```python
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
```

```python
data=pd.read_csv('spam.csv',encoding=('ISO-8859-1'))
data
```

|  | Category | Message | Unnamed: 2 | Unnamed: 3 | Unnamed: 4 |
|---|---|---|---|---|---|
| 0 | ham | Go until jurong point, crazy.. Available only ... | NaN | NaN | NaN |
| 1 | ham | Ok lar... Joking wif u oni... | NaN | NaN | NaN |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... | NaN | NaN | NaN |
| 3 | ham | U dun say so early hor... U c already then say... | NaN | NaN | NaN |
| 4 | ham | Nah I don't think he goes to usf, he lives aro... | NaN | NaN | NaN |
| ... | ... | ... | ... | ... | ... |
| 5567 | spam | This is the 2nd time we have tried 2 contact u... | NaN | NaN | NaN |
| 5568 | ham | Will Ì_ b going to esplanade fr home? | NaN | NaN | NaN |
| 5569 | ham | Pity, * was in mood for that. So...any other s... | NaN | NaN | NaN |
| 5570 | ham | The guy did some bitching but I acted like i'd... | NaN | NaN | NaN |
| 5571 | ham | Rofl. Its true to its name | NaN | NaN | NaN |

5572 rows × 5 columns

```python
data.columns
```

```
Index(['Category', 'Message', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], dtype='object')
```

```python
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   Category    5572 non-null   object
 1   Message     5572 non-null   object
 2   Unnamed: 2  50 non-null     object
 3   Unnamed: 3  12 non-null     object
 4   Unnamed: 4  6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB
```

**Dropped The Column Unnamed: 0**

```python
data.isna().sum()
```

```
Category         0
Message          0
Unnamed: 2    5522
Unnamed: 3    5560
Unnamed: 4    5566
dtype: int64
```

```python
data['Spam']=data['Category'].apply(lambda x:1 if x=='spam' else 0)
data.head(5)
```

|  | Category | Message | Unnamed: 2 | Unnamed: 3 | Unnamed: 4 | Spam |
|---|---|---|---|---|---|---|
| 0 | ham | Go until jurong point, crazy.. Available only ... | NaN | NaN | NaN | 0 |
| 1 | ham | Ok lar... Joking wif u oni... | NaN | NaN | NaN | 0 |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... | NaN | NaN | NaN | 1 |
| 3 | ham | U dun say so early hor... U c already then say... | NaN | NaN | NaN | 0 |

```
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(data.Message,data.Spam,test_size=0.25)


#CounterVectorizer Convert the text into matrics
from sklearn.feature_extraction.text import CountVectorizer
```
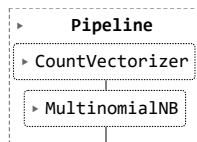
**Naive Bayes Have three Classifier(Bernouli,Multinominal,Gaussian) Here I use Multinominal Bayes Because here data in a discrete form discrete data(e.g movie ratings ranging 1 to 5 as each rating will have certain frequency to represent)**

```
from sklearn.naive_bayes import MultinomialNB


from sklearn.pipeline import Pipeline
clf=Pipeline([
    ('vectorizer',CountVectorizer()),
    ('nb',MultinomialNB())
])
```

## ▾ Tarining The Model

```
clf.fit(X_train,y_train)
```

```
┌─────────────────────────────┐
│  ▸      Pipeline             │
│  ┌─────────────────────────┐ │
│  │ ▸ CountVectorizer       │ │
│  └─────────────────────────┘ │
│    ┌───────────────────────┐ │
│    │ ▸ MultinomialNB       │ │
│    └───────────────────────┘ │
└─────────────────────────────┘
```

**Here I given Two email Two detect 1st One is looking good and the other one looking spam**

```
emails=[
    'Sounds great! Are you home now?',
    'Will u meet ur dream partner soon? Is ur career off 2 a flyng start? 2 find out free, txt HORO followed by ur star sign, e. g. HORO
]
```

**Predict Email**

```
clf.predict(emails)
```

```
array([0, 1])
```

## ▾ Prediction Of Model

```
clf.score(X_test,y_test)
```

```
0.990667623833453
```

✓ 0s    completed at 11:17 PM                                    ● ✕