

## Ham/Spam email detection

### Introduction

The purpose of this project is to use our data science and machine learning knowledge accrued over the last several months to train a Machine Learning model that can categorize emails as either Spam or Not Spam (from now on, Not Spam will be referred to as Ham). At the end of this project, I had built a simple applet in which the user can input the contents of an email, and the applet processes the email using a trained machine learning model to categorize it as Spam or Ham.

### Data Selection

Before a machine learning model can be trained, I was required to find an appropriate dataset with which to train it. I had chosen the Enron Spam Dataset, as it was a large and diverse enough dataset with over 33,000 entries, with a close balance of Ham and Spam email examples (16,614 spam entries and 16,493 ham entries).

Using an evenly balanced dataset was beneficial, reducing the amount of cleanup necessary during the data preparation phase, and leading to a more reliable final model when the training was done.

### Data Cleanup

The structure of the dataset initially was:

Message	Spam/Ham
The content of the email, stored as a text value.	Whether this email is ham or spam, stored as a text value in the form of the words: Ham Spam

Initially, the Spam/Ham column was easily mapped to a boolean representation. All Ham values were converted to a 0, and all spam values were converted to a 1.

Secondly, remove any “stop words” found in the email content. Stop words are common words (e.g., “the”, “a”, “is”, “and”) that are often removed during text preprocessing in Natural Language Processing (NLP) because they typically carry little semantic meaning and unnecessarily complicate our analysis.

[illegible]

index	Message	Spam	Bool
1	gari product high island larger block commencing saturday p gross carlo expect gross tomorrow vastar own gross product georg x forward georg weissman hou ect den jamer carlo j rodriguez hou ect ect cc georg weissman hou ect ect melissa grave hou ect ect subject vastar resourc inc carlo pleas call linda get everyth set go eslim come tomorrow increas follow day base convers bill fischer bmar forward den jamer hou ect ect enron north america corp georg weissman den jamer hou ect ect cc gari byran hou ect ect melissa grave hou ect ect subject vastar resourc inc denrre atnap appear nomin vastar resourc inc high island larger block previous enron refer well vastar expect well commencing product sometim tomorrow told linda harti get telephon number ga control provid notifi turn tomorrow linda number record vici fax would pleas see someone contact linda advs submit futur nomini via e mail fax vici want georg x forward georg weissman hou ect lindi harti georg weissman hou ect ect cc subject h i effect medt min ipi time hour hour hour hour hour hour hour hour hour hour hour hour hour hour hour hour hour		0
2	calpin dail ga nomin doc		0
3	fjyi see note alore dore stella forward stella i morri hou ect ect sheryn schumack stella i morri hou ect ect cc howard b camp hou ect ect subject issu stella already taken care yesterday thank howard b camp stella i morri hou ect ect cc sheryn schumack hou ect ect howard b camp hou ect ect stacey newseil hou ect ect den jamer hou ect ect subject issu stella work stacey denrre resolv hi forward howard b camp hou ect ect sheryn schumack pm howard b camp hou ect ect cc subject issu creat account arrang purchas unoc enrgi meter deal trck volum deal explr		0
4	fjyi forward lauri allen hou ect pm kimberli vaughn pm lauri allen hou ect ect cc mari smith hou ect ect subject meter nov alloc lauri pt stranga ga get contract denrre forward kimberli vaughn hou ect pm lauri allen pm kimberli vaughn hou ect ect anila luong hou ect ect cc howard b camp hou ect ect mari smith hou ect ect subject meter nov alloc kim anila volum nm show alloc reliant contract novemb reliant point novemb forwarder volum alloc contract pleas make sure volum move reliant contract prior novemb dore thanx		0
5	jacki sanc inlt river plant shut last day flow meter mcmullen ga divrt meter hpl bay residu ga ga teco vastar vintag tejon swift still set use activ deal meter path manag teco vastar vintag tejon swift also see ga schedul pop meter pleas advnc need resol		0

The alpha, stemmed messages only have the key information and lack any of the noise in the form of punctuation and formatting found in the original dataset. This simplifies them and decreases the number of features necessary for our machine learning algorithm to analyse.

The final step in the data cleaning step, is to convert these alpha values into numerical tokens so that an ML algorithm can read them properly. This was done using Sci-Kit Learn's CountVectorizer. After all data was cleaned, this was the final structure of the dataset table.

Message	Spam_Bool
The content of the email, stored as a numerical token value.	Whether this email is ham or spam, stored as a numerical value in the form of the numbers: 0 1

## Model Training

After cleaning all the data, a machine learning model needed to be trained. After trying three separate models, the Logistic Regression algorithm proved to be the most accurate, with the fewest false positive values and the fewest false negative values.

The final trained Logistic Regression model had the following evaluation:

### Confusion Matrix:

<b>True Positives:</b> 4876	<b>False Positives:</b> 87
<b>False Negatives:</b> 46	<b>True Negatives:</b> 4924

**Accuracy:** 0.9866

**Precision:** 0.9906

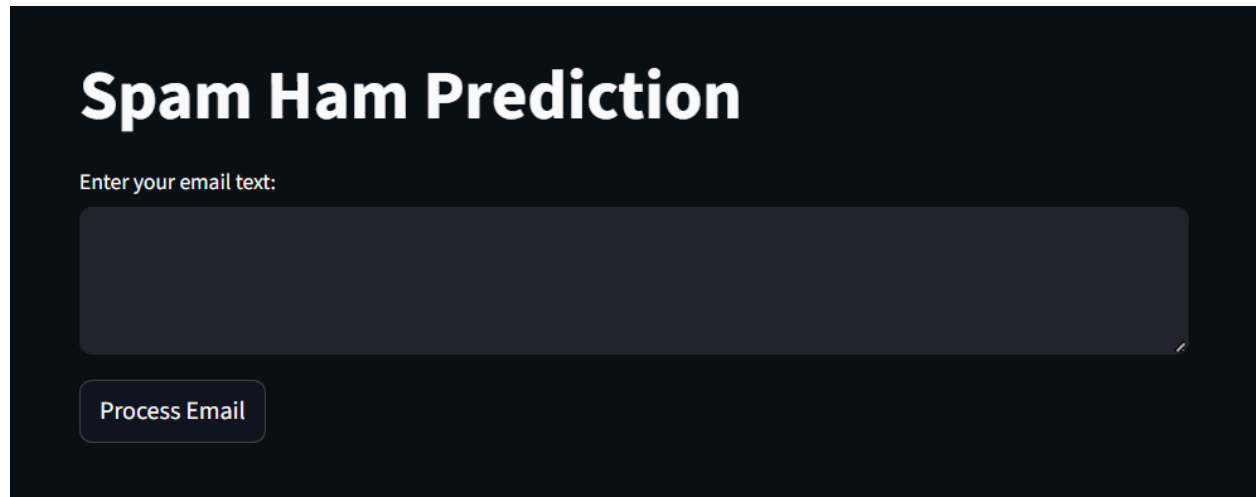
**Recall:** 0.9825

**f1-score:** 0.9865

This was a very satisfactory result, with a near-perfect Precision of over 99%. Satisfied with the model performance all that was left was to implement it inside an app.

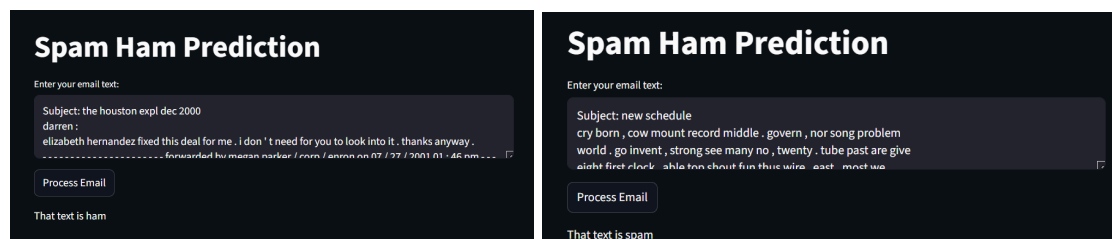
## Streamlit Application

The model was saved to a pickle file, and integrated into a simple streamlit applet.



The screenshot shows a Streamlit application titled "Spam Ham Prediction". It features a dark background with white text. Below the title, there is a label "Enter your email text:" followed by a large, empty text input box. Below the input box is a button labeled "Process Email".

The user can input any email text content into the textbox, and it'll categorize it into Spam or Ham using the previously-trained Machine Learning model. The sole complication I had not foreseen with the implementation into the streamlit applet, is that the same tokenizer used to initially tokenize the data had to be reused in the later applet. This was easily handled by saving the tokenizer to an external file in the transformation step, and loading it in the streamlit applet.



The image shows two side-by-side screenshots of the "Spam Ham Prediction" application. Both have a dark background and white text.

The left screenshot shows the input text: "Subject: the houston expl dec 2000  
darwin :  
elizabeth hernandez fixed this deal for me . i don ' t need for you to look into it . thanks anyway .  
..... forwarded by mason narker / csm / amon on 07 / 27 / 2003 01 : 46 pm ....". Below the input box is a button labeled "Process Email". At the bottom, it says "That text is ham".

The right screenshot shows the input text: "Subject: new schedule  
cry born , cow mount record middle , govern , nor song problem  
world . go invent , strong see many no , twenty . tube past are give  
eight first clock . able ten about fun this wire . east . most wa". Below the input box is a button labeled "Process Email". At the bottom, it says "That text is spam".

## **Reflection**

This project was an interesting exercise for me as it felt like my first time having to really do the entire ML pipeline with minimal guidance. I had to personally do the whole flow of finding a dataset, cleaning it, training a model with it, and then implementing it into a presentable app prototype. This project left me feeling more prepared for future Machine Learning projects, and steps that took longer this time around (most notably, data cleanup) I think will be much faster moving forward as the sense of confidence has seriously increased.