

Project No #7

Spam Email Detection



	Label	EmailText
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...

Importing needful libraries

```
: 1 import pandas as pd
  2 from sklearn.feature_extraction.text import CountVectorizer
  3 from sklearn.model_selection import train_test_split
  4 from sklearn.svm import SVC
  5 from sklearn.model_selection import GridSearchCV
  6 from sklearn.model_selection import KFold
  7 from sklearn.metrics import accuracy_score
  8 from imblearn.over_sampling import SMOTE
```

Importing data

```
: 1 data = pd.read_csv("D:\\Workshops\\Python for Data Science Comprehensive
  2 data.head()
```

	Label	EmailText
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...

class imbalance

```
: data["Label"].value_counts()
```

```
: ham      4825
   spam      747
   Name: Label, dtype: int64
```

Separating independent & dependent data

```
: 1 x = data["EmailText"]  
  2 y = data["Label"]
```

Creating a matrix with frequencies of email text

```
1 cvec = CountVectorizer()  
2 cx = cvec.fit_transform(x)
```

Using SMOTE for balancing the response data

```
1 smt=SMOTE()  
2 x_sm,y_sm=smt.fit_resample(cx,y)
```

```
3 y_sm.value_counts()  
:  
ham      4825  
spam     4825  
Name: Label, dtype: int64
```

Splitting training & testing data

```
: 1 x_train,x_test,y_train,y_test=train_test_split(x_sm,y_sm,test_size=0.2,random_state=0)
```

Grid Search for identifying best hyperparameters

```
: 1 params = {'kernel': ['rbf', 'linear']}  
  2 cv=KFold(n_splits=5)
```

```
: 1 model = GridSearchCV(SVC(), params, cv=cv)
```

```
: 1 model.fit(x_train, y_train)  
  2 print(model.best_params_)
```

```
{'kernel': 'linear'}
```

Defining the model with best hyperparameters

```
: 1 bmodel=SVC(kernel="linear")
```

Training the model with training data

```
1 bmodel.fit(x_train, y_train)
```

```
SVC(kernel='linear')
```

Predictions and accuracy of the model

```
: 1 y_pred=bmodel.predict(x_test)

: 1 y_pred
: array(['ham', 'spam', 'spam', ..., 'spam', 'spam', 'spam'], dtype=object)

: 1 y_test
: 1070      ham
   4488      ham
   8763     spam
   7372     spam
   7633     spam
   ...
   212      ham
   4546      ham
   6411     spam
   7916     spam
   6712     spam
Name: Label, Length: 1930, dtype: object
```

```
1 accuracy_score(y_test,y_pred)
```

0.9569948186528497

```
1 print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
ham	0.98	0.93	0.95	936
spam	0.94	0.98	0.96	994
accuracy			0.96	1930
macro avg	0.96	0.96	0.96	1930
weighted avg	0.96	0.96	0.96	1930

Checking the model with new sample emails

```
1 emails=["Hey, you have won a car !!!!!. Congrattzz","Dear applicant, Your CV has been recieved. Best regards"]  
1 yp=bmodel.predict(cvec.transform(emails))  
1 yp  
array(['spam', 'ham'], dtype=object)
```

using Previous objects

