

## RANDOM WALK TERM WEIGHTING FOR IMPROVED TEXT CLASSIFICATION

SAMER HASSAN\*, RADA MIHALCEA<sup>†</sup> and CARMEN BANE<sup>‡</sup>

*Department of Computer Science and Engineering  
University of North Texas, P.O. Box 311366  
Denton, TX 76203, USA*

*\*samer@unt.edu*

*†rada@cs.unt.edu*

*‡carmenb@unt.edu*

This paper describes a new approach for estimating term weights in a document, and shows how the new weighting scheme can be used to improve the accuracy of a text classifier. The method uses term **co-occurrence** as a measure of dependency between word features. **A random walk model is applied on a graph encoding words and co-occurrence dependencies**, resulting in scores that represent a quantification of how a particular word feature contributes to a given context. Experiments performed on three standard classification datasets show that the new random walk based approach outperforms the traditional term frequency approach of feature weighting.

*Keywords:* Random walk models; graph-based algorithms; TextRank; term weighting; text classification.

### 1. Introduction

Term frequency has long been used as a major factor for estimating the probabilistic distribution of features in a document, and it has been employed in a broad spectrum of tasks including language modeling [19], feature selection [30, 25], and term weighting [13, 21]. The main drawback associated with the term frequency method is the fact that it relies on a bag-of-words approach. It implies feature independence, and disregards any dependencies that may exist between words in the text. In other words, it defines a “random choice,” where the weight of the term is proportional to the probability of choosing the term randomly from the set of terms that constitute the text. Such an approach might be effective for capturing the relevance of a term in a local context, but it fails to account for the global effect that the term’s existence exerts on the entire text segment.

We argue that the bag-of-words model may not be the best technique to capture term importance. Instead, given that relations in the text could be preserved by maintaining the structural representation of the text, a method that takes into account the structural properties of the context could lead to a better term weighting

scheme. Previous work has shown that a higher but costly performance can be achieved by incorporating such dependencies [23].

In this paper we introduce a system that models the weighting problem as a “random walk” rather than “random choice.” We assume an imaginary reader (or “walker”) who steps through the text on a term by term basis. In this setting, the importance of the term is determined by the probability of the random walker to encounter the target term in the text during the walk.

The new measure of term weighting integrates both the locality of a term and its relation to the surrounding context. We model this local contribution using a co-occurrence relation in which terms that co-occur in a certain context are likely to share between them some of their importance (or significance). Note that in this model the relation between a given term and its context is not linear. A given term relates to a context, and the context, in turn, relates to a collection of terms. In order to model this recursive relation, we use a graph-based ranking algorithm, namely the PageRank random walk algorithm [2], and its TextRank adaptation to text processing [16]. In this paper, we show how TextRank can be used to model the probabilistic distribution of word features in a document. Through experiments performed on a text classification task, we show that the random walk scores outperform the traditional term frequencies, typically used to model the feature weights for this task.

In the following, we first overview the basic principles behind random walk algorithms, and briefly describe the TextRank application for text processing. We then show how these random walk models can be adapted to term weighting, and demonstrate that the new weighting scheme can be used to significantly improve the accuracy of a text classification system, as compared to the traditional term frequency weighting scheme. Finally, we conclude with a discussion and directions for future work.

## 2. Random Walk Algorithms

The basic idea implemented by a random walk algorithm is that of “voting” or “recommendation.” When one vertex links to another, it is basically casting a vote for this vertex. The higher the number of votes that are cast for a vertex, the higher the importance of the vertex. Moreover, the importance of the vertex casting a vote determines how important the vote itself is, and this information is also taken into account by the ranking algorithm. While there are several random walk algorithms that have been proposed in the past, we focus on only one such algorithm, namely PageRank [2], as it was previously found to be successful in a number of applications, including Web link analysis [2], social networks [8], citation analysis, and more recently in several text processing applications [16, 9].

Given a graph  $G = (V, E)$ , let  $In(V_a)$  be the set of vertices that point to vertex  $V_a$  (predecessors), and  $Out(V_a)$  be the set of vertices that vertex  $V_a$  points to (successors). The PageRank score associated with the vertex  $V_a$  is defined using a

recursive function that integrates the scores of its predecessors:

$$S(V_a) = (1 - d) + d * \sum_{V_b \in In(V_a)} \frac{S(V_b)}{|Out(V_b)|}. \quad (1)$$

where  $d$  is a parameter that is set between 0 and 1.<sup>a</sup>

The score of each vertex is recalculated upon each iteration based on the new weights that the neighboring vertices have accumulated. The algorithm terminates when the convergence point is reached for all the vertices, meaning that the error rate for each vertex falls below a pre-defined threshold.

This vertex scoring scheme is based on a random walk model, where a walker takes random steps on the graph, with the walk being modeled as a Markov process. Under certain conditions (the graph is aperiodic and irreducible), the model is guaranteed to converge to a stationary distribution of probabilities associated with the vertices in the graph [10]. Intuitively, the stationary probability associated with a vertex represents the probability of finding the walker at that vertex during the random walk, and thus it represents the importance of the vertex within the graph.

Particularly relevant for our work is the application of random walks to text processing, as done in the TextRank system [16]. TextRank has been successfully applied to three natural language processing tasks: document summarization, word sense disambiguation, and keyword extraction, with results competitive with those of state-of-the-art systems. The strength of the model lies in the global representation of the context and its ability to model how the co-occurrence between features might propagate across the context and affect other distant features. Our approach follows similar steps as used in the TextRank keyword extraction application, which derives term weights using a graph representation that accounts for the co-occurrence dependencies between words in the text. We are, however, incorporating a larger number of lexical units, and use different window sizes, as we will show in the following section.

### 3. Random Walks for Term Weighting

Starting with a given document, we determine a ranking over the words in the document by using the following models.

#### 3.1. Random walk models

In our work, we experimented with several variations of PageRank that incorporate additional information and variables into the traditional version shown in Eq. (1). We summarize the best PageRank-based term ranking models as follows:

$\vec{rw}_o$  : It represents the original model, as described in Eq. (1) in which we use an undirected graph with a constant damping factor that adheres strictly to the traditional formula of PageRank.

<sup>a</sup>The typical value for  $d$  is 0.85 [2], and this is the value we are using in our implementation.

$\overleftrightarrow{rw_{e.idf}}$  : This model represents an undirected graph approach that uses the **weighted edge version** of PageRank with a variable damping factor. The edge weight is calculated by the following formula:

$$E_{V_1, V_2} = tf.idf_{V_1} * tf.idf_{V_2}. \quad (2)$$

where  $E_{V_1, V_2}$  is the edge connecting  $V_1$  to  $V_2$ , and  $tf.idf$  represents the **term frequency** multiplied by the **inverse document frequency**.

The **damping factor** is expressed as a function of the incoming edges' weight, calculated as follows:

$$d_{E_{V_1, V_2}} = E_{V_1, V_2} / E_{max}. \quad (3)$$

where  $d_{E_{V_1, V_2}}$  is the damping function and  $E_{max}$  represents the highest weight for an edge in the graph. The resulting node ranking formula is:

$$S'(V_a) = \frac{(1-d)}{|N|} + \sum_{V_b \in In(V_a)} C * \frac{d_{E_{V_b, V_a}} * S(V_b)}{|Out(V_b)|}. \quad (4)$$

where  $N$  represents the total number of nodes in the graph,  $d$  is the damping constant, and  $C$  is a scaling constant.<sup>b</sup>

The model biases the random walker toward nodes with stronger (higher weight) edges, as opposed to nodes with weaker (lower weight) edges.

$\overleftrightarrow{rw_{e.oc}}$  : This model is similar to the above approach, however the damping factor for an edge is estimated in terms of the bigram co-occurrence frequency of the two nodes connected by the edge, Eq. (5). **For example, if the bigram "free software" occurred four times in a document**, then the weight of the edge connecting "free" and "software" is four.

$$E_{V_1, V_2} = tf(v_1 v_2). \quad (5)$$

### 3.2. Term weighting

Given a document, the following steps are applied to derive a weight associated with the words in the text:

First, the **document is tokenized** for punctuation, special symbols, and word abbreviations. **Common words are also removed, using a list of approximately 500 frequently used words.**<sup>c</sup>

Next, the resulting text is processed to extract both term frequency ( $tf$ ) and random walk ( $rw$ ) weights for each term in the document. To determine  $tf$ , we simply count the frequencies of each word in the document. To determine  $rw$ , all the

<sup>b</sup> $C$  is a scaling constant which is set to 0.95 in order to ensure an upper bound on the summation term.

<sup>c</sup>We use the list of common words distributed with the Smart system <ftp://ftp.cs.cornell.edu/pub/smart>.

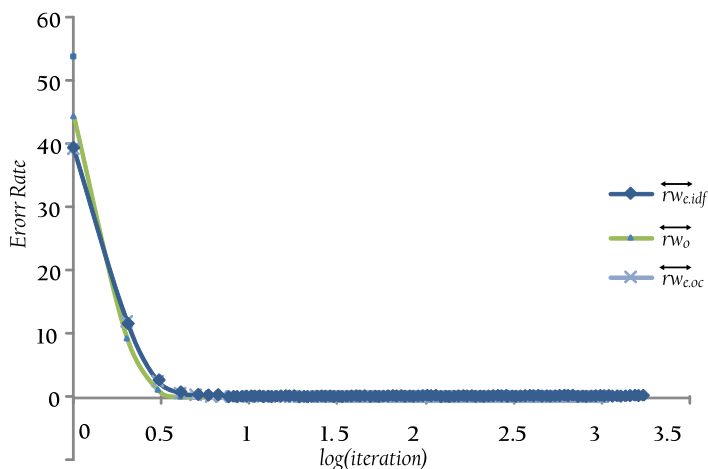


Fig. 1. Convergence graphs for the random walk models.

terms are added as vertices in a graph representing the document. A co-occurrence scanner is then applied to the text to relate the terms that co-occur within a given window size. For a given term, all the terms that fall in the vicinity of this term are considered dependent terms. This is represented by a set of edges that connect the term to all the other terms in the window. Experiments are performed for window sizes of 2, 4, 6, and 8. Once the graph is constructed and the edges are in place, the random walk algorithm is applied.<sup>d</sup> The result of the ranking is a list of all the input terms and their corresponding  $rw$  scores.

### 3.3. An example

To understand why the  $rw$  weights might be a good replacement for the traditional  $tf$  weights, consider the example in Fig. 2, which models a sample document. Starting with this text, a graph is constructed as follows. If a term has not been previously seen, then a node is added to the graph to represent this term. A term

---

*London-based sugar operator Kaines Ltd confirmed it sold two cargoes of white sugar to India out of an estimated overall sales total of four or five cargoes in which other brokers participated. The sugar, for April/May and April/June shipment, was sold at between 214 and 218 dls a tonne cif, it said.*

---

Fig. 2. Sample document.

<sup>d</sup>Unless otherwise stated, throughout this paper we refer to a random walk implementation where the damping factor is set to 0.85 as recommended in [2], and a tight convergence threshold of 0.0001. Each graph node is assigned an initial weight of 0.25.

can only be represented by one node in the graph. An undirected edge is drawn between two nodes if they co-occur within a certain window size. Figure 3 shows the graph constructed for this text, assuming a window size of 2, corresponding to two consecutive terms in the text (e.g. *London* is linked to *based*).

After the graph is constructed, the random walk model is applied on the graph, resulting in a set of scores associated with the vertices (words). Table 1 shows the *tf* and *rw* weights. By analyzing the *rw* weights, we can observe a non-linear correlation with the *tf* weights, with an emphasis given to terms surrounding important

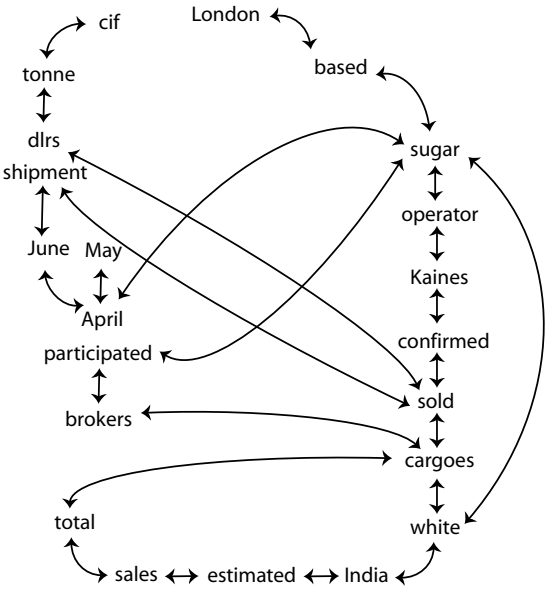


Fig. 3. Sample graph.

Table 1. *tf* & *rw* scores for a sample text.

Term	<i>rw</i>	<i>tf</i>	Term	<i>rw</i>	<i>tf</i>
sugar	16.88	3	participated	3.87	1
sold	14.15	2	april	3.87	2
based	7.39	1	india	1.00	1
confirmed	6.90	1	estimated	1.00	1
Kaines	6.81	1	sales	1.00	1
operator	6.76	1	total	1.00	1
London	4.14	1	brokers	1.00	1
cargoes	4.01	2	may	1.00	1
shipment	4.01	1	june	1.00	1
dlrs	4.01	1	tonne	1.00	1
white	3.87	1	cif	1.00	1

key terms such as *sugar* or *cargoes*. This spatial locality has resulted in higher ranks for terms like *operator* compared to other terms like *London*.<sup>e</sup>

## 4. Experimental Setup

To evaluate the random walk based approach to feature weighting, we integrate it in a text classification algorithm, and evaluate its performance on several standard text classification datasets.

Text classification is a problem typically formulated as a machine learning task, where a classifier learns how to distinguish between categories in a given set by using features automatically extracted from a collection of training documents.

### 4.1. Text classifiers

We compare the results obtained with three frequently used text classifiers — Rocchio, Naïve Bayes, and SVM, selected based on their performance and diversity of learning methodologies.

**Naïve Bayes.** The basic idea in a Naïve Bayes text classifier is to estimate the probability of a category given a document using joint probabilities of words and documents. Naïve Bayes text classifiers assume word independence, but despite this simplification, they were shown to perform surprisingly well [11, 24]. While there are several versions of Naïve Bayes classifiers (variations of multinomial and multivariate Bernoulli), we use the multinomial model [15], which was shown to be more effective.

**Rocchio.** The Rocchio text classification method uses standard term weighted vectors to represent documents, and builds a prototype vector for each category by summing up the vectors of the training documents in each category. Test documents are then assigned to the category with the closest prototype vector, based on cosine similarity. Classification experiments with different versions of the algorithm showed competitive results on standard benchmarks [11, 17].

**SVM.** Support Vector Machines (SVM) [28] is a state-of-the-art machine learning approach based on decision plans. The algorithm defines the best hyper-plan that separates the set of points associated with different class labels with a maximum-margin. The unlabeled examples are then classified by deciding on which side of the hyper-surface they reside. In our evaluations we use *SVMtorch* [5] with a linear kernel, since it was proved to be as powerful as other kernels in text classification experiments [29]. This SVM implementation is also observed to be the fastest when compared to *SVMlib* and *Weka's SMO* [12].

We use the *tf* and *rw* feature weights (and their alternatives, as described below) to create feature vectors for the SVM and the Naïve Bayes classifiers. Following

<sup>e</sup>All the missing words not shown in the graph, e.g. *Ltd*, *it*, are common words that were eliminated during pre-processing.

standard practice in term weighting for Rocchio classifiers, we use the *tf.idf* and *rw.idf* feature weights in our initial evaluation of the models (Tables 2, 3, 5 and 6).<sup>f</sup> In the final experiments (Tables 8 and 9) we report *tf* and *tf.idf* results separately. The results obtained using *tf* will act as a baseline for all the evaluations.

4.2. Datasets

We use three standard datasets: *WebKB*, *LingSpam*, and *20Newsgroups* — commonly used in text classification evaluations [27, 1, 24].

**WebKB**<sup>g</sup> is a data set collected from computer science departments of various universities. The dataset contains seven class labels: Project, Student, Department, Faculty, Staff, Course, and Other. The Other label was removed from the dataset for evaluation purposes. Most of the evaluations in the literature have been performed on only four of the categories (Project, Student, Faculty, and Course) since they represent the largest categories. However, since we wanted to see how our system behaves when only a few training examples are available, we also considered the Staff and the Department classes which have only a few training documents available. We performed our evaluations on two versions of *WebKB*: one with the four categories version (*WebKB*<sub>4</sub>) and one with the six categories (*WebKB*<sub>6</sub>).

**20Newsgroups**<sup>h</sup> is a collection of 20,000 messages from 20 newsgroups, corresponding to different topics or subjects. Each newsgroup has about 1000 messages split into 400 test and 600 training documents.

Table 2. Micro-average results for different random walk models.

Dataset	<i>tf</i>	$\overleftrightarrow{rw}_{e.oc}$	$\overleftrightarrow{rw}_{e.idf}$	$\overleftrightarrow{rw}_o$
<i>Naïve Bayes</i>				
<i>WebKB</i> <sub>4</sub>	84.2	<b>86.1*</b>	85.8*	86.1**
<i>WebKB</i> <sub>6</sub>	81.3	82.4	81.9	<b>83.3*</b>
<i>LSpam</i>	99.2	<b>99.3</b>	99.2	99.3
<i>20NG</i>	89.3	91.5**	<b>91.7**</b>	90.6**
<i>Rocchio<sup>idf</sup></i>				
<i>WebKB</i> <sub>4</sub>	84.3	87.5**	<b>87.7**</b>	86.9**
<i>WebKB</i> <sub>6</sub>	80.1	<b>84.3**</b>	84.2**	83.4**
<i>LSpam</i>	98.1	<b>98.5</b>	98.3	98.2
<i>20NG</i>	91.5	<b>94.3**</b>	93.8**	93.0**
<i>SVM</i>				
<i>WebKB</i> <sub>4</sub>	81.3	<b>94.1**</b>	93.2**	80.3
<i>WebKB</i> <sub>6</sub>	79.3	<b>90.7**</b>	90.4**	78.9
<i>LSpam</i>	93.6	98.9**	<b>99**</b>	92.1
<i>20NG</i>	90.1	<b>94.4**</b>	94.4**	86.9

<sup>f</sup>We refer to these results as *Rocchio<sup>idf</sup>*.

<sup>g</sup><http://www-2.cs.cmu.edu/afs/cs.cmu.edu/project/theo-20/www/data/>.

<sup>h</sup><http://people.csail.mit.edu/jrennie/20Newsgroups>.



Table 3. Macro-average results for different random walk models.

Dataset	$tf$	$\leftrightarrow rw_{e.oc}$	$\leftrightarrow rw_{e.idf}$	$\leftrightarrow rw_o$
<i>Naïve Bayes</i>				
<i>WebKB<sub>4</sub></i>	82.5	83.5*	82.9*	<b>84.2**</b>
<i>WebKB<sub>6</sub></i>	68.1	66.5	65.5	<b>70.0*</b>
<i>LSpam</i>	98.6	<b>98.8</b>	98.6	98.8
<i>20NG</i>	89.3	<b>91.4**</b>	90.6**	90.5**
<i>Rocchio<sup>idf</sup></i>				
<i>WebKB<sub>4</sub></i>	84.3	<b>86.7**</b>	86.7**	86.1**
<i>WebKB<sub>6</sub></i>	71.6	<b>75.8**</b>	75.3**	75.1**
<i>LSpam</i>	96.7	<b>97.3</b>	97.1	96.9
<i>20NG</i>	91.5	<b>94.3**</b>	93.7**	93.0**
<i>SVM</i>				
<i>WebKB<sub>4</sub></i>	77.4	<b>93.2**</b>	93.2**	80.3
<i>WebKB<sub>6</sub></i>	68.2	81.1**	<b>81.5**</b>	65.2
<i>LSpam</i>	86.3	97.9**	<b>98.1**</b>	81.2
<i>20NG</i>	90.6	<b>94.4**</b>	94.4**	88.1

Table 4. Running time in seconds for the processing of 1000 documents from *WebKB<sub>4</sub>*, on a Pentium-IV machine with 2Gb RAM.

$tf$	$\leftrightarrow rw_{e.oc}$	$\leftrightarrow rw_{e.idf}$	$\leftrightarrow rw_o$
112	138	138	165

**LingSpam**<sup>i</sup> is a spam corpus [1], consisting of email messages organized in 10 sets to allow for 10-fold cross validation. Each collection has roughly 300 spam and legitimate messages. There are four versions of the corpus standing for bare, stop-word filtered, lemmatized, and stop-word and lemmatized. We use the bare collection with a standard 10-fold cross validation.

## 5. Evaluation and Discussion

As a first step, we evaluate each of the random walk models presented in Sec. 3 ( $\leftrightarrow rw_o$ ,  $\leftrightarrow rw_{e.idf}$ , and  $\leftrightarrow rw_{e.oc}$ ). Tables 2 and 3 show the micro-average and macro-average accuracy figures for each model, classifier, and dataset for a window size of 2. These results are also shown graphically in Figs. 4 and 5. The  $tf$  column shows the results obtained using the term frequency weighting scheme, which, as stated before, acts as a baseline throughout all our experiments.

As seen in the tables, most of the models presented perform better than the  $tf$  baseline. Both the  $\leftrightarrow rw_{e.oc}$  and  $\leftrightarrow rw_{e.idf}$  models stand out as the best performing models with noticeable improvements, especially for the SVM and the Rocchio classifiers.

<sup>i</sup><http://boole.cs.iastate.edu/book/acad/bag/data/lingspam>.

<sup>i\*</sup>indicates a statistically significant result where  $0.05 > \rho > 0.001$ . The result is marked by \*\* when  $\rho \leq 0.001$ .

Table 5. Micro-average results for the  $rw_{e.oc}^{\leftrightarrow}$  random walk model for different window sizes.

Dataset	$tf$	$rw_2$	$rw_4$	$rw_6$	$rw_8$
<i>Naïve Bayes</i>					
<i>WebKB<sub>4</sub></i>	84.2	<b>86.1*</b>	85.8*	85.8*	85.7*
<i>WebKB<sub>6</sub></i>	81.3	<b>82.4*</b>	81.9*	81.8	76.6
<i>LSpam</i>	99.2	99.3	<b>99.3</b>	99.3	99.3
<i>20NG</i>	89.3	<b>91.5</b>	91.2	91.2	91.2
<i>Rocchio<sup>idf</sup></i>					
<i>WebKB<sub>4</sub></i>	84.3	87.5**	87.5**	87.4**	<b>87.6**</b>
<i>WebKB<sub>6</sub></i>	80.1	84.3**	84.0**	84.3**	<b>84.4**</b>
<i>LSpam</i>	98.1	<b>98.5</b>	98.3	98.4	98.3
<i>20NG</i>	91.5	94.3**	<b>94.3**</b>	94.2**	94.2**
<i>SVM</i>					
<i>WebKB<sub>4</sub></i>	81.3	<b>94.1**</b>	93.6**	93.5**	93.7**
<i>WebKB<sub>6</sub></i>	79.3	<b>90.7**</b>	90.6**	90.7**	90.6**
<i>LSpam</i>	93.6	98.9**	<b>99.1**</b>	99.1**	99.0**
<i>20NG</i>	90.1	94.4**	94.5**	94.5**	<b>94.6**</b>

Table 6. Macro-average results for the  $rw_{e.oc}^{\leftrightarrow}$  random walk model for different window sizes.

Dataset	$tf$	$rw_2$	$rw_4$	$rw_6$	$rw_8$
<i>Naïve Bayes</i>					
<i>WebKB<sub>4</sub></i>	82.5	<b>83.5*</b>	83.0*	82.9*	82.6*
<i>WebKB<sub>6</sub></i>	68.1	<b>66.5*</b>	66.0*	65.6	62.1
<i>LSpam</i>	98.6	98.8	<b>98.8</b>	98.8	98.8
<i>20NG</i>	89.3	<b>91.4</b>	91.2	91.1	91.1
<i>Rocchio<sup>idf</sup></i>					
<i>WebKB<sub>4</sub></i>	83.4	<b>86.7**</b>	86.6**	86.5**	86.7**
<i>WebKB<sub>6</sub></i>	71.6	75.8**	75.5**	75.7**	<b>76.2**</b>
<i>LSpam</i>	96.7	<b>97.3</b>	97.1	97.3	97.1
<i>20NG</i>	91.5	<b>94.3**</b>	94.2**	94.2**	94.2**
<i>SVM</i>					
<i>WebKB<sub>4</sub></i>	77.4	<b>93.2**</b>	92.9**	92.7**	92.9**
<i>WebKB<sub>6</sub></i>	68.2	81.1**	<b>81.5**</b>	81.3**	81.3**
<i>LSpam</i>	86.3	97.9**	98.3**	<b>98.4**</b>	98.2**
<i>20NG</i>	90.6	94.4**	94.5**	94.5**	<b>94.6**</b>

The  $rw_{e.idf}^{\leftrightarrow}$  and  $rw_{e.oc}^{\leftrightarrow}$  models redefine the random jump component of PageRank, by considering the damping factor as a function that can be estimated per edge. A highly connected node with relatively strong edges would tend to encourage the random walker to following its outgoing links rather than randomly jumping out. The consideration given to the relative weight of the edges signifies the encapsulation of global information in the biasing factor. This allows us, in a sense, to steer the random walker toward useful nodes more effectively, which is valuable in emphasizing the discriminative power of central features.

Table 7. Correlation between  $tf$  and  $rw$  feature weights, for different window sizes.

Dataset	$rw_2$	$rw_4$	$rw_6$	$rw_8$
<i>Naïve Bayes</i>				
<i>WebKB<sub>4</sub></i>	0.80	<b>0.81</b>	0.80	0.79
<i>WebKB<sub>6</sub></i>	<b>0.83</b>	0.80	0.81	0.80
<i>LSpam</i>	<b>0.84</b>	0.84	0.84	0.84
<i>20NG</i>	<b>0.85</b>	0.84	0.83	0.83
<i>Rocchio<sup>idf</sup></i>				
<i>WebKB<sub>4</sub></i>	<b>0.73</b>	0.71	0.69	0.69
<i>WebKB<sub>6</sub></i>	<b>0.74</b>	0.71	0.70	0.68
<i>LSpam</i>	<b>0.82</b>	0.82	0.82	0.82
<i>20NG</i>	<b>0.62</b>	0.54	0.56	0.59
<i>SVM</i>				
<i>WebKB<sub>4</sub></i>	<b>0.48</b>	0.43	0.41	0.42
<i>WebKB<sub>6</sub></i>	<b>0.53</b>	0.52	0.51	0.52
<i>LSpam</i>	0.71	<b>0.73</b>	0.72	0.71
<i>20NG</i>	0.62	0.66	<b>0.67</b>	0.66

In addition to accuracy, we also evaluated the efficiency of the new models. By comparing the processing time for 1000 *WebKB<sub>4</sub>* documents using the proposed models, we notice a small overhead of 53 seconds for the  $\vec{rw}_o$  model and 26 seconds for the  $\vec{rw}_{e.oc}$  and  $\vec{rw}_{e.idf}$ , as compared to the  $tf$  baseline (which also includes the time required for the tokenization and stopwords removal). This is due to the fast convergence of these models in approximately 15 iterations (Fig. 1). We believe that this small increase in processing time is a reasonable cost for achieving significantly higher accuracies.

### 5.1. Different window sizes

Among the various models, the  $\vec{rw}_{e.oc}$  model seems to consistently outperform the other models. To take a closer look at this model, we further analyze it under different window sizes. Tables 5 and 6 show the  $\vec{rw}_{e.oc}$  classification results for *WebKB<sub>4</sub>*, *WebKB<sub>6</sub>*, *LingSpam*, and *20Newsgroups* respectively. The  $rw_2$ ,  $rw_4$ ,  $rw_6$ , and  $rw_8$  represent the accuracies achieved using the  $\vec{rw}_{e.oc}$  weighting scheme under window sizes of 2, 4, 6, and 8 respectively.

The system displays consistent performance across different window sizes. By further analyzing the results using statistical t-tests we notice that windows of size 2 and 4 supply the most significant results across all the classifiers and the datasets.

Comparing the  $tf$  and  $rw$  weighting schemes for the *WebKB<sub>6</sub>* dataset, we found that both schemes failed to predict the class Staff. However, a significant improvement was obtained over the class Department, in which our  $rw$  model scores an accuracy of 47% compared to 4% when using  $tf$ . This could be due to the ability of the model to extract more realistic and smoother distribution of terms, hence reducing the feature bias imposed by the limited number of training examples.

Table 8. Micro-average results for different weighting schemes ( $\vec{rw}_{e.oc}$ ).

Dataset	<i>tf</i>	<i>rw</i>	<i>tf.idf</i>	<i>rw.idf</i>	$\sqrt{tf}$	$\sqrt{rw}$	$\log(tf)$	$\log(rw)$	<i>itf</i>	<i>irw</i>	$\log(tf).idf$	$\log(rw).idf$
<i>Naïve Bayes</i>												
<i>WebKB<sub>4</sub></i>	84.2	<b>86.1*</b>	81.6	<b>83.8*</b>	85.2	<b>85.9</b>	85.1	<b>85.9</b>	83.3	83.3	83.2	<b>83.8</b>
<i>WebKB<sub>6</sub></i>	81.3	<b>82.4</b>	78.8	<b>81.5*</b>	81.5	<b>81.9</b>	81.6	<b>81.9</b>	77.3	<b>77.7</b>	80.4	<b>81.7</b>
<i>LSpam</i>	99.2	<b>99.3</b>	99.1	99.1	99.3	99.3	99.3	<b>99.4</b>	99.3	<b>99.4</b>	99.0	99.1
<i>20NG</i>	89.3	<b>91.5**</b>	86.9	<b>88.6**</b>	90.8	<b>91.4</b>	77.3	<b>91.4*</b>	91.2	<b>91.5</b>	76.7	<b>88.4**</b>
<i>Rocchio</i>												
<i>WebKB<sub>4</sub></i>	74.9	<b>83.5**</b>	84.3	<b>87.5**</b>	81.2	<b>84.5**</b>	80.6	<b>84.2**</b>	82.9	<b>84.7*</b>	86.5	<b>87.8</b>
<i>WebKB<sub>6</sub></i>	70.1	<b>78.6**</b>	80.1	<b>84.3**</b>	76.3	<b>79.9**</b>	75.3	<b>79.5**</b>	78.1	<b>80.1*</b>	82.9	<b>84.3</b>
<i>LSpam</i>	96.5	<b>97.8**</b>	98.1	<b>98.5</b>	97.4	<b>98.1</b>	97.3	<b>97.9</b>	97.6	<b>98.1</b>	98.2	<b>98.5</b>
<i>20NG</i>	90.5	<b>94.8**</b>	91.5	<b>94.3**</b>	94.1	<b>94.8*</b>	67.0	<b>94.9**</b>	94.6	<b>94.7</b>	76.1	<b>94.4**</b>
<i>SVM</i>												
<i>WebKB<sub>4</sub></i>	81.3	<b>94.1**</b>	88.1	<b>90.4</b>	92.0	<b>93.8*</b>	91.4	<b>94.0*</b>	92.6	<b>92.9</b>	89.8	<b>90.6</b>
<i>WebKB<sub>6</sub></i>	79.3	<b>90.7**</b>	87.1	<b>88.7*</b>	89.7	<b>90.8</b>	89.9	<b>90.9*</b>	<b>89.9</b>	89.6	88.1	<b>88.6</b>
<i>LSpam</i>	93.6	<b>98.9**</b>	93.6	<b>93.9</b>	98.2	<b>99.1*</b>	97.6	<b>99.1**</b>	98.6	<b>98.7</b>	94.5	<b>94.7</b>
<i>20NG</i>	90.1	<b>94.4**</b>	92.5	<b>93.2*</b>	92.0	<b>93.9</b>	77.6	<b>94.5**</b>	94.6	<b>94.7</b>	93.1	<b>93.3</b>

Table 9. Macro-average results for different weighting schemes ( $\overrightarrow{rw_{e.oc}}$ ).

Dataset	$tf$	$rw$	$tf.idf$	$rw.idf$	$\sqrt{tf}$	$\sqrt{rw}$	$\log(tf)$	$\log(rw)$	$itf$	$irw$	$\log(tf).idf$	$\log(rw).idf$
<i>Naïve Bayes</i>												
<i>WebKB<sub>4</sub></i>	82.5	<b>83.5*</b>	78.9	<b>81.2*</b>	82.4	<b>83.1</b>	82.4	<b>83.2</b>	<b>78.1</b>	77.5	80.6	<b>81.1</b>
<i>WebKB<sub>6</sub></i>	<b>68.1</b>	66.5	68.5	<b>69.8*</b>	<b>66.3</b>	65.6	<b>66.9</b>	65.7	<b>52.6</b>	52.3	69.4	<b>70.2</b>
<i>LSpam</i>	98.6	<b>98.8</b>	98.4	98.4	98.7	<b>98.8</b>	98.8	<b>98.9</b>	98.9	98.9	98.3	<b>98.4</b>
<i>20NG</i>	89.3	<b>91.4**</b>	86.8	<b>88.5**</b>	90.7	<b>91.3</b>	77.1	<b>91.3*</b>	91.1	<b>91.4</b>	76.5	<b>88.4**</b>
<i>Rocchio</i>												
<i>WebKB<sub>4</sub></i>	74.4	<b>82.6**</b>	83.4	<b>86.7**</b>	80.3	<b>83.5**</b>	79.7	<b>83.2**</b>	82.0	<b>83.7*</b>	85.7	<b>86.9</b>
<i>WebKB<sub>6</sub></i>	62.7	<b>71.9**</b>	71.6	<b>75.8**</b>	69.6	<b>73.2**</b>	68.6	<b>72.7**</b>	71.5	<b>73.1*</b>	74.7	<b>75.9</b>
<i>LSpam</i>	94.1	<b>96.2**</b>	96.7	<b>97.3</b>	95.5	<b>96.7</b>	95.5	<b>96.5</b>	96.0	<b>96.7</b>	97.0	<b>97.4</b>
<i>20NG</i>	90.5	<b>94.8**</b>	91.5	<b>94.3**</b>	94.2	<b>94.8*</b>	68.1	<b>94.9**</b>	94.6	<b>94.8</b>	76.1	<b>94.4**</b>
<i>SVM</i>												
<i>WebKB<sub>4</sub></i>	77.4	<b>93.2**</b>	85.8	<b>88.6**</b>	90.8	<b>93.0*</b>	90.1	<b>93.2*</b>	91.7	<b>91.9</b>	87.9	<b>88.8</b>
<i>WebKB<sub>6</sub></i>	68.2	<b>81.1**</b>	75.9	<b>77.3*</b>	81.0	<b>82.0</b>	80.0	<b>81.7*</b>	<b>82.1</b>	81.3	77.1	<b>77.2</b>
<i>LSpam</i>	86.3	<b>97.9**</b>	86.1	<b>86.8</b>	96.6	<b>98.3*</b>	95.4	<b>98.3**</b>	97.6	<b>97.7</b>	88.4	<b>88.8</b>
<i>20NG</i>	90.6	<b>94.4**</b>	92.6	<b>93.1*</b>	90.8	<b>94.0</b>	77.6	<b>94.5**</b>	94.6	<b>94.7</b>	93.1	<b>93.3</b>

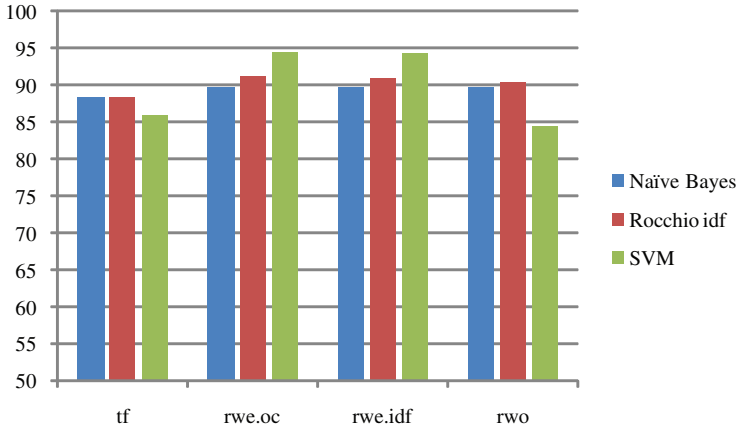


Fig. 4. Overall micro-average results for different random walk models.

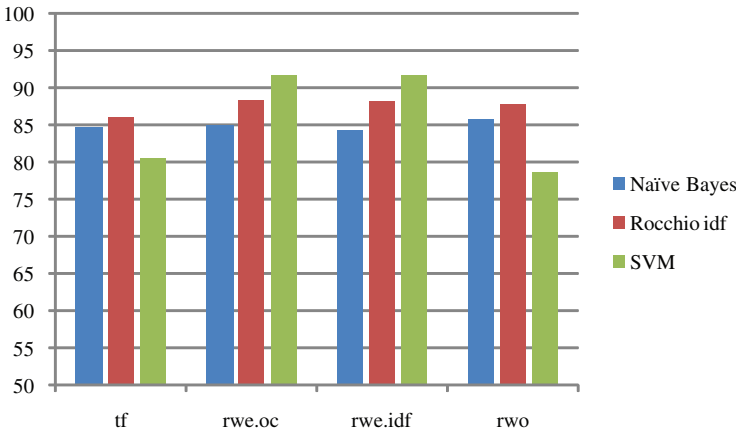


Fig. 5. Overall macro-average results for different random walk models.

We also notice a positive impact of using our model with the Rocchio classifier. In order to reach a deeper understanding of the correlation between the two models in the context of the Rocchio classifier, we macro-averaged the *tf.idf* and the *rw.idf* term scores over all the *WebKB4* documents, in a manner similar to the construction of prototype classification vectors, but taking into account the entire corpus rather than individual classes. The resulting term scores are plotted in Fig. 6. Polynomial approximations are used to visualize the trends of the plotted points, where each trendline represents a different window size.

By analyzing the graph, we can distinguish three interesting properties:

- (1) A clear non-linear correlation between the *tf* and *rw* models,
- (2) An increasing drift from the *tf* model as we increase the window size,

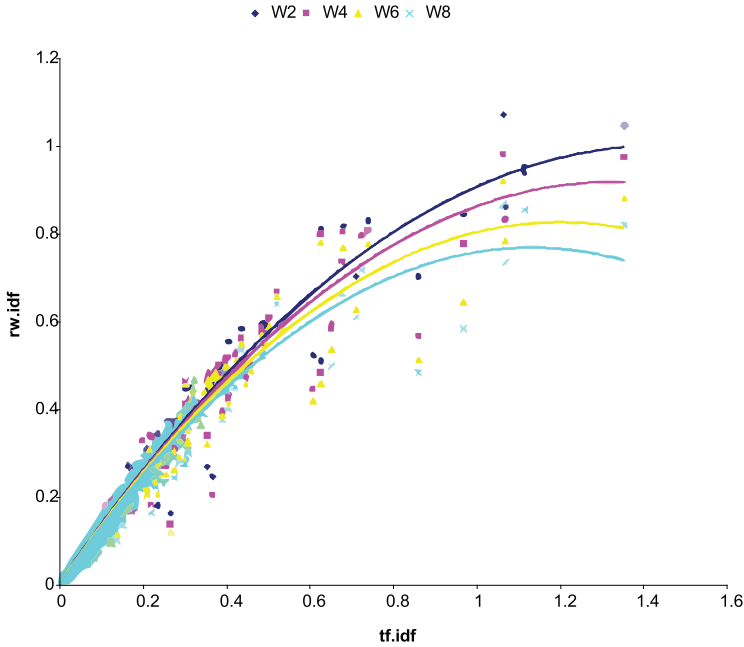


Fig. 6. Correlations of the  $tf$  and  $rw$  models for the WebKB4 collection.

- (3) A smoothing effect associated with the  $rw$  model, with a growth rate of the  $tf$  values clearly faster than  $rw$ .

We also examined the predictions' correlation between the  $tf$  and the  $rw$  models with respect to different classifiers and window sizes. The correlation, shown in Table 7, is calculated using the Pearson correlation metric, and indicates once again a weak correlation between the two models. The two models are generally more diverse and less correlated when using windows of size 6 and 8, than with windows of size 2 and 4. This is due to the increasing deviation from the feature independence assumption implied by  $tf$ . However, increasing the dependency is not always desirable, as seen in the reported results. We expect that at a certain window size the system performance will degrade to a binary weighting scheme, for a window size equal to the document size. In such a case, each term will depend on all the remaining terms resulting in an almost completely connected graph. Consequently, each feature contribution to the surrounding context will be identical, resulting in similar  $rw$  scores for all the features.

## 5.2. Other weighting scheme policies

We also compared our models to other reported state-of-the-art weighting schemes:

$tf.idf$ : since its introduction [26],  $tf.idf$  has been one of the most extensively studied weighting schemes [22, 18]. It served as a standard baseline in term weighting studies [4, 13] and proved hard to beat in [7].

*itf*: defined as one minus the inverse term frequency, this weighting scheme was found to have an excellent performance when used with an SVM linear kernel [14, 6].

*log(tf)*: first introduced in [3], this scheme was recently used with the purpose of smoothing term frequencies and hence minimizing the feature bias [20, 4].

*log(tf).idf*: this scheme was suggested in [3], and it showed superior performance in [4]. In this scheme, the smoothed term frequency is scaled by its idf to confer higher weights to domain relevant features.

$\sqrt{tf}$ : due to the interesting trends observed in Fig. 6, we introduced an approximation of the *rw/tf* correlation using the square root  $\sqrt{tf}$  function, which exhibits the general behavior of the curves plotted in Fig. 6.

For each of the schemes, we also introduce the *rw* alternative, by replacing the *tf* factor in all the schemes with the *rw* values calculated using our random walk model. For instance, for the *tf.idf* scheme, we introduce a *rw.idf* scheme; for *log(tf)*, we introduce *log(rw)*; and so forth.

The classification results obtained for the different weighting schemes on the three datasets are shown in Tables 8 and 9, and graphically illustrated in Figs. 7 and 8. Statistical significance tests were run to compare the performance of the *rw* and *tf* alternatives for each of the weighting schemes.

As seen in the tables, our random walk models clearly outperform the term frequency alternative in both micro and macro averages under all the datasets and all the classifiers. In the worst case, the random walk system performs as good as the term frequency baseline model. The superiority of our random walk models indicates

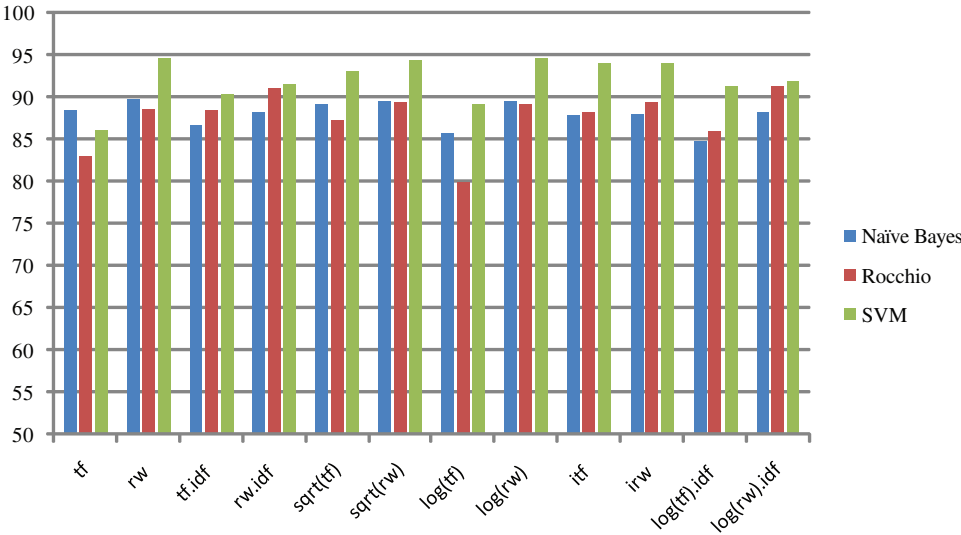


Fig. 7. Overall micro-average results for different weighting schemes.



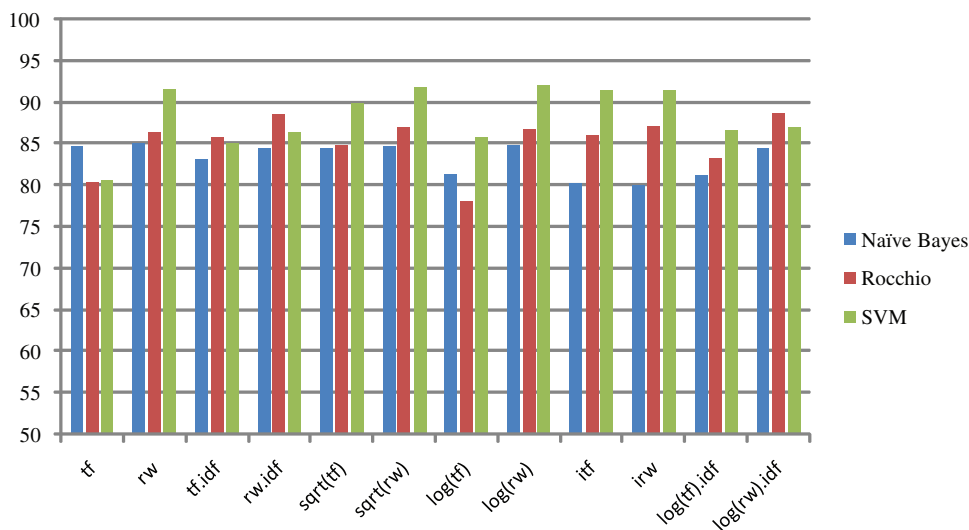


Fig. 8. Overall macro-average results for different weighting schemes.

that the use of dependencies between features can lead to significant improvements, and these improvements are consistent among different weighting schemes.

## 6. Conclusions and Future Work

In this paper, we introduced a random walk approach for term weighting that has the ability to capture term dependencies in a text by accounting for the structural properties of the text. Through experiments performed on a text classification task, we showed that the random walk model can achieve relative error rate reductions of 3.2–84.3%, as compared to the traditional term frequency based approach. The evaluation results have shown that the system's performance is consistent for various window sizes, and its running time is comparable to the *tf.idf* model.

Additionally, in experiments carried out using a variety of weighting scheme policies, the random walk term weighting was consistently found, superior as compared to the traditional term frequency weighting scheme.

We believe these results support our claim that random walk models can accurately estimate term weights by accounting for term dependencies, and can be used as a technique to model the probabilistic distribution of features in a document.

In future work we plan to extend the model and use it to define a formal language model, in which we can estimate the probability of longer n-gram sequences of words.

## Acknowledgments

This work was supported in part by a research grant from the Texas Advanced Research Program (#003594).

## References

- [1] I. Androutsopoulos, J. Koutsias, K. V. Chandrinos, G. Paliouras, and C. D. Spyropoulos, An evaluation of naive bayesian anti-spam filtering, in *Proceedings of the Workshop on Machine Learning in the New Information Age*, 2000.
- [2] S. Brin and L. Page, The anatomy of a large-scale hypertextual Web search engine, *Computer Networks and ISDN Systems* **30**(1-7) (1998).
- [3] C. Buckley, G. Salton, J. Allan, and A. Singhal, Automatic query expansion using smart: Trec 3, in *Proceedings of the Text Retrieval Conference*, 1994.
- [4] P. D. Ciya Liao and Shamim Alpha, Feature preparation in text categorization, in *Oracle Corporation*, 2002.
- [5] R. Collobert and S. Bengio, SVM Torch: Support vector machines for large-scale regression problems, *Journal of Machine Learning Research* **1** (2001) 143-160.
- [6] P. Dai, U. Iurgel, and G. Rigoll, A novel feature combination approach for spoken document classification with support vector machines, 2003.
- [7] F. Debole and F. Sebastiani, Supervised term weighting for automated text categorization, in *SAC '03: Proceedings of the 2003 ACM Symposium on Applied Computing* (New York, USA, ACM Press, 2003), pp. 784-788.
- [8] B. Dom, I. Eiron, A. Cozzi, and Y. Shang, Graph-based ranking algorithms for e-mail expertise analysis, in *Proceedings of the 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery* (San Diego, California, 2003).
- [9] G. Erkan and D. Radev, Lexrank: Graph-based centrality as salience in text summarization, *Journal of Artificial Intelligence Research* (December 2004).
- [10] G. Grimmett and D. Stirzaker, *Probability and Random Processes* (Oxford University Press, 1989).
- [11] T. Joachims, A probabilistic analysis of the Rocchio algorithm with TFIDF for text categorization, in *Proceedings of the 14th International Conference on Machine Learning* (Nashville, US, 1997).
- [12] A. Klautau, Speech recognition based on discriminative classifiers, in *Proceedings of the Simposio Brasileiro de Telecomunicacion-SBT* (Rio de Janeiro, Brazil, 2003).
- [13] M. Lan, C. Tan, H. Low, and S. Sungy, A comprehensive comparative study on term weighting schemes for text categorization with support vector machines, in *Proceedings of the 14th International Conference on World Wide Web*, 2005, pp. 1032-1033.
- [14] E. Leopold and J. Kindermann, Text categorization with support vector machines. how to represent texts in input space? In *Machine Learning*, Vol. 46 (Hingham, MA, USA, Kluwer Academic Publishers, 2002), pp. 423-444.
- [15] A. McCallum and K. Nigam, A comparison of event models for Naive Bayes text classification, in *Proceedings of AAAI-98 Workshop on Learning for Text Categorization*, 1998.
- [16] R. Mihalcea and P. Tarau, TextRank — bringing order into texts, in *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP 2004)* (Barcelona, Spain, 2004).
- [17] A. Moschitti, A study on optimal paramter tuning for Rocchio text classifier, in *Proceedings of the European Conference on Information Retrieval*, Pisa, Italy, 2003.
- [18] K. Papineni, Why inverse document frequency? In *NAACL '01: Second Meeting of the North American Chapter of the Association for Computational Linguistics on Language Technologies, 2001*, Association for Computational Linguistics (Morristown, NJ, USA, 2001), pp. 1-8.
- [19] J. M. Ponte and W. B. Croft, A language modeling approach to information retrieval, in *Research and Development in Information Retrieval*, 1998, pp. 275-281.

- [20] M. Radovanovic and M. Ivanovic, Document representations for classification of short web-page descriptions, in *DaWaK*, 2006, pp. 544–553.
- [21] R. Robertson and K. Sparck-Jones, Simple, proven approaches to text retrieval, Technical report, 1997.
- [22] S. Robertson, Understanding inverse document frequency: on theoretical arguments for idf, *Journal of Documentation* **5** (2004) 503–520.
- [23] M. Sahami, Learning limited dependence bayesian classifiers, in *Proceedings of the International Conference on Knowledge Discovery and Data Mining*, 1996, pp. 335–338.
- [24] K. Schneider, A new feature selection score for multinomial naive Bayes text classification based on kl-divergence, in *The Companion Volume to the Proceedings of 42nd Annual Meeting of the Association for Computational Linguistics*, Barcelona, Spain, July 2004.
- [25] H. Schutze, D. A. Hull, and J. O. Pedersen, A comparison of classifiers and document representations for the routing problem, in *Proceedings of the 18th Annual international ACM SIGIR Conference on Research and Development in Information Retrieval* (Seattle, Washington, 1995).
- [26] K. Sparck Jones, A statistical interpretation of term specificity and its application in retrieval, *Journal of Documentation* **28**(1) (1972) 11–21.
- [27] S. Tan, X. Cheng, M. M. Ghanem, B. Wang, and H. Xu, A novel refinement approach for text categorization, in *CIKM '05: Proceedings of the 14th ACM International Conference on Information and Knowledge Management* (Bremen, Germany, 2005), pp. 469–476.
- [28] V. Vapnik, *The Nature of Statistical Learning Theory* (Springer, New York, 1995).
- [29] Y. Yang and X. Liu, A reexamination of text categorization methods, in *Proceedings of the 22nd ACM SIGIR Conference on Research and Development in Information Retrieval*, 1999.
- [30] Y. Yang and J. O. Pedersen, A comparative study on feature selection in text categorization, in *Proceedings of the 14th International Conference on Machine Learning* (Nashville, US, 1997).