

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

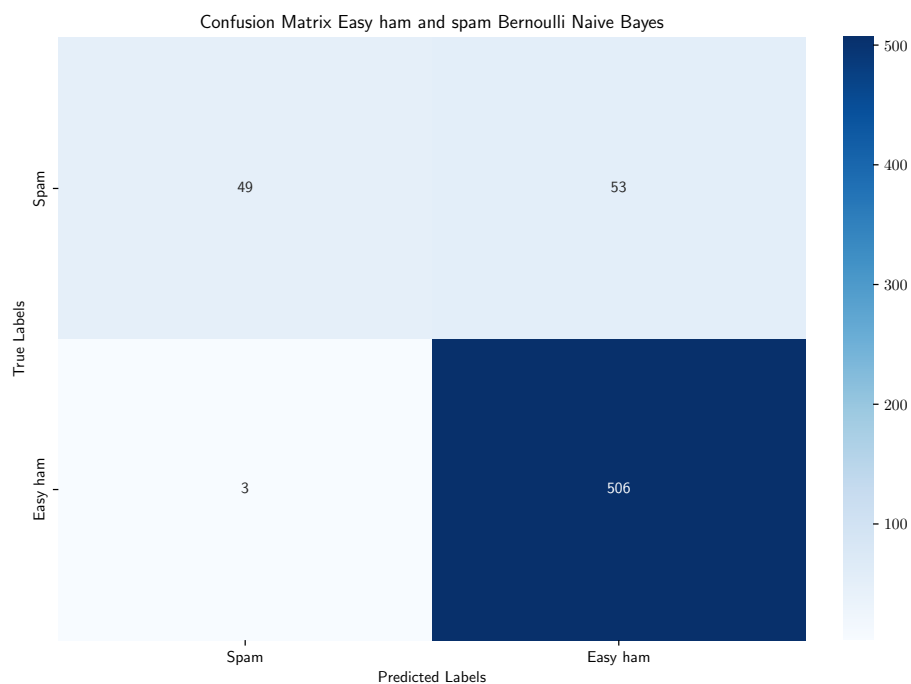
```
X_train, X_test, y_train, y_test =  
train_test_split(email_matrix, labels, test_size=0.2)
```

Problem 2: Preprocessing

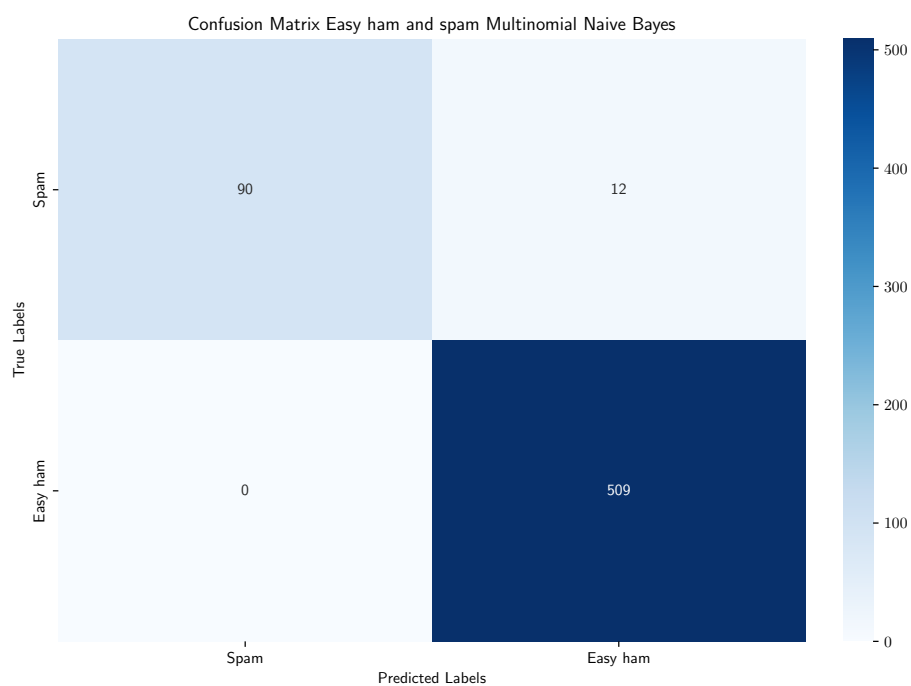
The "bag of words" model is a basic and intuitive way to analyze and compare documents based on their textual content. However, it does not consider the context or the order of words, which can limit its effectiveness in capturing the semantics and meaning of the text.

Problem 3: Easy Ham

```
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)
```



(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

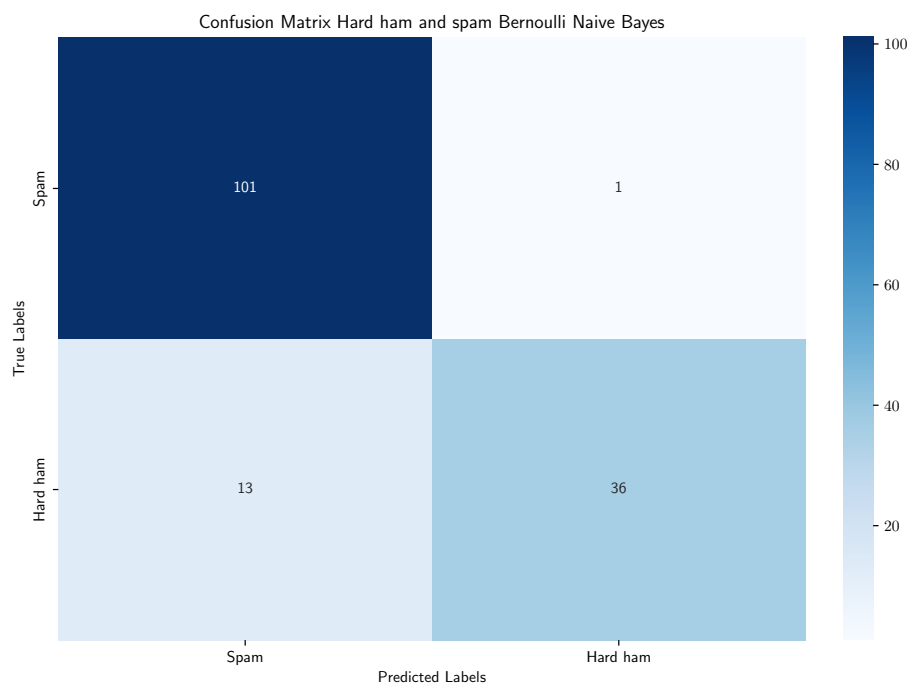
Model	accuracy	precision	recall	F1 score
Multinomial Naive Bayes	0.985	0.984	0.998	0.991
Bernoulli Naive Bayes	0.923	0.918	0.996	0.956

Table 1: Precision and accuracy for Easy Ham and Spam

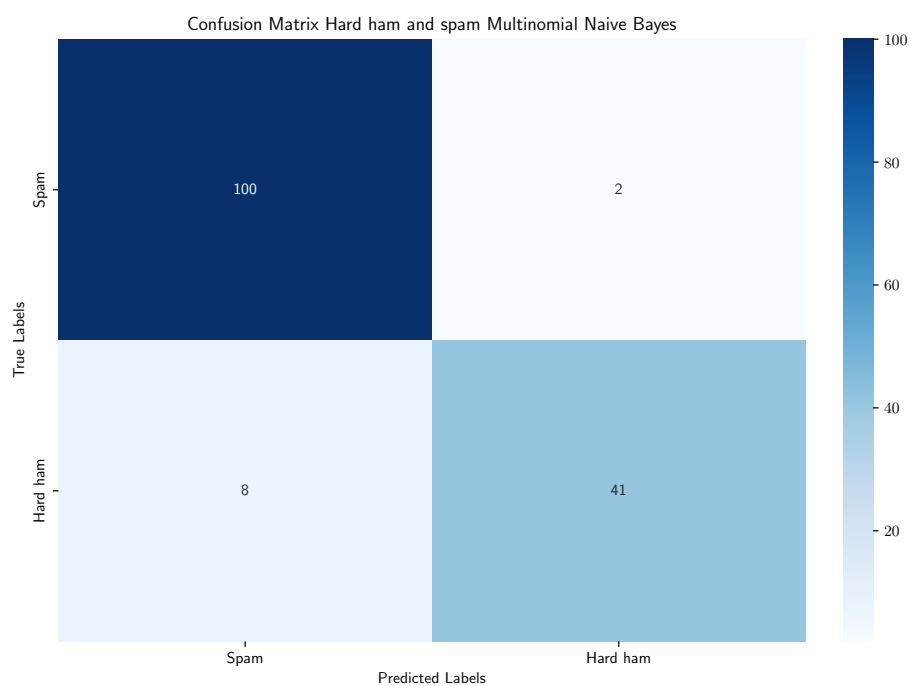
Model	accuracy	precision	recall	F1 score
Multinomial Naive Bayes	0.947	0.956	0.878	0.915
Bernoulli Naive Bayes	0.934	0.976	0.816	0.889

Table 2: Precision and accuracy for Hard Ham and Spam

Problem 3: Hard Ham



(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     accuracy = accuracy_score(y_test, y_pred)
40     precision = precision_score(y_test, y_pred)
41     recall = recall_score(y_test, y_pred)
42     print(title + "~and~spam~" + classifier + "~accuracy:",
43         ↪ accuracy)
44     print(title + "~and~spam~" + classifier + "~precision:",
45         ↪ precision)
46     print(title + "~and~spam~" + classifier + "~recall:", recall)
47     print(title + "~and~spam~" + classifier + "~F1-score:", 2 * (
48         ↪ precision * recall) / (precision + recall))
49
50 conf_matrix = confusion_matrix(y_test, y_pred)
```

```

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

```

```
100  classify_email(emails_hard_and_spam , labels_hard_and_spam , "Hard-  
    ↪ ham")
```