

An Improved Term Weighting Method Based on Relevance Frequency for Text Classification

Sichuan University

Wenqiang Li

Sichuan University

Zhong Tang

Sichuan University

Song Li

Sichuan University

Hai Xiang

Sichuan University

Research Article

Keywords: text classification, term weighting, relevance frequency, logarithm processing

Posted Date: August 6th, 2021

DOI: https://doi.org/10.21203/rs.3.rs-680515/v1

License: © 1) This work is licensed under a Creative Commons Attribution 4.0 International License.

Read Full License

An improved term weighting method based on relevance frequency for

text classification

Chuanxiao Li, Wenqiang Li*, Zhong Tang, Song Li, Hai Xiang School of Mechanical Engineering, Sichuan University, Chengdu 610065, China Innovation Method and Creative Design Key Laboratory of Sichuan Province, Chengdu 610065, China

Abstract: As a vital step of text classification (TC) task, the assignment of term weight has a great influence on the performance of TC. Currently, masses of term weighting schemes can be utilized, such as term frequency-inverse documents frequency (TF-IDF) and term frequency-relevance frequency (TF-RF), and they are all consisted of local part (TF) and global part (e.g., IDF, RF). However, most of these schemes adopt the logarithmic processing on their respective global parts, and it is natural to consider whether the logarithmic processing apply to all these schemes or not. Actually, for a specific term weighting scheme, due to its different ratio of local weight and global weight resulting from logarithmic processing, it usually shows diverse text clasification results on different text sets, which presents poor robustness. To explore the influence of logarithmic processing imposed on the global weight on the classification result of term weighting schemes, TF-RF is selected as the representative because it can achieve a better performance among these schemes adopting logarithmic processing. Then, two propositions along with corresponding methods about the relation between TF part and RF part are proposed based on TF-RF. In addition, two groups of experiments are conducted on the two methods. The first group of experiments proves that one method (denoted as TF-ERF) is more helpful to the improvement than the other one (denoted as ETF-RF). The second group of experiments shows that TF-ERF not only ourperforms TF-RF but also obtains better performance than other existing term weighting schemes.

Keywords: text classification, term weighting, relevance frequency, logarithm processing

1 Introduction

Nowadays, with the development of Internet technology, we have entered the information age quickly. Massive data is created and transformed to the virtual environment rapidly every day and most of them exists in a form of textual document [1,2,3]. Meanwhile, people's demand that improving the matching degree between retrieval words and provided documents is also raising [4]. Therefore, it is essential to classify these documents accurately according to their content. However, confronted with the continuously growing text data, it is inefficient and impractical to deal with a large amount of information only by manual work [5]. Luckily, text classification (TC) technology rising in recent years can be utilized to accomplish this target [6,7]. TC is such a task which aims at assigning corresponding category labels to the documents in accordance with their respective topics by specific classification algorithms [8,9,10]. This technology has been applied in many fields, including spam filtering [11], web page detection and function extraction of patent texts [12,13]. However, computers can't identify document texts like human beings do. Therefore, it is essential to transform the format of these texts into an appropriate one in order to make them recognized by computers and classifiers successfully, and this process of transformation is called text representation [14]. Currently, vector space model (VSM) is a widely-used text

E-mail address: liwenqiang@scu.edu.cn (W. Li), tangzscu@163.com (Z. Tang), sculisong1992@163.com (S. Li), 1270978696@qq.com (H. Xiang).

^{*} Corresponding author.

representation method in which documents are transformed into vectors weighted by specific measurements [15,16]. In this model, the document can be represented by the form of $d_k = (t_1, t_2, ..., t_n)$ with a corresponding weight of $w = (w_1, w_2, ..., w_n)$, in which n denotes the number of selected features and w is the set of all term weights [17,18]. In the whole process, the assignment of term weight is a vital step because the weight demonstrates the importance of a specific term and the contribution made by this term in classifying different kinds of texts [3,14,19].

In order to assign appropriate weights to terms, it is essential to choose a reasonable term weighting scheme. Generally, term weighting schemes can be mainly divided into two categories, supervised term weighting (STW) schemes and unsupervised term weighting (UTW) schemes according to whether the predefined category information is utilized or not [14,20,21,22]. Due to different term weighting schemes adopt various models, therefore, even if these schemes belong to the same UTW or STW scheme, there are still many differences between them. For example, the TF scheme only takes the term frequency into account when computing the term weight, and it holds the assumption that the higher the frequency of appearance of a term, the more important this term is. On the contrary, the IDF scheme is based on the assumption that the less documents in which a term appears, the more significant the term will be, which ignores the effect of term frequency completely [9,23]. Considering that both TF and IDF have respective defects as a single term weighting scheme, TF-IDF was then proposed by combining them together [24], which supports the idea that the term weight should be measured from both term frequency and document frequency comprehensively [25,26]. It should be noticed that the three schemes mentioned above are all UTW schemes as the available category information isn't utilized in the text training process. From the core ideas of the above schemes, we can see that the term weight measured by them only reflect the relationship between term and document as well as document and document. However, TC is a supervised learning task which aims at classifying various texts into different categories according to their content [14]. It seems that UTW schemes can't achieve the goal of TC with a satisfactory result.

In order to obtain a more reasonable term weighting scheme, it is natural to take the category information into account to match the TC task which is also a supervised learning process. Relative to the UTW scheme, the method which make use of the prior category information when computing the term weight, is called STW scheme [19]. Most STW schemes follow the pattern of TF-IDF which is consisted of local part and global part. Besides, it is publicly accepted that term frequency (TF) is an excellent representation of local weight, so the TF part is still retained while the IDF part is substituted by others in these new-proposed STW schemes. For example, considering that the methods Chi-square (CHI2), information gain (IG) and mutual information (MI) perform well in the feature selection procedure, there is probability that these methods also apply to measuring the term weight equally. Therefore, replacing the IDF with CHI2, IG and MI separately, three new STW schemes are proposed, which are named as TF-CHI2, TF-IG and TF-MI [3]. In view of that the number of categories containing the specific term may contain useful information for TC task, then two STW schemes named TF-ICF [20], TF-IDF-ICF [21] generated as a result. Chen et al. [22] proposed that the traditional TF-IDF is not fully effective for TC task, in order to make some improvements, based on a new statistic model, a new STW scheme named TF-IGM and its variants were proposed which claims that this method can make full use of the fine-grained term distribution across different classes of texts. In addition, there are also a variety of STW schemes based on

different models like TF-OR [27] and TF-PB [28].

Intuitively, the STW schemes should have performed better than UTW ones in terms of text classification performance because they make full use of the predefined category information. Actually, as the representative of UTW schemes, TF-IDF shows better performance than some STW schemes, this phenomenon is conflict with our intuition [14,29]. Aiming at this problem, we have analyzed the above listed supervised term weighting schemes and found that part of them become invalid under some special circumstances. For example, the category frequency (CF) part of TF-ICF and TF-IDF-ICF, it represents the total number of categories whose documents contain the chosen term. In the two models, the number of documents containing the chosen term in a specific category has no effect on the term weight, that is to say, one document or ten documents belong to the same category in which the chosen term appears is regarded as no difference, this is obviously unreasonable. To eliminate this defect, TF-IDF-ICSDF [21] was proposed implementing a new model named inverse class-space-density frequency. However, it will degenerate into TF-IDF when the number of documents in each category is the same [22]. It is not a unique instance, a similar situation will also happen to TF-IGM which makes it become invalid when the number of different kinds of documents meets certain conditions. Hence, in order to offset this shortcoming, a novel method named TF-IGM_{imp} [30] was proposed by adding a ratio to the initial TF-IGM. Among these STW schemes, term frequency-relevance frequency (TF-RF) proposed by Lan et al [14] is considered as an outstanding scheme with reasonable theoretical explanation and good classification performance. More importantly, similar failure circumstances will not happen in TF-RF. it is noticeable that most of the listed STW schemes adopt the same logarithmic processing borrowed from TF-IDF to their respective global parts, but it is not clear whether the logarithmic processing is beneficial to the performance of TC or not. To explore this problem, two improved methods named TF-ERF and ETF-RF are proposed by strengthen the RF part and TF part separately. As a result, the method TF-ERF is proved to be more helpful to the improvement of TC and it shows certain advantages over other term weighting schemes.

The rest of this paper is arranged as follows. Section 2 points out problems existing in the related work. Section 3 proposes improved methods based on TF-RF. Section 4 introduces the experimental settings. Section 5 analyzes the experimental data in detail. Section 6 conclusions this paper.

2 Analysis about current term weighting schemes

In this section, we conduct analysis about some current term weighting schemes. Meanwhile, the existing problems are also raised. For convenience, the notations utilized in this study are first presented in Table 1 and seven existing representative term weighting schemes are shown in Table 2. By the way, the mathematical forms of these term weighting schemes are all not normalized.

Table 1 Notations and descriptions

Notation	Description
A_{ij}	Number of documents that contain feature t_i in category c_i .
B_{ij}	Number of documents that do not contain feature t_i in category c_i .
C_{ij}	Number of documents that contain feature t _j but do not belong to
	category c _i .
D_{ij}	Number of documents that do not contain feature t_{j} and do not belong
	to category c _i .
N	Total number of documents in the training set, $N=A_{ij}+B_{ij}+C_{ij}+D_{ij}$.
$tf(t_j, d_k)$	The term frequency of feature t_j in document d_k .
$df(t_j)$	Number of documents containing t_{j} , and $df(t_{j}) = A_{ij} + C_{ij}$.

m	Total number of categories.
λ	λ is an adjustable coefficient and its default value is set to 7.0, and λ
	∈[5.0,9.0].
f_{ij}	Frequencies of t _j 's occurring in different categories, which are sorted
	in descending order with i (i=1,2,,m) being the rank.
$n_{c_i}(t_j)$	Number of documents containing the term t_j in a certain category c_i .
N_{c_i}	Total number of documents in a certain category $c_{\rm i}$.
k_{tf}	Coefficient added to the TF part of TF-RF
k_{rf}	Coefficient added to the RF part of TF-RF

Table 2 Term weighting schemes to be compared

Scheme		Name	Mathematical form
Unsupervised weighting	term	TF-IDF	$tf(t_j, d_k) \times \log\left(\frac{N}{df(t_j)}\right)$
Supervised weighting	term	TF-CHI2	$tf(t_j, d_k) \times \frac{N \times (A_{ij} \times D_{ij} - B_{ij} \times C_{ij})^2}{(A_{ij} + B_{ij}) \times (C_{ij} + D_{ij}) \times (A_{ij} + C_{ij}) \times (B_{ij} + D_{ij})}$
		TF-MI	$tf(t_{j}, d_{k}) \times \log \left(\frac{N \times A_{ij}}{(A_{ij} + B_{ij}) \times (A_{ij} + C_{ij})} \right)$
		TF-OR	$tf(t_j, d_k) \times \log \left(\frac{A_{ij} \times D_{ij}}{B_{ij} \times C_{ij}} \right)$
		TF-RF	$tf(t_j, d_k) \times \log_2\left(2 + \frac{A_{ij}}{\max(1, C_{ij})}\right)$
		TF-IDF-ICSDF	$tf(t_{j}, d_{k}) \times \left(1 + \log\left(\frac{N}{df(t_{j})}\right)\right) \times \left(1 + \log\left(\frac{m}{\sum_{i=1}^{m} \left(n_{c_{i}}(t_{j}) / N_{c_{i}}\right)}\right)\right)$
		TF-IGM	$tf(t_j, d_k) \times \left(1 + \lambda \times \frac{f_{j1}}{\sum_{i=1}^{m} f_{ji} \times i}\right)$

2.1 Term frequency (TF)

TF not only plays the role as the local part of a term weighting scheme but also exists as an independent weighting scheme which have been mentioned in Section 1. TF supports the assumption that the more frequent a term appears in a specific document, the greater contribution it makes to the representation of this document. That is to say, the term is more important to this document [31]. This conclusion can be intuitively obtained, but is it really reasonable? We can imagine a scenario like this, there is a document in which a term (t₁) appears 20 times while another term (t₂) appears only once, can we directly come to the conclusion that t₁ is 20 times more important than t₂? The answer is obviously not. By the way, we have mentioned that TC is a task whose goal is to classify a variety of documents into corresponding categories according to their content. It is intuitive that terms which possess the ability to distinguish between different types of documents ought to be distributed a higher term weight. However, the TF weight only reflects the ability of a term in representing the document containing this term [32], this is against with the target of TC task. Therefore, we can conclude that the single TF can't make meaningful contribution as a term weighting scheme to serve the TC task.

2.2 Inverse document frequency (IDF)

As mentioned in previous subsection, the single TF can't support the requirement of text classification as a term weighting scheme, because the importance of a term is measured only in the dimension of term-document. TC task in only limited in a certain text without considering the relationship between all texts, which results in a

bad performance of classification. Therefore, in order to make up the defect of TF, the IDF part is introduced to balance the excessive influence of TF on the final term weight [6]. With regard to a specific term, its IDF weight can be defined as

$$IDF(t_{j}) = \log\left(\frac{N}{df(t_{j})}\right) \tag{1}$$

It can be seen from Eq. (1) that the smaller the value of $df(t_i)$, the larger the value of $IDF(t_i)$, which

holds the assumption that the less documents in which a term appears, the more significant the term will be. By introducing the IDF part to the TF part, the scheme TF-IDF was generated as a result. Intuitively, this new scheme is more appropriate as it takes both local part and global part into account, and the actual TC performance is also consisted with the intuition. However, there are still some defects existing in this scheme, a simple example can be given to explain this. Assuming that there are four kinds of texts in the training text sets and they all consist of 60 documents. Four terms with the same term frequency are selected and their document frequency can be represented as $\{60, 0, 0, 0\}$, $\{30, 30, 0, 0\}$, $\{20, 20, 20, 0\}$ and $\{15, 15, 15, 15\}$ separately. It can be easily seen that the ability to discriminate different classes is ranked as $t_1 > t_2 > t_3 > t_4$. However, the four terms are assigned the same IDF value, which is contrary to our intuition. Actually, although the global factor is taken into account in the TF-IDF scheme, its incomplete application of global weight leads to these extreme cases occasionally which make the scheme invalid. This can be attributed to the absence of available category information in the process of computing term weight, so it is essential to measure the importance of a term in the document-category dimension.

2.3 Relevance frequency (RF)

Due to the TC task is a supervised learning process aiming at classifying different types of documents, it is natural to utilize STW schemes to match the TC process. However, there are some defects in most existing STW schemes, which will make the schemes invalid in some extreme situations. Apparently, TC task looks like a multi-label classification problem, but in fact, what we need to do is just separate the chosen category from others instead of taking every unrelated category into account. Following this thought, TF-RF was proposed which simplifies the multi-label classification problem into multiple independent binary classification problems [14]. Specifically speaking, in the training text corpus, when computing the weight of a specific term, the category of document containing this term is tagged as the positive category and the other categories are uniformly classified as the negative category [14]. TF-RF supports the idea that the more concentrated the chosen term is in the positive category than in the negative category, the greater ability it possesses to select a correct category for the documents containing it. Besides, this scheme also holds the assumption that the importance of a term is only related to the documents containing it [14]. Based on the two thoughts mentioned above, the term weighting scheme TF-RF was proposed which can be presented with the form shown in Table 2. It can be noticed that the logarithmic processing is also adopted on the global part of TF-RF.

Considering that logarithmic processing has certain restrictions on its parameters, we need to limit the parameters with a certain range to prevent failure. In terms of TF-RF, firstly, in order to avoid the RF part being meaningless when the chosen term doesn't appear in positive documents, i.e. a = 0, the minimum of the independent variable is limited to 2. Meanwhile, the base is also set to 2 to match the independent variable. With

this processing, the RF value becomes 1 when a = 0, it is logical that the term weight depends on the TF weight entirely in this case. Secondly, the minimum of the denominator is limited to 1 for the purpose that preventing the RF value becomes infinite. It seems that the processed formula of TF-RF shown in Table 2 solves the problems existing perfectly. However, due to the TF part depends on the frequency a certain term appears in a document, it is bound to be affected by the length of the document. Therefore, the normalization is always performed to eliminate the influence of different text lengths when computing the term weight. As a result, the normalized TF-RF formula can be defined as.

$$TFRF(t_{j}, d_{k}, c_{i}) = \frac{tf(t_{j}, d_{k}) \cdot \log_{2}(a / \max(1, c) + 2)}{\sqrt{\sum_{j=1}^{n} \left(tf(t_{j}, d_{k}) \cdot \log_{2}(a / \max(1, c) + 2) \right)^{2}}}$$
(2)

We have analyzed that the introduction of IDF is to balance the excessive influence of the single TF on the term weight. For the same purpose, the global part RF is introduced to generate a more reasonable term weighting model and the logarithmic processing borrowed from IDF part is adopted on the RF part directly. Although TF-RF presents outstanding performance in TC task, it is not clear whether the logarithmic processing adopted on the RF part contribute to the good performance or not. Maybe there is possibility that the logarithmic processing restricts TF-RF to achieve better TC performance. Therefore, a problem can be naturally proposed, that is, "does the logarithmic processing adopted on IDF also apply to RF?" In terms of this problem, related analysis will be carried out in the next section.

3 Proposed method

In this section, aiming at solving the problem mentioned in the previous section, two propositions are proposed along with corresponding strategy. In addition, the reliability of the two propositions are also discussed at the end of the section.

3.1 Motivation

Intuitively, if terms are assigned the same term weight, there is no doubt that they possess equivalent importance. But is the actual situation really as simple as it seems? As we know, most term weighting schemes consist of two parts, local part and global part, and the term weight is the product of their respective weights of the two parts. It seems that the term weight gives a comprehensive consideration from the two parts. Conversely, we can interpret it from another perspective, that is, the respective characteristics of the two parts can't be fully shown because the size relation between local weight and global weight is neglected when multiplying them. Here an example can be given to make an explanation about this. Assuming that there are two terms named t_1 and t_2 , the TF weights of t_1 and t_2 are 1 and 100 separately and the RF weights are 100 and 1 separately, it is obvious that their term weights are the same. However, there is great difference between the two terms because the RF weight has absolute dominance over term weight for t_1 while t_2 is just the opposite, which leads to the influence of the weak part on the term weight is suppressed by the strong part.

In terms of the TF-RF scheme we are studying, it has a relatively good performance for TC task, but we don't know whether the logarithmic processing adopted on the RF part borrowed from TF-IDF contributes to the excellent performance or not. Besides, there is even possibility that the influence of TF part and RF part on the term weight is out of balance due to the logarithmic processing, resulting that the characteristic of one part is concealed by the other part. As a result, TF-RF can't measure term weight from the two parts appropriately. Based

on this problem, two propositions can be proposed as following.

Proposition 1: the RF part is weakened too much through logarithmic processing. In other words, the impact of the TF part is overemphasized, which results in the RF part is suppressed by the TF part in contributing the discriminating power to the selected term.

Proposition 2: the weakening of RF part is not enough only by logarithmic processing. The RF part still occupies excessive dominance on the term weight, which leads to suppressing the effect of the TF part in contributing the discriminating power to the selected term.

The two propositions mentioned above can be presented as Fig. 1. It can be seen from the figure that some terms with different characteristics are assigned the same weight for initial TF-RF. Among these terms, some are bias to tf and others prefer rf. For proposition 1, due to the RF part has a greater dominance on the term weight which will restrain the expression of the characteristic of the TF part, so the dominance of TF part should be strengthened to balance the overemphasis on the RF part in order to make the scheme more reasonable. For proposition 2, it is just opposite to proposition 1. Therefore, the impact of TF part should be enhanced to balance the excessive influence caused by the RF part. Naturally, the result is also contrary to the result of proposition 1. In view of the two propositions, their respective strategy will be analyzed in detail in the next section.

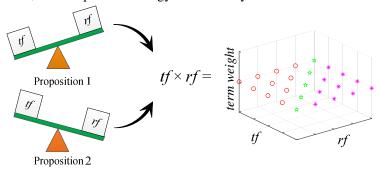


Fig. 1 Two propositions based on TF-RF

3.2 Improvement approaches

In view of the two propositions mentioned above, in order to gain a more reasonable scheme, the part which is suppressed by the other part for the dominance on the term weight, is ought to be strengthened. And there are two strategies shown in Fig. 2 which can be used to achieve this purpose.

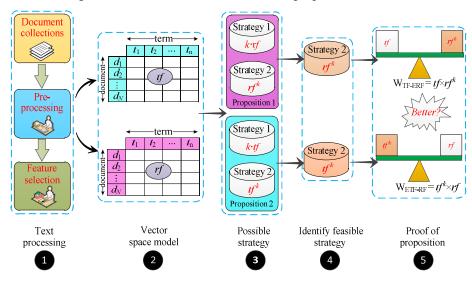


Fig. 2 Framework of the proposed method

Strategy 1: Adding a coefficient k as multiple to the certain part ($k \cdot \text{rf}$ or $k \cdot \text{rf}$, k > 1).

Strategy 2: Adding a coefficient k as power exponent to the certain part (tf^k or rf^k , k>1).

For proposition 1, adopting strategy 1 to the normalized TF-RF shown as Eq. (2). Then the deformed formula can be defined as Eq. (3), in which $k = k_{\rm rf}$.

$$TFRF(t_{j}, d_{k}, c_{i})^{*} = \frac{tf(t_{j}, d_{k}) \cdot k_{rf} \cdot \log_{2}(a / \max(1, c) + 2)}{\sqrt{\sum_{j=1}^{n} \left(tf(t_{j}, d_{k}) \cdot k_{rf} \cdot \log_{2}(a / \max(1, c) + 2) \right)^{2}}}$$

$$= \frac{tf(t_{j}, d_{k}) \cdot \log_{2}(a / \max(1, c) + 2)}{\sqrt{\sum_{j=1}^{n} \left(tf(t_{j}, d_{k}) \cdot \log_{2}(a / \max(1, c) + 2) \right)^{2}}} = TFRF(t_{j}, d_{k}, c_{i})$$
(3)

It can be noticed from Eq. (3) that the coefficient k_{rf} added to the RF part will be offset, so strategy 1 becomes meaningless. Then adopting strategy 2 to Eq. (2), the deformed formula can be defined as Eq. (4). Obviously, the failure will not happen in strategy 2. So after the identification, strategy 2 is chosen as a feasible to strengthen the RF part, and the improved method Eq. (4) is named as term frequency-exponential relevance frequency (TF-ERF).

$$TFERF(t_{j}, d_{k}, c_{i}) = \frac{tf(t_{j}, d_{k}) \cdot \left(\log_{2}(a / \max(1, c) + 2)\right)^{k_{rf}}}{\sqrt{\sum_{j=1}^{n} \left(tf(t_{j}, d_{k}) \cdot \left(\log_{2}(a / \max(1, c) + 2)\right)^{k_{rf}}\right)^{2}}}$$
(4)

For proposition 2, the same conclusion can be drawn as proposition 1. Strategy 1 will also become invalid due to the coefficient k_{tf} added to the TF part is offset. As a result, strategy 2 is selected to strengthen the TF part. Adopting approach 2 to Eq. (2), then the deformed formula can be defined as Eq. (5), which is named as exponential term frequency-relevance frequency (ETF-RF), in which $k = k_{tf}$.

$$ETFRF(t_{j}, d_{k}, c_{i}) = \frac{tf(t_{j}, d_{k})^{k_{ij}} \cdot \log_{2}(a / \max(1, c) + 2)}{\sqrt{\sum_{j=1}^{n} \left(tf(t_{j}, d_{k})^{k_{ij}} \cdot \log_{2}(a / \max(1, c) + 2)\right)^{2}}}$$
(5)

3.3 Qualitative analysis of two improved methods

In terms of the methods proposed in previous section, their actual effect performance can be presented as Fig. 3. Fig. 3 (a) and Fig. 3 (b) present the result for different values of tf and rf by adding a coefficient $k_{rf}(k_{rf}=1,2,3)$ or $k_{tf}(k_{rf}=1,2,3)$ to the corresponding RF or TF part separately. It can be seen from the figure that the slope of the surface shown in Fig. 3 (b) is steeper than that of Fig. 3 (a) for the same k_{tf} and k_{rf} . In addition, we can also notice that there is little difference between the partial derivative in direction tf and that of rf for Fig. 3 (a) though the coefficient k_{rf} is introduced to strengthen the RF part. While in Fig. 3 (b), the partial derivative in direction tf far exceeds that in direction rf due to the introduction of k_{tf} . This proves that the introduction of k_{tf} has greater influence on the term weight than k_{rf} . Due to the initial TF-RF is an excellent term weighting scheme which presents a good performance in TC task, the respective influence of the TF part and RF part on the term weight is reasonable to a certain extent. Therefore, the treatment to strengthen the influence of TF part on the term weight by adding a coefficient k_{tf} may greatly undermine this rationality and make the classification performance reduced.

By the way, similar analysis and conclusion also mentioned in [14], which is consistent with what we have analyzed. So if the initial TF-RF can be improved, there is a great possibility that the dominance of the RF part on the term weight ought to be enhanced, that is to say, TF-ERF may be more helpful to the improvement of TF-RF compared with ETF-RF.

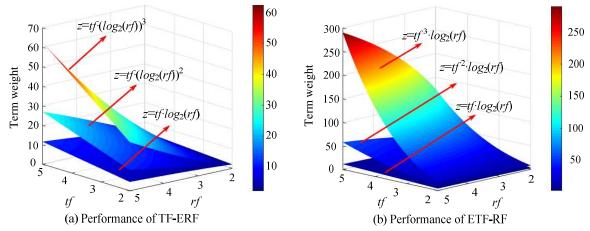


Fig. 3 Performance adding different k_{tf} and k_{rf} to TF part and RF part separately

4 Experimental setup

In this section, the experimental datasets, feature selection method, classification algorithms and evaluation of the performance measures used in our experiments are introduced successively.

4.1 Experimental datasets and pre-processing

In order to verify whether the improved methods proposed before is helpful to improve the performance of TC task, a series of experiments are conducted. Experiments are conducted on two datasets, Reuters-21587 corpus and WebKB corpus. Both the two datasets can be downloaded directly from the Internet.¹

4.1.1 Reuters-21587 corpus

The first dataset used in our experiments is the Reuters-21587 corpus, which is widely used in the text classification field. This English corpus contains 90 classes of news documents. The top 8 largest classes are selected for the text classification task. The reduced corpus contains 7674 documents which have been divided into a training set with 5485 documents and a test set with 2189 documents (shown in Table 3).

The pre-processing is an important step which will make an effect on the result of text classification to a certain degree. In this step, punctuation marks, numbers and other symbols are all removed. Furthermore, in order to reduce the size of the feature set, the terms which appear less than two times are discarded. At last, all letters are converted to lowercase and words are stemmed using Porter's stemmer [33].

Finally, a total of 8541 distinct terms left build the feature set. After the pre-processing stage, we acquire the document-term matrices of training set and test set, which are 5485×8541 and 2189×8541 respectively.

Table 3 Data Description on Reuters-21578

No.	Class label	Training data	Test data
1	acq	1596	696
2	crude	253	121
3	earn	2840	1083
4	grain	41	10
5	interest	190	81
6	money-fx	206	87

 $^{^1\,}https:\!/\!/ana.cachopo.org/datasets-for-single-label-text-categorization.$

7	ship	108	36
8	trade	251	75

4.1.2 WebKB corpus

The second dataset used in our experiments is the WebKB corpus, this English corpus contains 4199 documents which have been divided into a training set with 2803 documents and a test set with 1396 documents (shown in Table 4). In the step of pre-processing, the same pre-processing mentioned previous is also carried out on the WebKB corpus. Eventually, a total of 7061 distinct terms (features) are selected to build the feature set, the document-term matrices of training set and test set are 2803×7061 and 1396×7061 respectively.

Table 4 Data Description on WebKB

No.	Class label	Training data	Test data
1	project	336	168
2	course	620	310
3	faculty	750	374
4	student	1097	544

4.2 Feature selection

The initial feature set is achieved after the pre-processing of the datasets. However, the feature set can't be directly applied to experiments due to the fact that masses of features are meaningless for TC task. In addition, these invalid features may also cause harmful effect to the classifier, which will lead to a bad TC performance [34,35]. Therefore, it is essential to select the most effective features on the promise of not sacrificing the performance of TC task. As we know, feature selection (FS) is such a task aiming at building a more reasonable model for TC task by selecting relatively valuable terms. Presently, there are many methods that can be applied to feature selection. For example, three typical methods have been mentioned in section 1 named CHI2, IG and MI separately. Among them, CHI2 is recognized to be the most effective method due to its outstanding performance [36,37,38], so it is adopted to select features in this study.

4.3 Classification algorithms

Bayes classifier is the general name of a type of classification algorithm which are all based on Bayes law. Among them, Naive Bayes (NB) is the most common and simple one which is widely used in the field of TC [39]. NB treats all features as independent and no interaction and it regards the TC task as a probability problem that the class with the highest probability will be selected as the final category [40]. Based on this idea, the number of parameters to be estimated is greatly reduced, which simplifies the requirements of feature space and the calculation of solution to a great extent. As a result, the simplicity and efficiency will be greatly promoted when using NB classifier. In view of these advantages, NB classifier is utilized in our experiments and its algorithm interpretation is presented as Eq. (6).

Assuming that the document d_k consisting of a certain number of terms which can be represented as $d_k = (t_1, t_2, ..., t_n)$. The probability that the document d_k belongs to category c_i can be defined as

$$P(c_{i}|\mathbf{d}_{k}) = \frac{P(c_{i}) \prod_{j=1}^{n} P(t_{jk}|\mathbf{c}_{i})}{P(d_{k})}$$
(6)

4.4 Evaluation of the performance

In order to verify the effectiveness of the improved methods, the TC performance need to be evaluated with a certain standard of measurement. Among masses of evaluation index, precision and recall are two popular

measures for evaluating the performance of TC task. Precision denotes the proportion of correct assignments among all the test documents that should be assigned to the target category and recall represents the proportion of correct assignments among all the test documents assigned to the target category. However, neither of them can be directly used to evaluating the performance for the reason that the higher level of one indicator may be obtained at the expense of sacrificing the level of the other one. As a result, a new measure named F1 was proposed, in which precision and recall are combined together and assigned the same importance. The precision, recall and the F1 measure can be defined as follows and the explanations of corresponding notations in the formula are listed in the Table 5.

$$P(c_i) = \frac{TP(c_i)}{TP(c_i) + FP(c_i)}$$
(7)

$$R(c_i) = \frac{TP(c_i)}{TP(c_i) + FN(c_i)}$$
(8)

$$F(c_i) = \frac{2 \cdot P(c_i) \cdot R(c_i)}{P(c_i) + R(c_i)} = \frac{2 \cdot TP(c_i)}{2 \cdot TP(c_i) + FP(c_i) + FN(c_i)}$$
(9)

Table 5 Contingency table for category c_i .

	True label c _i	True not c _i
Predicted label ci	True Positive (TP)	False Positive (FP)
Predicted not ci	False Negative (FN)	True Negative (TN)

In general, the F1 measure is estimated from two ways, micro-averaged F1 (micro-F1) and macro-average F1 (macro-F1). The macro-F1 and micro-F1 can be defined as Eq.(10) and Eq.(11).

$$macro-F1 = \frac{1}{m} \sum_{i=1}^{m} F_1(c_i)$$
 (10)

$$micro-F1 = \frac{2 \cdot \sum_{i=1}^{m} TP(c_i)}{2 \cdot \sum_{i=1}^{m} TP(c_i) + \sum_{i=1}^{m} FP(c_i) + \sum_{i=1}^{m} FN(c_i)}$$
(11)

5 Experiment results and analysis

In this section, orthogonal experimental design will be firstly described, by which one of ETF-RF and TF-ERF will be proved to be more helpful to the improvement of the initial TF-RF. Then, the chosen method (i.e., TF-ERF or ETF-RF) will be compared with other existing term weighting schemes (listed in Table 2) in order to verify its effectiveness for the improvement of TC performance.

5.1 Performance comparisons between TF-ERF and ETF-RF

The first group of experiments is to distinguish which one of TF-ERF and ETF-RF is more helpful to the improvement of the initial TF-RF. Orthogonal experimental design is a kind of experimental design method to study multi-factors and multi-levels. According to orthogonality, some representative points are selected from the comprehensive test to carry out the test. The main tool of orthogonal test design is the orthogonal table, testers can seek the corresponding orthogonal table according to the requirements of the number of factors, the level of factors and whether there is interaction, and then select some representative points from the comprehensive test to test based on the orthogonality of the orthogonal table. In this study, the orthogonal table $L_{25}(5^6)$ is selected to arrange the combinations of parameters at different level and the five levels are set to 1,2,3,4,5 separately. The macro-F1 and micro-F1 are chosen to evaluate the performance. The results in different values of $k_{\rm tf}$ and $k_{\rm rf}$ are

plotted in the Fig. 4.

The figure presents the performance obtained on Reuters-21578 dataset and WebKB dataset separately. It can be seen from the figure that the classification performance will deteriorate rapidly with the increase of k_{tf} . In contrast, the change of k_{rf} doesn't have great influence on the performance of TC as k_{tf} shows, while the increase of k_{rf} presents a positive impact on the improvement of the classification performance. This proves that the improved method TF-ERF is beneficial to the improvement of initial TF-RF scheme, which is consistent with the analysis result mentioned in Section 3.3.

In terms of Fig. 4, taking both macro-F1 and micro-F1 into account, the best classification performances can be observed under $k_{tf} = 1$ and $k_{rf} = 5$ for the Reuters-21578 dataset as well as $k_{tf} = 1$ and $k_{rf} = 2$ for the WebKB dataset. Therefore, the parameters added to TF-ERF is determined, and the improved method TF-ERF will be compared with other term weighting schemes in the following experiments.

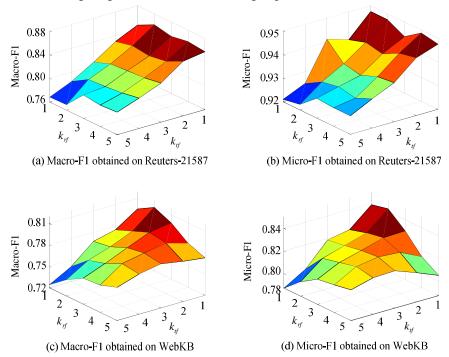


Fig. 4 Classification performance according to different parameters k_{tf} and r_{rf} under two corpus

5.2 Performance comparisons of existing methods

The second group of experiments is to verify the effectiveness of TF-ERF by comparing its performance with other term weighting schemes listed in Table 2. Fig. 5 shows the experiment performance of TC on the Reuters-21578 corpus. It can be seen from the figure that TF-IDF and TF-IDF-ICSDF present the worst performance in all feature sets for both macro-F1 and micro-F1. Especially at a small feature sets (less than 200), the performance of the two schemes is far worse than the other ones. Meanwhile, the rest schemes perform well even when the number of features is small.

In terms of macro-F1, almost all term weighting schemes reach their peaks at a feature set around 1000 and TF-ERF obtain the best performance at the peak compared with others. On the whole, TF-ERF is superior to other schemes when the number of features is less than 7000. In addition, it can be seen that TF-RF does not show advantages over other schemes and even inferior to TF-CHI2 and TF-MI in most feature sets. By contrast, it is obvious that TF-ERF is very effective in improving the performance of TC as an improved method of TF-RF.

In terms of micro-F1, these schemes don't reach the peak at the same feature set as the figure of macro-F1 shows. But there are also corresponding turning points at a feature set around 1000. After that, the micro-F1 of some schemes begin to decrease like TF-IDF and TF-IDF-ICSDF, while the rest schemes continue to maintain an increasing trend at a slower speed. It can be seen that TF-ERF also shows evident advantage over other schemes for most feature sets in addition to the situation when the number of features is 800. TF-ERF reaches its peak when the number of features is around 5000 and the peak presents the best performance of all the term weighting schemes for all the features.

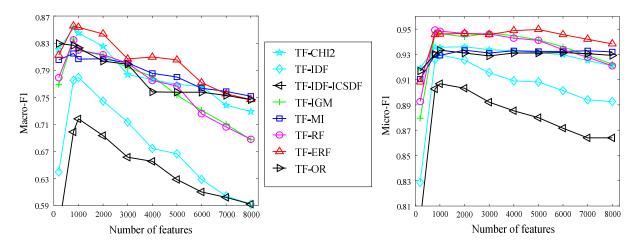


Fig. 5 Macro-F1 and Micro-F1 measure of the eight weighting approaches on the Reuters-21578 corpus with different numbers of features

Fig. 6 shows the experimental performance of TC on the WebKB corpus. It can be seen from the figure that most schemes show poor performance when the number of features is relatively small, which is different from what Fig. 5 shows. In addition, TF-ERF does not present great advantage over other schemes as it shows in Fig. 5. However, it can't be neglected that TF-ERF has a better performance in a certain feature sets and the improvement relative to the initial TF-RF.

In terms of macro-F1, it can be seen that TF-CHI2, TF-MI and TF-OR almost maintain the growing trend with the increase of the number of features, and other schemes start to decrease when reaching their respective peaks. As far as TF-ERF is concerned, it is superior to other schemes when the number of features falls in [200, 3000]. After the number of features exceeds 3000, the performance of TF-ERF begin to deteriorate and be surpassed by TF-MI, TF-OR and TF-CHI2 successively. However, we should notice that TF-ERF reaches its peak at a feature set around 2000 and the performance of this point is the best among all the schemes. In addition, as an improved method of the initial TF-RF, it also shows a better performance than TF-RF for the whole feature sets.

In terms of micro-F1, the curves become much tighter compared with the left figure which means that the gap of performance between different schemes has become smaller. From the figure we can see that TF-IGM, TF-RF and TF-ERF obtain better performance than the rest schemes for almost the whole feature sets. There is almost no difference in the performance of these three schemes when the number of features is less than 1500. After that they all present a descent trend one after another. Among them, due to the attenuation of TF-ERF is less than the other two schemes, so the performance of TF-ERF is superior to TF-RF and TF-IGM. TF-ERF reaches its peak at a feature set around 3000 and it also represents the best performance of all these eight schemes for all feature sets.

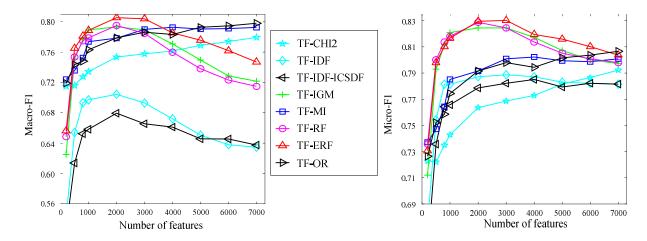


Fig. 6 Macro-F1 and Micro-F1 measure of the five weighting approaches on the WebKB corpus with different numbers of features

5.3 Performance improvements of proposed method over others

1000

0.9461

(-0.24%)

In order to understand the performance of the proposed scheme compared with other existing schemes more accurately, the specific experimental data of each scheme for Reuters-21578 dataset and WebKB dataset are listed in the Table 6 and Table 7 separately. The data in parentheses is the improvement of TF-ERF relative to other schemes. Among them, the data with the best performance in each chosen feature set is presented with boldface. And the data with the best performance in the whole feature sets is underlined extra.

Table 6 Macro-F1 and micro-F1 for Reuters-21578 dataset Classifier, TF-IDF-TF-ERF TF-CHI2 TF-IDF TF-IGM TF-OR TF-RF TF-MI -ICSDF #Features Macro-F1 0.5550 0.7790 0.8206 0.6398 0.76860.8053 0.8295 200 0.8123 (4.28%)(46.36%) (5.69%)(0.87%)(-2.07%) (-1.01%)(6.95%)0.8524 0.7752 0.8267 0.8348 0.6989 0.8361 0.8148 800 0.8555 (2.48%)(0.36%)(10.36%)(22.41%)(2.32%)(5.00%)(3.48%) 0.8197 0.8450 0.7793 0.7182 0.8190 0.8064 0.8231 0.8535 1000 (4.12%)(1.00%)(9.52%) (18.84%)(4.21%)(5.83%)(3.69%) 0.8132 0.8253 0.7444 0.6934 0.8122 0.8068 0.8034 2000 0.8437 (3.75%)(2.23%)(13.34%) (21.68%)(3.88%)(4.57%)(5.02%) 0.7835 0.80040.7136 0.6615 0.7936 0.8009 0.7990 3000 0.8066 (0.78%)(2.95%)(13.04%)(21.94%)(1.64%)(0.72%)(0.95%) 0.77480.77610.6744 0.6554 0.78120.7855 0.7581 4000 0.8095 (20.03%)(4.47%)(4.30%)(23.51%)(3.62%)(3.05%)(6.77%) 0.7668 0.7690 0.6667 0.6288 0.7533 0.7800 0.7575 5000 0.8053 (20.79%)(28.07%) (5.02%)(4.71%)(6.90%)(3.25%)(6.31%) 0.7256 0.7663 0.6287 0.6102 0.7309 0.7635 0.7576 6000 0.7717 (0.71%) (22.74%)(6.36%)(26.47%)(5.58%)(1.07%)(1.86%) 0.7063 0.7385 0.6042 0.6022 0.7104 0.7588 0.7528 7000 0.7551 (6.92%)(2.26%)(24.97%) (25.40%)(6.30%) (-0.48%) (0.31%) 0.6881 0.7290 0.5920 0.5922 0.7464 0.6877 0.7517 8000 0.7467 (2.43%)(26.09%) (0.04%)(8.52%)(26.14%)(8.58%)(-0.67%)Micro-F1 0.9191 0.9169 0.89260.8287 0.79760.8794 0.9105 0.9082200 (1.75%)(-1.19%)(9.59%)13.86%) (3.27%)(-0.26%)(-0.95%) 0.9470 0.9493 0.9356 0.9260 0.9027 0.9296 0.9301 800 0.9456 (2.12%) (-0.39%)(1.07%)(1.67%) (4.76%)(-0.14%)(1.73%)0.9484 0.9356 0.9296 0.9068 0.9292 0.9333 0.9466

(4.33%)

(1.12%)

(1.82%)

(1.37%)

(-0.05%)

2000	0.9466	0.9465	0.936	0.9255	0.9032	0.9438	0.9333	0.9310
2000	0.9400	(0.01%)	(1.13%)	(2.27%)	(4.80%)	(0.29%)	(1.42%)	(1.67%)
2000	0.9456	0.9461	0.9338	0.9155	0.8922	0.9470	0.9310	0.9287
3000	0.9436	(-0.05%)	(1.27%)	(3.29%)	(5.99%)	(-0.14%)	(1.57%)	(1.82%)
4000	0.9488	0.9429	0.9328	0.9091	0.8853	0.9447	0.9328	0.9310
4000	0.9488	(0.63%)	(1.72%)	(4.37%)	(7.18%)	(0.44%)	(1.72%)	(1.92%)
5000	0.9497	0.9411	0.9315	0.9082	0.8799	0.9415	0.9324	0.9310
3000	<u>0.9497</u>	(0.92%)	(1.96%)	(4.57%)	(7.94%)	(0.88%)	(1.86%)	(2.01%)
6000	0.9456	0.9338	0.9296	0.9013	0.8716	0.9360	0.9324	0.9315
0000	0.9430	(1.27%)	(1.73%)	(4.92%)	(8.49%)	(1.03%)	(1.42%)	(1.52%)
7000	0.0420	0.9278	0.9251	0.8940	0.8639	0.9292	0.9328	0.9306
/000	0.9420	(1.53%)	(1.82%)	(5.37%)	(9.04%)	(1.38%)	(0.98%)	(0.22%)
9000	0.0202	0.9210	0.921	0.8926	0.8639	0.9223	0.9319	0.9292
8000	0.9383	(1.88%)	(1.88%)	(5.12%)	(8.62%)	(1.74%)	(0.69%)	(0.98%)

Table 7 Macro-F1 and micro-F1 for WebKB dataset

Classifier, #Features	TF-ERF	TF-RF	TF-CHI2	TF-IDF	TF-IDF- -ICSDF	TF-IGM	TF-MI	TF-OR
Macro-F1								
200	0.6560	0.6485	0.7149	0.5433	0.5259	0.6250	0.7240	0.7190
200	0.6562	(1.18%)	(-8.21%)	(20.78%)	(24.78%)	(4.99%)	(-9.37%)	(-8.74%)
500	0.7651	0.7535	0.7161	0.654	0.6138	0.7511	0.7357	0.7437
500	0.7651	(1.54%)	(6.85%)	(16.99%)	(24.66%)	(1.87%)	(4.00%)	(2.88%)
800	0.7813	0.7751	0.7276	0.6933	0.6524	0.7768	0.7526	0.7481
800	0./013	(0.80%)	(7.38%)	(12.69%)	(19.75%)	(0.58%)	(3.81%)	(4.44%)
1000	0.7887	0.7788	0.7343	0.6971	0.6584	0.7893	0.7739	0.7630
1000	0.7007	(1.28%)	(7.41%)	(13.14%)	(19.80%)	(-0.07%)	(1.92%)	(3.37%)
2000	0.8055	0.7952	0.7535	0.7044	0.6794	0.7934	0.7784	0.7789
2000	<u>0.8055</u>	(1.30%)	(6.91%)	(14.36%)	(18.57%)	(1.53%)	(3.49%)	(3.42%)
3000	0.8040	0.7847	0.7578	0.6929	0.6654	0.7894	0.7900	0.7864\
3000	0.8040	(2.47%)	(6.10%)	(16.04%)	(20.83%)	(1.86%)	(1.78%)	(2.24%)
4000	0.7854	0.7599	0.7615	0.6721	0.6611	0.7708	0.7927	0.7833
4000	0.7834	(3.36%)	(3.140%)	(16.86%)	(18.80%)	(1.90%)	(-0.92%)	(0.27%)
5000	0.7758	0.7382	0.7690	0.6506	0.6456	0.7495	0.7904	0.7927
3000	0.7758	(5.09%)	(0.88%)	(19.24%)	(20.16%)	(3.51%)	(-1.85%)	(-2.13%)
6000	0.7619	0.7232	0.7741	0.6378	0.6454	0.7286	0.7912	0.7947
0000		(5.35%)	(-1.58%)	(19.46%)	(18.06%)	(4.57%)	(-3.70%)	(-4.13%)
7000	0.7471	0.7148	0.7795	0.6346	0.6374	0.7216	0.7926	0.7981
7000	0.7471	(4.52%)	(-4.16%)	(17.72%)	(17.21%)	(3.53%)	(-5.74%)	(-6.39%)
Micro-F1								
200	0.7307	0.7357	0.7228	0.6748	0.6562	0.7120	0.7371	0.7264
200		(-0.69%)	(1.09%)	(8.28%)	(11.35%)	(2.62%)	(-0.87%)	(0.59%)
500	0.7980	0.8001	0.7221	0.7543	0.7357	0.7930	0.7471	0.7521
300		(-0.26%)	(10.51%)	(5.79%)	(8.47%)	(0.63%)	(6.81%)	(6.10%)
800	0.8102	0.8138	0.735	0.7815	0.7643	0.8116	0.7643	0.7586
	0.0102	(-0.45%)	(10.23%)	(3.67%)	(6.00%)	(-0.18%)	(6.00%)	(6.80%)
1000	0.8166	0.8181	0.7428	0.7815	0.7658	0.8209	0.7851	0.7744
1000	0.0100	(-0.18%)	(9.94%)	(4.49%)	(6.64%)	(-0.52%)	(4.01%)	(5.45%)
2000	0.8295	0.8288	0.7636	0.7872	0.7787	0.8245	0.7915	0.7915
2000	J.0273	(0.09%)	(8.63%)	(5.38%)	(6.53%)	(0.61%)	(4.80%)	(4.80%)
3000	0.8302	0.8245	0.7686	0.7887	0.7822	0.8245	0.8009	0.798
	0.0302	(0.69%)	(8.02%)	(5.27%)	(6.14%)	(0.69%)	(3.66%)	(4.04%)
4000	0.8195	0.8138	0.7729	0.7872	0.7851	0.8173	0.8023	0.7944
7000	0.0175	(0.70%)	(6.03%)	(4.10%)	(4.38%)	(0.27%)	(2.14%)	(3.16%)
5000	0.8159	0.8052	0.7815	0.7830	0.7794	0.8073	0.7994	0.8016
5000	0.8159	(1.33%)	(4.40%)	(4.20%)	(4.68%)	(1.07%)	(2.06%)	(1.78%)
6000	0.8102	0.8009	0.7865	0.7822	0.7822	0.8001	0.7987	0.8037

		(1.16%)	(3.01%)	(3.58%)	(3.58%)	(1.26%)	(1.44%)	(0.81%)
7000	0.8037	0.7980	0.7923	0.7815	0.7815	0.7973	0.8009	0.8066
		(0.72%)	(1.44%)	(2.84%)	(2.84%)	(0.81%)	(0.35%)	(-0.36%)

It can be seen from Table 6 that most of the best performances of selected feature set focus on TF-ERF. For both macro-F1 and micro-F1, the best performance of TF-ERF is also the best performance of all the eight term weighting schemes for the whole feature sets. As an improvement method of TF-RF, the performance of TF-ERF is superior to that of TF-RF for each feature set. On the whole, the improvement growth relative to TF-RF raises with the increase of the number of features except for some special points. For Table 7, although TF-ERF doesn't show strong advantage over other schemes as Table 6 shows, in most cases, TF-ERF is still superior than other schemes. Besides, what has not changed is that the best performance of TF-ERF still represents the optimal level of all the selected term weighting schemes for the whole feature sets.

6 Conclusion

In this study, considering that the logarithmic processing of the global weight borrowed from TF-IDF may not adapt to other term weighting schemes, e.g., TF-RF. So two improved methods based on TF-RF were presented to explore this problem. First of all, two propositions were proposed to explain the problems faced by TF-RF. In terms of the two propositions, feasible improvement strategies were identified and then applied to them. As a result, two methods named TF-ERF and ETF-RF were proposed. Then, through orthogonal experimental design, the improved method TF-ERF, which holds the proposition that the TF part suppresses the impact of the RF part in contributing the discriminating power to the selected term, was proved to be more helpful to the improvement of TC than the other one,. Meanwhile, the parameters added to TF-ERF is also determined. After that, TF-ERF was compared with seven existing representative schemes to evaluate the text classification performance, the experimental results proved that TF-ERF overperforms the listed term weighting schemes including TF-RF.

7 Declarations

Funding

This work is supported by the National Natural Science Foundation of China (No. 52075350), the Sichuan Major Science and Technology Project (No.2019ZDZX0001).

Conflicts of interest

All authors declare that they have no conflict of interest.

Availability of data and material

The sources of relevant data have been described in this manuscript.

Code availability

The codes in this manuscript are programmed by Python. Part of the code have been listed in the appendix at the end of the manuscript.

Authors' contributions

The problems studied in this manuscript and the corresponding contributions are stated in the summary and conclusion respectively

Ethics approval

This article does not contain any studies with human participants or animals performed by any of the authors.

Consent to participate

All the authors agreed to participate.

Consent for publication

All the authors agreed to publish.

References

- [1] Sebastiani F. Machine learning in automated text categorization [J]. ACM Computing Surveys, 2002, 34(1): 1–47
- [2] Al-Mubaid H, Umair S A. A new text categorization technique using distributional clustering and learning logic [J]. IEEE Transactions on Knowledge and Data Engineering, 2006, 18(9): 1156-1165.
- [3] F Debole, F Sebastiani. Supervised term weighting for automated text categorization [J]. Proceedings of the 2003 ACM Symposium on Applied Computing, 2003, pp. 784–788.
- [4] Li W, Miao D, Wang W. Two-level hierarchical combination method for text classification [J]. Expert Systems with Applications, 2011, 38(3): 2030-2039.
- [5] Labani M, Moradi P, Ahmadizar F, et al. A novel multivariate filter method for feature selection in text classification problems [J]. Engineering Applications of Artificial Intelligence, 2018, 70: 25-37.
- [6] Shang C X, Li M, Feng S Z, et al. Feature selection via maximizing global information gain for text classification [J]. Knowledge-Based Systems, 2013, 54: 298-309.
- [7] Tellez E S, Moctezuma D, Miranda-Jiménez S, et al. An automated text categorization framework based on hyperparameter optimization [J]. Knowledge-Based Systems, 2018, 149: 110-123.
- [8] Shang W, Huang H, Zhu H, et al. A novel feature selection algorithm for text categorization [J]. Expert Systems with Applications, 2007, 33(1): 1-5.
- [9] Zhang W, Yoshida T, Tang X. A comparative study of TF*IDF, LSI and multi-words for text classification [J]. Expert Systems with Applications, 2011, 38(3): 2758-2765.
- [10] Haddoud M, Mokhtari A, Lecroq T, et al. Combining supervised term-weighting metrics for SVM text classification with extended term representation [J]. Knowledge and Information Systems, 2016, 49(3): 909-931.
- [11] Chaoliang Li, Shigang Liu. A comparative study of the class imbalance problem in Twitter spam detection[J]. Concurrency and Computation: Practice and Experience, 2018.
- [12] Li W, Li Y, Chen J, et al. Product functional information based automatic patent classification: Method and experimental studies [J]. Information Systems, 2017, 67: 71-82.
- [13] Liu L, Li Y, Xiong Y, et al. A new function-based patent knowledge retrieval tool for conceptual design of innovative products [J]. Computers in Industry, 2020, 115:103154.
- [14] Lan M, Tan C L, Su J, et al. Supervised and traditional term weighting methods for automatic text categorization [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2009, 31(4): 721-735.
- [15] Salton G, Wong A, Yang C S. A vector space model for automatic indexing [J]. Communications of the ACM, 1974, 18(11): 613–620.
- [16] Zong W, Wu F, Chu, L K, et al. A discriminative and semantic feature selection method for text categorization [J]. International Journal of Production Economics, 2015, 165: 215-222.
- [17] Sabbah T, Selamat A, Selamat M H, et al. Modified frequency-based term weighting schemes for text classification [J]. Applied Soft Computing, 2017, 58: 193–206.
- [18] Wu H, Gu X, Gu Y. Balancing between over-weighting and under-weighting in supervised term weighting [J]. Information Processing and Management, 2017, 53(2): 547–557.
- [19] Guru D S, Suhil M, Raju L N, et al. An alternative framework for univariate filter based feature selection for

- text categorization [J]. Pattern Recognition Letters, 2018, 103: 23-31.
- [20] Wang D, Zhang H. Inverse-Category-Frequency Based Supervised Term Weighting Schemes for Text Categorization [J]. Journal of information science and engineering, 2013, 29(2): 209-225.
- [21] Ren F, Sohrab M G. Class-indexing-based term weighting for automatic text classification [J]. Information Sciences, 2013, 236: 109-125.
- [22] Chen K, Zhang Z, Long J, et al. Turning from TF-IDF to TF-IGM for term weighting in text classification [J]. Expert Systems with Applications, 2016, 66: 245-260.
- [23] Sabbah T, Selamat A, Selamat M H, et al. Hybridized term-weighting method for Dark Web classification [J]. Neurocomputing, 2016, 173: 1908–1926.
- [24] Li Z, Lu W, Sun Z, et al. A parallel feature selection method study for text classification [J]. Neural Comput and Applications, 2016, 28: 513-524.
- [25] Spärck Jones K. A statistical interpretation of term specificity and its application in retrieval [J]. Journal of Documentation, 2004, 60(5): 493–502.
- [26] Salton G, Buckley C. Term-weighting approaches in automatic text retrieval [J]. Information Processing and Management, 1988, 24(5): 513–523.
- [27] Altınçay H, Erenel Z. Analytical evaluation of term weighting schemes for text categorization [J]. Pattern recognition letters, 2010, 31(11): 1310-1323.
- [28] Liu Y, Loh H T, Sun A. Imbalanced text classification: A term weighting approach [J]. Expert Systems with Applications, 2009, 36(1): 690-701.
- [29] Quan X, Wenyin L, Qiu B. Term weighting schemes for question categorization [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2011, 33(5): 1009-1021.
- [30] Dogan T, Uysal A K. Improved inverse gravity moment term weighting for text classification [J]. Expert Systems with Application, 2019, 130: 45-59.
- [31] Lakshmi R, Baskar S. Novel term weighting schemes for document representation based on ranking of terms and Fuzzy logic with semantic relationship of terms [J]. Expert Systems with Applications, 2019, 137: 493-503.
- [32] Zhang W, Li Y, Wang S. Learning document representation via topic-enhanced LSTM model [J]. Knowledge Based Systems, 2019, 174: 194-204.
- [33] Porter, M.F. An algorithm for suffix stripping [J]. Program Electronic Library & Information Systems, 2006, 14(3): 130 137.
- [34] Meng J N, Lin H F, Yu Y H. A two-stage feature selection method for text categorization [J]. Computers and Mathematics with Applications, 2011, 62(7): 2793-2800.
- [35] Wang S, Pedrycz W, Zhu Q, et al. Subspace learning for unsupervised feature selection via matrix factorization [J]. Pattern Recognition, 2015, 48(1): 10–19.
- [36] Liu H, Yu L. Toward integrating feature selection algorithms for classification and clustering [J]. IEEE Transactions on Knowledge and Data Engineering, 2005, 17(4): 491-502.
- [37] Taşcı Ş, Güngör T. Comparison of text feature selection policies and using an adaptive framework [J]. Expert Systems with Application, 2013, 40(12): 4871-4886.
- [38] Şahin D Ö, Kılıç E. Two new feature selection metrics for text classification [J]. Automatika, 2019,

- 60(2):162–171.
- [39] Yang Y M, Pedersen J O. A comparative study on feature selection in text categorization [C]. In Proceedings of the 14th International Conference on Machine Learning, Nashville, USA, 1997, pp. 412-420.
- [40] Friedman N, Geiger D, Goldszmidt M. Bayesian network classifiers [J]. Machine Learning, 1997, 29(2-3): 131-163.

Appendix

```
import sklearn
import nltk
import json
import math
import numpy as np
from math import log
from sklearn import metrics
from sklearn import svm, datasets
from sklearn.datasets import load_files
from sklearn.metrics import classification_report
from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
from sklearn.feature_selection import SelectKBest,chi2
from sklearn.naive_bayes import MultinomialNB, BernoulliNB, GaussianNB
from openpyxl import Workbook
from openpyxl import load_workbook#引入库
from sklearn.preprocessing import Normalizer
print("====读取 r8-train-stemmed 训练集文本=====")
with open('r8-train-stemmed.txt') as file_object:
    r8_train = file_object.read()
print("====读取 r8-test-stemmed 测试集文本=====")
with open('r8-test-stemmed.txt') as file_object:
    r8_test = file_object.read()
print("====建立数据集类别名列表=====")
r8_target_names=['acq\t','crude\t','earn\t','grain\t','interest\t','money-fx\t','ship\t','trade\t']#建立数据集类别名列表,
其中\t 表示制表符,即 tab 键,用于增加类别名的区分度
#因为用 word 打开 r8 数据集可以发现,每个类别名后面都有一个 tab 键,因此可以用类别名加 tab 键的形
式(增加了区分度)来计算训练集合测试集中各类别的文档数。因为类别名区分度不强,可能在文中也会
出现,例如类别名 acq 可能作为一个单词在文中出现
print("====查看训练集中各类别下的文档数=====")
train name numbers=[]#建立一个空列表存储文档数
for r8_target_name in r8_target_names:#依次从类别名列表中读取类别名
    train_name_numbers.append(r8_train.count(r8_target_name))#根据类别名计算其文档数
    print(r8_target_name, r8_train.count(r8_target_name))
print("用于训练的文档数: ",sum(train_name_numbers))
print("====查看测试集中各类别下的文档数=====")
test_name_numbers=[]#建立一个空列表存储文档数
for r8_target_name in r8_target_names:
    test_name_numbers.append(r8_test.count(r8_target_name))
```

print(r8_target_name, r8_test.count(r8_target_name))

print("用于测试的文档数: ",sum(test_name_numbers))

print("训练和测试的文档总数: ",sum(train_name_numbers)+sum(test_name_numbers))

print("====将训练集中的类别名全部替换为英文句号"点",以方便后面分句=====")

for r8_target_name in r8_target_names:#用 for 循环依次替换 20 个不同的类别名

f1 = open("r8-train-stemmed_sentences.txt", "r")#注意这里用循环来替换,因此 f1 的打开与关闭都要放到 for 循环中,否则只能对最后一个类别名进行替换,因为那样后面分句后的结果为 378

content = f1.read()

f1.close()

 $t = content.replace(r8_target_name,".")$ #将类别名全部替换为英文句号"点",以方便后面分句。注意"点"后面有一个空格,否则分句后的结果为 0,即 print(len(sentences))为 0

with open("r8-train-stemmed_sentences.txt","w") as f2:

f2.write(t)

f2.close()

print("====将训练文本中的回车符去掉=====")#将文本中的回车符"\n"去掉,即用""来替换,否则打印一个句子(相当于一篇文档)时,都是分了很多个段落。但是不去掉回车符也不影响后面的词频计算

f1 = open("r8-train-stemmed_sentences.txt","r")

content = f1.read()

f1.close()

t= content.replace("\n"," ")#将文中的回车符"\n"去掉,即用" "来替换,否则打印一个句子(相当于一篇文档)时,都是分了很多个段落的。

#注意这里是替换为空格,两个引号之间有一个空格。如果两个引号之间没有空格,分句之后的结果为 11213,可能是因为回车符恰好在一行的结尾,导致上下两行链接的中间没有空格。

#如果两个引号之间没有空格,则可以将上述语句放到分句之后才去掉回车符,其结果仍为11293.

with open("r8-train-stemmed_sentences.txt","w") as f2:

f2.write(t)

f2.close()

•••••