

Text mining for identification of biological entities related to antibiotic resistant organisms

Kelle Fortunato Costa¹, **Fabício Almeida Araújo**^{2,3}, **Jefferson Moraes**⁴, **Carlos Renato Lisboa Frances**¹, **Rommel T J Ramos**^{Corresp. 5}

¹ Programa de pós-graduação em Engenharia Elétrica, Universidade Federal do Pará, Belém, Pará, Brazil

² Biological Science Institute, Universidade Federal do Pará, Belém, Pará, Brazil

³ Universidade Federal Rural da Amazônia, Belém, Pará, Brazil

⁴ Universidade Federal do Pará, Belém, Pará, Brazil

⁵ Biological Science Institute, Universidade Federal do Para, Belém, Pará, Brazil

Corresponding Author: Rommel T J Ramos

Email address: rommelramos@ufpa.br

Antimicrobial resistance is a significant public health problem worldwide. In recent years, the scientific community has been intensifying efforts to combat this problem; many experiments have been developed, and many articles are published in this area. However, the growing volume of biological literature increases the difficulty of the biocuration process due to the cost and time required. Modern text mining tools with the adoption of artificial intelligence technology are helpful to assist in the evolution of research. In this article, we propose a text mining model capable of identifying and ranking prioritizing scientific articles in the context of antimicrobial resistance. We retrieved scientific articles from the PubMed database, adopted machine learning techniques to generate the vector representation of the retrieved scientific articles, and identified their similarity with the context. As a result of this process, we obtained a dataset labeled "Relevant" and "Irrelevant" and used this dataset to implement one supervised learning algorithm to classify new records. The model's overall performance reached 90% accuracy and the f-measure (harmonic mean between the metrics) reached 82% accuracy for positive class and 93% for negative class, showing quality in the identification of scientific articles relevant to the context. The dataset, scripts and models are available at <https://github.com/engbiopct/TextMiningAMR>.

Text Mining for Identification of Biological Entities Related to Antibiotic Resistant Organisms

Kelle Fortunato da Costa¹, Fabricio Almeida Araujo^{2,3}, Jefferson Moraes⁴, Renato Francês⁵ and Rommel Ramos⁶

¹ Universidade Federal do Pará, Belém, Pará, Brasil

² Computational Biology, Universidade Federal do Pará, Belém, Pará, Brasil

³ Universidade Federal do Rural da Amazônia, Belém, Pará, Brasil

⁴ Universidade Federal do Pará, Belém, Pará, Brasil

⁵ Universidade Federal do Pará, Belém, Pará, Brasil

⁶ Biological Science Institute, Universidade Federal do Pará, Belém, Pará, Brasil

Corresponding Author:

Rommel Ramos

R. Augusto Corrêa, 01 - Guamá, Belém - PA, 66075-110, Brasil

E-mail address: rommelramos@ufpa.br

Abstract

Antimicrobial resistance is a significant public health problem worldwide. In recent years, the scientific community has been intensifying efforts to combat this problem; many experiments have been developed, and many articles are published in this area. However, the growing volume of biological literature increases the difficulty of the biocuration process due to the cost and time required. Modern text mining tools with the adoption of artificial intelligence technology are helpful to assist in the evolution of research. In this article, we propose a text mining model capable of identifying and ranking prioritizing scientific articles in the context of antimicrobial resistance. We retrieved scientific articles from the PubMed database, adopted machine learning techniques to generate the vector representation of the retrieved scientific articles, and identified their similarity with the context. As a result of this process, we obtained a dataset labeled "Relevant" and "Irrelevant" and used this dataset to implement one supervised learning algorithm to classify new records. The model's overall performance reached 90% accuracy and the f-measure (harmonic mean between the metrics) reached 82% accuracy for positive class and 93% for negative class, showing quality in the identification of scientific articles relevant to the context. The dataset, scripts and models are available at <https://github.com/engbiopct/TextMiningAMR>.

Introduction

Antibiotics are the most successful drugs of the last 100 years, responsible for saving countless lives and enabling modern medical procedures that would otherwise be unthinkable. However, all antibiotics derived from secondary or fully synthetic microbial metabolism products are subject to resistance [1].

Antimicrobial resistance (AMR) has been increasingly recognized as an important public health problem worldwide, considering that infections caused by multidrug-resistant organisms (MDR) result in a significant increase in mortality and cause a tremendous economic burden [2]. In addition to the costs of rising hospital admission rates, it is estimated that by 2050 there will be an economic loss of \$100 trillion in global antibiotic production [3].

In recent years, the scientific community has intensified efforts to combat this problem by making available a wide range of public databases specific to AMR, such as: National Database of Antibiotic Resistant Organisms (NDARO) [4], Comprehensive Antibiotic Resistance Database (CARD) [5], Resfinder [6], ResfinderFG [7], Resfams [8], Antibiotic Resistance Genes Database (ARDB) [9], MEGARes [10], Antibiotic Resistance Gene Annotation (ARG-ANNOT) [11], Mustard [12], Functional Antibiotic Resistance Metagenomic Element (FARME database) [13], SARG (v2) [14], Lahey list of β -lactamases [15], β -Lactamase Database (BLDB) [16], Lactamase Engineering Database (LacED) [17, 18], Comprehensive β -Lactamase Molecular Annotation Resource (CBMAR) [19], among others, which are frequently used as reference databases, with gene sequences related to antimicrobial resistance and metadata that enrich the characterization of sequences. Due to the increase in the volume of literature related to biological and health sciences, the curation process has become challenging for researchers and biocurators who use these databases as a source of research, mainly due to the time required to locate relevant information about biological entities related to antibiotic-resistant organisms. Even queries in specialized databases in biomedical literature such as PubMed (scientific and medical abstracts/citations), PubMed Central (full-text journal articles), NLM Catalog (index of NLM collections), Books (books and reports), and MeSH (ontology used for PubMed indexing) [20], tends to make document selection difficult due to a large amount of retrieved items.

In this context, the adoption of text mining (TM) techniques are viable alternatives [21], as they can help in different stages of the standard biocuration workflow. According to [22], the steps are:

- Selection: search for articles relevant to the curation.
- Identification and standardization of bioentities: detection of mentions of bioentities relevant to the curatorship; for example, genes, proteins, or small molecules, linked to unique identifiers from databases such as UniProt, EntrezGene, or ChEBI.
- Detection of annotation events: identification and encoding of events, such as descriptions of protein-protein interactions, characterizations of gene products in terms of cellular location, molecular function, involvement in the biological process and phenotypic effect.
- Evidence qualifier association: association of experimental evidence that supports the annotation event performed due to biocuration efforts.
- Completion and verification of the database record.

TM technologies combine knowledge resources such as controlled vocabularies, taxonomies, and ontologies with linguistic analysis and machine learning to deal with language variations and extract not only terms from the text but also relationships between terms [23].

In the last decade, several applications in the biomedical area were developed [24 - 50] using one or more of the following TM steps: (i) retrieving textual resources relevant to a particular subject of interest, a process known as information retrieval (IR), (ii) detect the occurrence of specific keywords of interest and the relationships between these keywords and (iii) infer new relationships based on known facts, and this step is called knowledge discovery (KD) [51].

One of the most used machine learning techniques in knowledge discovery, especially in document screening, is text classification. However, supervised classification requires the prior labeling of a training set, a non-trivial task for human curators, as it requires a lot of time and effort.

In this sense, Suomela and Andrade [52] propose a methodology for automatic (binary) classification of large volumes of data, adopting the word counting technique, known as the bag of words (BOW) [53], a textual representation that composes the vector space model [54], where documents are converted into vectors of words. A weighting scheme is applied to each word, which can be a simple word count or a metric such as Term Frequency - Inverse Document Frequency (tf -idf) [53][55][56] and based on the arithmetic mean of the weights of these words, text summaries are classified as relevant or irrelevant.

This methodology was implemented from specific abstracts of an area of interest, extracted from Pubmed, and inspired the development of the MedlineRanker application [57] and served as a baseline for this study, which instead of the bag of words, adopts a representation approach based on neural networks [58], called Paragraph Vector – Distributed Memory (PV-DM) [59], capable of revealing semantic characteristics between documents, a property that makes this approach useful for many natural language processing (NLP) tasks and justifies its wide use in works involving natural language understanding [60; 61], machine translation [62; 63], image comprehension [64] and relational extraction [65].

The study of antimicrobial resistance genes is essential to public health. However, there is challenging to handle and extract the amount of available data of scientific and medical manuscripts from public databases without computational methods. Thus, this work proposes an unsupervised learning-based TM approach for ranking the relevance of articles on AMR context to generate a set of training, accurate enough to generalize new data, maximizing the efficiency of the supervised classifiers.

Materials & Methods

1. Labeling pipeline

Figure 1 shows the TM steps implemented in this work in order to label the data.

Figure 1. Proposed TM model. Steps (A) and (B) include retrieving the information. Steps (C) and (D) include the recognition of entities and the discovery of knowledge, resulting in a metric (cosine similarity) responsible for determining the binary classification performed in step (E).

1.1. Information retrieval

An Application Programming Interface (API) was implemented to retrieve a collection of relevant articles in the Drug Resistance, and Microbial domain through the Pubmed Central (PMC) database, which contains free full-text files of the Library's of Medicine and the US National Institutes of Health (NIH/NLM) biomedical literature are available. In the API, the E-Search and E-Fetch tools from the E-utilities package were used, which provide a structured interface for accessing the Entrez system, the NCBI database system, which currently includes 38 databases, covering a variety of biomedical data, including nucleotide and protein sequences, gene records, three-dimensional molecular structures, and biomedical literature [66]. Table 1 presents the set of parameters incorporated into the search (E-Search).

Table 1. Parameters *E-Search* PubMed Central

The terms of the MeSH hierarchy were adopted for antimicrobial resistance, considering that the terms MeSH is a controlled vocabulary of biomedical terms whose elements are assigned to a document by indexers (specialists in biomedical subjects) based on its context. They contain high-density document-wide information that cannot be deduced from the title or abstract that PubMed returns using keywords [67].

Then, a list of PMCIDs (unique identifiers provided by PubMed Central to each document) is generated to be used to access the full texts of articles through the E-Fetch utility.

1.2. Named entity recognition and Knowledge Discovery

The entity recognition and knowledge discovery process consist of two steps: in the first step (Figure 1-C), the Doc2Vec unsupervised learning algorithm from the Gensim library was used, which implements the Paragraph Vector – Distributed Memory model, to obtain the embedding of the retrieved documents (dense representation of a sequence of words). Table 2 displays the parameters used in the algorithm.

Table 2. Parameters Doc2Vec algorithm

In the second step (Figure 1-D), the pre-trained model was used, capable of predicting whether a set of documents $\{Doc_1, Doc_2, Doc_3...Doc_n\}$ belong to the context of a central document Doc_0 , to infer the similarity of the documents to the AMR context, represented by 4,290 terms extracted from CARD and the Gene Ontology Database [68], calculating the cosine distance [69] between them. The resulting value varies in the range between -1 and 1 where the higher the number, the greater the similarity with the context.

Finally, each of the scientific articles was automatically labeled (Figure 1-E), considering the arithmetic average of the cosine similarity, values above the average were defined as relevant and below the mean as irrelevant.

2. Evaluation of the proposed method

To evaluate the classification performance of the proposed TM approach (Figure 2 - A), the same dataset was labeled using the Bag of words text representation model (Figure 2 - B), adopting the specific vocabulary for AMR as a vector of features, and the documents were classified as relevant or irrelevant through the arithmetic mean of the weights of the words, obtained with the TF-IDF weighting method, similar to the methodology proposed in [52]. Finally, the two (automatically) labeled databases were compared with a test database labeled by experts (Figure

2 - C), who independently labeled the articles as relevant or irrelevant. Only the samples where the three experts converged on the labels were included in the test dataset.

Figure 2. Evaluation of the proposed method.

3. Predictions with automatically labeled data

To evaluate the efficiency of the proposed approach, which uses neural embeddings for labeling, the generated datasets (Figure 3 - Dataset_1 and Figure 3 - Dataset_2) were used as input to the supervised classifier SVM [70; 71]. The classification performance was evaluated through the analysis of the Precision, Recall, Accuracy, and F-Measure metrics [72], calculated from the test base labeled by experts and not used to train the SVM models (Figure 3).

Figure 3. Performance of predictions with automatically labeled data.

For the SVM classifier, the feature/attribute vector (AMR vocabulary with 4,290 words) was weighted using the TF-IDF technique, and cross-validation, with 5-folds (default value of the adopted algorithm), was applied to explore the combination of parameters for determining the best model, as the effectiveness of the SVM depends on the kernel selection, which is the function that will be used by the algorithm, in the margin parameter (C), which determines a balance between maximizing the margin and minimizing classification errors, and the Gamma parameter, when the chosen kernel is Gaussian (or RBF) [73], adopted in this experiment.

4. Data Availability

The dataset, script, and models generated by this work are available at <https://github.com/engbiopct/TextMiningAMR>, under CC BY 4.0 Copyright license, with information regards the workflow adopted, a short step-by-step guide to the readers reproduce this experiment and the complementary materials.

Results and Discussion

1. Labeling

A collection of 88,300 scientific articles on antimicrobial resistance was retrieved from Pubmed Central, using the terms of reference of the MeSH (Medical Subject Headings, developed at the National Library of Medicine) hierarchical vocabulary referring to the AMR domain (<https://meshb.nlm.nih.gov/record/ui?ui=D004352>).

The retrieved dataset was submitted to the PV-DM text representation model, with the embedding of the documents obtained. The similarity of the documents to the AMR context (Table S1) was inferred, and the label "relevant" was automatically assigned (class 0) to all articles whose cosine distance value was equal to or greater than the arithmetic average of the cosine distances of the entire corpus, and the label "irrelevant" (class 1) to all articles with a cosine distance value lower than that referred to average, resulting in 43,136 records labeled relevant and 45,164 labeled irrelevant (Table S2).

The same initial dataset was submitted to the Bag of words text representation model in order to obtain the weights of the words in the documents according to the AMR dictionary and thus

automatically assign the label "relevant" (class 0) to all articles with weight equal to or greater than the arithmetic average of the weights, and the label "irrelevant" (class 1) to all articles with a value lower than the mean, resulting in 45,946 records labeled as relevant and 42,354 labeled as irrelevant (Table S3).

With the two labeled datasets, the results were compared with a test dataset labeled by experts, consolidated with a total of 62 scientific articles, 15 labeled as relevant and 47 labeled as irrelevant (Table S4).

In the comparison, the proposed method labeled 44 articles according to the experts, which represents 71% of hits in general, with 80% of hits for the relevant label and 68% of hits for the irrelevant label. As for the baseline method, there were only 26 labels according to the experts, which represents 42% of overall performance, with 66% of correctness for the relevant label and 34% of correctness for the irrelevant label.

The proposed approach presents a superior performance about the baseline, which despite its simplicity, efficiency, and often surprising precision, does not take into account the order and semantics of the words, that is, the distances between them. This means that the words "mighty", "strong" and "Paris" are equally distant. Although semantically, "powerful" is closer to "strong" than "Paris" [59], characteristics present in the PV-DM model and fundamental human skills in manual data labeling tasks.

2. Classification

The two labeled datasets were submitted to the supervised SVM classifier, excluding the test dataset labeled by experts from training.

Figure 4 presents the confusion matrix of the SVM_1 classifier, trained with data from dataset 1 (Figure 3 - Dataset_1). There is a high degree of precision both in terms of true positives (relevant articles classified as relevant) and true negatives (non-relevant articles classified as non-relevant).

Only 1 article was incorrectly classified as relevant (false positive), and 5 articles were incorrectly classified as irrelevant (false negative).

Figure 4. SVM classifier confusion matrix for dataset_1 (PV-DM)

Figure 5 presents the confusion matrix of SVM_2, trained with data from dataset 2 (Figure 3 - Dataset_2), where a lower degree of precision is observed for both true positives and true negatives in relation to SVM_1. With this classifier, however, there was an increase in the number of incorrect classifications, with 33 articles incorrectly classified as relevant (false positives) and 6 articles incorrectly classified as irrelevant (false negatives).

Figure 5. SVM classifier confusion matrix for dataset_1 (Bag of Words)

Table 3 presents the results of the evaluation metrics: precision, recall, accuracy and the f-measure, obtained based on the results of the confusion matrix of the two classifiers. The results of the SVM_1 classifier were superior to the SVM_2 classifier in all evaluated metrics.

Table 3. Classifier Performance Assessment

Accuracy, a metric that represents the overall performance of a model, reached 90% of accuracy and the f-measure, which is a harmonic average between the precision and recall metrics and that

can be used as a single measure to represent the quality in the text mining [74] reached 82% success rate for the positive class and 93% for the negative class.

The results show that the best performance was obtained with the database labeled by the PV-DM model.

Table 4. Results of Labeling and Classification steps vs Experts

Table 4 presents the percentage of correct predictions, both in the labeling and in the classification stage, in comparison with the data labeled by experts and validates the hypothesis that the use of Paragraph Vector, Distributed Representations of Sentences and Documents associated with similarity with a specific context is able not only to perform the binary classification of large volumes of data but also to optimize the percentage of hits of supervised classifiers. The SVM_2 classifier showed a reduction in the number of hits compared to the Labeling step, although we adopted the same attribute vector and the same representation in both experiments (bag of words, weighted with TF-IDF).

Conclusions

The proposed TM approach proved capable of identifying and prioritizing documents in the AMR context, as well as predicting the relevance of new documents in the same context. For this, we used the TM steps summarized in [51], plus some specifics of the proposal such as (1) Use of the MeSH hierarchy; (2) Use of full text; (3) Use of domain-specific dictionaries (CARD and Gene Ontology), fundamental in this process, as it facilitated the detection of similarity of articles to the AMR context and (4) Adoption of unsupervised learning for better representation of texts. Additionally, we submitted the labeled bases to the SVM classifier to evaluate their performance in comparison to the test base labeled by human experts.

The proposed approach efficiently identifies scientific articles relevant to the AMR context. Therefore, it is a valuable tool to automate information capture processes from robust bibliographic reference databases such as Pubmed Central, as well as to accelerate the screening of documents with biocuration potential, facilitating the other stages of the biocuration process. This work presents a new set of pre-trained document embeddings in the AMR domain and a base labeled for relevance according to similarity to CARD and Gene Ontology, which, in future work, can be used as input to other algorithms of supervised learning and for biocuration tools aimed at the identification and normalization of bioentities, detection of annotation events and filling out specific databases for AMR. This proposal is limited to the previous existence of an accurate database representing the main terms related to the target.

Acknowledgements

References

- [1] G. D. Wright, Molecular mechanisms of antibiotic resistance. *Chem. Commun. (Camb.)* 47, 4055–061 (2011). 10.1039/c0cc05111j pmid:21286630 doi:10.1039/c0cc05111j.

- [2] Tran TT, Munita JM, Arias CA. Mechanisms of drug resistance: daptomycin resistance. *Ann N Y Acad Sci.* 2015; 1354:32–53.
- [3] Review on Antimicrobial Resistance. 2016. Antimicrobial resistance: TACKLING DRUG-Resistant infections globally: final report and recommendations. accessed march 11, 2020.
- [4] NLM Annual Report 2016 - National Library of Medicine – NIH.
- [5] Alcock et al . CARD 2020: antibiotic resistance surveillance with the comprehensive antibiotic resistance database, *Nucleic Acids Research*, Volume 48, Issue D1, 08 January 2020, Pages D517–D525, <https://doi.org/10.1093/nar/gkz935>.
- [6] Zankari E et al. Identification of acquired antimicrobial resistance genes. *J. Antimicrob. Chemother* 67, 2640–2644 (2012).
- [7] Munk P et al. Abundance and diversity of the faecal resistome in slaughter pigs and broilers in nine European countries. *Nat. Microbiol* 3, 898–908 (2018).
- [8] Gibson MK, Forsberg KJ & Dantas G Improved annotation of antibiotic resistance determinants reveals microbial resistomes cluster by ecology. *ISME J* 9, 207–216 (2015).
- [9] Liu B & Pop M ARDB—antibiotic resistance genes database. *Nucleic Acids Res* 37, D443–D447 (2009).
- [10] Lakin SM et al. MEGARes: an antimicrobial resistance database for high throughput sequencing. *Nucleic Acids Res* 45, D574–D580 (2017).
- [11] Gupta SK et al. ARG-ANNOT, a new bioinformatic tool to discover antibiotic resistance genes in bacterial genomes. *Antimicrob. Agents Chemother* 58, 212–220 (2014).
- [12] Ruppe E et al. Prediction of the intestinal resistome by a three-dimensional structure-based method. *Nat. Microbiol* 4, 112–123 (2019)
- [13] Wallace JC, Port JA, Smith MN & Faustman EM FARME DB: a functional antibiotic resistance element database. *Database (Oxford)* 2017, baw165 (2017).
- [14] Yin X et al. ARGs-OAP v2.0 with an expanded SARG database and hidden Markov models for enhancement characterization and quantification of antibiotic resistance genes in environmental metagenomes. *Bioinformatics* 34, 2263–2270 (2018).
- [15] Bush K & Jacoby GA Updated functional classification of β -lactamases. *Antimicrob. Agents Chemother* 54, 969–976 (2010).
- [16] Naas T, Oueslati S, Bonnin RA, Dabos ML, Zavala A, Dortet L, Retailleau P, Iorga BI. Beta-lactamase database (BLDB) - structure and function. *J Enzyme Inhib Med Chem.* 2017 Dec;32(1):917-919. doi: 10.1080/14756366.2017.1344235. PMID: 28719998; PMCID: PMC6445328.
- [17] Thai QK, Bos F & Pleiss J The lactamase engineering database: a critical survey of TEM sequences in public databases. *BMC Genomics* 10, 390 (2009).
- [18] Thai QK & Pleiss J SHV lactamase engineering database: a reconciliation tool for SHV beta-lactamases in public databases. *BMC Genomics* 11, 563 (2010).
- [19] Srivastava A, Singhal N, Goel M, Virdi JS & Kumar M CBMAR: a comprehensive beta-lactamase molecular annotation resource. *Database (Oxford)* 2014, bau111 (2014).

- [20] Sayers, Eric W et al. “Database resources of the National Center for Biotechnology Information.” *Nucleic acids research* vol. 47,D1 (2019): D23-D28. doi:10.1093/nar/gky1069
- [21] Hirschman L, Burns GA, Krallinger M, Arighi C, Cohen KB, Valencia A, Wu CH, Chatr-Aryamontri A, Dowell KG, Huala E, et al. Text mining for the biocuration workflow, *Database*, 2012, vol. 2012 April 18 (doi:10.1093/database/bas020; epub ahead of print).
- [22] Wei CH, et al. PubTator: a web-based text mining tool for assisting biocuration, *Nucleic Acids Res.*, 2013, vol. 41 Web server(pg. W518-W522).
- [23] Chaix Estelle, Deléger Louise, Bossy Robert, Nédellec Claire. Text mining tools for extracting information about microbial biodiversity in food. *Food Microbiology*. 2018 doi: 10.1016/j.fm.2018.04.011.
- [24] W.W. Fleuren et al., *Nucleic Acids Res.* 39 (Web Server issue) (2011) W450–W454.
- [25] P. Fontelo, F. Liu, M. Ackerman, *BMC Med. Inf. Decis. Mak.* 5 (2005) 5.
- [26] C. Perez-Iratxeta, P. Bork, M.A. Andrade, *Trends Biochem. Sci.* 26 (9) (2001) 573–575.
- [27] J. Lewis et al., *Bioinformatics* 22 (18) (2006) 2298–2304.
- [28] J.F. Fontaine et al., *Nucleic Acids Res.* 37 (Web Server issue) (2009) W141– W146.
- [29] D.J. States et al., *Bioinformatics* 25 (7) (2009) 974–976.
- [30] K.C. Huang et al., *J. Biomed. Inf.* 46 (5) (2013) 940–946.
- [31] K. Hokamp, K.H. Wolfe, *Nucleic Acids Res* 32 (Web Server issue) (2004) W16–W19.
- [32] M.V. Plikus, Z. Zhang, C.M. Chuong, *BMC Bioinf.* 7 (2006) 424.
- [33] K.G. Becker et al., *BMC Bioinf.* 4 (2003) 61.
- [34] S.M. Douglas, G.T. Montelione, M. Gerstein, *Genome Biol.* 6 (9) (2005) R80.
- [35] B. Brancotte et al., *Bioinformatics* 27 (8) (2011) 1187–1189.
- [36] S. De et al., *Physiol. Genomics* 42A (2) (2010) 162–167.
- [37] N.R. Smalheiser, W. Zhou, V.I. Torvik, *J. Biomed. Discov. Collab.* 3 (2008) 2.
- [38] H. Chen, B.M. Sharp, *BMC Bioinf.* 5 (2004) 147.
- [39] C. Li et al., *Database (oxford)* 2013 (2013) bat030.
- [40] R.W. Glynn, M.J. Kerin, K.J. Sweeney, *Br. J. Surg.* 97 (8) (2010) 1304–1308.
- [41] W. Xuan et al., *Comput. Syst. Bioinf. Conf.* 6 (2007) 359–369.
- [42] E. Giglia, *Eur. J. Phys. Rehabil. Med.* 47 (4) (2011) 687–690.
- [43] Y. Tsuruoka et al., *Bioinformatics* 27 (13) (2011) i111–i119.
- [44] J.M. Fernandez, R. Hoffmann, A. Valencia, *Nucleic Acids Res.* 35 (Web Server issue) (2007) W21–W26.
- [45] K. Raja, S. Subramani, J. Natarajan, *Database (Oxford)* 2013 (2013) bas052.
- [46] E. Pafilis et al., *Nat. Biotechnol.* 27 (6) (2009) 508–510.
- [47] D. Rebholz-Schuhmann et al., *Bioinformatics* 24 (2) (2008) 296–298.
- [48] C. Plake et al., *Bioinformatics* 22 (19) (2006) 2444–2445.
- [49] T.G. Soldatos et al., *Nucleic Acids Res.* 38 (1) (2010) 26–38.
- [50] A. Franceschini et al., *Nucleic Acids Res.* 41 (Database issue) (2013) D808–D815.
- [51] Fleuren, W. W. M. & Alkema, W. Application of text mining in the biomedical domain. *Methods* 74, 97–106 (2015).

- [52] Suomela BP, Andrade MA. Ranking the whole MEDLINE database according to a large training set using text indexing, BMC Bioinformatics, 2005, vol. 6 pg. 75
- [53] C. D. Manning, P. Raghavan, H. Schütze, Introduction to Information Retrieval, 1st Edition, Cambridge University Press, Cambridge, 2008.
- [54] SALTON, Gerard; WONG, Anita; YANG, Chung-Shu. A vector space model for automatic indexing. Communications of the ACM, v. 18, n. 11, p. 613-620, 1975.
- [55] J. H. Paik, A novel tf-idf weighting scheme for effective ranking, in: Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval, 2013, pp. 343–352.
- [56] C.-H. Chen, Improved tfidf in big news retrieval: An empirical study, Pattern Recognition Letters 93 (2017) 113–122.
- [57] Jean-Fred Fontaine, Adriano Barbosa-Silva, Martin Schaefer, Matthew R. Huska, Enrique M. Muro, Miguel A. Andrade-Navarro, MedlineRanker: ranking flexível de literatura biomédica, Nucleic Acids Research , Volume 37, Issue suppl_2, 1 de julho 2009, Pages W141 – W146, <https://doi.org/10.1093/nar/gkp353>
- [58] Bengio, Yoshua, Schwenk, Holger, Senécal, Jean-Sébastien, Morin, Frédéric, and Gauvain, Jean-Luc. A Neural probabilistic language models. In Innovations in Machine Learning, pp. 137–186. Springer, 2006.
- [59] LE, Quoc; MIKOLOV, Tomas. Distributed representations of sentences and documents. In: International Conference on Machine Learning. [S.l.: s.n.], 2014. p. 1188–1196.
- [60] Collobert, Ronan and Weston, Jason. A unified architecture for natural language processing: Deep neural networks with multitask learning. In Proceedings of the 25th International Conference on Machine Learning, pp. 160– 167. ACM, 2008.
- [61] Zhila, A., Yih, W.T., Meek, C., Zweig, G., and Mikolov, T. Combining heterogeneous models for measuring relational similarity. In NAACL HLT, 2013.
- [62] Mikolov, Tomas, Le, Quoc V., and Sutskever, Ilya. Exploiting similarities among languages for machine translation. CoRR, abs/1309.4168, 2013b.
- [63] Zou, Will, Socher, Richard, Cer, Daniel, and Manning, Christopher. Bilingual word embeddings for phrasebased machine translation. In Conference on Empirical Methods in Natural Language Processing, 2013.
- [64] Frome, Andrea, Corrado, Greg S., Shlens, Jonathon, Bengio, Samy, Dean, Jeffrey, Ranzato, Marc’Aurelio, and Mikolov, Tomas. DeViSE: A deep visual-semantic embedding model. In Advances in Neural Information Processing Systems, 2013.
- [65] Socher, Richard, Chen, Danqi, Manning, Christopher D., and Ng, Andrew Y. Reasoning with neural tensor networks for knowledge base completion. In Advances in Neural Information Processing Systems, 2013a.
- [66] NCIBI Homepage, <https://www.ncbi.nlm.nih.gov/books/NBK3827/>, acessado em 25/04/2021.
- [67] NLM: Medical Subject Headings (MeSH).
- [68] Ashburner et al. Gene ontology: tool for the unification of biology. Nat Genet. May 2000.

401 [69] H. Nguyen and L. Bai, “Cosine similarity metric learning for face verification,” in
 402 Computer Vision C ACCV 2010, ser. Lecture Notes in Computer Science, R. Kimmel, R. Klette,
 403 and A. Sugimoto, Eds. Springer Berlin Heidelberg, 2011, vol. 6493, pp. 709–720.
 404 [70] BOSER, B. E.; GUYON, I. M.; VAPNIK, V. N. A Training Algorithm for Optimal Margin
 405 Classifiers. In: ANNUAL WORKSHOP ON COMPUTACIONAL LEARNING, 5, 1992,
 406 Pittsburgh. ACM Press. Pittsburgh: Haussler D, jul 1992. p.144-152 .
 407 [71] DRUCKER, H.; BURGESS, C. J.; KAUFMAN, L.; SMOLA, A.; VAPNIK, V. Support
 408 vector regression machines. Advances in neural information processing systems, Morgan
 409 Kaufmann Publishers, p. 155–161, 1997.
 410 [72] PANG, B.; LEE, L. Opinion mining and sentiment analysis. Foundations and Trends in
 411 Information Retrieval, v. 2, n. 1-2, p. 1–135, 2008.
 412 [73] SYARIF, Iwan; PRUGEL-BENNETT, Adam; WILLS, Gary. SVM parameter optimization
 413 using grid search and genetic algorithm to improve classification performance. Telkomnika, v.
 414 14, n. 4, p. 1502, 2016.
 415 [74] Rodriguez-Esteban, R. (2009) Biomedical text mining and its applications. PLoS Comput.
 416 Biol. 5, e1000597.

Figure 1

Proposed TM model

Steps (A) and (B) include retrieving the information. Steps (C) and (D) include the recognition of entities and the discovery of knowledge, resulting in a metric (cosine similarity) responsible for determining the binary classification performed in step (E).

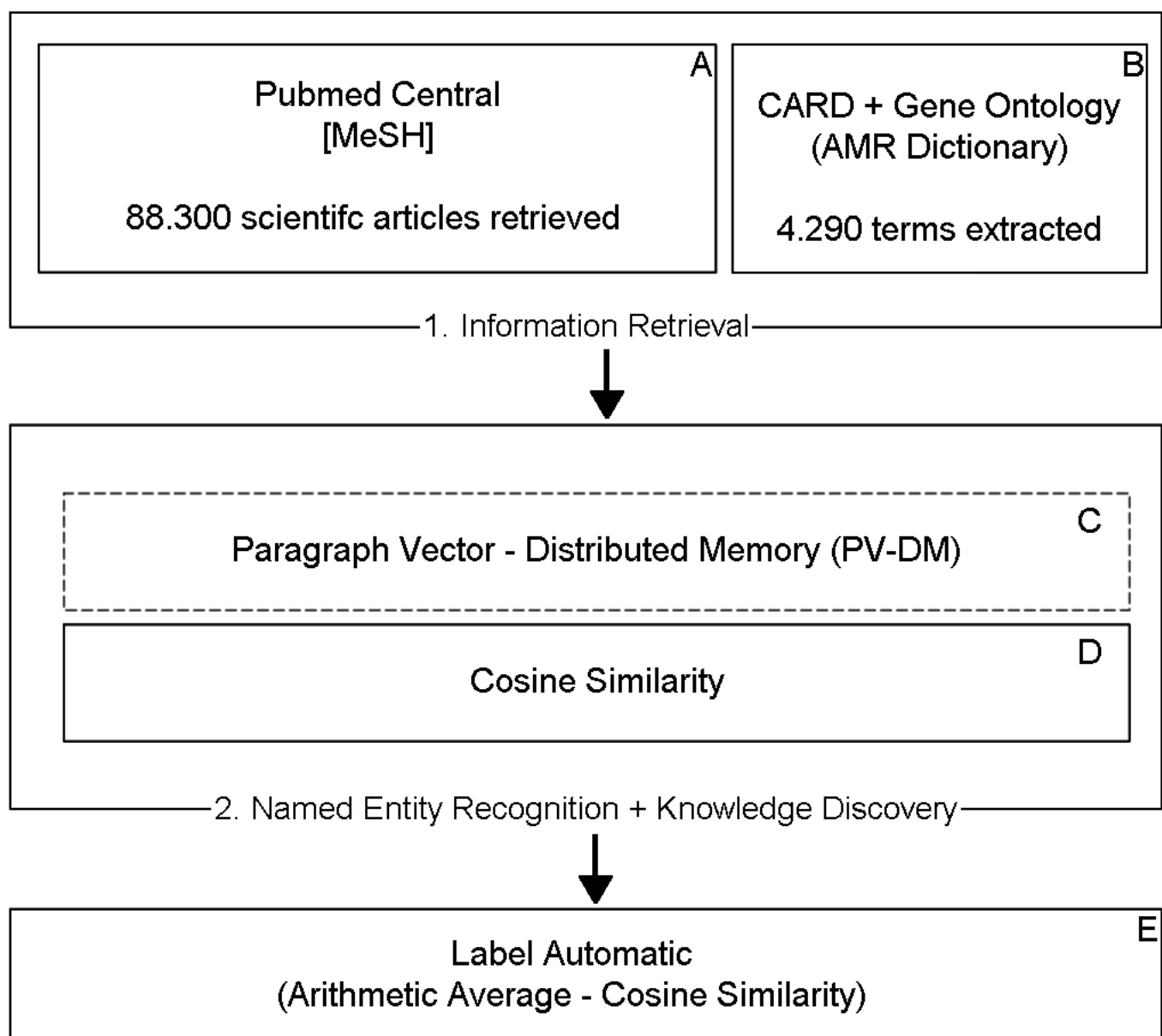


Figure 2

Evaluation of the proposed method

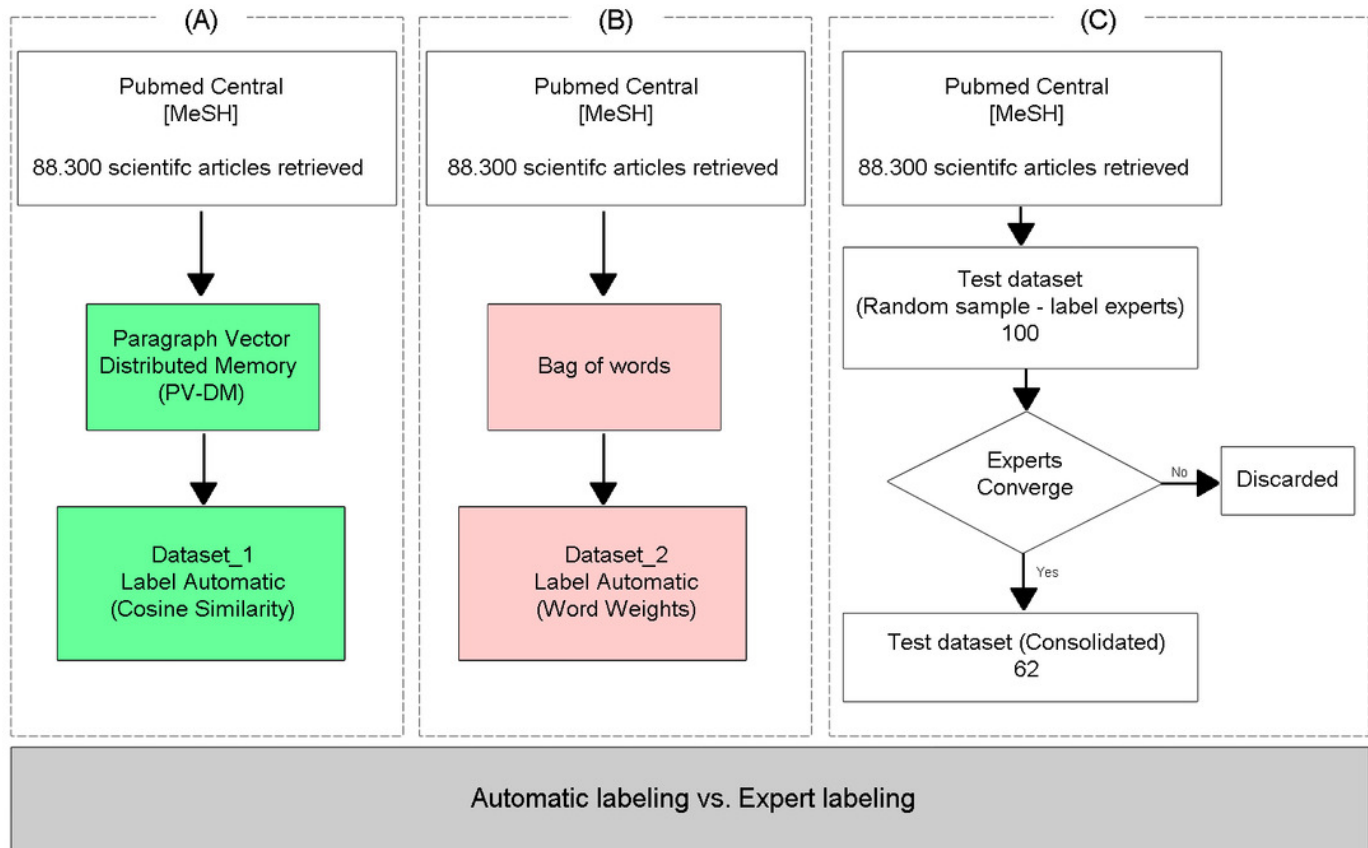


Figure 3

Performance of predictions with automatically labeled data



Figure 4

SVM classifier confusion matrix for dataset_1 (PV-DM)

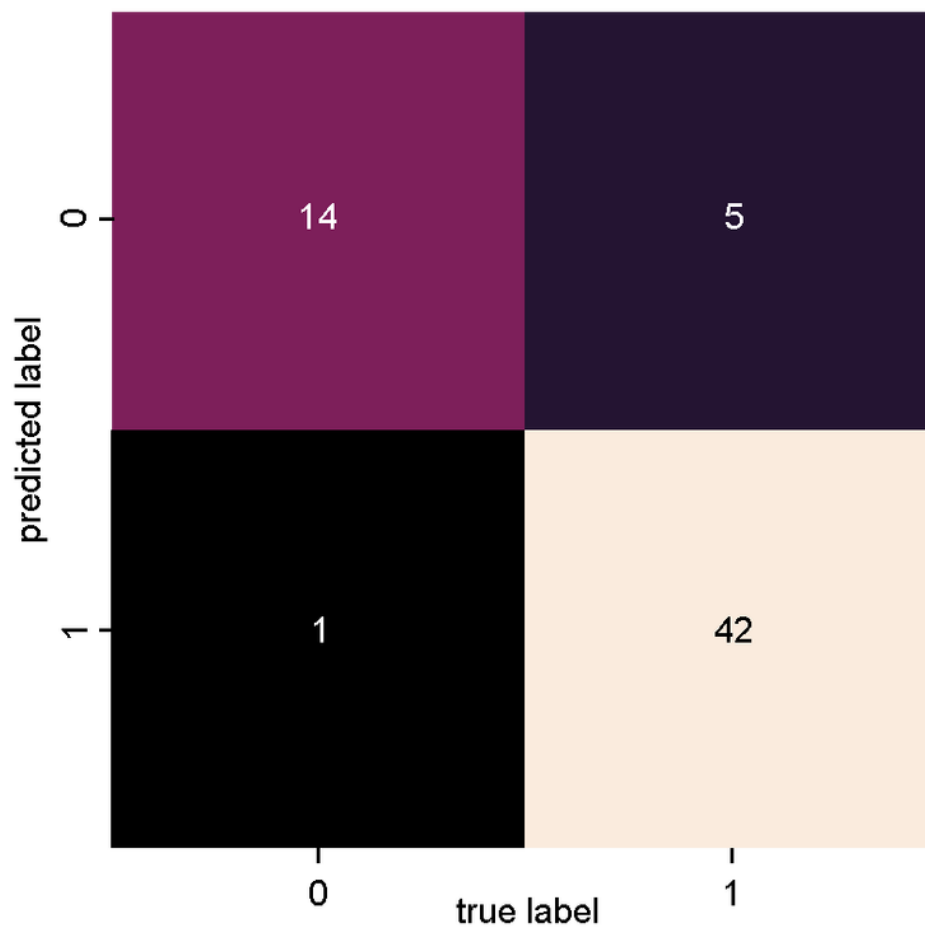


Figure 5

SVM classifier confusion matrix for dataset_1 (Bag of Words)

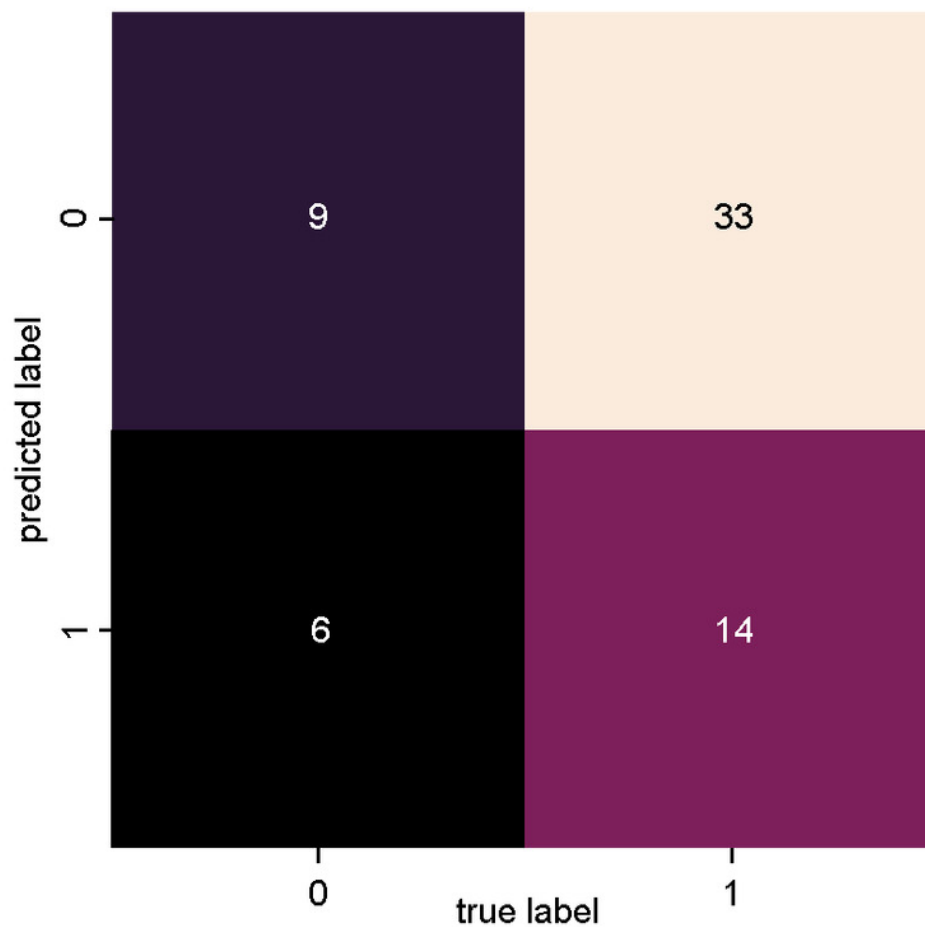


Table 1(on next page)

Parameters *E-Search* PubMed Central

Table 1. Parameters *E-Search* PubMed Central

Parameters	Value
URL	<a ""[filter]"="" "&term="" and="" href="https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=">https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=""&term="" AND ""[filter]
db	PMC (full text articles)
Term	("drug resistance, microbial"[MeSH Terms] OR ("drug"[All Fields] AND "resistance"[All Fields] AND "microbial"[All Fields]) OR "microbial drug resistance"[All Fields] OR ("drug"[All Fields] AND "resistance"[All Fields] AND "microbial"[All Fields]) OR "drug resistance, microbial"[All Fields])
Free text articles	Open access

Table 2 (on next page)

Parameters Doc2Vec algorithm

1 **Table 2.** Parameters Doc2Vec algorithm

Parameters	Value	Description
vector_size (<i>int, optional</i>)	300	Dimensionality of the feature vectors.
alpha (<i>float, optional</i>)	0,025	The initial learning rate.
min_alpha (<i>float, optional</i>)	0,00025	Learning rate will linearly drop to <i>min_alpha</i> as training progresses.
workers (<i>int, optional</i>)	18	Use these many worker threads to train the model (=faster training with multicore machines).
min_count (<i>int, optional</i>)	3	Ignores all words with total frequency lower than this.
epochs (<i>int, optional</i>)	30	Number of iterations (epochs) over the corpus.
dm (<i>{1,0}, optional</i>)	1	Defines the training algorithm. If <i>dm=1</i> , ‘distributed memory’ (PV-DM) is used. Otherwise, <i>distributed bag of words</i> (PV-DBOW) is employed.

2

Table 3(on next page)

Classifier Performance Assessment

1 **Table 3.** Performance evaluation of the classifiers SVM

	Class	Precision	Recall	F1-score	Support	Accuracy
SVM_1 (Doc2Vec+Mean)	0	0.74	0.93	0.82	15	0.90
	1	0.98	0.89	0.93	47	
SVM_2 (TF-IDF+Mean)	0	0.21	0.60	0.32	15	0.37
	1	0.70	0.30	0.42	47	

2

3

Table 4(on next page)

Results of Labeling and Classification steps vs Experts

1 **Table 4.** Labeling and Classification vs Experts

	Relevant	Irrelevant
Labeling		
Dataset_1	80%	68%
Dataset_2	66%	34%
Classification		
SVM_1	93%	89%
SVM_2	60%	29%

2