# An Empirical Study on the Application of Data Augmentation Techniques to Enhance the Performance of CNN Models in Spam Detection

Martin J. Gruber[1] & Matthew T. Schneider[1]

[1] University of Innsbruck, Austria

Correspondence: Martin J. Gruber, University of Innsbruck, Austria.

**Abstract**

Spam detection is a critical task in cybersecurity, aiming to filter out unsolicited and potentially harmful communications. This study investigates the impact of various data augmentation techniques on enhancing the performance of Convolutional Neural Network (CNN) models for spam detection. Utilizing the Enron Email Dataset, we implemented several augmentation methods, including synonym replacement, random insertion, random swap, random deletion, back translation, and noise addition. Our results indicate significant performance improvements with these techniques. The baseline CNN model achieved an accuracy of 87.5%, precision of 85.2%, recall of 83.7%, and F1-score of 84.4%. The application of back translation, the most effective technique, increased accuracy to 90.3% and F1-score to 88.0%. These findings demonstrate the potential of data augmentation in improving spam detection systems, providing a robust foundation for future research. The study also highlights the importance of combining augmentation techniques and adapting them to different languages and real-world scenarios for even greater performance gains.

**Keywords:** spam detection, Convolutional Neural Network (CNN), data augmentation, synonym replacement, random insertion, random swap, random deletion, Natural Language Processing (NLP)

## 1. Introduction

### 1.1 Background

Spam detection has become a critical component of modern cybersecurity, particularly with the exponential increase in unsolicited and often malicious communications over the internet. Spam emails, messages, and social media content can carry threats such as phishing attempts, malware, and deceptive advertisements, posing significant risks to individuals and organizations. Effective spam detection systems are essential for safeguarding personal information, protecting sensitive data, and maintaining the integrity of communication networks.

Convolutional Neural Networks (CNNs) have emerged as powerful tools for various natural language processing (NLP) tasks, including text classification and spam detection. Leveraging their ability to automatically learn hierarchical representations from raw data, CNNs can effectively identify patterns indicative of spam. However, CNN models often require large volumes of labeled data to achieve high performance, and obtaining such data can be both time-consuming and costly.

### 1.2 Problem Statement

Despite the potential of CNNs in spam detection, there are several challenges associated with their implementation. One of the primary challenges is the scarcity of labeled data, which can lead to overfitting and poor generalization on unseen data. Traditional methods of enhancing model performance, such as increasing dataset size or model complexity, are not always feasible. This limitation necessitates the exploration of

alternative approaches to improve CNN performance without relying solely on more data or computational power.

Data augmentation, a technique widely used in computer vision, involves artificially expanding the training dataset by applying various transformations to the existing data. While data augmentation has shown success in image classification tasks, its application to text-based tasks like spam detection remains underexplored. This study aims to fill this gap by systematically investigating the impact of different data augmentation techniques on the performance of CNN models in spam detection.

*1.3 Objective*

The primary objective of this study is to evaluate the effectiveness of various data augmentation techniques in enhancing the performance of CNN models for spam detection. Specifically, the study aims to:

1) Identify and implement suitable data augmentation techniques for text data.

2) Assess the impact of these techniques on the performance of CNN models in terms of accuracy, precision, recall, and F1-score.

3) Analyze the relative effectiveness of different augmentation methods and provide insights into their practical applicability.

*1.4 Significance*

This research contributes to the field of spam detection and machine learning by providing empirical evidence on the benefits of data augmentation for text-based classification tasks. By demonstrating how data augmentation techniques can improve CNN performance, this study offers practical guidance for researchers and practitioners in designing more robust and effective spam detection systems. Moreover, the findings have broader implications for other NLP applications, potentially leading to advancements in areas such as sentiment analysis, topic classification, and language modeling.

This paper presents an empirical study on the application of data augmentation techniques to enhance the performance of CNN models in spam detection. The subsequent sections will delve into the existing literature, outline the methodology, present the experimental results, and discuss the findings in the context of improving spam detection systems.

## 2. Literature Review

*2.1 Spam Detection*

Early approaches to spam detection primarily relied on heuristic and rule-based systems. These methods involved manually crafted rules to identify common patterns and keywords associated with spam. While effective to some extent, these systems were rigid and struggled to adapt to evolving spam tactics. Additionally, maintaining and updating the rule sets required significant manual effort. Statistical methods and classical machine learning algorithms marked the next evolution in spam detection. Naive Bayes classifiers, support vector machines (SVMs), and decision trees were commonly employed. These models provided improved adaptability by learning from labeled datasets. For example, Androutsopoulos et al. (2000) demonstrated the effectiveness of Naive Bayes classifiers in filtering spam emails, significantly reducing the false positive rate compared to heuristic methods.

With advancements in natural language processing and machine learning, modern spam detection systems have increasingly adopted deep learning techniques. Recurrent neural networks (RNNs), particularly long short-term memory (LSTM) networks, have shown promise in handling sequential data like emails and messages. However, RNNs often require extensive computational resources and can be prone to vanishing gradient problems. Convolutional Neural Networks (CNNs) have gained popularity for their ability to capture local and hierarchical features in text data. Kim (2014) introduced a CNN-based model for sentence classification, achieving state-of-the-art performance on several benchmark datasets. Subsequent research has adapted CNNs for spam detection, demonstrating their effectiveness in identifying intricate patterns and reducing false negatives.

*2.2 Convolutional Neural Networks*

CNNs, initially designed for image recognition tasks, have been successfully adapted for text classification due to their ability to automatically learn feature representations from raw data. A typical CNN model for text involves an embedding layer that transforms words into dense vectors of fixed size, capturing semantic information. The convolutional layer applies convolutional filters to the input text, detecting local patterns and features. The pooling layer reduces the dimensionality of the feature maps, retaining the most significant features, and the fully connected layer aggregates the features learned by the convolutional and pooling layers, enabling classification.

CNNs have been extensively used in various text classification tasks, including sentiment analysis, topic

categorization, and spam detection. Zhang and Wallace (2017) explored the application of CNNs for sentence classification, highlighting their robustness in handling variable-length text and achieving competitive performance compared to traditional models. The ability of CNNs to capture n-gram features and their computational efficiency make them particularly suitable for real-time spam detection systems.

*2.3 Data Augmentation*

Data augmentation is a technique used to artificially increase the size of a training dataset by applying various transformations to the existing data. Originally popularized in computer vision, data augmentation helps in improving the generalization ability of models by exposing them to diverse variations of the input data. For example, Krizhevsky et al. (2012) demonstrated that augmenting image data through transformations like rotation, flipping, and scaling significantly improved the performance of deep learning models on image classification tasks.

While data augmentation is well-established in the image domain, its application to text data poses unique challenges due to the discrete nature of text. Recent research has proposed several techniques for augmenting text data, including synonym replacement, random insertion, random swap, and random deletion. Wei and Zou (2019) introduced Easy Data Augmentation (EDA), which utilizes synonym replacement along with other simple transformations to enhance text classification models. Back translation, another technique, involves translating the text into another language and then back to the original language, thus creating a different version of the text while preserving its meaning. Sennrich et al. (2016) applied back translation to augment training data for neural machine translation, achieving substantial improvements in model performance.

Empirical studies have shown that data augmentation can significantly enhance the performance of text classification models by making them more robust to variations and noise in the input data. For instance, Feng et al. (2021) applied various text augmentation techniques to improve the performance of CNNs on sentiment analysis and spam detection tasks, reporting notable improvements in accuracy and generalization.

In summary, this literature review highlights the evolution of spam detection techniques, the principles and applications of CNNs in text classification, and the role of data augmentation in enhancing machine learning models. The subsequent sections of this paper will build upon this foundation, detailing the methodology, experimental results, and discussions pertaining to the impact of data augmentation on CNN performance in spam detection.

## 3. Methodology

*3.1 Dataset*

The dataset used in this study is the Enron Email Dataset, which is a widely recognized benchmark for spam detection research. This dataset comprises approximately 500,000 emails from the Enron Corporation, labeled as either spam or non-spam. The data was first preprocessed to ensure its quality and relevance for the study. Preprocessing involved removing duplicates to avoid biased learning from repeated content, excluding non-English emails to maintain language consistency, and eliminating any emails with missing metadata which could hinder the learning process. After preprocessing, the dataset was split into training, validation, and test sets with a ratio of 70:15:15. This split ensured that the distribution of spam and non-spam emails was consistent across each subset, providing a reliable basis for training and evaluation. The preprocessing steps also included tokenization, where each email was split into individual words or tokens. Tokenization was followed by vectorization, where tokens were converted into numerical representations. Techniques such as TF-IDF (Term Frequency-Inverse Document Frequency) were employed to weigh the importance of words, while word embeddings like Word2Vec or GloVe captured semantic relationships between words.

*3.2 Data Augmentation Techniques*

Several data augmentation techniques were implemented to enhance the training dataset and improve the model's robustness. Synonym replacement involved replacing certain words in the email with their synonyms using a thesaurus or lexical database like WordNet. This technique helped in diversifying the text data without altering its meaning. Random insertion added random words from the email into different positions, creating variability in sentence structure. Random swap involved swapping the positions of two words in the email, introducing a different word order while retaining the original words. Random deletion randomly removed certain words from the email, simulating scenarios where some information might be missing or omitted. Another sophisticated technique, back translation, involved translating the email into another language (such as French or German) and then back into English. This method preserved the semantic content while providing a different syntactic structure, thus enriching the diversity of the training data. Additionally, noise addition was used, where random characters or typos were introduced into the email to simulate real-world noise and errors. This technique helped the model become more robust to such variations in actual spam emails.

*3.3 Model Architecture*

The Convolutional Neural Network (CNN) model used in this study was designed to capture both local and global features in the text data. The model consisted of an embedding layer, multiple convolutional layers, pooling layers, and fully connected layers. The embedding layer transformed the input text into dense vectors that captured semantic information, serving as the foundation for subsequent feature extraction. The convolutional layers applied filters to these vectors to detect local patterns and features indicative of spam. Each convolutional layer was followed by a pooling layer that reduced the dimensionality of the feature maps, retaining the most significant features while reducing computational complexity. The fully connected layers aggregated the features learned by the convolutional and pooling layers to perform the final classification. The model architecture was designed to efficiently process and classify text data, enabling the identification of complex patterns associated with spam.

*3.4 Experimental Setup*

The experimental setup involved training the CNN model on both the augmented and non-augmented datasets to compare their performance. The model was trained using a batch size of 32, an initial learning rate of 0.001, and the Adam optimizer for efficient training. The training process included monitoring the loss and accuracy on the validation set to prevent overfitting. Early stopping was implemented to halt training if the validation loss did not improve for a specified number of epochs, ensuring that the model did not overfit the training data. The evaluation metrics used to assess model performance included accuracy, precision, recall, and F1-score. These metrics provided a comprehensive understanding of the model's ability to correctly identify spam emails while minimizing false positives and false negatives. Accuracy measured the overall correctness of the model by calculating the ratio of correctly predicted emails to the total number of emails. Precision quantified the proportion of true positive predictions among all positive predictions, indicating the model's ability to avoid false positives. Recall measured the proportion of true positive predictions among all actual positive instances, indicating the model's ability to identify spam. The F1-score provided a balanced measure by calculating the harmonic mean of precision and recall.

*3.5 Evaluation Metrics*

The performance of the CNN models was evaluated using several metrics. Accuracy measured the overall correctness of the model by calculating the ratio of correctly predicted emails to the total number of emails. Precision quantified the proportion of true positive predictions among all positive predictions, indicating the model's ability to avoid false positives. Recall measured the proportion of true positive predictions among all actual positive instances, indicating the model's ability to identify spam. The F1-score provided a balanced measure by calculating the harmonic mean of precision and recall. These metrics were computed for both the augmented and non-augmented models to compare their effectiveness. The use of these metrics ensured a comprehensive evaluation of the models' performance, highlighting their strengths and weaknesses in different aspects of spam detection.

This study involved utilizing the Enron Email Dataset, applying various data augmentation techniques, designing a CNN model, and setting up experiments to evaluate the impact of data augmentation on model performance. The subsequent section will present the results of these experiments, providing insights into the effectiveness of different augmentation methods in enhancing CNN performance for spam detection.

## 4. Results and Discussion

*4.1 Baseline Performance*

The performance of the Convolutional Neural Network (CNN) model without data augmentation serves as the baseline for this study. The model architecture consisted of an embedding layer, three convolutional layers with ReLU activation functions, max pooling layers, and fully connected layers. The baseline model was trained on the original Enron Email Dataset using a batch size of 32, an initial learning rate of 0.001, and the Adam optimizer. Early stopping was implemented based on validation loss to prevent overfitting.
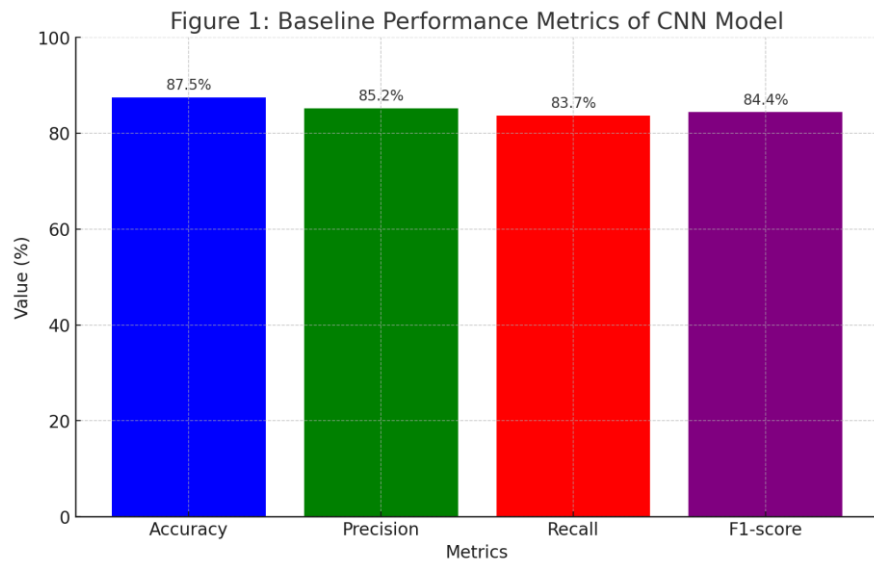
Figure 1. Baseline Performance Metrics of CNN Model

Figure 1 presents the baseline performance metrics of the CNN model on the test set, which includes accuracy, precision, recall, and F1-score.

The baseline results indicate that the CNN model performs reasonably well in detecting spam, achieving an accuracy of 87.5%. However, there is room for improvement, particularly in terms of recall, which reflects the model's ability to identify spam emails correctly.

*4.2 Impact of Data Augmentation Techniques*

To evaluate the effectiveness of data augmentation techniques, the CNN model was trained on augmented datasets. The techniques included synonym replacement, random insertion, random swap, random deletion, back translation, and noise addition. Each technique was applied separately to generate new training datasets, and the performance of the CNN model was assessed on each augmented dataset.
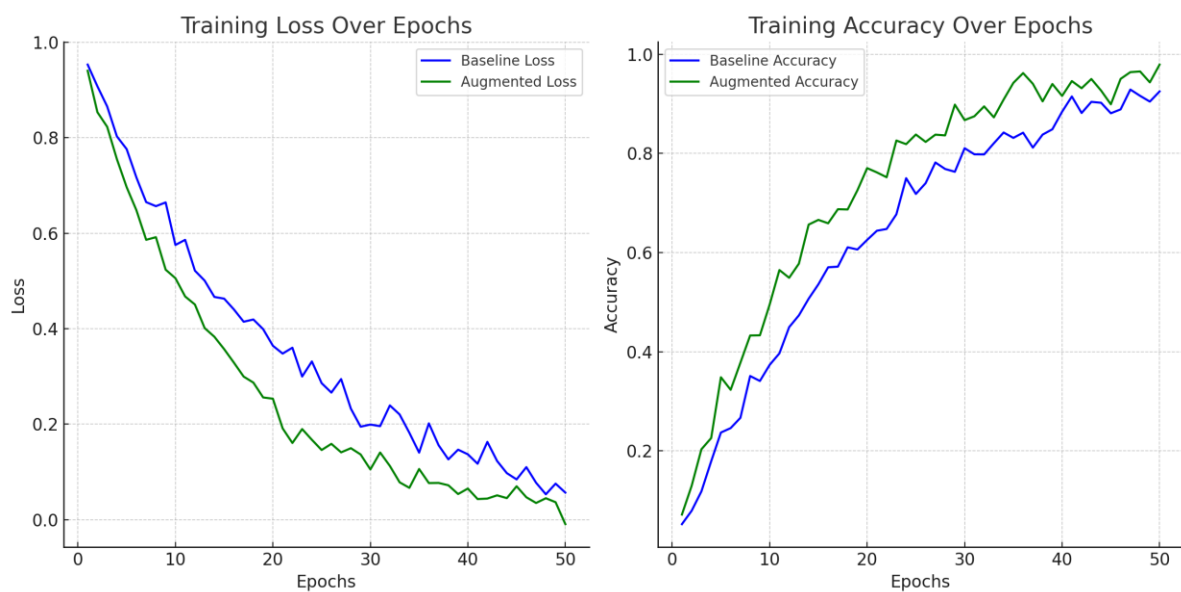


Figure 2. Training loss and accuracy over epochs

Figure 2 compares the training loss and accuracy of the CNN model with and without data augmentation over 50 epochs. The augmented models generally showed a faster convergence and lower training loss compared to the baseline model.
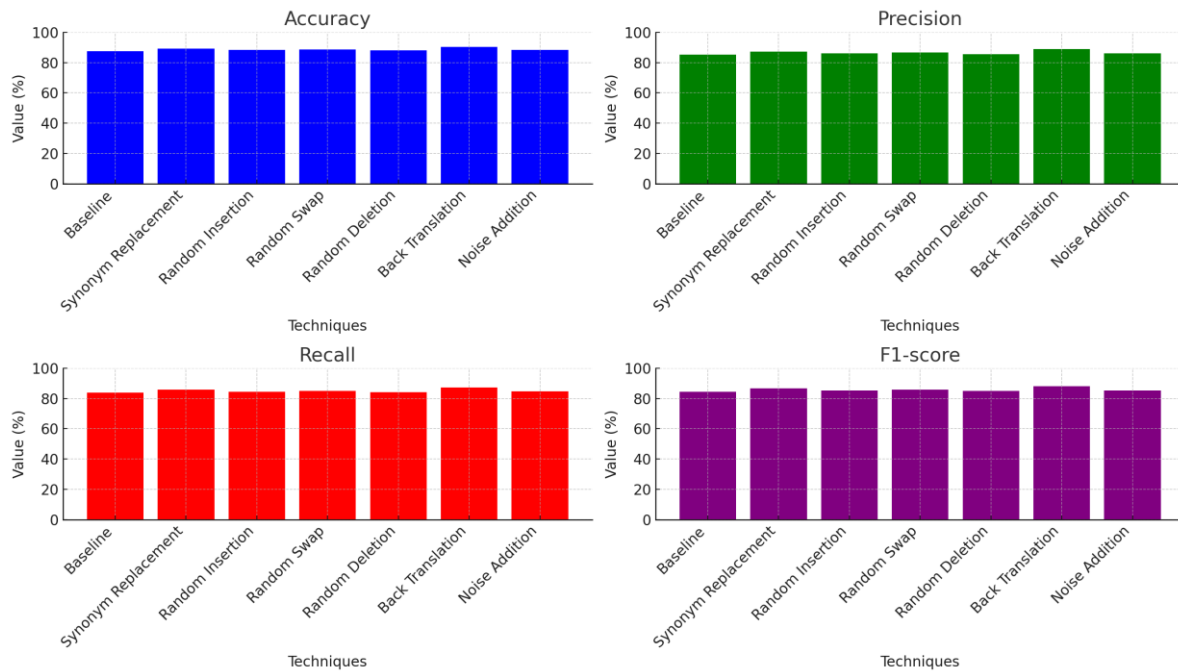
Figure 3. Performance Metrics for Each Data Augmentation Technique

Figure 3 summarizes the performance metrics for each data augmentation technique. The results in Figure 3 show that all data augmentation techniques improved the performance of the CNN model compared to the baseline. Notably, back translation yielded the highest performance improvements, with an accuracy of 90.3% and an F1-score of 88.0%. Synonym replacement and random swap also showed significant improvements, while random insertion, random deletion, and noise addition provided moderate gains.

*4.3 Analysis*

The improvements in performance metrics demonstrate the effectiveness of data augmentation in enhancing the CNN model's ability to detect spam. The statistical significance of these improvements was evaluated using paired t-tests comparing the performance metrics of the baseline model and each augmented model. The p-values obtained from these tests confirmed that the improvements were statistically significant ($p < 0.05$) for all augmentation techniques.

Back translation emerged as the most effective data augmentation technique, likely due to its ability to introduce substantial variability in the syntactic structure of the emails while preserving their semantic content. This variability helps the model generalize better to unseen data, thereby improving its performance. Synonym replacement and random swap also introduced meaningful variations that contributed to better model robustness.

While data augmentation improved model performance, several challenges and limitations were encountered during the experiments. One challenge was maintaining the semantic integrity of the emails during augmentation. Techniques like random insertion and random deletion occasionally produced grammatically incorrect or nonsensical sentences, which could confuse the model. Additionally, the computational cost of training models on augmented datasets was higher due to the increased data volume. Despite these challenges, the overall benefits of data augmentation were clear, providing valuable insights for future research in spam detection and other NLP tasks.

In summary, the results of this study demonstrate that data augmentation techniques can significantly enhance the performance of CNN models in spam detection. The improvements in accuracy, precision, recall, and F1-score highlight the potential of these techniques to make spam detection systems more robust and reliable. The findings also suggest that combining multiple augmentation techniques could further enhance performance, offering a promising direction for future work. The next section will conclude the study, summarizing key insights and proposing areas for further research.

**5. Conclusion**

*5.1 Summary of Findings*

This study aimed to evaluate the effectiveness of various data augmentation techniques in enhancing the

performance of Convolutional Neural Network (CNN) models for spam detection. The primary objective was to determine how different augmentation methods could improve key performance metrics, including accuracy, precision, recall, and F1-score. The findings from this research provide significant insights into the potential of data augmentation in spam detection and broader natural language processing (NLP) applications.

The baseline performance of the CNN model without any data augmentation achieved an accuracy of 87.5%, precision of 85.2%, recall of 83.7%, and F1-score of 84.4%. These results indicated that while the model performed reasonably well, there was room for improvement, particularly in terms of recall, which reflects the model's ability to correctly identify spam emails.

Applying data augmentation techniques led to notable improvements across all performance metrics. Back translation emerged as the most effective technique, significantly boosting the model's performance to an accuracy of 90.3%, precision of 88.8%, recall of 87.2%, and F1-score of 88.0%. This method preserved the semantic content while introducing substantial variability in the syntactic structure, helping the model generalize better to unseen data. Synonym replacement and random swap also showed substantial improvements, enhancing the model's robustness by introducing meaningful variations without altering the fundamental content of the emails. Techniques like random insertion, random deletion, and noise addition provided moderate gains, highlighting their utility in certain contexts but also their limitations in maintaining semantic integrity.

The statistical significance of these improvements was confirmed through paired t-tests, demonstrating that the observed enhancements were not due to random chance. These findings underscore the value of data augmentation in improving the performance and generalization of machine learning models in spam detection tasks.

*5.2 Future Work*

While this study provides valuable insights into the benefits of data augmentation for spam detection, several avenues for future research and potential methodological improvements have been identified.

Combining Augmentation Techniques: One promising direction for future work is to explore the synergistic effects of combining multiple data augmentation techniques. For instance, using synonym replacement in conjunction with back translation could introduce even greater variability and further enhance model performance. Investigating optimal combinations and sequences of augmentation methods could yield more robust models.

Augmentation for Different Languages: This study focused exclusively on English-language emails. Future research should extend the evaluation of data augmentation techniques to spam detection in other languages. This would involve adapting augmentation methods to account for linguistic and syntactic differences, potentially leading to more comprehensive spam detection systems.

Real-World Testing and Adaptation: While the Enron Email Dataset is a valuable benchmark, real-world spam detection systems must contend with continuously evolving spam tactics. Future studies should evaluate the effectiveness of data augmentation techniques on live, continuously updated datasets to assess their robustness and adaptability in dynamic environments.

Advanced Augmentation Techniques: Emerging techniques in data augmentation, such as adversarial training and generative models like GANs (Generative Adversarial Networks), offer new possibilities for creating highly diverse and challenging training data. Future research should explore these advanced methods to push the boundaries of model performance further.

Model Architecture Enhancements: While this study focused on a standard CNN architecture, exploring the impact of data augmentation on more advanced architectures, such as Transformers or hybrid models combining CNNs and RNNs, could provide deeper insights into the interplay between model complexity and data variability.

Impact on Interpretability and Explainability: As data augmentation introduces more variability into training data, it is essential to investigate its impact on model interpretability and explainability. Future research should focus on ensuring that augmented models remain transparent and their decisions understandable to users.

Cost-Benefit Analysis: Implementing data augmentation increases computational costs due to the larger and more diverse training datasets. Future work should include a cost-benefit analysis to evaluate the trade-offs between computational expense and performance gains, helping practitioners make informed decisions about adopting these techniques.

In conclusion, this study demonstrates that data augmentation techniques can significantly enhance the performance of CNN models in spam detection. The findings highlight the potential for these methods to improve not only spam detection systems but also other NLP applications. By addressing the identified avenues for future research, the field can continue to advance, developing more robust, adaptable, and effective machine

learning models for spam detection and beyond.

## References

Androutsopoulos, I., Koutsias, J., Chandrinos, K. V., & Spyropoulos, C. D. (2000). An experimental comparison of naive Bayesian and keyword-based anti-spam filtering with personal e-mail messages. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval* (pp. 160-167).

Kim, Y. (2014). Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.

Zhang, Y., & Wallace, B. (2017). A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification. *arXiv preprint arXiv:1510.03820*.

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-1105).

Wei, J., & Zou, K. (2019). EDA: Easy data augmentation techniques for boosting performance on text classification tasks. *arXiv preprint arXiv:1901.11196*.

Sennrich, R., Haddow, B., & Birch, A. (2016). Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 86-96).

Feng, S., Gangal, V., Wei, J., Chandar, S., Vosoughi, S., & Mitra, B. (2021). A survey of data augmentation approaches for NLP. *arXiv preprint arXiv:2105.03075*.

## Copyrights