

From sequences to Papers: an Information Retrieval Exercise

Célia Talma Gonçalves
*LIACC &
Faculdade de Engenharia
Universidade do Porto &
Instituto Superior de Contabilidade
e Administração do Porto
& CEISE-STI
Porto, Portugal
talma@fe.up.pt*

Rui Camacho
*LIAAD & DEI &
Faculdade de Engenharia
Universidade do Porto
Porto, Portugal
rcamacho@fe.up.pt*

Eugénio Oliveira
*LIACC & DEI &
Faculdade de Engenharia
Universidade do Porto
Porto, Portugal
eco@fe.up.pt*

Abstract—Whenever new sequences of DNA or proteins have been decoded it is almost compulsory to look at similar sequences and papers describing those sequences in order to both collect relevant information concerning the function and activity of the new sequences and/or know what is known already about similar sequences that might be useful in the explanation of the function or activity of the newly discovered ones.

In current web sites and data bases of sequences there are, usually, a set of paper references linked to each sequence. Those links are very useful because the papers describe useful information concerning the sequences. They are, therefore, a good starting point to look for relevant information related to a set of sequences. One way is to implement such approach is to do a blast with the new decoded sequences, and collect similar sequences. Then one looks at the papers linked with the similar sequences. Most often the number of retrieved papers is small and one has to search large data bases for relevant papers.

In this paper we propose a process of generating a classifier based on the initially set of relevant papers that are directly linked to the similar sequences retrieved and use that classifier to automatically enlarge the set of relevant papers by searching the MEDLINE using the automatically constructed classifier.

We have empirically evaluated our proposal and report very promising results.

Keywords-MEDLINE; classification; information retrieval system;

I. INTRODUCTION

Molecular Biology and Biomedicine scientific publications are available (at least the abstracts) in Medical Literature Analysis and Retrieval System On-line (MEDLINE). MEDLINE is the U.S. National Library of Medicine (NLM), premier bibliographic database: contains over 16 million references to journal articles in life sciences with a concentration on Biomedicine. A distinctive feature of MEDLINE is that the records are indexed with NLM's Medical Subject Headings (MeSH terms). MEDLINE is the major component of PubMed [1], a database of citations of the NLM. PubMed comprises more than 19 million citations for biomedical articles from MEDLINE and life science journals. The

PubMed database maintained by the National Center for Biotechnology Information (NCBI) is a key resource for biomedical science, and is our first base of work. The NCBI's PubMed system is a widely used method for accessing MEDLINE.

The result of a MEDLINE/PubMed search is a list of citations (including authors, title, journal name, paper abstract, keywords and MeSH terms) to journal articles. The results of such search is, quite often, a huge amount of documents, making it very hard for researchers to efficiently reach the most relevant documents. As this is a very relevant and actual topic of investigation we assess the use of Machine Learning-based text classification techniques to help in the identification of a reasonable amount of relevant documents in MEDLINE. The core of the reported work is to study the best way to construct the data sets and the classifiers from the starting set of sequences.

These experiences were done using a set of positive examples associated to the sequences/keywords given by the user and a set of negative examples which is the focuses of this paper. The negative examples were generated in three different ways and we intend to show which is the best approach for our classification purpose. In our experiments we have used several classification algorithms available in the WEKA [2] tool.

The rest of the paper is structured as follows. In Section II we present an architecture for an information retrieval system. Section III presents the local data base construction process. Section IV presents the related work. Section V details the construction of the data sets and and thereby Section VI describes the experiences and presents the results obtained. Finally Section VII concludes the paper.

II. AN ARCHITECTURE FOR AN INFORMATION RETRIEVAL SYSTEM

The overall goal of our work is to implement a web-based search tool that receives a set of genomic or proteomic sequences and returns an ordered set of papers relevant to

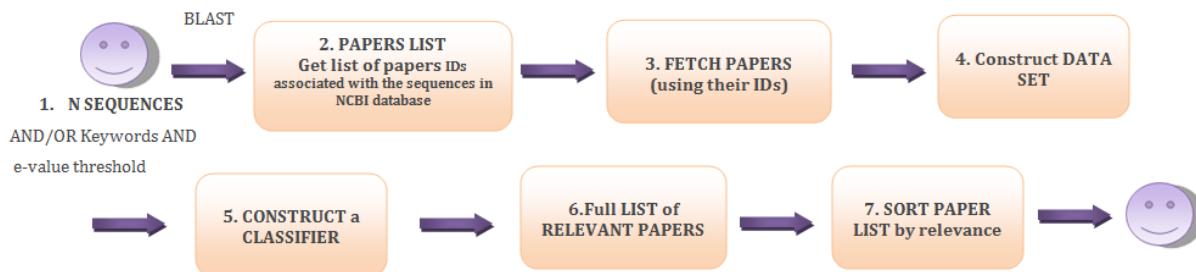


Figure 1. Sequence of steps executed by BioTextRetriever where a user provide a set of initial DNA/protein sequences.

the study of such sequences. The initial set of sequences is supplied by a biologist together with a set of relevant keywords and an e-value¹. These three items are the input for BioTextRetriever as can be seen in Figure 1. Figure 1 presents a summary of our approach that we will now describe in detail. In the following description we use NCBI as the sequence Data Base.

In Step 1, the user (a biologist researcher) provides an initial of sequences, optionally a list of keywords, and an e-value. With these three items (sequences, keywords and e-value) and using the NCBI BLAST tool we collect a set of similar sequences together with the paper references associated to them. We could also use Ensembl with the same inputs because Ensembl may return a different set of papers references. However for the proposed work we have only used the NCBI database.

With this list of paper citations we search for their abstracts in a local copy of MEDLINE (LDB - Local Database) (Step 3). For this we have previously preprocessed MEDLINE. Step 3 searches and collects the following information in the pre-processed local copy of MEDLINE: pmid, journal title, journal ISSN, article title, abstract, list of authors, list of keywords, list of MeSH terms and publication date.

For the scope of this paper we are considering only the paper citations that have the abstract available in MEDLINE. After Step 3 we have a data set of papers related to the sequences. We will take this set of papers as the positive examples for the full construction of the data set (Step 4) but we need to get some negative examples. To obtain the negative examples we have three possible approaches.

Thus this step is explained in detail in Section III.

The following step, Step 5, is one of the most important steps of our work which is to Construct a Classifier using Machine Learning techniques that is explained in the next section. As a result of this step we have a full list of articles considered relevant by our classifier (Step 6). However, we need to present them to the biologist in an ordered fashion way. So Step 7 presents an ordered list of relevant articles to

the biologist. Here we will develop and implement a ranking algorithm based on features such as the number of citations of the paper and the impact factor of the journal/conference where it was published.

This paper focuses on the construction of the data sets highlighting the research from the different approaches to obtain the negative examples. According to the figure and for the purpose of this paper we focus on Step 4, although we have made a set of experiences with some classifiers to conclude what was the best approach.

III. THE LOCAL DATA BASE

We have downloaded 80GB (617 XML files) of MEDLINE 2010 from the NCBI website. Each XML file has information characterizing one citation. Among these characteristics we have considered the following ones: PMID - the PubMed Identifier; the PubMed Date; the Journal Title; the Journal ISSN that corresponds to the ISI Web of Knowledge ISSN; the Title; the Abstract of the article if available; the list of the Authors; the MeSH Headings list and the Keywords list. After download the files were pre-processed as follows.

A. Pre-Processing MEDLINE XML files

An independent step of our tool is to maintain a local copy of MEDLINE, that we will call Local Data Base (LDB). The LDB will enable efficient search of the paper and will have that relevant information of each paper in format adequate, the algorithm that construct a chain as describe further in this paper.

The first step is to read the XML files and extract the relevant information to store in the LDB. Article's title and abstract are preprocessed with "traditional" text pre-processing techniques. Next section presents the pre-processing techniques applied.

B. Pre-processing Techniques

We have empirically [3] evaluate which are the best combination of pre-processing techniques to achieve a better accuracy. Based on this previous study and with some more research in the meanwhile we have used the following pre-processing techniques.

¹A e-value is a statistic to estimate the significance of a match between 2 sequences

Document Representation

For each paper with the information referred in the beginning of this section

However the text facts of a document (title and abstract) are filtered using text processing techniques and represented using the vector space model from Information Retrieval where the value of a term in a document is given by the standard term-frequency inverse document frequency (TFIDF=TF*IDF) function [4], to assign weights to each term in the document.

TF is the frequency of term in document
and

$$IDF = \log \frac{\text{number of documents in collection}}{\text{number of documents with term}} + 1$$

Named Entity Recognition (NER)

NER is the task of identifying terms that mention a known entity. We have used ABNER [5], which stands for A Biomedical Named Entity Recognition, that is a software tool for molecular biology that identifies entities in the biology domain : proteins, RNA, DNA, cell type and cell line. Although we have implemented these technique we have concluded that the identification of NER terms augments significantly the number of attributes instead of reducing them. We concluded that the use of NER increases strongly the number of terms which is a problem for the classifiers. Thus we did not use NER in the pre-processing phase.

Handling Synonyms

We handle synonyms using the WordNet [6] to search for similar terms, in the case of regular terms, and used Gene Ontology [7] to find biological synonyms. If two words mean the same then they are synonyms, so they could be replaced by one of them in the entire MEDLINE (title and abstract fields) without changing the semantic meaning of the term thus reducing the number of attributes. In this step we have replaced all the synonyms found by one synonym-term thus reducing the number of terms.

Dictionary Validation

A term is considered a valid term if it appears in available dictionaries. We have gathered several dictionaries for the common English terms (such as Ispell and WordNet) and for the medical and biological terms (BioLexicon [8], The Hosford Medical Terms Dictionary [9] and Gene Ontology [7]). The Hosford Medical Terms Dictionary consists of a file that contains a long list of medical terms. BioLexicon is a large-scale terminological resource developed to address text mining requirements in the biomedical domain. The BioLexicon is publicly available both as an XML-formatted term repository and as a relational database (MySQL) and it adheres to the LMF ISO standards for lexical resources. We have also used the Gene Ontology available files that are related to genes, enzymes, chemical resources, species and proteins. We have processed each of these resource files in order to have a simple text file with one term per line.

Our approach is in the sense that if a term appears in one of these dictionaries it is a valid term, otherwise, it is not a valid term, so we remove it from the collection of terms.

The application of these technique is fundamental in attribute reduction once a lot of terms that have no biology, medical and normal significance are discarded.

Stop Words Removal

Stop Words Removal removes words that are meaningless such as articles, conjunction and prepositions (e.g., a, the, at, etc.). These words are meaningless for the evaluation of the document content. We have used a set of 659 stop words file.

Tokenization

Tokenization is the process of breaking a text into tokens. A token is a non empty sequence of characters, excluding spaces and punctuation.

Special Characters Removal

Special character removal removes all the special characters (+, -, !, ?, ., ,, ;, :, {, }, =, &, #, %, \$, [,], /, <, >, \, “, ”, |) and digits.

Stemming

Stemming is the process of removing inflectional affixes of words reducing the words to their stem (the words computer, computing and computation are all transformed into comput, which means that three different terms are transformed into only one term thus reducing the number of attributes. We implemented the Porter's Stemmer Algorithm [10].

Pruning

Using pruning we discard in the documents collection terms that either appear too rarely or too frequently.

IV. RELATED WORK

There are some work being done on biological and biomedical document classification. Some of them applied to MEDLINE document document classification and other databases.

The work of [11] tries to automate the process of adding new information to TCDB database (Transport Classification Database) that is a web free access database (<http://www.tcdb.org>) about comprehensive information on transport proteins. The authors restricted themselves to the documents in MEDLINE. The main goal is to highlight the use of Machine Learning techniques outperforms rules created by hand by a human expert. To train the classifier they have used a set of MEDLINE documents referred TCDB as positive examples and have selected randomly also from MEDLINE a set of negative examples.

The authors in [12] describe a new model for text classification using estimating term weights which improves accuracy classification according to the authors experiences. Documents are represented as vectors of terms with their normalized global frequency. Global weights are functions that count how many times a term appears in the entire

collection and the normalization process compensates the discrepancies in the lengths of the documents. They have used 1000 documents from PubMed; 600 documents for the training data set and 400 for the test data set. All these documents belong to four categories with MeSH terms related to Diabetes melitus. The authors compare the different weighting methods: local-binary, local-log, local df and global relevant. They concluded in this study that global relevant weighting method achieves a higher precision. In our own work we have also used all normalized global frequency.

BioQSpace [13] is a GUI where users can query abstracts from PubMed using an embedded search facility. BioQSpace performs pairwise similarity calculations between all the abstracts based on a set of individual attributes namely: structure, function, disease and therapeutic compounds word list obtained from MeSH terms, word usage, PubMed related articles, publication date among others. These attributes are given more or less importance according to the weight attributed by users. A clustering algorithm is used to group abstracts that are very similar.

[14] describe a methodology to build an application capable of identifying and disseminating health care information using a Machine Learning approach. The main objective of their work is to study the best information representation model and what classification algorithms are suitable for classifying relevant medical information in short texts. They have used 6 different Machine Learning algorithms. The authors concluded that naive Bayes performed very well on short texts in the medical domain and that adaboost had the worst result.

In [15] the authors present an approach for classifying a collection of biomedical abstracts downloaded from MEDLINE database with the help of ontology alignment. Although this work classifies MEDLINE documents it is based on ontology alignment which is out of our scope.

LigerCat [16] stands for Literature and Genomic Electronic Resource Catalogue, and it is a system for exploring biomedical literature through the selection of terms within a MeSH cloud that is generated based on an initial query using journal, article, or gene data. The central idea of LigerCat is to create a tag cloud showing an overview of important concepts and trends associated to the MeSH descriptors. LigerCat aggregates multiple articles in PubMed, combining the associated MeSH descriptors into a cloud, weighted by frequency. LigerCat does not apply any Machine Learning techniques for paper classification as we present in our study.

V. AUTOMATIC CONSTRUCTION OF DATA SETS

A. Constructing the Data Sets

In order to solve our problem that is given a set of genomic or proteomic sequences return a set of related sequences and papers with relevant information for the study of such sequences, we need to obtain first of all the articles

associated with the given set of sequences and construct the data set to give to the classifier. Figure 2 illustrates sequence of steps included in BioTextRetriver. The empirical work repeated in the paper concerns the construction of a data set (Step 4).

The input of our work is a set of sequences given in the FASTA format. We use the netblast-2.2.22 tool that perform a remote blast search at the NCBI site. We have embedded this application into our code and automatically have access to both, the original sequences and the set of similar sequences retrieved by BLAST.

These results show us the similar sequences and the e-value associated with each of the retrieved sequence. The e-value is a statistic to estimate the significance of a match between two sequences. The e-value is an input that is given us by the biologist. We relax this threshold value in order to obtain the negative examples based on the e-value as we can see in Figure 3. The positive examples are the one's that are lower than the e-value previously specified by the biologist. We establish a "no man's land" zone and after that zone we collect the negative examples.

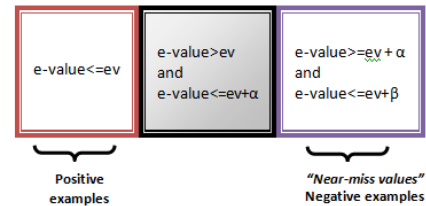


Figure 3. How positive and near-miss (negative) examples are obtained. ev is the e-value threshold to obtain the positive examples. α and β are parameters for the cut off of the negative examples.

The positive examples are the set of papers associated with the set of sequences with e-value below the respective threshold. In this study we have empirically evaluated three different ways of obtaining the negatives examples. We now explain the alternatives.

Near-Miss Values (NMV)

To obtain the Near-Miss Values (NMV) we collect the papers associated with the similar sequences that have e-value above the threshold but close to that. In Figure 3 there is a strip gray to better discriminate what are positive examples and negative examples. The examples in the right most box contains near-miss negative examples because they are not positives but have a certain degree of similarity with the sequences. This works on the examples that have a minimum number of negative examples. If we do not have any negative examples with this approach, or if the negative examples are few, we can follow one of the following approaches: to use MeSH Random Values or to use Random Values. In our experiments we have considered e-value = 0.001 and we have relaxed it to 1, 2 and 5 as we can see in sequences distribution tables in the next Section. We have

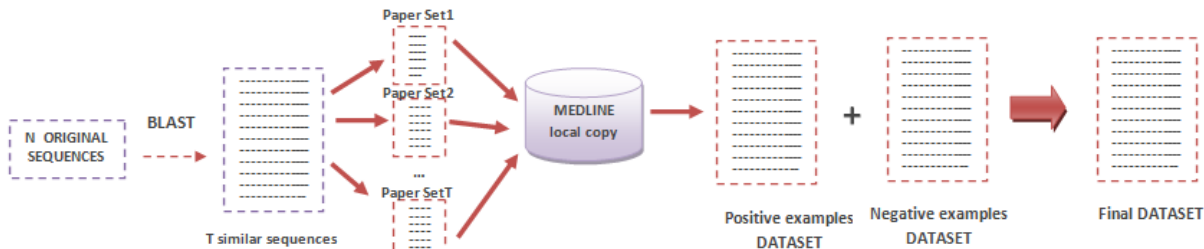


Figure 2. Dataset construction

relaxed to these different values to obtain more negative examples. The articles associated to the similar sequences with e-value less than 0.001 are considered positive; the articles associated to the similar sequences that have e-values greater than 0.001 and e-values less than 0.001 plus 10% ($\alpha = 10\%$) are considered in the gray strip so they are not considered positive or negative; the articles associated with to the similar sequences with e-values greater than 0.001 plus 10% are considered negative examples (near-miss values).

MeSH Random Values (MRV)

This alternative to generate negatives is adopted when we do not have sufficient number of negative examples for the classifier to learn. The negative examples are obtained combining the near miss values, if they exist, with some random examples generated from the LDB. But, these MRV examples must have the maximum number of MeSH terms from the positive examples. At the end the number of negative examples is equal to the number of positive examples.

Random Values (RV)

The last approach is to generate just randomly the negative examples from our LDB in a number equal to the number of positive examples. We guarantee that in this set there is no positive example.

VI. COMPARING THE ALTERNATIVES TO DATA SET CONSTRUCTION

A. Data Set Characterization

For, this study we have generated several data sets based on sequences that belong to six different classes, with the following distribution:

- RNASES: 2 sequences
- Escherichia Coli: 5 sequences
- Cholesterol: 5 sequences
- Hemoglobin: 5 sequences
- Blood Pressure: 5 sequences
- Alzheimer: 5 sequences

We have also used three different relaxation values for the e-value (1, 2 and 5). If the user enters an e-value of 0.001, then the positive examples are the ones that have e-value less or equal to 0.001. And the negative examples are the one's greater than 0.001. But as we can see in Figure 3 we leave

a gray strip to better separate the positive from the negative examples. This strip is also defined by the user. For these examples we have defined a strip of 10% of the number of not similar sequences. So the negative near-miss examples are the one's that are greater than 0.001 plus 10% of the number of not similar sequences and lower than e-value relaxation value (1, 2 or 5 in our examples).

The main idea of this study is to study the best way to construct the negative examples based on our experiences.

The distributions of positive and negative examples are show in the Appendix.

B. Experimental Results

In our experiences we have used a set of algorithms available in the WEKA [2] tools and that are listed Table I.

Table I
MACHINE LEARNING ALGORITHMS USED IN THE STUDY.

Acronym	Algorithm	Type
ZeroR	Majority predictor	Rule learner
smo	Sequential Minimal Optimization	Support Vector Machines
rf	Random Forest	Ensemble
ibk	K-nearest neighbors	Instance-based learner
BayesNet	Bayesian Network	Bayes learner
j48	Decision tree C4.5)	Decision Tree learner
dtmb	Decision table/naive bayes hybrid	Rule learner

The data sets used are characterised in Tables II and III. Table II characterizes data sets for which the negative examples are made only of near miss examples. Table III characterizes the data sets for which there were not enough negative examples and therefore we have used the MRV and RV strategies.

Tables IV, V and VI show the accuracy results obtained using the classifiers of Table I. Accuracy results were obtained performing a 10-fold Cross Validation.

The results tables show very promising results. Almost all values are weigh above the naive classifier of predicting the majority class. We can also say that the use of near miss values outperforms in most of the common data sets the other two strategies for generating negative examples. This finding is in the line of the use of near miss examples in Machine Learning.

Table II
CHARACTERIZATION OF DATA SETS WHERE THERE WERE ENOUGH
NEAR MISS EXAMPLES FOR LEARNING.

Data Sets	NA	Positive E.	Negative E.	Total E.
T11	539	18	10	28
EC15	276	53	53	106
EC45	313	110	82	192
EC55	162	57	4	61
BP12	470	164	81	245
BP25	1135	63	81	231
C35	583	20	39	59
C45	583	20	4	24

Table III
CHARACTERIZATION OF DATA SETS FOR WHICH THERE WERE NO, OR
NOT ENOUGH NEAR MISS EXAMPLES.

Data Sets	NA	Positive E.	Negative E.	Total E.
S12	1602	128	128	156
ALZ11	544	24	24	48
ALZ31	1485	114	114	228
C15	423	13	13	26
C21	354	8	8	16
C35	583	20	20	40
C45	583	20	20	40
C55	583	20	20	40
H11	1461	120	120	240
H21	396	12	12	24
H31	1408	130	130	260
H41	396	12	12	24
H51	1363	124	124	248

Table IV
NMV ACCURACY CLASSIFICATION RESULTS

Data Set	ZeroR	smo	rf	ibk	BayesNet	j48	dtnb
T11NMV	53.84	98.035 (5.3)	96.51 (8.49)	97.39 (6.65)	75.35 (11.7)	86.71 (13.3)	73.75 (12.35)
EC15NMV	46.72	95.63 (7.60)	95.00 (7.81)	96.60 (6.57)	77.10 (10.83)	80.87 (11.58)	82.31 (10.76)
EC45NMV	56.99	64.4 (9.47)	62.76 (9.57)	61.99 (9.63)	64.96 (7.91)	62.76 (9.25)	65.50 (7.08)
EC55NMV	93.44	89.26 (9.3)	90.73 (9.6)	89.26 (9.3)	89.26 (9.3)	89.26 (9.3)	91.73 (8.5)
BP12NMV	66.93	95.34 (4.65)	95.14 (4.33)	94.32 (4.73)	84.33 (6.81)	88.30 (6.55)	86.47 (6.82)
BP25NMV	64.93	91.76 (5.46)	92.67 (4.96)	91.42 (6.10)	80.30 (8.84)	85.75 (6.62)	83.38 (6.95)
C35NMV	57.60	92.05 (9.83)	92.55 (9.55)	93.22 (8.48)	80.17 (11.70)	83.27 (10.62)	81.28 (11.72)
C45NMV	92.98	93.33 (8.20)	90.73 (10.63)	93.33 (8.20)	93.33 (8.20)	93.33 (8.20)	93.33 (8.20)

VII. CONCLUSIONS

This paper focuses on data set construction for posteriori classification of MEDLINE documents. Our study highlights the impact of three different ways to construct the data sets for posteriori classification. These three different ways are applied only to construct the negative examples. The first one is based on the concept of near miss values (NMV), which are examples that although are negative examples are relatively close to the positive examples. The second approach is to use MeSH Random values (MRV) that is applied when we do not have enough negative near-miss

Table V
MRV ACCURACY CLASSIFICATION RESULTS

Data Set	smo	rf	ibk	BayesNet	j48	dtnb
T11MRV	70.00 (24.28)	78.33 (16.76)	70.00 (24.28)	73.33 (23.83)	65.00 (25.40)	70.83 (21.96)
S12MRV	88.37 (8.70)	91.82 (3.83)	54.78 (11.73)	98.05 (2.79)	98.05 (2.06)	97.29 (2.60)
Alz11MRV	60.50 (23.86)	80.00 (23.09)	52.50 (30.84)	83.50 (20.82)	94.00 (9.66)	93.50 (14.15)
Alz31MRV	78.12 (9.29)	89.98 (7.67)	53.97 (8.34)	98.26 (2.25)	94.76 (7.04)	97.83 (3.07)
C11MRV	43.33 (34.43)	58.33 (29.66)	50.00 (26.06)	33.33 (19.25)	76.67 (26.29)	61.67 (15.81)
C21MRV	65.00 (41.16)	50.00 (40.82)	55.00 (36.89)	85.00 (33.75)	75.00 (42.49)	85.00 (33.75)
C35MRV	45.00 (25.82)	60.00 (21.08)	50.00 (26.35)	70.00 (19.72)	70.00 (25.82)	67.50 (12.08)
C55MRV	52.50 (18.45)	65.00 (24.15)	52.50 (27.51)	80.00 (19.72)	82.50 (16.87)	75.00 (20.41)
H11MRV	89.17 (6.27)	95.83 (5.20)	64.58 (17.60)	96.25 (4.14)	94.58 (5.22)	98.75 (2.01)
H21MRV	38.33 (28.38)	66.67 (24.85)	50.00 (38.49)	81.67 (24.15)	78.33 (23.64)	81.67 (24.15)
H35MRV	82.69 (9.29)	88.08 (7.57)	67.69 (13.10)	91.15 (6.29)	91.54 (6.98)	95.00 (3.65)
H41MRV	48.33 (25.40)	55.00 (30.48)	50.00 (26.06)	38.33 (15.81)	75.00 (27.50)	80.00 (26.99)
H51MRV	93.92 (6.53)	95.55 (2.33)	64.12 (14.34)	96.35 (4.10)	96.78 (4.14)	98.78 (2.72)

Table VI
RV ACCURACY CLASSIFICATION RESULTS

Data Set	smo	rf	ibk	BayesNet	j48	dtnb
T11RV	84.17 (24.67)	94.17 (12.45)	53.33 (22.97)	94.17 (12.45)	94.17 (12.45)	97.50 (7.91)
S12RV	86.09 (10.30)	92.66 (4.24)	51.92 (7.69)	96.91 (3.98)	97.68 (3.72)	98.05 (3.33)
C11RV	76.67 (21.08)	68.33 (24.15)	50.00 (26.06)	90.00 (22.50)	90.00 (22.50)	93.33 (21.08)
C21RV	85.00 (24.15)	80.00 (34.96)	55.00 (36.89)	85.00 (33.75)	85.00 (33.75)	100.00 (0.00)
C35RV	67.50 (16.87)	87.50 (17.68)	50.00 (26.35)	97.50 (7.91)	100.00 (0.00)	97.50 (7.91)
C45RV	72.50 (21.89)	87.50 (13.18)	50.00 (26.35)	100.00 (0.00)	92.50 (16.87)	97.50 (7.91)
C55RV	60.00 (21.08)	82.50 (16.87)	50.00 (26.35)	97.50 (7.91)	95.00 (10.54)	97.50 (7.91)
H11RV	88.75 (6.53)	94.17 (5.27)	66.25 (17.51)	95.42 (4.99)	96.67 (5.12)	99.17 (1.76)
H21RV	75.00 (27.50)	91.67 (18.00)	60.00 (32.58)	88.33 (19.33)	91.67 (18.00)	85.00 (19.95)
H31RV	90.30 (7.48)	96.24 (3.06)	69.56 (11.56)	95.11 (4.77)	96.61 (2.79)	99.26 (2.34)
ALZ11RV	80.00 (26.67)	88.00 (21.50)	51.00 (30.62)	88.00 (21.50)	94.00 (13.50)	95.50 (9.56)
BP25RV	82.63 (8.73)	87.37 (6.42)	71.35 (11.62)	95.96 (4.26)	92.82 (8.89)	100.00 (0.00)

examples (when negative examples are less than the number of positive examples) and add to this few or none examples some random negative examples from our LDB. However these examples are not just random negative examples, they must have some MeSH terms common with the positive examples MeSH terms that were previously processed. The last approach is to generate randomly the negative examples from our LDB in equal number to the number of positive examples. However in the second and third approach we guarantee that in these random negative examples there

aren't any positive examples.

We have generated several data sets with these different techniques. We have presented comparison tables, for these three techniques under study for six different categories: RNASES, Escherichia Coli, Blood Pressure, Alzheimer, Hemoglobin and Cholesterol. The categories were chosen by a biologist expert. These tables present the accuracy obtained performing a 10-fold cross validation and using different classification algorithms and the three different approaches (NMV, MRV and RV).

From the results presented in the accuracy tables of Section VI-B, we can say that the use of near miss examples achieves better accuracy. We argue that the classifier achieves better results with the use of near miss examples because they establish a close boundary around the positive ones and so the classifier learns better to discriminate between positive and negative examples. We can also conclude that when we have enough number of near-miss examples we achieve accuracies very near to the Random Values accuracies as long as the data set is relatively small. We can see this, for example, in Table IV in example T11 and in Table V in example EC15.

REFERENCES

- [1] D. L. Wheeler, T. Barrett, D. A. Benson, S. H. Bryant, K. Canese, V. Chetvernin, D. M. Church, M. Dicuccio, R. Edgar, S. Federhen, L. Y. Geer, W. Helmsberg, Y. Kapustin, D. L. Kenton, O. Khovayko, D. J. Lipman, T. L. Madden, D. R. Maglott, J. Ostell, K. D. Pruitt, G. D. Schuler, L. M. Schriml, E. Sequeira, S. T. Sherry, K. Sirotkin, A. Souvorov, G. Starchenko, T. O. Suzek, R. Tatusov, T. A. Tatusova, L. Wagner, and E. Yaschenko. Database resources of the national center for biotechnology information. *Nucleic Acids Res*, 34(Database issue), January 2006.
- [2] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H. Witten. The weka data mining software: an update. *SIGKDD Explor. Newsl.*, 11(1):10–18, 2009.
- [3] Carlos Adriano Gonçalves, Célia Talma Gonçalves, Rui Camacho, and Eugénio C. Oliveira. The impact of pre-processing on the classification of medline documents. In Ana L. N. Fred, editor, *Pattern Recognition in Information Systems, Proceedings of the 10th International Workshop on Pattern Recognition in Information Systems, PRIS 2010, In conjunction with ICEIS 2010, Funchal, Madeira, Portugal, June 2010*, pages 53–61, 2010.
- [4] Wei Zhou, Neil R. Smalheiser, and Clement Yu. A tutorial on information retrieval: basic terms and concepts. *Journal of Biomedical Discovery and Collaboration*, 1:2+, March 2006.
- [5] Burr Settles. Abner: an open source tool for automatically tagging genes, proteins and other entity names in text. *Bioinformatics*, 21(14):3191–3192, 2005.
- [6] C. Fellbaum. *WordNet: An Electronical Lexical Database*. The MIT Press, Cambridge, MA, 1998.
- [7] M. Ashburner. Gene ontology: Tool for the unification of biology. *Nature Genetics*, 25:25–29, 2000.
- [8] D. Rebbholz-Schuhmann, P. Pezik, V. Lee J-J Kim, R. del Gratta, Y. Sasaki, J. McNaught, S. Montemagni, M. Monachini, N. Calzolari, and S. Ananiadou. Biolexicon: Towards a reference terminological resource in the biomedical domain. In *Proceedings of the of the 16th Annual International Conference on Intelligent Systems for Molecular Biology (ISMB-2008)*, 2008.
- [9] Hosford medical terms dictionary v3.0, 2004.
- [10] M. F. Porter. An algorithm for suffix stripping. pages 313–316, 1997.
- [11] Aditya Kumar Sehgal, Sanmay Das, Keith Noto, Milton H. Saier Jr., and Charles Elkan. Identifying relevant data for a biological database: Handcrafted rules versus machine learning. *IEEE/ACM Trans. Comput. Biology Bioinform.*, 8(3):851–857, 2011.
- [12] S.Sagar Imambi and T.Sudha. Classification of medline documents using global relevant weighing schema. *International Journal of Computer Applications*, 16(3):45–48, February 2011. Published by Foundation of Computer Science.
- [13] Anna Divoli, Rasmus Winter, Steve Pettifer, and Terri Attwood. 20. bioqspace: An interactive visualisation tool for clustering medline abstracts., 2005.
- [14] Oana Frunza, Diana Inkpen, and Thomas Tran. A machine learning approach for identifying disease-treatment relations in short texts. *IEEE Trans. Knowl. Data Eng.*, 23(6):801–814, 2011.
- [15] R. Dollah, Md. H. Seddiqui, and M. Aono. The effect of using hierarchical structure for classifying biomedical text abstracts. 2010.
- [16] Indra Neil N. Sarkar, Ryan Schenk, Holly Miller, and Catherine N. Norton. LigerCat: using "MeSH Clouds" from journal, article, or gene citations to facilitate the identification of relevant biomedical literature. *AMIA ... Annual Symposium proceedings / AMIA Symposium. AMIA Symposium*, 2009:563–567, 2009.

APPENDIX

In all the tables of this appendix SEQ means the type of sequence which can be one of six types (RNASE: starts with T or S, Escherichia Coli starts with EC, Hemoglobin starts with H, Cholesterol starts with C, Blood Pressure starts with BP and Alzheimer starts with Alz). S stands for "Similar" is the number of similar sequences; SP means "Similar Pmid" and represents the number of articles associated to the set of similar sequences; LS means "Less Similar" and is the total number of less similar sequences (e-value above the threshold); and LP stands for "Less Pmid" and represents the total number of articles associated with the less similar sequences.

Table VII
RNASE'S SEQUENCES DISTRIBUTION

SEQ	S	SP	LS	LP
T11	49	26	114	57
T12	49	26	146	77
T15	49	26	155	78
S11	250	172	0	0
S12	250	172	0	0
S15	250	172	0	0

Table VIII
ESCHERICHIA COLI SEQUENCES DISTRIBUTION

SEQ	S	SP	LS	LP
EC11	43	31	27	13
EC12	43	31	40	17
EC15	43	31	60	25
EC41	88	66	57	50
EC42	88	66	67	51
EC45	88	66	76	57
EC51	241	38	219	8
EC52	241	38	221	37
EC55	241	38	223	37

Table IX
BLOOD PRESSURE SEQUENCES DISTRIBUTION

SEQ	S	SP	LS	LP
BP11	69	94	104	56
BP12	69	94	166	106
BP15	69	94	183	116
BP21	71	94	71	40
BP22	71	94	112	58
BP25	71	94	180	117

Table X
ALZHEIMER SEQUENCES DISTRIBUTION

SEQ	S	SP	LS	LP
Alz11	33	39	0	0
Alz12	33	39	1	0
Alz15	33	39	5	1
Alz31	239	178	28	22
Alz32	239	178	36	27
Alz35	239	178	46	29

Table XI
CHOLESTEROL SEQUENCES DISTRIBUTION

SEQ	S	SP	LS	LP
C11	24	21	2	4
C12	24	21	2	4
C15	24	21	2	4
C21	24	21	3	4
C22	24	21	3	4
C25	24	21	3	4
C31	43	34	19	12
C32	43	34	24	17
C35	43	34	31	18
C41	36	34	7	2
C42	36	34	12	2
C45	36	34	33	16
C51	36	34	7	2
C52	36	34	12	2
C55	36	34	33	16

Table XII
CHOLESTEROL SEQUENCES DISTRIBUTION

SEQ	S	SP	LS	LP
H11	250	156	0	0
H12	250	156	0	0
H15	250	156	0	0
H21	250	194	0	0
H22	250	194	0	0
H25	250	194	0	0
H31	250	183	0	0
H32	250	183	0	0
H35	250	183	2	3
H41	250	194	0	0
H42	250	194	0	0
H45	250	194	0	0
H51	250	191	0	0
H52	250	191	0	0
H55	250	191	0	0