

K Nearest Neighbors Project

Welcome to the KNN Project! This will be a simple project very similar to the lecture, except you'll be given another data set. Go ahead and just follow the directions below.

Import Libraries

Import pandas,seaborn, and the usual libraries.

```
In [4]:  
import numpy as np  
import pandas as pd  
import seaborn as sns  
import matplotlib.pyplot as plt  
%matplotlib inline  
from sklearn.metrics import classification_report as cr,confusion_matrix as cm,accuracy_score as ac  
from sklearn.model_selection import train_test_split,cross_val_score  
from sklearn.preprocessing import StandardScaler  
  
import warnings  
warnings.filterwarnings('ignore')
```

Get the Data

Read the 'KNN_Project_Data csv file into a dataframe

```
In [2]:  
df=pd.read_csv('KNN_Project_Data.csv')  
df.head()  
  
Out[2]:  
      XVPY  GWHY  TRAT  TLLZ  IGGA  HYKR  EDFY  GUUB  MGJM  JH2C  TARGET CLASS  
0  1636.670614  817.988525  2565.995189  358.347163  550.417491  1618.870897  2147.641254  330.27893  1494.878631  845.136088  0  
1  1013.402760  577.587332  2644.141273  280.428203  1161.873391  2084.107872  853.404981  447.157619  1193.032521  861.081809  1  
2  1300.035501  820.518697  2025.854469  525.562292  922.206261  252.355407  818.676686  845.491492  1968.367513  1647.186291  1  
3  1059.347542  1066.866418  612.000041  480.827789  419.467495  685.666983  852.867810  341.664784  1154.391368  1450.935357  0  
4  1018.340526  1313.679056  950.622661  724.742174  843.065903  1370.554164  905.469453  658.18202  539.459350  1899.850792  0
```

Check the head of the dataframe.

```
In [7]:  
df.info()  
  
Out[7]:  
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 1000 entries, 0 to 999  
Data columns (total 11 columns):  
 #   Column          Non-Null Count  Dtype     
---  --  
 0   XVPY            1000 non-null   float64  
 1   GWHY            1000 non-null   float64  
 2   TRAT            1000 non-null   float64  
 3   TLLZ            1000 non-null   float64  
 4   IGGA            1000 non-null   float64  
 5   HYKR            1000 non-null   float64  
 6   EDFY            1000 non-null   float64  
 7   GUUB            1000 non-null   float64  
 8   MGJM            1000 non-null   float64  
 9   JH2C            1000 non-null   float64  
 10  TARGET CLASS  1000 non-null   int64  
dtypes: float64(10), int64(1)  
memory usage: 86.1 KB
```

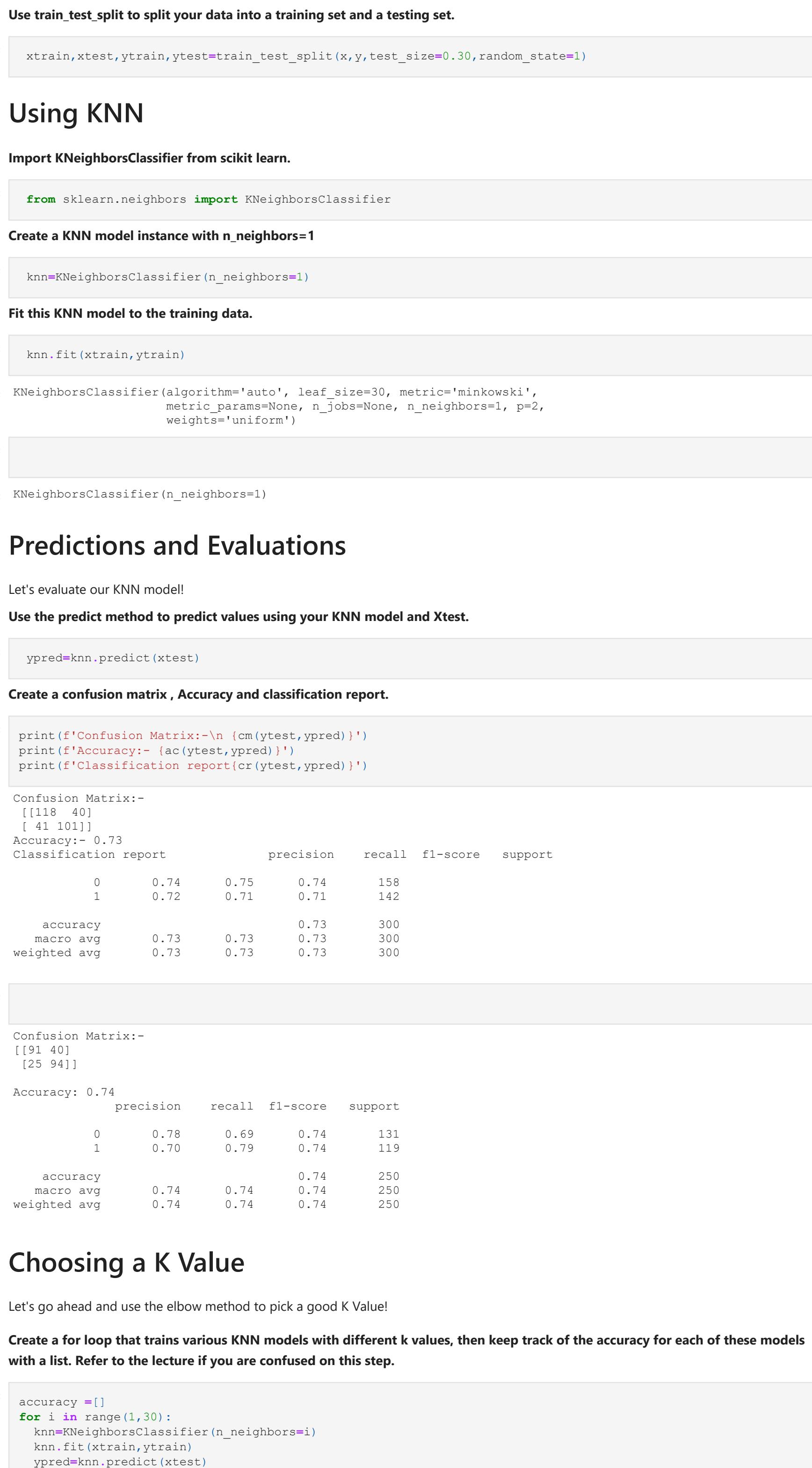
```
In [ ]:  
  
Out[ ]:  
      XVPY  GWHY  TRAT  TLLZ  IGGA  HYKR  EDFY  GUUB  MGJM  JH2C  TARGET CLASS  
0  1636.670614  817.988525  2565.995189  358.347163  550.417491  1618.870897  2147.641254  330.27893  1494.878631  845.136088  0  
1  1013.402760  577.587332  2644.141273  280.428203  1161.873391  2084.107872  853.404981  447.157619  1193.032521  861.081809  1  
2  1300.035501  820.518697  2025.854469  525.562292  922.206261  252.355407  818.676686  845.491492  1968.367513  1647.186291  1  
3  1059.347542  1066.866418  612.000041  480.827789  419.467495  685.666983  852.867810  341.664784  1154.391368  1450.935357  0  
4  1018.340526  1313.679056  950.622661  724.742174  843.065903  1370.554164  905.469453  658.18202  539.459350  1899.850792  0
```

EDA

Since this data is artificial, we'll just do a large pairplot with seaborn.

Use seaborn on the dataframe to create a pairplot with the hue indicated by the TARGET CLASS column.

```
In [3]:  
sns.pairplot(df,hue='TARGET CLASS')  
  
Out[3]: <seaborn.axisgrid.PairGrid at 0x7f271f60ca50>
```



Divide X & Y

```
In [5]:  
x=df.iloc[:, :-1]  
y=df.iloc[:, -1]
```

Standardize the Variables

Time to standardize the variables.

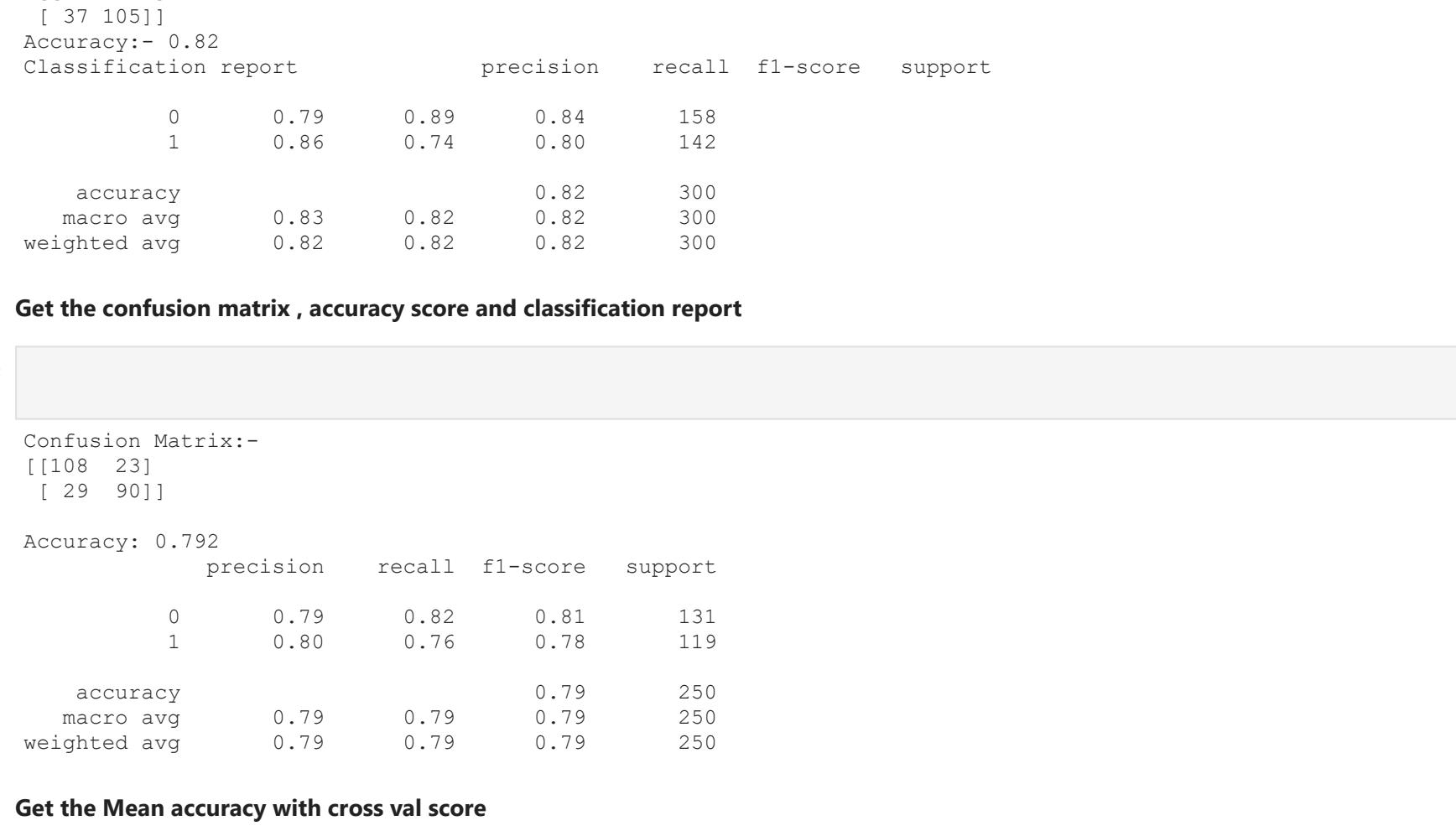
Import StandardScaler from Scikit learn.

```
In [6]:  
sc=StandardScaler()  
sc.fit_transform(x)
```

```
Out[6]: array([[ 1.56852168, -0.44343461,  1.61980773, ..., -0.93279392,  
    1.00831307, -1.03692723],  
   [-0.11237594, -1.05657361,  1.7419175, ..., -0.46186435,  
    0.25832069, -0.104154625],  
   [ 0.66064691, -0.43698145,  0.77579285, ...,  1.14929806,  
    2.1847836, -0.34281129],  
   ...,  
   [-0.35889496, -0.97901454,  0.83771499, ..., -1.51472604,  
    -0.27512225,  0.86428656],  
   [ 0.27507999, -0.99239881,  0.0303711, ..., -0.03623294,  
    0.43668516, -0.21245586],  
   [ 0.62589594,  0.79510909,  1.12180047, ..., -1.25156478,  
    -0.60352946, -0.87985868]])
```

Convert the scaled features to a dataframe and check the head of this dataframe to make sure the scaling worked.

```
In [ ]:  
  
Out[ ]: <seaborn.axisgrid.PairGrid at 0xle6d5923d0>
```



Train Test Split

Use train_test_split to split your data into a training set and a testing set.

```
In [53]:  
xtrain,xtest,ytrain,ytest=train_test_split(x,y,test_size=0.30,random_state=1)
```

Using KNN

Import KNeighborsClassifier from scikit learn.

```
In [54]:  
from sklearn.neighbors import KNeighborsClassifier
```

Create a KNN model instance with n_neighbors=1

```
In [55]:  
knn=KNeighborsClassifier(n_neighbors=1)
```

Fit this KNN model to the training data.

```
In [56]:  
knn.fit(xtrain,ytrain)
```

```
Out[56]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',  
    metric_params=None, n_jobs=None, n_neighbors=1, p=2,  
    weights='uniform')
```

```
In [ ]:  
  
Out[ ]:
```

Now create the following plot using the information from your for loop.

```
In [60]:  
plt.figure(figsize=(10,6))  
plt.plot(range(1,30),accuracy,color='black', linestyle='dashed', markerfacecolor='green', markeredgecolor='black', markersize=10)  
  
plt.title('Accuracy vs. K Value')  
plt.xlabel('K')  
plt.ylabel('Accuracy Rate')  
plt.grid(True)  
plt.show()
```



```
In [ ]:  
  
Out[ ]: plt.figure(figsize=(10,6))  
plt.plot(range(1,30),accuracy,color='blue', linestyle='dashed', markerfacecolor='red', markeredgecolor='red', markersize=10)  
  
plt.title('Accuracy vs. K Value')  
plt.xlabel('K')  
plt.ylabel('Accuracy Rate')  
plt.grid(True)  
plt.show()
```


Choosing a K Value

Let's go ahead and use the elbow method to pick a good K Value!

Create a for loop that trains various KNN models with different k values, then keep track of the accuracy for each of these models with a list. Refer to the lecture if you are confused on this step.

```
In [59]:  
accuracy =[]  
for i in range(1,30):  
    knn=KNeighborsClassifier(n_neighbors=i)  
    knn.fit(xtrain,ytrain)  
    ypred=knn.predict(xtest)  
    acc=accuracy(ytest,ypred)  
    accuracy.append(acc)
```

```
In [42]:  
len(accuracy)
```

```
Out[42]: 29
```

Now create the following plot using the information from your for loop.

```
In [60]:  
plt.figure(figsize=(10,6))  
plt.plot(range(1,30),accuracy,color='black', linestyle='dashed', markerfacecolor='black', markeredgecolor='black', markersize=10)  
  
plt.title('Accuracy vs. K Value')  
plt.xlabel('K')  
plt.ylabel('Accuracy Rate')  
plt.grid(True)  
plt.show()
```



```
In [ ]:  
  
Out[ ]: plt.figure(figsize=(10,6))  
plt.plot(range(1,30),accuracy,color='blue', linestyle='dashed', markerfacecolor='red', markeredgecolor='red', markersize=10)  
  
plt.title('Accuracy vs. K Value')  
plt.xlabel('K')  
plt.ylabel('Accuracy Rate')  
plt.grid(True)  
plt.show()
```


Get the mean accuracy with cross_val_score

```
In [80]:  
cv=cross_val_score(KNeighborsClassifier(n_neighbors=8),xtrain,ytrain,cv=6,scoring='accuracy')  
print(f'Accuracy: {cv.mean()*100}')  
print(f'Standard Deviation: {cv.std()}' )
```

```
Accuracy: 80.7165399351607  
Standard Deviation: 0.038975541917816096
```

```
In [81]:  
#given accuracy in file  
#Accuracy: 80.254592  
#Standard Deviation: 3.026876
```

```
In [ ]:  
  
Out[ ]:
```

Get the mean accuracy with cross_val_score

```
In [80]:  
cv=cross_val_score(KNeighborsClassifier(n_neighbors=8),xtrain,ytrain,cv=6,scoring='accuracy')  
print(f'Accuracy: {cv.mean()*100}')  
print(f'Standard Deviation: {cv.std()}' )
```

```
Accuracy: 80.7165399351607  
Standard Deviation: 0.038975541917816096
```

```
In [81]:  
#given accuracy in file  
#Accuracy: 80.254592  
#Standard Deviation: 3.026876
```

```
In [ ]:  
  
Out[ ]:
```

Get the mean accuracy with cross_val_score

```
In [80]:  
cv=cross_val_score(KNeighborsClassifier(n_neighbors=8),xtrain,ytrain,cv=6,scoring='accuracy')  
print(f'Accuracy: {cv.mean()*100}')  
print(f'Standard Deviation: {cv.std()}' )
```

```
Accuracy: 80.7165399351607  
Standard Deviation: 0.038975541917816096
```

```
In [81]:  
#given accuracy in file  
#Accuracy: 80.254592  
#Standard Deviation: 3.026876
```

```
In [ ]:  
  
Out[ ]:
```

