

# DIT407 Introduction to data science and AI

## Assignment 3

Reynir Siik  
reynir.siik@gmail.com

Franco Zambon  
guszamfr@student.gu.se

2024-02-05

## 1 Problem 1: Spam & Ham

### 1.1 Introduction

*Abbreviated citation from assignment:*

The three files `20021010 easy ham.tar.bz2`, `20021010 hard ham.tar.bz2`, and `20021010 spam.tar.bz2` contain three sets of emails: *easy ham*, *hard ham*, and *spam*[1]. “Ham” refers to the benign, good emails that we want to read. “Spam”, on the other hand, refers to unsolicited bulk email, also known as junk. We are going to adopt a bag of words approach to train a Naïve Bayesian Classifier to tell the two classes of email apart.

### 1.2 A. Data exploration

*Explore the three datasets. Use your judgment as a human being to describe what makes you able to tell spam apart from ham. Also explore the differences between the two classes of ham (easy and hard ham). What makes them different?*

#### 1.2.1 Easy ham

The *Easy ham* dataset is at a quick glance personal email sent from specific mail address and addressed to one, or a few addresses that are enclosed in the header. It can also be a bulk mail, although

#### 1.2.2 Hard ham

The *Hard ham* dataset is often bulk email that are presumably wanted. It could be from a mailing list we follow, or a company we do business with. It can also be a mail about something that the spam-assassin will mark as spam due to its content, ie. mail describing fraudulent mails, and what to look out for. The recipient address could be a fake address, ie. not my address. Many mails have html content with a lot of hyperlinks.

#### 1.2.3 Spam

The *Spam* dataset is mail we haven’t asked for, and don’t want. It is mostly bulk email with undisclosed recipients. Many emails are offering porn or an offer ”to good to be true”.

### 1.2.4 Conclusion

Overall it is not trivial to look at someone else's email, and just by looking at it decide whether it is spam or not. The truth lies in the eye of the beholder.

## 1.3 B. Data splitting

*Perform an appropriate train-test split on the datasets. We will use the training sets to train a classifier, and evaluate its performance against the test sets.*

The code for splitting the data sets into train and test set are listed in A.1. The data sets are stored in csv files after being converted to ISO-8859-1 character encoding. The Unicode character u2019 - Right Single Quotation Mark, is substituted with ASCII Single Quotation Mark. This is presumed to have a negligible impact on training result.

## 2 Problem 2: Preprocessing

### 2.1 Introduction

*Abbreviated citation from assignment:*

A *bag of words* maps text documents into vectors of word counts. That is, each word is assigned a token number, and we count the number of occurrences of each unique word in each document. The vector then consists of these counts.

### 2.2 Our preprocessing

What we did in this assignment was splitting each of the mail content into a vector of words by using the `CountVectorizer`. After doing that we applied it to the training split by using the function `fit_transform()` and to the test split with the `transform()` function. Finally we transformed the vectors into a matrix using the `toarray()` function.

## 3 Problem 3: Easy Ham

### 3.1 Introduction

*Abbreviated citation from assignment:*

We are going to try out two different Naïve Bayesian Classifiers:

- The Multinomial Naïve Bayesian Classifier, and <sup>1</sup>
- The Bernoulli Naive Bayesian Classifier<sup>3</sup>. <sup>2</sup>

For this problem, use the set easy ham as your ham. Use the training sets of ham and spam to train an instance of both classifiers. Then, evaluate the classifier against the test set.

---

<sup>1</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.CountVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)

<sup>2</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html)

### 3.2 What we did

First of all we converted the  $y$  variables into integer (previously they were object) and after that we trained the Multinomial Bayes Naïve Classifier by using the `fit()` function. After that we used it to predict on the test. The result was pretty good, since we obtained these statistics:

Accuracy: 0.9748908296943232

Precision: 0.9751227237126167

Recall: 0.9748908296943232

Confusion Matrix:

$$\begin{bmatrix} 127 & 21 \\ 2 & 766 \end{bmatrix} \begin{matrix} 148 \\ 768 \end{matrix}$$

129 787

After that we use the Bernoulli Naive Bayesian Classifier. From that we obtained:

Accuracy: 0.9104803493449781

Precision: 0.9135972782624029

Recall: 0.9104803493449781

Confusion Matrix:

$$\begin{bmatrix} 70 & 78 \\ 4 & 764 \end{bmatrix} \begin{matrix} 148 \\ 768 \end{matrix}$$

74 842

## 4 Problem 4: Hard Ham

### 4.1 Introduction

*Abbreviated citation from assignment:*

Repeat the experiment of the previous problem, but use hard ham instead. Report the same values. Discuss the differences in results.

### 4.2 What we did

The procedure for this exercise was the same of the previous one. In this case we obtained for the Multinomial Bayes Naïve Classifier:

Accuracy: 0.9646017699115044

Precision: 0.9648439683278993

Recall: 0.9646017699115044

Confusion Matrix:

$$\begin{bmatrix} 144 & 2 \\ 6 & 74 \end{bmatrix} \begin{matrix} 146 \\ 80 \end{matrix}$$

150 76

For the Bernoulli Naive Bayesian Classifier:

Accuracy: 0.8938053097345132

Precision: 0.9054539575976868

Recall: 0.8938053097345132

Confusion Matrix:

$$\begin{bmatrix} 145 & 1 \\ 23 & 57 \end{bmatrix} \begin{matrix} 146 \\ 80 \end{matrix}$$

168   58

## 5 Discussion

What we can conclude from this analysis is that the Multinomial Bayes Naive Classifier work better in both cases.

In general the prediction are pretty good since they exceed 95% precision.

## References

- [1] J. Mason, *Spam assassin public mail corpus*. Made available on Canvas, From assignment handout on Canvas, 2004. [Online]. Available: <https://spamassassin.apache.org/old/publiccorpus/>.

## A Source code

This is the complete listing of the source code.

### A.1 Problem 1

Code for splitting the datasets into training and test sets.

```
1 # Reynir Siik, Franco Zambon
2 # DIT407 lp3 2024-02-10
3 # Assignment 3
4 # Problem 1
5 # Spam & Ham
6 # Splitting the datasets into training set and test set
7 #####
8 import pandas as pd
9 import tarfile
10 import re
11 from sklearn.model_selection import train_test_split
12 ## Read data from tar file #####
13 def read_data(fn, ds):
14     dataset = list()
15     with tarfile.open(fn) as tf:
16         for i in tf:
17             if i.isfile():
18                 f = tf.extractfile(i)
19                 if f:
20                     with f:
21                         try:
22                             s = f.read().decode('iso-8859-1')
23                         except:
24                             s = f.read().decode('utf-8')
25                             s = re.sub(u"[\u2018|\u2019]", "'", s)
26                     dataset.append((ds, s))
27     return dataset
28 #####
29 trainingSet = list()
30 testingSet = list()
31 set_list = ["easy_ham", "hard_ham", "spam"]
32 for set in set_list:
33     fname = "Uppgift03\\\\00_Uppg_Data\\\\20021010_"+set+".tar.bz2"
34     print(fname)
35     df = read_data(fname, set)
36     tr, te = train_test_split(df, train_size=0.7, random_state=42) # 23145)
37     trainingSet.extend(tr)
38     testingSet.extend(te)
39 # Create the pandas DataFrame
40 DFtrainingSet = pd.DataFrame(trainingSet, columns=['Status', 'Message'])
41 DFtestingSet = pd.DataFrame(testingSet, columns=['Status', 'Message'])
42 # Save to csv file
43 DFtrainingSet.to_csv(r"Uppgift03\\easy_ham\\20021010_training.csv")
44 DFtestingSet.to_csv(r"Uppgift03\\easy_ham\\20021010_test.csv")
```

## A.2 Problem 2 to 4

Code for making a bag of words and calculate statistics.

```

1  # Reynir Siik, Franco Zambon
2  # DIT407 lp3 2024-02-10
3  # Assignment 3
4  # Problem 2, 3, 4
5  # Spam & Ham
6  # Preprocessing and Classifying
7  #####
8  import pandas as pd
9  from sklearn.feature_extraction.text import CountVectorizer
10 from sklearn.model_selection import train_test_split
11 from sklearn.naive_bayes import MultinomialNB
12 from sklearn.naive_bayes import BernoulliNB
13 from sklearn.metrics import accuracy_score, precision_score, recall_score
14 from sklearn.metrics import confusion_matrix
15 ## #####
16 filelocation = "Uppgift03\\\\"easy_ham\\\\"
17 ## #####
18
19 df1=pd.read_csv(filelocation + '20021010_training.csv')
20 df2=pd.read_csv(filelocation + '20021010_test.csv')
21 df = pd.merge(df1, df2, how='outer')
22 print(df['Status'].value_counts())
23
24 # Removing the hard_ham from the dataset and converting the Status to a
25 # boolean value
26 df_eh = df[df['Status'] != 'hard_ham']
27 df_eh.loc[df_eh["Status"]=="spam", "Status"]=0
28 df_eh.loc[df_eh["Status"]=="easy_ham", "Status"]=1
29 df_eh = df_eh.drop(columns=['Unnamed: 0']).copy()
30 print(df_eh.head())
31
32 ## PROBLEM 3 #####
33 ## Splitting the dataset, vectorizing the Message and counting each word
34 ## and finally converting it to an array
35 print('*****_Problem_3_*****')
36 X = df_eh['Message']
37 y = df_eh['Status']
38 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
39                                                    random_state=42)
40 vectorizer = CountVectorizer(lowercase=True, stop_words='english')
41 X_traincv = vectorizer.fit_transform(X_train)
42 X_testcv = vectorizer.transform(X_test)
43 X_traincv.toarray()
44 X_testcv.toarray()
45
46 # Applying the Multinomial Bayes Naive Classifier
47 mnb = MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True)
48 y_train=y_train.astype('int')
49 y_test=y_test.astype('int')
50 mnb.fit(X_traincv,y_train)
51 pred_Mult = mnb.predict(X_testcv)
52 # Accuracy
53 accuracy = accuracy_score(y_test, pred_Mult)
54 # Precision

```

```
55 precision = precision_score(y_test, pred_Mult, average='weighted')
56 # Recall
57 recall = recall_score(y_test, pred_Mult, average='weighted')
58 print("Accuracy:", accuracy)
59 print("Precision:", precision)
60 print("Recall:", recall)
61 # Confusion matrix
62 conf_matrix = confusion_matrix(y_test, pred_Mult)
63 print("Confusion_Matrix:")
64 print(conf_matrix)
65 # Applying the Bernulli Naive Bayesian Classifier
66 clf = BernoulliNB()
67 clf.fit(X_traincv, y_train)
68 pred_Bern = clf.predict(X_testcv)
69 # Accuracy
70 accuracy = accuracy_score(y_test, pred_Bern)
71 # Precision
72 precision = precision_score(y_test, pred_Bern, average='weighted')
73 # Recall
74 recall = recall_score(y_test, pred_Bern, average='weighted')
75 print("Accuracy:", accuracy)
76 print("Precision:", precision)
77 print("Recall:", recall)
78 # Confusion matrix
79 conf_matrix = confusion_matrix(y_test, pred_Bern)
80 print("Confusion_Matrix:")
81 print(conf_matrix)
82
83 ## PROBLEM 4 #####
84 ## Removing the easy_ham from the dataset and converting the Status to
85 ## a boolean value
86 print('*****_Problem_4_*****')
87 df_hh = df[df['Status'] != 'easy_ham']
88 df_hh.loc[df_hh["Status"]=="spam", "Status"]=0
89 df_hh.loc[df_hh["Status"]=="hard_ham", "Status"]=1
90 df_hh = df_hh.drop(columns=['Unnamed: 0']).copy()
91 df_hh.head()
92 # Splitting the dataset, vectorizing the Message and counting each word
93 # and finally converting it to an array
94 X = df_hh['Message']
95 y = df_hh['Status']
96 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
97                                                    random_state=42)
98 vectorizer = CountVectorizer(lowercase=True, stop_words='english')
99 X_traincv = vectorizer.fit_transform(X_train)
100 X_testcv = vectorizer.transform(X_test)
101 X_traincv.toarray()
102 X_testcv.toarray()
103 # Applying the Multinomial Bayes Naive Classifier
104 mnb = MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True)
105 y_train=y_train.astype('int')
106 y_test=y_test.astype('int')
107 mnb.fit(X_traincv, y_train)
108 pred_Mult = mnb.predict(X_testcv)
109 # Accuracy
110 accuracy = accuracy_score(y_test, pred_Mult)
111 # Precision
112 precision = precision_score(y_test, pred_Mult, average='weighted')
```

```
113 # Recall
114 recall = recall_score(y_test, pred_Mult, average='weighted')
115 print("Accuracy:", accuracy)
116 print("Precision:", precision)
117 print("Recall:", recall)
118 # Confusion matrix
119 conf_matrix = confusion_matrix(y_test, pred_Mult)
120 print("Confusion_Matrix:")
121 print(conf_matrix)
122 # Applying the Bernulli Naive Bayesian Classifier
123 clf = BernoulliNB()
124 clf.fit(X_traincv, y_train)
125 pred_Bern = clf.predict(X_testcv)
126 # Accuracy
127 accuracy = accuracy_score(y_test, pred_Bern)
128 # Precision
129 precision = precision_score(y_test, pred_Bern, average='weighted')
130 # Recall
131 recall = recall_score(y_test, pred_Bern, average='weighted')
132 print("Accuracy:", accuracy)
133 print("Precision:", precision)
134 print("Recall:", recall)
135 # Confusion matrix
136 conf_matrix = confusion_matrix(y_test, pred_Bern)
137 print("Confusion_Matrix:")
138 print(conf_matrix)
```