

DAT565/DIT407 Assignment 3

Ola Bratt
ola.bratt@gmail.com

Patrick Attimont
patrickattimont@gmail.com

2024-02-xx

This paper is addressing the assignment 3 study queries within the *Introduction to Data Science & AI* course, DIT407 at the University of Gothenburg and DAT565 at Chalmers. The main source of information for this project is derived from the lectures and Skiena [1]. Assignment 3 is about text classification and the use of correct data splitting and encoding handling.

Problem 1: Spam and Ham

A. Data exploration

B. Data splitting

Since we have a large dataset, we can use the `train_test_split` function from the `sklearn.model_selection`. With a smaller dataset it would be better to use cross-validation to avoid overfitting.

```
X_train, X_test, y_train, y_test =  
train_test_split(email_matrix, labels, test_size=0.2)
```

Problem 2: Preprocessing

The "bag of words" model is a basic and intuitive way to analyze and compare documents based on their textual content. However, it does not consider the context or the order of words, which can limit its effectiveness in capturing the semantics and meaning of the text.

Problem 3: Easy Ham

To calculate the precision, recall, accuracy and confusion matrix, we use the following code (These functions are available in the `sklearn.metrics` package):

```
accuracy = accuracy_score(y_test, y_pred)  
precision = precision_score(y_test, y_pred)  
recall = recall_score(y_test, y_pred)  
conf_matrix = confusion_matrix(y_test, y_pred)
```

Model	accuracy	precision	recall	F1 score
Multinomial Naive Bayes	0.985	0.984	0.998	0.991
Bernoulli Naive Bayes	0.923	0.918	0.996	0.956

Table 1: Metrics for Easy Ham and Spam

Accuracy measure the proportion of true results among the total number of cases examined, this is calculated according to Equation 1. Precision measures the proportion of true positive results among the total number of cases that were predicted to be positive, this is calculated according to Equation 2. Recall measures the proportion of true positive results among the total number of cases that were actually positive, this is calculated according to Equation 3. F1 score is the harmonic mean of precision and recall, this is calculated according to Equation 4. These metrics are used to evaluate the performance of the models. Value close to 1 indicates that a high percentage of the classifier's predictions are correct.

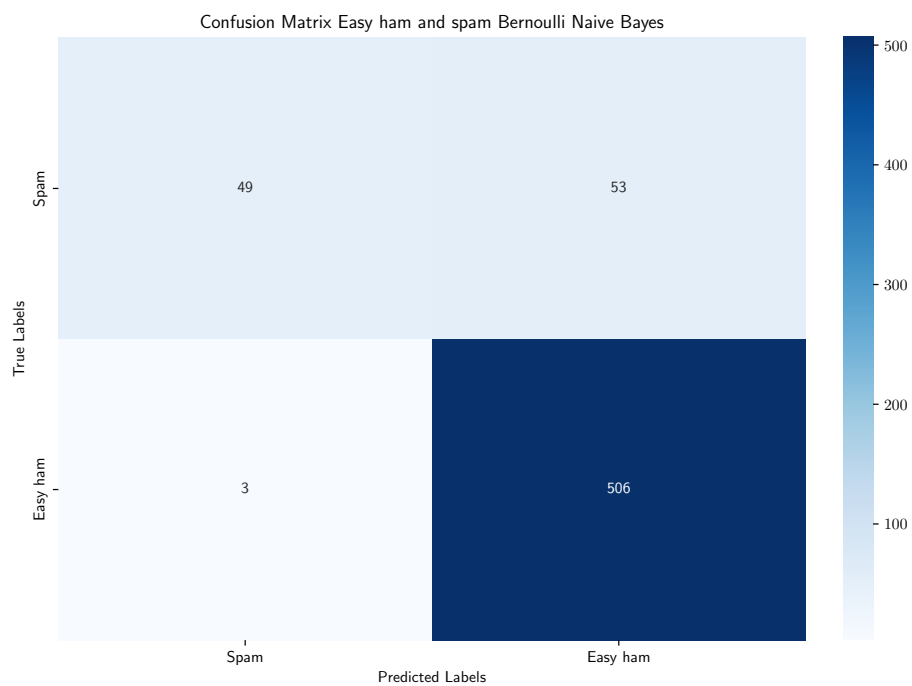
The accuracy, precision, recall, and F1 score for the easy ham and spam dataset are shown in Table 1. The confusion matrixes for the easy ham and spam dataset are shown in Figure 1.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

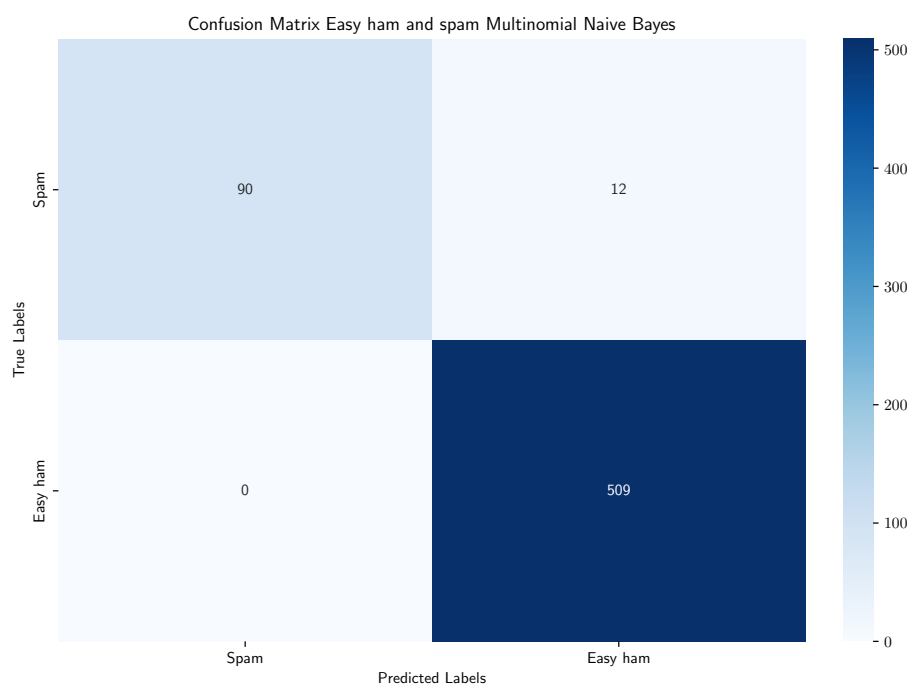
$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (4)$$



(a) Easy ham vs spam, Bernoulli Naive Bayes



(b) Easy ham vs spam, Multinomial Naive Bayes

Figure 1: Confusion matrixes of easy ham and spam

Model	accuracy	precision	recall	F1 score
Multinomial Naive Bayes	0.947	0.956	0.878	0.915
Bernoulli Naive Bayes	0.934	0.976	0.816	0.889

Table 2: Metrics for Hard Ham and Spam

Problem 3: Hard Ham

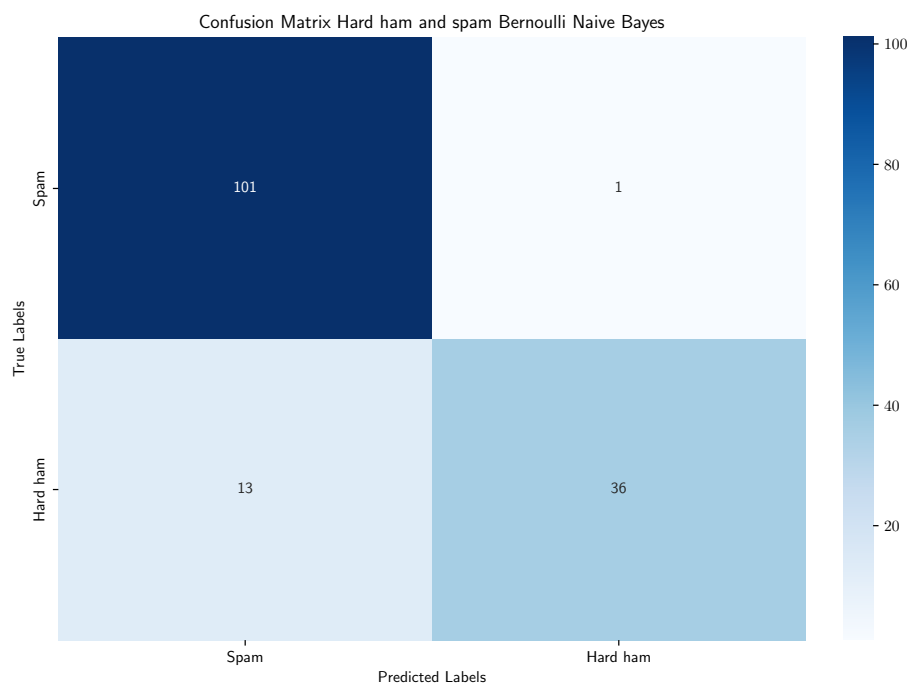
The accuracy, precision, recall, and F1 score for the hard ham and spam dataset are shown in Table 2. The confusion matrixes for the hard ham and spam dataset are shown in Figure 2.

There are 501 emails categorized as spam, 2551 emails categorized as easy ham, while only 250 emails are labeled as hard ham, making the dataset for easy ham significantly larger. Consequently, there is more data available to train the models for easy ham, suggesting that these models should perform better.

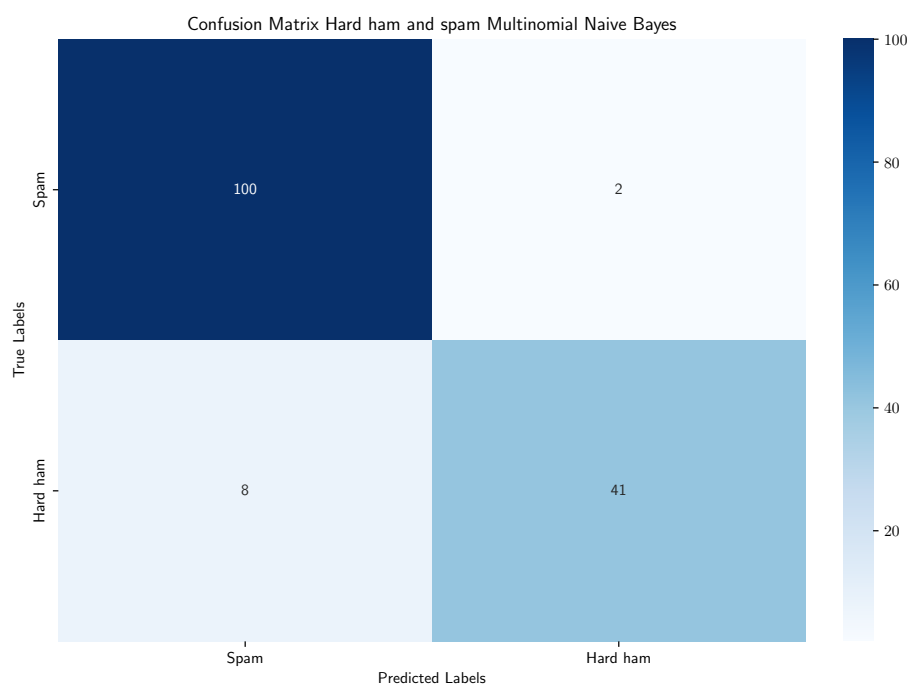
During the model training process, we split the data in the same manner for both easy and hard ham, with 20% allocated for testing and 80% for training. Given the substantial dataset difference between easy and hard ham, some may propose adjusting the data split, potentially allocating 30% for testing and 70% for training for easy ham.

Upon examining the metrics, it becomes evident that the models for easy ham generally exhibit superior performance. Specifically, the Multinomial Naive Bayes model outperformed the other model.

Interestingly, when comparing the confusion matrices, the models encountered difficulty in correctly identifying ham emails in the hard ham category, but demonstrated better success in identifying spam. Conversely, for easy ham, the trend was reversed: the models excelled at identifying ham emails but struggled more with identifying spam. This is likely due to the distribution of the different email types.



(a) Hard ham vs spam, Bernoulli Naive Bayes



(b) Hard ham vs spam, Multinomial Naive Bayes

Figure 2: Confusion matrixes of hard ham and spam

References

- [1] Steven S Skiena. *The Data Science Design Manual*. Retrieved 2024-01-20. 2024. URL: <https://ebookcentral.proquest.com/lib/gu/detail.action?docID=6312797>.

Appendix: Source Code

```
1 from matplotlib import pyplot
2 import tarfile
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.model_selection import train_test_split
5 from sklearn.feature_extraction.text import CountVectorizer
6 from sklearn.naive-bayes import MultinomialNB
7 from sklearn.naive-bayes import BernoulliNB
8 from sklearn.metrics import accuracy_score
9 from sklearn.metrics import confusion_matrix
10 from sklearn.metrics import precision_score, recall_score
11 import seaborn as sns
12
13 def decode_bytes(bytes, encodings=('utf-8', 'ascii', 'ISO-8859-1'))
14     ↪ :
15     for encoding in encodings:
16         try:
17             decoded_text = bytes.decode(encoding)
18             return decoded_text
19         except UnicodeDecodeError:
20             continue
21     return None
22
23 def parse_tar_bz2(file_path):
24     emails = []
25     try:
26         with tarfile.open(file_path, 'r:bz2') as tar:
27             for member in tar.getmembers():
28                 #print("File:", member.name)
29                 file = tar.extractfile(member)
30                 if file is not None:
31                     content = file.read()
32                     emails.append(decode_bytes(content))
33     except tarfile.TarError as e:
34         print("Error~occurred~while~processing~the~tar.bz2~file:",
35             ↪ e)
36     return emails
37
38 def evaluate_model(y_test, y_pred, title, classifier):
39     # Calculate accuracy, precision, recall, and F1 score
40     accuracy = accuracy_score(y_test, y_pred)
41     precision = precision_score(y_test, y_pred)
42     recall = recall_score(y_test, y_pred)
43     print(title + "~and-spam~" + classifier + "~accuracy:",
44         ↪ accuracy)
45     print(title + "~and-spam~" + classifier + "~precision:",
46         ↪ precision)
47     print(title + "~and-spam~" + classifier + "~recall:", recall)
48     print(title + "~and-spam~" + classifier + "~F1-score:", 2 * (
49         ↪ precision * recall) / (precision + recall))
```

```

48     # Create confusion matrix
49     conf_matrix = confusion_matrix(y_test, y_pred)
50     fig, ax = pyplot.subplots(figsize=(8, 6), layout='constrained')
51     sns.heatmap(conf_matrix, annot=True, cmap='Blues', fmt='d',
52                 xticklabels=['Spam', title],
53                 yticklabels=['Spam', title])
54     ax.set_xlabel('Predicted-Labels')
55     ax.set_ylabel('True-Labels')
56     ax.set_title('Confusion-Matrix-' + title + '-and-spam-' +
57                 classifier)
58     filename = title + '-and-spam-' + classifier + '
59                 _confusion-matrix.pdf'
60     filename = filename.replace('-', '_').lower()
61     fig.savefig(filename, bbox_inches='tight')
62
63 def classify_email(emails, labels, title):
64     vectorizer = CountVectorizer()
65
66     # Fit CountVectorizer object to email data and
67     # transform email data into a matrix of token counts
68     email_matrix = vectorizer.fit_transform(emails)
69     # Split data into training and test sets, with 20% of data
70     # reserved for testing
71     X_train, X_test, y_train, y_test = train_test_split(
72         email_matrix, labels, test_size=0.2)
73     print('Size-of-test-set-' + title + ':', len(y_test))
74
75     # Train classifier (Multinomial Naive Bayes and Bernoulli Naive
76     # Bayes)
77     classifierMNB = MultinomialNB()
78     classifierBNB = BernoulliNB()
79     classifierMNB.fit(X_train, y_train)
80     classifierBNB.fit(X_train, y_train)
81
82     # Predict labels for test set
83     y_predMNB = classifierMNB.predict(X_test)
84     y_predBNB = classifierBNB.predict(X_test)
85
86     # Evaluate the classifier
87     evaluate_model(y_test, y_predMNB, title, "Multinomial-Naive-
88         Bayes")
89     evaluate_model(y_test, y_predBNB, title, "Bernoulli-Naive-Bayes
90         ")
91
92     pyplot.rcParams['text.usetex'] = True
93     file_path_easy_ham = "../20021010_easy_ham.tar.bz2"
94     emails_easy_ham = parse_tar_bz2(file_path_easy_ham)
95     file_path_hard_ham = "../20021010_hard_ham.tar.bz2"
96     emails_hard_ham = parse_tar_bz2(file_path_hard_ham)
97     file_path_spam = "../20021010_spam.tar.bz2"
98     emails_spam = parse_tar_bz2(file_path_spam)
99
100    labels_easy_and_spam = [1] * len(emails_easy_ham) + [0] * len(
101        emails_spam)
102    emails_easy_and_spam = emails_easy_ham + emails_spam
103
104    labels_hard_and_spam = [1] * len(emails_hard_ham) + [0] * len(
105        emails_spam)

```

```

101 emails_hard_and_spam = emails_hard_ham + emails_spam
102
103 print("Number of easy ham emails:", len(emails_easy_ham))
104 print("Number of hard ham emails:", len(emails_hard_ham))
105 print("Number of spam emails:", len(emails_spam))
106
107 classify_email(emails_easy_and_spam, labels_easy_and_spam, "Easy -
    ↳ ham")
108 classify_email(emails_hard_and_spam, labels_hard_and_spam, "Hard -
    ↳ ham")

```