datasets.

	Nudging	PTSD	Software	ACE	Virus	Wilson
SVM + TF-IDF	60.2 (3.12)	98.6 (1.40)	99.0 (0.00)	86.2 (5.25)	73.4 (1.62)	90.6 (1.17)
NB + TF-IDF	65.3 (2.61)	99.6 (0.95)	98.2 (0.34)	90.5 (1.40)	73.9 (1.70)	87.3 (2.55)
RF + TF-IDF	53.6 (2.71)	94.8 (1.60)	99.0 (0.00)	82.3 (2.75)	62.1 (3.19)	86.7 (5.82)
LR + TF-IDF	62.1 (2.59)	99.8 (0.70)	99.0 (0.00)	88.5 (5.16)	73.7 (1.48)	89.1 (2.30)
SVM + D2V	67.3 (3.00)	97.8 (1.12)	99.3 (0.44)	84.2 (2.78)	73.6 (2.54)	91.5 (4.16)
RF + D2V	62.6 (5.47)	97.1 (1.90)	99.2 (0.34)	80.8 (5.72)	67.3 (3.19)	75.5 (14.35)
LR + D2V	67.5 (2.59)	98.6 (1.40)	99.0 (0.00)	81.7 (1.81)	70.6 (2.21)	90.6 (5.00)
median (MAD)	62.6 (3.89)	98.6 (1.60)	99.0 (0.00)	84.2 (3.71)	73.4 (0.70)	89.1 (2.70)

Additional files

Additional file 1 — The active learning cycle

additional-file-1-active-learning-cycle.pdf. Description: The active learning cycle for screening prioritization in systematic reviews.

Additional file 2 — Dynamic Resampling

additional-file-2-DR.pdf. Description: Algorithm describing how to rebalance training data by the Dynamic Resampling (DR) strategy.

Additional file 3 — Average Time to Discovery

additional-file-3-ATD.pdf. Description: Definition of the Average Time to Discovery (ATD), a metric to assess the model performance.

Additional file 4 — Recall curves

additional-file-4-recall-curves.pdf. Description: Various subsets of recall curves for the PTSD, Software, ACE, Virus, and Wilson datasets, like Figure 1 presents curves for the Nudging dataset.

References

- [1] PRISMA-P Group, Moher D, Shamseer L, Clarke M, Ghersi D, Liberati A, et al. Preferred Reporting Items for Systematic Review and Meta-Analysis Protocols (PRISMA-P) 2015 Statement. *Syst Rev.* 2015;4(1):1. Available from: <https://systematicreviewsjournal.biomedcentral.com/articles/10.1186/2046-4053-4-1>.
- [2] Gough D, Richardson M. Systematic Reviews. In: *Advanced Research Methods for Applied Psychology*. Routledge; 2018. p. 75–87.
- [3] Chalmers I. The Lethal Consequences of Failing to Make Full Use of All Relevant Evidence about the Effects of Medical Treatments: The Importance of Systematic Reviews. In: *Treating Individuals - from Randomised Trials to Personalised Medicine*. *Lancet*; 2007. p. 37–58.
- [4] Borah R, Brown AW, Capers PL, Kaiser KA. Analysis of the Time and Workers Needed to Conduct Systematic Reviews of Medical Interventions Using Data from the PROSPERO Registry. *BMJ Open.* 2017;7(2):e012545. Available from: <https://bmjopen-bmj-com.proxy.library.uu.nl/content/7/2/e012545>.
- [5] Lau J. Editorial: Systematic Review Automation Thematic Series. *Syst Rev.* 2019;8(1):70. Available from: <https://doi.org/10.1186/s13643-019-0974-z>.
- [6] Harrison H, Griffin SJ, Kuhn I, Usher-Smith JA. Software Tools to Support Title and Abstract Screening for Systematic Reviews in Healthcare: An Evaluation. *BMC Med Res Methodol.* 2020;20(1):7. Available from: <https://doi.org/10.1186/s12874-020-0897-3>.
- [7] O'Mara-Eves A, Thomas J, McNaught J, Miwa M, Ananiadou S. Using Text Mining for Study Identification in Systematic Reviews: A Systematic Review of Current Approaches. *Syst Rev.* 2015;4(1):5. Available from: <https://doi.org/10.1186/2046-4053-4-5>.
- [8] Cohen AM, Ambert K, McDonagh M. Cross-Topic Learning for Work Prioritization in Systematic Review Creation and Update. *J Am Med Inform Assoc.* 2009;16(5):690–704. Available from: <https://academic-oup-com.proxy.library.uu.nl/jamia/article/16/5/690/804676>.
- [9] Shemilt I, Simon A, Hollands GJ, Marteau TM, Ogilvie D, O'Mara-Eves A, et al. Pinpointing Needles in Giant Haystacks: Use of Text Mining to Reduce Impractical Screening Workload in Extremely Large

- Scoping Reviews. *Res Synth Methods*. 2014;5(1):31–49. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jrsm.1093>.
- [10] Yu Z, Menzies T. FAST2: An Intelligent Assistant for Finding Relevant Papers. *Expert Syst Appl*. 2019;120:57–71. Available from: <http://www.sciencedirect.com/science/article/pii/S0957417418307413>.
- [11] Yu Z, Kraft NA, Menzies T. Finding Better Active Learners for Faster Literature Reviews. *Empir Softw Eng*. 2018;23(6):3161–3186. Available from: <http://dx.doi.org/10.1007/s10664-017-9587-0>.
- [12] Miwa M, Thomas J, O’Mara-Eves A, Ananiadou S. Reducing Systematic Review Workload through Certainty-Based Screening. *J Biomed Inform*. 2014;51:242–253. Available from: <http://www.sciencedirect.com/science/article/pii/S1532046414001439>.
- [13] Cormack GV, Grossman MR. Evaluation of Machine-Learning Protocols for Technology-Assisted Review in Electronic Discovery. In: *Proceedings of the 37th International ACM SIGIR Conference on Research & Development in Information Retrieval. SIGIR ’14*. Association for Computing Machinery; 2014. p. 153–162. Available from: <https://doi.org/10.1145/2600428.2609601>.
- [14] Cormack GV, Grossman MR. Autonomy and Reliability of Continuous Active Learning for Technology-Assisted Review. *CoRR*. 2015;abs/1504.06868. Available from: <http://arxiv.org/abs/1504.06868>.
- [15] Wallace BC, Trikalinos TA, Lau J, Brodley C, Schmid CH. Semi-Automated Screening of Biomedical Citations for Systematic Reviews. *BMC Bioinform*. 2010;11(1):55. Available from: <https://doi.org/10.1186/1471-2105-11-55>.
- [16] Gates A, Johnson C, Hartling L. Technology-Assisted Title and Abstract Screening for Systematic Reviews: A Retrospective Evaluation of the Abstrackr Machine Learning Tool. *Syst Rev*. 2018;7(1):45. Available from: <https://systematicreviewsjournal.biomedcentral.com/articles/10.1186/s13643-018-0707-8>.
- [17] Settles B. Active Learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*. 2012;6(1):1–114. Available from: <http://www.morganclaypool.com/doi/abs/10.2200/S00429ED1V01Y201207AIM018>.
- [18] Settles B. Active learning literature survey. University of Wisconsin-Madison Department of Computer Sciences; 2009.
- [19] Singh G, Thomas J, Shawe-Taylor J. Improving active learning in systematic reviews. *arXiv preprint arXiv:180109496*. 2018;.

- [20] Carvallo A, Parra D. Comparing Word Embeddings for Document Screening based on Active Learning. In: BIRNDL@ SIGIR; 2019. p. 100–107.
- [21] Ma Y. Text classification on imbalanced data: Application to Systematic Reviews Automation. University of Ottawa (Canada). 2007;.
- [22] Wallace BC, Small K, Brodley CE, Lau J, Trikalinos TA. Deploying an Interactive Machine Learning System in an Evidence-Based Practice Center: Abstract. In: Proceedings of the 2nd ACM SIGHT International Health Informatics Symposium. IHI '12. Association for Computing Machinery; 2012. p. 819–824. Available from: <https://doi.org/10.1145/2110363.2110464>.
- [23] Cheng SH, Augustin C, Bethel A, Gill D, Anzaroot S, Brun J, et al. Using Machine Learning to Advance Synthesis and Use of Conservation and Environmental Evidence. *Conserv Biol*. 2018;32(4):762–764. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1111/cobi.13117>.
- [24] Ouzzani M, Hammady H, Fedorowicz Z, Elmagarmid A. Rayyan—a Web and Mobile App for Systematic Reviews. *Syst Rev*. 2016;5(1):210. Available from: <http://dx.doi.org/10.1186/s13643-016-0384-4>.
- [25] Przybył a P, Brockmeier AJ, Kontonatsios G, Pogam MAL, McNaught J, Erik von Elm, et al. Prioritising References for Systematic Reviews with RobotAnalyst: A User Study. *Res Synth Methods*. 2018;9(3):470–488. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jrsm.1311>.
- [26] Kilicoglu H, Demner-Fushman D, Rindfleisch TC, Wilczynski NL, Haynes RB. Towards Automatic Recognition of Scientifically Rigorous Clinical Research Evidence. *J Am Med Inform Assn*. 2009;16(1):25–31. Available from: <https://academic.oup.com/jamia/article-lookup/doi/10.1197/jamia.M2996>.
- [27] Aphinyanaphongs Y. Text Categorization Models for High-Quality Article Retrieval in Internal Medicine. *J Am Med Inform Assoc*. 2004;12(2):207–216. Available from: <https://academic.oup.com/jamia/article-lookup/doi/10.1197/jamia.M1641>.
- [28] Aggarwal CC, Zhai C. A Survey of Text Classification Algorithms. In: Aggarwal CC, Zhai C, editors. *Mining Text Data*. Springer US; 2012. p. 163–222. Available from: http://link.springer.com/10.1007/978-1-4614-3223-4_6.
- [29] Zhang W, Yoshida T, Tang X. A Comparative Study of TF*IDF, LSI and Multi-Words for Text Classification. *Expert Syst Appl*. 2011;38(3):2758–2765. Available from: <http://www.sciencedirect.com/science/article/pii/S0957417410008626>.
- [30] Le Q, Mikolov T. Distributed representations of sentences and documents. In: *International conference on machine learning*; 2014. p. 1188–1196.

- [31] Marshall IJ, Johnson BT, Wang Z, Rajasekaran S, Wallace BC. Semi-Automated Evidence Synthesis in Health Psychology: Current Methods and Future Prospects. *Health Psychol Rev.* 2020;14(1):145–158. Available from: <https://doi.org/10.1080/17437199.2020.1716198>.
- [32] Cohen AM, Hersh WR, Peterson K, Yen PY. Reducing Workload in Systematic Review Preparation Using Automated Citation Classification. *J Am Med Inform Assoc.* 2006;13(2):206–219. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1447545/>.
- [33] Appenzeller-Herzog C, Mathes T, Heeres MLS, Weiss KH, Houwen RHJ, Ewald H. Comparative Effectiveness of Common Therapies for Wilson Disease: A Systematic Review and Meta-Analysis of Controlled Studies. *Liver Int.* 2019;39(11):2136–2152. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1111/liv.14179>.
- [34] Kwok KTT, Nieuwenhuijse DF, Phan MVT, Koopmans MPG. Virus Metagenomics in Farm Animals: A Systematic Review. *Viruses.* 2020;12(1):107. Available from: <https://www.mdpi.com/1999-4915/12/1/107>.
- [35] Nagtegaal R, Tummers L, Noordegraaf M, Bekkers V. Nudging Healthcare Professionals towards Evidence-Based Medicine: A Systematic Scoping Review. *J Behav Public Adm.* 2019;2(2).
- [36] van de Schoot R, Sijbrandij M, Winter SD, Depaoli S, Vermunt JK. The GRoLTS-Checklist: Guidelines for Reporting on Latent Trajectory Studies. *Struct Equ Model Multidiscip J.* 2017;24(3):451–467. Available from: <https://doi.org/10.1080/10705511.2016.1247646>.
- [37] van de Schoot R, de Bruin J, Schram R, Zahedi P, de Boer J, Weijdema F, et al. ASReview: Open Source Software for Efficient and Transparent Active Learning for Systematic Reviews. *arXiv preprint arXiv:200612166.* 2020;.
- [38] Wynants L, Calster BV, Collins GS, Riley RD, Heinze G, Schuit E, et al. Prediction Models for Diagnosis and Prognosis of Covid-19: Systematic Review and Critical Appraisal. *BMJ.* 2020;369. Available from: <http://www.bmj.com/content/369/bmj.m1328>.
- [39] Tong S, Koller D. Support Vector Machine Active Learning with Applications to Text Classification. *J Mach Learn Res.* 2001;2:45–66.
- [40] Kremer J, Steenstrup Pedersen K, Igel C. Active Learning with Support Vector Machines. *WIREs Data Min Knowl Discov.* 2014;4(4):313–326. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/widm.1132>.

- [41] Zhang H. The Optimality of Naive Bayes. In: Proceedings of the Seventeenth International Florida Artificial Intelligence Research Society Conference, FLAIRS 2004. vol. 2; 2004. p. 562–567.
- [42] Breiman L. Random Forests. *Machine Learning*. 2001;45(1):5–32. Available from: <https://doi.org/10.1023/A:1010933404324>.
- [43] Ramos J, et al. Using Tf-Idf to Determine Word Relevance in Document Queries. In: Proceedings of the First Instructional Conference on Machine Learning. vol. 242. Piscataway, NJ; 2003. p. 133–142.
- [44] Shahmirzadi O, Lugowski A, Younge K. Text Similarity in Vector Space Models: A Comparative Study. In: 2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA); 2019. p. 659–666. Available from: [10.1109/ICMLA.2019.00120](https://doi.org/10.1109/ICMLA.2019.00120).
- [45] Yang Y, Loog M. A Benchmark and Comparison of Active Learning for Logistic Regression. *Pattern Recognition*. 2018;83:401–415. Available from: <http://www.sciencedirect.com/science/article/pii/S0031820318302140>.
- [46] Fu JH, Lee SL. Certainty-Enhanced Active Learning for Improving Imbalanced Data Classification. In: 2011 IEEE 11th International Conference on Data Mining Workshops. IEEE; 2011. p. 405–412. Available from: <http://ieeexplore.ieee.org/document/6137408/>.
- [47] van de Schoot R, de Bruin J, Schram R, Zahedi P, Kramer B, Ferdinands G, et al. ASReview: Active Learning for Systematic Reviews. Zenodo; 2020.
- [48] Bergstra J, Yamins D, Cox D. Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures. In: International Conference on Machine Learning; 2013. p. 115–123. Available from: <http://proceedings.mlr.press/v28/bergstra13.html>.
- [49] R Core Team. R Foundation for Statistical Computing, editor. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2019. Available from: <https://www.R-project.org/>.
- [50] Appenzeller-Herzog C. Data from Comparative Effectiveness of Common Therapies for Wilson Disease: A Systematic Review and Meta-analysis of Controlled Studies. Zenodo. 2020; Available from: <https://doi.org/10.5281/zenodo.3625931>.
- [51] Hall T, Beecham S, Bowes D, Gray D, Counsell S. A Systematic Literature Review on Fault Prediction Performance in Software Engineering. *IEEE Trans Softw Eng*. 2012;38(6):1276–1304.
- [52] Nagtegaal R, Tummers L, Noordegraaf M, Bekkers V. Nudging healthcare professionals towards

evidence-based medicine: A systematic scoping review. Harvard Dataverse. 2019; Available from: <https://doi.org/10.7910/DVN/WMGPGZ>.

[53] Mitchell TM. Does Machine Learning Really Work? *AI Mag.* 1997;18(3):11–11. Available from: <https://www.aaai.org/ojs/index.php/aimagazine/article/view/1303>.

[54] Rathbone J, Hoffmann T, Glasziou P. Faster Title and Abstract Screening? Evaluating Abstrackr, a Semi-Automated Online Screening Program for Systematic Reviewers. *Systematic Reviews.* 2015;4(1):80. Available from: <https://doi.org/10.1186/s13643-015-0067-6>.