

Assuming a set of documents that need to be classified, use the naive Bayesian classifier model to perform this task. Built-in java class / API can be used to write program. You can use java / python ml library class / API

```
import pandas as pd
msg = pd.read_csv('C:\Users\P r Bhat\Desktop\
csvfiles\sampladata.csv', name = ['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 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('In 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
      CountVecorizer
count_vec = CountVecorizer()
```

```
xtrain_dtm = count_vec.fit_transform(xtrain)
xtest_dtm = count_vec.transform(xtest)
print('In Total features extracted using CountVec-  
- total: ', xtrain_dtm.shape[1])
print('In Features for first 5 training instances are  
listed below')
df = pd.DataFrame(xtrain_dtm.toarray(),  
                  columns=count_vec.get_feature_names())
print(df[0:5])
from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB().fit(xtrain_dtm, ytrain)
predicted = clf.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('1.5 → 1.5' if (doc, pred))
from sklearn import metrics
print('In Summary metrics')
print('Accuracy of the classifier is ', metrics.accuracy-  
score(ytest, predicted))
print('Recall: ', metrics.recall_score(ytest, predicted),  
      'In Precision: ', metrics.precision_score(ytest, predicted))
print('Confusion matrix')
print(metrics.confusion_matrix(ytest, predicted))
```

Dataset :

I love this sandwich	pos
This is an amazing place	pos
I feel very good about these bars	pos
This is my best work	pos
what an awesome view	pos
I do not like this restaurant	neg
I'm tired of this stuff	neg
I can't deal with this	neg
He is my sworn enemy	neg
my boss is horrible	neg
This is an awesome place	pos
I do not like the taste of this juice	neg
I love to dance	pos
I am sick and tired of this place	neg
what a great holiday	pos
This is a bad locality to stay	neg
We will have good fun tomorrow	pos
I went to my enemy's house today	neg

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 proud about these beers, pos
 This is my best world, pos
 What an awesome view, pos

Dataset is split into Training and Testing samples

Total training instances: 13

Total testing instances: 5

Total features extracted using CountVecDfizer: 32

Features for first 5 training instances are listed below

am amazing an and awesome bad best con done...

0	0	0	1	0	1	0	0	0	0	...
1	0	0	0	0	0	0	0	0	0	...
2	0	0	0	0	0	0	0	0	0	...
3	0	0	1	0	1	0	0	0	0	...
4	0	1	1	0	0	0	0	0	0	...

the this 1

0 1

0 0

1 1

0 0

0 1

lined to today view wine what with world

0	0	0	0	0	0	0	0	0
1	0	1	1	0	1	0	0	0
2	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0
4	0	0	0	1	0	1	0	0
	0	0	0	0	0	0	0	0

[5 rows x 38 columns]

Classification results of testing samples are given below

I feel very good about these beers \rightarrow neg

my boss is horrible \rightarrow pos

He is my sworn enemy \rightarrow pos

We will have good fun tomorrow \rightarrow neg

what a great holiday \rightarrow pos

Accuracy metrics

Accuracy of the classifier is 0.2

Recall : 0.3333333333333333

Precision : 0.3333333333333333

confusion matrix

$\begin{bmatrix} 0 & 2 \end{bmatrix}$

$\begin{bmatrix} 2 & 1 \end{bmatrix}$