# ELG 5255: Applied Machine Learning

## Assignment 1

## Group23_HW1

## 1. Overview

The main objective of this assignment is to build some ML models and check their performance using many different packages and tools of python and learning new Techniques (OVR and OVO) to deal with muti-classification problems as it is binary classification problems.

## 2. Methodology

We followed some defined steps to obtain the aimed results:

1. **(a) Load the DUMD dataset and convert UNS column into numerical data**

```python
# Function to read the Dataset
@staticmethod
def readDataSet(DataSet_name, Sheet_Name):
    Extension = re.findall('((.csv)|(.xls)|(.xlsx))', DataSet_name)
    Extension = str(Extension)
    if '.csv' in Extension:
        Extension = '.csv'
    elif '.xls' in Extension:
        Extension = '.xls'
    elif '.xlsx' in Extension:
        Extension = '.xlsx'

    if Extension == ('.xls' or '.xlsx'):
        Assignment1.DataFrame = pd.read_excel(DataSet_name, sheet_name = Sheet_Name)
    elif Extension == '.csv':
        Assignment1.DataFrame = pd.read_csv(DataSet_name)
    return Assignment1.DataFrame
```
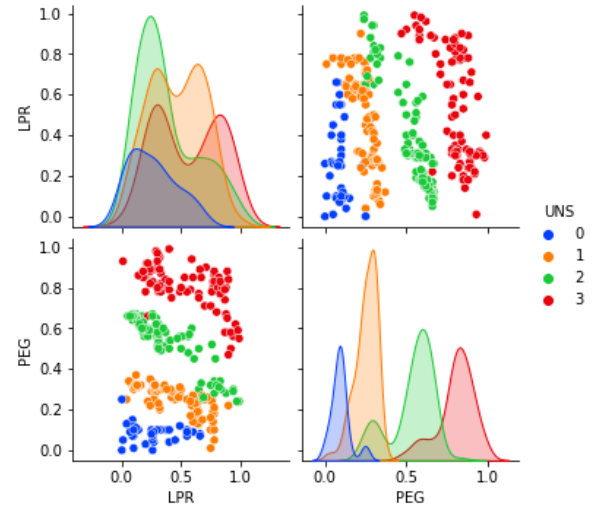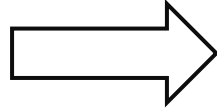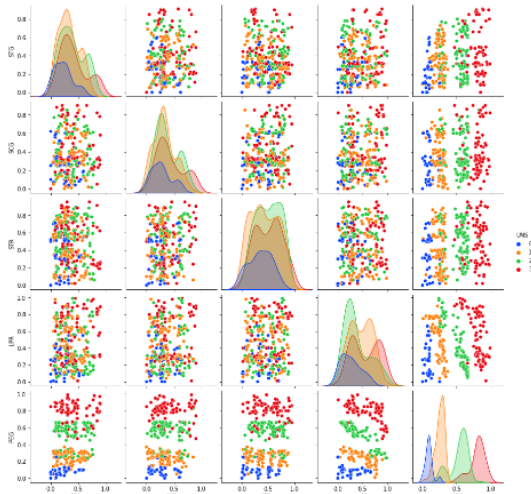
```python
39    # Function to Encode Labels of categorical data
      @staticmethod
40    def labelEncoder(DataFrame , Categorical_Column, ListOfClassesNames):
41        for i in range(len(ListOfClassesNames)):
42            DataFrame.loc[DataFrame[str(Categorical_Column)] == ListOfClassesNames[i], str(Categorical_Column)] = i
43        return DataFrame
44
```

Before and after.

| Index | STG | SCG | STR | LPR | PEG | UNS |
|---|---|---|---|---|---|---|
| 0 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | Very Low |
| 1 | 0.080 | 0.080 | 0.100 | 0.240 | 0.900 | High |
| 2 | 0.100 | 0.100 | 0.150 | 0.650 | 0.300 | Medium |
| 3 | 0.080 | 0.080 | 0.080 | 0.980 | 0.240 | Low |
| 4 | 0.090 | 0.150 | 0.400 | 0.100 | 0.660 | Medium |
| 5 | 0.100 | 0.100 | 0.430 | 0.290 | 0.560 | Medium |
| 6 | 0.200 | 0.140 | 0.350 | 0.720 | 0.250 | Low |
| 7 | 0.000 | 0.000 | 0.500 | 0.200 | 0.850 | High |
| 8 | 0.180 | 0.180 | 0.550 | 0.300 | 0.810 | High |
| 9 | 0.060 | 0.060 | 0.510 | 0.410 | 0.300 | Low |
| 10 | 0.200 | 0.200 | 0.700 | 0.300 | 0.600 | Medium |
| 11 | 0.120 | 0.120 | 0.750 | 0.350 | 0.800 | High |
| 12 | 0.050 | 0.070 | 0.700 | 0.010 | 0.050 | Very Low |

| Index | STG | SCG | STR | LPR | PEG | UNS |
|---|---|---|---|---|---|---|
| 0 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0 |
| 1 | 0.080 | 0.080 | 0.100 | 0.240 | 0.900 | 3 |
| 2 | 0.100 | 0.100 | 0.150 | 0.650 | 0.300 | 2 |
| 3 | 0.080 | 0.080 | 0.080 | 0.980 | 0.240 | 1 |
| 4 | 0.090 | 0.150 | 0.400 | 0.100 | 0.660 | 2 |
| 5 | 0.100 | 0.100 | 0.430 | 0.290 | 0.560 | 2 |
| 6 | 0.200 | 0.140 | 0.350 | 0.720 | 0.250 | 1 |
| 7 | 0.000 | 0.000 | 0.500 | 0.200 | 0.850 | 3 |
| 8 | 0.180 | 0.180 | 0.550 | 0.300 | 0.810 | 3 |
| 9 | 0.060 | 0.060 | 0.510 | 0.410 | 0.300 | 1 |
| 10 | 0.200 | 0.200 | 0.700 | 0.300 | 0.600 | 2 |
| 11 | 0.120 | 0.120 | 0.750 | 0.350 | 0.800 | 3 |
| 12 | 0.050 | 0.070 | 0.700 | 0.010 | 0.050 | 0 |

1. **(b) Choose two features based on PairPlot (the two most separable features that classification model can classify the different classes from each other easily) and "ExtraTreesClassifier"**



```
3        from sklearn.ensemble import ExtraTreesClassifier
```

```
44
45         # Function to detemine which features will be eliminated
46         @staticmethod
47         def FeatureSelection(X,Y):
48             Selector = ExtraTreesClassifier(n_estimators=2)
49             Selector = Selector.fit(X, Y)
50             return Selector.feature_importances_
51
```

**Output :** feature importance score for each feature.

[0.05059389, 0.02684576, 0.06839334, 0.15636901, 0.697798 ]

After choosing two features.



| Index | LPR | PEG | UNS |
|---|---|---|---|
| 0 | 0.000 | 0.000 | 0 |
| 1 | 0.240 | 0.900 | 3 |
| 2 | 0.650 | 0.300 | 2 |
| 3 | 0.980 | 0.240 | 1 |
| 4 | 0.100 | 0.660 | 2 |
| 5 | 0.290 | 0.560 | 2 |
| 6 | 0.720 | 0.250 | 1 |
| 7 | 0.200 | 0.850 | 3 |
| 8 | 0.300 | 0.810 | 3 |
| 9 | 0.410 | 0.300 | 1 |
| 10 | 0.300 | 0.600 | 2 |
| 11 | 0.350 | 0.800 | 3 |

## 1. (c) Apply SVM (rbf) and Perceptron

```python
# SVM Classifier
@staticmethod
def SVM(X ,Y ,GeneralizationTerm):
    ClassifierSVM = SVC(kernel="rbf", C = GeneralizationTerm, probability=True)
    ClassifierSVM.fit(X,Y)
    return ClassifierSVM

# PERCEP Classifier
@staticmethod
def PERCEP(X,Y,LearningRate,Epoch):
    ClassifierPERCEP = Perceptron(eta0=LearningRate, max_iter=Epoch)
    ClassifierPERCEP.fit(X,Y)
    return ClassifierPERCEP
```

## 1. (c) Classify testing data by using SVM and Perceptron classifiers. Provide accuracies, confusion matrix and decision boundaries for both classifiers.

| | SVM | Perceptron |
|---|---|---|
| **Accuracies on test** | 95.0 | 85.0 |
| **Confusion matrix** |  |  |
| **Decision boundaries** |  |  |

## 2. (a) Obtain the binarized labels (OvR)
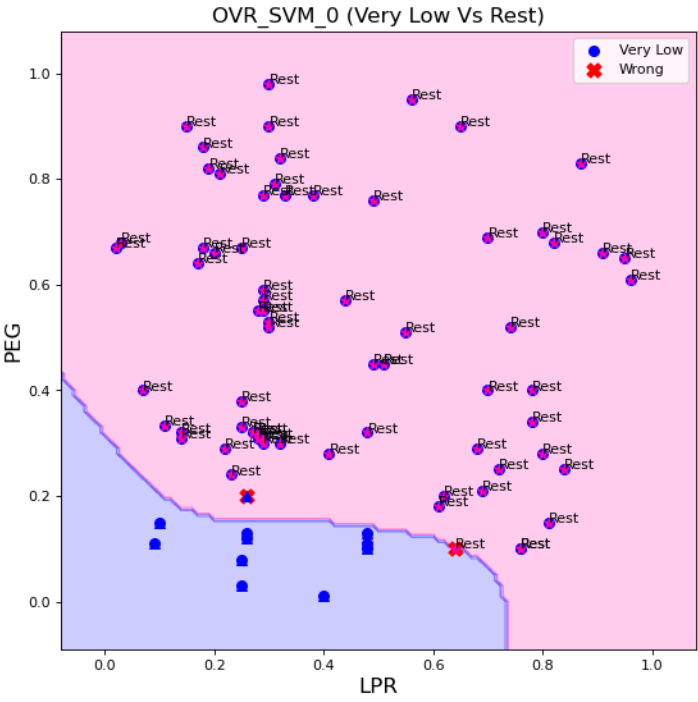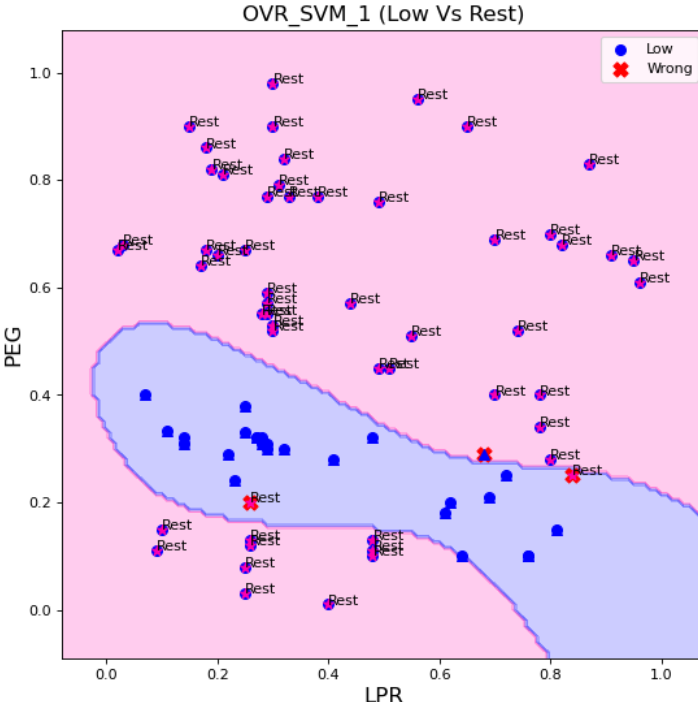
```
150        # One Vs Rest Spliting
151        @staticmethod
152        def OvR(DataSet,RequiredClass, Categorical_Column):
153            BinaryDF = DataSet.copy()
154            if RequiredClass == 0:
155                BinaryDF.loc[BinaryDF[str(Categorical_Column)] != RequiredClass, str(Categorical_Column)] = 2
156                BinaryDF.loc[BinaryDF[str(Categorical_Column)] == RequiredClass, str(Categorical_Column)] = 1
157                BinaryDF.loc[BinaryDF[str(Categorical_Column)] == 2, str(Categorical_Column)] = 0
158            else:
159                BinaryDF.loc[BinaryDF[str(Categorical_Column)] != RequiredClass, str(Categorical_Column)] = 0
160                BinaryDF.loc[BinaryDF[str(Categorical_Column)] == RequiredClass, str(Categorical_Column)] = 1
161            return BinaryDF
```

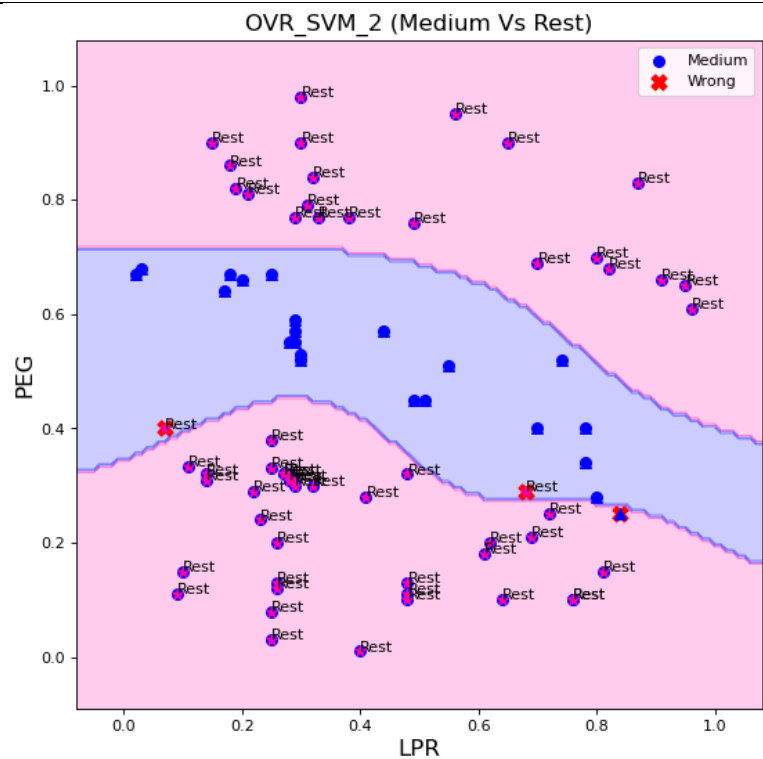**The required class = 1, and the rest = 0.**

**Y_train for each of 4 models.**

| OvR_DF_0_TrainY - Se | | OvR_DF_1_TrainY - Se | | OvR_DF_2_TrainY - Se | | OvR_DF_3_TrainY - Series | |
|---|---|---|---|---|---|---|---|
| Index | UNS | Index | UNS | Index | UNS | Index | UNS |
| 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 |
| 2 | 0 | 2 | 0 | 2 | 1 | 2 | 0 |
| 3 | 0 | 3 | 1 | 3 | 0 | 3 | 0 |
| 4 | 0 | 4 | 0 | 4 | 1 | 4 | 0 |
| 5 | 0 | 5 | 0 | 5 | 1 | 5 | 0 |
| 6 | 0 | 6 | 1 | 6 | 0 | 6 | 0 |
| 7 | 0 | 7 | 0 | 7 | 0 | 7 | 1 |
| 8 | 0 | 8 | 0 | 8 | 0 | 8 | 1 |
| 9 | 0 | 9 | 1 | 9 | 0 | 9 | 0 |
| 10 | 0 | 10 | 0 | 10 | 1 | 10 | 0 |
| 11 | 0 | 11 | 0 | 11 | 0 | 11 | 1 |
| 12 | 1 | 12 | 0 | 12 | 0 | 12 | 0 |

| Classes | SVM's accuracies on test | SVM's decision boundary | Comments |
|---|---|---|---|
| Very low vs Rest | 97.5 | OVR_SVM_0 (Very Low Vs Rest) | This model correctly classifies very low class from the rest.<br><br>And there are only 2 wrong points. |
| low vs Rest | 96.25 | OVR_SVM_1 (Low Vs Rest) | This model correctly classifies low class from the rest.<br><br>And there are only 3 wrong points. |

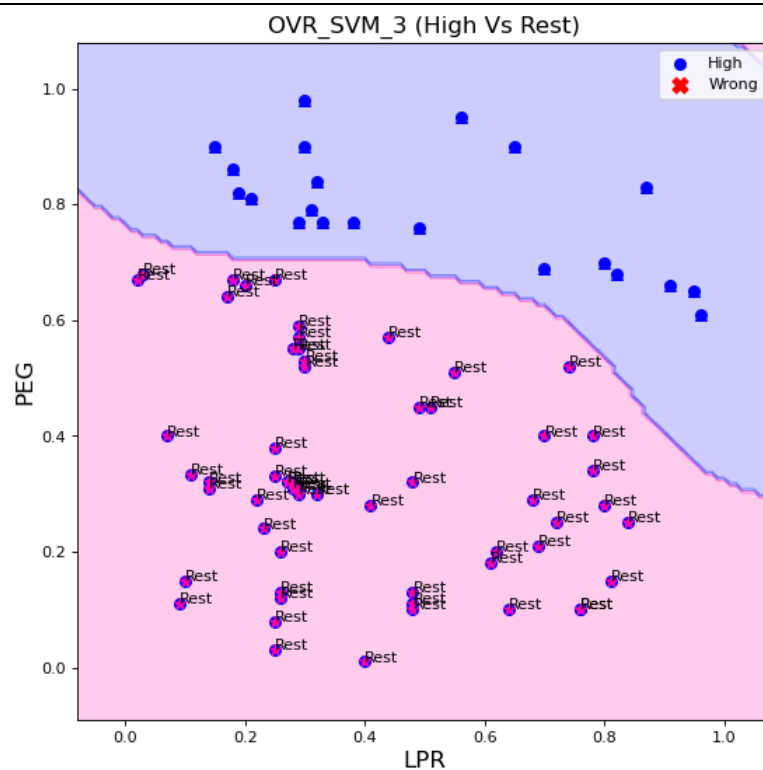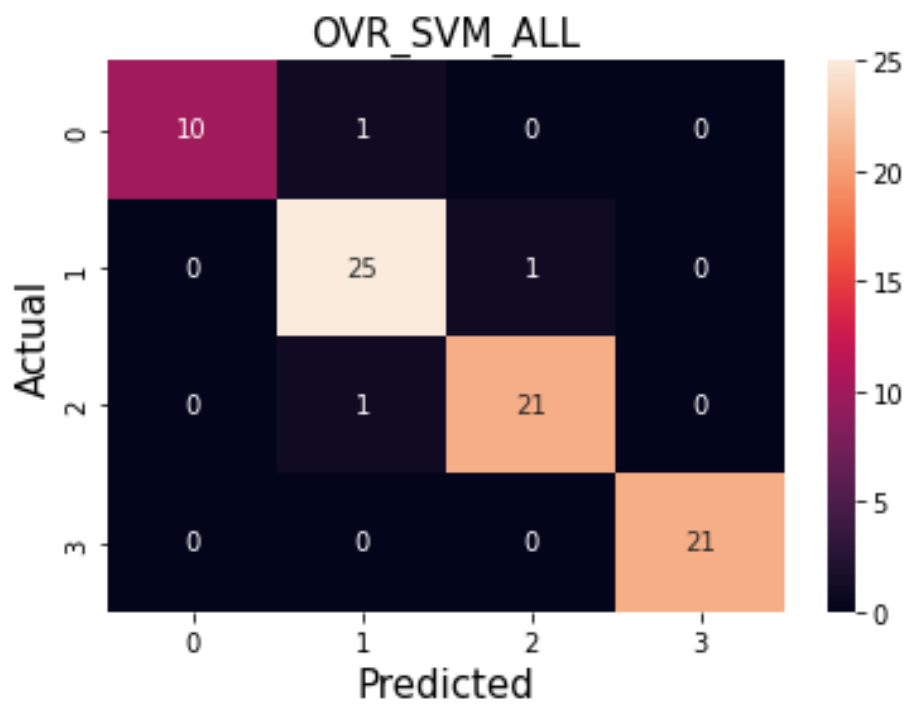| | | | |
|---|---|---|---|
| **Medium vs Rest** | 96.25 |  | **This model correctly classifies medium class from the rest.**<br><br>**And there are only 3 wrong points.** |
| **High VS Rest** | 100.0 |  | **This model correctly classifies very low class from the rest.**<br><br>**And there are no wrong points detected.** |

## 2. (b) OvR Overall accuracy (after aggregation) is = 96.25
We applied numpy.argmax() to these probabilities.

| Y_PredSVM_0_Prop - N | Y_PredSVM_1_Prop - | Y_PredSVM_2_Prop - | Y_PredSVM_3_Prop - NumPy object array |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0: 0.022 | 0: 0.009 | 0: 0.8660 | 0: 0.046 |
| 1: 0.000 | 1: 0.000 | 1: 0.9864 | 1: 0.001 |
| 2: 0.053 | 2: 0.013 | 2: 0.9490 | 2: 0.036 |
| 3: 0.866 | 3: 0.019 | 3: 0.0000 | 3: 0.003 |
| 4: 0.061 | 4: 0.000 | 4: 0.2165 | 4: 0.880 |
| 5: 0.002 | 5: 0.000 | 5: 0.0000 | 5: 1.000 |
| 6: 0.000 | 6: 0.002 | 6: 0.9905 | 6: 0.000 |
| 7: 0.050 | 7: 0.000 | 7: 0.0780 | 7: 0.952 |
| 8: 0.001 | 8: 0.965 | 8: 0.0018 | 8: 0.000 |
| 9: 0.003 | 9: 0.000 | 9: 0.9882 | 9: 0.001 |
| 10: 0.000 | 10: 0.001 | 10: 0.9949 | 10: 0.000 |

## To obtain this (OvR Y_Predict) and compare it with Y_Actual.

YProbALL_OvR - NumPy object array

| | 0 |
|---|---|
| 0 | 2 |
| 1 | 2 |
| 2 | 2 |
| 3 | 0 |
| 4 | 3 |
| 5 | 3 |
| 6 | 2 |
| 7 | 3 |
| 8 | 1 |
| 9 | 2 |
| 10 | 2 |

OVR_SVM_ALL

There are 3 wrong points predicted.

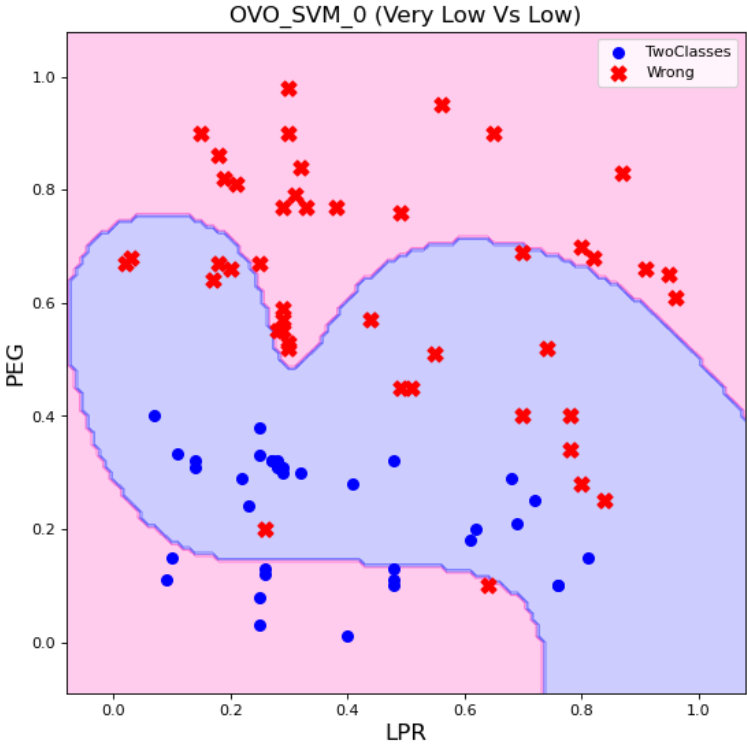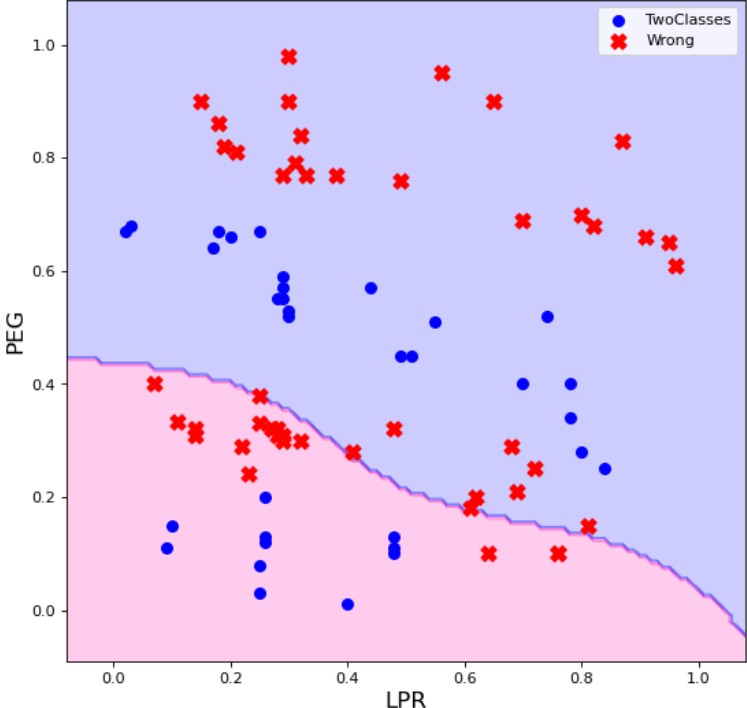### 3. (a) Obtain the binarized labels (OvR)

```
163        # One Vs Rest Spliting
164        @staticmethod
165        def OvO(FristClass, SecondClass, DataSet, Categorical_Column):
166            TwoClassesDF = DataSet.copy()
167            TwoClassesDF = TwoClassesDF[(TwoClassesDF[str(Categorical_Column)] == FristClass) | (TwoClassesDF[str(Categorical_
168            TwoClassesDF = TwoClassesDF.reset_index(drop=True)
169            TwoClassesDF.loc[TwoClassesDF[str(Categorical_Column)] == FristClass, str(Categorical_Column)] = FristClass
170            TwoClassesDF.loc[TwoClassesDF[str(Categorical_Column)] == SecondClass, str(Categorical_Column)] = SecondClass
171            return TwoClassesDF
```
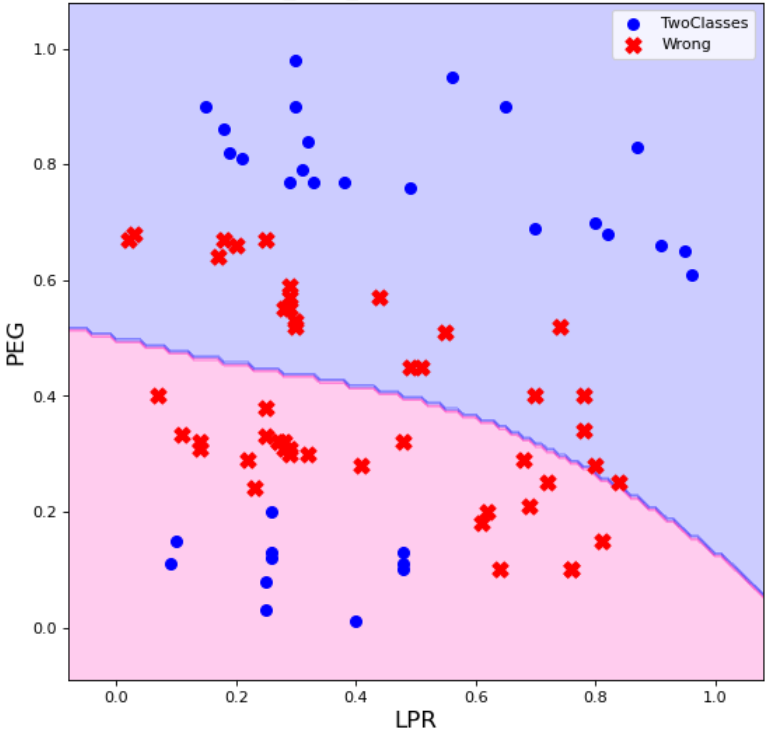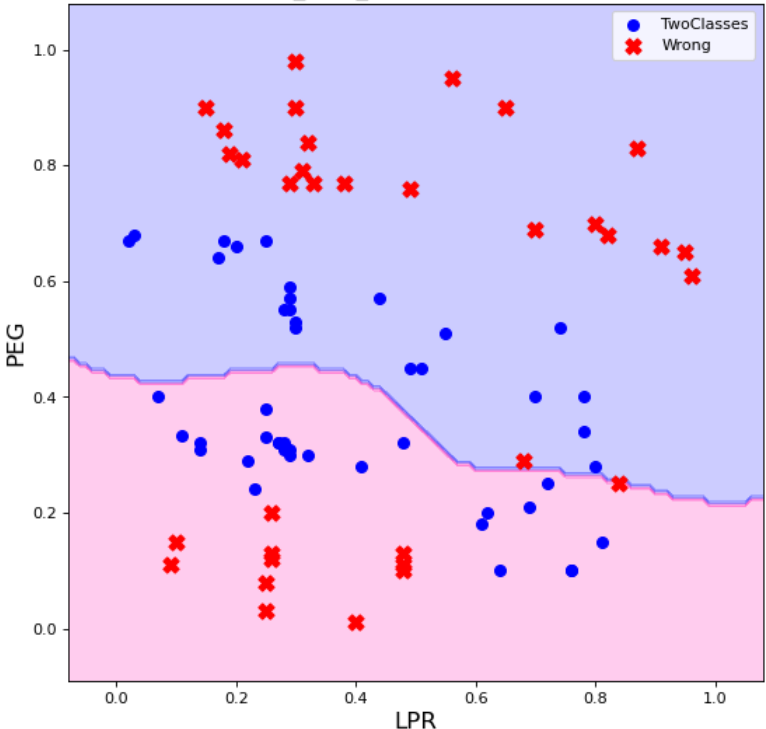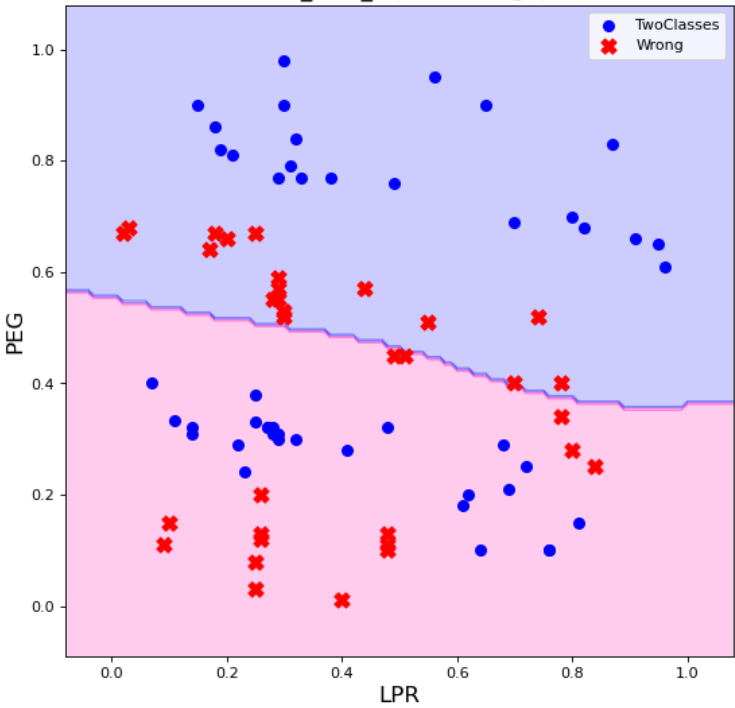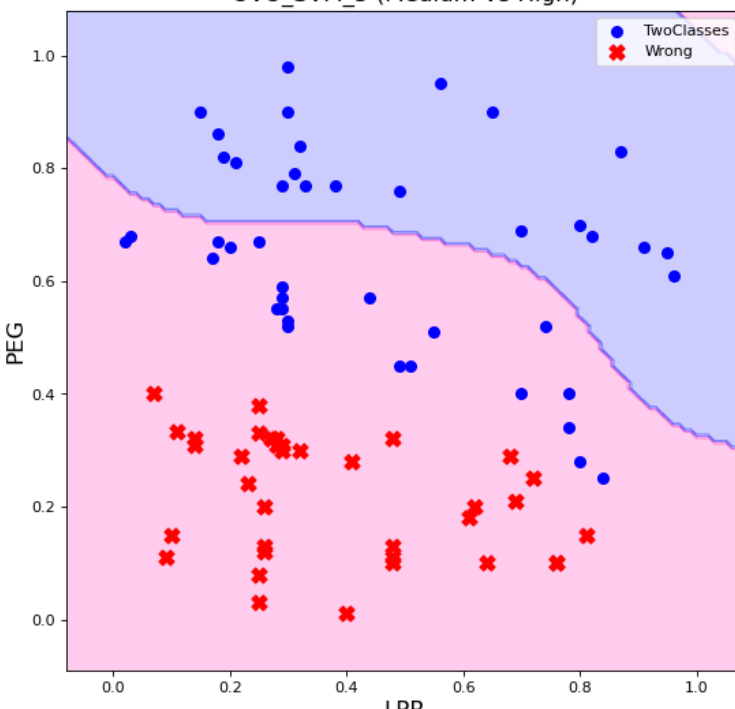
**Very Low -> 0, Low -> 1, Medium -> 2, High -> 3**

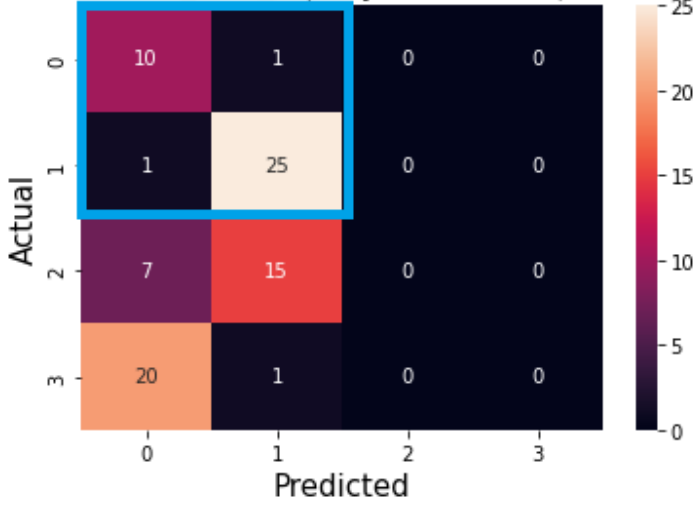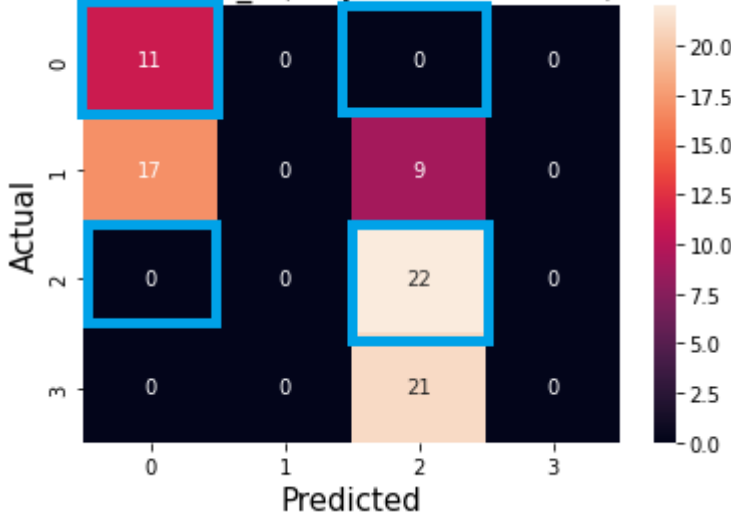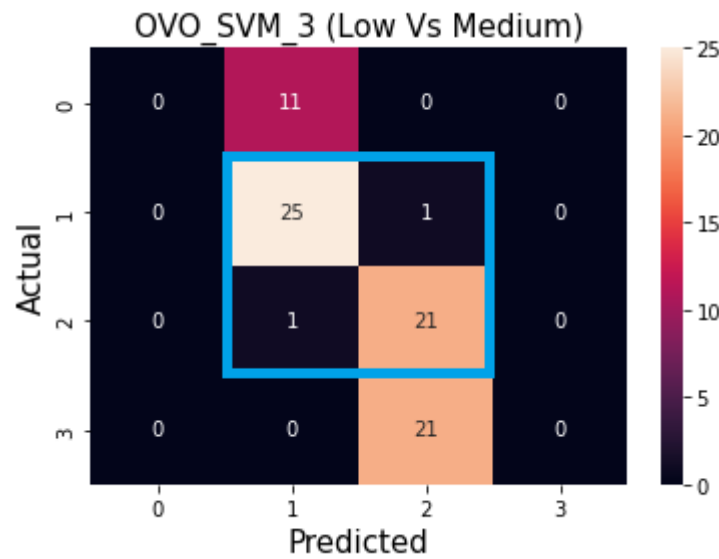**Y_train for each of 6 models.**

| OvO_DF_0_TrainY | | OvO_DF_1_TrainY | | OvO_DF_2_TrainY | | OvO_DF_3_TrainY | | OvO_DF_4_TrainY | | OvO_DF_5_TrainY - Series | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Index | UNS | Index | UNS | Index | UNS | Index | UNS | Index | UNS | Index | UNS |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 3 | 0 | 3 |
| 1 | 1 | 1 | 2 | 1 | 3 | 1 | 1 | 1 | 1 | 1 | 2 |
| 2 | 1 | 2 | 2 | 2 | 3 | 2 | 2 | 2 | 1 | 2 | 2 |
| 3 | 1 | 3 | 2 | 3 | 3 | 3 | 2 | 3 | 3 | 3 | 2 |
| 4 | 0 | 4 | 2 | 4 | 3 | 4 | 1 | 4 | 3 | 4 | 3 |
| 5 | 1 | 5 | 0 | 5 | 0 | 5 | 1 | 5 | 1 | 5 | 3 |
| 6 | 1 | 6 | 2 | 6 | 3 | 6 | 2 | 6 | 3 | 6 | 2 |
| 7 | 1 | 7 | 2 | 7 | 0 | 7 | 1 | 7 | 1 | 7 | 3 |
| 8 | 1 | 8 | 2 | 8 | 3 | 8 | 1 | 8 | 1 | 8 | 2 |
| 9 | 1 | 9 | 0 | 9 | 3 | 9 | 2 | 9 | 1 | 9 | 2 |
| 10 | 0 | 10 | 2 | 10 | 0 | 10 | 1 | 10 | 1 | 10 | 2 |
| 11 | 1 | 11 | 0 | 11 | 3 | 11 | 2 | 11 | 1 | 11 | 3 |
| 12 | 1 | 12 | 2 | 12 | 0 | 12 | 2 | 12 | | 12 | 3 |

| Classes | SVM's accuracies on test | SVM's decision boundary |
|---|---|---|
| **Very Low Vs Low** | **94.5 for the two classes only** **And** **43.75 (Compared to all classes)** |  |
| **Very Low Vs Medium** | **100 for the two classes only** **And** **41.25 (Compared to all classes)** |  |

| | | |
|---|---|---|
| **Very low Vs High** | **100 for the two classes only**<br><br>**And**<br><br>**40.0 (Compared to all classes)** | OVO_SVM_2 (Very Low Vs High) |
| **Low Vs Medium** | **95.8 for the two classes only**<br><br>**And**<br><br>**57.49 (Compared to all classes)** | OVO_SVM_3 (Low Vs Medium) |

| | | |
|---|---|---|
| **Low Vs High** | **100 for the two classes only**<br><br>**And**<br><br>**58.75 (Compared to all classes)** | <br>OVO_SVM_4 (Low Vs High) |
| **Medium Vs High** | **100 for the two classes only**<br><br>**And**<br><br>**53.75 (Compared to all classes)** | <br>OVO_SVM_5 (Medium Vs High) |

| Classes | SVM's Confusion matrix and comments |
|---|---|
| **Very Low Vs Low** |  This model correctly classifies Very low class from low class correctly. And there are 2 wrong points. |
| **Very Low Vs Medium** |  This model correctly classifies Very low class from medium class correctly. And there are no wrong points detected. |
| **Very low Vs High** |  This model correctly classifies Very low class from high class correctly. And there are no wrong points detected. |

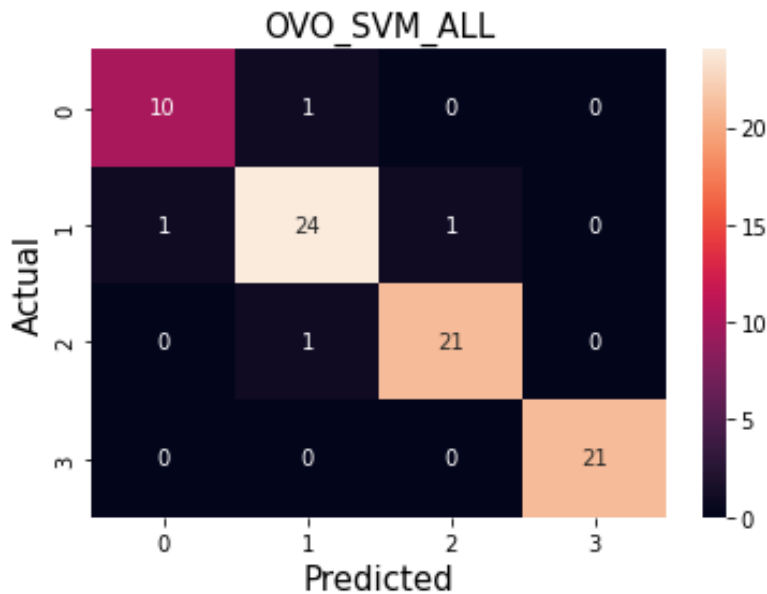| | | |
|---|---|---|
| **Low Vs Medium** |  OVO_SVM_3 (Low Vs Medium) | **This model correctly classifies low class from medium class correctly. And there are 2 wrong points.** |
| **Low Vs High** |  OVO_SVM_4 (Low Vs High) | **This model correctly classifies low class from high class correctly. And there are no wrong points detected.** |
| **Medium Vs High** |  OVO_SVM_5 (Medium Vs High) | **This model correctly classifies medium class from high class correctly. And there are no wrong points detected.** |

1. (b) **OvO Overall accuracy (after aggregation) is = 95.0**
   **We applied numpy.argmax() to these probabilities.**

| | VeryLow_Aggregation 0 | | Low_Aggregation - Nu 0 | | Medium_Aggregation 0 | | High_Aggregation - NumPy object array 0 |
|---|---|---|---|---|---|---|---|
| 0 | 0.107213 | 0 | 0.309074 | 0 | 0.936275 | 0 | 0.647438 |
| 1 | 0.0652639 | 1 | 0.592218 | 1 | 0.998943 | 1 | 0.343575 |
| 2 | 0.09403 | 2 | 0.273229 | 2 | 0.975206 | 2 | 0.657535 |
| 3 | 0.937252 | 3 | 0.716295 | 3 | 0.333967 | 3 | 0.012486 |
| 4 | 0.203998 | 4 | 0.14569 | 4 | 0.693825 | 4 | 0.956487 |
| 5 | 0.216133 | 5 | 0.144854 | 5 | 0.643092 | 5 | 0.995921 |
| 6 | 0.0699496 | 6 | 0.550776 | 6 | 0.985792 | 6 | 0.393482 |
| 7 | 0.216478 | 7 | 0.153331 | 7 | 0.650732 | 7 | 0.97946 |
| 8 | 0.410991 | 8 | 0.993057 | 8 | 0.559245 | 8 | 0.0367066 |
| 9 | 0.0896021 | 9 | 0.288878 | 9 | 0.995859 | 9 | 0.625661 |
| 10 | 0.0627806 | 10 | 0.533418 | 10 | 0.993407 | 10 | 0.410395 |

**To obtain this (OvO Y_Predict) and compare it with Y_Actual.**

YProbALL_OvO - NumPy object array

| | 0 |
|---|---|
| 0 | 2 |
| 1 | 2 |
| 2 | 2 |
| 3 | 0 |
| 4 | 3 |
| 5 | 3 |
| 6 | 2 |
| 7 | 3 |
| 8 | 1 |
| 9 | 2 |
| 10 | 2 |

**There are 4 wrong points predicted.**

## 4. (a) Conclusion

We have learned many new things during this assignment, and we have discovered some useful techniques like OvR and OvO. We have gotten familiar with new libraries. We have learnt how to select features based on PairPlot and discover which features is more important.

Now, we can say that we are capable of dealing with different types of SVM (We have tried rbf and linear) and we discover that the Perceptron algorithm is based on neural network.

We learnt how to divide a multi-classification problem into small Binary classification problems by using OVO and OVR techniques.

Those approaches divide a big classification problem into small ones by changing the labels of the target variables in the train dataset for each of binary classification problem (for OvR) and train one model on each dataset, for instance -> 1 for Very low class and 0 for others and so on….

Or by splitting our train dataset into different datasets for each model to train on it (for OvO) the number of binary datasets will be equal to (NumClasses * (NumClasses – 1)) / 2.

After that, we aggregate our results to obtain final accuracy for OvR and OvO and, we have obtained an accuracy of 96.25 from OvR which is higher then the normal SVM (rbf) approach that we tried in point number 1 (c).