

# Supervised Web Page Classification: Towards a Rocchio-Based Semantic Classification

**Abstract.** This article focuses on supervised classification presenting and comparing several methods for web page classification. Rocchio, the method we choose for its efficiency, is tested on the reference corpus "20newsGroups", adopting different similarity measures. Results illustrate some limitations mainly related to ignoring text semantics. In order to overcome these limitations, this work proposes to extend the original Rocchio through vector conceptualization replacing terms with related concepts and by adopting semantic similarity measures, all using semantic resources as ontologies, towards semantic Rocchio-based classification method.

## 1 INTRODUCTION

Actually, traditional keyword-based indexing techniques seem to be unable to respond efficiently to user queries through existing search engines. This is often related to the explosive increase in published information on the web, and also to the indexing techniques neglecting search context [1]. In order to save Internet users' time spent on checking traditional search engines' answers, it appears appropriate to apply classification techniques to these answers considering their contents. Web page classification is currently a challenging research topic, particularly in areas such as information retrieval, recommendation, personalization, user profiles etc..

Comparing the heterogeneous structure of web pages to plain text documents, web page classification can be considered as a particular case of text classification as many features can be extracted from different parts of a web page's HTML code (title, metadata, header, URL, ...) [2, 3]. Despite these differences, the principles of plain text classification also apply to web page classification.

At the beginning, text classification was completely a manual task realized by experts. Then, it was automated by the use of rules that generally bind the occurrence of certain keywords or "features" in a document to its association to a specific class. However, rule implementation and maintenance demand a lot of time and effort from experts, in addition to their limited adaptability to their original context dynamics and for each new context [2].

Consequently, learning-based techniques appeared, introducing new methods for classification. Generally based on a supervised learning, these methods use training corpus to learn decision criteria in order to be able to discriminate relevant classes. These criteria are often crystallized in induced rules, or statistical

estimations. Such supervised methods require training corpus preparation through manual tagging, that associates its documents to their relevant classes. Even if this preparation effort is significant, it is nevertheless much smaller than the effort required for rule implementation [4].

Second section concerns the most commonly used supervised methods for plain text classification (Naïve Bays, SVM, Rocchio, KNN) and also relevant for web page classification. Comparing these methods, Rocchio seems to be an efficient baseline classifier for the rest of this work and so it is evaluated, in section three, through experimentations realized on the reference corpus "20NewsGroups" using several similarity measures. Analyzing statistical results, we can relate the majority of Rocchio method's limitations to omitting semantics throughout classification process. Section four presents an extension of Rocchio overcoming these limits, so document is represented using vector of concepts instead of terms, introducing semantic classification. Finally, we conclude with an assessment of our work, followed by different research perspectives.

## 2 SUPERVISED CLASSIFICATION METHODS

For text document classification, many supervised methods are proposed in the literature. This section focuses on the most popular methods for text classification: Naïve Bayes Classifier (NB), Support Vector Machines (SVMs), Rocchio, and K Nearest Neighbor (KNN). NB classifier uses a particular document representation that models only the absence or the presence of a feature in the text using a binary vector. Other methods apply the Vector Space Model (VSM) [5] in order to represent treated text documents taking into account feature occurrence frequencies as well by means of weighting schemes (like TF/IDF).

*Naïve Bayes (NB) Classifier* [6] is a probabilistic learning method that, applying the simple *Bays' theorem*, calculates the probability of a document  $d$  occurring in class  $C$ . Once the probability model is learned, new test documents might appear to be classified. Using the learned model, the best class is assigned to the treated document; the most likely class with maximum probability that the document occurs in it.

Despite its attractive simplicity, this classifier, also called "The Binary Independence Model", has many critical weaknesses. First of all, unrealistic independence hypothesis of this model considers each feature independently for calculating their occurrence probabilities related to a class. Second, binary vectors used for document representation neglect information that can be derived from terms' frequencies in the processed document or even its

length [6]. Thus, many works propose different variations of this model to overcome its limitations [7].

*Support Vector Machines (SVMs)* [8-10] is a supervised classification method based on the assumption that classes' examples in training corpus can be linearly separated in feature space using a hyper plane defining decision boundary between classes. During learning phase, SVMs try to find this separation that maximizes margins between the examples of two classes in order to minimize classification error.

As for text classification using SVMs, the number of features characterizing documents is crucial to learning efficiency as it can significantly increment its complexity. So it is essential to this method to eliminate noisy and irrelevant features that might have negative influence on complexity and also on classification results [4]. Consequently, SVMs is considered a time and memory consuming method for text classification where class discrimination needs a considerable set of features [4].

*Rocchio classifier*, sometimes called *Centroid-based classifier*, [11] is based on the idea that documents of each class constitute a sphere in the feature space, so its centroid is considered the class prototype (or model), and so at the end of training phase this classifier calculates a set of centroids as a classification model. New document (represented in the VSM) is then compared to each of these centroids during test phase, and so the class having the most similar centroid is assigned to the document. This similarity is estimated according to a particular similarity measure.

*K Nearest Neighbor (KNN)* [12], compared to Rocchio, has no prior assumption on regularities concerning training documents' distribution in feature space. Indeed, and during test phase, the algorithm looks for the K closest documents to the processed document choosing the class of the majority of its neighbors. The absence of a distinct training phase and deploying the complete training corpus to classify each document might slow down KNN-based classification, especially when using a big corpus. In addition, as decision making depends only on the K nearest documents in the training corpus to the treated document, a wrong classification decision can be made having some noisy examples (incorrectly tagged documents) nearby [4].

Compared to other methods for text classification, Rocchio (or centroid-based classifier) has many advantages [11]. First, learned classification model summarizes the characteristics of each class through a centroid vector, even if these characteristics aren't all present simultaneously in all documents. This summarization is relatively absent in other classification methods except for NB that learns term-probability distribution functions summing up their occurrences in different classes. Another advantage is the use of similarity measure that compares a document to class centroids taking into account summarization result as well as term occurrences in the document in order to classify it. NB uses learned probability distribution only to estimate the occurrence probability of each term independently to other terms in a class summarization or to document co-occurring terms.

Vector-based representation (binary or TF/IDF) used by all methods permits semantic enrichments. Conceptualization is the process of enriching these vectors by concepts related to their features that can be retrieved using a certain background knowledge base [13]. In addition, both KNN and Rocchio enable the use of knowledge bases in decision making through new semantic similarity functions [14].

As a conclusion, we consider Rocchio an adequate baseline classifier for text and web page classification. Its efficiency,

simplicity, and extendibility with semantic resources in addition to other advantages lead us to choose Rocchio for the rest of this work. In next section, Rocchio is evaluated using different similarity measures on 20NewsGroups corpus.

### 3 EVALUATING ROCCHIO CLASSIFIER USING DIFFERENT SIMILARITY MEASURES

In previous section, Rocchio classifier has been chosen for its efficiency and semantic extendibility. This section presents an experimental study of Rocchio-based classification of text documents referring to some implementation details. Then we introduce five well known similarity measures used for decision making in our system (Cosine, Jaccard, Pearson, Averaged Kullback-Leibler Divergence, and Levenshtein distance). Using these five measures separately in experimentations enables us to evaluate Rocchio's performance independently to similarity calculation. Afterwards, results of different system configurations (each one with a different similarity measure) applied to three different variations of 20NewsGroups corpus (original corpus, reorganized corpus, six chosen classes only) are compared and analyzed. Finally, we discuss certain limitations in these results and relate some of them to the absence of semantic aspects in classification process.

#### 3.1 Rocchio Implementation details

Rocchio or centroid based classification [11] for text documents is widely used in Information Retrieval tasks, in particular for relevance feedback [15]. In Rocchio, VSM is adopted for document representation through applying four preprocessing steps (Tokenization, Stemming, StopWord Removal and Weighting). Multiple weighting schemes might be used to represent the corresponding importance of each term in a document [16] like idf, idf-prob, Odds Ratio,  $\chi^2$  etc. According to TF/IDF, the most popular schema, the score of a term  $t_j$  in document  $d_i$  is estimated as follows:

$$w_{ij} = tf_{ij} * \log(N/df_j) \quad (1)$$

$tf_{ij}$ : Frequency of term  $t_j$  in document  $d_i$ .

N: Number of documents.

$df_j$ : Number of documents that contain term  $t_j$ .

The result of applying vector space modeling to a text document is a weighted vector of features:

$$d_i = (w_{i1}, w_{i2}, w_{i3}, \dots, w_{in}) \quad (2)$$

For centroid-based classification, each class is represented by a vector positioned at the center of the sphere delimited by training documents related to this class. This vector is so called the class's centroid as it summarized all features of the class as collected during learning phase through vectors representing training documents following the VSM as detailed earlier. Having  $n$  classes in the training corpus,  $n$  centroid vectors  $\{C_1, C_2, \dots, C_n\}$  are calculated through the training phase. In order to classify a new document  $x$ , first we use the TF/IDF weighting scheme to calculate the vector representing this document in the space. Then, resulting vector is compared to all centroids of  $n$  candidate classes using a

similarity measure. So the class of the document  $x$  is the one represented by the most similar centroid.

$$\arg \operatorname{Max}_{i=1,2,\dots,n} (\operatorname{SimFun}(x, C_i)) \quad (3)$$

The most commonly used similarity measure is the Cosine measure. We present it among other similarity measures in next section.

## 3.2 Similarity Measures

Many similarity measures were used for both document classification and document clustering [17] to estimate the similarity between a document and a class prototype. Using VSM, this similarity is calculated to compare a document vector with the vector representing a class or the centroid. Next, are introduced five similarity measures (Cosine, Jaccard, Pearson, Kullback Leibler, and Levenshtein) all used in experimentations with Rocchio.

### 3.2.1 Cosine

Cosine is the most popular similarity measure and largely used in information retrieval, document clustering, and document classification research domains.

Having two vectors  $A(a_1, a_2, \dots, a_n)$ ,  $B(b_1, b_2, \dots, b_n)$ , the similarity between these vector is estimated using the cosine of the angle they delimit:

$$\operatorname{Sim}(A, B) = \cos(t) = \frac{A \cdot B}{|A| * |B|} \quad (4)$$

Where:  $A \cdot B = \sum a_i * b_i$   $|A|^2 = \sum a_i^2$   
 $i \in [0, n-1]$ ;  $n$ : the number of features in vector space.

In systems using this similarity measure, changing documents' length has no influence on the result as the angle they delimit is still the same.

### 3.2.2 Jaccard

Jaccard (sometimes called Tanimoto) estimates the similarity to the division of the intersection by the union. Having two vectors  $A(a_1, a_2, \dots, a_n)$ ,  $B(b_1, b_2, \dots, b_n)$ , according to Jaccard the similarity between A and B is by definition:

$$\operatorname{Sim}(A, B) = \frac{A \cdot B}{|A|^2 + |B|^2 - A \cdot B} \quad (5)$$

Where:  $A \cdot B = \sum a_i * b_i$   $|A|^2 = \sum a_i^2$   
 $i \in [0, n-1]$ ;  $n$ : the number of features in the vector space.

### 3.2.3 Pearson correlation coefficient

Given two vectors  $A(a_1, a_2, \dots, a_n)$ ,  $B(b_1, b_2, \dots, b_n)$ , Pearson calculates the correlation between these vectors. Deriving their centric vectors:  $A(a_1 - \bar{a}, \dots, a_n - \bar{a})$  and  $B(b_1 - \bar{b}, \dots, b_n - \bar{b})$

Where:  $\bar{a}$  is the average of all A's features,  $\bar{b}$  is the average of all B's features.

Pearson correlation coefficient is by definition the cosine of the angle  $\alpha$  between the centric vectors as follows:

$$r_{ab} = \frac{n \sum a_i b_i - \sum a_i \sum b_i}{\sqrt{[n \sum a_i^2 - (\sum a_i)^2][n \sum b_i^2 - (\sum b_i)^2]}} \quad (6)$$

This form represents also Pearson similarity measure.

### 3.2.4 Averaged Kullback-Leibler divergence

According to probability and information theory, Kullback-Leibler divergence is a measure estimating dis-similarities between two probability distributions. In the particular case of text processing, this measure calculates the divergence between feature distributions in documents. Given vectors' representations of their features distribution  $A(a_1, a_2, \dots, a_n)$ ,  $B(b_1, b_2, \dots, b_n)$ , the divergence is calculated as follows

$$D_{AvgKL}(\vec{t}_a || \vec{t}_b) = \sum_{t=1}^n (\pi_1 * D(w_{t,a} || w_t) + \pi_2 * D(w_{t,b} || w_t)) \quad (7)$$

Where:  $\pi_1 = \frac{w_{t,a}}{w_{t,a} + w_{t,b}}$ ,  $\pi_2 = \frac{w_{t,b}}{w_{t,a} + w_{t,b}}$ ,  
 $w_t = \pi_1 * w_{t,a} + \pi_2 * w_{t,b}$

### 3.2.5 Levenshtein

Levenshtein is used to compare two strings. A possible extension for vector comparison can be derived as the following equation:

$$\operatorname{Sim}(A, B) = 1 - (\operatorname{Distance} / \operatorname{Max}) \quad (8)$$

Where:

$$\operatorname{Distance}(A, B) = \sum |a_i - b_i|, \operatorname{Max}(A, B) = \sum \operatorname{Max}(a_i, b_i)$$

## 3.3 20NewsGroups corpus

20NewsGroups corpus [18] is a collection of 20,000 newsgroups documents almost evenly divided in twenty news classes according to their content topic assigned by authors. This collection is divided according to the percentages (60:40) into training corpus and test corpus respectively. Corpus organization is illustrated in Table 1. Some classes cover similar topics for example (comp.sys.ibm.pc.hardware & comp.sys.mac.hardware), whereas others concern relatively different ones as (rec.autos & sci.crypt).

**Table 1.** Twenty actuality classes of 20NewsGroups corpus

comp.graphics comp.os.ms- windows.misc comp.sys.ibm.pc.hardware comp.sys.mac.hardware comp.windows.x	rec.autos rec.motorcycles rec.sport.baseball rec.sport.hockey	sci.crypt sci.electronics sci.med sci.space
misc.forsale	talk.politics.misc talk.politics.guns talk.politics.mideast	talk.religion.misc alt.atheism soc.religion.christian

Being frequently used in evaluating state of the art supervised classification and clustering techniques, 20NewsGroups corpus seems to us a wise choice to evaluate Rocchio classifier as well.

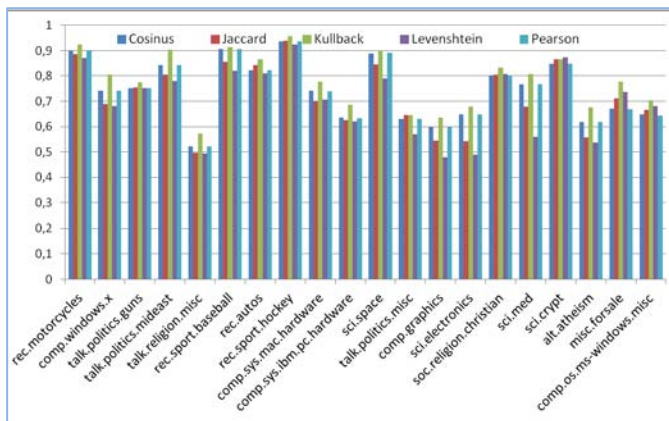
## 3.4 Experimentations and results

In these experimentations, three variations of the 20NewsGroups corpus are used: (i) the original corpus, (ii) six chosen classes, and finally (iii) the original corpus reorganized in more general six

classes as they are assembled in Table 1. Training is realized on each of these variations and so class centroids are calculated. As for test, on each variation of the corpus, five experimentations are executed applying five similarity measures (see section 3.2). For most classification tasks, classifier's accuracy [19] exceeded 90%. Thus, we use  $F_1$ -Measure [19] for performance comparison.

### 3.4.1 Experimentations on the original corpus

Our first test concerns the original corpus with twenty classes. As illustrated in **Figure. 1**, similarity measures' performance varies according to the treated class. For instance, the class (talk.religion.misc) is vast compared to other religious classes so it can draw in their documents (False negative) resulting in its relatively low value of F-measure. This observation is so called: general class issue. Classes related to computers seem to use similar vocabulary so the classifier cannot be able to distinguish them properly having similar centroids (similar class issue). On the other hand, distinct classes like (rec.sport.hockey, rec.sport.baseball) are well distinguished by all classifiers. Finally, through this diagram we can observe that the similarity measure Kullback-Leibler outperforms other similarity measures. In addition, Cosine and Pearson similarities show similar behaviors.



**Figure. 1.** Evaluating five similarity measures on the original 20NewsGroups corpus (F1-Measure).

After detailed results analysis, it is observed that at least (50%) of incorrectly classified documents (False Negative) are classified in a similar class. Indeed, similar classes, using similar vocabularies, usually have their centroids close to each other in the feature space. This implies some classification difficulties in order to distinguish classes' boundaries affecting overall performance. In addition, document contents might be related to multiple classes making classifier's task tricky.

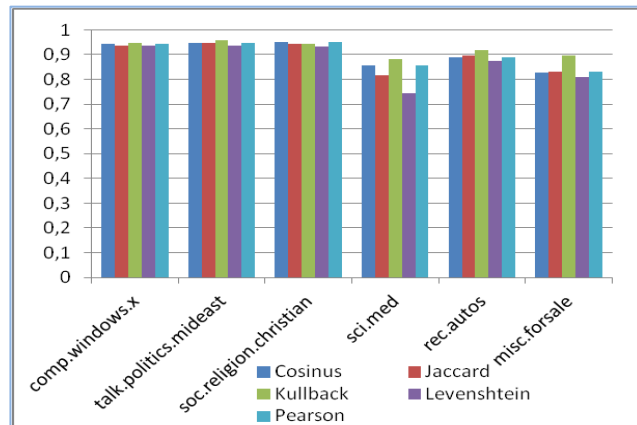
In order to support our observations, we present next two supplementary experimentations, using in the first six classes chosen from the original corpus, and in the second the original corpus reorganized in six meta-classes.

### 3.4.2 Experimentations on six chosen classes

In these experimentations, six classes, relatively distinct, of the twenty classes of the original corpus are chosen for both training and test. Classifier is first trained and then tested on the following

classes: comp.windows.x, misc.forsale, rec.auto, sci.med, soc.religion.christian, talk.politics.mideast.

Even though (sci.med) is isolated from other scientific classes, classifier's performance is still relatively poor, compared to other classes. This is due to the large distribution of medical documents in feature space so no learned centroid can be sufficiently representative. For other classes, classifier shows better performance as they are distincts, and therefore positions of their centroids are well dispersed in the feature space. Kullback-Leibler seems to outperform other similarity measures in these experimentations as well. Results are illustrated in **Figure. 2**.



**Figure. 2.** Evaluating five similarity measures on six classes of 20NewsGroups corpus (F1-Measure)

### 3.4.3 Experimentations on the corpus after reorganization

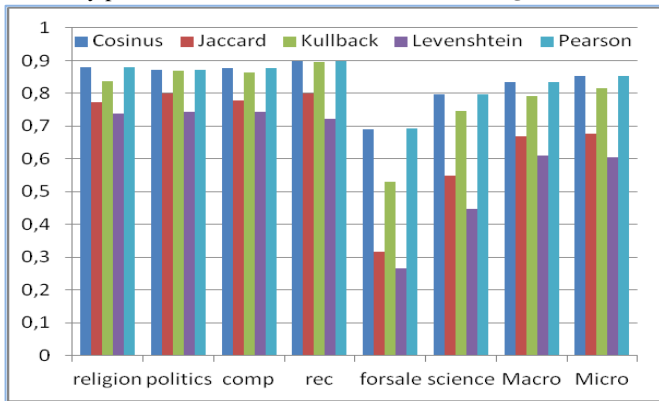
Original 20NewsGroups corpus classes are reorganized depending on initial class similarities, so documents of similar classes are gathered in a more general class resulting in six new classes in total: "comp", "rec", "science", "forsale", "politics", "religion" as presented in Table 2. Rocchio based classifier is trained on the training set and then the learned model is tested using five different similarity measures on the test set.

**Table 2.** The reorganization of 20NewsGroups corpus in six classes

comp	rec	science
comp.graphics		sci.crypt
comp.os.ms-windows.misc	rec.autos	sci.electronics
comp.sys.ibm.pc.hardware	rec.motorcycles	sci.med
comp.sys.mac.hardware	rec.sport.baseball	sci.space
comp.windows.x	rec.sport.hockey	
forsale	politics	religion
	talk.politics.misc	talk.religion.misc
misc.forsale	talk.politics.guns	alt.atheism
	talk.politics.mideast	soc.religion.christian

According to evaluation results illustrated in **Figure. 3**, classifier's performance is relatively high for most classes applying at least for one of the similarity measures. In fact these classes assemble similar original classes like (religion) or well specified classes like (rec). Classifier show some difficulties classifying (science) as the classes it assembles contain diverse information (heterogeneous class issue). In fact, one centroid for such

heterogeneous class is not very representative justifying the relatively poor value of f-measure for this class in **Figure 3**.



**Figure 3.** Evaluating five similarity measures on 20NewsGroups corpus reorganized in six classes (F1-Measure)

### 3.4.4 Conclusion

Throughout previous experimentations, several limitations seem to affect Rocchio's performance particularly in dealing with similarities among classes, general classes and heterogeneous classes. These limitations are mainly related to class representation and similarity calculations. Limitations observed with similar classes can be overcome by means of semantic resources. Indeed, redefining centroids using concepts instead of terms might limit intersections between spheres of similar classes in concept space. Consequently, ambiguities between classes using similar vocabulary can be resolved at representation level through vector conceptualization using semantic resources or ontologies. This idea and other related semantic solutions are introduced in next section.

## 4 TOWARDS A ROCCHIO-BASED SEMANTIC CLASSIFICATION OF TEXT

In spite of being considered the most popular text representation method, VSM suffers from certain limitations [20, 21] 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. Thus, text classification using conceptualized vectors is called "Semantic Classification".

This section presents first different strategies for vector conceptualization. Then, the use of conceptualized vectors during classification process is discussed according to the chosen strategy for conceptualization. These two steps constitute a semantic extension to the original Rocchio, thanks to semantic resources.

### 4.1 Vector conceptualization

Conceptualization is the process of mapping terms literally occurring in treated text to their semantically corresponding concepts or senses in semantic resources that might permit better classification. As an example of semantic resources that might be used for conceptualization: Wordnet, Wikipedia and other domain specific ontologies usually called domain ontologies.

In general, vector conceptualization is realized in two steps: search for corresponding concepts related to vector's terms and then the integration of these concepts in the vector producing the final conceptualized vector. Three different strategies have been proposed for conceptualization [13] :

- *Adding Concepts*: Where the original vector is extended and corresponding concepts are added.
- *Partial Conceptualization*: Where terms are substituted by corresponding concepts. Terms having no related concepts are held in the vector.
- *Complete Conceptualization*: Similarly to Partial Conceptualization, terms are substituted by concepts whereas remaining terms are eliminated from the final vector.

Integrated concepts are assigned new scores derived from the frequencies of their related terms.

Second strategy seems to be the most appropriate as it removes no term without replacing it with a related concept, so no original feature is removed from the vector (compared to the third one), and no extra feature is added (compared to the first one) resulting in minimized efficiency effects. Yet, the classifier has to be adapted to hybrid (concepts + terms) representation.

While searching concepts corresponding to a polysemic term in semantic resources, multiple matches are detected introducing some ambiguities in final document representation. For example: the term "Book" signifies in English a book and also a reservation (Ticket, accommodations...). According to [13], three strategies for disambiguation deploying WordNet can be used:

- *All*: Accept all candidate concepts as matches for the considered term.
- *First*: Accept the most frequently used concept among candidates using document language statistics.
- *Context*: Accept the candidate concepts having the most similar semantic context compared to the term's context in the document.

First strategy, despite being the simplest, is the least reliable as it accepts all candidate concepts without choosing a specific meaning for the term. In cases where a term is used in the document signifying its rarely used sense, second strategy chooses an incorrect concept. Despite its complexity, last strategy seems to be more accurate as concept context can be derived from semantic resources using: its definition, its descriptive terms or from text corpus.

### 4.2 Using conceptualized vectors in Rocchio

Vectors resulting from the conceptualization step are applied to decision-making criteria in order to classify corresponding documents. Most semantic classification systems neglect semantic relations between relevant concepts and apply the same pure mathematic similarity measures for decision-making. Indeed, the exact concept must occur in both compared vectors to influence decision-making through its score. Similar or related concepts are rarely taken into consideration.

Using semantic resources, it is important to incorporate related concepts to those already adopted through conceptualization step enriching document representation and eventually improving classification results. For example, some works propose integrating super concepts to the conceptualized vector and demonstrate some

ameliorations in classification results as multiple superior levels of ontology are considered [13].

New semantic similarity measures considering semantic relations among ontology concepts are also being developed making more advances towards a semantic classification. These measures permit similar concepts to be compared and contribute to vector comparison beside common concepts, [14]. For an overview on different semantic similarity measures see [22].

Integrating similarity measures in semantic classification depends on adopted conceptualization strategy. Precedent measures can be directly applied to vectors resulting from "Complete Conceptualization" (see section 4.1). Considering other conceptualization strategies that produce hybrid (concepts+terms) document representation, both mathematic and semantic similarity measures must be applied to terms and concepts respectively.

New representation models using parts of ontology hierarchy are also proposed. These parts constitute semantic trees or forests where each concept is assigned an importance score. As for decision making in systems using these models, similarity measures between semantic hierarchies is considered as the accumulation of similarities between their concepts peer to peer [23]. Semantic similarity between two concepts is related to their scores and their positions in the hierarchy.

## 5 CONCLUSION AND PERSPECTIVES

In the past decade, publishing millions of pages has resulted in the explosive growth of the web. The performance of many search engines is still limited to meet the needs of users leading to a challenging need for new filtering and ranking techniques. In this context, we presented in this paper and compared four traditional methods for text classification that can be applied to web page classification as well.

We retain Rocchio that shows a good performance compared to its relatively minimal complexity. Moreover, it may provide feedback on the relevance of classification results permitting better result understanding and so potential classification improvements. Rocchio method has been tested using different similarity measures on the 20newsGroups corpus. However, Rocchio results illustrate several limitations in text classification, which could be surpassed by means of semantic resources taking meaning into consideration in text classification.

In this work, we propose a new method for semantic classification of web pages based on Rocchio. Indeed, Rocchio allows integrating semantics in conceptualization and during decision making through different semantic similarity measures. Taking into account information stored in HTML tags, document vectorization followed by the conceptualization using a knowledge base, helps to complete VSM approach with semantics. Appearing ambiguities can be resolved using the context strategy. Finally, semantic similarity measures can also be used implying similar concepts in calculating vector similarities. The next step of our work consists in developing this method and then in applying and evaluating it on classification tasks using other web pages corpus.

## REFERENCES

[1] Asirvatham, A.P. and K.K. Ravi, *Web page classification based on document structure. Awarded second prize in National Level Student Paper Contest conducted by IEEE India Council.*, in 2001.

[2] Pierre, J.M., *On the Automated Classification of Web Sites*. Linköping Electronic Articles in Computer and Information Science, 2001. 6(1).

[3] Qi, X. and B.D. Davison, *Web page classification: Features and algorithms*. ACM Comput. Surv., 2009. 41(2): p. 1-31.

[4] Manning, C.D., P. Raghavan, and H. Schütze, *Introduction to Information Retrieval* 2008, New York, NY, USA: Cambridge University Press. 496.

[5] Salton, G., A. Wong, and C.S. Yang, *A vector space model for automatic indexing*. Commun. ACM, 1975. 18(11): p. 613-620.

[6] Lewis, D.D., *Naive (Bayes) at Forty: The Independence Assumption in Information Retrieval*, in *Proceedings of the 10th European Conference on Machine Learning* 1998, Springer-Verlag. p. 4-15.

[7] Sebastiani, F., *Machine learning in automated text categorization*. ACM Comput. Surv., 2002. 34(1): p. 1-47.

[8] Vapnik, V., *Statistical learning theory* 1998 NY: Springer-Verlag.

[9] Vapnik, V.N., *The nature of statistical learning theory* 1995, New York, NY, USA: Springer-Verlag New York, Inc.

[10] Burges, C.J.C., *A Tutorial on Support Vector Machines for Pattern Recognition*. Data Min. Knowl. Discov., 1998. 2(2): p. 121-167.

[11] Han, E.-H. and G. Karypis, *Centroid-Based Document Classification: Analysis and Experimental Results*, in *Proceedings of the 4th European Conference on Principles of Data Mining and Knowledge Discovery* 2000, Springer-Verlag. p. 424-431.

[12] Soucy, P. and G.W. Mineau, *A Simple KNN Algorithm for Text Categorization*, in *Proceedings of the 2001 IEEE International Conference on Data Mining* 2001, IEEE Computer Society. p. 647-648.

[13] Hotho, A., S. Staab, and G. Stumme, *Text clustering based on background knowledge*, 2003.

[14] Guisse, A., K. Khelif, and M. Collard, *PatClust : une plateforme pour la classification sémantique des brevets, in Conférence d'Ingénierie des connaissances* 2009: Hammamet, Tunisie.

[15] Salton, G., *The SMART Retrieval System-Experiments in Automatic Document Processing* 1971: Prentice-Hall, Inc.

[16] Lan, M., et al., *Supervised and Traditional Term Weighting Methods for Automatic Text Categorization*. IEEE Trans. Pattern Anal. Mach. Intell., 2009. 31(4): p. 721-735.

[17] Huang, A., *Similarity measures for text document clustering*. Proceedings of the Sixth New Zealand Computer Science Research Student Conference (NZCSRSC2008), Christchurch, New Zealand, 2008: p. 49-56.

[18] Rennie, J. *Home Page for 20 Newsgroups Data Set*. Mon Jan 14 2008; Available from: <http://people.csail.mit.edu/jrennie/20Newsgroups>.

[19] Sokolova, M. and G. Lapalme, *A systematic analysis of performance measures for classification tasks*. Information Processing & Management, 2009. 45(4): p. 427-437.

[20] Bloehdorn, S. and A. Hotho, *Boosting for text classification with semantic features*, in *Proceedings of the 6th international conference on Knowledge Discovery on the Web: advances in Web Mining and Web Usage Analysis* 2006, Springer-Verlag: Seattle, WA. p. 149-166.

[21] Wang, P. and C. Domeniconi, *Building semantic kernels for text classification using wikipedia*, in *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining* 2008, ACM: Las Vegas, Nevada, USA. p. 713-721.

[22] Pirro, G., *A semantic similarity metric combining features and intrinsic information content*. Data Knowl. Eng., 2009. 68(11): p. 1289-1308.

[23] Lee, J.-J., et al., *Novel web page classification techniques in contextual advertising*, in *Proceedings of the eleventh international workshop on Web information and data management* 2009, ACM: Hong Kong, China. p. 39-47.