

CS5691: Pattern Recognition and Machine Learning

Assignment 3- Spam Classifier

Name: Anukul Dewangan

Roll No: CS22M016

SPAM or HAM?

In this assignment, you will build a spam classifier from scratch. No training data will be provided. You are free to use whatever training data that is publicly available/does not have any copyright restrictions (You can build your own training data as well if you think that is useful). You are free to extract features as you think will be appropriate for this problem. The final code you submit should have a function/procedure that when invoked will be able to automatically read a set of emails from a folder titled test in the current directory.

Each file in this folder will be a test email and will be named 'email#.txt' ('email1.txt', 'email2.txt', etc). For each of these emails, the classifier should predict +1 (spam) or 0 (non-spam). You are free to use whichever algorithm learned in the course to build a classifier (or even use more than one). The algorithms (except SVM) need to be coded from scratch. Your report should clearly detail information relating to the data set chosen, the features extracted and the exact algorithm/procedure used for training including hyperparameter tuning/kernel selection if any. The performance of the algorithm will be based on the accuracy of the test set.

Solution: I have used Naïve Bayes Algorithm for solving the Spam Classification Problem since a Naïve Bayes text Classifier is based on the Bayes theorem which helps us compute the conditional probabilities of occurrence of two events based on the probabilities of occurrence of each element. Therefore, those probabilities are extremely useful

Data set used for training and testing Naïve Bayes Algorithm for spam classification: 'spam.csv'.

Percentage of ham/spam emails in this dataset:

```
ham      0.865937
spam     0.134063
Name: Label, dtype: float64
```

We can see that 86.57 percent of our emails are ham and 13.4 percent of emails are spam.

Spam Classifier Accuracy using naïve bayes algorithm:

```
Accuracy on test Data: 97.4789237668162
Number of Incorrect Prediction: 28
number of correct Prediction: 1072
```

My Approach:

Step 1: Data pre-processing/ Data cleaning: Remove punctuations and convert every email into lowercase.

Step 2: Now, We are going to split out the data set into a training set and a testing set. I have used 80% dataset for training and 20% for testing.

Step 3: Creating a dictionary for training emails, testing emails, spam emails, and non-spam emails that consist of only unique words. We transform each email in the message column into a list by splitting the string at the space character and appending the list unique words.

Step 4: calculate prior Probabilities of $p(\text{spam} = 1)$ and $p(\text{spam} = 0)$

$$P(\text{spam} = 1) = \text{number of emails flagged 'spam'} / \text{total number of emails}$$
$$P(\text{ham} = 0) = \text{number of emails flagged 'ham'} / \text{total number of emails}.$$

Step 4: Take the input string, first do data cleaning, and create a dictionary that contains words that are present in the input string and our training dictionary, ignore words that are not present in our training dictionary.

Example: Input string “lucky win lucky joy kick “

Calculate $p(\text{lucky}/\text{label}=1)$, $p(\text{win}/\text{label}=1)$, $p(\text{joy}/\text{label}=1)$, $p(\text{kick}/\text{label}=1)$
 $P(\text{luck}/\text{label}=0)$, $p(\text{win}/\text{label} = 0)$, $p(\text{ joy}/\text{label} = 0)$, $p(\text{kick}/\text{label} = 0)$

$$P(\text{label}=1) = p(\text{label} = 1) * p(\text{lucky}/\text{label}=1) * p(\text{win}/\text{label}=1) * p(\text{joy}/\text{label}=1) * p(\text{kick}/\text{label}=1)$$
$$P(\text{label} = 0) = p(\text{label} = 0) * p(\text{lucky}/\text{label}=0) * p(\text{win}/\text{label}=0) * p(\text{joy}/\text{label}=0) * p(\text{kick}/\text{label}=0)$$

If $P(\text{label}=1) > p(\text{label} = 0)$:
 return ‘This email is a Spam’
else:
 return ‘This email is not a spam mail’

’.

Below is the result of some emails :

```
Email: SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply
Reply HL 4 info
This Email is a Spam
Email: I'm gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough today.
This Email is Not a Spam mail
```