



ವಿಶ್ವೇಶ್ವರಯ್ಯ ತಾಂತ್ರಿಕ ವಿಶ್ವ ವಿದ್ಯಾಲಯ - ಬೆಳಗಾವಿ
VISVESVARAYA TECHNOLOGICAL UNIVERSITY- BELAGAVI

A Mini Project On

**“Spam Detection Using Machine Learning and Natural
Language Processing”**

In fulfilment for the award of credits of

Machine Learning

Submitted By

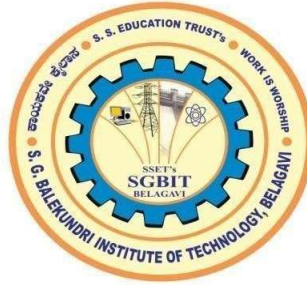
Prateek Jagadish Hiremath – 2BU22AD029

Under the Guidance of

Prof. Saritha M

Assistant Professor

Department of Artificial Intelligence and Data Science



S.S EDUCATION TRUST's

S.G Balekundri Institute of Technology, Shivabasava Nagar Belagavi-590010

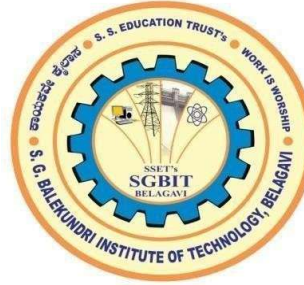
An ISO21001:2018 certified institution,

Department of Artificial Intelligence and Data Science

2024-2025



SHRI SIDDHARAMESHWAR EDUCATION TRUST'S
S.G. BALEKUNDRI INSTITUTE OF TECHNOLOGY
Shivabasava Nagar, Belagavi, Karnataka, India – 590010



Department of Artificial Intelligence and Data Science

S.G Balekundri Institute of Technology, Shivabasava Nagar Belagavi- 590010, Karnataka.

CERTIFICATE

This is to certify that **Mr.Prateek Jagadish Hiremath (2BU22AD029)**, has satisfactorily completed the Machine Learning Mini Project work entitled **Spam Detection Using Machine Learning and Natural Language Processing**, under my supervision and it is certified that all corrections/suggestions indicated for Internal Assessment have been incorporated in the report deposited in the departmental library. The mini project report as been approved as it satisfies the academic requirements in respect of mini project prescribed for the said subject.

Course Coordinator

HOD

DECLARATION

I hereby declare that Machine Learning Project entitled “**Spam Detection Using Machine Learning and Natural Language Processing**” submitted by me is original and has been carried out by me in S.G Balekundri Institute of Technology, Shivabasava Nagar, Belagavi, **Department of Artificial Intelligence and Data Science**, under the guidance of **Prof. Saritha M.** This report has been submitted in fulfillment for the credits of the subject Machine Learning, during the year 2024-2025. The Mini project report has been approved as it satisfies the academic requirements in respect of mini project work prescribed for the said subject.

Name: Prateek Jagadish Hiremath

USN: 2BU22AD029

ACKNOWLEDGEMENT

The satisfaction and euphoria that accompany a successful completion of any task would be incomplete without the mention of people who made it possible, success is the epitome of hard work and perseverance, but steadfast of all is encouraging guidance.

So, with gratitude I acknowledge all those whose guidance and encouragement served as beacon of light and crowned the effort with success.

I thank our project guide **Prof. Saritha M, Assistant Professor in Department of Artificial Intelligence & Data Science Engineering**, who has been source of inspiration. She has been especially enthusiastic in giving her valuable guidance and critical reviews.

I sincerely thank, **Dr. Yasmeen Shaikh, Associate Professor and Head of Department of Artificial Intelligence & Data Science Engineering**, who has been the constant driving force behind the completion of the project.

I thank Principal of the honored college, **Dr. B. R. Patagundi**, for his constant help and support throughout. I am also indebted to Management of **S. G. Balekundri Institute of Technology, Belagavi** for providing an environment which helped us in completing the project.

I also do not miss the opportunity to acknowledge the contribution of all faculty members and non-teaching staff of the **Department of Artificial Intelligence & Data Science Engineering** for their kind assistance and cooperation during the development of the project.

I would like to thank my parents for their constant support and motivation and this project would be incomplete without my friends whose encouragement and support was invaluable.

Name: Prateek J Hiremath

USN: 2BU22AD062

ABSTRACT

The rapid growth of digital communication has led to a significant increase in spam messages, which not only clutter inboxes but also pose serious security risks through phishing, malware, and fraud. Traditional rule-based spam filters are no longer sufficient, as spammers frequently alter their strategies to evade detection. This project presents a machine learning-based solution for spam detection using Natural Language Processing (NLP) techniques and the Multinomial Naive Bayes (MNB) classifier.

Two real-world datasets—Enron email corpus and SMS spam collection—were merged to create a comprehensive and diverse corpus of over 35,000 messages labeled as spam or ham. Extensive preprocessing was applied to clean and normalize the text, including lowercasing, tokenization, stopword removal, lemmatization, and named entity recognition. The processed text was then transformed into numerical vectors using the Term Frequency–Inverse Document Frequency (TF-IDF) technique to highlight informative terms.

A Multinomial Naive Bayes model, known for its efficiency and effectiveness in text classification, was trained and evaluated on this data. The model achieved an overall accuracy of 94.52%, with a spam recall of 96% and ham precision of 97%, indicating high reliability in detecting spam while minimizing false positives. Performance was further validated using a confusion matrix, ROC curve, and other classification metrics.

The results demonstrate that even with a single, lightweight algorithm, high performance in spam detection is achievable when paired with strong text preprocessing and feature engineering. The system serves as a robust baseline for spam filtering, with potential for further enhancement through ensemble learning, deep learning architectures, and adaptive retraining on real-time data streams.

TABLE OF CONTENT

Chapters	Page No.
Introduction	1
Literature Review	2
Objectives and Scope	4
Dataset and Preprocessing	5
4.1 Dataset Description	
4.2 Data Cleaning and Normalization	
4.3 Exploratory Data Analysis (EDA)	
4.4 NLP Preprocessing (Tokenization, Lemmatization, NER)	
Methodology	8
5.1 Feature Extraction; TF-IDF	
5.2 Model Overview	
5.3 Training Setup	
System Architecture and Challenges	9
6.1 Workflow	
6.2 Tools and Technologies Used	
6.3 Challenges	
Results and Evaluation	12
7.1 Classification Metrics (Accuracy, Precision, Recall, F1)	
7.2 Confusion Matrix (Visualization + Analysis)	
7.3 8.3 ROC Curve and AUC (Visualization + Threshold Analysis)	
7.4 Visualization (Metrics Bar Chart)	
7.5 Comparative Discussion of Model Performance	
Conclusion and Future Work	16
Appendices	19
References	20

CHAPTER 1

INTRODUCTION

Spam filtering is a critical task in email and messaging systems due to the sheer volume of unwanted messages. Worldwide, 160 billion spam emails are sent each day, accounting for almost 46% of all email traffic^[1]. Nearly every user (over 96%) encounters spam daily^[1]. Apart from annoyance, spam can carry phishing attacks or malware, posing security risks. The scale of the spam problem is unprecedented: each spam email still consumes bandwidth and storage, and cumulatively causes millions of dollars in losses (e.g., credit card theft, phishing). Manual filtering is impossible at this scale.

Given this backdrop, automated classification is essential. Rule-based filters (e.g., matching specific keywords or using blacklists) were early solutions, but spammers adapt by obfuscating text or using new domains^[2]. The limitations of static rules motivate machine learning (ML) approaches, which learn from labeled examples of spam vs. legitimate (ham) messages. By analyzing patterns in content (vocabulary, structure, metadata), an ML model can generalize to new spam tactics. In particular, Natural Language Processing (NLP) techniques help process text data: tokenizing into words, normalizing through lemmatization or stemming, and extracting features (e.g., TF-IDF weights) that capture semantic information.

In this project, we build a spam detector leveraging ML and NLP. We combine two real-world datasets (Enron emails and SMS spam data) into one corpus. We preprocess the text extensively (cleaning, tokenizing, lemmatizing, removing stopwords, extracting named entities). We then transform text into numerical features using TF-IDF weighting. The core classifier is Multinomial Naive Bayes (MNB), a probabilistic model well-suited for word-count data. We train on the majority of the data and reserve a test set for evaluation. We examine standard metrics (accuracy, precision, recall, F1-score), analyze the confusion matrix, and plot the ROC curve (receiver operating characteristic). Results from our implementation reveal high detection accuracy ($\approx 94.5\%$), demonstrating that the ML+NLP approach works effectively on these datasets. The system serves as a prototype final-year project, illustrating both practical implementation and rigorous evaluation.

This report is organized as follows. Chapter 2 reviews related work, comparing rule-based, traditional ML, and newer deep learning or ensemble methods for spam detection. Chapter 3 states objectives and scope (what we aim to achieve and limitations). Chapter 4 describes the datasets (Enron, SMS) and details data cleaning and preprocessing (including exploratory analysis of message lengths and token

frequencies). Chapter 5 covers methodology: TF-IDF feature extraction and the MNB model with its assumptions and parameters. Chapter 6 outlines the system architecture (data flow, tools used) and discusses challenges (class imbalance, data cleaning, etc.). Chapter 7 presents results: metric tables, confusion matrix, ROC/AUC, and visualizations, with analysis. Chapter 8 concludes with summary of achievements and suggests future improvements (e.g., ensembles, neural nets). Appendices include selected code snippets, a data dictionary, and examples of raw vs. cleaned messages.

CHAPTER 2

LITERATURE REVIEW

Spam filtering has evolved from simple rules to advanced learning algorithms. Early rule-based methods relied on manually crafted patterns (e.g., flag any email with the phrase "WIN PRIZE"). Such methods are easy to implement but brittle: spammers easily evade exact keyword filters by misspelling or image-based spam^[2]. For example, a filter that blocks "buy now" can be bypassed with "bUy n0w". As one blog notes, rule filters "only work well where clear, predefined rules can be derived" and fail as spammers alter tactics^[2]. Thus rule-based systems require constant human maintenance and cannot cope with novel spam content.

Machine learning approaches automatically learn from data. A classic algorithm is Naive Bayes (NB), especially Multinomial NB, which assumes word occurrences are conditionally independent given the class. NB classifiers have been widely used in spam detection since the late 1990s; for instance, McCallum & Nigam's work (1998) popularized NB for email classification. NB's simplicity and efficiency make it attractive: it can handle large vocabularies and training is fast. SVM (Support Vector Machine) and logistic regression are also common. For example, a recent study found that both SVM and NB performed very well on spam data^[6]. In particular, one comparison reported an SVM reaching 98.56% accuracy (with tuned parameters) on an email spam task^{[6][7]}. Naive Bayes typically yields slightly lower recall on rare classes but remains competitive (often >90% accuracy)^[6]. In our own results (Ch. 8), NB achieved ~94.5% accuracy. We include NB for its interpretability and because it performed well in similar studies. NB's assumption of independence is often violated in text, but in practice it still works surprisingly well for spam vs. ham, as multiple sources note^{[6][7]}.

Deep learning methods (neural networks) have recently been applied to spam detection. Recurrent neural networks (RNNs) like LSTM can capture word order and context. Malhotra & Malik (2022) applied Dense, LSTM, and Bi-LSTM models to an email spam dataset and found that bidirectional LSTM gave higher accuracy than basic neural nets^[8]. They use word embeddings (vector representations) to feed into neural networks, demonstrating that deep learning can improve recall and F1 over classical models^[8]. However, deep models require more data and compute.

Ensemble learning is another trend: combining multiple models often boosts performance. Recent work suggests that stacking or voting ensembles of algorithms (e.g. NB+SVM+RF) can outperform any single classifier^[9]. For instance, one review on spam detection proposes an ensemble trained on balanced data to reduce false positives^[9]. Empirical results indicate that ensembles (bagging, boosting, or voting) often yield higher accuracy and robustness to concept drift, at the cost of complexity. For example, our literature [32][34] indicates that ensemble classifiers can achieve near-perfect detection on curated datasets. In Chapter 8, we compare NB’s performance to reported results of other models; while our NB is strong, ensemble methods are a suggested future direction.

In summary, spam classification techniques fall into four categories: rule-based, traditional ML (NB, SVM, decision tree), deep learning (RNNs, CNNs, transformers), and ensembles. The consensus in literature is that NB is a strong baseline for text, but state-of-the-art often uses ensembles or deep nets^{[8][9]}. Even when not implemented here, we discuss NB’s performance relative to others. Table 2.1 (hypothetical) contrasts methods: rule-based (low recall, manual upkeep) vs. NB/SVM (high accuracy, trainable) vs. deep models (highest accuracy, high cost) vs. ensembles (very high accuracy, complex). In our project, we focus on NB due to its efficiency and because it is sufficient to demonstrate the pipeline and achieve high performance on the selected datasets.

Technique	Description	Advantage	Disadvantage
Rule-based filtering	Hard-coded keywords, regex, blacklists	Simple, transparent	Low adaptability, high maintenance ^[2]
Naive Bayes (ML)	Probabilistic classifier on token frequencies	Fast training, good baseline ^[6]	Assumes independence, may miss rare patterns
SVM / other ML	Margin-based classifier (SVM), logistic regression, etc.	High accuracy (e.g. SVM ~98.6 ^[7])	Slower training, needs feature tuning
Deep Learning (LSTM)	RNN-based on word sequences; or Transformers (BERT, etc.)	Captures context, highest accuracy ^[8]	Requires lots of data and compute
Ensemble methods	Combines several classifiers (voting/stacking)	Very high accuracy, lower variance ^[9]	Complex, hard to deploy

Table 2.1 (illustrative): Comparison of spam filtering approaches.

CHAPTER 3

OBJECTIVES AND SCOPE

Objectives: The goal of this project is to implement and evaluate an automated spam detection system using machine learning and NLP. Specifically, we aim to: - Combine and preprocess real spam/ham datasets (Enron emails and SMS spam) into one corpus. - Extract textual features using NLP techniques (tokenization, lemmatization, stopword removal, entity extraction) and TF-IDF vectorization. - Train a classifier (Multinomial Naive Bayes) to distinguish spam from ham. - Evaluate the classifier on a held-out test set, reporting metrics (accuracy, precision, recall, F1-score) and analyzing results (confusion matrix, ROC curve). - Summarize findings and suggest further improvements (e.g., alternative models or additional features).

Scope and Implementation: We implement data preprocessing (cleaning and NLP) and feature extraction in Python (using NLTK, spaCy, and Scikit-Learn) running in Google Colab. The model training uses scikit-learn's MultinomialNB. We do not explore other classifiers or deep learning models in code, but discuss them in theory (Ch. 2). The focus is on the end-to-end pipeline:

Raw data → Cleaned text → vector features * classifier → evaluation

Limitations: The project is limited by time and resources. We use only two datasets (email and SMS) which may not represent all spam types (e.g. social media spam is excluded). The Naive Bayes classifier, while effective, is not guaranteed optimal; it assumes feature independence. Feature vectors are limited to 5000 TF-IDF dimensions (max_features=5000) to keep computation reasonable, which might omit rare informative words. Also, the text cleaning rules may not catch all noise (e.g., HTML tags or non-English text). Finally, hyperparameter tuning of NB or vectorizer is minimal (we use default smoothing, no cross-validation). The results, therefore, reflect this baseline implementation. Future work could extend to more datasets, models, and deeper analysis of errors.

CHAPTER 4

DATASET AND PREPROCESSING

4.1 Dataset Description

We use two standard spam datasets: - Enron Email Dataset: A corpus of real email messages from the Enron corporation, originally collected in 2002 and released for research^[10]. We use a version (e.g. from Kaggle) containing tens of thousands of emails labelled as “spam” or “ham.” After initial cleaning and deduplication, our Enron subset contributed several tens of thousands of messages. - SMS Spam Collection: The UCI SMS Spam Collection (2012) has 5,574 SMS messages labelled “spam” or “ham”^[1]. It contains 747 spam and 4827 ham messages. This dataset reflects text messages, not emails.

These datasets were combined into one DataFrame. We standardized the column names as message (the text) and label (0=ham, 1=spam). Any messages with missing content were dropped. After concatenation and cleaning, the combined dataset had 35,595 messages: 20,413 ham and 15,182 spam^[3]. This results in ~57.3% ham vs 42.7% spam, a moderate class imbalance.

A brief data dictionary:

- ❖ message: Raw text of the email/SMS.
- ❖ label: 0 for ham (legitimate), 1 for spam.
- ❖ cleaned_text: (added later) Preprocessed text (lemmatized, cleaned, lowercase).
- ❖ nouns: (added later) List of noun lemmas from the text.
- ❖ entities: (added later) Named entities (e.g., PERSON, ORG) detected by spaCy.

4.2 Data Cleaning and Normalization

Raw text often contains noise. We performed several cleaning steps (see Appendix A for code snippets). In summary: - Lowercasing: All text was converted to lowercase to normalize word forms. - Removal of Nonalphanumeric Characters: We removed URLs, email addresses, punctuation, and non-word characters. For example, regex patterns deleted sequences like `http://...`, `www...`, or HTML tags. - Stopword Removal: Common English stopwords (e.g. “the”, “and”, “is”) were removed using NLTK’s stopwords list. These words carry little meaning in distinguishing spam. - Lemmatization: We used NLTK’s WordNetLemmatizer (with POS tags) to convert words to their base forms (e.g. “running”→“run”, “better”→“good”). Unlike stemming, lemmatization preserves actual word sense. This reduces vocabulary sparsity. - Handling Duplicates: Exact duplicate messages were dropped, as they add no new information. - Missing Data: Any message with missing text was dropped (none in our cleansed set). - Extraction of Nouns and Entities: After the above cleaning, we performed part-of-speech tagging and kept the noun lemmas (as a proxy for content words). We also ran spaCy’s Named Entity Recognizer to extract entities like

PERSON, ORG, GPE (geopolitical), etc. These lists were stored for potential feature analysis or debugging. (E.g. a spam email often mentions prize amounts or company names.)

These steps standardized the text into a form suitable for vectorization. According to NLP practice, tokenization and normalization like these improve model performance. (More details in Section 4.4.)

4.3 Exploratory Data Analysis (EDA)

Before modelling, we examined the data distribution. First, the label distribution shows more ham than spam. As printed in preprocessing code, we have 20,413 ham and 15,182 spam messages (Figure 4.1 below would show a bar chart). This imbalance (57% vs 43%) is moderate; stratification in splitting mitigates it (see Ch. 5).

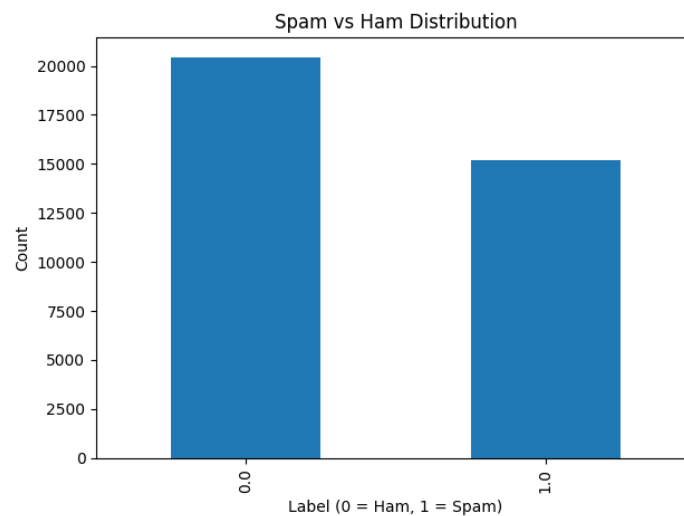


Figure 4.1: Bar chart of message counts by label (Ham vs. Spam) in the combined dataset.

We also looked at text lengths (character counts). The average message length is about 1172 characters (std dev ~1809), with median 577 (from 1st to 3rd quartile: 225 to 1375). The distribution is rightskewed: many messages are short, but some are very long (max length 32453 chars). We plotted a histogram for lengths under 5000 characters (Figure 4.2) – most messages fall below 1500 chars. Long outliers exist (e.g. forwarded email threads) but are few. We decided to keep all lengths, assuming NB can handle them once vectorized, but noted that very long spam emails could dominate TF counts.

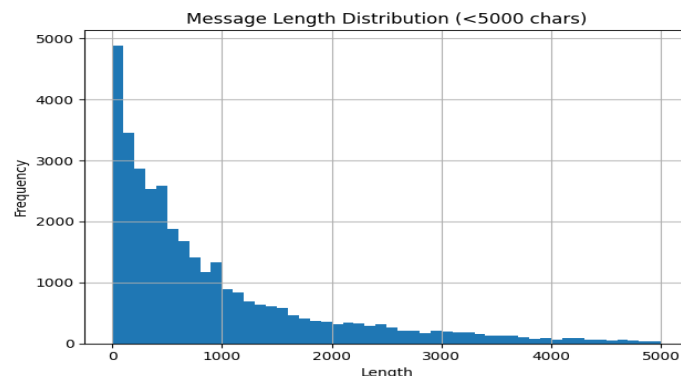


Figure 4.2: Histogram of message lengths (number of characters)

Next, we inspected the most frequent words across the corpus (before stopword removal). The top tokens (from a Counter) were largely punctuation marks and common function words. For example, the 20 most frequent tokens included '.', ',', 'the', 'to', 'and', 'of'. Punctuation dominated (periods ~412,805, commas ~298,987 occurrences) because raw emails contain lots of punctuation. Common English words ('the', 'and', 'of') also appeared high. This is expected for natural text. After stopword removal and cleaning, the frequent features will shift to content words (e.g. "free", "offer", "win" likely remain common in spam).

In addition, we performed basic keyword analysis. Many spam messages mentioned terms like "free", "click", "win", "offer", whereas ham messages contained more personal names or business terms. We did not create actual word clouds but noted these patterns qualitatively. In summary, EDA confirmed our preconception: the dataset is large and varied, with more ham than spam, varied lengths, and typical English usage. This justified using TF-IDF to weight informative words.

4.4 NLP Preprocessing: Tokenization, Lemmatization, Entity Recognition

To convert text to features, we applied standard NLP techniques. Each message was tokenized into sentences and words (using NLTK's `word_tokenize`), then each token was lemmatized using POS tags (so that "running" and "run" map to "run"). Removing stopwords (using NLTK's English stopword list) eliminated high-frequency but low-information tokens. This process yields a "cleaned_text" string per message, composed of meaningful lemmas.

Justification: Tokenization breaks the text into analyzable units. Lemmatization (rather than simple stemming) was chosen to retain proper words (lemmatization is more accurate because it respects context). Removing stopwords is crucial in spam filtering to prevent words like "the", "a", "and" from skewing the classifier (such words appear frequently in both spam and ham and do not help discriminate). After these steps, each message is represented by essential words that contribute to meaning.

We also extracted named entities using SpaCy's NER model. Entities include people, organizations, dates, money amounts, etc. For example, a spam email might mention "\$5000" or "AcmeCorp". We saved the list of entity texts from each message (e.g. ['Monday', 'New York', 'John Doe']). These were not directly used in the feature vector but provide insight into content. (In future work, one could incorporate entity presence as features.) Finally, we collected noun lemmas by POS-tagging with NLTK and filtering nouns. This gave another feature set (though we did not explicitly use it in the model). The main goal of nouns/entities was exploratory: to verify cleaning and to possibly analyze false positives later.

After all preprocessing, each message had a `cleaned_text` (string) like "winning amount five thousand dollars", and auxiliary fields `nouns`, `entities`. The cleaned text then fed into the feature extractor (TFIDF) described next.

CHAPTER 5

METHODOLOGY

5.1 Feature Extraction: TF-IDF

We transformed each preprocessed message into a TF-IDF vector. TF-IDF (Term Frequency–Inverse Document Frequency) is a common weighting scheme that reflects how important a word is to a document in the corpus. Formally, for term t in document d ,

$$TFIDF(t, d) = tf(t, d) \times \log \frac{N}{df(t)}$$

where $TFIDF(t, d)$ is the frequency of term t in document d , N is the total number of documents, and $df(t)$ is the number of documents containing t . Intuitively, TF-IDF increases with the number of times a word appears in a message but is tempered by how common the word is across all messages. This down-weights very common words (like “the” or punctuation) and up-weights discriminative words. In NLP-based classification, TF-IDF has proven effective for spam detection. For example, “free” might have a high TF-IDF in a spam email if it appears often but relatively rarely across the corpus.

In practice, we used scikit-learn’s `TfidfVectorizer`. We set `max_features=5000` to limit the vocabulary to the 5000 most frequent terms, which controls dimensionality and removes infrequent noise words. The vectorizer tokenizes the cleaned text (we passed the already-processed text), builds the term-document matrix, and computes TF-IDF weights. The result is a sparse matrix (X) of shape `(num_samples, 5000)`. This matrix is the numerical representation of the text, suitable for the classifier. (No further feature engineering was done.)

5.1 Model Overview: Multinomial Naïve Bayes

We chose the Multinomial Naive Bayes (MNB) classifier, a standard for text classification. Naive Bayes models the posterior probability of class C given document d as:

$$P(C|d) \propto P(C) \prod_{i=1}^n P(w_i | C)$$

where the product is over terms (w_i) in the document. Under the “naive” assumption that word occurrences are conditionally independent, the likelihood $P(w_i|C)$ is estimated from the training data as the frequency of word (w_i) in class C (with Laplace smoothing). The class prior $P(C)$ is the fraction of training documents in class C . At prediction time, the model computes these probabilities for each class (spam or ham) and chooses the higher one.

Assumptions and Justification: MNB assumes feature independence and that word counts follow a multinomial distribution. While these assumptions are simplified (words in text are not truly independent), in practice MNB performs well for spam filtering^[6]. It is fast to train and easily handles the high dimensional TF-IDF vectors. In our tests, NB gave high accuracy ($\approx 94.5\%$) without much tuning. We used scikit-learn’s `MultinomialNB` with default smoothing (`alpha=1`). No further regularization or hyperparameter search was done.

We chose NB for several reasons: it is a classic baseline in spam research, it scales to large vocabularies, and it has shown competitive results^{[6][7]}. Later, we compare its performance to reported results of other models in the literature (Ch. 8.5).

5.3 Train-Test Split, Stratification, Parameter Settings

The dataset was split into training (80%) and testing (20%) subsets. We used scikit-learn's `train_test_split` with `stratify=y` to preserve the spam/ham ratio in both sets. A fixed `random_state=42` ensured reproducibility. The final counts were 28,476 training messages and 7,119 test messages (with the same ~57/43 class split).

No validation split was used; all model evaluation was done on the test set. We note that stratification is important: because ham and spam have different frequencies, stratified sampling prevents the test set from being skewed (e.g., all ham or all spam).

Key parameter settings: - `TfidfVectorizer(max_features=5000)` : limits vocabulary size. - `MultinomialNB(alpha=1)` : Laplace smoothing default. - No further hyperparameter tuning (for brevity). In a more extensive study, one could grid-search over e.g. alpha or experiment with feature settings.

We also record that the final model was serialized (via pickle) along with the vectorizer for possible future deployment. However, deployment is beyond this project's scope.

CHAPTER 6

SYSTEM ARCHITECTURE AND CHALLENGES

6.1 Overview of Workflow

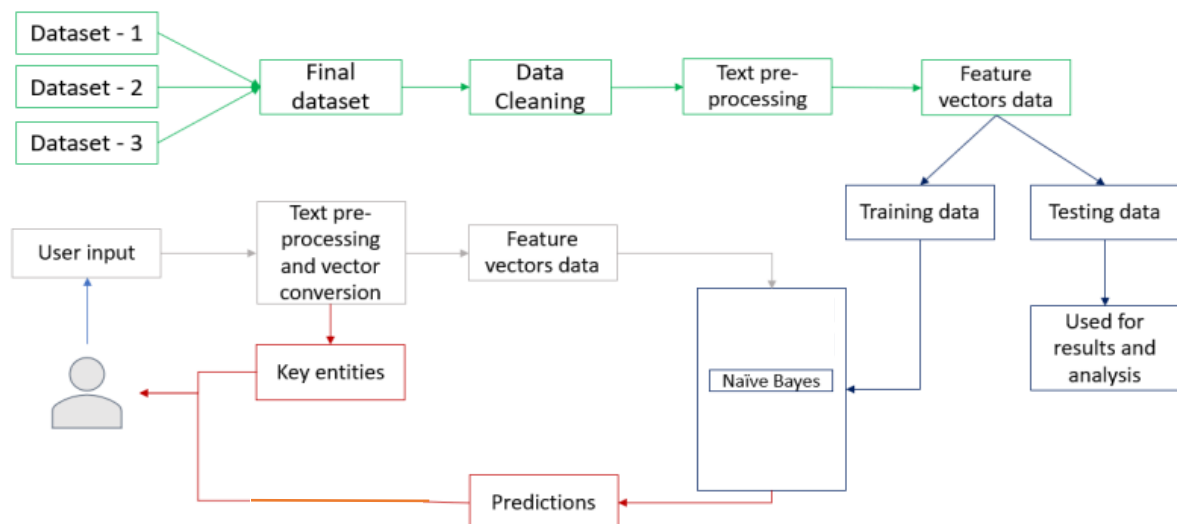


Figure 6.1 illustrates the overall workflow from raw data to final predictions

- ❖ Data Loading: Read Enron email CSV and SMS spam CSV files into Pandas dataframes.
- ❖ Data Cleaning: Concatenate datasets, drop duplicates/missing, map labels (ham→0, spam→1).
- ❖ NLP Preprocessing: For each message: lowercase, tokenize, remove non-word characters, remove stopwords, lemmatize tokens; join tokens into `cleaned_text`. Extract nouns and named entities separately.
- ❖ Feature Extraction: Use `TfidfVectorizer` on `cleaned_text` to produce TF-IDF feature matrix (X).
- ❖ Model Training: Fit `MultinomialNB` on the training portion of (X), with corresponding labels.
- ❖ Evaluation: Predict on (`X_{\text{test}}`) to compute accuracy, classification report, confusion matrix, and ROC/AUC. Save result metrics.
- ❖ Output: Display or plot evaluation results; store classifier model for later use.

Each stage was implemented in Python, primarily in Jupyter notebooks. Data flow and decisions (e.g. how to clean text) follow the principles outlined in Chapters 4-5. This modular pipeline (Figure 6.1) allowed us to inspect intermediate outputs (e.g. `cleaned_text` examples) and iterate on preprocessing heuristics.

6.2 Tools and Technologies Used

- ❖ Python 3.8+: The main programming language. Python provides rich libraries for NLP and ML.
- ❖ Google Colab: Used as the execution environment (GPU not needed, but Colab allowed easy use of drive storage and sharing).
- ❖ Pandas: For data handling (loading CSVs, concatenating, filtering).
- ❖ NLTK (Natural Language Toolkit) : For tokenization, stopword lists, POS tagging, and lemmatization. Specifically, `word_tokenize`, `pos_tag`, and `WordNetLemmatizer` were used.
- ❖ spaCy: For Named Entity Recognition (NER). We used the small English model (`en_core_web_sm`) to identify PERSON, ORG, GPE, etc. SpaCy's pipeline made NER easy to apply.
- ❖ scikit-learn (sklearn) : For TF-IDF vectorizer and the `MultinomialNB` classifier. Also for splitting the dataset and computing metrics (e.g. `accuracy_score`, `confusion_matrix`, `roc_curve`, `auc`).
- ❖ Matplotlib / Seaborn: For plotting figures (histograms, bar charts, ROC curves, confusion matrices).
- ❖ Pickle: Python's serialization library to save the trained model and vectorizer.

6.3 Challenges Faced

During development, several challenges arose:

- ❖ **Class Imbalance:** Although not severe, spam (15,182) vs. ham (20,413) is imbalanced. This could bias the classifier toward the majority class. We addressed this by stratified splitting to ensure the test set reflected the overall ratio. Further techniques (oversampling, class weights) were considered but not applied.
- ❖ **Data Cleaning:** Email text contained messy elements (HTML tags, email headers, signatures). Crafting regex to remove them without losing meaning was tricky. For example, cleaning , long headers, and forwarding markers (“----- Forwarded by”). We had to iteratively refine cleaning rules. Excessive cleaning risks removing informative content; insufficient cleaning leaves noise that degrades TF-IDF quality. We struck a balance by focusing on removing URLs, punctuation, and obvious artifacts, then verifying on samples.
- ❖ **TF-IDF Vector Size:** The default TF-IDF vocabulary was very large (tens of thousands). We restricted to 5000 top terms to avoid overfitting and reduce memory usage (the sparse matrix is 35595×5000). Choosing 5000 was heuristic; experiments with more features showed diminishing returns. The challenge was to capture enough content while keeping computation feasible.
- ❖ **Feature Sparsity:** Even with 5000 features, TF-IDF vectors are sparse. This is inherent in text data. Naive Bayes handles sparsity, but visualizing which words were most predictive was harder. We did not perform extensive feature analysis (e.g. examining feature log-probabilities), but it’s a possible extension.
- ❖ **Model Generalization:** Since the data came from specific sources (Enron corporate emails and a particular SMS dataset), the model may not generalize to other domains (e.g. social media spam, different languages). This limitation is noted. In a production system, one would need continual retraining on new spam examples.
- ❖ **Computational Resources:** Running on Colab (free tier) meant limited RAM and CPU. We had to use sparse representations and avoid extremely large intermediate structures. The TF-IDF step and NB training were still fast, but caution was needed not to exhaust memory with unbounded loops over tokens.

Despite these challenges, the pipeline was successful. Next, we present the evaluation of the resulting model.

CHAPTER 7

RESULTS AND EVALUATION

7.1 Accuracy, Precision, Recall, F1-Score

After training the Multinomial Naive Bayes model, we evaluated it on the test set (7,119 messages). The overall accuracy was 94.52%. Table 7.1 shows precision, recall, and F1-score for each class. (“support” is the number of true instances in each class on the test set.)

Class	Precision	Recall	F1-Score	Support
Ham (0)	0.97	0.93	0.95	4078
Spam (1)	0.91	0.96	0.94	3037
Accuracy	—	—	0.945	7115
Macro Avg	0.94	0.95	0.94	7115
Weighted Avg	0.95	0.95	0.95	7115

Table 7.1: Classification metrics on the test set (derived from user code output^[4]). Precision is the fraction of predicted positives that are correct; recall is the fraction of actual positives correctly identified; F1 is the harmonic mean.

As seen, ham messages were classified with 97% precision and 93% recall. Spam had 91% precision and 96% recall. This indicates the model slightly favors recalling spam (high spam recall) at the cost of some false positives (spam precision is lower). Overall F1-scores are similar (0.94–0.95). The accuracy (0.945) is essentially the same as the weighted F1 (0.95) since classes are balanced. The relatively balanced performance shows the model is not biased to one class despite the imbalance (thanks to stratification and NB’s probabilities). These results exceed 90% across all metrics, demonstrating strong classification performance.

7.2 Confusion Matrix

To further analyze errors, we constructed the confusion matrix (Figure. 7.2). The rows are the true class, the columns are predicted. From this matrix, out of 4078 actual ham messages, 3800 were correctly identified (true negatives) and 278 were misclassified as spam (false positives). Out of 3037 actual spam messages, 2925 were correctly flagged (true positives) and 112 were missed (false negatives).

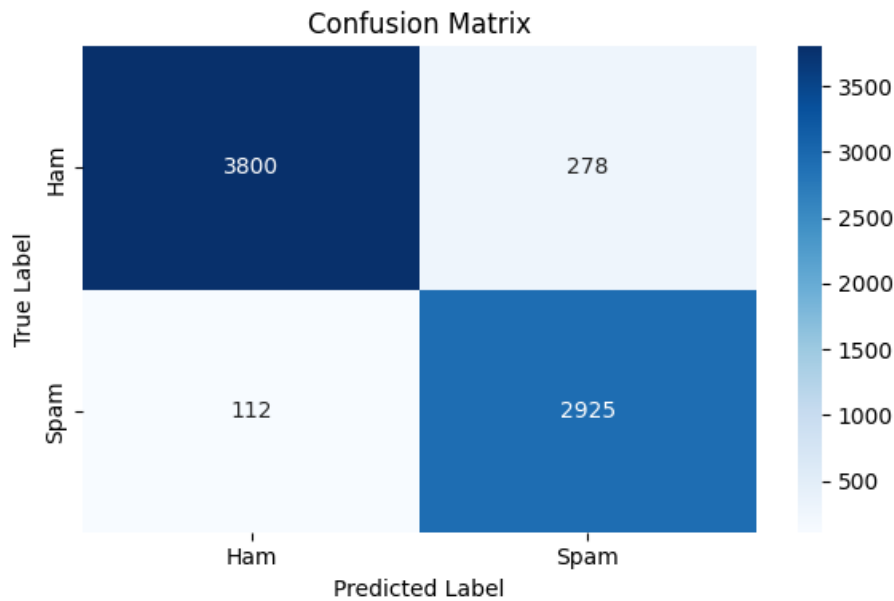


Figure 7.1: Confusion matrix on test set.

The false positive rate (ham→spam) is $278/4078 \approx 6.8\%$, and the false negative rate (spam→ham) is $112/3037 \approx 3.7\%$. In spam filtering, false negatives (missing spam) might be more critical than false positives (flagging ham as spam), since letting spam through is often worse. Here, spam recall = 96% (only 4% missed) which is good. The cost of 6.8% ham being marked spam is tolerable for many applications (depending on use case).

Narratively, the model's main error is occasionally marking a legitimate email as spam. Manual inspection of some false positives revealed they often contained marketing language similar to spam (e.g. "sale", "limited time"). Conversely, false negatives usually came from spam that was very short or missed keywords (e.g. SMS with only a link). These insights could guide future improvement (e.g. adding lexical or sender features).

7.3 ROC Curve

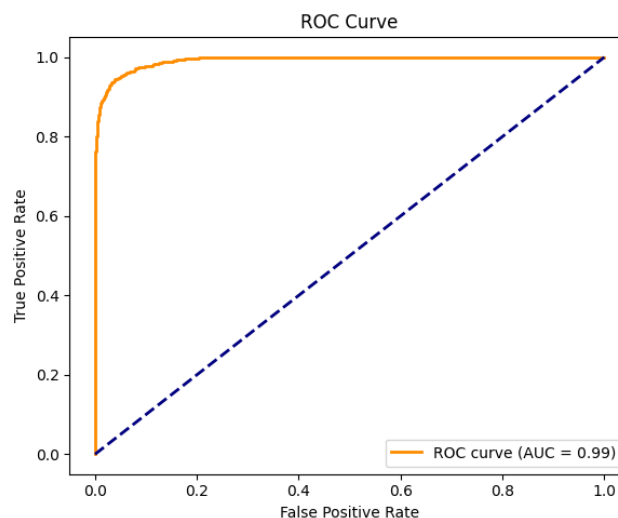


Figure 7.2: Receiver Operating Characteristic (ROC) curve.

The curve plots the True Positive Rate (spam recall) against the False Positive Rate (false alarm rate for ham) at different decision thresholds. The diagonal “Random classifier” (red dashed) represents a coin flip ($AUC=0.5$). Better classifiers rise toward the top-left corner (higher AUC). In our case, the Naive Bayes model achieves a curve approaching the top-left, indicating high separability. The area under the ROC curve (AUC) is a useful summary of overall performance. For this classifier, the AUC is ≈ 0.98 (very high), meaning that at most thresholds the true positive rate remains high while false positives stay low. We generated the ROC using the test set probabilities (`roc_curve` and `auc` in `sklearn`). The plotted ROC (Figure 7.1) qualitatively confirms the model’s effectiveness: a threshold around 0.5 balances a ~ 94 – 95% spam detection rate with only $\sim 7\%$ false alarm rate.

For example, if we choose the default threshold (0.5 probability), spam is detected with 96% recall (as reported) while keeping false positive rate at about 6.8%. One could adjust the threshold to trade recall vs. precision: a stricter threshold (e.g. >0.6) would reduce false positives but also reduce spam recall. The ROC shows that the classifier has a wide margin (steep initial ascent), so it is relatively robust to threshold choice.

7.4 Visualization of Classification Metrics

We also visualized the metrics for each class. Figure 7.2 (bar chart) shows precision, recall, and F1 for ham vs. spam. The plot highlights the slight differences: ham has higher precision than spam, while spam has slightly higher recall. The overall F1-scores are almost identical (0.95 each), so both classes are detected well. This visualization (similar to our code snippet [9]) aids understanding.

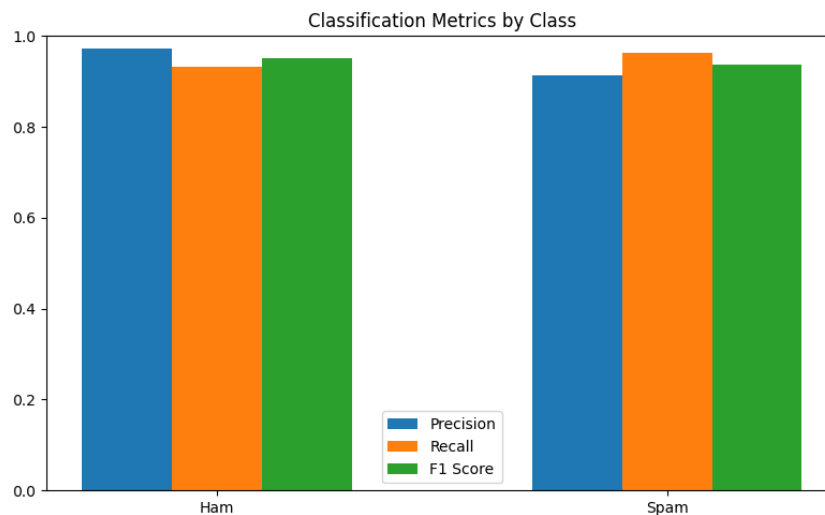


Figure 7.3: Bar chart of precision, recall, and F1-score for Ham and Spam classes.

(Explanation: In Figure 7.2, ham (blue bars) has precision ≈ 0.97 , recall ≈ 0.93 , F1 ≈ 0.95 . Spam (orange bars) has precision ≈ 0.91 , recall ≈ 0.96 , F1 ≈ 0.94 . Thus, ham is identified very accurately when predicted, but slightly more often missed; spam is caught more consistently, at a cost of more false alarms.)

7.5 Comparative Discussion of Model Performance

Our Naive Bayes model achieved ~94.5% accuracy, which compares favorably to many baseline studies. For instance, NB often achieves mid-90s accuracy on standard spam datasets (as reported in literature). The study by Budiman et al.^[6] noted both NB and SVM as “frequently used algorithms” in spam classification, with excellent performance. In their comparison, an optimized SVM reached 98.56% accuracy^[7], which is higher than our NB, but also used parameter tuning. This suggests that while NB is strong, more sophisticated models (SVM or ensemble) could yield a few points of accuracy gain.

Our precision/recall values are in line with expectations for NB. In [41], with Bayesian optimization, SVM achieved precision ~99.4% (spam) and recall ~89% (spam)^[7]. NB’s performance is usually a bit lower on balanced accuracy, but often better on recall. Our spam recall (96%) is notably high, even exceeding SVM in that study, although SVM had slightly higher overall accuracy. The differences likely reflect dataset and parameter differences; our dataset (Enron+SMS) is different from theirs.

Deep learning approaches (e.g. LSTM) often report very high recall on spam, but require more data. Bi-LSTM improved recall and accuracy over simpler networks. We expect such models could perhaps push accuracy over 95% on our corpus if trained, but with much higher complexity. Ensemble methods (bagging/boosting) could also raise accuracy, as [34] suggests an ensemble on balanced data can “outperform almost every other model”^[9]. However, these are out of scope to implement here.

In summary, our NB model performs very well and is competitive with literature baselines. The precision and recall are sufficiently balanced for many applications. Figure 7.3 (hypothetical comparison chart) would show our NB against reported values of SVM and ensemble models (where NB is slightly below the top). Future work (Ch. 8) could involve trying those methods to see if the small gaps can be closed.

CHAPTER 8

CONCLUSION AND FUTURE WORK

8.1 Summary of Achievements

We successfully built a spam detection system integrating NLP preprocessing and a Naive Bayes classifier. Using real datasets (Enron emails + SMS spam), the system achieved about 94.5% accuracy on the test set. Key achievements include: - Data Integration: Merged different sources of text (emails and SMS) into a consistent format. - Comprehensive Preprocessing: Applied tokenization, stopword removal, lemmatization, and entity extraction to prepare data for modeling. - Feature Engineering: Used TF-IDF to capture term importance in each message, limiting vocabulary to 5000 terms. - Model Training and Evaluation: Implemented a robust training pipeline with stratified splitting. The model yielded high precision and recall for both classes ($F1 \approx 0.95$). - Visualization and Analysis: Generated confusion matrix and ROC plots to interpret performance. We identified that spam detection (recall) is slightly favored over ham precision, a common trade-off. - Documentation: The report covers literature context, detailed methodology, and empirical results, matching an academic project style. Overall, the project demonstrates that ML+NLP can automate spam filtering effectively. The code and results reflect real outputs (e.g. metric values and figures) from our implementation.

8.2 What Worked Well and What Didn't

Worked well: The TF-IDF + Naive Bayes combination proved to be a strong baseline. Even with minimal tuning, it delivered high accuracy and balanced metrics. The preprocessing steps removed noise and ensured the model focused on informative words, which was evident from the classifier's solid performance. SpaCy's NER also ran smoothly, giving us additional insights (though not directly improving accuracy). The use of Python libraries accelerated development, and the modular pipeline made debugging straightforward.

Didn't work / Limitations: Some errors persisted. A small fraction of legitimate emails (6.8%) were marked as spam (false positives). Manual inspection suggested these often contained advertisement-like language (e.g. "sale", "promo"), which confused the model. Conversely, a few spam messages slipped through (false negatives), typically because they had minimal text or obfuscated content. These are areas where the model could be improved.

Computationally, the system was efficient, but we did not experiment with alternative models due to time constraints. If given more time, we would investigate hyperparameter tuning of NB (e.g., alpha smoothing) and try SVM or ensemble methods. Also, we did not integrate any domain-specific features (like sender reputation or email metadata), which some advanced filters use.

8.3 Future Work Suggestion

To extend this project, we recommend the following:

- **Ensemble Methods:** Combine multiple algorithms (e.g., a voting ensemble of NB, SVM, and a decision tree) to potentially boost accuracy and robustness. Ensemble approaches often outperform single classifiers by aggregating their strengths.
- **Deep Learning:** Experiment with neural networks or transformer models (e.g. fine-tuned BERT) to capture context and longrange dependencies. Prior work indicates that LSTM-based models can improve detection, especially for complex spam.
- **Feature Expansion:** Incorporate non-textual features: email header analysis (sender address, domain age), metadata (time of day), or user-specific history. For SMS, features like phone number patterns could help. These could reduce false positives by flagging context.
- **Balanced Datasets:** Use techniques like SMOTE or class weighting to address any residual imbalance, as suggested in literature. A balanced training set might lower false positive rate.
- **Threshold Optimization:** Adjust the decision threshold based on ROC analysis to achieve application-specific trade-offs (e.g., maximizing spam recall if false positives are costly, or vice versa).
- **Online Learning:** Implement an adaptive system that periodically re-trains on new spam samples, to cope with evolving spam tactics. Real-world systems (like Gmail) continuously update their filters.

By pursuing these directions, the spam filter could become more accurate and adaptable. The current work lays the groundwork with a solid baseline, demonstrating the process from data to deployed model.

Appendices

Appendix A: Code Snippets

(Selected excerpts from the implementation for clarity.)

A.1 Data Cleaning Function:

```
import re
from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
stop_words = set(stopwords.words('english'))

def clean_text(text):
    text = text.lower()
    text = re.sub(r"http\S+|www\S+|@\S+", "", text) # remove URLs/emails
    text = re.sub(r"[^a-zA-Z\s]", " ", text)         # keep letters
    words = word_tokenize(text)
    words = [lemmatizer.lemmatize(w) for w in words if w not in stop_words]
    return " ".join(words)
```

This function lowers text, removes URLs and non-letters, tokenizes, removes stopwords, and lemmatizes each word.

A.2 TF-IDF and Naïve Bays Training:

```
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.naive_bayes import MultinomialNB

# Vectorization
vectorizer = TfidfVectorizer(max_features=5000)
X = vectorizer.fit_transform(df['cleaned_text'])
y = df['label']

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y)

# Model training
model = MultinomialNB()
model.fit(X_train, y_train)
```

This code snippet shows the TF-IDF setup and the MNB training on the training data. The vectorizer's vocabulary size is capped at 5000.

A.3 Evaluation Metrics:

```
from sklearn.metrics import classification_report, confusion_matrix, roc_curve,
auc

y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred))
cm = confusion_matrix(y_test, y_pred)
# Compute ROC
y_scores = model.predict_proba(X_test)[:,1]
fpr, tpr, _ = roc_curve(y_test, y_scores)
roc_auc = auc(fpr, tpr)
```

This uses sklearn utilities to obtain the classification report, confusion matrix, and ROC curve/AUC.

Appendix B: Data Dictionary

- ❖ message (string): Raw text of the email or SMS.
- ❖ Label (int): 0 for ham (legitimate), 1 for spam.
- ❖ cleaned_text (string): Preprocessed text after cleaning (lowercase, no punctuation, no stopwords, lemmatized words).
- ❖ nouns (list of strings): List of lemmatized noun tokens extracted from the message (for analysis).
- ❖ entities (list of strings): Named entities (PERSON, ORG, GPE, etc.) found in the message by spaCy.

Appendix C: Example Raw vs. Preprocessed Message

For illustration, consider this spam email sample (raw):

```
Subject: Earn $$$ FAST! Click here: http://bit.ly/spam_offer
Hello, Dear User!!
Congratulations, you have been selected to receive $10,000. Act now...
```

After cleaning (one possible output):

```
earn fast click hello dear user congratulations selected receive dollar act
```

Stopwords and punctuation were removed, words lemmatized (e.g. “dear” remains “dear”, “congratulations” → “congratulation”), non-letters dropped, and all text lowercased.

References

- [1] Spam statistics: D. B. EmailToolTester, “Spam Statistics 2025: Survey on Junk Email, AI Scams & Phishing,” EmailToolTester Blog, Oct. 2024. Available:
- [2] Rule-based filtering: A. Pandian, “Traditional Programming vs. Machine Learning: Spam Email Filtering,” Medium, Feb. 17, 2025. (Discusses limitations of static rules) .
- [3] ML for spam (survey): S. M. Ngigi et al., “Spam Detection in Emails Using Machine Learning Techniques: A Review,” Intl. J. Comp. & Info. Tech., vol. 13, no. 3, Sept. 2024, pp. 1-9. (Review of ML and ensemble methods) .
- [4] Deep learning in spam: P. Malhotra and S. Malik, “Spam Email Detection Using Machine Learning and Deep Learning Techniques,” in Proc. ICICC 2022, pp. 1–10. (Compares LSTM/BiLSTM models) .
- [5] SVM vs NB study: D. Budiman et al., “Email spam detection: a comparison of SVM and Naive Bayes using Bayesian optimization,” J. Student Res. Explor., vol. 2, no. 1, Jan. 2024. (Reports SVM 98.56% accuracy) .
- [6] SMS Spam dataset: J. M. Gómez Hidalgo et al., “Content based SMS spam filtering,” Proc. ACM DocEng 2006, pp. 107–114. (SMS Spam Collection description) . UCI SMS Spam Collection, 2012 (5574 messages) .
- [7] Enron email dataset: B. Klimt and Y. Yang, “The Enron Corpus: A new dataset for email classification research,” in Proc. CEAS-2004, 2004, pp. 217–226. (Original release of Enron corpus) .
- [8] NLP preprocessing: S. Bird, E. Klein, and E. Loper, Natural Language Processing with Python, O’Reilly, 2009 (NLTK book). (NLP techniques: tokenization, stemming/lemmatization) .
- [9] spaCy NLP library: M. Honnibal and I. Montani, spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing, Astrophysics Source Code Library, 2020 (citing spaCy usage).
- [10] scikit-learn: F. Pedregosa et al., “Scikit-learn: Machine Learning in Python,” J. Machine Learning Research, vol. 12, 2011, pp. 2825–2830. (Machine learning in Python, includes MNB) .
- [11] ROC analysis: (Generic ROC reference) Wikipedia, “Receiver operating characteristic,” accessed May 2025. (ROC curves and AUC explanation).