

HW 4

This assignment covers Linear Classification methods

DO NOT ERASE MARKDOWN CELLS AND INSTRUCTIONS IN YOUR HW submission

- **Q** - QUESTION
- **A** - Where to input your answer

Instructions

Keep the following in mind for all notebooks you develop:

- Structure your notebook.
- Use headings with meaningful levels in Markdown cells, and explain the questions each piece of code is to answer or the reason it is there.
- Make sure your notebook can always be rerun from top to bottom.
- Please start working on this assignment as soon as possible. If you are a beginner in Python this might take a long time. One of the objectives of this assignment is to help you learn python and scikit-learn package.
- See [README.md \(README.md\)](#) for homework submission instructions

Related Tutorials

Refreshers

- [Intro to Machine Learning w scikit-learn \(https://scikit-learn.org/stable/tutorial/basic/tutorial.html\)](https://scikit-learn.org/stable/tutorial/basic/tutorial.html)
- [A tutorial on statistical-learning for scientific data processing \(https://scikit-learn.org/stable/tutorial/statistical_inference/index.html#stat-learn-tut-index\)](https://scikit-learn.org/stable/tutorial/statistical_inference/index.html#stat-learn-tut-index)

Classification Approaches

- [Logistic Regression with Sklearn \(https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)
- [KNN with sklearn \(https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html)

Modeling

- [Cross-validation \(https://scikit-learn.org/stable/modules/cross_validation.html\)](https://scikit-learn.org/stable/modules/cross_validation.html)
- [Plot Confusion Matrix with Sklearn \(https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html\)](https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html)
- [Confusion Matrix Display \(https://scikit-learn.org/stable/modules/generated/sklearn.metrics.ConfusionMatrixDisplay.html#sklearn\)](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.ConfusionMatrixDisplay.html#sklearn)

Data Processing

Q1 Get training data from the dataframe

1. Load mobile_data.csv from ""data" folder into the dataframe
2. Assign values of price_range column to y
3. Drop 'price_range' column from data frame,
4. Assign remaining df column values to x
5. Print the head of the dataframe

A1 Replace ??? with code in the code cell below

```
In [5]: import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from matplotlib import pyplot as plt
from sklearn.preprocessing import StandardScaler

#Read the mobile_data.csv file using the prrropriate separator as input
df = pd.read_csv('../data/mobile_data.csv')

df.head()
```

Out [5]:

	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_
0	842	0	2.2	0	1	0	7	0.6	188	
1	1021	1	0.5	1	0	1	53	0.7	136	
2	563	1	0.5	1	2	1	41	0.9	145	
3	615	1	2.5	0	0	0	10	0.8	131	
4	1821	1	1.2	0	13	1	44	0.6	141	

5 rows × 21 columns

In [6]: `df.head()`

Out[6]:

	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_
0	842	0	2.2	0	1	0	7	0.6	188	
1	1021	1	0.5	1	0	1	53	0.7	136	
2	563	1	0.5	1	2	1	41	0.9	145	
3	615	1	2.5	0	0	0	10	0.8	131	
4	1821	1	1.2	0	13	1	44	0.6	141	

5 rows × 21 columns

In [7]: `y= df['price_range']`
`x= df.drop(columns='price_range')`

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

Out[8]:

	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_
0	842	0	2.2	0	1	0	7	0.6	188	
1	1021	1	0.5	1	0	1	53	0.7	136	
2	563	1	0.5	1	2	1	41	0.9	145	
3	615	1	2.5	0	0	0	10	0.8	131	
4	1821	1	1.2	0	13	1	44	0.6	141	

5 rows × 21 columns

Q2:

1. Check number of null values per column in the x dataframe.

A2 Replace ??? with code in the code cell below

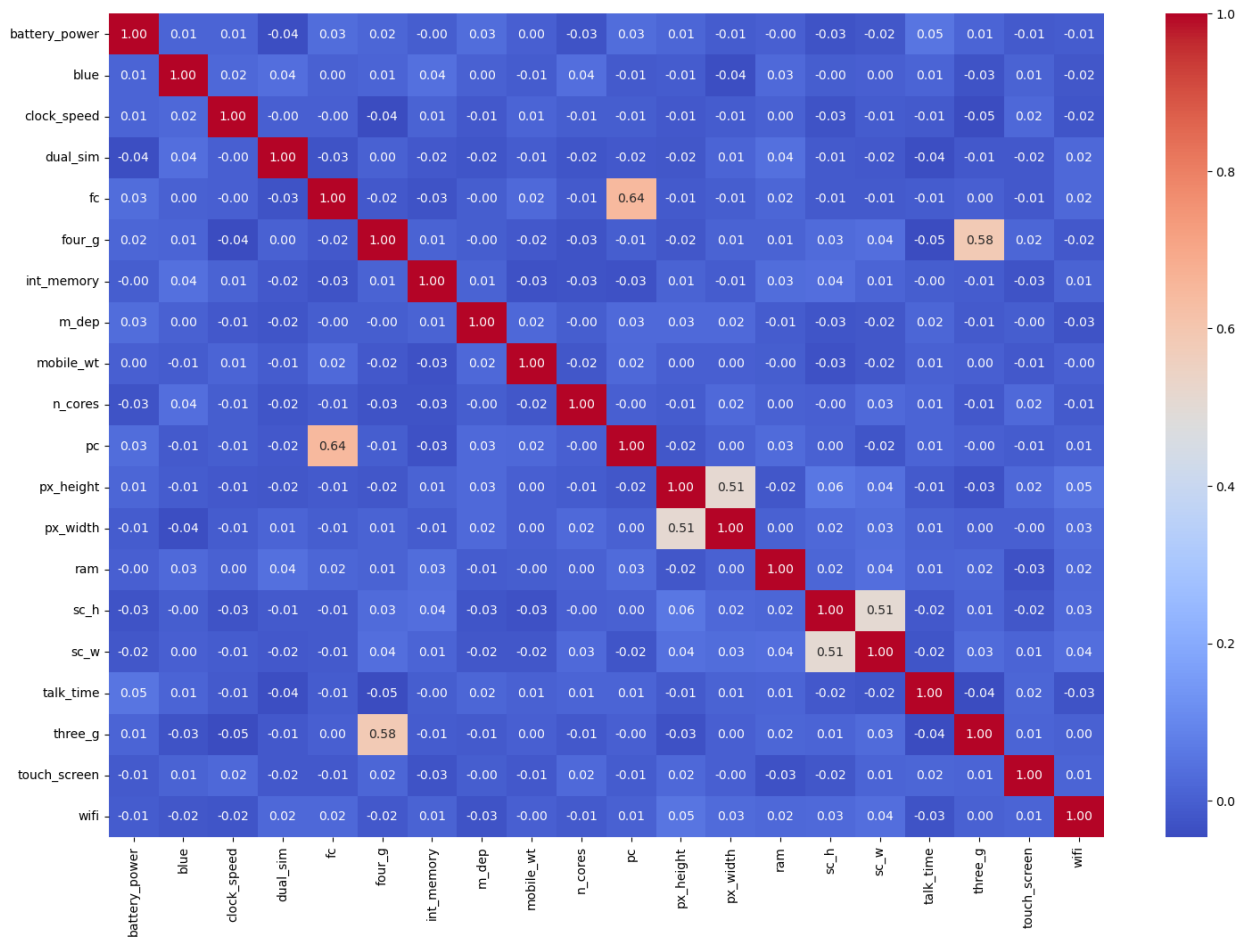
```
In [9]: x.isnull().sum()
```

```
Out[9]: battery_power    0  
blue                    0  
clock_speed             0  
dual_sim                0  
fc                      0  
four_g                  0  
int_memory              0  
m_dep                   0  
mobile_wt               0  
n_cores                 0  
pc                      0  
px_height               0  
px_width                0  
ram                     0  
sc_h                    0  
sc_w                    0  
talk_time               0  
three_g                 0  
touch_screen            0  
wifi                    0  
dtype: int64
```

Q3.1 Use seaborn heatmap chart to visualize the correlations between the columns. Replace ??? with code in the code cell below

A3.1

```
In [11]: import seaborn as sns
plt.figure(figsize=(18,12))
sns.heatmap(x.corr(), annot=True, fmt=".2f", cmap='coolwarm')
plt.show()
```



Q3.2 List columns that correlate the most with the 'price_range' column.

Note: For this dataset any column that has correlation factor over or near 0.1 can be considered as a good predictor/ feature.

A3.2 Answer here

Q3.3 Update the 'x' dataframe defined earlier in Q1 with your selected features/columns for 'price_range'.

A3.3 Replace ??? with code in the code cell below

In [12]: `# A3 Part 3:`

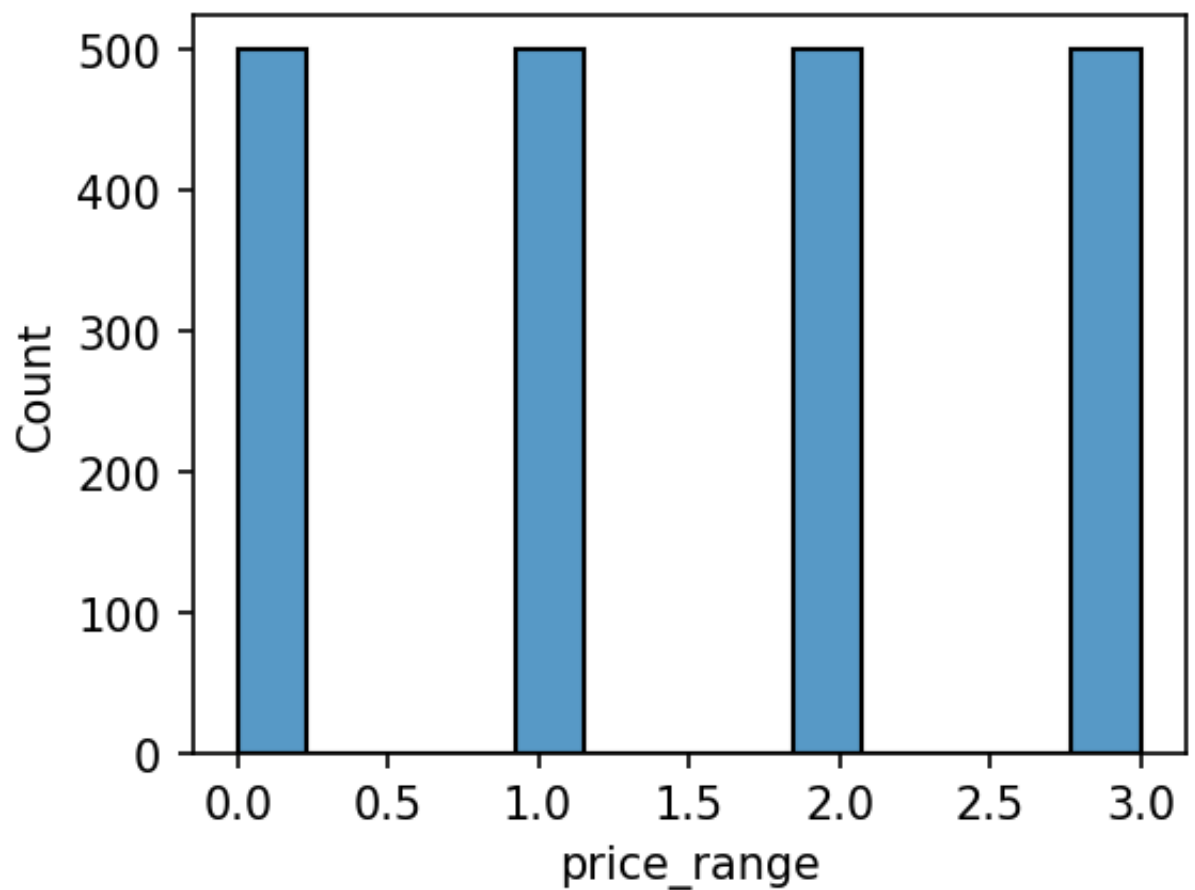
```
x = x[['ram', 'battery_power', 'n_cores']]
```

Q4: Use seaborn *histplot* to plot a distribution graph for the *price_range* column

A4 Replace ??? with code in the code cell below

In [13]: `plt.figure(figsize=(4,3),dpi=150)
sns.histplot(data=y)`

Out[13]: `<Axes: xlabel='price_range', ylabel='Count'>`

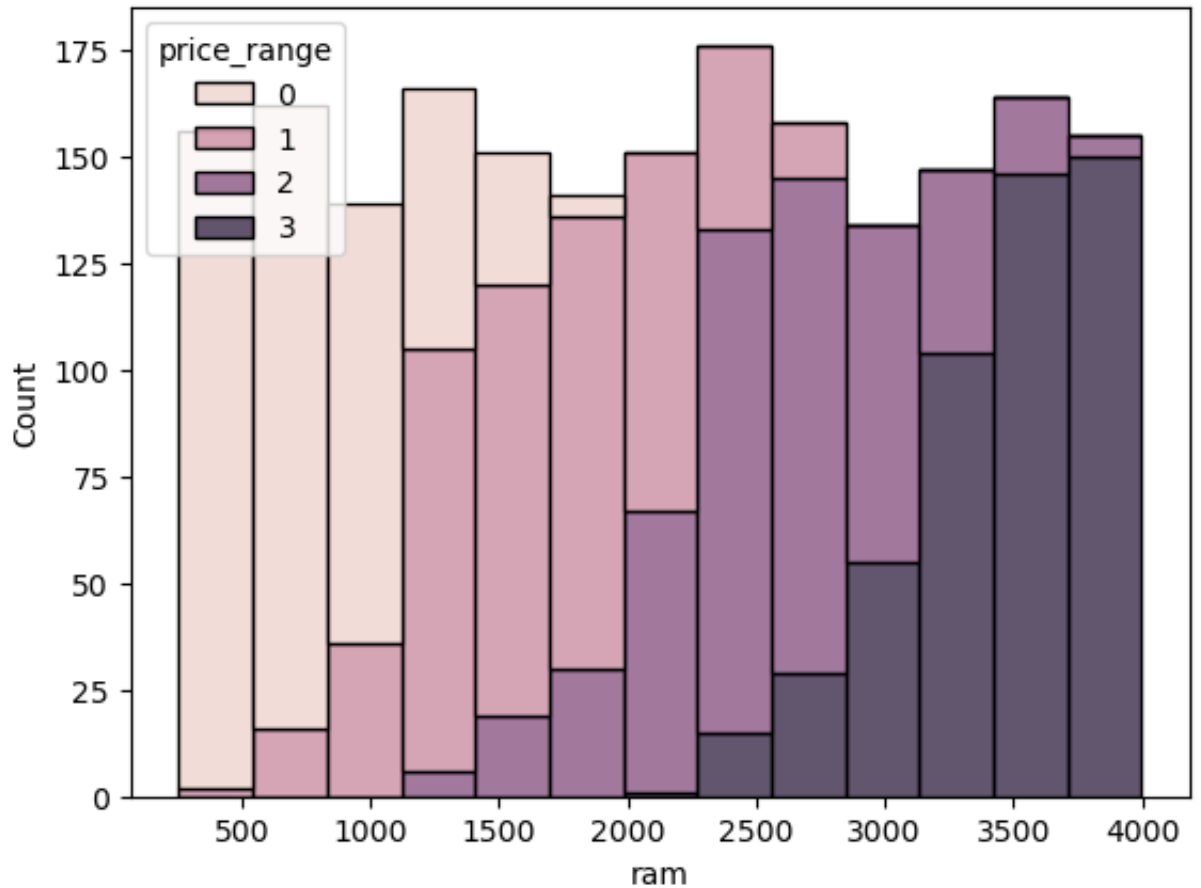


Q5: Use seaborn *histplot* to present the relation between *price_range* and the *ram* of a mobile

A5 Replace ??? with code in the code cell below

```
In [14]: sns.histplot(data=df, x='ram', hue='price_range', multiple="stack")
```

```
Out[14]: <Axes: xlabel='ram', ylabel='Count'>
```



Q6:

1. Use StandardScaler from sklearn to transform the x dataframe.
2. Split dataset into train and test data use train_test_split with test_size = 0.2 and random_state = 42
3. Check the number of instance in the train and test set.
4. Check the number of instance per class in train and test set using ytrain and ytest

A6 Replace ??? with code in the code cell below

```
In [15]: scaler = StandardScaler()
scaler.fit(x)
x = scaler.transform(x)
```

```
In [16]: xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size=0.2, r
```

```
In [17]: print(len(xtrain))  
print(len(xtest))
```

```
1600  
400
```

```
In [18]: ytrain.value_counts()
```

```
Out[18]: price_range  
1      409  
2      408  
0      395  
3      388  
Name: count, dtype: int64
```

```
In [19]: ytest.value_counts()
```

```
Out[19]: price_range  
3      112  
0      105  
2       92  
1       91  
Name: count, dtype: int64
```

Classification Model 1: Logistic Regression

Here, we fit Logistic Regression model to the train dataset using K-fold cross validation

Q7 Train Logistic Regression Model

1. Create a logistic regression model using sklearn [linear_model \(https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html) library.
2. Fit the model with the train data
3. Get the score from the model using test data
4. Plot confusion matrix using [ConfusionMatrixDisplay \(https://scikit-learn.org/stable/modules/generated/sklearn.metrics.ConfusionMatrixDisplay.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.ConfusionMatrixDisplay.html), see [Visualization with Display Objects \(https://scikit-learn.org/stable/auto_examples/miscellaneous/plot_display_object_visualization.html\)](https://scikit-learn.org/stable/auto_examples/miscellaneous/plot_display_object_visualization.html) example.

A7 Replace ??? with code in the code cell below

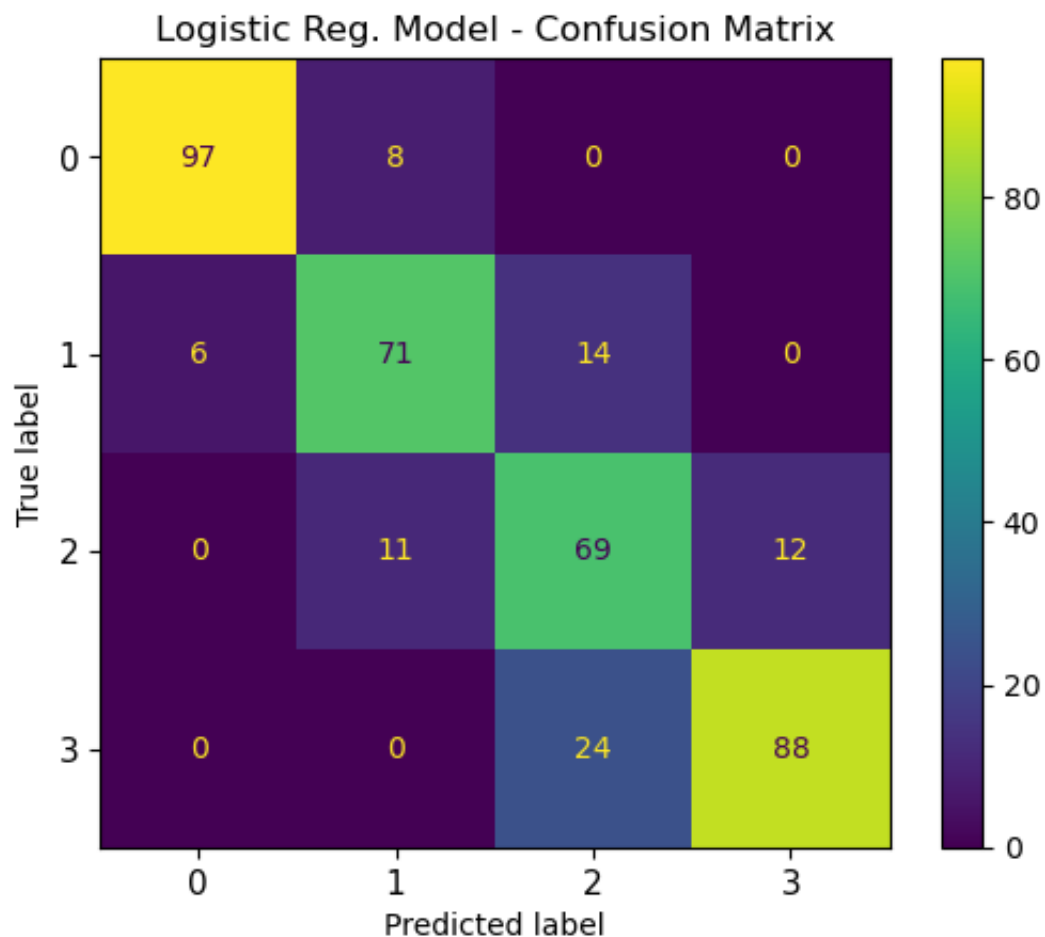

```
In [20]: from sklearn.metrics import ConfusionMatrixDisplay
import matplotlib.pyplot as plt

# Create a logistic regression model using sklearn library
clf=LogisticRegression()
clf.fit(xtrain, ytrain)

#print score for test data
print(clf.score(xtest, ytest))

0.8125
```

```
In [21]: cm = ConfusionMatrixDisplay.from_estimator(clf, xtest, ytest)
plt.title("Logistic Reg. Model - Confusion Matrix")
plt.xticks(ticks=np.arange(len(np.unique(y))), labels=np.unique(y), fo
plt.yticks(ticks=np.arange(len(np.unique(y))), labels=np.unique(y), fo
plt.show()
```



Q8: Train Logistic Regression Model using cross-validation on *xtrain*, *ytrain* data.

- Apply K fold cross validation technique for the model training (cross_val_score), and set K to 5 or 10.
- Print the different scores from different folds

A8: Replace ??? with code in the code cell below

```
In [22]: from sklearn.model_selection import cross_val_score

# Use sklearn for 5 fold cross validation
scores_log= cross_val_score(clf, x, y, cv=5)

# print the scores from different folds
print(scores_log)

[0.82    0.8175 0.8425 0.82    0.815 ]
```

Classification Model 2: K Nearest Neighbor Classifier

Here, we learn how to fit KNN on the train dataset using k-fold cross validation, and evaluate its classification accuracy on the train dataset using confusion matrix.

Q9 Build a KNN Classification Model for the dataset as following:

1. Create a KNN model using sklearn library, and initialize n_neighbors as described in [documentation \(https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html).
2. Fit the model with the train data
3. Predict the values from test data
4. Print out the score from training and test data
5. Repeat Step 1.- 4. for a range of n_neighbors values (k in kNN) from 1 to 30.

A9 Replace ??? with code in the code cell below

In [23]: `from sklearn.neighbors import KNeighborsClassifier`

```
# Define KNN model
for k in range(1,31):

    knn = KNeighborsClassifier(n_neighbors=k)

    #Fit KNN model on xtrain, ytrain from above
    knn.fit(xtrain, ytrain)
    #predict y values from xtest
    y_pred=knn.predict(xtest)

    #print score for test data
    print("K: ",k,"Train Score: ",knn.score(xtrain, ytrain), "Test Score: ")
```

```
K: 1 Train Score: 1.0 Test Score: 0.745
K: 2 Train Score: 0.89 Test Score: 0.7525
K: 3 Train Score: 0.88625 Test Score: 0.7725
K: 4 Train Score: 0.86 Test Score: 0.7925
K: 5 Train Score: 0.86875 Test Score: 0.805
K: 6 Train Score: 0.855625 Test Score: 0.7925
K: 7 Train Score: 0.86 Test Score: 0.795
K: 8 Train Score: 0.84875 Test Score: 0.7875
K: 9 Train Score: 0.85 Test Score: 0.805
K: 10 Train Score: 0.840625 Test Score: 0.79
K: 11 Train Score: 0.845625 Test Score: 0.8
K: 12 Train Score: 0.844375 Test Score: 0.7925
K: 13 Train Score: 0.84375 Test Score: 0.79
K: 14 Train Score: 0.839375 Test Score: 0.7925
K: 15 Train Score: 0.8375 Test Score: 0.8025
K: 16 Train Score: 0.84125 Test Score: 0.7975
K: 17 Train Score: 0.843125 Test Score: 0.7875
K: 18 Train Score: 0.844375 Test Score: 0.795
K: 19 Train Score: 0.838125 Test Score: 0.8
K: 20 Train Score: 0.840625 Test Score: 0.805
K: 21 Train Score: 0.835 Test Score: 0.7975
K: 22 Train Score: 0.84 Test Score: 0.7975
K: 23 Train Score: 0.83625 Test Score: 0.8
K: 24 Train Score: 0.836875 Test Score: 0.805
K: 25 Train Score: 0.83625 Test Score: 0.81
K: 26 Train Score: 0.836875 Test Score: 0.81
K: 27 Train Score: 0.84125 Test Score: 0.8175
K: 28 Train Score: 0.838125 Test Score: 0.8125
K: 29 Train Score: 0.84 Test Score: 0.8175
K: 30 Train Score: 0.839375 Test Score: 0.8075
```

Q9 Part 2:

What is the best `n_neighbors` ? Why?

A9 19 because it provides one of the highest accuracies on the test set without being too complex.

Q10.

1. Create a KNN Classifier model using the best value of `k` found from previous question.
2. Train the model using `xtrain`, `ytrain` values.
3. Plot confusion matrix for the `xtest` and `ytest`, using [ConfusionMatrixDisplay](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.ConfusionMatrixDisplay.html) (<https://scikit-learn.org/stable/modules/generated/sklearn.metrics.ConfusionMatrixDisplay.html>), see [Visualization with Display Objects](https://scikit-learn.org/stable/auto_examples/miscellaneous/plot_display_object_visualization.html) (https://scikit-learn.org/stable/auto_examples/miscellaneous/plot_display_object_visualization.html) example.

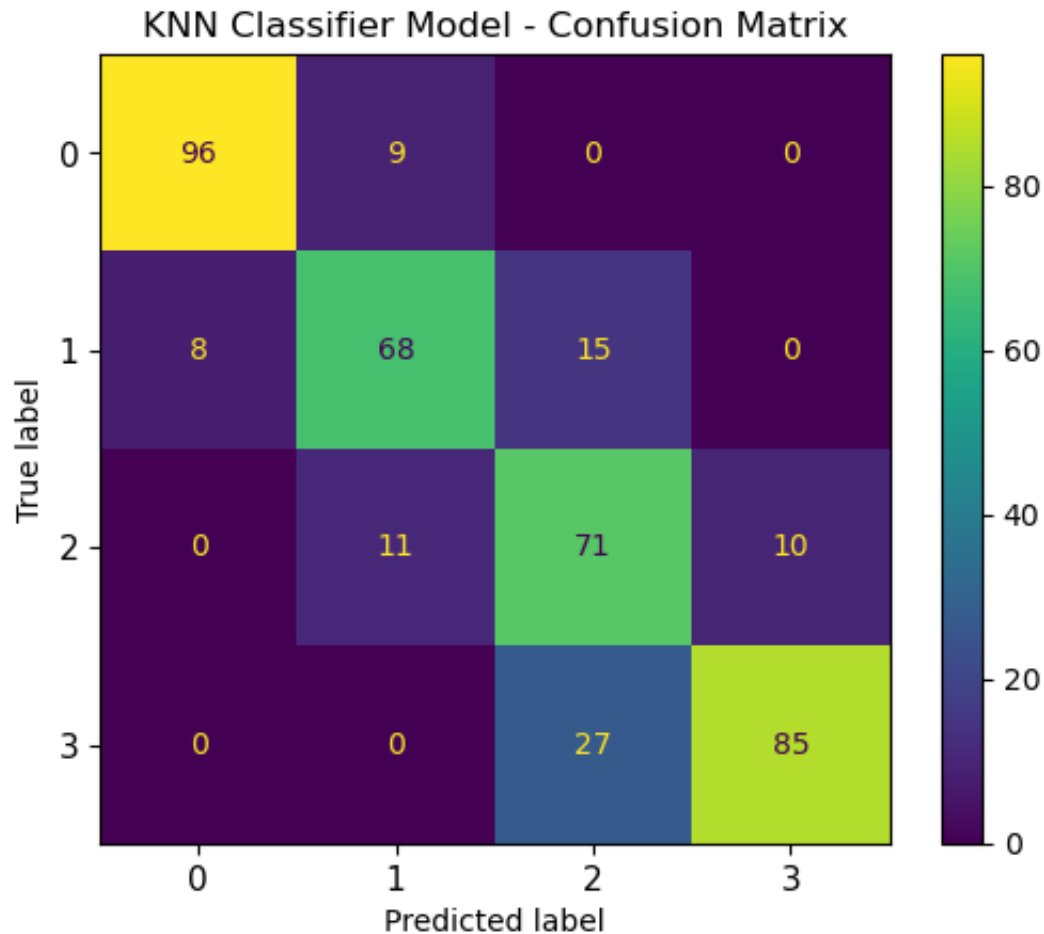
A10 Replace ??? with code in the code cell below

```
In [25]: knn_best = KNeighborsClassifier(n_neighbors=19)

knn_best.fit(xtrain, ytrain)
```

```
Out[25]: KNeighborsClassifier
KNeighborsClassifier(n_neighbors=19)
```

```
In [26]: from sklearn.metrics import ConfusionMatrixDisplay
cm = ConfusionMatrixDisplay.from_estimator(knn_best, xtest, ytest)
#plt.figure()
#plot_confusion_matrix(cm, figsize=(12,8), hide_ticks=True, cmap=plt.cm
plt.title("KNN Classifier Model - Confusion Matrix")
plt.xticks(ticks=np.arange(len(np.unique(ytest))), labels=np.unique(yt
plt.yticks(ticks=np.arange(len(np.unique(ytest))), labels=np.unique(yt
plt.show()
```



Q11 Train KNN classifier using cross-validation approach, [sklearn.cross_validation](https://scikit-learn.org/stable/modules/cross_validation.html) (https://scikit-learn.org/stable/modules/cross_validation.html) tutorial.

Note:

Try a range of `n_neighbors` values (k in kNN) from 1 to 30.

A11 Replace ??? with code in the code cell below **

```
In [27]: # Define KNN model
from sklearn.model_selection import cross_val_score

for k in range(1, 31):
    #Define KNN model
    knn_crossval = KNeighborsClassifier(n_neighbors=k)

    # Use sklearn for 5 fold cross validation
    scores_cv=cross_val_score(knn_crossval, x, y, cv=5)

    # print the scores from different folds
    print(scores_cv)
```

```
[0.7525 0.76    0.7425 0.7575 0.705 ]
[0.77    0.7675 0.755   0.7425 0.715 ]
[0.7775 0.79    0.77    0.7575 0.7525]
[0.7875 0.805   0.7925 0.7775 0.77   ]
[0.7975 0.7825 0.7925 0.785   0.7675]
[0.79    0.81    0.8     0.78   0.775]
[0.7975 0.8     0.7975 0.78    0.7775]
[0.7975 0.8125 0.7975 0.78    0.7625]
[0.7875 0.8     0.8     0.78    0.7675]
[0.8     0.8125 0.8075 0.79    0.7675]
[0.8025 0.7925 0.8     0.805   0.7875]
[0.81    0.795   0.7975 0.79    0.7925]
[0.805   0.7925 0.8025 0.8075 0.7975]
[0.8125 0.805   0.8     0.8175 0.7975]
[0.815   0.805   0.8     0.8275 0.8075]
[0.82    0.8025 0.8     0.8175 0.81   ]
[0.8225 0.8     0.815   0.8175 0.8125]
[0.81    0.81    0.815   0.825   0.8025]
[0.815   0.81    0.8125 0.815   0.815   ]
[0.8175 0.8075 0.815   0.8175 0.8     ]
[0.815   0.8025 0.815   0.815   0.8075]
[0.8175 0.8125 0.8075 0.8075 0.8175]
[0.8175 0.8125 0.8125 0.81    0.815   ]
[0.8125 0.8075 0.8125 0.8125 0.815   ]
[0.8125 0.81    0.815   0.81    0.8225]
[0.815   0.81    0.8175 0.8125 0.81    ]
[0.815   0.81    0.8225 0.8125 0.8225]
[0.82    0.8125 0.825   0.8175 0.8275]
[0.8175 0.815   0.8175 0.82    0.825   ]
[0.8225 0.82    0.81    0.8175 0.825   ]
```

Comparison

Q12 Compare the two models (trained using xtrain,ytrain) in terms of score.

- Train two different models on Train data
- Predict xtest using the trained models
- Make a correlation matrix between ytest and predicted ytest values from the two Models
- Your resulting matrix should be 3x3 correlation matrix for xtest, ytest data
 - The matrix is symmetric
 - It will provide the correlation between two models predictions plus ytest
 - Hint: You can create a new dataframe using these values and use corr() function for creating the correlation matrix. Use meaningful column name while creating the dataframe.

A12 Replace ??? with code in the code cell below

```
In [29]: import matplotlib.pyplot as plt
import seaborn as sns

# Predict Train dataset using logistic reg
clf= LogisticRegression()
clf.fit(xtrain, ytrain)
ypred_clf= clf.predict(xtest)

# Predict Train dataset using KNN
knn= KNeighborsClassifier(n_neighbors=29)
knn.fit(xtrain, ytrain)
ypred_knn= knn.predict(xtest)

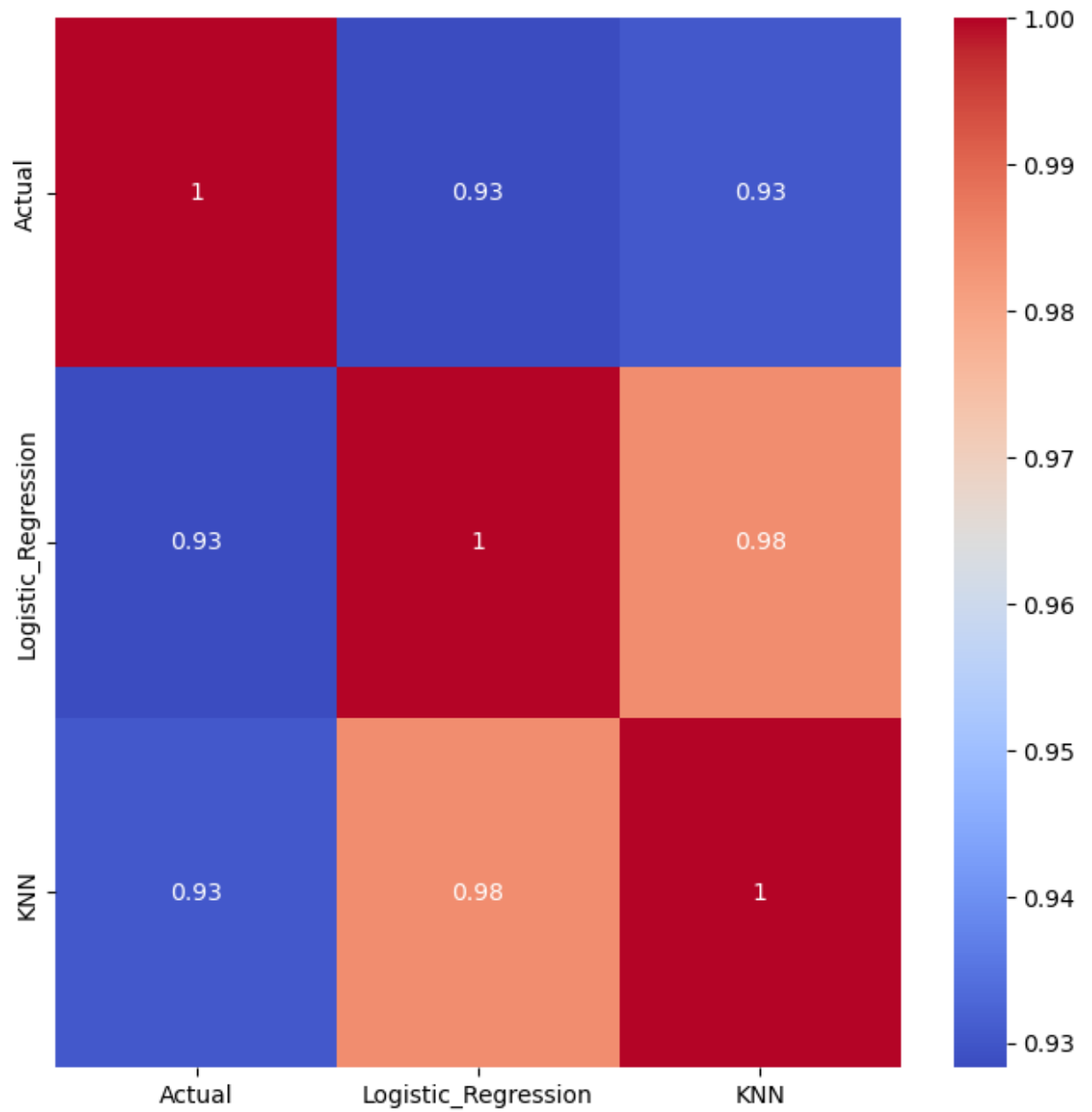
print(ytest.shape, ypred_clf.shape, ypred_knn.shape)
# Create a dataframe using the predicted results from the models
df = pd.DataFrame({
    'Actual': ytest,
    'Logistic_Regression': ypred_clf,
    'KNN': ypred_knn
})

#compute correlation

# Now use seaborn library to plot the heatmap correlation matrix
correlation_matrix = df.corr()
plt.figure(figsize=(8,8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
```

(400,) (400,) (400,)

Out[29]: <Axes: >



In []: