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

→ import pandas as pd

```
msg = pd.read_csv('lab6.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("\n The message and its label
        of first 5 instances are listed
        below")
x5, y5 = x[0:5], msg.label[0:5]
for x, y in zip(x5, y5):
        print(x, '.', y)
```

```
from sklearn.model.selection import
train_test_split
xtrain, xtest, ytrain, ytest = train_test
                                _split(x,y)
print("Dataset is split into Training
        and Testing samples")
print("Total training instances :",
                xtrain.shape[0])
print("Total testing instances :",
                xtest.shape[0])


from sklearn.feature_extraction.text
        import countVectorizer.
count_vect = countVectorizer()
xtrain_dbm = count_vect.fit_transform
                (xtrain)
xtest_dbm = count_vect.transform
                (xtest)
print("\nTotal features extracted
        using countVectorizer:")
                xtrain_dbm.shape[1])
print("\nFeatures for first 5 training
instances are listed below")
df = pd.Dataframe(xtrain_dbm.
toarray(), columns = count_vect.
        get_feature_names())
```

```python
print(df[0:5])

from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB().fit(X_train_dtm, y_train)
predicted = clf.predict(X_test_dtm)
print("\n classification results of
             testing samples are given
             below ")
for doc, p in zip(X_test, predicted):
    pred = 'pos' if p == 1 else 'neg'
    print("%s --> %s" % (doc, pred))


from sklearn import metrics
print("\n Accuracy metrics")
print("\n Accuracy of the classifier
is", metrics.accuracy_score(y_test,
             predicted))
print("Recall:", metrics.recall_score(
             y_test, predicted))
print("Precision:", metrics.precision
             _score(y_test, predicted))
print("Confusion matrix")
print(metrics.confusion_matrix(
             y_test, predicted))
```

output:
Total instances in the dataset: 18

The message and 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 video, pos

Dataset is split into Training & Testing samples:
Total training instances: 13
Total testing instances: 5

Total features extracted using count vectorizer: 46
Features for first 5 training instances are listed below:

| | about | am | an | awesome | beers | best | boss | can | deal | do...td |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1...0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0...1 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0...0 |
| 3 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0...0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0...0 |

| | tomorrow | very | view | we | went | whats | will | with | wor |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 1 | 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 x 46 columns]

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 stay → neg

Accuracy metrics

Accuracy of the classifier is 1.0

Recall : 1.0

Precision : 1.0

confusion matrix :

[[2 0]
 [0 3]]