

```
1 # Fetching the Dataset for training...
2 path = "https://raw.githubusercontent.com/shakil1819/CSE442-Machine-Learning-Sessional/main/Week%201%2B2/Knee-Torque-ZDataSet.csv"
3 df = pd.read_csv(path)
4 df.tail()
```

	Body Weight (kg)	Body Height (m)	Weight-Height Ratio (R)	Internal Moment (M)	Required Torque (N-m)	Torque Category
52	62	1.62	NaN	NaN	9.302	Low
53	65	1.50	NaN	NaN	8.624	Low
54	87	1.87	NaN	NaN	7.012	Very Low
55	77	1.88	NaN	NaN	2.514	Very Low
56	67	1.73	NaN	NaN	20.000	Very High

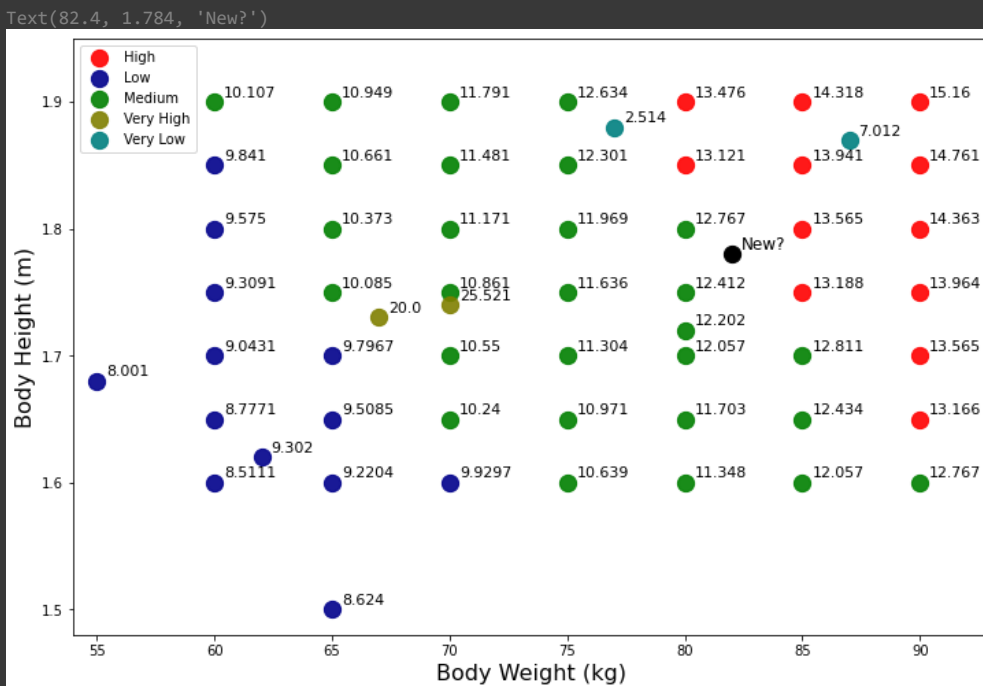
Scatter plot showing Body Height (m) on the Y-axis versus Body Weight (kg) on the X-axis. The data points are labeled with subject numbers, indicating a positive correlation between weight and height.

Subject	Body Weight (kg)	Body Height (m)
8.001	55	1.68
8.5111	60	1.60
8.624	65	1.50
8.7771	60	1.65
9.2204	65	1.60
9.302	62	1.62
9.3091	60	1.75
9.5085	65	1.65
9.575	60	1.80
9.9297	70	1.60
10.0431	60	1.70
10.085	65	1.75
10.107	60	1.90
10.24	70	1.65
10.251	70	1.75
10.301	75	1.70
10.373	65	1.80
10.639	75	1.60
10.661	65	1.85
10.861	70	1.75
11.055	70	1.70
11.171	70	1.80
11.202	80	1.72
11.304	75	1.70
11.348	80	1.60
11.481	70	1.85
11.636	75	1.75
11.791	70	1.90
11.969	75	1.80
12.057	80	1.70
12.057	85	1.60
12.202	80	1.72
12.267	80	1.80
12.301	75	1.85
12.348	80	1.60
12.412	80	1.75
12.434	85	1.65
12.514	78	1.88
12.634	75	1.90
12.767	90	1.60
12.811	85	1.70
12.851	87	1.87
13.057	90	1.70
13.121	80	1.85
13.166	90	1.65
13.363	90	1.80
13.476	80	1.90
13.565	85	1.80
13.624	87	1.82
13.841	60	1.85
13.941	85	1.85
13.964	90	1.75
14.012	87	1.87
14.318	85	1.90
14.363	90	1.80
14.761	90	1.85
15.16	90	1.90

```

1 T_CATEGORY = df['Torque Category'].values
2 T_CATEGORY_ = np.unique(T_CATEGORY)
3 COLORS = ["#FF0000", "#00008b", "#008000", "#808000", "#008080"]
4
5 fig, ax = plt.subplots(figsize=(12,8))
6 plt.axis([54, 93, 1.48, 1.95]) # plt.axis([x-min, x-max, y-min, y-max])
7 for category, color in zip(T_CATEGORY_, COLORS):
8     idxs = np.where(T_CATEGORY == category)
9     # No legend will be generated if we do not pass label=category...
10    ax.scatter(
11        w[idxs], h[idxs], label=category, s=150, color=color, alpha=0.9
12    )
13 ax.legend()
14
15 font = {'family': 'sans serif', 'color': 'black',
16         'weight': 'normal', 'style': 'normal', 'size': 11,
17         }
18 i=0
19 for v in df["Required Torque (N-m)"].to_numpy(): #df.iloc[:, 4].values:
20     plt.text(w[i]+.4, h[i]+.004, v, fontdict=font)
21     i=i+1
22
23 # https://matplotlib.org/stable/tutorials/text/text_props.html
24 font = {'family': 'sans serif', 'color': 'black',
25         'weight': 'normal', 'style': 'normal', 'size': 16,
26         }
27 ax.set_xlabel("Body Weight (kg)", fontdict=font, fontsize=16)
28 ax.set_ylabel("Body Height (m)", fontdict=font, fontsize=16)
29
30 plt.scatter(82, 1.78, s=150, c='#000000')
31 plt.text(82+.4, 1.78+.004, "New?", fontname='sans serif', fontsize=12)

```



```

1 # df # df will show all the records with values...
2 # df.isnull() # This will show all the records with True or False (bool) values...
3 df.isnull().sum() # [df.isna().sum()] Shows the features with NULL values count...

```

```

Body Weight (kg)      0
Body Height (m)       0
Weight-Height Ratio (R)  8
Internal Moment (M)   8
Required Torque (N-m)  0
Torque Category       0
dtype: int64

```

```

1 # =====
2 # Dropping the NULL records...
3 # =====
4 # df1 = df
5 # avg1 = df['Weight-Height Ratio (R)'].sum()/df['Weight-Height Ratio (R)'].count()
6 # avg2 = df['Internal Moment (M)'].sum()/df['Internal Moment (M)'].count()
7
8 ##### df.dropna() or df.fillna(100) or df.replace(np.nan, 0) or df.interpolate() ...
9 # df1['Weight-Height Ratio (R)'] = df['Weight-Height Ratio (R)'].replace(np.nan, avg1)
10 # df1['Internal Moment (M)'] = df['Internal Moment (M)'].replace(np.nan, avg2)
11
12 df1 = df.dropna()
13
14 df1.tail(10)

```

	Body Weight (kg)	Body Height (m)	Weight-Height Ratio (R)	Internal Moment (M)	Required Torque (N-m)	Torque Category
39	85	1.80	47.222	1.3841	13.565	High
40	85	1.85	45.946	1.4226	13.941	High
41	85	1.90	44.737	1.4610	14.318	High
42	90	1.60	56.250	1.3027	12.767	Medium
43	90	1.65	54.545	1.3434	13.166	High
44	90	1.70	52.941	1.3841	13.565	High
45	90	1.75	51.429	1.4249	13.964	High
46	90	1.80	50.000	1.4656	14.363	High
47	90	1.85	48.649	1.5063	14.761	High
48	90	1.90	47.368	1.5470	15.160	High

```

1 y = df1['Required Torque (N-m)']
2 X = df1.drop(['Weight-Height Ratio (R)', 'Internal Moment (M)',
3             'Required Torque (N-m)', 'Torque Category'], axis=1)
4 X.describe()

```

	Body Weight (kg)	Body Height (m)
count	49.00000	49.000000
mean	75.00000	1.750000
std	10.10363	0.101036
min	60.00000	1.600000
25%	65.00000	1.650000
50%	75.00000	1.750000
75%	85.00000	1.850000
max	90.00000	1.900000

```

1 # Splitting Train and Test data (75% & 25%, respectively) ...
2 from sklearn.model_selection import train_test_split
3 # SEED = 42
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25)
5
6 len(X)      # 57
7 len(X_train) # 42
8 len(X_test)  # 15
9
10 X_test

```

	Body Weight (kg)	Body Height (m)
30	80	1.70
5	60	1.85
43	90	1.65
21	75	1.60

```
1 X_test2 = X_test[:1]
2 X_test2
```

	Body Weight (kg)	Body Height (m)
30	80	1.7
18	70	1.80

```
1 # Feature Scaling for KNN Regression...
2 from sklearn.preprocessing import StandardScaler
3
4 scaler = StandardScaler()
5
6 # Fit only on X_train...
7 scaler.fit(X_train)
8
9 X_t = X_test
10 X_t2 = X_test2
11
12 # Scale both X_train and X_test...
13 X_train = scaler.transform(X_train)
14 X_test = scaler.transform(X_test)
15 X_test2 = scaler.transform(X_test2)
```

```
1 X_test2
```

```
array([[ 0.54403307, -0.66422963]])
```

```
1 col_names=['Body Weight (kg)', 'Body Height (m)'] # , 'Weight-Height Ratio (R)', 'Internal Moment (M)']
2 scaled_df = pd.DataFrame(X_train, columns=col_names)
3 scaled_df.describe().T
```

	count	mean	std	min	25%	50%	75%	max
Body Weight (kg)	36.0	4.502571e-16	1.014185	-1.464704	-0.96252	0.041849	1.046217	1.548402
Body Height (m)	36.0	2.800229e-15	1.014185	-1.703893	-0.66423	-0.144398	0.895266	1.415098

```
1 # Training and Predicting K-NN Regression...
2 from sklearn.neighbors import KNeighborsRegressor
3 K = 5
4 regressor = KNeighborsRegressor(n_neighbors=K)
5 regressor.fit(X_train, y_train)
6
7 y_pred = regressor.predict(X_test)
8 print(X_test), print(y_test), print(y_pred)
```

```
array([[ 0.54403307, -0.66422963],
       [-1.46470441,  0.89526602],
       [ 1.54840181, -1.18406151],
       [ 0.0418487 , -1.7038934 ],
       [-1.46470441,  1.41509791],
       [ 1.04621744, -1.7038934 ],
       [-0.46033567, -0.66422963],
       [ 0.54403307, -1.18406151],
       [-0.46033567,  0.37543414],
       [-0.96252004,  0.89526602],
       [ 1.54840181, -0.66422963],
       [ 1.04621744, -1.18406151],
       [ 0.54403307, -1.7038934 ]])
```

```
1 y_pred2 = regressor.predict(X_test2)
2 y_pred2
```

```
array([12.2702])
```

```
1 col_names=['Required Torque (N-m)']
2 scaled_y_pred = pd.DataFrame(y_pred2, columns=col_names)
3 scaled_y_pred
```

Required Torque (N-m)

```
1 X_t2
```

Body Weight (kg) Body Height (m)

30	80	1.7
----	----	-----

```
1 #col_names=['Body Weight (kg)', 'Body Height (m)', 'Weight-Height Ratio (R)', 'Internal Moment (M)']
2 pred_df = X_t2 # pd.DataFrame(X_test2, columns=col_names)
3 pred_df
```

Body Weight (kg) Body Height (m)

30	80	1.7
----	----	-----

```
1 pred_df['Required Torque (N-m)'] = y_pred2
2 pred_df
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
"""Entry point for launching an IPython kernel.
```

Body Weight (kg) Body Height (m) Required Torque (N-m)

30	80	1.7	12.2702
----	----	-----	---------

```
1 # Evaluating the Algorithm for KNN Regression...
2 from sklearn.metrics import mean_absolute_error, mean_squared_error
3
4 mae = mean_absolute_error(y_test, y_pred)
5 mse = mean_squared_error(y_test, y_pred)
6 rmse = mean_squared_error(y_test, y_pred, squared=False)
7
8 print(f'mae: {mae}')
9 print(f'mse: {mse}')
10 print(f'rmse: {rmse}')
```

```
mae: 0.2714353846153841
mse: 0.10830680212307671
rmse: 0.32909998803262924
```

```
1 # R2 can be calculated...
2 regressor.score(X_test, y_test)
```

```
0.9105616219325816
```

```
1 y.describe()
```

```
count    49.000000
mean     11.636300
std       1.707826
min       8.511100
25%      10.240000
50%      11.636000
75%      12.811000
max       15.160000
Name: Required Torque (N-m), dtype: float64
```

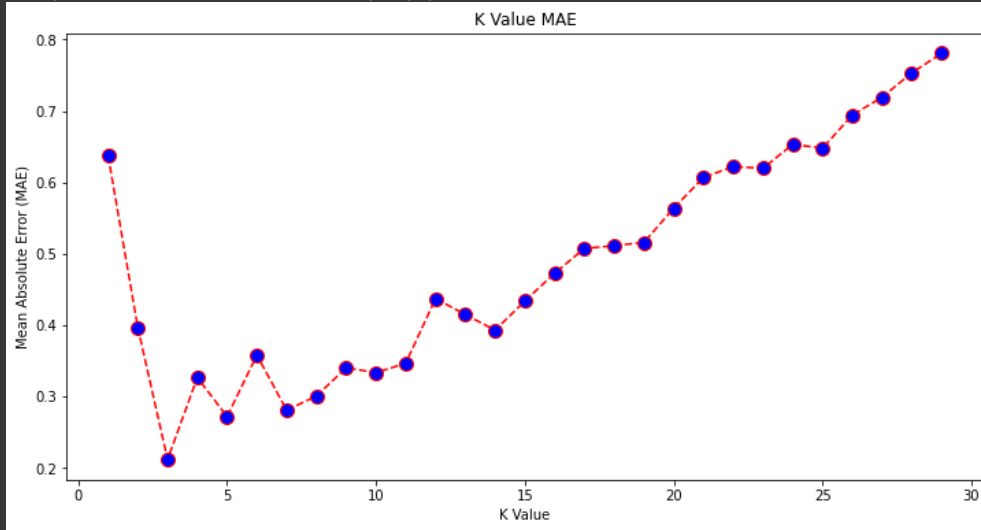
```
1 # =====
2 # Finding the Best K for KNN Regression...
3 # =====
4 error = []
5
6 # Calculating MAE error for K values between 1 and 29
7 for i in range(1, 30):
8     knn = KNeighborsRegressor(n_neighbors=i)
9     knn.fit(X_train, y_train)
10    pred_i = knn.predict(X_test)
11    mae = mean_absolute_error(y_test, pred_i)
12    error.append(mae)
```

```

1 import matplotlib.pyplot as plt
2
3 plt.figure(figsize=(12, 6))
4 plt.plot(range(1, 30), error, color='red',
5          linestyle='dashed', marker='o',
6          markerfacecolor='blue', markersize=10)
7
8 plt.title('K Value MAE')
9 plt.xlabel('K Value')
10 plt.ylabel('Mean Absolute Error (MAE)')

```

Text(0, 0.5, 'Mean Absolute Error (MAE)')



```

1 # =====
2 # Predicting with minimum error position as K...
3 # =====
4 import numpy as np
5
6 # print(min(error))
7 # print(np.array(error).argmin())
8
9 print(min(error[2:]))          # 0.14875230769230813
10 print(np.array(error[2:]).argmin()+2) # 4
11 K_neigh = np.array(error[2:]).argmin()+2
12
13 knn_reg = KNeighborsRegressor(n_neighbors=4)
14 knn_reg.fit(X_train, y_train)
15 y_pred_reg = knn_reg.predict(X_test)
16 r2_reg = knn_reg.score(X_test, y_test)
17
18 mae_reg = mean_absolute_error(y_test, y_pred_reg)
19 mse_reg = mean_squared_error(y_test, y_pred_reg)
20 rmse_reg = mean_squared_error(y_test, y_pred_reg, squared=False)
21 print(f'r2: {r2_reg}, \nmae: {mae_reg} \nmse: {mse_reg} \nrmse: {rmse_reg}')

```

```

0.21149487179487172
2
r2: 0.8485152716705489,
mae: 0.3261134615384617
mse: 0.18344279995192306
rmse: 0.42830222968357645

```

```

1 # =====
2 # Classification using K-Nearest Neighbors with Scikit-Learn...
3 # =====
4
5 dfc = df.dropna()
6
7 # Creating 4 categories and assigning them to a NewColumn...
8 # dfc["NewColumn"] = pd.qcut(dfc["ClassificatinColumn"], 4, retbins=False,
9 #                             labels=[1, 2, 3, 4])
10 y = dfc['Torque Category']
11 X = dfc.drop(['Weight-Height Ratio (R)', 'Internal Moment (M)',
12              'Required Torque (N-m)', 'Torque Category'], axis = 1)
13

```

```

1 # Splitting Data into Train and Test Sets...
2 from sklearn.model_selection import train_test_split
3
4 # SEED = 42
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25) #, random_state=SEED)

```

```

1 # Feature Scaling for Classification...
2 from sklearn.preprocessing import StandardScaler
3
4 scaler = StandardScaler()
5 scaler.fit(X_train)
6
7 X_train = scaler.transform(X_train)
8 X_test = scaler.transform(X_test)

```

```

1 # Training and Predicting for Classification...
2 from sklearn.neighbors import KNeighborsClassifier
3
4 classifier = KNeighborsClassifier() # n_neighbors = 5
5 classifier.fit(X_train, y_train)
6
7 y_pred = classifier.predict(X_test)
8 y_pred

```

```

array(['High', 'High', 'High', 'Medium', 'Medium', 'Medium', 'Medium',
       'Medium', 'High', 'High', 'High', 'Medium', 'Medium'], dtype=object)

```

```

1 # Evaluating KNN for Classification...
2 acc = classifier.score(X_test, y_test)
3 print(acc) # 0.8461538461538461

```

```

0.8461538461538461

```

```

1 # Visualize using a heatmap....
2 from sklearn.metrics import classification_report, confusion_matrix
3 #importing Seaborn's to use the heatmap
4 import seaborn as sns
5
6 # Adding classes names for better interpretation
7 classes_names = ['High','Medium']
8 cm = pd.DataFrame(confusion_matrix(y_test, y_pred),
9                   columns=classes_names, index = classes_names)
10
11 # Seaborn's heatmap to better visualize the confusion matrix
12 sns.heatmap(cm, annot=True, fmt='d');
13
14 print(classification_report(y_test, y_pred))

```

	precision	recall	f1-score	support
High	0.67	1.00	0.80	4
Medium	1.00	0.78	0.88	9
accuracy			0.85	13
macro avg	0.83	0.89	0.84	13
weighted avg	0.90	0.85	0.85	13



```

1 print(f'Test: {np.array(y_test)} \n')
2 print(f'Pred: {y_pred}')

```

```

Test: ['High' 'Medium' 'Medium' 'Medium' 'Medium' 'Medium' 'Medium' 'Medium'
       'High' 'High' 'High' 'Medium' 'Medium']

```

```

Pred: ['High' 'High' 'High' 'Medium' 'Medium' 'Medium' 'Medium' 'Medium'
       'High' 'High' 'Medium' 'Medium']

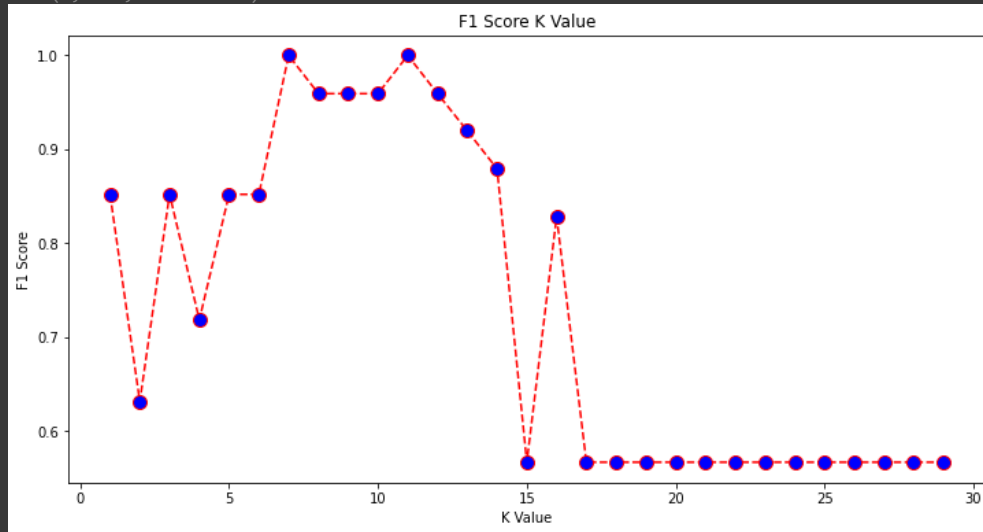
```

```

1 # =====
2 # Finding the Best K for KNN Classification...
3 # =====
4 from sklearn.metrics import f1_score
5
6 f1s = []
7
8 # Calculating f1 score for K values between 1 and 30
9 for i in range(1, 30):
10     knn = KNeighborsClassifier(n_neighbors=i)
11     knn.fit(X_train, y_train)
12     pred_i = knn.predict(X_test)
13     # average='weighted' to calculate a weighted average for the classes
14     f1s.append(f1_score(y_test, pred_i, average='weighted'))
15
16 plt.figure(figsize=(12, 6))
17 plt.plot(range(1, 30), f1s, color='red', linestyle='dashed', marker='o',
18         markerfacecolor='blue', markersize=10)
19 plt.title('F1 Score K Value')
20 plt.xlabel('K Value')
21 plt.ylabel('F1 Score')

```

Text(0, 0.5, 'F1 Score')



```
1 np.array(f1s)
```

```

array([0.85192308, 0.63116371, 0.85192308, 0.71828172, 0.85192308,
        0.85192308, 1.          , 0.95927602, 0.95927602, 0.95927602,
        1.          , 0.95927602, 0.91960671, 0.87912088, 0.56643357,
        0.82820513, 0.56643357, 0.56643357, 0.56643357, 0.56643357,
        0.56643357, 0.56643357, 0.56643357, 0.56643357, 0.56643357,
        0.56643357, 0.56643357, 0.56643357, 0.56643357, 0.56643357])

```

```

1 # The f1-score is the highest when the value of the K is 7 or 11...
2 # Retrain the classifier with 7 neighbors...
3 classifier7 = KNeighborsClassifier(n_neighbors=7)
4 classifier7.fit(X_train, y_train)
5 y_pred7 = classifier7.predict(X_test)
6 print(classification_report(y_test, y_pred7))

```

	precision	recall	f1-score	support
High	1.00	1.00	1.00	4
Medium	1.00	1.00	1.00	9
accuracy			1.00	13
macro avg	1.00	1.00	1.00	13
weighted avg	1.00	1.00	1.00	13

```

1 # Retrain the classifier with 11 neighbors...
2 classifier11 = KNeighborsClassifier(n_neighbors=11)
3 classifier11.fit(X_train, y_train)
4 y_pred11 = classifier11.predict(X_test)
5 print(classification_report(y_test, y_pred11))

```

	precision	recall	f1-score	support
High	1.00	1.00	1.00	4
Medium	1.00	1.00	1.00	9
accuracy			1.00	13

macro avg	1.00	1.00	1.00	13
weighted avg	1.00	1.00	1.00	13

```
1 # Retrain the classifier with 5 neighbors...
2 classifier5 = KNeighborsClassifier(n_neighbors=5)
3 classifier5.fit(X_train, y_train)
4 y_pred5 = classifier5.predict(X_test)
5 print(classification_report(y_test, y_pred5))
```

	precision	recall	f1-score	support
High	0.67	1.00	0.80	4
Medium	1.00	0.78	0.88	9
accuracy			0.85	13
macro avg	0.83	0.89	0.84	13
weighted avg	0.90	0.85	0.85	13

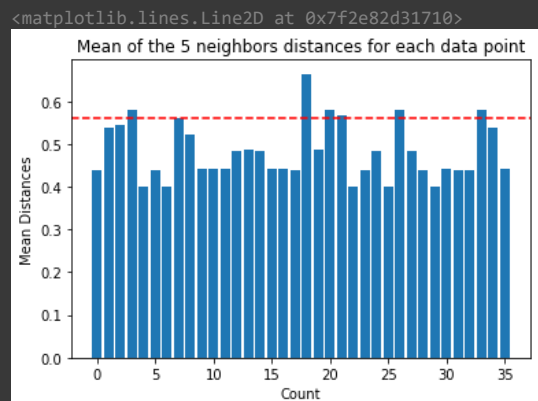
```
1 # =====
2 # Implementing KNN for Outlier Detection with Scikit-Learn...
3 # =====
4 from sklearn.neighbors import NearestNeighbors
5
6 nbrs = NearestNeighbors(n_neighbors = 5)
7 nbrs.fit(X_train)
8 # Distances and indexes of the 5 neighbors...
9 distances, indexes = nbrs.kneighbors(X_train)
10
11 # 5 distances for each data point (distance between itself and 5 neighbors)...
12 distances[:3], distances.shape
```

```
(array([[0.      , 0.48989795, 0.48989795, 0.51034318, 0.70742502],
        [0.      , 0.48989795, 0.51034318, 0.70742502, 0.9797959 ],
        [0.      , 0.48989795, 0.51034318, 0.70742502, 1.02068636]]),
      (36, 5))
```

```
1 # Look at the neighbors' indexes for 3 rows...
2 indexes[:3], indexes[:3].shape
```

```
(array([[ 0, 23, 27,  4, 25],
        [ 1, 23, 10, 25,  0],
        [ 2, 24, 11,  6, 30]]), (3, 5))
```

```
1 # calculate the mean of the 5 distances and plot a graph that counts
2 # each row on the X-axis and displays each mean distance on the Y-axis...
3 dist_means = distances.mean(axis=1)
4 plt.bar(np.array(range(0, 36)), dist_means)
5 plt.title('Mean of the 5 neighbors distances for each data point')
6 plt.xlabel('Count')
7 plt.ylabel('Mean Distances')
8
9 plt.axhline(y = 0.56, color = 'r', linestyle = '--')
```



```
1 import numpy as np
2
3 # Visually determine cutoff values > 0.56
4 outlier_index = np.where(dist_means > 0.56)
5 outlier_index
```

```
(array([ 3,  7, 18, 20, 21, 26, 33]),)
```

```
1 # Filter outlier values (locate them in the dataframe)...
2 outlier_values = df1.iloc[outlier_index]
3 outlier_values
```

	Body Weight (kg)	Body Height (m)	Weight-Height Ratio (R)	Internal Moment (M)	Required Torque (N-m)	Torque Category
3	60	1.75	34.286	0.94990	9.3091	Low