

# DAT565/DIT407 Assignment 3

David Duong  
thuan@chalmers.se

Savinjith Walisadeera  
savwal@chalmers.se

October 1, 2024

## Problem 1

A)

We can differentiate between easy ham, hard ham and spam by reading the email. Easy ham is usually a email conversation between two people. The difference between hard ham and easy ham is that hard ham emails are not necessarily written by a human. There can be part or all of it which is copied or a generated response. Spam emails are all generic responses or emails not specified to the receiver, they often contain words that are uncommon in conversations and links.

## Problem 2

We split the data with `train-test-split` into training and testing data with the ratio 75:25. We organize all the training emails into vectors using `fit-transform()` and the test emails with `transform()`.

## Problem 3

Multinomial	Pred. Positive	Pred. Negative	Sum
Actual Positive	641	1	642
Actual Negative	27	94	121
Sum	668	95	

Table 1: Easy ham confusion matrix for multinomial naive Bayesian classifier

Bernoulli	Pred. Positive	Pred. Negative	Sum
Actual Positive	642	0	642
Actual Negative	69	52	121
Sum	711	52	

Table 2: Easy ham confusion matrix for Bernoulli naive bayesian classifier

	Accuracy	Precision	Recall
Multinomial	0.9633	0.9595	0.9984
Bernoulli	0.9095	0.9029	1.0000

Table 3: Classifier metrics for easy ham + spam

## Problem 4

Multinomial	Pred. Positive	Pred. Negative	Sum
Actual Positive	56	8	64
Actual Negative	3	121	124
Sum	59	129	

Table 4: Hard ham confusion matrix for Multinomial naive Bayesian classifier

Bernoulli	Pred. Positive	Pred. Negative	Sum
Actual Positive	43	21	64
Actual Negative	2	122	124
Sum	45	143	

Table 5: Hard ham confusion matrix for Bernoulli naive Bayesian classifier

	Accuracy	Precision	Recall
Multinomial	0.9414	0.9491	0.8750
Bernoulli	0.8776	0.9556	0.6719

Table 6: Classifier metrics for hard ham + spam

After changing all easy ham emails to hard ham the accuracy, precision and recall values all decrease. This is because hard ham contains words that are more used in spam emails, which makes it harder for the classifiers to differentiate between them. This causes the metrics to decrease between the two.

## A Code

```

1 import numpy as np
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import
  CountVectorizer
4 import os
5 from pathlib import Path
6
7 #Extract the files
8 file_list_easy = []
9 file_list_hard = []
10 types_list_easy = []

```

```

11 types_list_hard = []
12 tar_list_easy = []
13 tar_list_hard = []
14
15 directory = './easy_ham'
16 for file in os.listdir(directory):
17     content = Path(directory, file).read_text(encoding
18         = 'latin-1')
19     file_list_easy.append(content)
20     types_list_easy.append('ham')
21
22 directory = './hard_ham'
23 for file in os.listdir(directory):
24     content = Path(directory, file).read_text(encoding
25         = 'latin-1')
26     file_list_hard.append(content)
27     types_list_hard.append('ham')
28
29 directory = './spam'
30 for file in os.listdir(directory):
31     content = Path(directory, file).read_text(encoding
32         = 'latin-1')
33     file_list_easy.append(content)
34     types_list_easy.append('spam')
35     file_list_hard.append(content)
36     types_list_hard.append('spam')
37
38 #Split the easy data 75-25
39 file_train, file_test, type_train, type_test =
40     train_test_split(file_list_easy, types_list_easy,
41         test_size=0.25, random_state=6)
42
43 #Vectorize
44 vectorizer = CountVectorizer()
45 text_vector_train = vectorizer.fit_transform(
46     file_train)
47 text_vector_test = vectorizer.transform(file_test)
48
49 #Train our models
50 from sklearn.naive_bayes import MultinomialNB
51 model_multinomial = MultinomialNB()
52 model_multinomial.fit(text_vector_train, type_train)
53
54 from sklearn.naive_bayes import BernoulliNB
55 model_bernoulli = BernoulliNB()
56 model_bernoulli.fit(text_vector_train, type_train)
57
58 #Evaluate the models
59 multinomial_TP = 0
60 multinomial_FP = 0
61 multinomial_TN = 0

```

```

55 multinomial_FN = 0
56
57 for i in range(np.shape(text_vector_test)[0]):
58     prediction = model_multinomial.predict(
59         text_vector_test[i])[0]
60     answer = type_test[i]
61
62     if prediction == 'ham':
63         if answer == 'ham':
64             multinomial_TP += 1
65         elif answer == 'spam':
66             multinomial_FP += 1
67     elif prediction == 'spam':
68         if answer == 'spam':
69             multinomial_TN += 1
70         elif answer == 'ham':
71             multinomial_FN += 1
72
73
74 multinomial_accuracy = (multinomial_TN+multinomial_TP)
75                        /(multinomial_TP+multinomial_FP+multinomial_FN+
76                        multinomial_TN)
77
78 multinomial_precision = multinomial_TP/(multinomial_TP
79 +multinomial_FP)
80
81 multinomial_recall = multinomial_TP/(multinomial_TP+
82 multinomial_FN)
83
84
85 bernoulli_TP = 0
86 bernoulli_FP = 0
87 bernoulli_TN = 0
88 bernoulli_FN = 0
89
90 for i in range(np.shape(text_vector_test)[0]):
91     prediction = model_bernoulli.predict(
92         text_vector_test[i])[0]
93     answer = type_test[i]
94
95     if prediction == 'ham':
96         if answer == 'ham':
97             bernoulli_TP += 1
98         elif answer == 'spam':
99             bernoulli_FP += 1
100     elif prediction == 'spam':
101         if answer == 'spam':
102             bernoulli_TN += 1
103         elif answer == 'ham':
104             bernoulli_FN += 1
105
106

```

```

198 bernoulli_accuracy = (bernoulli_TN+bernoulli_TP)/(
    bernoulli_TP+bernoulli_FP+bernoulli_FN+bernoulli_TN
    )
199 bernoulli_precision = bernoulli_TP/(bernoulli_TP+
    bernoulli_FP)
200 bernoulli_recall = bernoulli_TP/(bernoulli_TP+
    bernoulli_FN)
201
202
203 print('accuracy, precision, recall')
204 print('Multinomial:', multinomial_accuracy,
    multinomial_precision, multinomial_recall)
205 print('Bernoulli:', bernoulli_accuracy,
    bernoulli_precision, bernoulli_recall)
206
207 print(multinomial_TP, multinomial_FN, multinomial_FP,
    multinomial_TN)
208 print(bernoulli_TP, bernoulli_FN, bernoulli_FP,
    bernoulli_TN)
209
210 #Repeat experiment with hard_ham
211 file_train, file_test, type_train, type_test =
    train_test_split(file_list_hard, types_list_hard,
    test_size=0.25, random_state=6)
212
213 vectorizer = CountVectorizer()
214 text_vector_train = vectorizer.fit_transform(
    file_train)
215 text_vector_test = vectorizer.transform(file_test)
216
217 #Train our models
218 from sklearn.naive_bayes import MultinomialNB
219 model_multinomial = MultinomialNB()
220 model_multinomial.fit(text_vector_train, type_train)
221
222 from sklearn.naive_bayes import BernoulliNB
223 model_bernoulli = BernoulliNB()
224 model_bernoulli.fit(text_vector_train, type_train)
225
226 #Evaluate the models
227 multinomial_TP = 0
228 multinomial_FP = 0
229 multinomial_TN = 0
230 multinomial_FN = 0
231
232 for i in range(np.shape(text_vector_test)[0]):
233     prediction = model_multinomial.predict(
        text_vector_test[i])[0]
234     answer = type_test[i]
235

```

```

136     if prediction == 'ham':
137         if answer == 'ham':
138             multinomial_TP += 1
139         elif answer == 'spam':
140             multinomial_FP += 1
141     elif prediction == 'spam':
142         if answer == 'spam':
143             multinomial_TN += 1
144         elif answer == 'ham':
145             multinomial_FN += 1
146
147
148
149 multinomial_accuracy = (multinomial_TN+multinomial_TP)
150                        /(multinomial_TP+multinomial_FP+multinomial_FN+
151                          multinomial_TN)
152 multinomial_precision = multinomial_TP/(multinomial_TP
153 +multinomial_FP)
154 multinomial_recall = multinomial_TP/(multinomial_TP+
155 multinomial_FN)
156
157
158 bernoulli_TP = 0
159 bernoulli_FP = 0
160 bernoulli_TN = 0
161 bernoulli_FN = 0
162
163 for i in range(np.shape(text_vector_test)[0]):
164     prediction = model_bernoulli.predict(
165         text_vector_test[i])[0]
166     answer = type_test[i]
167
168     if prediction == 'ham':
169         if answer == 'ham':
170             bernoulli_TP += 1
171         elif answer == 'spam':
172             bernoulli_FP += 1
173     elif prediction == 'spam':
174         if answer == 'spam':
175             bernoulli_TN += 1
176         elif answer == 'ham':
177             bernoulli_FN += 1
178
179
180 bernoulli_accuracy = (bernoulli_TN+bernoulli_TP)/(
181     bernoulli_TP+bernoulli_FP+bernoulli_FN+bernoulli_TN
182 )
183 bernoulli_precision = bernoulli_TP/(bernoulli_TP+
184 bernoulli_FP)
185 bernoulli_recall = bernoulli_TP/(bernoulli_TP+
186 bernoulli_FN)
187

```

```

177
178 print('accuracy,precision,recall')
179 print('Multinomial:', multinomial_accuracy,
        multinomial_precision, multinomial_recall)
180 print('Bernoulli:', bernoulli_accuracy,
        bernoulli_precision, bernoulli_recall)
181
182 print(multinomial_TP, multinomial_FN, multinomial_FP,
        multinomial_TN)
183 print(bernoulli_TP, bernoulli_FN, bernoulli_FP,
        bernoulli_TN)

```