# 02_ForestProject-Modeling&Presentation

August 19, 2020

# 1 Prediciton (Classification) the Types of Trees.

### 1.0.1 This Jupyter Notebook contains;

- Classification Models for Predicting the Types of Trees,
- Visualization of the Result.

### 1.0.2 What I Plan to implement in terms of ML models :

- SVM (I will use `LinearSVC` model from `sklearn.svm` module),
- XGBoost (I will use `XGBClassifier` model from `xgboost` module)
- Additionally :
  - Decision Tree (I will use `DecisionTreeClassifier` model from `sklearn.tree` module)
  - KNN (I will use `KNeighborsClassifier` model from `sklearn.neighbors` module)
  - LGBM (I will usem `LGBMClassifier` model from `lightgbm` module).

### 1.0.3 I will use `yellowbrick`, `seaborn` and `matplotlib` modules to visualize the model results.

### 1.0.4 Importing `covtype1.csv` dataset for modelling and required libraries.

```python
[58]: import pandas as pd
      import numpy as np
      import matplotlib.pyplot as plt
      import seaborn as sns
      from sklearn.preprocessing import StandardScaler
      from sklearn.decomposition import PCA
      from sqlalchemy import create_engine
      import warnings
      from IPython.core.pylabtools import figsize
      from scipy.stats import zscore
      from scipy import stats
      from numpy import percentile
      from sklearn.metrics import accuracy_score
      from sklearn.model_selection import train_test_split
      from sklearn.model_selection import cross_val_score
```

```python
from statsmodels.formula.api import ols
from scipy.stats import zscore
from sklearn.tree import DecisionTreeClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVC, LinearSVC
from sklearn.ensemble import RandomForestClassifier , GradientBoostingClassifier
import seaborn as sns
from xgboost import XGBClassifier
from lightgbm import LGBMClassifier
from sklearn.model_selection import TimeSeriesSplit
from yellowbrick.classifier import ClassificationReport
from yellowbrick.datasets import load_occupancy
from sklearn.metrics import f1_score
font_title = {'family': 'times new roman', 'color': 'darkred',
              'weight': 'bold', 'size': 14}


warnings.filterwarnings('ignore')
sns.set_style("whitegrid")


plt.rcParams['figure.dpi'] = 100
```

### 1.0.5 Reading dataset with pandas

```python
[135]: df = pd.read_csv("covtype1.csv")
```

```python
[8]: df.head()
```

```
[8]:    Elevation  Aspect  Slope  Horizontal_Distance_To_Roadways  Hillshade_9am  \
    0       2596      51      3                              510            221
    1       2590      56      2                              390            220
    2       2804     139      9                             3180            234
    3       2785     155     18                             3090            238
    4       2595      45      2                              391            220

       Hillshade_Noon  Horizontal_Distance_To_Fire_Points  Wilderness_Area1  \
    0             232                                6279                 1
    1             235                                6225                 1
    2             238                                6121                 1
    3             238                                6211                 1
    4             234                                6172                 1

       Wilderness_Area2  Wilderness_Area3  …  Soil_Type33  Soil_Type34  \
    0                 0                 0  …            0            0
    1                 0                 0  …            0            0
```

```
2                    0           0 …              0              0
3                    0           0 …              0              0
4                    0           0 …              0              0

    Soil_Type35  Soil_Type38  Soil_Type39  Soil_Type40  Cover_Type  \
0            0            0            0            0           5
1            0            0            0            0           5
2            0            0            0            0           2
3            0            0            0            0           2
4            0            0            0            0           5

    Square_Hypo_Distance  Average_Dist_Road_Hydro  Average_Elevation_Hydro
0                  66564                      384                     1298
1                  44980                      301                     1292
2                  76049                     1724                     1434
3                  72488                     1666                     1451
4                  23410                      272                     1297

[5 rows x 46 columns]
```

### 1.0.6  Modeling

```python
[23]: X = df.drop("Cover_Type", axis = 1)
```

```python
[24]: y = df["Cover_Type"]
```

- Splitting the data set into two pieces : Test split - Train split

```python
[25]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,␣
      ↪random_state=101)
```

### 1.0.7  XGBoost Classifer

```python
[26]: xgb_classifier = XGBClassifier()
      xgb_classifier.fit(X_train , y_train)
```

```
[26]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                    colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                    importance_type='gain', interaction_constraints='',
                    learning_rate=0.300000012, max_delta_step=0, max_depth=6,
                    min_child_weight=1, missing=nan, monotone_constraints='()',
                    n_estimators=100, n_jobs=0, num_parallel_tree=1,
                    objective='multi:softprob', random_state=0, reg_alpha=0,
                    reg_lambda=1, scale_pos_weight=None, subsample=1,
                    tree_method='exact', validate_parameters=1, verbosity=None)
```

```
[27]: y_predicted = xgb_classifier.predict(X_test)
```

```
[31]: y_predicted
```

```
[31]: array([7, 3, 3, …, 1, 3, 2], dtype=int64)
```

```
[44]: xgb_accuracy = accuracy_score(y_test, y_predicted)
```

```
[45]: # Very good!

      xgb_accuracy
```

```
[45]: 0.8714599602298091
```

**Now, I would like to show three plots for the results.** - Class Prediction Error Bar Plot - Confusion Matrix - Classification Report
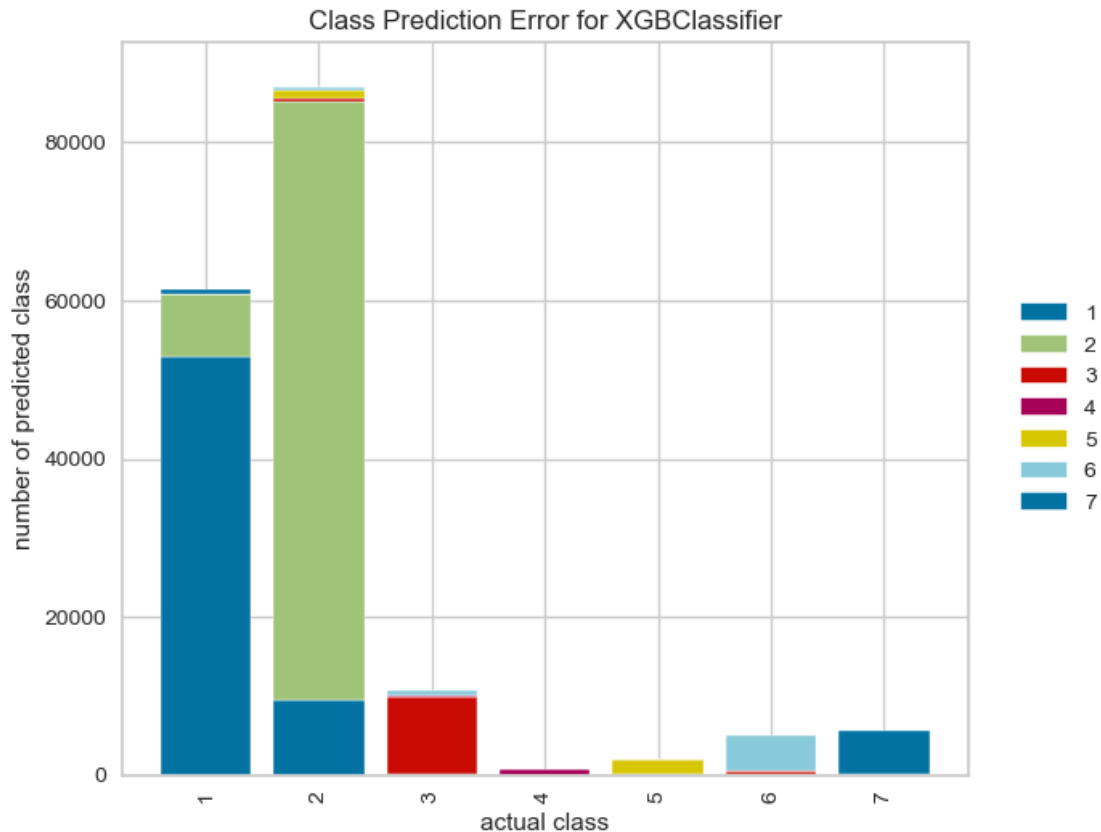
```
[42]: from sklearn.datasets import make_classification
      from sklearn.model_selection import train_test_split
      from sklearn.ensemble import RandomForestClassifier
      from yellowbrick.classifier import ClassPredictionError

      visualizer = ClassPredictionError(xgb_classifier)

      # Fit the training data to the visualizer
      visualizer.fit(X_train, y_train)

      # Evaluate the model on the test data
      visualizer.score(X_test, y_test)

      # Draw visualization
      visualizer.show()
```

## Class Prediction Error for XGBClassifier



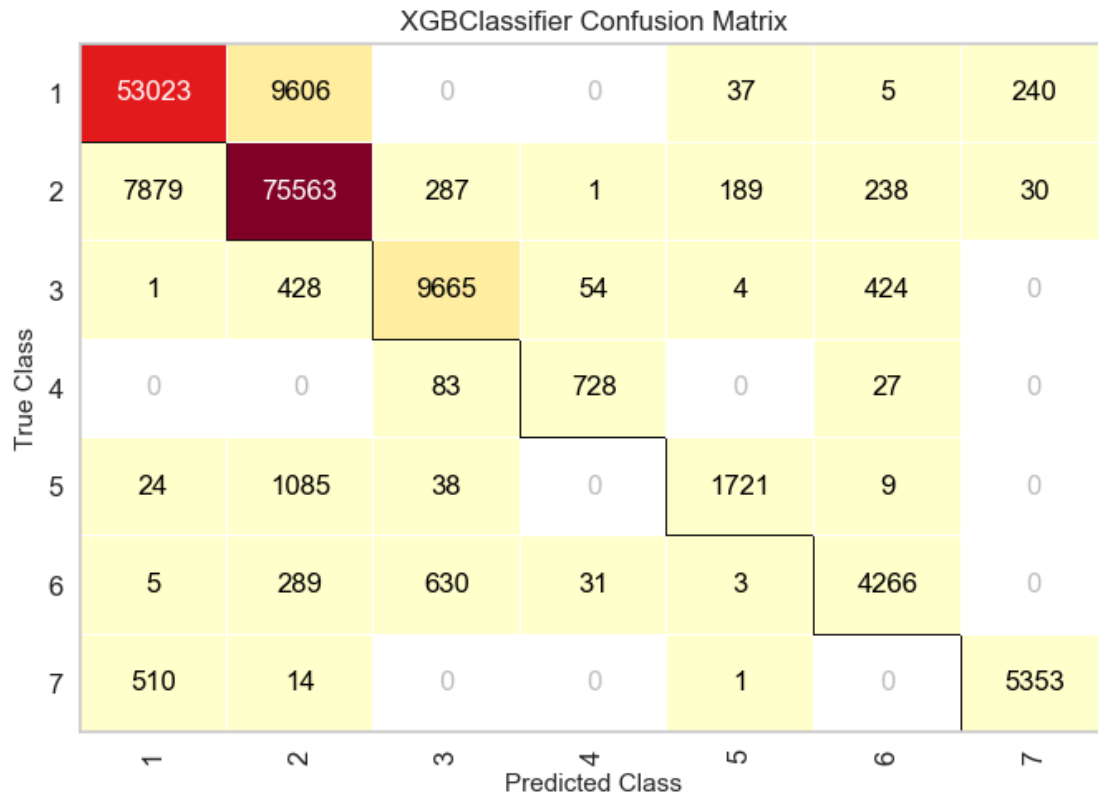[42]: `<matplotlib.axes._subplots.AxesSubplot at 0x2030637a048>`

```python
[41]: from sklearn.datasets import load_digits
      from sklearn.model_selection import train_test_split as tts
      from sklearn.linear_model import LogisticRegression
      from yellowbrick.classifier import ConfusionMatrix

      # The ConfusionMatrix visualizer taxes a model
      cm = ConfusionMatrix(xgb_classifier)

      # Fit fits the passed model. This is unnecessary if you pass the visualizer a␣
       ↪pre-fitted model
      cm.fit(X_train, y_train)

      # To create the ConfusionMatrix, we need some test data. Score runs predict()␣
       ↪on the data
      # and then creates the confusion_matrix from scikit-learn.
      cm.score(X_test, y_test)

      cm.show()
```
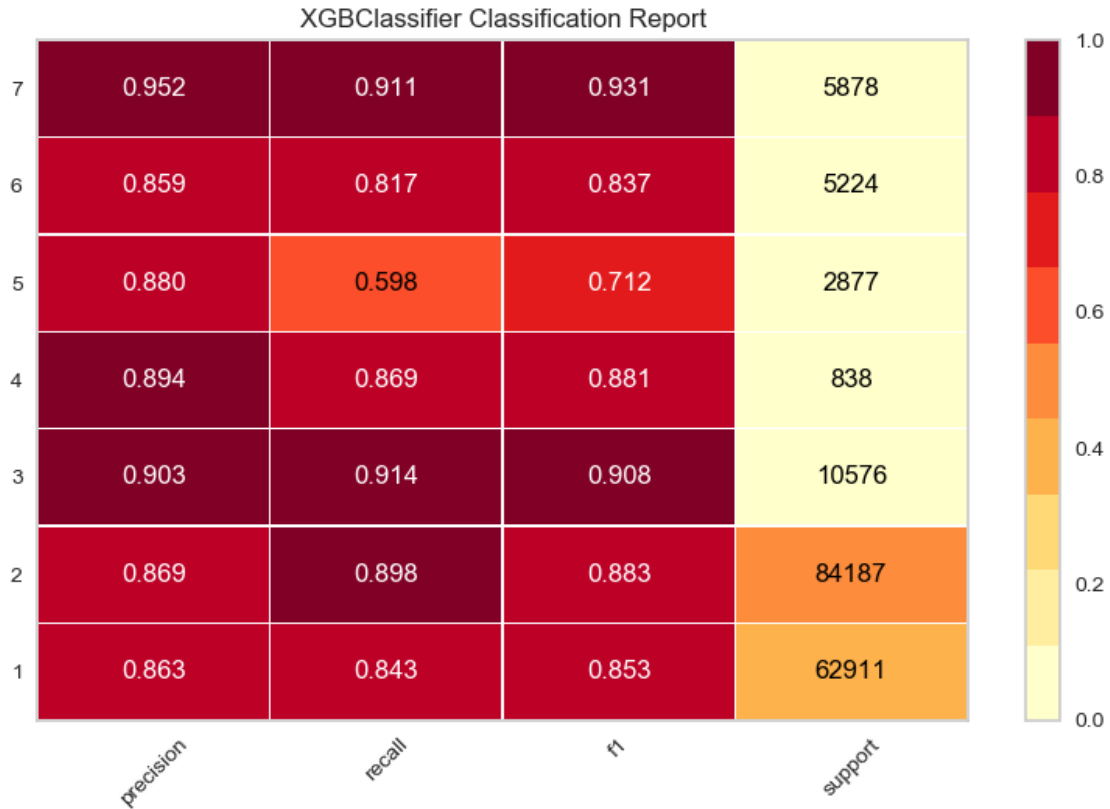
## XGBClassifier Confusion Matrix

| True Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | 53023 | 9606 | 0 | 0 | 37 | 5 | 240 |
| 2 | 7879 | 75563 | 287 | 1 | 189 | 238 | 30 |
| 3 | 1 | 428 | 9665 | 54 | 4 | 424 | 0 |
| 4 | 0 | 0 | 83 | 728 | 0 | 27 | 0 |
| 5 | 24 | 1085 | 38 | 0 | 1721 | 9 | 0 |
| 6 | 5 | 289 | 630 | 31 | 3 | 4266 | 0 |
| 7 | 510 | 14 | 0 | 0 | 1 | 0 | 5353 |

[41]: <matplotlib.axes._subplots.AxesSubplot at 0x203062dc390>

```python
[43]: from sklearn.model_selection import TimeSeriesSplit
      from sklearn.naive_bayes import GaussianNB

      from yellowbrick.classifier import ClassificationReport
      from yellowbrick.datasets import load_occupancy

      visualizer = ClassificationReport(xgb_classifier, support=True)

      visualizer.fit(X_train, y_train)        # Fit the visualizer and the model
      visualizer.score(X_test, y_test)        # Evaluate the model on the test data
      visualizer.show()
```

## XGBClassifier Classification Report



| | precision | recall | f1 | support |
|---|---|---|---|---|
| 7 | 0.952 | 0.911 | 0.931 | 5878 |
| 6 | 0.859 | 0.817 | 0.837 | 5224 |
| 5 | 0.880 | 0.598 | 0.712 | 2877 |
| 4 | 0.894 | 0.869 | 0.881 | 838 |
| 3 | 0.903 | 0.914 | 0.908 | 10576 |
| 2 | 0.869 | 0.898 | 0.883 | 84187 |
| 1 | 0.863 | 0.843 | 0.853 | 62911 |

[43]: `<matplotlib.axes._subplots.AxesSubplot at 0x203065f2668>`

- There is a little bit low `recall` score of the class 5. I conclueded that it is caused by the values in the data set. %59.8 of the class 5 is predicted true. Although the `recall` score level of class 5 is a little bit low, `precision` and `f1` score is quite well.

### 1.0.8 LinearSVC

[111]: `modelSVM = LinearSVC()`

[112]: `modelSVM.fit(X_train , y_train)`

[112]: 
```
LinearSVC(C=1.0, class_weight=None, dual=True, fit_intercept=True,
          intercept_scaling=1, loss='squared_hinge', max_iter=1000,
          multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
          verbose=0)
```

[114]: `pred = modelSVM.predict(X_test)`

[115]: `SVM_accuracy = accuracy_score(pred, y_test)`

7

```
[116]: SVM_accuracy
```

```
[116]: 0.23831968044709578
```

- It seems the model is failed. I decided to drop additional columns ('Average_Dist_Road_Hydro', 'Average_Elevation_Hydro') and re-split the dataset.

```
[68]: df['Square_Hypo_Distance'] = np.sqrt(df['Square_Hypo_Distance'])
      df1 = df.drop(['Average_Dist_Road_Hydro','Average_Elevation_Hydro'], axis = 1)
      X1 = df1.drop("Cover_Type", axis = 1)
      y1 = df1["Cover_Type"]
```

- Splitting the data set into two pieces : Test split - Train split

```
[96]: X1_train, X1_test, y1_train, y1_test = train_test_split(X1, y1, test_size=0.3,␣
      ↪random_state=101)
```

```
[118]: modelSVM1 = LinearSVC()
```

```
[119]: modelSVM1.fit(X1_train , y1_train)
```

```
[119]: LinearSVC(C=1.0, class_weight=None, dual=True, fit_intercept=True,
                 intercept_scaling=1, loss='squared_hinge', max_iter=1000,
                 multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
                 verbose=0)
```

```
[120]: pred1 = modelSVM1.predict(X1_test)
```

```
[121]: SVM_accuracy1 = accuracy_score(pred1, y1_test)
```

```
[122]: SVM_accuracy1
```

```
[122]: 0.5566145480054032
```

- I've doubled the **accuracy score** but it's still insufficient. SVM Classifier ;
  - SVM has been widely used in finance. For example, predicting stock price via SVM has been a acknowledged application in the industry.
  - In classification of text and handwritten objects, SVM performs well.
  - It may not be very successful in datasets with more than 100,000 data.

**Now, I would like to show three plots for the results.** - Class Prediction Error Bar Plot - Confusion Matrix - Classification Report
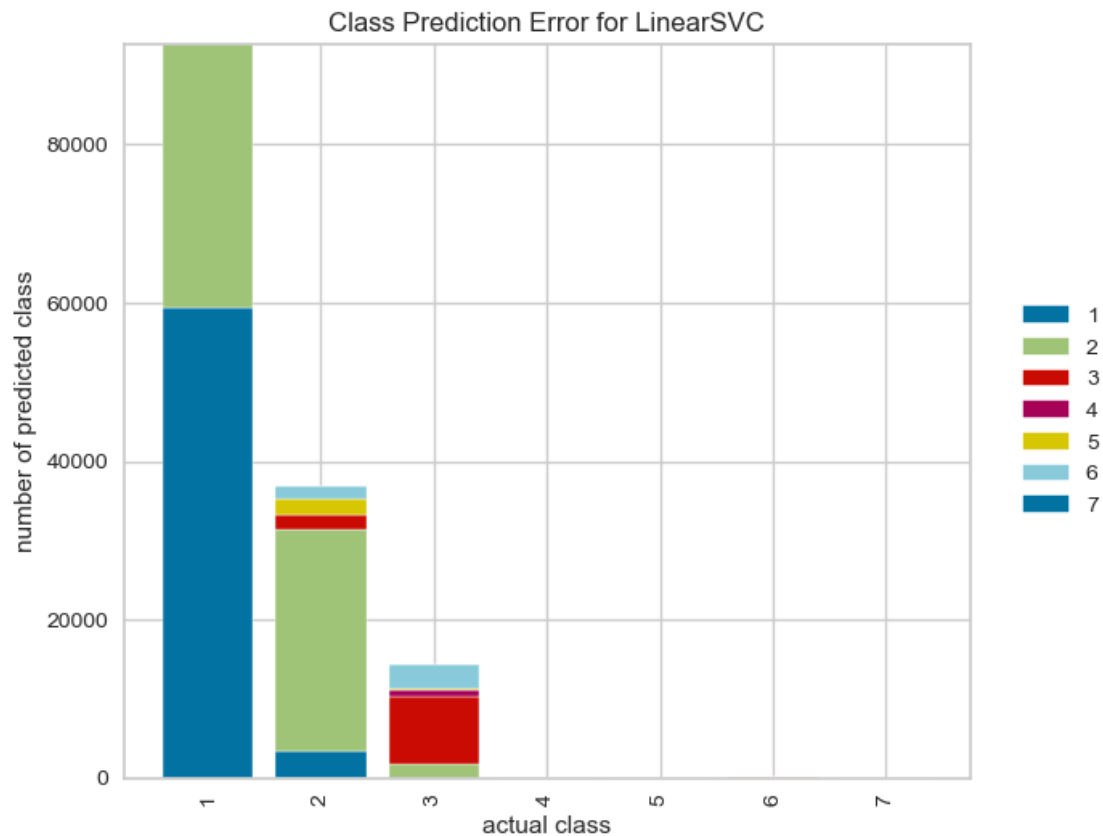
```
[123]: visualizer = ClassPredictionError(modelSVM1)

       # Fit the training data to the visualizer
       visualizer.fit(X1_train, y1_train)
```

```
# Evaluate the model on the test data
visualizer.score(X1_test, y1_test)

# Draw visualization
visualizer.show()
```
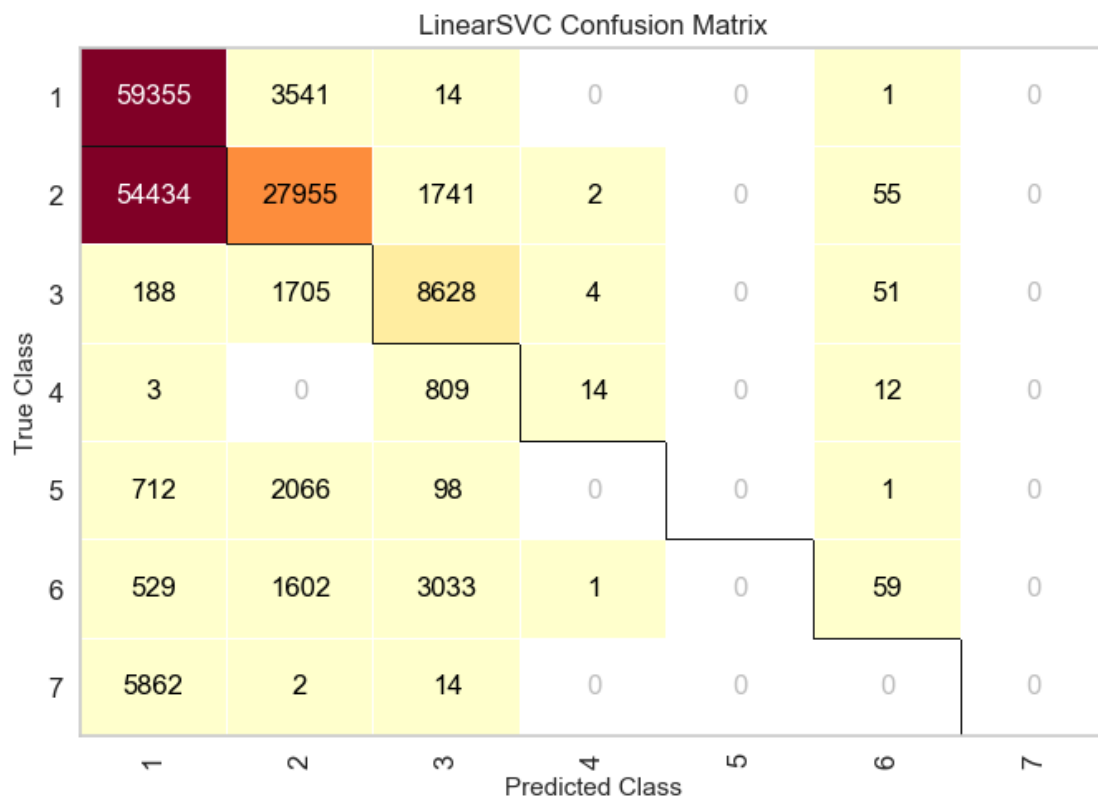
### Class Prediction Error for LinearSVC

[123]: `<matplotlib.axes._subplots.AxesSubplot at 0x2030b33fe10>`

[125]:
```
# The ConfusionMatrix visualizer taxes a model
cm = ConfusionMatrix(modelSVM1)

# Fit fits the passed model. This is unnecessary if you pass the visualizer a
 ↪pre-fitted model
cm.fit(X1_train, y1_train)

# To create the ConfusionMatrix, we need some test data. Score runs predict()
 ↪on the data
# and then creates the confusion_matrix from scikit-learn.
cm.score(X1_test, y1_test)
```
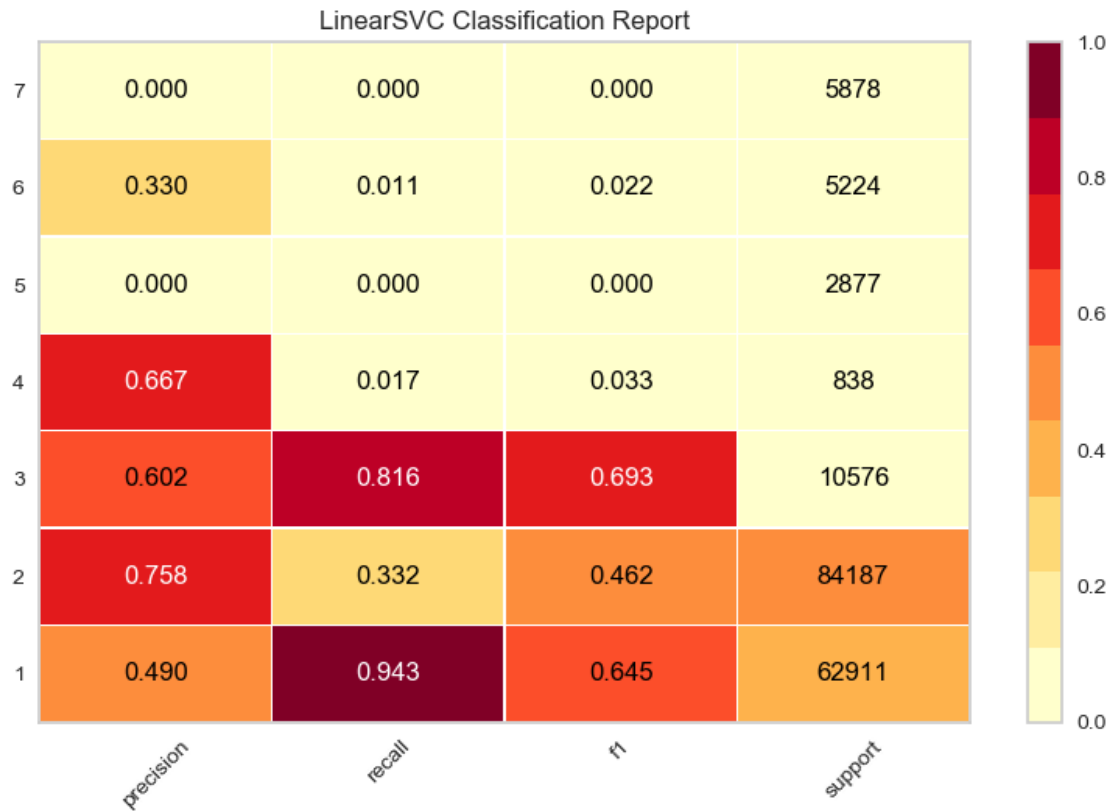
```
cm.show()
```



LinearSVC Confusion Matrix

| True Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | 59355 | 3541 | 14 | 0 | 0 | 1 | 0 |
| 2 | 54434 | 27955 | 1741 | 2 | 0 | 55 | 0 |
| 3 | 188 | 1705 | 8628 | 4 | 0 | 51 | 0 |
| 4 | 3 | 0 | 809 | 14 | 0 | 12 | 0 |
| 5 | 712 | 2066 | 98 | 0 | 0 | 1 | 0 |
| 6 | 529 | 1602 | 3033 | 1 | 0 | 59 | 0 |
| 7 | 5862 | 2 | 14 | 0 | 0 | 0 | 0 |

[125]: <matplotlib.axes._subplots.AxesSubplot at 0x2030c4c5198>

```
[126]: visualizer = ClassificationReport(modelSVM1, support=True)

       visualizer.fit(X1_train, y1_train)        # Fit the visualizer and the model
       visualizer.score(X1_test, y1_test)        # Evaluate the model on the test data
       visualizer.show()
```

LinearSVC Classification Report

|   | precision | recall | f1 | support |
|---|-----------|--------|-----|---------|
| 7 | 0.000 | 0.000 | 0.000 | 5878 |
| 6 | 0.330 | 0.011 | 0.022 | 5224 |
| 5 | 0.000 | 0.000 | 0.000 | 2877 |
| 4 | 0.667 | 0.017 | 0.033 | 838 |
| 3 | 0.602 | 0.816 | 0.693 | 10576 |
| 2 | 0.758 | 0.332 | 0.462 | 84187 |
| 1 | 0.490 | 0.943 | 0.645 | 62911 |

[126]: <matplotlib.axes._subplots.AxesSubplot at 0x2030c5a59e8>

- We can see from the plots that more than half of the classes are predicted correctly. But the model predicted 5 classes although there are 7. Class 5 and 7 could not detected.

### 1.0.9 DecisionTreeClassifier

```
[47]: modelTree = DecisionTreeClassifier()
```

```
[48]: modelTree.fit(X_train , y_train)
```

```
[48]: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                             max_depth=None, max_features=None, max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=1, min_samples_split=2,
                             min_weight_fraction_leaf=0.0, presort='deprecated',
                             random_state=None, splitter='best')
```

```
[49]: pred = modelTree.predict(X_test)
```

11

```
[50]: tree_accuracy = accuracy_score(pred, y_test)
```

```
[51]: # Quite well!

      tree_accuracy
```
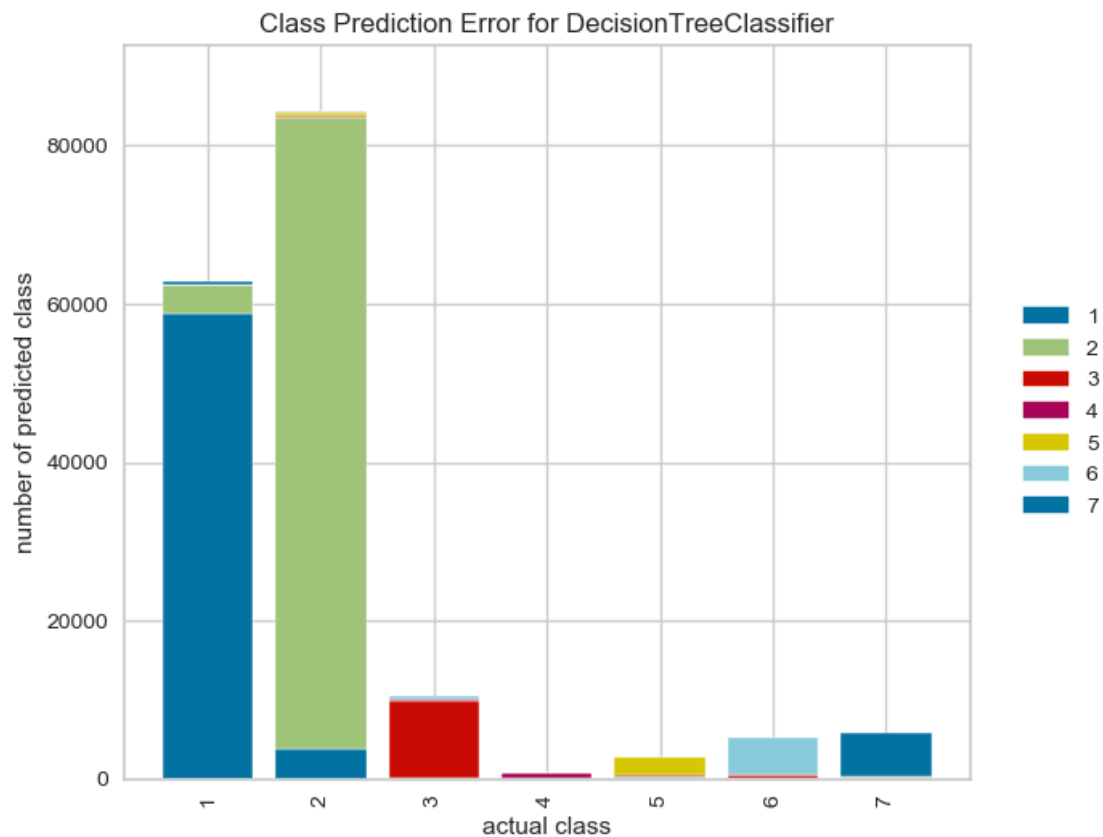
```
[51]: 0.9354517047266234
```

**Now, I would like to show three plots for the results.** - Class Prediction Error Bar Plot - Confusion Matrix - Classification Report

```
[54]: visualizer = ClassPredictionError(modelTree)

      # Fit the training data to the visualizer
      visualizer.fit(X_train, y_train)

      # Evaluate the model on the test data
      visualizer.score(X_test, y_test)

      # Draw visualization
      visualizer.show()
```

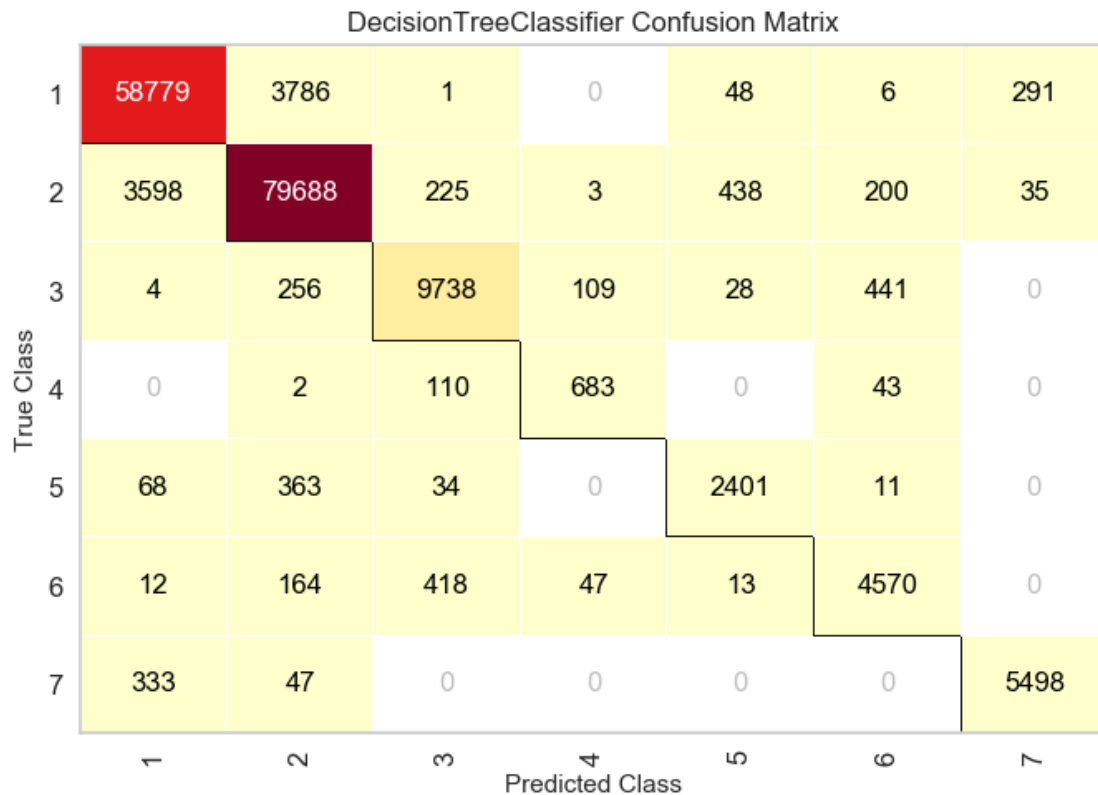[54]: `<matplotlib.axes._subplots.AxesSubplot at 0x203097713c8>`

[53]:
```python
# The ConfusionMatrix visualizer taxes a model
cm = ConfusionMatrix(modelTree)

# Fit fits the passed model. This is unnecessary if you pass the visualizer a
 ↪pre-fitted model
cm.fit(X_train, y_train)

# To create the ConfusionMatrix, we need some test data. Score runs predict()
 ↪on the data
# and then creates the confusion_matrix from scikit-learn.
cm.score(X_test, y_test)

cm.show()
```
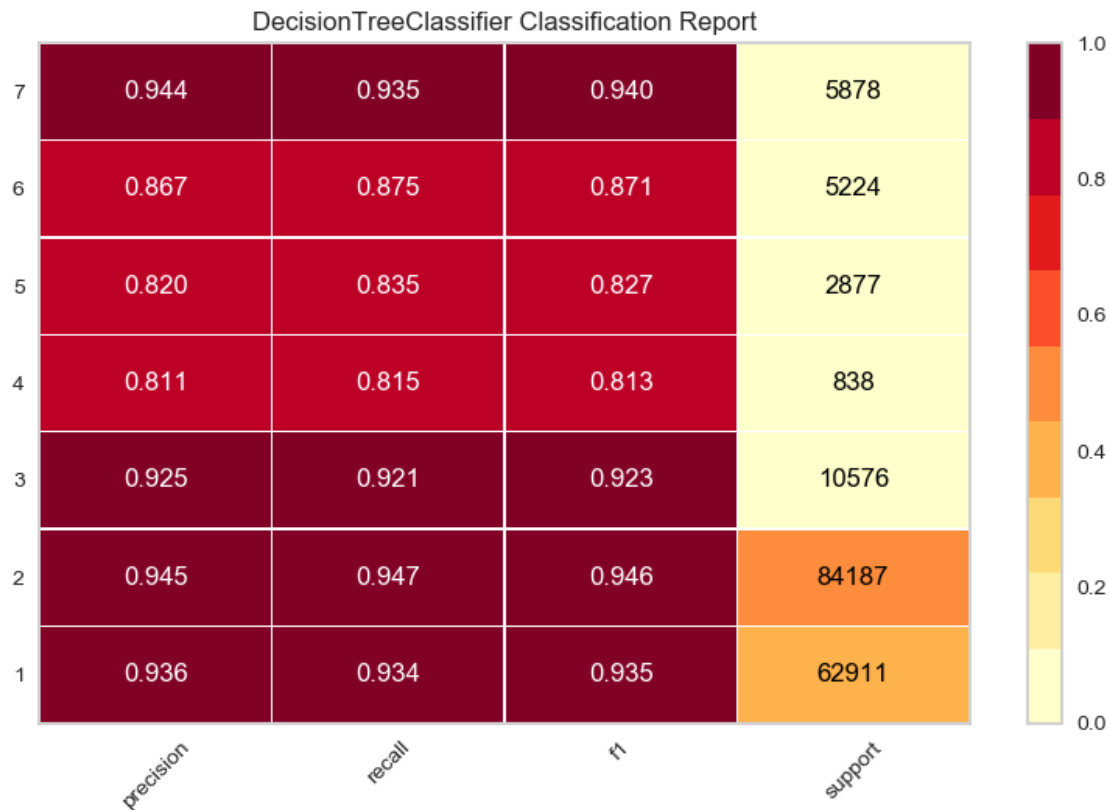
DecisionTreeClassifier Confusion Matrix

| True Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | 58779 | 3786 | 1 | 0 | 48 | 6 | 291 |
| 2 | 3598 | 79688 | 225 | 3 | 438 | 200 | 35 |
| 3 | 4 | 256 | 9738 | 109 | 28 | 441 | 0 |
| 4 | 0 | 2 | 110 | 683 | 0 | 43 | 0 |
| 5 | 68 | 363 | 34 | 0 | 2401 | 11 | 0 |
| 6 | 12 | 164 | 418 | 47 | 13 | 4570 | 0 |
| 7 | 333 | 47 | 0 | 0 | 0 | 0 | 5498 |

[53]: `<matplotlib.axes._subplots.AxesSubplot at 0x20309c15a20>`

[52]:
```python
visualizer = ClassificationReport(modelTree, support=True)
```

```
visualizer.fit(X_train, y_train)          # Fit the visualizer and the model
visualizer.score(X_test, y_test)          # Evaluate the model on the test data
visualizer.show()
```

DecisionTreeClassifier Classification Report

| | precision | recall | f1 | support |
|---|---|---|---|---|
| 7 | 0.944 | 0.935 | 0.940 | 5878 |
| 6 | 0.867 | 0.875 | 0.871 | 5224 |
| 5 | 0.820 | 0.835 | 0.827 | 2877 |
| 4 | 0.811 | 0.815 | 0.813 | 838 |
| 3 | 0.925 | 0.921 | 0.923 | 10576 |
| 2 | 0.945 | 0.947 | 0.946 | 84187 |
| 1 | 0.936 | 0.934 | 0.935 | 62911 |

[52]: <matplotlib.axes._subplots.AxesSubplot at 0x20309c15940>

### 1.0.10  KNeighborsClassifer

- Deciding the number of neighbors

```
[55]: neighbors = np.arange(1, 7)
      train_accuracy =np.empty(len(neighbors))
      test_accuracy = np.empty(len(neighbors))

      for i,k in enumerate(neighbors):
          #Setup a knn classifier with k neighbors
          knn = KNeighborsClassifier(n_neighbors = k)

          #Fit the model
          knn.fit(X_train, y_train)
```

14

```
    #Compute accuracy on the training set
    train_accuracy[i] = knn.score(X_train, y_train)

    #Compute accuracy on the test set
    test_accuracy[i] = knn.score(X_test, y_test)
```
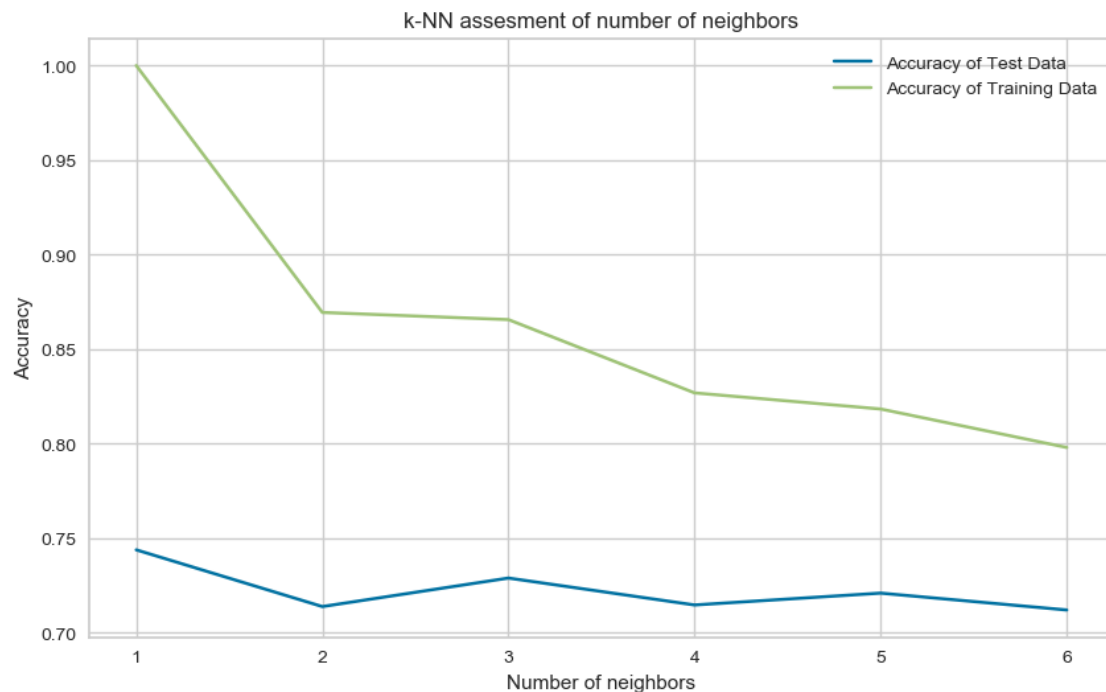
```
[56]: plt.figure(figsize=(10,6))
      plt.title('k-NN assesment of number of neighbors')
      plt.plot(neighbors, test_accuracy, label='Accuracy of Test Data')
      plt.plot(neighbors, train_accuracy, label='Accuracy of Training Data')
      plt.legend()
      plt.xlabel('Number of neighbors')
      plt.ylabel('Accuracy')
      plt.show()
```



- **The graph lines stabilize around 5. So, Let's try 5 neighbors**. Let's prepare the train and test data for KNN Model.

```
[101]: knn5 = KNeighborsClassifier(n_neighbors = 5)
```

```
[102]: knn5.fit(X1_train,y1_train)
```

```
[102]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                            metric_params=None, n_jobs=None, n_neighbors=5, p=2,
```

```
                      weights='uniform')
```

[103]: `knn_accuracy = knn5.score(X1_test,y1_test)`

[104]:
```python
# Excellent!

knn_accuracy
```
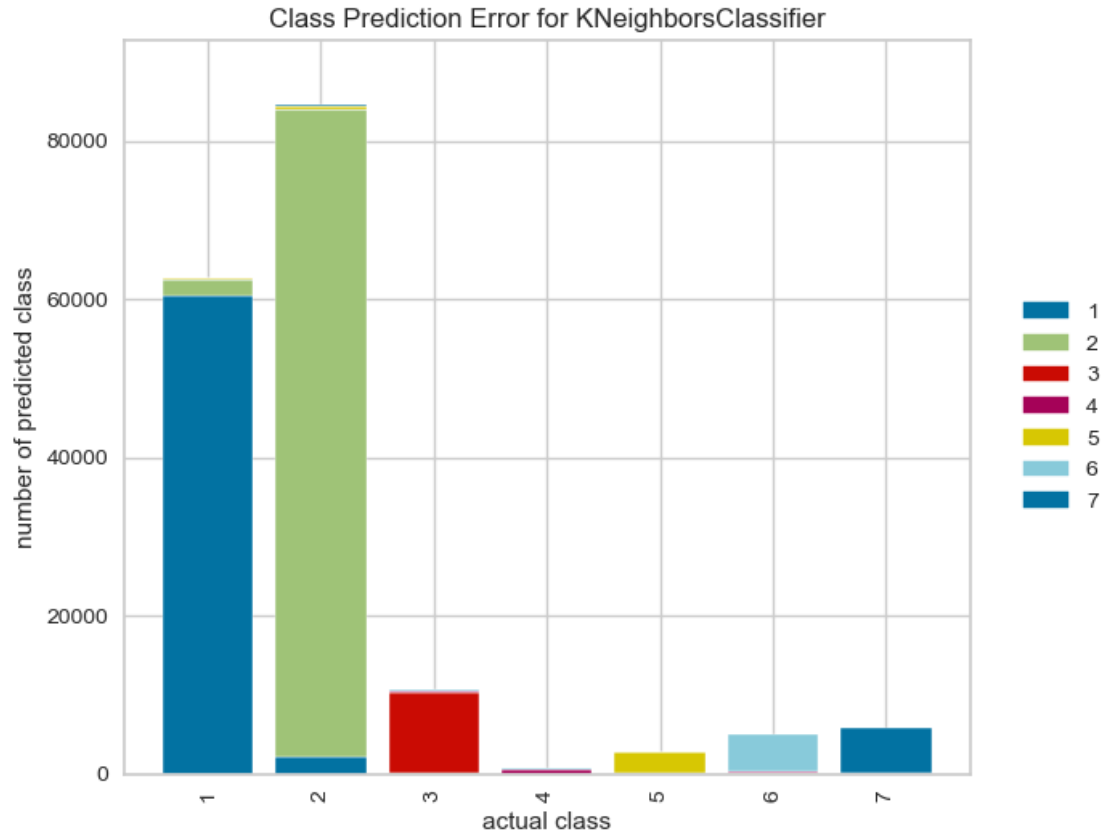
[104]: `0.962745882393864`

**Now, I would like to show three plots for the results.** - Class Prediction Error Bar Plot - Confusion Matrix - Classification Report

[108]:
```python
visualizer = ClassPredictionError(knn5)

# Fit the training data to the visualizer
visualizer.fit(X1_train, y1_train)

# Evaluate the model on the test data
visualizer.score(X1_test, y1_test)

# Draw visualization
visualizer.show()
```
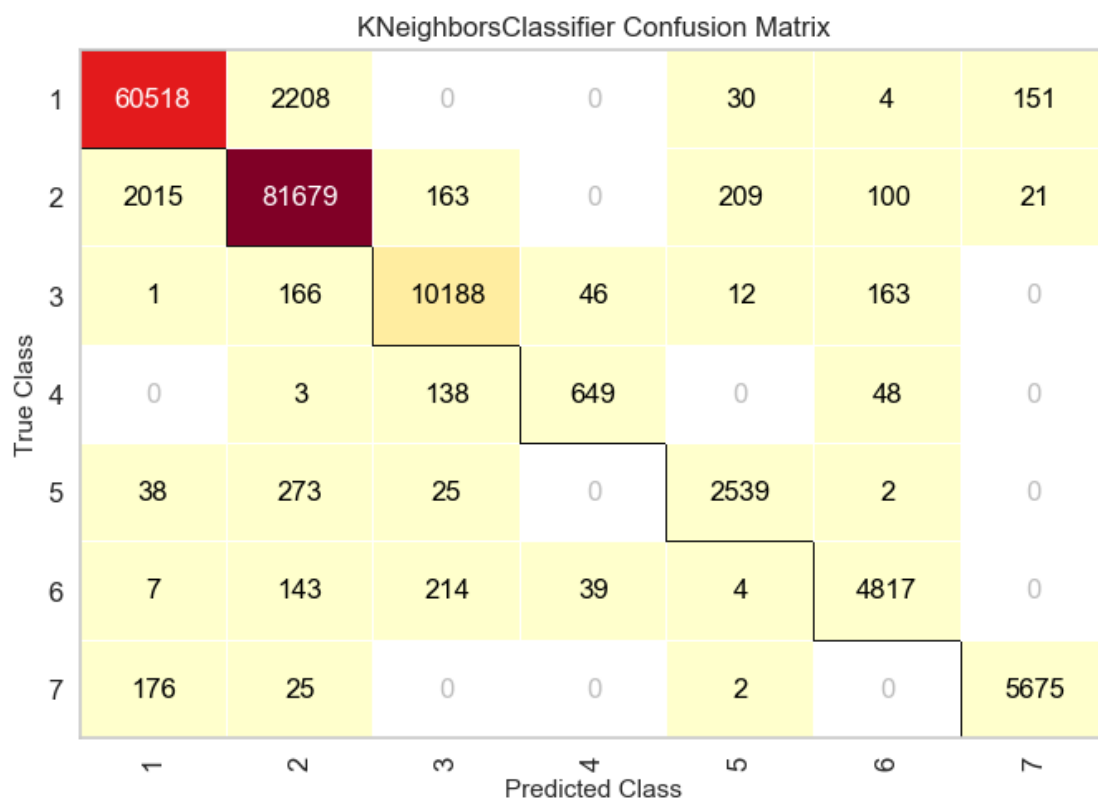
```
[108]: <matplotlib.axes._subplots.AxesSubplot at 0x20306be8358>
```

```
[109]: # The ConfusionMatrix visualizer taxes a model
       cm = ConfusionMatrix(knn5)

       # Fit fits the passed model. This is unnecessary if you pass the visualizer a
        ↪pre-fitted model
       cm.fit(X1_train, y1_train)

       # To create the ConfusionMatrix, we need some test data. Score runs predict()
        ↪on the data
       # and then creates the confusion_matrix from scikit-learn.
       cm.score(X1_test, y1_test)

       cm.show()
```
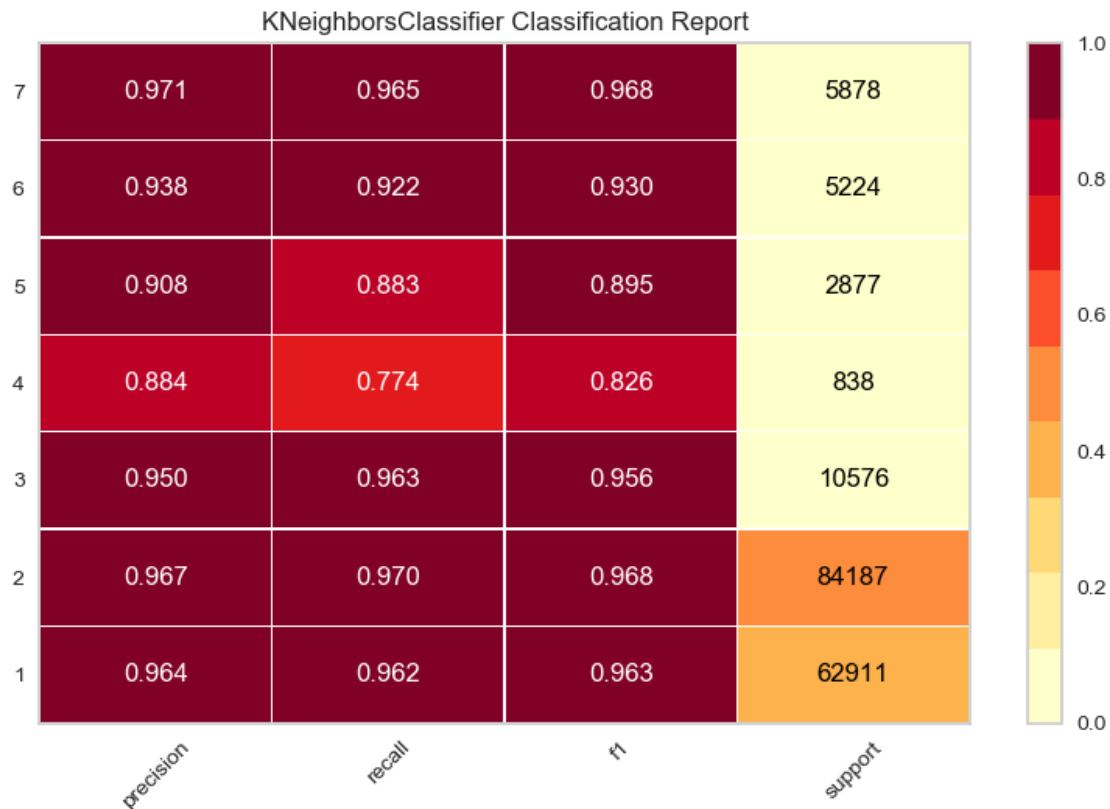
KNeighborsClassifier Confusion Matrix

| True Class | Predicted 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | 60518 | 2208 | 0 | 0 | 30 | 4 | 151 |
| 2 | 2015 | 81679 | 163 | 0 | 209 | 100 | 21 |
| 3 | 1 | 166 | 10188 | 46 | 12 | 163 | 0 |
| 4 | 0 | 3 | 138 | 649 | 0 | 48 | 0 |
| 5 | 38 | 273 | 25 | 0 | 2539 | 2 | 0 |
| 6 | 7 | 143 | 214 | 39 | 4 | 4817 | 0 |
| 7 | 176 | 25 | 0 | 0 | 2 | 0 | 5675 |

Predicted Class

```
[109]: <matplotlib.axes._subplots.AxesSubplot at 0x203000f5358>
```

```python
[110]: visualizer = ClassificationReport(knn5, support=True)

       visualizer.fit(X1_train, y1_train)          # Fit the visualizer and the model
       visualizer.score(X1_test, y1_test)          # Evaluate the model on the test data
       visualizer.show()
```

KNeighborsClassifier Classification Report

|   | precision | recall | f1 | support |
|---|-----------|--------|------|---------|
| 7 | 0.971 | 0.965 | 0.968 | 5878 |
| 6 | 0.938 | 0.922 | 0.930 | 5224 |
| 5 | 0.908 | 0.883 | 0.895 | 2877 |
| 4 | 0.884 | 0.774 | 0.826 | 838 |
| 3 | 0.950 | 0.963 | 0.956 | 10576 |
| 2 | 0.967 | 0.970 | 0.968 | 84187 |
| 1 | 0.964 | 0.962 | 0.963 | 62911 |

```
[110]: <matplotlib.axes._subplots.AxesSubplot at 0x2030708e400>
```

### 1.0.11 LGBMClassifier

```python
[77]: from lightgbm import LGBMClassifier
```

```python
[80]: lgbm_classifier = LGBMClassifier()
```

```python
[81]: lgbm_classifier.fit(X_train, y_train)
      y_predicted = lgbm_classifier.predict(X_test)
```

```python
[127]: lgbm_accuracy = accuracy_score(y_test, y_predicted)
```

```
[128]:  # Very good..

         lgbm_accuracy
```

[128]: 0.8332782579960694

**Now, I would like to show three plots for the results.** - Class Prediction Error Bar Plot - Confusion Matrix - Classification Report

```
[130]:  visualizer = ClassPredictionError(lgbm_classifier)

         # Fit the training data to the visualizer
         visualizer.fit(X_train, y_train)

         # Evaluate the model on the test data
         visualizer.score(X_test, y_test)

         # Draw visualization
         visualizer.show()
```
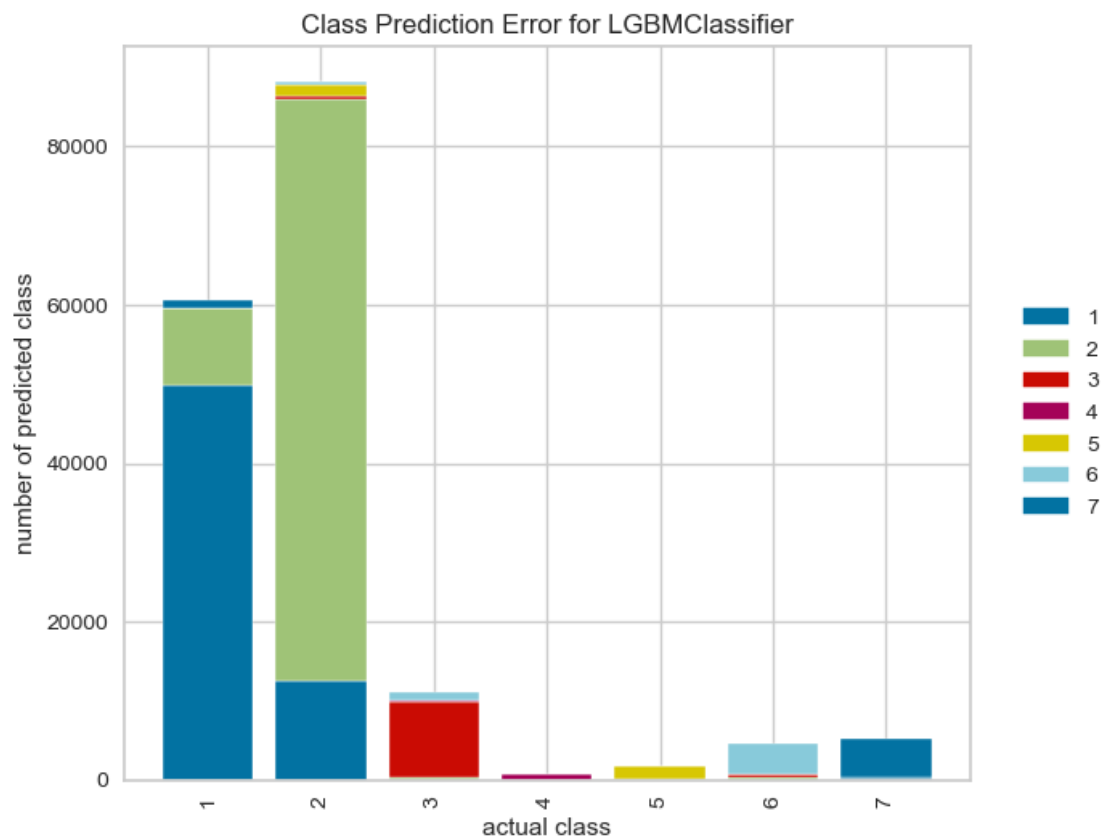


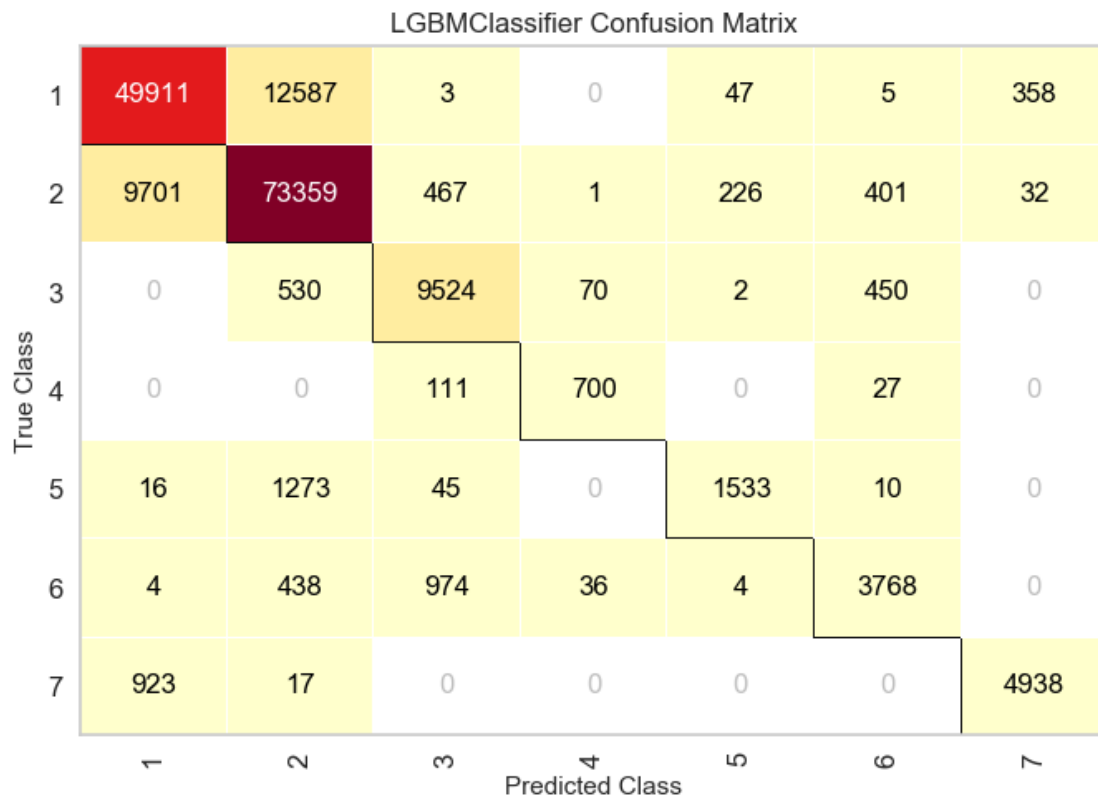[130]: <matplotlib.axes._subplots.AxesSubplot at 0x2030c5a5b38>

```
[131]: # The ConfusionMatrix visualizer taxes a model
       cm = ConfusionMatrix(lgbm_classifier)

       # Fit fits the passed model. This is unnecessary if you pass the visualizer a
       ↪pre-fitted model
       cm.fit(X_train, y_train)

       # To create the ConfusionMatrix, we need some test data. Score runs predict()
       ↪on the data
       # and then creates the confusion_matrix from scikit-learn.
       cm.score(X_test, y_test)

       cm.show()
```



LGBMClassifier Confusion Matrix

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | 49911 | 12587 | 3 | 0 | 47 | 5 | 358 |
| 2 | 9701 | 73359 | 467 | 1 | 226 | 401 | 32 |
| 3 | 0 | 530 | 9524 | 70 | 2 | 450 | 0 |
| 4 | 0 | 0 | 111 | 700 | 0 | 27 | 0 |
| 5 | 16 | 1273 | 45 | 0 | 1533 | 10 | 0 |
| 6 | 4 | 438 | 974 | 36 | 4 | 3768 | 0 |
| 7 | 923 | 17 | 0 | 0 | 0 | 0 | 4938 |

True Class / Predicted Class
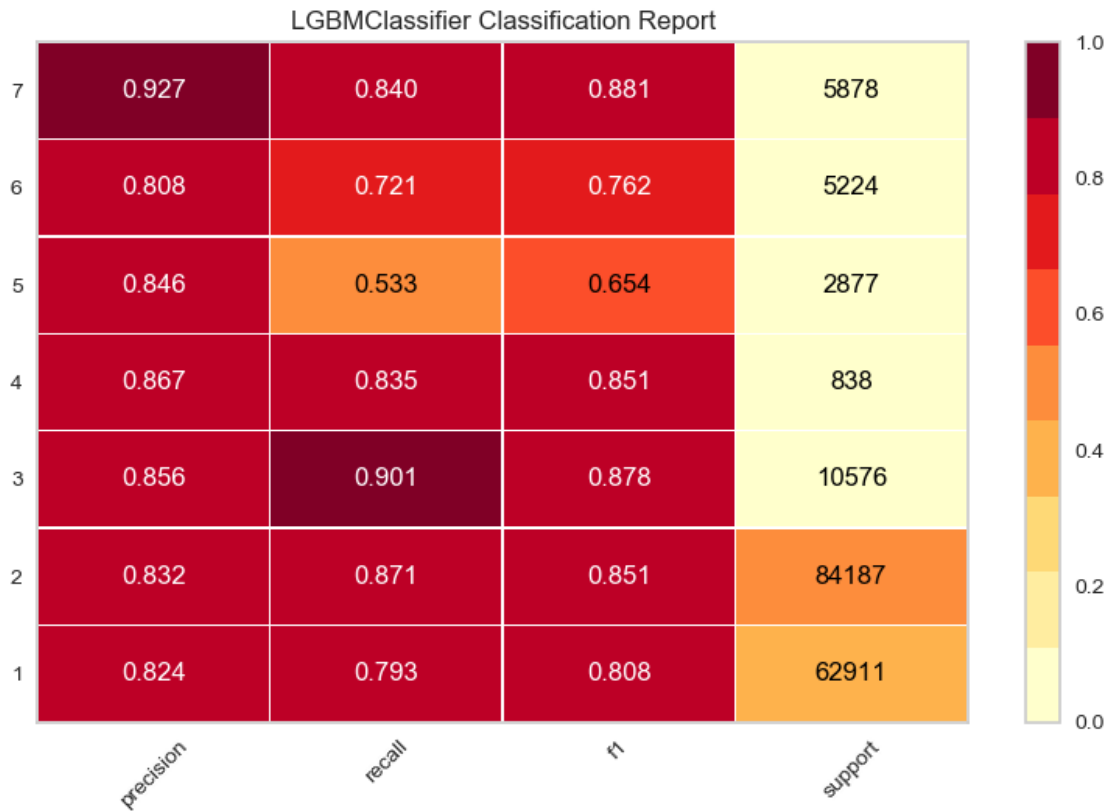
```
[131]: <matplotlib.axes._subplots.AxesSubplot at 0x203096c2278>
```

```
[132]: visualizer = ClassificationReport(lgbm_classifier, support=True)

       visualizer.fit(X_train, y_train)        # Fit the visualizer and the model
       visualizer.score(X_test, y_test)        # Evaluate the model on the test data
       visualizer.show()
```
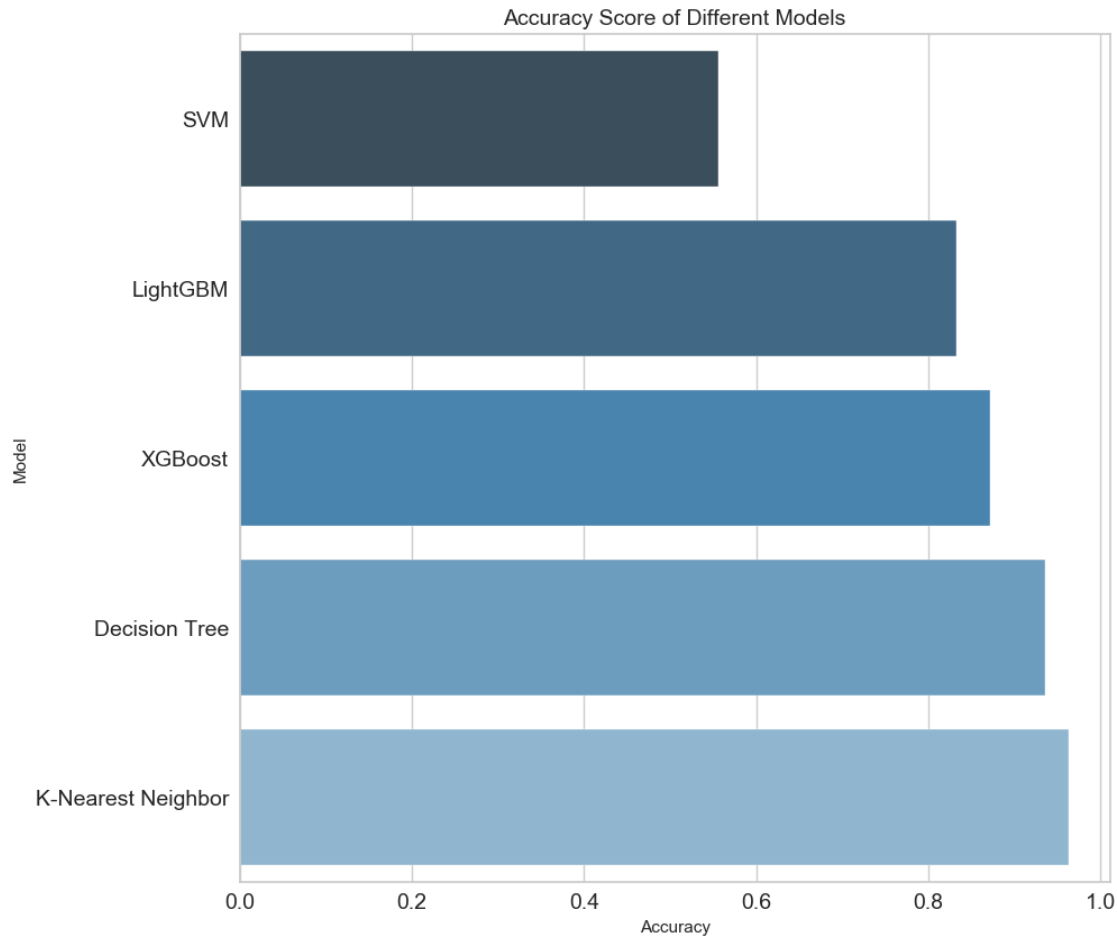
## LGBMClassifier Classification Report

| | precision | recall | f1 | support |
|---|---|---|---|---|
| 7 | 0.927 | 0.840 | 0.881 | 5878 |
| 6 | 0.808 | 0.721 | 0.762 | 5224 |
| 5 | 0.846 | 0.533 | 0.654 | 2877 |
| 4 | 0.867 | 0.835 | 0.851 | 838 |
| 3 | 0.856 | 0.901 | 0.878 | 10576 |
| 2 | 0.832 | 0.871 | 0.851 | 84187 |
| 1 | 0.824 | 0.793 | 0.808 | 62911 |

[132]: <matplotlib.axes._subplots.AxesSubplot at 0x2032c4635f8>

[134]:
```python
compare = pd.DataFrame({"Model": ["K-Nearest Neighbor", "LightGBM", "XGBoost",
 →"Decision Tree", "SVM"],
                        "Accuracy": [knn_accuracy, lgbm_accuracy, xgb_accuracy,
 →tree_accuracy, SVM_accuracy1]})

compare = compare.sort_values(by="Accuracy", ascending=True)

plt.figure(figsize=(10,10))
ax = sns.barplot(x="Accuracy", y="Model", data=compare, palette="Blues_d")

plt.yticks(size = 14)
plt.xticks(size = 14)
plt.title("Accuracy Score of Different Models", size=14)
```

[134]: Text(0.5, 1.0, 'Accuracy Score of Different Models')

Accuracy Score of Different Models



### 1.0.12 Overall Results :

- The overall evaluation in terms of `accuracy scores` is displayed above.
- The model with the least performance is SVM.
  - SVM has been widely used in finance. For example, predicting stock price via SVM has been a acknowledged application in the industry.
  - In classification of text and handwritten objects, SVM performs well.
  - It may not be very successful in datasets with more than 100,000 data (We have half a million).
  - It also doesn't perform well against unstable data (`Cover_Type`).
- `XGBoost` has performed well enough. As in `linear regression` models, it has proven its success in `clasification` once again.
  - There is a little bit low `recall` score of the class 5.
  - I conclueded that it is caused by the values in the data set. %59.8 of the class 5 is predicted true.
  - Although the `recall` score level of class 5 is a little bit low, `precision` and `f1` score is quite well.

- All the other Models (`K-Nearest Neighbor`, `LightGBM`, `Decision Tree`) are quite well. They all have performed well and given high accuracy scores.
- Multi Class Classification is `KNN`'s job...

Best Regards..