7. Assuming a set of documents that need to be classified, use naive Bayesian classifier model to perform this task. Built in java class / API can be used to write the program calculate the accuracy precision & recall of your dataset.

```
import pandas as pd
msg = pd.read-csv ('c: /users/ lenovo/
          Desktop / 4MT16CS060 - Prajwal /
          lab 6.csv ') names = ['message', 'label'])
print ( 'Total instances in the dataset:'
                  msg. shape [0])
msg ['labelnum'] = msg. label .map
          [ {' pos : 1, 'neg' : 0 })
x = msg. message
y = msg. labelnum
print ("In The message & its label
    of first 5 instances are listed
             below ")
xs, ys = x[0:5] / msg. label [0:5]
for x, y in zip (xs, ys):
    print (x, ',', y)
```

```
from sklearn.model.selection import
    train-test-split
X train, X test, y train, y test = train-test-split
                              (x,y).
print("Dataset is split into Training &
    Testing samples")
print("Total training instances :', x train.
                     shape [0])
print("Total Testing instances :', x test.
                     shape [0])
from sklearn.feature-extraction.text
    import Count vectorizer
count_vect = Count vectorizer()
x train-dtm = count-vect.fit-transform
                (x train).
x test - dtm = count-vect.transform (x test)
print("\n Total features extracted using
    countvectorizer:", x train-dtm.shape(i))
print("\n Features for first 5
    training instances are listed below")
df = pd.DataFrame (x train-dtm.
    toarray(), columns = count-vect.
        get-feature-names())
print(df [0:5])
```

```
from sklearn.naive.bayes import MultiNominalB
df = MultiNominalNB().fit(x train = dtm,
                                      ytrain)
predicted = df.predict(xtest = dtm)
print("In classification results of testing
          samples are given below")
for doc, p in zip(xtest, predicted):
     pred = 'pos' if p == 1 else 'neg'
     print("%s -> %s".%(doc,pred))

from sklearn import metrices
print("In Accuracy metrices")
print("In Accuracy of the classifier is",
          matrices.accuracy_score
                    (ytest, predicted))
print("Recall:", metrices.recall_score
                    (ytest, predicted))
print("Precision:", metrices.precision_
          score(ytest, predicted))
print("confusion matrix")
print(metrices.confusion_matrix
          (ytest, predicted))
```

## Output:

Total instances in the dataset : 18
The message & its label of first 5 instances
are listed below

    I love this sandwich, pos
    This is an amazing place, pos
    I feel very good about these beers, pos
    This is my best work, pos
    what an awesome view, pos

Dataset is split into Training & Testing Samples
    Total training Instances : 13
    Total testing Instances : 5

Total features extracted using countvec clarizer:
                                                    46
features for first 5 training instances are listed

| | about | am | an | awesome | beers | best | bets | com | deal |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| 3 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

|   | do | ... | today |
|---|----|-----|-------|
| 0 | 1  | ... | 0     |
| 1 | 0  | ... | 1     |
| 2 | 0  | ... | 0     |
| 3 | 0  | ... | 0     |
| 4 | 0  | ... | 0     |

|   | tomorrow | very | view | we | went | what | will | with | work |
|---|----------|------|------|----|------|------|------|------|------|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

[5 rows × 46 colum]

classification results of testing samples are given below.

    I love to dance → pos

    I am sick and tired of this place → neg

    This is an amazing place → pos

    what a great holiday → pos

    This is a bad locality to study → neg

Accuracy metrics
Accuray of the classifier is 1.0
Recall : 1.0
Precision : 1.0
Confusion matrix
$$\{ \{2 \ 0\}$$
$$[0 \ 3]\}$$