

Final Project

INFO 4330 S50

Group I

Ali Manghat

James Tariga

## Project Topic:

For this Project we decided to search for a data set that had data on a topic we found interesting and implement some of the important concepts we have covered in this course such as supervised learning which includes machine learning models such as linear regression, K-Nearest Neighbors, and logistical regression. We also wanted to test the topic of unsupervised learning which contains the topic of clustering. We found a data set that provides a detailed overview of the recent pandemic and an overview of covid-19 statistics for various different regions around the world. After finding this dataset we performed many of the machine learning techniques we have learned in this class.

## Linear Regression:

Step 1: Importing the dataset in the variable df.

Step 2: We decided to use 'Confirmed', 'Active', 'Recovered / 100 Cases' as the independent variables, and 'Death / 100 cases' as the dependent variable.

Step 3: We split the data set for the first 131 rows for training and the remaining for testing. There are a total of 188 rows.

```
1172] df = pd.read_csv('country_wise_latest.csv')

1173] x = df[['Confirmed', 'Active', 'Recovered / 100 Cases']]
      y = df['Deaths / 100 Cases']

1174] x_train = x[:131]
      x_test = x[131:]

      y_train = y[:131]
      y_test = y[131:]
```

Step 4: We created a linear regression model.

Step 5: We trained the model using the training data.

Step 6: We make predictions on the model using the testing data.

```
linearRegModel = LinearRegression()
linearRegModel.fit(x_train, y_train)
```

[5] ✓ 0.0s

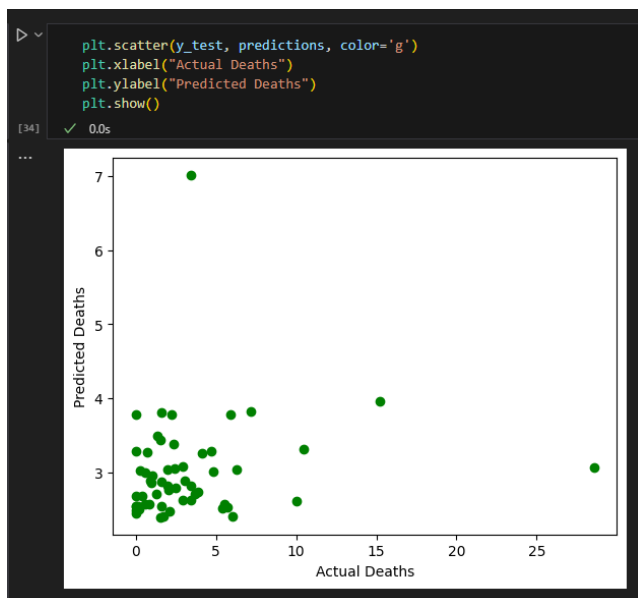
▼ LinearRegression

LinearRegression()

```
predictions = linearRegModel.predict(x_test)
```

[6] ✓ 0.0s

Step 7: Creating a scatter plot to show the actual and predicted values of the deaths / 100 cases.



Step 8: Calculating MSE. The MSE is 19.78 which is relatively low. This means that the model is pretty accurate. The predicted values are like the actual values. We can see the max value for

the dependent variable 'Deaths / 100 cases' is 28.56 and the lowest is 0. This further proves that the MSE is relatively accurate.

```
mse = mean_squared_error(y_test, predictions)
print(mse)
35] ✓ 0.0s
.. 19.77795136072621

max_value = df['Deaths / 100 Cases'].max()
print("Max Value: ", max_value)

min_value = df['Deaths / 100 Cases'].min()
print("Min Value: ", min_value)
36] ✓ 0.0s
.. Max Value: 28.56
   Min Value: 0.0
```

## Clustering:

Step 1: We decided to make two clustering models. The first with 2 clusters, and the second with 3 clusters. W

```
km1 = KMeans(n_clusters=2, n_init=20)
km1.fit(x_train)
1180]

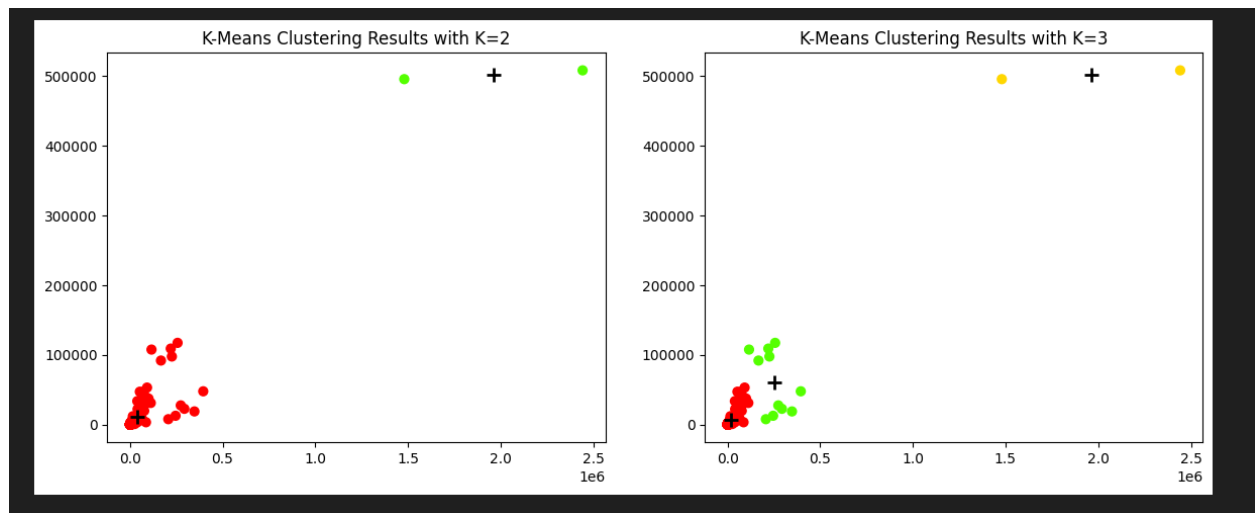
... KMeans
   KMeans(n_clusters=2, n_init=20)

km2 = KMeans(n_clusters=3, n_init=20)
km2.fit(x_train)
1181]

... KMeans
   KMeans(n_clusters=3, n_init=20)
```

Step 2: Plotting the clustered data. Here we can see the two models plotted. The graph on the left represents the first model which is grouping the data into two clusters and on the right the model is

grouping the data in three clusters. We can see covid-19 cases/deaths grouped by different regions. First world countries have a much lower number of deaths, the number increases as the countries become less developed. The centroids have also been plotted to show the 'mean' for each of the clusters.



### Comparison:

Both Linear regression and clustering have their use cases. Choosing which model to use comes down to what you are trying to predict. In our scenario linear regression makes sense because we are trying to map the features confirmed cases, recovered cases, and active cases to the dependent variable deaths. With clustering we are able to see what regions had the highest death rates but it does not provide us with a way to predict the deaths based on the features we are using. Therefore, linear regression is more appropriate for this specific scenario but clustering is also useful if you want to break down the dataset into clusters to visualize similarities in the data.

## KNN and Logistic Regression:

Step 1: Initialized a feature called xKNN and a target called yKNN to train the model. Then the models are trained using the fit method on the provided training data.

```
xKNN = df[['Confirmed', 'Active', 'Recovered / 100 Cases', 'Deaths / 100 Cases']]
yKNN = df['Americas']

logregmodel = LogisticRegression()
knn = KNeighborsClassifier(5)

logregmodel.fit(xKNN, yKNN)
knn.fit(xKNN, yKNN)
```

Step 2: We then make predictions on the same training data using the trained models.

```
logmodelpreds = logregmodel.predict(xKNN)
knnmodelpreds = knn.predict(xKNN)
```

Step 3: Confusion matrices are computed for both models using the true labels (yKNN) and predicted labels (logmodelpreds, knnmodelpreds).

```
logconf = confusion_matrix(yKNN, logmodelpreds, labels=knn.classes_)
print("Logistical Regression: \n", logconf)

KNNconf = confusion_matrix(yKNN, knnmodelpreds, labels=knn.classes_)
print("KNN: \n", KNNconf)
```

Logistical Regression:

```
[[148  4]
 [ 32  3]]
```

KNN:

```
[[148  4]
 [ 27  8]]
```

Step 4: The accuracy scores are calculated for both models using the true labels and predicted labels.

```
logacc = accuracy_score(yKNN, logmodelpreds)
print("Logistical Regression: ", logacc)

KNNacc = accuracy_score(yKNN, knnmodelpreds)
print("KNN: ", KNNacc)
```

[126] ✓ 0.0s

... Logistical Regression: 0.8074866310160428  
KNN: 0.8342245989304813

Step 5: Precision, recall, and F1-score are computed using the `precision_recall_fscore_support` function for both models, considering a binary classification scenario.

This essentially trains two classifiers, makes predictions, evaluates their performance using confusion matrices, accuracy scores, and precision-recall-f1 metrics. It provides a comprehensive overview of how well each model is performing on the given dataset.

```
precision_recall_fscore_support(yKNN, logmodelpreds, average='binary')
```

[127] ✓ 0.0s

.. (0.42857142857142855, 0.08571428571428572, 0.14285714285714285, None)

```
precision_recall_fscore_support(yKNN, knnmodelpreds, average='binary')
```

[128] ✓ 0.0s

.. (0.6666666666666666, 0.22857142857142856, 0.3404255319148936, None)

## Comparison:

Based off the output below, we can conclude that the Logistic Regression model correctly predicts the majority class ('not Americas' in this case) well but struggled with the minority class ('Americas').

**True Negative: 148**

**False Positive: 4**

**False Negative: 32**

**True Positive: 3**

**KNN:**

On the other hand, the KNN model shows a slight improvement over Logistic Regression in predicting the minority class, with more true positives and fewer false negatives.

**True Negative: 148**

**False Positive: 4**

**False Negative: 27**

**True Positive: 8**

These values indicate the overall proportion of correct predictions made by each model. KNN performs slightly better in this regard.

**Logistical Regression: 0.8074866310160428**

**KNN: 0.8342245989304813**

Both models are more inclined towards predicting the majority class correctly while struggling with the minority class. We believe the reason as to why this occurs is because our dataset is imbalanced, and this is a common issue with imbalanced datasets. To improve model performance, techniques like resampling the dataset, trying different classification algorithms, or using different feature selection methods might alter the result. Additionally, adjusting the decision threshold and employing methods to address class imbalance can be considered to improve the recall without significantly compromising precision.



**Citation of data set:**

Devakumar, K.P. (n.d.). Corona Virus Report. Kaggle. Retrieved from

<https://www.kaggle.com/datasets/imdevskp/corona-virus-report>