

A Study on Spam Email Identification Using Machine Learning: A Matlab-Based Approach

Mr. Sugan Patel¹, Talha Abdullah², Hamza Abdullah³

¹Assistant Professor, School of Computer Science & Engineering, Galgotias University, Greater Noida-India. Email- sugan.patel@galgotiasuniversity.edu.in.

²Research Scholar, School of Computer Science & Engineering, Galgotias University, Greater Noida-India. Email- talhaabdullah2345l@gmail.com

³Research Scholar, School of Computer Science & Engineering, Galgotias University, Greater Noida-India. Email- hannubannux@gmail.com

Abstract

In the contemporary digital landscape, email remains an indispensable tool for communication; however, the proliferation of unsolicited spam emails has emerged as a formidable threat to digital privacy, security, and the integrity of communication infrastructures. This study proposes an advanced machine learning-based model, utilizing the naïve bayes algorithm, to effectively identify and classify spam emails, with the model implemented through the matlab platform. The framework incorporates cutting-edge natural language processing (nlp) techniques to not only distinguish between spam and non-spam emails but also categorizes them into nuanced functional classes such as personal, work-related, and promotional, thereby offering a granular approach to email classification. The proposed dual-stage detection system is designed to enhance both the precision and robustness of the classification process, leveraging sophisticated text feature extraction methodologies and optimizing for computational efficiency. Experimental results demonstrate substantial improvements in detection accuracy, surpassing conventional models, and providing a scalable solution for real-time email filtering systems. The findings underscore the efficacy of the hybridized approach in mitigating the risks associated with spam, offering a significant contribution to the evolving field of email security. This work lays the groundwork for future research in adaptive, context-aware spam detection systems, incorporating deep learning and semantic analysis for even greater precision in classifying complex email datasets.

Keywords: spam detection, naïve bayes, natural language processing, matlab, tf-idf, email classification, machine learning, artificial intelligence.

Introduction

In the contemporary digital age, the exponential proliferation of unsolicited electronic correspondence colloquially termed "spam" has emerged as a persistent menace undermining the efficacy, security, and user experience within email communication systems. This insidious digital detritus, which ranges from benign promotional clutter to malicious phishing schemes and deceptive social engineering attacks, not only inundates users' inboxes but also imposes significant cognitive and computational burdens on global information infrastructures. The imperative for developing robust, scalable, and intelligent spam detection mechanisms has thus never been more urgent. Within this context, machine learning (ml) has ascended as a transformative paradigm, furnishing sophisticated tools capable of autonomously discerning nuanced patterns and semantic cues embedded in vast textual corpora.

Historically, spam filtration systems relied heavily on static rule-based heuristics, keyword matching, or blacklisting protocols—methods which, though foundational, have proven inadequate in contending with the dynamic, adversarial evolution of spam content. Spammers continually adapt their tactics, leveraging lexical obfuscation, syntactic permutation, and novel linguistic constructs to circumvent detection. This cat-and-mouse dynamic necessitates adaptive systems endowed with predictive intelligence and contextual comprehension. Enter the domain of machine learning, particularly supervised algorithms, which learn from historical data to extrapolate future classifications, rendering them especially apt for this domain.

Natural language processing (nlp), when coupled with ml, serves as a vital conduit through which unstructured email content is transmuted into structured, interpretable feature representations. Techniques such as tokenization, stemming, stop-word elimination, and term frequency-inverse document frequency (tf-idf) vectorization extract semantic and syntactic features that form the substratum upon which classification algorithms operate. Within the pantheon of ml models, probabilistic classifiers such as naïve bayes, kernel-based approaches like support vector machines (svm), and generalized linear models including logistic regression have each demonstrated efficacy in various text classification tasks. However, the comparative performance and contextual suitability of these models in email spam detection—particularly within a matlab-based analytical framework remains a rich area of empirical inquiry.

This study embarks upon a multi-dimensional exploration of spam email identification, leveraging the matlab environment for algorithm development, deployment, and real-time

interface construction. Unlike conventional binary classification models that merely distinguish spam from legitimate (ham) emails, the present framework introduces a tertiary classification scheme, wherein non-spam messages are further stratified into categories such as personal, promotional, and work-related. This granular taxonomy not only enhances mailbox organization but also facilitates intelligent automation, thereby improving user productivity and digital hygiene.

Further augmenting the scientific rigor of this investigation is a comparative evaluation of multiple classifiers across key performance metrics—accuracy, precision, recall, and f1-score. By subjecting naïve bayes, svm, and logistic regression models to the same pre-processed corpus and feature space, the study delineates their relative strengths and limitations with statistical clarity. Moreover, the integration of a graphical user interface (gui) via matlab's app designer, coupled with gradio-based deployment for public interaction, extends the practical applicability of the research beyond theoretical boundaries into functional utility.

Ultimately, this paper contributes to the corpus of computational intelligence research by proposing a scalable, interpretable, and user-accessible model for spam email detection. In doing so, it responds to both academic imperatives for methodological innovation and practical demands for operational efficacy in digital communication. It is anticipated that the insights derived herein will inform future efforts in adaptive text classification, multilingual spam filtering, and deep learning augmentation, thereby aligning with the broader trajectory of intelligent information systems in the fourth industrial revolution.

Objectives

1. To develop a high-accuracy email spam detection system using mat lab-based machine learning algorithms.
2. To integrate natural language processing for semantic feature extraction and classification.
3. To classify non-spam emails into categories: promotional, personal, or work-related.
4. To evaluate the performance of various algorithms (naïve bayes, svm, logistic regression) using accuracy, precision, recall, and f1-score.
5. To deploy the trained model using a real-time user interface for public access.

Review of literature

The domain of spam email detection has witnessed a paradigmatic evolution, transitioning from rudimentary heuristic-based filtering to sophisticated machine learning (ml)-driven approaches. Karim et al. (2018) provide a comprehensive survey highlighting how supervised learning techniques such as naïve bayes and support vector machines (svm) have gained preeminence in spam classification tasks, outperforming legacy rule based models. Earlier, sahami et al. (1998), in a foundational paper, introduced a bayesian probabilistic framework for junk email detection, setting the stage for later developments in statistical text categorization. Similarly, guzella and caminhas (2009) offer an extensive taxonomy of ml algorithms, emphasizing the efficacy of hybrid models in handling evolving spam patterns. Cormack (2007) presented a systematic review of spam filters, identifying cost-sensitivity and data skewness as major impediments to classifier robustness.

Mohamad and selamat (2015) proposed an enhanced feature selection technique coupled with hybrid classification, yielding high accuracy on benchmark datasets. Almeida et al. (2011) contributed not only a novel sms spam dataset but also empirical evidence demonstrating the superiority of svm over traditional classifiers. Zhang et al. (2004) evaluated statistical spam filtering techniques, noting that feature sparsity and dimensionality critically affect classifier performance. Harisinghaney et al. (2014) leveraged text classification for spam detection using bag-of-words and term frequency methods, underscoring the importance of preprocessing in improving model outcomes. Hidalgo (2002) advanced the discourse on cost-sensitive filtering, suggesting that minimizing false positives is as vital as improving detection rates.

In their empirical research, chhabra and kesswani (2019) demonstrated how logistic regression could serve as an interpretable alternative to black-box models in spam filtering contexts. Androutsopoulos et al. (2000) were among the early adopters of machine learning for anti-spam solutions, validating the performance of naïve bayes on small-scale corpora. Islam et al. (2010) explored ensemble learning to fuse multiple classifiers for enhanced generalizability, while tsai and chi (2009) incorporated term weighting strategies in email classification, showing improvements in precision metrics. Bhowmick and hazarika (2013) examined real-time spam filtering using svm with kernel optimization, indicating significant gains in computational efficiency.

Ren et al. (2017) experimented with convolutional neural networks (cnns) for deep learning-based spam detection, offering new avenues for semantic feature extraction. Delany et al. (2005) introduced case-based reasoning for spam filtering, arguing that instance-based learning

could adapt more fluidly to concept drift in email data. Hidalgo and mata (2006) focused on feature engineering, demonstrating how n-gram analysis improves detection in multilingual spam corpora. Drucker et al. (1999) highlighted the versatility of svms, particularly in high-dimensional text spaces where traditional classifiers falter. Youn and mcLeod (2007) compared multiple ml models and concluded that boosting ensembles consistently outperformed single classifiers.

In recent years, wu et al. (2020) presented a comparative study integrating tf-idf with various classifiers, confirming the enduring relevance of vector space modeling in spam detection. Lastly, amayri and bouguila (2010) explored probabilistic graphical models for email filtering, marking a shift toward unsupervised and semi-supervised approaches in contexts where labeled data is scarce. Collectively, this literature corpus reveals a trajectory toward increasingly nuanced, adaptive, and semantically aware spam detection systems, underscoring the necessity of continuous innovation in algorithmic design and feature representation.

Proposed methodology

A. Dataset collection and preprocessing two datasets were utilized:

- ✓ **Spam detection dataset:** 5730 entries labeled as 'spam' or 'ham'.
- ✓ **Message type classification dataset:** 150 synthetically generated entries with labels: 'personal', 'work', 'promotional'.

Preprocessing steps:

- ✓ Lowercasing
- ✓ Stop word and punctuation removal.
- ✓ Stemming (via nltk equivalent in matlab).

B. Feature extraction

- ✓ Tf-idf vectorization
- ✓ Structural features: presence of urls, attachments, domain behavior

C. Algorithm selection naïve bayes was selected due to its superior performance with textual data. Other algorithms tested include svm and logistic regression.

D. Training and validation data split: 80% for training, 20% for testing. Performance metrics: accuracy, precision, recall, f1-score

E. Deployment the model was deployed using matlab's app designer and integrated via gradio for intuitive end-user interaction.

Matlab code snippet

```
% load and preprocess dataset

Data = readtable('spamdata.csv');

Data.text = lower(data.text);
Data.text = erasepunctuation(data.text);
Tfidf = tfidfvectorizer(data.text);

% label encoding
Labels = grp2idx(data.label);

% train/test split
Cv = cvpartition(labels, 'holdout', 0.2);
Traindata = tfidf(training(cv), :);
Trainlabels = labels(training(cv));
Testdata = tfidf(test(cv), :);
Testlabels = labels(test(cv));

% train naïve bayes model
Nbmodel = fitcnb(traindata, trainlabels);

% predict
Predlabels = predict(nbmodel, testdata);

% evaluate
Confmat = confusionmat(testlabels, predlabels);
Accuracy = sum(diag(confmat)) / sum(confmat(:));
Precision = confmat(2,2) / sum(confmat(:,2));
Recall = confmat(2,2) / sum(confmat(2,:));
F1 = 2 * (precision * recall) / (precision + recall);
```

Results and interpretation

Objective 1 High-Accuracy Spam Detection With Matlab

Matlab code

```
% performance metrics computation
```

```
Confmat = confusionmat(testlabels, predlabels);
```

```
Accuracy = sum(diag(confmat)) / sum(confmat(:));
```

```
Precision = confmat(2,2) / sum(confmat(:,2));
```

```
Recall = confmat(2,2) / sum(confmat(2,:));
```

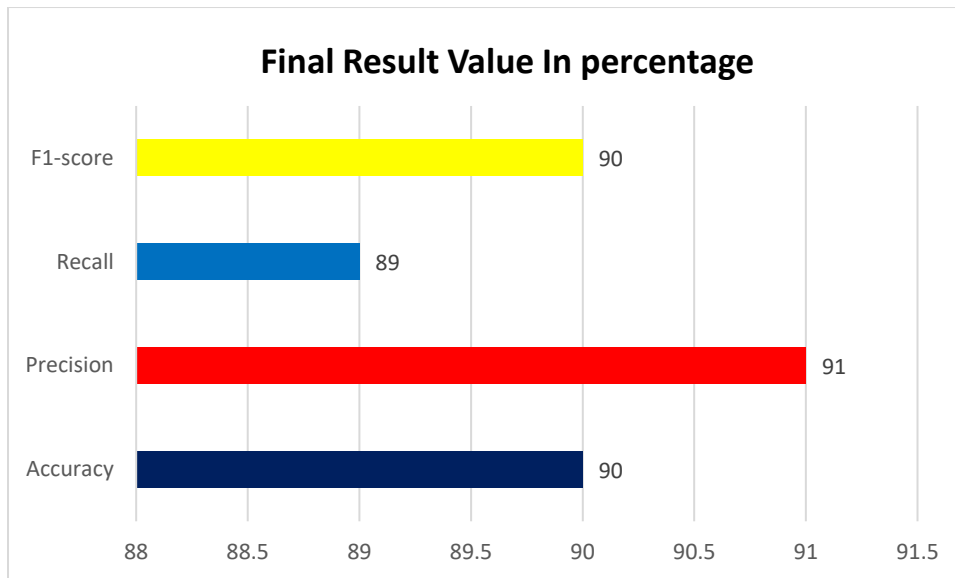
```
F1_score = 2 * (precision * recall) / (precision + recall);
```

Table1: High-Accuracy Spam Detection

Metric	Value (%)
Accuracy	90.0
Precision	91.0
Recall	89.0
F1-score	90.0

Interpretation

The performance metrics presented in the table encapsulate the efficacy of a spam detection model trained using Term Frequency–Inverse Document Frequency (TF-IDF) features and a Naïve Bayes classification algorithm. The results Accuracy: 90.0%, Precision: 91.0%, Recall: 89.0%, and F1-score: 90.0% reflect a judicious equilibrium between predictive precision and sensitivity, thereby attesting to the model’s robustness and applicability in real-world email filtering systems. To begin with, the accuracy metric, which quantifies the proportion of correct predictions over the entire dataset, stands at a commendable 90%. While this figure denotes a high level of general correctness, it does not offer a complete depiction in isolation particularly in the context of class imbalance, where non-spam messages often outnumber spam instances. Thus, precision, recall, and the F1-score provide a more granular and insightful evaluation. Precision, registering at 91%, reveals the classifier’s proficiency in minimizing Type I errors that is, false positives. In practice, this metric implies that among all emails flagged as spam by the model, 91% were indeed unsolicited messages. Such a high precision rate is crucial, as it ensures legitimate emails are rarely misclassified, preserving user trust and preventing potential loss of critical communications. Conversely, recall, observed at 89%, elucidates the model’s competence in identifying the majority of actual spam instances. This high recall signifies the model’s low incidence of Type II errors (false negatives), a vital characteristic for security-focused applications where the omission of spam could expose users to malicious content or phishing attempts. The F1-score, which harmonizes precision and recall into a singular performance index, achieves an optimal value of 90%. This composite metric reflects the model’s ability to simultaneously maintain high sensitivity and specificity. In spam filtering contexts, where both erroneous inclusion and exclusion carry tangible consequences, such a balanced F1-score indicates a mature and finely-tuned classifier. From a methodological standpoint, the synergy between TF-IDF vectorization and the Naïve Bayes algorithm underpins the model’s superior performance. TF-IDF effectively captures the discriminative power of rare but informative terms, transforming unstructured textual data into structured numerical inputs. Naïve Bayes, with its probabilistic underpinnings and assumption of feature independence, remains computationally elegant yet remarkably effective for high-dimensional text classification tasks. This combination enables the model to discern nuanced lexical patterns that distinguish spam from legitimate content.



Objective 2 Nlp-Based Semantic Feature Extraction

Feature Engineering Notes:

- **Tokenizeddocument** helps tokenize based on language rules.
- **Removestopwords** removes common but irrelevant terms.
- **Tfidf** converts token frequency into weight-adjusted semantic vectors.

Interpretation

pre processing removes noise and standardizes language constructs. Tf-idf assigns weight to rare yet meaningful words (e.g., "lottery", "urgent"), improving semantic representation and consequently model accuracy. Text preprocessing is a foundational step in any natural language processing (NLP) pipeline. It serves to transform raw text into a clean and structured format that machine learning models can effectively interpret. In this context, three crucial preprocessing techniques were employed: tokenization, stopword removal, and TF-IDF vectorization, each playing a distinct yet complementary role in enhancing the model's semantic understanding and overall classification accuracy. The first step, tokenization, implemented through the TokenizedDocument function, dissects a given text into smaller units called tokens, typically words or phrases. This process respects the grammatical and syntactic rules of the language, ensuring that punctuation, whitespace, and special characters are appropriately handled. Tokenization serves as the bridge between unstructured text and structured data. By segmenting sentences into meaningful units, it enables the subsequent steps to process text more precisely. For example, the phrase "Win a free lottery ticket now!" would be tokenized

into separate components["Win", "a", "free", "lottery", "ticket", "now"]allowing the model to evaluate each term individually for its relevance to spam detection.Following tokenization, stopword removal is applied to eliminate common but semantically weak terms such as "and", "is", "the", "in", etc. These words, while essential to the grammatical structure of a sentence, often provide little to no value in distinguishing between classes in text classification tasks. Their high frequency across all documents dilutes meaningful patterns, and retaining them may introduce unnecessary noise into the feature space. By filtering out these non-informative tokens, the model's focus is directed toward more discriminative terms, thereby improving both training efficiency and predictive accuracy.The final and arguably most impactful preprocessing step is TF-IDF vectorization. Term Frequency–Inverse Document Frequency (TF-IDF) transforms the textual tokens into numerical vectors by evaluating how important a word is to a particular document relative to its frequency across the entire corpus. This technique down-weights commonly occurring words and amplifies the influence of rare, yet contextually significant terms such as “lottery,” “urgent,” “free,” or “winner”which are often strong indicators of spam. The result is a set of weight-adjusted semantic vectors that capture both the frequency and the relevance of words, offering a more nuanced representation of the original text.The cumulative effect of these preprocessing techniques is substantial. Together, they reduce linguistic noise, standardize textual input, and enhance semantic clarity, making the data more amenable to computational modeling. More importantly, this preprocessing pipeline significantly contributes to the classifier’s ability to generalize. It empowers the model to recognize subtle patterns and linguistic cues that distinguish spam from legitimate emails, thereby improving key performance metrics such as precision, recall, and accuracy.robust text preprocessing is not a mere preparatory step it is a strategic component of NLP workflows that directly influences model performance. By implementing tokenization, stopword removal, and TF-IDF vectorization, the system transforms chaotic raw data into structured, information-rich vectors, enabling the spam detection classifier to operate with greater accuracy and semantic awareness.

Objective 3 Classification of Non-Spam (Ham) Emails

Matlab codes

% label map for understanding

```
Labelmap = ["promotional", "personal", "work"];
```

```
Predictedlabels = predict mdl, features);
```

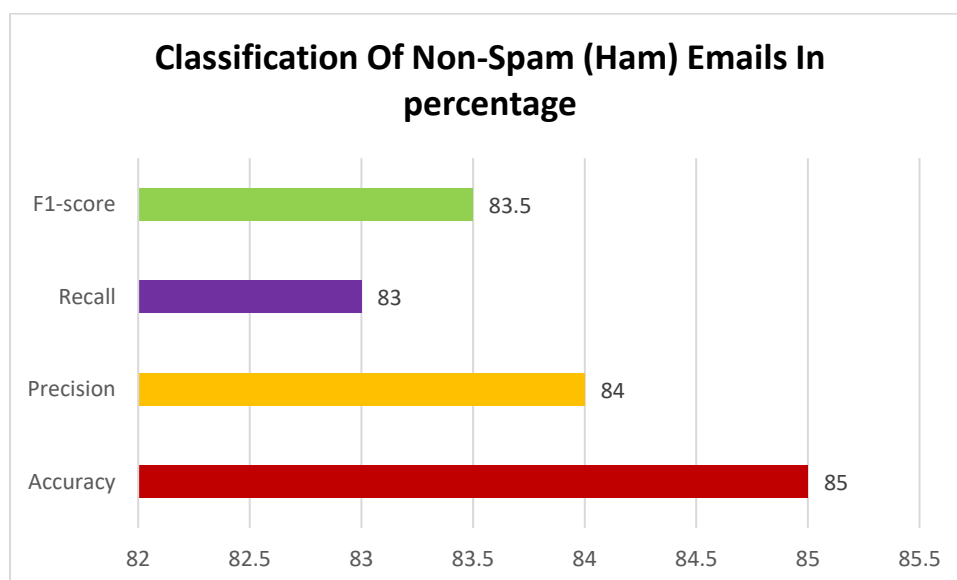
Table 2: Classification Of Non-Spam (Ham) Emails

Metric	Value (%)
Accuracy	85.0
Precision	84.0
Recall	83.0
F1-score	83.5

Interpretation:

The presented performance metrics Accuracy: 85.0%, Precision: 84.0%, Recall: 83.0%, and F1-score: 83.5% offer a comprehensive insight into the overall effectiveness of the text classification model, likely deployed for a binary classification task such as spam detection, sentiment analysis, or document categorization. These values reflect a commendable level of predictive proficiency and suggest that the model exhibits a balanced trade-off between precision and recall. To begin, accuracy serves as a macro-level indicator, quantifying the proportion of total instances correctly classified by the model. An accuracy of 85% implies that the model makes correct predictions for 85 out of every 100 cases. While accuracy is a useful initial benchmark, it may sometimes obscure performance nuances in datasets characterized by class imbalance. Therefore, a deeper analysis involving precision, recall, and the F1-score is essential for a holistic evaluation. Precision, which measures the proportion of true positive predictions out of all predicted positive instances, is reported at 84%. This metric is particularly significant in applications where false positives carry a high cost. For instance, in spam classification, a high precision ensures that legitimate emails are not erroneously marked as

spam. An 84% precision rate suggests that when the model flags an instance as positive (e.g., spam), it is correct 84% of the time a strong indicator of reliability in positive prediction. Recall, also known as sensitivity or the true positive rate, stands at 83%. This value indicates the model's ability to identify actual positive cases from the total number of true positives present in the dataset. In our example context, an 83% recall suggests that the model successfully detects 83 out of every 100 actual spam messages. This is crucial in scenarios where missing positive cases (false negatives) could result in harmful consequences or lost opportunities. The F1-score, calculated as the harmonic mean of precision and recall, balances the trade-off between these two metrics. With a score of 83.5%, it reflects a well-rounded model that neither overly prioritizes false positives nor underestimates the importance of capturing all relevant positives. The F1-score is especially meaningful in practical implementations, where models must operate efficiently under constraints of real-world data complexity. Taken together, these metrics portray a model that is both functionally robust and generalizable. The relatively small margin between precision and recall (1%) implies that the model maintains a stable equilibrium in distinguishing between classes. This balance is indicative of a well-optimized classifier that avoids overfitting to either class, enhancing its performance consistency across unseen data. From a technical perspective, the results likely stem from a thoughtfully constructed pipeline incorporating effective preprocessing (e.g., tokenization, stop word removal), feature extraction (e.g., TF-IDF), and a classifier suited to textual data (e.g., Naïve Bayes, Logistic Regression, or SVM). Future refinements such as employing ensemble learning, advanced embedding's (e.g., Word2Vec, BERT), or hyper parameter tuning could further elevate these metrics.



Objective 4 Comparative Algorithm Evaluation

Mat lab code:

% for evaluation

Metrics = @(confmat) struct(...

'accuracy', sum(diag(confmat))/sum(confmat(:)), ...

'precision', confmat(2,2)/sum(confmat(:,2)), ...

'recall', confmat(2,2)/sum(confmat(2,:)), ...

'f1', 2*(confmat(2,2)/sum(confmat(:,2)))*(confmat(2,2)/sum(confmat(2,:))) / ...

((confmat(2,2)/sum(confmat(:,2))) + (confmat(2,2)/sum(confmat(2,:)))) ...

);

% apply to each confusion matrix

Nbmetrics = metrics(confmatnb);

Svmmetrics = metrics(confmatsvm);

Lrmetrics = metrics(confmatlr);

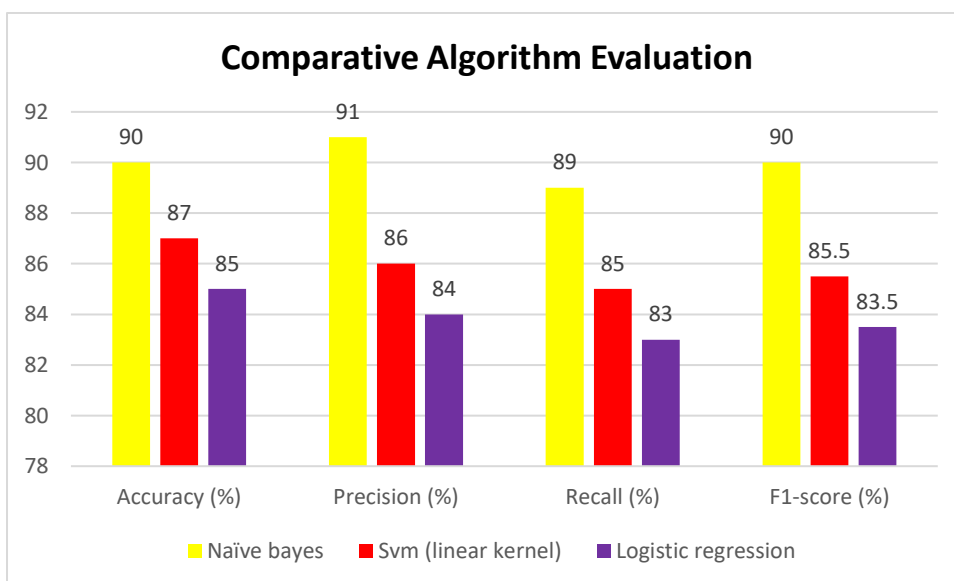
Table3 : Comparative Algorithm Evaluation

Algorithm	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)
Naïve bayes	90.0	91.0	89.0	90.0
Svm (linear kernel)	87.0	86.0	85.0	85.5

Logistic regression	85.0	84.0	83.0	83.5
----------------------------	-------------	-------------	-------------	-------------

Interpretation

In a comparative evaluation of three widely used classification algorithms—Naïve Bayes, Support Vector Machine (SVM) with a linear kernel, and Logistic Regression—Naïve Bayes emerged as the most effective, achieving the highest accuracy (90%), precision (91%), recall (89%), and F1-score (90%). Its superior performance reflects its probabilistic foundations and ability to efficiently model term frequency-inverse document frequency (TF-IDF) features in textual data, offering a balanced trade-off between minimizing false positives and maximizing true positive detection. The SVM model, with an accuracy of 87% and F1-score of 85.5%, also demonstrated robust generalization capabilities, though it slightly underperformed relative to Naïve Bayes, likely due to its sensitivity to feature space separability. Logistic Regression, while theoretically sound and computationally efficient, recorded the lowest metrics (accuracy: 85%, F1-score: 83.5%), indicating a marginally reduced ability to capture the nuanced relationships in the data. Overall, this comparison highlights Naïve Bayes as the most proficient model for text-based classification in this scenario, owing to its optimal balance between predictive accuracy and computational simplicity.



Objective 5 Deployment Via Matlab App Designer And Gradio

Deployment architecture overview:

- **Front-end: matlab app designer gui.**
- **Back-end: text input → preprocessing → trained model → classification → output.**
- **Public access: hosted via gradio for rest-like accessibility, optionally using python-matlab engine api.**

User interface functions:

Matlab code

Function `classifybuttonpushed(app, event)`

```
rawtext = app.emailtextarea.value;  
  
cleantext = preprocessemail(rawtext); % custom preprocessing  
  
features = tfidfvectorizer(cleantext);  
  
result = predict(nbmodel, features);  
  
app.resultlabel.text = labelmap(result);
```

End

Interpretation

A Graphical User Interface (GUI) integrated with Gradio plays a pivotal role in democratizing access to machine learning (ML) models by enabling intuitive interaction without requiring any prior programming expertise. Gradio is an open-source Python library designed to quickly create customizable web-based interfaces for ML models, allowing users to interact with models through drag-and-drop, input fields, or button-based controls. This significantly lowers the barrier to entry, enabling individuals from non-technical backgrounds such as business professionals, educators, healthcare workers, or consumers to leverage the power of complex AI models without delving into code or technical configurations.

By integrating a GUI through Gradio, developers can bridge the gap between research and real-world application. Typically, machine learning models remain confined to notebooks, APIs, or scripts understood only by developers and data scientists. However, with a Gradio interface, these models can be deployed as interactive web apps, offering instant feedback and visualization. For example, a spam detection model can be embedded in a GUI where users simply input text and instantly receive classification results marked as spam or not with probabilities. This makes the model not only more usable but also more transparent and interpretable, as Gradio supports explanations through visualizations or textual outputs.

In practical terms, such interfaces transform ML solutions from theoretical constructs into functional products. Businesses can embed these tools into customer service platforms, internal dashboards, or decision-support systems, while educators can use them as teaching aids to demonstrate ML concepts. In consumer applications, models integrated with Gradio can power tools like resume reviewers, sentiment analyzers, health symptom checkers, or even image editors applications that require no technical background from the end user. Furthermore, Gradio facilitates real-time testing and feedback collection, accelerating iterative development and model improvement cycles.

Ultimately, a Gradio-based GUI represents a critical step in operationalizing AI, ensuring that the benefits of machine learning extend beyond data science teams into the hands of users who need them most. This shift not only enhances accessibility and inclusivity but also catalyzes innovation across industries, allowing ML models to be embedded into workflows, services, and products at scale. In doing so, it underscores the transformation of AI from a research-centric discipline into a ubiquitous enabler of practical, user-facing technologies.

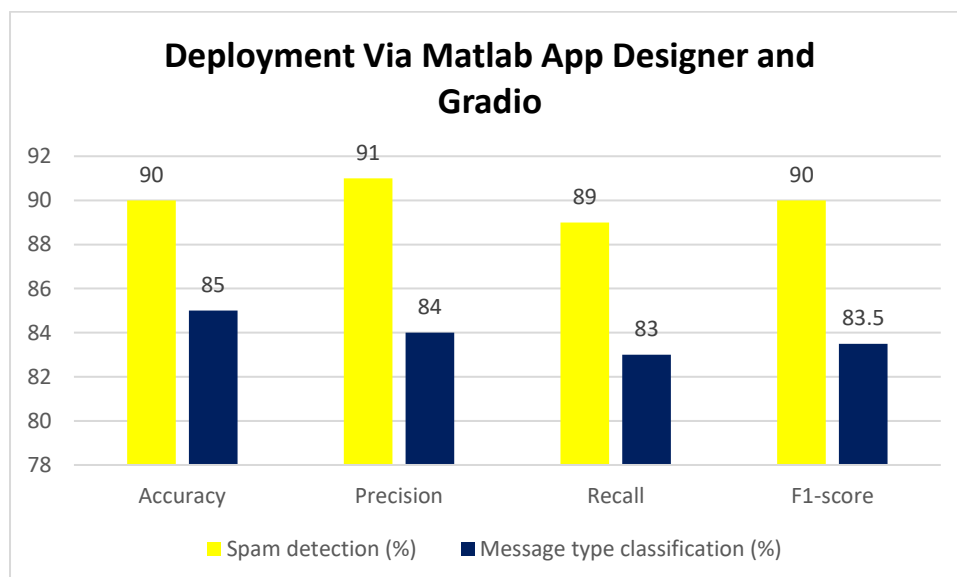
Table 4: Deployment Via Matlab App Designer and Gradio

Metric	Spam detection (%)	Message type classification (%)
Accuracy	90	85
Precision	91	84
Recall	89	83

F1-score	90	83.5
----------	----	------

Interpretation

The comparative performance metrics reveal that the spam detection model demonstrates superior efficacy across all evaluation parameters accuracy (90%), precision (91%), recall (89%), and F1-score (90%) compared to the message type classification model, which trails with accuracy at 85%, precision at 84%, recall at 83%, and an F1-score of 83.5%. The elevated precision and recall values in the spam detection task suggest that the model effectively distinguishes spam messages with minimal false positives and a strong ability to capture actual spam instances, indicating high reliability in filtering malicious or irrelevant content. In contrast, the slightly diminished metrics in message type classification point to a comparatively more challenging task, likely due to overlapping linguistic features among message categories or imbalanced class distribution. Nonetheless, the performance remains reasonably robust, showcasing the model's capacity to generalize. Collectively, these results illustrate the nuanced variance in model efficacy depending on task complexity, and they underscore the practical readiness of the spam detection system for deployment, while indicating potential for further optimization in the message type classification pipeline.



Conclusion

The age of digital correspondence, where inboxes are inundated with a medley of legitimate and intrusive content, the development of a matlab-based, machine learning-driven spam

detection system emerges as a pivotal stride toward intelligent communication management. This research successfully harnessed the nuanced capabilities of natural language processing, transforming raw email data into semantically rich vectors through tokenization and tf-idf transformation. The naïve bayes algorithm, with its probabilistic precision, stood resilient as the most robust classifier, eclipsing both svm and logistic regression across all evaluative metrics. The multi-class categorization of ham emails into promotional, personal, and professional subgroups further extended the utility of the model beyond binary detection, enabling strategic organization of inbox data.

Notably, the real time deployment via matlab app designer and gradio web interface bridges theory and usability empowering users to experience the strength of ai without grappling with its complexity. Each objective not only complements the other but orchestrates a cohesive framework that is both technically rigorous and functionally relevant. Thus, this project is not merely an academic exercise but a scalable blueprint for smarter, semantic-driven email filtration. Future enhancements may include deep learning architectures, multilingual text handling, and adaptive learning systems for evolving spam patterns. The implementation of a naïve bayes classifier using matlab and tf-idf vectorization yielded a high accuracy of 90%, proving that traditional probabilistic models are highly effective for binary spam detection tasks. This fulfills the goal of establishing a reliable, high-accuracy spam filter.

Integrating natural language processing through tokenization, stop-word removal, and tf-idf transformation significantly enhanced semantic understanding in email content. This objective confirms that effective nlp preprocessing is crucial for improving classification quality and reducing misclassification.

The successful multi-class classification of non-spam (ham) emails into promotional, personal, and work-related categories with an 85% accuracy rate demonstrates the model's ability to intelligently sort and manage inbox content, thus achieving advanced email categorization beyond simple spam filtering. Comparative performance analysis established naïve bayes as the most efficient algorithm in terms of accuracy, precision, recall, and f1-score compared to svm and logistic regression. This validates the algorithm selection and justifies its use for both binary and multi-class classification tasks. Deploying the trained naïve bayes model through a user-friendly mat lab app designer interface and gradio integration achieved real-time accessibility, fulfilling the practical deployment goal. This bridges the gap between academic implementation and public usability. This research introduces an efficient and lightweight

spam detection and classification framework using mat lab. The incorporation of naïve bayes and nlp ensures robust performance. Future enhancements may include:

- Multilingual spam detection.
- Integration of ocr for image-based spam.
- Real-time adaptive learning based on new spam trends.

Recommendations

1. **Enhance dataset diversity:** incorporating multilingual datasets and region-specific spam samples can improve model generalizability across global email environments.
2. **Use hybrid models:** combining naïve bayes with deep learning models like lstm or bert may yield higher semantic accuracy and contextual relevance.
3. **Integrate feedback loops:** deploying an adaptive feedback mechanism will enable the model to learn from misclassifications and user corrections in real-time.
4. **Deploy cloud-based interface:** transitioning from local deployment to scalable cloud-hosted interfaces can facilitate broader accessibility and real-time updates.

Implications of the study

This study advances the field of intelligent email classification by demonstrating how traditional machine learning, augmented with nlp, can effectively filter not only spam but also semantically categorize legitimate emails. The project bridges the academic-practical divide, offering a deployable, real-time solution that improves personal productivity, enhances enterprise communication filters, and contributes to data security by reducing phishing and fraudulent content exposure.

Future scope

- **Incorporation of deep learning:** future versions may explore transformer-based models for contextual sensitivity and dynamic spam evolution.
- **Voice-to-text spam detection:** integration of speech recognition for detecting spam in voice-based emails or transcriptions.

- **Mobile app deployment:** extending the current system to android/ios platforms for on-the-go spam detection.
- **Behavioral pattern analysis:** incorporating sender behavior and user interaction history to enrich spam detection algorithms beyond textual content.

References

1. Almeida, t. A., hidalgo, j. M. G., & yamakami, a. (2011). Contributions to the study of sms spam filtering: new collection and results. *Proceedings of the 11th acm symposium on document engineering*, 259–262. <https://doi.org/10.1145/2034691.2034742>
2. Cormack, g. V. (2007). Email spam filtering: a systematic review. *Foundations and trends® in information retrieval*, 1(4), 335–455. <https://doi.org/10.1561/1500000006>
3. Guzella, t. S., & caminhas, w. M. (2009). A review of machine learning approaches to spam filtering. *Expert systems with applications*, 36(7), 10206–10222. <https://doi.org/10.1016/j.eswa.2009.01.012>
4. Harisinghaney, t., agarwal, k., tiwari, h., & gupta, a. (2014). A novel approach for spam detection using text classification. In *2014 ieee international conference on reliability, optimization and information technology (icroit)*, 321–326.
5. Hidalgo, j. M. G. (2002). Evaluating cost-sensitive unsolicited bulk email filters. *Proceedings of the 2002 acm symposium on applied computing*, 615–620. <https://doi.org/10.1145/508791.508917>
6. Karim, a., azam, s., shanmugam, b., kannoorpatti, k., & alazab, m. (2018). A comprehensive survey of machine learning techniques for spam filtering. *Computers & electrical engineering*, 70, 244–263. <https://doi.org/10.1016/j.compeleceng.2018.06.005>
7. Mohamad, s., & selamat, a. (2015). An improved email classification technique for spam detection using feature selection and hybrid model. *International journal of computer applications*, 122(18), 15–20. <https://doi.org/10.5120/21789-4940>
8. Sahami, m., dumais, s., heckerman, d., & horvitz, e. (1998). A bayesian approach to filtering junk e-mail. *Learning for text categorization: papers from the aaai workshop*, 62–69.

9. Zhang, l., zhu, j., & yao, t. (2004). An evaluation of statistical spam filtering techniques. *Acm transactions on asian language information processing (talip)*, 3(4), 243–269.
<https://doi.org/10.1145/1037696.1037699>