

# Table of Contents

- [LAB3 tutorial for Machine Learning Naïve Bayes Classifier Algorithm](#)
- [Objective](#)
- [1 The Naïve Bayes algorithm](#)
  - [1.1 Bayes' Theorem:](#)
  - [1.2 Working of Naïve Bayes' Classifier:](#)
  - [1.3 naive assumption](#)
  - [1.4 Gaussian Naive Bayes classifier](#)
  - [1.5 Zero Conditional Probability](#)
- [2 LAB Assignment](#)
  - [Exercise1 Spam filtering](#)
  - [1\) Preparing the text data](#)
  - [2\) Creating word dictionary](#)
  - [3\) Feature Extraction Process](#)
  - [4\) Training the Classifiers](#)
  - [5\) Implementation of the Naive Bayes algorithm](#)
  - [5\) Checking the results on test set](#)
  - [Exercise 2 Compare your Naïve Bayes algorithm with GassianNB from Sklearn, which one does better? Where is the gap?](#)
  - [Exercise 3 Questions](#)

## LAB3 tutorial for Machine Learning Naïve Bayes Classifier Algorithm



Written by Jla Yanhong in 2022. Sept. 16th

### Objective

- Understand one of the most popular and simple machine learning classification algorithms, the Naive Bayes algorithm
- Learn how to implement the Naive Bayes Classifier in Python
- Complete the LAB assignment and submit it to BB.

### 1 The Naïve Bayes algorithm

The Naïve Bayes algorithm is comprised of two words Naïve and Bayes, Which can be described as:

- **Naïve:** It is called Naïve because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the bases of color, shape, and taste, then red, spherical, and sweet fruit is recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.

- **Bayes:** It is called Bayes because it depends on the principle of [Bayes' Theorem](https://www.javatpoint.com/bayes-theorem-in-artificial-intelligence) (<https://www.javatpoint.com/bayes-theorem-in-artificial-intelligence>).

## 1.1 Bayes' Theorem:

- Bayes' theorem is also known as **Bayes' Rule** or **Bayes' law**, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability.
- The formula for Bayes' theorem is given as:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Where,

**P(A|B) is Posterior probability:** Probability of hypothesis A on the observed event B.

**P(B|A) is Likelihood probability:** Probability of the evidence given that the probability of a hypothesis is true.

**P(A) is Prior Probability:** Probability of hypothesis before observing the evidence.

**P(B) is Marginal Probability:** Probability of Evidence.

$$posterior = \frac{likelihood \times prior}{evidence}$$

In practice, there is interest only in the numerator of that fraction, because the denominator does not depend on C and the values of the features  $x_i$  are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

## 1.2 Working of Naïve Bayes' Classifier:

Working of Naïve Bayes' Classifier can be understood with the help of the below example:

Suppose we have a dataset of **weather conditions** and corresponding target variable "**Play**". So using this dataset we need to decide that whether we should play or not on a particular day according to the weather conditions. So to solve this problem, we need to follow the below steps:

1. Convert the given dataset into frequency tables.
2. Generate Likelihood table by finding the probabilities of given features.
3. Now, use Bayes theorem to calculate the posterior probability.

**Problem:** If the weather is sunny, then the Player should play or not?

**Solution:** To solve this, first consider the below dataset:

Day	Outlook	Play
0	Rainy	Yes
1	Sunny	Yes
2	Overcast	Yes
3	Overcast	Yes
4	Sunny	No
5	Rainy	Yes
6	Sunny	Yes
7	Overcast	Yes

Day	Outlook	Play
8	Rainy	No
9	Sunny	No
10	Sunny	Yes
11	Rainy	No
12	Overcast	Yes
13	Overcast	Yes

#### Frequency table for the Weather Conditions:

Weather	Yes	No
Overcast	5	0
Rainy	2	2
Sunny	3	2
Total	10	5

#### Likelihood table weather condition:

Weather	No	Yes	
Overcast	0/4	5/10	5/14= 0.35
Rainy	2/4	2/10	4/14=0.29
Sunny	2/4	3/10	5/14=0.35
All	4/14=0.29	10/14=0.71	

#### Applying Bayes'theorem:

$$P(\text{Yes}|\text{Sunny}) = P(\text{Sunny}|\text{Yes}) \cdot P(\text{Yes}) / P(\text{Sunny})$$

$$P(\text{Sunny}|\text{Yes}) = 3/10 = 0.3$$

$$P(\text{Sunny}) = 0.35$$

$$P(\text{Yes}) = 0.71$$

So

$$P(\text{Yes}|\text{Sunny}) = 0.3 \cdot 0.71 / 0.35 = *0.60**$$

$$P(\text{No}|\text{Sunny}) = P(\text{Sunny}|\text{No}) \cdot P(\text{No}) / P(\text{Sunny})$$

$$P(\text{Sunny}|\text{No}) = 2/4 = 0.5$$

$$P(\text{No}) = 0.29$$

$$P(\text{Sunny}) = 0.35$$

So

$$P(\text{No}|\text{Sunny}) = 0.5 \cdot 0.29 / 0.35 = *0.41**$$

So as we can see from the above calculation that  $P(\text{Yes}|\text{Sunny}) > P(\text{No}|\text{Sunny})$

**Hence on a Sunny day, Player can play the game.**

## 1.3 naive assumption

To start with, let us consider a dataset.

Consider a fictional dataset that describes the weather conditions for playing a game of golf. Given the weather conditions, each tuple classifies the conditions as fit("Yes") or unfit("No") for playing golf.

Here is a tabular representation of our dataset.

Day	OUTLOOK	TEMPERATURE	HUMIDITY	WINDY	PLAY GOLF
0	Rainy	Hot	High	False	No
1	Rainy	Hot	High	True	No
2	Overcast	Hot	High	False	Yes
3	Sunny	Mild	High	False	Yes
4	Sunny	Cool	Normal	False	Yes
5	Sunny	Cool	Normal	True	No
6	Overcast	Cool	Normal	True	Yes
7	Rainy	Mild	High	False	No
8	Rainy	Cool	Normal	False	Yes
9	Sunny	Mild	Normal	False	Yes
10	Rainy	Mild	Normal	True	Yes
11	Overcast	Mild	High	True	Yes
12	Overcast	Hot	Normal	False	Yes
13	Sunny	Mild	High	True	No

The dataset is divided into two parts, namely, **feature matrix** and the **response vector**.

- Feature matrix contains all the vectors(rows) of dataset in which each vector consists of the value of **dependent features**. In above dataset, features are 'Outlook', 'Temperature', 'Humidity' and 'Windy'.
- Response vector contains the value of **class variable**(prediction or output) for each row of feature matrix. In above dataset, the class variable name is 'Play golf'.

Now, its time to put a naive assumption to the Bayes' theorem, which is, **independence** among the features. So now, we split **evidence** into the independent parts.

Now, if any two events A and B are independent, then,

$$P(A, B) = P(A)P(B)$$

Hence, we reach to the result:

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y) \dots P(x_n|y)P(y)}{P(x_1)P(x_2) \dots P(x_n)}$$

which can be expressed as:

$$P(y|x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i|y)}{P(x_1)P(x_2) \dots P(x_n)}$$

Now, as the denominator remains constant for a given input, we can remove that term:

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

Now, we need to create a classifier model. For this, we find the probability of given set of inputs for all possible values of the class variable  $y$  and pick up the output with maximum probability. This can be expressed mathematically as:

$$y = \operatorname{argmax}_y P(y) \prod_{i=1}^n P(x_i|y)$$

So, finally, we are left with the task of calculating  $P(y)$  and  $P(x_i|y)$ .

Please note that  $P(y)$  is also called **class probability** and  $P(x_i|y)$  is called **conditional probability**.

The different naive Bayes classifiers differ mainly by the assumptions they make regarding the distribution of  $P(x_i|y)$ .

Let us try to apply the above formula manually on our weather dataset. For this, we need to do some precomputations on our dataset.

We need to find  $P(x_i|y_i)$  for each  $x_i$  in  $X$  and  $y_i$  in  $y$ . All these calculations have been demonstrated in the tables below:

**Outlook**

	Yes	No	P(yes)	P(no)
Sunny	2	3	2/9	3/5
Overcast	4	0	4/9	0/5
Rainy	3	2	3/9	2/5
<b>Total</b>	<b>9</b>	<b>5</b>	<b>100%</b>	<b>100%</b>

**Temperature**

	Yes	No	P(yes)	P(no)
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cool	3	1	3/9	1/5
<b>Total</b>	<b>9</b>	<b>5</b>	<b>100%</b>	<b>100%</b>

**Humidity**

	Yes	No	P(yes)	P(no)
High	3	4	3/9	4/5
Normal	6	1	6/9	1/5
<b>Total</b>	<b>9</b>	<b>5</b>	<b>100%</b>	<b>100%</b>

**Wind**

	Yes	No	P(yes)	P(no)
False	6	2	6/9	2/5
True	3	3	3/9	3/5
<b>Total</b>	<b>9</b>	<b>5</b>	<b>100%</b>	<b>100%</b>

Play		P(Yes)/P(No)
Yes	9	9/14
No	5	5/14
<b>Total</b>	<b>14</b>	<b>100%</b>

So, in the figure above, we have calculated  $P(x_i|y_i)$  for each  $x_i$  in  $X$  and  $y_i$  in  $y$  manually in the tables 1-4. For example, probability of playing golf given that the temperature is cool, i.e  $P(\text{temp.} = \text{cool} \mid \text{play golf} = \text{Yes}) = 3/9$ .

Also, we need to find class probabilities ( $P(y)$ ) which has been calculated in the table 5. For example,  $P(\text{play golf} = \text{Yes}) = 9/14$ .

So now, we are done with our pre-computations and the classifier is ready!

Let us test it on a new set of features (let us call it today):

*today = (Sunny, Hot, Normal, False)*

So, probability of playing golf is given by:

$$P(Yes|today) = \frac{P(Sunny Outlook|Yes)P(Hot Temperature|Yes)P(Normal Humidity|Yes)P(No Wind|Yes)}{P(today)}$$

and probability to not play golf is given by:

$$P(No|today) = \frac{P(Sunny Outlook|No)P(Hot Temperature|No)P(Normal Humidity|No)P(No Wind|No)}{P(today)}$$

Since, P(today) is common in both probabilities, we can ignore P(today) and find proportional probabilities as:

$$P(Yes|today) \propto \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{6}{9} \cdot \frac{6}{9} \cdot \frac{9}{14} \approx 0.0141$$

and

$$P(No|today) \propto \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{2}{5} \cdot \frac{5}{14} \approx 0.0068$$

Now, since

$$P(Yes|today) + P(No|today) = 1$$

These numbers can be converted into a probability by making the sum equal to 1 (normalization):

$$P(Yes|today) = \frac{0.0141}{0.0141 + 0.0068} = 0.67$$

and

$$P(No|today) = \frac{0.0068}{0.0141 + 0.0068} = 0.33$$

Since

$$P(Yes|today) > P(No|today)$$

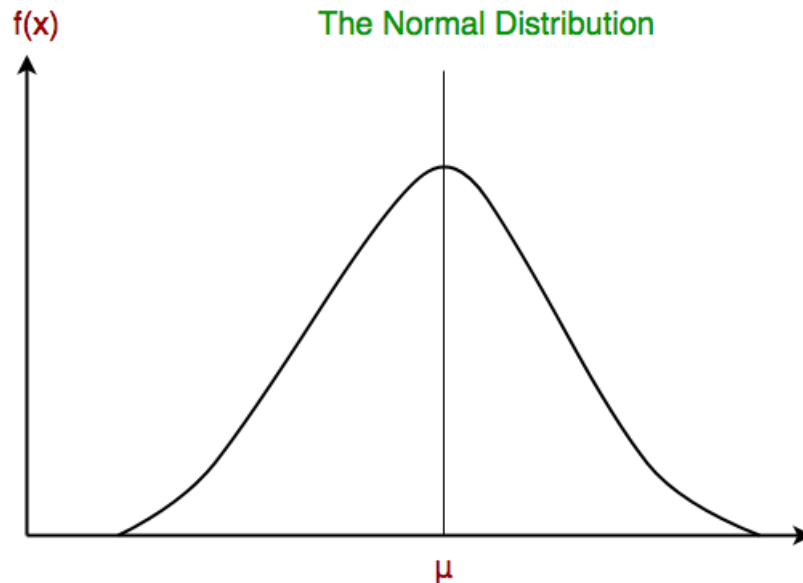
So, prediction that golf would be played is 'Yes'.

The method that we discussed above is applicable for discrete data. In case of continuous data, we need to make some assumptions regarding the distribution of values of each feature. The different naive Bayes classifiers differ mainly by the assumptions they make regarding the distribution of  $P(x_i | y)$ .

Now, we discuss one of such classifiers here.

## 1.4 Gaussian Naive Bayes classifier

In Gaussian Naive Bayes, continuous values associated with each feature are assumed to be distributed according to a **Gaussian distribution**. A Gaussian distribution is also called [Normal distribution](https://en.wikipedia.org/wiki/Normal_distribution) ([https://en.wikipedia.org/wiki/Normal\\_distribution](https://en.wikipedia.org/wiki/Normal_distribution)). When plotted, it gives a bell shaped curve which is symmetric about the mean of the feature values as shown below:



The likelihood of the features is assumed to be Gaussian, hence, conditional probability is given by:

$$P(x_i|y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right)$$

## Example: Continuous-valued Features

- Temperature is naturally of continuous value

Yes: 25.2, 19.3, 18.5, 21.7, 20.1, 24.3, 22.8, 23.1, 19.8

No: 27.3, 30.1, 17.4, 29.5, 15.1

- Estimate mean and variance for each class

$$\mu = \frac{1}{N} \sum_{n=1}^N x_n, \sigma^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$$

$$\mu_{Yes} = 21.64, \sigma_{Yes} = 2.35$$

$$\mu_{No} = 23.88, \sigma_{No} = 7.09$$

- Learning Phase: output two Gaussian models

$$\hat{P}(x|Yes) = \frac{1}{2.35\sqrt{2\pi}} \exp\left(-\frac{(x-21.64)^2}{2 \times 2.35^2}\right)$$

$$\hat{P}(x|No) = \frac{1}{7.09\sqrt{2\pi}} \exp\left(-\frac{(x-23.88)^2}{2 \times 7.09^2}\right)$$

## 1.5 Zero Conditional Probability



Sometimes, we face a zero conditional probability problem during test.

$$\hat{P}(x_1|c_i) \cdots \hat{P}(a_{jk}|c_i) \cdots \hat{P}(x_n|c_i) = 0 \text{ for } x_i = a_{jk}, \hat{P}(a_{jk}|c_i) = 0$$

Let's test another example:

$$today = (Overcast, Hot, Normal, False)$$

and probability to not play golf is given by:

$$\begin{aligned} P(No|today) &\propto P(Overcast Outlook|No)P(Hot Temperature|No)P(Normal Humidity|No)P(No W \\ &\propto \frac{0}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{2}{5} \cdot \frac{5}{14} = 0 \end{aligned}$$

In this case, the model becomes ineffective!!!

For a remedy, class conditional probabilities reestimated with:

$$\hat{P}(a_{jk}|c_i) = \frac{n_c + mp}{n + m}$$

$n_c$ : number of training examples for which  $x_j = a_{jk}$  and  $c = c_i$

$n$ : number of training examples for which  $c = c_i$

$p$ : prior estimate (usually,  $p = \frac{1}{t}$  for  $t$  possible values of  $x_j$ )

$m$ : weight to prior (number of "virtual" examples,  $m \geq 1$ ) + Temperature is naturally of continuous value

## • Example: P(outlook=overcast|no)=0 in the play-golf dataset

- Adding  $m$  virtual examples
  - In this dataset, training examples for the `no` class is 5.
  - We can only add  $m=1$  virtual example in our m-estimate remedy.
- The `outlook` feature can takes only 3 values. So  $p=1/3$ .
- Re-estimate  $P(\text{outlook}|\text{no})$  with the m-estimate

$$p(\text{overcast}|\text{no}) = \frac{0+1*\frac{1}{3}}{5+1} = \frac{1}{18}$$

$$p(\text{sunny}|\text{no}) = \frac{3+1*\frac{1}{3}}{5+1} = \frac{5}{9}$$

$$p(\text{rain}|\text{no}) = \frac{2+1*\frac{1}{3}}{5+1} = \frac{7}{18}$$

## 2 LAB Assignment

Text mining (deriving information from text) is a wide field which has gained popularity with the huge text data being generated. Automation of a number of applications like sentiment analysis, document classification, topic classification, text summarization, machine translation, etc., has been done using machine learning models.

### Exercise1 Spam filtering

In this Exercise, you are required to write your spam filter by using naïve Bayes method. This time you should not use 3rd party libraries including scikit-learn.

Spam filtering is a beginner's example of the document classification task which involves classifying an email as spam or non-spam (a.k.a. ham) mail. An email dataset will be provided. We will use the following steps to build this application:

1) Preparing the text data 2) Creating a word dictionary 3) Feature extraction 4) Implementation of the Naive Bayes algorithm 5) Training the classifier 6) Checking the results on the test set

## 1) Preparing the text data

The data-set used here, is split into a training set and a test set containing 702 mails and 260 mails respectively, divided equally between spam and ham mails. You will easily recognize spam mails as it contains `spmsg` in its filename.

In any text mining problem, text cleaning is the first step where we remove those words from the document which may not contribute to the information we want to extract. Emails may contain a lot of undesirable characters like punctuation marks, stop words, digits, etc which may not be helpful in detecting the spam email. The emails in Ling-spam corpus have been already preprocessed in the following ways:

1. **Removal of stop words** – Stop words like “and”, “the”, “of”, etc are very common in all English sentences and are not very meaningful in deciding spam or legitimate status, so these words have been removed from the emails.
2. **Lemmatization** – It is the process of grouping together the different inflected forms of a word so they can be analyzed as a single item. For example, “include”, “includes,” and “included” would all be represented as “include”. The context of the sentence is also preserved in lemmatization as opposed to stemming (another buzz word in text mining which does not consider meaning of the sentence)

We still need to remove the non-words like punctuation marks or special characters from the mail documents. There are several ways to do it. Here, we will remove such words after creating a dictionary, which is a very convenient method to do so since when you have a dictionary; you need to remove every such word only once.

## 2) Creating word dictionary

We will only perform text analytics on the content to detect the spam mails. As the first step, we need to create a dictionary of words and their frequency. For this task, a training set of 700 mails is utilized. This python function will create the dictionary for you.

```
from collections import Counter
def make_Dictionary(train_dir):
    emails = [os.path.join(train_dir,f) for f in os.listdir (train_dir)]
    all_words = []
    for mail in emails:
        with open (mail) as m:
            for i,line in enumerate (m) :
                if i == 2:
                    words = line.split()
                    all_words += words
    dictionary = Counter(all_words)
    # Write code for non-word removal here

    return dictionary
```

Once the dictionary is created we can add just a few lines of code written below to the above function to remove non-words. Absurd single characters in the dictionary which are irrelevant here are also removed. Do not forget to insert the below code in the function of `make_Dictionary`:

```
list_to_remove = list(dictionary.keys())
for item in list_to_remove:
    if item.isalpha() == False:
        del dictionary[item]
    elif len(item) == 1:
        del dictionary[item]
dictionary = dictionary.most_common(3000)
```

The dictionary can be seen by the command “print dictionary”. You may find some absurd word counts to be high but don’t worry, it’s just a dictionary and you always have a chance to improve it later. If you use the provided dataset, make sure your dictionary has some of the entries given below as most frequent words. Here 3000 most frequently used words are chosen in the dictionary.

In [ ]:

```
1 # To show the most frequent words in train-mails
2 dictionary = make_Dictionary('ling-spam/train-mails')
3 dictionary
```

### 3) Feature Extraction Process

Once the dictionary is ready, we can extract word count vector (our feature here) of 3000 dimensions for each email of the training set. Each **word count vector** contains the frequency of 3000 words in the training file. Of course you might have guessed by now that most of them will be zero. Let us take an example. Suppose we have 500 words in our dictionary. Each word count vector contains the frequency of 500 dictionary words in the training file. Suppose the text in the training file is “Get the work done, work done”, then it will be encoded as

[0, 0, 0, 0, 0, ... .0, 0, 2, 0, 0, 0, ... ., 0, 0, 1, 0, 0, ... 0, 0, 1, 0, 0, ... . 2, 0, 0, 0, 0, 0]

Here, all the word counts are placed at the 296th, 359th, 415th, 495th elements of the word count vector in the length of 500 and the rest are zero.

The below python code will generate a feature vector matrix whose rows denote 700 files of the training set and columns denote 3000 words of the dictionary. The value at index  $ij$  will be the number of occurrences of the  $j^{th}$  word of the dictionary in the  $i^{th}$  file

In [ ]:

```

1 import os
2 import numpy as np
3 def extract_features(mail_dir):
4     files = [os.path.join(mail_dir, fi) for fi in os.listdir(mail_dir)]
5     features_matrix = np.zeros((len(files), 3000))
6     docID = 0
7     _i = 0
8     print(len(files))
9     for fil in files:
10         _i+=1
11         with open(fil) as fi:
12             for i, line in enumerate(fi):
13                 if i == 2:
14                     words = line.split()
15                     for word in words:
16                         wordID = 0
17                         for i, d in enumerate(dictionary):
18                             if d[0] == word:
19                                 wordID = i
20                                 features_matrix[docID, wordID] += 1
21             docID = docID + 1
22         print('\r', 'done {} files'.format(_i), flush=True, end='')
23     return features_matrix

```

In [ ]:

```

1 # Create a dictionary of words with its frequency
2 features_matrix = extract_features("ling-spam/train-mails")
3 features_matrix

```

## 4) Training the Classifiers

Here you should write your Naïve Bayes classifiers after fully understanding its principle.

In [ ]:

```

1 ##### Write Your Code Here #####
2 # Prepare feature vectors per training mail and its labels
3
4 X_train = [[]]
5
6 y_train = []
7
8 #####

```

In [ ]:

```

1 ##### Write Your Code Here #####
2 # Prepare feature vectors per testing mail and its labels
3 X_test = [[]]
4 y_test = []
5
6
7
8 #####

```

## 5) Implementation of the Naive Bayes algorithm

Complete the code for naive Bayes algorithm in `predict` function.

In [ ]:

```

1
2 ##### Write Your Code Here #####
3 class NaiveBayes():
4
5     def fit(self, X, y):
6         y = y.astype(int)
7         self.X = X
8         self.y = y
9         self.classes = np.unique(y)
10        self.parameters = {}
11        for i, c in enumerate(self.classes):
12            # Calculate the mean, variance, prior probability of each class
13            X_Index_c = X[np.where(y == c)]
14            X_index_c_mean = np.mean(X_Index_c, axis=0, keepdims=True)
15            X_index_c_var = np.var(X_Index_c, axis=0, keepdims=True)
16            parameters = {"mean": X_index_c_mean, "var": X_index_c_var, "prior": X_Index_c.shape[0]}
17            self.parameters["class" + str(c)] = parameters
18
19
20
21        def predict(self, X):
22            # return class with highest probability
23            prediction = 0
24
25            # Complete code for naive Bayes algorithm
26
27
28
29
30        return prediction
31
32
33
34 #####

```

In [ ]:

```

1 from sklearn.metrics import confusion_matrix
2
3 #Call the Naive Bayes algorithm, which we wrote ourselves
4 model = NaiveBayes()
5 model.fit(X_train, y_train)
6 result = model.predict(X_test)
7 print(confusion_matrix(y_test, result))

```

## 5) Checking the results on test set

The test set contains 130 spam emails and 130 non-spam emails. Please compute accuracy, recall, F-1 score to evaluate the performance of your spam filter.

In [ ]:

```
1 ##### Write Your Code Here #####
2 from sklearn.metrics import accuracy_score, recall_score, f1_score
3
4
5 #####
```

**Exercise 2 Compare your Naïve Bayes algorithm with GassianNB from Sklearn, which one does better? Where is the gap?**

In [ ]:

```
1 ##### Write Your Code Here #####
2 from sklearn.naive_bayes import GaussianNB
3
4 #####
```

### Exercise 3 Questions

1. Describe another real-world application where the naïve Bayes method can be applied
2. What are the strengths of the naïve Bayes method; when does it perform well?
3. What are the weaknesses of the naïve Bayes method; when does it perform poorly?
4. What makes the naïve Bayes method a good candidate for the classification problem, if you have enough knowledge about the data?