

# Conceptualization Effects on MEDLINE Documents Classification Using Rocchio Method

Shereen Albitar

Aix Marseille University LSIS UMR  
CNRS 7296

Domaine Universitaire de St Jérôme  
13397 Marseille Cedex 20, France  
shereen.albitar@lsis.org

Sebastien Fournier

Aix Marseille University LSIS UMR  
CNRS 7296

Domaine Universitaire de St Jérôme  
13397 Marseille Cedex 20, France  
sebastien.fournier@lsis.org

Bernard Espinasse

Aix Marseille University LSIS UMR  
CNRS 7296

Domaine Universitaire de St Jérôme  
13397 Marseille Cedex 20, France  
bernard.espinasse@lsis.org

**Abstract-** The aim of this paper is to propose a supervised text classification method for the biomedical domain using semantic resources. We choose the traditional text classification method, Rocchio, for its scalability and extendibility with semantic knowledge. This paper proposes to integrate semantic aspects into Rocchio through a conceptualization task. This conceptualization is realized by mapping terms that are extracted from text to their corresponding concepts in the UMLS® Metathesaurus® in order to take meaning into consideration during text classification. The proposed classifier is tested on the Ohsumed text corpus, which is composed of abstracts of biomedical articles retrieved from the MEDLINE® database. The effects of Conceptualization on Rocchio's performance are discussed according to different standard similarity measures and to a variety of conceptualization strategies.

**Keywords-Text Classification, Semantic classification, Information retrieval, Rocchio, Similarity measures, conceptualization.**

## I. INTRODUCTION

Nowadays and due to the explosive increase in published information, existing search engines seem to be unable to respond effectively to users' requests. This is often related to the traditional keyword-based indexing techniques that neglect search context [1]. Aiming at more effective and less time expensive search, it seems adequate to involve classification techniques in order to consider the contents of answers provided by search engines, applying thorough filtering and ranking.

Text classification is currently a challenging research topic, particularly in areas such as information retrieval, recommendation, personalization, user profiles etc.

Generally, text classification methods use syntactical and statistical models for text document representation. This applies to the most popular text classification methods such as: Naïve Bayes Classifier (NB), Support Vector Machines (SVMs), Rocchio, and K Nearest Neighbor (KNN). These models suffer the lack of sense in resulting representations, ignoring all semantics that reside in the original text that can help in text classification.

Vector-based (binary or TF/IDF) representation used by preceding methods permits semantic integration or

"Conceptualization" that enriches document representation model using a certain background knowledge base [2, 3]. In addition, both KNN and Rocchio enable involving knowledge bases in decision making through semantic similarity functions [4].

In this work, we mainly focus on improving text classification in the biomedical domain using domain specific knowledge bases, particularly UMLS®. Many works have already tried to improve biomedical text representation for better classification using UMLS® [5, 6] or MeSH® [7, 8]. To the best of our knowledge, few works investigated in details the gain of integrating these semantic resources in classifying biomedical text.

Making a comparative study on the different traditional classification methods is out of the scope of this paper, for detailed comparisons please refer to [9]. We consider Rocchio an adequate baseline text classifier for its scalability and its extendibility with semantic resources at both levels: text representation and similarity calculation for decision making. Most of other traditional classification methods, such as SVM and NB, allow the integration of semantics essentially in text representation. Nevertheless, deploying semantic similarity functions allows a full exploitation of semantic resources (concept properties, relations between concepts, etc.) during decision making.

Next section proposes appropriate semantic solutions in order to improve text representation and classification. In third section, we apply Rocchio method to conceptualized Ohsumed corpus [10] using UMLS® (Unified Medical Language System) Metathesaurus® and MetaMap® tool. The forth section presents some preliminary results using these conceptualized documents. Conceptualization effects on Rocchio's classification according to different standard similarity measures and conceptualization strategies are discussed. Finally, we conclude with an assessment of our work, followed by research perspectives.

## II. TOWARDS A ROCCHIO-BASED SEMANTIC CLASSIFICATION

In spite of being considered the most popular text representation method, VSM representation suffers from certain limitations [11, 12] especially for processing composed words, synonyms, polysemy, etc. In order to overcome these limitations, semantic resources (like thesaurus & ontologies) can be used to replace term-based

representation by concept-based one through "conceptualization". As an example of semantic resources that might be used for conceptualization: WordNet, Wikipedia and other domain specific resources usually called domain ontologies such as UMLS® thesaurus in the medical domain.

In general, vector conceptualization is realized in two steps: (i) search for corresponding concepts related to vector's terms and then (ii) the integration of these concepts in the vector producing the final conceptualized vector. Three different strategies have been proposed for conceptualization [2]: (i) *Adding Concepts*: Where the original vector is extended and corresponding concepts are added. (ii) *Partial Conceptualization*: Where terms are substituted by corresponding concepts. Terms having no related concepts are kept in the vector. (iii) *Complete Conceptualization*: Similarly to Partial Conceptualization, terms are substituted by concepts whereas remaining terms are eliminated from the final vector. The integrated concepts are assigned scores derived from related terms' frequencies.

While searching concepts that correspond to a polysemic term in semantic resources, multiple matches are detected. This introduces some ambiguities in document representation. For example: the term "Book" signifies in English a book and also a reservation (Ticket, accommodations...). According to [2], three strategies for disambiguation deploying WordNet can be used: (i) *All*: accept all candidate concepts as matches for the considered term. (ii) *First*: Accept the most frequently used concept among candidates using document language statistics. (iii) *Context*: Accept the candidates having the most similar semantic context compared to the term's context in text.

### III. APPLYING ROCCHIO TO CONCEPTUALIZED CORPUS

This section presents details concerning the integration of conceptualization process in the Rocchio classifier (figure 1). Two resources, a knowledge base and a text-to-concept mapping utility, are needed in order to transform a plain text into a conceptualized one. We chose to use the UMLS® Metathesaurus® and the MetaMap® tool for mapping Ohsumed text to UMLS® concepts.

After introducing these resources, we present the new classification process that integrates a conceptualization task using these resources. Finally we introduce the different conceptualization strategies used in our experimentations.

#### A. Resources

In this section we present UMLS® and MetaMap® that are both used in our system during the conceptualization process. Both were developed by the National Library of Medicine (NLM) aiming at facilitating the development of sophisticated medical information systems.

##### 1) Unified Medical Language System® (UMLS®)

The Unified Medical Language System® (UMLS®) [13] was developed to model the language of biomedicine and health. UMLS' sources enhance the development of information systems in the biomedical domain. The UMLS® knowledge base consists of three main resources: the Metathesaurus, the Semantic Network and the

SPECIALIST Lexicon. The Metathesaurus is a multilingual vocabulary database of biomedical concepts, their names, their attributes and the relations among them. This database organizes concepts of various source vocabularies (like MeSH®, SNOMED-CT®, etc.) according to their senses grouping common concepts together. Concepts and relations among them are assigned at least one type from the Semantic Network. Indeed, the Semantic Network provides a higher level of abstraction through concept and relation categorization in inter-related Types constituting a network of 133 semantic types and 54 relationships. The SPECIALIST lexicon contains a large variety of general terms as well as medical terms and words. Text conceptualization using UMLS® as a semantic resource, if compared to other more generic semantic resources such as WordNet, is more relevant in the medical domain. This allows the mapping between text and the related specific concepts in UMLS®.

##### 2) MetaMap® tool

In addition to the UMLS® knowledge resources, many tools are developed and provided by NLM in order to facilitate deploying these sources by medical information system developers. In this work we are interested particularly in MetaMap® [14]. The major goal of MetaMap's developers was to improve biomedical text retrieval using UMLS® sources. Indeed, MetaMap® can discover the links between biomedical text and the knowledge in the Metathesaurus. This mapping is the result of a rigorous linguistic analysis of each phrase of the text: First the text is tokenized and phrase boundaries are identified, then part-of-speech-tags are added. Second, the Specialist lexicon and the shallow parser are used to analyze these phrases syntactically. Finally, different candidates are identified in the Metathesaurus and then final mappings combining these candidates are evaluated resulting in confidence scores for each mapping. In cases where ambiguities are detected, MetaMap® keeps the mappings that are the most semantically similar to the surrounding text following the context strategy (section 2).

#### B. Classification Process

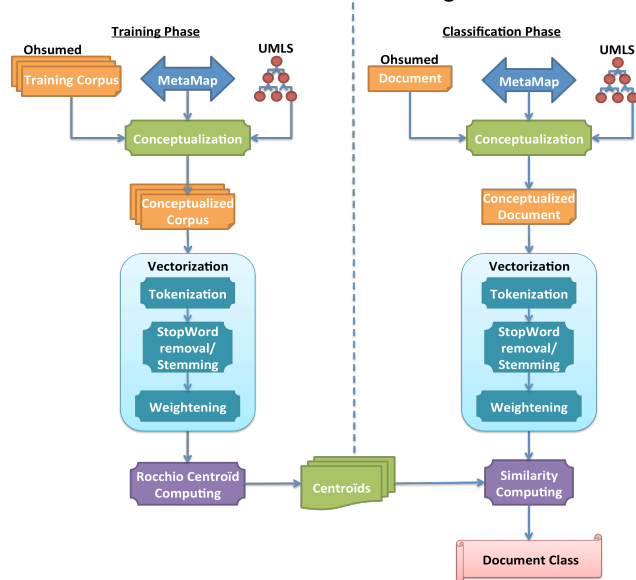
Figure 1 illustrates the classification process that is realized through two phases: Training Phase and Classification Phase. The *Conceptualization* task is introduced before the *Vectorization* task using both MetaMap® and UMLS®.

During the training phase, all documents of the training corpus are conceptualized using UMLS® by means of MetaMap®; the system transforms all training corpus into conceptualized text documents. This *Conceptualization task* is realized on text documents directly and not on their vectors in order to keep all information in text. This information is needed by MetaMap® for realizing a better text to concept mappings. Inspired by vector conceptualization strategies described in the previous section, text conceptualization can also be done using several strategies. Then, the *Vectorization task* converts these conceptualized documents into vectors according to the well known VSM. Finally the *Rocchio Centroid Computing task* computes classes' centroids using vectors resulting from conceptualized training corpus vectorization.

During the classification phase, the new document to classify is first conceptualized and then vectorized through preceding *Conceptualization* and *Vectorization* tasks. Then, through a *Similarity Computing* task, the resulting class is the one having the most similar centroid according to one of the five similarity measures (Cosine, Jaccard, Kullback-Leibler, Levenshtein and Pearson).

### C. Conceptualization step

During the conceptualization task, different strategies can be implemented as described in section 2 (adding concepts, partial conceptualization and complete conceptualization). According to MetaMap® text-to-concept matching results, two strategies can be chosen: *Best concept strategy* that takes the best concept among several candidate concepts matched to the text according to a matching score computed by MetaMap® [14], or *All concepts strategy* that keeps all candidates. Candidates resulting from matching have many properties. We chose to use the *concept name* or the *concept ID*. In fact, during the tokenization step, the *concept ID* is considered as a single token so it stays intact. *Concept names*, being sometimes compound words, are cut during tokenization when applied to a conceptualized text according to concept name strategy. In this work, conceptualization is done using 12 combinations of the previous strategies. Classification step is executed using the same five similarity measures that we used in the original system without conceptualization. Next section shows a selection of the most relevant results using F-measure.



**Figure 1. Conceptualization process: Rocchio applied to conceptualized corpus**

## IV. EXPERIMENTATION AND RESULTS

These experimentations were realized in order to evaluate the effect of integrating MetaMap® into the classification process; Rocchio is tested on the Ohsumed before and after conceptualization. This integration enables mapping text to UMLS® concepts.

### A. Results analysis

Table 1 shows the F1-measure obtained from applying Rocchio method to conceptualized Ohsumed corpus, for each of the five most frequent categories {C06, C23, C04, C14, C20}. These results are obtained first from using one of the five similarity measures that are also tested and compared on the original corpus without conceptualization. Then the classifier is tested with each similarity measure on the conceptualized corpus according to one of the twelve different conceptualization strategies. The two last columns present Micro and Macro averaged F-measure obtained for each similarity measure and conceptualized strategy. In Micro-averaging, F-measure is computed globally over all documents, whereas in Macro-averaging it is equal to the average of locally calculated F-measure for each class. As illustrated in the table, in most cases the original system outperforms the new system integrating conceptualization phase. In fact this applies to approximately 70% of results evaluated using the Micro-averaged F-measure. However, after a thorough look into the results, it seems clear that the system using the similarity measure of KullbackLeibler shows some amelioration. Indeed, results are improved in two thirds of the cases after conceptualization. Considering classes independently in the results, we can observe that conceptualization improves the outcome in about 70% of cases for the class "C23". The original system, showed the worst performance in treating this class. This difference is significant according to the McNemar [15] test on classification results of C23 documents before and after conceptualization ( $p \ll 0.01$ ).

Concerning *conceptualization strategies*, one of the 12 conceptualization strategies tested in these experimentations seems to provide the maximal improvement. Indeed, *Best Concept* strategy outperforms others as it retrieves the best candidate concept provided by MetaMap® for the mapping result. *Best concept names* rather than *IDs* are used in text conceptualization. However this strategy does not present a significant improvement over the original system without conceptualization. Indeed, the gain is each time less than 1%. The highest increases are obtained when using conceptualization strategies taking into account all the candidates found by MetaMap® and not only the best. Thus, in some cases this increase exceeds 10%. Furthermore, largest increase in the Micro-averaged F-Measure values are attained using the strategy of adding Concepts applied to the classifier with KullbackLeibler similarity measure. Previous improvements seem to be significant according to McNemar [15] test having ( $p \ll 0.01$ ).

Concerning *similarity measures*, the least improvement can be observed for the classifier using the similarity measure Levenshtein. In fact, among all Micro-averaged F-measure values for this classifier, only one surpasses its value when compared to system results using Levenshtein without conceptualization. As we noticed during experiments on the original corpus, Rocchio classifiers that use Cosine and Pearson have relatively similar behaviors. Indeed, both conceptualization strategies that add the names and the IDs of the best concepts seem to improve classification results for both similarity measures. Since the Micro-averaged F-measure's increases are in most cases less than 0.5%, this increase does not seem very significant.



concepts can participate in vectorization. Terms that are shared among concept with different IDs are excluded from vectors even if they had an important influence on results.

Second, when the system performance has a good F1-measure value (i.e. exceeds 60%), no significant effect can be observed for the integration of the conceptualization task into the system. In fact, as the same similarity measures are used for both cases with/without conceptualization, results' improvements were limited.

Third, when the system performance using a specific similarity measure has a low F1-measure value, as it the case for the class "C23", introducing conceptualization can significantly improve this value with a maximum gain reaching (10%) in some cases. Indeed, the class "C23" is very large compared to others and so enriching class representation by semantics might result in a better identification of this class and also in better results.

Fourth, the best conceptualization strategy is *Add concept* adding the *Best concept* among mapped candidates into the text. In fact, best mappings retrieved by MetaMap® are added into text in order to enrich text representation with semantics avoiding any information loss.

Finally, even if the results are still preliminary, it seems useful to introduce semantic enrichments to the Rocchio method in order to ameliorate classification results. Nevertheless, the exploitation of semantic resources was limited in this work ignoring all relations (like Subsumption and Transversal relations) among concepts identified during the conceptualization task. Thus, it seems adequate to deploy these relations during classification through the introduction of new semantic similarity measures which can be coupled with traditional similarity measures.

## V. CONCLUSION AND PERSPECTIVES

Due to the explosive growth of published data, many search engines demonstrate limited performance to meet the needs of users. This leads to a challenging need for effective filtering and ranking techniques. This paper concerns text classification in the biomedical domain, demonstrating the effects of involving semantic resources in text representation on classification effectiveness.

We choose Rocchio that demonstrates a good performance compared to its relatively minimal complexity in addition to its extendibility with semantics. Moreover, it can provide relevance feedback on classification results permitting better result understanding and potential improvements in classification. Moreover, some of its limitations could be overcome by means of semantic resources considering meaning in text classification. These extensions are promising.

Thus, in this research we have proposed to extend the text classification method, Rocchio, using semantic resources in order to improve its performance. Indeed, semantics can be integrated into Rocchio through conceptualization and also during decision making through different semantic similarity measures. Document conceptualization using knowledge bases helps to complete VSM approach with semantics. We realized some experiments using different conceptualization strategies and UMLS® on the Ohsumed corpus with standard similarity measures according to Rocchio's classification process.

These experiments show in some cases considerable performance improvements. However we expect better improvements, through deploying semantic similarity measures that can be calculated by aid of semantic resources. These measures can be combined with standard similarity measures, already used in this work, permitting the development of effective semantic text classification methods. Moreover, we intend to realize comparative studies on the effects of conceptualization on other traditional classification algorithms such as SVM and NB.

## VI. REFERENCES

- [1] A. P. Asirvatham and K. K. Ravi, *Web page classification based on document structure*, 2001.
- [2] A. Hotho, S. Staab, and G. Stumme, "Text clustering based on background knowledge," ed, 2003.
- [3] E. Ferretti, M. Errecalde, and P. Rosso, "Does Semantic Information Help in the Text Categorization Task?," *Journal of Intelligent Systems*, vol. 17, pp. 91-107, 2008.
- [4] A. Guisse, K. Khelif, and M. Collard, "PatClust: une plateforme pour la classification sémantique des brevets," in *IC 2009 Conference*, Hammamet, Tunisie, 2009.
- [5] V. N. Garla and C. Brandt, "Ontology-guided feature engineering for clinical text classification," *J Biomed Inform*, in press.
- [6] M. Yetisgen-Yildiz and W. Pratt, "The effect of feature representation on MEDLINE document classification," in *AMIA Annu Symp*, 2005, pp. 849-853.
- [7] X. Zhang, L. Jing, X. Hu, M. Ng, and X. Zhou, "A comparative study of ontology based term similarity measures on PubMed document clustering," in *12th international conference on Database systems for advanced applications*, Bangkok, Thailand, 2007, pp. 115-126.
- [8] F. Camous, S. Blott, and A. F. Smeaton, "Ontology-based MEDLINE document classification," in *1st international conference on Bioinformatics research and development*, Berlin, Germany, 2007, pp. 439-452.
- [9] T. Joachims, "Text categorization with support vector machines: learning with many relevant features," in *Proceedings of ECML-98, 10th European Conference on Machine Learning*, Chemnitz, Germany, 1998, pp. 137-142.
- [10] W. Hersh, C. Buckley, T. J. Leone, and D. Hickam, "OHSUMED: an interactive retrieval evaluation and new large test collection for research," in *17th annual international ACM SIGIR conference on Research and development in IR*, Dublin, Ireland, 1994, pp. 192-201.
- [11] S. Bloehdorn and A. Hotho, "Boosting for text classification with semantic features," in *6th international conference on Knowledge Discovery on the Web: advances in Web Mining and Web Usage Analysis*, Seattle, WA, 2006, pp. 149-166.
- [12] P. Wang and C. Domeniconi, "Building semantic kernels for text classification using wikipedia," in *14th ACM SIGKDD international conference on Knowledge discovery and data mining*, Las Vegas, Nevada, USA, 2008, pp. 713-721.
- [13] *Unified Medical Language System (UMLS®)*. Available: <http://www.nlm.nih.gov/research/umls/>
- [14] A. R. Aronson and F. M. Lang, "An overview of MetaMap: historical perspective and recent advances," *J Am Med Inform Assoc*, vol. 17, pp. 229-36, May-Jun 2010.
- [15] T. G. Dietterich, "Approximate statistical tests for comparing supervised classification learning algorithms," *Neural Comput.*, vol. 10, pp. 1895-1923, 1998.