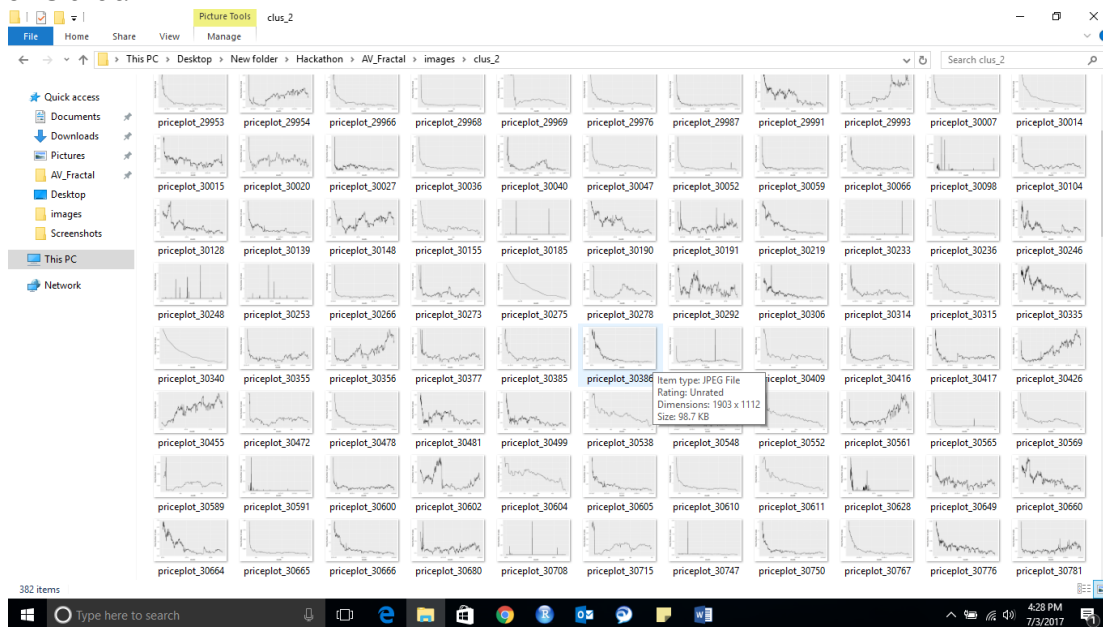# FRACTAL ANALYTICS HIRING HACKATHON
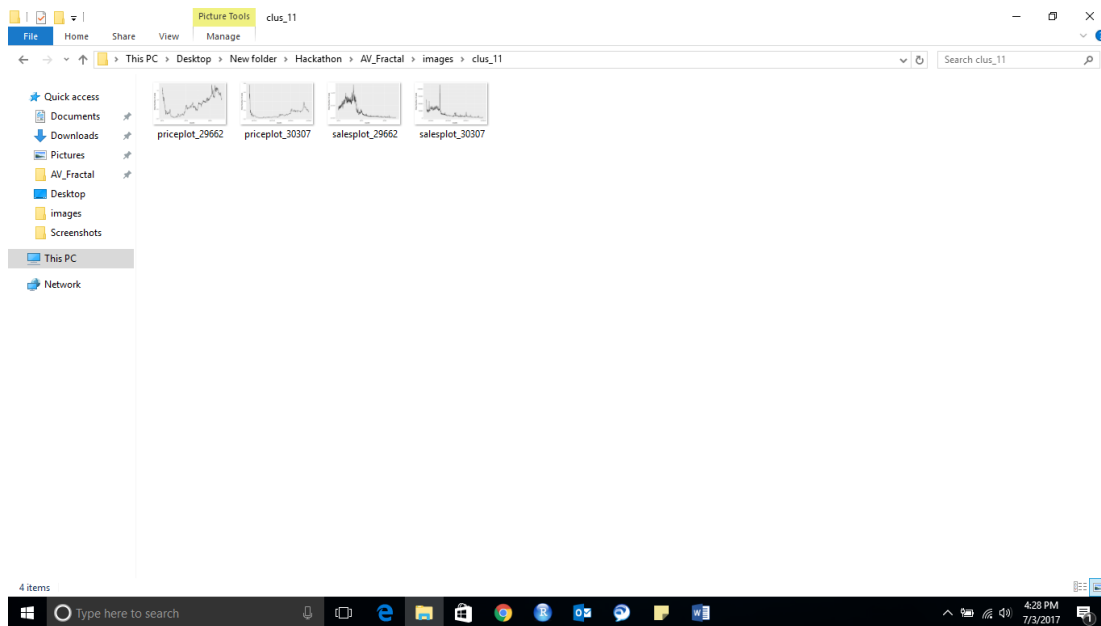## HARSHIT SAXENA

## Approach for the challenge:

- Programming Language: **R**, Visualization: **ggplot2**
- We need to build a model to predict **Price** and **Number_Of_Sales** for next 6 months (Huge Number) for all the Items.
- To find out any relation in distributions of each item represented as Item_ID.
- Blindly predicting/forecasting on all data gives random and very bad results which is a proof of not having a same behavior for all Items >> data cannot be combined to build one forecasting or linear model.
- Next step was to cluster all the items so that each cluster has its own distribution/behavior for **Price** and **Number_Of_Sales** and separate models can be created for the same. This also gave bad variations in each cluster. Screenshot below shows that:



[ All the images show **Price** distribution for all the items in one cluster: cluster 2]
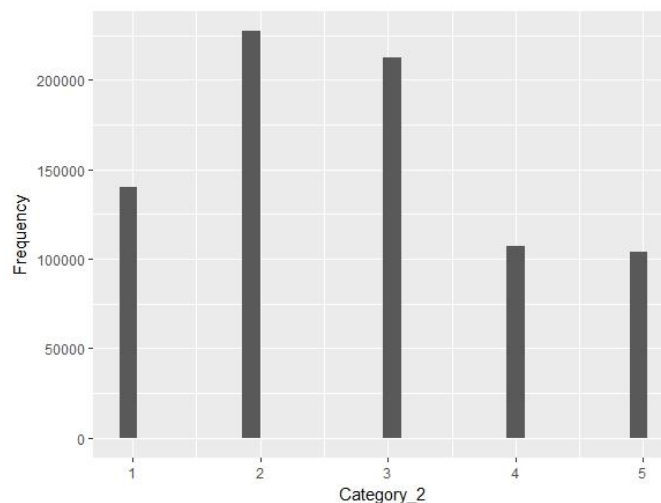
[ All the images show **Price** and **Number_Of_Sales** distribution for all the items in one cluster: cluster 11]

- These clusters didn't perform well so I approached for one vs one approach and build separate models for each Item_ID.
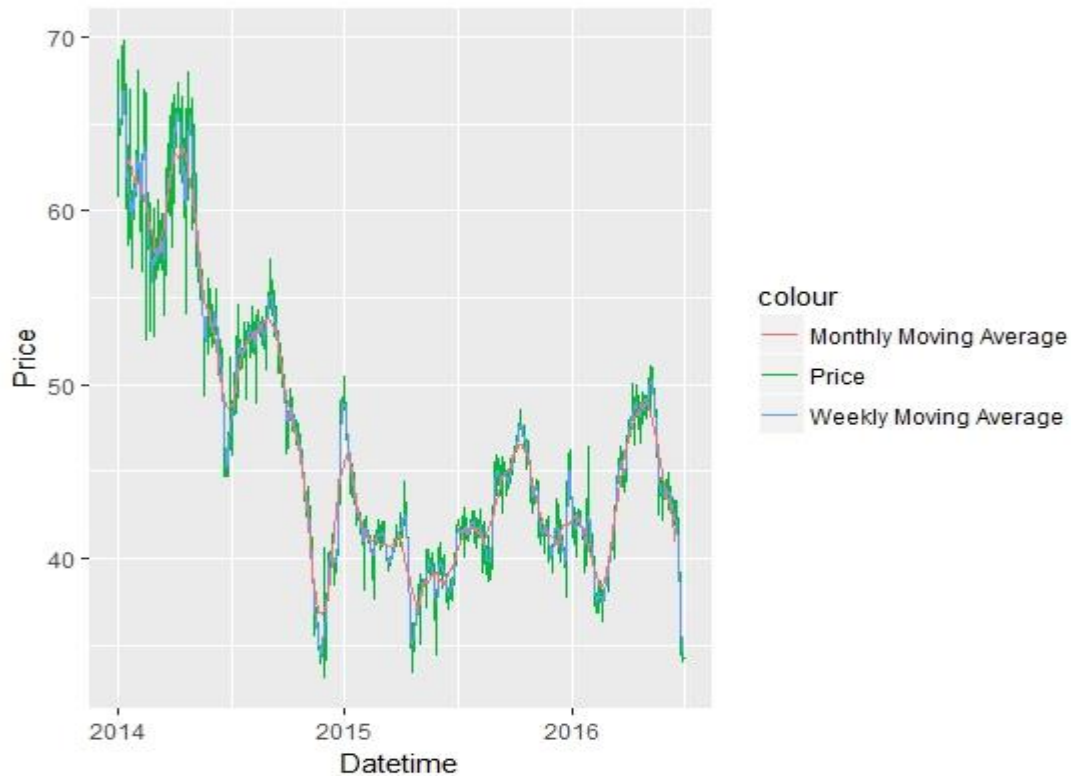- Model Selections is given below under algorithm selection section.

# Data preprocessing/ cleaning steps:
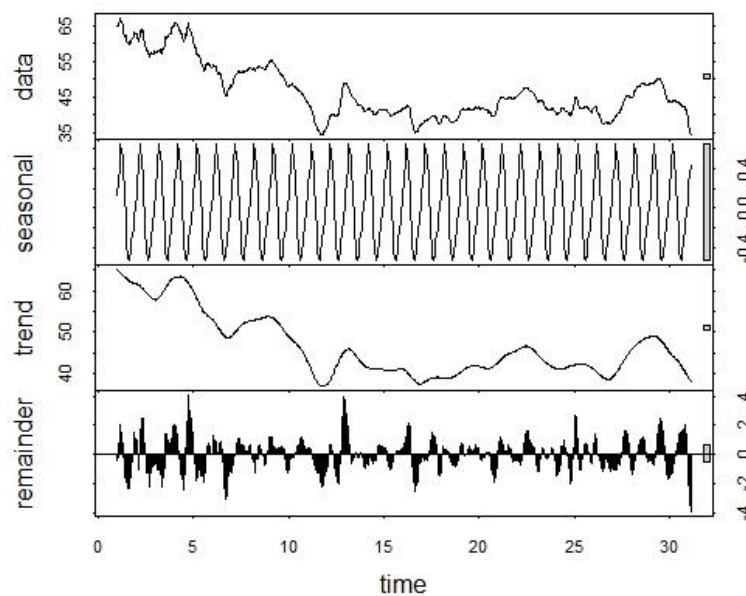
- Imputed Missing values in Category_2 with value = 2



- Creating **Price & Number_Of_Sales** a time series object and removing outliers from the distribution using tsclean () function.
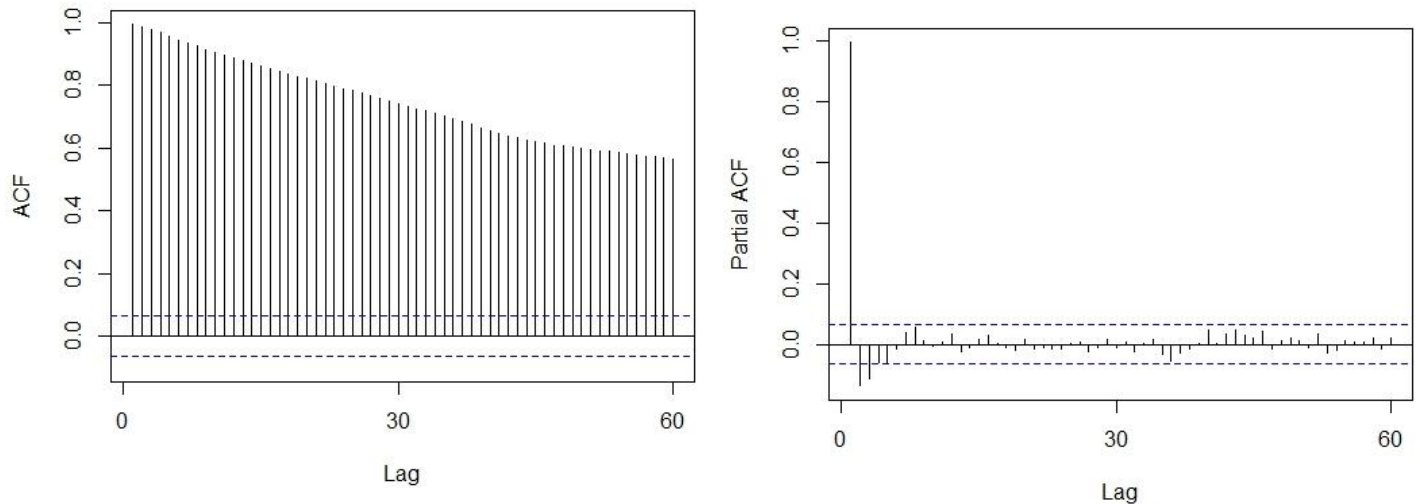
- Created features like year, month, date, weekday, weekend, diff_days (difference of days from day 1 which is 01-01-2014)
- Looked onto distribution using ACF and PACF plots which tells pattern with every 7 days.



[Distribution of Price for Item_ID = "30375" with weekly and monthly moving averages]
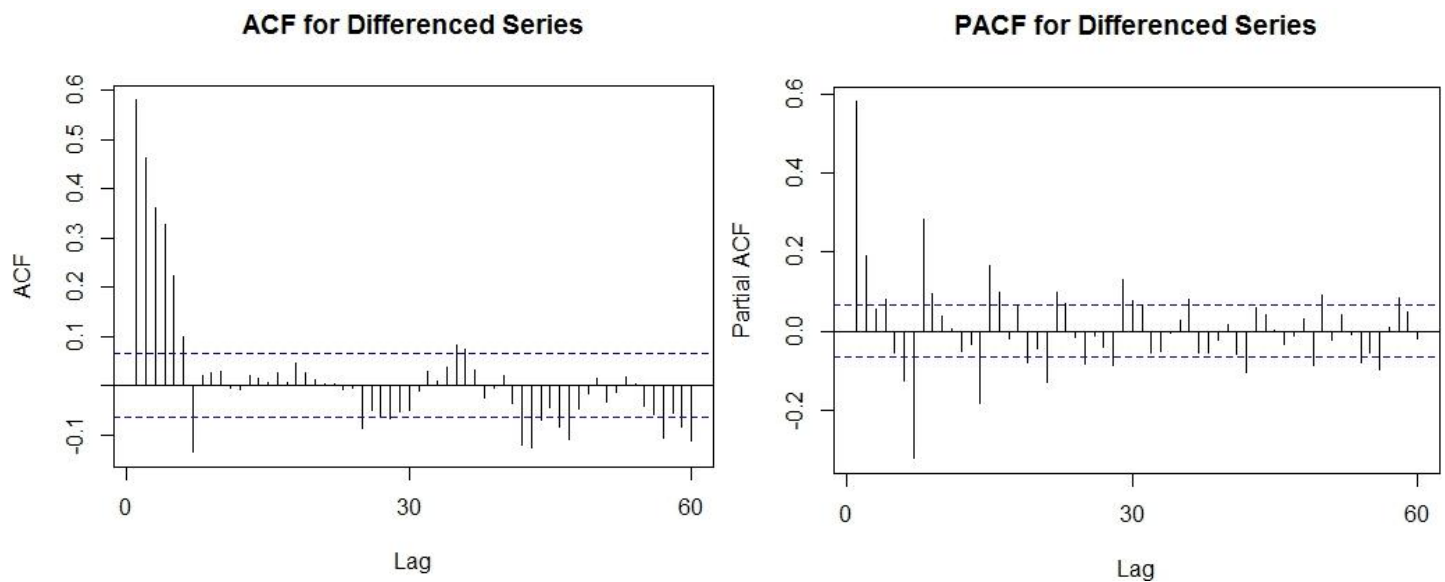


[checking for trend, seasonality in our distribution]

[ACF and PACF plots for our price distribution]

- We need to make our series stationary for ARIMA models and with differencing of 1 gives us stationary series from Dickey-Fuller Test.

**ACF for Differenced Series**                    **PACF for Differenced Series**



- We can see PACF clearly showing 7$^{th}$ lag out of blue line which tells weekly pattern.

## Algorithm Selection:

- Firstly, I tried ARIMA on each item which gave good results. It was making sense to me as well but what about smaller datasets like Item_ID = "30908" has just five entries.

- So, I introduced simple exponential smoothing forecasts for smaller datasets which further improved the results.
- Then I introduced linear regression for datasets which gave very bad results on leaderboard dataset
- Even XGBoost didn't improve my score because of less data and large forecasting range of 6 months.
- Then I did parameter pruning and found SES (simple exponential smoothing) to outperform ARIMA models and used it for final submission.

HARSHIT SAXENA

LEADERBOARD RANK – 19

I have attached code also in the mail.