

PROJECT REPORT ON WINE QUALITY CLASSIFICATION

Programming Language: Python

1. Separating raw data into Training data and Testing Data (three times each Wine).

Read data:

```
1 red_wine = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\winequality-red.csv',header = 0)
2 white_wine = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\winequality-white.csv',header = 0)
```

Separating Training (90%) and Testing data three time for each:

export csv file

```
1 X = red_wine.drop('quality',axis = 1)
2 y = red_wine['quality']
3 for i in range(3):
4     i = i + 1
5     X_train,X_test,y_train,y_test = train_test_split(X,y,train_size = 0.9)
6     train = pd.concat([X_train,y_train],axis=1)
7     test = pd.concat([X_test,y_test],axis=1)
8     train.to_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\train\red_train_{}.csv'.format(i), index = False)
9     test.to_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\test\red_test_{}.csv'.format(i), index = False)
```

```
1 X = white_wine.drop('quality',axis = 1)
2 y = white_wine['quality']
3 for i in range(3):
4     i = i + 1
5     X_train,X_test,y_train,y_test = train_test_split(X,y,train_size = 0.9)
6     train = pd.concat([X_train,y_train],axis=1)
7     test = pd.concat([X_test,y_test],axis=1)
8     train.to_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\train\white_train_{}.csv'.format(i), index = False)
9     test.to_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\test\white_test_{}.csv'.format(i), index = False)
```

Determine Independent and Dependent variable:

```

1 red_train_1 = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\train\red_train_1.csv')
2 red_train_2 = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\train\red_train_2.csv')
3 red_train_3 = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\train\red_train_3.csv')
4 red_test_1 = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\test\red_test_1.csv')
5 red_test_2 = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\test\red_test_2.csv')
6 red_test_3 = pd.read_csv(r'C:\Users\Do Anh Luyen\Creative Cloud Files\Desktop\Task 2\test\red_test_3.csv')

```

```

1 X_red_train_1=red_train_1.drop(['quality'],axis=1)
2 X_red_train_2=red_train_2.drop(['quality'],axis=1)
3 X_red_train_3=red_train_3.drop(['quality'],axis=1)
4 X_red_test_1=red_test_1.drop(['quality'],axis=1)
5 X_red_test_2=red_test_2.drop(['quality'],axis=1)
6 X_red_test_3=red_test_3.drop(['quality'],axis=1)
7 y_red_train_1=red_train_1['quality']
8 y_red_train_2=red_train_2['quality']
9 y_red_train_3=red_train_3['quality']
10 y_red_test_1=red_test_1['quality']
11 y_red_test_2=red_test_2['quality']
12 y_red_test_3=red_test_3['quality']

```

2. KMeans Model

- Data preprocessing:

There is no missing value or noise data.

Using StandardScaler for scaling data:

```

1 #scaler
2 std = StandardScaler()
3 X_red_train_1 = std.fit_transform(X_red_train_1)
4 X_red_test_1 = std.transform(X_red_test_1)

```

- Determine number of clusters for Red and White wine:
- Group Dependent variable:

0 will be represented for low quality (quality from 0 to 7).

1 will be represented for high quality (quality from 7 to 10).

```

y_red_train_1= [ 0 if ( i < 7 ) else 1 for i in y_red_train_1]
y_red_train_2= [ 0 if ( i < 7 ) else 1 for i in y_red_train_2]
y_red_train_3= [ 0 if ( i < 7 ) else 1 for i in y_red_train_3]
y_red_test_1= [ 0 if ( i < 7 ) else 1 for i in y_red_test_1]
y_red_test_2= [ 0 if ( i < 7 ) else 1 for i in y_red_test_2]
y_red_test_3= [ 0 if ( i < 7 ) else 1 for i in y_red_test_3]

y_white_train_1 = [0 if ( i < 7 ) else 1 for i in y_white_train_1]
y_white_train_2 = [0 if ( i < 7 ) else 1 for i in y_white_train_2]
y_white_train_3 = [0 if ( i < 7 ) else 1 for i in y_white_train_3]
y_white_test_1 = [0 if ( i < 7 ) else 1 for i in y_white_test_1]
y_white_test_2 = [0 if ( i < 7 ) else 1 for i in y_white_test_2]
y_white_test_3 = [0 if ( i < 7 ) else 1 for i in y_white_test_3]

```

⇒ For Red_Wine we will have n_clusters = 2.

⇒ For White_Wine we will have n_clusters = 2.

- Using Kmeans:

Red_Wine:

- Training data:

```
#Kmeans
km = KMeans(n_clusters = 2, random_state = 42)
km.fit(X_red_train_1)
clusters = km.predict(X_red_train_1)
```

+ Using PCA to 2-Dimensional for Visualization:

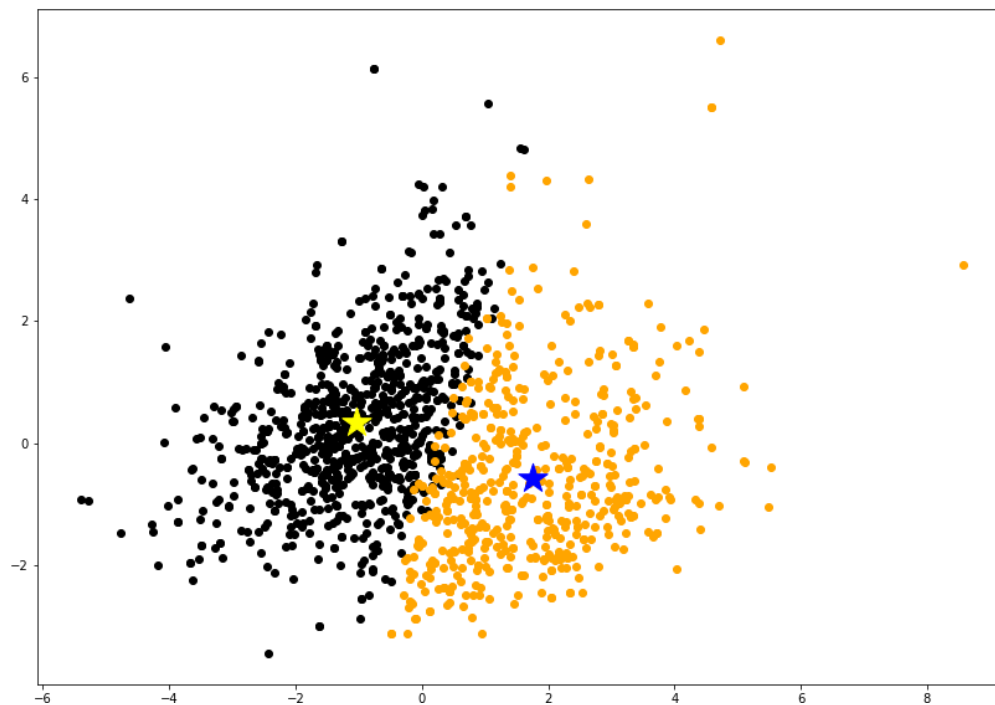
```
X_centre = km.cluster_centers_
```

```
pca = PCA(n_components = 2, random_state = 42)
X_red_train_1 = pd.DataFrame(pca.fit_transform(X_red_train_1), columns = ['PCA1', 'PCA2'])
X_red_train_1['cluster'] = clusters
X_centre = pd.DataFrame(pca.transform(X_centre), columns = ['PCA1', 'PCA2'])
```

+ Visualization:

```
1 plt.figure(figsize = (14,10))
2 plt.scatter(X_red_train_1[X_red_train_1['cluster'] == 0].loc[:, 'PCA1'], X_red_train_1[X_red_train_1['cluster'] == 0].loc[:, 'PCA2'])
3 plt.scatter(X_red_train_1[X_red_train_1['cluster'] == 1].loc[:, 'PCA1'], X_red_train_1[X_red_train_1['cluster'] == 1].loc[:, 'PCA2'])
4 plt.scatter(X_centre['PCA1'].loc[0], X_centre['PCA2'].loc[0], color = 'yellow', marker = '*', s = 600)
5 plt.scatter(X_centre['PCA1'].loc[1], X_centre['PCA2'].loc[1], color = 'blue', marker = '*', s = 600)
```

<matplotlib.collections.PathCollection at 0x245871c9940>



Yellow star is the central of 0 group.

Blue star is the central of 1 group.

- Tesing Data:

```

1 #for testing
2 clusters = km.predict(X_red_test_1)
3 X_centre = km.cluster_centers_
4
5 X_red_test_1 = pd.DataFrame(pca.transform(X_red_test_1), columns = ['PCA1', 'PCA2'])
6 X_red_test_1['cluster'] = clusters
7 X_centre = pd.DataFrame(pca.transform(X_centre), columns = ['PCA1', 'PCA2'])

```

+ Confusion matrix:

```

1 classificationSummary(y_red_test_1, clusters)

```

Confusion Matrix (Accuracy 0.6687)

	Prediction	
Actual	0	1
0	98	43
1	10	9

```

1 precision = (9/(9+43))
2 print("Precision : {}".format(precision))

```

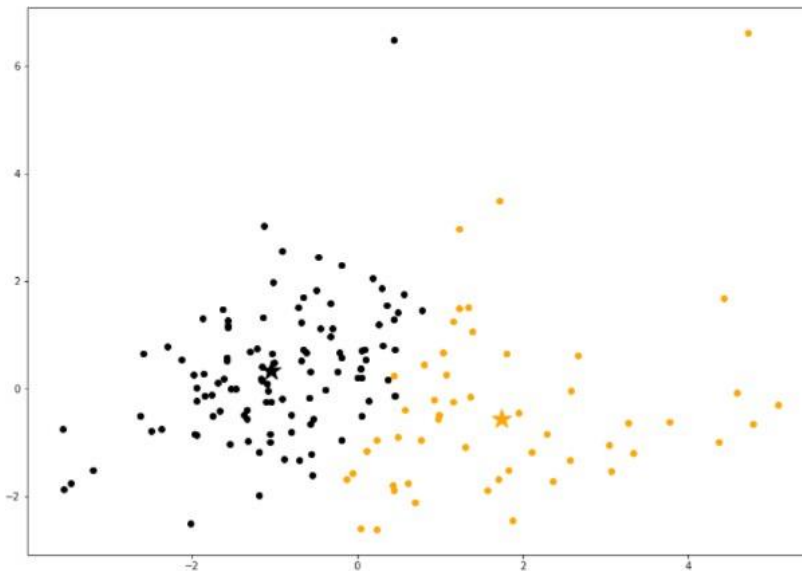
Precision : 0.17307692307692307

Accuracy = 0.6687,

True positive = 9, False positive = 43, TN = 98, FN = 10,

Precision = 0.173.

+ Visualization:



- For another testing set, we will have score below:

```
1 classificationSummary(y_red_test_2,clusters)
```

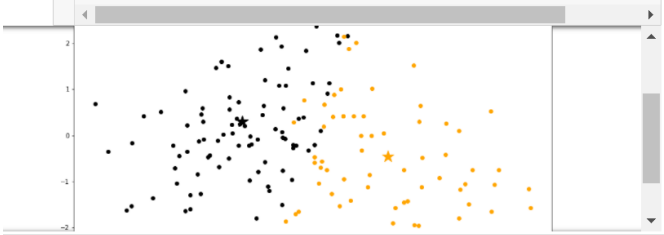
Confusion Matrix (Accuracy 0.6375)

	Prediction
Actual 0	1
0	87 53
1	5 15

```
1 precision = (15/(15+53))
2 print("Precision : {}".format(precision))
```

Precision : 0.22058823529411764

```
In [23]: 1 plt.figure(figsize = (14,10))
2 plt.scatter(X_red_test_2[X_red_test_2['cluster'] == 0].loc[:, 'PCA1'],X_red_test_2[X_red_test_2['cluster'] == 0].loc[:, 'PCA2']
3 plt.scatter(X_red_test_2[X_red_test_2['cluster'] == 1].loc[:, 'PCA1'],X_red_test_2[X_red_test_2['cluster'] == 1].loc[:, 'PCA2']
4
5 plt.scatter(X_centre['PCA1'].loc[0],X_centre['PCA2'].loc[0],color = 'black',marker = '*',s = 400)
6 plt.scatter(X_centre['PCA1'].loc[1],X_centre['PCA2'].loc[1],color = 'orange',marker = '*',s = 400)
7
```



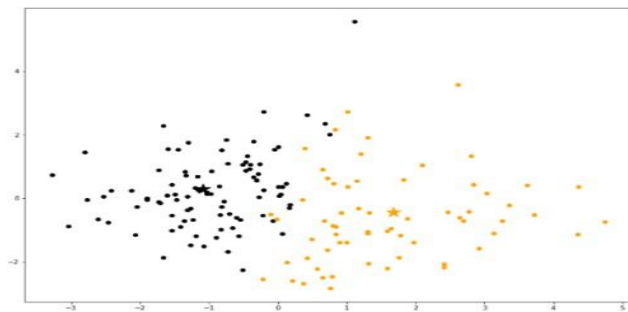
```
1 classificationSummary(y_red_test_3,clusters)
```

Confusion Matrix (Accuracy 0.6188)

	Prediction
Actual 0	1
0	85 52
1	9 14

```
1 precision = (14/(14+52))
2 print("Precision : {}".format(precision))
```

Precision : 0.21212121212121213



- Calculate Mean of 3 testing:

```
1 Mean = (a+b+c)/3
2 print("Mean_Confusion_matrix_Red_wine:\n {}".format(Mean))
3 mean = (a_1 + b_1 + c_1)/3
4 print("Mean_accuracy_Red_wine: {}".format(mean))
```

Mean_Confusion_matrix_Red_wine:

```
[[90.      49.33333333]
 [ 8.      12.66666667]]
```

Mean_accuracy_Red_wine: 0.6416666666666666

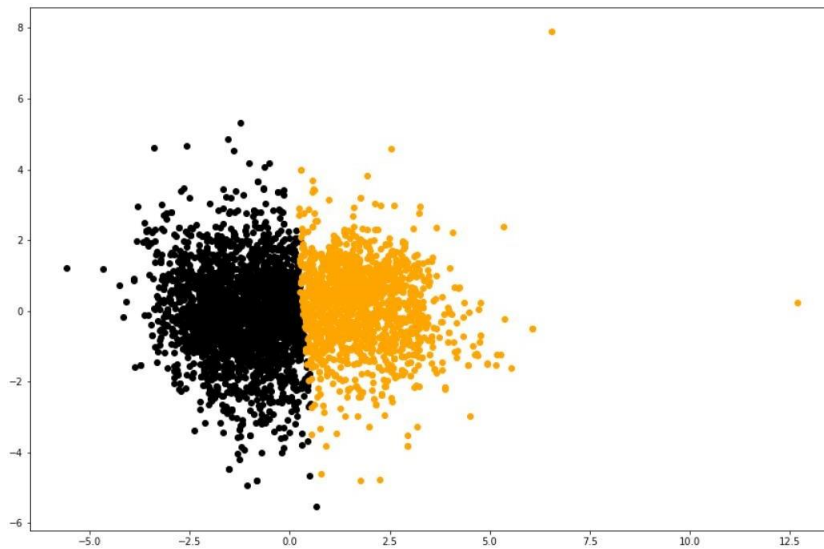
⇒ In conclusion, we have mean_accuracy equals 0.64

White_Wine:

- Training data set:

```
#train_1
std_1 = StandardScaler()
X_white_train_1 = std_1.fit_transform(X_white_train_1)
X_white_test_1 = std_1.transform(X_white_test_1)
km = KMeans(n_clusters = 2,random_state=42)
km.fit(X_white_train_1)
clusters = km.predict(X_white_train_1)
X_centre = km.cluster_centers_
pca_1 = PCA(n_components = 2,random_state=42)
X_white_train_1 = pd.DataFrame(pca_1.fit_transform(X_white_train_1),columns = ['PCA1','PCA2'])
X_white_train_1['cluster'] = clusters
X_centre = pd.DataFrame(pca_1.transform(X_centre),columns = ['PCA1','PCA2'])
```

+ Visualization:



- Testing data set:

```
#test1
clusters = km.predict(X_white_test_1)
X_white_test_1 = pd.DataFrame(pca_1.transform(X_white_test_1), columns = ['PCA1', 'PCA2'])
X_white_test_1['cluster'] = clusters
```

+ Confusion Matrix:

```
1 classificationSummary(y_white_test_1, clusters)
```

Confusion Matrix (Accuracy 0.4551)

	Prediction	
Actual	0	1
0	206	192
1	75	17

```
1 precision = (17/(17+192))
2 print("Precision : {}".format(precision))
```

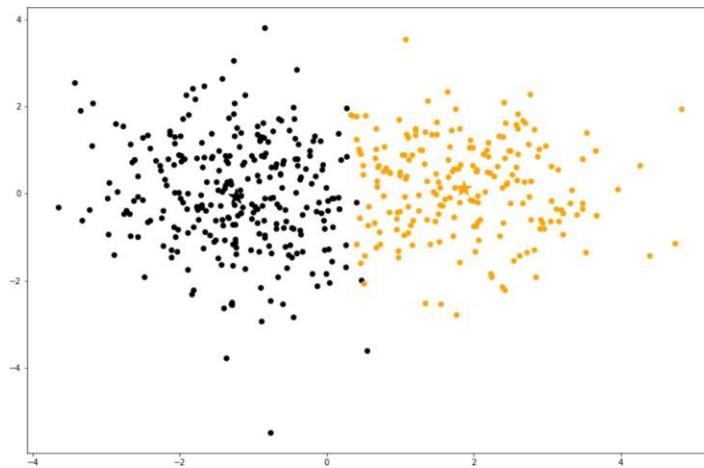
Precision : 0.08133971291866028

Accuracy = 0.455,

TP = 17, FP = 192, TN = 206, FN = 75,

Precision = 0.08.

+ Visualization:



- For another testing sets:

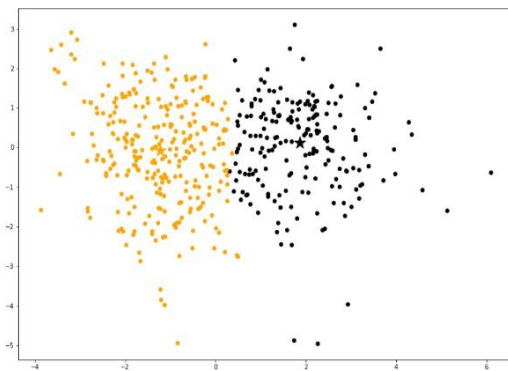
```
1 classificationSummary(y_white_test_2,clusters)
```

Confusion Matrix (Accuracy 0.6020)

	Prediction	
Actual	0	1
0	197	187
1	8	98

```
1 precision = (98/(98+187))
2 print("Precision : {}".format(precision))
```

Precision : 0.34385964912280703



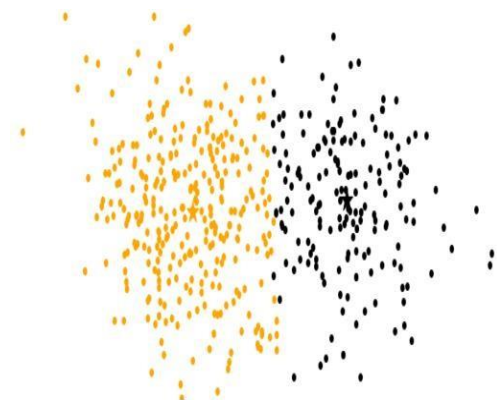
```
1 classificationSummary(y_white_test_3,clusters)
```

Confusion Matrix (Accuracy 0.5327)

	Prediction	
Actual	0	1
0	163	217
1	12	98

```
1 precision = (98/(98+217))
2 print("Precision : {}".format(precision))
```

Precision : 0.3111111111111111



- Calculate Mean of 3 testing:

```
1 Mean1 = (d+e+f)/3
2 print("Mean_Confusion_matrix_White_wine:\n {}".format(Mean1))
3 mean1 = (d_1 + e_1 + f_1)/3
4 print("Mean_accuracy_White_wine: {}".format(mean1))
```

Mean_Confusion_matrix_White_wine:
[[188.66666667 198.66666667]
[31.66666667 71.]]
Mean_accuracy_White_wine: 0.5299319727891156

⇒ In conclusion, Mean_accuracy for White_wine is 0.53.

3. Regression model:

- Preprocessing data (StandardScaler):

```
mm = MinMaxScaler()
X_red_train_1 = mm.fit_transform(X_red_train_1)
X_red_test_1 = mm.transform(X_red_test_1)
```

```
1 X_white_train_1 = mm.fit_transform(X_white_train_1)
2 X_white_test_1 = mm.transform(X_white_test_1)
```

- Linear Regression:

Red_Wine:

```
lr = LinearRegression().fit(X_red_train_1,y_red_train_1)
```

- Confusion matrix:

```
1 confusion_matrix = confusion_matrix(y_red_test_1,y_prediction )
2 print(confusion_matrix)
```

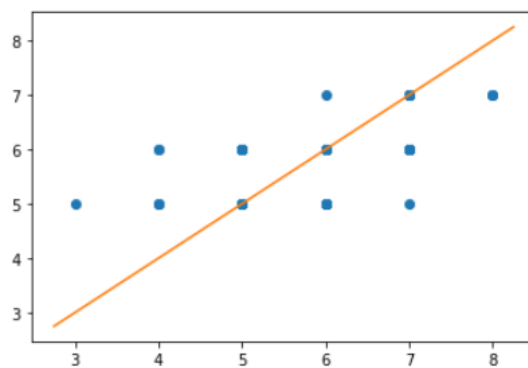
```
[[ 0  0  1  0  0  0]
 [ 0  0  3  2  0  0]
 [ 0  0 41 18  0  0]
 [ 0  0 32 43  1  0]
 [ 0  0  1  8  7  0]
 [ 0  0  0  0  3  0]]
```

- False positive, false negative, true positive, true negative, precision, accuracy.

	3	4	5	6	7	8
FP	0.0	0.0	37.000000	28.000000	4.000000	0.0
FN	1.0	5.0	18.000000	33.000000	9.000000	3.0
TP	0.0	0.0	41.000000	43.000000	7.000000	0.0
TN	159.0	155.0	64.000000	56.000000	140.000000	157.0
Pecision	NaN	NaN	0.525641	0.605634	0.636364	NaN

Accuracy: Accuracy: 0.2905520217191234

- Visualization:



White_Wine:

```
lr = LinearRegression().fit(X_white_train_1,y_white_train_1)
```

- Confusion matrix:

```
1 confusion_matrix = confusion_matrix(y_white_test_1,y_prediction )
2 print(confusion_matrix)
```

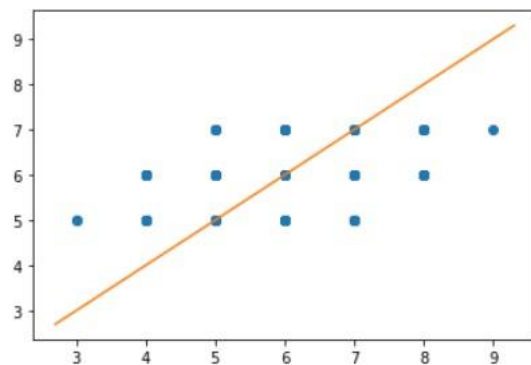
```
[[ 0  0  2  0  0  0  0]
 [ 0  0 12 10  0  0  0]
 [ 0  0 55 89  4  0  0]
 [ 0  0 31 172 23  0  0]
 [ 0  0  5 48 21  0  0]
 [ 0  0  0 10  7  0  0]
 [ 0  0  0  0  1  0  0]]
```

- False positive, false negative, true positive, true negative, precision, accuracy.

	3	4	5	6	7	8	9
FP	0.0	0.0	50.00000	157.000000	35.000	0.0	0.0
FN	2.0	22.0	93.00000	54.000000	53.000	17.0	1.0
TP	0.0	0.0	55.00000	172.000000	21.000	0.0	0.0
TN	488.0	468.0	292.00000	107.000000	381.000	473.0	489.0
Precision	NaN	NaN	0.52381	0.522796	0.375	NaN	NaN

Accuracy: 0.2655541307739012

- Visualization:



• Logistic Regression:

Red_Wine:

```
In [68]: 1 lg = LogisticRegression(max_iter = 10000).fit(X_red_train_1,y_red_train_1)
2
3 print("Accuracy: {}".format(lg.score(X_red_test_1,y_red_test_1)))
4 y_prediction = lg.predict(X_red_test_1)
5 y_true = y_red_test_1

Accuracy: 0.4875

In [69]: 1 confusion_matrix = confusion_matrix(y_red_test_1,y_prediction )
2 print(confusion_matrix)

[[ 0  0  1  0  0  0]
 [ 0  0  2  3  0  0]
 [ 0  0 44 15  0  0]
 [ 0  0 43 32  1  0]
 [ 0  0  1 13  2  0]
 [ 0  0  0  1  2  0]]

In [70]: 1 X = a(y_true,y_prediction)
2 pd.DataFrame(X,columns = sorted(y_red_test_1.unique()), index = ['FP','FN','TP','TN','Prcision'])

<ipython-input-58-87a397bf5c65>:8: RuntimeWarning: invalid value encountered in true_divide
Precision = TP/(TP+FP)

Out[70]:
```

	3	4	5	6	7	8
FP	0.0	0.0	47.000000	32.0	3.0	0.0
FN	1.0	5.0	15.000000	44.0	14.0	3.0
TP	0.0	0.0	44.000000	32.0	2.0	0.0
TN	159.0	155.0	54.000000	52.0	141.0	157.0
Precision	NaN	NaN	0.483516	0.5	0.4	NaN

In [69] is the confusion matrix.

In [70] is the False positive, false negative, true positive, true negative, precision, accuracy.

In [68] is the logistic model application and Accuracy = 0.4875.

White_Wine:

```
In [79]: 1 lg = LogisticRegression(max_iter = 10000).fit(X_white_train_1,y_white_train_1)
2 print("Accuracy: {}".format(lg.score(X_white_test_1,y_white_test_1)))
3 y_prediction = lg.predict(X_white_test_1)
4 y_true = y_white_test_1

Accuracy: 0.5204081632653061

In [80]: 1 confusion_matrix = confusion_matrix(y_white_test_1,lg.predict(X_white_test_1))
2 print(confusion_matrix)

[[ 0  0  2  0  0  0  0]
 [ 0  0 15  7  0  0  0]
 [ 0  0 69 78  1  0  0]
 [ 0  0 42 173 11  0  0]
 [ 0  0  5  56 13  0  0]
 [ 0  0  1  12  4  0  0]
 [ 0  0  0  0  1  0  0]]

In [81]: 1 X = a(y_true,y_prediction)
2 pd.DataFrame(X,columns = sorted(y_white_test_1.unique()), index = ['FP','FN','TP','TN','Precision'])

<ipython-input-58-87a397bf5c65>:8: RuntimeWarning: invalid value encountered in true_divide
Precision = TP/(TP+FP)

Out[81]:
```

	3	4	5	6	7	8	9
FP	0.0	0.0	65.000000	153.000000	17.000000	0.0	0.0
FN	2.0	22.0	79.000000	53.000000	61.000000	17.0	1.0
TP	0.0	0.0	69.000000	173.000000	13.000000	0.0	0.0
TN	488.0	468.0	277.000000	111.000000	399.000000	473.0	489.0
Precision	NaN	NaN	0.514925	0.530675	0.433333	NaN	NaN

In [80] is confusion matrix.

In [81] is the False positive, false negative, true positive, true negative, precision, accuracy.

In [79] is the logistic model application and Accuracy is 0.52.

4. Comparing:

Before separating Dependent variable (quality) into two group 0 and 1 (low and high quality) and run Kmeans with n_clusters = 2 for both

Red_wine and White_wine, my team has already run Kmeans without separating Dependent variable which means $n_clusters = 6$ for Red_wine and $n_clusters = 7$ for White_wine, we saw that our accuracy was very bad just 20% or less. This score is lower than LinearRegression model (around 30%) because Kmeans is an unsupervised learning model. This means Kmeans didn't learn from the dependent variable, it just clusters based on the distance between each point, higher clusters will lead to lower accuracy, while Linear Regression learns from both Independent and Dependent variables. Otherwise, Logistic Regression has a higher score than Linear Regression because the dependent variable in this exercise is a categorical variable so a classification model will be the best.

Eg: Logistic Regression, Decision Tree, Random Forest, Support Vector Machine.