

1 The use of text-mining and machine learning
2 algorithms in systematic reviews: reducing
3 workload in preclinical biomedical sciences and
4 reducing human screening error

5 Bannach-Brown, A.^{1,5} (<https://orcid.org/0000-0002-3161-1395>), Przybyła, P.,² Thomas, J.³
6 (<https://orcid.org/0000-0003-4805-4190>), Rice, A.S.C.⁴ (<https://orcid.org/0000-0002-1237-4769>)
7 , Ananiadou, S.², Liao, J.¹, Macleod, M.R.¹ (<https://orcid.org/0000-0001-9187-9839>)

8

9 1 Centre for Clinical Brain Sciences, University of Edinburgh.

10 2 National Centre for Text Mining, School of Computer Science, University of Manchester.

11 3 EPPI-Centre, Department of Social Science, University College London.

12 4 Pain Research, Department of Surgery and Cancer, Imperial College.

13 5 Translational Neuropsychiatry Unit, Aarhus University.

14

15

16 Abstract:

17 Background: In this paper we outline a method of applying machine learning (ML) algorithms to aid
18 citation screening in an on-going broad and shallow systematic review, with the aim of achieving a
19 high performing algorithm comparable to human screening.

20 Methods: We tested a range of machine learning algorithms. We applied ML algorithms to
21 incremental numbers of training records and recorded the performance on sensitivity and specificity
22 on an unseen validation set of papers. The performance of these algorithms was assessed on measures
23 of recall, specificity, and accuracy. The classification results of the best performing algorithm was
24 taken forward and applied to the remaining unseen records in the dataset and will be taken forward
25 to the next stage of systematic review. ML was used to identify potential human errors during
26 screening by analysing the training and validation datasets against the machine-ranked score.

27 Results: We found that ML algorithms perform at a desirable level. Classifiers reached 98.7%
28 sensitivity based on learning from a training set of 5749 records, with an inclusion prevalence of 13.2%
29 (see below). The highest level of specificity reached was 86%. Human errors in the training and
30 validation set were successfully identified using ML scores to highlight discrepancies. Training the ML
31 algorithm on the corrected dataset improved the specificity of the algorithm without compromising
32 sensitivity. Error analysis sees a 3% increase or change in sensitivity and specificity, which increases
33 precision and accuracy of the ML algorithm.

34 Conclusions: The technique of using ML to identify human error needs to be investigated in more
35 depth, however this pilot shows a promising approach to integrating human decisions and automation
36 in systematic review methodology.

37

38 Key-words: machine learning, systematic review, analysis of human error, citation screening,
39 automation tools

40 Introduction:

41 The rate of publication is increasing exponentially within biomedical research [1]. Researchers are
42 finding it increasingly difficult to keep up with new findings and discoveries even within one
43 biomedical domain, an issue that has been emerging for a number of years [2]. Synthesising research
44 – either informally or through systematic reviews to provide an unbiased summary also becomes an
45 increasingly resource intensive task as search strings retrieve larger and larger corpuses of potentially
46 relevant papers for reviewers to screen for relevance to the research question at hand.

47 Within the animal model literature we see this increase in rate of publication. In an update to a
48 systematic review of animal models of neuropathic pain, 11,880 further unique records were retrieved
49 in 2015, in addition to the 33,184 unique records identified initially in the search in 2012. In the field
50 of animal models of depression, the amount of unique records retrieved from the search has increased
51 in 1 year from 70,365 (May 2016) to 76,679 (August 2017).

52 The use of text-mining tools and machine learning (ML) algorithms to aid systematic review is
53 becoming an increasingly popular tool to reduce human burden and monetary resources required and
54 increase efficiency to complete systematic reviews as well as increase the speed at which results are
55 produced [3; 4; 5]. ML algorithms are primarily employed at the screening stage in the systematic
56 review process. This screening stage involves categorising records identified from the search into
57 ‘Relevant’ or ‘Not-Relevant’ to the research question, typically carried out to the gold standard by two
58 independent human reviewers. This decision is typically made using the title and abstract of an article
59 in the first instance. In previous experience at CAMARADES, screening a preclinical systematic review
60 with 33,184 independent hits took 18 person months in total, which can be divided between the
61 number of reviewers available on a project. With two reviewers working concurrently, it took 9
62 months to screen 33,184 records. Based on this, we have estimated that a systematic review with
63 roughly 10,000 publications retrieved takes an estimated 40 weeks. In clinical systematic reviews,
64 Borah and colleagues [6] showed the average clinical systematic review registered on PROSPERO takes

65 an average 67.3 weeks to complete. ML algorithms are employed to learn this categorisation ability,
66 based on training instances that have been screened by human reviewers. This approach is an example
67 of using supervised machine learning for automatic classification in systematic review screening,
68 where the algorithm learns the classification decision based on example information given by human
69 screeners. Other machine learning approaches for screening in systematic reviews include supervised
70 learning for document prioritisation, often used for rapid evidence assessments, where the most
71 potentially relevant documents are presented to reviewers first for screening [7; 5]. Further, active
72 learning, which involves the algorithm selecting data that it is trained on, as opposed to a random
73 sample. The concept being that the classifier presents the expert reviewer with the next record that
74 it ‘thinks’ it will learn most from [8; 9]. Different approaches can be taken on how to use the machine
75 learning decision, whether, with high enough performance, review teams accept the decision of the
76 machine, or whether to use the machine as a second screener.

77 In this paper we outline the approach taken to screen a novel corpus for a broad and shallow
78 systematic review question, non-human animal models of depression, a corpus of 70,365 records
79 retrieved from two online biomedical databases. In the context of systematic review, ML algorithms
80 have often been tested in projects, where a large corpus of ‘already-classified’ data is available, for
81 example in the case of updating a systematic review [3]. *Here, the aim was to identify the amount of*
82 *training data required for an algorithm to achieve the level of performance equivalent to the gold*
83 *standard of two independent human screeners.* The goal performance for the ML algorithms was
84 performance at the human gold standard of 95% sensitivity and to maximise performance for
85 specificity.

86 Sena and colleagues developed guidelines for the appraisal of systematic reviews of animal studies
87 [10]. These guidelines consider dual extraction or two independent human reviewers as a feature of
88 a high quality review. The inter-screener agreement in systematic reviews at CAMARADES is estimated
89 to be between 95% and 99%. This leaves a 1-5% human error rate. This can include both random error

90 due to fatigue or distraction of human screeners, or – more consequentially - a systematic error,
91 which, if included in a training set, might be propagated into a ML algorithm. Sources of systematic
92 error at the screening stage could be if there are differences in the way screeners interpret inclusion
93 criteria. As far as we are aware, the nature of this 5% residual human error in systematic review
94 methodology has not been formally investigated. The training data used for ML categorisation is based
95 on training instances that has been screened by two independent human screeners. *We therefore*
96 *aimed to explore the use of ML algorithms as part of a systematic review framework at the*
97 *classification stage, to investigate if the ML algorithms could be used to improve the human gold*
98 *standard by identifying human screening errors and thus improve the overall performance of ML.*

99

100 **Methods:**

101 The methodology of this paper is outlined as follows. Firstly, different machine learning algorithms
102 were applied to assist in the screening of a large systematic review (> 70,000 records identified).
103 Three sub-sets of the 70,365 records were screened by two human reviewers. These sub-sets were
104 used to train the algorithms. Performance of the classifiers was assessed on a validation set of
105 unseen documents, and a number of different metrics were used. Secondly, the best performing
106 algorithm was used to identify human error in the training and validation sets. The error analysis is
107 assessed on the net reclassification index.

108

109 **Step 1: Application of ML tools to screening of a large preclinical systematic review.**

110 **Training Sets:**

111 70,365 potentially relevant records were identified from Pubmed and EMBASE (for search strings see
112 [11]). The training sets were chosen at random from the 70,365 by assigning each record a random
113 number using the RAND function in excel and ranking them from smallest to largest. The first 1993

114 records were chosen for the first training set. The second training set consisted of adding the next 996
115 records, bringing the training set size to 2989. The third training set consisted of the adding the next
116 2760 records, bringing the total to 5749 records. These datasets were intended to be round numbers
117 (i.e.2000 and 1000), however data disposition errors occurred when transferring the records between
118 databases, hence the non-even numbers of records in each sub-dataset. The total number of training
119 records was 5749. The training sets were screened by two independent human screeners. (Datasets
120 are available on Zenodo, see “Availability of Data & Materials” below). Performance was assessed at
121 each level on a validation set of unseen records. All training and validation sets were selected
122 consecutively from the initial random ordering. The first training set had a validation set of 996
123 records. The second training set had a validation set of 1010 records. And for the full training set of
124 5749 records, the validation set was the subsequent 1251 records. This validation set had more than
125 150 “included” records, which can give reasonably precise 95% confidence intervals for sensitivity and
126 specificity.

127

128 **Feature Generation:**

129 First, documents in the training set were transformed into a representation that is appropriate for the
130 machine learning algorithms. Namely, every case (document) is represented by a fixed number of
131 features, i.e., numerical quantities describing certain properties that could be used by the classifier to
132 extract rules and make predictions about it. The classifiers described below used similar approaches,
133 though there were differences.

134 The “bag-of-words” model was used to characterise document features – their titles and abstracts –
135 in both classifiers. This approach to translating text to a numeric representation is commonly used in
136 document classification and, at its most basic, is simply the frequency count of each word in a
137 document. In order to account not only for the relative importance of words within a document, ‘Term
138 Frequency – Inverse Document Frequency – is often used. This helps the classifier to focus on terms

139 which help to distinguish between documents, rather than on terms which occur frequently [12].
140 Technically, each word is known as a ‘gram’ – and the simple indexing of individual words would be
141 described as being a ‘uni-gram’ bag of words. ‘n-gram’ bag of words are usually used – where ‘n’
142 indicates the number of words that are included (e.g. ‘bi-gram’ bag of words approaches index every
143 word, and also every pair of words). Lastly ‘stopword’ lists are used, to remove frequently occurring
144 words which provide little relevant information for classification purposes.

145 Bag of words representation is straightforward and easy to implement, but has a serious drawback: it
146 generates as many unigram features as there are words in the collection (typically at least several
147 thousand), and much more when using higher-order n-grams. This is known to be a problem for
148 general-purpose classifiers, which have to be extended with additional feature selection stages.
149 Moreover, this multi-dimensionality is related to high redundancy, as quite often different words refer
150 to the same underlying concept, e.g. in case of synonymy or morphological variants. That is why
151 several techniques have been developed for representing textual data in a more efficient way.

152 One such option is Latent Semantic Indexing (LSI) [13]. In LSI the training set is represented as a matrix,
153 where rows correspond to documents, columns to terms (words or n-grams), while cells contain
154 frequency or TF/IDF score of a given term in a given document. The matrix is then decomposed using
155 a general matrix factorisation technique: Singular Value Decomposition (SVD) and truncated to the
156 first n dimensions. Thanks to the properties of SVD, the new features will be such linear combinations
157 of features of the old space that minimise the differences between the original and the transformed
158 space. In case of textual data it means that those words that frequently occur in the same documents
159 (probably because of the similar meaning) will be treated in the same way. The n is set a-priori to a
160 reasonably low value – usually a few hundred.

161 Latent Dirichlet Allocation (LDA) is another method aiming to provide lower-dimensional
162 representation of documents by exploiting distributional similarities between words, but based on
163 explaining document contents using a Bayesian network [14]. This method is based on the premise

164 that every document is a mixture of topics, which in turn consist of related words. The correspondence
165 between documents and topics and between topics and words could be inferred via Gibbs sampling
166 process. As a result, similarly to LSI, every document is represented by a sequence of n numbers,
167 indicating how related it is to every topic [15]. As well as being used to identify distinct (and
168 overlapping) clusters of documents, topic model representation can be included as feature space for
169 document classification tasks, including screening [9]. Other related methods exploiting vocabulary
170 similarities for document representation include paragraph vectors and descriptive clustering [16; 17].

171

172

173 **Classifiers:**

174 Following the transformations made in feature selection, the documents are then used to train the
175 machine learning classifier. One of the most commonly used document classifiers is the Support
176 Vector Machine (SVM). SVM is a supervised learning algorithm, learning to classify new documents
177 based on a training set of labelled documents [18]. This algorithm represents training documents as
178 points in a feature space and seeks the optimal hyper-plane separating positive and negative cases,
179 i.e. included and excluded documents. The new, unseen, documents are then ranked by their distance
180 from the optimal hyper-plane. Logistic regression is a similar linear classifier, which instead of
181 hyperplane, seeks such coefficients of a linear combination of feature values that will give high values
182 for positive cases (included documents) and low for negative (excluded documents). Both of these
183 approaches could be enriched with feature selection elements to mitigate the problems with
184 multitude of features.

185

186 After testing a number of machine learning approaches, these two algorithms (below) performed the
187 best for this dataset of 70,365 records, on the broad topic of preclinical animal models of depression.

188

189 **Linear classifiers:**

190 Classifier 1, 3, & 5:

191 Classifier one, three, and five used a tri-gram ‘bag-of-words’ model for feature selection and used a
192 simple linear classifier. The classifier selected applied Stochastic Gradient Descent as implemented in
193 the SciKit-Learn python library [19]. This classifier was chosen as previous internal evaluations showed
194 its results to be indistinguishable from SVM; moreover, it is efficient, scales well to large numbers of
195 records, and provides an easily interpretable list of probability estimates when predicting class
196 membership (i.e. scores for each document lying between 0 and 1). Efficiency and interpretability are
197 important, as this classifier is deployed in a large systematic review platform [20], and any deployed
198 algorithm therefore needs not to be too computationally demanding, and its results understood by
199 users who are not machine learning specialists. The tri-gram feature selection approach without any
200 additional feature engineering also reflects the generalist need of deployment on a platform used in
201 a wide range of reviews: the algorithm needs to be generalisable across disciplines and literatures,
202 and not ‘over-fitted’ to a specific area. This approach aims to give the best compromise between
203 reliable performance across a wide range of domains and that achievable from a workflow that has
204 been highly tuned to a specific context.

205

206 Classifier 2 & 4:

207 Classifier two and four used a regularised logistic regression model built on LDA and SVD features.
208 Namely, the document text (consisting of title and abstract) was first lemmatised with GENIA tagger
209 [21] and then converted into bag of words representation of unigrams, which was then used to create
210 two types of features. First, the word frequencies were converted into a matrix TF/IDF scores, which
211 was then decomposed via SVD implemented in scikit-learn library and truncated to the first 300

212 dimensions. Second, an LDA model was built using MALLET library [22], setting 300 as a number of
213 topics. As a result each document was represented by 600 features, and an L1-regularised logistic
214 regression model was built using glmnet package [23] in R statistical framework [24].

215 Thanks to this procedure every document is represented with a constant, manageable number of
216 features, no matter what corpus or vocabulary size is. As a result, we can use a relatively simple
217 classification algorithm and obtain good performance with short processing time even for very large
218 collections. This feature is particularly useful when running the procedure numerous times in cross-
219 validation mode for error analysis (see below).

220

221 Algorithms assign each record a score of predicted probability of being relevant and an optimal cut-
222 off score that gives the best performance is calculated.

223

224 **Assessing Machine Learning Performance:**

225 The facets of a machine learning algorithm performance that would be beneficial to this field of
226 research are high sensitivity, comparable to the current gold standard, two independent human
227 screeners, which we estimate at 95%. We accept an error rate of 5%, missing out on 5% of relevant
228 studies, is permissible. We therefore aim to achieve at least 95% sensitivity, including the lower 95%
229 confidence bound. Once the level of sensitivity has been reached, the aim is to maximise specificity,
230 in order to reduce the amount of irrelevant records included by an algorithm.

231 *Performance metrics:*

232 Performance was assessed using sensitivity (or recall), specificity, precision, accuracy, and Work Saved
233 over Sampling (WSS) (see table 1), carried out in R (R version 3.4.2; [24]) using the 'caret' package [25].
234 95% Confidence Intervals were calculated using the efficient-score method [26].

235

Table 1. Equations used to assess performance of machine learning algorithms

Sensitivity or Recall	$TP / (TP+FN)$
Specificity	$TN / (TN+FP)$
Precision	$TP / (TP+FP)$
Accuracy	$(TP+TN) / (TP+FP+FN+TN)$
WSS@95%	$((TN+FN) / N) - (1.0 - 0.95)$

All equations from [5].

236

237

238 **Step 2: Application of ML tools to training datasets to identify human error.**

239 **Error Analysis Methods:**

240 The methodology for the error analysis was outlined in an *a priori* protocol, published on the
241 CAMARADES website 18th December 2016 [27]. To generate the machine learning scores for training
242 set records, the non-exhaustive cross-validation method, k-fold validation, was used. This method
243 involved randomly partitioning the test set into *k* number of equal sized subsamples. One subsample
244 was set aside for validation, and the remaining *k* -1 subsamples were used to train the algorithm [28].
245 Here *k* = 5. This process was carried out until all records in the training set had a machine learning
246 algorithm assigned score. The scores were used to draw attention to discrepancies or disagreements
247 between machine decision and human decision, by ranking the machine assigned labels in order of
248 predictive probability, from most likely to be relevant to least likely to be relevant, and highlighting
249 potential false positive and potential false negative cases, i.e. errors in the human decision. In order
250 to avoid reassessing the full 5749 record dataset for false positive and false negative results, a stopping
251 rule was established. If the initial human decision was correct five consecutive times, further false
252 positive, or false negative, records were not reassessed.

253 After the errors in the training set were investigated and corrected, a new model was built on the
254 updated training data. The outcome of error analysis is presented as reclassification tables and the
255 net reclassification index [29] is used to compare the performance of the classifier built on the

256 updated training data with the performance of the classifier built on the original training data. The
257 following equation was used:

$$258 \text{NRI}_{\text{binary outcomes}} = (\text{Sensitivity} + \text{Specificity})_{\text{second test}} - (\text{Sensitivity} + \text{Specificity})_{\text{first test}}$$

259 [30]

260 Further, we applied the same technique as above to identify human screening errors in the validation
261 dataset. Due to the small number of records in the validation set (1251 records), it was assumed that
262 every error would be likely to impact performance, therefore, the manual screening of the validation
263 set involved revisiting every record where the human and machine decision were incongruent. The
264 number of reclassified records was noted.

265

266 Results:

267 In this section we display the performance from the ML algorithms. We display the results from the
268 analysis of human error. And we show the performance of the ML algorithm after human errors in the
269 training and validation set have been corrected.

270

271 Performance of Machine Learning Algorithms

272 Table 2 shows the performance of the two best machine learning approaches from the SLIM
273 collaboration. The desired sensitivity of 95% has been reach by both classifiers. Both classifiers
274 reached 98.7% sensitivity based on learning from a training set of 5749 records, with an inclusion
275 prevalence of 13.2% (see below). Classifier 5 reached a higher specificity level of 86%.

276

277

Table 2. Performance of machine learning approaches on depression training dataset.

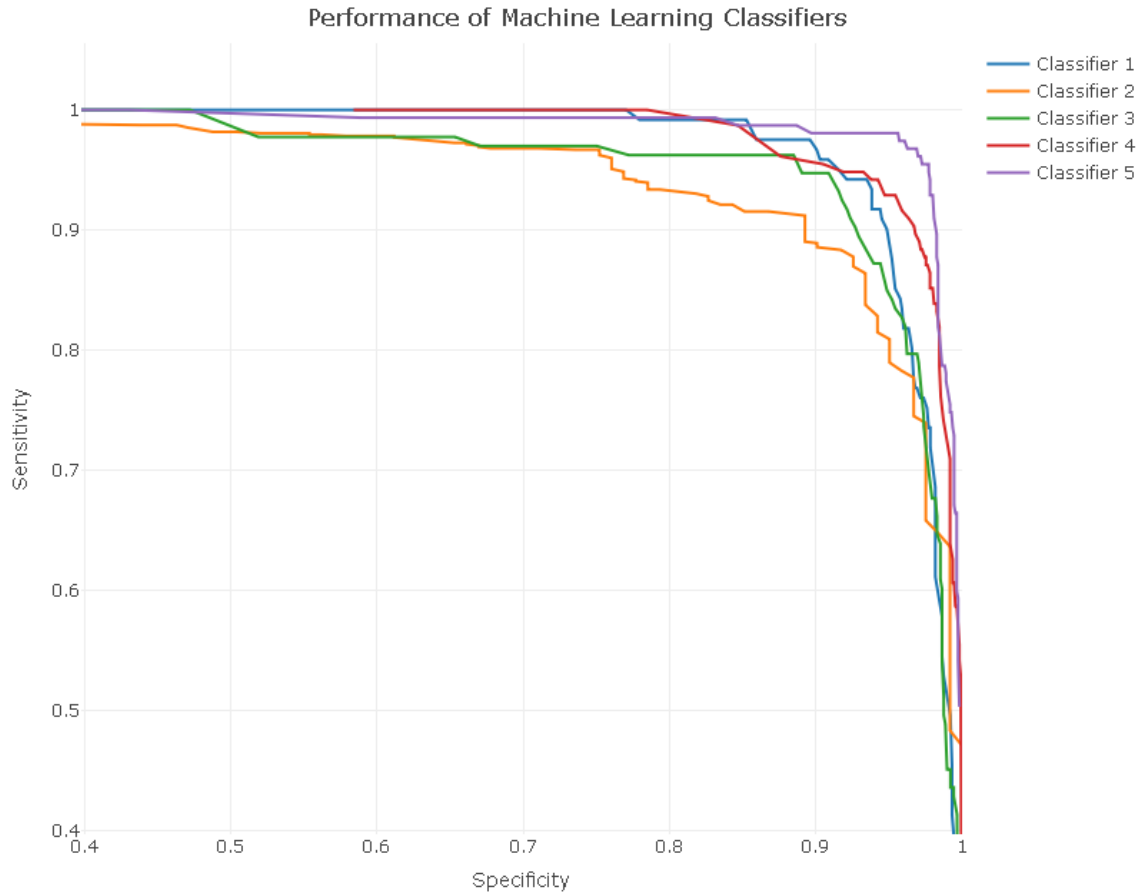
	Classifier 5	Classifier 4
Training Set Size	5749	5749
Optimal Cut-Off Score	0.1	0.07
Sensitivity	98.7%	98.7%
Upper 95% CI	0.997	0.997
Lower 95% CI	0.949	0.949
Specificity	86.0%	84.7%
Precision	50%	47.66%
Accuracy	1096/1251 = 87.6%	1081/1251= 86.4%
WSS@95%	0.705	0.693

278

279 The performance of the classifiers altered with increasing numbers of training records, and can be
280 visualised as a learning curve (figure 1). Here we see a trade-off between sensitivity and specificity.
281 With increasing amounts of records, specificity is increased for Classifier 5, the sensitivity increases
282 for Classifier 4 to reach above the desired 95%.

283

284 Figure 1. Performance of classifiers with increasing amounts of training records. Classifiers 1 & 2 are
285 trained on the first training set of 1993 records. Classifier 3 is trained on the second training set of
286 2989 records. Classifier 4 & 5 are trained on the third and full dataset of 5749 records. For the
287 interactive version of this plot with cut-off values, see code and data at
288 [https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-](https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-systematic-reviews/blob/master/Performance%20of%20Machine%20Learning%20Classifiers.html)
289 [systematic-reviews/blob/master/Performance%20of%20Machine%20Learning%20Classifiers.html](https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-systematic-reviews/blob/master/Performance%20of%20Machine%20Learning%20Classifiers.html)



290

291

292 Error Analysis & Reclassification

293 To assess whether machine learning algorithms can identify human error and therefore improve the

294 training data, error analysis was conducted. Seventy-five papers out of 5749 papers were identified

295 by the machine learning algorithm as potential human errors (i.e. where ML classification was correct,

296 human was wrong) and were rescreened by a human. Out of 75 rescreened papers, the machine

297 corrected the human decision 47 times. The machine was wrong, or the initial human decision was

298 correct, 28 times. The validation set was also rescreened. Ten papers out of the 1251 records were

299 identified as potential human errors. Out of 10 errors, the machine corrected 8 human decisions.

300 These 8 records were all falsely excluded by the human and were now included. The initial human

301 decision was correct twice.

302 To calculate human error in the training set, the number of errors identified (47) out of the training
303 set (5749 records) was calculated to be 0.8%. Of the 47 records reclassified, 11 records were falsely
304 included in the original screening process and were now correctly excluded, and 36 records were
305 falsely excluded in the original screening process and were now correctly included. The machine
306 correctly identified human screening errors, which were calculated to be just under 1% of the dual
307 screened training set. By looking at the prevalence of inclusion in this training set, which is 13.2% (760
308 out of the 5749), it highlights that it essential to correctly identify relevant papers. Therefore any
309 missing relevant papers, false negatives, or including irrelevant papers, false positives, will have a large
310 impact on the learning of the ML algorithm. Forty-seven papers out of 760 were 'correctly' reclassified,
311 6% of the included papers.

312 Similarly, the human error rate in the validation set (1251 records) was 0.6%. Again looking at the
313 prevalence of inclusion in this dataset (155/1251), which is 12.4%, the 8 records of out the now 163
314 were correctly reclassified which is 4.9% reclassified. All 8 records we falsely excluded in the original
315 screening process and are now correctly included.

316

317 Test 1: $98.7\% + 86\% = 184.7\%$

318 Test 2: $98.2\% + 89.3\% = 187.5\%$

319 **NRI = 3.2%**

320

321 We consider the updated validation set to be the new gold standard as 8 records were now included.
322 The confusion matrix for the performance of the machine learning algorithm after the error analysis
323 update on the training records is displayed below in table 3.

324

325

Table 3. Reclassification of records in validation after error analysis.

Test 1 – Original Machine Learning Algorithms results				
Test 2 – Post-error analysis ML results		In	Out	Total
	In	153	153	306
		160	116	276
	Out	2	943	945
	3	972	975	
Total	155	1096	1251	
	163	1088		

326

327 By analysing the human errors identified by the machine learning algorithm, correcting for these
328 errors and re-teaching the algorithm, you can see an increase in performance of the algorithm,
329 particularly sensitivity and thus precision. This can save considerable human time in the screening
330 stage of a systematic review. Consider the remaining approximately 64,000 papers, if the ML algorithm
331 results are 3% more accurate, that is approximately 2000 papers that are correctly 'excluded' that a
332 human reviewer does not need to screen through.

333

334 **After Error Analysis: Improving Machine Learning**

335 Using the error analysis technique above, of the 47 errors identified in the full training dataset of 5749
336 records, 0.8% were corrected. We retrained classifier 5 on the corrected training set and measured
337 performance on the corrected validation set of 1251 records as we consider this to be the 'new' gold
338 standard. The performance of the original classifier 5 and classifier 6 was assessed on the corrected
339 validation set of 1251 records. The performance of this retrained algorithm in comparison to the
340 performance of the original classifier 5 on the updated validation set is shown in table 4.

341

342

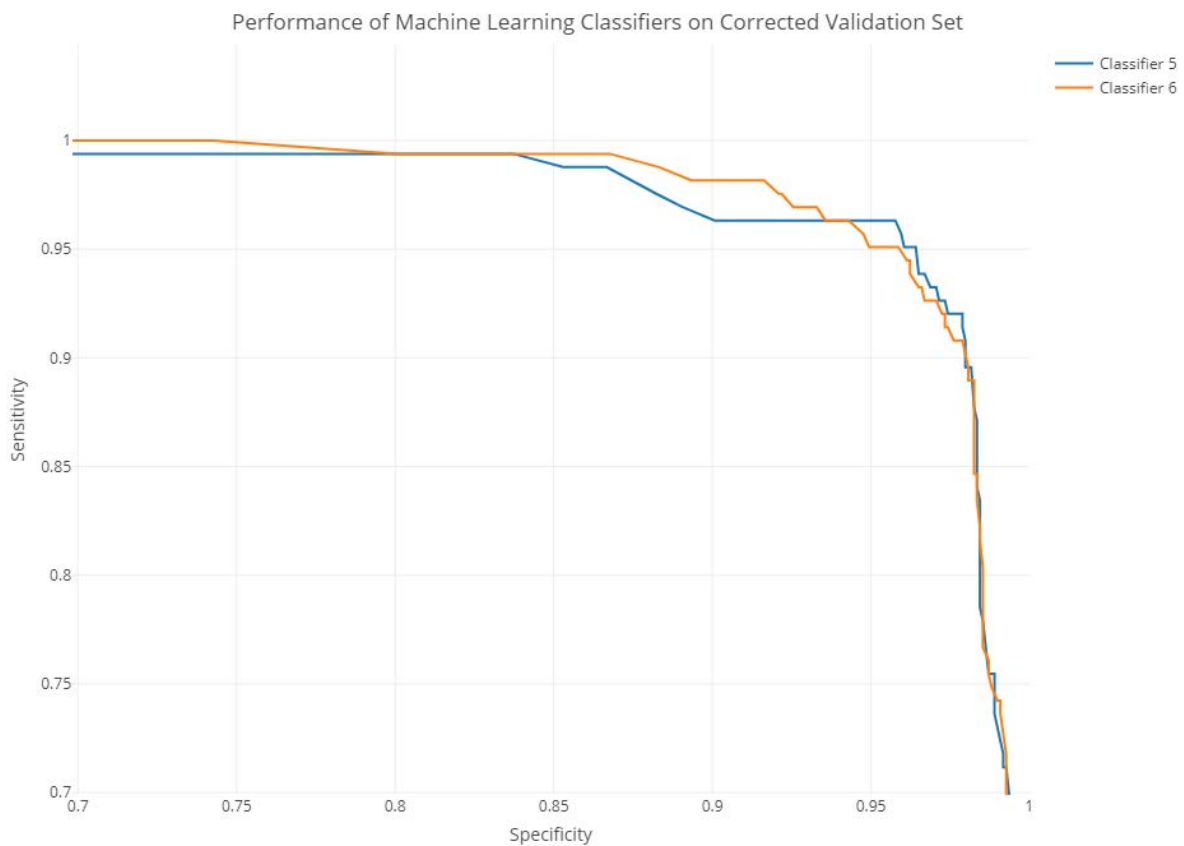
343

Table 4. Performance of machine learning classifier after error analysis.

	Classifier 6	Classifier 5
Cut-Off	0.09	0.10
Sensitivity	98.7%	98.7%
Upper 95% CI of Sensitivity	0.997	0.997
Lower 95% CI of Sensitivity	0.949	0.949
Specificity	88.3%	86.7%
Precision	55.9%	52.61%
Accuracy	89.7%	88.2%
WSS@95%	961/ 1251 – (0.05) = 0.718	945/1251 – (0.05) = 0.705

344

345 Figure 2. Performance of classifiers after error analysis. Classifier 6 is classifier 5 retrained on the
 346 corrected training set after error analysis correction. Performance on both Classifier 5 and 6 is
 347 measured on the corrected validation set (with error analysis correction). For the interactive version
 348 of this plot with exact cut-off values, see code and data at [https://github.com/abannachbrown/The-](https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-systematic-reviews/blob/master/Performance-after-error-analysis.html)
 349 [use-of-text-mining-and-machine-learning-algorithms-in-systematic-](https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-systematic-reviews/blob/master/Performance-after-error-analysis.html)
 350 [reviews/blob/master/Performance-after-error-analysis.html](https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-systematic-reviews/blob/master/Performance-after-error-analysis.html)



351

352

353 We compared the area under the ROC curve for classifier 5 and 6. The AUC for Classifier 5 was
354 0.9272 (95% CI calculated using DeLong method; 0.914-0.9404). The AUC for Classifier 6 was 0.9355
355 (95% CI calculated using DeLong method; 0.9227-0.9483). DeLong's test to compare the AUC
356 between the ROC of the two classifiers was applied ('pROC' package in R; [31]), $Z = -2.3685$, $p =$
357 0.0178.

358 Discussion:

359 Document Classification:

360 Overall, machine learning algorithms are shown here to have high levels of performance, with a
361 sensitivity at least comparable to two independent human screeners. The decision making process for
362 selecting classifiers in this project was to maximise sensitivity, in order that the minimum of potentially
363 relevant papers or research are missed. Thereafter, algorithms were then chosen for
364 improved/highest specificity. This was to reduce the subsequent human time required to sort through
365 and assess papers.

366 The two best performing classifiers have similar performance. The slight differences may reflect the
367 method of feature generation. These algorithms have high performance on this specific topic of animal
368 models of depression. As demonstrated previously, the performance of various classifiers can alter
369 depending on the topic and specificity of the research question [3].

370 In a similar broad preclinical research project, neuropathic pain, it took 18 person months to screen
371 33,814 unique records – based on these numbers it would take an estimated 40 person months to
372 screen 70,365 unique records. Performance of machine learning tools demonstrated in this paper can
373 greatly reduce the amount of human resource needed for initial title and abstract screening of a large
374 corpus of records retrieved from a broad search.

375 We have applied the algorithm to the full dataset (remaining 63,365 records) and are in the process
376 of full-text screening. Following this process, it will allow a more in depth learning of the machine that
377 it can apply to any updates to the search.

378

379 **Error Analysis:**

380 By using the ML algorithm to classify the likelihood of inclusion for each record in the training set, we
381 ranked and highlighted discrepancies between the human inclusion or exclusion decision and the
382 machine decision. Using this technique we identified human errors, which were then corrected to
383 update the training set.

384 We have successfully identified human screening errors which were calculated to be just under 1% of
385 the training set which was dual screened by two independent human reviewers. This error analysis
386 results in a 3% increase or change in sensitivity and specificity, which has increased precision,
387 accuracy, and work saved over sampling of the algorithm. We observed an increase in specificity of
388 1.6% without compromise to specificity. In a systematic review with this number of records, or larger,
389 this saves considerable human resources as the number of records required to screen reduces by at
390 least 1125.

391 This was an initial pilot with stopping criteria where if the initial human decision was correct five
392 consecutive times, further records were not reassessed. However with a more in-depth analysis of the
393 training dataset, investigating every instance where the human and machine decision were
394 incongruent, it is possible that more errors could be identified and correcting them would increase
395 the precision and accuracy of machine learning classifiers, further reducing human resources required
396 for this stage of systematic review

397

398 **Limitations & Future Directions:**

399 Here we display the best performing algorithms for this dataset with a broad research question. Other
400 dissimilar research questions or topics may require different levels of training data to achieve the same
401 levels of performance, or may require different topic modelling approaches or classifiers. The best
402 performing algorithm, outlined in this paper, is being applied in an ongoing research project, therefore
403 the ‘true’ inclusion and exclusion results for the remaining 63365 records is not yet known. The ‘true’
404 results will unfold with the fullness of time.

405 These machine learning algorithms are deployed in an existing systematic review online platform,
406 EPPI-Reviewer [20] and are in the process of being integrated into the SyRF tool, which is focused on
407 the preclinical domain (www.app.syrf.org). This will improve the ease of use of machine learning
408 functions for systematic reviewers, increase the usage of machine learning algorithms for systematic
409 review and significantly reduce the amount of human resources required to conduct systematic review
410 across a range of topics. By allowing a degree of user control over which classifiers and the levels of
411 performance are required for each specific research project. With a broad collaboration such as SLIM
412 we aim to test many ML algorithms across a range of research topics to identify which classifiers
413 perform best under which circumstances, to be able to provide recommendations to users of SyRF.

414

415 This paper outlines a pilot approach to using machine learning algorithms to identify human errors in
416 current systematic review methodology. Future research can investigate this concept more
417 thoroughly by setting up a more comprehensive experimental design. After further investigation into
418 the extent of human error in dual reviewing, the picture will be clearer as to the scale of human error
419 and to what extent a machine learning algorithm can identify and aid in rectifying this. These tools can
420 could be integrated into systematic review platforms, such as SyRF (www.app.syrf.org), and may
421 provide feedback to the systematic reviewer during screening, and could ultimately flag incorrectly
422 screened records as the human screens them for inclusion in a dataset for machine training.

423

424 **Conclusions:**

425 We have demonstrated that machine learning techniques can be successfully applied to an ongoing,
426 broad pre-clinical systematic review. We have demonstrated that machine learning techniques can be
427 used to identify human errors in the training and validation datasets. We have demonstrated that
428 updating the learning of the algorithm after error analysis improves performance. This error analysis
429 technique requires further detailed elucidation and validation. These machine learning techniques are
430 in the process of being integrated into existing systematic review applications to enable more wide-
431 spread use. In future, machine learning and error analysis techniques that are optimised for different
432 types of review topics and research questions can be applied seamlessly within the existing
433 methodological framework.

434 **References:**

435

436 [1] Bornmann, L., & Mutz, R. (2015). Growth rates of modern science: A bibliometric analysis based
437 on the number of publications and cited references. *Journal of the Association for Information*
438 *Science and Technology*, 66(11), 2215-2222.

439 [2] Cohen, A. M., Adams, C. E., Davis, J. M., Yu, C., Yu, P. S., Meng, W., ... & Smalheiser, N. R. (2010,
440 November). Evidence-based medicine, the essential role of systematic reviews, and the need for
441 automated text mining tools. In *Proceedings of the 1st ACM international Health Informatics*
442 *Symposium* (pp. 376-380). ACM.

443 [3] Howard, B.E., Phillips, J., Miller, K., Tandon, A., Mav, D., Shah, M.R., Holmgren, S., Pelch, K.E.,
444 Walker, V., Rooney, A.A. and Macleod, M., 2016. SWIFT-Review: a text-mining workbench for
445 systematic review. *Systematic reviews*, 5(1), p.87.

446 [4] Tsafnat, G., Glasziou, P., Choong, M.K., Dunn, A., Galgani, F. and Coiera, E., 2014. Systematic
447 review automation technologies. *Systematic reviews*, 3(1), p.74.

448 [5] O'Mara-Eves, A., Thomas, J., McNaught, J., Miwa, M., & Ananiadou, S. (2015). Using text mining
449 for study identification in systematic reviews: a systematic review of current approaches. *Systematic*
450 *reviews*, 4(1), 5.

451 [6] Borah, R., Brown, A.W., Capers, P.L., *et al.* (2017). Analysis of the time and workers needed to
452 conduct systematic reviews of medical interventions using data from the PROSPERO registry. *BMJ*
453 *Open*;7:e012545. doi: 10.1136/bmjopen-2016-012545

454 [7] Thomas, J., McNaught, J., Ananiadou, S. (2011). Applications of text mining within systematic
455 reviews. *Res Synth Methods*,2(1):1–14. 10.1002/jrsm.27

- 456 [8] Wallace, B. C., Small, K., Brodley, C. E., & Trikalinos, T. A. (2010). Active learning for biomedical
457 citation screening. In Proceedings of the 16th ACM SIGKDD International conference on Knowledge
458 Discovery and Data mining (pp. 173-182). ACM.
- 459 [9] Miwa, M., Thomas, J., O'Mara-Eves, A. and Ananiadou, S., 2014. Reducing systematic review
460 workload through certainty-based screening. *Journal of biomedical informatics*, 51, pp.242-253.
- 461 [10] Sena, E. S., Currie, G. L., McCann, S. K., Macleod, M. R., & Howells, D. W. (2014). Systematic
462 reviews and meta-analysis of preclinical studies: why perform them and how to appraise them
463 critically. *Journal of Cerebral Blood Flow & Metabolism*, 34(5), 737-742.
- 464 [11] Bannach-Brown, A., Liao, J., Wegener, G., & Macleod, M.R. (2016). Understanding in vivo
465 modelling of depression in non-human animals: a systematic review protocol. *Evidence-based*
466 *Preclinical Medicine*, 3(2), 20-27.
- 467 [12] Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to Information Retrieval,
468 Cambridge University Press: USA.
- 469 [13] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). Indexing by
470 latent semantic analysis. *Journal of the American society for information science*, 41(6), 391.
- 471 [14] Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of machine*
472 *Learning research*, 3(Jan), 993-1022.
- 473 [15] Mo, Y., Kontonatsios, G., & Ananiadou, S. (2015). Supporting systematic reviews using LDA-
474 based document representations. *BMC Systematic Reviews*, 4(1), 1.
- 475 [16] Hashimoto, K., Kontonatsios, G., Miwa, M., & Ananiadou, S. (2016). Topic detection using
476 paragraph vectors to support active learning in systematic reviews. *Journal of biomedical*
477 *informatics*, 62, 59-65.

- 478 [17] Mu, T., Goulermas, J. Y., Korkontzelos, I., & Ananiadou, S. (2016). Descriptive document
479 clustering via discriminant learning in a co-embedded space of multilevel similarities. *Journal of the*
480 *Association for Information Science and Technology*, 67(1), 106-133.
- 481 [18] Mertsalov, K., & McCreary, M. (2009). Document classification with support vector machines.
482 Rational Enterprise: White Paper. Accessed from:
483 [http://www.rationalenterprise.com/assets/content/files/Classification_with_Support_Vector_Machi](http://www.rationalenterprise.com/assets/content/files/Classification_with_Support_Vector_Machines.pdf)
484 [nes.pdf](http://www.rationalenterprise.com/assets/content/files/Classification_with_Support_Vector_Machines.pdf) , on: 05/09/2016.
- 485 [19] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P,
486 Weiss R, Dubourg V, Vanderplas J. (2011). Scikit-learn: Machine learning in Python. *Journal of*
487 *machine learning research*.12(Oct):2825-30.
- 488 [20] Thomas, J., Brunton, J., Graziosi, S., (2010). EPPI-Reviewer 4.0: software for research synthesis.
489 EPPI-Centre Software. London: Social Science Research Unit, Institute of Education.
- 490 [21] Tsuruoka, Y., Tateishi, Y., Kim, J. D., Ohta, T., McNaught, J., Ananiadou, S., & Tsujii, J. I. (2005).
491 Developing a robust part-of-speech tagger for biomedical text. In *Panhellenic Conference on*
492 *Informatics* (pp. 382-392). Springer, Berlin, Heidelberg.
- 493 [22] McCallum, Andrew Kachites. (2002). "MALLET: A Machine Learning for Language Toolkit."
494 <http://mallet.cs.umass.edu>.
- 495 [23] Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization paths for generalized linear
496 models via coordinate descent. *Journal of statistical software*, 33(1), 1.
- 497 [24] R Core Team, "R: A Language and Environment for Statistical Computing." Vienna, Austria, 2013.
- 498 [25] Kuhn, M., (2017) "The caret package". <https://topepo.github.io/caret/>
- 499 [26] Newcombe, R.G. (1998)."Two-Sided Confidence Intervals for the Single Proportion: Comparison
500 of Seven Methods," *Statistics in Medicine*, 17, 857-872

- 501 [27] Bannach-Brown, A., Thomas, J., Przybyła, P., Liao, J., (2016). “Protocol for Error Analysis:
502 Machine learning and text mining solutions for systematic reviews of animal models of depression”.
503 Published on CAMARADES Website. www.CAMARADES.info. Direct Access:
504 <https://drive.google.com/file/d/0BxckMffc78BYTm0tUzJJZkc1alk/view>
- 505 [28] Rodriguez, J. D., Perez, A., & Lozano, J. A. (2010). Sensitivity analysis of k-fold cross validation in
506 prediction error estimation. *IEEE transactions on pattern analysis and machine intelligence*, 32(3),
507 569-575.
- 508 [29] Kerr, K. F., Wang, Z., Janes, H., McClelland, R. L., Psaty, B. M., & Pepe, M. S. (2014). Net
509 Reclassification Indices for Evaluating Risk-Prediction Instruments: A Critical Review. *Epidemiology*
510 *(Cambridge, Mass.)*, 25(1), 114–121. <http://doi.org/10.1097/EDE.0000000000000018>
- 511 [30] Pencina, M.J., D'Agostino, R.B. and Vasan, R.S., (2008). Evaluating the added predictive ability of
512 a new marker: from area under the ROC curve to reclassification and beyond. *Statistics in*
513 *medicine*, 27(2), 157-172.
- 514 [31] Robin, X. (2017). “pROC” Package. <https://cran.r-project.org/web/packages/pROC/pROC.pdf>

515

516 Declarations:

517

518 Availability of Data & Materials:

519 The training and validation datasets, error analysis datasheets, as well as all the records in the
520 depression systematic review are available on Zenodo: DOI [10.5281/zenodo.60269](https://doi.org/10.5281/zenodo.60269)

521 The protocol for the systematic review of animal models of depression is available from:

522 <http://onlinelibrary.wiley.com/doi/10.1002/ebm2.24/pdf>

523 The protocol for the Error Analysis is available via the CAMARADES website and can be accessed
524 directly from this link: <https://drive.google.com/file/d/0BxckMffc78BYTm0tUzJJZkc1alk/view>

525 The results of the classification algorithms and the R code used to generate the results is available on
526 GitHub: [https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-](https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-systematic-reviews)
527 [algorithms-in-systematic-reviews](https://github.com/abannachbrown/The-use-of-text-mining-and-machine-learning-algorithms-in-systematic-reviews).

528

529 Competing Interests:

530 The authors declare that they have no competing interests.

531

532 Funding:

533 This work is supported by a grant from the Wellcome Trust & Medical Research Council (Grant
534 Number: MR/N015665/1). ABB is supported by a scholarship from the Aarhus-Edinburgh Excellence
535 in European Doctoral Education Project.

536

537 Authors' Contributions:

538 ABB screened and analysed the datasets. JT & PB conducted feature selection and built the
539 classifiers. ABB, JT & PB wrote the manuscript. ABB, JT, PB, MRM, JL, AR & SA devised the study. JL,
540 MRM & SA supervised the study. All authors edited and approved the final manuscript.

541

542 Acknowledgements:

543 Thank you to Kaitlyn Hair & Paula Grill for their assistance in second screening the training and
544 validation datasets.