

Article

Not peer-reviewed version

---

# Evaluation and Comparison of SVM, Deep Learning, and Naïve Bayes Performances for Natural Language Processing Text Classification Task

---

[Destiny Ogaga](#)<sup>\*</sup> and Abiodun Olalere

Posted Date: 23 November 2023

doi: 10.20944/preprints202311.1462.v1

Keywords: text classification; natural language processing; multinomial naïve bayes; support vector machine; deep learning



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

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

*Article*

# Evaluation and Comparison of SVM, Deep Learning, and Naïve Bayes Performances for Natural Language Processing Text Classification Task

Destiny Ogaga <sup>1,\*</sup> and Abiodun Olalere <sup>2</sup>

<sup>1</sup> Department of Computer Science, University of Hertfordshire, Hatfield, United Kingdom

<sup>2</sup> Department of Computer Science, Espam Formation University, Cotonou, Benin; Shervie20@gmail.com

\* Correspondence: destiny@schoolnigeria.com.ng

**Abstract:** Text classification is one of the most important task in natural language processing, In this research, we carried out several experimental research on three (3) of the most popular Text classification NLP classifier in Convolutional Neural Network (CNN), Multinomial Naive Bayes (MNB), and Support Vector Machine (SVN). In the presence of enough training data, Deep Learning CNN works best in all parameters for evaluation with 77% accuracy, followed by SVM with accuracy of 76%, and multinomial Bayes with least performance of 69% accuracy. CNN has the best performance in the presence of large enough training dataset because of the presence of filter/kernels which help to indentify patterns in text data regardless of their position in the sentence. We repeated the training again with just one-third of our data, at this point SVM comes with the best performance, the performance of CNN noticeably drops but still better than Multinomial Naive Bayes, the reason why SVM works best when we reduce the training data was because of its ability to look for a hyper-plane that creates a boundary between different classes of data so as to properly classify them, so we believed that getting the hyper-plane was more efficient when we reduce the dataset, hence reason for the good performance. Multinomial Naive Bayes have the least performance which we attributed to its assumption of independency between the features which sometimes does not hold true. We concluded that availability of data should be an important factor when choosing classifier for Natural Language Processing Text Classification task. CNN should be use in the presence of enough dataset, and SVM should be use when data is not enough. Multinomial Naive Bayes must not be trusted with state of the art NLP task due to its assumption of independency between the features.

Keywords: text classification; natural language processing; multinomial naïve bayes; Support Vector Machine; deep learning

## I. Introduction

One of the most important problems in natural language processing is text classifications which are applicable in so many areas of modern technology and innovation. With text as one of the most common categories of data available, more than 80% of data are unstructured. This makes it difficult and challenging to comprehensively analyze them, hence many businesses are unable to exploit its full potential for business benefit [7]. The inability to exploit full potential of billions of unstructured data brings about the suitability of machine learning algorithm for training model for text classification and sentiment analysis in areas like web search [4], spam detection [5], topic classification [6], news classification, sentiment classification [7,22] to the spotlight, as machine learning algorithms can be use in Natural Language Processing for text classification and sentiment [13–16] related activities such as spam filtering and other use of text classification as a core technique for machine learning process.

The aim of this research is evaluate and compare performances of each of the three of the most important classifier commonly used for state of the art natural language processing text classifier and

sentiment analysis. We experimented this by implementing each of Linear Support Vector Machine, deep learning through Convolutional Neural Network (CNN), and multinomial variant of Bayesian classifier for text classification, after evaluation and comparison of the accuracy and performances of each classifier, we also dived into causes of variation and differences in their performances.

For all implementations in this research, we use stack overflow dataset which is publicly available at: <https://storage.googleapis.com/tensorflow-workshop-examples/stack-overflow-data.csv> dataset and can also be directly query from Google BigQuery platform.

## II. RELATED WORK

### 1) Convolutional Neural Network (deep learning based)

Artificial neural network works to mimic functionality of human brain, input in artificial neural network represents the dendrites find in human brain while the different axon terminals then represent output from the neural network. In deep learning, a typical network contains one or more several hidden layers called deep neural network (Figure 1), and it works by performing computation on input data fed to the network to give an output.

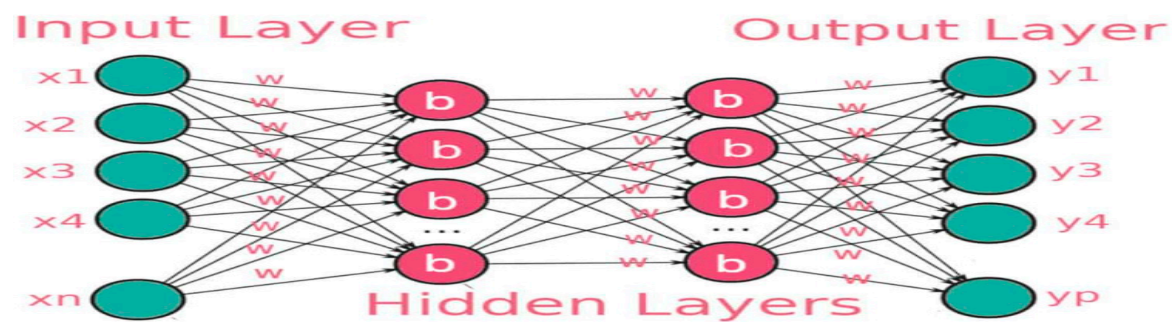


Figure 1. Hidden Layers of Convolutional Neural Network (CNN).

Convolution Neural Networks(CNNs) are complex multi-layered artificial neural networks with the ability to detect complex features such as extraction of features in images and text in dataset. They are very efficient in computer vision task and so are commonly used for image segmentation, object detection, image classification, and text classification tasks [19,20] for natural language processing.

Convolutional neural network has convolution layer and pooling layer (Figure 2) as the two major layers for separate purposes, and while the convolution layer obtains features from data, pooling layer reduces sizes of the feature map from the data.

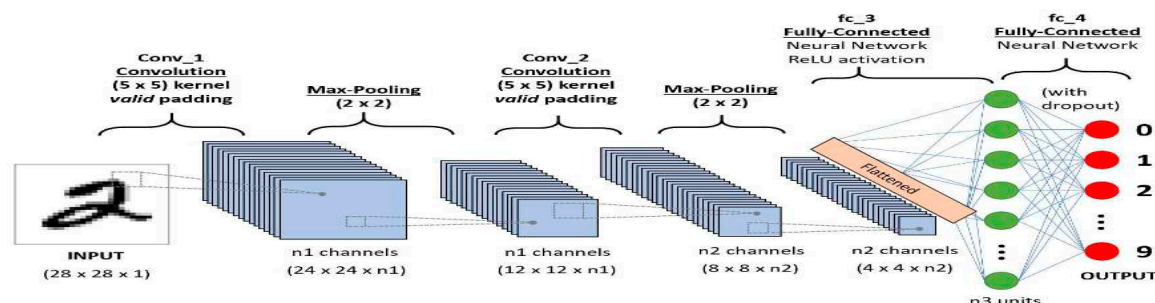


Figure 2. Connectivity of the Pooling and Kernels in CNN.

During convolution, features are obtained from data, features that are obtained during convolution operation are fed to CNN. Outputs from this convolution operation are the most important features and they are called feature map, a convolved feature, or an activation map. The output is computed by applying feature detector which is also known as kernel or filter to the input data. The next stage of the process is computation of the feature map by multiplication of the kernel

and the matrix representation of the input data together, this process ensures that the feature map that is passed to the CNN is smaller but contains all essential features. This process is done by filter by going step by step process known as strides through every element in the input data.

## 2) SUPPORT VECTOR MACHINE (SVM)

Support Vector Machine (SVM) is a supervised machine algorithm used for both classification and regression task [17,18,21]. It works by looking for a hyper-plane (Figure 3) that creates a boundary between two classes of data so as to properly classify them, it determines the best decision boundary between categories, hence they can be applied to vector which can encode data, and so text are classified into vector during text classification tasks by SVM algorithm.

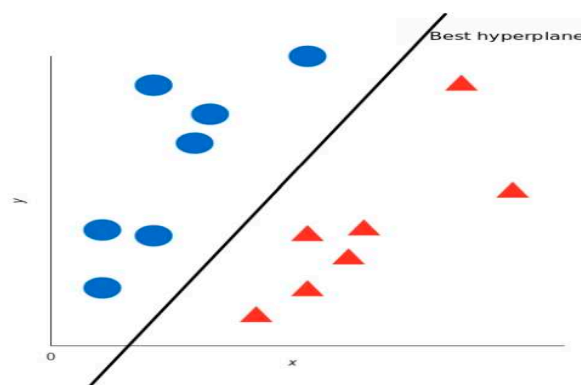


Figure 3. Hyper-plane mapping region of separation.

Once the algorithm determined the decision boundary for each of the category in the dataset for which we want to analyze, we can proceed to obtain the representations of each of all the texts we want to classify in our NLP [9–12] and check for the side of the boundary that those representations fall into.

### Theories supporting suitability of SVM for text classification

**High dimensional input space:** when training a classifier for text classification task, it is common to deal with several features, but the fact that SVM tries to prevent overfitting by using overfitting protective measure which are independent on the number of features in the data gives SVM the potential to handle feature spaces in dataset, and hence their suitability for tasks relating to natural language processing.

**Linearly separable properties of text categorization:** Virtually, all of categories involved in text classification are linearly separable. And so, the fact that SVM works by finding linear separator between each of the categories in the textual dataset makes SVM suitable for tasks relating to NLP such as text classification, sentiment analysis, cyber troll and so on.

## 3) Multinomial Naïve Bayes (MNB)

Multinomial variant of Bayesian classifier is an algorithm commonly used for text classification related tasks and also problems with multiple classes [8]. In order to understand how Bayesian classifier works, it is important to understand the basic concept of Bayes theorem.

Tosin Ige, and Sikiru Adewale [10] successfully use multinomial Naïve Bayes algorithm to developed a machine learning based model along with an automated chatbot that can identify and intercept bully messages in an online chat conversation to protect potential victim with 92% accuracy. The accuracy of their model is a big discovery in Natural language processing field of artificial intelligence.

In Bayes theorem, the probability of an event occurring based on the prior knowledge of conditions related to an event is calculated based on the formula:

$$P(A|B) = P(A) * P(B|A)/P(B).$$

Where we are calculating the probability of class A when predictor B is already provided.

$P(B)$  = prior probability of B

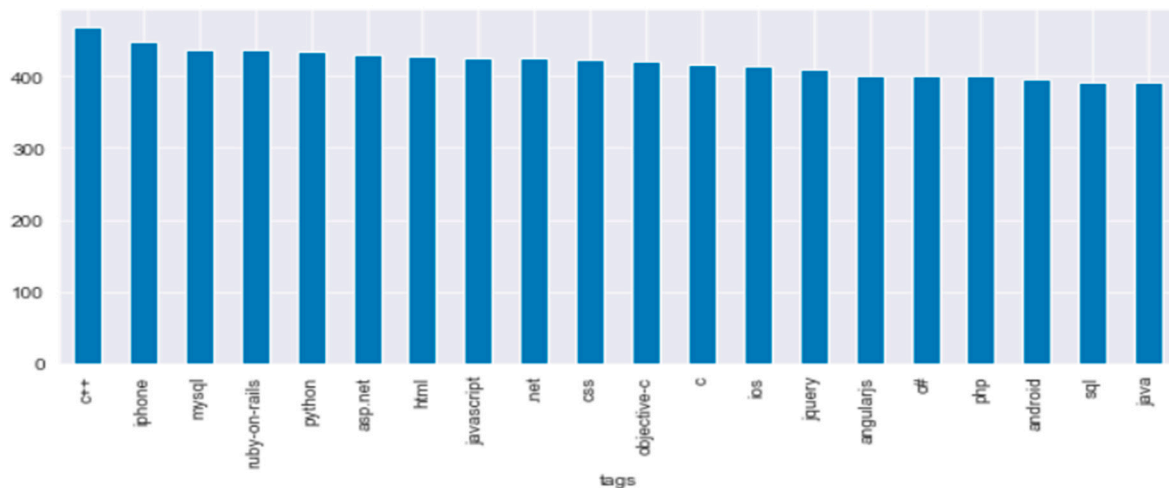
$P(A)$  = prior probability of class A

$P(B|A)$  = occurrence of predictor B given class A probability

Base on the principle of this calculation, we can automate the computation of tags in text and categorize them.

### III. Research Methodology

1) *Dataset*: For each of the implementations in this research, we use stack Overflow questions and tags dataset which is publicly available at: <https://storage.googleapis.com/tensorflow-workshop-examples/stack-overflow-data.csv> and can also be query from Google BigQuery platform. The dataset contains two main columns, the first column contains the question for all non-deleted Stack Overflow questions in the database while the second column is Tag which contains the tags on each of these questions.



**Figure 4.** Tag Classification Visualization from Stackoverflow question and answer dataset.

2) *Pre-processing and Feature Engineering*: Since our dataset is a collection of several posts and tags from stackoverflow dataset, it is imperative to clean the data of several unwanted characters. This necessitated the need for pre-processing and feature engineering before we could actually use the data. As part of the pre-processing step, we removed unwanted characters from text, remove punctuation mark, stopwords, search for and decoding HTML, and so on, and then we finally split the dataset into training and validation dataset as part of our final data pre-processing steps.

After pre-processing followed feature engineering during which we converted each of the text to matrix of token count by CountVectorizer. The count matrix is further converted to a normalized tf-idf representation also known as tf-idf transformer. Haven completed both pre-processing and feature engineering process, we then proceeded to train our model on SVM, Multinomial, and CNN classifiers.

#### 3) IMPLEMENTATION

We used pipeline which serves as a compound classifier in scikit-learn, each unique word was assigned an index while we used tokenizer to count. Parameter for number of words was passed to the tokenizer to ensure that our vocabulary is limited to only top words. Haven done this; we were able to use our tokenizer along with texts\_to\_matrix method to create proper training data that could be passed to the model. After feeding our model with one hot encoded vector, and the transformation



of features and labels to a format such that it could be read by keras deep learning python library. We trained the model by passing training data and labels, batch size and epochs to the fit() method. As for the SVM implementation, we use CountVectorizer() method, TfidfTransformer() method, after which we called SGDClassifier with the following configuration argument

```
loss='hinge',
penalty='l2',
alpha=1e-3,
random_state=42,
max_iter=5,
tol=None
```

In the case of training our model with Convolutional Neural Network which requires larger amount of training dataset for optimal performance, we set the validation\_split to 0.1 to ensure that 90% of the dataset is use for training while the remaining 10% is use for validation, we set the batch size to 4 and number of epochs to 50, and then monitor both the result and performance of each of the models.

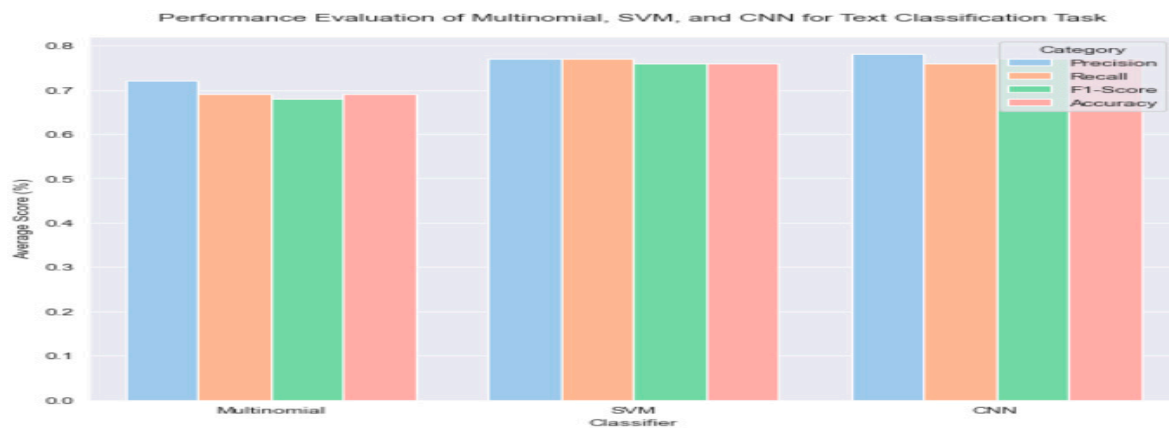
#### IV. Result and Discussion

First and foremost, we observed some notable differences in training time for each of the classifiers. We got different performances and results for each of the experimental implementation of the three classifiers. Our first attention was drawn to the performance of Multinomial Naïve Bayes (Table 1) as it has the worst performance for each of the evaluation criteria among the three (3). In order to investigate this we, we used another dataset with few data to train the model, it was at this stage that we got a better result. This clearly proof to us that Multinomial variants of Bayesian classifier is good for small dataset, the reason for this is not far fetch, since multinomial naïve Bayes assume independency between features in the dataset, as the data gets bigger, the assumption of independency tends not to be true for some of the features in the dataset, hence the drop in performances with large dataset.

**Table 1.** Metric Scores from Different Evaluation Criteria.

<i>Classifier</i>	<i>Average Precision</i>	<i>Average Recall</i>	<i>Average F1-Score</i>	<i>Average Accuracy</i>
<b>Multinomial Naïve Bayes (MNB)</b>	0.72	0.69	0.68	0.69
<b>Support Vector Machine (SVM)</b>	0.77	0.77	0.76	0.76
<b>Deep Learning (CNN)</b>	0.78	0.76	0.77	0.77

As for Support Vector Machine (SVM) classifier whose performance is in-between that of Multinomial and CNN classifier (Figure 4). SVM works better than Multinomial for two reasons, the first reason is because it is not based on any independence assumption among the features, the second reason is because of its ability to look for a hyperplane that creates a boundary between two classes of data so as to properly classify them, as this enables it to determine the best decision boundary between the categories. Although, SVM works better than Multinomial Naïve Bayes for Natural Language Processing text classification tasks, it is also not suitable for a very large dataset due to complexity in training.



**Figure 5.** Comparing performances of CNN, SVM, and Multinomial across different evaluation metrics.

For each of all the evaluation parameters for our experimental implementation, Convolutional Neural Network (CNN) has the best overall performance, with average precision (78%), average recall (76%), average F1-Score (76%), average accuracy (77%). We believed that deep learning convolutional neural network work best among the three (3) for natural language processing text classification task because it contains filter/ kernels which can help to identify patterns in text data, also the fact that these filters are translational invariant means that they can detect patterns regardless of their position in the sentence. Here, the convolutional architecture can identify -gram that could be use for prediction with the need to pre-specify any embedding vector for ngram.

## V. Evaluation

Since each of the evaluation techniques have its own drawback, so in order to have better evaluation of our model and to ensure effective comparison of performances of each of the classifier on Natural Language Processing text classification task, we combined different evaluation techniques to effectively evaluate the performance of each of the models and compare result. We used average precision score, average recall score, average F1-Score, and average accuracy as our evaluation parameters for comparison;

### ACCURACY

This is the actual ratio of the number of correct predictions to the total number of predictions. It represents the most fundamental of all the evaluating metrics used to evaluate model performance. The formula is given by.

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

As fundamental as it is, it performs poorly in the presence of an imbalance dataset. Suppose a model classifies that most of the data belongs to the major class label. It yields higher accuracy. But in general, the model cannot classify on minor class labels and has poor performance.

### PRECISION AND RECALL

Precision is the ratio of true positives to the summation of true positives and false positives. It basically analyses the positive predictions.

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

The drawback of Precision is that it does not consider the True Negatives and False Negatives.

Recall is the ratio of true positives to the summation of true positives and false negatives. It analyses the number of correct positive samples.

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

The drawback of Recall is that often it leads to a higher false positive rate.

### F1 SCORE

This is the harmonic mean of precision and recall. It is well known that there is precision-recall trade-off such that if we increase the precision, recall decreases and vice versa. F1 score evaluation

technique combines both the precision and recall score to give harmonic mean to better evaluate model performance.

$$F1 \text{ score} = (2 \times \text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$$

## VI. Conclusion

In this research, we confirmed that of the three (3) of the most popular classifier for NLP text classification task, Convolutional Neural Network work best and even with all parameters for metric evaluation when we have enough training dataset, CNN is good for text classification in the presence of enough data because due to the presence of filter/ kernels which help to indentify patterns in text data regardless of their position in the sentence.

In the absence of enough training data, Support Vector Machine (SVM) works best which we believed to be due to its ability to look for a hyperplane that creates a boundary between different classes of data so as to properly classify them. Multinomial Naive Bayes works based on independent assumption between features which are sometimes not valid in real data, we believe this is the reason for its least performance among the three, and we believed that Multinomial Bayes classifiers must not be trusted for state of the Natural Language Processing (NLP) text classification task.

## VII. Future Research Direction

Bayesian classifier works best when two conditions are met. The first is the generic condition of independence that all features are independent of each other which rarely holds in real life. The fact that this condition rarely holds in real life limits the applicability of naïve Bayes in real-world use cases as they become unsuitable where there is any full or partial dependency between any of the features in the data. The second condition is based on individual assumption of each variant of Bayesian classifier which does not always hold. The workings of each of the existing variants of Bayesian classifier are based on different assumption which is the single most important factor that influences their performance and accuracy. This explains why each of the existing variants of Bayesian classifier is suitable for different classification tasks depending on the nature of distribution of the dataset if it agrees with their assumption. Each of the current variants of Bayesian classifiers works well if only the distribution of the data agrees with their assumption and the generic condition of independence holds.

In order to improve on the performance of Naïve bayes classifier, additional research work is needed to lower each of the two levels aof assumption of the Bayesian variants to the barest minimum. Lowering these assumptions will significantly improve the performances of naïve bayes.

## References

1. Ige, T., & Kiekintveld, C. (2023). Performance Comparison and Implementation of Bayesian Variants for Network Intrusion Detection. *arXiv preprint arXiv:2308.11834*.
2. Sentiment Analysis of Internet Posts," 2019 IEEE Eurasia Conference on IOT, Communication and Engineering (ECICE), Yunlin, Taiwan, 2019, pp. 154-155. <https://doi.org/10.1109/ECICE47484.2019.8942736>.
3. J. Zhao, L. Dong, J. Wu, K. Xu. MoodLens: An Emoticon-Based Sentiment Analysis System for Chinese Tweets in Weibo. KDD 2012.).
4. Y. Chen and Z. Zhang, "Research on text sentiment analysis based on CNNs and SVM," 2018 13th IEEE Conference on Industrial Electronics and Applications (ICIEA), Wuhan, China, 2018, pp. 2731-2734. <https://doi.org/10.1109/ICIEA.2018.8398173>.
5. Park D S, Chan W, Zhang Y, et al. Specaugment: A simple data augmentation method for automatic speech recognition[J]. *arXiv preprint arXiv:1904.08779*, 2019.
6. Faris H, Ala'M A Z, Heidari A A, et al. An intelligent system for spam detection and identification of the most relevant features based on evolutionary random weight networks[J]. *Information Fusion*, 2019, 48:67-83.
7. Watanabe K, Zhou Y. Theory-driven analysis of large corpora: Semisupervised topic classification of the UN speeches[J]. *Social Science Computer Review*, 2020: 0894439320907027.
8. Gao Z, Feng A, Song X, et al. Target-dependent sentiment classification with BERT[J]. *IEEE Access*, 2019, 7: 154290-154299.



9. Ige, T., & Adewale, S. (2022a). Implementation of data mining on a secure cloud computing over a web API using supervised machine learning algorithm. *International Journal of Advanced Computer Science and Applications*, 13(5), 1–4. <https://doi.org/10.14569/IJACSA.2022.0130501>
10. Ige, T., & Adewale, S. (2022b). AI powered anti-cyber bullying system using machine learning algorithm of multinomial naïve Bayes and optimized linear support vector machine. *International Journal of Advanced Computer Science and Applications*, 13(5), 5–9. <https://doi.org/10.14569/IJACSA.2022.0130502>
11. Park D S, Chan W, Zhang Y, et al. SpecAugment: A simple data augmentation method for automatic speech recognition[J]. *arXiv preprint arXiv:1904.08779*, 2019.
12. Faris H, Ala'M A Z, Heidari A A, et al. An intelligent system for spam detection and identification of the most relevant features based on evolutionary random weight networks[J]. *Information Fusion*, 2019, 48: 67-83.
13. Watanabe K, Zhou Y. Theory-driven analysis of large corpora: Semisupervised topic classification of the UN speeches[J]. *Social Science Computer Review*, 2020: 0894439320907027.
14. Gao Z, Feng A, Song X, et al. Target-dependent sentiment classification with BERT[J]. *IEEE Access*, 2019, 7: 154290-154299.
15. Yong Z, Youwen L, Shixiong X. An improved KNN text classification algorithm based on clustering[J]. *Journal of computers*, 2009, 4(3): 230-237.
16. Sang-Bum Kim, Kyoung-Soo Han, Hae-Chang Rim and Sung Hyon Myaeng, "Some Effective Techniques for Naive Bayes Text Classification," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 18, no. 11, pp. 1457-1466, Nov. 2006. <https://doi.org/10.1109/TKDE.2006.180>.
17. Nigam K, Lafferty J, McCallum A. Using maximum entropy for text classification[C]//IJCAI-99 workshop on machine learning for information filtering. 1999, 1(1): 61-67.
18. Wang Z Q, Sun X, Zhang D X, et al. An optimal SVM-based text classification algorithm[C]//2006 International Conference on Machine Learning and Cybernetics. IEEE, 2006: 1378-1381.
19. Berger M J. Large scale multi-label text classification with semantic word vectors[J]. Technical report, Stanford University, 2015.
20. Wang S, Huang M, Deng Z. Densely connected CNN with multi-scale feature attention for text classification[C]//IJCAI. 2018: 4468-4474.
21. Guo B, Zhang C, Liu J, et al. Improving text classification with weighted word embeddings via a multi-channel TextCNN model[J]. *Neurocomputing*, 2019, 363: 366-374.
22. Tao P, Sun Z, Sun Z. An improved intrusion detection algorithm based on GA and SVM[J]. *Ieee Access*, 2018, 6: 13624-13631.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.