

Gaussian estimate conditional PDFs

In []:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn
5 import warnings
6 from sklearn import model_selection
7 import gensim
8
9
10 fil_data=pd.read_excel('filtered_data.xlsx')
11 X_train, X_test, Y_train, Y_test = model_selection.train_test_split(fil_data[['text']]
12                                                                    fil_data[
13 print(X_train.shape)
```

In []:

```
1 from gensim.test.utils import common_texts
2 from gensim.models.doc2vec import Doc2Vec, TaggedDocument
3
4 documents = [TaggedDocument(doc, [i]) for i, doc in enumerate(common_texts)]
5 model = Doc2Vec(documents, vector_size=100, window=1, min_count=1, workers=4, dm=1)
```

In []:

```
1 from gensim.test.utils import get_tmpfile
2
3 fname = get_tmpfile("my_doc2vec_model")
4
5 model.save(fname)
6 model = Doc2Vec.load(fname) # you can continue training with the loaded model!
```

In []:

```
1 X_train_no=np.zeros((87144,500))
2
3 for i in range(87144):
4     X_train_no[i,:]= model.infer_vector(np.asarray(X_train)[i])
```

In []:

```
1 Y_train = np.array(Y_train);
2 Y_test = np.array(Y_test);
```

In []:



```
1 Y_trainno=np.zeros((Y_train.size))
2 for i in range(Y_train.size):
3     if Y_train[i]=="BUSINESS":
4         Y_trainno[i]=1
5     elif Y_train[i]=="COMEDY":
6         Y_trainno[i]=2
7     elif Y_train[i]=="ENTERTAINMENT":
8         Y_trainno[i]=3
9     elif Y_train[i]=="FOOD & DRINK":
10        Y_trainno[i]=4
11    elif Y_train[i]=="HEALTHY LIVING":
12        Y_trainno[i]=5
13    elif Y_train[i]=="PARENTING":
14        Y_trainno[i]=6
15    elif Y_train[i]=="POLITICS":
16        Y_trainno[i]=7
17    elif Y_train[i]=="QUEER VOICES":
18        Y_trainno[i]=8
19    elif Y_train[i]=="SPORTS":
20        Y_trainno[i]=9
21    elif Y_train[i]=="STYLE & BEAUTY":
22        Y_trainno[i]=10
23    elif Y_train[i]=="TRAVEL":
24        Y_trainno[i]=11
25    elif Y_train[i]=="WELLNESS":
26        Y_trainno[i]=12
27
```

In []:



```
1 x1_train=[]
2 x2_train=[]
3 x3_train=[]
4 x4_train=[]
5 x5_train=[]
6 x6_train=[]
7 x7_train=[]
8 x8_train=[]
9 x9_train=[]
10 x10_train=[]
11 x11_train=[]
12 x12_train=[]
13 x1_test=[]
14 x2_test=[]
15 x3_test=[]
16 x4_test=[]
17 x5_test=[]
18 x6_test=[]
19 x7_test=[]
20 x8_test=[]
21 x9_test=[]
22 x10_test=[]
23 x11_test=[]
24 x12_test=[]
25 y1_train=[]
26 y2_train=[]
27 y3_train=[]
28 y4_train=[]
29 y5_train=[]
30 y6_train=[]
31 y7_train=[]
32 y8_train=[]
33 y9_train=[]
34 y10_train=[]
35 y11_train=[]
36 y12_train=[]
37 y1_test=[]
38 y2_test=[]
39 y3_test=[]
40 y4_test=[]
41 y5_test=[]
42 y6_test=[]
43 y7_test=[]
44 y8_test=[]
45 y9_test=[]
46 y10_test=[]
47 y11_test=[]
48 y12_test=[]
49
```

In []:



```
1 for i in range(Y_train.size):
2     if Y_train[i]=="BUSINESS":
3         x1_train.append(X_train_no[i,:])
4         y1_train.append(Y_trainno[i])
5     elif Y_train[i]=="COMEDY":
6         x2_train.append(X_train_no[i,:])
7         y2_train.append(Y_trainno[i])
8     elif Y_train[i]=="ENTERTAINMENT":
9         x3_train.append(X_train_no[i,:])
10        y3_train.append(Y_trainno[i])
11    elif Y_train[i]=="FOOD & DRINK":
12        x4_train.append(X_train_no[i,:])
13        y4_train.append(Y_trainno[i])
14    elif Y_train[i]=="HEALTHY LIVING":
15        x5_train.append(X_train_no[i,:])
16        y5_train.append(Y_trainno[i])
17    elif Y_train[i]=="PARENTING":
18        x6_train.append(X_train_no[i,:])
19        y6_train.append(Y_trainno[i])
20    elif Y_train[i]=="POLITICS":
21        x7_train.append(X_train_no[i,:])
22        y7_train.append(Y_trainno[i])
23    elif Y_train[i]=="QUEER VOICES":
24        x8_train.append(X_train_no[i,:])
25        y8_train.append(Y_trainno[i])
26    elif Y_train[i]=="SPORTS":
27        x9_train.append(X_train_no[i,:])
28        y9_train.append(Y_trainno[i])
29    elif Y_train[i]=="STYLE & BEAUTY":
30        x10_train.append(X_train_no[i,:])
31        y10_train.append(Y_trainno[i])
32    elif Y_train[i]=="TRAVEL":
33        x11_train.append(X_train_no[i,:])
34        y11_train.append(Y_trainno[i])
35    elif Y_train[i]=="WELLNESS":
36        x12_train.append(X_train_no[i,:])
37        y12_train.append(Y_trainno[i])
```

In []:



```
1 Y_testno=np.zeros((Y_test.size))
2 for i in range(Y_test.size):
3     if Y_test[i]=="BUSINESS":
4         Y_testno[i]=1
5     elif Y_test[i]=="COMEDY":
6         Y_testno[i]=2
7     elif Y_test[i]=="ENTERTAINMENT":
8         Y_testno[i]=3
9     elif Y_test[i]=="FOOD & DRINK":
10        Y_testno[i]=4
11    elif Y_test[i]=="HEALTHY LIVING":
12        Y_testno[i]=5
13    elif Y_test[i]=="PARENTING":
14        Y_testno[i]=6
15    elif Y_test[i]=="POLITICS":
16        Y_testno[i]=7
17    elif Y_test[i]=="QUEER VOICES":
18        Y_testno[i]=8
19    elif Y_test[i]=="SPORTS":
20        Y_testno[i]=9
21    elif Y_test[i]=="STYLE & BEAUTY":
22        Y_testno[i]=10
23    elif Y_test[i]=="TRAVEL":
24        Y_testno[i]=11
25    elif Y_test[i]=="WELLNESS":
26        Y_testno[i]=12
27
```

In []:



```
1 total=len(x1_train)+len(x2_train)+len(x3_train)+len(x4_train)+len(x5_train)+len(x6_tr
2 print(total)
3 total=len(y1_train)+len(y2_train)+len(y3_train)+len(y4_train)+len(y5_train)+len(y6_tr
4 print(total)
```

In []:

```

1
2 mu1_es=(1/len(x1_train))*(np.sum(x1_train,axis=0));
3 mu2_es=(1/len(x2_train))*(np.sum(x2_train,axis=0));
4 mu3_es=(1/len(x3_train))*(np.sum(x3_train,axis=0));
5 mu4_es=(1/len(x4_train))*(np.sum(x4_train,axis=0));
6 mu5_es=(1/len(x5_train))*(np.sum(x5_train,axis=0));
7 mu6_es=(1/len(x6_train))*(np.sum(x6_train,axis=0));
8 mu7_es=(1/len(x7_train))*(np.sum(x7_train,axis=0));
9 mu8_es=(1/len(x8_train))*(np.sum(x8_train,axis=0));
10 mu9_es=(1/len(x9_train))*(np.sum(x9_train,axis=0));
11 mu10_es=(1/len(x10_train))*(np.sum(x10_train,axis=0));
12 mu11_es=(1/len(x11_train))*(np.sum(x11_train,axis=0));
13 mu12_es=(1/len(x12_train))*(np.sum(x12_train,axis=0));
14
15 #Estimated variance
16 cov1_es=(1/len(x1_train))*np.dot((np.transpose(x1_train-mu1_es)),(x1_train-mu1_es));
17 cov2_es=(1/len(x2_train))*np.dot((np.transpose(x2_train-mu2_es)),(x2_train-mu2_es));
18 cov3_es=(1/len(x3_train))*np.dot((np.transpose(x3_train-mu3_es)),(x3_train-mu3_es));
19 cov4_es=(1/len(x4_train))*np.dot((np.transpose(x4_train-mu4_es)),(x4_train-mu4_es));
20 cov5_es=(1/len(x5_train))*np.dot((np.transpose(x5_train-mu5_es)),(x5_train-mu5_es));
21 cov6_es=(1/len(x6_train))*np.dot((np.transpose(x6_train-mu6_es)),(x6_train-mu6_es));
22 cov7_es=(1/len(x7_train))*np.dot((np.transpose(x7_train-mu7_es)),(x7_train-mu7_es));
23 cov8_es=(1/len(x8_train))*np.dot((np.transpose(x8_train-mu8_es)),(x8_train-mu8_es));
24 cov9_es=(1/len(x9_train))*np.dot((np.transpose(x9_train-mu9_es)),(x9_train-mu9_es));
25 cov10_es=(1/len(x10_train))*np.dot((np.transpose(x10_train-mu10_es)),(x10_train-mu10_
26 cov11_es=(1/len(x11_train))*np.dot((np.transpose(x11_train-mu11_es)),(x11_train-mu11_
27 cov12_es=(1/len(x12_train))*np.dot((np.transpose(x12_train-mu12_es)),(x12_train-mu12_
28

```

In []:

```

1 X_test_no=np.zeros((42923,500))
2
3 for i in range(42923):
4     X_test_no[i,:]= model.infer_vector(np.asarray(X_test)[i])

```

In []:

```

1 prob_prior=np.zeros((12))
2 prob_prior[0]=len(x1_train)/total
3 prob_prior[1]=len(x2_train)/total
4 prob_prior[2]=len(x3_train)/total
5 prob_prior[3]=len(x4_train)/total
6 prob_prior[4]=len(x5_train)/total
7 prob_prior[5]=len(x6_train)/total
8 prob_prior[6]=len(x7_train)/total
9 prob_prior[7]=len(x8_train)/total
10 prob_prior[8]=len(x9_train)/total
11 prob_prior[9]=len(x10_train)/total
12 prob_prior[10]=len(x11_train)/total
13 prob_prior[11]=len(x12_train)/total
14 print(prob_prior)

```

In []:



```
1 from sklearn.metrics import confusion_matrix, accuracy_score
2 from scipy.stats import multivariate_normal
3 y_pred=np.zeros((42923))
4 for i in range(42923):
5     p1=multivariate_normal.pdf(X_test_no[i,:],mu1_es,cov1_es)*prob_prior[0]
6     p2=multivariate_normal.pdf(X_test_no[i,:],mu2_es,cov2_es)*prob_prior[1]
7     p3=multivariate_normal.pdf(X_test_no[i,:],mu3_es,cov3_es)*prob_prior[2]
8     p4=multivariate_normal.pdf(X_test_no[i,:],mu4_es,cov4_es)*prob_prior[3]
9     p5=multivariate_normal.pdf(X_test_no[i,:],mu5_es,cov5_es)*prob_prior[4]
10    p6=multivariate_normal.pdf(X_test_no[i,:],mu6_es,cov6_es)*prob_prior[5]
11    p7=multivariate_normal.pdf(X_test_no[i,:],mu7_es,cov7_es)*prob_prior[6]
12    p8=multivariate_normal.pdf(X_test_no[i,:],mu8_es,cov8_es)*prob_prior[7]
13    p9=multivariate_normal.pdf(X_test_no[i,:],mu9_es,cov9_es)*prob_prior[8]
14    p10=multivariate_normal.pdf(X_test_no[i,:],mu10_es,cov10_es)*prob_prior[9]
15    p11=multivariate_normal.pdf(X_test_no[i,:],mu11_es,cov11_es)*prob_prior[10]
16    p12=multivariate_normal.pdf(X_test_no[i,:],mu12_es,cov12_es)*prob_prior[11]
17    prob=np.column_stack((p1,p2,p3,p4,p5,p6,p7,p8,p9,p10,p11,p12))
18    print(prob)
19    try:
20        y_pred[i]=(np.where(prob[0,:] == np.amax(prob[0,:]))[0])+1
21    except:
22        y_pred[i]=1
23 print(confusion_matrix(Y_testno, y_pred))
24 print(accuracy_score(Y_testno, y_pred))
```

In []:



```
1 print(accuracy_score(Y_testno, y_pred))
```