

7-26-2024

A new dynamic classifier selection method for text classification

İSMAİL TERZİ

ALPER KÜRŞAT UYSAL

Follow this and additional works at: <https://journals.tubitak.gov.tr/elektrik>



Part of the [Computer Engineering Commons](#), [Computer Sciences Commons](#), and the [Electrical and Computer Engineering Commons](#)

Recommended Citation

TERZİ, İSMAİL and UYSAL, ALPER KÜRŞAT (2024) "A new dynamic classifier selection method for text classification," *Turkish Journal of Electrical Engineering and Computer Sciences*: Vol. 32: No. 4, Article 10. <https://doi.org/10.55730/1300-0632.4092>

Available at: <https://journals.tubitak.gov.tr/elektrik/vol32/iss4/10>



This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

This Research Article is brought to you for free and open access by TÜBİTAK Academic Journals. It has been accepted for inclusion in Turkish Journal of Electrical Engineering and Computer Sciences by an authorized editor of TÜBİTAK Academic Journals. For more information, please contact pinar.dundar@tubitak.gov.tr.

A new dynamic classifier selection method for text classification

İsmail TERZİ¹ , Alper Kürşat UYSAL^{2,*} 

¹Department of Computer Engineering, Faculty of Engineering, Zonguldak Bülent Ecevit University, Zonguldak, Türkiye

²Department of Computer Engineering, Faculty of Engineering, Alanya Alaaddin Keykubat University, Antalya, Türkiye

Received: 25.08.2023

Accepted/Published Online: 11.05.2024

Final Version: 26.07.2024

Abstract: The primary objective of employing multiple classifier systems (MCS) in pattern recognition is to enhance classification accuracy. Dynamic classifier selection (DCS) and dynamic ensemble selection (DES) are two purposeful forms of multiple classifier systems. While DES involves the selection of a classifier set followed by decision combination, DCS opts for the choice of a single competent classifier, eliminating the necessity for classifier combination. As a consequence, DCS methods exhibit superior efficiency in terms of processing time and memory usage compared to DES methods. Moreover, a substantial performance gap exists between the performance of Oracle and both DES and DCS methods. In this study, we introduce a novel method termed dynamic classifier selection technique-decision quotient (DCS-DQ) for text classification based on dynamic classifier selection. Our experimental investigation encompasses four distinct text datasets, with classification accuracy and macro F1-score serving as the primary evaluation criteria. The proposed DCS-DQ method is subjected to comparison with seven state-of-the-art DCS methods. Based on our empirical findings, the DCS-DQ method outperforms the other seven DCS methods in terms of classification accuracy across the majority of feature sizes. Notably, in the Reuters dataset, the classification accuracy of DCS-DQ surpasses that of other DCS methods for all feature sizes except when the feature size is 100. Similarly, in the Ohsumed dataset, the DCS-DQ method demonstrates significant performance improvement, with an accuracy value of 77.02% for 3000 features compared to the maximum accuracy value of 72.74% achieved by the DCS method MCB. Additionally, the performance of the proposed DCS-DQ method closely aligns with the oracle performance compared to the other methods. In conclusion, our proposed DCS-DQ method holds promise for significantly improving classification accuracy in text classification literature.

Key words: Text classification, dynamic classifier selection, multiple classifier systems, DCS-DQ

1. Introduction

The volume of text documents in the databases of companies and on the Internet is increasing day by day. Consequently, there has been a growing inclination towards the field of text classification. E-mail messages, articles on web pages, research articles, tweets, medical reports, customer correspondence, blogs, customer reviews on shopping sites are composed of text messages. People not only save documents but also discover some useful patterns within them. Since such an amount of text data is overwhelming for individuals to analyze, they need useful tools to deal with such number of text documents. Classification of text documents is the process of determining classes for text documents based on their content. The foremost objective within the process of text classification is assigning the text document to the appropriate class. Numerous machine learning methods

*Correspondence: alper.uysal@alanya.edu.tr

have been employed in the field of text classification so far, including; naïve Bayes (NB), logistic regression (LR), support vector machines (SVM), random forest (RF), k-nearest neighbor (kNN), decision tree (DT) classifiers, and Rocchio algorithm (RA) [1]. Text classification techniques are utilized in many applications to simplify people's lives. Document categorization [2], document routing application [3], author recognition [4], opinion mining and sentiment analysis [5, 6], question answering systems [7], and detection of spam SMS messages and social spam [8] were performed using text classification techniques. Ensembles of classifiers represent a widely discussed area in the domain of machine learning. According to Dietterich [9], ensembles of classifiers are the leading research direction in machine learning and they can improve the accuracy of classification. In the literature, authors frequently refer to the ensembles of classifiers as multiple classifier systems [10]. In this work, we use multiple classifier systems (MCS). MCS have better performance than traditional single classifier systems [9, 10]. MCS, as illustrated in Figure 1, consists of three sequential components: the initial phase entails generation, followed by selection, and culminating in integration. The selection component is categorized into two distinct groups: dynamic selection (DS) and static selection (SS). In DS, there are two approaches, dynamic ensemble selection (DES) and dynamic classifier selection (DCS). SS methods employ a singular classifier or an ensemble of classifiers during the training phase, subsequently using the identical selected classifiers to predict outcomes for all unknown samples. DS methods select only a single classifier or a combination of classifiers for every unknown samples. As DCS methods employ a single classifier, the requirement for an integration phase is obviated. The absence of a requirement for integration renders DCS methods more efficient than DES methods in processing time and memory usage.

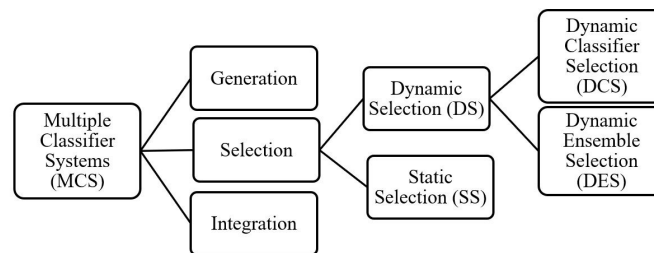


Figure 1. Multiple classifier systems.

Oracle [11] stands as another significant concept within the domain of DCS. Oracle is an abstract method used to identify the classifier that can accurately classify the text instance among the available classifiers, given the existence of such a classifier. The key aspect of Oracle is that at least one of the classifiers in the pool should be capable of correctly classifying the unknown sample. The performance of Oracle serves as the upper bound for DCS methods [12] and there is a substantial performance gap between Oracle performance and DCS methods. In this study, the proposed DCS-DQ method can contribute to closing these gaps. Performance gap between the existing methods and oracle performance is shown in Figure 2 on two different datasets. When the number of features is 3000, the classification accuracy of the DCS method is 80%. However, Oracle performance is 97% with the same number of features. The discrepancy is significant. Studies aimed at addressing this gap will yield significant contributions to the DCS literature. In all datasets, a substantial performance gap is observable between existing methods and the performance exhibited by the Oracle. In the following sections, we will demonstrate how our proposed method, namely DCS-DQ, effectively narrows this significant gap.

The main motivation of this work is to:

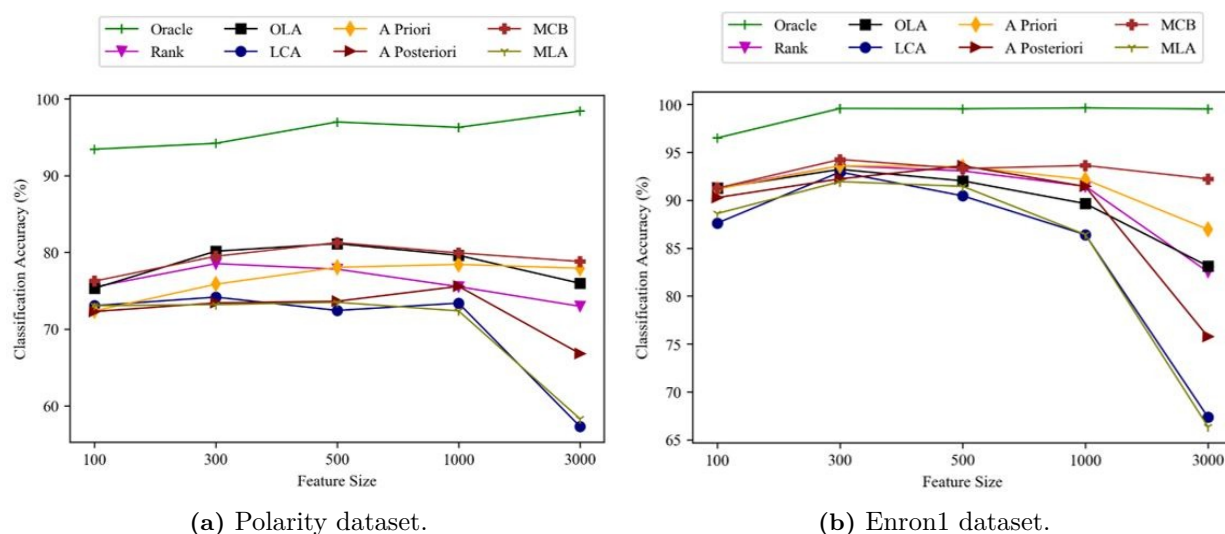


Figure 2. Oracle performance for Polarity and Enron1 datasets.

1. Propose a new DCS method namely DCS-DQ for text classification.
2. Make contribution for reducing the gap between DCS and Oracle performance.
3. Show that the proposed DCS-DQ outperforms current state-of-the-art techniques.
4. Analyze the effectiveness of proposed DCS-DQ method on text datasets with different characteristics.

The organization of this paper is as follows: Works related to applications using multiple classification systems are explained in Section 2, multiple classifier systems are briefly explained in Section 3, existing DCS methods in the literature is given in Section 4, the proposed method DCS-DQ is presented in Section 5, experimental studies are presented in Section 6, the experimental findings are given in Section 7, and conclusion about our work is presented in Section 8.

2. Related works

In the literature, a significant number of articles utilizing DCS methods have been published recently. These methods are being applied to a wide range of real-world problems. Credit scoring [21], face recognition systems [22], biometric verification [23], signature verification [24], and customer classification [25] are applications of DCS methods. Wen et al. utilized dynamic classifier selection techniques to identify the ball screw degradation [26]. Groccia et al. [27] introduced a framework that integrates several classification algorithms by dynamically selecting the most proficient classifier. In the literature, many classification datasets are unbalanced and proposing techniques for unbalanced data is more valuable. Roy et al. [28] have experimentally shown that dynamic selection methods have the potential to achieve superior performance compared to static ensembles in unbalanced classification problems. In experimental studies, they used DS methods, LCA, and Rank. With the rise in the use of Android smartphones today, the number of malicious applications posing threats to the Android platform has also increased. Feng et al. [29] proposed an ensemble-based Android malware detection method called EnDroid to protect Android platforms from malware. Various machine learning algorithms were employed, including NB, kNN, SVM, DT, boosted tree, and RF. Credit scoring is another critical issue for

financial institutions. Junior et al. [30] investigated the suitability of dynamic selection techniques for credit scoring and introduced the reduced minority k-nearest neighbors (RMkNN) method. Their proposed approach improves the delineation of local regions in dynamic selection techniques for imbalanced credit scoring datasets. Another ensemble based study on credit scoring was published by Feng et al. [31]. In their study, a dynamically weighted ensemble method is proposed for credit scoring. They used a Markov Chain to dynamically weight the classifiers in the classifier pool for each sample in the test set and combine the classifiers' decisions. Martins et al. [32] published a research on forest species recognition. In this research, they used dynamic classifier selection methods such as MCB, OLA, LCA, A Priori, and A Posteriori. The best result in this work is 93.03%. The best result is observed when integrating probabilistic information into a dynamic classifier selection method based on MCB. DS techniques are used in time series forecasting Sergio et al. [35] proposed a dynamic selection of regressors for time series forecasting. The authors developed an algorithm inspired by the dynamic classifier selection method MCB to predict the competence of each combiner. The technique, termed dynamic selection of forecast combiners (DS-FC), is a heuristic approach designed to choose an optimal ensemble from a provided pool of classifiers [36]. The proposed algorithm is a pruning algorithm based on accuracy and diversity. It evaluates both the accuracy of individual classifiers and the pairwise diversity among them. Cruz et al. [37] showed that DCS methods offer a substantial increase in classification accuracy compared to kNN. Dogo et al. [38] used dynamic selection techniques, including Rank, LCA, and OLA to address the issue of water quality anomaly detection problem. Groccia et al. used dynamic classifier selection techniques for clinical diagnosis [39]. Text data is frequently used in the medical field, and sorting medical texts into clinical texts, clinical notes, prescriptions and examination requests is an important task. Sousa et al. [33] analyzed the success of classifier ensemble approaches in classifying medical texts. In their analysis, they obtained results that can automatically and accurately classify clinical texts with higher accuracy than individual approaches. There are also studies that use dynamic selection methods in feature selection. Li et al. [34] proposed a dynamic feature selection method for extracting semantic features from agricultural texts.

3. Multiple classifier systems (MCS)

Throughout this paper the following notations are used.

- $D = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_m, y_m)\}$ is a text dataset containing m documents and $(x_1, x_2, x_3, \dots, x_m)$ are documents, $(y_1, y_2, y_3, \dots, y_m)$ are corresponding class labels of $(x_1, x_2, x_3, \dots, x_m)$.
- $\Omega = (\omega_1, \omega_2, \omega_3, \dots, \omega_n)$ are the corresponding class labels and $y_m \in \Omega$.
- x_i is a sample text document with unknown class label in test set.
- \check{D} is the dynamic selection data (DSEL).
- $P = (C_1, C_2, C_3 \dots C_T)$ is the pool composed number of T base classifiers and C_i is the most competent classifier for x_i selected by DCS method.
- $\delta = (\partial_{C_1}, \partial_{C_2}, \partial_{C_3}, \dots, \partial_{C_T})$ is the set of accuracy value of the base classifiers on \check{D} and ∂_{C_t} is the accuracy value of a base classifier C_t on \check{D} .
- $\theta_{x_i} = (\dot{x}_1, \dot{x}_2, \dot{x}_3 \dots \dot{x}_k)$ is the k nearest neighbor of x_i in \check{D} and \dot{x}_k is a neighbor of x_i in θ_{x_i} .

- θ_{x_i} ; competence level of a base classifier C_t for x_i .
- $\Psi_{x_i} = \{\sigma_{c_1}, \sigma_{c_2}, \sigma_{c_3}, \dots, \sigma_{c_T}\}$ is the set of accuracy value of the base classifiers on θ_{x_i} and σ_{c_t} is the accuracy value of a base classifier C_t on θ_{x_i} .
- $W_i = \frac{1}{d_i}$; d_i is distance between x_i and \hat{x}_i .
- \tilde{x}_i ; output profile of x_i

MCS systems encompass three primary stages. The first step is creation of a classifier pool P, the second step is selection of a classifier C_i or a subset of classifiers C' , and the last step is the combination of the classification results of various classifiers [13]. DCS methods do not have a combination step. The absence of a combination step can be considered one of the strengths of DCS methods. The pool of classifiers, denoted as P, can be generated in six distinct methods [20]. These involve different initialization, different feature sets, different parameters, different classifier models, different architectures, and different training sets. Additionally, bagging and boosting methods can be utilized for creating the pool [14]. Different subset of training set is used in bagging method [15]. After creating a pool of classifiers, the selection process is performed. SS methods select one competent classifier or a set of classifiers which is also named as ensemble of classifiers at the training stage and anticipate all unknown samples x_i by using the same classifier or a set of classifiers. The most competent classifier for x_i is the classifier that classifies all the samples in θ_{x_i} with the highest accuracy. In the DCS strategy, base classifier C_i is selected on the fly. Different base classifiers are selected for each x_i . In order to yield more accurate results for the multiple classifier system, each classifier in the pool must be accurate and diverse [16]. A classifier can be considered accurate if it yields an error rate lower than that of random guessing for unknown test samples x_i . If two classifiers exhibit different errors on the same test sample x_i , then these two classifiers can be deemed diverse [17].

The key idea of using DCS technique is that the competence of each base classifier must be determined. The performance of DCS methods is very dependent on the detection of this competence [18]. The rationale behind this explanation is that each base classifier specializes in a distinct region of the feature space [13]. Determining θ_{x_i} for a given sample x_i is the fundamental concept. Once the region is determined, the classifier C_i that has the highest classification accuracy on θ_{x_i} is chosen. k-nearest neighbors technique is mostly used to determine this region [19]. Given a test instance x_i , DCS methods select the most capable classifier C_i . C_i assigns the instance x_i to a class ω_l . The diagram depicting the operation of the DCS method is presented in Figure 3 below.

As shown in Figure 3, vital task of DCS methods is to reveal the most capable classifier C_i for each x_i . DCS methods use a labeled validation or dynamic selection data (DSEL) [20] when deciding which classifier to be chosen. The DSEL data is initially segregated from the dataset and does not overlap with the test and train data.

4. DCS methods

In this study, proposed DCS method has been compared with the most commonly used DCS methods in the literature. In this section, DCS methods will be briefly explained. In our experimental study, these methods were employed with their default parameters, as described in [44].

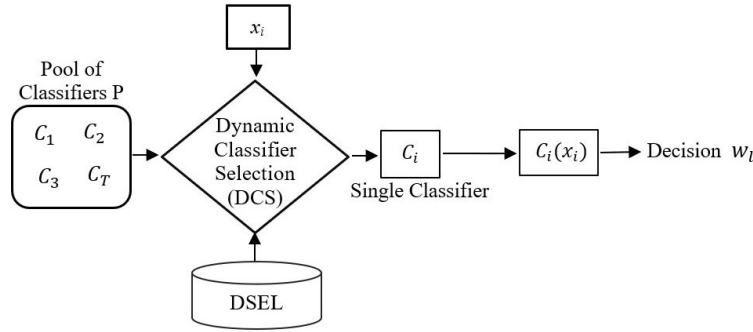


Figure 3. Dynamic classifier selection (DCS).

4.1. Modified classifier rank (DCS-Rank):

In this approach, the ranking of a base classifier C_t is determined by the consecutive correct classifications within the neighborhood θ_{x_i} . "Consecutive" in this context refers to classifications from the nearest neighbor to the farthest one. A classifier C_t that correctly classifies the greatest number of consecutive nearest samples is assigned the highest "rank". Consequently, the classifier with the highest rank is selected, and x_i is classified accordingly [40].

4.2. Overall local accuracy (OLA):

In this method, the classifier C_t that has the highest classification accuracy on θ_{x_i} is the most competent classifier; therefore, x_i is classified by C_t . Classifier competency $\emptyset_{t_{x_i}}$ is calculated by Equation (1) [40].

$$\emptyset_{t_{x_i}} = \frac{1}{k} \sum_{i=1}^k P(\omega_l | \hat{x}_i \in \omega_l, C_t) \quad (1)$$

4.3. Local classifier accuracy (LCA):

In this method, first, the θ_{x_i} of the test sample x_i is formed. Following this step, the proficiency of the base classifier C_t is determined based on its classification accuracy, focusing solely on the samples from the class ω_l within this neighborhood. Here, ω_l represents the class predicted by the base classifier C_t , for x_i . The competency level $\emptyset_{t_{x_i}}$ is assessed using Equation (2). If multiple base classifiers achieve identical competency levels, the first one encountered is chosen [40].

$$\emptyset_{t_{x_i}} = \frac{\sum_{\hat{x}_i \in \omega_l} P(\omega_l | \hat{x}_i, C_t)}{\sum_{i=1}^k P(\omega_l | \hat{x}_i, C_t)} \quad (2)$$

4.4. A Priori:

$\emptyset_{t_{x_i}}$ is predicted based on the probability of true classification of the base classifier C_t , taking into account all samples in θ_{x_i} . This method weights the impact of each training sample as per its distance W_i to x_i . Competence level of a classifier $\emptyset_{t_{x_i}}$ is calculated by Equation (3). If multiple base classifiers achieve identical competency levels, the first one encountered is chosen [41].

$$\emptyset_{t_{x_i}} = \frac{\sum_{i=1}^k P(\omega_l | \dot{x}_i \in \omega_l, C_t) W_i}{\sum_{i=1}^k W_i} \quad (3)$$

4.5. A Posteriori:

$\emptyset_{t_{x_i}}$ is estimated based on the probability of correct classification of the base classifier C_t , for each neighbor \dot{x}_k belonging to a specific class ω_l . In this case, ω_l is the class predicted by the base classifier C_t , for x_i . This method weights the impact of each sample in θ_{x_i} according to its distance W_i to x_i . $\emptyset_{t_{x_i}}$ is calculated by Equation (4) [41].

$$\emptyset_{t_{x_i}} = \frac{\sum_{\dot{x}_i \in \omega_l} P(\omega_l | \dot{x}_i \in \omega_l, C_t) W_i}{\sum_{i=1}^k P(\omega_l | \dot{x}_i \in \omega_l, C_t) W_i} \quad (4)$$

4.6. Multiple classifier behavior (MCB):

The region of competence is estimated considering the feature space and the decision space (using the behavior knowledge space (BKS) method [52]). First, θ_{x_i} is formed for x_i . Then, the similarity in the behavior knowledge space between x_i and θ_{x_i} are estimated using Equation (5),

$$S(\tilde{x}_i, \tilde{\tilde{x}}_i) = \frac{1}{M} \sum_{i=1}^M T(x_i, \dot{x}); \quad T(x_i, \dot{x}) = \begin{cases} 1 & \text{if } C_t(x_i) = C_t(\dot{x}_i) \\ 0 & \text{if } C_t(x_i) \neq C_t(\dot{x}_i) \end{cases} \quad (5)$$

where $S(\tilde{x}_i, \tilde{\tilde{x}}_i)$ denotes the similarity between x_i and \dot{x}_i according to BKS method. M represents the number of base classifiers in the classifier Pool. Instances with a similarity below a predetermined threshold are excluded from the θ_{x_i} . The competence level of the base classifiers $\emptyset_{t_{x_i}}$ is predicted based on their classification accuracy within the last region of competence θ_{x_i} [41].

4.7. Modified local accuracy (MLA):

The competence level $\emptyset_{t_{x_i}}$ of C_t is determined according to its classification accuracy, considering solely the samples associated with a specific class ω_l . Here, ω_l denotes the class predicted by the base classifier C_t , for x_i . This approach evaluates the significance of each training sample based on its proximity to the query instance. The competence level of a classifier $\emptyset_{t_{x_i}}$ is computed using Equation (6) [42].

$$\emptyset_{t_{x_i}} = \sum_{i=1}^k P(\omega_l | \dot{x}_i \in \omega_l, C_t) W_i \quad (6)$$

5. Proposed method

In this study, DCS methods that are prominent in terms of their popularity and number of citations have been selected. The main shortcoming of the existing DCS methods is the way of deciding the competence level of a base classifier. Existing DCS methods decide competence level of the base classifiers by using k-nearest neighbors of unknown sample x_i in \check{D} . A base classifier possessing the highest level of competence within the k-nearest neighbors of an unknown sample in \check{D} is selected and the selected classifier classifies x_i . Most of the

time it is challenging to find *k-nearest neighbors* with a high degree of similarity for x_i since datasets are very sparse [43]. Deciding the competence level of the base classifiers using k-unlike neighbors is often inaccurate. The most important limitation of existing methods is that an example to be classified is classified by one of the classifiers selected from the pool even if it has no similar neighbors.

The proposed method has been developed in order to eliminate the mentioned shortcomings of the existing methods. The basic principle of the proposed method is the k-nearest neighbor (k-NN) algorithm, which is a widely recognized machine learning technique. The basic philosophy underlying the k-NN algorithm is that any test instance is similar to the nearest instance whose label is known in training set. With the same inference:

- *Claim1* : Any test instance x_i is similar to the nearest samples in DSEL(\check{D}) data. This similarity can also be used to infer whether the test instance is easily classifiable or hardly classifiable.
- *Claim2* : If every classifier in the classifier pool can classify the nearest neighbors of x_i with high accuracy, then x_i is said to be easily classified. Conversely, if the nearest neighbors of the test sample x_i are hard to classify, then x_i will also hard to be classified.

Figure 4 below presents the proposed DCS-DQ method, where $(x_1, x_2, x_3, \dots, x_n)$ are test documents in test data. $\theta_{x_i} = (\check{x}_1, \check{x}_2, \check{x}_3, \dots, \check{x}_k)$ are the k nearest neighbors for x_i in DSEL(\check{D}), P is the classifier pool. Most important part of the proposed DCS-DQ method is to construct a $(\theta_{x_i} \times P)$ Matrix.

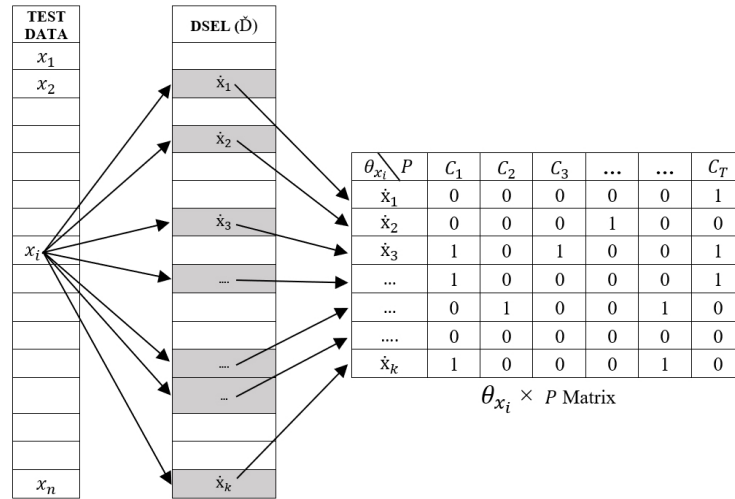


Figure 4. Proposed dynamic classifier selection method DCS-DQ.

For each x_i in test data, a $(\theta_{x_i} \times P)$ matrix is constructed. Rows of $(\theta_{x_i} \times P)$ matrix represent each neighbor \check{x}_i for x_i in \check{D} . Columns of $(\theta_{x_i} \times P)$ matrix represent number of T base classifiers $(C_1, C_2, C_3, \dots, C_T)$ in the pool P. Cell of $(\theta_{x_i} \times P)$ matrix represents decision of C_t for \check{x}_i . Value of the cell is equal to 1 if $\check{x}_i \in w_n$ and $C_t(\check{x}_i) \in w_n$; otherwise, it is 0. The function $\varphi(C_t, \check{x}_i)$ in Equation (7) determines the values of the cells in $(\theta_{x_i} \times P)$ matrix.

$$\varphi(C_t, \check{x}_i) = \begin{cases} 1 & \text{if } \check{x}_i \in w_n \text{ and } C_t(\check{x}_i) \in w_n \\ 0 & \text{if } \check{x}_i \in w_n \text{ and } C_t(\check{x}_i) \notin w_n \end{cases} \quad (7)$$

In Equation (8), a function $\lambda(\theta_{x_i})$ is defined. $\lambda(\theta_{x_i})$ gives the summation of the values of all the cells in $(\theta_{x_i} \times P)$ matrix.

$$\lambda(\theta_{x_i}) = \sum_{i=1}^k \sum_{t=1}^T \varphi(C_t, \dot{x}_i) \quad (8)$$

Let T be the number of base classifiers in P and the number of nearest neighbors for x_i in \check{D} is k , then maximum value of the function $\varphi(C_t, \dot{x}_i)$ is equal to $T * k$. This situation is possible if all base classifiers correctly classify all samples in θ_{x_i} . Minimum value of the function $\lambda(\theta_{x_i})$ is equal to 0. This situation is possible if none of the base classifiers can correctly classify any of the samples in θ_{x_i} .

Decision quotient (DQ) of a test instance x_i is calculated by Equation (9).

$$DQ(x_i) = \frac{\lambda(\theta_{x_i})}{T * k} \quad (9)$$

Range of the value of $DQ(x_i)$ is between 0 and 1. The value closer to 1 means that x_i is easy to classify since its neighbors are easily classified. The value closer to 0 means that x_i is hard to classify since its neighbors are hardly classified. $DQ(x_i)$ varies between 0 and 1. It is important to define a threshold for making a decision about a test instance x_i . This decision is about if a test instance x_i is hard to classify or easy to classify (Equation 10). During experimental study, this threshold can be taken as $\min(\partial_{C_1}, \partial_{C_2}, \partial_{C_3}, \dots, \partial_{C_T})$. As stated before, $\delta = (\partial_{C_1}, \partial_{C_2}, \partial_{C_3}, \dots, \partial_{C_T})$ are accuracy values of the base classifiers on \check{D} . Accuracy values of $(\partial_{C_1}, \partial_{C_2}, \partial_{C_3}, \dots, \partial_{C_T})$ are also between 0 and 1.

$$x_i = \begin{cases} \text{easy to classify} & \text{if } DQ(x_i) \geq \min(\delta) \\ \text{hard to classify} & \text{if } DQ(x_i) < \min(\delta) \end{cases} \quad (10)$$

$\Psi'_{x_i} = \{\sigma_{c_1}, \sigma_{c_2}, \sigma_{c_3}, \dots, \sigma_{c_T}\}$ are the accuracy values of the base classifiers on θ_{x_i} , $\delta = (\partial_{C_1}, \partial_{C_2}, \partial_{C_3}, \dots, \partial_{C_T})$ are accuracy values of the base classifiers on \check{D} . In Equation (11), we define $\mathcal{L}G_{x_i} = \Psi_{x_i} + \delta$ which are local and global accuracies of the base classifiers for x_i . \mathcal{L} stands for local and \mathcal{G} stands for global. Local accuracy is the accuracy of the base classifiers on θ_{x_i} , and global accuracy is the accuracy of the base classifiers on \check{D} .

$$\mathcal{L}G_{x_i} = \Psi_{x_i} + \delta = (\sigma_{c_1} + \partial_{C_1}, \sigma_{c_2} + \partial_{C_2}, \sigma_{c_3} + \partial_{C_3}, \dots, \sigma_{c_T} + \partial_{C_T}) \quad (11)$$

There are two ways to select the most capable classifiers:

1. If x_i is easy to classify, then an arbitrary classifier from P can classify it. However, in this situation, we select the classifier that classifies DSEL(\check{D}) with the highest accuracy. This classifier is called the single best [13].
2. If x_i is hard to classify, then an arbitrary classifier from P cannot classify it easily. In this situation, the classification accuracies of the base classifiers on \check{D} and θ_{x_i} are summed, i.e. as stated in Equation (11). $\max(\mathcal{L}G_{x_i})$ is calculated, and the base classifier C_i that achieves the highest value is selected.

Figure 5 depicts the DCS-DQ method. However, an example scenario for DCS-DQ method is also presented below.

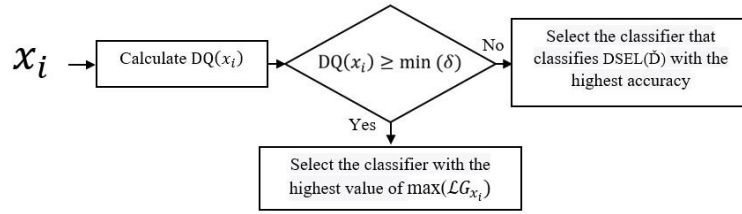


Figure 5. DCS-DQ method.

An example for DCS-DQ method: In this scenario, we have four different test items, namely $x_1, x_2, x_3,$ and x_4 and five base classifiers, namely $C_1, C_2, C_3, C_4,$ and C_5 in classifier pool P. We set $k = 7$ and $T = 5$ for all x_1, x_2, x_3, x_4

Step 1: Form $(\theta_{x_i} \times P)$ matrices for x_1, x_2, x_3, x_4 . These matrices are given in Figures 6a–6d.

$\theta_{x_1} \backslash P$	C_1	C_2	C_3	C_4	C_5	Total
\dot{x}_1	1	1	1	1	1	5
\dot{x}_2	1	1	1	1	1	5
\dot{x}_3	1	1	1	1	0	4
\dot{x}_4	1	1	1	1	1	5
\dot{x}_5	1	1	1	1	1	5
\dot{x}_6	1	1	1	1	1	5
\dot{x}_7	1	1	0	1	1	4
Total	7	7	6	7	6	33

(a) $(\theta_{x_1} \times P)$ Matrix

$\theta_{x_2} \backslash P$	C_1	C_2	C_3	C_4	C_5	Total
\dot{x}_1	1	1	0	1	0	3
\dot{x}_2	1	1	1	1	1	5
\dot{x}_3	1	1	1	1	0	4
\dot{x}_4	1	0	1	1	1	4
\dot{x}_5	0	0	1	1	1	3
\dot{x}_6	1	1	1	1	1	5
\dot{x}_7	1	1	1	0	1	4
Total	6	5	6	6	5	28

(b) $(\theta_{x_2} \times P)$ Matrix

$\theta_{x_3} \backslash P$	C_1	C_2	C_3	C_4	C_5	Total
\dot{x}_1	1	0	1	0	1	3
\dot{x}_2	0	0	0	1	0	1
\dot{x}_3	0	1	0	1	0	2
\dot{x}_4	0	0	1	0	1	2
\dot{x}_5	1	1	0	1	0	3
\dot{x}_6	0	1	0	0	1	2
\dot{x}_7	1	0	0	0	1	2
Total	3	3	2	3	4	15

(c) $(\theta_{x_3} \times P)$ Matrix

$\theta_{x_4} \backslash P$	C_1	C_2	C_3	C_4	C_5	Total
\dot{x}_1	0	0	0	0	0	0
\dot{x}_2	0	0	0	1	0	1
\dot{x}_3	0	1	0	1	0	2
\dot{x}_4	0	0	0	0	1	1
\dot{x}_5	1	0	0	0	0	1
\dot{x}_6	0	1	0	0	1	2
\dot{x}_7	0	0	0	0	1	1
Total	1	2	0	2	3	8

(d) $(\theta_{x_4} \times P)$ Matrix

Figure 6. $(\theta_{x_i} \times P)$ matrices (a, b, c, and d).

Step 2: Calculate $\lambda(\theta_{x_i})$.

$$\lambda(\theta_{x_1}) = \sum_{i=1}^7 \sum_{t=1}^5 \varphi(C_t, \dot{x}_i) = 33, \quad \lambda(\theta_{x_2}) = \sum_{i=1}^7 \sum_{t=1}^5 \varphi(C_t, \dot{x}_i) = 28$$

$$\lambda(\theta_{x_3}) = \sum_{i=1}^7 \sum_{t=1}^5 \varphi(C_t, \dot{x}_i) = 15, \quad \lambda(\theta_{x_4}) = \sum_{i=1}^7 \sum_{t=1}^5 \varphi(C_t, \dot{x}_i) = 8$$

Step 3: Calculate DQ (x_i) .

$$DQ(x_1) = \frac{33}{35} = 0.94, \quad DQ(x_2) = \frac{28}{35} = 0.80, \quad DQ(x_3) = \frac{15}{35} = 0.43, \quad DQ(x_4) = \frac{8}{35} = 0.22$$

Step 4: Determine if x_i is easy or hard to classify. Suppose the accuracy value of the base classifiers $C_1, C_2, C_3, C_4,$ and C_5 is $\delta = (0.75, 0.90, 0.85, 0.95, 0.85)$. In this case, $\min(\delta) = 0.75$. This means that 0.75 is the minimum classification accuracy of the base classifiers on D. As a result x_i with $DQ(x_i) \geq 0.75$ is easy to classify; on the contrary, x_i with $DQ(x_i) < 0.75$ is hard to classify. In our example, $DQ(x_1) = 0.94 \geq 0.75$,

$DQ(x_2) = 0.80 \geq 0.75$, $DQ(x_3) = 0.42 < 0.75$, and $DQ(x_4) = 0.22 < 0.75$. Therefore, x_1 and x_2 are easy to classify, conversely, x_3 and x_4 are hard to classify. An arbitrary classifier from P can classify x_1 and x_2 . We select the classifier that classifies DSEL (D) with highest accuracy. In this situation, $\max(\delta) = 0.95$. Thus, C_4 is the most capable classifier for x_1 and x_2 . Consequently, DCS-DQ selects the classifier C_4 for x_1 and x_2 . Since $DQ(x_3)$ and $DQ(x_4)$ are less than 0.75; therefore, x_3 and x_4 are hard to classify. Suppose the accuracy value of the base classifiers on θ_{x_3} and θ_{x_4} are $\Psi_{x_3} = \{0.65, 0.85, 0.90, 0.75, 0.95\}$, $\Psi_{x_4} = \{0.70, 0.95, 0.75, 0.85, 0.90\}$.

$$\mathcal{LG}_{x_3} = \Psi_{x_3} + \delta = (0.65 + 0.75, 0.85 + 0.90, 0.90 + 0.85, 0.75 + 0.95, 0.95 + 0.85) = (1.40, 1.75, 1.75, 1.70, 1.80)$$

$$\mathcal{LG}_{x_4} = \Psi_{x_4} + \delta = (0.70 + 0.75, 0.95 + 0.90, 0.75 + 0.85, 0.85 + 0.95, 0.90 + 0.85) = (1.45, 1.85, 1.60, 1.80, 1.75)$$

The maximum values \mathcal{LG}_{x_3} and \mathcal{LG}_{x_4} are as follows: $\max(\mathcal{LG}_{x_3}) = 1.80$ and $\max(\mathcal{LG}_{x_4}) = 1.85$. These values belong to C_5 and C_2 , respectively. Thus, C_5 and C_2 are the most capable classifiers for x_3 and x_4 , respectively. Consequently, DCS-DQ selects the classifier C_5 from the classifier pool to classify x_3 , in the same way DCS-DQ selects the classifier C_2 from the classifier pool to classify x_4 .

6. Experimental study

In this section, the proposed DCS method, namely DCS-DQ, is compared with seven DCS techniques across four different text datasets. By using the same experimental protocol, we compared the DCS-DQ with seven state-of-the-art DCS techniques empirically. Details of the experimental part of this study are expressed in the following subsections. The experiments are carried out using seven DCS techniques given in Section 4. DCS methods are implemented using DESlib library in Python [44]. We used the default parameter values of all DCS methods in DESlib library. The pseudo-code for the DCS techniques is provided in the study of Britto et al. [45].

6.1. Classification scheme

The main research objective of this study is to propose a new DCS method, namely DCS-DQ. Experiments were carried out using four different benchmark text data sets and five state-of-the-art individual diverse base classifiers. The flow of the classification scheme utilized in this study is as follows: document collection, preprocessing, feature extraction, feature weighting, feature selection, classification and analysis of the results. Classification part includes pool generation and classifier selection dynamically. The text classification scheme illustrating the flow of the experiments is pictured in Figure 7. The following section of the paper explains the flow in detail.

6.2. Generating pool of classifier

To generate the pool of classifiers, we employed various classifier models. Five state-of-the-art classifiers have been selected. The selected classifiers are commonly employed in the text classification domain due to their high classification accuracy. Furthermore, these classifiers are heterogeneous, meaning they have diverse structures, thereby establishing diversity. Selected classifiers are kNN, DT, SVM, LR, and MNB.

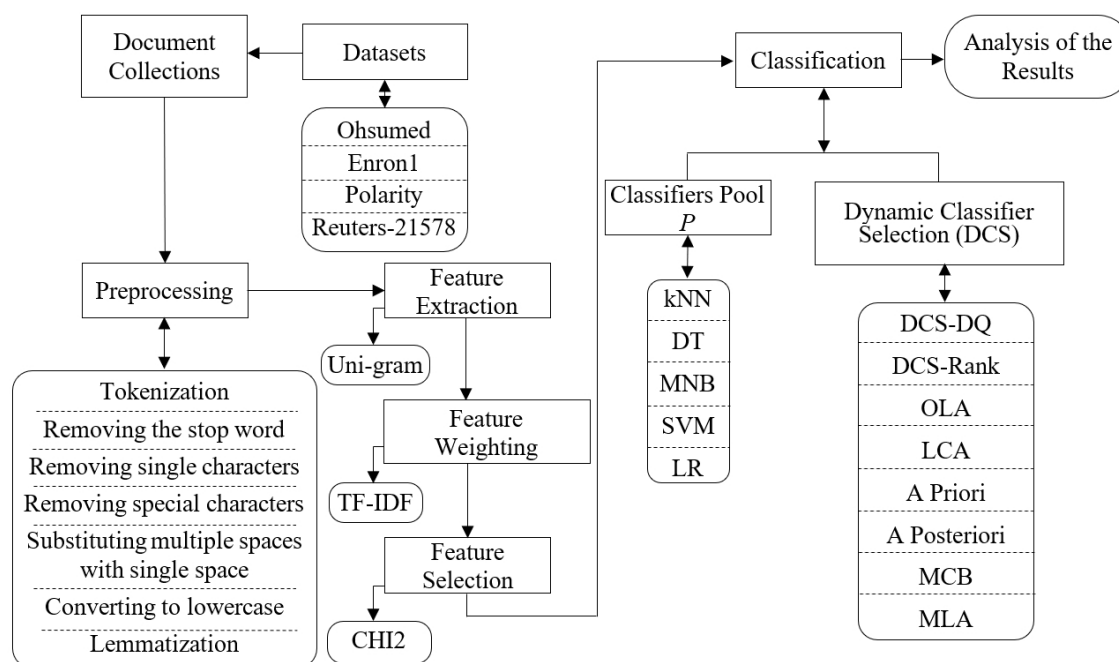


Figure 7. Classification scheme utilized in this study.

6.3. Datasets

Four distinct benchmark datasets have been utilized in this study. Reuters-21578 [48] dataset is widely used in text classification studies in the literature. Reuters-21578, known as ModApte split, contains the top 10 classes with the highest number of documents. Ohsumed [46] is a multiclass-unbalanced dataset. Ohsumed dataset is generated from a subset of the Medline database. The top 10 classes are used in the experiments. Enron1 [47] is an e-mail dataset consisting of two unbalanced classes, namely “Legitimate” and “Spam”. Polarity [50] is a two-class balanced dataset including 1000 positive and 1000 negative processed reviews about movies. The total number of documents and class labels related to the datasets used in the experimental studies are presented in Table 1.

6.4. Text preprocessing steps

Text preprocessing is essential in text classification due to the presence of noise and redundant words in documents within a Corpus. Therefore, preprocessing steps such as tokenization, stop word removal, removal of single characters, removal of special characters, substitution of multiple spaces with a single space, lowercase conversion, and lemmatization have been applied. Classifiers cannot directly operate on raw text data. Initially, text documents must be converted into numerical values. Various methods are employed for this conversion, among which the most popular model is the Bag of Words. In this approach, the text dataset is initially transformed into a matrix, where the rows represent documents from the Corpus and the columns represent features or words. The numerical values within the cells of the matrix are weighted utilizing term frequency-inverse document frequency (TF-IDF) [49] values.

Table 1. Properties of datasets.

Ohsumed		Reuters-21578- ModApte	
Class label	Number of documents	Class label	Number of documents
Neoplasms	2513	earn	3964
Digestive system diseases	837	acq	2369
Respiratory tract diseases	634	money-fx	717
Urologic and male genital diseases	842	grain	582
Nervous system diseases	1328	crude	578
Cardiovascular diseases	2876	trade	486
Nutritional and metabolic diseases	815	interest	478
Immunologic diseases	1060	ship	286
Disorders of environmental origin	1283	wheat	283
Pathological conditions, signs and symptoms	1924	corn	237
Enron1		Polarity	
Legitimate	3672	Positive	1000
Spam	1500	Negative	1000

6.5. Feature extraction and selection

The total number of features obtained after applying the preprocessing steps described in the previous section is presented in Table 2. The total number of features is relatively high. However, utilizing a large number of features during the classification process can lead to reduced accuracy and classifier performance. Therefore, the chi-square (CHI2) [50] feature selection method, which is commonly used, is employed to select the most appropriate features. Feature sizes of 100, 300, 500, 1000, and 3000 are used for all datasets.

Table 2. The numbers of documents and features in datasets.

Dataset	Number of documents	Number of features
Ohsumed	14,112	13,169
Enron1	5172	9190
Polarity	2000	12638
Reuters-ModApte	9980	9942

6.6. Evaluation

Although three out of four datasets used are unbalanced, there is not a significant imbalance between the classes. Taking this into consideration, classification accuracy and macro-averaged F-measure [51] are used as success measures. Macro F1-measure is particularly suitable for imbalanced data. Macro F1 is calculated for each class and then averaged across all classes. This ensures that equal weight is given to each class, irrespective of the number of documents in the classes.

Although three out of four datasets used are unbalanced, there is not a great imbalance between the classes. Considering this imbalance we used classification accuracy and macro-averaged F-measure [51] as success measures. Macro F1-measure is a suitable success measure for imbalanced data. Macro F1 is calculated for each class and averaged across all classes. In this manner, each class is assigned equal importance, irrespective of the number of documents within each class. The calculation of macro F1 can be formulated as follows (Equation 12):

$$MacroF1 = \frac{\sum_{k=1}^n F_k}{n}, F_k = \frac{2 \cdot p_k \cdot r_k}{p_k + r_k} \quad (12)$$

In Equation (12); p_k is precision value for class k , r_k is the recall value for class k , and n is the number of classes in datasets.

Classification accuracy is the second success measure for this study. It is defined as the proportion of accurately classified samples to the total number of samples to be classified. The equation for classification accuracy is given by Equation (13).

$$\text{Classification accuracy} = \frac{\text{correctly classified samples}}{\text{total number of samples}} * 100\% \quad (13)$$

7. Experimental results

After the implementation of the feature selection method, the resulting matrix was randomly split into 50% for training, 25% for testing, and 25% for dynamic selection (DSEL). The base classifiers were trained employing identical training data and tested using identical testing data. Due to the random partitioning of the test and training data, experiments were conducted 10 times, and the resulting classification accuracies were averaged. In each of the 10 scenarios, training data, testing data, and dynamic selection (DSEL) data were generated randomly. In order to evaluate the performance of DCS-DQ method, experiments were conducted using feature sizes of 100, 300, 500, 1000, and 3000. Experimental results are presented using line graph. In the graphs, the results of the proposed method DCS-DQ are presented in red color. The Oracle scores are shown in green color. In the graphs representing experimental results, DCS methods, namely Rank, OLA, LCA, A Priori, A Posteriori, MCB, and MLA are shown in different colors. Oracle performance is presented in Figures 8 and 9. As explained in the introduction, there is a big performance gap between oracle performance and the DCS methods. The proposed DCS-DQ method helps to fill the performance gap between the oracle performance and existing DCS methods in all datasets.

7.1. The results of the experiments conducted on Polarity and Enron1 datasets

Polarity and Enron1 are two-class datasets. While the Polarity dataset is balanced, the Enron1 dataset is imbalanced. The classification performances of DCS-DQ and existing methods are presented in Figure 8.

When analyzing the experimental results in the Polarity dataset, DCS-DQ method outperforms other DCS methods for all feature sizes. Additionally, with an increasing number of features, the performance of the DCS-DQ method improves, whereas the performance of other DCS methods increases up to 500 features and then starts to decline. Upon reviewing the experimental outcomes of the Enron1 dataset, the proposed method demonstrates superior performance compared to the pool mean for all feature sizes. The DCS-DQ method also outperforms other DCS methods for all feature sizes. Moreover, as the feature size increases, the performance of the DCS-DQ method improves, while the performance of other DCS methods decreases. From Figure 8, it can be inferred that the proposed method performs well on both the Enron1 and Polarity datasets.

7.2. The results of the experiments conducted on Ohsumed and Reuters datasets

Classification accuracies for Ohsumed and Reuters datasets are shown in Figure 9.

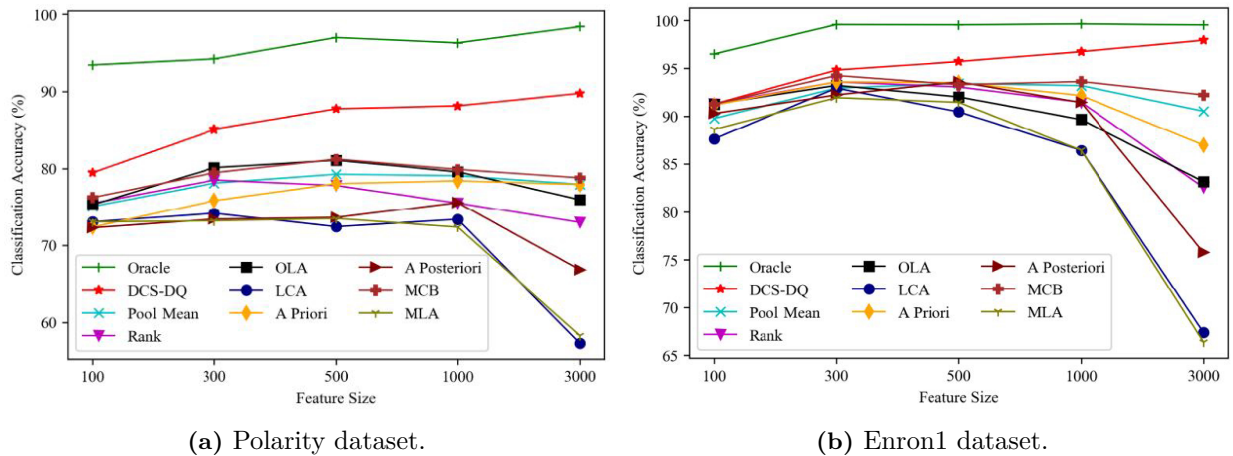


Figure 8. Classification accuracies on Polarity and Enron1 datasets.

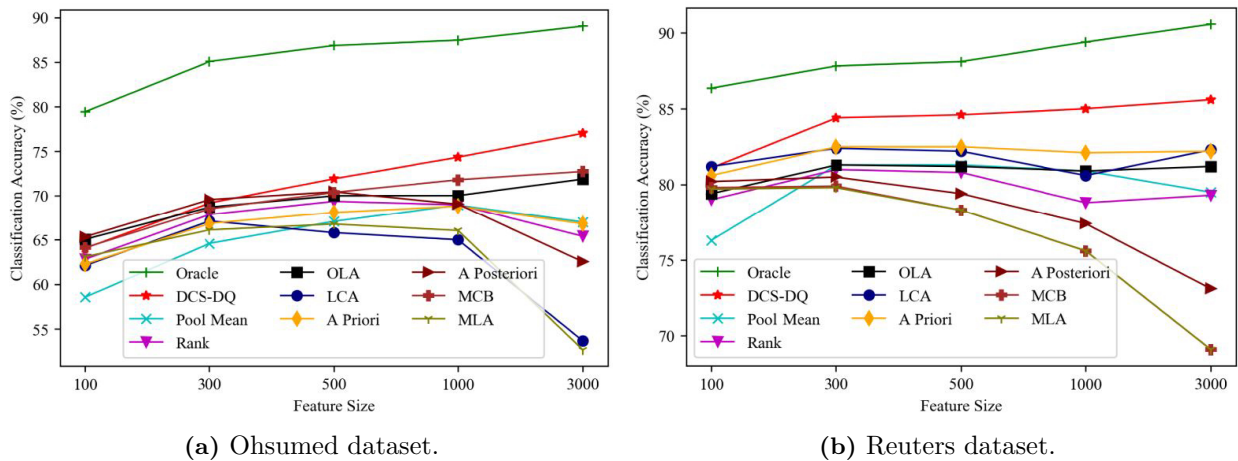


Figure 9. Classification accuracies on Ohsumed and Reuters datasets.

In accordance with the experimental results of the Ohsumed dataset, for feature size 100, OLA and A Posteriori methods perform well than DCS-DQ. Besides, only A Priori method has better performance than DCS-DQ when the number of feature is 300. For feature sizes which are 500, 1000, and 3000, DCS-DQ method’s classification accuracy is much better than others. Upon scrutinizing the experimental results of the Reuters dataset, it is observed that for all feature sizes except for 100, the classifiers selected by the DCS-DQ method yield higher classification accuracy than those selected by other methods. According to the results presented in Figure 9, it is evident that the DCS-DQ method is more efficient than other methods on both the Ohsumed and Reuters datasets. Ohsumed is a collection of medical texts. As depicted in Figure 9a, the classification accuracy of all DCS methods is low. Based on this outcome, it can be inferred that the Ohsumed dataset poses challenges in classification. Moreover, Ohsumed exhibits the lowest classification accuracy among all datasets for the proposed method.

Table 3. Macro F1-scores for Polarity and Enron 1 datasets.

Datasets	Polarity					Enron1				
	100	300	500	1000	3000	100	300	500	1000	3000
kNN	0.712	0.701	0.707	0.622	0.519	0.873	0.871	0.858	0.809	0.667
SVM	0.768	0.832	0.848	0.862	0.890	0.877	0.925	0.949	0.961	0.973
DT	0.670	0.669	0.662	0.651	0.650	0.888	0.923	0.924	0.926	0.921
MNB	0.807	0.858	0.873	0.876	0.889	0.809	0.899	0.926	0.947	0.959
LR	0.803	0.841	0.854	0.857	0.877	0.869	0.914	0.939	0.951	0.963
Rank	0.731	0.782	0.796	0.761	0.750	0.896	0.924	0.922	0.913	0.806
OLA	0.734	0.785	0.814	0.802	0.772	0.895	0.918	0.904	0.898	0.803
LCA	0.710	0.729	0.739	0.708	0.628	0.872	0.913	0.894	0.858	0.663
A Priori	0.725	0.754	0.768	0.779	0.766	0.899	0.925	0.920	0.923	0.829
A Posteriori	0.713	0.732	0.748	0.756	0.733	0.885	0.917	0.922	0.919	0.735
MCB	0.747	0.792	0.808	0.803	0.796	0.895	0.935	0.926	0.932	0.904
MLA	0.710	0.729	0.739	0.708	0.628	0.872	0.913	0.895	0.858	0.663
DCS-DQ	0.775	0.849	0.881	0.880	<u>0.894</u>	0.894	0.930	0.947	0.971	<u>0.977</u>

7.3. The analysis of the macro F1-scores of the DCS methods and base classifiers

The macro F1-scores of the DCS methods and base classifiers in the pool are presented in Tables 3 and 4. The highest score is highlighted in bold for each corresponding feature size. We have also denoted the highest score with both bold and underlined formatting for clarity. As shown in Tables 3 and 4, there is no single base classifier or DCS method that achieves the highest macro F1-score for all dimensions across different datasets. However, it is evident that the proposed method, DCS-DQ, consistently outperforms all other methods in terms of macro F1-score. The range of macro F1-scores for the Polarity dataset varies between 0.519 and 0.894. Notably, the highest macro F1-score achieved for feature sizes 500, 1000, and 3000 in the Polarity dataset is attributed to the DCS-DQ method. Consequently, it can be inferred that the DCS-DQ method demonstrates remarkable performance compared to all other methods in the Polarity dataset. Moreover, the highest macro F1-score for feature sizes 1000 and 3000 in the Enron 1 dataset, 500 and 1000 in the Ohsumed dataset, and 300, 500, 1000, and 3000 in the Reuters dataset is achieved by the DCS-DQ method. This further reinforces the effectiveness of the DCS-DQ method across multiple datasets and feature sizes.

7.4. The analysis of the overall performance of DCS-DQ method on all datasets

To elucidate the results depicted in Figures 8 and 9, the names of the methods corresponding to the peak points for each feature size are provided in Table 5. In the last column of the table, the performance ratio of the proposed DCS-DQ method is presented. This ratio indicates the number of peak points reached by the DCS-DQ method for each classifier on each dataset. For instance, the DCS-DQ method attains the peak point for all feature sizes on the Polarity dataset, resulting in a performance ratio of 1. Similarly, the DCS-DQ method outperforms other methods for 4 out of 5 feature sizes, leading to a performance ratio of 0.8 for DCS-DQ on the Enron1 dataset.

Table 5 demonstrates that the proposed method is more effective than existing state-of-the-art methods.

Table 4. Macro F1-scores for Ohsumed and Reuters datasets.

Datasets	Ohsumed					Reuters				
	100	300	500	1000	3000	100	300	500	1000	3000
kNN	0.613	0.664	0.672	0.639	0.482	0.637	0.600	0.593	0.410	0.610
SVM	0.578	0.653	0.689	0.720	0.761	0.599	0.635	0.605	0.557	0.652
DT	0.577	0.624	0.627	0.624	0.600	0.554	0.649	0.609	0.575	0.664
MNB	0.394	0.492	0.551	0.592	0.600	0.474	0.654	0.601	0.581	0.669
LR	0.596	0.667	0.696	0.718	0.744	0.384	0.660	0.590	0.553	0.664
Rank	0.615	0.681	0.694	0.679	0.632	0.722	0.722	0.714	0.691	0.680
OLA	0.642	0.686	0.698	0.695	0.690	0.730	0.732	0.725	0.714	0.708
LCA	0.627	0.654	0.667	0.648	0.524	0.712	0.688	0.671	0.620	0.548
A Priori	0.611	0.663	0.678	0.678	0.650	0.728	0.740	0.737	0.722	0.727
A Posteriori	0.647	0.693	0.706	0.690	0.622	0.728	0.709	0.699	0.651	0.592
MCB	0.637	0.686	0.703	0.708	0.701	0.729	0.737	0.737	0.724	0.728
MLA	0.627	0.654	0.667	0.648	0.525	0.711	0.688	0.671	0.620	0.548
DCS-DQ	0.610	0.685	0.709	0.731	0.756	0.725	0.746	0.752	0.753	0.755

Table 5. DCS methods producing the highest classification accuracies on all datasets for each feature size.

Datasets	Feature size					Ratio of DCS-DQ
	100	300	500	1000	3000	
Polarity	DCS-DQ	DCS-DQ	DCS-DQ	DCS-DQ	DCS-DQ	1.0
Enron1	OLA	DCS-DQ	DCS-DQ	DCS-DQ	DCS-DQ	0.8
Ohsumed	A Posteriori	A Posteriori	DCS-DQ	DCS-DQ	DCS-DQ	0.6
Reuters	LCA	DCS-DQ	DCS-DQ	DCS-DQ	DCS-DQ	0.8

7.5. The statistical analysis of the overall performance of DCS-DQ method

The paired samples t-test [53] was employed to analyze the experimental findings. This statistical test compares the means of a variable observed in two different situations. In our study, we compared the classification accuracy of the DCS-DQ method with that of existing methods. Our analysis revealed a significant difference between the DCS-DQ method and the other seven methods. The results of this comparison are presented in Table ??.

Hypotheses for paired samples t-test:

Null Hypothesis (H0): With 95% confidence, there is no statistically significant difference between the mean classification accuracy before and after the experiment. ($M1 = M2$)

Alternative Hypothesis (H1): With 95% confidence, there is a statistically significant difference between the mean classification accuracy before and after the experiment. ($M1 \neq M2$)

The decision about statistical significance is given using the p values on Table ?. Since the p values are less than 0.05 ($0.000 < 0.05$) for all seven pairs in Table ?, the null hypothesis (H0) is rejected. This indicates that there is a statistically significant difference between the mean classification accuracy before and after the experiment. Therefore, the alternative hypothesis (H1) is accepted. Consequently, we can infer for all pairs that the DCS-DQ method is statistically significant with 95% confidence. In other words, the DCS-DQ method is effective in improving classification accuracy.

Table 6. Time analysis of DCS methods.

Dataset	Polarity					Ohsumed				
	100	300	500	1000	3000	100	300	500	1000	3000
DCS-DQ	3.4	8.5	18.0	37.0	101.8	14.3	51.1	101.8	213.4	761.8
DCS-Rank	2.3	7.2	15.0	30.0	76.2	10.4	48.3	102.7	208.4	556.0
OLA	2.4	8.3	17.5	35.8	92.0	10.9	43.5	96.1	203.6	547.4
LCA	2.2	7.1	14.7	29.3	71.7	9.6	39.2	87.7	190.0	550.7
A Priori	2.8	8.5	17.7	36.3	97.9	11.3	43.5	98.0	213.0	583.4
A Posteriori	2.5	8.2	17.0	34.5	95.3	11.6	44.1	100.1	218.5	594.6
MCB	2.1	6.9	13.8	29.0	74.0	9.8	37.9	87.3	186.3	483.3
MLA	2.2	7.2	15.1	30.7	78.3	10.9	44.5	96.3	211.9	640.5

7.6. Time analysis of DCS methods

We compared the time analysis of the proposed method with other methods, and the results of the comparison are presented in Table 6. The table presents the time analysis for the Polarity and Ohsumed datasets. All values are expressed in milliseconds and denote the average classification time for only one test sample in the respective dataset. The best performing DCS method for any number of feature size is highlighted in bold. For the Polarity dataset, the MCB method is the fastest performing DCS method for 100 attributes. Interestingly, the MCB method also stands out as the fastest performing DCS method for both the Polarity and Ohsumed datasets overall. Although the running time of our proposed method is not better than that of other DCS methods, it is not significantly worse. Given its strong performance in terms of classification accuracy, the slight disadvantage in running time can be considered negligible.

8. Conclusion and future works

The purpose of the current study is to propose a new DCS method, namely DCS-DQ. DCS methods have been shown to improve classification accuracy in many classification problems. However, according to our research, this is the first time that DCS methods have been used in a text classification problem. Four different benchmark text datasets are used in the experimental study. In all datasets used in experimental studies, classification accuracies were improved when compared to existing methods. The highest improvement in classification accuracy was observed in the Polarity dataset. The proposed method outperforms other DCS methods in terms of macro F1-score. The DCS-DQ method demonstrates superior performance compared to other DCS methods in the Polarity dataset for feature sizes of 500, 1000, and 3000. Similarly, in the other three datasets, the proposed method outperforms other DCS methods based on the macro F1-score. We also demonstrated that the DCS-DQ method is statistically significant with a 95% confidence. Classification accuracy is one of the most important success measure used in classification problems and the proposed method has been shown to have high classification accuracy. In the light of all the findings, we can infer that the proposed DCS-DQ method can make a significant contribution to the text classification literature. The difference between the classification accuracy of DCS methods and Oracle classification accuracy is still quite large for all datasets. For example, when the number of attributes in Ohsumed dataset is 3000, the Oracle classification accuracy is 89.04%. For the same number of attributes, the classification accuracy of the DCS-DQ method is 77.02%. In future studies, developing different DCS methods to close this gap can contribute to the classification literature. In this study, we used different classifier model. In the experiments, we used the CHI2 as a feature selection method. In our future studies, we will investigate the effects of alternative feature selection methods on classification accuracy. However, the proposed DCS-DQ method can be adapted and applied to other pattern recognition problems, such as credit scoring, face recognition systems, and music genre classification.

References

- [1] Kowsari K, Meimandi KJ, Heidarysafa M, Mendu S, Barnes L et al. Text classification algorithms: a survey. *Information* 2019; 10 (4): 150. <https://doi.org/10.3390/info10040150>
- [2] Abooraig R, Al-Zu'bi S, Kanan T, Hawashin B, Al Ayoub M et al. Automatic categorization of Arabic articles based on their political orientation. *Digital Investigation* 2018; 25: 24-41. <https://doi.org/10.1016/j.diin.2018.04.003>
- [3] Iyer RD, Lewis DD, Schapire RE, Singer Y, Singhal A. Boosting for document routing. In: *Proceedings of the Ninth International Conference on Information and Knowledge Management*; McLean, VA, USA; 2000. pp. 70-77. <https://doi.org/10.1145/354756.354794>
- [4] Kale SD, Prasad RS. Influence of language-specific features for author identification on Indian literature in Marathi. In: Reddy V, Prasad V, Wang J, Reddy K (editors). *Soft Computing and Signal Processing. IC-SCSP 2019. Advances in Intelligent Systems and Computing*, vol 1118. Singapore: Springer, 2020, pp. 639-652. https://doi.org/10.1007/978-981-15-2475-2_59
- [5] Uysal D, Uysal AK. Automatic classification of EFL Learners' self-reported text documents along and effective continuum. *Advanced Education* 2022; 9 (20): 4-14. <https://doi.org/10.20535/2410-8286.248091>
- [6] Zhu L, Zhu Z, Zhang C, Xu Y, Kong X. Multimodal sentiment analysis based on fusion methods: a survey. *Information Fusion* 2023; 95: 306-325. <https://doi.org/10.1016/j.inffus.2023.02.028>
- [7] Ghosh S, Razniewski S, Weikum G. Answering count questions with structured answers from text. *Journal of Web Semantics* 2023; 76, 100769. <https://doi.org/10.1016/j.websem.2022.100769>
- [8] Rao S, Verma AK, Bhatia T. Hybrid ensemble framework with self-attention mechanism for social spam detection on imbalanced data. *Expert Systems with Applications* 2023; 217: 119594. <https://doi.org/10.1016/j.eswa.2023.119594>
- [9] Dietterich TG. Machine-learning research. *AI magazine* 1997; 18 (4): 97-136. <https://doi.org/10.1609/aimag.v18i4.1324>
- [10] Polikar R. Ensemble based systems in decision making. *IEEE Circuits and Systems Magazine* 2006; 6 (3): 21-45. <https://doi.org/10.1109/MCAS.2006.1688199>
- [11] Kuncheva LI. A theoretical study on six classifier fusion strategies. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 2002; 24 (2): 281-286. <https://doi.org/10.1109/34.982906>
- [12] Cruz RMO, Sabourin R, Cavalcanti GDC, Ren TI. META-DES: a dynamic ensemble selection framework using meta-learning. *Pattern Recognition* 2015; 48 (5): 1925-1935. <https://doi.org/10.1016/j.patcog.2014.12.003>
- [13] Cruz RMO, Sabourin R, Cavalcanti GDC. Dynamic classifier selection: recent advances and perspectives. *Information Fusion* 2018; 41: 195-216. <https://doi.org/10.1016/j.inffus.2017.09.010>
- [14] Xiao H, Xiao Z, Wang Y. Ensemble classification based on supervised clustering for credit scoring. *Applied Soft Computing* 2016; 43: 73-86. <https://doi.org/10.1016/j.asoc.2016.02.022>
- [15] Breiman L. Bagging predictors. *Machine Learning* 1996; 24 (2): 123-140. <https://doi.org/10.1007/BF00058655>
- [16] Woźniak M, Graña M, Corchado E. A survey of multiple classifier systems as hybrid systems. *Information Fusion* 2014; 16: 3-17. <https://doi.org/10.1016/j.inffus.2013.04.006>
- [17] Dietterich TG. Ensemble Methods in Machine Learning. In: *Multiple Classifier Systems. MCS 2000. Lecture Notes in Computer Science*, vol 1857. Berlin, Heidelberg: Springer Berlin Heidelberg, 2000. https://doi.org/10.1007/3-540-45014-9_1
- [18] Cruz RMO, Sabourin R, Cavalcanti GDC. A DEEP analysis of the META-DES framework for dynamic selection of ensemble of classifiers. *arXiv preprint* 2015; 1509.00825v2. <https://doi.org/10.48550/arXiv.1509.00825>
- [19] Oliveira DVR, Cavalcanti GDC, Sabourin R. Online pruning of base classifiers for dynamic ensemble selection. *Pattern Recognition* 2017; 72: 44-58. <https://doi.org/10.1016/j.patcog.2017.06.030>

- [20] Cavalin PR, Sabourin R, Suen CY. Dynamic selection approaches for multiple classifier systems. *Neural Computing and Applications* 2013; 22 (3): 673-688. <https://doi.org/10.1007/s00521-011-0737-9>
- [21] Melo L, Macêdo JF, Nardini FM, Renso C. RMkNN and KNORA-IU: combining imbalanced dynamic selection techniques for credit scoring. In: 2021 IEEE 33rd International Conference on Tools with Artificial Intelligence (ICTAI); Washington, DC, USA; 2021. pp. 823-830. <https://doi.org/10.1109/ICTAI52525.2021.00131>
- [22] Bashbaghi S, Granger E, Sabourin R, Bilodeau GA. Dynamic ensembles of exemplar-SVMs for still-to-video face recognition. *Pattern Recognition* 2017; 69: 61-81. <https://doi.org/10.1016/j.patcog.2017.04.014>
- [23] Porwik P, Doroz R, Wrobel K. An ensemble learning approach to lip-based biometric verification, with a dynamic selection of classifiers. *Expert Systems with Applications* 2019; 115: 673-683. <https://doi.org/10.1016/j.eswa.2018.08.037>
- [24] Batista L, Granger E, Sabourin R. Dynamic ensemble selection for off-line signature verification. In: Sansone C, Kittler J, Roli F (editors). *Multiple Classifier Systems. MCS 2011. Lecture Notes in Computer Science*, vol 6713. Berlin, Germany: Springer, 2011. https://doi.org/10.1007/978-3-642-21557-5_18
- [25] Xiao J, Xie L, He C, Jiang X. Dynamic classifier ensemble model for customer classification with imbalanced class distribution. *Expert Systems with Applications* 2012; 39 (3): 3668-3675. <https://doi.org/10.1016/j.eswa.2011.09.059>
- [26] Wen J, Gao H, Liu Q, Hong X, Sun Y. A new method for identifying the ball screw degradation level based on the multiple classifier system. *Measurement* 2018; 130, 118-127. <https://doi.org/10.1016/j.measurement.2018.08.005>
- [27] Groccia MC, Guido R, Conforti D. Multi-classifier approaches for supporting clinical decision making. *Symmetry* 2020; 12 (5): 699. <https://doi.org/10.3390/sym12050699>
- [28] Roy A, Cruz RMO, Sabourin R, Cavalcanti GDC. A study on combining dynamic selection and data preprocessing for imbalance learning. *Neurocomputing* 2018; 286: 179-192. <https://doi.org/10.1016/j.neucom.2018.01.060>
- [29] Feng P, Ma J, Sun C, Xu X, Ma Y. A novel dynamic android malware detection system with ensemble learning. *IEEE Access* 2018; 6: 30996-31011. doi: <https://doi.org/10.1109/ACCESS.2018.2844349>
- [30] Junior LM, Nardini FM, Renso C, Trani R, Macedo JA. A novel approach to define the local region of dynamic selection techniques in imbalanced credit scoring problems. *Expert Systems with Applications* 2020; 152: 113351. <https://doi.org/10.1016/j.eswa.2020.113351>
- [31] Feng X, Xiao Z, Zhong B, Dong Y, Qiu J. Dynamic weighted ensemble classification for credit scoring using Markov Chain. *Applied Intelligence* 2019; 49: 555-568. <https://doi.org/10.1007/s10489-018-1253-8>
- [32] Martins JG, Oliveira LS, Britto AS, Sabourin R. Forest species recognition based on dynamic classifier selection and dissimilarity feature vector representation. *Machine Vision and Applications* 2015; 26: 279-293. <https://doi.org/10.1007/s00138-015-0659-0>
- [33] Sousa OLV, da Silva DP, Campelo VES, e Silva RRV, Magalhães DMV. Ensemble of classifiers for multilabel clinical text categorization in Portuguese. In: Abraham A, Pillana S, Casalino G, Ma K, Bajaj A (editors). *Intelligent Systems Design and Applications. ISDA 2022. Lecture Notes in Networks and Systems*, vol 715. Cham, Switzerland: Springer, 2023, p. 42. https://doi.org/10.1007/978-3-031-35507-3_5
- [34] Li Y, Zhang S, Lai C. Agricultural text classification method based on dynamic fusion of multiple features. *IEEE Access* 2023; 11: 27034-27042. <https://doi.org/10.1109/ACCESS.2023.3253386>
- [35] Sergio AT, de Lima TPF, Ludermir TB. Dynamic selection of forecast combiners. *Neurocomputing* 2016; 218: 37-50. <https://doi.org/10.1016/j.neucom.2016.08.072>
- [36] Bhatnagar V, Bhardwaj M, Sharma S, Haroon S. Accuracy-diversity based pruning of classifier ensembles. *Progress in Artificial Intelligence* 2014; 2: 97-111. <https://doi.org/10.1007/s13748-014-0042-9>
- [37] Cruz RMO, Zakane HH, Sabourin R, Cavalcanti GDC. Dynamic ensemble selection vs k-nn: why and when dynamic selection obtains higher classification performance? In: 2017 Seventh International Confer-

- ence on Image Processing Theory, Tools and Applications (IPTA); Montreal, QC, Canada; 2017. pp. 1-6. <https://doi.org/10.1109/IPTA.2017.8310100>
- [38] Dogo EM, Nwulu NI, Twala B, Aigbavboa C. Accessing imbalance learning using dynamic selection approach in water quality anomaly detection. *Symmetry* 2021; 13 (5): 818. <https://doi.org/10.3390/sym13050818>
- [39] Groccia MC, Guido R, Conforti D. Multi-classifier approaches for supporting clinical diagnosis. In: Sforza A, Sterle C (editors). *Optimization and Decision Science: Methodologies and Applications*. ODS 2017. Springer Proceedings in Mathematics & Statistics, vol 217. Cham, Switzerland: Springer, 2017. https://doi.org/10.1007/978-3-319-67308-0_13
- [40] Woods K, Kegelmeyer WP, Bowyer K. Combination of multiple classifiers using local accuracy estimates. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 1997; 19 (4): 405-410. <https://doi.org/10.1109/34.588027>
- [41] Giacinto G, Roli F. Methods for dynamic classifier selection. In: *Proceedings 10th International Conference on Image Analysis and Processing*; Venice, Italy; 1999. <https://doi.org/10.1109/ICIAP.1999.797670>
- [42] Smits PC. Multiple classifier systems for supervised remote sensing image classification based on dynamic classifier selection. *IEEE Transactions on Geoscience and Remote Sensing* 2002; 40 (4): 801-813. <https://doi.org/10.1109/TGRS.2002.1006354>
- [43] Liu G, Guo J. Bidirectional LSTM with attention mechanism and convolutional layer for text classification. *Neurocomputing* 2019; 337: 325-338. <https://doi.org/10.1016/j.neucom.2019.01.078>
- [44] Cruz RMO, Hafemann LG, Sabourin R, Cavalcanti GDC. DESlib: a dynamic ensemble selection library in Python. *Journal of Machine Learning Research* 2020; 21: 1-5. <https://doi.org/10.48550/arXiv.1802.04967>
- [45] Britto Jr AS, Sabourin R, Oliveira LES. Dynamic selection of classifiers—a comprehensive review. *Pattern Recognition* 2014; 47 (11): 3665-3680. <https://doi.org/10.1016/j.patcog.2014.05.003>
- [46] Hersh W, Buckley C, Leone TJ, Hickam D. OHSUMED: an interactive retrieval evaluation and new large test collection for research. In: Croft BW, van Rijsbergen CJ (editors). *SIGIR'94*. London, UK: Springer, 1994, pp. 192-201. https://doi.org/10.1007/978-1-4471-2099-5_20
- [47] Klimt B, Yang Y. The Enron corpus: a new dataset for email classification research. In: Boulicaut JF, Esposito F, Giannotti F, Pedreschi D (editors). *Machine Learning: ECML 2004*. Lecture Notes in Computer Science, vol 3201. Berlin, Germany: Springer, 2004, pp. 217-226. https://doi.org/10.1007/978-3-540-30115-8_22
- [48] Asuncion A, Newman D. *UCI Machine Learning Repository* 2007; Irvine, CA, USA.
- [49] Uysal AK, Gunal S. The impact of preprocessing on text classification. *Information Processing & Management* 2014; 50 (1): 104-112. <https://doi.org/10.1016/j.ipm.2013.08.006>
- [50] Uysal AK, Gunal S. A novel probabilistic feature selection method for text classification. *Knowledge-Based Systems* 2012; 36: 226-235. <https://doi.org/10.1016/j.knosys.2012.06.005>
- [51] Lever J, Krzywinski M, Altman N. Classification evaluation. *Nature Methods* 2016; 13 (8): 603-604. <https://doi.org/10.1038/nmeth.3945>
- [52] Huang YS, Suen CY. The behavior-knowledge space method for combination of multiple classifiers. In: *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*; New York, NY, USA; 1993. <https://doi.org/10.1109/CVPR.1993.1626170>
- [53] Ross A, Willson VL. Paired samples T-test. In: *Basic and Advanced Statistical Tests*. Rotterdam, the Netherlands: SensePublishers, 2017, pp. 17-19. https://doi.org/10.1007/978-94-6351-086-8_4