

Report on Spam Message Detection

1. Stakeholder:

Our primary stakeholder is an email service provider aiming to enhance its spam detection system. The company's objective is to improve user satisfaction by significantly reducing the number of spam messages that reach users' inboxes while ensuring that legitimate emails are not incorrectly classified as spam.

2. Problem Statement:

The company seeks to develop a more accurate spam detection system to classify incoming messages as either spam or ham (non-spam). The challenge is to achieve high detection accuracy while minimizing false positives, ensuring that legitimate emails are not mistakenly flagged as spam.

3. Dataset:

The dataset used for this project is the SMS Spam Collection Dataset. This dataset contains a total of 5,572 SMS messages categorized into two classes: 'ham' (legitimate messages) and 'spam' (unwanted messages). The dataset distribution includes 4,827 'ham' messages and 747 'spam' messages, highlighting an imbalanced nature that poses a challenge for classification. Here is the link of dataset:

<https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset/data>

4. Model Selection:

In this project, we experimented with several models to identify the best fit for our spam detection system. The selection of these models was based on their ability to handle text data efficiently and their potential to achieve high accuracy and precision:

Multinomial Naive Bayes: This model is known for its simplicity and efficiency in text classification tasks. It assumes that the features (words) are conditionally independent given the class, which often holds true for text data. Naive Bayes is particularly effective for high-dimensional data and provides fast training and prediction times.

XGBoost: XGBoost is a powerful gradient boosting algorithm that combines the predictions of multiple weak models to provide high performance and accuracy. It is known for its ability to handle large datasets and complex data structures. XGBoost also includes several hyperparameters that can be tuned to optimize performance, making it a versatile choice for classification tasks.

LSTM with GloVe Embeddings: Long Short-Term Memory (LSTM) networks are a type of recurrent neural network (RNN) effective for sequential data and capturing context in text. We used pre-trained GloVe embeddings to capture semantic relationships between words, enhancing the model's ability to

understand the text's meaning. LSTMs are particularly suited for tasks where the sequence of words is important, such as spam detection.

Each model was selected based on its strengths in handling text data and its potential to achieve high accuracy and precision.

5. Feature Engineering:

Feature engineering plays a critical role in preparing the text data for model training. The following steps were undertaken to transform the raw text into meaningful features:

Text Cleaning: We performed several text cleaning steps to preprocess the data. This included converting text to lowercase, removing punctuation, HTML tags, URLs, and numbers. Text cleaning helps focus on meaningful words and reduces noise in the data, making it easier for models to learn from the text.

Stopwords Removal: Stopwords are common words (such as "and", "the", "is") that do not contribute significant meaning to the text. We removed these stopwords to reduce the dimensionality of the data and focus on more informative words. This step helps in capturing the essential information while ignoring common but uninformative words.

Stemming: Stemming involves reducing words to their root forms, consolidating similar words (e.g., "running" and "runs" to "run"). This helps in reducing the feature space and capturing the core meaning of the words. Stemming is useful in grouping words with similar meanings, thereby enhancing the model's ability to generalize.

Vectorization:

We used two techniques for vectorizing the text data:

CountVectorizer: Converts the text data into a bag-of-words representation, where each word is represented by its frequency in the document. This method is straightforward and provides a simple representation of the text.

TF-IDF Transformer: Adjusts the frequency of words based on their importance across documents. TF-IDF (Term Frequency-Inverse Document Frequency) helps in emphasizing important words that are frequent in a document but rare across the corpus, making it a more informative representation.

These features were chosen based on their ability to capture the essential information in text while reducing noise and dimensionality.

6. Model Evaluation:

To evaluate the performance of our models, we used several metrics that provide a comprehensive view of the model's effectiveness:

Accuracy: Measures the proportion of correctly classified messages out of the total messages. It provides a general sense of the model's performance but can be misleading in imbalanced datasets.

Precision: Indicates the proportion of true positive predictions out of all positive predictions. Precision is crucial for minimizing false positives, ensuring that legitimate emails are not incorrectly classified as spam. High precision is essential in applications where the cost of false positives is high.

Recall: Measures the proportion of true positive predictions out of all actual positives. Recall is important for capturing as many spam messages as possible, ensuring that spam is effectively filtered. High recall is essential in applications where missing a positive instance (false negative) is costly.

F1-Score: The F1-score is the harmonic mean of precision and recall, providing a balanced measure of the model's performance. It is especially useful when the dataset is imbalanced, as it considers both false positives and false negatives. The F1-score provides a single metric that balances the trade-off between precision and recall.

These metrics were selected to ensure a comprehensive evaluation of the models, focusing on both precision (to avoid false positives) and recall (to catch most spam messages).

7. Results and Discussion:

The models were trained and evaluated using the aforementioned metrics. Below are the detailed results:

Multinomial Naive Bayes: This model achieved an accuracy of approximately 97.8%. The precision and recall were also high, making it a strong baseline model. The confusion matrix indicated that it was able to correctly classify most of the spam and ham messages, with minimal misclassifications.

XGBoost: The XGBoost model showed improved performance with an accuracy of around 96.4%. The model provided a good balance between precision and recall, indicating its robustness in handling complex data structures. The confusion matrix showed that XGBoost was able to correctly classify most messages, with fewer false positives and false negatives compared to Naive Bayes.

LSTM with GloVe Embeddings: The LSTM model, although computationally intensive, demonstrated promising results. It achieved an accuracy comparable to the other models but required more resources and fine-tuning. The model's ability to capture the sequential nature of text data provided a deeper understanding of the context, leading to accurate predictions.

8. Future Work and Recommendations:

Given more time, the following improvements could be made:

Hyperparameter Tuning: Conduct a more extensive search for optimal hyperparameters to further improve model performance. This includes experimenting with different parameter values and

combinations to find the best settings for each model. Grid search and random search are commonly used techniques for hyperparameter tuning.

Additional Features: Explore additional text features such as n-grams (combinations of consecutive words), part-of-speech tags, and sentiment analysis. These features can provide more context and improve the model's understanding of the text. N-grams, for instance, can capture phrases and word pairs that carry more information than individual words.

Ensemble Methods: Combine multiple models (ensemble methods) to leverage their strengths and improve overall accuracy. Ensemble methods, such as stacking or voting, can enhance performance by integrating the predictions of several models. These methods often result in better generalization and robustness.

We recommend the company use the XGBoost model due to its high accuracy and balanced performance in terms of precision and recall. While the LSTM model with GloVe embeddings shows promise, it requires more computational resources and fine-tuning. XGBoost offers a good trade-off between performance and resource requirements.

9. Deployment:

To deploy the model, the following steps will be taken:

Model Export: Save the trained model using a suitable format (e.g., Pickle for traditional models or SavedModel for TensorFlow). This allows for easy loading and use in different environments. The model export step ensures that the model can be reused and integrated into the deployment pipeline without retraining.

API Integration: Develop an API using Flask or FastAPI to serve the model predictions. This API will handle incoming messages, process them through the model, and return the classification result (spam or ham). API integration allows for real-time predictions and easy access to the model's functionality.

Continuous Monitoring: Implement logging and monitoring to track the model's performance in real-time. Regular updates and retraining with new data will ensure the model remains effective and adapts to new spam tactics. Continuous monitoring helps in identifying any performance degradation and allows for timely interventions.

These steps ensure the model can be seamlessly integrated into the company's existing infrastructure and continuously improved over time.

10. Conclusion:

This report outlines the rationale behind our choices and the methodologies applied throughout the project. We focused on achieving a robust spam detection system with a balanced performance suitable

for real-world applications. The detailed exploration, feature engineering, model selection, and evaluation processes underscore the effort to create an effective solution for spam detection.

By leveraging advanced machine learning techniques and comprehensive feature engineering, we have developed a spam detection model that meets the stakeholder's requirements. The XGBoost model, in particular, stands out as a reliable and efficient solution for the company's spam detection needs. The deployment plan ensures that the model can be integrated and maintained effectively, providing long-term value to the company.