

International Research Journal of Modernization in Engineering Technology and Science

( Peer-Reviewed, Open Access, Fully Refereed International Journal ) Volume:06/Issue:02/February-2024 Impact Factor- 7.868 ww

www.irjmets.com

# **ENHANCING COMMUNICATION SECURITY THROUGH MACHINE**

# **LEARNING AND STREAMLIT**

# Lokesh<sup>\*1</sup>, Kaviya<sup>\*2</sup>, Yuvashree<sup>\*3</sup>, Vimalesh<sup>\*4</sup>

\*1,2,3,4Student, Department Of Computer Science And Engineering, Gojan School Of Business And Technology, Chennai, Tamil Nadu, India.

DOI : https://www.doi.org/10.56726/IRJMETS49790

# ABSTRACT

The effectiveness of many machine learning methods for identifying spam in SMS and email interactions is examined in this study's abstract. Our study assesses each method's performance using Multinomial Naive Bayes, Random Forest, Support Vector Classifier, and Extra Tree Classifier. Preprocessing labelled datasets and extracting relevant features for model training are part of the study. Additionally, interaction with the Google Spreadsheet API and Google Drive API is implemented to help with real-time analysis and decision-making. Comprehensive testing and assessment are used to provide insights into the scalability, computational efficiency, and algorithm performance. The results further the progress of digital environment communication security and spam detection.

**Keywords:** Spam Detection, Machine Learning, Multinomial Naive Bayes, Support Vector Classifier, Google Spreadsheet API, Google Drive API.

# I. INTRODUCTION

The project uses an interactive machine learning pipeline built using Python modules like Streamlit and Google APIs to solve the ongoing problem of spam in communication technologies, such as email and SMS. Strong models are trained on pre-labeled datasets by the system using methods such as Multinomial Naive Bayes, Support Vector Classifiers, Random Forests, and Extra Trees Classifiers to achieve accurate spam classification. An interactive online interface with Streamlit is part of the implementation, which makes it simple for users to test and contrast various classifiers. The smooth import and export of data between Google Spreadsheets API and Google Drive enhances user accessibility. The research highlights the importance of spam identification in contemporary communication systems and how it protects users against a range of unsolicited messages, including malicious information and advertisements. Through the utilization of machine learning methodologies, including supervised learning and natural language processing, the project facilitates the examination of email and SMS content, sender information, and additional variables to ensure efficient spam identification. The online application development platform Streamlit is marketed as being easy to use, enabling users to enter messages for analysis and build interactive spam detection algorithms. To sum up, the initiative makes using the platform more enjoyable while protecting personal information and offering a useful tool to fight spam on real-world communication channels.

#### **OVERVIEW OF A DATASET:**

## II. METHODOLOGY

Collections of labeled messages (spam and ham) are easily accessible through public datasets, which are typically derived from users who voluntarily provided their information. Easy access and pre-existing labels are benefits. Still, they might not accurately represent unique user experiences or current spam patterns.

Type of Messages	Total Percentage	Total Attributes	
Spam	87%	2	
Ham	13%	2	

#### Table 1: Classification of Dataset

#### DATA CLEANING & PRE-PROCESSING:

This study uses a Kaggle dataset to examine the significance of data cleaning for improving machine learning models' accuracy and performance. To illustrate their effect on computing efficiency and forecast accuracy, different data cleaning methods are used. The findings underscore the critical role that preprocessing processes



# International Research Journal of Modernization in Engineering Technology and Science

(Peer-Reviewed, Open Access, Fully Refereed International Journal) Volume:06/Issue:02/February-2024 Impact Factor- 7.868 ww

www.irjmets.com

play in machine learning workflows by demonstrating a strong link between clean data and effective models. Increasing the accuracy of machine learning models requires effective data preprocessing. In order to improve model reliability, this research suggests methods for removing duplicate data and noise. We address these issues by talking about feature-based approaches, normalization strategies, and statistical methodologies, and we provide implementation and impact analysis. Case studies from the real world show how effective our method is in different fields.

#### FEATURE EXTRACTION:

Transforming unstructured text input into formats appropriate for machine learning is known as feature extraction, and it is used in spam detection. Word frequencies and significance are measured, respectively, by techniques like TF-IDF and Bag-of-Words. Word embeddings represent semantic linkages; N-grams capture sequential word patterns. Metadata characteristics offer contextual signals, and text preprocessing techniques standardise data. By combining these techniques, significant characteristics are taken out and used to train models that can differentiate between spam and real messages in email and SMS correspondence.

#### TF-IDF:

In natural language processing, TF-IDF (Term Frequency-Inverse Document Frequency) is a feature extraction method. It evaluates a word's significance in a document by taking into account both how frequently it appears in the document (TF) and how uncommonly it occurs in the corpus (IDF). Higher weights are assigned to words that are common inside a document but uncommon throughout the corpus. TF-IDF works well at identifying the importance of words unique to a given document while undervaluing phrases that are used frequently. This approach helps find important phrases that set texts apart from each other. It is extensively employed in tasks where a precise understanding of word importance is essential for analysis and categorization, such as spam detection, document classification, and information retrieval.

## III. MODELING AND ANALYSIS

#### MULTINOMINAL NAIVE BAYES(MNB):

Through probabilistic concepts and feature independence assumptions, this paper examines the Multinomial Naive Bayes method for spam identification. Extensive testing shows that the system performs well in correctly categorising spam messages. The study demonstrates Multinomial Naive Baye.

## RANDOM FOREST(RF):

The present investigation explores the Random Forest algorithm for spam detection through decision tree aggregation and ensemble learning. Extensive testing shows that the system performs well in correctly categorising spam messages. The study highlights Random Forest's potential for reliable spam detection systems.

## **SUPPORT VECTOR CLASSIFIER(SVC):**

It investigates the use of kernel functions and margin maximisation in Support Vector Classifiers for spam identification. The algorithm demonstrates its accuracy in classifying spam messages after thorough study. The study demonstrates how SVC can be used to create reliable spam detection system.

## EXTRA TREE CLASSIFIER(ETC):

The present investigation uses randomised decision trees and feature selection to investigate the Extra Tree Classifier for spam detection. After a comprehensive analysis, the programme demonstrates that it can reliably identify spam messages. The study emphasises how the Extra Tree Classifier may be used to create reliable spam detection systems.

## TESTING THE CLASSIFICATION RESULTS:

It addresses the spam detection training phase by using labelled data to train algorithms such as Random Forests and Support Vector Machines. The models learn to classify spam messages accurately through iterative optimisation. The study emphasises how useful customised model training is for strong spam detection systems.

It tests machine learning models during the spam detection testing phase by employing methods such as holdout validation and cross-validation. The model's efficacy in categorising messages is evaluated using metrics like accuracy and precision. The study emphasises how crucial comprehensive testing is to accurate



# International Research Journal of Modernization in Engineering Technology and Science (Peer-Reviewed, Open Access, Fully Refereed International Journal)

Volume:06/Issue:02/February-2024 Impact Factor- 7.868 www.irjmets.com

performance evaluation of spam detection systems.

#### FRAMEWORK:

Users can quickly determine whether or not a communication is spam with our Streamlit-powered spam SMS and email detection tool. They are able to view the outcomes and determine how to handle each message. Streamlit makes it simple for us to add new features and enhance the tool in response to user feedback so that we can combat spam more effectively.

#### **IMPLEMENTING APIs:**

We use a google sheets Api as a Database when we click a button it automatically stores the our input Messages and also store a whether message is spam or ham.

<> /	app.py	×
73		# Display Result
		result message = f" <div class="result-message"> Analysis Complete: {'Spam' if result == 1 else 'Not Spam'}!</div> "
75		st.markdown(result message, unsafe allow html=True)
76		
77		# Write to Google Sheets
		<pre>sheet url = "https://docs.google.com/spreadsheets/d/12MrdL6UMnFW7GSw5mHvPUAKiejXy409WUd50pCjh390/edit?usp=sharing"</pre>
79		<pre>sheet_name = "SpamPredictions"</pre>
		<pre>sheet = client.open_by_url(sheet_url).sheet1</pre>
81		
82		# Append only the necessary information to columns A and B
83		<pre>sheet.append_row([input_sms, "Spam" if result == 1 else "Not Spam"], value_input_option='USER_ENTERED')</pre>
84		
85		
		<pre>sheet.resize(rows=len(sheet.get_all_values()) + 1)</pre>
87		
88		# Send email using smtplib
89		<pre>subject = f"Prediction Result: {'Spam' if result == 1 else 'Not Spam'}"</pre>
90		body = f"The input message is classified as: {'Spam' if result == 1 else 'Not Spam'}."
91		
92		<pre>msg = MIMEText(body)</pre>
93		<pre>msg['Subject'] = subject</pre>
94		msg['From'] = GMAIL_USER
95		<pre>msg['To'] = 'lokeshkcse314@gmail.com' # Replace with the recipient's email address</pre>
97		with smtplib.SMTP_SSL('smtp.gmail.com', 465) as server:
98		server.login(GMAIL_USER, GMAIL_PASSWORD)
99		<pre>server.sendmail(GMAIL_USER, ['lokeshkcse314@gmail.com'], msg.as_string()) # Use a list for multiple recipients</pre>
100		
101		
102		new_page_content = f"""
103		<h2 class="text-center">Prediction Results</h2>
104		<pre>The input message is classified as: {'Spam' if result == 1 else 'Not Spam'}.</pre>
105		
106		st.markdown(new_page_content, unsafe_allow_html=True)
107		

# Figure 1: Connecting to Sheets API IV. RESULTS AND DISCUSSION

#### **CLASSIFICATION REPORT:**

While precision evaluates the accuracy of positive predictions, accuracy gauges the overall soundness of a model's predictions. Precision precisely assesses the ratio of true positives to the sum of true positives and false positives, whereas accuracy computes the percentage of correctly classified samples out of the total. These metrics are essential for assessing how well spam detection models perform in distinguishing between spam and authentic messages.

#### **PERFORMANCE ANALYSIS:**

We compare 11 algorithms graph namely KNN, MNB, RF, SVC, ETC, LR, ADA-BOOST, XGB, GBDT and DT. We got a higher accuracy on Multinominal Naïve Bayes.

www.irjmets.com



# International Research Journal of Modernization in Engineering Technology and Science

Table 2: Comparison of Algorithms

( Peer-Reviewed, Open Access, Fully Refereed International Journal ) Volume:06/Issue:02/February-2024 Impact Factor- 7.868 ww

www.irjmets.com

)	new_	df_scaled.	merge(temp	_df,on=' <mark>Alg</mark>	orithm')					
		Algorithm	Ассигасу	Precision	Accuracy_scaling_x	Precision_scaling_x	Accuracy_scaling_y	Precision_scaling_y	Accuracy_num_chars	Precision_num_chars
	0	KN	0.905222	1.000000	0.905222	1.000000	0.905222	1.000000	0.905222	1.000000
	1	NB	0.970986	1.000000	0.970986	1.000000	0.970986	1.000000	0.970986	1.000000
	2	RF	0.975822	0.982906	0.975822	0.982906	0.975822	0.982906	0.975822	0.982906
	3	SVC	0.975822	0.974790	0.975822	0.974790	0.975822	0.974790	0.975822	0.974790
	4	ETC	0.974855	0.974576	0.974855	0.974576	0.974855	0.974576	0.974855	0.974576
	5	LR	0.958414	0.970297	0.958414	0.970297	0.958414	0.970297	0.958414	0.970297
	6	AdaBoost	0.960348	0.929204	0.960348	0.929204	0.960348	0.929204	0.960348	0.929204
	7	xgb	0.967118	0.926230	0.967118	0.926230	0.967118	0.926230	0.967118	0.926230
	8	GBDT	0.946809	0.919192	0.946809	0.919192	0.946809	0.919192	0.946809	0.919192
	9	BgC	0.958414	0.868217	0.958414	0.868217	0.958414	0.868217	0.958414	0.868217
	10	DT	0.930368	0.836735	0.930368	0.836735	0.930368	0.836735	0.930368	0.836735



Figure 2: Accuracy and Precision of Algorithms

## **RESULTS:**

Users can promptly evaluate and respond to incoming messages for efficient spam management as the output shows a confidence score and whether the communication is categorised as spam or Ham.





International Research Journal of Modernization in Engineering Technology and Science (Peer-Reviewed, Open Access, Fully Refereed International Journal)

Volume:06/Issue:02/Febru	ary-2024 Impact Factor- 7.868	www.irj	mets.	com
← → C ③ localhost8501		☆	. 0	:
IAbout         This app predicts whether an input message is spam or not using a pre-trained model.	<ul> <li>Image: Image: Im</li></ul>	, ,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Deploy	

#### Figure 4: Output For Not Spam

We have effectively used the API to combine our input messages with Google Sheets. To give you an idea, picture the spreadsheet as having rows of nicely arranged data, each row representing a message whether it is Spam or Ham.

 	25 docs.google	.com/spreads	sheets/d/12MrdL6UMnFW7GSwSmHvPUAKiejXy4Q9WUd5C	DpCjh390/edit#gid=0		☆ 🛛 🕒
BamP File Ed	redictions ☆ t View Insert	Format [	Data Tools Extensions Help		U =	CH - Share -
९ ५ २ १५	🖨 🗣 100%	6 👻   E	% .0, .00 123   Defaul ▼   - 10 +   B	I ↔ A A 5 + E 5 +	E • ↓ • IqI • A •   •	⇔⊥шү⊑∙Σ ∧
A	B	c	D	E	F	G H I
1			SPAM MESSAGES	SPAM OR HAM C	LASSIFIER	
2			"Go until jurong point, crazy Available only in bugis n great	t world la e buffet C Not Spam		
3			"Go until jurong point, crazy Available only in bugis n great	t world la e buffet C Not Spam		
4			"FreeMsg Hey there darling it's been 3 week's now and no v	"FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like sor Not Spam		
5			*FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like sor Not Spam			
6			WINNER! As a valued network customer you have been se	lected to receivea å£. Spam		
7						
8						



#### V. CONCLUSION

We get a high accuracy rate of 97% with our spam SMS and email detection system, which uses the Multinomial Naive Bayes (Mnb) algorithm and connects effortlessly with the Sheets API. This technology helps users efficiently recognise and remove spam communications from their email and SMS inboxes. We guarantee effective data handling and administration by utilising Mnb in conjunction with the Sheets API, improving the user experience overall in the fight against spam mails.

## ACKNOWLEDGEMENTS

We also express our thanks to our Head of the Department Mr. P. Senthil. M.Tech., who has been a constant source of inspiration and guidance in the course of the project, We record our sincere thanks to our Supervisor Mrs. S. Divya Assistant Professor for being instrumental in the completion of our project with his exemplary guidance.

## VI. REFERENCES

E. G. Dada, J. S. Bassi, H. Chiroma, S. M. Abdulhamid, A. O. Adetunmbi, and O. E. Ajibuwa, "Machine learning for email spam filtering: Review, approaches and open research problems," Heliyon, vol. 5, no. 6, Jun. 2019, Art. no. e01802, doi: 10.1016/j.heliyon.2019.e01802



#### International Research Journal of Modernization in Engineering Technology and Science (Peer-Reviewed, Open Access, Fully Refereed International Journal)

Volume:06/Issue:02/February-2024	Impact Factor- 7.868	www.irjmets.com

- [2] W. Awad and S. ELseuofi, "Machine learning methods for spam E-Mail classification," Int. J. Comput. Sci. Inf. Technol., vol. 3, no. 1, pp. 173–184, Feb. 2011, doi: 10.5121/ijcsit.2011.3112.
- [3] A. Wijaya and A. Bisri, "Hybrid decision tree and logistic regression classifier for email spam detection," in Proc. 8th Int. Conf. Inf. Technol. Electr. Eng. (ICITEE), Oct. 2016, pp. 1–4, doi: 10.1109/ICITEED.2016.7863267.
- [4] O. Saad, A. Darwish, and R. Faraj, "A survey of machine learning techniques for Spam filtering," Int. J. Comput. Sci. Netw. Secur., vol. 12, no. 2, p. 66, Feb. 2012.
- [5] Y. Cohen, D. Gordon, and D. Hendler, "Early detection of spamming accounts in large-Scale service provider networks," Knowl.-Based Syst., vol. 142, pp. 241–255, Feb. 2018.
- [6] S. Smadi, N. Aslam, and L. Zhang, "Detection of online phishing email using dynamic evolving neural network based on reinforcement learning," Decis. Support Syst., vol. 107, pp. 88–102, Mar. 2018.
- [7] M. Yesilyurt and Y. Yalman, "Security threats on mobile devices and their effects: estimations for the future", International Journal of Security and Its Applications, vol. 10, no. 2, pp. 13-26, 2016.
- [8] F. Ahmed, "The impact of SMS marketing on consumer behavior", The Business Management Review, vol. 10, no. 1, pp. 115-125, 2018.
- [9] R. Patel and P. Thakkar, "Opinion spam detection using feature selection", 2014 International Conference on Computational Intelligence and Communication Networks, pp. 560-564, 2014, November.
- [10] P.P. Chan, C. Yang, D. S. Yeung and W. W. Ng, "Spam filtering for short messages in adversarial environment", Neurocomputing, vol. 155, pp. 167-176, 2015.