

# Progress Report

## Baseline Model (M1) implementation

With the Naive Bayes Classifier as our model of choice, we have successfully developed a baseline model for our spam detection task [1]. The Naive Bayes method was chosen because it is straightforward, effective, and has a track record of success in text categorization applications like spam detection.

*The following were the main steps in the implementation:*

1. **Data preparation:** The 'SMSSpamCollection' dataset was imported and prepared. The dataset was preprocessed by randomizing it, converting it to lowercase, deleting non-word characters, and tokenizing (dividing the text into individual words).
2. **Training/Test split:** Data was split into training and test datasets, with 80% of the data going to training and 20% to testing.
3. **Feature Extraction:** To transform the text input into numerical features, we employed a bag-of-words model. The dataset's individual words were each handled as a feature.
4. **Model Training:** Using the training set of data, the Naive Bayes classifier was trained. In order to determine the conditional probabilities of each word given a class (spam or ham), the model was implemented.
5. **Testing and Evaluation of the Model:** The accuracy of the trained model was evaluated using test data that had not been seen.

## Research and Performance Evaluation

With an accuracy of about 86.44% on the test dataset, the Naive Bayes Classifier did reasonably well at detecting spam. Although this performance shows promise, there is still potential for development.

## Stretch Goals

Future plans call for:

1. **Implement Additional Classifiers:** In order to enhance performance, we will test additional classifiers such K-Nearest Neighbors (K-NN). It has been demonstrated that K-NN works well for text classification issues [2].
2. **Integrate Performance Metrics:** In addition to assessing the model's accuracy, we will also monitor its efficiency in terms of runtime and memory consumption. In real-world applications, where efficiency is frequently a crucial factor [3], these measures are significant.

## References

[1] Rish, Irina. "An empirical study of the naive Bayes classifier." IJCAI 2001 workshop on empirical methods in artificial intelligence. Vol. 3. 2001.

[2] Soucy, Pierre, and Guy W. Mineau. "A simple KNN algorithm for text categorization." Proceedings 2001 IEEE International Conference on Data Mining. IEEE, 2001.

# Progress Report

[3] Witten, Ian H., et al. "Text mining: practical prediction of text." More data mining with Weka (2016): 9-17.