

Program 6

Assuming a set of documents that need to be classified, use the naive Bayesian classifier model to perform this task. Built-in java classes/API can be used to write the program. Calculate the accuracy, precision, and recall for your dataset.

```
import pandas as pd
msg = pd.read_csv('data6.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 and its label of first 5 instances are listed below')
x5, y5 = x[0:5], msg['labelnum'][0:5]
for x, y in zip(x5, y5):
    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('In Dataset is split into Training and Testing Samples')
print('In Total Training instances:', x_train.shape[0])
print('In Total testing instances:', x_test.shape[0])
from sklearn.feature_extraction.text import CountVectorizer
count_vect = CountVectorizer()
x_train_dtm = count_vect.fit_transform(x_train)
x_test_dtm = count_vect.transform(x_test)
```

Teacher's Signature

```
print('In Total features extracted using CountVecorizer:',  
      xtrain_dtm.shape[1])  
print('In Features for first 5 training instances are listed  
      below')  
df = pd.DataFrame(xtrain_dtm.toarray(), columns=word_vec.  
                  get_feature_names())  
print(df[0:5])  
from sklearn.naive_bayes import MultinomialNB  
clf = MultinomialNB().fit(xtrain_dtm, ytrain)  
predicted = clf.predict(xtrain_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 metrics  
print('In Accuracy metrics')  
print('Accuracy of the classifier is', metrics.accuracy_score  
      (ytest, predicted))  
print('Recall', metrics.recall_score(ytest, predicted), 'In Precision:  
      , metrics.precision_score):  
  
print('Confusion Matrix')  
print(metrics.confusion_matrix(test, predicted))
```

Teacher's Signature _____

Dataset

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

I do not like 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

That is 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 good about these beers, pos

This is my best work, 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 countvectorizer: 46

Features for first 5 training instances are listed below

	about	am	an	awesome	beers	but	boss	can	did	do	total
0	0	0	0	0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	1	1	0	0
3	0	0	0	1	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0

	tomorrow	my	view	we	went	what	will	with	when
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 10 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]]