

Entropy-based Term Weighting Schemes for Text Categorization in VSM

Tao Wang*, Yi Cai*, Ho-fung Leung[†], Zhiwei Cai* and Huaqing Min*

*School of Software Engineering, South China University of Technology, China

[†] Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong
 {wtgmme, cai.zhiwei}@gmail.com, {ycai, hqmin}@scut.edu.cn, lhf@cuhk.edu.hk

Abstract—Term weighting schemes have been widely used in information retrieval and text categorization models. In this paper, we first investigate into the limitations of several state-of-the-art term weighting schemes in the context of text categorization tasks. Considering that category-specific terms are more useful to discriminate different categories, and these terms tend to have smaller entropy with respect to these categories, we then explore the relationship between a term's discriminating power and its entropy with respect to a set of categories. To this end, we propose two entropy-based term weighting schemes (i.e., *tf-dc* and *tf-bdc*) which measure the discriminating power of a term based on its global distributional concentration in the categories of a corpus. To demonstrate the effectiveness of the proposed term weighting schemes, we compare them with seven state-of-the-art schemes on a long-text corpus and a short-text corpus respectively. Our experimental results show that the proposed schemes outperform the state-of-the-art schemes in text categorization tasks with KNN and SVM.

Keywords—Term Weighting; Entropy; Text Categorization;

I. INTRODUCTION

As the increasing growth of digital documents, text categorization (TC) which automatically classifies documents into some pre-defined categories has become an effective technique to organize these materials [5]. To handle a TC task, an important step is document representation. In the Vector Space Model (VSM), documents are represented as vectors of terms so as to be processed by classifiers. Since different terms have different importance degrees to indicate the semantics of a document, term weighting schemes, which assign appropriate weights to terms, are widely used in document representation to boost TC. Although some classifiers (e.g., SVM) can learn weights for terms, term weighting in document representation is to map documents into proper positions in the vector space [1]. A high-quality mapping can help classifiers to categorize documents more effectively and achieve a better performance.

Term weighting in TC can be interpreted to measure the utility of a term in discriminating different categories. Currently, term weighting schemes are broadly classified into *unsupervised* ones and *supervised* ones according to whether the schemes exploit the category information of training documents [8]. Most unsupervised schemes are derived from Information Retrieval (IR), such as *tf*, *tf-idf* [6] and some variants [10]. As the category labels of training

documents are neglected, many unsupervised schemes, e.g., *tf-idf* and BM25 [15], are actually to measure the utility of a term in distinguishing a document from other documents rather than distinguishing a category from other categories. Thus, unsupervised schemes may be insufficient to measure the discriminating power of terms in TC.

In contrast, the category labels of training documents are exploited in supervised schemes to weight terms [8]. However, most supervised schemes weight a term based on the term's occurrences in positive category (*PC*) and negative category (*NC*), which may have the following limitations. First, in multi-class cases, *PC* is a single class while *NC* is a combination of many classes. The terms' statistics (e.g., terms' occurrences and absences) in *NC* tend to dominate terms' weights in some supervised schemes. However, these statistics are often less useful to reflect terms' discriminating power. Second, multiple categories are combined together as a single *NC* straightforwardly in most supervised schemes. However, the specific occurrences of a term in different categories of *NC* are neglected. This leads to information loss on the terms' categorical distributions, so that more specific differences of terms cannot be reflected. Third, these schemes weight terms with a given *PC*, while test documents have no class labels. Thus, a challenge is how to properly represent test documents with these schemes.

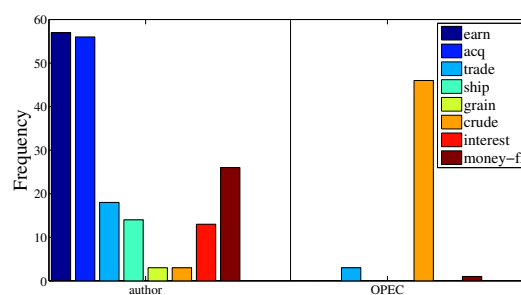


Figure 1. Categorical distributions of “author” and “OPEC”

To address these limitations, we propose two new supervised term weighting schemes in this paper. In our observation, some terms only occur in certain categories (called category-specific terms) while some terms scatter across many categories (called category-general terms) in a corpus. As category-specific terms are relatively better

indicators of categories, these terms are intuitively more useful to discriminate different categories (i.e., have more discriminating power) than category-general terms. On the other hand, a term more concentrated in a fewer categories has a smaller entropy [4]. Hence, the entropy of a term can denote the specificity that a term distributes in categories and reflect its discriminating power in TC. For example, “author” and “OPEC” are two terms in Reuters-21578. Figure 1 shows their frequencies in 8 categories of Reuters. We see that “author” scatters across each category and has a larger entropy. In contrast, “OPEC” mainly concentrates in a single category and has a smaller entropy. Obviously, “OPEC” has more discriminating power than “author”. Based on this observation, we propose two entropy-based schemes to weight terms for TC tasks. The contributions of our work are summarized as follows.

- We first reveal and elaborate the limitations of the state-of-the-art term weighting schemes in TC. To address these limitations, we then propose an entropy-based scheme *tf-dc* which weights terms based on the global concentration degree that a term distributes in the categories of a corpus. Due to no *PC/NC* splits on categories, the occurrences of a term in each category can be taken care of by *tf-dc*. Moreover, *tf-dc* weights are category-independent, and hence can properly represent test documents although they have no class labels.
- To boost the performance of *tf-dc* in the corpus with a skewed categorical distribution, we replace the absolute term frequencies in standard entropy formula with the proportions of a term in its relevant categories, and propose another new entropy-based scheme *tf-bdc*.
- To explore the answer to “Is it a better option to adopt entropy to measure the discriminating power of a term than other schemes?”, the proposed schemes are compared with seven state-of-the-art schemes based on a long-text corpus and a short-text corpus respectively. We present a detailed discussion on the experimental results and draw conclusions from different aspects.

II. INVESTIGATION ON TERM WEIGHTING SCHEMES

Some detailed introductions of term weighting schemes have been given in [8], [13]. In this paper, we focus on revealing the limitations of some widely used term weighting schemes in TC. For intuitive representation, we explore term weighting schemes through some running examples in this section. These examples are obtained from a subset of a search snippets corpus. This subset contains 1,000 documents from 4 categories of *Business (Bus)*, *Computers (Com)*, *Education (Edu)* and *Engineering (Eng)*. Table I shows document frequencies of four terms in each category.

A. Unsupervised Schemes

Most unsupervised term weighting schemes are originally proposed in IR. These schemes weight terms relying on the

Table I
TERM WEIGHTING EXAMPLES

Term	PC	NC			idf	rf	idf-qf-icf	dc
	Bus	Com	Edu	Eng				
IPO	23	0	0	0	5.442	4.644	57.938	1.0
AMD	12	8	0	0	5.644	1.807	33.102	0.515
private	12	4	3	0	5.718	1.893	25.864	0.344
market	169	10	3	6	2.411	3.446	17.865	0.691

occurrence counts of terms but neglecting the category labels of training documents. Some widely used unsupervised schemes include *binary*, *tf*, *tf-idf* and other variants [10].

Term frequency (tf) of a term in a document is one of the most intuitive term weighting schemes. Currently, *tf* has various variants, e.g., $\log(tf)$, $\log(tf + 1)$ and $\log(tf) + 1$, while previous studies [8] report that there is no significant difference in these variants. Thus, we only study the raw *tf* in this paper. Replacing any *tf* which is greater than 0 with 1, we can reduce *tf* to the *binary* scheme. Generally, *rf* serves as a local weight of a term in a document and often combines with other global factors together to weight terms.

Inverse document frequency (idf) [6] is a widely used global factor in IR. Unlike unsupervised IR tasks, the categorical membership of training documents is known beforehand in TC tasks. High-frequency terms concentrated in a category are good discriminators in TC [13] and should be given larger weights. However, *idf* cannot reflect the strength of category-term relevance due to its unsupervised nature. Let us consider an example in Table I.

Example 1. Since the document frequency (*df*) of “IPO” is 23 and “AMD” is 20, *idf* assigns a larger weight to “AMD” than “IPO”. However, “IPO” only occurs in *Bus* while “AMD” scatters across *Bus* and *Com*. As a good discriminator of *Bus*, “IPO” intuitively has more discriminating power than “AMD”. Therefore, *idf* is insufficient to reflect terms’ discriminating power in TC.

In fact, *tf-idf* is a measure to reflect the utility of a term in distinguishing a document from other documents rather than distinguishing a category from other categories. Although many variants of *tf-idf* have been proposed, such as BM25 [15] and *tff-idf* [12], these schemes share the same foundation as *tf-idf* and have the same limitation of *tf-idf* in TC.

Dumais replaces *idf* with an entropy-based factor and reports that the new scheme outperforms *tf-idf* [2]. However, Dumais calculates the entropy value of a term over each document, which is a computationally expensive task, especially in a corpus with a large number of documents. Besides, the prior category information of training documents has not been considered in this scheme. Therefore, similar to other unsupervised schemes, this scheme may be insufficient to reflect the discriminating power of terms in TC.

B. Supervised Schemes

Recently, some researchers have exploited the category labels of training documents to weight terms and proposed some supervised term weighting schemes [8], [7]. Most existing supervised schemes are category-specific, i.e., each

term has a weight vector in which each value associates with a specific category. For example, to estimate the discriminating power of a term t , each category c_i is alternately treated as positive category (PC) and other categories \bar{c}_i are combined together as negative category (NC). Table II is a category-term contingency table that shows the numbers of documents which contain and do not contain t in c_i and \bar{c}_i , and $N = a + b + c + d$. Then, the discriminating power of t for c_i is weighted based on this contingency table. The intuition behind these schemes is that terms with more discriminating power are the ones distributed more differently in PC and NC [16]. To maintain representing consistency for the running examples in Table I, we suppose the category *Bus* to be PC in these examples.

Table II
CONTINGENCY TABLE OF CATEGORY c_i AND TERM t

	c_i (PC)	\bar{c}_i (NC)
t	a	c
\bar{t}	b	d

1) *Feature-selection based Schemes*: Currently, most supervised schemes are based on feature-selection methods that use different metrics, e.g., *chi-square* (χ^2), *information gain* (ig) [1], *mutual information* (mi) [19], *gain ratio* (gr), G^2 test (*Log-Likelihood Ratio*) [14] and *eccd* [9] etc. As G^2 test is reported to approximate to χ^2 test for large samples [14], we only include χ^2 in our experiments. The basic idea in these schemes is that a term more correlated with PC is assumed to have more discriminating power.

As represented in [8], [13], most feature-selection schemes weight terms based on contingency tables and split a set of categories into PC and NC . However, in multi-class cases, PC is a single class while NC is a combination of many classes. This may lead to a high imbalance between PC and NC . As a result, item d in Table II tends to be so large (compared with a , b and c) as to dominate the results of these feature-selection schemes. However, d has less significance than a and c to indicate the discriminating power of a term [8]. Hence, these feature-selection schemes may reduce effectiveness to reflect the discriminating power of a term in multi-class cases. Note that, although *eccd* [9] uses the entropy of a term in categories as a weighting component, this scheme is also based on contingency tables and includes d in another weighting component.

2) *State-of-the-art Schemes*: To prevent b and d from dominating the weighting results, *relevance frequency* (rf) uses the ratio of a and c to estimate a term's discriminating power [8]. However, rf has another limitation of PC/NC -split based schemes. The occurrences of a term in different categories of NC are combined together as a single number c . The distributional information of a term in NC , such as the number of categories in which a term occurs, has been lost. As a result, the discriminating power of a term may be reflected inappropriately. Let us consider an example.

Example 2. “AMD” and “private” in Table I have the same occurrence in PC (i.e., a), but with different occurrences in NC (i.e., c), namely, 8 for “AMD” and 7 for “private”. Measured by rf , “AMD” has less discriminating power than “private”. However, “AMD” occurs in a smaller number of categories than “private”. Intuitively, a term which occurs in a smaller number of categories can indicate the category of a document more certainly. Thus, “AMD” has more discriminating power than “private”, which is opposite to the results of rf .

Although some variants of rf have been proposed, e.g., a logarithmically scaled rf (named vr) [13] and a probabilistic rf (named trr) [7], these variants have the same foundation as rf , as well as the same limitation of rf in TC.

$iqf \cdot qf \cdot icf$ is another scheme recently proposed in [13]. Compared with rf , the number of categories in which a term occurs, called *category frequency* (cf), is additionally considered in $iqf \cdot qf \cdot icf$. However, a single number of cf is insufficient to reflect the categorical distribution of a term, such as the probability of a term distributed in each category. Let us consider another example.

Example 3. As “private” in Table I occurs in a smaller number of categories (i.e., a smaller cf) than “market”, $iqf \cdot qf \cdot icf$ assigns a larger weight to “private”. However, comparing their categorical distributions, “market” has a higher probability to occur in *Bus*. The certainty of “market” to indicate *Bus* is relatively higher than “private”. Thus, “market” seems to have more discriminating power than “private”, which is against $iqf \cdot qf \cdot icf$ results.

In addition, a scheme based on statistical confidence intervals has been proposed in [17]. Although this supervised scheme is reported to outperform $tf \cdot idf$ and $tf \cdot gr$, it is more complicated and difficult to implement than other schemes.

III. ENTROPY-BASED TERM WEIGHTING SCHEMES

A. Intuition behind Entropy-based Schemes

As discussed above, apart from the occurrences of a term in PC , the specific occurrences of the term in each category of NC is also useful to reflect the term's discriminating power. As shown in the examples of Table I, “IPO” which only occurs in *Bus* is a good indicator of *Bus*. Thus, “IPO” has more discriminating power than other terms. Although “market” occurs in multiple categories, most occurrences of “market” concentrate in *Bus*. Compared with “AMD” and “private”, “market” has relatively more certainty to indicate *Bus*. Hence, we consider that “market” has more discriminating power than “AMD” and “private”. Similarly, “AMD” has more certainty to indicate the category of a document than “private”, since “private” is a more category-general term. Thus, “AMD” has more discriminating power than “private”. Accordingly, we have:

Observation 1. Given two terms t_1 and t_2 in a document, if t_1 only concentrates within a smaller number of categories in a corpus while t_2 scatters across a larger number of categories more uniformly, then t_1 has more certainty to indicate the document's category than t_2 ; t_1 has more discriminating power than t_2 in TC.

According to Observation 1, we assume that the discriminating power of a term depends on the global concentration degree that the term distributes in categories. A term with a higher concentration has more discriminating power. On the other hand, a term with a more concentrated distribution in categories often has smaller entropy. Thus, we then have:

Assumption 1. *Given two terms t_1 and t_2 in a document, if t_1 has a smaller entropy value than t_2 in the categories of a corpus, we assume that t_1 has more discriminating power than t_2 in TC.*

According to Assumption 1, the discriminating power of a term in TC depends on its entropy with respect to a set of categories. However, to the best of our knowledge, this assumption has not been verified in previous studies. To verify our assumption, we propose two entropy-based term weighting schemes in this paper.

B. Distributional Concentration

Since a smaller entropy denotes a higher concentration and more discriminating power, we define our first supervised scheme named *distributional concentration* (dc) as:

$$dc(t) = 1 - \frac{H(t)}{\log(|C|)} = 1 + \frac{\sum_{i=1}^{|C|} \frac{f(t, c_i)}{f(t)} \log \frac{f(t, c_i)}{f(t)}}{\log(|C|)}, \quad (1)$$

where $H(t)$ denotes the entropy of term t in the categories of a corpus, and $|C|$ is the number of categories. $f(t, c_i)$ denotes the frequency of t in category c_i , which is calculated according to the category labels of training documents. $f(t)$ denotes the sum of frequencies of t in all categories. As $H(t) \in [0, \log(|C|)]$, the values of dc can be normalized into $[0, 1]$ after the division of $\log(|C|)$. This ensures the dc values being comparable across terms in the same document.

As shown in Table I, dc weights are more reasonable and intuitive to reflect the discriminating power of these terms. Since the search snippets are quite short, each term almost occurs once per document. For intuitive comparisons with other schemes, we assume the $f(t, c_i)$ used here to be the $df(t, c_i)$ values in Table I. However, this assumption only works in these examples and $f(t, c_i)$ is calculated according to its definition above in practice.

Unlike unsupervised schemes, the proposed dc is a supervised one that exploits the category information of training documents. Unlike calculating entropy over documents in [2], we calculate entropy over categories. As the number of categories in a corpus is often much smaller than the number of documents, the computational expense of calculating entropy over categories is much less than over documents. Unlike most feature-selection schemes, dc does not involve the values of b and d , which prevents these values from dominating weighting results and reducing the effectiveness of results. Although ig involves the calculation of entropy, ig weights a term t by measuring the information change of categorical distribution in a corpus after setting t as a separator, and entropy in ig is to quantify the information of states that before and after setting separator t . In contrast, we weight

a term based on the categorical specificity/concentration of the term and use entropy to measure these specificity degrees directly. Unlike category-specific schemes, dc is a category-independent one which weights terms based on the global distributional concentration of a term in categories. That is, dc assigns weights without requiring a pre-specified PC . Thus, dc can properly represent test documents even though they have no class labels.

C. Balanced Distributional Concentration

Previous work [4] reports that the best probability distribution to represent the semantic topics of a corpus is the one with maximum entropy. However, due to lack of prior category information in IR tasks, the principle of maximum-entropy often reduces to a single constraint that the sum of prior probabilities must be one. Under this constraint, the maximum-entropy distribution is the uniform distribution. That is, the prior probability of each category is assumed to be equal in TC tasks. However, in most real-world TC tasks, different categories have different sizes (namely contain different amounts of documents). Generally, terms have more occurrences in larger categories (containing many documents) than in smaller categories (containing a few documents). This uniform assumption may result in the significance of a term to indicate larger categories being over-estimated, and cause a bias toward larger categories in classification. Let us consider an example in Table III.

Example 4. Table III shows the occurrences of term “interest” in two different-size categories *earn* and *interest* of Reuters. $f(t, c_i)$ denotes the frequency of term t in category c_i . $f(c_i)$ denotes the frequency sum of all terms in category c_i . As a larger category often has a larger $f(c_i)$, we use $f(c_i)$ to denote the size of a category. In the standard entropy of Eq. 1, “interest” has more significance to indicate *earn* than to indicate *interest*, since $f(t, c_i)$ in *earn* is larger than that in *interest*. However, relative to the category sizes (reflected by $f(c_i)$), the significance of $f(t, c_i)$ in the larger category *earn* is actually less than in the smaller category *interest*.

Table III
TERM “interest” IN CATEGORIES OF *earn* AND *interest*

c_i	$f(t, c_i)$	$f(c_i)$	$p(t c_i)$
<i>earn</i>	213	120,219	0.002
<i>interest</i>	144	12,408	0.012

To balance the disproportion of category sizes and avoid the bias toward larger categories, the frequency of a term in each category should be normalized. Hence, we replace the absolute frequency (i.e., $f(t, c_i)$) in standard entropy with the proportion of a term in its relevant categories (i.e., $p(t|c_i) = \frac{f(t, c_i)}{f(c_i)}$). As shown in Table III, “interest” is more significant to indicate *interest* than *earn* measured by $p(t|c_i)$, which is more reasonable. Hence, another scheme named

balanced distributional concentration (*bdc*) is defined as:

$$bdc(t) = 1 - \frac{BH(t)}{\log(|C|)} = 1 + \frac{\sum_{i=1}^{|C|} \frac{p(t|c_i)}{\sum_{i=1}^{|C|} p(t|c_i)} \log \frac{p(t|c_i)}{\sum_{i=1}^{|C|} p(t|c_i)}}{\log(|C|)}, \quad (2)$$

where $BH(t)$ denotes the balanced entropy value of term t .

The normalization in *bdc* is not necessary in probabilistic models, e.g., Naïve Bayes and sLDA, since these models can inherently handle prior probabilities of categories. However, probabilistic models use a sequence of terms or a binary vector to represent documents and rarely exploit term weighting schemes in document representation. Our focus in this work is to study term weighting schemes widely used in VSM.

IV. EXPERIMENTS

Next, we conduct experiments to support our analysis and verify the effectiveness of our proposed schemes. According to the number of category labels that a document has, TC can be classified into multi-label TC and single-label TC. Since a multi-label task can be transformed into a set of single-label tasks to be handled, for simplicity, we choose single-label TC as the benchmark in our experiments.

A. Datasets

To examine the adaptability of different schemes on difference kinds of datasets and compare our results with the results presented in previous works, a long-text corpus and a short-text corpus, which have been widely used in previous studies [8], [20], are used in our experiments.

Reuters corpus: As the task we considered is single-label classification, the documents with less than or more than one label in Reuters-21578 are eliminated. We use 7,674 documents from the 8 largest categories. Adopted the ModApte split, we partition the dataset into a training set with 5,485 documents and a test set with 2,189 documents. The categorical distribution in this dataset is highly skewed. The largest category (*earn*) has 51.7% of training documents and 75% of categories have less than 4.6% of documents.

Snippets corpus: This dataset [20] consists of 12,340 search snippets which are the results of web search transaction using pre-defined phrases of 8 different categories. For each query phrase put into Google search engine, the top 20 (for queries in training data) or 30 (for queries in test data) snippets from the search results are used to construct the dataset. Compared with the documents in Reuters corpus, the documents in Snippets corpus are much shorter. The average length of snippets is 14.6.

B. Methodology

We adopt seven schemes as baseline methods. Some of them (e.g., *tf* and *tf-idf*) are the widely used unsupervised schemes and some of them (e.g., *tf-chi*, *tf-ig*, *tf-eccd*, *tf-rf* and *iqf-qf-icf*) are the state-of-the-art supervised schemes. To represent test documents with category-specific schemes, e.g., *tf-chi*, *tf-ig*, *tf-eccd*, *tf-rf* and *iqf-qf-icf*, we adopt a

popular method in previous studies [1], [7] that assigning the maximum value among $|C|$ estimated weights to each term in test documents.

Since KNN and SVM are considered the two best performing VSM based classifiers [5], [18], we investigate the performance of different term weighting schemes using these classifiers. In our implement of KNN, we use the Cosine similarity as the distance measure between test documents and training documents. In predicting a category for a test document, the vote of each neighbor is weighted by its similarity to the test document [1]. The SVM classifier we used is the LIBLINEAR package [3], in which a multi-class classification task is partitioned into n binary classification tasks in one-vs.-the-rest manner. We adopt linear SVM as evaluating classifier since linear SVM is reported to outperform no-linear SVM in TC tasks [18].

The *microcoverage* and *macrocoverage* of F_1 are used in TC evaluation [17]. As the classification performance on smaller categories is more emphasized by *MacroF1* but less by *MicroF1*, the results of two metrics may have some difference. Higher F_1 values denote better performance.

C. Results

1) *Comparison in KNN:* The results of different weighting schemes in combination with KNN are shown in Figure 2. Since the number of neighbors k is an important parameter in KNN, we present the results in different values of k .

Figure 2(a) and Figure 2(b) show the *MicroF1* and *MacroF1* results on Reuters corpus. The highest *MicroF1* result is achieved by *dc* with 95.84% while the highest *MacroF1* result is achieved by *tf-bdc* with 91.81%. Since the categories in Reuters are imbalanced, we can find that the results of *MicroF1* and *MacroF1* vary greatly. For example, the performance of *iqf-qf-icf* in *MicroF1* results is much better than that in *MacroF1* results. This illustrates that *iqf-qf-icf* improves classification performance on larger categories more effectively while less on smaller categories. Consistent with the results in [8], feature-selection schemes, such as *tf-chi*, *tf-ig*, and *tf-eccd*, perform poorly. Another finding is that KNN in combination with *tf-idf* performs better gradually as k increases. The reason is that *tf-idf* tends to assign large weights to the terms with a low document frequency (i.e., rare terms in a corpus). However, some of these rare terms are noise. More neighbors help to reduce the influence of noise terms. Thus, the performance of KNN with *tf-idf* improves as the increase of k .

Figure 2(c) and Figure 2(d) show the *MicroF1* and *MacroF1* results on Snippets corpus. Both the highest *MicroF1* and *MacroF1* results are achieved by *tf-bdc* with 77.84% and 77.37% respectively. As the categories in Snippets are approximately balanced, the trends of each scheme in *MicroF1* and *MacroF1* results are much similar. Since *tf-rf* is originally proposed for long document categorization, the performance of *tf-rf* on short-text Snippets is not as good

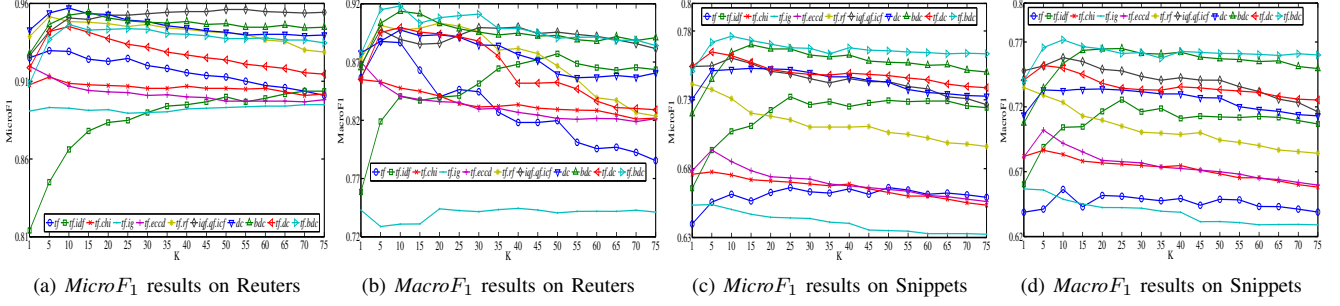


Figure 2. Comparison of term weighting schemes with KNN in different values of K

as that on long-text Reuters. Similarly, feature-selection schemes *tf-chi*, *tf-ig* and *tf-eccd* perform worse on Snippets. The reason is that the short-text snippets have sparser vector representations. This leads to *d* in Table II being larger, so that feature-selection schemes are more likely to be dominated by *d* and lose effectiveness to reflect terms' discriminating power. Furthermore, *bdc* adopts a normalization to avoid the bias toward larger categories in classification and improve the performance on smaller categories. Hence, we can find that *bdc* and *tf-bdc* perform better than *dc* and *tf-dc* in most cases on two datasets. These improvements are more obvious in *MacroF1* results, since the performance on smaller categories is more emphasized by *MacroF1*.

Table IV
COMPARISON OF TERM WEIGHTING SCHEMES WITH SVM

Scheme	Reuters		Snippets	
	<i>MicroF1</i>	<i>MacroF1</i>	<i>MicroF1</i>	<i>MacroF1</i>
<i>tf</i>	97.17%	93.48%	70.18%	70.14%
<i>tf-idf</i>	97.85%	94.88%	70.98%	70.81%
<i>tf-chi</i>	94.20%	89.68%	68.34%	69.05%
<i>tf-ig</i>	92.14%	75.97%	65.96%	66.16%
<i>tf-eccd</i>	95.20%	92.00%	68.60%	69.21%
<i>tf-rf</i>	97.30%	94.28%	71.50%	71.35%
<i>iqr-qf-icf</i>	96.85%	92.29%	70.89%	70.04%
<i>dc</i>	97.21%	93.78%	71.11%	69.83%
<i>bdc</i>	96.98%	92.44%	71.59%	70.32%
<i>tf-dc</i>	97.49%	94.51%	71.90%	71.66%
<i>tf-bdc</i>	97.58%	95.06%	72.38%	72.02%

2) *Comparison in SVM*: Table IV shows the performance of eleven term weighting schemes in combination with SVM classifier on Reuters and Snippets. Since SVM has self-optimizing to learn term weights [1], the comparison results of SVM are slightly different from those of KNN. For example, *tf-idf* performs better in SVM than in KNN, and it achieves the best *MicroF1* on Reuters in SVM. However, the proposed *tf-bdc* performs best in *MacroF1* results. In addition, the proposed *tf-bdc* performs best on Snippets in both *MicroF1* and *MacroF1* results. Comparing the results of *dc*, *bdc*, *tf-dc* and *tf-bdc*, we can find that *tf-dc* and *tf-bdc* which combine with *tf* factor perform better than single *dc* and *bdc* consistently. Hence, we assume that the local weight of a term in a document (reflected by *tf* factor) is more emphasized in SVM. However, *iqr-qf-icf* does not include *tf* factor and its performance in SVM is inferior to that in

KNN. In contrast, *tf-rf* performs better in SVM than in KNN. Another finding is that KNN outperforms SVM on Snippets. A possible reason is the sparse vector representation of short-text snippets. Previous study [11] reports that SVM performs poorly in sparse datasets.

3) *Comparison in Binary Classification*: To examine the effectiveness of our proposed schemes in binary classification cases, we draw two subsets of Reuters and Snippets datasets respectively as binary classification benchmarks. The first dataset includes the *crude* and *trade* categories from Reuters corpus, with 374 documents from *crude* category (253 for training and 121 for test) and 326 documents from *trade* (251 for training and 75 for test). Another dataset is extracted from the *Business* and *Computers* categories of Snippets corpus. Both *Business* and *Computers* contain 1,500 documents (1,200 for training and 300 for test).

Table V
RESULTS OF REUTERS AND SNIPPETS SUBSETS

Scheme	KNN		SVM	
	<i>MicroF1</i>	<i>MacroF1</i>	<i>MicroF1</i>	<i>MacroF1</i>
<i>tf</i>	99.49%	99.46%	98.98%	98.93%
<i>tf-idf</i>	97.96%	97.85%	98.98%	98.93%
<i>tf-chi</i>	98.98%	98.93%	99.49%	99.46%
<i>tf-ig</i>	98.98%	98.93%	99.49%	99.46%
<i>tf-eccd</i>	98.98%	98.93%	99.49%	99.46%
<i>tf-rf</i>	100.00%	100.00%	100.00%	100.00%
<i>iqr-qf-icf</i>	99.49%	99.46%	98.47%	98.39%
<i>dc</i>	100.00%	100.00%	100.00%	100.00%
<i>bdc</i>	100.00%	100.00%	100.00%	100.00%
<i>tf-dc</i>	99.49%	99.46%	100.00%	100.00%
<i>tf-bdc</i>	99.49%	99.46%	100.00%	100.00%
<i>tf</i>	87.67%	87.67%	88.17%	88.12%
<i>tf-idf</i>	90.83%	90.83%	88.50%	88.45%
<i>tf-chi</i>	84.83%	84.72%	86.33%	86.29%
<i>tf-ig</i>	85.17%	85.09%	87.00%	86.97%
<i>tf-eccd</i>	86.83%	86.81%	86.67%	86.63%
<i>tf-rf</i>	90.50%	90.49%	90.00%	89.97%
<i>iqr-qf-icf</i>	90.67%	90.66%	88.17%	88.10%
<i>dc</i>	93.00%	93.00%	89.67%	89.64%
<i>bdc</i>	93.00%	93.00%	89.67%	89.64%
<i>tf-dc</i>	93.00%	92.99%	90.67%	90.63%
<i>tf-bdc</i>	93.00%	92.99%	91.17%	91.14%

Table V shows the best results of eleven weighting schemes on two subsets with KNN and SVM. We can see that *tf-rf*, *dc* and *bdc* perform best in both KNN and SVM, and *tf-dc* and *tf-bdc* work well in SVM. Apart from the good performance in multi-class cases, the proposed schemes can achieve competitive performance in binary

cases as well. In our schemes, the terms only occurring in a single category are considered to be good discriminators. These terms are given large weights and have significant power to determine which category a document belongs to in classification. The relevance between documents and categories can be indicated by these single-category occurring terms even though our weights are not category-specific. On the other hand, since two categories in this binary classification are more balanced, *PC* and *NC* are more balanced when feature-selection schemes weigh terms. This alleviates the imbalance between *PC* and *NC*, and prevents *d* from dominating weighting results in multi-class cases. Thus, feature-selection schemes, e.g., *tf·chi*, *tf·ig* and *tf·eccd*, perform better in binary classification.

V. DISCUSSIONS

A. Performance on Representing Test Documents

Previous supervised schemes weight terms depending on a pre-specified *PC*. However, test documents have no class labels. This brings a challenge on appropriately representing test documents and may reduce TC performance. Let us consider an example in Snippets weighted by *tf·chi*.

Example 5. Two documents in Snippets corpus belong to different categories, $d_1 < \text{math:}0.7, \text{homework:}0.2 > \in \text{Education}$ and $d_2 < \text{math:}0.1, \text{programmer:}0.5 > \in \text{Computers}$. If we take d_2 as a test document, we can obtain a test vector of $\hat{d}_2 < \text{math:}0.7, \text{programmer:}0.5 >$ by adopting the maximum-weight representing method in [1], [7]. Calculated the Cosine similarities of these documents, we can find that $\cos(d_1, \hat{d}_2) = 0.78$ is larger than $\cos(d_2, \hat{d}_2) = 0.73$. Thus, \hat{d}_2 will be classified into Education, while \hat{d}_2 actually belongs to Computers.

Table VI
MATCHING ACCURACY OF TERM WEIGHTING SCHEMES

Corpus	<i>tf·chi</i>	<i>tf·ig</i>	<i>tf·eccd</i>	<i>tf·rf</i>	others
Reuters	93.78%	92.84%	65.76%	99.91%	100.00%
Snippets	86.12%	86.93%	69.49%	99.92%	100.00%

To illustrate the bias of category-specific schemes in representing test documents, we use the training data of two datasets (used in the above experiments) as test sets respectively, and examine whether the most similar document of each test document could match itself. Table VI shows the matching accuracy of eleven weighting schemes using KNN with $k = 1$. As shown in Table VI, most category-specific schemes, e.g., *tf·chi*, *tf·ig*, *tf·eccd* and *tf·rf*, have lower accuracy than the category-independent schemes, e.g., *tf*, *tf·idf*, *tf·dc* and *tf·bdc* etc.

Although category-specific weights are intuitively reasonable in representing training documents, previous schemes have no proper strategies to properly represent test documents. This may instead lead to poor classification performance. In contrast, the proposed schemes weight a term based on its global concentration in categories, which is

category-independent. Thus, test documents can be properly represented with our schemes, so as to improve TC.

B. Quality of Term Weights

In this subsection, we compare the quality of term weights obtained using different schemes. Figure 3(a) shows the correlation of *dc* and *rf* weights of each term in Snippets corpus, and Figure 3(b) shows the correlation of *dc* and *iqf·qf·icf*. The *rf* and *iqf·qf·icf* weights are the maximum values among $|C|$ estimated weights of each term.

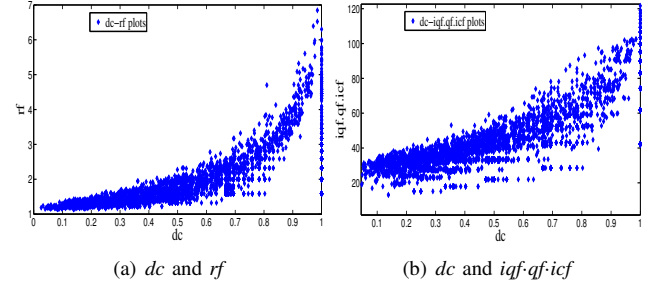


Figure 3. Correlation of weights in Snippets

Clearly, the proposed *dc* is highly correlated to the state-of-the-art schemes *rf* and *iqf·qf·icf*. These schemes thus have comparable performance in experiments. However, for the terms with $dc < 0.5$ (i.e., the terms occurring in multiple categories), the *rf* weights are concentrated in $[1, 2]$ and *iqf·qf·icf* in $[20, 40]$. Relative to the ranges of *rf* ($[1, 7]$) and *iqf·qf·icf* ($[0, 120]$) weights, these weights are so close that the different discriminating power of these terms cannot be reflected significantly. The reason is that these schemes combine multiple categories as a *NC* straightforwardly, but neglect the distribution information of a term in different categories of *NC*. This causes information loss on terms' distributions in *NC*, so that more specific difference of terms is not reflected by these schemes. In contrast, the probability of a term in each category is considered in our *dc*. The different degrees of terms' discriminating power are reflected more significantly by *dc*, and hence *dc* weights are more distinguishable. This is one reason why our schemes have better performance.

Another finding is that the *rf* (as well *iqf·qf·icf*) weights of the terms with $dc = 1$ (i.e., the terms only occurring in one category) vary greatly. The reason is that *rf* and *iqf·qf·icf* use the absolute document frequencies of terms to give weights. This causes that some category-specific terms, such as "Adidas" for the *Sports* category in Snippets corpus, are given small weights due to their low document frequencies. However, these category-specific terms are good discriminators of categories and should be given large weights. In our schemes, we use the probability of a term distributed in each category rather than the absolute frequency of a term in each category, which avoids the unfairness for category-specific terms in *rf* and *iqf·qf·icf*. This is another reason why our proposed schemes have a superior performance.

VI. CONCLUSIONS AND FUTURE WORK

In this paper, we explore the limitations of several state-of-the-art term weighting schemes in TC and propose two entropy-based supervised weighting schemes *tf-dc* and *tf-bdc*. Unlike most existing supervised schemes that often convert a multi-class case into multiple binary cases, the proposed schemes inherently permit multi-class cases since we weight a term by its global distributional concentration with respect to the categories in a corpus. Our experimental results show the proposed schemes outperform the state-of-the-art schemes in TC tasks. Thus, the answer to the question in Section 1 is that entropy can reflect the discriminating power of a term in TC more effectively than most other schemes. Furthermore, we have the following findings. First, rare terms as well as noise terms are over-weighted by *tf-idf*, and hence more neighbors in KNN are needed to alleviate the influence of noise and achieve optimal performance. Second, the feature-selection based schemes (e.g., *tf-chi* and *tf-ig*) perform better in binary classification than in multi-class classification. Third, compared with *dc* factor, *bdc* factor improves the classification performance on small categories significantly. Finally, the local weight of a term in a document (reflected by *tf*) is more emphasized in SVM, and the schemes with *tf* (e.g., *tf-dc* and *tf-bdc*) perform better than those without *tf* (e.g., *dc* and *bdc*) in SVM.

In the future, we will verify the effectiveness of the proposed schemes in other datasets and classifiers.

ACKNOWLEDGMENT

This work is supported by National Natural Science Foundation of China (Grant NO. 61300137), the Guangdong Natural Science Foundation, China (NO. S2013010013836), Science and Technology Planning Project of Guangdong Province China NO. 2013B010406004, the Fundamental Research Funds for the Central Universities, SCUT(NO. 2014ZZ0035).

REFERENCES

- [1] Iyad Batal and Milos Hauskrecht. Boosting knn text classification accuracy by using supervised term weighting schemes. In *CIKM*, pages 2041–2044. ACM, 2009.
- [2] Susan T Dumais. Improving the retrieval of information from external sources. *Behavior Research Methods, Instruments, & Computers*, 23(2):229–236, 1991.
- [3] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. Liblinear: A library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874, 2008.
- [4] Edwin T Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620, 1957.
- [5] Thorsten Joachims. *Text categorization with support vector machines: Learning with many relevant features*. Springer, 1998.
- [6] Karen Sparck Jones. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21, 1972.
- [7] Youngjoong Ko. A new term-weighting scheme for text classification using the odds of positive and negative class probabilities. *Journal of the Association for Information Science and Technology*, 2015.
- [8] Man Lan, Chew Lim Tan, Jian Su, and Yue Lu. Supervised and traditional term weighting methods for automatic text categorization. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 31(4):721–735, 2009.
- [9] Christine Largeron, Christophe Moulin, and Mathias Géry. Entropy based feature selection for text categorization. In *Proceedings of the 2011 ACM Symposium on Applied Computing*, pages 924–928. ACM, 2011.
- [10] Edda Leopold and Jörg Kindermann. Text categorization with support vector machines. how to represent texts in input space? *Machine Learning*, 46(1-3):423–444, 2002.
- [11] Xinghua Lu, Bin Zheng, Atulya Velivelli, and ChengXiang Zhai. Enhancing text categorization with semantic-enriched representation and training data augmentation. *Journal of the American Medical Informatics Association*, 13(5):526–535, 2006.
- [12] Jiaul H Paik. A novel tf-idf weighting scheme for effective ranking. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*, pages 343–352. ACM, 2013.
- [13] Xiaojun Quan, Wenyin Liu, and Bite Qiu. Term weighting schemes for question categorization. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(5):1009–1021, 2011.
- [14] Paul Rayson, Damon Berridge, and Brian Francis. Extending the cochrane rule for the comparison of word frequencies between corpora. In *7th International Conference on Statistical analysis of textual data (JADT 2004)*, pages 926–936, 2004.
- [15] Stephen Robertson. Understanding inverse document frequency: on theoretical arguments for idf. *Journal of documentation*, 60(5):503–520, 2004.
- [16] Fabrizio Sebastiani. Machine learning in automated text categorization. *ACM computing surveys (CSUR)*, 34(1):1–47, 2002.
- [17] Pascal Soucy and Guy W Mineau. Beyond tfidf weighting for text categorization in the vector space model. In *IJCAI*, volume 5, pages 1130–1135, 2005.
- [18] Yiming Yang and Xin Liu. A re-examination of text categorization methods. In *SIGIR*, pages 42–49. ACM, 1999.
- [19] Yiming Yang and Jan O Pedersen. A comparative study on feature selection in text categorization. In *ICML*, volume 97, pages 412–420, 1997.
- [20] Shitao Zhang, Xiaoming Jin, Dou Shen, Bin Cao, Xuetao Ding, and Xiaochen Zhang. Short text classification by detecting information path. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pages 727–732. ACM, 2013.