

## Research Article

# Text Classification Using Novel Term Weighting Scheme-Based Improved TF-IDF for Internet Media Reports

Zhiying Jiang <sup>1,2</sup>, Bo Gao <sup>1,2</sup>, Yanlin He <sup>1,2</sup>, Yongming Han <sup>1,2</sup>, Paul Doyle <sup>3</sup>  
and Qunxiong Zhu <sup>1,2</sup>

<sup>1</sup>College of Information Science & Technology, Beijing University of Chemical Technology, Beijing 100029, China

<sup>2</sup>Engineering Research Center of Intelligent PSE, Ministry of Education of China, Beijing 100029, China

<sup>3</sup>School of Computer Science within the College of Science and Health, Technological University Dublin, Dublin, Ireland

Correspondence should be addressed to Qunxiong Zhu; [zhuqx@mail.buct.edu.cn](mailto:zhuqx@mail.buct.edu.cn)

Received 6 December 2020; Revised 18 January 2021; Accepted 29 January 2021; Published 5 March 2021

Academic Editor: Nianyin Zeng

Copyright © 2021 Zhiying Jiang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

With the rapid development of the internet technology, a large amount of internet text data can be obtained. The text classification (TC) technology plays a very important role in processing massive text data, but the accuracy of classification is directly affected by the performance of term weighting in TC. Due to the original design of information retrieval (IR), term frequency-inverse document frequency (TF-IDF) is not effective enough for TC, especially for processing text data with unbalanced distributions in internet media reports. Therefore, the variance between the DF value of a particular term and the average of all DFs ( $\overline{DF}$ ), namely, the document frequency variance (ADF), is proposed to enhance the ability in processing text data with unbalanced distribution. Then, the normal TF-IDF is modified by the proposed ADF for processing unbalanced text collection in four different ways, namely, TF-IADF, TF-IADF<sup>+</sup>, TF-IADF<sub>norm</sub>, and TF-IADF<sub>norm</sub><sup>+</sup>. As a result, an effective model can be established for the TC task of internet media reports. A series of simulations have been carried out to evaluate the performance of the proposed methods. Compared with TF-IDF on state-of-the-art classification algorithms, the effectiveness and feasibility of the proposed methods are confirmed by simulation results.

## 1. Introduction

Due to the rapid development of internet technology and information infrastructure construction, the volume of text data which can be obtained online has increased dramatically. As the China Internet Network Information Center (CNNIC) stated, the number of netizens in China had increased to 828.51 million by the end of 2018 [1]. The internet has become the main channel for Chinese people to obtain information. The content of internet media is the most important data source, in which textual documents are the main one. And it is increasingly important to effectively analyze massive textual documents such as classification, indexing, and clustering. As a consequence, text classification (TC) is concerned by many researchers working in the field. Based on the previous studies, many applications based

on TC technology have been developed, such as author identification [2, 3], spam e-mail filtering [4], medical documents' classification [5], management of customer relationship, and classification of web pages [6, 7]. Text classification (TC) is a task that assigns textual documents to predefined classes based on knowledge extracted from their content. The process of TC is as follows [8]:

- (i) Given a set of  $k$  different discrete class label values  $C = \{C_1, \dots, C_k\}$  and training data and a set of documents  $D = \{D_1, \dots, D_n\}$ , each document of which is labeled with a specific value from set  $C$
- (ii) Calculate text representations for documents in set  $D$
- (iii) Build a classification model based on the training data, which indicates the relationship between the

features in the underlying document and one of the classes

- (iv) Predict class labels for class-unknown documents using the trained model

Calculating text representation, training classification models, and predicting class labels for class-unknown documents are the main steps of text classification. The entire steps, factors, and the way they organize in TC are shown in Figure 1.

As shown in Figure 1, before documents can be analyzed by a classification model, documents need to be pre-processed in a specific way such as be represented by vectors with numerical values. These values relate to predefined classes that the classification model can understand. This process is called text representation, and it is an essential prerequisite for TC tasks [9]. There are many methods proposed for text representation among which the vector space model (VSM) is the most commonly used one [10]. VSM is a feature vector that consists of numerical values which are also called term weights for representing a document. Components of this kind of model can be of different types such as words, sentences, and phrases [11]. These components are also called terms which are extracted from a document to form a bag of words (BOW) [12]. The abilities of these terms in distinguishing different documents are represented by numerical values (weights) related to the terms [13]. For example, a document can be represented as a vector of weighted features (or terms)  $d_k = (t_1, t_2, \dots, t_n)$  and a corresponding weight vector  $w_k = (w_1, w_2, \dots, w_n)$ , where  $n$  is the number of selected features (terms) and  $w_1, w_2, \dots, w_n$  are the weights of  $t_1, t_2, \dots, t_n$ . Then, a collection of documents (corpus) can be represented as shown in Figure 2, where element  $w_{i,j}$  represents the weight of  $t_j$  from  $d_i$ .

As we can see in the matrix, each term in each document can only be assigned to one weight at the same time in VSM. It is obviously crucial to assign appropriate weights to terms for the performance of text classification. Therefore, many methods which are called term weighting scheme (TWS) are proposed to determine the weights for terms of documents. Different TWSs generate different vectors for the same document, thus attributing to the document with different representations. ‘‘Good’’ term weighting methods are of fundamental importance for guaranteeing good TC performance. So far, there are two main categories of TWSs in the literature: semantic-based TWSs and statistics-based TWSs [14].

The semantic-based TWSs focus on the semantic relationships between terms and documents which are hidden behind the extracted features (words) as well as focus on the meaning of words [15]. For example, Rao et al. proposed a new model based on a neural network which captures semantics of continuous text representations [16]. Based on the distributed hypothesis in which meanings of terms from documents with similar meanings will also be similar, a neural network was used to embed words into a continuous vector space (Word2Vec) for capturing the semantic information of words [17]. Doc2Vec, based on Word2Vec, was extended from the word level to the document level by fully

using the information of the word sequence [18]. Because the above methods all have limitations when being used on their own, Kim et al. combined the BOW and Doc2Vec and proposed the bag-of-concepts method for overcoming the limitations [19], and Kim et al. also tried different methods, namely, TF-IDF, LDA, and Doc2Vec, for text representation [20]. More recently, Wu et al. proposed a novel phrase-based text representation called Phrase2Vec, which includes skip phrase, CBOP, and GloVeFP. They applied the novel method to text analysis research, and results show that Phrase2Vec can improve the performance of TC and clustering tasks [21]. Jaeyoung Kim et al. first applied capsule networks, which achieved success in image classification, in the TC task and demonstrated comparable performance to well-known schemes at the time [22].

For non-English languages, semantic analysis is also widely used for TC. Ye-wang Chen et al. proposed a novel method using the biggest open and free internet knowledge-based Baidu Baike to capture the semantic relationships of words to categories to enhance the performance of TC in Chinese text [23]. Ashraf Elnagar et al. introduced two new datasets for Arabic TC tasks, namely, SANAD and NADiA, both of which are freely available for research studies. In their experiments of extensive comparisons among several deep learning (DL) models for Arabic TC on SANAD and NADiA, their method outperformed others because of no requirement of a preprocessing stage and being completely based on deep learning models [24].

However, TWSs based on semantic analysis are more complex than statistical counterparts in analyzing and calculating the process. Furthermore, for semantic-based methods, performance cannot be significantly improved. Therefore, statistics-based TWSs are still major topics in the field of text classification [25]. Normally, most statistics-based TWSs all depend on the following philosophies:

- (i) Terms with higher occurrences in a document relate to the document better, which is the basic idea of ‘‘term frequency’’ (TF)
- (ii) Terms with occurrence in fewer documents relate to the documents where they occur better, which is the basic idea of ‘‘inverse document frequency’’ (IDF)

According to these principles, there are two main factors in a statistics-based TWS as it is shown in the following equation:

$$\text{TWS} = \text{term frequency factor} * \text{collection frequency factor.} \quad (1)$$

Many methods are proposed in the literature based on different implementations of equation (1). Some popular examples are shown in Table 1, where values of ‘‘NONE’’ indicate there is no corresponding method for the specific parameter. As the table shows, some methods focus on modifying the term frequency factor (i.e., LogTF-RF [11] and SQRT\_TF-IGM [25]), while some focus on developing novel methods as the collection frequency factor (i.e., TF-IDF [26], TF-CHI2 [27], TF-IEF [14], and TF-IGM [25]). Nevertheless, TF-IDF is still one of the most preferred methods.

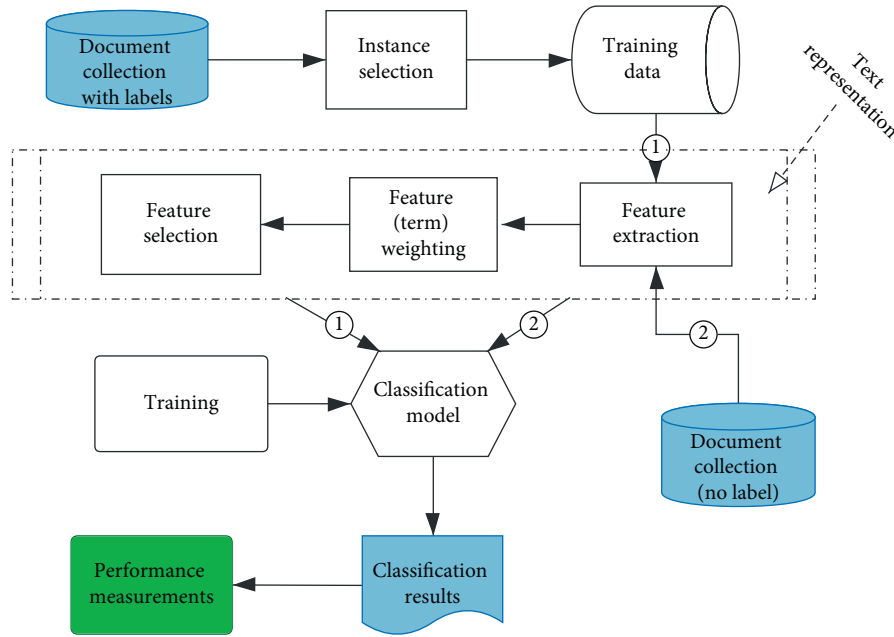


FIGURE 1: The entire steps, factors, and the way they organize in text classification.

$$\begin{bmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,j} \\ w_{2,1} & w_{2,2} & \dots & w_{2,j} \\ \vdots & \vdots & \vdots & \vdots \\ w_{i,1} & w_{i,2} & \dots & w_{i,j} \end{bmatrix}$$

FIGURE 2: An example of the weight matrix for a collection of documents in the vector space model.

TABLE 1: Some popular examples of statistics-based term weighting schemes.

TWS	Term frequency factor	Collection frequency factor
TF	Term frequency (TF)	NONE
DF	NONE	Document frequency (DF)
TF-IDF	Term frequency (TF)	Inverse document frequency (IDF)
TF-CHI2	Term frequency (TF)	Globalized chi-square score (CHI2)
LogTF-RF	Logarithm of term frequency	Relevance frequency (RF)
TF-IGM	Term frequency (TF)	Inverse gravity moment (IGM)
SQRT_TF-IGM	Square root of term frequency	Inverse gravity moment (IGM)
TF-IEF	Term frequency (TF)	Inverse exponential frequency (IEF)

In this paper, enriching the collection frequency factor of statistics-based TWS is concerned, i.e., TF-IDF, for handling situations of imbalanced data distribution. A new formula is designed by using the variance between the DF of a specific term and the average value of all DFs ( $\overline{DF}$ ) instead of original DF in TF-IDF. Based on the new formula, a novel method named TF-IADF and three other TWSs based on the same idea are proposed to enhance the TC performance in the imbalanced situation of internet media reports.

The remainder of this paper is organized as follows. An overview of background study on statistics-based term weighting schemes is given in Section 2. The main idea of our novel methods is described in Section 3. Section 4 briefly introduces the experimental settings and datasets, including data preprocessing, the classifiers, and the measurements.

Experimental results and the analysis are presented in Section 5. The final conclusion is given in Section 6.

## 2. Background Study

When looking at the studies related to TWS in the literature, TF-IDF, originally designed for information retrieval (IR), may be at the top of the list. However, as Chen et al. stated, due to its original design, TF-IDF is not effective enough in the text classification domain [25]. Thus, they proposed a new statistics-based model named inverse gravity moment (IGM) to describe the inter-category distribution. Based on IGM, TF-IGM and sqrt\_TF-IGM (RTF) are proposed. In their demonstration on popular classifiers, namely, SVM and kNN, the

proposed methods had better performance in measurements such as micro-F1 and macro-F1 than existing TWSs (TF, TF-IDF, TFIDF-ICSDF, TF-CHI, TF-PB, and TF-RF). However, Turgut Dogan et al. revealed that, for each case where the term document frequency changes, the term with the same weight is given by TF-IGM. This means terms with different distinguishing abilities obtain the same weights from the standard IGM method which is unreasonable [28]. In their studies, two novel TWSs, namely, SQRT\_TF-IGM<sub>imp</sub> and TF-IGM<sub>imp</sub>, are proposed deriving from IGM to overcome its limitations. In other aspects, Zhong Tang et al. described two deficiencies from which TF-IDF suffers, namely, collection frequency factor being undefined (division by zero) or being equal to zero in some special cases. They proposed a novel method, namely, term frequency-inverse exponential frequency (TF-IEF), to overcome these drawbacks [14]. The proposed methods replaced the IDF with a global weighting factor IEF, and a log-like method is used to characterize the collection frequency factor. It greatly reduced the influence caused by terms with high TF values, which helped in generating a more representative vector of terms. The experiments stated that the novel methods had an improved performance than compared schemes. The knowledge about Chinese language and Chinese culture provided by Baidu Baike is learned and organized by Chinese language-speaking people and professional employees of Baidu company. Therefore, Baidu Baike is used for optimizing TC on Chinese text a couple of times in the Chinese language aspect [23, 29]. However, both Baidu Baike-based methods are based on semantic analysis, and huge calculations are required for processing.

However, most of these methods are based on the assumption that the dataset is relatively balanced in distribution. In fact, the imbalanced distribution of the dataset occurs frequently in the TC domain [30]. Furthermore, the classification performance is heavily affected by the imbalanced distribution of the dataset in TC [31, 32]. Many studies have been proposed to address this problem, such as [33, 34]. In these proposed studies, two common ways are used to solve the problem of data imbalance, namely, the data-driven methods and the algorithm-driven methods. The data-driven method is to adjust the proportion of data categories by under-sampling, oversampling, or a combination of under-sampling and oversampling. The algorithm-driven method is to adjust the classification algorithm to achieve the effect of promoting learning without changing the dataset. The simulation results of these proposed methods show that the more unbalanced the proportion of categories is, the lower the overall performance of TC becomes. One of the main reasons is that some less-common terms in large-scale categories are weighted even higher than some more-common terms in small-scale categories due to their frequencies of occurrence.

Document classification of Chinese media reports on the internet which is also a TC problem with imbalanced dataset is researched in this paper. And a more representative model

in cases of imbalance data is tried to create by modifying the term weighting method.

### 3. Novel Term Weighting Methods Based on Improved TF-IDF

TF-IDF is the most widely used TWS proposed by Karen Spärck Jones [26]. In this section, a new TWS based on TF-IDF, namely TF-IADF, and its variants proposed in this paper are described in specific.

*3.1. Overview to TF-IDF.* TF-IDF [35] is a combination of term frequency (TF) and inverse document frequency (IDF). Since the original value of term frequency in a document is used directly, the TF representation is one of the simplest TWSs. TF is based on the assumption that a term with a higher term frequency value is regarded to be more important than that with a lower term frequency value. It only depends on the number of occurrences of a specific term in a local document. Therefore, the capacity of TF for distinguishing all relevant documents from other irrelevant documents is very low due to its ignorance of collection frequency. To address this problem, the inverse document frequency (IDF) was proposed with a concern of collection frequency which enhanced the discriminative capacity of a term for text classification [36]. IDF extends from document frequency (DF) which means the number of documents where a term occurs. It is proposed based on the assumption that a term which occurs in fewer documents is regarded to be more important than that which occurs in more documents [11]. The IDF value of a specific term can be obtained as shown in the following:

$$\text{IDF}(t, d, D) = \log \frac{|D|}{\text{DF}(t, D)}. \quad (2)$$

In equation (2),  $\text{DF}(t, D)$  represents the DF value of term  $t$  in corpus  $D$ . The symbol in equation (2) represents the total number of documents in corpus  $D$ . To avoid infinity of some extreme cases, the formula is sometimes optimized as shown in the following:

$$\text{IDF}(t, d, D) = \log \frac{|D| + 1}{\text{DF}(t, D) + 1}. \quad (3)$$

After that, Jones extended the IDF method by adding the TF value into calculation [26]. The proposed combination with TF and IDF is the most well-known term weighting method, namely, TF-IDF. Similar with IDF, TF-IDF is also a global statistical measure. The classical structure of TF-IDF is shown as

$$\text{TF-IDF}(t, d, D) = \text{TF}(t, d) * \text{IDF}(t, d, D). \quad (4)$$

In equation (4),  $\text{TF-IDF}(t, d, D)$  represents the weight of term  $t$  of document  $d$  in corpus  $D$ , while  $\text{TF}(t, d)$  represents the TF value of term  $t$  in document  $d$ .

As we introduced in Section 2, TF-IDF is not effective enough in the text classification domain due to its original

design. And many research studies have been deployed in optimizing term weighting methods based on TF-IDF from different perspectives. Some of them developed new methods replacing the term frequency factor or document frequency factor of TF-IDF, while some of them modified the existing method of TF-IDF. This paper focuses on modifying IDF to improve the TC performance, especially for Chinese internet media content.

**3.2. Proposed Methods.** When looking into the formula of calculating the value of IDF as shown in equations (2) and (3), we notice that when the corpus is not very balanced which means the size of different categories in a corpus varies from each other, terms from categories with larger size will be assigned smaller values than terms from other categories. This is obviously not in line with the real situation. Moreover, for some low document-frequency terms, the value of IDF is much higher than others even when those low document-frequency terms are meaningless, which is not in line with the true situation either. To address this kind of problems, we focus on the deviation of the DF value between a specific term and the overview of all terms in the whole corpus since when the deviation between the DF value of a specific term and the average of all DF values is large, its discriminative ability is weak. This factor should be considered in the term weighting process.

**Definition 1** (average document frequency (ADF)). It is the variance between the DF value of a specific term and the average of all DF values in a corpus.

We modified the collection frequency factor by adding ADF into calculation to address the problems mentioned above. In this study, the average of all DF values in the corpus is represented as  $\overline{DF}$ , while the ADF value of term  $t$  in document  $D$  is represented as  $A_{DF}(t, D)$ . Equations (5) and (6) show how they are calculated, where  $n$  is the number of terms:

$$\overline{DF} = \frac{\sum DF(t, D)}{n}, \quad (5)$$

$$A_{DF}(t, D) = \frac{(DF(t, D) - \overline{DF})^2}{n}. \quad (6)$$

As ADF is extended from DF, the simplest way of optimizing IDF is to replace DF by ADF in the formula. Then, we get a novel formula of collection frequency which is shown as follows:

$$IADF(t, D) = \log \frac{|D| + 1}{A_{DF}(t, D) + 1}. \quad (7)$$

In fact, the IDF method is successful enough in most cases; what we need to do is just to modify it for some extreme cases. Then, we get another novel formula as shown in equation (8), where ADF is used to reduce the weight of the terms with extremely high or extremely low DF value.

$$IADF^+(t, D) = \log \frac{|D| + 1}{DF(t, D) + 1} * \frac{1}{\log(A_{DF}(t, D) + 1) + 1}. \quad (8)$$

The two ADF-based methods can improve the TC performance in some cases we mentioned. However, there are still limitations due to the variance itself that when the size is too large or too small, the variance will be relatively too small or too large. Extreme values for terms will obviously impact the TC performance. Therefore, we further optimized the formula by normalizing the ADF to reduce the effect caused by the extreme value of terms. First,  $A_{DF}(t, D)$  is modified as shown in equation (9) and then using the normalization formula as shown in equation (10):

$$A'_{DF}(t, D) = \log \frac{1}{(A_{DF}(t, D) + 1)} + 1, \quad (9)$$

$$A''_{DF}(t, D) = \frac{A'_{DF}(t, D) - \min(A'_{DF}(t, D))}{\max(A'_{DF}(t, D)) - \min(A'_{DF}(t, D))}. \quad (10)$$

Based on  $A_{DF}''$ , another two novel formulas are designed as shown in equations (11) and (12), where  $\alpha$  (default value is 1) is used as an optional weight proportion to adjust the importance of  $A_{DF}''$  in different cases.

$$IADF_{\text{norm}}(t, D) = \log \frac{|D| + 1}{A_{DF}''(t, D) + 1}, \quad (11)$$

$$IADF_{\text{norm}}^+(t, D) = \log \frac{|D| + 1}{DF(t, D) + 1} * (A_{DF}''(t, D) * \alpha). \quad (12)$$

Based on the above four proposed formulas of collection frequency based on IDF, we get four novel term weighting methods which are shown in equations (13)–(16):

$$\text{TF-IADF}(t, d, D) = \text{TF}(t, d) * IADF(t, D), \quad (13)$$

$$\text{TF-IADF}^+(t, d, D) = \text{TF}(t, d) * IADF^+(t, D), \quad (14)$$

$$\text{TF-IADF}_{\text{norm}}(t, d, D) = \text{TF}(t, d) * IADF_{\text{norm}}(t, D), \quad (15)$$

$$\text{TF-IADF}_{\text{norm}}^+(t, d, D) = \text{TF}(t, d) * IADF_{\text{norm}}^+(t, D). \quad (16)$$

As a result, the optimized text representation model of processing internet media reports is shown in Figure 3. Four new calculation formulas are used to replace IDF of TF-IDF. And four novel term weighting methods are obtained to enhance the performance of processing unbalanced text collection.

## 4. Case Study

To evaluate our proposed TWSs, experiments are carried out by using proposed methods in state-of-the-art classification algorithms on both Chinese and English corpora. In this section, datasets used in experiments are briefly described. Then, algorithms utilized for the classification process and the measurements used for performance evaluation in this

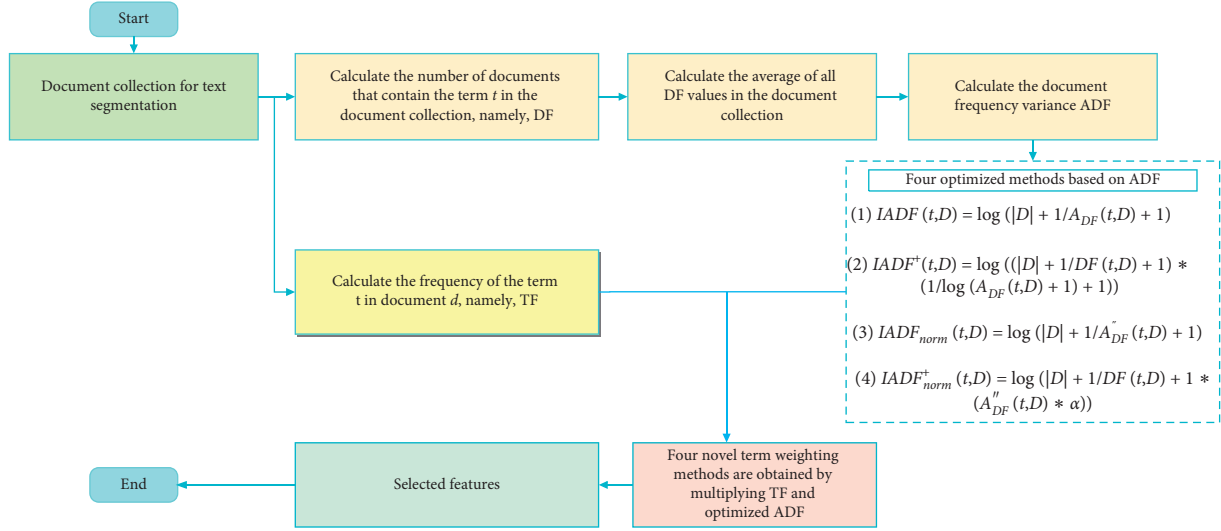


FIGURE 3: The flowchart of the text classification method based on the improved TF-IDF.

study are introduced. Finally, the experiment settings are also presented.

**4.1. The Data Source.** This study carried out experiments on three different datasets, i.e., standard dataset of English text, namely, Reuters-21578 corpus, classic dataset of Chinese text, namely, Fudan corpus, and a collection of Chinese internet media reports named Internet corpus, which were crawled off web and transformed into forms of Chinese textual document.

**4.1.1. Reuters-21578 Corpus.** The Reuters-21578 corpus contains top-10 categories of Reuters-ModApte separately split which is most preferred in the TC domain [37]. In this study, multilabeled samples are removed since single-label-classification is focused. So, only 8 categories of 5607 training samples and 2270 test samples in Reuters-21578 were used in our experiments. The detail of data distribution of this corpus is shown in Figure 4 and Table 2.

**4.1.2. Fudan Corpus.** The Fudan University TC corpus is from the Chinese NLP group in Department of Computer Information and Technology, Fudan University of China. There are 20 categories of which the data distributions are shown in Figure 5 and Table 3. Similar to Reuters-21578, Fudan corpus is also an unbalanced dataset, but in Chinese language.

**4.1.3. Internet Corpus.** To test the performance of our proposed methods on Chinese internet media reports, some reports from the web are crawled, and this corpus is formed. There are six categories, namely, sport, education, tourism, traffic, tech, finance, and food. The data format is shown in Figure 6. There are three parts in each instance of the test data which are the category index, the article content, and

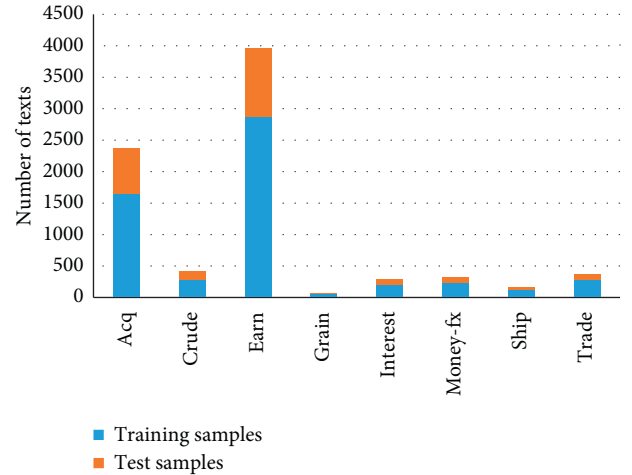


FIGURE 4: Data distribution of Reuters-21578 (unbalanced).

TABLE 2: Reuters-21578 corpus.

No.	Class label	Training samples	Testing samples
1	Acq	1634	728
2	Crude	277	131
3	Earn	2859	1086
4	Grain	41	10
5	Interest	198	87
6	Money-fx	214	93
7	Ship	115	43
8	Trade	269	92

the total number of words. In this study, both balanced and unbalanced data distributions of this corpus are tried.

**4.2. Classification Algorithms Used for Experiments and Measurements.** Three popular classification algorithms, namely, naïve Bayes (NB), support vector machine (SVM), and random forests (RF), are utilized using our proposed methods and existing methods for a brief comparison.

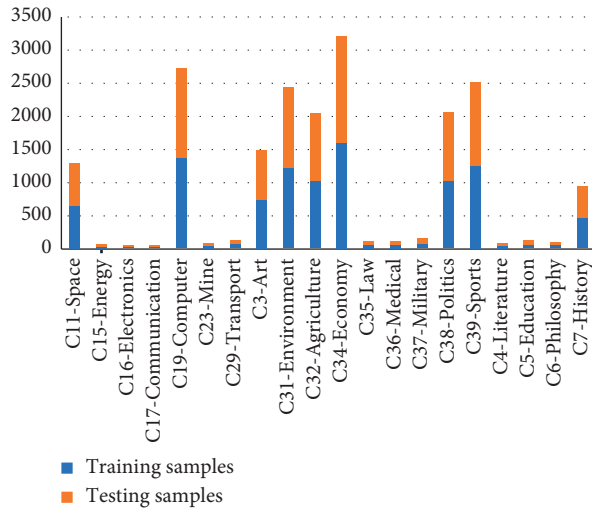


FIGURE 5: Data distribution of Fudan corpus (unbalanced).

TABLE 3: Fudan corpus.

No.	Class label	Training samples	Testing samples
1	C11-Space	640	642
2	C15-Energy	32	33
3	C16-Electronics	27	28
4	C17-Communication	25	27
5	C19-Computer	1357	1358
6	C23-Mine	33	34
7	C29-Transport	57	59
8	C3-Art	740	742
9	C31-Environment	1217	1218
10	C32-Agriculture	1021	1022
11	C34-Economy	1600	1601
12	C35-Law	51	52
13	C36-Medical	51	53
14	C37-Military	74	76
15	C38-Politics	1024	1026
16	C39-Sports	1253	1254
17	C4-Literature	33	34
18	C5-Education	59	61
19	C6-Philosophy	44	45
20	C7-History	466	468

,分类,文章,字数

0,娱乐,《青蛇》造型师默认新《红楼梦》额妆抄袭 (图)凡是看过电影《青蛇》的人,都不会忘记青白二蛇的经典造型:飘逸的身材,妩媚的妆容,而最独特的就是取材于京剧女角头型的精华之处"打片子"(也被称之为额妆).....由于吴宝玲在香港出道很早,所以无论是对造型设计本身还是对文化的钻研,都有独到之处,而且往往能把两者结合得很好,她甚至感觉新版《红楼梦》的设计缺乏文学性,"《红楼梦》本身就让人很期待,假如我来设计造型的话,因为片中将出现很多生活场景,我可能会尽量做得不那么沉重,而且人物角色的感觉应该有所区分,薛宝钗就是薛宝钗,而不会分不清楚她和林黛玉。"最后吴宝玲表示,她还会观察一段时间,看后期的设计是否真的完全会对《青蛇》有所模仿,必要的时候她可能向法律寻找帮助。昨晚,记者致电新《红楼梦》导演李少红,但她的手机一直无人接听。(责任编辑:CL).1519

FIGURE 6: Style of test data.

Introductions about these algorithms and measurements for evaluation are given as follows.

**4.2.1. Naïve Bayes.** Naïve Bayes algorithm [38] is a well-known TC classifier based on Bayes' assumption that the features are regarded to be independent from each other. In the TC process, document  $d_k$  can be represented as a vector of terms  $(t_1, t_2, \dots, t_n)$ . The probability that  $d_k$  belongs to a specific category  $c_i$  can be calculated using equation (17). More details about the NB classifier can be accessed in [39]. In this study, a NB classifier is used for evaluating the text weighting performance.

$$P(c_i|d_k) = P(c_i) \frac{\prod_{j=1}^n P(t_{jk}|c_i)}{P(d_k)}. \quad (17)$$

**4.2.2. Support Vector Machine.** SVM [40, 41] is one of the most preferred algorithms for TC and many other pattern recognition problems. Since it is a learning algorithm, it can handle problems with high dimensions well. The main principle of the SVM is to create linear or nonlinear hyperplanes to separate positive and negative samples. SVM uses some samples in the training set (called support vectors) to create hyperplanes at locations maximizing margins between negative and positive samples. In this study, a classic SVM classifier is used for evaluating the text weighting performance.

**4.2.3. Random Forests.** The random forest (RF) algorithm [42] is a parallelizable integration method which is one of the most preferred classifiers in the field of TC [43]. RF is composed of multiple decision trees. It is used to build a forest in a random way, which consists of many decision trees (DT). There is no correlation between each decision tree in the RF. After the RF is obtained, for each sample input, each decision tree in the forest is judged to see which category this sample belongs to, which category ultimately gets the most results, and which type of input prediction is. The architecture of RF is shown in Figure 7, where DT refers to the decision tree. In this study, a RF classifier is used for evaluating the text weighting performance.

**4.2.4. Measurements for TC Performance.** To evaluate the classification performance on the aforementioned datasets, accuracy, precision, recall, and  $F_1$  score are calculated for validation according to equations (18)–(21), where  $c_i$  refers to a predefined category, while  $TP(c_i)$  refers to the number of documents which belong to  $c_i$  resulting in  $c_i$ ,  $FN(c_i)$  refers to the number of documents which do not belong to  $c_i$  resulting in  $c_i$ ,  $FP(c_i)$  refers to the number of documents which belong to  $c_i$  not resulting in  $c_i$ , and  $TN(c_i)$  refers to the number of documents which do not belong to  $c_i$  not resulting in  $c_i$ . The relationship between them and the classification result are shown in Table 4.

TABLE 4: Meanings of TP, TN, FP, and FN.

An instance of documents in the corpus ( $d_k$ )		
Sample $d_k$	Result in $c_i$	Not a result in $c_i$
Belongs to $c_i$	TP ( $c_i$ )	FP ( $c_i$ )
Does not belong to $c_i$	FN ( $c_i$ )	TN ( $c_i$ )

$$\text{Accuracy}(c_i) = \frac{TP(c_i) + TN(c_i)}{TP(c_i) + TN(c_i) + FP(c_i) + FN(c_i)}, \quad (18)$$

$$\text{precision}(c_i) = \frac{TP(c_i)}{TP(c_i) + FP(c_i)}, \quad (19)$$

$$\text{recall}(c_i) = \frac{TP(c_i)}{TP(c_i) + FN(c_i)}, \quad (20)$$

$$F_1(c_i) = \frac{2 * \text{precision}(c_i) * \text{recall}(c_i)}{\text{precision}(c_i) + \text{recall}(c_i)}. \quad (21)$$

In multiclass classification problems, the overall performance can be measured by averaging the evaluation methods. Microaverage and macroaverage are used widely for this purpose. In this study, the microaveraged  $F_1$  (micro- $F_1$ ) and macroaveraged  $F_1$  (macro- $F_1$ ) measurements are also calculated to evaluate the experimental methods. The definition of macro- $F_1$  and micro- $F_1$  is as shown in equations (22) and (23):

$$\text{macro} - F_1 = \frac{1}{m} \sum_{i=1}^m F_1(c_i), \quad (22)$$

$$\text{micro} - F_1 = \frac{2 * \sum_{i=1}^m TP(c_i)}{2 * \sum_{i=1}^m TP(c_i) + \sum_{i=1}^m FP(c_i) + \sum_{i=1}^m FN(c_i)}. \quad (23)$$

In cases of unbalanced distribution, it is better to use micro- $F_1$  than macro- $F_1$  since the data size of categories is not considered in micro- $F_1$  score calculation.

**4.3. Experiment Settings.** In this study, we carried out three experiments on the aforementioned datasets. All experiments were implemented on a 64 bit Windows 10 computer with 8 GB internal storage. The experimental code was written in Python language using Scikit-learn (sk-learn). sk-learn is a commonly used third-party module in machine learning which encapsulates many commonly used machine learning algorithms such as regression, dimension reduction, clustering [44], and classification. For each dataset, in preprocessing, term weighting, and term extraction, term representation and classification were utilized. In preprocessing, all documents were segmented into words by the open-source tool Jieba, and stop words were removed in this process. After that, a vector space model (VSM) was used for term representation using words as terms. Term weighting methods including TF-IDF and proposed methods, i.e., TF-IADF, TF-IADF<sub>norm</sub>, TF-IADF<sup>+</sup>, and TF-IADF<sup>+</sup><sub>norm</sub> were used here to form a final representation for each document. Finally, NB classifier, SVM classifier, and RF classifier were utilized for TC purpose. Combinations of different term



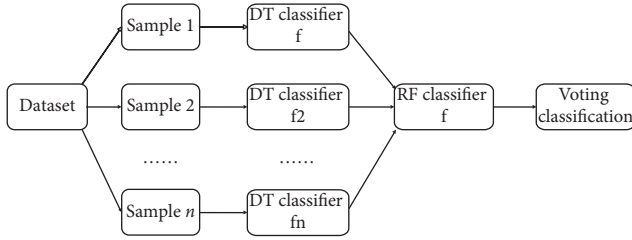


FIGURE 7: The architecture of random forests.

weighting methods and different classification algorithms on different datasets were compared for a brief analysis.

## 5. Results' Analysis and Discussion

**5.1. Results on the Internet Corpus.** In this part, the experiments on both a balanced dataset and an unbalanced dataset of internet media reports in the Chinese language are carried out.

**5.1.1. Balanced Dataset.** In this dataset, there are 1000 training samples and 1000 test samples in each category. This is an exactly balanced dataset. SVM, NB, and RF classifiers are utilized as TWS for a comparison with TF-IDF. The overall performance of the proposed TF-IADF outperformed all other methods in SVM and RF classifiers as shown in Figure 8. The details of experimental results are shown in Table 5 in that the proposed TF-IADF<sup>+</sup><sub>norm</sub> demonstrates better performance than TF-IDF in all cases. Furthermore, all proposed methods outperformed the TF-IDF, in some cases, respectively.

For the SVM classifier, the overall classification effect is better than the other two classifiers, and the micro-*F1* value is over 94%. TF-IADF has achieved the best effect of 94.45% (increased by 0.31% than TF-IDF). For the RF classifier, TF-IADF achieves the best effect. In addition, TF-IADF<sub>norm</sub> and TF-IADF<sup>+</sup><sub>norm</sub> also come with some improvement. The micro-*F1* value of TF-IDF is 84.51%, while that of TF-IADF<sup>+</sup><sub>norm</sub> is 85.52%, and that of TF-IADF<sub>norm</sub> is 85.26%. The micro-*F1* value of TF-IADF reaches to 86.37%, which is the best among all methods, an improvement of 1.86% from TF-IDF. For the NB classifier, both TF-IADF<sup>+</sup> and TF-IADF<sup>+</sup><sub>norm</sub> show improvements for the micro-*F1* value of TF-IADF<sup>+</sup> and TF-IADF<sup>+</sup><sub>norm</sub> reaching 92.26% and 92.33%, while that of TF-IDF is only 92.13%. TF-IADF<sup>+</sup><sub>norm</sub> achieves the best effect with an increase of 1.2% over TF-IDF. In the classification results, we notice that the precision rate and recall rate of finance are relatively low, at only 80% and 85%, respectively. The reason may be that the terms of this category are not obvious enough, and the scope involved is relatively wide, which may cover some contents from tourism and traffic categories, resulting in the poor classification effect of the whole category.

The results show that, in the case of balanced dataset which is not focused by our design, the proposed methods can improve the classifiers to a certain extent even though the improvement range is not so obvious.

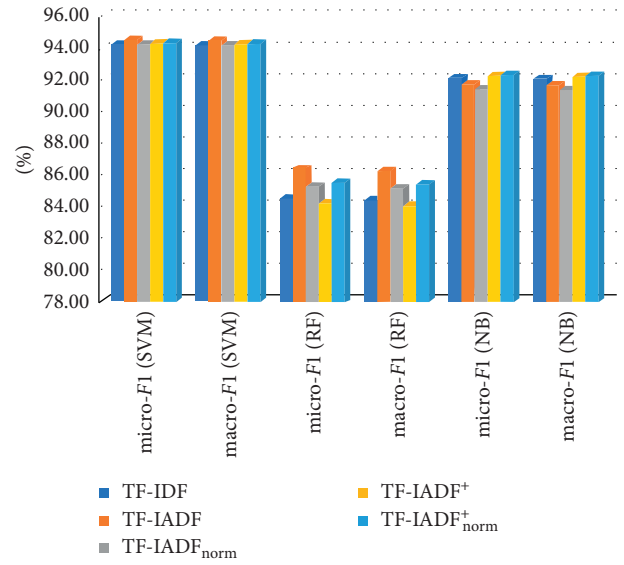


FIGURE 8: Data distribution of the Fudan corpus (unbalanced).

**5.1.2. Unbalanced Dataset.** In this section, we investigate how TC performance is impacted by unbalanced datasets. Dataset here is specially designed. The size of food is increased to 5000, and the sizes of other categories are kept the same as in the balanced dataset (1000) except sport. For sport, the size is gradually reduced from 400 to 50, decreasing 50 at a time.

**Definition 2 (balance ratio).** It is the proportion between the sizes of the category with the smallest size and the category with the largest size.

In this section, we call the proportion between the size of sport and that of food as balance ratio which drops from 8% to 1% with the decrease of sport's size. Experiments by using the SVM classifier with TF-IDF and our proposed methods are carried out. The experimental results show that performances on categories other than sport and food are basically the same with the balance ratio changing. However, the recall and *F1* score of sport and the precision of food are impacted heavily by decreasing the balance ratio. The details are shown in Tables 6–8. As shown, the balance ratio is decreasing, with the recall and *F1* score of sport also decreasing in all methods. However, the proposed TF-IADF and TF-IADF<sub>norm</sub> outperformed TF-IDF in all cases. Especially, when the balance ratio decreased to an extreme value (1%), TF-IADF<sub>norm</sub> came with a recall of 0.7168 which is almost 170% of that of TF-IDF. Due to the balance ratio changes, precision of the category with a relatively large size (food) was also impacted which can be seen from Table 8.

**Definition 3 (decline ratio).** It is the ratio of the current value to the initial value in a declining trend.

To make the relationship between the balance ratio and the performance clear, the decline ratio is calculated. It is the ratio of the value for each balance ratio to its initial value (value at 8%). This can be seen in Figure 9, where the ordinate refers to the decline ratio of the corresponding

TABLE 5: Experiments on the balanced dataset of the internet corpus.

Balanced dataset	TF-IDF (%)	TF-IADF (%)	TF-IADF <sub>norm</sub> (%)	TF-IADF <sup>+</sup> (%)	TF-IADF <sup>+</sup> <sub>norm</sub> (%)
micro-F1 (SVM)	94.23	<b>94.54</b>	<u>94.23</u>	<u>94.32</u>	<u>94.35</u>
macro-F1 (SVM)	94.16	<b>94.48</b>	<u>94.18</u>	<u>94.25</u>	<u>94.28</u>
micro-F1 (RF)	84.51	<b>86.37</b>	<u>85.26</u>	84.23	<u>85.52</u>
macro-F1 (RF)	84.41	<b>86.26</b>	<u>85.16</u>	84.08	<u>85.41</u>
micro-F1 (NB)	92.13	91.73	91.42	<u>92.26</u>	<u>92.33</u>
macro-F1 (NB)	92.05	91.66	91.36	<u>92.18</u>	<u>92.24</u>

TABLE 6: Recall (sport).

Balance ratio (%)	TF-IDF	TF-IADF	TF-IADF <sub>norm</sub>	TF-IADF <sup>+</sup>	TF-IADF <sup>+</sup> <sub>norm</sub>
8	0.9479	<b>0.9581</b>	<u>0.9581</u>	0.9315	0.9162
7	0.9376	<u>0.9458</u>	<b>0.9479</b>	0.8988	0.8886
6	0.9182	<u>0.9397</u>	<b>0.9417</b>	0.8824	0.8691
5	0.9018	<u>0.9346</u>	<b>0.9407</b>	0.8650	0.8507
4	0.8763	<u>0.9243</u>	<b>0.9315</b>	0.8262	0.8129
3	0.8272	<u>0.8937</u>	<b>0.9141</b>	0.7464	0.7301
2	0.7423	<u>0.8252</u>	<b>0.8507</b>	0.6012	0.5910
1	0.4274	<u>0.6626</u>	<b>0.7168</b>	0.2567	0.2413

TABLE 7: F1 score (sport).

Balance ratio (%)	TF-IDF	TF-IADF	TF-IADF <sub>norm</sub>	TF-IADF <sup>+</sup>	TF-IADF <sup>+</sup> <sub>norm</sub>
8	0.9732	<b>0.9786</b>	<u>0.9776</u>	0.9645	0.9562
7	0.9678	<u>0.9722</u>	<b>0.9732</b>	0.9467	0.9410
6	0.9574	<u>0.9689</u>	<b>0.9700</b>	0.9375	0.9300
5	0.9484	<u>0.9662</u>	<b>0.9694</b>	0.9276	0.9193
4	0.9341	<u>0.9607</u>	<b>0.9645</b>	0.9048	0.8968
3	0.9054	<u>0.9438</u>	<b>0.9551</b>	0.8548	0.8440
2	0.8521	<u>0.9042</u>	<b>0.9193</b>	0.7510	0.7429
1	0.5989	<u>0.7971</u>	<b>0.8350</b>	0.4085	0.3888

TABLE 8: Precision (food).

Balance ratio (%)	TF-IDF	TF-IADF	TF-IADF <sub>norm</sub>	TF-IADF <sup>+</sup>	TF-IADF <sup>+</sup> <sub>norm</sub>
8	0.8726	<u>0.9185</u>	<b>0.9305</b>	0.8310	0.7788
7	0.8711	<u>0.9160</u>	<b>0.9288</b>	0.8241	0.7716
6	0.8688	<u>0.9160</u>	<b>0.9296</b>	0.8213	0.7668
5	0.8650	<u>0.9143</u>	<b>0.9296</b>	0.8146	0.7580
4	0.8583	<u>0.9126</u>	<b>0.9279</b>	0.8041	0.7444
3	0.8480	<u>0.9101</u>	<b>0.9253</b>	0.7875	0.7170
2	0.8275	<u>0.8849</u>	<b>0.9092</b>	0.7517	0.6728
1	0.7813	<u>0.8725</u>	<b>0.8976</b>	0.6728	0.5691

performance factor. As it is shown, with the balance ratio decreasing, the growth rate of decline ratio becomes faster and faster. It is obvious that, for datasets with extreme categories such as food (relatively too large) and sport (relatively too small), performances of TF-IADF and TF-IADF<sub>norm</sub> are more stable than TF-IDF.

The overall performance on this dataset is impacted by the balance ratio also. The details of micro-F1 and macro-F1 are shown in Tables 9 and 10, respectively. The proposed TF-IADF and TF-IADF<sub>norm</sub> outperformed TF-IDF in all cases. To make it clear how balance ratio impacts the performance, we calculated the decline ratio which is shown in Figure 10. As it is shown, with the decrease of balance ratio, both micro-F1 and macro-F1 decrease

gradually. For example, when looking at micro-F1, TF-IADF<sub>norm</sub> was with a decrease of just 3.65%, while TF-IDF was with a decrease of 8.06% which is more than twice of that of the proposed TF-IADF<sub>norm</sub>, meaning the performance of TF-IADF<sub>norm</sub> is much more stable than that of TF-IDF in this dataset.

Even though TF-IADF<sup>+</sup> and its variance do not achieve improvement in this experiment, TF-IADF and its variance outperformed TF-IDF significantly. Furthermore, it can be seen from Figure 10 that our proposed TF-IADF and TF-IADF<sub>norm</sub> are not only numerically better but also more stable than TF-IDF. Therefore, considering ADF in the term weighting method can actually improve the performance of text classification considerably in unbalanced cases.

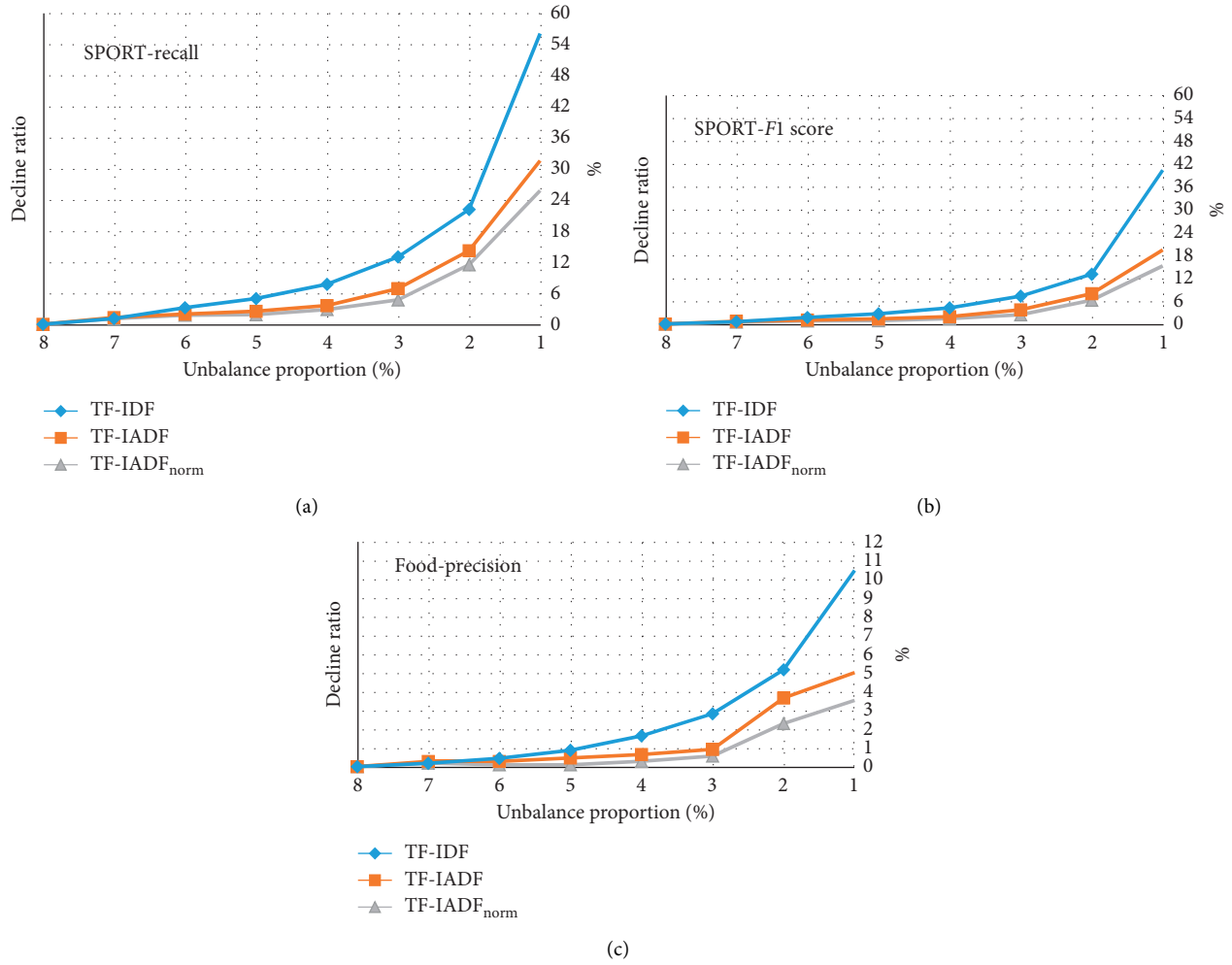


FIGURE 9: The decline ratio of the performance with balance ratio changing.

TABLE 9: Micro-F1.

Balance ratio (%)	TF-IDF	TF-IADF	TF-IADF <sub>norm</sub>	TF-IADF <sup>+</sup>	TF-IADF <sup>+</sup> <sub>norm</sub>
8	0.9292	<b><u>0.9372</u></b>	<b><u>0.9372</u></b>	0.9226	0.9140
7	0.9278	<b><u>0.9354</u></b>	<b><u>0.9360</u></b>	0.9179	0.9094
6	0.9250	<b><u>0.9345</u></b>	<b><u>0.9352</u></b>	0.9154	0.9063
5	0.9226	<b><u>0.9338</u></b>	<b><u>0.9351</u></b>	0.9129	0.9032
4	0.9189	<b><u>0.9323</u></b>	<b><u>0.9338</u></b>	0.9073	0.8977
3	0.9119	<b><u>0.9279</u></b>	<b><u>0.9313</u></b>	0.8960	0.8859
2	0.8994	<b><u>0.9179</u></b>	<b><u>0.9222</u></b>	0.8753	0.8658
1	0.8542	<b><u>0.8946</u></b>	<b><u>0.9029</u></b>	0.8258	0.8156

TABLE 10: Macro-F1.

Balance ratio (%)	TF-IDF	TF-IADF	TF-IADF <sub>norm</sub>	TF-IADF <sup>+</sup>	TF-IADF <sup>+</sup> <sub>norm</sub>
8	0.9287	<b><u>0.9366</u></b>	<b><u>0.9366</u></b>	0.9225	0.9147
7	0.9273	<b><u>0.9349</u></b>	<b><u>0.9355</u></b>	0.9179	0.9102
6	0.9246	<b><u>0.9341</u></b>	<b><u>0.9348</u></b>	0.9154	0.9072
5	0.9223	<b><u>0.9333</u></b>	<b><u>0.9346</u></b>	0.9130	0.9042
4	0.9186	<b><u>0.9319</u></b>	<b><u>0.9333</u></b>	0.9074	0.8987
3	0.9115	<b><u>0.9275</u></b>	<b><u>0.9309</u></b>	0.8956	0.8868
2	0.8985	<b><u>0.9174</u></b>	<b><u>0.9217</u></b>	0.8726	0.8649
1	0.8443	<b><u>0.8922</u></b>	<b><u>0.9014</u></b>	0.8055	<b><u>0.7987</u></b>

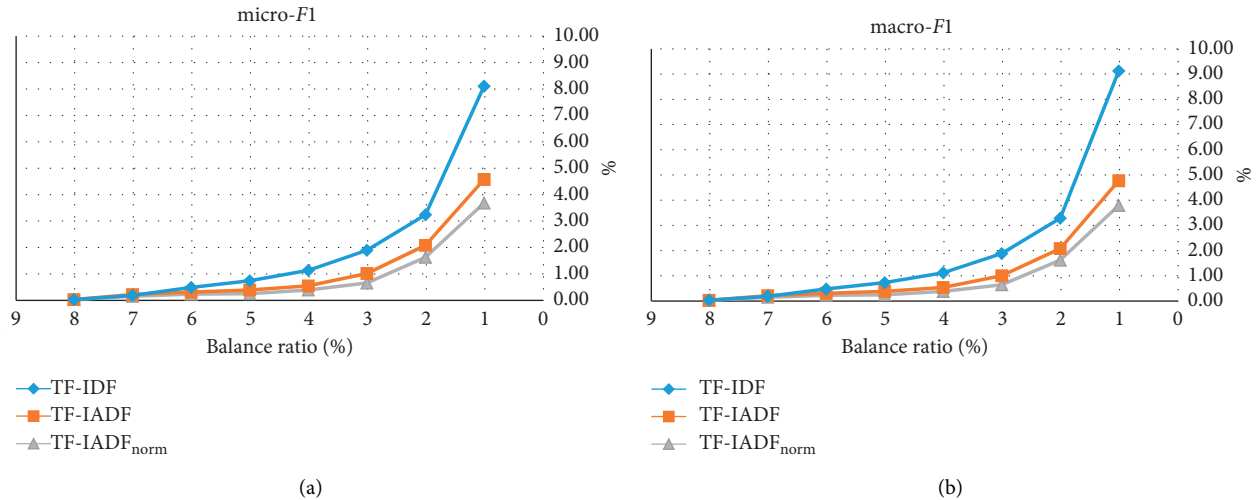


FIGURE 10: The decline ratio of the overall performance with balance ratio changing.

**5.1.3. Analysis.** For this corpus, the TF-IADF<sup>+</sup> and TF-IADF<sub>norm</sub><sup>+</sup> methods are more suitable for the NB classifier, and TF-IADF<sub>norm</sub><sup>+</sup> is the best. For the RF and SVM classifiers, TF-IADF<sub>norm</sub><sup>+</sup> and TF-IADF are more suitable. In addition, TF-IADF is better in case of a balanced dataset, while TF-IADF<sub>norm</sub> is better in case of an unbalanced dataset. This means that the processed TF-IADF<sub>norm</sub> method is more sensitive and stable in the case of unbalanced datasets, while TF-IADF improves the classification effect more when the datasets are relatively balanced. For several mathematical models proposed by different algorithms, it can be concluded that the formula suitable for different algorithms may be different, and for the improvement of the corresponding algorithm effect, it can also be concluded that the ADF index has improved the effect of text classification, which confirms the conjecture; especially, when the dataset is not evenly distributed, the effect of text classification is more stable.

## 5.2. Results on the Fudan Corpus

**5.2.1. Results' Analysis.** Table 11 and Figure 11 show the micro-F1 and macro-F1 scores obtained on Fudan corpus using SVM, RF, and NB algorithms with different TWSs, while Tables 12–14 show the detailed results.

For the NB classifier, TF-IADF<sub>norm</sub><sup>+</sup> has the best performance, which is 89.38%, an increase of 0.39% compared to TF-IDF. Meanwhile, TF-IADF<sup>+</sup> also gets an increase of 0.32%. This means for Chinese text datasets, these two models are more suitable for a NB classifier and can achieve an improvement on the overall classification effect. When comparing the specific measurements which are shown in Figure 12 and more details can be seen in Table 14, TF-IADF<sub>norm</sub><sup>+</sup> has achieved a lot of the highest performance-score items. Furthermore, the difference between performance scores of TF-IDF and those of TF-IADF<sub>norm</sub><sup>+</sup> is relatively small in categories, where TF-IADF<sub>norm</sub><sup>+</sup> is not as good as TF-IDF. For example, in C32 and C35, where TF-

IDF achieves the best precision score, the difference between the precision score of TF-IDF and that of TF-IADF<sub>norm</sub><sup>+</sup> is less than 1%. However, in categories such as C17, where TF-IADF<sub>norm</sub><sup>+</sup> achieves the highest precision score, the difference between that and precision score of TF-IDF is more than 6% which is six times of the difference occurring in categories where TF-IDF achieves the higher score. In fact, in C17, when comparing TF-IADF<sub>norm</sub><sup>+</sup> to TF-IDF, the precision is improved by 6.62%, and the recall score of C17 is also increased by 7.41%, which is obvious. In an overall view of the F1 scores, it can be noticed that, in all the 20 categories except C35 and C5, the scores of TF-IADF<sub>norm</sub><sup>+</sup> are not lower than those of TF-IDF. Especially, in cases of TF-IADF<sub>norm</sub><sup>+</sup> with higher elevation such as C16, the F1 score is increased by nearly 10%.

For the RF classifier, the micro-F1 score obtained by TF-IADF is 81.95%, which is 1.63% higher than that obtained by TF-IDF. Meanwhile, TF-IADF<sup>+</sup> and TF-IADF<sub>norm</sub> also have achieved improvements in different extents. It is known that RF algorithm has certain randomness, but our proposed TF-IADF method is more stable and has achieved better results which can be seen in Figure 13. When looking into detailed results as shown in Table 13, in some categories, such as C29 and C5, the precision score has been improved by 25.98% and 30.16%, respectively, and the recall score has remained basically unchanged. For the performance of F1 score, in some categories, such as C16 and C36, the score obtained by TF-IADF is about 10% higher than that obtained by TF-IDF. Furthermore, there are only three categories where F1 of TF-IADF is not as good as TF-IDF.

For the SVM classifier, it is still the case that TF-IADF and TF-IADF<sub>norm</sub> have achieved improved performances. When comparing with TF-IDF in the micro-F1 score, TF-IADF has increased by 0.73%, while TF-IADF<sub>norm</sub> has increased by 0.63%. It can be seen in Figure 14 that, for the SVM classifier, these two improved methods are more stable and have achieved some improvements. As shown in Table 12, there are many categories where the precision scores are very high, even up to 1. It can easily be seen in Figure 14

TABLE 11: Overall performances on the Fudan corpus.

Fudan corpus	TF-IDF	TF-IADF	TF-IADF <sub>norm</sub>	TF-IADF <sup>+</sup> (%)	TF-IADF <sup>+</sup> <sub>norm</sub> (%)
micro-F1 (SVM)	50.02	<b>54.57</b>	54.93	48.85	47.90
macro-F1 (SVM)	90.31	<b>91.04</b>	90.94	90.24	89.82
micro-F1 (RF)	48.57	<b>51.43</b>	48.80	49.78	51.64
macro-F1 (RF)	80.32	<b>81.95</b>	80.83	81.20	79.85
micro-F1 (NB)	65.44	62.28	59.94	66.39	<b>66.97</b>
macro-F1 (NB)	88.99	88.20	87.55	89.31	<b>89.38</b>

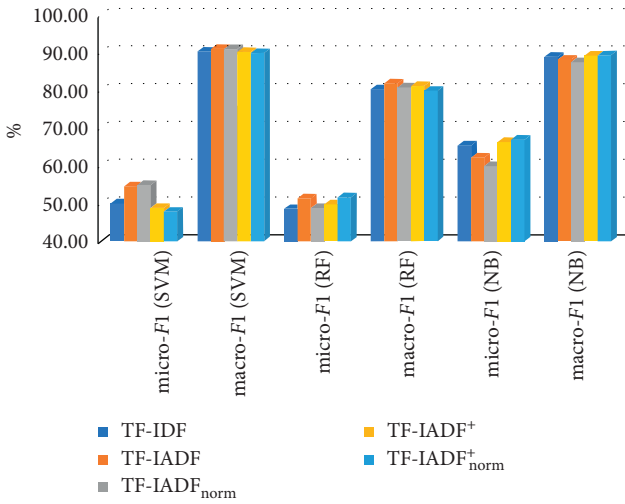


FIGURE 11: Overall performances of classification on the Fudan corpus.

that the size (data amount) of those categories is very small. For example, there are only 32 training samples in C15. The reason is that due to the small eigenvalues, no other categories being considered as this category are responsible for the high precision score. However, it can also be concluded from the detailed results that, for most of those categories with a high precision but a small size, the recall score is significantly reduced; that is to say, many documents are assigned to the wrong categories, which may be due to the impact caused by a small number of features. For example, the recall score of C15 is only 0.09. TF-IADF and TF-IADF<sub>norm</sub> perform better in most of the categories, especially in C29, where the size is small. Although the precision scores obtained by the two methods are decreased by about 4%, the recall scores are improved by about 30% which is significant. In addition, there are only two categories where *F1* scores of TF-IADF are lower than those of TF-IDF, while for most categories, TF-IADF is better.

**5.2.2. Further Discussion.** Due to the unbalanced distribution of this corpus, all categories can be roughly divided into three groups by their size for specific analysis. The first group refers to categories with size in the range of 0 to 100, the second refers to categories with size in the range of 101 to 1000, and the third group refers to categories with size over 1000. First of all, for the first group, due to the insufficient size of training data, the same property is reflected in all term weighting methods, i.e., high

precision and low recall scores, such as C16. Taking SVM combined with TF-IDF as an example, the precision score reaches 100%, while the recall score is only 3.57%, resulting in poor classification effect. In short, it assigns a few samples correctly, while a large number of test samples are assigned to the wrong categories. And the performance is similar in the RF and NB classifiers. This is the impact caused by the unbalanced distribution and the small size of the training data. For the second group, the number of training sets has been improved to a higher level. For example, in C11 when using the SVM classifier with TF-IDF, the precision score and recall score are 96.10% and 92.06%, respectively, which is a significant improvement compared to the first group. For the third category, such as C19, the precision score and recall score are 95.64% and 98.53%, respectively. It can be concluded that the precision score in categories with large training set size is relatively low, while the recall rate is relatively high. In categories with a small training set size, it will have better precision but very low recall score. In these kinds of conditions, our proposed methods will demonstrate a similar but more stable performance compared to TF-IDF. Taking C29 as an example, there are only 57 training samples. Comparing TF-IDF with TF-IADF using the SVM classifier, the precision scores are 100% to 96.15%, while the recall scores are 8.48% to 42.37%. It can be seen that although the accuracy rate of our algorithm is slightly reduced, the recall rate is greatly increased by nearly 500%, and the *F1* score is also greatly improved, from 15.63% to 58.82%, showing a very obvious improvement.

In the Fudan corpus, a similar phenomenon with the internet corpus can be seen. For example, the TF-IADF<sub>norm</sub> method is the best one in the NB classifier, while TF-IADF<sup>+</sup> is also better than TF-IDF. And with RF and SVM classifiers, both TF-IADF and TF-IADF<sub>norm</sub> achieve a relatively stable performance. The difference is that TF-IADF achieves the best effect in the Fudan unbalanced dataset. It is also noticed that although the micro-*F1* score of the SVM classifier is the best, the macro-*F1* score, which is about 50%, is not as good as that of the NB classifier which is over 65%. That is to say, although the overall accuracy of the SVM classifier is high, the effect of the NB classifier is better when each category is regarded as equally important.

**5.3. Results on Reuters-21578.** The overall performance on Reuters-21578 of all methods is shown in Table 15, while Figure 15 shows a brief comparison between all proposed

TABLE 12: Performance on the Fudan corpus using the SVM classifier.

SVM	F1 score											
	Precision						Recall					
	TF-IDF	TF-IADF	TF-IADF <sup>norm</sup>	TF-IDF	TF-IADF	TF-IADF <sup>norm</sup>	TF-IDF	TF-IADF	TF-IADF <sup>norm</sup>	TF-IDF	TF-IADF	TF-IADF <sup>norm</sup>
C11	96.10	95.98	<u>96.13</u>	92.06	<u>92.99</u>	92.84	91.75%	<u>94.46</u>	94.03	<u>94.45</u>	93.94%	93.94%
C15	<b>100.00</b>	<b>100.00</b>	<b>100.00%</b>	<b>9.09</b>	<b>9.09</b>	<b>9.09</b>	<b>9.09%</b>	<b>16.67</b>	<b>16.67</b>	<b>16.67</b>	<b>16.67%</b>	<b>16.67%</b>
C16	<b>100.00</b>	<b>100.00</b>	<b>100.00%</b>	3.57	3.57	3.57	3.57%	6.90	6.90	6.90	6.90%	6.90%
C17	<b>100.00</b>	<b>100.00</b>	<b>100.00%</b>	18.52	<u>22.22</u>	<u>22.22</u>	18.52%	<b>36.36</b>	31.25	<b>36.36</b>	31.25%	31.25%
C19	95.64	<u>95.93</u>	95.56%	98.53	<u>98.97</u>	<u>99.12</u>	98.31%	<u>97.43</u>	97.06	<u>97.29</u>	96.99%	96.92%
C23	<b>100.00</b>	<b>100.00</b>	<b>100.00%</b>	<b>5.88</b>	<b>5.88</b>	<b>5.88</b>	<b>5.88%</b>	<b>11.11</b>	<b>11.11</b>	<b>11.11</b>	<b>11.11%</b>	<b>11.11%</b>
C29	<b>100.00</b>	96.15	95.65	8.48	<u>42.37</u>	37.29	5.09%	<b>58.82</b>	15.63	53.66	9.68%	6.56%
C3	87.10	88.85	<b>89.17</b>	<b>95.55</b>	<b>95.55</b>	95.42	94.74%	<u>92.08</u>	91.13	<b>92.19</b>	91.05%	90.65%
C31	95.61	<u>96.97</u>	<u>95.80%</u>	96.63	97.21	96.96	97.04%	<u>97.09</u>	96.12	<u>97.36</u>	<u>96.54%</u>	<u>96.45%</u>
C32	94.17	94.10	<u>94.37%</u>	96.38	<u>96.67</u>	96.58	95.21%	<u>95.37</u>	95.26	95.00	<u>95.30%</u>	<u>94.79%</u>
C34	90.20	90.18	<b>90.23</b>	94.88	<b>95.19</b>	94.63	94.88%	<b>92.62</b>	92.48	92.38	92.20%	90.73%
C35	<b>100.00</b>	95.24	<b>100.00%</b>	23.08	<u>38.46</u>	<u>46.15</u>	11.54%	<u>54.80</u>	37.50	<b>61.54</b>	20.69%	14.29%
C36	<b>100.00</b>	<b>100.00</b>	<b>100.00%</b>	1.89	<u>13.21</u>	<u>15.09</u>	1.89%	<u>23.33</u>	3.70	<u>26.23</u>	3.70%	3.70%
C37	<b>100.00</b>	<b>100.00</b>	<b>100.00%</b>	<b>5.26</b>	<b>5.26</b>	<b>5.26</b>	<b>5.26%</b>	<b>10.00</b>	<b>10.00</b>	<b>10.00</b>	<b>10.00%</b>	<b>10.00%</b>
C38	82.65	82.36	82.78	<b>95.61</b>	94.64	94.64	95.52%	88.07	88.66	88.31	<b>88.93%</b>	88.50%
C39	85.36	<b>88.69</b>	84.86%	96.73	<u>96.97</u>	96.49	<u>96.97%</u>	<b>92.65</b>	90.69	91.88	90.51%	<u>91.08%</u>
C4	<b>100.00</b>	<b>100.00</b>	<b>100.00%</b>	<b>5.88</b>	<b>5.88</b>	<b>5.88</b>	<b>5.88%</b>	<b>11.11</b>	<b>11.11</b>	<b>11.11</b>	<b>11.11%</b>	<b>11.11%</b>
C5	<b>66.67</b>	<b>66.67</b>	<b>66.67%</b>	<b>3.28</b>	<b>3.28</b>	<b>3.28</b>	<b>3.28%</b>	<b>6.25</b>	<b>6.25</b>	<b>6.25</b>	<b>6.25%</b>	<b>6.25%</b>
C6	88.89	88.89	88.89%	17.78	17.78	<b>20.00</b>	17.78%	29.63	29.63	<b>32.73</b>	29.63%	29.63%
C7	<b>82.78</b>	81.30	80.99	68.80	<u>72.44</u>	<u>73.72</u>	68.16%	76.61	75.15	<b>77.18</b>	74.62%	72.90%

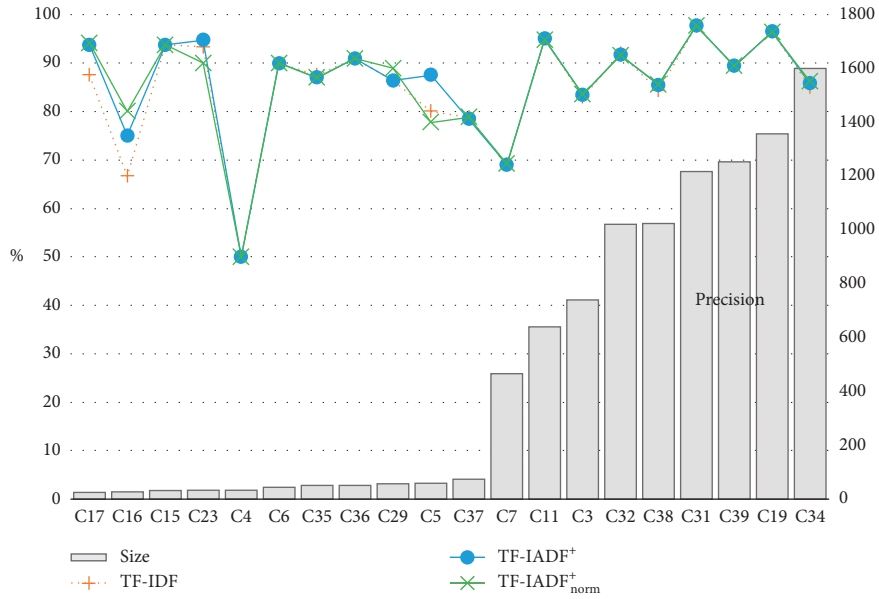
TABLE 13: Performance on the Fudan corpus using the RF classifier.

RF	Precision				Recall				F1 score				
	TF-IDF (%)	TF-IDF IADF (%)	TF-IDF IADF <sup>+</sup> (%)	TF-IDF IADF <sup>norm</sup> (%)	TF-IDF (%)	TF-IDF IADF (%)	TF-IDF IADF <sup>norm</sup> (%)	TF-IDF IADF <sup>norm</sup> (%)	TF-IDF (%)	TF-IDF IADF (%)	TF-IDF IADF <sup>norm</sup> (%)	TF-IDF IADF <sup>norm</sup> (%)	TF-IDF IADF <sup>norm</sup> (%)
C11	82.94	<b>90.85</b>	<u>87.61</u>	<u>86.50</u>	82.56	<b>85.05</b>	79.13	78.19	<u>83.80</u>	82.75	<b>87.85</b>	82.20	82.63
C15	33.33	80.00	<u>66.67</u>	<b>83.33</b>	<b>15.15</b>	12.12	12.12	12.12	<b>15.15</b>	20.83	21.05	19.05	20.51
C16	40.00	50.00	<u>30.77</u>	<b>57.14</b>	7.14	<b>17.86</b>	14.29	14.29	14.29	12.12	<b>26.32</b>	22.86	19.51
C17	55.56	50.00	<u>38.71</u>	<b>61.91</b>	37.04	37.04	44.44	44.44	<b>48.15</b>	44.44	42.55	41.38	35.90
C19	89.90	<b>90.55</b>	86.31	88.24	96.32	<u>96.61</u>	96.32	<b>97.05</b>	<u>96.69</u>	93.00	<b>93.48</b>	91.76	91.37
C23	43.75	38.46	<u>33.33</u>	<b>50.00</b>	20.59	14.71	8.82	<b>23.53</b>	14.71	28.00	21.28	13.95	<b>28.07</b>
C29	43.59	<b>69.57</b>	67.74	<u>60.87</u>	28.81	27.12	22.03	<u>35.59</u>	<b>47.46</b>	34.69	39.02	30.59	<u>46.67</u>
C3	<b>75.77</b>	74.89	<u>75.14</u>	71.73	86.39	<b>90.43</b>	<b>90.43</b>	<u>89.62</u>	<u>87.87</u>	80.73	<b>81.93</b>	81.04	<u>81.75</u>
C31	92.04	<u>93.17</u>	91.77	<b>93.49</b>	90.23	<b>91.79</b>	91.54	<u>90.72</u>	88.42	91.13	<b>92.47</b>	<u>91.66</u>	<u>91.78</u>
C32	84.10	87.01	<u>86.30</u>	81.96	84.34	86.50	83.86	<b>86.79</b>	80.92	84.22	86.75	85.06	<b>87.43</b>
C34	74.49	<b>75.85</b>	<u>75.03</u>	72.61	88.26	<b>90.26</b>	89.82	<u>89.69</u>	88.07	80.79	<b>82.43</b>	80.90	81.71
C35	66.67	<b>77.78</b>	<u>68.75</u>	60.00	7.69	13.46	<b>25.00</b>	<u>21.15</u>	17.31	13.79	<u>22.95</u>	<b>34.67</b>	<u>32.35</u>
C36	57.14	<u>66.67</u>	<b>75.00</b>	51.52	22.64	30.19	<u>26.42</u>	<u>28.30</u>	<b>32.08</b>	32.43	<b>41.56</b>	35.00	<u>41.10</u>
C37	42.86	42.86	<u>13.64</u>	<b>66.67</b>	11.84	<b>15.79</b>	6.58	3.95	10.53	18.56	<b>23.08</b>	11.49	6.12
C38	65.66	66.45	<u>66.86</u>	66.81	76.22	<u>77.97</u>	<b>80.51</b>	80.02	76.51	70.55	<u>71.75</u>	<b>74.89</b>	72.85
C39	84.79	86.29	<u>87.54</u>	<u>86.67</u>	84.45	<b>85.33</b>	84.05	82.62	81.42	84.62	<b>85.81</b>	<u>85.76</u>	<u>85.06</u>
C4	<b>100.00</b>	75.00	<b>100.00</b>	<b>100.00</b>	5.88	<b>8.82</b>	5.88	5.88	<b>8.82</b>	11.11	<u>15.79</u>	11.11	11.11
C5	55.56	<b>85.71</b>	62.50	77.78	8.20	9.84	<b>11.48</b>	8.20	<b>11.48</b>	14.29	<u>17.65</u>	18.67	14.49
C6	<b>90.91</b>	84.62	<u>63.64</u>	90.00	22.22	<b>24.44</b>	22.22	15.56	20.00	35.71	<b>37.93</b>	35.09	25.00
C7	69.54	70.48	<u>73.04</u>	65.05	25.86	25.00	17.95	<b>27.14</b>	25.86	37.70	36.91	28.82	<b>40.25</b>

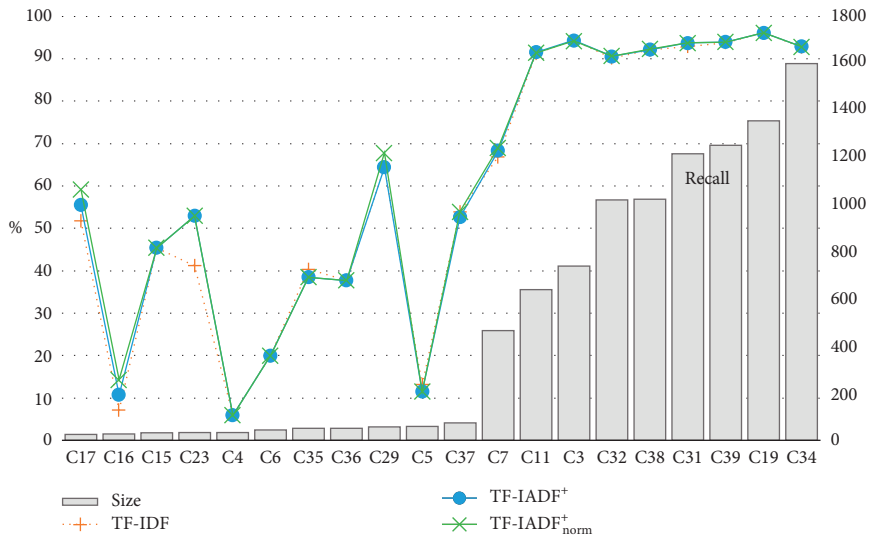
TABLE 14: Performance on the Fudan corpus using the NB classifier.

NB	Precision				Recall				F1 score					
	TF-IDF (%)	TF-IADF (%)	TF-IADF <sup>norm</sup> (%)	TF-IADF <sup>+</sup> (%)	TF-IDF (%)	TF-IADF (%)	TF-IADF <sup>norm</sup> (%)	TF-IADF <sup>+</sup> (%)	TF-IDF (%)	TF-IADF (%)	TF-IADF <sup>norm</sup> (%)	TF-IADF <sup>+</sup> (%)	TF-IADF <sup>norm</sup> (%)	TF-IADF <sup>+</sup> (%)
C11	94.98	94.57	94.53	<b>94.99</b>	94.83	91.28	89.56	88.79	<b>91.59</b>	<b>91.43</b>	93.09	92.00	91.57	<b>93.26</b>
C15	93.75	<b>100.00</b>	<b>100.00</b>	93.75	93.75	<b>45.46</b>	30.30	27.27	<b>45.46</b>	<b>45.46</b>	<b>61.22</b>	46.51	42.86	<b>61.22</b>
C16	66.67	<b>100.00</b>	<b>100.00</b>	75.00	80.00	7.14	3.57	3.57	10.71	<b>14.29</b>	12.90	6.90	6.90	<b>24.24</b>
C17	87.50	86.67	83.33	93.75	<b>94.12</b>	51.85	48.15	37.04	55.56	<b>59.26</b>	65.12	61.91	51.28	<b>72.73</b>
C19	96.38	96.02	95.87	<b>96.46</b>	96.38	<b>96.17</b>	96.02	95.80	<b>96.17</b>	<b>96.17</b>	96.28	96.02	95.84	96.28
C23	93.33	92.31	<b>100.00</b>	94.74	90.00	41.18	35.29	20.59	<b>52.94</b>	<b>52.94</b>	57.14	51.06	34.15	<b>67.93</b>
C29	86.36	87.88	87.10	86.36	<b>88.89</b>	64.41	49.15	45.76	64.41	<b>67.80</b>	73.79	63.04	60.00	<b>76.92</b>
C3	82.94	81.81	81.19	83.43	<b>83.51</b>	94.34	95.15	<b>95.42</b>	94.34	94.21	88.27	87.98	87.73	<b>88.55</b>
C31	97.59	<b>97.89</b>	97.78	97.69	<b>97.69</b>	92.94	91.38	90.48	<b>93.76</b>	<b>93.76</b>	95.21	94.52	93.99	<b>95.69</b>
C32	<b>91.92</b>	90.70	90.35	91.68	91.60	90.22	88.75	87.97	90.61	<b>90.71</b>	91.06	89.71	89.14	<b>91.15</b>
C34	85.05	82.86	81.32	85.76	<b>86.20</b>	92.76	92.69	92.44	<b>92.94</b>	<b>92.88</b>	88.74	87.50	86.52	<b>89.42</b>
C35	<b>87.50</b>	86.96	86.96	86.96	<b>86.96</b>	<b>40.39</b>	38.46	<b>40.39</b>	38.46	38.46	<b>55.26</b>	53.33	<b>55.26</b>	53.33
C36	90.91	88.89	<b>93.33</b>	90.91	90.91	<b>37.74</b>	30.19	26.42	<b>37.74</b>	<b>37.74</b>	<b>53.33</b>	45.07	41.18	<b>53.33</b>
C37	78.85	<b>84.09</b>	81.40	78.43	78.85	<b>53.95</b>	48.68	46.05	52.63	<b>53.95</b>	<b>64.06</b>	61.67	58.82	<b>64.06</b>
C38	84.30	83.12	82.16	85.46	<b>85.55</b>	92.11	<b>92.59</b>	92.01	92.20	92.30	88.03	87.60	86.81	<b>88.80</b>
C39	89.30	88.29	87.30	<b>89.39</b>	<b>89.39</b>	93.86	93.22	92.66	<b>94.02</b>	<b>94.02</b>	91.52	90.69	89.90	<b>91.64</b>
C4	50.00	<b>100.00</b>	<b>100.00</b>	50.00	50.00	<b>5.88</b>	<b>5.88</b>	<b>5.88</b>	<b>5.88</b>	<b>5.88</b>	10.53	<b>11.11</b>	<b>11.11</b>	10.53
C5	80.00	75.00	71.43	<b>87.50</b>	77.78	<b>13.12</b>	9.84	8.20	11.48	11.48	<b>22.54</b>	17.39	14.71	20.00
C6	<b>90.00</b>	<b>90.00</b>	<b>90.00</b>	<b>90.00</b>	<b>90.00</b>	<b>20.00</b>	<b>20.00</b>	<b>20.00</b>	<b>20.00</b>	<b>20.00</b>	<b>32.73</b>	<b>32.73</b>	<b>32.73</b>	<b>32.73</b>
C7	69.10	71.76	<b>72.49</b>	68.97	69.17	66.88	66.24	64.74	68.38	<b>69.02</b>	67.97	<b>68.89</b>	68.40	<b>69.09</b>



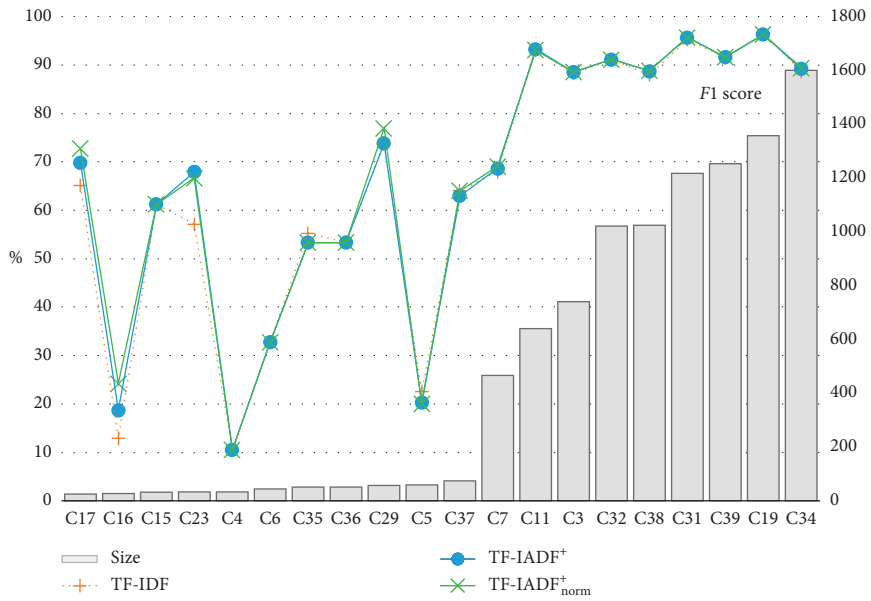


(a)



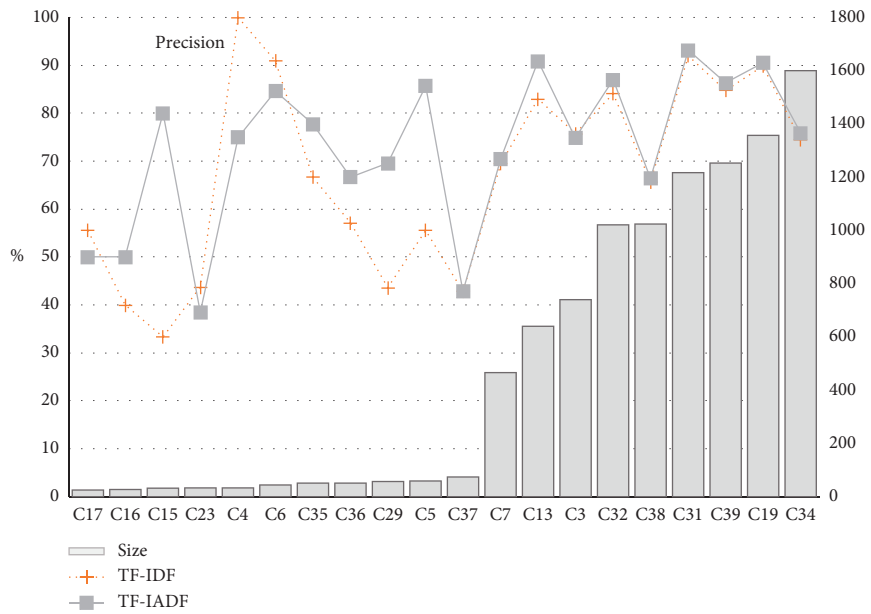
(b)

FIGURE 12: Continued.



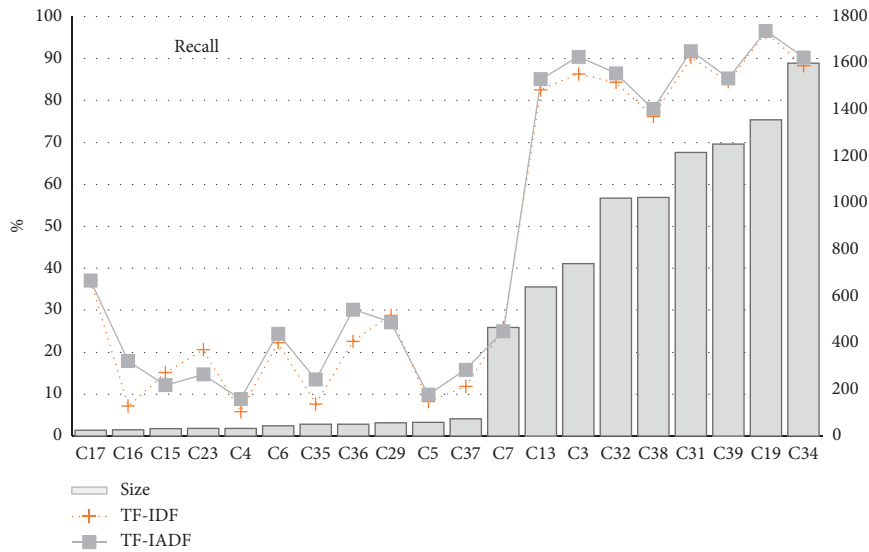
(c)

FIGURE 12: Detailed performances on the Fudan corpus using the NB classifier.

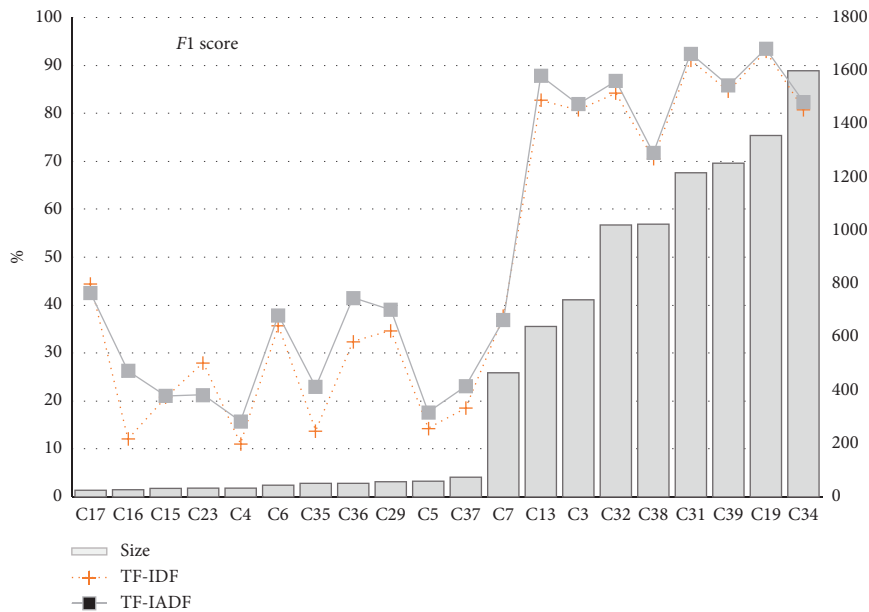


(a)

FIGURE 13: Continued.



(b)



(c)

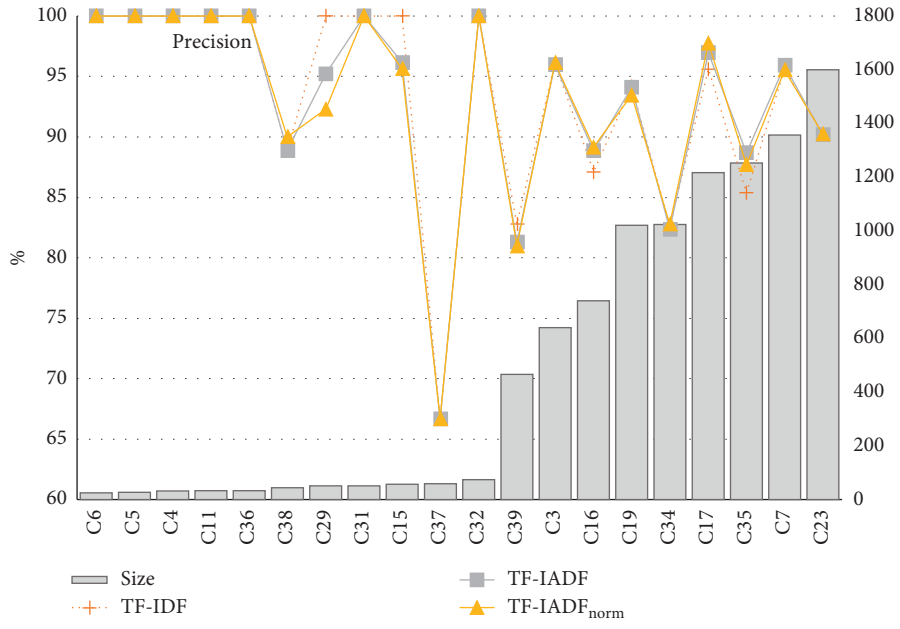
FIGURE 13: Detailed performances on the Fudan corpus using the RF classifier.

methods and the well-known TF-IDF method when combining with different classifiers. More detailed results are shown in Tables 16–18.

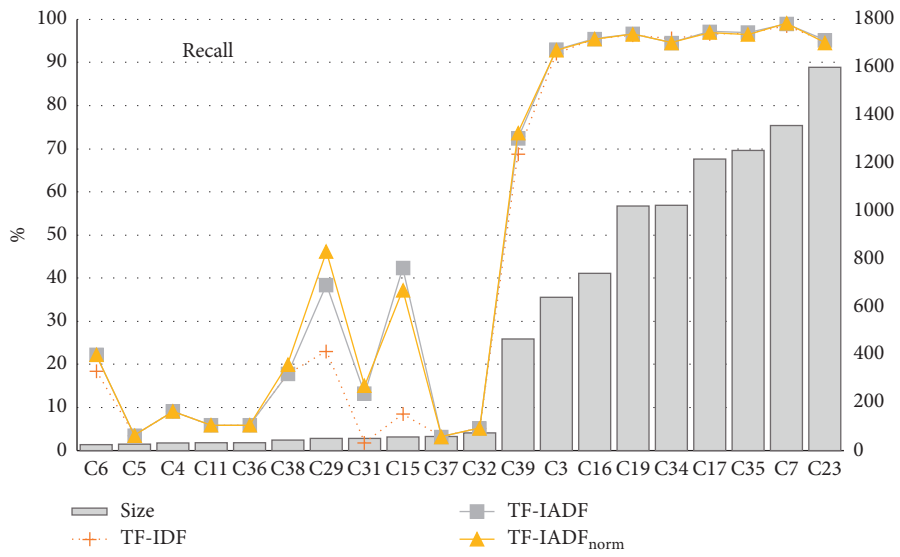
For the SVM classifier,  $TF-IADF_{norm}$  comes with the best performance which is 0.86% higher than TF-IDF in the micro-F1 score. According to the detailed results shown in Table 15 and Figure 16, the improvement of  $TF-IADF_{norm}$  is mainly contributed by the recall score, which is greatly improved in the category of crude, grain, interest, and money-fx. Especially in the category of money-fx, the recall score has increased by 10.35%. The other categories also improved in different extents, namely, crude increased by 3.31%, interest increased by 4%, and grain increased by 10%. From the perspective of the F1 score, in the total eight categories, the max F1 score of five categories is obtained by

$TF-IADF_{norm}$ , and the largest increase is obtained in the category of grain, where the F1 score increased from 57.14% of TF-IDF to 66.67% of  $TF-IADF_{norm}$ , close to 10%.

For the NB classifier, the results are different from those on the Chinese dataset, where  $TF-IADF^+$  and  $TF-IADF^+_{norm}$  achieved better performance, and performances of  $TF-IADF^+$  and  $TF-IADF^+_{norm}$  were worse on this corpus. However, the proposed TF-IADF and  $TF-IADF_{norm}$  showed an improvement, the micro-F1 score of which is 0.59% and 0.5% higher than that of TF-IDF, respectively. As shown in Table 16 and Figure 17, TF-IADF has achieved the highest F1 score among all methods in all categories of this corpus, and the maximum improvement, 3.83%, occurs in the category of interest. Specifically, in terms of precision, we see that the four categories of cloud, interest, ship, and trade have



(a)



(b)

FIGURE 14: Continued.

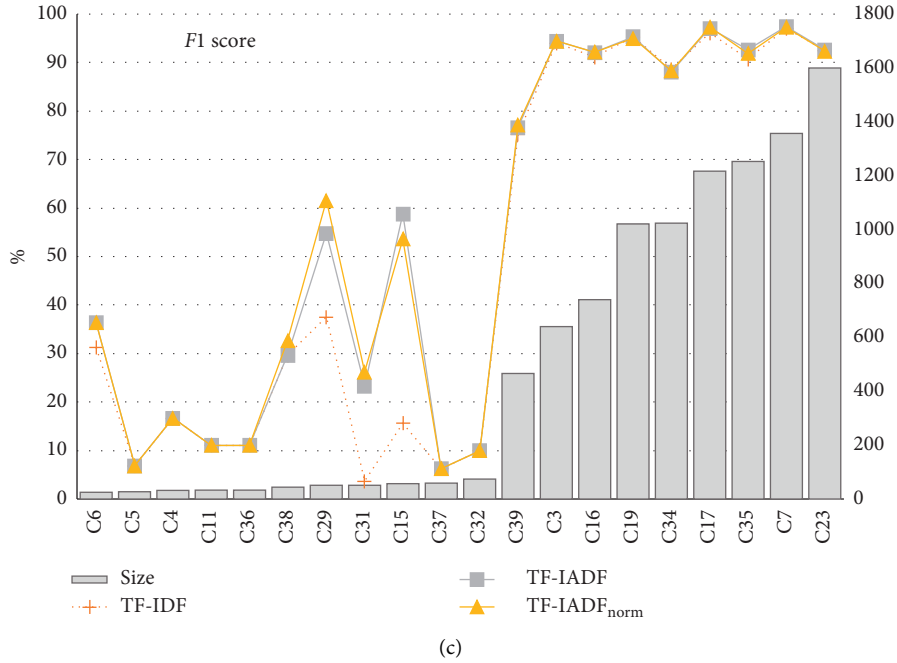


FIGURE 14: Detailed performances on the Fudan corpus using the SVM classifier.

TABLE 15: Overall performances on Reuters-21578.

	TF-IDF (%)	TF-IADF (%)	TF-IADF <sub>norm</sub> (%)	TF-IADF <sup>+</sup> (%)	TF-IADF <sup>+</sup> <sub>norm</sub> (%)
micro-F1 (SVM)	91.00	<b>91.41</b>	<u>91.86</u>	90.59	90.23
macro-F1 (SVM)	79.63	<b>80.33</b>	<u>81.97</u>	78.42	75.77
micro-F1 (RF)	87.45	<b>85.23</b>	<u>87.77</u>	87.41	<u>88.14</u>
macro-F1 (RF)	69.90	<b>59.03</b>	65.43	67.91	<u>67.76</u>
micro-F1 (NB)	89.82	<u>90.41</u>	<u>90.32</u>	89.41	<b>89.36</b>
macro-F1 (NB)	76.22	<u>77.86</u>	<u>77.77</u>	76.21	<b>75.95</b>

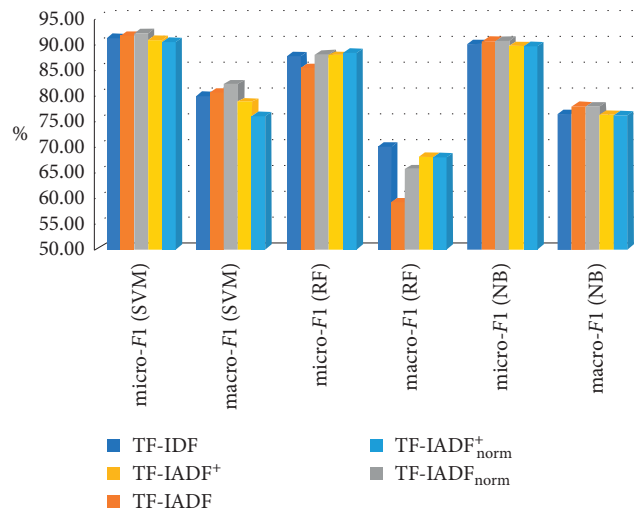


FIGURE 15: Performance of the NB classifier on Reuters-21578 (precision, recall, and F1 score).

improved significantly, with an improvement of 2.36%, 4.35%, 6.82%, and 3.75%, respectively. In terms of the recall score, the increase in the money-fx category is 4.60%.

For the RF algorithm, TF-IADF<sup>+</sup><sub>norm</sub> and TF-IADF<sub>norm</sub> have improved the effect, micro-F1 score of which is 0.69% and 0.32% higher than that of TF-IDF, respectively. TF-

TABLE 16: Performance on Reuters-21758 using the SVM classifier.

SVM	Precision				Recall				F1 score			
	TF-IDF (%)	TF-IDF IADF (%)	TF-IDF IADF <sup>+</sup> (%)	TF-IDF IADF <sup>+</sup> norm (%)	TF-IDF (%)	TF-IDF IADF (%)	TF-IDF IADF <sup>+</sup> (%)	TF-IDF IADF <sup>+</sup> norm (%)	TF-IDF (%)	TF-IDF IADF (%)	TF-IDF IADF <sup>+</sup> (%)	TF-IDF IADF <sup>+</sup> norm (%)
Acq	90.17	91.43	90.06	90.42	92.24	91.95	92.10	92.24	91.19	91.69	91.21	91.32
Crude	90.39	89.72	<b>92.63</b>	92.55	77.69	79.34	<b>80.99</b>	71.90	83.56	84.21	81.48	80.93
Earn	92.06	92.00	91.35	90.21	97.42	97.69	97.69	<b>97.88</b>	94.66	94.76	94.33	93.89
Grain	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	40.00	40.00	<b>50.00</b>	30.00	57.14	57.14	57.14	46.15
Interest	91.30	90.00	88.24	90.48	56.00	<b>60.00</b>	<b>60.00</b>	50.67	69.42	<b>72.00</b>	68.33	64.96
Moneyfx	79.71	<b>84.06</b>	80.00	79.03	63.22	66.67	<b>73.56</b>	56.32	70.51	74.36	68.83	65.77
Ship	<b>100.00</b>	<b>100.00</b>	95.83	95.24	<b>63.89</b>	61.11	<b>63.89</b>	55.56	<b>77.97</b>	75.86	74.58	70.18
Trade	90.63	89.80	88.78	<b>92.47</b>	94.57	<b>95.65</b>	94.57	93.48	92.55	<b>92.63</b>	91.58	<b>92.97</b>

TABLE 17: Performance on Reuters-21758 using the RF classifier.

RF	Precision				Recall				F1 score				
	TF-IDF (%)	TF-IDF IADF (%)	TF-IDF IADF <sub>norm</sub> (%)	TF-IDF IADF <sup>+</sup> (%)	TF-IDF IADF (%)	TF-IDF IADF <sub>norm</sub> (%)	TF-IDF IADF <sup>+</sup> (%)	TF-IDF IADF <sub>norm</sub> (%)	TF-IDF IADF (%)	TF-IDF IADF <sub>norm</sub> (%)	TF-IDF IADF <sup>+</sup> (%)	TF-IDF IADF <sub>norm</sub> (%)	TF-IDF IADF <sup>+</sup> (%)
Acq	86.52	82.25	85.41	81.77	90.37	90.52	92.53	94.11	88.41	88.83	87.51	89.06	89.06
Crude	92.68	<b>93.24</b>	84.88	91.25	62.81	57.03	60.33	60.33	<b>74.88</b>	70.53	72.64	74.53	74.53
Earn	89.89	87.78	90.48	<b>93.11</b>	96.86	96.86	96.58	94.83	93.24	93.44	<b>93.96</b>	93.55	93.55
Grain	<b>100.00</b>	<b>100.00</b>	33.33	<b>100.00</b>	<b>30.00</b>	10.00	10.00	20.00	<b>46.15</b>	15.39	33.33	16.67	16.67
Interest	75.51	74.29	78.00	72.73	49.33	34.67	52.00	<b>53.33</b>	59.68	62.40	61.54	<b>62.71</b>	62.71
Moneyfx	75.47	73.17	<b>82.69</b>	72.41	45.98	34.48	<b>49.43</b>	<b>48.28</b>	57.14	<b>61.87</b>	<b>57.93</b>	55.94	55.94
Ship	90.00	90.00	92.31	<b>100.00</b>	50.00	25.00	33.33	41.67	64.29	48.98	58.82	<b>66.67</b>	66.67
Trade	72.73	78.21	<b>84.88</b>	80.23	78.26	66.30	79.35	75.00	75.39	82.02	77.53	<b>82.98</b>	82.98





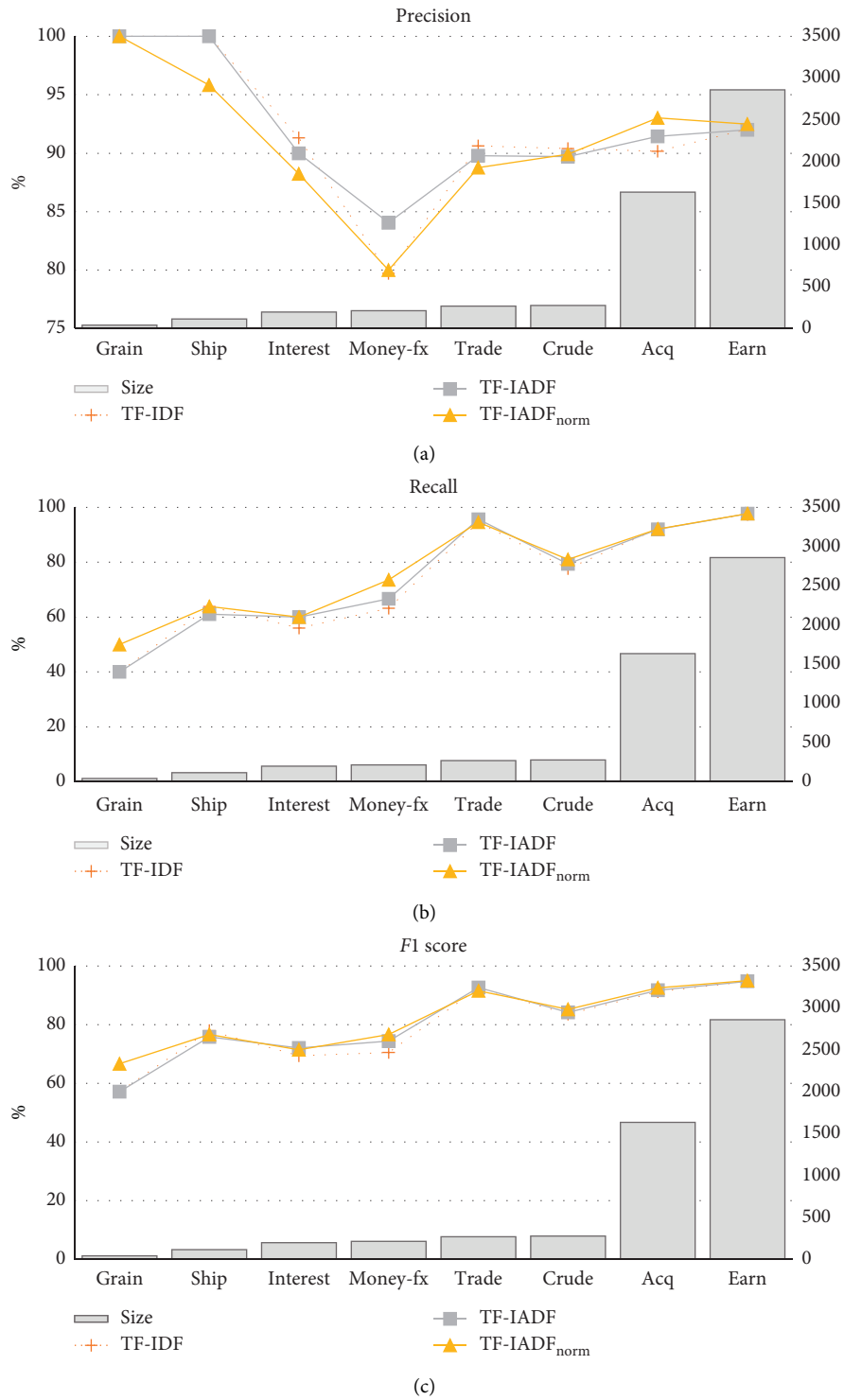


FIGURE 16: Performance of the SVM classifier on Reuters-21578 (precision, recall, and F1 score).

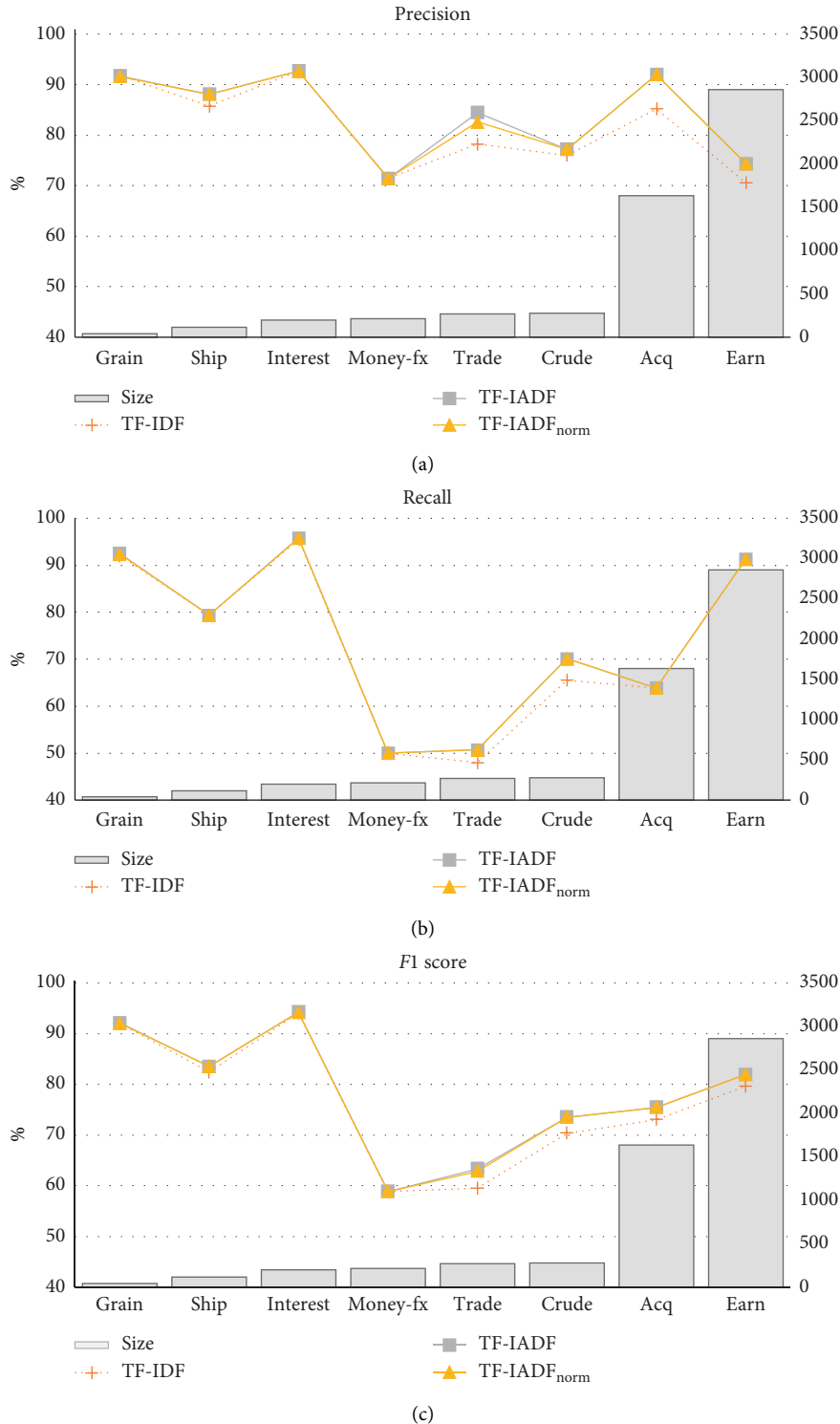


FIGURE 17: Performance of the NB classifier on Reuters-21578 (precision, recall, and *F1* score).

IADF, which has better performance in the Chinese dataset, has worse performance on this corpus. It can be seen in Table 17 and Figure 18 that the better performance of TF-IADF<sup>+</sup><sub>norm</sub> is mainly due to the higher recall score. Comparing with TF-IDF, the recall score obtained in the trade category is 6.52% higher, and in the crude category, it is 2.48% higher. In terms of precision, it has achieved good

performance in the categories of interest and trade, with an increase of 10.54% and 8.52%, respectively. On the side of *F1* score, the trade category improved significantly, up to 7.59%. However, we also see a decrease in the performance on the categories of grain and money-fx. This may be due to the small size of the test dataset, which changes greatly, and so, it is not easy to draw more accurate conclusions.

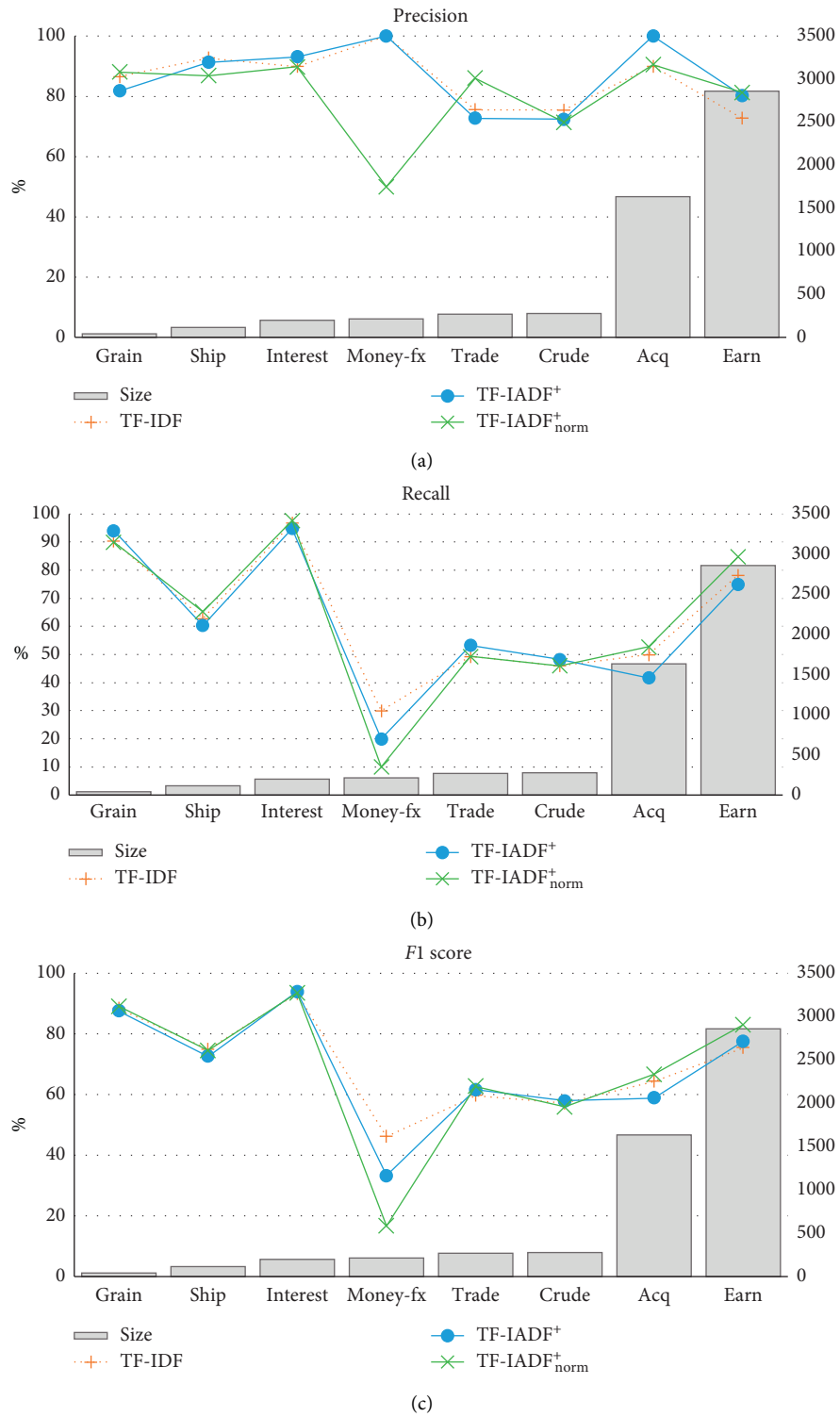


FIGURE 18: Performance of the RF classifier on Reuters-21578 (precision, recall, and F1 score).

TABLE 19: The best combinations of TWSs with classifiers.

Dataset/classifier	NB	RF	SVM
Balanced Chinese dataset	TF-IADF <sup>+</sup> <sub>norm</sub>	TF-IADF	TF-IADF
Unbalanced Chinese dataset	TF-IADF <sup>+</sup> <sub>norm</sub>	TF-IADF, TF-IADF <sub>norm</sub>	TF-IADF, TF-IADF <sub>norm</sub>
Unbalanced English dataset	TF-IADF	TF-IADF <sup>+</sup> <sub>norm</sub>	TF-IADF <sub>norm</sub>

According to all experimental results on Reuters-21578, the proposed TF-IADF<sub>norm</sub> outperformed TF-IDF in nearly all conditions except the micro-F1 score of the RF classifier. In fact, the RF classifier's performance was worse among all classifiers, suggesting that the RF classifier may not be suitable for this corpus.

**5.4. Discussion.** The experimental results obtained on the English dataset (Reuters) are somehow different from those on the Chinese dataset, and the most effective combination of a term weighting method and classification algorithm is different. For example, TF-IADF<sup>+</sup><sub>norm</sub> is suitable for the Chinese internet corpus using the NB algorithm, whereas TF-IADF performs best in the English unbalanced dataset. It can be inferred from the experimental results that all of the proposed algorithms can generally always be combined with a suitable mathematical model that shows a better performance than the original TF-IDF. For example, TF-IADF<sup>+</sup><sub>norm</sub> performs better in the RF algorithm, while TF-IADF<sub>norm</sub> performs better in the SVM. The best combinations concluded from the experiments are shown in Table 19. For the two Chinese datasets in the experiment, we draw the following conclusions: (1) for the NB algorithm, TF-IADF<sup>+</sup> and TF-IADF<sup>+</sup><sub>norm</sub> are more suitable, and TF-IADF<sup>+</sup><sub>norm</sub> can achieve better performance whether in balanced or unbalanced datasets; (2) for the RF algorithm and SVM algorithm, TF-IADF and TF-IADF<sub>norm</sub> can achieve a relatively stable improvement effect. In the experiments with the internet corpus, it can be concluded that TF-IADF has a better improvement effect when the dataset is relatively balanced, and TF-IADF<sub>norm</sub> has a better classification effect when the dataset is unbalanced. However, in the Fudan corpus (unbalanced), TF-IADF has achieved better performance than TF-IADF<sub>norm</sub>, while TF-IADF<sub>norm</sub> has also improved. It is possible that TF-IADF is more suitable for RF and SVM algorithms when there are many categories which are unbalanced in the corpus.

## 6. Conclusions

In this paper, an improved TF-IDF with novel term weighting schemes is proposed to greatly reduce the impact from the unbalanced distribution of datasets. It is easy to be observed that all unbalanced corpuses that are categorized with a larger amount of data will always have an impact on the classification effect due to the larger amount of feature words. Meanwhile, the precision score decreases, while documents from other categories will be easily mistaken into this category. This is due to the

increase in the training set; some feature words with strong performance capabilities in other categories have also appeared in this category which causes errors. Meanwhile, categories with smaller amount of data will be reflected in the significant reduction in the recall score due to the insufficient collection of feature words. The training process cannot classify these categories well without sufficient representation in the training process. As a result, many documents are assigned to wrong categories. In these cases, the proposed methods can increase the weight of those feature words which have a document frequency close to the average value, while reducing the weight of low-frequency and high-frequency words in order to obtain better results.

The simulation results show that the proposed methods with ADF are more effective than the original TF-IDF, although different mathematical models may be needed for improvement when utilizing different classification algorithms. Especially in experiments specifically designed in which the size of data in sport decreases while keeping other conditions the same, the results proved that the proposed methods are with a better performance and more stable than the well-known TF-IDF on the unbalanced corpus. It can also be concluded that our proposed methods come with a better performance in the balanced dataset when compared with TF-IDF. Document frequency of specific words may vary across categories, even in cases where the training sets appear roughly the same. Methods with ADF can weight those words more reasonably and form a more representative model.

However, for the purpose of TC on internet media reports, this paper just focuses on the term weighting scheme under unbalanced distribution but ignoring linguistic characteristics, which might be helpful in the term extraction process. Therefore, our next study will focus on enhancing the TC performance by combining the proposed methods with language characteristics.

## Data Availability

The data that support the findings of this study are available from the corresponding author upon reasonable request.

## Conflicts of Interest

The authors declare that they have no conflicts of interest.

## Acknowledgments

This research was supported by the National Natural Science Foundation of China under Grant nos. 61973024 and 61703027 and Fundamental Research Funds for the Central Universities under Grant nos. JD1914 and XK1802-4.

## References

- [1] CNNIC (China Internet Network Information Center), "The 43<sup>rd</sup> China statistical report on internet development," 2019, [http://www.cac.gov.cn/2019-02/28/c\\_1124175677.htm](http://www.cac.gov.cn/2019-02/28/c_1124175677.htm).
- [2] A. Sboev, T. Litvinova, D. Gudovskikh, R. Rybka, and I. Moloshnikov, "Machine learning models of text categorization by author gender using topic-independent features," *Procedia Computer Science*, vol. 101, pp. 135–142, 2016.
- [3] J. S. Li, L. C. Chen, J. V. Monaco, P. Singh, and C. C. Tappert, "A comparison of classifiers and features for authorship authentication of social networking messages," *Concurrency and Computation: Practice and Experience*, vol. 29, no. 14, p. e3918, 2017.
- [4] J. R. Méndez, T. R. Cotos-Yañez, and D. Ruano-Ordás, "A new semantic-based feature selection method for spam filtering," *Applied Soft Computing*, vol. 76, pp. 89–104, 2019.
- [5] B. Parlak and A. K. Uysal, "The impact of feature selection on medical document classification," in *Proceedings of the 2016 11th Iberian Conference on Information Systems and Technologies (CISTI)*, pp. 1–5, IEEE, Las Palmas, Spain, June 2016.
- [6] H. Lu, D. Zhan, L. Zhou, and D. He, "An improved focused crawler: using web page classification and link priority evaluation," *Mathematical Problems in Engineering*, vol. 2016, Article ID 6406901, 10 pages, 2016.
- [7] T. Sabbah, A. Selamat, M. H. Selamat, R. Ibrahim, and H. Fujita, "Hybridized term-weighting method for dark web classification," *Neurocomputing*, vol. 173, pp. 1908–1926, 2016.
- [8] C. C. Aggarwal and C. Zhai, "A survey of text classification algorithms," in *Mining Text Data*, pp. 163–222, Springer, Boston, MA, USA, 2012.
- [9] D. Wang and H. Zhang, "Inverse-category-frequency based supervised term weighting scheme for text categorization," 2010, <https://arxiv.org/abs/1012.2609>.
- [10] P. Soucy and G. W. Mineau, "Beyond TFIDF weighting for text categorization in the vector space model," in *Proceedings of the IJCAI*, pp. 1130–1135, Trento, Italy, July 2005.
- [11] T. Sabbah, A. Selamat, M. H. Selamat et al., "Modified frequency-based term weighting schemes for text classification," *Applied Soft Computing*, vol. 58, pp. 193–206, 2017.
- [12] N. Passalis and A. Tefas, "Entropy optimized feature-based bag-of-words representation for information retrieval," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 7, pp. 1664–1677, 2016.
- [13] M. Haddoud, A. Mokhtari, T. Lacroq, and S. Abdeddaim, "Combining supervised term-weighting metrics for SVM text classification with extended term representation," *Knowledge and Information Systems*, vol. 49, no. 3, pp. 909–931, 2016.
- [14] Z. Tang, W. Li, and Y. Li, "An improved term weighting scheme for text classification," *Concurrency and Computation: Practice and Experience*, vol. 32, p. e5604, 2019.
- [15] B. Altinel and M. C. Ganiz, "Semantic text classification: a survey of past and recent advances," *Information Processing & Management*, vol. 54, no. 6, pp. 1129–1153, 2018.
- [16] G. Rao, W. Huang, Z. Feng, and Q. Cong, "LSTM with sentence representations for document-level sentiment classification," *Neurocomputing*, vol. 308, pp. 49–57, 2018.
- [17] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013, <https://arxiv.org/abs/1301.3781>.
- [18] Q. Le and T. Mikolov, "Distributed representations of sentences and documents," in *Proceedings of the International Conference on Machine Learning*, pp. 1188–1196, Beijing, China, January 2014.
- [19] H. K. Kim, H. Kim, and S. Cho, "Bag-of-concepts: comprehending document representation through clustering words in distributed representation," *Neurocomputing*, vol. 266, pp. 336–352, 2017.
- [20] D. Kim, D. Seo, S. Cho, and P. Kang, "Multi-co-training for document classification using various document representations: TF-IDF, LDA, and Doc2Vec," *Information Sciences*, vol. 477, pp. 15–29, 2019.
- [21] Y. Wu, S. Zhao, and W. Li, "Phrase2Vec: phrase embedding based on parsing," *Information Sciences*, vol. 517, pp. 100–127, 2020.
- [22] J. Kim, S. Jang, E. Park, and S. Choi, "Text classification using capsules," *Neurocomputing*, vol. 376, pp. 214–221, 2020.
- [23] Y.-W. Chen, Q. Zhou, W. Luo, and J.-X. Du, "Classification of Chinese texts based on recognition of semantic topics," *Cognitive Computation*, vol. 8, no. 1, pp. 114–124, 2016.
- [24] A. Elnagar, R. Al-Debsi, and O. Einea, "Arabic text classification using deep learning models," *Information Processing & Management*, vol. 57, no. 1, Article ID 102121, 2020.
- [25] K. Chen, Z. Zhang, J. Long, and H. Zhang, "Turning from TF-IDF to TF-IGM for term weighting in text classification," *Expert Systems with Applications*, vol. 66, pp. 245–260, 2016.
- [26] K. S. Jones, "A statistical interpretation of term specificity and its application in retrieval," *Journal of Documentation*, vol. 28, 1972.
- [27] F. Debole and F. Sebastiani, "Supervised term weighting for automated text categorization," in *Text Mining and its Applications*, pp. 81–97, Springer, Berlin, Germany, 2004.
- [28] T. Dogan and A. K. Uysal, "Improved inverse gravity moment term weighting for text classification," *Expert Systems with Applications*, vol. 130, pp. 45–59, 2019.
- [29] F. Ren and C. Li, "Hybrid Chinese text classification approach using general knowledge from baidu baike," *IEEE Transactions on Electrical and Electronic Engineering*, vol. 11, no. 4, pp. 488–498, 2016.
- [30] R. Longadge and S. Dongre, "Class imbalance problem in data mining review," 2013, <https://arxiv.org/abs/1305.1707>.
- [31] M. N. Anwar, *Complexity Measurement for Dealing with Class Imbalance Problems in Classification Modelling*, Massey University, Palmerston, New Zealand, 2012.
- [32] M. Buda, A. Maki, and M. A. Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural Networks*, vol. 106, pp. 249–259, 2018.
- [33] J. M. Johnson and T. M. Khoshgoftaar, "Survey on deep learning with class imbalance," *Journal of Big Data*, vol. 6, no. 1, p. 27, 2019.
- [34] F. Thabtah, S. Hammoud, F. Kamalov, and A. Gonsalves, "Data imbalance in classification: experimental evaluation," *Information Sciences*, vol. 513, pp. 429–441, 2020.
- [35] A. Onan, "Sentiment analysis on product reviews based on weighted word embeddings and deep neural networks," *Concurrency and Computation: Practice and Experience*, p. e5909, 2020.
- [36] M. Lan, C. L. Tan, J. Su, and Y. Lu, "Supervised and traditional term weighting methods for automatic text categorization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 4, pp. 721–735, 2008.
- [37] M. Lichman, *UCI Machine Learning Repository*, University of California, School of Information and Computer Science, Irvine, CA, USA, 2017.

- [38] A. Onan, S. Korukoğlu, and H. Bulut, "Ensemble of keyword extraction methods and classifiers in text classification," *Expert Systems with Applications*, vol. 57, pp. 232–247, 2016.
- [39] D. M. Farid, L. Zhang, C. M. Rahman, M. A. Hossain, and R. Strachan, "Hybrid decision tree and naive bayes classifiers for multi-class classification tasks," *Expert Systems with Applications*, vol. 41, no. 4, pp. 1937–1946, 2014.
- [40] G. Bianchi, R. Bruni, and F. Scalfati, "Identifying e-commerce in enterprises by means of text mining and classification algorithms," *Mathematical Problems in Engineering*, vol. 2018, Article ID 7231920, 8 pages, 2018.
- [41] Z. Geng, Y. Zhang, and Y. Han, "Joint entity and relation extraction model based on rich semantics," *Neurocomputing*, vol. 429, 2020.
- [42] X. Wu, Y. Gao, and D. Jiao, "Multi-label classification based on random forest algorithm for non-intrusive load monitoring system," *Processes*, vol. 7, no. 6, p. 337, 2019.
- [43] Q. Wu, Y. Ye, H. Zhang, M. K. Ng, and S.-S. Ho, "Foretexter: an efficient random forest algorithm for imbalanced text categorization," *Knowledge-Based Systems*, vol. 67, pp. 105–116, 2014.
- [44] A. Onan, "Two-stage topic extraction model for bibliometric data analysis based on word embeddings and clustering," *IEEE Access*, vol. 7, pp. 145614–145633, 2019.