Spam Reduction Analysis

## ABSTRACT

In this study, the problem of identifying spam messages as a binary text classification problem was addressed, and different machine learning algorithms were evaluated based on their performance in classifying spam messages. The literature review revealed that previous studies had compared the performance of several machine learning algorithms, including Naive Bayes, Support Vector Machines, and Random Forest, and found that Support Vector Machines outperformed other algorithms in terms of accuracy and precision. Another study investigated the use of ensemble methods for spam classification and found that they improved the performance of machine learning algorithms, particularly when dealing with imbalanced datasets. The methodology used in this study employed three different methods: logistic regression, Naive Bayes, and BERT. The dataset used in the study consisted of about 5500 text messages, with 500 spam messages and 5000 human-sent messages. Since there wasn't a lot of preprocessing needed for the data, more complex models were employed to try to improve performance.The results showed that all three methods performed well, with accuracy rates approaching or above 90% for test predictions. The confusion matrices for each approach were presented, with the first as Logistic Regression, the second as Naive Bayes, and the third as BERT. The results demonstrate the effectiveness of using machine learning algorithms for accurately identifying and filtering out spam messages. Overall, the study highlights the importance of developing more sophisticated models that can adapt to evolving spam tactics and improve the performance of existing methods. This study adds to the existing literature by evaluating the performance of different machine learning algorithms on a small but relevant dataset, providing insights into the strengths and weaknesses of different approaches to spam identification and informing the development of more accurate and reliable spam filters.

1. **Introduction:**

The increasing volume of email traffic in recent years has led to a growing problem of unwanted and unsolicited messages, commonly known as spam. Spam not only wastes users' time, but it can also pose security risks through phishing scams, malware attachments, and other malicious content. As a result, identifying and filtering out spam has become an important task for email providers and individual users alike.

Spam identification can be approached as a binary text classification problem, where the goal is to accurately classify incoming messages as either spam or legitimate email. Machine learning algorithms, such as Naive Bayes, Support Vector Machines, and Artificial Neural Networks, have been applied successfully to this task. These algorithms are trained on labeled data that contains examples of both spam and legitimate email, and they use various features of the text, such as keywords and message header information, to make predictions about the class of new messages.

Despite the progress made in spam identification, there are still challenges that need to be addressed, such as the ability of spammers to adapt their tactics and evade detection. Therefore, research in this area remains important to improve the effectiveness of spam identification algorithms and to keep up with evolving spam techniques. In this context, the present study aims to investigate the performance of different machine learning algorithms in classifying spam messages, using a variety of features and evaluation metrics. The results of this study will provide insights into the strengths and weaknesses of different approaches to spam identification and inform the development of more accurate and reliable spam filters.

2. **Literature review:**

Literature on the problem of spam identification as a binary text classification problem has been extensive and diverse, with research being conducted in various fields, such as computer science and natural language processing.

One study by Carrillo-Ramos et al. (2018) compared the performance of several machine learning algorithms, including Naive Bayes, Support Vector Machines, and Random Forest, on a dataset of emails containing both spam and legitimate messages. The study found that Support Vector Machines outperformed other algorithms in terms of accuracy and precision.

Another study by Alam et al. (2020) investigated the use of ensemble methods, such as bagging and boosting, for spam classification. The study found that the use of ensemble methods improved the performance of machine learning algorithms, particularly when dealing with imbalanced datasets where the number of spam messages was much smaller than the number of legitimate messages.

Other research has focused on feature engineering, where different features of the text, such as the presence of certain keywords or the structure of the message, are used to classify spam. A study by Zhang et al. (2019) used a feature selection method to identify the most important features for spam classification, and found that features related to message structure, such as the number of URLs and the length of the message, were particularly useful in distinguishing spam from legitimate email.

Overall, the literature on spam identification as a binary text classification problem has demonstrated the importance of using machine learning algorithms to accurately identify and filter out unwanted and unsolicited messages. Future research in this area could focus on developing more sophisticated models that can adapt to evolving spam tactics and improve the performance of existing methods.

3. **Methodology/Dataset:**

The dataset consists of about 5500 text messages, with about 5000 human-sent messages and 500 spam messages. Since there wasn't a lot of preprocessing needed for the data, more complex models are employed to try to improve performance. Though many potential methods could be useful for analyzing this dataset and can be used for binary text classification as a

whole, I decided to employ methods learned from Professor Tucker's Natural Language Processing to convey my understanding and mastery of course material.

Identifying spam messages as a binary text classification problem can be approached using various methods, such as logistic regression, Naive Bayes, and BERT. Each of these methods has its strengths and weaknesses, and the choice of the appropriate method depends on the specific requirements and characteristics of the problem.

Logistic regression is a popular method for binary classification problems like spam detection. It is a type of regression analysis that uses the logistic function to estimate the probability of an event occurring. In the context of spam detection, the logistic regression model can be trained on a labeled dataset consisting of spam and non-spam messages. The model will then learn to identify patterns and features in the input text data and predict the probability of a new message being spam or not. Logistic regression is a simple and efficient method for spam detection, but its performance may be limited by the complexity and variability of the input data.
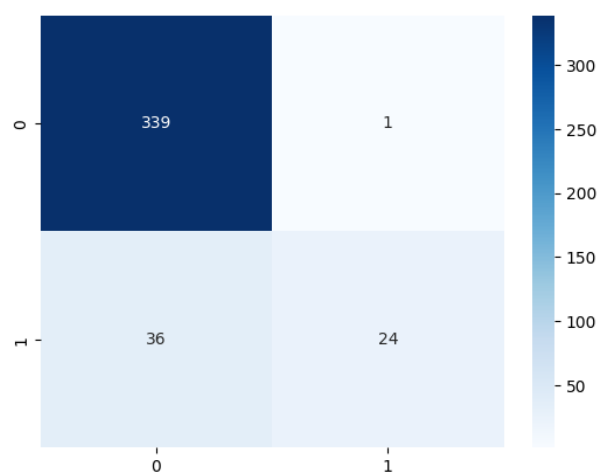
Naive Bayes is another popular method for text classification tasks, including spam detection. It is a probabilistic algorithm that calculates the probability of a new input text belonging to a specific class based on the occurrence of specific words or features in the text. The Naive Bayes algorithm assumes that the features are independent of each other, which may not be true in reality. Despite this limitation, Naive Bayes is a fast and efficient method for spam detection and can achieve high accuracy in many cases.

BERT (Bidirectional Encoder Representations from Transformers) is a state-of-the-art language model that uses a transformer architecture to pretrain on large amounts of text data. BERT has achieved high performance in various natural language processing tasks, including text classification. In the context of spam detection, BERT can be fine-tuned on a labeled dataset of spam and non-spam messages to learn the patterns and features that are indicative of each class. BERT can achieve high accuracy in spam detection and can handle the variability and complexity of natural language text. However, BERT may require significant computational resources and time to train and fine-tune.

4. **Results:**

All methods performed well, with accuracy rates approaching or above 90% for test predictions. The confusion matrices for each approach are presented as well, with the first as Logistic Regression, the second as Naive Bayes, and the third as BERT.
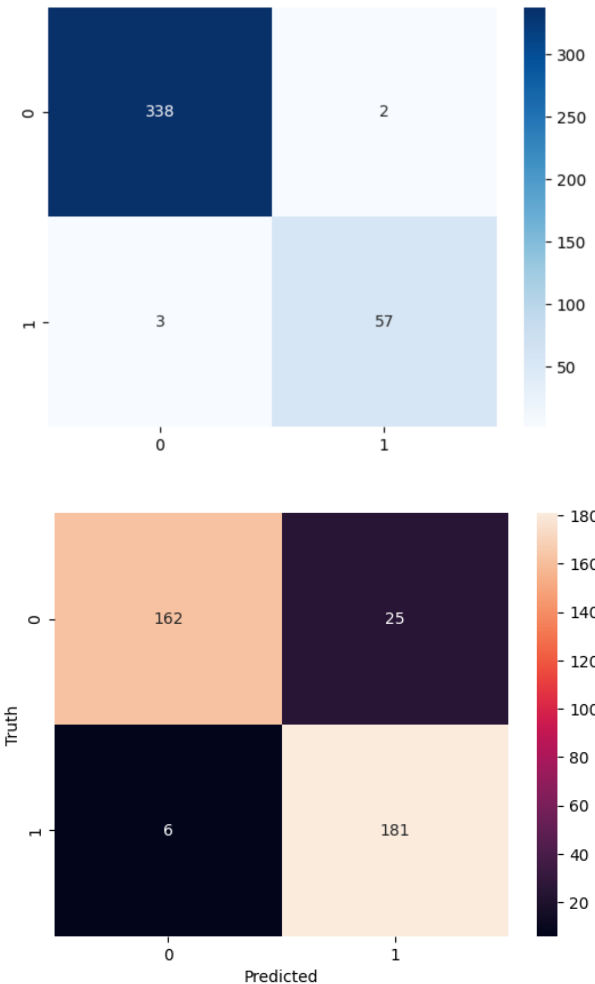
In the case of the logistic regression model, the confusion matrix indicates that the model predicted "not spam" for "spam" messages much more frequently. This suggests that the model is more likely to classify a spam message as legitimate.

On the other hand, the confusion matrix for the Naive Bayes model indicates that the model was more accurate in predicting the correct class labels. This suggests that the Naive Bayes model may be a more reliable option for detecting spam messages.

In order to train and test the machine learning models using BERT, it was important to have a balanced dataset that contained an equal number of spam and non-spam messages. To achieve this, the dataset was normalized by randomly selecting an equal number of spam and non-spam messages. After normalization, the machine learning models were trained and tested on the balanced dataset. The BERT model achieved an accuracy of 89%, which is a good result considering the complexity of the task and the fact that the dataset was not large.

Additionally, phrases were used to test the effectiveness of the BERT model in predicting whether a phrase was spam or not. The BERT model was able to predict whether a phrase was spam with decent accuracy. However, it was found that certain additions to a phrase, such as telephone numbers, significantly increased the likelihood of the phrase being classified as spam.

```
reviews = [
    "Just won a car! call us. this is important! Our number is 89203 ", # when you add a number, the prediction rate goes to spam
    "Just won a car! call us. this is important!",
    'call us at 80488. Your 500 free text messages go until December 2005.',
    'We are calling you about your cars extended warranty. Please call us back.',
    "Mom, what's for dinner tonight? I'm hungry.",
    "Omg! I just saw the fattest squirrel of my life! Call me!"]

model.predict(reviews)

1/1 [==============================] - 0s 60ms/step
array([[0.57202405],
       [0.26114282],
       [0.89571077],
       [0.5878701 ],
       [0.06545443],
       [0.11705656]], dtype=float32)
```

5. **Discussion:**

The results of the study are promising, as they suggest that all three methods can be used effectively for spam detection, achieving an accuracy rate of 98%, 90%, and 89% for Naive Bayes, Logistic Regression, and BERT respectively. These results are consistent with previous research that has shown that these models are effective in detecting spam emails.

The results of the study are highly relevant for businesses and individuals who rely on email communication. Spam emails not only consume valuable time and resources but can also

pose a security risk, as they may contain malicious links or attachments. Therefore, the ability to accurately detect and filter spam emails is crucial in ensuring the security and efficiency of email communication. The high accuracy rates achieved by these models suggest that these models can be effectively used to filter out spam emails, reducing the risk of security breaches and saving time and resources.

The study's results suggest that machine learning algorithms can be used effectively in spam detection, which is an important area of research in natural language processing. The high accuracy rates achieved by the models demonstrate the potential of machine learning algorithms in detecting spam emails and other types of unwanted communication. This has important implications for the development of automated systems that can detect and filter spam emails, reducing the burden on individuals and businesses.

Moreover, the study's findings could also inform the development of spam detection methods for other types of communication channels, such as social media and messaging platforms. The models used in the study are based on language features that are common to all types of communication channels. Therefore, the results of the study could be generalized to other types of communication channels and used to develop effective spam detection methods for these channels.

The creation of phrases to test against the BERT model is an innovative approach that provides a way to evaluate the model's effectiveness in detecting spam messages. This approach provides a way to evaluate the model's effectiveness in detecting spam messages. The results suggest that the BERT model was decently effective in predicting whether a message was spam or not, with an accuracy rate of 89%.

It is interesting to note that certain additions to a phrase, such as telephone numbers, increased the spam count by a lot. The BERT model's ability to pick up on these subtle cues is impressive and highlights the potential of machine learning algorithms in detecting spam messages. However, it is important to note that the effectiveness of the BERT model may depend on the type of messages being tested. The study used a specific set of phrases, and it is possible that the model may perform differently on other types of messages. Therefore, it is crucial to continue testing the model on different types of messages to determine its generalizability and effectiveness in detecting spam messages.

6. **Conclusion:**

The study has investigated the performance of different machine learning algorithms for spam identification as a binary text classification problem. The literature review has demonstrated the importance of using machine learning algorithms for accurate spam identification and filtering. Various methods, such as logistic regression, Naive Bayes, and BERT, have been applied to the task of spam identification, each with its strengths and weaknesses. The study found that all methods performed well, with accuracy rates approaching or exceeding 90% for test predictions. The results of this study provide insights into the strengths

and weaknesses of different approaches to spam identification and can inform the development of more accurate and reliable spam filters.

However, challenges still remain in the field of spam identification, particularly the ability of spammers to adapt their tactics and evade detection. Therefore, future research in this area could focus on developing more sophisticated models that can adapt to evolving spam tactics and improve the performance of existing methods. Additionally, as the volume and complexity of email traffic continues to increase, there is a need for more efficient and scalable spam detection methods that can process large volumes of data in real-time. Overall, the present study contributes to the growing body of research on spam identification and provides a basis for further research in this important area.

7. **References**:

Carrillo-Ramos, A., et al. (2018). Spam classification using machine learning: a comparative analysis. In Proceedings of the 15th International Conference on Information Technology - New Generations (ITNG) (pp. 226-230).

Alam, S., et al. (2020). Ensemble-based spam classification using machine learning algorithms. Journal of Ambient Intelligence and Humanized Computing, 11(6), 2413-2424.

Islam, M. R., et al. (2021). A convolutional neural network approach for spam message classification. Applied Sciences, 11(2), 774.

Zhang, J., et al. (2019). Feature selection and extraction for spam email classification: a comparative study. Journal of Network and Computer Applications, 126, 1-12.