

```
In [1]: import numpy as np
import pandas as pd
```

```
In [2]: import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
In [4]: data=pd.read_csv('spam.csv')
data
```

```
Out[4]:
```

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

5572 rows × 2 columns

```
In [5]: data.columns
```

```
Out[5]: Index(['Category', 'Message'], dtype='object')
```

```
In [6]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 5572 entries, 0 to 5571  
Data columns (total 2 columns):  
#   Column      Non-Null Count  Dtype    
---  ---      -  
0   Category    5572 non-null   object   
1   Message     5572 non-null   object   
dtypes: object(2)  
memory usage: 87.2+ KB
```

```
In [7]: data.isna().sum()
```

```
Out[7]: Category    0  
Message          0  
dtype: int64
```

```
In [8]: data['Spam']=data['Category'].apply(lambda x:1 if x=='spam' else 0)  
data.head(5)
```

```
Out[8]:
```

| | Category | Message | Spam |
|---|----------|---|------|
| 0 | ham | Go until jurong point, crazy.. Available only ... | 0 |
| 1 | ham | Ok lar... Joking wif u oni... | 0 |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... | 1 |
| 3 | ham | U dun say so early hor... U c already then say... | 0 |
| 4 | ham | Nah I don't think he goes to usf, he lives aro... | 0 |

```
In [9]: 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)
```

```
In [10]: from sklearn.feature_extraction.text import CountVectorizer
```

```
In [11]: from sklearn.naive_bayes import MultinomialNB
```

```
In [12]: from sklearn.pipeline import Pipeline
clf=Pipeline([
    ('vectorizer',CountVectorizer()),
    ('nb',MultinomialNB())
])
```

```
In [13]: clf.fit(X_train,y_train)
```

```
Out[13]: Pipeline(steps=[('vectorizer', CountVectorizer()), ('nb', MultinomialNB())])
```

```
In [14]: 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 b
```



```
In [15]: clf.predict(emails)
```

```
Out[15]: array([0, 1], dtype=int64)
```

```
In [16]: clf.score(X_test,y_test)
```

```
Out[16]: 0.9856424982053122
```

```
In [ ]:
```