

# EMAIL SPAM DETECTION USING PYTHON & MACHINE LEARNING

**PRESENTED BY:**  
Srividya Polneni  
Sindhuja Yerramalla



# PROBLEM STATEMENT

- These days, all official and sensitive communication is made through emails and Spam emails are a major issue on the internet. It is easy to send an email that contains spam messages by spammers.
- Spam fills our inbox with several irrelevant emails. Spammers can steal our sensitive information from our devices like files, and contact. Even though we have the latest technology, it is challenging to detect spam emails.





# ABSTRACT

- Automatic e-mail filtering seems to be the most effective method for countering spam at the moment.
- Only several years ago most of the spam could be reliably dealt with by blocking e-mails coming from certain addresses or filtering out messages with certain subject lines.
- Spammers began to use several tricky methods to overcome the filtering methods like using random sender addresses and/or appending random characters at the beginning or the end of the message subject line.
- Machine learning techniques are being used to automatically filter spam e-mail at a very successful rate.

# CONT..

- First, data collection and representation are mostly problem specific (i.e. e-mail messages), second, e-mail feature selection and feature reduction attempt to reduce the dimensionality (i.e. the number of features).
- Finally, the e-mail classification phase of the process finds the actual mapping between the training set and the testing set.
- Machine Learning approach includes lots of algorithms that can be used in e-mail filtering like Naïve Bayes, K-nearest neighbour, Support Vector Machine, and classifiers. In conclusion, we try to summarize the performance results of the few machine learning methods in terms of spam precision and accuracy.



# IMPLEMENTED METHODS

- In this system, we are implementing Natural Language Processing (NLP) like **TF-IDF** is one of the simple and robust methods to understand the context of a text.
- Term Frequency and Inverse Document Frequency (TF-IDF) are used to find the related content and important words and phrases in a larger text. Implementing TF-IDF analysis is very easy using Python.
- Computers cannot understand the meaning of a text, but they can understand numbers. The words can be converted to numbers which is called feature extraction.
- Later these features are trained with machine learning techniques such as **Support Vector Machine, K-Nearest Neighbors, and Naive Bayes Classifiers** for Spam email detection.

AutoSave Off spam\_train • Saved

File Home Insert Draw Page Layout Formulas Data Review View Automate Help

Undo Clipboard Font Alignment

F5569

	A	B	C	D	E	F	G	H	I	J	K
1	Class	Emails									
2	ham	Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat...									
3	ham	Ok lar... Joking wif u oni...									
4	spam	Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entry quest									
5	ham	U dun say so early hor... U c already then say...									
6	ham	Nah I don't think he goes to usf, he lives around here though									
7	spam	FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for it still? Tb									
8	ham	Even my brother is not like to speak with me. They treat me like aids patent.									
9	ham	As per your request 'Melle Melle (Oru Minnaminunginte Nuringu Vettam)' has been set as your callertune fo									
10	spam	WINNER!! As a valued network customer you have been selected to receivea 900 prize reward! To claim ca									
11	spam	Had your mobile 11 months or more? U R entitled to Update to the latest colour mobiles with camera for F									
12	ham	I'm gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough tod									
13	spam	SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days,									
14	spam	URGENT! You have won a 1 week FREE membership in our 100,000 Prize Jackpot! Txt the word: CLAIM to N									
15	ham	I've been searching for the right words to thank you for this breather. I promise i wont take your help for gr									
16	ham	I HAVE A DATE ON SUNDAY WITH WILL!!									
17	spam	XXXMobileMovieClub: To use your credit, click the WAP link in the next txt message or click here>> http://w									
18	ham	Oh k...i'm watching here:)									
19	ham	Eh u remember how 2 spell his name... Yes i did. He v naughty make until i v wet.									
20	ham	Fine if that is the way u feel. That is the way its gota b									
21	spam	England v Macedonia - dont miss the goals/team news. Txt ur national team to 87077 eg ENGLAND to 8707									
22	ham	Is that seriously how you spell his name?									
23	ham	I am going to try for 2 months ha ha only joking									
24	ham	So pay first lar... Then when is da stock comin...									
25	ham	Aft i finish my lunch then i go str down lor. Ard 3 smth lor. U finish ur lunch already?									
26	ham	Ffffffffff. Alright no way I can meet up with you sooner?									
27	ham	Just forced myself to eat a slice. I'm really not hungry tho. This sucks. Mark is getting worried. He knows I'm									
28	ham	Lol your always so convincing.									

spam\_train

Ready Accessibility: Unavailable

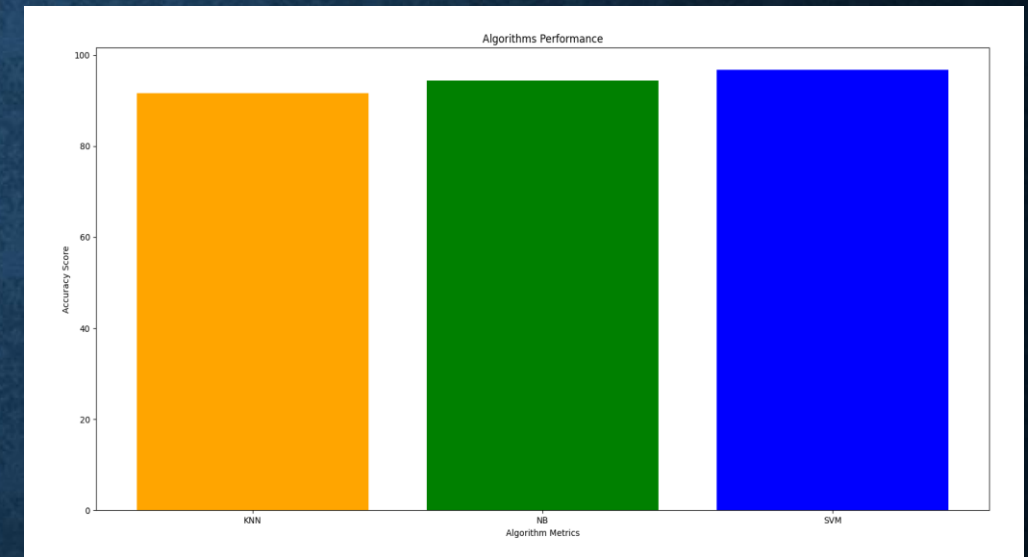
# DATA SETS USED

- In this system, we are using the Email Spam Classification dataset which is accessed from the Kaggle repository <https://www.kaggle.com/code/harshsinhal234/email-spam-classification-nlp/data>.
- In this system, we are using the text format dataset which is not understood by the ML techniques. So, it is a challenge to convert the text format dataset into a numerical dataset using the NLP technique and ML algorithms
- This dataset contains multiple columns like text and spam. The **text** column will contain the email subjects and the **spam** column will contain 0 or 1. The '0' value indicates the **Ham** and the '1' value indicated the **Spam**.



# RESULTS & CONCLUSION

	Algorithm	Accuracy	Precision	Recall	F1_Score
1	KNN	91.63179916317992	63.703703703703695	47.175141242937855	51.343823979947075
2	NB	94.38135086670651	64.61739699149771	53.672316384180796	57.53936122357175
3	SVM	96.77226539151226	64.42606859556011	60.14813283726057	62.0542543588769



- So by the above results, we can conclude that the Algorithm performance is better in the SVM classification technique in all the features like Accuracy, Precision, Recall etc.,

## EMAIL SPAM DETECTION USING PYTHON & MACHINE LEARNING



### Admin Login

User Name

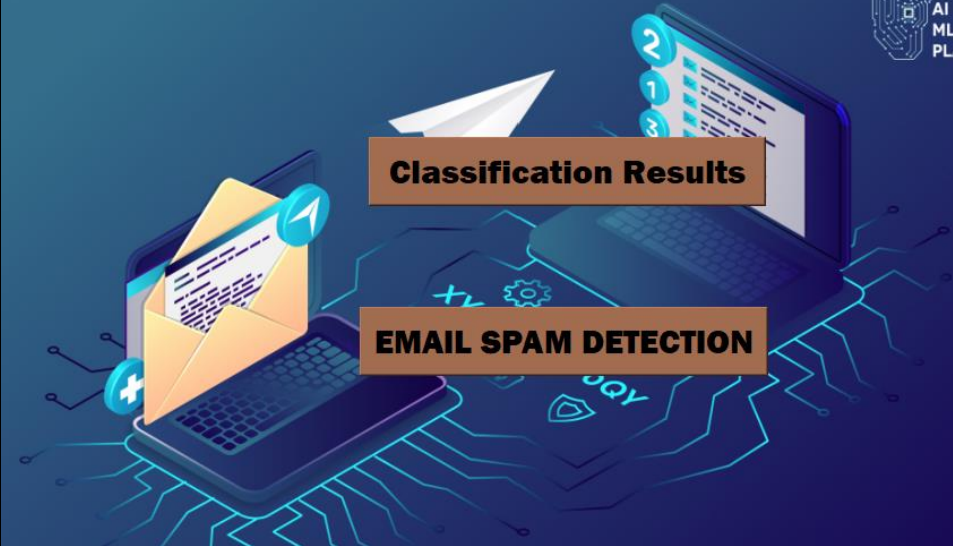
Password

LOGIN



### Classification Results

### EMAIL SPAM DETECTION



Enter Email Subject:

URGENT! You have won a 1 week FREE membership in our 100,000 Prize Jackpot! Txt the word: CLAIM to No: 81010 T&C www.dbuk.net LCCLTD POBOX 4403LDNW1A7RW18



Detection

Result : **SPAM**

Enter Email Subject:

Yup... Ok i go home look at the timings then i msg \_ again... Xuhui going to learn on 2nd may too but her lesson is at 8am



Detection

Result : **HAM**



# TEAM WORK!

- Contribution to work by each Individual

Team member: Srividya Polneni

- Dataset Selection

- Training Data

Team member: Sindhuja Yerramalla

- UI development

- Performance of Algorithms

**THANK YOU!!**