

In []:

```
1 Program:6
2 Assuming a set of documents that need to be classified, use the naïve Bayesian Classifier
3 model to perform this task. Built-in Java classes/API can be used to write the program.
4 Calculate the accuracy, precision, and recall for your data set.
```

In [13]:

```
1 import pandas as pd
```

In [14]:

```
1 msg=pd.read_csv('document.csv',names=['message','label'])
2 print('Total instances in the dataset:',msg.shape[0])
3 msg['labelnum']=msg.label.map({'pos':1,'neg':0})
4 X=msg.message
5 Y=msg.labelnum
6 print('\nThe message and its label of first 5 instances are listed below')
7 X5, Y5 = X[0:5], msg.label[0:5]
8 for x, y in zip(X5,Y5):
9     print(x,',',y)
```

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 view , pos
```

In [15]:

```
1 # Splitting the dataset into train and test data
2 from sklearn.model_selection import train_test_split
3 xtrain,xtest,ytrain,ytest=train_test_split(X,Y)
4 print('\nDataset is split into Training and Testing samples')
5 print('Total training instances :', xtrain.shape[0])
6 print('Total testing instances :', xtest.shape[0])
```

Dataset is split into Training and Testing samples

Total training instances : 13

Total testing instances : 5

In [16]:

```
1 # Output of count vectoriser is a sparse matrix
2 # CountVectorizer - stands for 'feature extraction'
3 from sklearn.feature_extraction.text import CountVectorizer
4 count_vect = CountVectorizer()
5 xtrain_dtm = count_vect.fit_transform(xtrain) #Sparse matrix
6 xtest_dtm = count_vect.transform(xtest)
7 print('\nTotal features extracted using CountVectorizer:',xtrain_dtm.shape[1])
8 print('\nFeatures for first 5 training instances are listed below')
9 df=pd.DataFrame(xtrain_dtm.toarray(),columns=count_vect.get_feature_names())
10 print(df[0:5])#tabular representation
11
```

Total features extracted using CountVectorizer: 42

Features for first 5 training instances are listed below

	am	amazing	an	and	awesome	best	boss	can	dance	deal	...	the	\
0	0	0	0	0	0	0	0	0	0	0	...	0	
1	1	0	0	1	0	0	0	0	0	0	...	0	
2	0	0	0	0	0	0	0	0	0	0	...	0	
3	0	1	1	0	0	0	0	0	0	0	...	0	
4	0	0	0	0	0	0	0	0	1	0	...	0	

	this	tired	to	tomorrow	we	what	will	with	work
0	0	0	0	0	0	0	0	0	0
1	1	1	0	0	0	0	0	0	0
2	0	0	0	1	1	0	1	0	0
3	1	0	0	0	0	0	0	0	0
4	0	0	1	0	0	0	0	0	0

[5 rows x 42 columns]

In [17]:

```
1 #print(xtrain_dtm) #Same as above but sparse matrix representation
2 # Training Naive Bayes (NB) classifier on training data.
3 from sklearn.naive_bayes import MultinomialNB
4 clf = MultinomialNB().fit(xtrain_dtm,ytrain)
5 predicted = clf.predict(xtest_dtm)
6 print('\nClassstification results of testing samples are given below')
7 for doc, p in zip(xtest, predicted):
8     pred = 'pos' if p==1 else 'neg'
9     print('%s -> %s ' % (doc, pred))
10
```

Classstification results of testing samples are given below
I am tired of this stuff -> neg

In [18]:

```
1
2
3 #printing accuracy metrics
4 from sklearn import metrics
5 print('\nAccuracy metrics')
6 print('Accuracy of the classifier is',metrics.accuracy_score(ytest,predicted))
7 print('Recall :',metrics.recall_score(ytest,predicted),
8       '\nPrecison :',metrics.precision_score(ytest,predicted))
9 print('Confusion matrix')
10 print(metrics.confusion_matrix(ytest,predicted))
11
```

Accuracy metrics
Accuracy of the classifer is 0.8
Recall : 1.0
Precison : 0.6666666666666666
Confusion matrix
[[2 1]
 [0 2]]

In []:

```
1
```