

# **CUSTOMER LIFE TIME VALUE (CLV) ANALYSIS**

**HOPMONK Gaming Company**

1



• Problem Statement

• Data Pre Processing and Feature Engineering

• Model Building

• Error Metric Evaluation

• Data Insights

• Summary

# AGENDA

# OVER VIEW

## Objective

Using machine learning techniques analyze the data and predict the **Customer Life Time Value CLV** that will enable Hopmonk to target and acquire customers based on the net potential as profit.

## Hopmonk Storage and the Approach

- Hopmonk currently has a million customers on their Enterprise Data Warehouse (EDW) where the data is spread across various tables.
- Hopmonk has implemented both Oracle and Cognos solutions designed to enable their business users to extract and analyze their data.
- Hopmonk feels they have locked away valuable details on consumer behavior, segmentation, demographics, and more.

# DATA PRE-PROCESSING AND FEATURE ENGINEERING

## Data Pre-Processing

7 Data frames are merged together based on the business problem and a final master data frame is made.

Based on the Variables business importance the missing NA values are been imputed.

## Feature Engineering

A new Features are been extracted based on the existing variables.

New Features are:

Frequency, Recency, Average Order Value, Purchase-frequency and Customer life time value.

Frequency=Sum(App+LF)

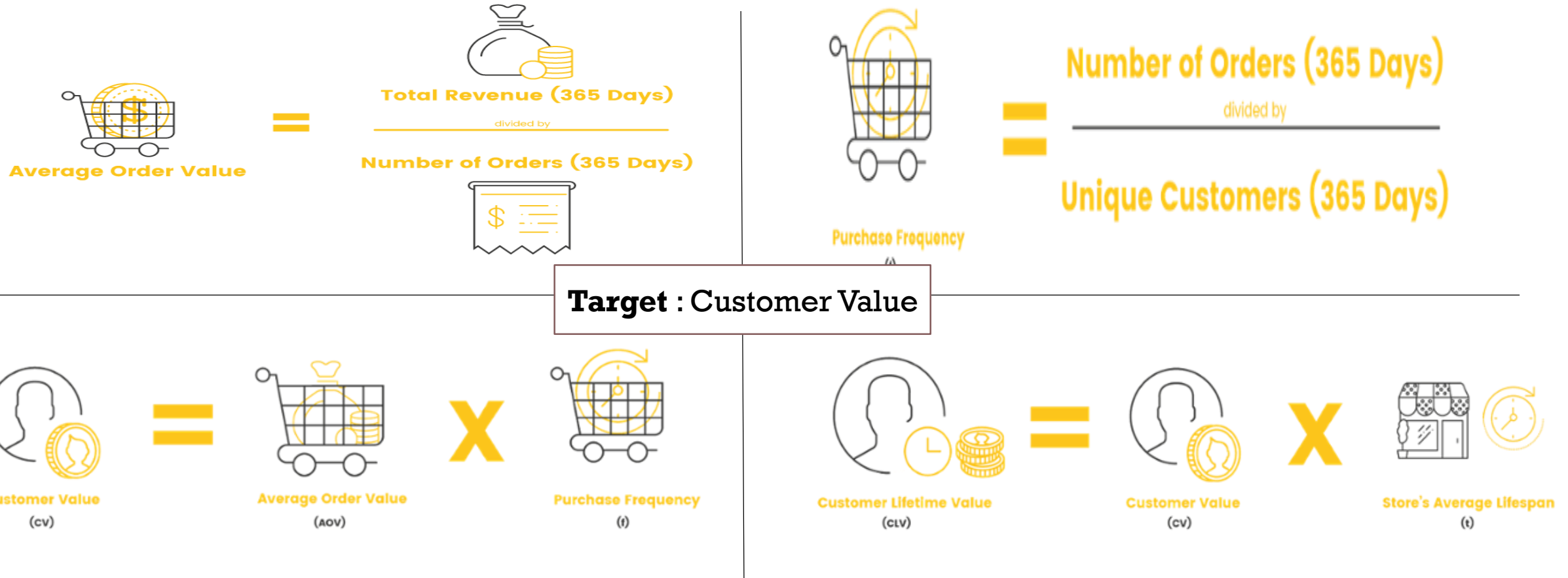
Recency=Recent value(App and LF)

## Variables Extraction

By using Hetcor function the variables which are highly corelated are been removed.

And using the Step Aic and VIF some of the variables are been removed.

# FEATURE ENGINEERING



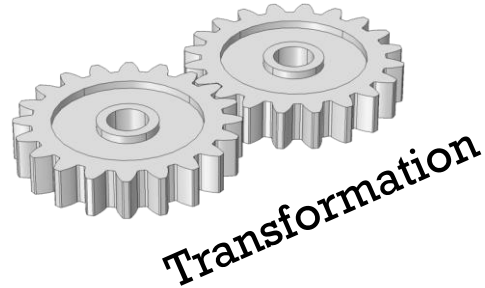
**Variables used:**  
 Total Revenue Generated  
 Units bought  
 No. of Unique Customer.

# DEMO GRAPHS OF THE FINAL DATA FRAME

Pre-



l Data Frame With  
0 by 103



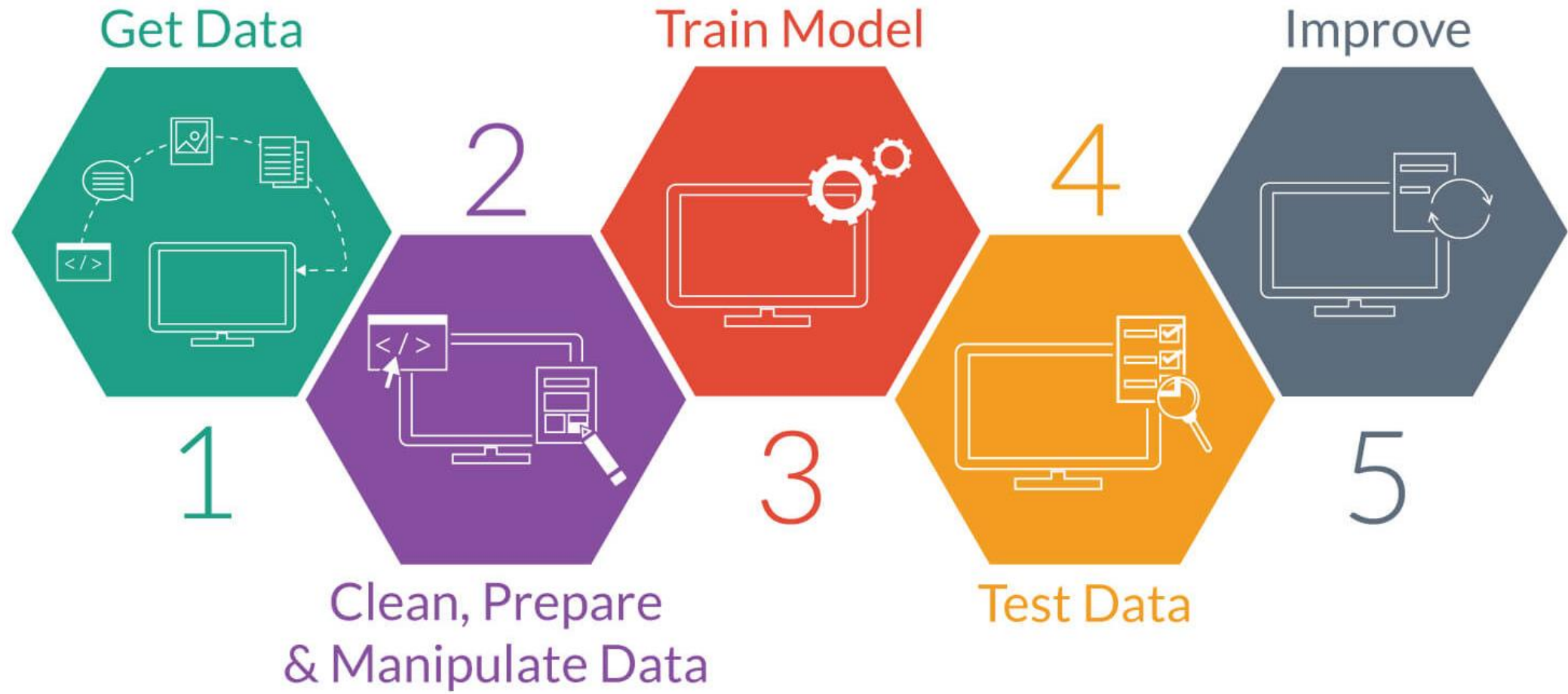
Feature engineered Data Frame With  
29271 by 83

| Customer Centric   | Business Centric  | Customer Demo graphs |
|--------------------|-------------------|----------------------|
| Games Played,      | Favourite Source  | Age                  |
| Time               | Favourite Channel | Country              |
| Units              |                   | No of Children       |
| Recency, Frequency |                   |                      |

## Hopmonk Dataset Summary:

**Total Observations** : 29271  
**Total variables** : 24  
**Split** : 70-30  
**Target Variable** : Customer\_value  
**Scaling** : Based on the model

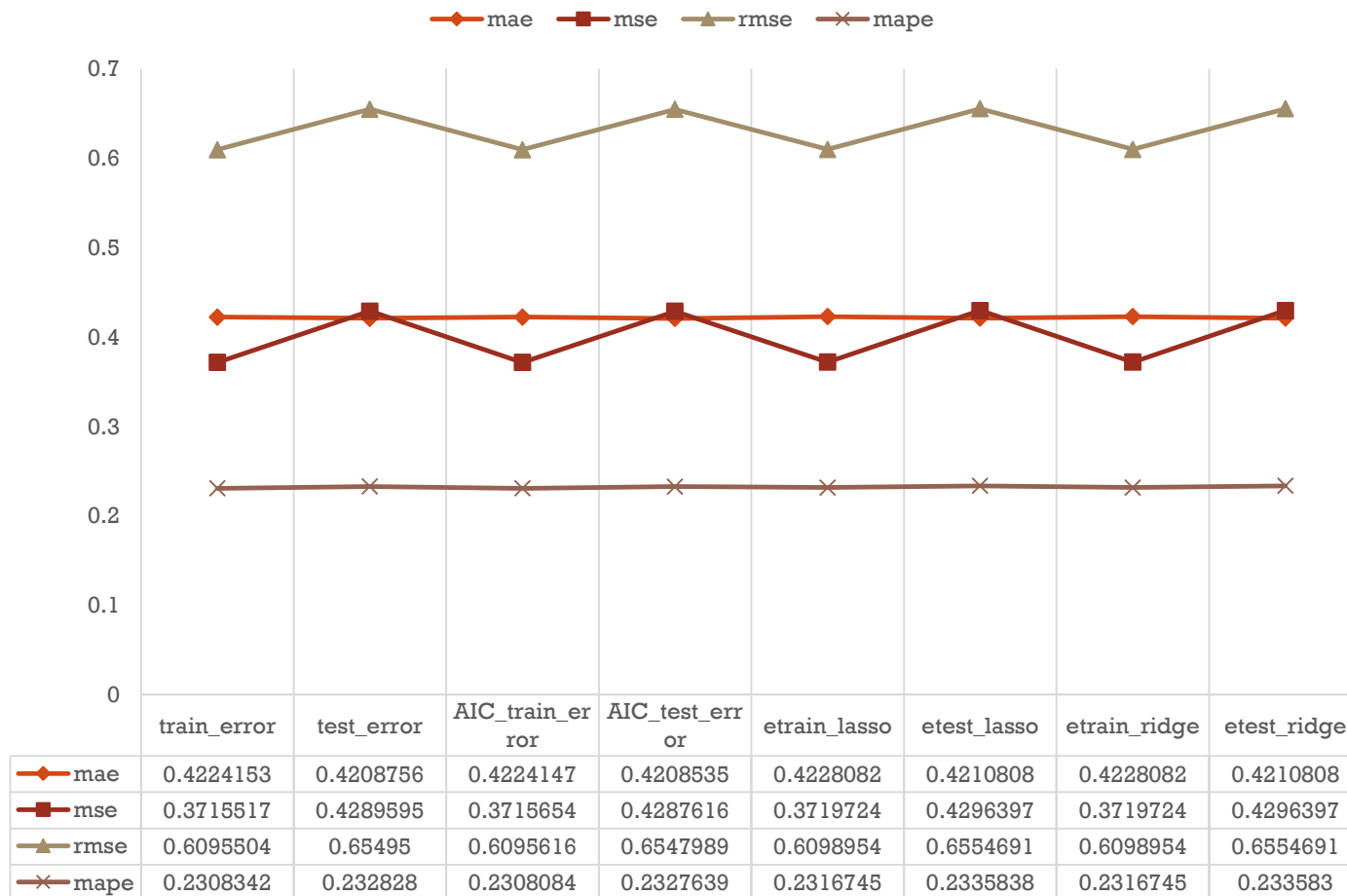
# MODEL BUILDING





# SIMPLE LINEAR REGRESSION

SIMPLE REGRESSION ERRORS



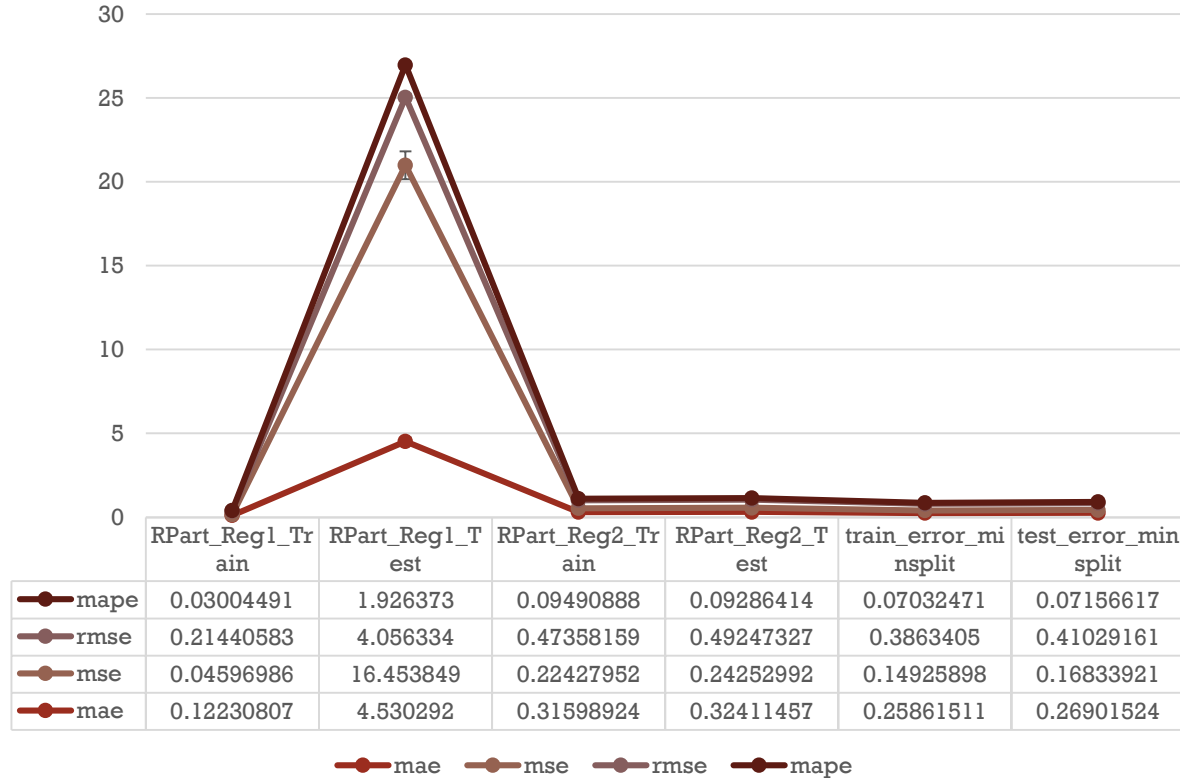
## Simple Linear Regression:

- The errors were computed based on the simple linear regression.
- MAPE stood almost same for various types of regression models.
- The Adjusted R2 is Around 0.90 so even the model is good.
- There is some non-linearity in the data as the data is heavily skewed.
- Did the transformation on the target variable such as  $\sinh(x)$ ,  $\log(x)$ ,  $x^2$  and  $\sqrt{x}$ .
- The model gave optimum values over the transformation of  $\sqrt{x}$  over other transformation.



# DECISION TREE

Decision Tree Errors

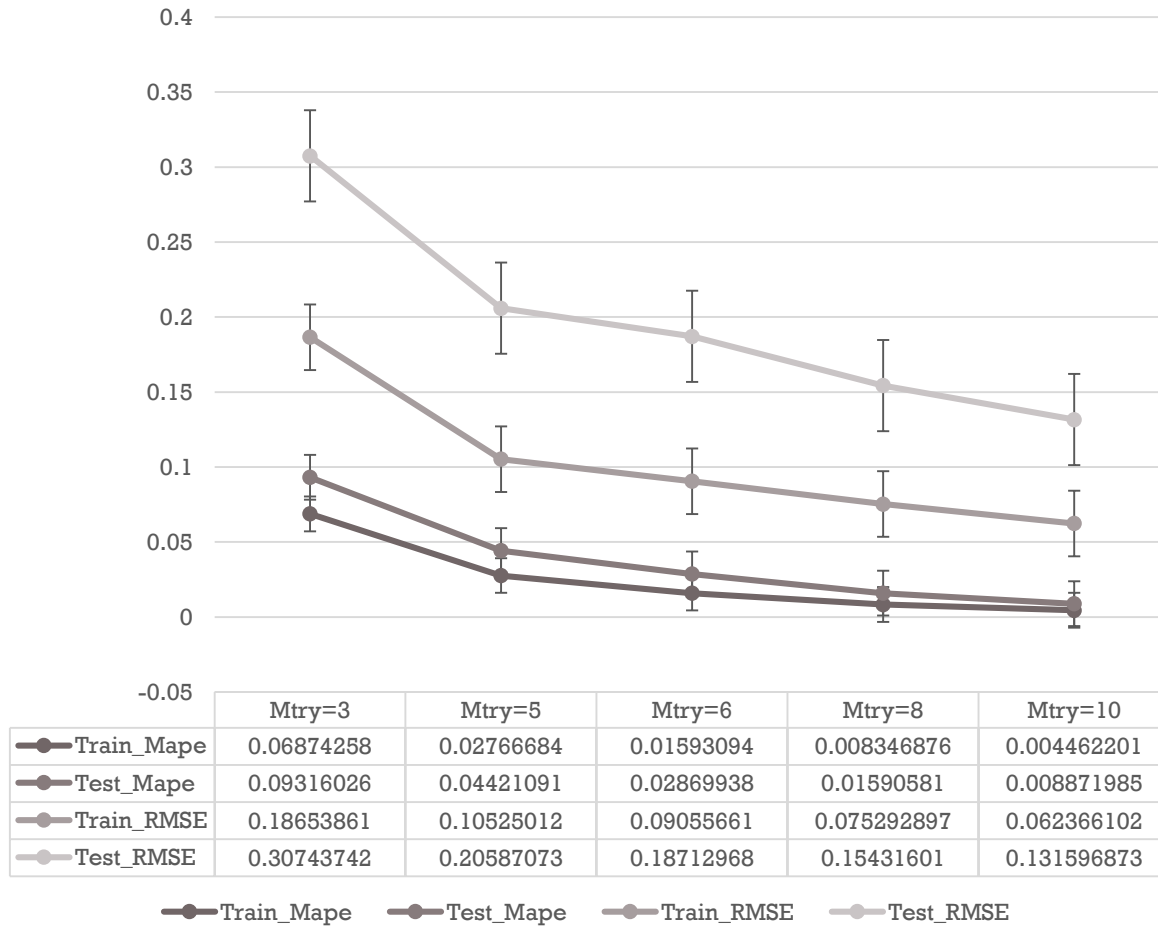


## Decision Tree Model:

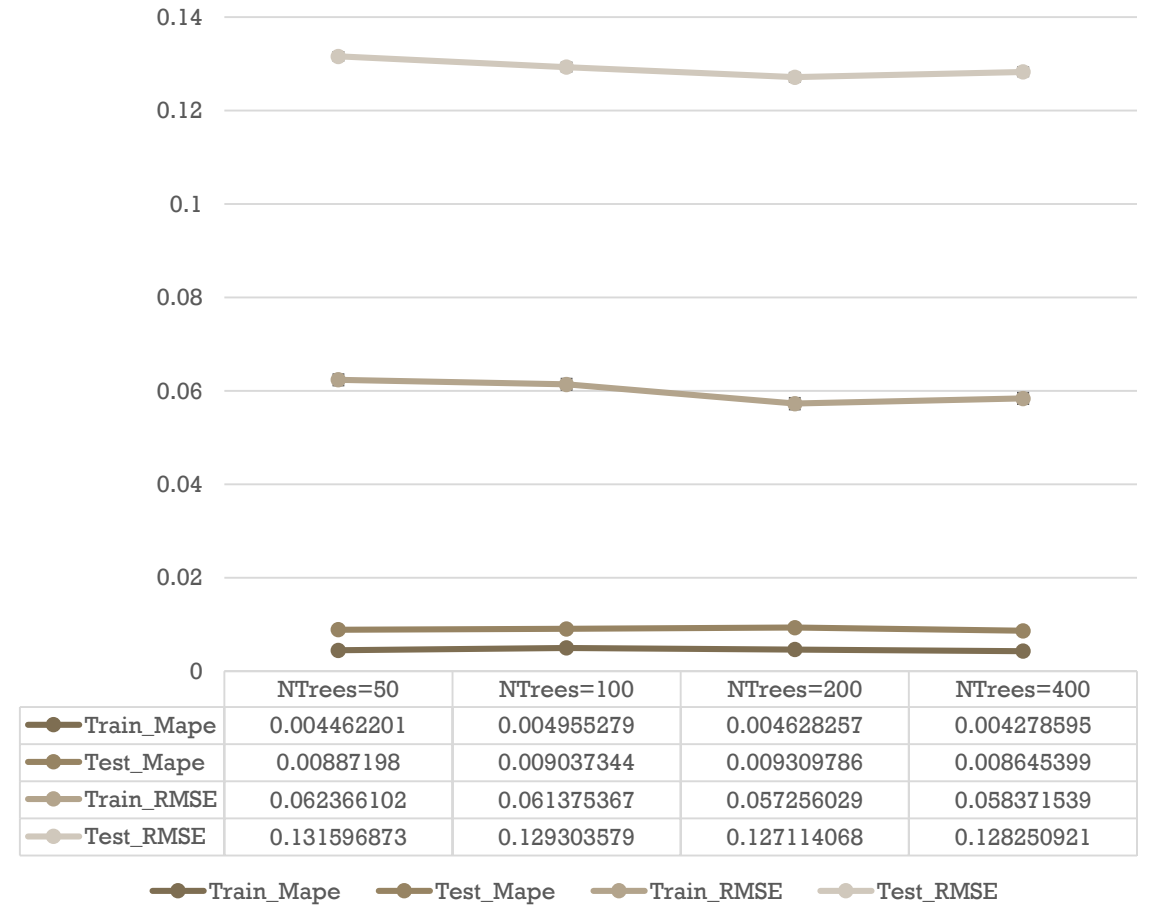
- The tree was made to grow without any constraints due to which the model got overfit.
- Then went for the optimal cp value based on the x–relative error.
- Then tuning the number of minimum split gave the least rmse and mape and stable over both train and test.
- As decision tree uses works with the greedy approach some sort of information might be left to get so going for the other models.

# RANDOM FOREST

Errors at Various Mtrys and Ntrees Kept Constant

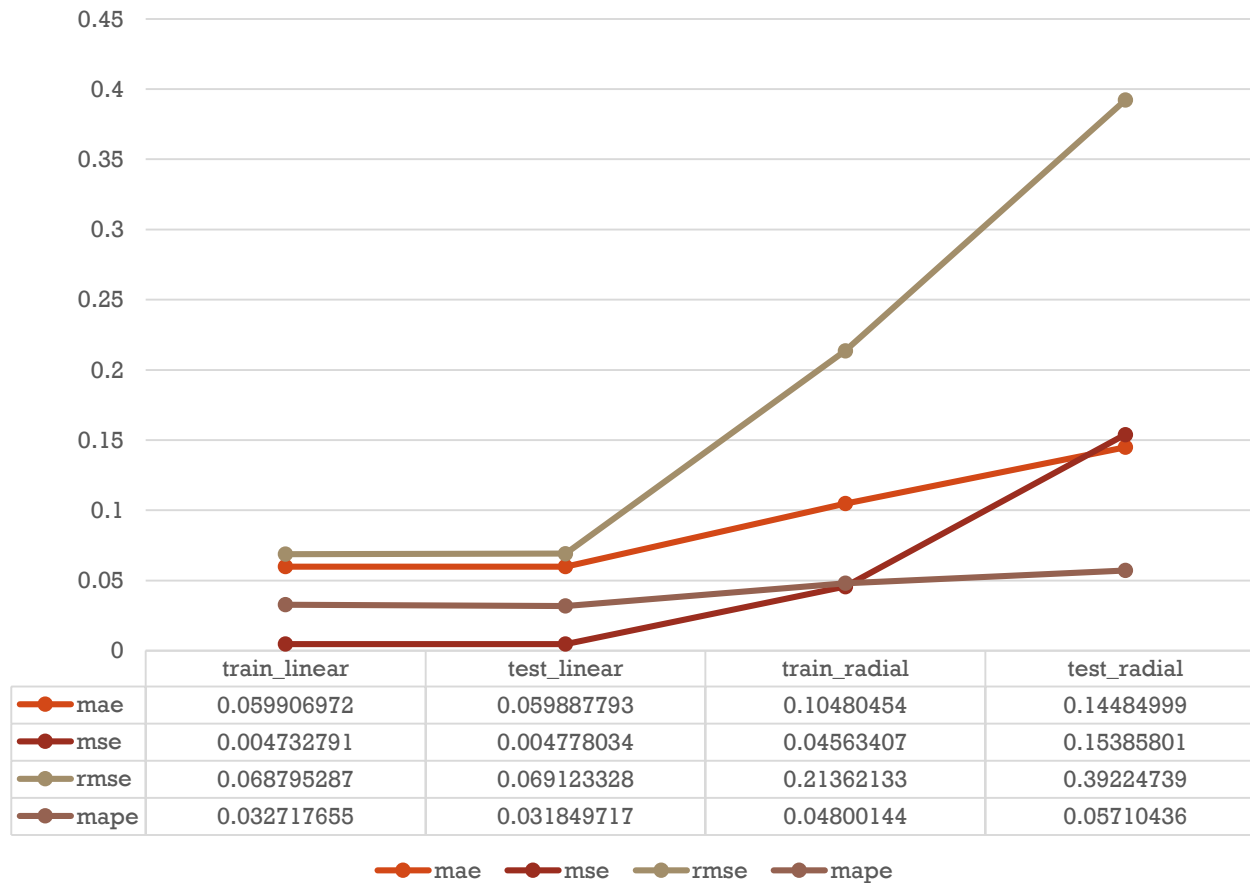


Errors at Various Ntrees and Mtrys Kept Constant



# SVM

Support Vector Machines Errors



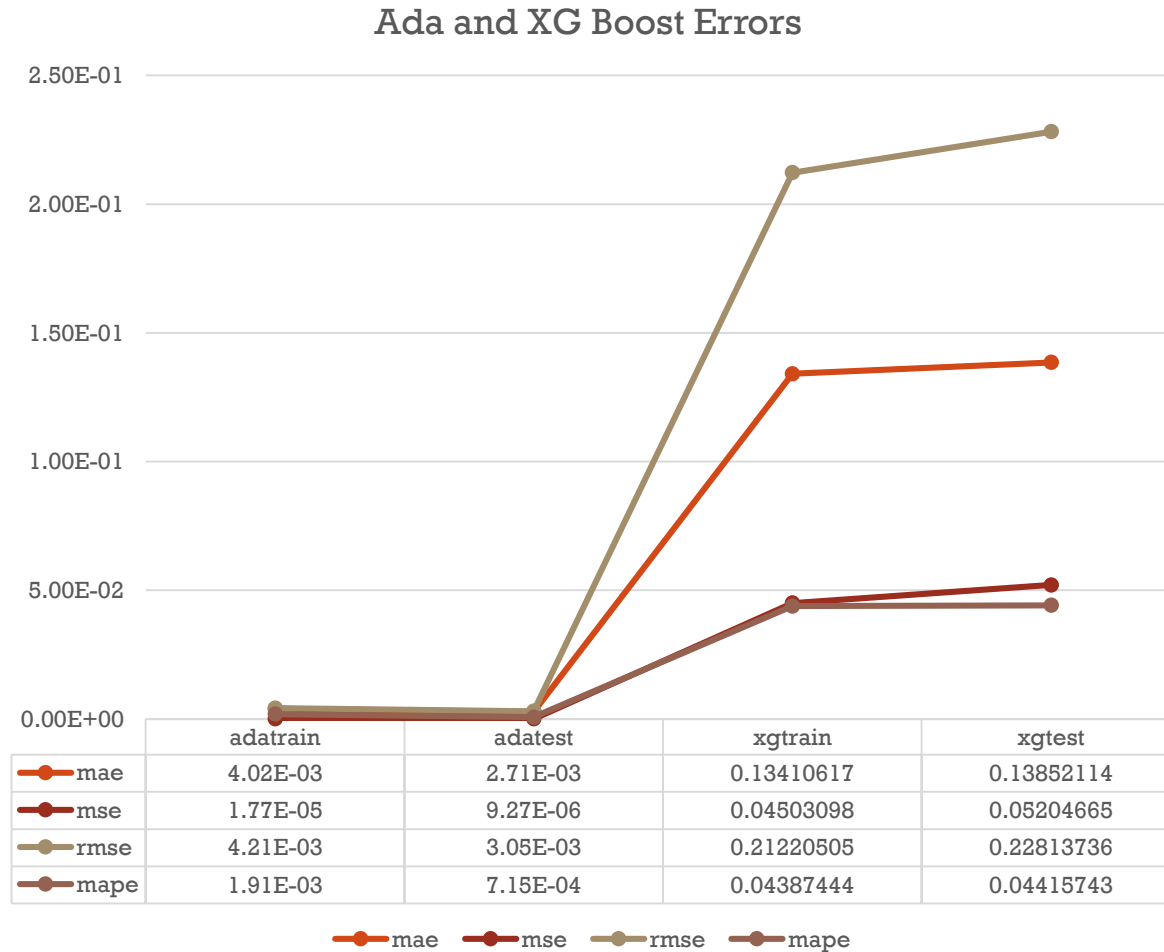
## SVM Regression:

- The values are separable well in the linear kernel and even the rmse of train and test is stable.
- The Tuning Parameter were not taken in account at the cost of time to tune in spite of errors being stable over train and test.
- But, in radial in hyperplane its not generalising well as the data might be linearly separable in linear SVM

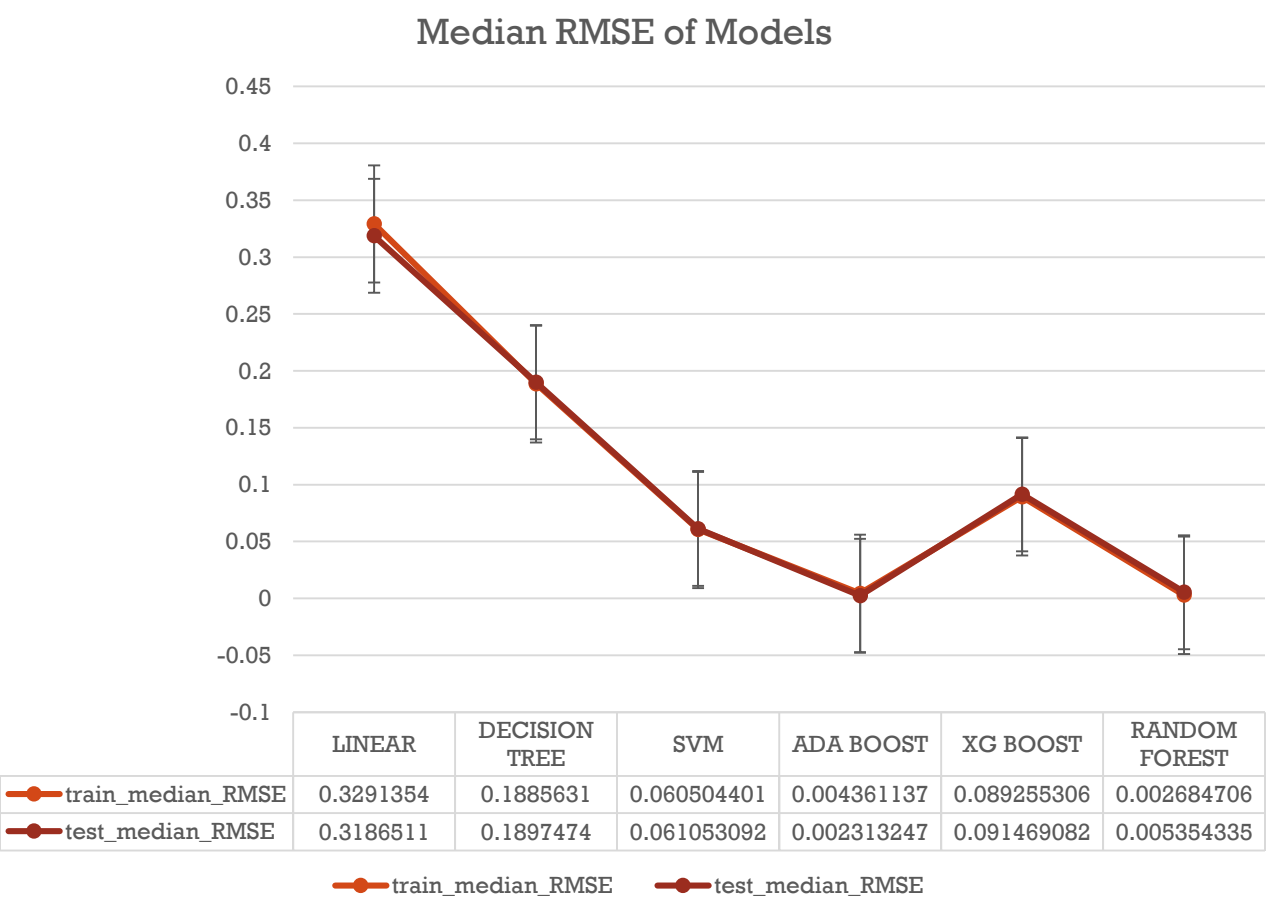
# ADA AND XG BOOST

## ADA and XG Boost:

- As ADA Boost take the class vectors did the range normalisation and predicted the probabilities.
- Then did denormalization of the probabilities and got the train and test error nearly about zero.
- XG Boost errors mape is generalised on train and test but the error deviation is high on both train and test due to which rmse is high when compared to other models.

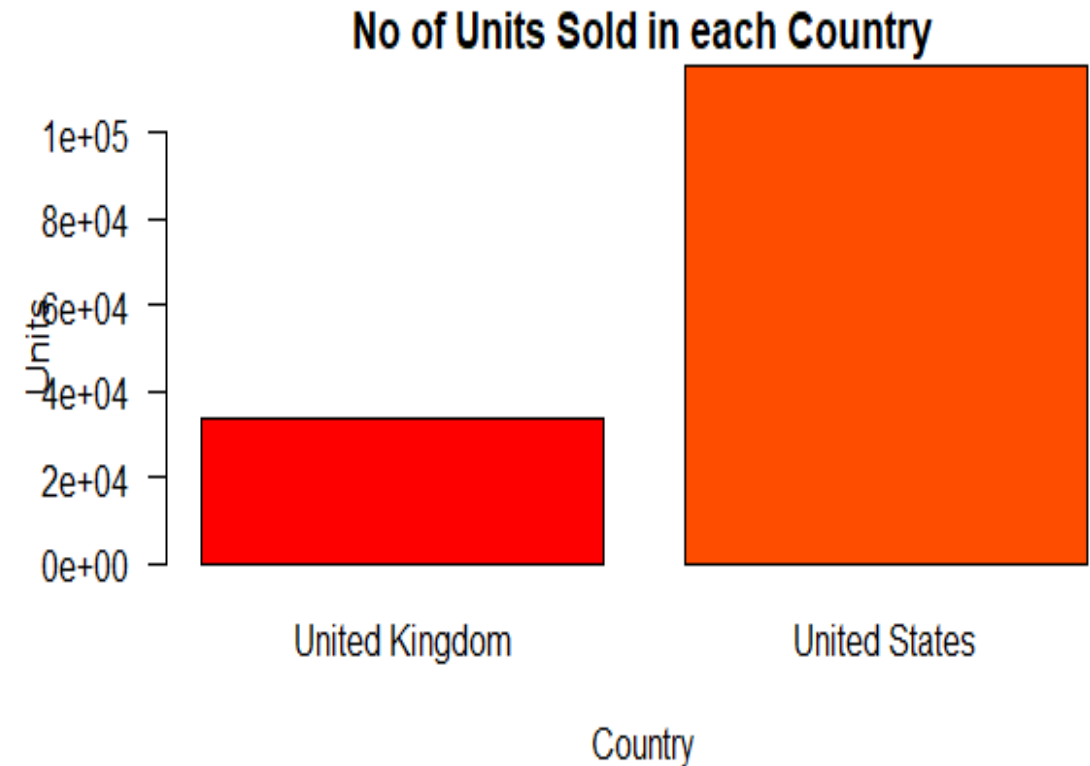
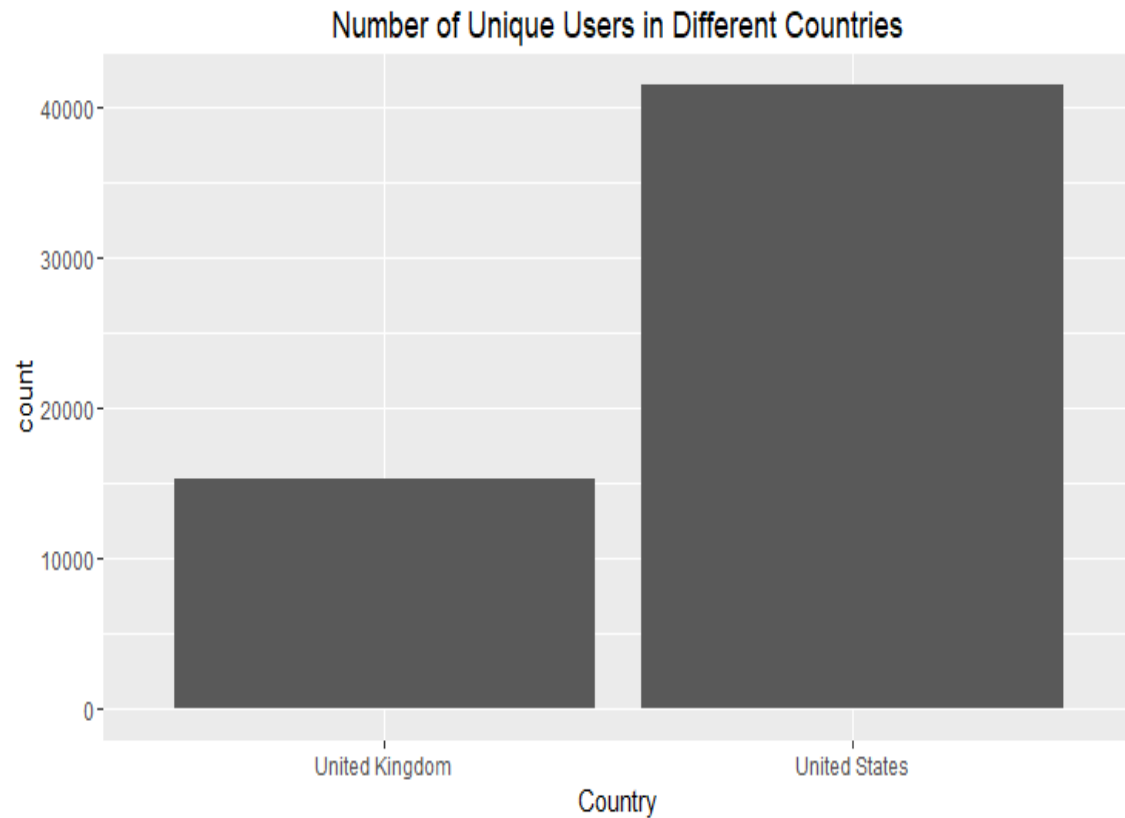


# Error Metric Evaluation Based on Median RMSE



# SUMMARY AND DATA INSIGHTS

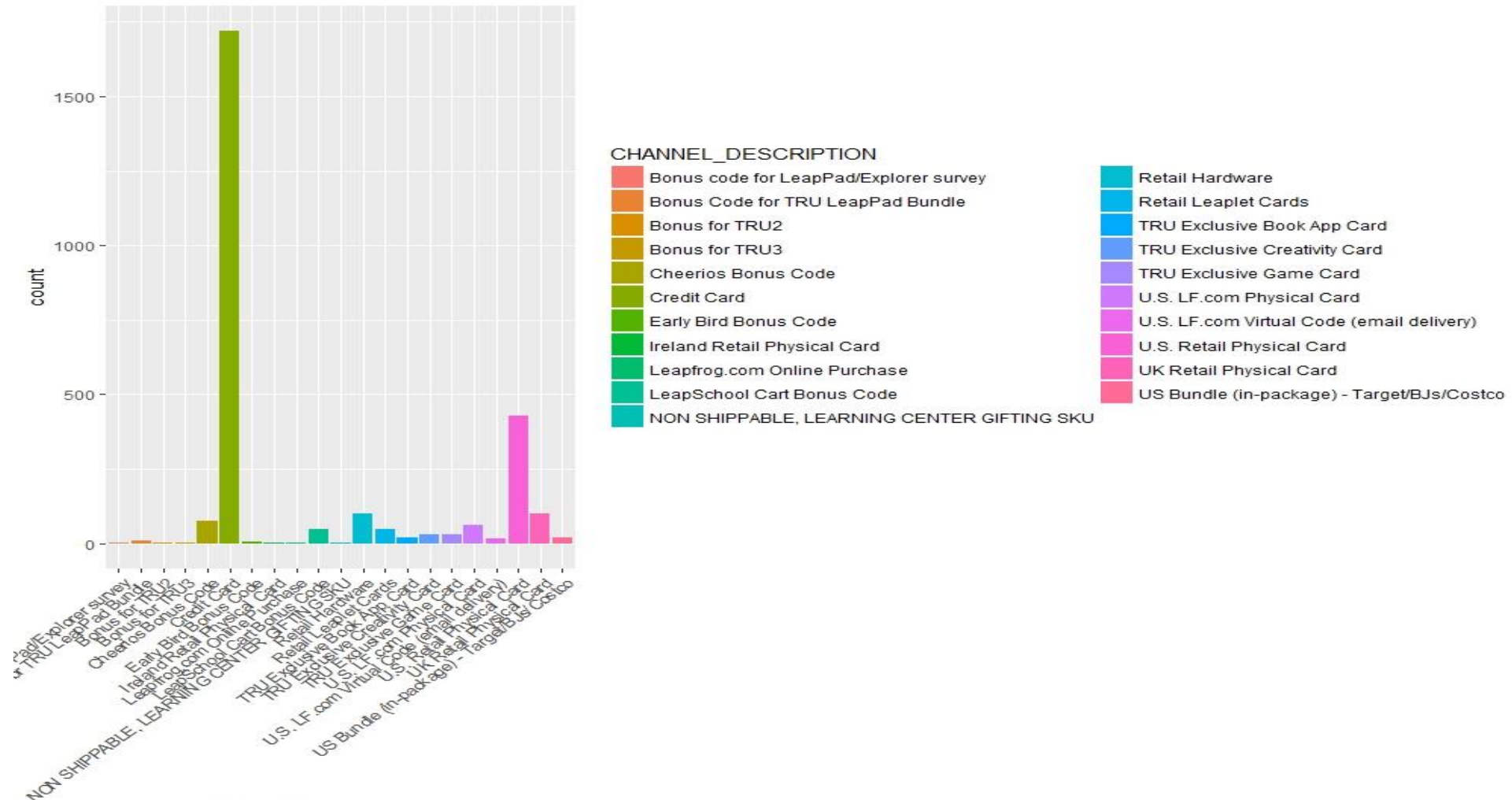
# DATA INSIGHTS BEFORE PRE-PROCESSING



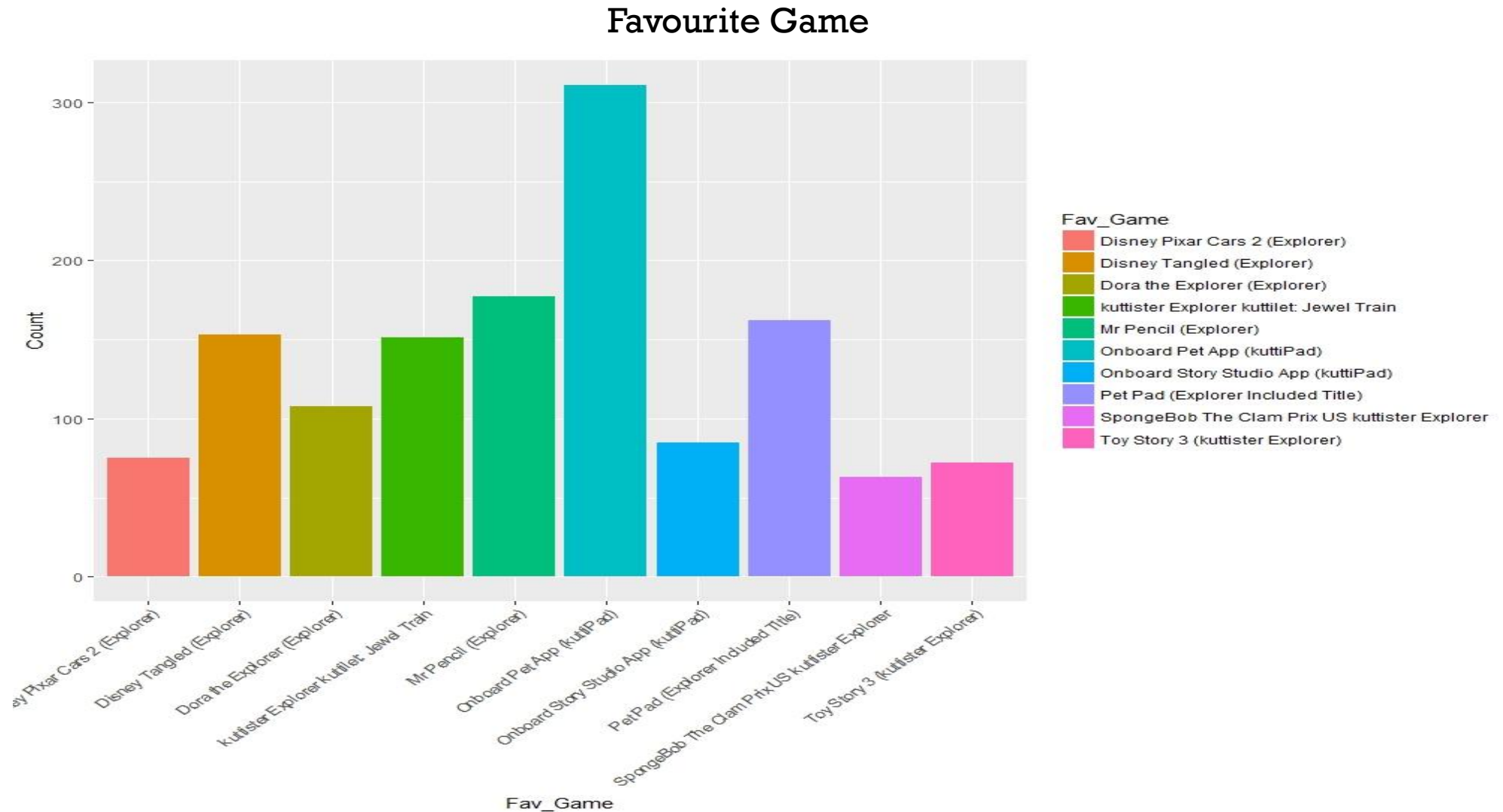


# Data Insights Before Pre-Processing

## Favourite Channel of Transactions



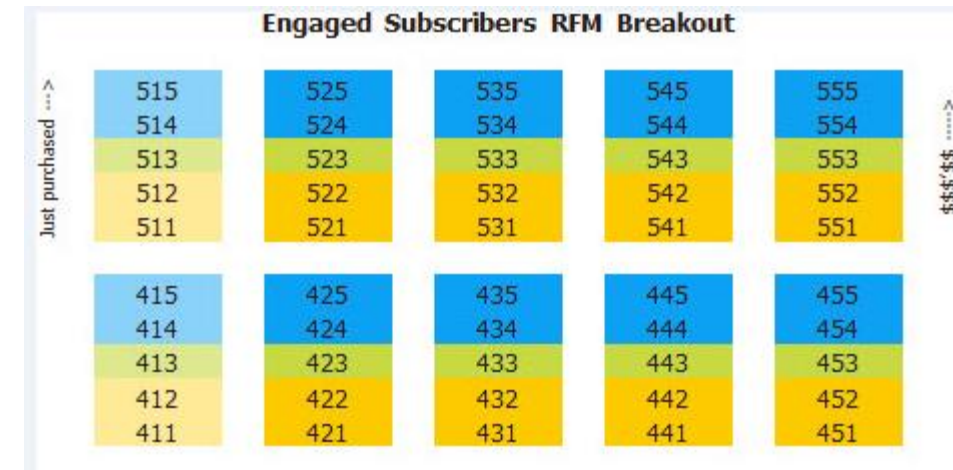
# DATA INSIGHTS BEFORE PRE-PROCESSING



# DATA INSIGHTS BEFORE PRE-PROCESSING

**RFM Analysis Table**

| RFM555    | RFM455   | RFM554   | RFM545   | RFM515    | RFM245   |
|-----------|----------|----------|----------|-----------|----------|
| 87890060  | 89620432 | 90600609 | 89435870 | 87838511  | 88431015 |
| 88686514  |          | 96136780 | 90511634 | 87975165  |          |
| 88721453  |          | 96934816 | 90863436 | 88089825  |          |
| 90343573  |          | 98825541 | 90896541 | 88388798  |          |
| 91409769  |          |          | 91003460 | 88422776  |          |
| 91730072  |          |          | 91577505 | 88650948  |          |
| 91838002  |          |          | 96392554 | 89017661  |          |
| 92139318  |          |          | 96571486 | 89165084  |          |
| 106253219 |          |          |          | 89171459  |          |
|           |          |          |          | 89515450  |          |
|           |          |          |          | 90823467  |          |
|           |          |          |          | 91004998  |          |
|           |          |          |          | 96149752  |          |
|           |          |          |          | 96633987  |          |
|           |          |          |          | 99678091  |          |
|           |          |          |          | 99802660  |          |
|           |          |          |          | 101022250 |          |



# SUMMARY

- As per the business constraints the Customer Life Time value has to be predicted and the SVM model is stable over others when compared in terms of Errors.
- As there is huge difference in the markets of US and UK the marketing strategies has to be improve a lot in UK market to scale up the profitability.
- As more customers are been inclined to the credit card as the channel of transaction company should make the portal highly secured.
- As there is a game named Disney Pixar car2 which is played a lot when compared to others a gamed around that themes will make the company to generate more revenue.
- There should be customers centric marketing strategies to retain the them and increase the revenue.

# Q&A

# THANK YOU