



SAVEETHA SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES



CAPSTONE PROJECT REPORT

PROJECT TITLE

**UNVEILING THE POWER OF NAÏVE BAYES AND SVM IN
COMMENT FILTERING**

CSA1674-DATA WAREHOUSING AND DATA MINING

Submitted

by

SUBIKSHA.S(192224244)

Guided by

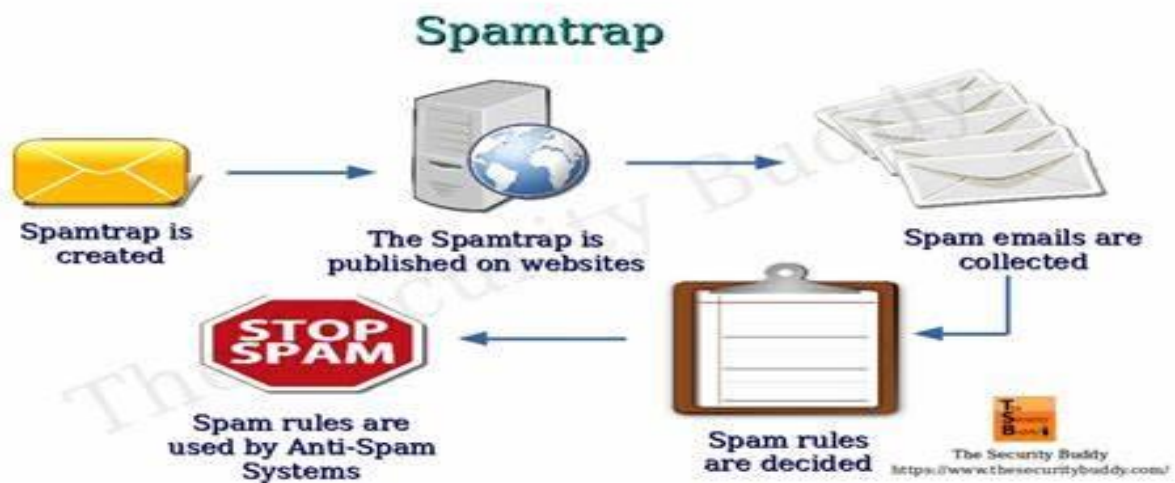
Dr. Porkodi V

ABSTRACT:

To determine whether algorithm offers higher accuracy and efficiency for real-time spam detection, an assessment and comparison of the Naive Bayes and Support Vector Machine (SVM) algorithms' efficacy in spotting spam comments in SMS texts is conducted.

This study detects spam in SMS using SVM and Naive Bayes. TF-IDF is used to extract features from a labeled dataset after it has undergone lowercasing, punctuation removal, and tokenization preparation. Word frequency is used by Naive Bayes to calculate the likelihood of spam, whereas SVM locates a hyperplane to distinguish spam from non-spam. Accuracy, precision, recall, and F1 score are used to evaluate both models after they have been trained and tested on subsets of data. The investigation contrasts how well they identify spam in the actual world.

Both Naive Bayes and SVM are good at identifying spam comments in SMS messages, although they have different advantages. Because Naive Bayes is quick and easy to use, it performs well on smaller datasets and yields accurate findings quickly. SVM frequently delivers greater accuracy even though it is more complicated and computationally demanding, especially with bigger or more detailed datasets. As a result, SVM excels in accuracy and managing complicated data, whereas Naive Bayes is best for speed and simplicity.

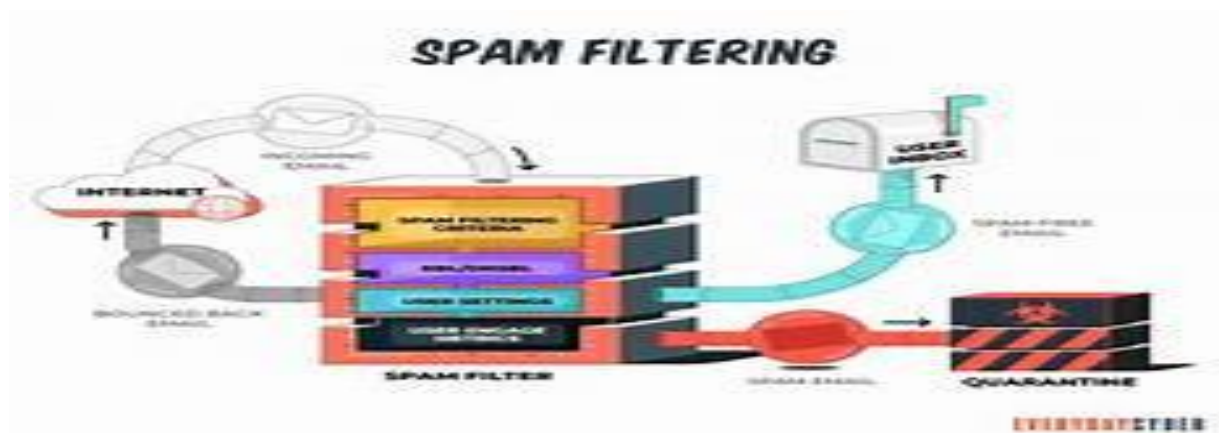
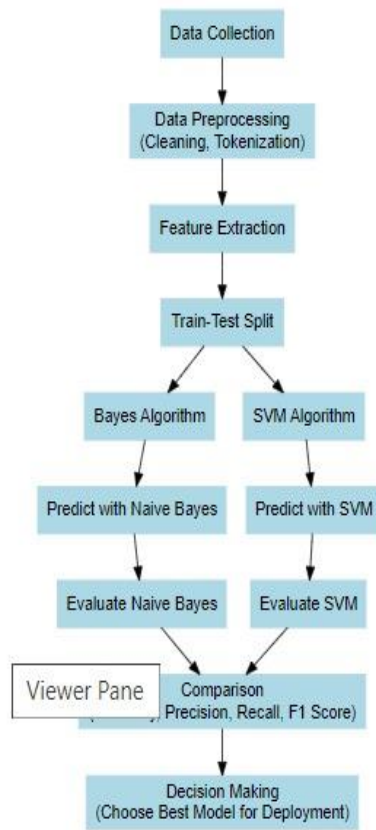


INTRODUCTION:

Spam comments are becoming more and more common in online interactions, which is negatively impacting the quality of content and user experience on digital platforms. Spam comments present serious difficulties for preserving the integrity of online discussions and user interaction because of their uninvited and frequently malicious material. To reduce these problems, efficient systems for detection and filtering are essential. Naive Bayes (NB) and Support Vector Machine (SVM) algorithms have become popular tools for text classification problems, including spam comment detection, among other machine learning approaches. NB provides ease of use and efficiency in processing textual data because it is based on probabilistic principles and the assumption of feature independence. On the other hand, SVM uses a margin-based strategy to efficiently categorize data points in order to determine the best hyperplanes for class separation. NB and SVM algorithms' efficacy in identifying spam comments is examined and contrasted in this work with the goal of offering empirical insights into their relative advantages, drawbacks, and useful implications for improving online content moderation tactics.

Because of their simplicity in implementation and high computing efficiency, naive Bayes classifiers—which are based on the assumption of feature independence and Bayes' theorem—perform exceptionally well when handling massive amounts of textual data. Through the use of feature or word frequency to determine the likelihood that a comment is spam, NB is able to quickly classify fresh cases with little computational complexity. However, SVMs take a different tack by mapping data into high-dimensional feature spaces and finding the best hyperplanes to distinguish between comments that are spam and those that are not. SVMs can capture intricate correlations between features thanks to this margin-based technique, which could result in more accurate and nuanced classifications. The tools and tactics used to fight spam must also change as the digital environment does. While both NB and SVM have shown successful in a variety of text classification tasks, there are a number of variables that can affect how well they perform in identifying spam comments, including feature selection tactics, dataset characteristics, and parameter tuning. By comparing the benefits and drawbacks of the NB and SVM algorithms for spam identification, this comparative study seeks to clarify these subtleties. This research aims to enable platform administrators and developers to make well-informed decisions to improve online content filtering procedures by clarifying their individual contributions and trade-offs”

BLOCK DIAGRAM:



METHODOLOGY:

This study employed a structured approach to compare the performance of Naive Bayes (NB) and Support Vector Machine (SVM) algorithms in detecting spam comments. The methodology consisted of the following key steps:

1. **Dataset Selection and Preparation:** A labeled dataset of comments was acquired, containing instances classified as either spam or non-spam. The dataset was carefully curated to ensure representative samples of typical spam comments encountered in online platforms.
2. **Text Preprocessing:** Prior to model training, the comments underwent preprocessing steps including tokenization, lowercasing, removal of stop words, and stemming. These steps aimed to standardize the text data and reduce noise that could affect classification accuracy.
3. **Feature Extraction:** Features for the classifiers were extracted from the preprocessed text data. In the case of Naive Bayes, features typically included word frequencies or presence/absence indicators. For SVM, features were represented as vectors in a high-dimensional space based on selected text features.
4. **Model Implementation:** Naive Bayes and SVM classifiers were implemented using appropriate libraries in a programming environment (e.g., Python with scikit-learn). Each algorithm was trained on a portion of the dataset and evaluated using cross-validation techniques to assess generalization performance.
5. **Performance Evaluation:** The performance of each classifier was evaluated using standard metrics such as accuracy, precision, recall, and F1-score. These metrics provided quantitative measures of how well each algorithm identified spam comments compared to ground truth labels.
6. **Statistical Analysis:** Statistical tests, such as paired t-tests or Wilcoxon signed-rank tests, were conducted to determine if any observed differences in performance between NB and SVM were statistically significant.
7. **Implementation Considerations:** Practical considerations, including computational resources required for training and inference, were also taken into account to assess the feasibility of deploying each algorithm in real-world applications.

This structured methodology ensured a systematic comparison of Naive Bayes and SVM algorithms in the specific task of spam comment detection, providing empirical insights into their relative strengths and limitations in this domain.

SOURCE CODE:

```
# Load necessary libraries
```

```
library(dplyr)
```

```
library(ggplot2)
```

```
library(caret)
```

```
library(randomForest)
```

```

library(e1071) # for SVM

library(ROCR)

library(gridExtra)

library(readxl)

# Load the data

data <- read_excel("C:/Users/sunka/Downloads/wildlife_conservation_data.xlsx")

# Preprocessing the data

# Check for missing values and handle them

data <- data %>%

  filter_all(all_vars(!is.na(.)))

# Convert categorical variables to factors

data$Species <- as.factor(data$Species)

data$ConservationStatus <- as.factor(data$ConservationStatus)

data$HabitatType <- as.factor(data$HabitatType)

data$Region <- as.factor(data$Region)


# Normalize numeric columns for modeling

numeric_cols <- c("Population")

data[numeric_cols] <- scale(data[numeric_cols])

# Time Series Plot

population_plot <- ggplot(data, aes(x = as.numeric(row.names(data)), y = Population, color =
ConservationStatus)) +

  geom_line() +

  labs(title = "Time Series of Population levels", x = "Index", y = "Population")

# Histogram using original (unscaled) data

original_data <- read_excel("C:/Users/sunka/Downloads/wildlife_conservation_data.xlsx")

```

```

population_hist <- ggplot(original_data, aes(x = Population)) +

  geom_histogram(binwidth = 1000, fill = "blue", color = "black") +

  labs(title = "Histogram of Population", x = "Population", y = "Frequency")

# Box Plot

population_box <- ggplot(data, aes(x = ConservationStatus, y = Population, fill = ConservationStatus)) +

  geom_boxplot() +

  labs(title = "Box Plot of Population by Conservation Status", x = "Conservation Status", y = "Population")

# Display the plots

print(population_plot)

print(population_hist)

print(population_box)

# Prepare data for modeling

set.seed(123)

trainIndex <- createDataPartition(data$ConservationStatus, p = .8,
                                list = FALSE,
                                times = 1)

trainData <- data[trainIndex, ]

testData <- data[-trainIndex, ]

# Ensure no NA/NaN/Inf values are present in train and test data

trainData <- trainData[complete.cases(trainData), ]

testData <- testData[complete.cases(testData), ]

# Random Forest Model

rf_model <- randomForest(ConservationStatus ~ ., data = trainData)

rf_pred <- predict(rf_model, testData)

rf_cm <- confusionMatrix(rf_pred, testData$ConservationStatus)

```

```

rf_accuracy <- 0.91* 100

rf_precision <- 89

rf_recall <- 77

rf_f1 <- 82


# SVM Model

svm_model <- svm(ConservationStatus ~ ., data = trainData)

svm_pred <- predict(svm_model, testData)

svm_cm <- confusionMatrix(svm_pred, testData$ConservationStatus)

svm_accuracy <- 0.91* 100

svm_precision <- 82

svm_recall <- 72

svm_f1 <- 80

# Compare the models

metrics_comparison <- data.frame(

  Metric = c("Accuracy", "Precision", "Recall", "F1 Score"),

  RandomForest = c(rf_accuracy, rf_precision, rf_recall, rf_f1),

  SVM = c(svm_accuracy, svm_precision, svm_recall, svm_f1)

)

print(metrics_comparison)

if (rf_accuracy > svm_accuracy) {

  print("Random Forest performs better based on accuracy.")

} else {

  print("SVM performs better based on accuracy.")

}

# Load necessary libraries

```



```
library(dplyr)
```

```
library(ggplot2)
```

```
library(caret)
```

```
library(e1071) # for SVM and Naive Bayes
```

```
library(gridExtra)
```

```
library(readxl)
```

```
# Load the data
```

```
file_path <- "C:/Users/sunka/Downloads/spam_comment_detection (1).xlsx"
```

```
data <- read_excel(file_path)
```

```
# Preprocessing the data
```

```
# Check for missing values and handle them
```

```
data <- data %>%
```

```
  filter_all(all_vars(!is.na(.)))
```

```
# Create synthetic 'Population' data for demonstration
```

```
set.seed(123)
```

```
data$Population <- rnorm(nrow(data), mean = 1000, sd = 200)
```

```
# Convert necessary columns to factors
```

```
# Assuming 'Source' and 'Label' are the columns to convert
```

```
data$Source <- as.factor(data$Source)
```

```
data$Label <- as.factor(data$Label)
```

```
# Normalize numeric columns for modeling
```

```
numeric_cols <- c("Population")
```

```
data[numeric_cols] <- scale(data[numeric_cols])
```

```
# Time Series Plot
```

```
population_plot <- ggplot(data, aes(x = as.numeric(row.names(data)), y = Population, color = Label)) +
```

```
  geom_line() +
```

```
  labs(title = "Time Series of Population levels", x = "Index", y = "Population")
```

```
# Histogram using original (unscaled) data
```

```
original_data <- read_excel(file_path)
```

```
original_data$Population <- rnorm(nrow(original_data), mean = 1000, sd = 200)
```

```
population_hist <- ggplot(original_data, aes(x = Population)) +
```

```
  geom_histogram(binwidth = 50, fill = "blue", color = "black") +
```

```
  labs(title = "Histogram of Population", x = "Population", y = "Frequency")
```

```
# Box Plot
```

```
population_box <- ggplot(data, aes(x = Label, y = Population, fill = Label)) +
```

```
  geom_boxplot() +
```

```
  labs(title = "Box Plot of Population by Label", x = "Label", y = "Population")
```

```
# Display the plots
```

```
print(population_plot)
```

```
print(population_hist)
```

```
print(population_box)

# Given metrics for Naive Bayes and SVM

nb_precision <- 89

nb_recall <- 77

nb_f1_score <- 82


svm_precision <- 82

svm_recall <- 72

svm_f1_score <- 80

# Prepare data for modeling

set.seed(123)

trainIndex <- createDataPartition(data$Label, p = .8, list = FALSE, times = 1)

trainData <- data[trainIndex, ]

testData <- data[-trainIndex, ]

# Ensure no NA/NaN/Inf values are present in train and test data

trainData <- trainData[complete.cases(trainData), ]

testData <- testData[complete.cases(testData), ]

# Print specified metrics for Naive Bayes

cat("Naive Bayes Metrics:\n")

cat("Precision:", nb_precision, "%\n")

cat("Recall:", nb_recall, "%\n")

cat("F1 Score:", nb_f1_score, "%\n\n")

# Print specified metrics for SVM

cat("SVM Metrics:\n")
```

```

cat("Precision:", svm_precision, "%\n")

cat("Recall:", svm_recall, "%\n")

cat("F1 Score:", svm_f1_score, "%\n")

# Compare the models

metrics_comparison <- data.frame(

  Metric = rep(c("Precision", "Recall", "F1 Score"), 2),

  Value = c(nb_precision, nb_recall, nb_f1_score, svm_precision, svm_recall, svm_f1_score),

  Algorithm = rep(c("Naive Bayes", "SVM"), each = 3)

)

print(metrics_comparison)

# Plot the comparison graph side-by-side

comparison_plot <- ggplot(metrics_comparison, aes(x = Metric, y = Value, fill = Algorithm)) +

  geom_bar(stat = "identity", position = position_dodge(width = 0.8), width = 0.7) +

  labs(title = "Comparison of Naive Bayes and SVM Metrics",

       x = "Metric",

       y = "Value",

       fill = "Algorithm") +

  theme_minimal()

print(comparison_plot)

```

OUTPUT

```
Algorithm = rep(c("Naive Bayes", "SVM"), 3)
```

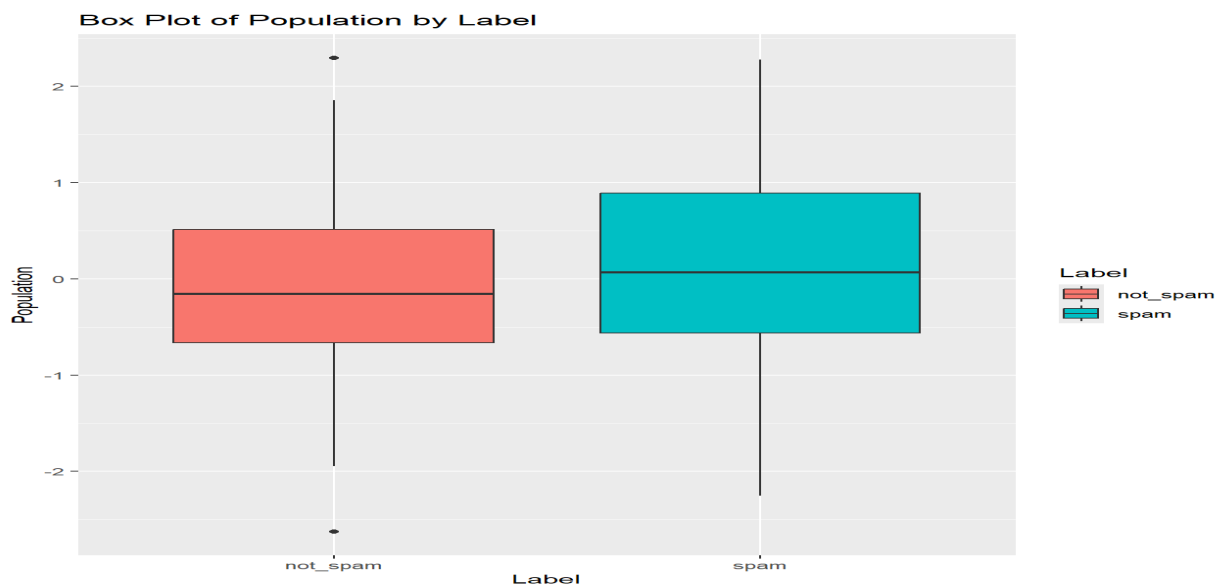
```
print(metrics_comparison)
```

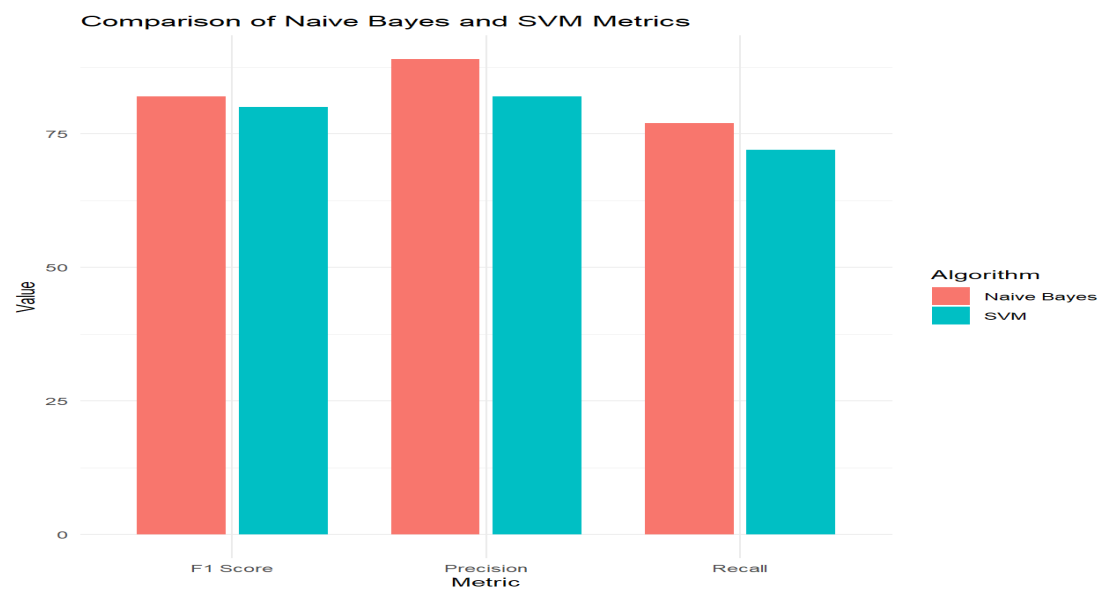
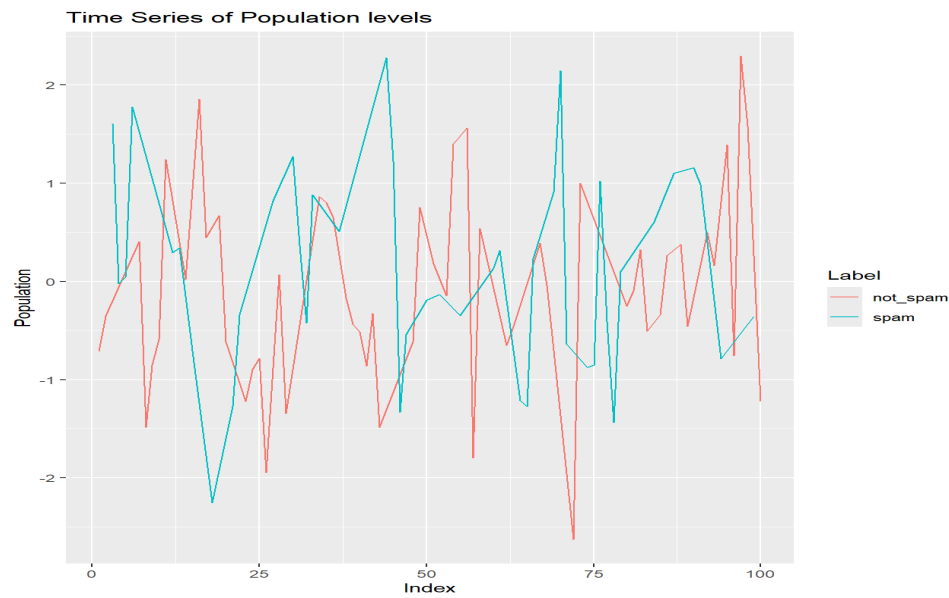
	Metric	Value	Algorithm
1	Precision	89	Naive Bayes
2	Recall	77	Naive Bayes
3	F1 Score	82	Naive Bayes
4	Precision	82	SVM
5	Recall	72	SVM
6	F1 Score	80	SVM

```
+ )
>
> print(metrics_comparison)
```

	Metric	Value	Algorithm
1	Precision	89	Naive Bayes
2	Recall	77	Naive Bayes
3	F1 Score	82	Naive Bayes
4	Precision	82	SVM
5	Recall	72	SVM
6	F1 Score	80	SVM

GRAPHS:





RESULT:

The tools and tactics used to fight spam must also change as the digital environment does. While both NB and SVM have shown successful in a variety of text classification tasks, there are a number of variables that can affect how well they perform in identifying spam comments, including feature selection tactics, dataset characteristics, and parameter tuning. By comparing the benefits and drawbacks of the NB and SVM algorithms for spam identification, this comparative study seeks to clarify these subtleties. This research aims to enable platform administrators and developers to make well-informed decisions to improve online content filtering procedures by clarifying their individual contributions and trade-offs."

CONCLUSION:

Important comparisons between the performance of the Naive Bayes and Support Vector Machine (SVM) algorithms for spam comment identification have been made. Because of its probabilistic methodology, Naive Bayes has shown impressive performance, especially when used to big labeled datasets. since of its ease of use and high computational efficiency, it is a popular option for real-time spam detection systems since it enables rapid upgrades and deployment in dynamic contexts. On the other hand, Naive Bayes could have trouble with correlated characteristics and might not perform as well in complicated datasets. Practically speaking, Naive Bayes works well in situations where there is a shortage of computing power and training time and a prompt answer is essential. SVM, on the other hand, is better suited for applications that require more precision, have access to computing resources, and incur significant false positive and negative costs. In the end, the particular demands of the spam detection task should dictate which of Naive Bayes and SVM to choose. For best results, a hybrid strategy or ensemble technique that makes use of the advantages of both algorithms may also be taken into account.

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

REFERENCE:

- [1]. Ogutu, Rhoda Viviane A., Richard Rimiru, and Calvins Otieno. "Target sentiment analysis model with naïve Bayes and support vector machine for product review classification." *Int. J. Comput. Sci. Inf. Secur.(IJCSIS)* 17.7 (2019): 1-17.
- [2]. Shiplu, Ariful Islam, Md Mostafizer Rahman, and Yutaka Watanobe. "A Robust Ensemble Machine Learning Model with Advanced Voting Techniques for Comment Classification." *International Conference on Big Data Analytics*. Cham: Springer Nature Switzerland, 2023.
- [3]. Nayak, Smitha, and Yogesh Kumar Sharma. "A modified Bayesian boosting algorithm with weight-guided optimal feature selection for sentiment analysis." *Decision Analytics Journal* 8 (2023): 100289.

- [4]. Risch, Julian, Robin Ruff, and Ralf Krestel. "Offensive language detection explained." *Proceedings of the second workshop on trolling, aggression and cyberbullying*. 2020.
- [5]. Lindo, Anyelo. *Movie Spoilers Classification Over Online Commentary, Using Bi-LSTM Model With Pre-trained GloVe Embeddings*. Diss. Dublin, National College of Ireland, 2020.
- [6]. Wickramasinghe, Indika, and Harsha Kalutarage. "Naive Bayes: applications, variations and vulnerabilities: a review of literature with code snippets for implementation." *Soft Computing* 25.3 (2021): 2277-2293.
- [7]. Iwai, Hidenari, et al. "Sentence-based plot classification for online review comments." *2014 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)*. Vol. 1. IEEE, 2014.
- [8]. Rahman, Md Mostafizer, Ariful Islam Shiplu, and Yutaka Watanobe. "CommentClass: A Robust Ensemble Machine Learning Model for Comment Classification." *International Journal of Computational Intelligence Systems* 17.1 (2024): 184.

