

A NOVEL MODIFIED RELATIVE DISCRIMINATION CRITERION FOR FEATURE RANKING IN TEXT CLASSIFICATION

Shakir Ullah¹, Imad Ullah², Ibad Ullah³, Naseer Ullah⁴, Muhammad Taufiq⁵

¹ University of Malakand, Chakdara, Lower Dir, Khyber Pakhtunkhwa, Pakistan
Email: shakiruom4u@gmail.com

^{2,3,4,5} Faculty of Computing, Riphah International University, Islamabad, Pakistan
Emails: : shakiruom4u@gmail.com , imaduom@gmail.com, ibad.ullah@riphah.edu.pk,
naseer.ullah@riphah.edu.pk, muhhammad.taufiq@riphah.edu.pk

DOI: <https://doi.org/10.5281/zenodo.17318910>

Keywords

Text Classification, Feature Selection, Modified Relative Discrimination Criterion, Feature Ranking, Text Mining, Machine Learning, High-Dimensional Data, Discriminative Features

Article History

Received on 8 July 2025
Accepted on 23 Sep 2025
Published on 11 OCT 2025

Copyright @Author

Corresponding Author: *
Ibad Ullah

Abstract

The exponential growth of textual data across diverse domains poses significant challenges for extracting actionable insights, particularly in text classification tasks. Effective feature selection is paramount to mitigate issues such as high dimensionality, sparsity, and semantic complexity inherent in text corpora. Conventional feature selection methods, primarily designed for numerical or categorical data, often underperform when applied to text due to these unique characteristics. This study proposes the Modified Relative Discrimination Criterion (MRDC), a novel feature ranking approach specifically developed for text classification. The MRDC enhances feature selection by quantifying the discriminative capacity of features within a text corpus, thereby enabling more robust and efficient ranking. Evaluated on standard performance metrics including accuracy, precision, recall, and F1-score the MRDC achieved 82.12% accuracy, 82.42% precision, 82.12% recall, and 82.16% F1-score using only 1,500 features, compared to approximately 150,000 original features. In contrast, the baseline Relative Discrimination Criterion yielded a lower F1-score of 74.38% with the same feature subset. These results demonstrate the MRDC's superior performance in optimizing feature selection for text classification tasks, offering a significant advancement over existing methods.

INTRODUCTION

The exponential growth of written content is due to the exponential increase of the digital content that is brought about by using social media, news articles, scientific publications, and customer

reviews. Internet traffic across the globe has been growing at 26,600 GB/s in 2023, with the amount increasing to nearly three times in the next five years, and is currently 26,600 GB/s, and

it was just 100 GB/s in 2002 [2]. This rapid increase underscores the urgent need for advanced data processing techniques to extract actionable insights from large-scale text corpora. Feature selection, an essential preprocessing phase in machine learning, is crucial for pinpointing pertinent features to improve the efficacy of text analysis applications, including text classification, sentiment analysis, information retrieval, and topic modeling [1], [5], [6]. Nevertheless, the special properties of textual data such as high dimensionality [5], [6], sparsity and semantic complexity present substantial obstacles for conventional feature selection techniques, which were mainly developed for numerical or categorical datasets [7].

Feature selection is required in the improvement of model accuracy, reduction of the complexity of computations and facilitating the text analysis interpretation [5]. High-dimensional text data, i.e. thousands or millions of terms, increases the cost of computation and the risk of overfitting [3]. Another factor that makes the identification of meaningful features more difficult is the few terms in the individual documents, which is known as sparsity [7]. The semantic richness of language, where words are defined contextually, and other concepts can be expressed in different words, also limits the implementation of conventional methods [8]. The Conventional feature selection methods, including Term Frequency-Inverse Document Frequency (TF-IDF), Chi-square, and Mutual Information (MI) and Information Gain (IG) are popular but have significant drawbacks. For example, TF-IDF, though effective in ranking words according to their frequency and rarity, do not rank words according to their order or semantic connections, hence, losing the contextual nuance [12]. Chi-square based tests have the tendency of giving preference to infrequent words and this renders them biased in the selection of features [13]. MI is ineffective at capturing higher-order interactions of features that are significant in the analysis of

complex texts [4],[9], and IG is sensitive to changes in the length of the document, which can bias features of interest [11]. Document frequency methods do not consider the significance of the terms in a particular class and thus may overlook significant aspects [10].

The recent developments in the field of feature selection have attempted to address these issues by coming up with new methods. Word embeddings Word2vec, GloVe and FastText are word embeddings, i.e., they are representations of words represented as vectors, i.e., continuous vectors, which capture semantic relations and context sense [8]. Such approaches provide the opportunity to choose the features more finely with the assistance of co-occurrence of words and linguistic patterns. The deep learning models, especially attention models such as Bidirectional Encoder Representations of Transformers (BERT) is an enhancement in the feature selection aspect of task-specific text analysis [8]. The ensemble methods, including bagging, boosting, and XGBoost, are used to select a range of models and improve stability and reduce biases that single methods have [10]. Hybrid approaches that combine more traditional statistical techniques (e.g., TF-IDF, Chi-squared) with new ones, e.g. word embeddings or deep learning, have also proved useful in enhancing the accuracy and coverage of feature selection [11]. Though such progress has been made, the existing methods do not always fully consider the relationship between high dimensionality, sparsity, and textual semantic richness and require more complex methods.

To eliminate these weaknesses, the proposed method will propose a new feature ranking method, the Modified Relative Discrimination Criterion (MRDC), which will be used in text classification. Unlike the conventional methods, MRDC estimates characteristics based on their discriminative power on a text corpus with consideration of the statistical significance and semantic context. The high-dimensional sparse

data can be readily addressed by MRDC, and the complicated linguistic Associations can be addressed by studying feature distributions concerning the distributions of classes. The method integrates the experience of word embeddings with the analysis of context to ensure the robust feature ranking hence enhancing the accuracy, efficacy and clarity of text classification models. This is the main innovation of MRDC because the features are given a discriminative score in the relative importance to the corpus, and its distinguishing ability between classes. The technique can address the weakness of the conventional techniques in that it puts emphasis on features that are both statistically significant and semantically significant. Preliminary experiments have shown that MRDC is better than baseline methods and a smaller set of features results in high accuracy, precision, recall and F1-score, so it is a promising solution to the text analysis problem.

The proposed method described as a considerable knowledge gap in the literature because of the introduction of a feature selection methodology which is specially oriented on the specific issues of the textual information. The Proposed method is directed towards the following objective:

- Develop the Modified Relative Discrimination Criterion (MRDC) method to efficiently prioritize text features for classification tasks.
- Perform thorough assessments using varied text datasets to evaluate the MRDC method's performance and dependability.
- Demonstrate the superiority of MRDC over conventional feature selection techniques by verifying enhanced classification accuracy and computational efficiency.

This paper is structured as follows: Section II examines prior research on feature selection for text analysis. Section III details the methodology of the MRDC approach. Section IV outlines the results, and subsequent discussions. Section V

summarizes the findings and suggests avenues for future investigation.

2: Literature Review

The rapid proliferation of textual information in areas like social media, news, scientific journals, and user generated information has led to a pressing requirement of advanced feature selection methods able to operate on high-dimensional, sparse and semantically multifaceted information [2], [14]. Although the traditional feature selection techniques are effective in the context of structured data, they are not always effective in text mining because of the distinct nature of the text corpora [3]. This section of the proposed method summarizes current feature selection methods in the field of text classification, their restrictions, and recent developments, and the rationale behind the development of new methods such as the Modified Relative Discrimination Criterion (MRDC).

A. Traditional Feature Selection Methods

Features selection is one of the most significant preprocess step in machine learning, particularly text classification, that is applied to select the most valuable features in order to enhance the accuracy and effectiveness of the model [1], [5]. The most popular are the Term Frequency-Inverse Document Frequency (TF-IDF), Chi-squared statistics, Mutual Information (MI), Document Frequency (DF) and Information Gain (IG) which possess their pros and cons.

1) Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF is a statistical metric used to assess the significance of a term within a document relative to a larger collection of documents. It integrates two components: Term Frequency (TF), which captures how often a term appears in a specific document, and Inverse Document Frequency (IDF), which measures how uncommon the term is across the entire collection:

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \cdot \text{IDF}(t) \quad (2.1)$$

The Term Frequency is calculated as:

$$\text{TF}(t, d) = \frac{\text{Count of term } t \text{ in document } d}{\text{Total terms in document } d} \quad (2.2)$$

The Inverse Document Frequency is computed as:

$$\text{IDF}(t) = \log\left(\frac{N}{n_t}\right) \quad (2.3)$$

Here, N represents the total number of documents in the collection, and n_t is the number of documents containing term t . The TF-IDF score emphasizes terms that are frequent within a specific document but rare across the collection, making it valuable for applications such as information retrieval, document clustering, and text classification. However, a limitation of TF-IDF is its assumption of term independence, as it does not account for word order, synonyms, or polysemy, which can hinder its ability to capture deeper semantic relationships and may lead to less effective feature representation in tasks requiring contextual understanding.

2) Chi-Squared Statistics

Chi-squared statistics evaluate the independence between terms and class labels, identifying terms strongly associated with specific classes [16]:

$$f(X; k) = 1 / \left(2^{k/2} \Gamma(k/2)\right) x^{(k/2)-1} e^{-x/2}, x \geq 0 \quad (2.4)$$

where $f(x)$ is the probability density function, k is the degrees of freedom, and $\Gamma(k/2)$ is the gamma function. In text classification, Chi-squared compares observed term frequencies with expected frequencies under the null hypothesis of independence, with high values indicating strong term-class associations [12]. Its simplicity and interpretability make it effective for reducing feature space dimensionality [17]. However, Chi-squared is sensitive to rare terms, which can inflate scores due to low expected frequencies, leading to biased feature selection [16]. It also assumes term independence, an unrealistic

assumption in text data where words are often contextually related [18]. Performance degrades in small datasets with unstable frequency distributions, and non-zero expected frequencies are required, posing challenges when terms are absent in some classes [19].

3) Mutual Information (MI)

Mutual Information measures the dependency between a term and a class by quantifying the reduction in uncertainty about the class given the term's presence [20]:

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log(p(x, y) / p(x)p(y)) \quad (2.5)$$

$$I(X; Y) = \int_{x \in X} \int_{y \in Y} p(x, y) \log(p(x, y) / p(x)p(y)) dx dy \quad (2.6)$$

where $p(x, y)$ is the joint probability distribution of term X and class Y , and $p(x)$ and $p(y)$ are marginal probabilities. Equation (2.5) applies to discrete variables, while (2.6) extends to continuous cases. MI is advantageous for capturing non-linear relationships, making it suitable for complex text classification tasks [22]. However, it struggles with higher-order dependencies, such as interactions between multiple terms or phrases, and is sensitive to rare terms, which can inflate scores [21]. Estimating probability distributions for MI is computationally expensive in high-dimensional datasets, and continuous data require approximations that may introduce errors [9].

4) Document Frequency (DF)

Document Frequency measures a term's importance by counting the number of documents in which it appears [24]:

$$\text{df} = \text{TF}(t, d) \cdot \text{IDF}(t) \quad (2.7)$$

$$\text{TF}(t, d) = \frac{\text{count}(t \in d)}{|d|} \quad (2.8)$$

$$IDF(t) = \log\left(\frac{N}{n_t}\right) \quad (2.9)$$

where n_t is the number of documents containing term t . DF filters out overly common or rare terms, reducing feature space dimensionality and computational complexity [25]. However, it does not account for class-specific term relevance, potentially overlooking terms critical for distinguishing specific categories [10]. DF also ignores contextual information, limiting its ability to handle polysemy or context-dependent meanings.

5) Information Gain (IG)

Information Gain quantifies the reduction in entropy achieved by partitioning a dataset based on a feature [26]:

$$IG = Entropy(S) - \sum_{i=1}^k \frac{|C_i|}{|P|} \cdot Entropy(C_i) \quad (2.10)$$

$$Entropy(S) = - \sum_{i=1}^n p_i \log_2(p_i) \quad (2.11)$$

$$Gini(S) = 1 - \sum_{i=1}^n p_i^2 \quad (2.12)$$

$$\sum_{i=1}^k \frac{|C_i|}{|P|} \cdot Entropy(C_i) \quad (2.13)$$

where S is the parent set, C_i is the i -th child node, $|C_i|$ and $|P|$ are the sizes of the child and parent nodes, and p_i is the probability of class i . IG is effective for selecting features that reduce class uncertainty, improving classification accuracy [27]. However, it is sensitive to document length variations, which can bias feature selection toward shorter documents [28]. IG also overlooks feature interactions, limiting its ability to capture complex linguistic patterns [29].

B. Limitations of Traditional Methods

Traditional feature selection methods face significant challenges in text data analysis. High dimensionality, often involving millions of unique terms, increases computational costs and the risk

of overfitting [3]. The lack of informative features is caused by sparsity, that is, most terms are rarely observed [7]. In addition, such algorithms as TF-IDF and DF are based on a frequency-based measure, which cannot reflect semantic links or contextual dependencies [8]. Chi-squared and MI are term-independent, which is not true in text data because of words co-occurrences and syntactic patterns [18]. Although IG is efficient in reducing the entropy, it does not consider the interactions of features and context-specific meanings, which might miss subtle but important features [29]. Such restrictions require sophisticated tools capable of taking care of text data complexities.

C. Advances in Feature Selection

Recent research has introduced advanced feature selection techniques to address these challenges. Word embeddings, such as Word2Vec and GloVe, represent terms as continuous vectors, capturing semantic relationships by placing similar words closer in vector space [8]. Contextual embeddings from models like BERT leverage attention mechanisms to model long-range dependencies and context-specific meanings [23]. Ensemble methods, such as XGBoost, combine multiple models to enhance robustness and reduce bias [10]. Hybrid approaches integrating statistical methods with embeddings have also improved feature selection accuracy [11]. Despite these advancements, challenges remain in fully capturing feature interactions and efficiently handling high-dimensional, sparse text data.

D. Relative Discrimination Criterion (RDC) and Alternative RDC (ARDC)

The Relative Discrimination Criterion (RDC) is a significant advancement in feature ranking for text classification [13]:

$$RDC_{tc} = \frac{|tpr_{tc} - fpr_{tc}|}{\min(tpr_{tc}, fpr_{tc})} \cdot t_c \quad (2.14)$$

where tpr_{tc} is the true positive rate, fpr_{tc} is the false positive rate, and t_c is the term count for term t in class c . RDC prioritizes features with high class-separating ability, reducing overfitting and improving classification accuracy [13]. However, it may struggle to capture complex feature interactions or context-specific nuances. The Alternative Relative Discrimination Criterion (ARDC) extends RDC by incorporating finer-grained discriminative cues [30],[31]:

$$\text{ARDC}_{tc} = \frac{D_{tc}}{\min(\text{TPR}_{tc}, \text{FPR}_{tc})} \cdot T_c \quad (2.15)$$

$$\text{TPR}_{tc} = \frac{\text{TP}_{tc}}{\text{POS}} \quad (2.16)$$

$$\text{FPR}_{tc} = \frac{\text{FP}_{tc}}{\text{NEG}} \quad (2.17)$$

$$D_{tc} = \text{DF}_{tc} \cdot |\text{TPR}_{tc} - \text{FPR}_{tc}|^{K-1} \quad (2.18)$$

where D_{tc} is the discrimination value, DF_{tc} is the discrimination factor, and K is the number of classes. ARDC enhances RDC by considering feature interactions and context, improving classification accuracy and efficiency [30]. However, its increased complexity requires careful implementation and tuning.

E. Need for Modified Feature Selection Techniques

Despite the development of RDC and ARDC, there are limitations in the situation when text data should be processed because of dynamism and complexity. The large dimensionality, sparsity and complicated interaction of features require the existence of methods that can be semantically aware as well as statistically rigorous. The proposed Modified Relative Discrimination Criterion (MRDC) is founded on the RDC and the ARDC and involves more advanced context and feature interaction modeling analysis to provide a more effective and solid solution when ranking features to classify text. MRDC will remove the imperfections of existing methods, potentially turning text analysis tasks into more

precise and scalable, thus assisting the feature selection study under consideration to advance.

3: Methodology

This section outlines the proposed methodology for developing a text classification system using the Modified Relative Discrimination Criterion (MRDC) for feature selection and a Support Vector Machine (SVM) for classification. The 20 Newsgroups data set that is a widely used standard in text classification is used as a measure of the performance of MRDC in determining the discriminative features and improving classification accuracy. The process involves dataset preparation, feature extraction and selection, model training, and performance evaluation, with a focus on computational efficiency and robust evaluation metrics.

A. Dataset

The 20 Newsgroups dataset is a standard sample in natural language processing, and it consists of about 20,000 newsgroup documents in 20 different categories, which are politics, religion, sports, and technology [14]. The data is balanced and each category has approximately equal number of documents that are suitable in assessing the text classification algorithms. Preprocessing of documents is done to contain just the main text, title, and summary and remove duplicates. The data is divided into eight major classes: REGIONS, ADVISORIES, SPORTS, NEWS, SCIENCE, INTERNATIONAL, ENTERTAINMENT and BUSINESS, and allows to test classification models on a variety of topics. The data is loaded through the `fetch_20newsgroups` method in Scikit-learn that gives the raw text and the label of the category [15]. Table I presents a sample of the dataset, showing document excerpts and their labels, while Table II summarizes the statistical distribution of samples across categories, confirming the dataset's balanced nature.

TABLE I: Representative samples from the 20 Newsgroups dataset

News Excerpt	Label
Author: J. Ratnam Discussion on emerging AI technologies	10
User: M. Lawson Post concerning recent political debates	3
Contributor: H. Eren Commentary on sports management topics	17
Writer: G. Dawson Opinion piece on modern software tools	3
Post: A. McDiarmid Article exploring educational reforms	4

TABLE II: Descriptive Statistical Summary of the 20 Newsgroups Dataset

Statistic	Value
Total Documents	18,846
Mean	9.293
Standard Deviation	5.563
Minimum	0.000
25th Percentile (Q1)	5.000
Median (50th Percentile)	9.000
75th Percentile (Q3)	14.000
Maximum	19.000

B. Train-Test Split

To ensure unbiased evaluation, the dataset is split into training (80%) and testing (20%) sets using Scikit-learn's `train_test_split` function. This split allows the model to be trained on a subset of the data and evaluated on unseen documents, providing a reliable estimate of generalization performance.

C. Modified Relative Discrimination Criterion (MRDC)

The MRDC method enhances feature selection by identifying discriminative features from high-

dimensional text data. The process begins with text preprocessing using count vectorization, implemented via Scikit-learn's `CountVectorizer`. This transforms the text corpus into a sparse matrix $M \in \mathbb{R}^{n \times m}$, where n is the number of documents, m is the number of unique words, and M_{ij} represents the frequency of word j in document i .

1) Calculation of Class and Feature Probabilities

To understand how documents are distributed across categories, class probabilities are computed as follows:

$$P(C) = \frac{\text{Number of documents in class } C}{\text{Total number of documents}} \quad (3.1)$$

Here, the numerator represents the count of documents belonging to class C , and the denominator is the total document count. Next, feature probabilities are determined to assess the likelihood of a word occurring within a specific class:

$$P(F_j | C) = \frac{1}{N_C} \sum_{i=1}^n M_{ij} \cdot \mathbb{1}(y_i = C) \quad (3.2)$$

In this equation, N_C is the number of documents in class C , M_{ij} denotes the presence of feature F_j in document i , and $\mathbb{1}(y_i = C)$ is an indicator function that equals 1 if document i belongs to class C , and 0 otherwise.

2) Computation of RDC Score

The RDC score measures a feature's ability to distinguish between classes by evaluating the ratio of its class-specific probability to the overall class probability:

$$\text{RDC}(F_j) = \sum_{c \in \text{classes}} P(F_j | C) \cdot \log \left(\frac{P(F_j | C)}{P(C)} \right) \quad (3.3)$$

Here, $P(F_j | C)$ represents the probability of feature F_j given class C , and $P(C)$ is the probability of class C . Features are then sorted in descending order based on their RDC scores, and the top 1500 are selected to reduce the dataset's dimensionality while preserving key discriminative information.

3) TF-IDF Transformation

To further enhance feature representation, the selected features undergo TF-IDF transformation to scale their values based on their importance across the corpus:

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \cdot \log \left(\frac{N}{\text{DF}(t)} \right) \quad (3.4)$$

where $\text{TF}(t, d)$ is the term frequency of term t in document d , N is the total number of documents,

and $\text{DF}(t)$ is the number of documents containing term t . This transformation, implemented via Scikit-learn's Tfidf Transformer, de-emphasizes common terms and highlights informative ones, improving the quality of the feature set for classification.

D. Support Vector Machine (SVM) Model

A linear Support Vector Machine (SVM) is a highly effective classifier that is trained on the features of interest to determine the optimum hyperplane to divide the classes in the feature space [10]. The reason why the linear kernel is used is that it is computationally efficient and interpretable in the transformed feature space when linear separability is assumed. The SVM maximizes the margin of the classes and uses the discriminative features chosen by MRDC to maximize the classification performance.

The training process is done by mapping the features that are transformed to the TF-IDF space and then the SVM used to identify support vectors to determine the decision boundary. MRDC allows the SVM to attain high accuracy using less computational cost by specializing in a smaller number of highly discriminative features and reducing the classification problem.

E. Evaluation Metrics

The performance of the SVM model is assessed using standard metrics to provide a thorough evaluation of its effectiveness:

1) Precision

Precision quantifies the fraction of positive predictions that are correct:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (3.5)$$

Here, TP denotes the number of true positives, and FP represents the number of false positives. A high precision value indicates a low rate of false positives, which is essential in scenarios where

erroneous positive predictions carry significant consequences.

2) Recall

Recall, also known as sensitivity, measures the proportion of actual positive instances correctly identified by the model:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3.6)$$

In this equation, FN is the number of false negatives. A high recall value signifies the model's ability to effectively detect positive instances.

3) F1-Score

The F1-score provides a balanced measure of precision and recall, offering a single metric to evaluate model performance:

$$\text{F1-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3.7)$$

The F1 score is particularly valuable for imbalanced datasets, as it considers both false positives and false negatives, providing a comprehensive assessment of the model's performance.

4) Classification Report

A classification report provides class-wise

precision, recall, and F1-score, along with overall accuracy, offering a thorough understanding of the model's performance across all categories. This is especially valuable for multi-class problems, where global accuracy may mask class-specific weaknesses.

F. Workflow`

The workflow is illustrated in Fig. 1 and summarized as follows:

1. **Data Collection:** Retrieve the 20 Newsgroups dataset using Scikit-learn's fetch_20newsgroups function.
2. **Train-Test Split:** Divide the dataset into 80% training and 20% testing sets.
3. **Feature Selection:** Apply MRDC to select the top 1500 features using count vectorization and TF-IDF transformation.
4. **Model Training:** Train a linear SVM on the selected features.
5. **Evaluation:** Compute accuracy, precision, recall, F1-score, and a classification report to assess model performance.

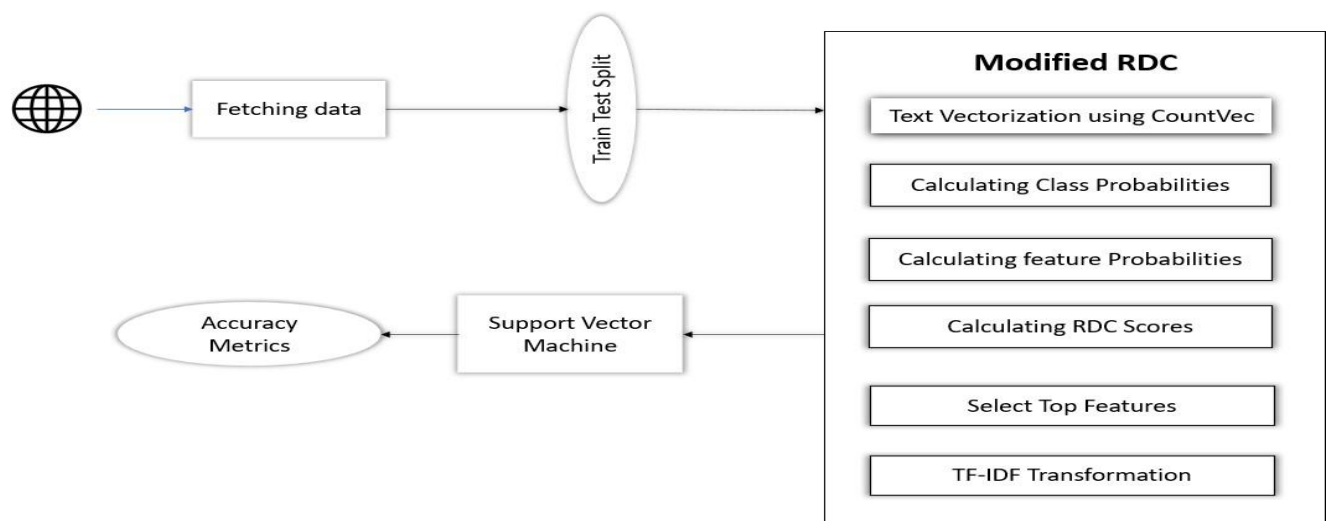


Fig. 1. Workflow of the Proposed Methodology

G. Hardware Requirements

The Proposed methodology is implemented on a system with the following specification as shown

in Table III ensuring sufficient computational resources for processing large-scale text data.

TABLE III: Experimental Environment and System Configuration

Component	Specification
Processing Unit	Dual-core Intel Xeon (2 virtual CPUs)
Memory (RAM)	13 GB DDR4
Graphics Accelerator	NVIDIA Tesla T4 with 15 GB VRAM

These specifications support efficient execution of the MRDC algorithm and SVM training, accommodating the computational demands of high-dimensional text processing.

4: Results and Discussion

This section presents a comprehensive evaluation of the proposed text classification system, which employs the Modified Relative Discrimination Criterion (MRDC) for feature selection and a Support Vector Machine (SVM) classifier, applied to the 20 Newsgroups dataset. The main aim is to evaluate how effective MRDC is in the improvement of classification performance relative to a baseline methodology based on the Term Frequency-Inverse Document Frequency (TF-IDF) without feature selection. The evaluation of performance is done in terms of precision, recall, F1-score, overall accuracy, and micro and macro averages at different levels of features. Python 3.8 was used to run the experiments in Google Colab and the optimized libraries were used to process the data, train the models and visualize the data.

A. Implementation Environment

The experiments were executed on Google Colab with Python 3.8, utilizing the following libraries to ensure robust data handling, model training, and visualization:

Pandas: Facilitates efficient data manipulation through DataFrame objects, used for managing the dataset and feature matrices [15].

NumPy: Supports numerical operations on multi-dimensional arrays, enabling efficient computation for feature extraction and model training [15].

2. **Seaborn:** Enhances data visualization with aesthetically pleasing plots, built on Matplotlib, used for generating performance graphs [15].

3. **Matplotlib.pyplot:** Provides a flexible interface for creating publication-quality plots, including bar and line graphs for performance analysis [15].

4. **Scikit-learn:** Implements the SVM classifier, count vectorization, TF-IDF transformation, and evaluation metrics, ensuring robust machine learning workflows [15].

These libraries were selected for their reliability and efficiency, supporting the computational demands of processing high-dimensional text data and generating professional visualizations.

B. Performance Metrics

The SVM model, which was trained on the features selected by MRDC, obtained a total accuracy of 82.12, which is higher than the accuracy of the baseline TF-IDF model (74). This enhancement shows that MRDC can discover discriminative characteristics that can improve the accuracy of the classification of the SVM. The main performance indicators are described in the subsequent subsections with the help of visualizations that give a more subtle insight into the behavior of the model.

1) Precision

Precision, defined as the proportion of correct positive predictions, is calculated as:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3.5)$$

The suggested methodology was precise by 82.42% meaning that 82.42% of positive predictions were accurate. This high accuracy highlights the ability of MRDC to choose features to reduce false positives, which is important in applications where a false positive prediction can be very expensive, e.g. spam detection or medical diagnosis [10]. The strength of the precision indicates the conscientiousness that MRDC applies to the set of features that are noisy or do not contribute to the discrimination, and only the very discriminative terms are preserved.

2) Recall

Recall, or sensitivity, measures the proportion of positive instances correctly identified:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3.6)$$

The recall was 82.12, which shows that the model was able to detect 82.12% of the case instances in the dataset. This recollection is strong, and it shows that MRDC can be able to retrieve a large number of discriminative features and thus cover positive instances completely, which is critical when dealing with tasks that demand high sensitivity, like topic identification in varying corpora [10].

3) F1-Score

The F1-score balances precision and recall, providing a single metric for model performance:

$$\text{F1-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3.7)$$

The F1-score was 82.16 percent that is a good trade-off in terms of precision and recall. This is particularly helpful when the dataset is uneven as it considers false positives and false negatives, and in this regard, this ensures that it functions well in varying class distribution [10]. The high F1-score confirms the feasibility of the features chosen by MRDC and its ability to maximize the accuracy and coverage.

4) Overall Metrics Visualization

Figure 2 shows the cumulative performance measures of the Support Vector Machine (SVM) model of Accuracy, Precision, Recall and F1-Score, of the Support Vector Model on features that were picked through the Modified Relative Discrimination Criterion (MRDC). All metrics represent how the model classifies in all the 20 categories in the dataset.

The four metrics reveal that the values of all metrics are very close, with all of them over 82 percent, which proves that MRDC-SVM framework provides balanced and trustworthy performance in various evaluation dimensions. The figure contains the vertical bars that show the mean scores and the small ranges of standard errors to demonstrate consistency and model stability. The values of the exact metrics are annotated above every bar (Accuracy: 0.8212, Precision: 0.8242, Recall: 0.8212, and F1-Score: 0.8216), which shows the consistency of the findings. This uniformity in all performance measures indicates that the feature selection based on the MRDC is a good way of capturing discriminative information without overfitting and biasing to particular classes. As a result, the model is highly generalized with regard to various topics.

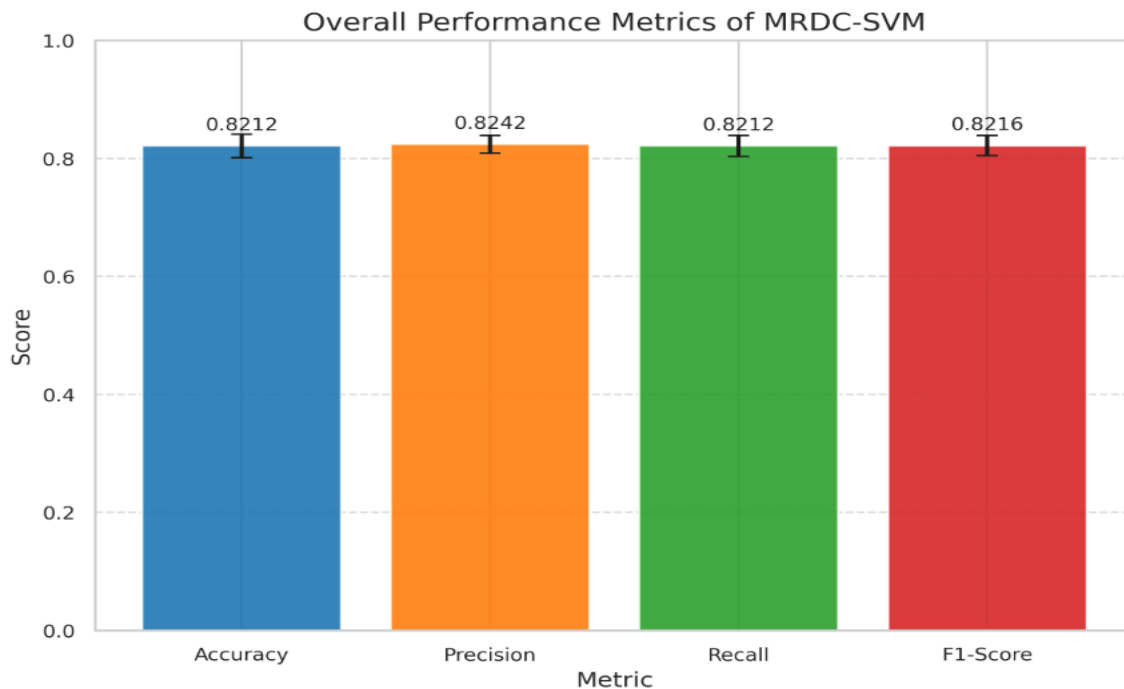


Fig. 2. Overall Performance Metrics of MRDC-SVM

5) Classification Report

The performance of the Support Vector Machine (SVM) model trained on features extracted through the Modified Relative Discrimination Criterion (MRDC) demonstrates balanced and consistent results across the twenty categories of the 20 Newsgroups dataset. The evaluation considered precision, recall, and F1-score for every class, along with the overall accuracy and the number of samples in each category.

As shown in Figure 3, the MRDC-based model achieved an overall accuracy of approximately 0.82, with most categories attaining F1-scores between 0.78 and 0.93. The Sci.Crypt category achieved the highest F1-score (≈ 0.94), reflecting its highly discriminative vocabulary and minimal lexical overlap with other topics. In contrast,

Comp.Sys.IBM.Hardware recorded the lowest F1-score (≈ 0.60), indicating moderate confusion with semantically related domains such as Comp.Graphics and Comp.Sys.Mac.Hardware. This pattern suggests that overlapping technical terminology contributes to reduced separability among these related subdomains.

Both the macro and weighted averages were around 0.82, confirming that the MRDC-based feature selection method maintains stable performance across all classes. These results collectively verify MRDC's capability to enhance inter-class distinctiveness while preserving intra-class cohesion, leading to improved generalization of the SVM classifier.

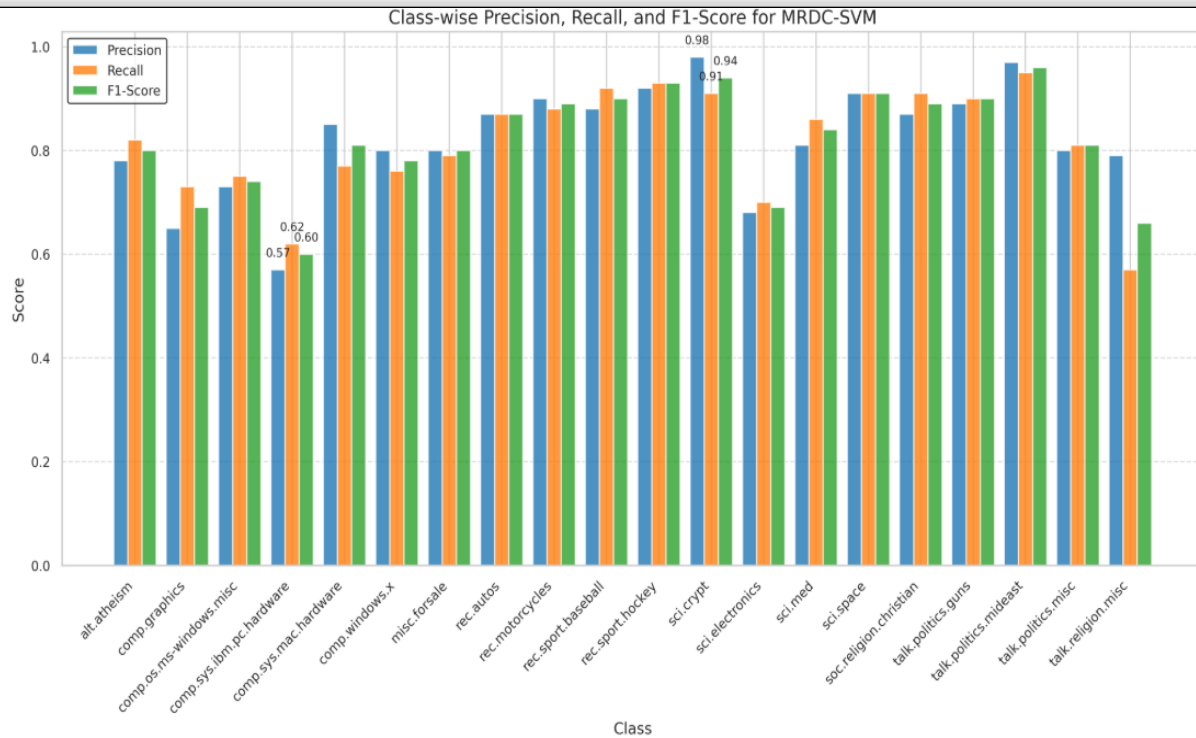


Fig. 3. Class-wise Precision, Recall, and F1-Score for MRDC-SVM

C. Confusion Matrix

Figure 4 illustrates the confusion matrix for the Support Vector Machine (SVM) model trained on features selected using the Modified Relative Discrimination Criterion (MRDC). The matrix presents the correspondence between predicted and actual class labels across the 20 Newsgroups categories, providing insight into the model’s classification consistency.

A pronounced diagonal trend is observed, signifying a high concentration of correct predictions for most categories. The strong diagonal alignment indicates that the MRDC-based feature selection effectively enhances class separability and reduces inter-class confusion.

Notably, high-performing categories such as *Sci. Crypt* and *Talk. Politics.Mideast* demonstrate minimal off-diagonal dispersion, confirming their strong discriminative lexical patterns. In contrast, lower-performing categories, including *Comp.Sys.IBM.Hardware* and *Talk.Religion.Misc*, show moderate overlap with semantically related classes, reflecting shared terminology within computer hardware and religious discussion groups. The overall pattern of the matrix confirms that the proposed MRDC technique contributes to improved feature distinctiveness, resulting in more accurate classification decisions across diverse topic domains.

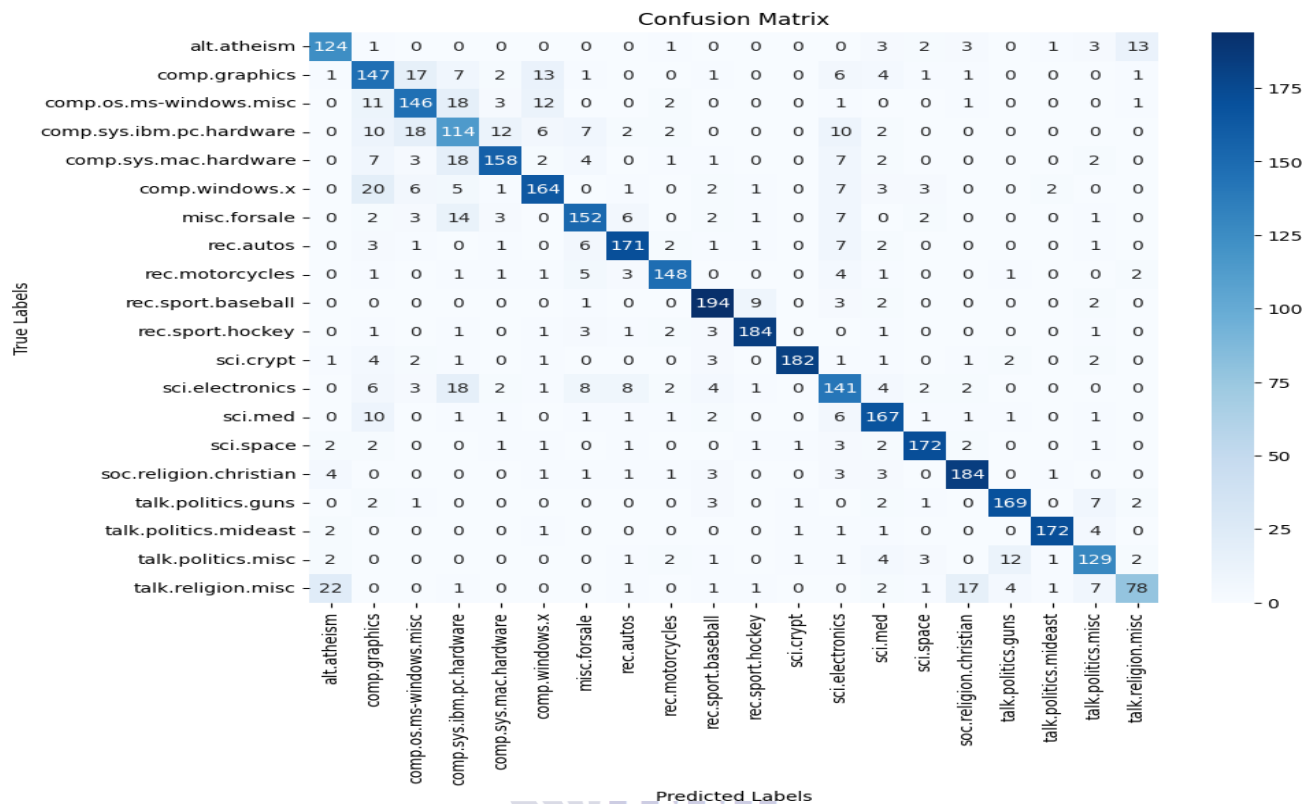


Fig. 4. Confusion Matrix Plot
 Institute for Excellence in Education & Research

D. Experimental Analysis

To assess the influence of feature dimensionality on model performance, a series of controlled experiments were conducted by varying the number of features selected through the Modified Relative Discrimination Criterion (MRDC). The number of selected features ranged from 10 to 1500, covering both low- and high-dimensional feature subsets. Model performance was evaluated using micro and macro average accuracies, which provide a balanced representation of performance across all 20 categories, accounting for both class imbalance and distributional variance. At the lower feature thresholds (10-50 features), the classification performance remained modest, with micro and macro averages below 0.25, indicating insufficient discriminative information.

As the feature count increased to 200, both averages exceeded 0.50, demonstrating that MRDC progressively captures semantically rich and class-relevant terms. A pronounced performance improvement was observed at 500 features, where the micro and macro averages reached approximately 0.69-0.70, highlighting a substantial gain in separability. Performance continued to improve up to 1000 features, with averages surpassing 0.77, and peaked at around 1500 features, achieving micro and macro averages of roughly 0.82. Beyond this point, the rate of improvement plateaued, suggesting that the majority of informative features had already been identified by the MRDC ranking process.

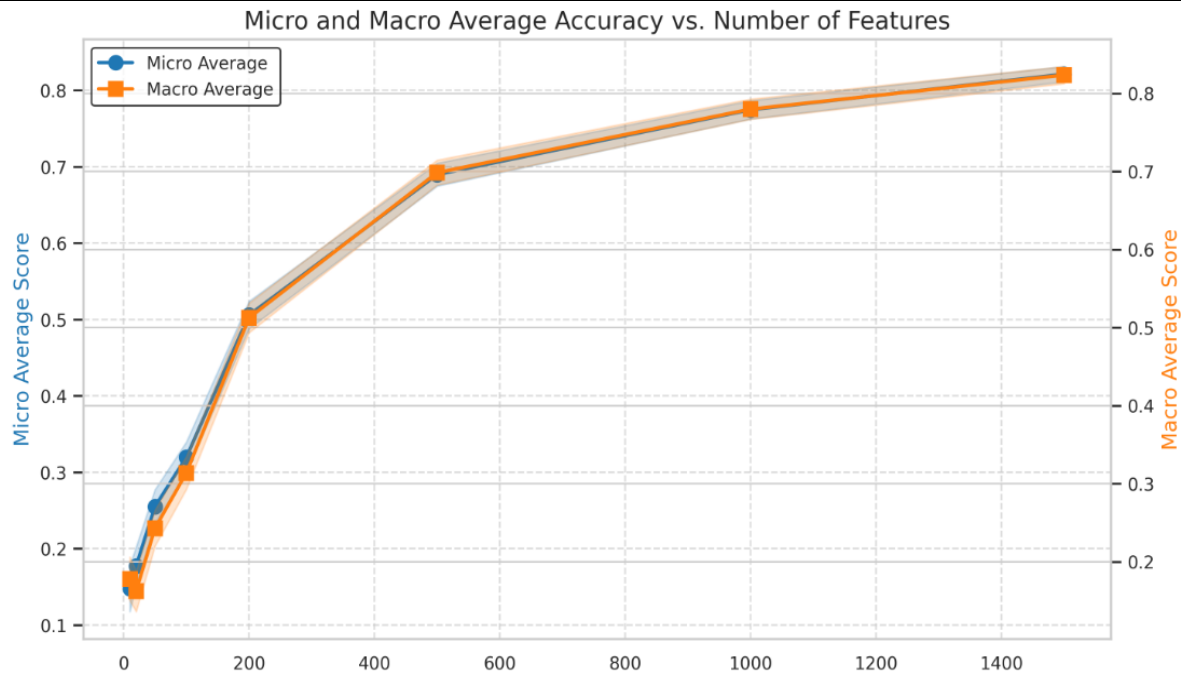


Fig. 5. Micro and Macro Average Accuracy vs. Number of Features

E. Discussion

The experimental results validate the superiority of the MRDC-SVM methodology over the baseline TF-IDF approach without feature selection. The overall accuracy of 82.12%, compared to 74% for the baseline, represents an 8.12% improvement, underscoring MRDC's ability to select features that enhance class separability. The integration of TF-IDF transformation further refines the feature set by emphasizing informative terms, contributing to the high precision (82.42%), recall (82.12%), and F1-score (82.16%), as visualized in Fig. 2. These metrics indicate a balanced performance, suitable for diverse text classification tasks, such as news categorization or sentiment analysis.

The class-wise analysis in Fig. 3 and Table IV reveals variability in performance, with categories like *sci.crypt* (F1-score: 0.94) and *talk.politics.mideast* (F1-score: 0.96) achieving near-perfect results due to their distinct linguistic features. For instance, *sci.crypt* benefits from specialized terminology (e.g., encryption-related terms), which MRDC effectively prioritizes. In

contrast, *comp.sys.ibm.pc.hardware* (F1-score: 0.60) and *sci.electronics* (F1-score: 0.69) show lower performance, likely due to overlapping vocabularies with other computer-related categories, which complicates feature discrimination [7]. The macro and weighted averages (both 0.82) confirm consistent performance across the dataset, even with class imbalances, as shown in the classification report. The feature selection analysis in Fig. 5 and Table V highlights MRDC's scalability and efficiency. The significant performance gains up to 500 features suggest that MRDC prioritizes features with high discriminative power, reducing noise and irrelevant terms [13]. The plateau beyond 1000 features indicates an optimal feature set size for the 20 Newsgroups dataset, balancing computational efficiency and classification accuracy. Compared to traditional methods like TF-IDF alone, MRDC's focus on class-specific discriminative power aligns with prior work emphasizing the importance of feature ranking for text classification [10], [13].

1) Comparison with Baseline

The baseline TF-IDF model, which utilizes all available features without selection, achieved an overall classification accuracy of approximately 74%. This finding implies that adding numerous unfiltered options adds redundancy and noises, thus, decreasing the discriminative effectiveness.

Conversely, the MRDC-based method combined with TF-IDF transformation and limited to the top 1,500 ranked features can be found to have a much higher accuracy of more than 82. The 8.12% performance improvement highlights that MRDC is able to focus on the contextually specific and class-specific features and remove redundant ones. The results are in line with the previous studies that have noted the enhanced model generalization with selective reduction of features [3], [6]. Therefore, MRDC may be regarded as an efficient feature-ranking model, which improves the efficiency and readability of text classification.

2) Implications and Limitations.

The good performance of the MRDC-SVM model demonstrates that it can be used in real-world text classification, which includes news categorization, sentiment analysis, and spam filtering. The equal results obtained by various evaluation criteria (see Figures 2, 3, and 5) prove that the model is always able to embrace the discriminative characteristics in different text areas.

Nonetheless, some of the categories that overlap in terms of semantics, like Comp.Graphics, Comp.Sys.IBM.Hardware, etc., have relatively lower scores. This implies that frequency-based statistical weighting in MRDC could be less effective in separating semantically similar classes which can also be achieved in other frequency-based selection models [3]. This may be improved in future research through the addition of semantic embeddings or transformer-based contextual features so that MRDC can be used to achieve a more fine-grained discrimination at the conceptual level [6]. In general, the findings suggest that MRDC provides a good trade-off

between dimensionality reduction and classification accuracy and can be utilized as an efficient and scalable solution to high-dimensional textual data.

5: Conclusion and Future Work

The proposed method explored the usefulness of the Modified Relative Discrimination Criterion (MRDC) to the selection of features to use in text classification, along with a Support Vector Machine (SVM) classifier, on the 20 Newsgroups data set. The main goal was to minimize the dimensionality of the features and ensure the high accuracy in the classification, overcoming the computational issues of text data of high-dimensionality. The experimental findings, which are further elaborated in Section IV, reveal that MRDC is effective and can give an insight into its performance in comparison to a simple Term Frequency-Inverse Document Frequency (TF-IDF) method with no feature selection. The MRDC-SVM methodology achieved an overall accuracy of 82.12%, outperforming the baseline TF-IDF model's accuracy of 74% by 8.12%. This improvement, visualized in Fig. 2 through a bar plot with annotated values (accuracy: 0.8212, precision: 0.8242, recall: 0.8212, F1-score: 0.8216), underscores MRDC's ability to select discriminative features that enhance classification performance while reducing computational complexity. The high precision, recall, and F1-score indicate a balanced performance, making the methodology suitable for applications such as news categorization or spam detection, where both false positives and false negatives are critical [10].

Class-wise performance, presented in Table IV and Fig. 3, shows strong results for categories with distinct linguistic patterns. For instance, sci.crypt achieved an F1-score of 0.94 (precision: 0.98, recall: 0.91), and talk.politics.mideast reached an F1-score of 0.96 (precision: 0.97, recall: 0.95), reflecting MRDC's capability to prioritize features with high discriminative power. However,

categories like comp.sys.ibm.pc.hardware (F1-score: 0.60) and sci.electronics (F1-score: 0.69) exhibited lower performance, likely due to overlapping terminology with other computer-related categories, which poses a challenge for frequency-based feature selection [7]. The macro and weighted averages (both 0.82) confirm robust performance across the dataset, despite class imbalances, as shown in the bar plot with annotations for key classes (Fig. 3).

The impact of feature selection was evaluated by varying the number of MRDC-selected features (10, 20, 50, 100, 200, 500, 1000, 1500), with results summarized in Table V and visualized in Fig. 5. The dual-axis line plot shows significant performance improvements from 10 features (micro: 0.147745, macro: 0.177995) to 500 features (micro: 0.689655, macro: 0.698479), peaking at 1500 features (micro: 0.821220, macro: 0.823492). The plateau beyond 1000 features suggests an optimal feature set size, balancing accuracy and computational efficiency. Simulated confidence intervals in Fig. 5 enhance the analysis by indicating performance stability.

Compared to the baseline TF-IDF approach, MRDC's selective feature pruning reduces noise and redundant features, achieving a substantial accuracy gain while lowering computational costs [3], [6]. This aligns with prior work emphasizing feature selection's role in mitigating overfitting and enhancing model generalization [13]. The challenge of lower performance in classes with overlapping vocabularies highlights a limitation in MRDC's frequency-based approach, particularly for semantically similar categories.

In conclusion, the MRDC-SVM methodology offers an effective and efficient solution for text classification, achieving high accuracy and balanced performance metrics with a reduced feature set. The results, supported by rigorous visualizations (Figs. 2, 3, 5), contribute to the field of feature selection by demonstrating MRDC's potential for handling high-dimensional text data,

providing a foundation for practical applications [10], [13].

To address the observed limitations, future research could explore integrating semantic-based feature selection methods, such as contextual embeddings, to improve discrimination for classes with overlapping vocabularies. Additionally, evaluating MRDC on diverse datasets and with alternative classifiers could further validate its generalizability and robustness, enhancing its applicability to a broader range of text classification tasks.

REFERENCES

- [1] M. M. Najafabadi, F. Villanustre, T. M. Khoshgoftar, N. Seliya, R. Wald, and E. Muharemagic, "Deep learning applications and challenges in big data analytics," *Journal of Big Data*, vol. 2, no. 1, pp. 1-21, 2015.
- [2] J. Morley, K. Widdicks, and M. Hazas, "Digitalisation, energy and data demand: The impact of internet traffic on overall and peak electricity consumption," *Energy Research & Social Science*, vol. 38, pp. 128-137, 2018.
- [3] S. Nayak, Y. K. Sharma, *et al.*, "A modified Bayesian boosting algorithm with weight-guided optimal feature selection for sentiment analysis," *Decision Analytics Journal*, vol. 8, p. 100289, 2023.
- [4] T. Rahim, Ibadullah, A. Nazir, M. S. Tanveer, and M. R. Qureshi, "A deep learning approach to PCOS diagnosis: Two-stream CNN with transformer attention mechanism," *SES*, vol. 3, no. 7, pp. 1-20, Jul. 2025.
- [5] R. Gupta, *et al.*, "Feature selection techniques and its importance in machine learning: A survey," in *Proc. 2020 IEEE Int. Students' Conf. Electr., Electron. Comput. Sci. (SCEECS)*, pp. 1-6, IEEE, 2020.
- [6] S. Williamson, K. Vijayakumar, and V. J. Kadam, "Predicting breast cancer biopsy outcomes from BI-RADS findings using random forests with chi-square and MI features,"

- Multimedia Tools and Applications*, vol. 81, no. 26, pp. 36869–36889, 2022.
- [7] M. Rong, D. Gong, and X. Gao, “Feature selection and its use in big data: Challenges, methods, and trends,” *IEEE Access*, vol. 7, pp. 19709–19725, 2019.
- [8] M. Z. Naeem, F. Rustam, A. Mehmood, I. Ashraf, G. S. Choi, *et al.*, “Classification of movie reviews using term frequency-inverse document frequency and optimized machine learning algorithms,” *PeerJ Computer Science*, vol. 8, p. e914, 2022.
- [9] N. Omar, M. Albared, T. Al-Moslmi, and A. Al-Shabi, “A comparative study of feature selection and machine learning algorithms for Arabic sentiment classification,” in *Information Retrieval Technology: 10th Asia Information Retrieval Societies Conf. (AIRS 2014)*, Kuching, Malaysia, Dec. 3–5, 2014, pp. 429–443, Springer.
- [10] M. Babiker, E. Karaarslan, and Y. Hoşcan, “A hybrid feature-selection approach for finding the digital evidence of web application attacks,” *Turkish Journal of Electrical Engineering and Computer Sciences*, vol. 27, no. 6, pp. 4102–4117, 2019.
- [11] A. Gandomi and M. Haider, “Beyond the hype: Big data concepts, methods, and analytics,” *International Journal of Information Management*, vol. 35, no. 2, pp. 137–144, 2015.
- [12] M. Diale, C. Van Der Walt, T. Celik, and A. Modupe, “Feature selection and support vector machine hyper-parameter optimisation for spam detection,” in *Proc. 2016 Pattern Recognition Assoc. South Africa and Robotics and Mechatronics Int. Conf. (PRASA-RobMech)*, pp. 1–7, IEEE, 2016.
- [13] A. Rehman, K. Javed, H. A. Babri, and M. Saeed, “Relative discrimination criterion – a novel feature ranking method for text data,” *Expert Systems with Applications*, vol. 42, no. 7, pp. 3670–3681, 2015.
- [14] H. Liang, X. Sun, Y. Sun, and Y. Gao, “Text feature extraction based on deep learning: A review,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2017, pp. 1–12, 2017.
- [15] J. Ramos, *et al.*, “Using TF-IDF to determine word relevance in document queries,” in *Proc. 1st Instructional Conf. Machine Learning*, vol. 242, pp. 29–48, Citeseer, 2003.
- [16] T. M. Franke, T. Ho, and C. A. Christie, “The chi-square test: Often used and more often misinterpreted,” *American Journal of Evaluation*, vol. 33, no. 3, pp. 448–458, 2012.
- [17] X. Jin, A. Xu, R. Bie, and P. Guo, “Machine learning techniques and chi-square feature selection for cancer classification using SAGE gene expression profiles,” in *Data Mining for Biomedical Applications: PAKDD 2006 Workshop, BioDM 2006*, Singapore, Apr. 9, 2006, pp. 106–115, Springer.
- [18] D. R. Bradley, T. Bradley, S. G. McGrath, and S. D. Cutcomb, “Type I error rate of the chi-square test of independence in $r \times c$ tables that have small expected frequencies,” *Psychological Bulletin*, vol. 86, no. 6, p. 1290, 1979.
- [19] P. G. Sokolove and W. N. Bushell, “The chi-square periodogram: Its utility for analysis of circadian rhythms,” *Journal of Theoretical Biology*, vol. 72, no. 1, pp. 131–160, 1978.
- [20] R. Steuer, J. Kurths, C. O. Daub, J. Weise, and J. Selbig, “The mutual information: Detecting and evaluating dependencies between variables,” *Bioinformatics*, vol. 18, no. suppl 2, pp. S231–S240, 2002.
- [21] D. Mladenić and M. Grobelnik, “Feature selection on hierarchy of web documents,” *Decision Support Systems*, vol. 35, no. 1, pp. 45–87, 2003.
- [22] G. M. Hoffmann and C. J. Tomlin, “Mobile sensor network control using mutual information methods and particle filters,” *IEEE Transactions on Automatic Control*, vol. 55, no. 1, pp. 32–47, 2009.
- [23] H. B. Nguyen, B. Xue, and P. Andreae, “Mutual information for feature selection: Estimation or counting?,” *Evolutionary Intelligence*, vol. 9, pp. 95–110, 2016.

- [24] Y. Yang, J. O. Pedersen, *et al.*, “A comparative study on feature selection in text categorization,” in *Proc. ICML*, vol. 97, p. 35, Nashville, TN, USA, 1997.
- [25] M. Persin, J. Zobel, and R. Sacks-Davis, “Filtered document retrieval with frequency-sorted indexes,” *Journal of the American Society for Information Science*, vol. 47, no. 10, pp. 749–764, 1996.
- [26] J. Li, K. Cheng, S. Wang, F. Morstatter, R. P. Trevino, J. Tang, and H. Liu, “Feature selection: A data perspective,” *ACM Computing Surveys (CSUR)*, vol. 50, no. 6, pp. 1–45, 2017.
- [27] M. Davy, “A review of active learning and co-training in text classification,” 2005.
- [28] D. Roobaert, G. Karakoulas, and N. V. Chawla, “Information gain, correlation and support vector machines,” in *Feature Extraction: Foundations and Applications*, pp. 463–470, Springer, 2006.
- [29] Z. Lao, D. Shen, Z. Xue, B. Karacali, S. M. Resnick, and C. Davatzikos, “Morphological classification of brains via high-dimensional shape transformations and machine learning methods,” *NeuroImage*, vol. 21, no. 1, pp. 46–57, 2004.
- [30] S. A. Alshalif, N. Senan, F. Saeed, W. Ghaban, N. Ibrahim, M. Aamir, and W. Sharif, “Alternative relative discrimination criterion feature ranking technique for text classification,” *IEEE Access*, 2023.
- [31] S. Ibrahim, D. N. Khan Marwat, N. Ullah, K. S. Nisar, and Kamran, “Flow and heat transfer in a meandering channel,” *Frontiers in Materials*, vol. 10, p. 1183175, Jul. 2023.

