# DAT565/DIT407 Assignment 3

Connell Hagen
connellh@student.chalmers.se

Måns Redin
mansre@student.chalmers.se

2024-09-19

## 1 Problem 1: Spam & Ham

### 1.1 Data Exploration

The spam emails frequently used sentences that started with a command. For example "Get this..., Experience this..., etc." The spam emails frequently uses exclamation marks, uppercase letters, and also had a lot of links in them. The ham emails are generally classified by them being directed to a specific person, involving sentimental expressions. The spam emails are often very explicit, and usually want you to buy something.

## 2 Problem 2: Preprocessing

The method applied to this problem was from the CountVectorizer library in Python. The three different operations were the fit, transform and the fit_transform. For the training dataset the fit_transform operation was applied, which is a combination of the fit and the transform operation. The fit operation goes though the training dataset and creates an array that has features for the different words in the emails. The transform operation is going through the dataset and counting the amount of appearances of each specific word, and stores it in the array constructed by the fit operation. Since the fit_transform operation was already made on the training dataset, which is a subset of the total dataset, the only operation needed for the testing dataset was the transform operation. The same applies to the type of emails, but for the labelEncoder operation. The labelEncoder operation makes binary classification columns for the different types of emails.

## 3 Problem 3: Easy Ham

In this part of the problem the accuracy, precision, and recall was desired. Equation 1, 2 and 3 defines these three measurements.

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \tag{1}$$

$$Precision = \frac{TP}{TP + FP} \tag{2}$$

$$Recall = \frac{TP}{TP + FN} \tag{3}$$

These were computed for the trained model, using both a Multinomial Naive Bayes classifier and a Bernoulli Naive Bayes classifier. Table 1 shows the accuracy, precision and recall for the two classifiers for the "easy ham" dataset.

|  | Multinomial | Bernoulli |
|---|---|---|
| Accuracy | 0.969 | 0.917 |
| Precision | 0.966 | 0.914 |
| Recall | 0.998 | 0.995 |

Table 1: The accuracy, precision and recall for the Multinomial and Bernoulli Naive Bayes classifier for the "easy ham" dataset

A confusion matrix was desired for the different Naive Bayes classifiers. Table 2 shows the scores of the Bernoulli Naive Bayes classifier.

|  | Predicted positive | Predicted negative | Total |
|---|---|---|---|
| Actual positive | 597 | 3 | 600 |
| Actual negative | 56 | 51 | 107 |
| Total | 653 | 54 | 707 |

Table 2: Confusion matrix for the Bernoulli Naive Bayes classifier for the "easy ham" dataset

Table 3 shows the results for the Multinomial Naive Bayes classifier.

|  | Predicted positive | Predicted negative | Total |
|---|---|---|---|
| Actual positive | 599 | 1 | 600 |
| Actual negative | 21 | 86 | 107 |
| Total | 620 | 87 | 707 |

Table 3: Confusion matrix for the Multinominal Naive Bayes classifier for the "easy ham" dataset

# 4 Problem 4: Hard Ham

The same approach as in the previous exercise, but now for the "hard ham" dataset instead of the "easy ham". Table 4 shows the accuracy, precision and recall values for the "hard ham" dataset.

The confusion matrices for the two classifiers for the "hard ham" dataset can be seen in table 5 and table 6.

## 4.1 Discussion

The overall accuracy and recall for the hard ham dataset were worse than for the easy ham dataset on both model types. It was harder for the models to tell what

|  | Multinomial | Bernoulli |
|---|---|---|
| Accuracy | 0.950 | 0.875 |
| Precision | 0.979 | 0.972 |
| Recall | 0.870 | 0.648 |

Table 4: The accuracy, precision and recall for the Multinomial and Bernoulli Naive Bayes classifier for the "hard ham" dataset

was hard ham than easy ham because the hard ham had features like HTML and advertisements within commonly. The recall (See equation 3) especially was lower, which can be explained by the fact that there were a lot of false negatives. This is probably because the hard ham was frequently newsletters or other mass distributed emails that share many characteristics with spam emails. The precision, in contrast, actually was better for the hard ham dataset. This is because there were not many false positives. This can be explained by the fact that there are many spam emails that are obviously spam, but among the hard ham emails, there were not many that were obviously not spam. This causes the mistakes to shift towards false negatives instead of false positives when it comes to identifying spam.

|  | Predicted positive | Predicted negative | Total |
|---|---|---|---|
| Actual positive | 35 | 19 | 54 |
| Actual negative | 1 | 105 | 106 |
| Total | 36 | 122 | 160 |

Table 5: Confusion matrix for the Bernoulli Naive Bayes classifier for the "hard ham" dataset

|  | Predicted positive | Predicted negative | Total |
|---|---|---|---|
| Actual positive | 47 | 7 | 54 |
| Actual negative | 1 | 105 | 106 |
| Total | 48 | 112 | 160 |

Table 6: Confusion matrix for the Multinomial Naive Bayes classifier for the "hard ham" dataset

# A Code

```
1  import numpy as np
2  import pandas as pd
3  import os
4  import re
5  from sklearn.model_selection import train_test_split
6  from sklearn.feature_extraction.text import
       CountVectorizer
7  from sklearn.preprocessing import LabelEncoder
8  from sklearn.naive_bayes import BernoulliNB
9  from sklearn.naive_bayes import MultinomialNB
10
11 # takes a directory of files, and a label, and returns
       a
12 # list of 2-tuples of each file's contents, and the
       given label
13 def get_email_content(dir, label):
14     files = []
15     for filename in os.listdir(dir):
16         path = os.path.join(dir, filename)
17         if os.path.isfile(path):
18             try:
19                 with open(path, encoding="utf-8") as f
                      :
20                     files += [(clean_email(f.read()),
                          label)]
21             except:
22                 pass
23     return files
24
25 # if we want to clean the emails we can modify this
       function
26 # however, our proposed way of removing the headers
       actually
27 # decreased the accuracy of the models
28 def clean_email(content):
29     # remove_header = '\\n\\n'
30     # email_parts = re.split(remove_header, content)
31     # if (len(email_parts) == 1):
32     #     return email_parts[0]
33     # else:
34     #     return email_parts[1]
35     return content
36
37 def get_confusion_values(y_test_inv, y_pred_inv, label
       ):
38     tp = ((y_test_inv == label) & (y_pred_inv == label
           )).sum()
```

```
39        fp = ((y_test_inv != label) & (y_pred_inv == label
             )).sum()
40        fn = ((y_test_inv == label) & (y_pred_inv != label
             )).sum()
41        tn = ((y_test_inv != label) & (y_pred_inv != label
             )).sum()
42        return tp, fp, fn, tn
43
44   def get_apr_values(tp, fp, fn, tn):
45        accuracy = (tp + tn) / (tp + fp + fn + tn)
46        precision = tp / (tp + fp)
47        recall = tp / (tp + fn)
48        return accuracy, precision, recall
49
50   SEED = 12345678
51   easy_ham_dir = 'assignment_3/easy_ham'
52   hard_ham_dir = 'assignment_3/hard_ham'
53   spam_dir = 'assignment_3/spam'
54
55   easy_ham_files = get_email_content(easy_ham_dir, "Ham"
         )
56   hard_ham_files = get_email_content(hard_ham_dir, "Ham"
         )
57   spam_files = get_email_content(spam_dir, "Spam")
58
59   # to switch between the easy_ham and hard_ham dataset,
60   # swap out either easy_ham_files or hard_ham_files for
         the
61   # other here
62   records = hard_ham_files + spam_files
63
64   df = pd.DataFrame.from_records(records, columns=["
         Email", "Type"])
65   df_train, df_test = train_test_split(df, random_state=
         SEED, test_size=0.25)
66
67   cv = CountVectorizer()
68   x_train = cv.fit_transform(df_train["Email"])
69   x_test = cv.transform(df_test["Email"])
70
71   le = LabelEncoder()
72   y_train = le.fit_transform(df_train["Type"])
73   y_test = le.transform(df_test["Type"])
74   y_test_inv = le.inverse_transform(y_test)
75
76   bnb = BernoulliNB()
77   bnb.fit(x_train, y_train)
78   y_pred_bnb = bnb.predict(x_test)
79   y_pred_bnb_inv = le.inverse_transform(y_pred_bnb)
```

```
80  acc_bnb , pre_bnb , rec_bnb = get_apr_values (*
        get_confusion_values ( y_test_inv , y_pred_bnb_inv , "
        Ham"))
81
82  mnb = MultinomialNB ()
83  mnb.fit( x_train , y_train )
84  y_pred_mnb = mnb.predict ( x_test )
85  y_pred_mnb_inv = le.inverse_transform ( y_pred_mnb )
86  acc_mnb , pre_mnb , rec_mnb = get_apr_values (*
        get_confusion_values ( y_test_inv , y_pred_mnb_inv , "
        Ham"))
```