

# Report on Programming Assignment Naive Bayesian Text Classification

Amy de Lange - s1534041

November 2019

## Introduction

The goal of the assignment is to program a classifier that can distinguish emails in the categories 'spam' and 'ham' (not-spam). The tools that were provided the assignment description manual and a corpus of emails to train and test with. I installed Python 2 and Atom (code editor) myself and followed an online course to retrieve the basic knowledge about the language of Python.

The purpose of this assignment is to learn the basic concepts of programming and machine learning. This will be done by building a Naive Bayesian Text Classifier, that implements a machine learning method for classifying texts based on the statistics about the occurrences of words in texts of the different categories. The first step would be to prepare the corpus into readable tokens where the classifier can be trained with. In the spam classifier, the occurrence of specific words in a text will be leading in classifying the text into one of the two classes. The specific words will be labelled in the training part where words will be counted in both classes to calculate later on their probability of being either spam or ham.

To improve the classifier, chi-square calculations will be first done on all the unique words of all the emails combined. These will generate the most distinguishing words from the corpus. This will be provided to the classifier to look for these specific words to calculate, at the end, the probability of a certain email for being spam or ham.

## Structure of the model

The model consists out of 6 parts. They will be all called upon in the main module. The parts are: Preparing the corpus for usage (1), creating two dictionaries with all the unique words from every class email set (2), counting word occurrences in the corpus per classified email (3), building a chi-square list (4), calculating word probabilities (5) and finally testing the the classifier by classifying the emails from the test corpus (6).

## Reading, Normalizing, Tokenizing

The corpus received for this assignment consist out of 10 parts, the first 9 for training the classifier and part 10 to test the classifier. The corpus consisted out of spam and ham emails. They could be distinguished by their filenames; spam consisted out of 'spms' and ham consisted out of 'msg'. To read the corpus, I made use of the module 'glob' to import files from my laptop directory. To categorize the right files to spamfiles and hamfiles, I made use of regular expressions to add upon the directory 'dir' that matches the filename characteristics of ham and spam.

By tokenizing the words in every email and normalizing them, we will get a dictionary of only lowercase words. The dictionaries of spam and ham will be filled up by a for-in loops their all files of that category.

```
1  """READING, NORMALIZING, TOKENIZING"""
2  def reader(dir):
3      ham_files = glob.glob(dir+"/*msg?.txt")
4      spam_files = glob.glob(dir+"/*spms*.txt")
5
6      ham_dictionary = {file: tokenizer(open(file).read()) for file
7                          in ham_files} # make a dictionary of all ham files
8      ham_tokens = ham_dictionary.values()
9      for ham_token in ham_tokens:
10         for e in ham_token:
11             ham_data.append(e)
12
13     spam_dictionary = {file: tokenizer(open(file).read()) for file
14                         in spam_files} # make a dictionary of all spam files
15     spam_tokens = spam_dictionary.values()
16     for spam_token in spam_tokens:
17         for e in spam_token:
18             spam_data.append(e)
19
20     return ham_data, spam_data
21
22 def normalise(word):
23     return re.sub(r'\W+', '', word.lower())
24
25 def tokenizer(str):
26     tokens = [normalise(word) for word in str.split(" ") if
27               normalise(word).isalpha()]
28     allTokens = tokens
29     return allTokens
```

Listing 1: Making dictionaries of all appearing words for each category

## Unique dictionaries

Later on in the assignment, we would like to retrieve a chi-square word list to train and test the classifier with. To generate that list, we would like to loop through a list of all unique words that occur in the corpus. Therefore, we make unique words dictionaries for each category. This will be done by an easy if statement to get a dictionary with the keys as words and the values as the times of occurrences in that class (ham/spam).

```
1  """UNIQUE WORDS"""
2
3  def UniqueWordsMaker():
4      print "=== spam_data: " + str(len(spam_data)) + " ==="
5      for word in spam_data:
6          if word not in spam_data_pr.keys():
7              spam_data_pr.update({word:1})
8
9          if word in spam_data_pr.keys():
10             value = spam_data_pr.get(word) + 1
11             spam_data_pr.update({word: value})
12
13     print "=== ham_data: " + str(len(ham_data)) + " ==="
14
15     for word in ham_data:
16         if word not in ham_data_pr.keys():
17             ham_data_pr.update({word:1})
18
19         if word in ham_data_pr.keys():
20             value = ham_data_pr.get(word) + 1
21             ham_data_pr.update({word: value})
```

Listing 2: Making dictionary with unique words

## Chi-square calculations

For the classifier to better distinguish between two classes, specific words are more distinguishing than others. For example, the word "a" will be probably used in both classes because most sentences are made up with that word without any specific meaning attached. To know which words are a good base to train the model with, we made use of the chi-square test.

To structure the calculations, I made use of pandas which generates a table with columns and rows. The retrieved value will be saved in the "chi-square" columns. The data is saved as "data" and will be sorted ascending to view the most distinguishing words on top of the list.

```
1  """BUILDING CHI SQUARE LIST"""
2
3  def chi_square():
4      print "chi square calculations starting..."
5      #data = np.array([ham_data_pr.keys(), "h",0])
6      spam_data_pd = pd.DataFrame(columns=["data", "n_words", '
7      chi_score', "label" ])
8      spam_data_pd["data"] = spam_data_pr.keys()
9      spam_data_pd["label"] = "s"
10     spam_data_pd["chi_score"] = 0
11     spam_data_pd["n_words"] = spam_data_pr.values()
12
13     ham_data_pd = pd.DataFrame(columns=["data", "n_words", '
14     chi_score', "label" ])
15     ham_data_pd["data"] = ham_data_pr.keys()
16     ham_data_pd["label"] = "h"
17     ham_data_pd["chi_score"] = 0
18     ham_data_pd["n_words"] = ham_data_pr.values()
19     data = spam_data_pd.append(ham_data_pd)
20
21     n_ham_data = len(data[data['label']=="h"])
22     n_spam_data = len(data[data['label']=="s"])
23     data['1_1'] = np.where(data["data"].isin(data[data['label']=="h"
24     ""]["data"]),data["n_words"],0)
25     data['1_2'] = np.where(data["data"].isin(data[data['label']=="s"
26     ""]["data"]),data["n_words"],0)
27     data['2_1'] = n_ham_data - data['1_1']
28     data['2_2'] = n_spam_data - data['1_2']
29     data['C1'] = data['1_1'] + data['1_2']
30     data['C2'] = data['1_2'] + data['2_2']
31     data['W1'] = data['1_1'] + data['1_2']
32     data['W2'] = data['2_1'] + data['2_2']
33     data['N'] = data['C1'] + data['C2']
34     data['E1_1'] = (data['W1'] * data['C1']) / data['N']
35     data['E1_2'] = (data['W1'] * data['C2']) / data['N']
36     data['E2_1'] = (data['W2'] * data['C1']) / data['N']
37     data['E2_2'] = (data['W2'] * data['C2']) / data['N']
38
39     cols = ['1_1','1_2','2_1','2_2']
40     cols2 = ['E1_1','E1_2','E2_1','E2_2']
41     data['chi_score'] = sum(((data[c] - data[e])**2)/data[e] for c
42     in cols for e in cols2)
```

```
39     data = data.sort_values(by=["chi_score"], ascending=True)
40
41     data = data.reset_index()
42     return data
```

Listing 3: Chi-square calculations

## Discussion

During meeting with Randy, we agreed to finalize this assignment by this report. We had several meetings where Randy and I worked on bugs that appeared in the code. Those meetings were really helpful and taught me a lot about coding and how to build a programming structure.

I came quite far in the assignment. However still there is a part missing in the final part of this assignment. I struggled quite a lot with the probability part of the assignment (part 5 and 6). First couldn't fix the counter for the words to be able to calculate the word probabilities. After that, I already planned that the word probabilities would be saved into a list that would be called upon in the probability calculations of the emails. The classification would then be done by simply looking at which probability is higher, spam or ham. To be able to conclude things upon the accuracy of the classifier, I would compare the filename of the testing emails to the outcome of the classification. This would give me a final percentage of the accuracy. By using different sizes of chi-square lists, I would have expected different outcomes of the accuracy.

## Listings

1	Making dictionaries of all appearing words for each category . . .	2
2	Making dictionary with unique words . . . . .	3
3	Chi-square calculations . . . . .	4