# DAT565/DIT407 Assignment 3

Saif Sayed
gussayedfa@student.gu.se

Gona Ibrahim Abdulrahman
gusibrgo@student.gu.se

2024-04-18

This is a report for assignment 3 for the course *Introduction to Data Science & AI* from Chalmers and Gothenburg University.

## Problem 1: Spam and Ham

## A    Data exploration

When analyzing email contents, it's crucial to differentiate between spam (unsolicited or promotional emails) and ham (legitimate personal or work-related correspondence). Spam emails often contain exaggerated language, excessive capitalization, and phrases like "Congratulations!" or "Don't miss your chance.". They originate from unfamiliar email addresses and may have grammatical errors. In contrast, ham emails are straightforward, relevant, and come from known contacts. There are two types of ham: easy ham, which includes simple work-related or personal messages, and hard ham, which involves more specialized content. Remember that spam expects no specific action, while ham often requires a reply or some form of follow-up.

**Differences Between Easy Ham and Hard Ham:**

**Easy Ham:** These straightforward emails are likely related to work, personal matters, or subscriptions. They have minimal noise or clutter and may include relevant attachments. Usually written in first person.

**Hard Ham:** Hard ham emails seem to be closer to spam as it contains promotions and newsletter from websites and companies the user has subscribed to. Hard ham emails tend to be presented in a more complex format and may require more careful evaluation to distinguish them from actual spam.

## B    Data splitting

We extracted email data from three archives: *easy_ham*, *hard_ham*, and *spam*. And then we split the data by performing a train-test split on each dataset (easy/hard ham and spam) and used the training sets to train the model and evaluated its performance against test sets.

**To model Training and Evaluation:** We selected an appropriate classification model (e.g., Naive Bayes ) and trained the model using training data $(X\_train, X1\_train)$ and corresponding labels $(Y\_train, Y1\_train)$. Then evaluated model performance on test data $(X\_test, X1\_test)$ and corresponding labels $(Y\_test, Y1\_test)$ using metrics like accuracy and precision.

By splitting the data and training a classifier, we aim to create an effective email classification model. Evaluation results will guide us in understanding how well the model performs across different email types.

Our source code can be found in Appendix A of this document.

# Problem 2: Preprocessing

Our goal is to transform a collection of emails into a matrix format. Each row in this matrix corresponds to an email, and each column represents the count of unique words within that email. This process allows us to map text documents (emails) into numerical vectors based on word frequencies .

We combined emails from two categories: Easy Ham and Spam. This simplification streamlines subsequent processes (problem 3).

We created an instance of CountVectorizer, as a tool for converting text data into numerical representations and then using the fit_transform method, we learned the vocabulary (unique words) from the entire email dataset and transformed the emails into a matrix representation of word counts.

By using the CountVectorizer, we successfully converted the emails into a matrix of word counts and this matrix will serve as input for training and evaluating our email classification model.

# Problem 3: Easy Ham

We combined emails from two categories: Easy Ham and Spam and then assigned labels to each dataset.

Later we instantiated two classifiers: Multinomial Naive Bayes (MNB) and Bernoulli Naive Bayes (BNB).

And then trained both classifiers using the training data and evaluated their performance against the test set mentioned in **Problem 1B**. Now both classifiers have been successfully trained and evaluated and the confusion matrices provide valuable insights into their performance.

**Table 1** shows the accuracy, precision and recall score of the Multinomial and Bernoulli Naive Bayes Classifier.

**Figure 1** shows the confusion matrix of Easy ham and Spam emails using

the Multinomial Naive Bayes Classifier.

**Figure 2** shows the confusion matrix of Easy ham and Spam emails using the Bernoulli Naive Bayes Classifier.

Table 1: Classification Report (Easy & Spam emails)

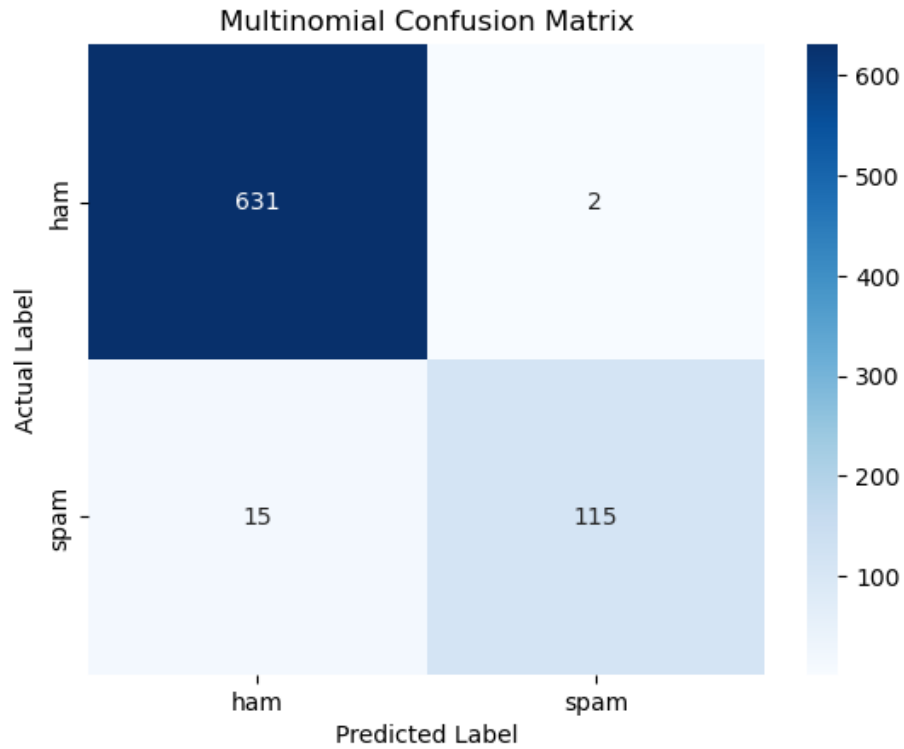| Stats | Multinomial | Bernoulli |
|---|---|---|
| Easy Ham emails | 2551 | 2551 |
| Spam emails | 501 | 501 |
| Test size | 763 | 763 |
| Accuracy | 0.9777195281782438 | 0.9043250327653998 |
| Precision (ham) | 0.9767801857585139 | 0.9011461318051576 |
| Recall Score (ham) | 0.9968404423380727 | 0.9936808846761453 |



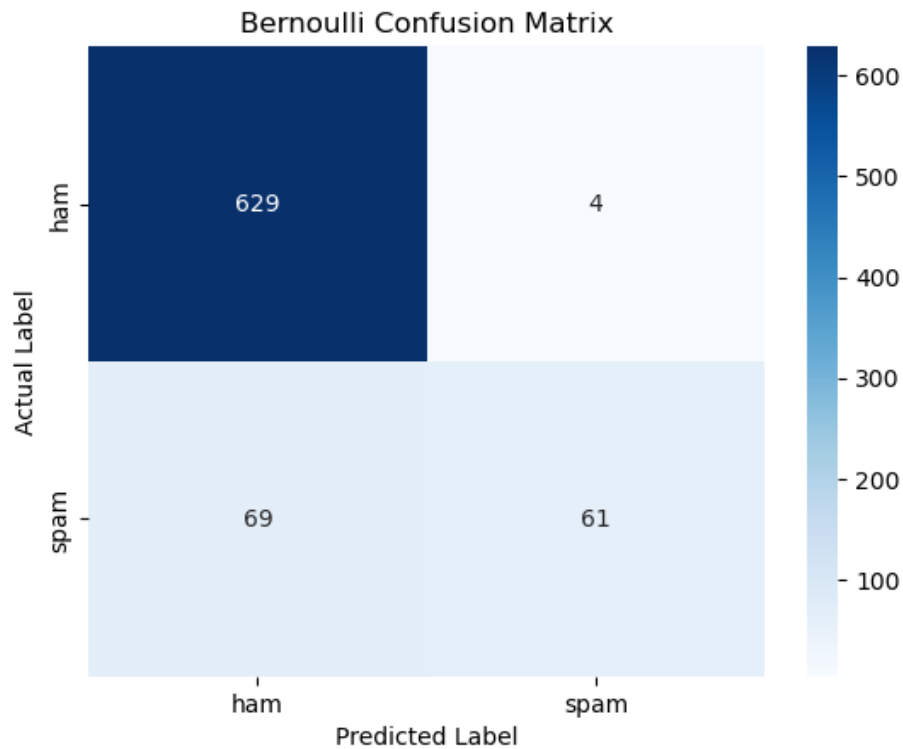Figure 1: Multinomial Confusion Matrix of Easy Ham and Spam emails

Figure 2: Bernoulli Confusion Matrix of Easy Ham and Spam emails

## Problem 4: Hard Ham

We applied the same solution as in problem 3, but this time we changed it for 'Hard ham' instead of 'Easy ham'.

**Table 2** shows the accuracy, precision and recall score of the Multinomial and Bernoulli Naive Bayes Classifier.

**Figure 3** shows the confusion matrix of Hard ham and Spam emails using the Multinomial Naive Bayes Classifier.

**Figure 4** shows the confusion matrix of Hard ham and Spam emails using the Bernoulli Naive Bayes Classifier.

Table 2: Classification Report (Hard & Spam emails)

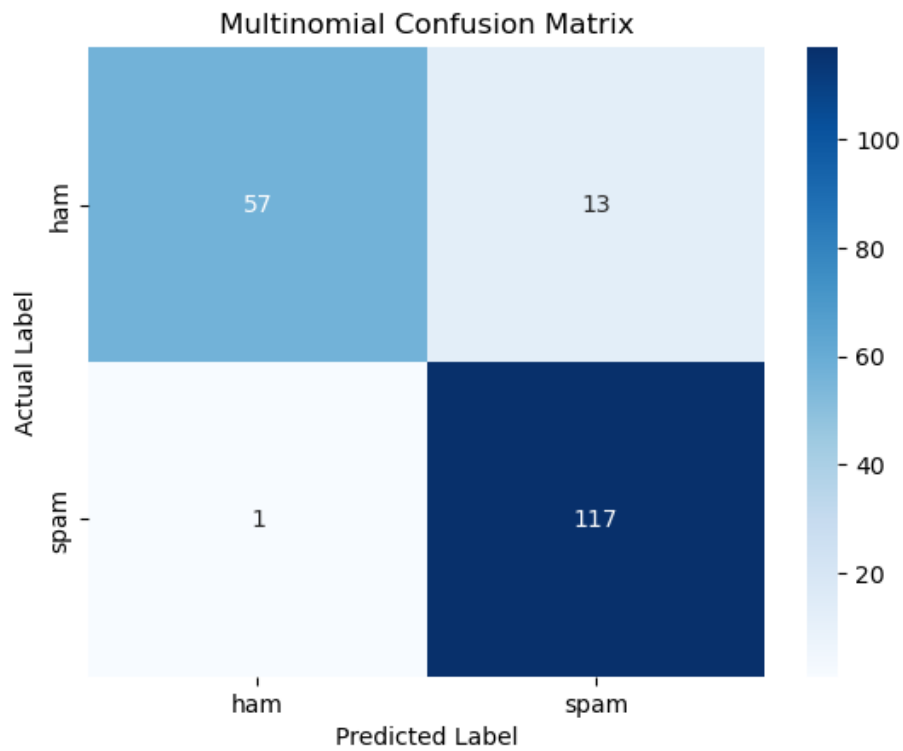| Stats | Multinomial | Bernoulli |
|---|---|---|
| Hard Ham emails | 250 | 250 |
| Spam emails | 501 | 501 |
| Test size | 188 | 188 |
| Accuracy | 0.925531914893617 | 0.8723404255319149 |
| Precision (ham) | 0.9827586206896551 | 0.96 |
| Recall Score (ham) | 0.8142857142857143 | 0.6857142857142857 |



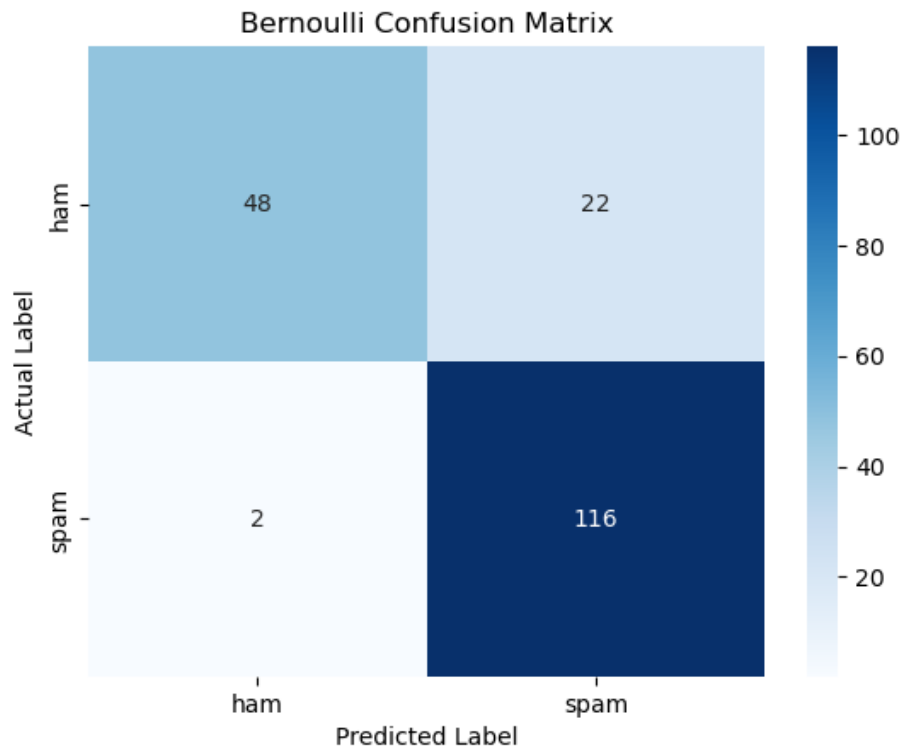Figure 3: Multinomial Confusion Matrix of Hard Ham and Spam emails

Figure 4: Bernoulli Confusion Matrix of Hard Ham and Spam emails

## Discussion

Both the multinomial and Bernoulli classifiers have higher accuracy and precision in classifying 'Easy ham' and 'Spam' emails compared to classifying 'Hard ham' and 'Spam' emails. The recall score for 'Hard ham' and 'Spam' emails indicates that both classifiers have a lower proportion of actual ham emails predicted as 'Hard ham' emails, in contrast to the recall score for 'Easy ham' and 'Spam' emails. Consequently, both classifiers yield more false positive results for 'Hard ham' emails. Although the results might be biased since the classifiers were trained with a larger dataset of 'Easy ham' emails compared to the dataset of 'Hard ham' emails. Overall, the Multinomial Naive Bayes classifier performs well in terms of accuracy, precision, and recall for both datasets.

## A  Python code

This is the code we used to classify between Easy ham emails, Hard ham emails and Spam emails .

```
1  import matplotlib.pyplot as plt
2  import tarfile
3  import os
```

```python
4   import chardet # to identify encoding format
5   from sklearn.model_selection import train_test_split
6   from sklearn.feature_extraction.text import
        CountVectorizer
7   from sklearn.naive_bayes import MultinomialNB,
        BernoulliNB
8   from sklearn.metrics import accuracy_score,
        precision_score, recall_score, confusion_matrix
9   import seaborn as sns
10
11  # B. Data splitting and Preprocessing
12
13  current_directory = os.getcwd()
14  extract_dir = current_directory
15
16  tar_easy_ham = os.path.join(current_directory, "
        20021010_easy_ham.tar.bz2")
17  tar_hard_ham = os.path.join(current_directory, "
        20021010_hard_ham.tar.bz2")
18  tar_spam = os.path.join(current_directory, "20021010
        _spam.tar.bz2")
19
20  with tarfile.open(tar_easy_ham, "r") as tar:
21      tar.extractall(path=extract_dir)
22  with tarfile.open(tar_hard_ham, "r") as tar:
23      tar.extractall(path=extract_dir)
24  with tarfile.open(tar_spam, "r") as tar:
25      tar.extractall(path=extract_dir)
26
27  easy_ham_dir = os.path.join(extract_dir, "easy_ham")
28  hard_ham_dir = os.path.join(extract_dir, "hard_ham")
29  spam_dir = os.path.join(extract_dir, "spam")
30
31  def read_emails_from_directory(directory):
32      emails = []
33      for filename in os.listdir(directory):
34          file_path = os.path.join(directory, filename)
35          with open(file_path, "rb") as file: # open in
                binary mode to detect encoding
36              raw_content = file.read()
37              detected_encoding = chardet.detect(
                    raw_content)['encoding']
38              try:
39                  content = raw_content.decode(
                        detected_encoding) # decode
                        according to detected encoding
40              except UnicodeDecodeError:
41                  content = raw_content.decode(
                        detected_encoding, errors='replace'
                        ) # else identify the error and
```

```
                       decode using another encoding.
42              emails.append(content)
43      return emails
44
45  easy_ham_emails = read_emails_from_directory(
        easy_ham_dir)
46  hard_ham_emails = read_emails_from_directory(
        hard_ham_dir)
47  spam_emails = read_emails_from_directory(spam_dir)
48
49  easy_and_spam_emails = easy_ham_emails + spam_emails
50  hard_and_spam_emails = hard_ham_emails + spam_emails
51
52  vectorizer = CountVectorizer()
53  easy_spam_email_vectors_train = vectorizer.
        fit_transform(easy_and_spam_emails)
54  hard_spam_email_vectors_train = vectorizer.
        fit_transform(hard_and_spam_emails)
55
56  label_easy_spam = ['ham'] * len(easy_ham_emails) + ['
        spam'] * len(spam_emails)
57  label_hard_spam = ['ham'] * len(hard_ham_emails) + ['
        spam'] * len(spam_emails)
58
59  #split data for easy ham and spam combination ( 25%
        test emails and 75% emails for training)
60  X_train, X_test, Y_train, Y_test = train_test_split(
        easy_spam_email_vectors_train, label_easy_spam,
        test_size=0.25, random_state=42)
61  #split data for hard ham and spam combination ( 25%
        test emails and 75% emails for training)
62  X1_train, X1_test, Y1_train, Y1_test =
        train_test_split(hard_spam_email_vectors_train,
        label_hard_spam, test_size=0.25, random_state=42)
63
64  # Problem 3: Easy Ham
65
66  multinomial_classifier = MultinomialNB()
67  bernoulli_classifier = BernoulliNB()
68
69  multinomial_classifier.fit(X_train, Y_train)
70
71  bernoulli_classifier.fit(X_train, Y_train)
72
73  multinomial_predictions = multinomial_classifier.
        predict(X_test)
74  bernoulli_predictions = bernoulli_classifier.predict(
        X_test)
75
```

```
76  multinomial_accuracy = accuracy_score(Y_test,
        multinomial_predictions)
77  bernoulli_accuracy = accuracy_score(Y_test,
        bernoulli_predictions)
78
79  multinomial_precision = precision_score(Y_test,
        multinomial_predictions, pos_label='ham')
80  bernoulli_precision = precision_score(Y_test,
        bernoulli_predictions, pos_label='ham')
81
82  multinomial_recall = recall_score(Y_test,
        multinomial_predictions, pos_label='ham')
83  bernoulli_recall = recall_score(Y_test,
        bernoulli_predictions, pos_label='ham')
84
85  multinomial_confusion = confusion_matrix(Y_test,
        multinomial_predictions)
86  bernoulli_confusion = confusion_matrix(Y_test,
        bernoulli_predictions)
87
88  print("Size_of_easy-ham_emails:", len(easy_ham_emails)
        )
89  print("Size_of_hard_ham_emails:", len(hard_ham_emails)
        )
90  print("Size_of_spam_emails:", len(spam_emails))
91  print()
92  # Print the results
93  print("Easy_Ham_and_Spam")
94  print("-------------------")
95  print("Size_of_test_set:", len(Y_test))
96  print("Multinomial␣Naive␣Bayes␣Classifier:")
97  print("Accuracy:", multinomial_accuracy)
98  print("Precision:", multinomial_precision)
99  print("Recall:", multinomial_recall)
100 print("Confusion_Matrix:")
101 # Create a heatmap of the confusion matrix
102 sns.heatmap(multinomial_confusion, annot=True, cmap='
        Blues', fmt='d', xticklabels=['ham', 'spam'],
103              yticklabels=['ham', 'spam'])
104 # Set the title and axis labels
105 plt.title('Multinomial_Confusion_Matrix')
106 plt.xlabel('Predicted_Label')
107 plt.ylabel('Actual-Label')
108 plt.show()
109 print()
110 print("Size_of_test_set:", len(Y_test))
111 print("Bernoulli_Naive_Bayes_Classifier:")
112 print("Accuracy:", bernoulli_accuracy)
113 print("Precision:", bernoulli_precision)
114 print("Recall:", bernoulli_recall)
```

```python
115  print("Confusion_Matrix:")
116
117  sns.heatmap(bernoulli_confusion, annot=True, cmap='
         Blues', fmt='d', xticklabels=['ham', 'spam'],
118                  yticklabels=['ham', 'spam'])
119  plt.title('Bernoulli_Confusion_Matrix')
120  plt.xlabel('Predicted_Label')
121  plt.ylabel('Actual_Label')
122  plt.show()
123  print()
124
125  # Problem 4: Hard Ham
126
127  multinomial_classifier.fit(X1_train, Y1_train)
128  bernoulli_classifier.fit(X1_train, Y1_train)
129
130  multinomial_predictions = multinomial_classifier.
         predict(X1_test)
131  bernoulli_predictions = bernoulli_classifier.predict(
         X1_test)
132
133  multinomial_accuracy = accuracy_score(Y1_test,
         multinomial_predictions)
134  bernoulli_accuracy = accuracy_score(Y1_test,
         bernoulli_predictions)
135
136  multinomial_precision = precision_score(Y1_test,
         multinomial_predictions, pos_label='ham')
137  bernoulli_precision = precision_score(Y1_test,
         bernoulli_predictions, pos_label='ham')
138
139  multinomial_recall = recall_score(Y1_test,
         multinomial_predictions, pos_label='ham')
140  bernoulli_recall = recall_score(Y1_test,
         bernoulli_predictions, pos_label='ham')
141
142  multinomial_confusion = confusion_matrix(Y1_test,
         multinomial_predictions)
143  bernoulli_confusion = confusion_matrix(Y1_test,
         bernoulli_predictions)
144  # Print the results
145  print("Hard_Ham_and_Spam")
146  print("------------------")
147  print("Size_of_test_set:", len(Y1_test))
148  print("Multinomial_Naive_Bayes_Classifier:")
149  print("Accuracy:", multinomial_accuracy)
150  print("Precision:", multinomial_precision)
151  print("Recall:", multinomial_recall)
152  print("Confusion_Matrix:")
153  # Create a heatmap of the confusion matrix
```

```
154  sns.heatmap(multinomial_confusion, annot=True, cmap='
         Blues', fmt='d', xticklabels=['ham', 'spam'],
155                     yticklabels=['ham', 'spam'])
156  plt.title('Multinomial_Confusion_Matrix')
157  plt.xlabel('Predicted_Label')
158  plt.ylabel('Actual_Label')
159  plt.show()
160  print()
161  print("Size_of_test_set:", len(Y1_test))
162  print("\nBernoulli_Naive_Bayes_Classifier:")
163  print("Accuracy:", bernoulli_accuracy)
164  print("Precision:", bernoulli_precision)
165  print("Recall:", bernoulli_recall)
166  print("Confusion_Matrix:")
167  sns.heatmap(bernoulli_confusion, annot=True, cmap='
         Blues', fmt='d', xticklabels=['ham', 'spam'],
168                     yticklabels=['ham', 'spam'])
169  plt.title('Bernoulli_Confusion_Matrix')
170  plt.xlabel('Predicted_Label')
171  plt.ylabel('Actual_Label')
172
173  plt.show()
```