

Aim: Implement Naive Bayes classifier for Spambased dataset.

```
In [171]: 1 import numpy as np
          2 from sklearn.cross_validation import train_test_split
```

```
In [172]: 1 data=open('C:/Users/cglab/Desktop/spam.txt','r')
```

```
In [173]: 1 d=[]
          2 for line in data:
          3     line=[float(element) for element in line.rstrip('\n').split(',')]
          4     d.append(np.asarray(line))
          5
```

```
In [174]: 1 num_features=48
          2 x=[d[i][:num_features] for i in range(len(d))]
          3 y=[int(d[i][-1]) for i in range(len(d))]
```

```
In [175]: 1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25,random_stat
```

```
In [176]: 1 #calculating likelihood estimations
          2 x_train_class0=[x_train[i] for i in range(len(x_train)) if y_train[i]==0]
          3 x_train_class1=[x_train[i] for i in range(len(x_train)) if y_train[i]==1]
```

```
In [177]: 1 likelihood_class0=np.mean(x_train_class0,axis=0)/100.0
          2 likelihood_class1=np.mean(x_train_class1,axis=0)/100.0
```

```
In [178]: 1 def cal_log_likelihood(feature_vector,Class):
          2     assert len(feature_vector)==num_features
          3     log_likelihood=0.0
          4     if Class==0:
          5         for feature_index in range(len(feature_vector)):
          6             if feature_vector[feature_index]==1:
          7                 log_likelihood += np.log10(likelihood_class0[feature_index])
          8             elif feature_vector[feature_index] ==0:
          9                 log_likelihood+=np.log10(1.0-likelihood_class0[feature_index])
         10     elif Class==1:
         11         for feature_index in range(len(feature_vector)):
         12             if feature_vector[feature_index]==1:
         13                 log_likelihood+=np.log10(likelihood_class1[feature_index])
         14             elif feature_vector[feature_index]==0:
         15                 log_likelihood+=np.log10(1.0-likelihood_class1[feature_index])
         16     else:
         17         raise ValueError("Class takes integer values 0 or 1")
         18     return log_likelihood
```

```
In [179]: 1 num_class0=float(len(x_train_class0))
2 num_class1=float(len(x_train_class1))
3
4 prior_prob_class0=num_class0/(num_class0+num_class1)
5 prior_prob_class1=num_class1/(num_class0+num_class1)
6
7 log_prior_class0=np.log10(prior_prob_class0)
8 log_prior_class1=np.log10(prior_prob_class1)
```

```
In [180]: 1 def cal_class_posterior(feature_vector):
2     log_likelihood_class0=cal_log_likelihood(feature_vector,Class=0)
3     log_likelihood_class1=cal_log_likelihood(feature_vector,Class=1)
4
5     log_posterior_class0=log_likelihood_class0+log_prior_class0
6     log_posterior_class1=log_likelihood_class1+log_prior_class1
7
8     return log_posterior_class0,log_posterior_class1
9
```

```
In [181]: 1 def classify_spam(doc_vector):
2     feature_vector=[int(element>0.0) for element in doc_vector]
3     log_posterior_class0,log_posterior_class1=cal_class_posterior(feature_vector)
4     if log_posterior_class0>log_posterior_class1:
5         return 0
6     else:
7         return 1
8
```

```
In [182]: 1 predictions=[]
2 for email in x_test:
3     predictions.append(classify_spam(email))
```

```
In [183]: 1 def accuracy(predictions,y):
2     correct_count=0.0
3     for item in range(len(predictions)):
4         if predictions[item]==y[item]:
5             correct_count+=1.0
6     acc=correct_count/len(predictions)
7     return acc
```

```
In [184]: 1 accuracy(predictions,y_test)*100
```

Out[184]: 89.22675933970461