

```
In [1]: #preprocess the data
#select the necessary libraries
import numpy as np
import matplotlib.pyplot as plt
from matplotlib import style
style.use("ggplot")
from sklearn.cluster import KMeans
import pandas as pd
from sklearn import preprocessing
from sklearn.preprocessing import Normalizer
from sklearn.decomposition import PCA
from sklearn.metrics import accuracy_score

In [2]: #select the data source
#the target data
path='D:/coding/divorce_csv'
data=pd.read_csv(path)
print("Selected Data Source:\n",Data)
```

```

Selected Data Source:
  Atr1  Atr2  Atr3  Atr4  Atr5  Atr6  Atr7  Atr8  Atr9  Atr10  ...  Atr46 \
0      2      2      4      1      0      0      0      0      0      0      ...  2
1      1      4      4      4      4      4      0      0      4      4      ...  2
2      2      2      2      2      1      3      2      1      1      2      ...  3
3      3      2      3      2      3      3      3      3      3      3      ...  2
4      2      2      1      1      1      1      0      0      0      0      ...  2
... ..
165     0      0      0      0      0      0      0      0      0      0      ...  1
166     0      0      0      0      0      0      0      0      0      0      ...  4
167     1      1      0      0      0      0      0      0      0      0      ...  3
168     0      0      0      0      0      0      0      0      0      0      ...  3
169     0      0      0      0      0      0      0      1      0      0      ...  3

  Atr47  Atr48  Atr49  Atr50  Atr51  Atr52  Atr53  Atr54  Class
0      1      3      3      3      2      2      3      2      1
1      2      3      4      4      4      4      2      2      1
2      2      3      1      1      1      2      2      2      1
3      2      3      3      3      3      2      2      2      1
4      1      2      3      2      2      2      1      0      1
... ..
165     0      4      1      1      1      2      2      2      0
166     1      2      2      2      2      2      3      2      2
167     0      2      0      1      1      3      0      0      0
168     3      2      2      3      2      4      3      1      0
169     4      4      0      1      3      3      3      1      0

[170 rows x 55 columns]

In [3]: #Transform the target to numeric array
Encoder=preprocessing.LabelEncoder()
Encoded_Data=Data.apply(preprocessing.LabelEncoder().fit_transform)
print("Transformed Data:\n",Encoded_Data)
Numeric_Array=Encoded_Data.values
print("Numeric Array:\n", Numeric_Array)

Transformed Data:
  Atr1  Atr2  Atr3  Atr4  Atr5  Atr6  Atr7  Atr8  Atr9  Atr10  ...  Atr46 \
0      2      2      4      1      0      0      0      0      0      0      ...  2
1      1      4      4      4      4      4      0      0      4      4      ...  2
2      2      2      2      2      1      3      2      1      1      2      ...  3
3      3      3      2      3      2      3      3      3      3      3      ...  2
4      2      2      1      1      1      1      0      0      0      0      ...  2
... ..
165     0      0      0      0      0      0      0      0      0      0      ...  1
166     0      0      0      0      0      0      0      0      0      0      ...  4
167     1      1      0      0      0      0      0      0      0      0      ...  3

```

```

167      1      0      0      0      0      0      0      0      1 ...    3
168      0      0      0      0      0      0      0      0      0 ...    3
169      0      0      0      0      0      0      0      0      0 ...    3

      Attr47  Attr48  Attr49  Attr50  Attr51  Attr52  Attr53  Attr54  Class
0          1      2      3      3      3      2      3      2      1
1          2      3      4      4      4      4      2      2      1
2          2      3      1      1      1      2      2      2      1
3          2      3      3      3      3      2      2      2      1
...
4          1      2      3      2      2      2      1      0      1
...
... .. . .. . .. . .. . .. . .. . .. . .. .
165      0      4      1      1      1      4      2      2      2
166      1      2      2      2      2      3      2      2      0
167      0      2      0      1      1      3      0      0      0
168      3      2      2      3      2      4      3      1      0
169      4      4      0      1      3      3      3      1      0

[170 rows x 95 columns]
Numeric Array
[[2 4 ... 2 1 1]
 [4 4 4 ... 2 2 1]
 [2 2 2 ... 2 2 1]
 ...
 [1 0 ... 0 0 0]
 [0 0 0 ... 3 1 0]
 [0 0 0 ... 3 1 0]]

In [4]: #select test sample
split=s0
Training_Sample=Numeric_Array[:,split,:]
print(("Training Sample:\n", Training_Sample))

('Training Sample:\n', array([[2, 2, 4, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 1, 1, 2, 1, 2, 1, 2, 0, 1, 2, 1, 3, 3, 2, 1, 1, 2,
3, 2, 1, 3, 3, 3, 2, 3, 2, 1, 1],
 [4, 4, 4, 4, 4, 0, 0, 4, 4, 4, 4, 4, 3, 4, 6, 4, 4, 4, 4, 4, 3, 2, 1, 1,
0, 2, 2, 1, 2, 0, 1, 1, 0, 4, 2, 4, 3, 0, 2, 3, 4, 2, 4, 2, 4, 2, 3, 4,
2, 2, 2, 3, 4, 4, 4, 4, 2, 2, 1],
 [2, 2, 2, 2, 1, 3, 2, 1, 1, 2, 3, 4, 2, 3, 3, 3, 3, 3, 3, 3, 3, 2, 1, 0,
1, 2, 2, 2, 2, 3, 2, 3, 3, 1, 1, 1, 1, 2, 1, 3, 3, 3, 3, 3, 2, 3,
2, 3, 2, 3, 1, 1, 1, 2, 2, 2, 1],
 [5, 2, 3, 2, 3, 3, 3, 3, 3, 3, 3, 4, 3, 3, 4, 3, 3, 3, 3, 3, 3, 4, 1, 1,
1, 1, 2, 1, 1, 1, 1, 3, 2, 3, 2, 2, 1, 1, 3, 3, 4, 4, 2, 2, 3, 2,
3, 2, 2, 3, 3, 3, 2, 2, 2, 1],
 [2, 2, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 2, 1, 2, 1, 0,
0, 0, 0, 2, 1, 2, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 2, 1, 0, 2, 3, 0])
```

[illegible]

```

2 1 3 3 3 3 2 3 2 3 1 1 1 2 2 2]
[3 2 3 2 3 3 3 3 3 4 3 3 4 3 3 3 3 3 4 1 1 1 1 2 1 1 1 1 1 3 2 3 2 2 1 1
3 3 4 2 2 2 3 2 2 3 2 3 3 3 2 2 2]
[2 2 1 1 1 1 0 0 0 0 0 0 1 1 1 1 1 1 1 1 0 0 0 0 2 1 1 1 1 1 1 1 1 1 1 1 0 0
0 0 2 1 0 2 3 0 2 2 1 3 2 2 2 1 0]
[0 0 0 1 0 0 0 0 1 0 2 1 0 2 0 1 0 1 0 0 0 0 2 2 0 0 0 0 4 1 1 1 1 1
1 2 0 2 2 1 2 3 0 2 2 1 1 1 2 0]
[3 3 3 2 1 3 4 3 2 2 2 2 2 3 3 3 3 2 3 3 3 3 2 3 3 2 2 2 1 2 2 1 1 2
2 2 2 3 3 3 3 4 3 2 3 2 3 2 3 2 2]
[2 1 2 2 2 1 0 3 2 4 3 2 3 4 3 2 3 2 1 2 1 2 3 3 2 2 3 1 1 0 2 2 1
4 4 4 4 4 3 2 0 6 1 2 2 2 1 1 0]
[2 1 2 0 0 4 1 3 3 3 3 3 3 3 3 3 3 3 2 2 2 3 2 3 2 3 1 1 1 1 1 1 1 1
1 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1]
[1 1 1 1 2 0 2 2 3 0 0 2 1 0 1 2 1 0 0 0 0 1 1 1 1 1 1 1 1 0 1 0 0
1 1 2 1 2 3 2 2 0 2 2 4 3 3]]

In [7]: %select output attribute from training sample
Ytrain_Sample=training.Sample[:, -1]
print("Output(Dependent) Attributes of training sample\n", Ytrain_Sample)

Output(Dependent) Attributes of training sample
[1 1 1 1 1 1 1 1 1]

In [8]: %select input attributes from test sample
Xtest_Sample=Test_Sample[:, :-1]
print("Independent Attributes of test sample\n\n", XTest_Sample)

Independent Attributes of test sample

[[4 4 4 ... 4 4 4]
 [4 4 4 ... 4 4 4]
 [3 4 3 ... 4 4 4]
 ...
 [1 1 0 ... 3 0 0]
 [0 0 0 ... 4 3 1]
 [0 0 0 ... 3 3 1]]

In [9]: %select output attributes from test sample

```

```
print("Dependent Attributes of test sample:\n\n", Actual_YTest_Sample)

Dependent Attributes of Test sample

[[ 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0]]

In [10]: #Compress input attributes data training sample into two attributes
# Defining PCA object for compressing the attributes
pca = PCA(n_components=2)
#use PCA objects to compress the input attributes of the Training Sample
Compressed_XTrain_Sample = pca.fit_transform(XTrain_Sample)
print("(Compressed input attributes of Training Sample:\n", Compressed_XTrain_Sample)

Compressed input attributes of Training Sample:

[[ 6.36181355 -4.35778468]
 [-6.9591194  6.97974282]
 [-1.86934224 -2.66948614]
 [-5.12634652  1.94824966]
 [ 6.75631959  8.08437835]
 [ 7.378128  -1.64959095]
 [-4.48735944 -1.46795007]
 [-3.66241826 -2.41908773]
```

```
[ -2.476263135 -5.28877 ]
[ 4.33489369  0.9637898 ] ]

In [11]: Compressed_XTest_Sample = pca_fit.transform(XTest_Sample)
print("Compressed input attributes of Training Sample:\n", Compressed_XTest_Sample)

Compressed input attributes of Training Sample:
[[ [ 15.01521765  0.2461208 ]
   [ 15.01521765  0.2461208 ]
   [ 15.68871096  0.40729336 ]
   [ 15.68871096  0.40729336 ]
   [ 15.64125871  0.58312961 ]
   [ 15.68998274  -0.62677004 ]
```

[13.]	68998274	-0.62677604	
[13.]	57965234	-1.23662623	
[13.]	14518801	-0.26826504	
[13.]	57965234	-1.23662623	
[13.]	12497241	-1.06981971	
[13.]	12497241	-1.06981971	
[12.]	8316997	2.11120323	
[12.]	47182686	-0.5885911	
[12.]	55962608	-1.51286938	
[12.]	47182686	-0.5885911	
[12.]	39856914	-0.42766145	
[12.]	37761977	-0.19524271	
[12.]	38496343	-0.69810622	
[12.]	47398985	-1.6224745	
[12.]	38496343	-0.69810622	
[12.]	36276825	-0.24092177	
[12.]	34128987	-0.47334065	
[12.]	26866696	-0.0724291	
[12.]	27463393	-1.16529095	
[11.]	6196521	-1.14115569	
[11.]	87438832	-0.35105671	
[11.]	64063963	-0.40315851	
[11.]	64063963	-0.40315851	
[11.]	68946711	-0.87931120	
[11.]	16013567	-1.19094717	
[11.]	26021455	-2.63187896	
[11.]	77632654	-1.45217051	

[11.16897933	- 1.69381968]
[11.25021455	- 2.61377896]
[11.14659196	- 1.46139194]
[11.16897933	- 1.69381968]
[11.14659196	- 1.46139194]
[11.25021455	- 2.61377896]
[11.16897933	- 1.69381968]
[11.25021455	- 2.61377896]
[11.67327224	- 1.36955229]
[11.61604507	- 0.39855033]
[10.71139121	- 0.89912643]
[10.94272624	- 0.23877817]
[19.09169262	- 0.18727176]

[10	493138	-2.97533832]
[10	493138	-2.97533832]
[10	38438779	-0.15065846]
[9	95119524	-2.68951514]
[9	95119524	-2.68951514]
[9	92970786	-2.4570964]
[9	95119524	-2.68951514]
[9	95119524	-2.68951514]
[9	94061505	-2.29086791]
[9	73251403	-0.79330525]
[9	71162685	-0.47680852]
[9	79015187	-1.39525483]
[9	63770693	-0.31094487]
[8	60215655	-3.85139575]
[9	49807188	-0.11518115]
[8	52032421	-2.52855173]
[8	62430314	-1.42897707]
[8	24330761	-1.86340753]
[8	06533586	-0.87955803]
[7	67139419	-3.1352999]
[5	58746826	-6.57715092]
[5	45109088	-1.31705767]
[5	63478512	-6.74343646]
[5	62784416	-5.27116594]
[4	23201682	-0.25563089]
[9	95809184	-3.84944742]
[9	70888613	-4.47558975]

```

-2.84602474 -0.20756240
-11.73217678 -2.29381528
-11.06102786 -0.24590244
-11.266818 -0.75588764
-11.37667236 -1.17658902
-11.01061583 -4.63737638
-10.71966982 -2.19928271
-11.02751695 -0.54373976
-10.90586627 -0.57698286
-10.13092371 -0.83636777
-10.07917386 -1.33131265
-9.9150317 -2.58229108
-10.77339074 -0.38447131
-10.11028757 -2.34561243
-10.10682272 -1.51871386
-10.53251987 -0.89479547
-10.84344696 -1.190844
-10.11699831 -0.35795755
-10.12721967 -0.118351
-10.14034718 -2.24473956
-9.9616849 -1.07289662
-10.15297121 -2.07833403
-9.65264669 -0.8394796
-10.1192188 -2.7538597
-9.64095546 -0.41503953
-9.65604686 -2.0255622
-9.49741134 -1.0743855
-9.48118265 -0.90229741
-9.38702931 -1.5406265
-9.51451876 -2.56622635
-9.03623688 -4.10912181
-9.08531027 -0.30787029
-9.66781809 -2.27676154
-9.67195772 -2.489888254
-9.04974507 -1.83082126
-9.4222482 -0.395662
-9.49937675 -0.52134363
-9.39943381 -0.41209569
-8.89124312 -4.05621579
-9.10955393 -0.02814904
-9.44524683 -0.53299767
-9.47198831 -1.44264706
-8.75349155 -1.47702453
-8.75349155 -1.47702453
-9.03589981 -2.60233179
-8.26287577 -2.18455125
-9.20723335 -2.90534591
-8.92933314 -0.98224272
-8.78382322 -0.25779692
-9.3261147 -0.47574662
-8.97952672 -0.80265764
-9.0957554 -0.06934773
-9.0102426 -1.9049662
-8.65457638 -0.11842163
-9.17542043 -2.43265738
-9.17542043 -2.43265738
-8.74086158 -0.20845271
-8.8340794 -0.66389650
-8.64095128 -2.32923657
-8.99491834 -4.45103769
-8.00299014 -1.35445751
-8.69114341 -0.32827108
-9.046822 -3.8450136
-8.65106639 -0.59021183
-8.78087287 -3.66874861
-8.684139 -1.26379765
-8.9092272 -4.62002304
-8.9064893 -3.83259997
-8.38769674 -0.79899572
-8.68949986 -0.58580734
-8.44821052 -1.8720181
-8.19472387 -1.61459195
-8.30067113 -1.79808061
-8.29402152 -1.43898808
-8.3279063 -1.23196773
-8.13819332 -1.20867867
-8.09864411 -0.93644064
-8.12051972 -0.7641283
-8.0281425 -3.31973116
-8.07039022 -4.03402034
-8.19226721 -0.83524766
-8.30940794 -2.40723665
-8.73397451 -1.47764654
-8.2363581 -3.10968905
-8.12098504 -0.44136881
-8.14974944 -2.90298564
-8.31713534 -1.76697272]

```

```
KNNClassifier_Object = KNeighborsClassifier(n_neighbors=5)
#Train KNN classifier model
Trained_KNNClassifier_Object=KNNClassifier_Object.fit(Compressed_XTrain_Sample, YTrain_Sample)
print(Trained_KNNClassifier_Object)

KNeighborsClassifier()

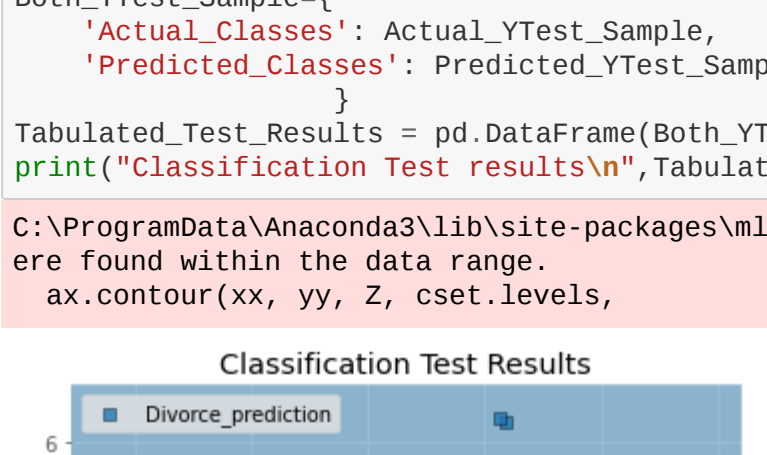
In [13]: from mlxtend.plotting import plot_decision_regions
from sklearn.metrics import confusion_matrix, plot_confusion_matrix
import matplotlib.pyplot as plt

#test it can classify new data and visualize the classifier
#use predict() function to predict classes of Test Sample
Predicted_YTest_Sample = Trained_KNNClassifier_Object.predict(Compressed_XTest_Sample)
#visualize the testing results using decision boundaries
from mlxtend.plotting import plot_decision_regions

plot_decision_regions(Compressed_XTest_Sample, Predicted_YTest_Sample, clf=Trained_KNNClassifier_Object, legend=2)
plt.title('Classification Test Results')
legend=plt.legend()
legend.get_texts()[0].set_text('Divorce_prediction')

plt.show()

#visualize Testing results using a Data Frame
Both_YTest_Samples=
```



```

Class_0 Actual_Classes Predicted_Classes
0      1      1
1      1      1
2      1      1
3      1      1
4      1      1
5      1      1
... ..
155     0      1
156     0      1
157     0      1
158     0      1
159     0      1

[160 rows x 2 columns]

In [34]: !pip install pandas_ml

Requirement already satisfied: pandas_ml in c:\programdata\anaconda3\lib\site-packages (0.6.1)Note: you may need to restart the kernel to use updated packages.
Requirement already satisfied: anaconda in c:\programdata\anaconda3\lib\site-packages (from pandas_ml) (1.1.10)
Requirement already satisfied: pandas>=0.19.0 in c:\programdata\anaconda3\lib\site-packages (from pandas_ml) (1.0.5)
Requirement already satisfied: pytz>=2017.2 in c:\programdata\anaconda3\lib\site-packages (from pandas>=0.19.0->pandas_ml) (2020.5)
Requirement already satisfied: python-dateutil>=2.6.1 in c:\programdata\anaconda3\lib\site-packages (from pandas>=0.19.0->pandas_ml) (2.8.1)
Requirement already satisfied: numba>=1.13.3 in c:\programdata\anaconda3\lib\site-packages (from pandas>=0.19.0->pandas_ml) (0.53.1)

```

```
as_ml) (1.18.5)
Requirement already satisfied: sklearn<1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil[>=2.6.1-> pandas==0.19.0->pandas_ml) (1.15.0)

In [15]: !pip install jaccard_score

Note: you may need to restart the kernel to use updated packages.

ERROR: Could not find a version that satisfies the requirement jaccard_score (from versions: none)
ERROR: No matching distribution found for jaccard_score

In [16]: from sklearn.metrics import accuracy_score
```

```

# Fit and validate the classifier
evaluate the model to ensure its performance
from sklearn.metrics import confusion_matrix
Test_Results = pd.DataFrame(Both_Test_Sample)
from sklearn.metrics import accuracy_score
Confusion_Matrix = confusion_matrix(Test_Results['Actual_Classes'], Test_Results['Predicted_Classes'])
print(Confusion_Matrix)

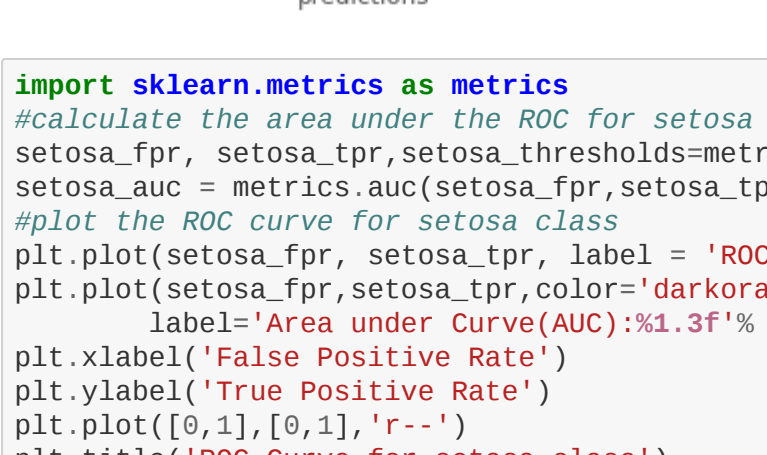
import seaborn as sns
visualize the validation results using a heat map
confusion_matrix=pd.crosstab(Test_Results['Actual_Classes'],
                             Test_Results['Predicted_Classes'],
                             rownames=['Actual_Classes'],

```

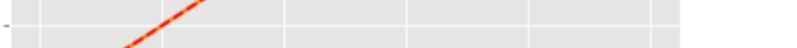
```
ax=plt.subplot()
sns.heatmap(confusion_matrix(), annot=True, ax=ax);
ax.set_xlabel('Predictions');ax.set_ylabel('Actual cases');
ax.set_title('Confusion Matrix Heatmap');
ax.xaxis.set_ticklabels(['No', 'Yes']);ax.yaxis.set_ticklabels(['No', 'Yes']);
plt.show()
```

```
[[ 0 86]
 [ 0 74]]
```

	No (Actual)	Yes (Actual)
No (Predicted)	0	86
Yes (Predicted)	0	74



```
plt.title("ROC Curve for Setosa class")
plt.legend(loc="upper right")
plt.show()
```



The figure shows an ROC curve for the 'versicolor' color. The x-axis is labeled 'False Positive Rate' and ranges from 0.0 to 1.0. The y-axis is labeled 'True Positive Rate' and ranges from 0.0 to 1.0. A red diagonal line represents the baseline for a random classifier. A blue curve, representing the model's performance, starts at (0,0) and rises above the diagonal line, indicating better-than-random performance. The curve reaches a True Positive Rate of approximately 0.8 at a False Positive Rate of about 0.2, then levels off towards the top right corner.

```
[29]: #validation of the ROC Curve
#calculate are under the ROC curve for versicolor color
versicolor_fpr,versicolor_tpr,versicolor_thresholds = metrics.roc_curve(Actual_YTest_Sample, Predicted_YTest_Sample,
pos_label=1)
versicolor_auc=metrics.auc(versicolor_fpr,versicolor_tpr)
```

```

#plot ROC curve
plot.roc.curve for setosa class
plt.plot(versicolor_fpr, versicolor_tpr, label = 'ROC for setosa')
plt.plot(versicolor_fpr, versicolor_tpr, color='darkorange',
        label='Area under Curve(AUC):%.3f'% versicolor_auc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.plot([0,1],[0,1], 'r--')
plt.title('ROC Curve for versicolor class')
plt.legend(loc='lower right')
plt.show()

ROC Curve for versicolor class

```

```

0.0    0.2    0.4    0.6    0.8    1.0
False Positive Rate

In [31]: #validation of the ROC curve
#calculate area under the ROC curve for versicolor color
virginica_fpr,virginica_tpr,virginica_thresholds = metrics.roc_curve(Actual_Ytest_Sample, Predicted_Ytest_Sample,pos_label=1)
virginica_auc=metrics_auc(virginica_fpr,virginica_tpr)
#plot the ROC curve for setosa class
plt.plot(virginica_fpr, virginica_tpr, label = 'ROC Virginica')
plt.plot(virginica_fpr, virginica_tpr, color='darkorange',
         label='%Area under Curve(AUC)=%.3f'% virginica_auc)
plt.plot([0,1],[0,1], color='blue',label='Baseline')

```

```
plt.ylabel('True Positive Rate')
plt.plot([0,1],[0,1], 'r-')
plt.title('ROC Curve for virginica class')
plt.legend(loc='lower right')
plt.show()
```



```

In [29]: #plot the roc curve for all the classes
plt.plot(setosa_fpr, setosa_tpr, linestyle='--', label='setosa AUC %1.3f'%setosa_auc)

```

```
plt.plot(versicolor_tpr, versicolor_fpr, marker='o', label='versicolor AUC: 0.97')
plt.plot(virginica_tpr, virginica_fpr, marker='v', label='virginica AUC: 0.88')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.plot([0,1],[0,1], 'r--')
plt.title('ROC Curve for all classes')
plt.legend()
plt.show()
```

ROC curve for the virginica dataset. The plot shows True Positive Rate (Y-axis) versus False Positive Rate (X-axis). A red diagonal line represents random performance, and a blue curve represents the model's performance. The Area Under the Curve (AUC) is 0.960.