

# DAT565/DIT407 Assignment 3

Henning Anevik      Joar Forsberg

2024-10-08

## Problem 1: Spam and Ham

### Data exploration

Starting by extracting the emails by using the library tarfile. We extracted the files from their bz2 format and now we are going to export the emails to individual text files A to be able to analyze their content and so one can supply an answer to question A. A sample of a mail can be seen in Figure 1

```
Subject: Rats' intestines and pigs' teeth From: boingboing rssfeeds@example.com To:
yyyy@example.com Body: URL: http://boingboing.net/#85497383 Date: Not supplied

This is the headline of the month, possibly the year: "Doctors Grow Pig Teeth in Rat
Intestines." Do we even need to read the story to understand it? It's like a freaking haiku of
near-singularity, future-shocky wonderment!

U.S. doctors said on Thursday they have managed to grow living pig teeth in
rats, a feat of biotechnology that experts said could spark a dental
revolution.

Researchers at Boston's Forsyth Institute said their successful experiment
suggests the existence of dental stem cells, which could one day allow a
person to replace a lost tooth with an identical one grown from his or her
own cells.

"The ability to identify, isolate and propagate dental stem cells to use in
biological replacement tooth therapy has the potential to revolutionize
dentistry," said Dominick DePaola, president and CEO of the institute that
focuses on oral and facial science.

Link[1] Discuss[2] (Thanks, Dave[3]!)
[1] http://story.news.yahoo.com/news?tmpl=story&ncid=585&e=1&cid=585&u=/nm/20020926/sc\_nm/health\_teeth\_dc [2] http://www.quicktopic.com/boing/H/88cUs6CumjiX [3] http://www.remtullaaurorscg.com
```

Figure 1: Sample of a mail in .txt file

Starting by looking at the readme.html provided by spamassassin its stated that the easy\_ham emails don't frequently contain spammish signatures which can be lots of capital letters, HTML and keyword like "adult-only" etc. Hard\_ham on the other hand have several aspects close to spam [3]. So some of the signatures of spam can be found in the emails, which makes them harder to detect as non spam. And the spam file includes spam.

## Data splitting

Start with creating python lists that include spam and easy-ham and the other will also include spam and hard-ham. The following code have resembles to the code seen in "Problem 1" but for more structured code this was the approach used. After the mails are put into the list the data splitting is performed. As mentioned in the lecture the trainset consist of 75% of the mails and only 25% for the test B.

## Problem 2: Preprocessing

Started with importing the library from Scikit-Learn that had the class CountVectorizer. Then applied fit\_transform to the train data. The fit\_transform perform two actions. It learns the models scaling parameters on the training set where it learns the different parameters and their connections. The fit function part calculates the mean and variance of the different features found in the data. Secondly it transform the data into a matrix that is used later by the two classifiers in problem 3 & 4. The transform function is applied on the test dataset [2].

The fit method that comes with the library sklearn is used to train the model on the dataset we provide. If we would have used fit on the test data it would have lead to the actual test becoming over fitted and resulted in an invalid evaluation. The transform function simply put use what it would have learn from the training data and apply the knowledge on the test data.

The preprocessing can be seen in C where CountVectorizer, fit\_transform and transform is used separately on Easy/Hard Ham.

## Problem 3: Easy Ham

With the test and training data gathered from the preprocessing we apply the two different Naïve Bayesian classifiers. The Multinomial Naïve Bayesian & the Bernoulli Naïve Bayesian Classifier seen in D.

Table 1: The Multinomial Naïve Bayesian Classifier (Easy Ham)

<b>Metric</b>	<b>Value</b>	
Accuracy	96.85%	
Precision	98.10%	
Recall	82.40%	
Confusion Matrix	636	2
	22	103

Table 2: The Bernoulli Naive Bayesian Classifier (Easy Ham)

<b>Metric</b>	<b>Value</b>	
Accuracy	90.69%	
Precision	100.00%	
Recall	43.20%	
Confusion Matrix	638	0
	71	54

Both classifiers performed well but had some differences when looking the their metrics. Looking at 1 & 2 one can see that their precision where close but the metrics that was most different was accuracy and recall. MultinomialNB had an accuracy of 96.85% which is around 6% better then BernoulliNB that had an accuracy of 90.68%. The reason behind this difference is that MultinomialNB is better regarding text classification where word count and/or frequency is important, while BernoulliNB works a bit more binary focus on a single keyword. So it can miss the bigger picture which the MultinomialNB captures [1]. Looking at recall MultinomialNB had a value of 82.40% & BernoulliNB only had 43.20%. The recall showcases how many time and actual spam was correctly identified. Meaning that MultinomialNB performed way better then the other classifier.

## Problem 4: Hard Ham

Same as done in problem 3. Only difference is that the Hard Ham is used instead of the Easy Ham E.

Table 3: The Multinomial Naive Bayesian Classifier (Hard Ham)

<b>Metric</b>	<b>Value</b>	
Accuracy	92.55%	
Precision	93.02%	
Recall	96.00%	
Confusion Matrix	54	9
	5	120

Table 4: The Bernoulli Naive Bayesian Classifier (Hard Ham)

<b>Metric</b>	<b>Value</b>	
Accuracy	84.04%	
Precision	81.46%	
Recall	98.40%	
Confusion Matrix	35	28
	2	123

Looking at 3 & 4 we can directly see that MultinomialNB wins this time as well but performed worse in the recall metric. The reason why we can probalby see the recall win for BernoulliNB in recall is because of the feature where its

better when looking into single keywords, which can be found more often in the Hard Ham (more spam signatures). The MultinomialNB achieves a better balance between that resulted in an higher overall performance. This trend can be seen when also looking at the data generated from the hard ham test. The classifiers on the hard ham performed a little worse which points that the ham is somewhat more complex. The Bernoulli Naive Bayes classifier has high recall but are having problems with many false positives. Both classifiers perform better on easy ham in terms of precision vs the hard ham.

## References

- [1] Valentin Calomme. *Difference between Bernoulli and Multinomial Naive Bayes*. Retrieved 2024-10-09. 2018. URL: <https://datascience.stackexchange.com/questions/27624/difference-between-bernoulli-and-multinomial-naive-bayes>.
- [2] Chetna Khanna. *What and why behind fit\_transform() and transform() in scikit-learn!*. Retrieved 2024-10-09. 2020. URL: <https://towardsdatascience.com/what-and-why-behind-fit-transform-vs-transform-in-scikit-learn-78f915cf96fe>.
- [3] Justin Mason. *Spam Assassin public mail corpus*. Retrieved 2023-11-10. 2004. URL: <https://spamassassin.apache.org/old/publiccorpus/>.

## A Data exploration

```
1 import tarfile
2 import os
3
4 for file in os.listdir("/content/"):
5     if file.endswith(".bz2"):
6         with tarfile.open(file, "r:bz2") as tar:
7             tar.extractall(path = "tmp")
8
9 extracted_folders = os.listdir("/content/tmp")
10 print("Extracted_folders:", extracted_folders)
11
12 import email
13 from email import policy
14 from email.parser import BytesParser
15 from email.header import decode_header
16
17 for folder in extracted_folders:
18     path = "/content/tmp/" + folder
19     file_list = os.listdir(path)
20
21 #Opening one folder at the time and then reading and
    ↳ writing each file to a .txt based on the folder
    ↳ name
```

```

22 #In the end genereates three .txt that can be read
    ↳ and analyzed for question A
23 for file_name in file_list:
24     with open(path + "/" + file_name, "rb") as file:
25         with open(folder + ".txt", "a",
    ↳ encoding="utf-8") as f:
26             msg =
    ↳ BytesParser(policy=policy.default).parse(file)
27
28             if msg["subject"] is None:
29                 continue
30             else:
31                 f.write(f"Subject:_{msg['subject']}\n")
32                 f.write(f"From:_{msg['from']}\n")
33
34             if msg['to'] is None:
35                 continue
36             else:
37                 f.write(f"To:_{msg['to']}\n")
38                 body = msg.get_body(preferencelist=('plain',))
39
40             if body is not None:
41                 try:
42                     content = body.get_content()
43                     charset = body.get_content_charset() or
    ↳ "utf-8"
44                     #Attempt to decode using the detected
    ↳ charset if not fallback will apply
45                     content = content.encode("ascii",
    ↳ "ignore").decode(charset,
    ↳ errors="replace")
46                     f.write(f"Body:\n{content}\n")
47                 except LookupError:
48                     f.write(f"Body:[Could not decode the
    ↳ body due to unknown encoding]\n")
49                 #Email seperator for easy reading
50                 f.write("#####\n")

```

## B Data splitting

```

1 from sklearn.feature_extraction.text import
    ↳ CountVectorizer
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score,
    ↳ classification_report, precision_score,
    ↳ recall_score, confusion_matrix
4
5 #Create 2 list which is gonna hold mail from easy and

```

```

    ↪ hard, both the list will also contain spam for
    ↪ the data splitting.
6 lst_ez_ham = []
7 lst_hd_ham = []
8 lst_spam = []
9
10 encodings = ["utf-8", "ascii", "iso-8859-1", "7-bit",
    ↪ "8-bit"]
11
12 for folder in extracted_folders:
13     for file in os.listdir("/content/tmp/" + folder +
    ↪ "/"):
14         for encoding in encodings:
15             try:
16                 with open("/content/tmp/" + folder + "/" +
    ↪ file, "rb") as f:
17                     data = f.readlines()
18                     data = [line.decode(encoding,
    ↪ errors='ignore') for line in data]
19                     content = ''.join(data) # Join lines into
    ↪ a single string
20                     break
21             except (UnicodeDecodeError, LookupError):
22                 continue # Try the next encoding if current
    ↪ fails
23
24         if folder == "easy_ham":
25             lst_ez_ham.append(content)
26         elif folder == "hard_ham":
27             lst_hd_ham.append(content)
28         else:
29             lst_spam.append(content)
30
31 #Rule of thumb: if we do a simple train-test split,
    ↪ use 75% of data for the training set and 25%
    ↪ for the test set
32
33 #Easy Ham Split datasplit
34 X_ez = lst_ez_ham + lst_spam
35 y_ez = [0] * len(lst_ez_ham) + [1] * len(lst_spam)
    ↪ #Labels
36 X_ez_train, X_ez_test, y_ez_train, y_ez_test =
    ↪ train_test_split(X_ez, y_ez, test_size=0.25,
    ↪ random_state=64, stratify=y_ez)
37
38 #Hard Ham Split datasplit
39 X_hd = lst_hd_ham + lst_spam
40 y_hd = [0] * len(lst_hd_ham) + [1] * len(lst_spam)
    ↪ #Labels
41 X_hd_train, X_hd_test, y_hd_train, y_hd_test =

```

```

↪ train_test_split(X_hd, y_hd, test_size=0.25,
↪ random_state=64, stratify=y_hd)

```

## C Preprocessing

```

1 from sklearn.feature_extraction.text import
  ↪ CountVectorizer
2 vectorizer = CountVectorizer()
3 x_ez_train = vectorizer.fit_transform(X_ez_train)
4 x_ez_test = vectorizer.transform(X_ez_test)
5 x_ez_train = vectorizer.fit_transform(X_ez_train)
6 x_ez_test = vectorizer.transform(X_ez_test)

```

## D Easy Ham

```

1 from sklearn.naive_bayes import MultinomialNB,
  ↪ BernoulliNB
2
3 clf = MultinomialNB()
4 clf.fit(x_ez_train, y_ez_train)
5 testA_predictions = clf.predict(x_ez_test)
6
7 accuracy = accuracy_score(y_ez_test,
  ↪ testA_predictions)
8 precision = precision_score(y_ez_test,
  ↪ testA_predictions)
9 recall = recall_score(y_ez_test, testA_predictions)
10 print("The Multinomial Naive Bayesian Classifier")
11 print(f"Accuracy: {accuracy*100:.2f}%")
12 print(f"Precision: {precision*100:.2f}%")
13 print(f"Recall: {recall*100:.2f}%")
14
15 conf_matrix = confusion_matrix(y_ez_test,
  ↪ testA_predictions)
16 print("Confusion Matrix:")
17 print(conf_matrix)
18 print("")
19 #####
20
21 clfB = BernoulliNB()
22 clfB.fit(x_ez_train, y_ez_train)
23 testB_predictions = clfB.predict(x_ez_test)
24
25 accuracy = accuracy_score(y_ez_test,
  ↪ testB_predictions)
26 precision = precision_score(y_ez_test,
  ↪ testB_predictions)

```

```

27 recall = recall_score(y_ez_test, testB_predictions)
28 print("The Bernoulli Naive Bayesian Classifier")
29 print(f"Accuracy: {accuracy*100:.2f}%")
30 print(f"Precision: {precision*100:.2f}%")
31 print(f"Recall: {recall*100:.2f}%")
32
33 conf_matrix = confusion_matrix(y_ez_test,
    ↪ testB_predictions)
34 print("Confusion Matrix:")
35 print(conf_matrix)

```

## E Hard Ham

```

1 clf = MultinomialNB()
2 clf.fit(x_hd_train, y_hd_train)
3 testA_predictions = clf.predict(x_hd_test)
4 accuracy = accuracy_score(y_hd_test,
    ↪ testA_predictions)
5 precision = precision_score(y_hd_test,
    ↪ testA_predictions)
6 recall = recall_score(y_hd_test, testA_predictions)
7
8 print("The Multinomial Naive Bayesian Classifier")
9 print(f"Accuracy: {accuracy*100:.2f}%")
10 print(f"Precision: {precision*100:.2f}%")
11 print(f"Recall: {recall*100:.2f}%")
12
13 conf_matrix = confusion_matrix(y_hd_test,
    ↪ testA_predictions)
14 print("Confusion Matrix:")
15 print(conf_matrix)
16 print("")
17 #####
18
19 clfB = BernoulliNB()
20 clfB.fit(x_hd_train, y_hd_train)
21 testB_predictions = clfB.predict(x_hd_test)
22 accuracy = accuracy_score(y_hd_test,
    ↪ testB_predictions)
23 precision = precision_score(y_hd_test,
    ↪ testB_predictions)
24 recall = recall_score(y_hd_test, testB_predictions)
25
26 print("The Bernoulli Naive Bayesian Classifier")
27 print(f"Accuracy: {accuracy*100:.2f}%")
28 print(f"Precision: {precision*100:.2f}%")
29 print(f"Recall: {recall*100:.2f}%")
30

```



```
31 conf_matrix = confusion_matrix(y_hd_test,  
    ↪ testB_predictions)  
32 print("Confusion_Matrix:")  
33 print(conf_matrix)
```