

```
In [30]: 1 import numpy as np
        2 from sklearn.datasets import fetch_20newsgroups
```

```
In [31]: 1 from sklearn.feature_extraction.text import CountVectorizer
```

```
In [32]: 1 from sklearn.naive_bayes import BernoulliNB, MultinomialNB
        2 #input features are binary
        3 #input features are multiple in multinomial (check frequency)
```

```
In [34]: 1 newsgroups=fetch_20newsgroups(subset='all')
```

```
In [35]: 1 vectorizer1=CountVectorizer(binary=True) #words are there are not
        2 vectorizer2=CountVectorizer(binary=False) #words are repeated or not
```

```
In [36]: 1 X1=vectorizer1.fit_transform(newsgroups.data)
        2 X2=vectorizer2.fit_transform(newsgroups.data)
```

```
In [37]: 1 y=newsgroups.target
```

```
In [38]: 1 from sklearn.model_selection import train_test_split
        2 xtrain1,xtest1,ytrain,ytest=train_test_split
        3 (X1,y,test_size=0.25,random_state=1) #X1 binary values
        4 xtrain2,xtest2,ytrain,ytest=train_test_split
        5 (X2,y,test_size=0.25,random_state=1) #X2 multivalues
```

```
In [39]: 1 bnb=BernoulliNB()
        2 bnb
        3 mnb=MultinomialNB()
        4 mnb
```

```
Out[39]: ▾ MultinomialNB
          MultinomialNB()
```

```
In [40]: 1 bnb.fit(xtrain1,ytrain)
```

```
Out[40]: ▾ BernoulliNB
          BernoulliNB()
```

In [41]: `1 mnb.fit(xtrain2,ytrain)`

Out[41]: `▼ MultinomialNB`
`MultinomialNB()`

In [42]: `1 y_pred1=mnb.predict(xtest1)`
`2 y_pred1`

Out[42]: `array([16, 9, 18, ..., 13, 7, 14])`

In [43]: `1 y_pred2=mnb.predict(xtest2)`
`2 y_pred2`

Out[43]: `array([16, 19, 18, ..., 13, 7, 14])`

In [44]: `1 from sklearn.metrics import accuracy_score`

In [45]: `1 accuracy_score(ytest,y_pred1)`

Out[45]: `0.681239388794567`

In [46]: `1 accuracy_score(ytest,y_pred2)`

Out[46]: `0.8384974533106961`

Conclusion: Multinomial is better than bernoulli for this corpus because it checks the frequency.

In [47]: `1 from sklearn.feature_extraction.text import TfidfVectorizer`

In [48]: `1 from sklearn.pipeline import make_pipeline`
`2 model=make_pipeline(TfidfVectorizer(),MultinomialNB())`

In [49]: `1 test_data=fetch_20newsgroups(subset='test')`
`2 train_data=fetch_20newsgroups(subset='train')`

In [50]: `1 model.fit(train_data.data,train_data.target)`

Out[50]: `► Pipeline`
`► TfidfVectorizer`
`► MultinomialNB`

```
In [52]: ▶ 1 predictions_tf=model.predict(test_data.data)
          2 predictions_tf
```

```
Out[52]: array([ 7, 11,  0, ...,  9,  3, 15])
```

```
In [53]: ▶ 1 accuracy_score(test_data.target,predictions_tf)
```

```
Out[53]: 0.7738980350504514
```