# DATA CLUSTTERING AND FITTING FOR DATASETS

## INTRODUCTION

Through this poster we are going through various levels of statistical study using datasets taken from Worldbank data.
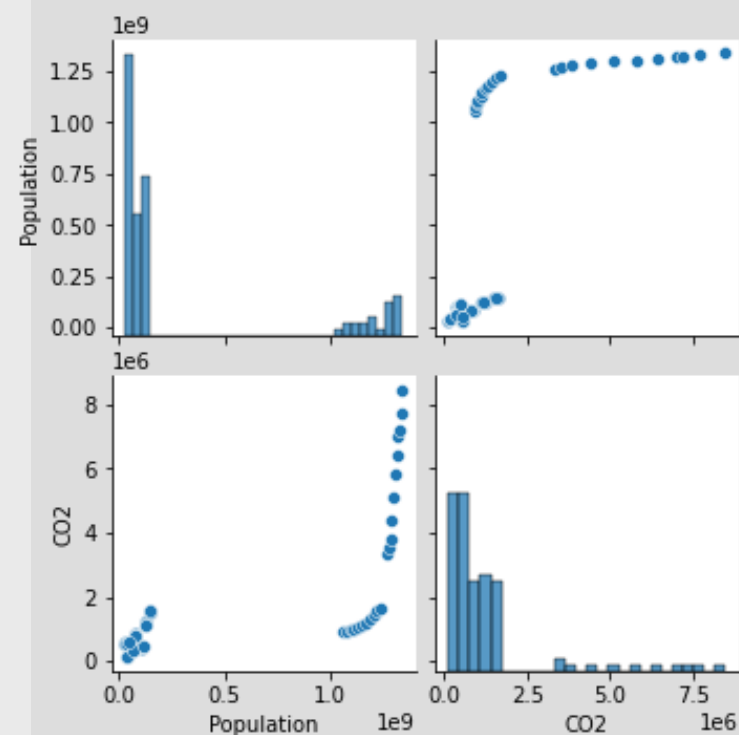At first we have taken the data regarding the total CO2 emissions from a period of 2000-2010 for 10 countries.



On the other hand we have to choose the dataset which contains the data of total population over the period 2000-2010 for 10 countries.
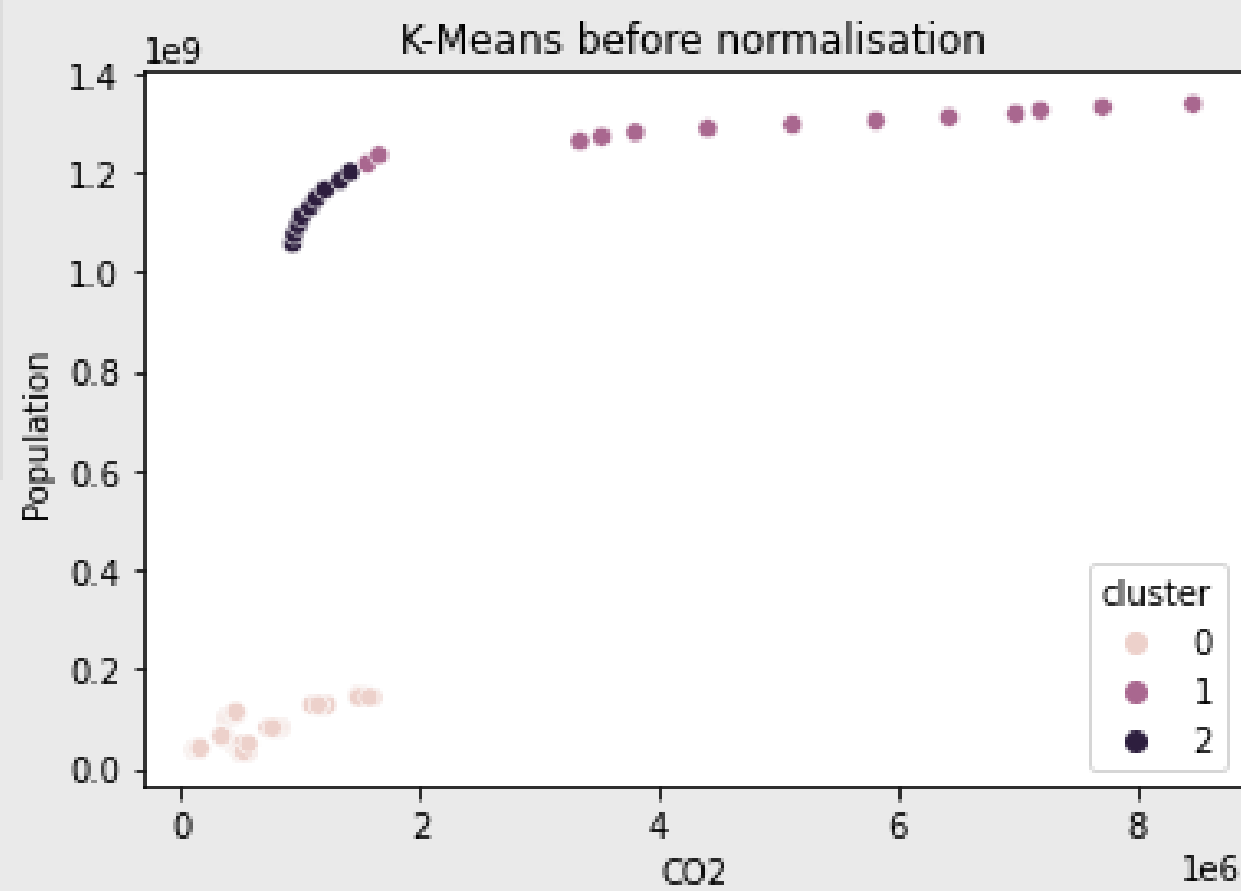


We have simply plotted a pair plot comparing both datasets and given as follow:



Primarly,we are looking onto the clustering using K means.Before that we have to know what is clusttering and what is the use of clustering in statistics.

- **process of making a group of abstract objects into classes of similar objects**.
- **Collection of an array of data points with similar traits and arrange them in a same cluster**.
- Here we have done the clusttering by K means.
- We have imported the libraries and done the whole clusttering process.
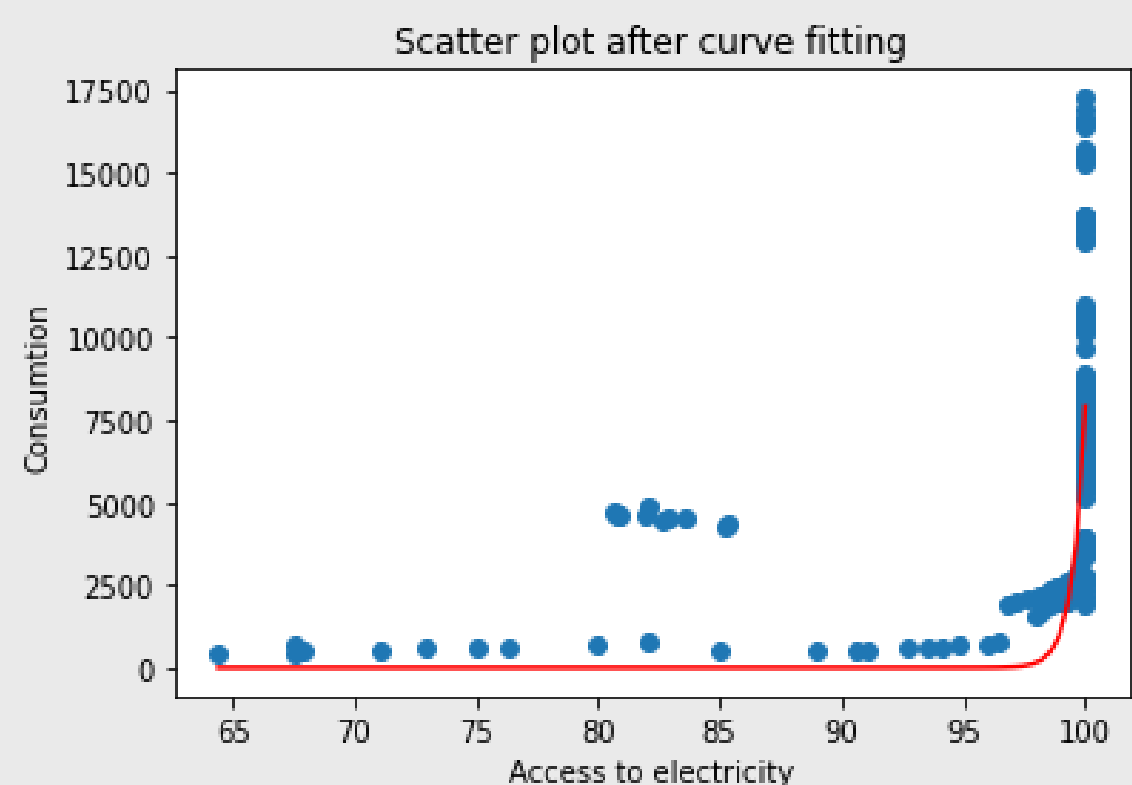- The K means clustering graph before and normalization is given below for further study:





Another interesting feature is data curve fitting.
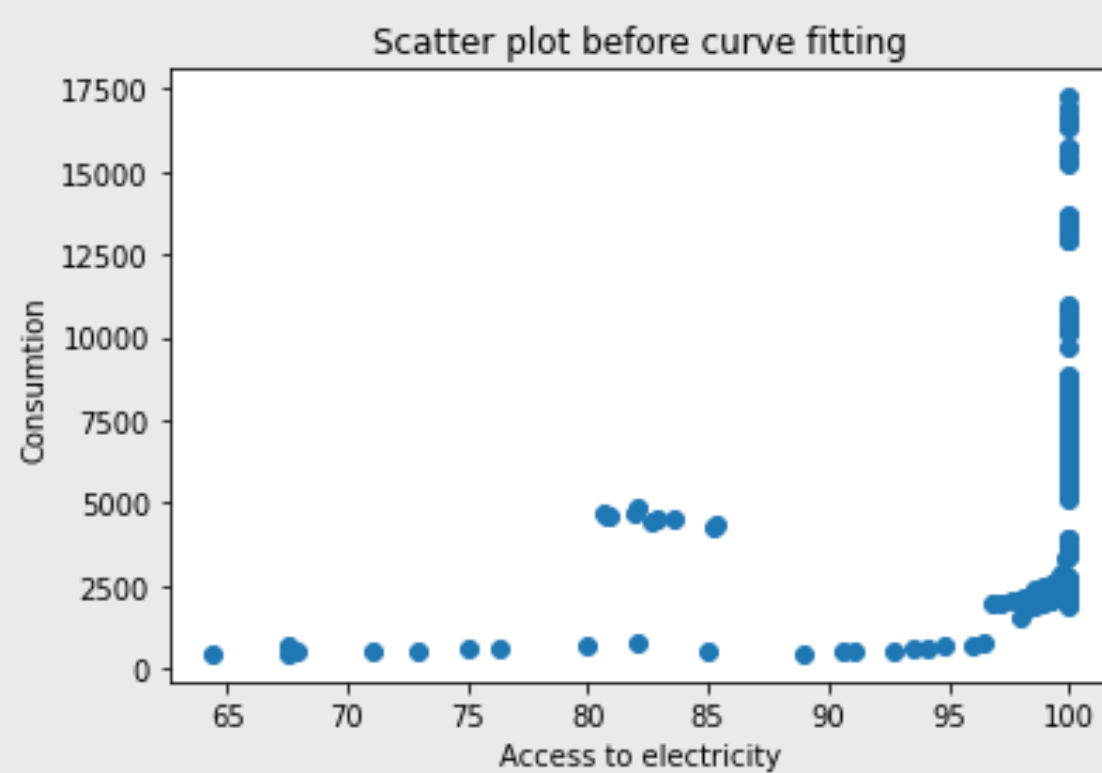What is curve fitting?
- A process of finding a function that fits perfectly for a data.
- For this fitting we have selected two different datasets.
- One having data containing the total access to electricity and another having data regarding the total energy consumption.
- Clustered data by both the datasets are tabled below:

| | consumption | access |
|---|---|---|
| 0 | 95.783287 | 2088.807630 |
| 1 | 100.000000 | 17037.072281 |
| 2 | 97.021797 | 992.943385 |
| 3 | 100.000000 | 6635.421406 |
| 4 | 100.000000 | 7224.526985 |
| ... | | |
| 105 | 76.300003 | 640.394607 |
| 106 | 100.000000 | 8594.909034 |
| 107 | 100.000000 | 9716.126081 |
| 108 | 99.236694 | 2018.827437 |
| 109 | 100.000000 | 6409.894365 |

110 rows × 2 columns

The graphs formed shown below:



The process of finding the curve fitting is mentioned below:
- The data is selected for 10 countries.
- Plotted a scattered plot.
- An exponential function generated and runned.
- Thus the curve is plotted in the scattered plot.



## CONCLUSION

Thus we have gone through total four types of data for total clustering and fitting.The findings are plotted and shown. A research regarding these findings are done and studied.

## ABSTRACT

- A study regarding two datasets taken in a certain period for selected countries to find and plot clustering
- A study for curve fitting is done using two datasets in a certain period for the countries selected to plot the curve in the data fitted.

## CONTACT

SREEJITH ANILKUMAR
Student ID: 21032113