# EXPERIMENT-5

**Problem Statement**

Use any machine learning method to classify the email dataset

**Algorithm**

Naive Bayes classifiers are a collection of classification algorithms based on Bayes' Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e. every pair of features being classified is independent of each other.

Naive Bayes classifier calculates the probability of an event in the following steps:

**Step 1:** Calculate the prior probability for given class labels
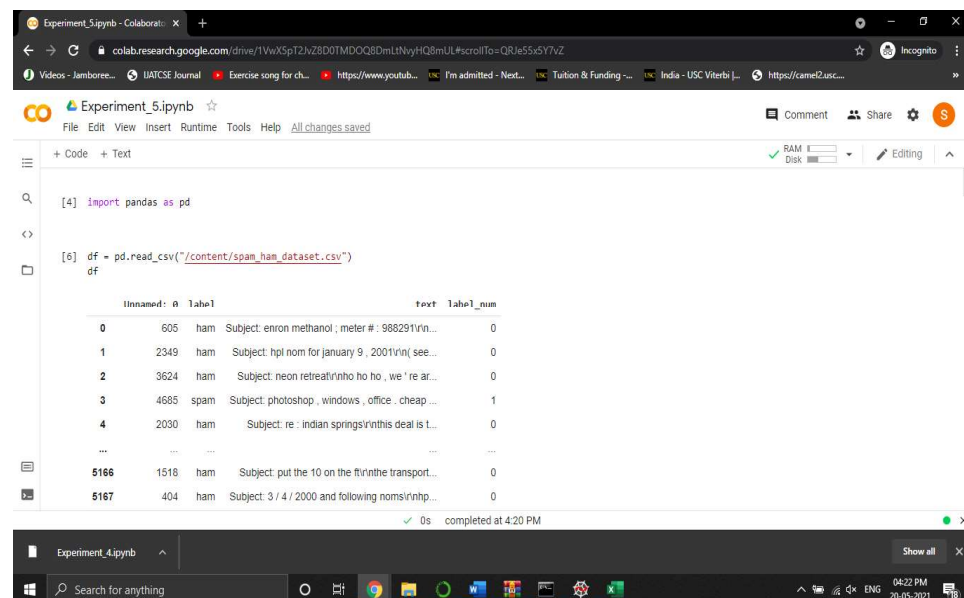
**Step 2:** Find Likelihood probability with each attribute for each class

**Step 3:** Put these values in Bayes Formula and calculate posterior probability.

**Step 4:** See which class has a higher probability, given the input belongs to the higher probability class.

**Program Snippet**

**Reading CSV**

Experiment_5.ipynb

File  Edit  View  Insert  Runtime  Tools  Help   All changes saved

+ Code   + Text

[6]  **5170**      4807    spam    Subject: important online banking alert\r\ndea...          1

5171 rows × 4 columns

[7]  df.head()

|   | Unnamed: 0 | label | text | label_num |
|---|---|---|---|---|
| 0 | 605 | ham | Subject: enron methanol ; meter # : 988291\r\n... | 0 |
| 1 | 2349 | ham | Subject: hpl nom for january 9 , 2001\r\n( see... | 0 |
| 2 | 3624 | ham | Subject: neon retreat\r\nho ho ho , we ' re ar... | 0 |
| 3 | 4685 | spam | Subject: photoshop , windows , office . cheap ... | 1 |
| 4 | 2030 | ham | Subject: re : indian springs\r\nthis deal is t... | 0 |

[8]  df.tail()

|   | Unnamed: 0 | label | text | label_num |
|---|---|---|---|---|
| 5166 | 1518 | ham | Subject: put the 10 on the ft\r\nthe transport... | 0 |

---

Experiment_5.ipynb

File  Edit  View  Insert  Runtime  Tools  Help   All changes saved

+ Code   + Text

df.tail()

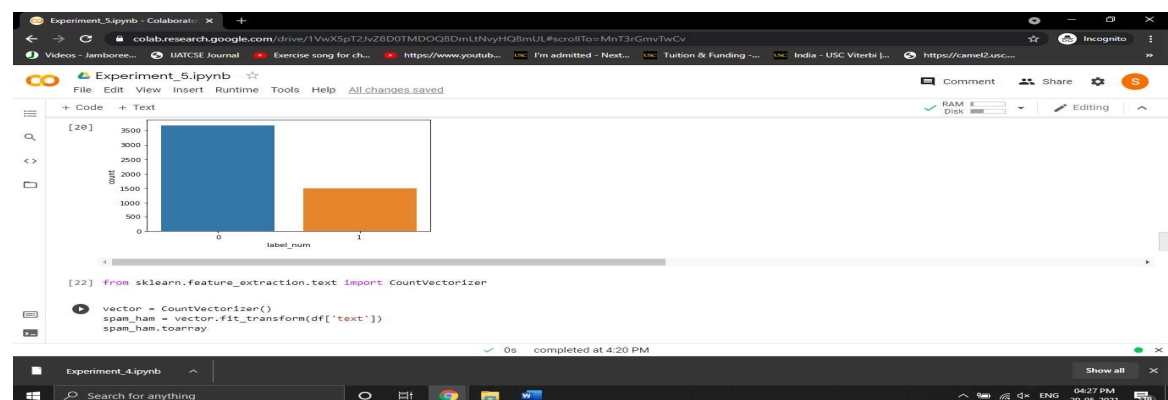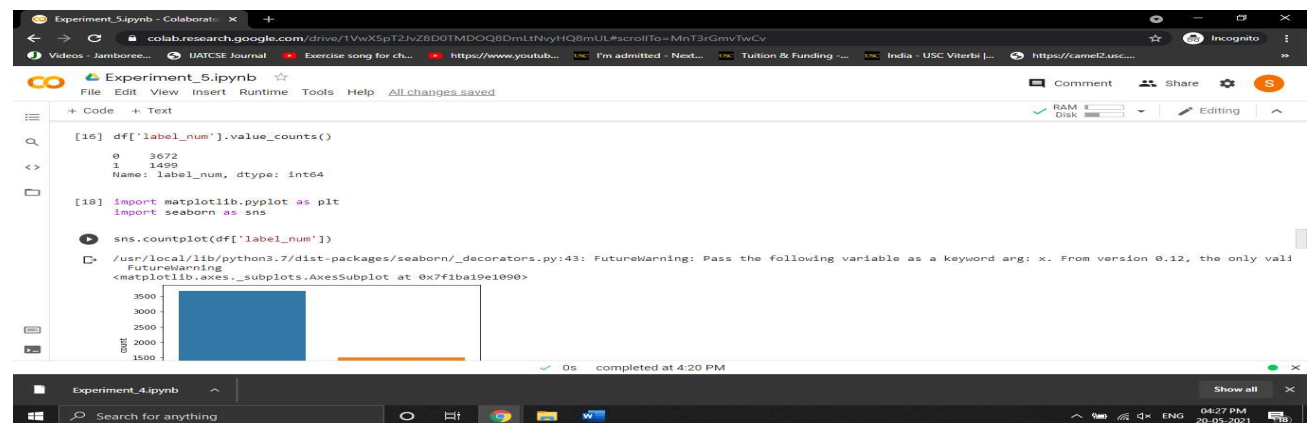|   | Unnamed: 0 | label | text | label_num |
|---|---|---|---|---|
| 5166 | 1518 | ham | Subject: put the 10 on the ft\r\nthe transport... | 0 |
| 5167 | 404 | ham | Subject: 3 / 4 / 2000 and following noms\r\nhp... | 0 |
| 5168 | 2933 | ham | Subject: calpine daily gas nomination\r\n>\r\n... | 0 |
| 5169 | 1409 | ham | Subject: industrial worksheets for august 2000... | 0 |
| 5170 | 4807 | spam | Subject: important online banking alert\r\ndea... | 1 |

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5171 entries, 0 to 5170
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   Unnamed: 0  5171 non-null   int64
 1   label       5171 non-null   object
 2   text        5171 non-null   object
```

**Estimating the Relation**

# Visualization





# Training and Testing the Dataset

```
[31] from sklearn.naive_bayes import MultinomialNB
     nb = MultinomialNB()
     nb.fit(xtrain,ytrain)

     MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True)

[32] ypred = nb.predict(xtrain)
     ypred

     array([0, 0, 0, ..., 1, 0, 0])

     ypredtest = nb.predict(xtest)
     ypredtest

     array([0, 1, 0, ..., 1, 0, 0])

[37] from sklearn.metrics import classification_report , confusion_matrix, accuracy_score
     cmtest = confusion_matrix( ytest, ypredtest)
     cmtrain = confusion_matrix (ytrain, ypred)
     cmtest
```

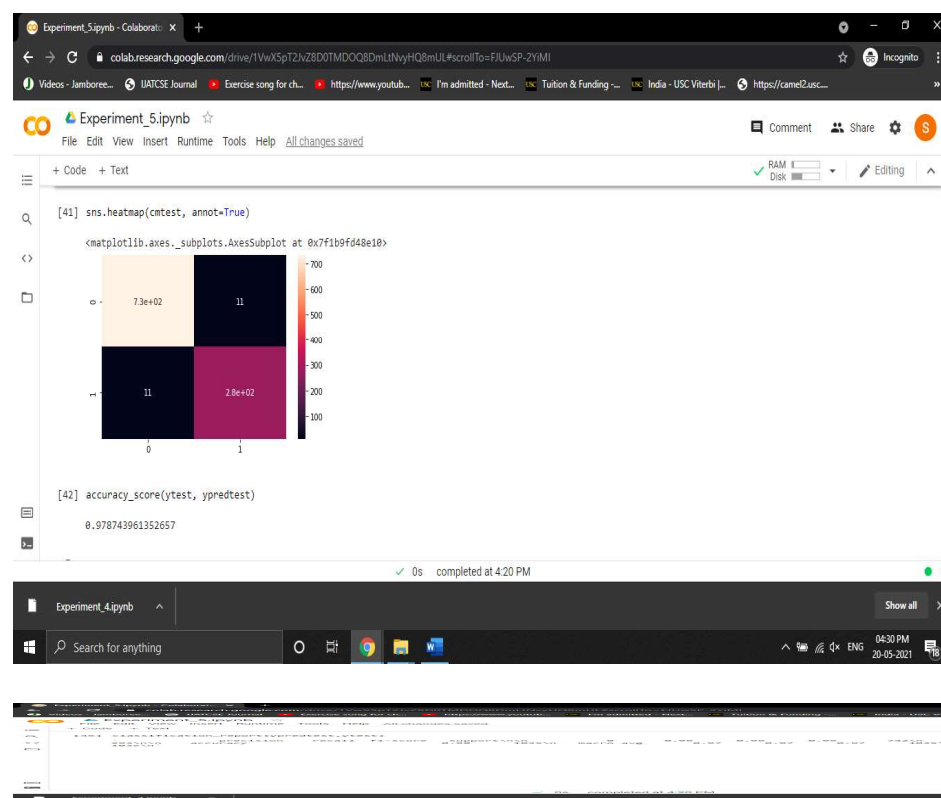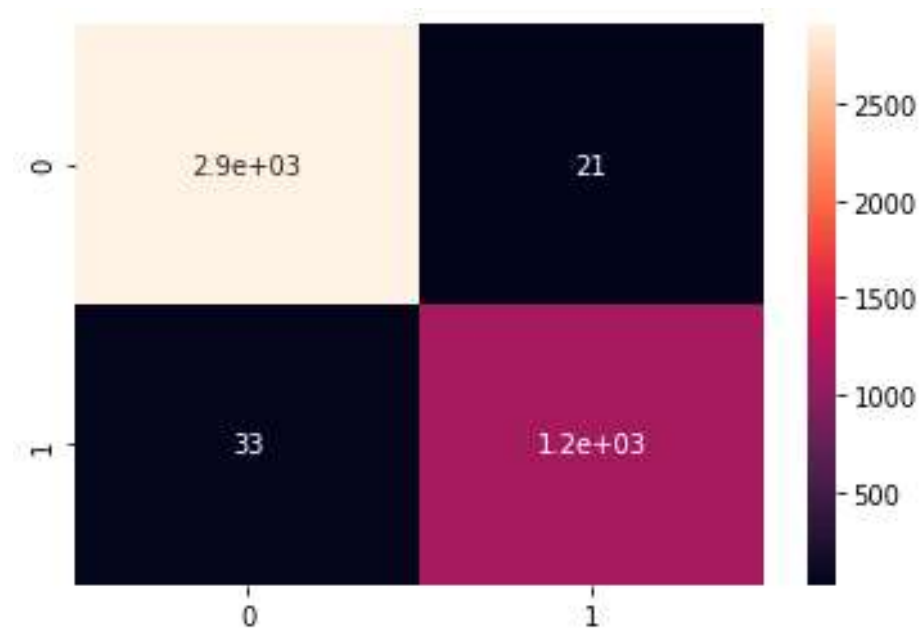**Classification Report and Confusion Matrix**



```
[37] from sklearn.metrics import classification_report , confusion_matrix, accuracy_score
     cmtest = confusion_matrix( ytest, ypredtest)
     cmtrain = confusion_matrix (ytrain, ypred)
     cmtest

     array([[731,  11],
            [ 11, 282]])

     cmtrain

     array([[2909,   21],
            [  33, 1173]])

[40] sns.heatmap(cmtrain, annot=True)

     <matplotlib.axes._subplots.AxesSubplot at 0x7f1b9fd87990>
```

**Github Link**

https://github.com/chanpreet1999/ML-Assignment/blob/master/Exp5/Machine%20Learning%20Experiment%205.ipynb