



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 [51]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

Get the Data

Read the 'KNN_Project_Data csv file into a dataframe

```
In [52]: df=pd.read_csv("KNN_Project_Data")
```

Check the head of the dataframe.

```
In [53]: df.head()
```

```
Out[53]:
```

	XVPM	GWYH	TRAT	TLLZ	IGGA	HYKR	EDF
0	1636.670614	817.988525	2565.995189	358.347163	550.417491	1618.870897	2147.64125
1	1013.402760	577.587332	2644.141273	280.428203	1161.873391	2084.107872	853.40498
2	1300.035501	820.518697	2025.854469	525.562292	922.206261	2552.355407	818.67668
3	1059.347542	1066.866418	612.000041	480.827789	419.467495	685.666983	852.86781
4	1018.340526	1313.679056	950.622661	724.742174	843.065903	1370.554164	905.46945



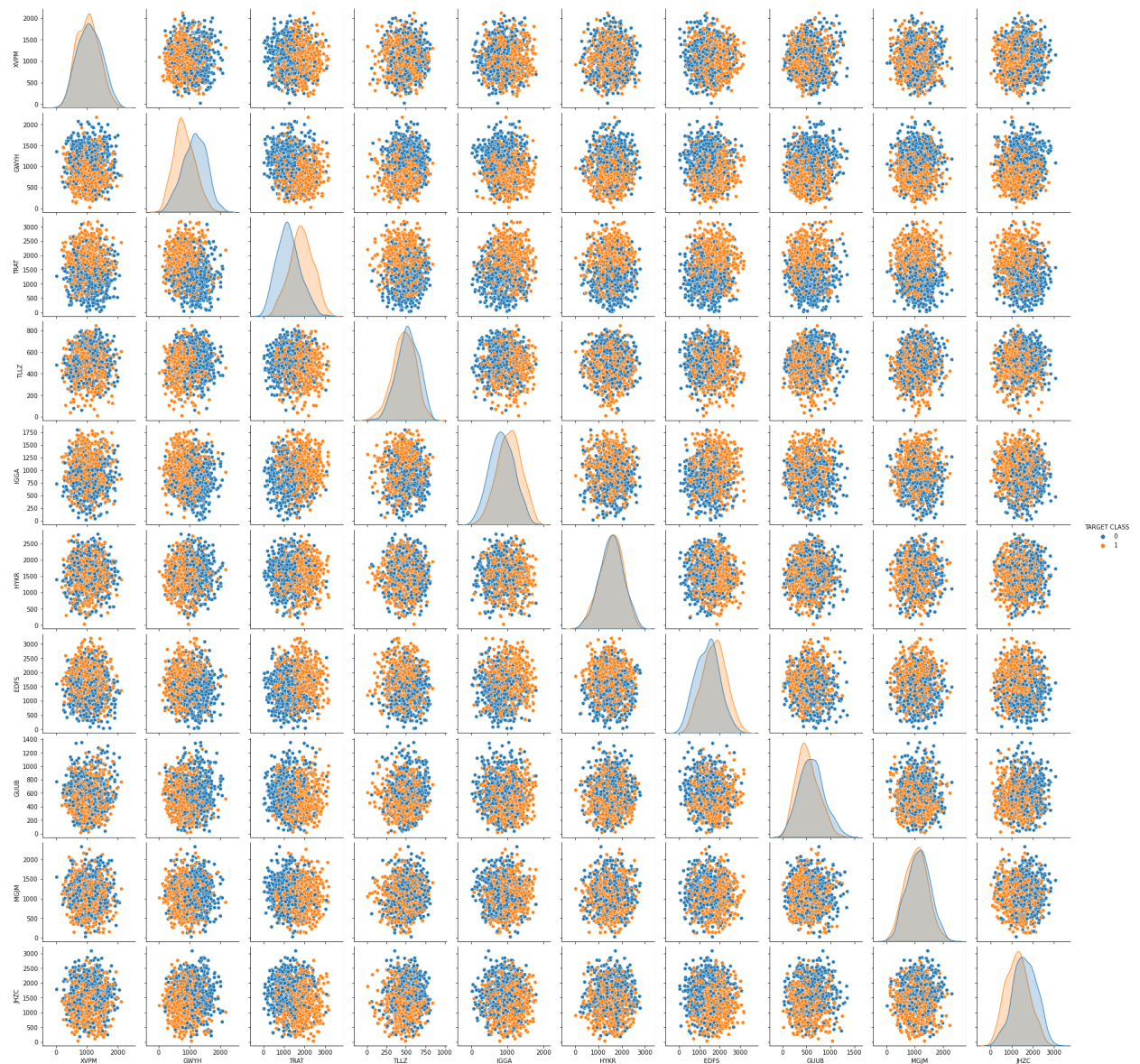
EDA

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 [54]: sns.pairplot(df, hue='TARGET CLASS')
```

```
Out[54]: <seaborn.axisgrid.PairGrid at 0x2c8880be9a0>
```



Standardize the Variables

Time to standardize the variables.

Import StandardScaler from Scikit learn.

```
In [55]: from sklearn.preprocessing import StandardScaler
```

Create a StandardScaler() object called scaler.

```
In [56]: scaler=StandardScaler()
```

Fit scaler to the features.

```
In [57]: scaler.fit(df.drop('TARGET CLASS',axis=1))
```

```
Out[57]: StandardScaler()
```

Use the .transform() method to transform the features to a scaled version.

```
In [58]: scaled_features = scaler.transform(df.drop('TARGET CLASS',axis=1))
```

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

```
In [59]: df_feat=pd.DataFrame(scaled_features,columns=df.columns[:-1])
df_feat.head()
```

```
Out[59]:
```

	XVPM	GWYH	TRAT	TLLZ	IGGA	HYKR	EDFS	GUUB	
0	1.568522	-0.443435	1.619808	-0.958255	-1.128481	0.138336	0.980493	-0.932794	1.0
1	-0.112376	-1.056574	1.741918	-1.504220	0.640009	1.081552	-1.182663	-0.461864	0.0
2	0.660647	-0.436981	0.775793	0.213394	-0.053171	2.030872	-1.240707	1.149298	2.0
3	0.011533	0.191324	-1.433473	-0.100053	-1.507223	-1.753632	-1.183561	-0.888557	0.0
4	-0.099059	0.820815	-0.904346	1.609015	-0.282065	-0.365099	-1.095644	0.391419	-1.0

Train Test Split

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

```
In [60]: from sklearn.model_selection import train_test_split
```

```
In [61]:
```

```
X_train,X_test,y_train,y_test=train_test_split(scaled_features,df['TARGET CLAS'])
```

Using KNN

Import KNeighborsClassifier from scikit learn.

```
In [62]: from sklearn.neighbors import KNeighborsClassifier
```

Create a KNN model instance with n_neighbors=1

```
In [63]: knn=KNeighborsClassifier(n_neighbors=1)
```

Fit this KNN model to the training data.

```
In [64]: knn.fit(X_train,y_train)
```

```
Out[64]: KNeighborsClassifier(n_neighbors=1)
```

Predictions and Evaluations

Let's evaluate our KNN model!

Use the predict method to predict values using your KNN model and X_test.

```
In [65]: predictions=knn.predict(X_test)
```

Create a confusion matrix and classification report.

```
In [66]: from sklearn.metrics import classification_report,confusion_matrix
```

```
In [67]: print(confusion_matrix(y_test,predictions))
```

```
[[113  47]
 [ 22 118]]
```

```
In [68]: print(classification_report(y_test,predictions))
```

	precision	recall	f1-score	support
0	0.84	0.71	0.77	160
1	0.72	0.84	0.77	140

accuracy			0.77	300
macro avg	0.78	0.77	0.77	300
weighted avg	0.78	0.77	0.77	300

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 error_rate for each of these models with a list. Refer to the lecture if you are confused on this step.

```
In [69]: error_rate = []

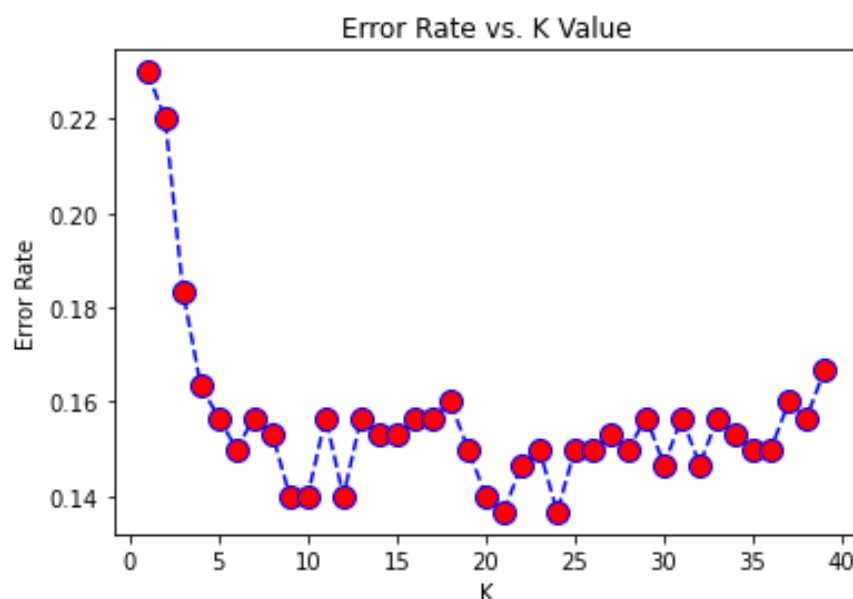
# Will take some time
for i in range(1,40):

    knn = KNeighborsClassifier(n_neighbors=i)
    knn.fit(X_train,y_train)
    pred_i = knn.predict(X_test)
    error_rate.append(np.mean(pred_i != y_test))
```

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

```
In [70]: #plt.figure(figsize=(10,6))
plt.plot(range(1,40),error_rate,color='blue', linestyle='dashed', marker='o',
         markerfacecolor='red', markersize=10)
plt.title('Error Rate vs. K Value')
plt.xlabel('K')
plt.ylabel('Error Rate')
```

```
Out[70]: Text(0, 0.5, 'Error Rate')
```



Retrain with new K Value

Retrain your model with the best K value (up to you to decide what you want) and re-do the classification report and the confusion matrix.

In [76]:

```
knn = KNeighborsClassifier(n_neighbors=35)

knn.fit(X_train,y_train)
pred = knn.predict(X_test)

print('WITH K=15')
print('\n')
print(confusion_matrix(y_test,pred))
print('\n')
print(classification_report(y_test,pred))
```

WITH K=15

```
[[121  39]
 [  6 134]]
```

	precision	recall	f1-score	support
0	0.95	0.76	0.84	160
1	0.77	0.96	0.86	140
accuracy			0.85	300
macro avg	0.86	0.86	0.85	300
weighted avg	0.87	0.85	0.85	300

Great Job!