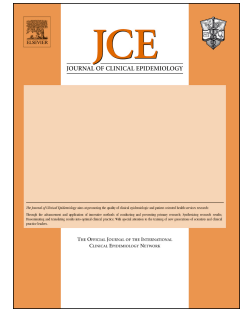


# Accepted Manuscript

Extending PubMed Searches to ClinicalTrials.gov Through a Machine Learning Approach for Systematic Reviews

Corrado Lanera, Clara Minto, Abhinav Sharma, Dario Gregori, Paola Berchiulla, Ileana Baldi



PII: S0895-4356(18)30085-4

DOI: [10.1016/j.jclinepi.2018.06.015](https://doi.org/10.1016/j.jclinepi.2018.06.015)

Reference: JCE 9692

To appear in: *Journal of Clinical Epidemiology*

Received Date: 26 January 2018

Revised Date: 19 June 2018

Accepted Date: 29 June 2018

Please cite this article as: Lanera C, Minto C, Sharma A, Gregori D, Berchiulla P, Baldi I, Extending PubMed Searches to ClinicalTrials.gov Through a Machine Learning Approach for Systematic Reviews, *Journal of Clinical Epidemiology* (2018), doi: 10.1016/j.jclinepi.2018.06.015.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

## Extending PubMed Searches to ClinicalTrials.gov Through a Machine Learning Approach for Systematic Reviews

Corrado Lanera<sup>1</sup>, Clara Minto<sup>1</sup>, Abhinav Sharma<sup>2\*</sup>, Dario Gregori<sup>3</sup>, Paola Berchialla<sup>4</sup>, Ileana Baldi<sup>5</sup>

<sup>1</sup>Researcher, Unit of Biostatistics, Epidemiology and Public Health, Department of Cardiac, Thoracic and Vascular Sciences, University of Padova, Via Loredan 18, 35131 Padova, Italy

<sup>2</sup>Student, Department of Biological Sciences and Bioengineering (BSBE), IIT Kanpur, India

(\*) The work was performed during an internship at the Unit of Biostatistics, Epidemiology and Public Health, Department of Cardiac, Thoracic and Vascular Sciences, University of Padova, Via Loredan 18, 35131 Padova, Italy

<sup>3</sup>Full Professor, Unit of Biostatistics, Epidemiology and Public Health, Department of Cardiac, Thoracic and Vascular Sciences, University of Padova, Via Loredan 18, 35131 Padova, Italy

<sup>4</sup>Assistant Professor, Department of Clinical and Biological Sciences, University of Torino, Via Santena 5bis, 10126 Torino, Italy

<sup>5</sup>Associate Professor, Unit of Biostatistics, Epidemiology and Public Health, Department of Cardiac, Thoracic and Vascular Sciences, University of Padova, Via Loredan 18, 35131 Padova, Italy

### Correspondence to:

Dr. Ileana Baldi,

Unit of Biostatistics, Epidemiology and Public Health

Department of Cardiac, Thoracic and Vascular Sciences, University of Padova

Via Loredan, 18 - 35131 Padova, Italy

Email: ileana.baldi@unipd.it

Phone: +39 049 8275403

Fax: +39 02 700445089

**Abstract**

**Aims.** Despite their essential role in collecting and organizing published medical literature, indexed search engines are unable to cover all relevant knowledge. Hence, current literature recommends the inclusion of clinical trial registries in systematic reviews. This study aims to provide an automated approach to extend a search on PubMed to the ClinicalTrials.gov database, relying on text mining and machine learning techniques.

**Study Design and Setting.** The procedure starts from a literature search on PubMed. Next, it considers the training of a classifier that can identify documents with a comparable word characterization in the ClinicalTrials.gov clinical trial repository. Fourteen systematic reviews, covering a broad range of health conditions, are used as case studies for external validation. A cross-validated support-vector machine model was used as the classifier.

**Results.** The sensitivity was 100% in all systematic reviews except one (87.5%), and the specificity ranged from 97.2 to 99.9%. The ability of the instrument to distinguish on-topic from off-topic articles ranged from an AUC of 93.4 to 99.9%.

**Conclusion.** The proposed machine learning instrument has the potential to help researchers identify relevant studies in the systematic review process by reducing workload, without losing sensitivity and at a small price in terms of specificity.

**Keywords:** Systematic Review; Meta-Analysis; Clinical Trial Registry; Indexed Search Engine; Machine Learning; Text Mining

**Running title:** A Machine Learning Approach for Comprehensive Searches

**Word count:** Abstract: 200; Manuscript: 3271

**What is new:**

- **Key findings:** Our study proposes a new classifier that extends PubMed searches to clinical trial registries with high discrimination ability and sensitivity.
- **What this adds to what is known:** Text mining and machine learning techniques can be used to support comprehensive systematic reviews.
- **What is the implication, what should change now:** The proposed machine learning instrument can help researchers identify relevant studies in the systematic review process by reducing workload without losing sensitivity and at a small price in terms of specificity.

## 1. Introduction

In medical practice and research, the highest level of evidence is represented by systematic reviews (SRs) [1]. An SR is the synthesis and evaluation of all relevant literature on a specific topic, aimed to make the available knowledge more accessible to physicians, care providers and decision makers [2]. Conducting an SR is not an easy task since it must follow specific guidelines and protocols to ensure reproducibility of the methods. After the definition of review questions, researchers should accurately identify evidence from articles, studies, and any other relevant documentation. This selection process consists of an active search through online and offline literature repositories and the identification of evidence from a large amount of irrelevant information [3]. In the search phase, researchers use keyword combinations to create queries that are able to filter documentations in large medical databases. This operational step is prone to potential bias related to the source of information, specificity, and completeness of search strings. After application of queries, researchers manually complete the study selection process by a screening of titles, abstracts, and full texts and assess the papers' eligibility. Finally, they describe the process using a PRISMA flow diagram [4].

The increasing number of web repositories and the development of new scientific topics makes the SR process even more complex [5]. Researchers can retrieve information using search engines, such as PubMed or Embase, that are organized in hierarchical branching structures (MeSH and EmTree, respectively), facilitating paper categorization and specific searches. This logical and hierarchical structure has important implications in the literature search process. First, it facilitates article retrieval by reducing or eliminating potential bias related to the differences in wording, language, and brand names. Second, even if not exhaustive, MeSH or EmTree structures are useful for limiting the number of records to the relevant ones, especially when the study topic is broad.

Despite their essential role in collecting and organizing published medical literature, indexed search engines are often unable to cover all relevant knowledge. A meta-analysis based only on their sources may provide biased estimates due to the exclusion of relevant unpublished information [6]. Furthermore, trial findings can influence the probability of publication and the presence of selective reporting outcomes [7]. The World Health Organization stated that

unreported studies could lead to a misleading picture of the risks and benefits of a treatment, leading to the use and consumption of ineffective or harmful products [8]. For this reason, SRs should be based on a wide literature dataset, which is essential in order for clinicians and patients to have a reliable and complete picture of their condition and make informed decisions. Among alternative informative sources, current literature recommends the inclusion of clinical trial registries such as ClinicalTrials.gov [9,10]. ClinicalTrials.gov is an international web-based platform organized by the US National Library of Medicine providing access to more than 263,373 clinical trials from 202 countries. Studies are registered and regularly updated by the principal investigator, and records are never removed from the site. On ClinicalTrials.gov, clinicians and patients can retrieve complete information about the disease, intervention, study design and phase, location, and contacts, as well as the links to published papers. Some records also include the results of the study, such as the main characteristics of the population, incidence of adverse events, and collected outcomes. Clinical trial registries are important literary sources contributing to an updated evidence-based medical practice and may contain data that cannot be found in published papers [11]. It has been estimated that 50% of results reported in ClinicalTrials.gov were not initially available elsewhere, while some other information on serious adverse events was not always reported in the corresponding publication [12,13]. In a recent study, Baudard and colleagues found that adding clinical trial registries to the search base of SRs that did not originally search such registries identified an additional 122 trials for 41 SRs, which affected the strength of evidence of the SRs [10]. Despite their relevant role, clinical trial registries are seldom used as sources of studies for systematic reviews, probably due to difficulties in record management. The main limitations are related to the absence of hierarchical order, poor interfaces, a limited number of synonyms and the impossible combination of different queries. In ClinicalTrials.gov, the search strategy is based only on retrieval of one or more text words in the fields Condition/Disease, Title, Brief Description, Interventions, Locations, and Country. Text word variations include a limited number of synonyms, but no hierarchical order or subcategories are used. Recently, the Clinical Trials Transformation Initiative (CTTI) proposed a solution to improve the usability of data included in ClinicalTrials.gov by creating a database for aggregate analysis (AACT) and categorization of clinical trials based on clinical specialty. However, this classification is limited to the definition

of Disease/Condition and is not consistent with the original MeSH classification, which does not allow for differentiation between clinical specialties.

This study aims to 1) provide an instrument based on text mining and machine learning (ML) techniques that can perform an automated literature search on clinical trial registries; and 2) evaluate the usability and effectiveness of the proposed instrument. To reach our objectives, we present a case study based on results reported by Baudard et al [10].

## **2. Materials and Methods**

### **2.1 Data sources**

To create and test the instrument for automated literature search, we used two different data sources. First, we used information reported in the article *Impact of searching clinical trials registries in systematic reviews of pharmaceutical treatments: methodological systematic review and reanalysis of meta-analyses* [10]. That study identified additional trials not included in original SRs, through a manual search of the International Clinical Trials Registry Platform (ICTRP). Specifically, the authors adapted and applied to ICTRP the search strings of fourteen SRs on the effectiveness of pharmacological treatments for several health conditions (i.e., atrial fibrillation, psoriasis, colorectal cancer, gastric cancer, Alzheimer's disease, Parkinson's disease, diabetes, rheumatoid arthritis, and hypertension). Then, they verified the consistency of the retrieved records with the inclusion criteria listed in the original papers and included relevant trials in a final estimation of treatment effectiveness. For our purpose, we used the same fourteen SRs listed in [10]. This information allowed us to recreate search strings for PubMed and compare the results of automated searches with those reported by the authors of each SR. Second, we used the full database of ClinicalTrials.gov downloaded from the website of the Clinical Trials Transformation Initiative. The database was organized in pipe-delimited files with data on each single study, such as identifier (NCT number), location, start date, sample size, etc. Data could be reported as a number, string (i.e., text), date, or Boolean (i.e., true and false).

### **2.2 Training datasets**

We created a training dataset for each of the fourteen SRs described above. Each training dataset included positive and negative records. Positive records were papers included in the original SRs, while negative records were a sample of papers off topic. Positive records were identified by running the original query in PubMed. When the search strategy did not allow us to retrieve

all relevant papers, missing citations were manually included in the training set. On the other hand, negative records were retrieved by adding the Boolean operator NOT to the original query. In other words, we identified off-topic papers by subtracting records of the original search strategy from the complete PubMed database. Negative records were filtered by “Text availability: abstract”, “Article types: Clinical trial”, “Species: Humans” and “Languages: English”. Since PubMed allows users to download up to 200 citations at a time, “Sort by: Best Match” option was selected to avoid any potential bias in the selection of papers based on Entrez Date. Then, negative records were downloaded in groups of 200 every time to achieve a ratio of at least twenty negative records to each positive one. The description of search strings and retrieved records is reported briefly in Table 1 (a more detailed description is reported as supplementary material Table S1). Finally, the first author, year, title, and abstract from each positive and negative paper were collected and included in the training set.

### 2.3 Testing datasets

A snapshot of the whole ClinicalTrials.gov was taken on January 5, 2017. This was composed of a set of pipe-delimited files from which we extracted the following information:

unique identifier (NCT number);

brief title;

official title;

brief summary;

detailed description;

study type (nature of investigation, such as interventional or observational);

overall recruitment status;

month and year of study start (enrolment of first participant);

month and year of primary completion (examination of final participant);

allocation;

number of arms;

study phase;



minimum age for participant eligibility;

interventional study model (otherwise, the strategy for assigning interventions to participants);

inclusion of drug product subject to the US FDA (Federal Food, Drug and Cosmetic Act).

We used the brief title, official title, and detailed description as textual information to perform our testing search. The other information was used to identify trials (NCT numbers) and to include filters similar to those applied in [10]. Specifically, Baudard and colleagues limited ICTRP results to clinical trials whose overall status was either completed or terminated. Moreover, we applied additional filters using fields consistently with the inclusion and exclusion criteria described in the fourteen SRs, replicating the selection filters used in [10] and in the original SRs (see *2.5 Procedure Workflow* for further details). Overall, 233,609 trials were finally included in the testing dataset.

## **2.4 Text Mining**

The text mining strategy consisted of (i) text preprocessing, (ii) training of the ML classifier, and (iii) estimation of the performance of the classifier on the testing dataset. We also considered an option to handle the unbalanced data in the training set. Text preprocessing steps converted the textual data into numbers. A support-vector machine (SVM) model, which is one of the most widely used classifiers for text mining [14], was chosen as the classifier and was trained using 5-fold cross-validation. In each of the training datasets, the ratio of positive to negative samples was at least 1:20 by construction. Data of this type are known as unbalanced data. Hence, on the side of the straight application of the defined procedure, we also used the data handling strategy random undersampling (RUS), which randomly removes cases from the majority samples (in our case, the negative samples) to make the classes more balanced [15]. We applied the RUS strategy to obtain a final positive:negative ratio of the class samples of 35:65 according to [16]. In this way, we had 28 datasets overall, two for each SR, i.e., the original one and the one after the application of the RUS.

## **2.5 Procedure workflow**

For each of the fourteen SRs, the title and the abstract of the retrieved records were merged, and text preprocessing steps were applied in the following order: conversion to lowercase, removing non-words, stemming words, stripping white space, and building the sequences of every two adjacent words from the original text (bi-grams). Further, a document-term matrix (DTM) was

created with this collection of tokens (i.e., a unit of textual information), and the matrix was filled with a term frequency (TF) weighting scheme. The sparsity of all 14 DTMs was very high, ranging from 99 to 100%. The top 4% of the features were selected according to term frequency-inverse document frequency (TF-IDF) rank as a tribute to (a double application of) Pareto's rule, i.e., that 80% of the effects come from 20% of the causes. These selected features were retained. The SVM was 5-fold cross-validated, and within the cross-validation step, the balancing strategy of RUS and the 35:65 positive:negative ratio were applied. Next, reweighting with TF-IDF was applied.

The testing ClinicalTrials.gov dataset went through the same text preprocessing strategy in the same order, and then DTM was created with the TF weighting scheme initially. Further, it was adapted with the same features retained from the training dataset and was reweighted with the TF-IDF weighing scheme with the same retained IDF weights of the corresponding training DTM, which were retained when applied on the whole training dataset.

Each cross-validated SVM model was applied to the corresponding testing dataset for each SR. The procedure workflow is briefly described in Figure 1. Analyses were carried out in R version 3.4.2 [17] with the packages *caret*, *tm*, *stringr*, and *unbalanced* [18–21].

To compare the consistency between the manual search in [10] and this automated search, we replicated the selection filters used in [10] and in the original SRs. Thus, positive citations identified by automated search were limited by adding all the following filters: 1) recruitment status, defined as completed or terminated; 2) interventional design; 3) start date before the search on ICTRP; 4) primary completion date before the search on ICTRP as reported in [10]; and 5) specific filters based on inclusion criteria reported in the original SRs. The goodness and robustness of our results were evaluated by verifying the inclusion of the additional clinical trials previously identified by Baudard and colleagues.

### 3. Results

The performance results of the most suitable filter are reported in Table 2. The sensitivity was 100% in all SRs except one (87.5%), and the specificity ranged from 97.2% to 99.9%. The AUC, which measures the ability of the instrument to distinguish relevant articles from off-topic articles, ranged from 93.4% to 99.9%. The performance of the procedures in which an RUS strategy was implemented was similar (data are not shown). Table 3 reports the numbers of

predicted positive citations before and after the application of a selection of filters. It also compares our results with the results of Baudard and colleagues' manual search results on ICTRP. As shown in the table, filters progressively reduced the number of citations (predicted positives), without excluding additional clinical trials identified in [10] (true positives).

The only false negative (1 out of 8 positives) pertained to an SR on the role of biological therapy in metastatic colorectal cancer [22] and referred to the study with ClinicalTrials.gov identifier NCT00079066.

Notably, the total number of records from our automated search (predicted positives) was lower than the number of records from the manual search in half the cases, with a mean of 472 and a maximum of 2119 records compared with 572 and 2680, respectively, retrieved in [10].

#### 4. Discussion

The time requirement and the need for the involvement of different professionals make an SR a very labour-intensive process [23]. The quality of the results depends on the extent to which the identified literature is accurate and comprehensive of all available knowledge on a specific topic. Also, the reliability of an SR is determined by the inclusion of up-to-date contents [24]. Our study proposes a classifier that can extend PubMed searches to clinical trials registries. This tool reduces the effort and time expenditure of an SR without losing accuracy and sensitivity.

Other researchers have highlighted how ML could make the standard SR process more efficient [25]. They focused on a “living” SR, considering as the starting point the existence of an initial SR provided by humans. Accordingly, we have provided an instrument that is also usable for the “living” step of updating an SR dataset, but it is specially tailored to the more complex and tricky step of contributing to the base dataset definition/extraction for new sources of data (work left to humans in [25]). Our procedure showed high performance in detecting true-positive citations of interest in completely different sources of data from the original one regarding the way meta-data are stored, the way information is accessed, and the structure of the information. It left out only 1 of 133 human-detected positive citations from fourteen independent SRs. From this starting point, we have also highlighted how, with simple and quick filtering, the number of false positives can be easily and drastically reduced without affecting the sensitivity of the procedure. In this way, the work left to humans can be reduced and quite limited on the first run of the

“living” update of the SR, i.e., the part of dataset definition that was completely based on human effort until now.

Other studies have shown how an ML approach for the classification of information based on clinical text could be very effective [26], including when tested on databases different (and not subsampled) from the original one [27]. On the other hand, to our knowledge, no other study was conducted on this wide a range of differentiated datasets with hundreds of thousands of entries.

The strict procedure that we followed allowed all the test sets to be blinded both from the training ones at every stage and from all the training procedures, making us confident in the quality of the results themselves. In an SR, both very specific positive and very specific negative sets can be selected to create a high-quality training set. This characteristic together with the ability of the SVMs to distinguish the well-separated type of data and the high proportion of (few) positive records against a huge number of negative ones have led to the nearly perfect results in sensitivity, which is the main characteristic of interest in this endeavour.

Our study demonstrates the usefulness of ML when scientific literature is not reported in indexed search engines. This is the case of clinical trial registries such as ClinicalTrials.gov, whose interfaces are usually not sophisticated. Their limited functionality has an important impact on process workload and often requires the application of long search strings, multiple searches and the screening of a high number of non-specific records. Moreover, when a researcher wants to use the same query on different search engines and registries, he must adapt each singular term and string according to the specific requirement of each platform. In the case of registries, an adaptation from common search engines (PubMed or Embase) is even more complex due to the frequent absence of text functionalities such as truncation or brackets. The use of ML could allow a more accurate and easier translation of queries by reducing the number of irrelevant records.

The main strength of the study is the robustness of the training and testing procedure, which was designed to be stable and unbiased. Furthermore, an R package and a companion GUI are under development (preliminary version publicly available at <https://github.com/UBESP-DCTV/costumer>). They are intended to be a user-friendly tool for healthcare researchers, who will only have to provide a) the set of citations finally retained, b) a personalized set of negative citations or the search string used on PubMed (to automatically identify a suitable set of random

negative citations), and c) an optional set of false positives already known from a previous run or directly the set of filters to be applied on non-textual meta-data. The first part of feature c) highlights the usability of the package for a very quick update of the SR, e.g., after the first run (for which the false positives must be manually identified).

Our study has some limitations. First, we adopted a defined ML algorithm and used only one strategy for managing the unbalanced data. We acknowledge that other techniques, such as convolutional neural networks (CNNs), are effective at achieving slightly better F-scores [28] over more traditional approaches to biomedical text classification, such as SVM, especially when there is significant label imbalance. Nevertheless, CNNs typically take at least an order of magnitude more time than traditional classifiers, especially compared with SVM [29]. Hence, we decided to start our investigation by considering SVM only. We are already working on testing both a wider range of ML techniques and more methods for unbalanced datasets. Nevertheless, the performance with the choice adopted in terms of the number of positives, number of true positives and number of negatives, as well as in terms of computational speed, is already good, and we do not expect more improvement, though small relative increases in specificity can still have a big impact on absolute numbers of false positives. Moreover, filters were manually applied after automatic searches and were not yet included in this ML instrument. The reason for this choice was that inclusion/exclusion criteria are rarely reported in the title, abstract or description. Thus, it was not possible to make a more accurate automatic selection of trials. That said, similar studies were able to reach a very high level of sensitivity at the cost of a discrete specificity [30]. Explanations for our results lie in the choice of the training set and of the task itself, that is, positives and negatives were highly separated by design: the set of positives is the (human-filtered) output of a PubMed query string and not of a “rule-free” human selection from the whole of PubMed, as was done, e.g., for the classification task performed in [30]. The same applies to the negatives, which were sampled from the output of the negative search with the PubMed query string used for positives. As a result, positives shared a similar word characterization, which is easily identified by SVM and can lead to a near-perfect sensitivity and an excellent specificity.

Following the recommended paradigm for model validation [31,32], this predictive tool underwent internal validation through cross-validation and external validation on an independent

data source. This aspect, in conjunction with the broad range of health conditions analysed, strongly argues in favour of the credibility of the proposed instrument.

**Contributions:** IB and PB designed the study; CM managed the systematic reviews; CL and AS performed the analysis; CM and AS wrote the manuscript; all authors contributed to result interpretation and approved the final manuscript.

**Funding:** This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

Conflicts of Interest: none.

## References

- [1] Howick J, Chalmers I, Glasziou P, Greenhalgh T, Heneghan C, Liberati A, et al. Explanation of the 2011 Oxford Centre for Evidence-Based Medicine (OCEBM) Levels of Evidence (Background Document) 2016. <https://www.cebm.net/2016/05/ocebml-levels-of-evidence/> (accessed January 13, 2018).
- [2] CRD's guidance for undertaking reviews in health care 2009. [https://www.york.ac.uk/media/crd/Systematic\\_Reviews.pdf](https://www.york.ac.uk/media/crd/Systematic_Reviews.pdf).
- [3] Hirschman L, Burns GAPC, Krallinger M, Arighi C, Cohen KB, Valencia A, et al. Text mining for the biocuration workflow. Database J Biol Databases Curation 2012;2012:bas020. doi:10.1093/database/bas020.
- [4] Liberati A, Altman DG, Tetzlaff J, Mulrow C, Gøtzsche PC, Ioannidis JPA, et al. The PRISMA Statement for Reporting Systematic Reviews and Meta-Analyses of Studies That Evaluate Health Care Interventions: Explanation and Elaboration. PLOS Med 2009;6:e1000100. doi:10.1371/journal.pmed.1000100.
- [5] Balan PF, Gerits A, Vanduffel W. A practical application of text mining to literature on cognitive rehabilitation and enhancement through neurostimulation. Front Syst Neurosci 2014;8:182. doi:10.3389/fnsys.2014.00182.
- [6] Guyatt GH, Oxman AD, Montori V, Vist G, Kunz R, Brozek J, et al. GRADE guidelines: 5. Rating the quality of evidence--publication bias. J Clin Epidemiol 2011;64:1277–82. doi:10.1016/j.jclinepi.2011.01.011.
- [7] Hopewell S, Loudon K, Clarke MJ, Oxman AD, Dickersin K. Publication bias in clinical trials due to statistical significance or direction of trial results. Cochrane Database Syst Rev 2009;MR000006. doi:10.1002/14651858.MR000006.pub3.
- [8] WHO. Major research funders and international NGOs to implement WHO standards on reporting clinical trial results 2017. <http://www.who.int/mediacentre/news/releases/2017/clinical-trial-results/en/> (accessed January 13, 2018).



- [9] Hughes S, Cohen D, Jaggi R. Differences in reporting serious adverse events in industry sponsored clinical trial registries and journal articles on antidepressant and antipsychotic drugs: a cross-sectional study. *BMJ Open* 2014;4:e005535. doi:10.1136/bmjopen-2014-005535.
- [10] Baudard M, Yavchitz A, Ravaud P, Perrodeau E, Boutron I. Impact of searching clinical trial registries in systematic reviews of pharmaceutical treatments: methodological systematic review and reanalysis of meta-analyses. *BMJ* 2017;356:j448.
- [11] Halfpenny NJ, Thompson JC, Quigley JM, Scott DA. Clinical Trials Registries For Systematic Reviews – An Alternative Source For Unpublished Data. *Value Health* 2015;18:A12. doi:10.1016/j.jval.2015.03.078.
- [12] Zarin DA, Tse T, Williams RJ, Califf RM, Ide NC. The ClinicalTrials.gov Results Database — Update and Key Issues. *N Engl J Med* 2011;364:852–60. doi:10.1056/NEJMsa1012065.
- [13] Tang E, Ravaud P, Riveros C, Perrodeau E, Dechartres A. Comparison of serious adverse events posted at ClinicalTrials.gov and published in corresponding journal articles. *BMC Med* 2015;13:189. doi:10.1186/s12916-015-0430-4.
- [14] Jindal R, Malhotra R, Jain A. Techniques for text classification: Literature review and current trends. *Webology* 2015;12:1.
- [15] Liu AC. The effect of oversampling and undersampling on classifying imbalanced text datasets. Univ Tex Austin 2004.
- [16] Khoshgoftaar TM, Seiffert C, Van Hulse J, Napolitano A, Folleco A. Learning with limited minority class data. *Mach. Learn. Appl. 2007 ICMLA 2007 Sixth Int. Conf. On, IEEE;* 2007, p. 348–353.
- [17] R Development Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria: R Foundation for Statistical Computing; 2008.
- [18] Wing MKC from J, Weston S, Williams A, Keefer C, Engelhardt A, Cooper T, et al. caret: Classification and Regression Training. 2017.
- [19] Feinerer I, Hornik K. tm: Text Mining Package. 2017.
- [20] Wickham H. stringr: Simple, Consistent Wrappers for Common String Operations. 2017.
- [21] Pozzolo AD, Caelen O, Bontempi G. unbalanced: Racing for Unbalanced Methods Selection. 2015.
- [22] Segelov E, Chan D, Shapiro J, Price TJ, Karapetis CS, Tebbutt NC, et al. The role of biological therapy in metastatic colorectal cancer after first-line treatment: a meta-analysis of randomised trials. *Br J Cancer* 2014;111:1122–31. doi:10.1038/bjc.2014.404.
- [23] Khabsa M, Elmagarmid A, Ilyas I, Hammady H, Ouzzani M. Learning to identify relevant studies for systematic reviews using random forest and external information. *Mach Learn* 2016;102:465–82. doi:10.1007/s10994-015-5535-7.
- [24] CRD's guidance for undertaking reviews in health care. Centre for Reviews and Dissemination, University of York; 2009.
- [25] Thomas J, Noel-Storr A, Marshall I, Wallace B, McDonald S, Mavergames C, et al. Living systematic reviews: 2. Combining human and machine effort. *J Clin Epidemiol* 2017;91:31–7.
- [26] Rochefort CM, Verma AD, Eguale T, Lee TC, Buckeridge DL. A novel method of adverse event detection can accurately identify venous thromboembolisms (VTEs) from narrative electronic health record data. *J Am Med Inform Assoc JAMIA* 2015;22:155–65. doi:10.1136/amiajnl-2014-002768.

- [27] Connolly B, Matykiewicz P, Bretonnel Cohen K, Standridge SM, Glauser TA, Dlugos DJ, et al. Assessing the similarity of surface linguistic features related to epilepsy across pediatric hospitals. *J Am Med Inform Assoc JAMIA* 2014;21:866–70. doi:10.1136/amiajnl-2013-002601.
- [28] Rios A, Kavuluru R. Convolutional Neural Networks for Biomedical Text Classification: Application in Indexing Biomedical Articles. *ACM-BCB ACM Conf Bioinforma Comput Biol Biomed* 2015;258–67. doi:10.1145/2808719.2808746.
- [29] Majumder S, Balaji N, Brey K, Fu W, Menzies T. 500+ Times Faster Than Deep Learning (A Case Study Exploring Faster Methods for Text Mining StackOverflow). *ArXiv180205319 Cs Stat* 2018.
- [30] Marshall IJ, Noel Storr A, Kuiper J, Thomas J, Wallace BC. Machine learning for identifying Randomized Controlled Trials: An evaluation and practitioner's guide. *Res Synth Methods* 2018;0. doi:10.1002/jrsm.1287.
- [31] Steyerberg EW, Harrell FE. Prediction models need appropriate internal, internal-external, and external validation. *J Clin Epidemiol* 2016;69:245–7. doi:10.1016/j.jclinepi.2015.04.005.
- [32] Altman DG, Royston P. What do we mean by validating a prognostic model? *Stat Med* 2000;19:453–73. doi:10.1002/(SICI)1097-0258(20000229)19:4<453::AID-SIM350>3.0.CO;2-5.
- [33] Yang Q, Qi X, Li Y. The preventive effect of atorvastatin on atrial fibrillation: a meta-analysis of randomized controlled trials. *BMC Cardiovasc Disord* 2014;14:99. doi:10.1186/1471-2261-14-99.
- [34] Meng Y, Dongmei L, Yanbin P, Jinju F, Meile T, Binzhu L, et al. Systematic review and meta-analysis of ustekinumab for moderate to severe psoriasis. *Clin Exp Dermatol* 2014;39:696–707. doi:10.1111/ced.12390.
- [35] Li D-H, Pan Z-K, Ye F, An H-X, Wu J-X. S-1-based versus 5-FU-based chemotherapy as first-line treatment in advanced gastric cancer: a meta-analysis of randomized controlled trials. *Tumour Biol J Int Soc Oncodevelopmental Biol Med* 2014;35:8201–8. doi:10.1007/s13277-014-2099-2.
- [36] Lv Z-C, Ning J-Y, Chen H-B. Efficacy and toxicity of adding cetuximab to chemotherapy in the treatment of metastatic colorectal cancer: a meta-analysis from 12 randomized controlled trials. *Tumour Biol J Int Soc Oncodevelopmental Biol Med* 2014;35:11741–50. doi:10.1007/s13277-014-2227-z.
- [37] Wang J, Yu J-T, Wang H-F, Meng X-F, Wang C, Tan C-C, et al. Pharmacological treatment of neuropsychiatric symptoms in Alzheimer's disease: a systematic review and meta-analysis. *J Neurol Neurosurg Psychiatry* 2015;86:101–9. doi:10.1136/jnnp-2014-308112.
- [38] Zhou C-Q, Zhang J-W, Wang M, Peng G-G. Meta-analysis of the efficacy and safety of long-acting non-ergot dopamine agonists in Parkinson's disease. *J Clin Neurosci Off J Neurosurg Soc Australas* 2014;21:1094–101. doi:10.1016/j.jocn.2013.10.041.
- [39] Liu X, Xiao Q, Zhang L, Yang Q, Liu X, Xu L, et al. The long-term efficacy and safety of DPP-IV inhibitors monotherapy and in combination with metformin in 18,980 patients with type-2 diabetes mellitus--a meta-analysis. *Pharmacoepidemiol Drug Saf* 2014;23:687–98. doi:10.1002/pds.3586.
- [40] Douxfils J, Buckinx F, Mullier F, Minet V, Rabenda V, Reginster J-Y, et al. Dabigatran etexilate and risk of myocardial infarction, other cardiovascular events, major bleeding, and



- all-cause mortality: a systematic review and meta-analysis of randomized controlled trials. *J Am Heart Assoc* 2014;3:e000515. doi:10.1161/JAHA.113.000515.
- [41] Kourbeti IS, Ziakas PD, Mylonakis E. Biologic therapies in rheumatoid arthritis and the risk of opportunistic infections: a meta-analysis. *Clin Infect Dis Off Publ Infect Dis Soc Am* 2014;58:1649–57. doi:10.1093/cid/ciu185.
- [42] Li ECK, Heran BS, Wright JM. Angiotensin converting enzyme (ACE) inhibitors versus angiotensin receptor blockers for primary hypertension. *Cochrane Database Syst Rev* 2014;CD009096. doi:10.1002/14651858.CD009096.pub2.
- [43] Cavender MA, Sabatine MS. Bivalirudin versus heparin in patients planned for percutaneous coronary intervention: a meta-analysis of randomised controlled trials. *Lancet Lond Engl* 2014;384:599–606. doi:10.1016/S0140-6736(14)61216-2.
- [44] Chatterjee S, Sardar P, Giri JS, Ghosh J, Mukherjee D. Treatment discontinuations with new oral agents for long-term anticoagulation: insights from a meta-analysis of 18 randomized trials including 101,801 patients. *Mayo Clin Proc* 2014;89:896–907. doi:10.1016/j.mayocp.2014.01.030.
- [45] Funakoshi T, Latif A, Galsky MD. Safety and efficacy of addition of VEGFR and EGFR-family oral small-molecule tyrosine kinase inhibitors to cytotoxic chemotherapy in solid cancers: a systematic review and meta-analysis of randomized controlled trials. *Cancer Treat Rev* 2014;40:636–47. doi:10.1016/j.ctrv.2014.02.004.

**Table 1.** Results of PubMed search strategies for the fourteen systematic reviews included in [10]. Final training datasets included the sum of positive and negative citations. Details of search strings are available in Table S1 (supplementary web appendix).

Systematic Review	Health condition	Positive records	Negative records
Yang Q et al. 2014 [33]	Atrial fibrillation	18	400
Meng Y et al. 2014 [34]	Psoriasis	9	200
Segelov E et al. 2014 [22]	Colorectal cancer	13	400
Li DH et al. 2014 [35]	Gastric cancer	6	200
Lv ZC et al. 2014 [36]	Colorectal cancer	12	400
Wang J et al. 2015 [37]	Alzheimer's disease	32	800
Zhou CQ et al. 2014 [38]	Parkinson's disease	9	200
Liu X et al. 2014 [39]	Type 2 diabetes mellitus	23	600
Douxflis J et al. 2014 [40]	Venous thromboembolic events	13	400
Kourbeti IS et al. 2014 [41]	Rheumatoid arthritis	75	1600
Li EC et al. 2014 [42]	Primary hypertension	9	200
Cavender MA et al. 2014[43]	Venous thromboembolic events	14	400
Chatterjee S et al. 2014 [44]	Venous thromboembolic events	18	400
Funakoshi T et al. 2014 [45]	Solid cancers	43	1000

**Table S1.** Replication of PubMed search strategies for the fourteen systematic reviews included in [10]. Final training datasets included the sum of positive and negative citations reported in bold characters.

Systematic Review	Positive search strategy <i>Sorted by Most recent</i>	Negative search strategy <i>Filter for Abstract &amp; Clinical trial &amp; Humans &amp; English Sorted by Best match</i>	Positive records	Negative records
<b>Yang Q et al. 2014 [33]</b>	(atorvastatin) AND atrial fibrillation AND ("0001/01/01"[PDat]: "2014/04/30"[PDat])	((("0001/01/01"[Date - Publication] : "2014/04/30"[Date - Publication])) NOT (atorvastatin) AND atrial fibrillation)  <b>Citations finally included in the main database</b>	n=76 total records n=5 manually added  <b>n=18 positives</b>	n=563037 total records  <b>n=400 negatives</b>
<b>Meng Y et al. 2014 [34]</b>	(((((ustekinumab) OR CNTO-1275) OR interleukin 12 23 monoclonal antibody) OR stelara) AND ((psoriasis) OR (pustulosis of palms and soles))) AND randomized) AND ("0001/01/01"[PDat] : "2013/08/01"[PDat] )	((("0001/01/01"[Date - Publication] : "2013/08/01"[Date - Publication]) AND ("0001/01/01"[PDat] : "2013/08/01"[PDat] ))) NOT ((((((ustekinumab) OR CNTO-1275) OR interleukin 12 23 monoclonal antibody) OR stelara) AND ((psoriasis) OR (pustulosis of palms and soles))) AND randomized))  <b>Citations finally included in the main database</b>	n=91 total records n=0 manually added  <b>n=9 positives</b>	n=538257 total records  <b>n=200 negatives</b>
<b>Segelov E et al. 2014 [22]</b>	(((((("Antibodies, Monoclonal"[Mesh]) OR "Antineoplastic Combined Chemotherapy Protocols"[Mesh]) OR "Antineoplastic Agents"[Mesh])) AND (((("Bevacizumab"[Mesh]) OR "Camptothecin"[Mesh]) OR "Fluorouracil"[Mesh]) OR "Leucovorin"[Mesh])) AND (((("Colorectal Neoplasms"[Mesh]) OR "Adenocarcinoma"[Mesh]) AND (advanced OR metastatic OR metastases OR metastasis))) AND "Humans"[Mesh]) AND "Randomized Controlled Trial" [Publication Type]) AND ("0001/01/01"[PDat] : "2012/05/31"[PDat] )	((("0001/01/01"[Date - Publication] : "2012/05/31"[Date - Publication])) NOT ((((((("Antibodies, Monoclonal"[Mesh]) OR "Antineoplastic Combined Chemotherapy Protocols"[Mesh]) OR "Antineoplastic Agents"[Mesh])) AND (((("Bevacizumab"[Mesh]) OR "Camptothecin"[Mesh]) OR "Fluorouracil"[Mesh]) OR "Leucovorin"[Mesh])) AND (((("Colorectal Neoplasms"[Mesh]) OR "Adenocarcinoma"[Mesh]) AND (advanced OR metastatic OR metastases OR metastasis))) AND "Randomized Controlled Trial" [Publication Type]))))  <b>Citations finally included in the main database</b>	n=913 total records n=10 manually added  <b>n=13 positives</b>	n=499438 total records  <b>n=400 negatives</b>
<b>Li DH et al. 2014 [35]</b>	(((((stomach cancer) OR gastric cancer)) AND S-1) AND fluorouracil) AND ("0001/01/01"[PDat] : "2014/02/20"[PDat] )	((("0001/01/01"[Date - Publication] : "2014/02/20"[Date - Publication])) NOT (((stomach cancer) OR gastric cancer) AND S-1) AND fluorouracil)  <b>Citations finally included in the main database</b>	n=1248 total records n=2 manually added  <b>n=6 positives</b>	n=557386 total records  <b>n=200 negatives</b>
<b>Lv ZC et al. 2014 [36]</b>	(((((("Colorectal Neoplasms"[Mesh]) OR ((colorectal) AND neoplasms)) OR colorectal neoplasms)) AND ((("Cetuximab"[Mesh]) OR cetuximab) AND ("Clinical Trial" [Publication Type]) AND "Humans"[Mesh])) AND ("0001/01/01"[PDat] : "2014/02/16"[PDat] )	((("0001/01/01"[Date - Publication] : "2014/02/16"[Date - Publication])) NOT (((("Colorectal Neoplasms"[Mesh]) OR ((colorectal) AND neoplasms)) OR colorectal neoplasms) AND ((("Cetuximab"[Mesh]) OR cetuximab)))  <b>Citations finally included in the main database</b>	n=201 total records n=0 manually added  <b>n=12 positives</b>	n=557094 total records  <b>n=400 negatives</b>
<b>Wang J et al. 2015 [37]</b>	(((((alzheimer's disease[Title/Abstract]) OR (alzheimer[Title/Abstract] OR (AD[Title/Abstract]) AND ((cholinesterase inhibitors[Title/Abstract]) OR (donepezil[Title/Abstract]) OR (galantamine[Title/Abstract]) OR (rivastigmine[Title/Abstract]) OR (metrifonate[Title/Abstract]) OR (tacrine[Title/Abstract]) OR (antipsychotics[Title/Abstract]) OR (haloperidol[Title/Abstract]) OR (thioridazine[Title/Abstract]) OR (thiothixene[Title/Abstract]) OR (chlorpromazine[Title/Abstract]) OR (acetophenazine[Title/Abstract]) OR (clozapine[Title/Abstract]) OR (olanzapine[Title/Abstract]) OR (risperidone[Title/Abstract]) OR (quetiapine[Title/Abstract]) OR (aripiprazole[Title/Abstract]) OR (antidepressants[Title/Abstract]) OR (sertraline[Title/Abstract]) OR (fluoxetine[Title/Abstract]) OR (citalopram[Title/Abstract]) OR (trazodone[Title/Abstract]) OR (mood stabilizers[Title/Abstract]) OR (valproate[Title/Abstract]) OR (carbamazepine[Title/Abstract]) OR (lithium[Title/Abstract]) OR (anticonvulsants[Title/Abstract]) OR (benzodiazepines[Title/Abstract]) OR (memantine[Title/Abstract]) OR (psychotropic drugs[Title/Abstract])) AND ((behavioural and psychological symptoms of dementia) OR (BPSD) OR (neuropsychiatric symptoms) OR (behaviour))) AND ("0001/01/01"[PDat] : "2013/11/30"[PDat] ))) AND English[lang])	((("0001/01/01"[Date - Publication] : "2013/11/30"[Date - Publication])) NOT (((alzheimer's disease[Title/Abstract]) OR (alzheimer[Title/Abstract] OR (AD[Title/Abstract]) AND ((cholinesterase inhibitors[Title/Abstract]) OR (donepezil[Title/Abstract]) OR (galantamine[Title/Abstract]) OR (rivastigmine[Title/Abstract]) OR (metrifonate[Title/Abstract]) OR (tacrine[Title/Abstract]) OR (antipsychotics[Title/Abstract]) OR (haloperidol[Title/Abstract]) OR (thioridazine[Title/Abstract]) OR (thiothixene[Title/Abstract]) OR (chlorpromazine[Title/Abstract]) OR (acetophenazine[Title/Abstract]) OR (clozapine[Title/Abstract]) OR (olanzapine[Title/Abstract]) OR (risperidone[Title/Abstract]) OR (quetiapine[Title/Abstract]) OR (aripiprazole[Title/Abstract]) OR (antidepressants[Title/Abstract]) OR (sertraline[Title/Abstract]) OR (fluoxetine[Title/Abstract]) OR (citalopram[Title/Abstract]) OR (trazodone[Title/Abstract]) OR (mood stabilizers[Title/Abstract]) OR (valproate[Title/Abstract]) OR (carbamazepine[Title/Abstract]) OR (lithium[Title/Abstract]) OR (anticonvulsants[Title/Abstract]) OR (benzodiazepines[Title/Abstract]) OR (memantine[Title/Abstract]) OR (psychotropic drugs[Title/Abstract])) AND ((behavioural and psychological symptoms of dementia) OR (BPSD) OR (neuropsychiatric symptoms) OR (behaviour))))))  <b>Citations finally included in the main database</b>	n=1091 total records n=5 manually added  <b>n=32 positives</b>	n=547357 total records  <b>n=800 negatives</b>
<b>Zhou CQ et al. 2014 [38]</b>	(((((pramipexole extended release) OR ropinirole prolonged release) OR rotigotine transdermal patch)) AND (((parkinson's disease) OR parkinson's) OR PD))) AND ("0001/01/01"[PDat] : "2013/02/10"[PDat] )	((("0001/01/01"[Date - Publication] : "2013/02/10"[Date - Publication])) NOT (((pramipexole extended release) OR ropinirole prolonged release) OR rotigotine transdermal patch)) AND (((parkinson's disease) OR parkinson's) OR PD)))  <b>Citations finally included in the main database</b>	n=107 total records n=2 manually added  <b>n=9 positives</b>	n=523622 total records  <b>n=200 negatives</b>
<b>Liu X et al. 2014 [39]</b>	((("Diabetes Mellitus, Type 2"[Mesh]) AND ((((((dpp-iv inhibitors) OR vildagliptin) OR sitagliptin) OR saxagliptin) OR alogliptin) OR linagliptin) OR metformin) OR sulfonylureas))) AND Randomized Controlled Trial[ptyp] AND ("0001/01/01"[PDat] : "2013/01/31"[PDat] ) AND Humans[Mesh] AND English[lang]	((("0001/01/01"[Date - Publication] : "2013/01/31"[Date - Publication])) NOT (((("Diabetes Mellitus, Type 2"[Mesh]) AND ((((((dpp-iv inhibitors) OR vildagliptin) OR sitagliptin) OR saxagliptin) OR alogliptin) OR linagliptin) OR dutogliptin) OR metformin) OR sulfonylureas))) AND Randomized Controlled Trial[ptyp]))  <b>Citations finally included in the main database</b>	n=1427 total records n=0 manually added  <b>n=23 positives</b>	n=521120 total records  <b>n=600 negatives</b>

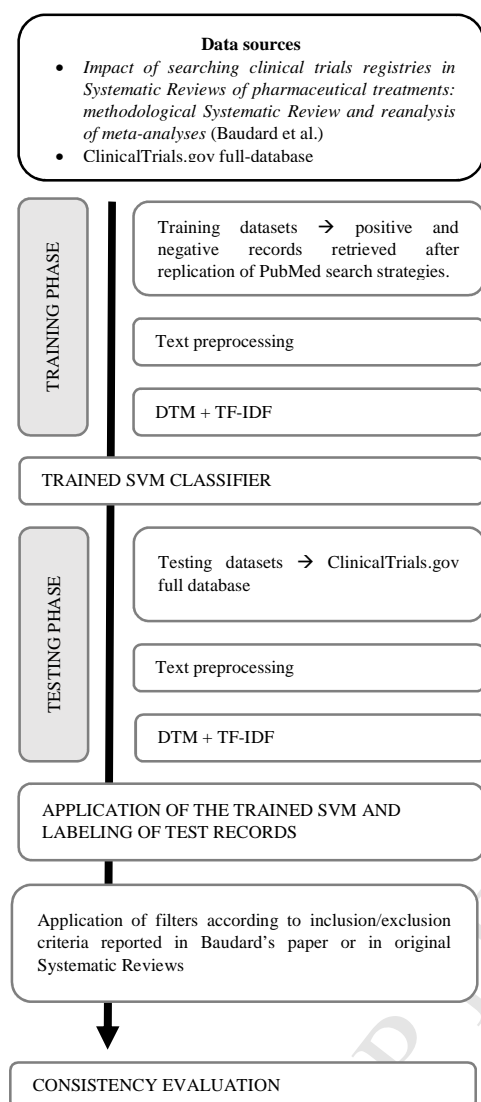
Systematic Review	Positive search strategy <i>Sorted by Most recent</i>	Negative search strategy <i>Filter for Abstract &amp; Clinical trial &amp; Humans &amp; English Sorted by Best match</i>	Positive records	Negative records
<b>Douxfls J et al. 2014 [40]</b>	(((((dabigatran) OR dabigatran etexilate) OR BIBR 1048)) AND ((((((randomized controlled trial) OR randomized clinical trial) OR randomized trial) OR randomised controlled trial) OR randomised clinical trial) OR randomised trial)) AND ( "0001/01/01"[PDat] : "2013/12/08"[PDat] ) AND English[lang]	((("0001/01/01"[Date - Publication] : "2013/12/08"[Date - Publication])) NOT ((((((dabigatran) OR dabigatran etexilate) OR BIBR 1048)) AND ((((((randomized controlled trial) OR randomized clinical trial) OR randomized trial) OR randomised controlled trial) OR randomised clinical trial) OR randomised trial)))	n=276 total records n=3 manually added	n=548647 total records
		<b>Citations finally included in the main database</b>	<b>n=13 positives</b>	<b>n=400 negatives</b>
<b>Kourbeti IS et al. 2014 [41]</b>	((((((rheumatoid) AND arthritis)) AND randomized) AND (((((((infiximab) OR etanercept) OR adalimumab) OR certolizumab) OR golimumab) OR anakinra) OR abatacept) OR tocilizumab) OR rituximab)) AND ( "0001/01/01"[PDat] : "2013/06/24"[PDat] ) AND English[lang])) AND ( "0001/01/01"[PDat] : "2013/06/24"[PDat] )	((("0001/01/01"[Date - Publication] : "2013/06/24"[Date - Publication])) NOT (((((((rheumatoid) AND arthritis)) AND randomized) AND (((((((infiximab) OR etanercept) OR adalimumab) OR certolizumab) OR golimumab) OR anakinra) OR abatacept) OR tocilizumab) OR rituximab)))	n=827 total records n=4 manually added	n=534353 total records
		<b>Citations finally included in the main database</b>	<b>n=75 positives</b>	<b>n=1600 negatives</b>
<b>Li EC et al. 2014 [42]</b>	((((((("Angiotensin Receptor Antagonists"[Mesh]) OR (((((((((((abitesartan) OR azilsartan) OR candesartan) OR elisartan) OR embusartan) OR eprosartan) OR forasartan) OR irbesartan) OR losartan) OR milfasartan) OR olmesartan) OR saprisartan) OR tasosartan) OR telmisartan) OR valsartan) OR zolasartan))) AND (((("Angiotensin-Converting Enzyme Inhibitors"[Mesh]) OR angiotensin converting enzyme inhibit*) OR (((((((((((((((((((((((acei) OR alacepril) OR altiopril) OR ancovenin) OR benazepril*) OR captopril) OR ceranapril) OR ceronapril) OR cilazapril*) OR deacetylalacepril) OR delapril) OR enalapril*) OR epicaptopril) OR fasidotril*) OR foroxymithine) OR fosinopril*) OR gemopatrilat) OR idapril) OR imidapril*) OR indolapril) OR libenzapril) OR lisinopril) OR moexipril*) OR moveltipril) OR omapatrilat) OR pentopril*) OR perindopril*) OR pivopril) OR quinapril*) OR ramipril*) OR rentiapril) OR saralasin) OR s nitrosocaptopril) OR spirapril*) OR temocapril*) OR teprotide) OR trandolapril*) OR utibapril*) OR zabicipril*) OR zofenopril*)) AND (((hypertension) OR hypertens*) OR "Blood Pressure"[Mesh])) AND (((("Randomized Controlled Trial" [Publication Type]) OR "Controlled Clinical Trial" [Publication Type]) OR randomi*[Title/Abstract]) OR placebo[Title/Abstract]) OR "Clinical Trials as Topic"[Mesh]) OR randomly[Title/Abstract]) OR trial[Title])) AND "Humans"[Mesh]) AND ( "0001/01/01"[PDat] : "2014/02/15"[PDat] )	((("0001/01/01"[Date - Publication] : "2014/02/15"[Date - Publication])) NOT (((((((("Angiotensin Receptor Antagonists"[Mesh]) OR (((((((((((abitesartan) OR azilsartan) OR candesartan) OR elisartan) OR embusartan) OR eprosartan) OR forasartan) OR irbesartan) OR losartan) OR milfasartan) OR olmesartan) OR saprisartan) OR tasosartan) OR telmisartan) OR valsartan) OR zolasartan))) AND (((("Angiotensin-Converting Enzyme Inhibitors"[Mesh]) OR angiotensin converting enzyme inhibit*) OR (((((((((((((((((((((((acei) OR alacepril) OR altiopril) OR ancovenin) OR benazepril*) OR captopril) OR ceranapril) OR ceronapril) OR cilazapril*) OR deacetylalacepril) OR delapril) OR enalapril*) OR epicaptopril) OR fasidotril*) OR foroxymithine) OR fosinopril*) OR gemopatrilat) OR idapril) OR imidapril*) OR indolapril) OR libenzapril) OR lisinopril) OR moexipril*) OR moveltipril) OR omapatrilat) OR pentopril*) OR perindopril*) OR pivopril) OR quinapril*) OR ramipril*) OR rentiapril) OR saralasin) OR s nitrosocaptopril) OR spirapril*) OR temocapril*) OR teprotide) OR trandolapril*) OR utibapril*) OR zabicipril*) OR zofenopril*)) AND (((hypertension) OR hypertens*) OR "Blood Pressure"[Mesh])) AND (((("Randomized Controlled Trial" [Publication Type]) OR "Controlled Clinical Trial" [Publication Type]) OR randomi*[Title/Abstract]) OR placebo[Title/Abstract]) OR randomly[Title/Abstract]) OR trial[Title]))	n=1441 total records n=0 manually added	n=556668 total records
		<b>Citations finally included in the main database</b>	<b>n=9 positives</b>	<b>n=200 negatives</b>
<b>Cavender MA et al. 2014[43]</b>	(((((bivalirudin) OR Angiomax) OR Hirulog)) AND ((((((stent) OR percutaneous coronary intervention) OR acute coronary syndromes) OR st-elevation myocardial infarction) OR non-ST-elevation myocardial infarction) OR unstable angina)) AND ( "0001/01/01"[PDat] : "2014/04/09"[PDat] )	((("0001/01/01"[Date - Publication] : "2014/04/09"[Date - Publication])) NOT (((((((bivalirudin) OR Angiomax) OR Hirulog)) AND ((((((stent) OR percutaneous coronary intervention) OR acute coronary syndromes) OR st-elevation myocardial infarction) OR non-ST-elevation myocardial infarction) OR unstable angina))))	n=745 total records n=1 manually added	n=561617 total records
		<b>Citations finally included in the main database</b>	<b>n=14 positives</b>	<b>n=400 negatives</b>
<b>Chatterjee S et al. 2014 [44]</b>	((((((("Rivaroxaban"[Mesh]) OR dabigatran) OR "apixaban" [Supplementary Concept]) OR "new oral anticoagulants" OR "oral thrombin inhibitors" OR "oral factor Xa inhibitors") AND ( "2001/01/01"[PDat] : "2013/09/15"[PDat] ) AND English[lang]	((("0001/01/01"[Date - Publication] : "2001/01/01"[Date - Publication])) NOT (((((((("Rivaroxaban"[Mesh]) OR dabigatran) OR "apixaban" [Supplementary Concept]) OR "new oral anticoagulants" OR "oral thrombin inhibitors" OR "oral factor Xa inhibitors"))	n=2034 total records n=0 manually added	n=223263 total records
		<b>Citations finally included in the main database</b>	<b>n=18 positives</b>	<b>n=400 negatives</b>
<b>Funakoshi T et al. 2014 [45]</b>	((((((((((axitinib) OR cabozantinib) OR erlotinib) OR gefitinib) OR lapatinib) OR pazopanib) OR regorafenib) OR sorafenib) OR sunitinib) OR vandetanib)) AND "Randomized Controlled Trial" [Publication Type]) AND ( "1966/01/01"[PDat] : "2013/03/31"[PDat] ) AND English[lang]	((("0001/01/01"[Date - Publication] : "2013/03/31"[Date - Publication])) NOT (((((((((((axitinib) OR cabozantinib) OR erlotinib) OR gefitinib) OR lapatinib) OR pazopanib) OR regorafenib) OR sorafenib) OR sunitinib) OR vandetanib)) AND "Randomized Controlled Trial" [Publication Type]))	n=418 total records n=5 manually added	n=527068 total records
		<b>Citations finally included in the main database</b>	<b>n=43 positives</b>	<b>n=1000 negatives</b>

**Table 2.** Number of training (PubMed) and testing (ClinicalTrials.gov) positive and negative records as well as the number of predicted positives and the relevant statistics for each systematic review considered (AUC = area under the receiver operator characteristic curve; PREV = prevalence of positive in ClinicalTrials.gov; PPV = positive predictive value; SENS = sensitivity; SPEC = specificity; LR+ = positive likelihood ratio LR- = negative likelihood ratio).

Systematic Review	Training positives	Training negatives	Testing positives	Testing negatives	Predicted positives	AUC	PPV	SENS	SPEC	LR+	LR-
Yang Q [33]	18	400	5	233604	1718	0.9963	0.0029	1	0.9927	136.9863	0
Meng Y [34]	9	200	4	233605	462	0.9990	0.0087	1	0.9980	500.0000	0
Segelov E [22]	13	400	8	233601	1595	0.9341	0.0044	0.875	0.9932	128.6765	0.1259
Li DH [35]	6	200	3	233606	1635	0.9965	0.0018	1	0.9930	142.8571	0
Lv ZC [36]	12	400	3	233606	1429	0.9969	0.0021	1	0.9939	163.9344	0
Wang J [37]	32	800	5	233604	1901	0.9959	0.0026	1	0.9919	123.4568	0
Zhou CQ [38]	9	200	3	233606	1011	0.9978	0.0030	1	0.9957	232.5581	0
Liu X [39]	23	600	30	233579	2178	0.9954	0.0138	1	0.9908	108.6957	0
Douxfls J [40]	13	400	10	233599	378	0.9992	0.0265	1	0.9984	625.0000	0
Kourbeti IS [41]	75	1600	25	233584	1843	0.9961	0.0136	1	0.9922	128.2051	0
Li EC [42]	9	200	2	233607	6558	0.9860	0.0003	1	0.9719	35.5872	0
Cavender MA [43]	14	400	7	233602	149	0.9997	0.0470	1	0.9994	1666.6667	0
Chatterjee S [44]	18	400	17	233592	771	0.9984	0.0220	1	0.9968	312.5000	0
Funakoshi T [45]	43	1000	11	233598	3851	0.9918	0.0029	1	0.9836	60.9756	0

**Table 3.** The number of predicted positives and true positives in manual and automated searches after filter application. Records of the manual search are those retrieved on ICTRP by Baudard and colleagues [10]. Records of the automated search are those retrieved on ClinicalTrials.gov using our ML instrument. Predicted positives are a pool of citations resulting from manual search strings or from automated searches. True positives are clinical trials added by Baudard and colleagues to each systematic review. The description of filters reports data element entries and the number of retrieved records. Filters were applied sequentially from Filter 0 to Filter 5.

Systematic Review	Manual search		Automated search												
	Predicted positives	True positives	Filter 0		Filter 1		Filter 2		Filter 3		Filter 4		Filter 5		All
			none	Predicted positives	study_type	Predicted positives	overall_status	Predicted positives	start_before	Predicted positives	primary_completion before/within	Predicted positives	specific filters	Predicted positives	True positive
Yang Q [33]	12	1	-	1718	interventional	1341	completed OR terminated	759	April 2014	705	April 2014	588	allocation = randomized number_of_arms ≠ 1	457	1
Meng Y [34]	26	1	-	462	interventional	399	completed OR terminated	282	August 2013	243	August 2013	202	allocation = randomized number_of_arms ≠ 1	144	1
Segelov E [22]	684	2	-	1595	interventional	1432	completed OR terminated	836	May 2012	770	May 2012	588	allocation = randomized number_of_arms ≠ 1	274	2
Li DH [35]	201	1	-	1635	interventional	1545	completed OR terminated	376	February 2014	837	February 2014	695	allocation = randomized number_of_arms ≠ 1 phase = 2 OR 3	289	1
Lv ZC [36]	665	1	-	1429	interventional	1294	completed OR terminated	727	February 2014	716	February 2014	583	allocation = randomized number_of_arms ≠ 1 minimum_age ≥ 18 years	243	1
Wang J [37]	227	1	-	1901	interventional	1690	completed OR terminated	1191	December 2013	1118	December 2013	972	allocation = randomized number_of_arms ≠ 1 intervention_model = parallel OR crossover	729	1
Zhou CQ [38]	3	1	-	1011	interventional	793	completed OR terminated	534	February 2013	468	February 2013	372	allocation = randomized number_of_arms ≠ 1	263	1
Liu X [39]	1661	21	-	2178	interventional	2028	completed OR terminated	1587	January 2013	1317	January 2013	1112	allocation = randomized number_of_arms ≠ 1 minimum_age ≥ 18 years phase ≥ 3	622	21
Douxifls J [40]	76	1	-	378	interventional	270	completed OR terminated	190	December 2013	174	December 2013	150	allocation = randomized number_of_arms ≠ 1	116	1
Kourbeti IS [41]	581	4	-	1843	interventional	1449	completed OR terminated	1023	June 2013	941	June 2013	756	allocation = randomized number_of_arms ≠ 1 is_fda_regulated = true	409	4
Li EC [42]	909	2	-	6558	interventional	5483	completed OR terminated	3629	January 2014	3412	January 2014	2771	allocation = randomized number_of_arms ≠ 1	2119	2
Cavender MA [43]	71	1	-	149	interventional	130	completed OR terminated	87	April 2014	84	April 2014	74	allocation = randomized number_of_arms ≠ 1	60	1
Chatterjee S [44]	217	1	-	771	interventional	509	completed OR terminated	279	March 2014	263	January 2001 - March 2014	207	allocation = randomized number_of_arms ≠ 1	169	1
Funakoshi T [45]	2680	2	-	3851	interventional	3762	completed OR terminated	2147	February 2014	2111	January 2004 - February 2014	1699	allocation = randomized number_of_arms ≠ 1 phase = 2 OR 3	711	2

**Figure 1.** General procedure workflow.

**What is new:**

- **Key findings:** Our study proposes a new classifier that extends PubMed searches to clinical trial registries with high discrimination ability and sensitivity.
- **What this adds to what is known:** Text mining and machine learning techniques can be used to support comprehensive systematic reviews.
- **What is the implication, what should change now:** The proposed machine learning instrument can help researchers identify relevant studies in the systematic review process by reducing workload without losing sensitivity and at a small price in terms of specificity.