

# Intelligent Email Filtering: Exploring Naive Bayes for Spam and Subtle Ham Detection

2025-06-23

## Problem 1: Spam and Ham Classification

### Introduction

This project investigates email classification by distinguishing between two primary categories:

- **Ham:** Legitimate (non-spam) emails
- **Spam:** Unsolicited and potentially malicious bulk emails

The datasets used include:

- **Easy Ham:** Clearly benign emails
- **Hard Ham:** Legitimate emails that are more difficult to distinguish from spam
- **Spam:** Emails containing spam characteristics

Two Naive Bayes classifiers—Multinomial and Bernoulli—were implemented to analyze their effectiveness in classifying these emails.

### A. Data Exploration

Initial inspection reveals key differences:

- **Easy Ham:** Shorter messages, clear structure, and fewer HTML or marketing terms.
- **Hard Ham:** More complex messages, occasional links, and vocabulary overlapping with spam.
- **Spam:** Rich in HTML content, includes numbers in email addresses, many links, and terms like “free”, “offer”, “win”.

### B. Data Splitting

Each dataset was labeled appropriately. After labeling, the datasets were split into training and testing sets separately for each category. The feature matrix  $X$  contained the email content, and the target vector  $y$  contained the corresponding labels.

## Problem 2: Preprocessing

We applied text preprocessing using `CountVectorizer` to tokenize and vectorize the text data:

- Unique tokens were assigned to words in the dataset.
- We used `fit_transform()` for training data and `transform()` for test data.
- Default settings were retained for simplicity.

## Problem 3: Easy Ham Classification

### Multinomial Naive Bayes

#### 1. Accuracy:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

- *TP*: Correctly predicted spam
- *TN*: Correctly predicted easy ham
- *FP*: Ham misclassified as spam
- *FN*: Spam misclassified as ham

Calculated using:

$$\text{accuracy\_score}(y\_test\_eh\_spam, \text{predictMN\_eh})$$

#### 2. Precision:

$$\text{Precision} = \frac{TP}{TP + FP}$$

#### 3. Recall:

$$\text{Recall} = \frac{TP}{TP + FN}$$

Both metrics are computed using Scikit-learn functions.

#### Confusion Matrix:

```
eh_MN_confusion_matrix = confusion_matrix(y_test_eh_spam, predictMN_eh)
```

Easy ham Accuracy: 0.9777  
 Easy ham Precision: 0.9910  
 Easy ham Recall: 0.8730

Easy ham Confusion Matrix(Multinomial):

	Positive prediction	Negative prediction	Total
Actual Positive	637	1	638
Actual Negative	16	110	126
Total	653	111	764

Figure 1: Multinomial Naive Bayes for Easy Ham vs Spam

## Bernoulli Naive Bayes

Easy ham Accuracy: 0.9123  
 Easy ham Precision: 0.9683  
 Easy ham Recall: 0.4841

Easy ham Confusion Matrix(Bernoulli):

	Positive prediction	Negative prediction	Total
Actual Positive	636	2	638
Actual Negative	65	61	126
Total	701	63	764

Figure 2: Bernoulli Naive Bayes for Easy Ham vs Spam

## Problem 4: Hard Ham Classification

### Multinomial Naive Bayes

Hard ham Accuracy: 0.9259  
Hard ham Precision: 0.9308  
Hard ham Recall: 0.9603

Hard ham Confusion Matrix(Multinomial):

	Positive prediction	Negative prediction	Total
Actual Positive	54	9	63
Actual Negative	5	121	126
Total	59	130	189

Figure 3: Multinomial Naive Bayes for Hard Ham vs Spam

### Bernoulli Naive Bayes

Hard ham Accuracy: 0.8995  
Hard ham Precision: 0.8741  
Hard ham Recall: 0.9921

Hard ham Confusion Matrix(Bernoulli):

	Positive prediction	Negative prediction	Total
Actual Positive	45	18	63
Actual Negative	1	125	126
Total	46	143	189

Figure 4: Bernoulli Naive Bayes for Hard Ham vs Spam

## Discussion

The experiments show that:

- Both classifiers perform better on easy ham compared to hard ham.
- The hard ham dataset likely reduces accuracy due to similarities with spam and smaller sample size.
- The Multinomial Naive Bayes classifier consistently outperforms Bernoulli in both cases.

- Bernoulli ignores word frequency, which limits its effectiveness for text-heavy data like emails.

Multinomial Naive Bayes benefits from capturing term frequency, making it more suitable for spam filtering in this context.

## References

1. [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html)
2. [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.precision\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.precision_score.html)
3. [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.recall\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.recall_score.html)
4. [https://en.wikipedia.org/wiki/Confusion\\_matrix](https://en.wikipedia.org/wiki/Confusion_matrix)

# Appendix

## Setting up

```
1 #Importing necessary packages
2 import numpy as np
3 import pandas as pd
4 import os
5 import tarfile
6
7 from sklearn.model_selection import train_test_split
8 from sklearn.feature_extraction.text import
   CountVectorizer
9 from sklearn.naive_bayes import BernoulliNB
10 from sklearn.naive_bayes import MultinomialNB
11 from sklearn.metrics import accuracy_score,
   precision_score, recall_score, confusion_matrix
12
13 #Connect to drive to be able to import files
14 from google.colab import drive
15 drive.mount('/content/drive')
16
17 #File path
18 file_path_easy_ham = '/content/drive/MyDrive/20021010
   _easy_ham.tar.bz2'
19 file_path_hard_ham = '/content/drive/MyDrive//20021010
   _hard_ham.tar.bz2'
20 file_path_spam = '/content/drive/MyDrive//20021010
   _spam.tar.bz2'
```

## Extracting tarfiles to working directory

```
1 #Extracting all tarfiles to current woring directory
2
3 #All tarfiles
4 tarfile_list = [file_path_easy_ham, file_path_hard_ham
   , file_path_spam]
5
6 for file in tarfile_list:
7     #Open the tarfiles we have recieved for this
       assignment
8     with tarfile.open(file, "r") as tf:
9         #Extracting files to current working directory
10        tf.extractall()
```

## Reading the files

```
1 #Defining generic function to read and decode the
   files
2 def read_files(folder_name):
3
4     data = []
```

```

5
6 for filename in sorted(os.listdir(folder_name)):
7     #Joining the file path together
8     filepath = os.path.join(folder_name, filename)
9
10    #Try reading with encoding utf-8
11    try:
12        with open(filepath, 'r', encoding='utf-8')
13    as file:
14        content = file.read()
15        data.append(content)
16    except UnicodeDecodeError:
17        #If utf-8 decoding fails, try iso-8859-1
18        try:
19            with open(filepath, 'r', encoding='iso
20            -8859-1') as file:
21                content = file.read()
22                data.append(content)
23            except UnicodeDecodeError:
24                #If that also fails, try ascii
25                try:
26                    with open(filepath, 'r', encoding=
27                    'ascii') as file:
28                        content = file.read()
29                        data.append(content)
30                    except Exception as e:
31                        #If none of these decodings work,
32                        show errors
33                        print(f"Failed to decode {filename}:
34                        {e}")
35
36    return data
37
38
39 #Reading the files
40 easy_ham = read_files('easy_ham')
41 hard_ham = read_files('hard_ham')
42 spam = read_files('spam')
43 print('Reading the files is done')
44
45
46 #Look at the lengths of the data
47 print(f'Length of Easy ham data: {len(easy_ham)}')
48 print(f'Length of Hard ham data: {len(hard_ham)}')
49 print(f'Length of Spam data: {len(spam)}')
50
51
52 #Looking at data
53 print(easy_ham[1])
54 print(hard_ham[7])
55 print(spam[7])

```

SPAM and HAM

```

1 #List of data frames and corresponding labels
2 df_list = [easy_ham, hard_ham, spam]
3 labels = ['easy_ham', 'hard_ham', 'spam']
4
5 #Loop to create a dataframe
6 df_eh, df_hh, df_spam = [
7     #Converting to dataframe by renaming the 1st column as
        email and assigning a new column to each dataframe
8     pd.DataFrame(df).rename(columns={0: 'email'}).
        assign(label=label)
9     for df, label in zip(df_list, labels)]

1 #Looking at the data, it is now divided into email and
    a corresponding label
2 print("Easy Ham DataFrame:\n", df_eh.head(100))
3 print("Hard Ham DataFrame:\n", df_hh.head(100))
4 print("Spam DataFrame:\n", df_spam.head(100))

1 #Creating CSV files if we want to export the data
2 df_eh.to_csv("eh_content.csv")
3 df_hh.to_csv("hh_content.csv")
4 df_spam.to_csv("spam_content.csv")

```

## Data Splitting

```

1 #Data splitting
2 #X contains the full emails, while y contains the
    corresponding label
3 #can test with some different text sizes and random
    states
4
5 #Train-test split for easy ham
6 X_train_eh, X_test_eh, y_train_eh, y_test_eh =
    train_test_split(df_eh['email'], df_eh['label'],
        random_state=50)
7
8 #Train-test split for hard ham
9 X_train_hh, X_test_hh, y_train_hh, y_test_hh =
    train_test_split(df_hh['email'], df_hh['label'],
        random_state=50)
10
11 #Train-test split for spam
12 X_train_spam, X_test_spam, y_train_spam, y_test_spam =
    train_test_split(df_spam['email'], df_spam['label'],
        random_state=50)

1 #Creating vectorizer for assignment 3 (easy ham) and
    assignment 4 (hard ham)
2 count_eh = CountVectorizer()
3 count_hh = CountVectorizer()

```



```

4
5 #Training data
6
7 #Concatenating the spam training data with the two
  different ham datas
8 X_train_eh_spam = pd.concat([X_train_spam, X_train_eh
  ],ignore_index=True)
9 X_train_hh_spam = pd.concat([X_train_spam, X_train_hh
  ],ignore_index=True)
10
11 y_train_eh_spam = pd.concat([y_train_spam, y_train_eh
  ],ignore_index=True)
12 y_train_hh_spam = pd.concat([y_train_spam, y_train_hh
  ],ignore_index=True)
13
14 #CountVectorizer to transform the train data
15 X_train_eh_spam = count_eh.fit_transform(
  X_train_eh_spam)
16 X_train_hh_spam = count_hh.fit_transform(
  X_train_hh_spam)
17
18 #Testing data
19
20 #Apply the same method on the test data
21 X_test_eh_spam = pd.concat([X_test_spam, X_test_eh],
  ignore_index=True)
22 X_test_hh_spam = pd.concat([X_test_spam, X_test_hh],
  ignore_index=True)
23
24 y_test_eh_spam = pd.concat([y_test_spam, y_test_eh],
  ignore_index=True)
25 y_test_hh_spam = pd.concat([y_test_spam, y_test_hh],
  ignore_index=True)
26
27 #CountVectorizer to transform the test data
28 X_test_eh_spam = count_eh.transform(X_test_eh_spam)
29 X_test_hh_spam = count_hh.transform(X_test_hh_spam)

1 #Print the shapes of the transformed test data
2 print("\nShape of X_test_eh_spam matrix:",
  X_test_eh_spam.shape)
3 print("Shape of X_test_hh_spam matrix:",
  X_test_hh_spam.shape)
4
5 #Print the first rows of the transformed test data for
  easy ham an hard ham
6 print("First 5 rows of the transformed test data (
  X_test_eh_spam):")
7 print(X_test_eh_spam.toarray()[:5])
8

```

```

9 print("First 5 rows of the transformed test data (
    X_test_hh_spam):")
10 print(X_test_hh_spam.toarray()[:5])

```

### MultinomialNB for Easy Ham and Spam

```

1 #MultinomialNB - Easy ham and spam
2
3 #clf short for classifier
4 clfMN_eh = MultinomialNB()
5 clfMN_eh.fit(X_train_eh_spam, y_train_eh_spam)
6
7 predictMN_eh = clfMN_eh.predict(X_test_eh_spam)
8 print(predictMN_eh)
9 print()
10
11 #Accuracy
12 #accuracy_score(y_true, y_predicted)
13 eh_MN_accuracy = accuracy_score(y_test_eh_spam,
    predictMN_eh)
14 print(f"Easy ham Accuracy: {eh_MN_accuracy:.4f}")
15
16 #Precision
17 #Can also use label 'easy_ham'
18 eh_MN_precision = precision_score(y_test_eh_spam,
    predictMN_eh, pos_label='spam')
19 print(f"Easy ham Precision: {eh_MN_precision:.4f}")
20
21 #Recall
22 #Can also use label 'easy_ham'
23 eh_MN_recall = recall_score(y_test_eh_spam,
    predictMN_eh, pos_label='spam')
24 print(f"Easy ham Recall: {eh_MN_recall:.4f}")
25
26 #Confusion matrix
27 eh_MN_confusion_matrix = confusion_matrix(
    y_test_eh_spam, predictMN_eh)
28 #Creating a confusion matrix that is easier to read
29 rows = ['Actual Positive', 'Actual Negative']
30 columns = ['Positive prediction', 'Negative prediction',
    '']
31 confusion_MN_eh = pd.DataFrame(eh_MN_confusion_matrix,
    index=rows, columns=columns)
32 #Calculating the sums
33 confusion_MN_eh_total_columns = confusion_MN_eh.sum()
34 confusion_MN_eh_total_rows = confusion_MN_eh.sum(axis
    =1)
35 #Adding sums to matrix
36 confusion_MN_eh["Total"] = confusion_MN_eh_total_rows
37 confusion_MN_eh.loc["Total"] =

```

```

        confusion_MN_eh_total_columns
38 confusion_MN_eh.loc["Total", "Total"] =
        confusion_MN_eh.iloc[: -1, : -1].sum().sum()
39
40 print()
41 print("Easy ham Confusion Matrix:")
42 confusion_MN_eh.astype(int)

```

### BernoulliNB for Easy Ham and Spam

```

1 #Bernoulli - Easy ham and spam
2
3 clfBN_eh = BernoulliNB()
4 clfBN_eh.fit(X_train_eh_spam, y_train_eh_spam)
5
6 predictBN_eh = clfBN_eh.predict(X_test_eh_spam)
7 print(predictBN_eh)
8 print()
9
10 #Accuracy
11 #accuracy_score(y_true, y_predicted)
12 eh_BN_accuracy = accuracy_score(y_test_eh_spam,
    predictBN_eh)
13 print(f"Easy ham Accuracy: {eh_BN_accuracy:.4f}")
14
15 #Precision
16 #Can also use label 'easy_ham'
17 eh_BN_precision = precision_score(y_test_eh_spam,
    predictBN_eh, pos_label='spam')
18 print(f"Easy ham Precision: {eh_BN_precision:.4f}")
19
20 #Recall
21 #Can also use label 'easy_ham'
22 eh_BN_recall = recall_score(y_test_eh_spam,
    predictBN_eh, pos_label='spam')
23 print(f"Easy ham Recall: {eh_BN_recall:.4f}")
24
25 #Confusion matrix
26 eh_BN_confusion_matrix = confusion_matrix(
    y_test_eh_spam, predictBN_eh)
27
28 #Creating a confusion matrix that is easier to read
29 rows = ['Actual Positive', 'Actual Negative']
30 columns = ['Positive prediction', 'Negative prediction',
    '']
31 confusion_BN_eh = pd.DataFrame(eh_BN_confusion_matrix,
    index=rows, columns=columns)
32
33 #Calculating the sums
34 confusion_BN_eh_total_columns = confusion_BN_eh.sum()

```

```

35 confusion_BN_eh_total_rows = confusion_BN_eh.sum(axis
    =1)
36 #Adding sums to matrix
37 confusion_BN_eh["Total"] = confusion_BN_eh_total_rows
38 confusion_BN_eh.loc["Total"] =
    confusion_BN_eh_total_columns
39 confusion_BN_eh.loc["Total", "Total"] =
    confusion_BN_eh.iloc[:-1, :-1].sum().sum()
40
41 print()
42 print("Easy ham Confusion Matrix:")
43 confusion_BN_eh.astype(int)

```

### MultinomialNB for Hard Ham and Spam

```

1 #MultinomialNB - Hard ham and spam
2
3 #clf short for classifier
4 clfMN_hh = MultinomialNB()
5 clfMN_hh.fit(X_train_hh_spam, y_train_hh_spam)
6
7 predictMN_hh = clfMN_hh.predict(X_test_hh_spam)
8 print(predictMN_hh)
9 print()
10
11 #Accuracy
12 #accuracy_score(y_true, y_predicted)
13 hh_MN_accuracy = accuracy_score(y_test_hh_spam,
    predictMN_hh)
14 print(f"Hard ham Accuracy: {hh_MN_accuracy:.4f}")
15
16 #Precision
17 #Can also use label 'hard_ham'
18 hh_MN_precision = precision_score(y_test_hh_spam,
    predictMN_hh, pos_label='spam')
19 print(f"Hard ham Precision: {hh_MN_precision:.4f}")
20
21 #Recall
22 #Can also use label 'hard_ham'
23 hh_MN_recall = recall_score(y_test_hh_spam,
    predictMN_hh, pos_label='spam')
24 print(f"Hard ham Recall: {hh_MN_recall:.4f}")
25
26 #Confusion matrix
27 hh_MN_confusion_matrix = confusion_matrix(
    y_test_hh_spam, predictMN_hh)
28
29 #Creating a confusion matrix that is easier to read
30 rows = ['Actual Positive', 'Actual Negative']
31 columns = ['Positive prediction', 'Negative prediction']

```

```

    ,]
32 confusion_MN_hh = pd.DataFrame(hh_MN_confusion_matrix,
    index=rows, columns=columns)
33
34 #Calculating the sums
35 confusion_MN_hh_total_columns = confusion_MN_hh.sum()
36 confusion_MN_hh_total_rows = confusion_MN_hh.sum(axis
    =1)
37 #Adding sums to matrix
38 confusion_MN_hh["Total"] = confusion_MN_hh_total_rows
39 confusion_MN_hh.loc["Total"] =
    confusion_MN_hh_total_columns
40 confusion_MN_hh.loc["Total", "Total"] =
    confusion_MN_hh.iloc[: -1, : -1].sum().sum()
41
42 print()
43 print("Hard ham Confusion Matrix:")
44 confusion_MN_hh.astype(int)

```

### BernoulliNB for Hard Ham and Spam

```

1 BernoulliNB - Hard ham and spam
2
3 #clf short for classifier
4 clfBN_hh = BernoulliNB()
5 clfBN_hh.fit(X_train_hh_spam, y_train_hh_spam)
6
7 predictBN_hh = clfBN_hh.predict(X_test_hh_spam)
8 print(predictBN_hh)
9 print()
10
11 #Accuracy
12 #accuracy_score(y_true, y_predicted)
13 hh_BN_accuracy = accuracy_score(y_test_hh_spam,
    predictBN_hh)
14 print(f"Hard ham Accuracy: {hh_BN_accuracy:.4f}")
15
16 #Precision
17 #Can also use label 'hard_ham'
18 hh_BN_precision = precision_score(y_test_hh_spam,
    predictBN_hh, pos_label='spam')
19 print(f"Hard ham Precision: {hh_BN_precision:.4f}")
20
21 #Recall
22 #Can also use label 'hard_ham'
23 hh_BN_recall = recall_score(y_test_hh_spam,
    predictBN_hh, pos_label='spam')
24 print(f"Hard ham Recall: {hh_BN_recall:.4f}")
25
26 #Confusion matrix

```

```

27 hh_BN_confusion_matrix = confusion_matrix(
    y_test_hh_spam, predictBN_hh)
28
29 #Creating a confusion matrix that is easier to read
30 rows = ['Actual Positive', 'Actual Negative']
31 columns = ['Positive prediction', 'Negative prediction
    ']
32 confusion_BN_hh = pd.DataFrame(hh_BN_confusion_matrix,
    index=rows, columns=columns)
33
34 #Calculating the sums
35 confusion_BN_hh_total_columns = confusion_BN_hh.sum()
36 confusion_BN_hh_total_rows = confusion_BN_hh.sum(axis
    =1)
37 #Adding sums to matrix
38 confusion_BN_hh["Total"] = confusion_BN_hh_total_rows
39 confusion_BN_hh.loc["Total"] =
    confusion_BN_hh_total_columns
40 confusion_BN_hh.loc["Total", "Total"] =
    confusion_BN_hh.iloc[:,-1, :-1].sum().sum()
41
42 print()
43 print("Hard ham Confusion Matrix:")
44 confusion_BN_hh.astype(int)

```