# Capstone Project - 3

## Online Retail Customer Segmentation

## ML Unsupervised Clustering

**Individual Project:**
**Soumya Ranjan Mishra**

# Contents

- Problem Statement
- Data Summary
- Data Analysis
- Analysis Details
- RFM Table for Customer ID
- Data Preparation
- Implementing Various Clustering Algorithms
- RFM Table for Cluster ID
- Challenges
- Conclusions

# Problem Statements

- **Identify major customer segments on UK Based online retail dataset.**
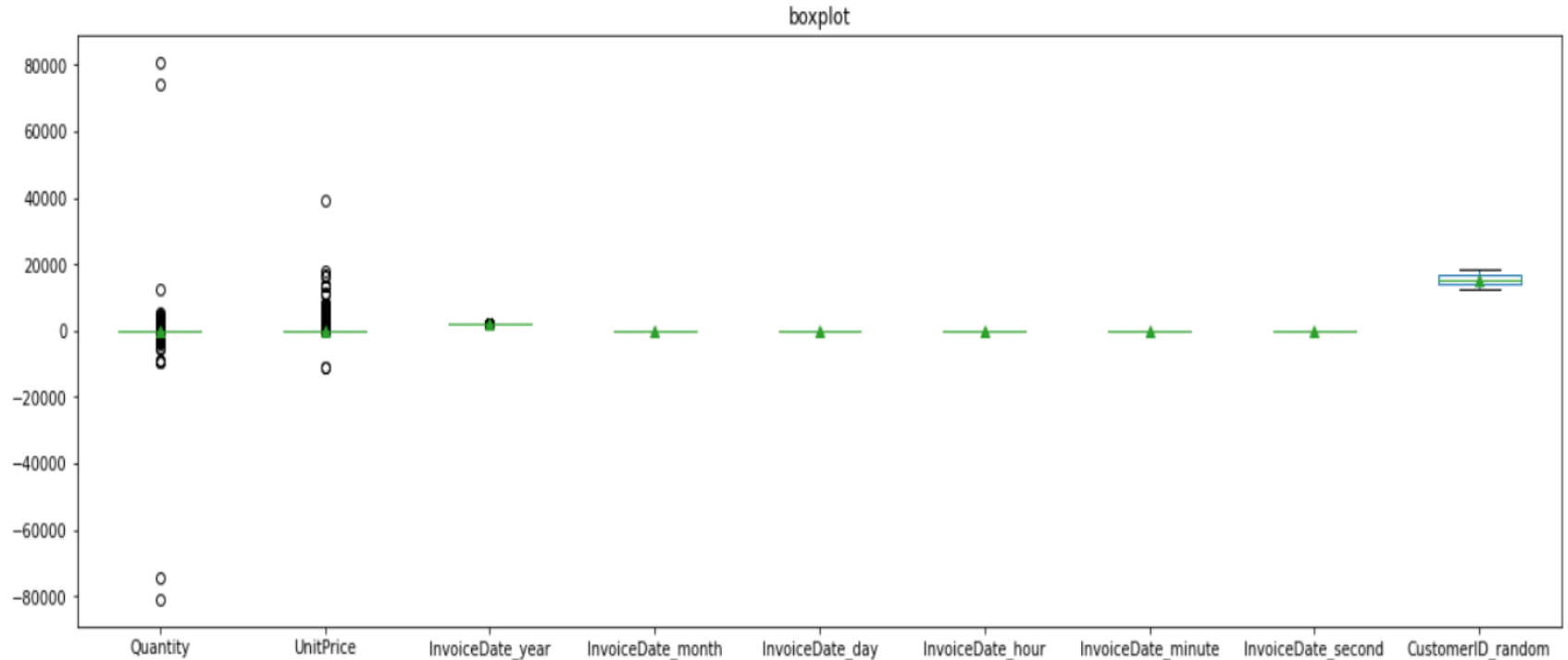- **Create RFM Table.**

# Data Summary

- **InvoiceNo:** Invoice number. Nominal, a 6-digit integral number uniquely assigned to each transaction. If this code starts with letter 'c', it indicates a cancellation.
- **StockCode:** Product (item) code. Nominal, a 5-digit integral number uniquely assigned to each distinct product.
- **Description:** Product (item) name. Nominal.
- **Quantity:** The quantities of each product (item) per transaction. Numeric.
- **InvoiceDate:** Invoice Date and time. Numeric, the day and time when each transaction was generated.
- **UnitPrice:** Unit price. Numeric, Product price per unit in sterling.
- **CustomerID:** Customer number. Nominal, a 5-digit integral number uniquely assigned to each customer.
- **Country:** Country name. Nominal, the name of the country where each customer resides.
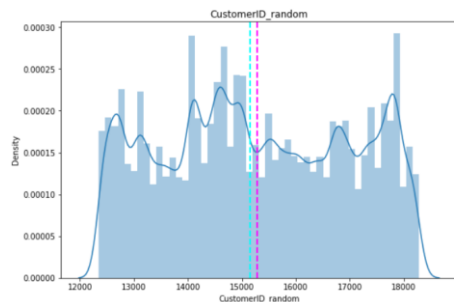
4

# Basic Data Exploration

**AI**

● **The dataset has** *541909* **rows and 8 features(columns).**

●**Four categorical features 'InvoiceNo', 'StockCode', &**

**'Description', ' Country'.**

● **One Datetime[ns] features 'InvoiceDate'.**

● **Outliers present only in "Quantity" & "UnitPrice"column.**

● **Missing Values on** *Description & CustomerID columns***.**
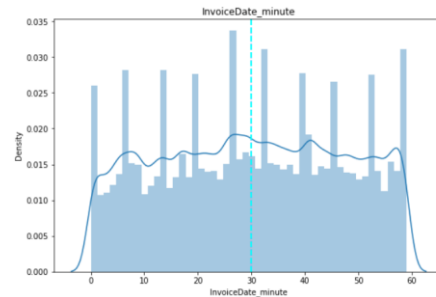
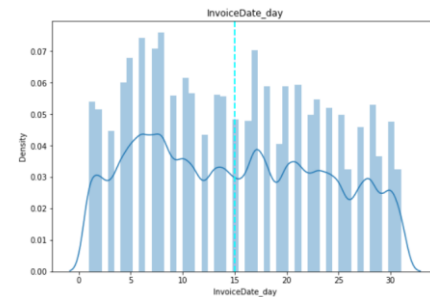● **Duplicated values present.**

# Outliers in the features



boxplot

# Mean Distribution of Features
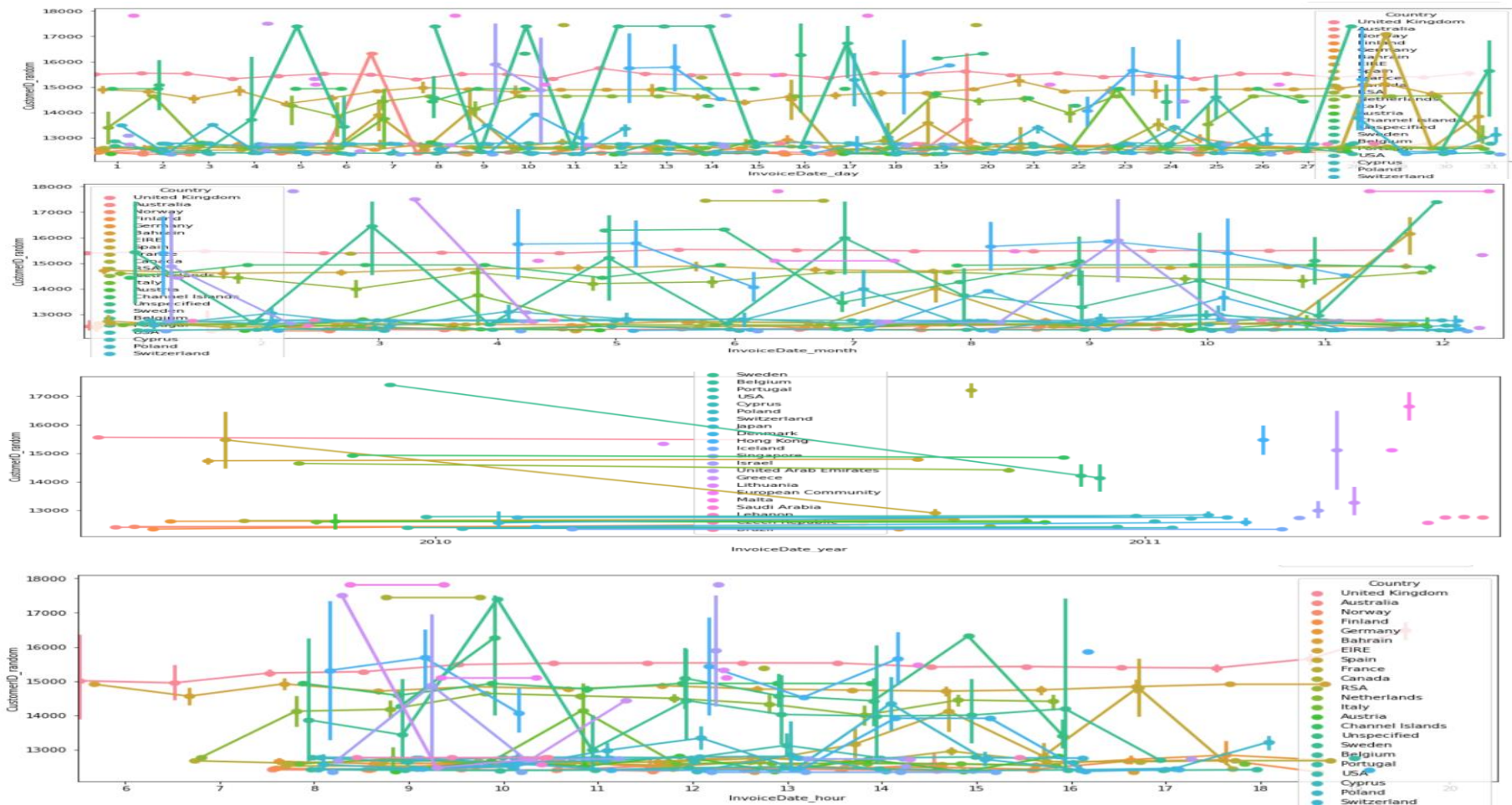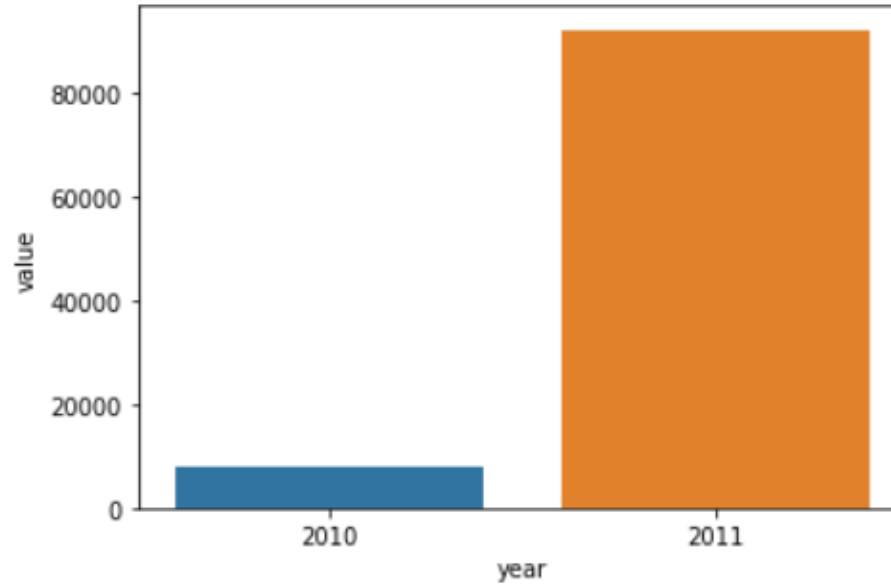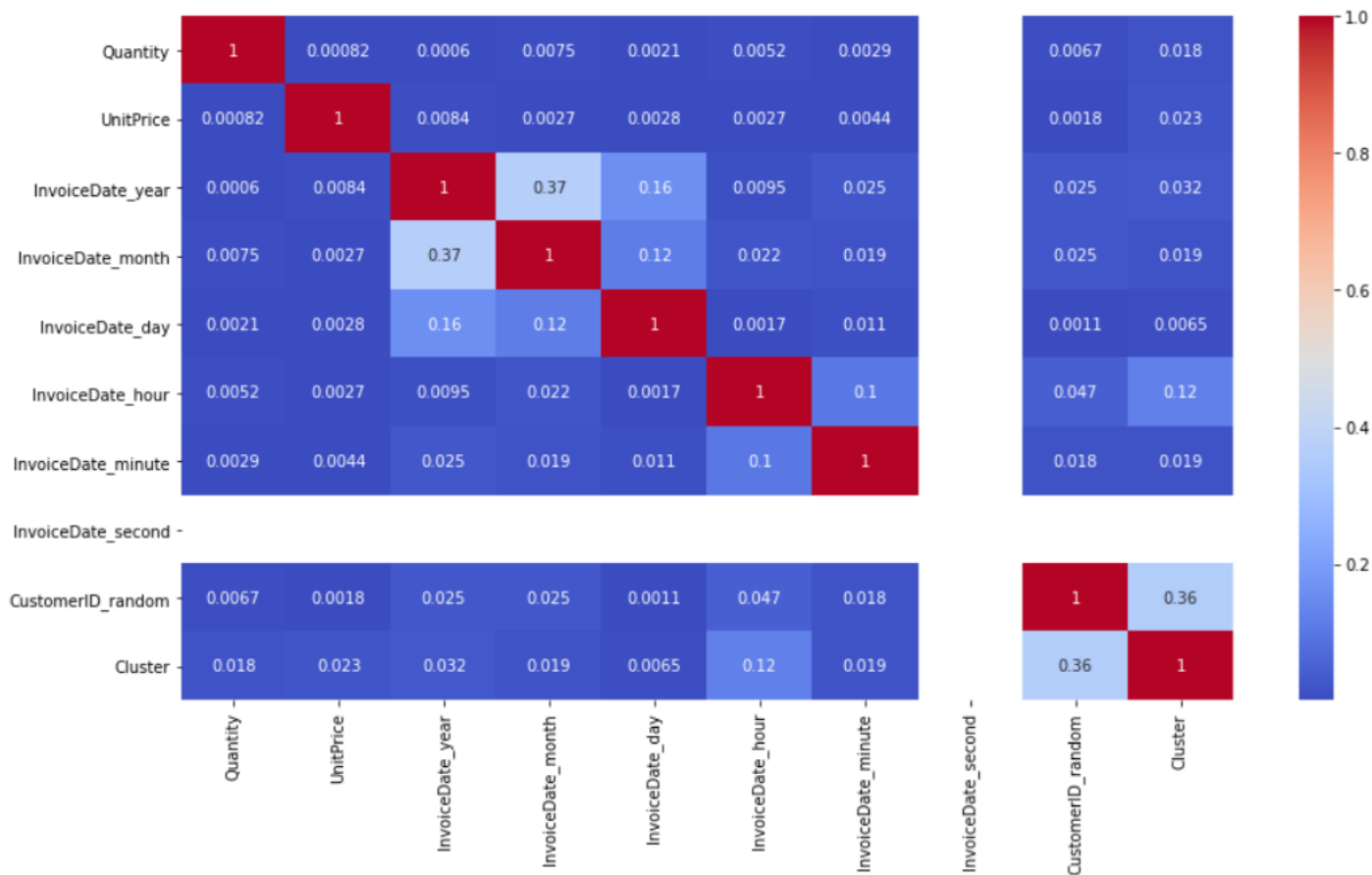
# Spread over time

# Spread over time and Country
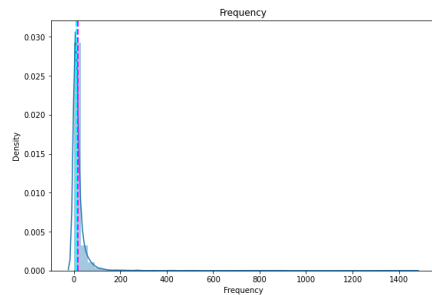


9

# Distribution of Number of reviews
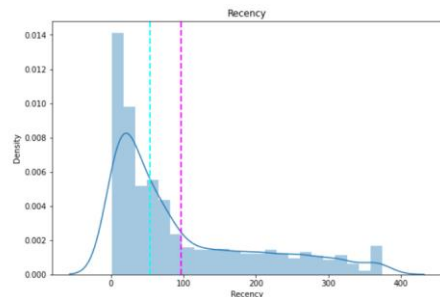
# Correlation Matrix

# RFM Table for Customer ID

- R (Recency): Number of days since last purchase
- F (Frequency): Number of tracsactions
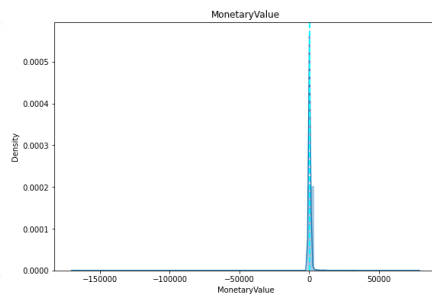- M (Monetary): Total amount of transactions (revenue contributed)

# Data Preparation

- **One Hot Encoding**

- **Outlier Treatment**

- **Standard Scaler Scaling**

- **Principal Component Analysis (** n_components = 4 **)**



(-7.028006510533737, 9.383883035465741, -39.92151455621802, 43.0757823245935)

# K Means Clustering

For n_clusters = 3 The average silhouette_score is : 0.6194838287845018

**Hyper parameters**
{n_clusters=3,
max_iter=1000,
random_state=10}

# RFM For Cluster

# K-Means Clustering with Silhouette

For n_clusters = 2 The average silhouette_score is : 0.611432364435861
For n_clusters = 3 The average silhouette_score is : 0.6194838287845018
For n_clusters = 4 The average silhouette_score is : 0.3372942137064119
For n_clusters = 5 The average silhouette_score is : 0.2748597164906843
For n_clusters = 7 The average silhouette_score is : 0.2809492799812412
For n_clusters = 8 The average silhouette_score is : 0.25979222243476296
For n_clusters = 10 The average silhouette_score is : 0.22728886796555026

## Hyper parameter
{ n_clusters=[2,3,4,5,7,8,10],
max_iter=1000,
random_state=10 }



16

# K-Means Clustering with Elbow method



The Elbow Method



KElbowVisualizer(ax=<matplotlib.axes._subplots.AxesSubplot object at 0x7f09590699d0>,
        k=None, metric=None, model=None, timings=False)

## Hyper parameter
{ n_clusters=[1,10],
 init='k-means++',
random_state=0}

# Hierarchical Clustering

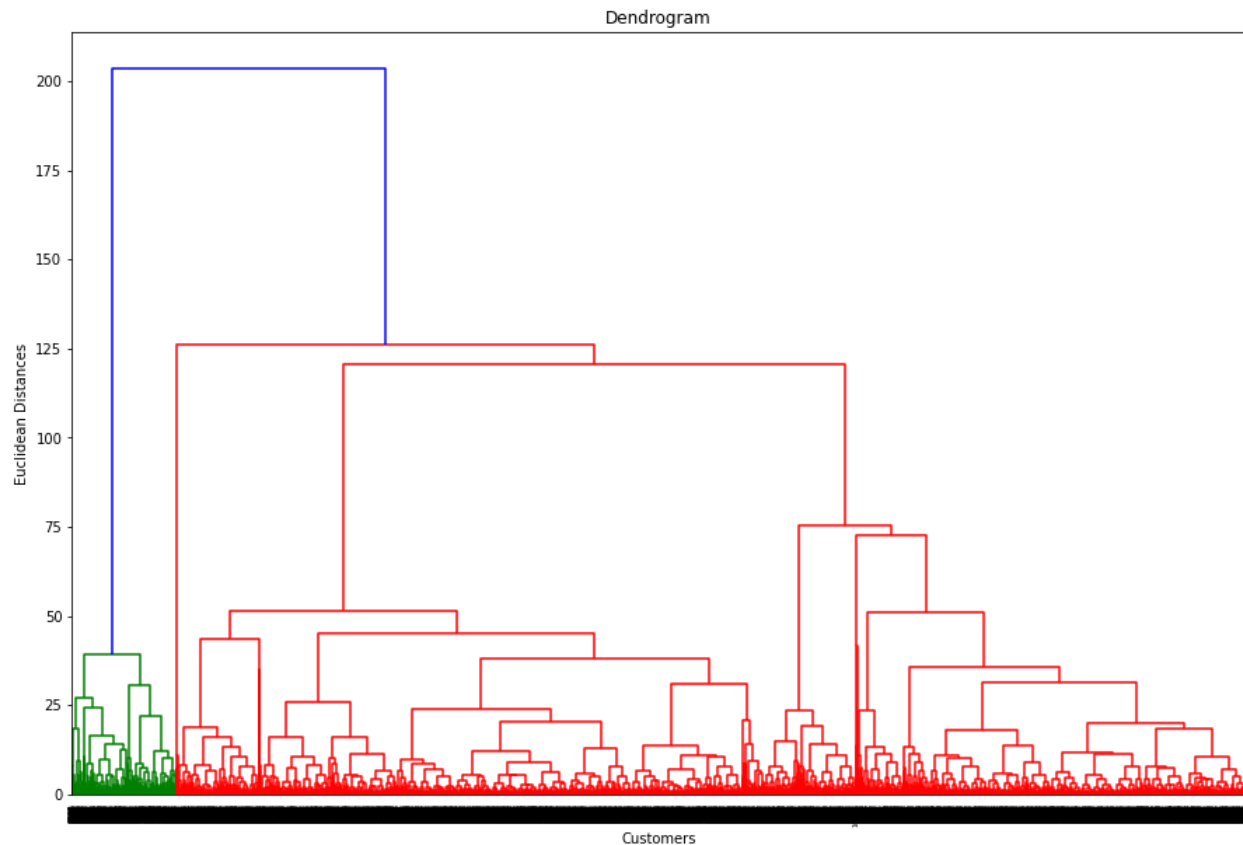**Hyper parameter**

AgglomerativeClustering
{ n_clusters = 3,
affinity = 'euclidean',
linkage = 'ward'}

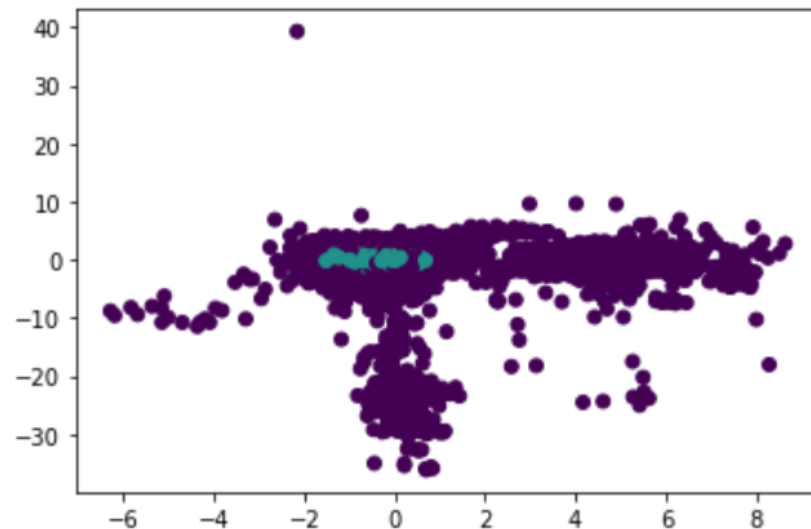# Density-Based Spatial Clustering Of Applications With Noise (DBSCAN)

## Hyper parameter

{ eps=0.3,
min_samples=100 }

```
Estimated number of clusters: 2
Estimated number of noise points: 39095
Homogeneity: 0.293
Completeness: 0.125
V-measure: 0.176
Adjusted Rand Index: 0.117
Adjusted Mutual Information: 0.176
Silhouette Coefficient: 0.148
```



Estimated number of clusters: 2

# Challenges

- **Large Dataset to handle.**

- **Needs to plot lot of Graphs to analyse.**

- **Lot of NaN values.**

- **Continuous Runtime and RAM Crash due to large dataset.**

- **Carefully tuned Hyper parameters .**

# Conclusion

- **K-Means Clustering with Silhouette gives the highest score of 61.9% for number of clusters 3.**

- **Sales has been increased from 2010 to 2011.**

- **RFM for Cluster ID box plots tells well about Cluster detail.**

- **We can deploy this model.**

# THANK YOU

## Q & A