# INSURANCE CLAIM

Prediction & Classification

ARUN RAJESH B.

# AGENDA

# Project Background & Objectives

Insurance Prediction dataset has 1.8L observations – 6 Numeric columns and 12 factor columns.

1. Explore the data & deliver key business Insights

2. Fit a regression model with highest adjusted R-square and least RMSE

3. Perform classification on the same dataset & achieve high accuracy scores

4. Perform clustering on the dataset and arrive at optimal clusters

5. Deploy the model for production ready

# Project Approach/Framework

| Data Loading | Data Preparation | Feature Selection | Model Building | Model Deployment |
|---|---|---|---|---|
| • **Data Import**<br>PostgreSQL<br>R –csv import<br>Spark<br><br>• **Data Exploration**<br>Visualization using<br>Metabase<br>PowerBI | • **Cleaning of data**<br>Char to factor conv.<br>Missing Values<br>one hot coding<br><br>• **Descriptive Stats**<br>Relation of Y& X's<br>Box Plot<br>Correlation Plot<br>auto - EDA | • **Relevant X only**<br>Removal of Unique Col<br><br>• **Relevant Transformation**<br>Log Transform Y<br>Normalization<br><br>• **New Feature Added**<br>Established Years | • **Regression**<br>Linear Regression<br>Random Forest<br>Regularized Model<br>Neural Networks<br><br>• **Classification**<br>Naïve Bayes<br>K-Nearest Neighbour<br><br>• **Clustering**<br>Kmeans | • **R Shiny**<br>Random Forest<br><br>• **Plumber API**<br>Linear Regression |

## H2o Auto ML Validation & Benchmark Setting**

**
H2o Flow: Distributed Random Forest Model   : R2: 0.998, RMSE: 3231.59
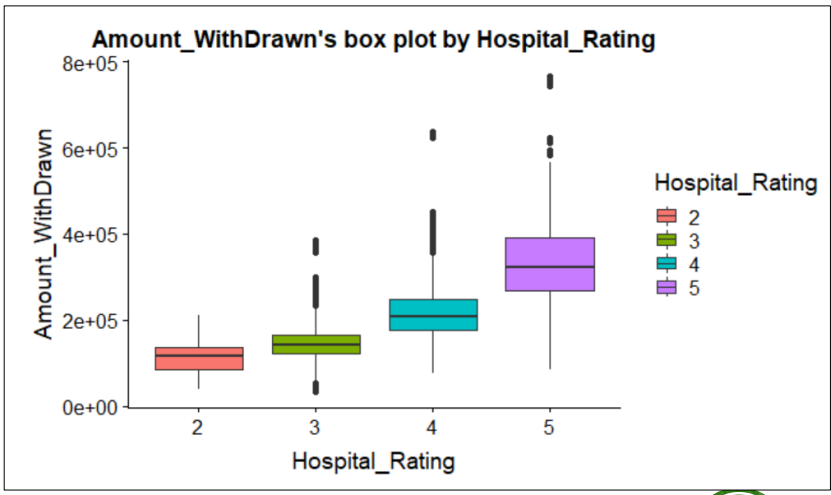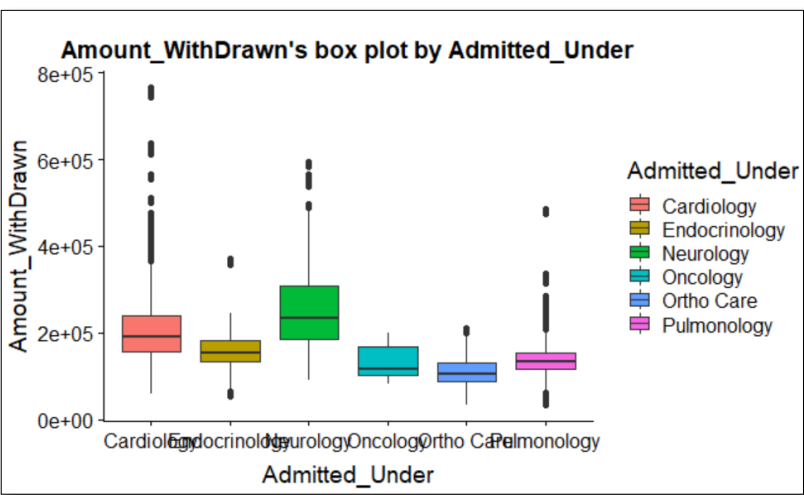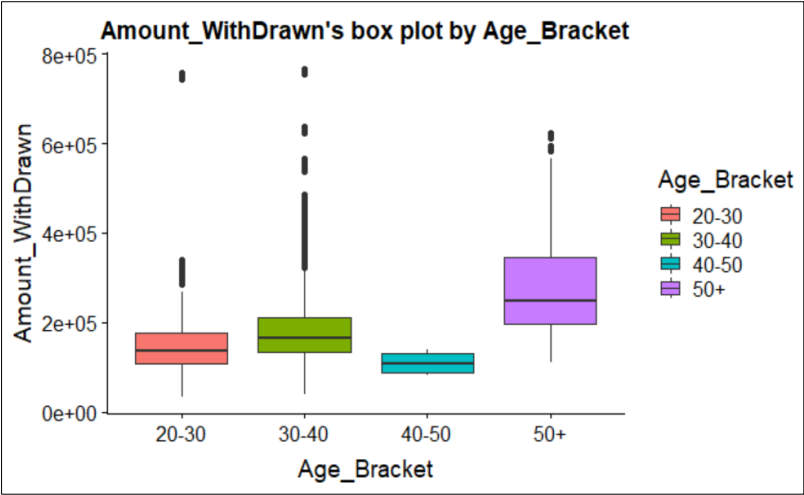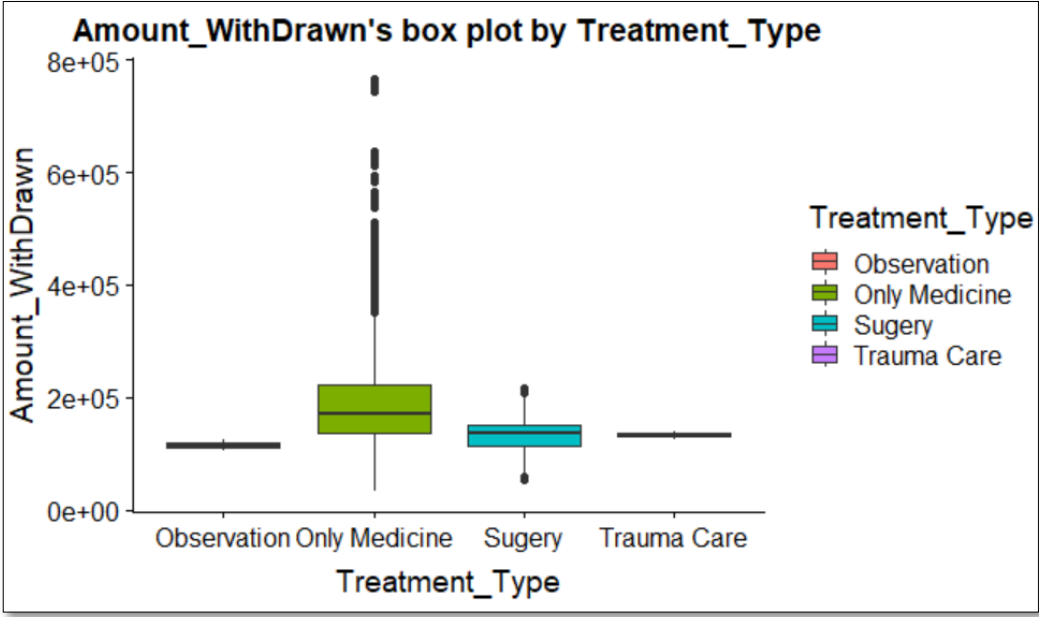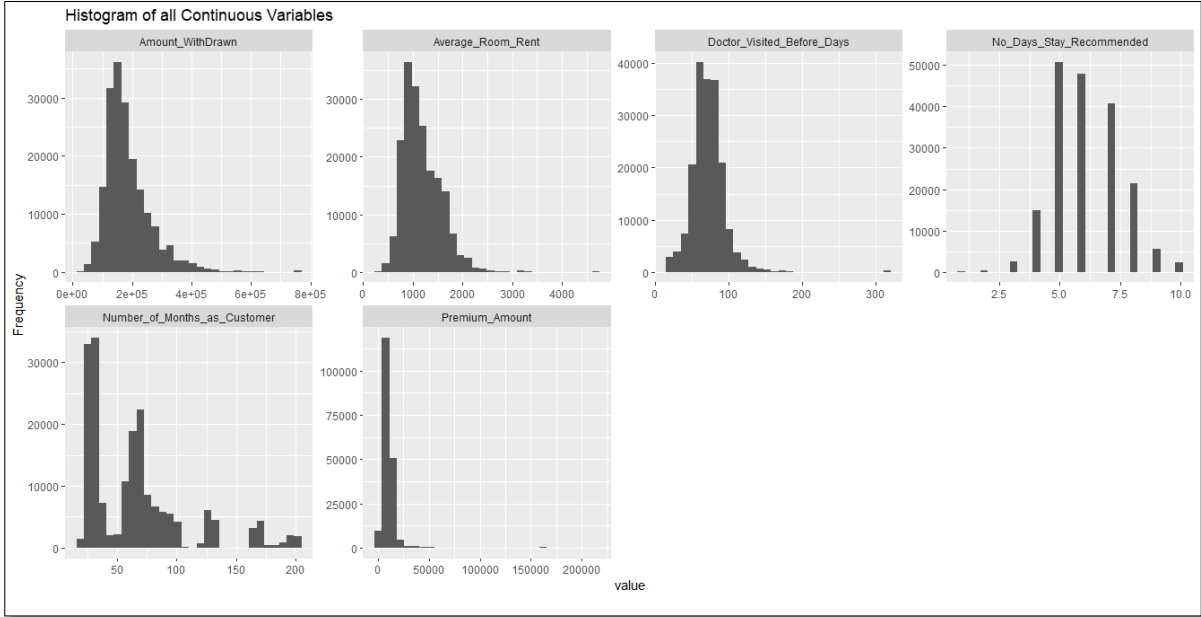H2o Flow: Gradient Boosting Machine Model  : R2:0.945, RMSE: 18,801

# Key Business Insights

- Average amount withdrawn under Polo4 for is highest mostly coming from semi urban region and urban regions in popular multi-speciality hospitals

- Most withdrawals are made in Urban, followed by town (<50% of Urban), semi urban and bigger town are very less withdrawals **(scope for expansion in semi urban and bigger towns)**

- Across all policy types, most withdrawals were made for medicine treatment type. For POLO5, comparatively higher withdrawal made for surgery treatment type

- Average amount withdrawn is highest for patients admitted under neurology followed by cardiology and endocrinology **(Scope for more hospitals to come / Health awareness to be made on neurology related disorders)**

- Most withdrawals are made by people in age bracket **30 to 40 (surprising as not old)** across all hospital locations. 50 + withdraw most in bigger towns. 40 to 50 is the only age bracket which is making withdrawals for surgery and trauma care. People in other brackets are only doing it for Medicine. Average amount paid is highest for 50+ group for "medicine" treatment type.

# Few Exploratory analysis in R (auto – EDA)

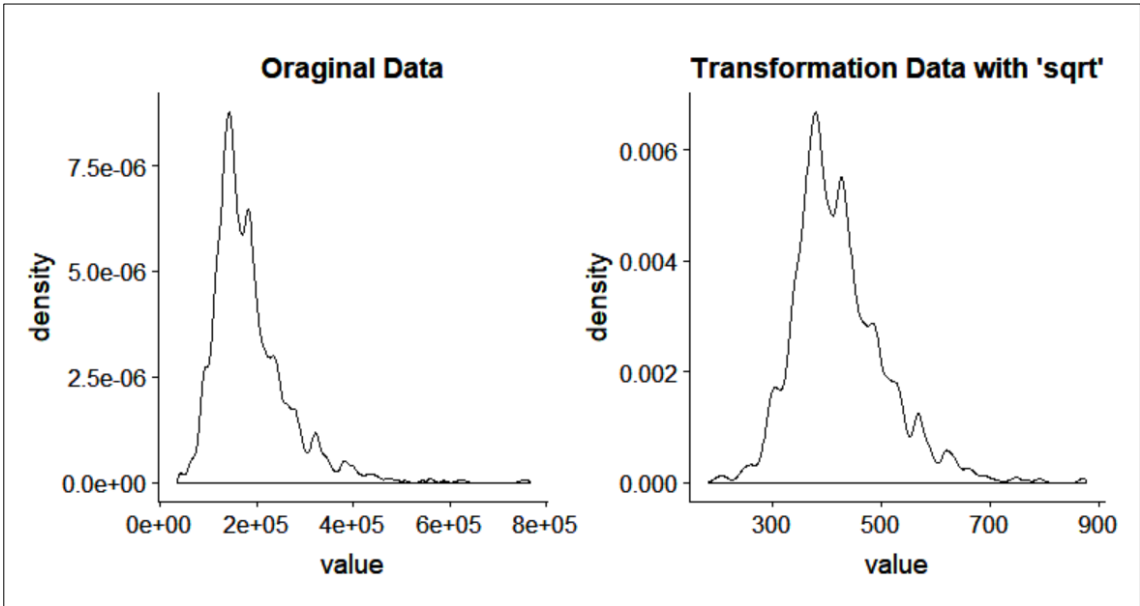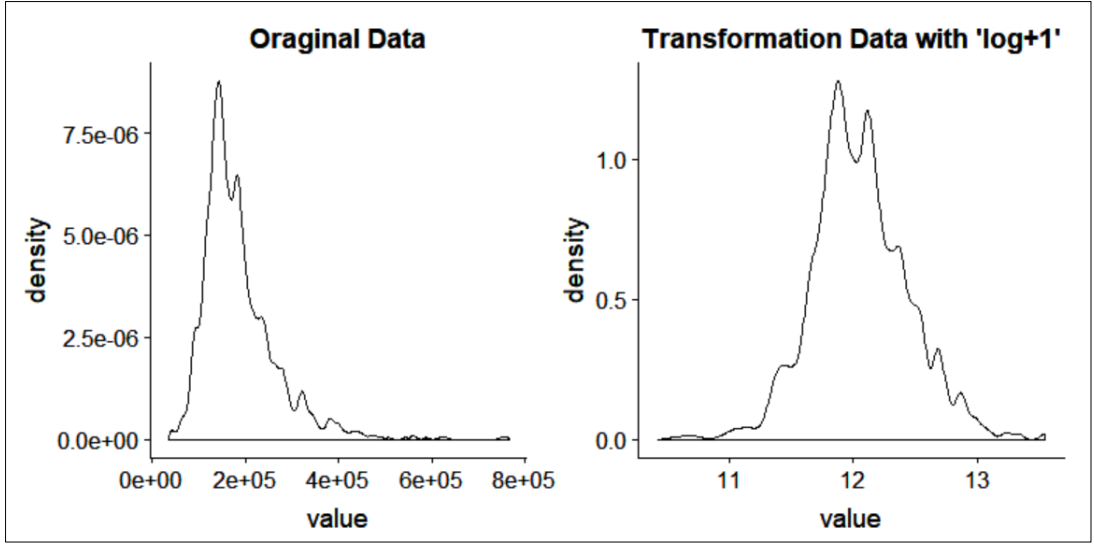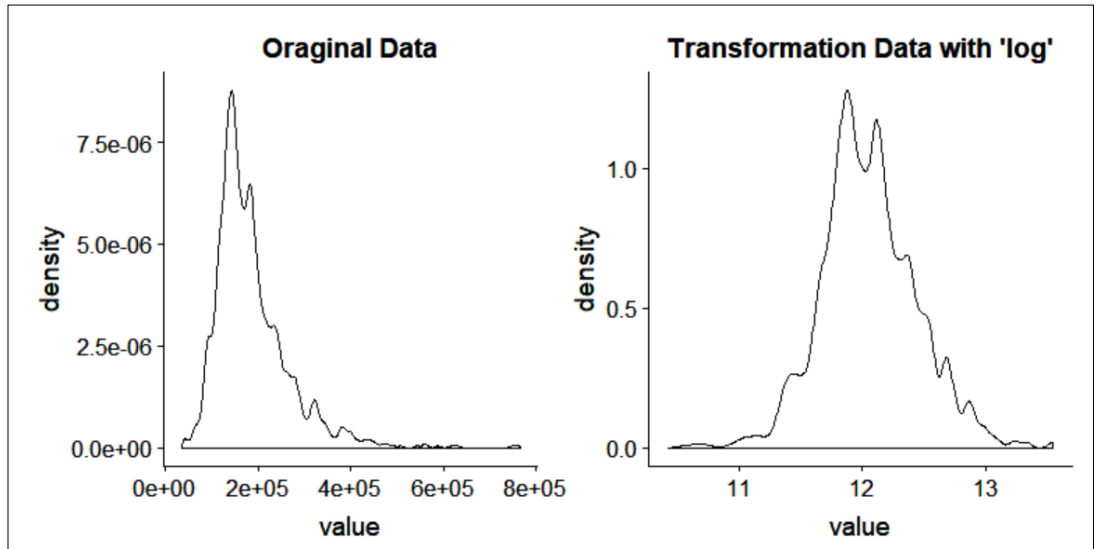# Few Transformations(auto – EDA)

Table of correlation coefficients (0.5 or more)

| Variable1 | Variable2 | Correlation Coefficient |
|---|---|---|
| No_Days_Stay_Recommended | Amount_WithDrawn | 0.7902981 |
| Average_Room_Rent | Amount_WithDrawn | 0.6053716 |

# Regression Algorithm Results I

Target Variable : Amount With Drawn

| S.No | Regression Algorithm | Input Pre-processing & Other | Train Model Accuracy (Percentage) | Hyper Parameter &Others | Test Model Accuracy (Percentage) |
|------|---------------------|------------------------------|-----------------------------------|-------------------------|----------------------------------|
| 1 | Linear Regression (Base Model) | - Scaling , With/Without Cross Validation | Res.SE : 0.4785<br>R -squared: 0.7703 (adjusted) | - CSV file output Generated for test data | RMSE : 0.4738<br>R-squared : 0.7740<br>MAE : 0.3178<br>MAPE : 4.1097 |
| 2 | Linear Regression (Log Transform of Y Variable) | - Skewness was observed in Y Parameter.<br>- Log Transform, Scaling ,<br>- With/Without Cross Validation | Res.SE : 0.1645<br>R -squared: 0.8187 (adjusted) | - CSV file output Generated for test data | RMSE : 0.1629<br>R-squared : 0.824<br>MAE : 0.1218<br>MAPE : 0.0101 |
| | Linear Regression (Log Transform of Y and Continuous X Variables) | - Skewness was observed in Y Parameter.<br>- Log Transform, Scaling ,<br>- With/Without Cross Validation | Res.SE : 0.1603<br>R -squared: 0.8278 (adjusted) | | RMSE : 0.1599<br>R-squared : 0.8305<br>MAE : 0.1183<br>MAPE : 0.0098 |

# Learnings | Log Transformation of Skewed data gives better results

# Regression Algorithm Results II

Target Variable : Amount With Drawn

| S.No | Regression Shrinkage Models | Input Pre-processing & Other | Train Model Performance (Percentage) | Hyper Parameter &Others | Test Model Accuracy (Percentage) |
|---|---|---|---|---|---|
| 4 | Lasso Regression | - Scaling ,Center<br>- With Cross Validation | R -squared: 0.8264 | - Model best tune<br>Fraction : 0.9 | RMSE : 0.1606<br>R-squared : 0.8294<br>MAE : 0.1185<br>MAPE : 0.0098 |
| 5 | Ridge Regression | - Scaling , Center<br>- With Cross Validation | R -squared: 0.8276 | - Model best tune<br>lambda : 0 | RMSE : 0.1599<br>R-squared : 0.8305<br>MAE : 0.1183<br>MAPE : 0.0098 |
| 6 | Elastic Net Regression | - Scaling , Center<br>- With Cross Validation | R -squared: 0.8276 | - Model best tune<br>lambda : 0<br>Fraction : 1 | RMSE : 0.1599<br>R-squared : 0.8305<br>MAE : 0.1183<br>MAPE : 0.0098 |

# Learnings

Shrinkage Models didn't improve model performance. They yield the same result as normal Linear models for this dataset
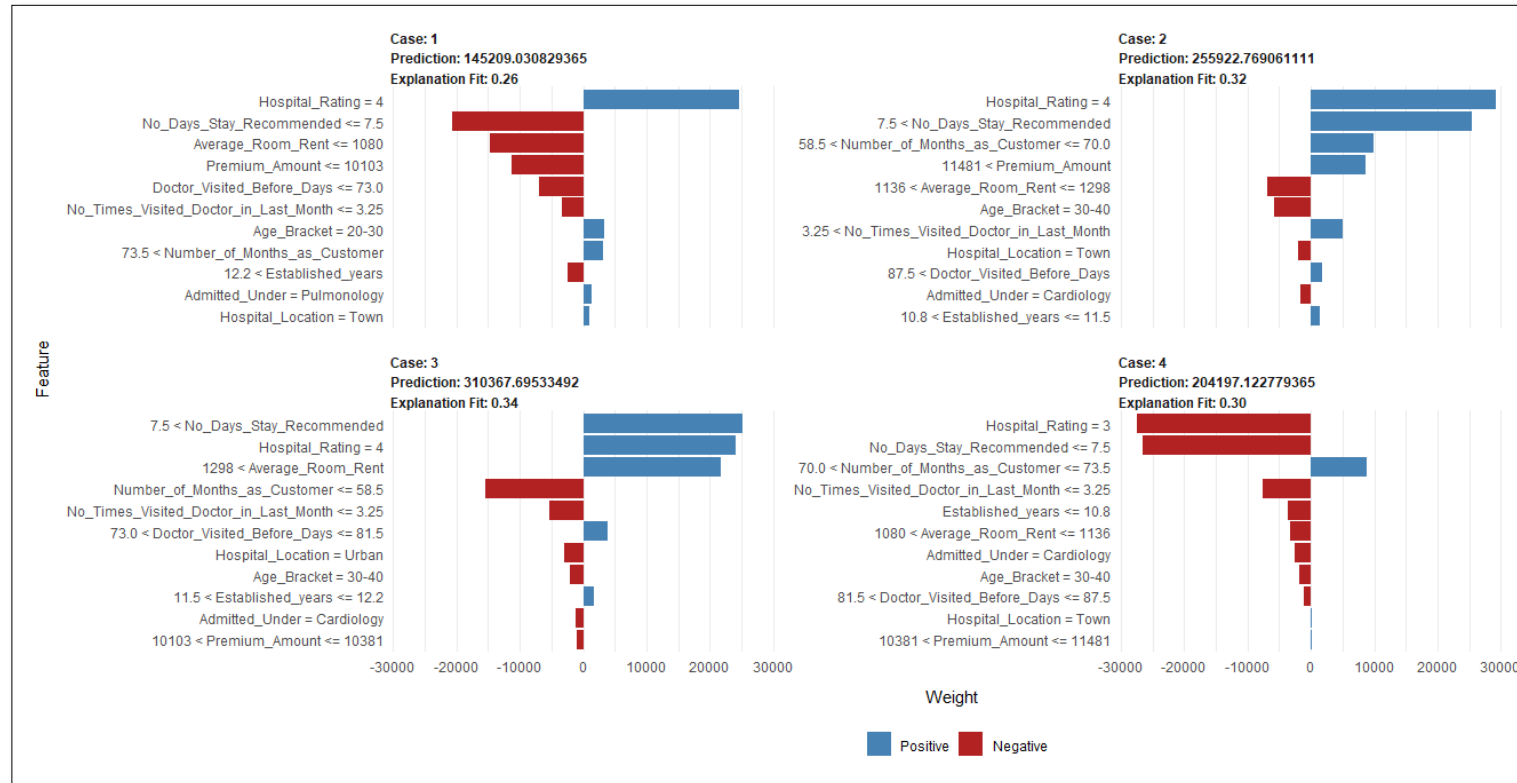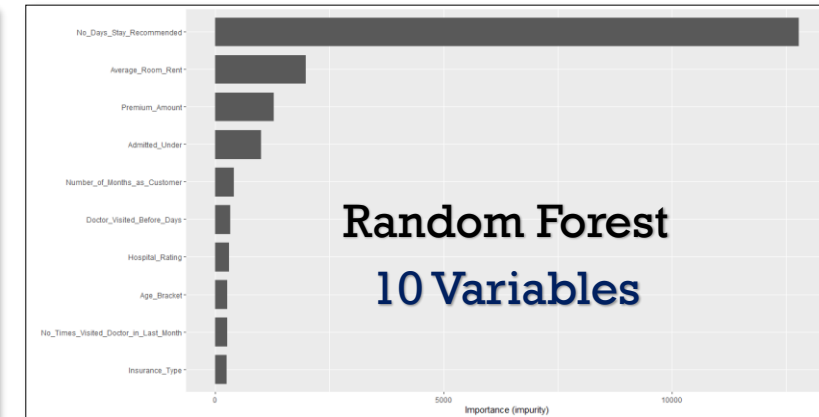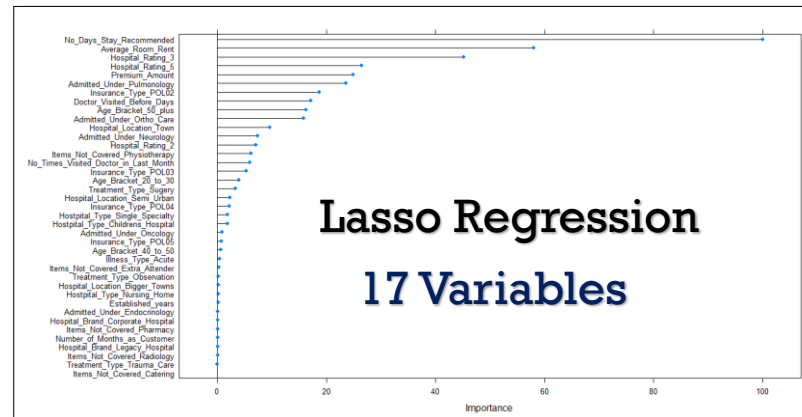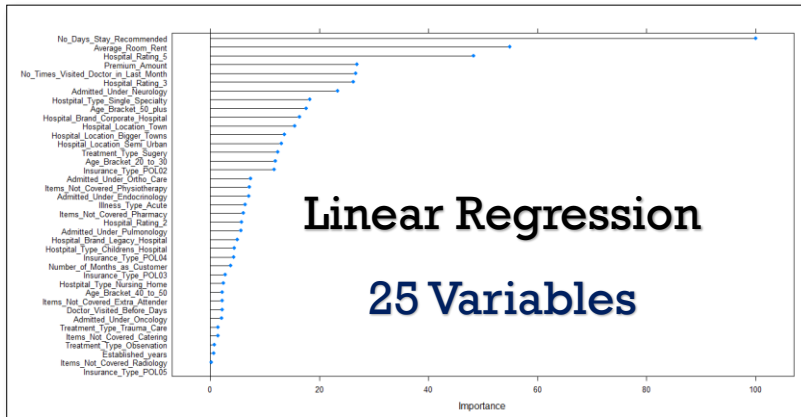
# Regression Algorithm Results III

| S.No | Regression - Tree Based Models | Input Pre-processing & Other | Train Model Performance (Percentage) | Hyper Parameter &Others | Test Model Accuracy (Percentage) |
|---|---|---|---|---|---|
| 7 | Random Forest (Spark – ML Lib through Sparklyr Package) | - Log Transform of Y Parameter | RMSE: 0.05376 R2:0.9805 | Num. trees =80 Max_depth =12 | RMSE : 0.0537 R-squared : 0.9811 MAE : 0.0376 MAPE : 0.0032 |
| 8 | Random Forest (Ranger Package) | - Scaling , <br> - With Cross Validation | OOB prediction error (MSE): 0.01809 R squared (OOB): 0.98184 | - Num. trees =500 <br> - Mtry =6 <br> - Node. Size =5 | RMSE : 0.1377 R-squared : 0.9851 MAE : 0.0935 MAPE : 1.3804 |
| | Random Forest (Ranger Package) | - No Scaling of X Parameters <br> - Log Transform of Y Parameter | OOB prediction error (MSE): 0.00048 R squared (OOB): 0.99677 | - Mtry= 11:17 <br> - Min. Node. Size = 3:9 <br> - Splitrule="variance" <br> - CSV file output Generated <br> Optimum Value <br> Mtry =17 <br> Min. Node. Size =9 | RMSE : 0.2212 R-squared : 0.9967 MAE : 0.0167 MAPE : 0.0014 |

# Learnings
Tree models perform better than regression models when hyperparameter tuning is involved

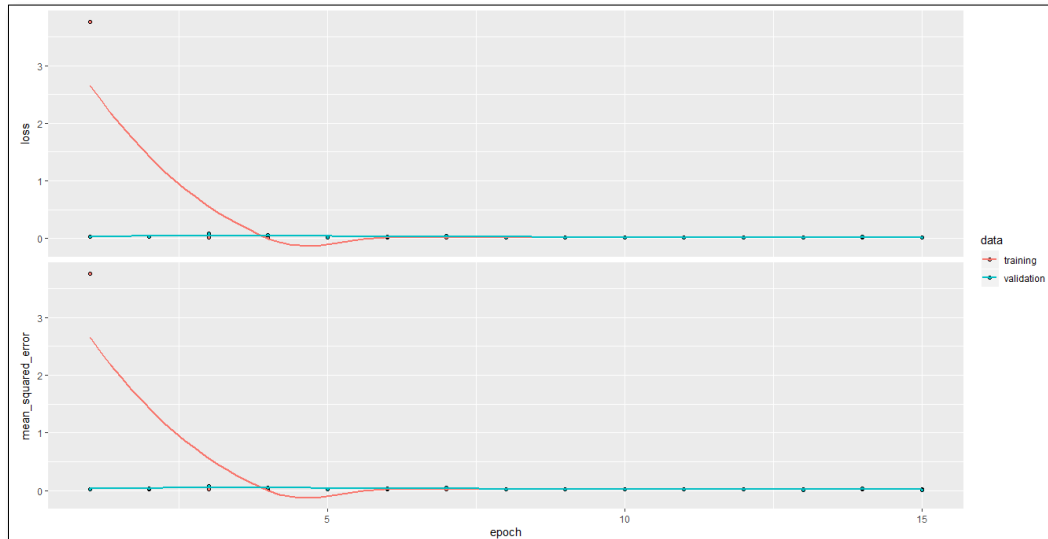# Variable Importance & Model Interpretation (LIME)



Linear Regression — 25 Variables

Lasso Regression — 17 Variables

Random Forest — 10 Variables

Case: 1
Prediction: 145209.030829365
Explanation Fit: 0.26

Case: 2
Prediction: 255922.769061111
Explanation Fit: 0.32

Case: 3
Prediction: 310367.69533492
Explanation Fit: 0.34

Case: 4
Prediction: 204197.122779365
Explanation Fit: 0.30

- local interpretations of ML models using 4 Observations of training dataset (LIME Package)
- 11 influential variables are plotted. Explains the support of variables for prediction
- No. of Days stay recommended, Hospital Rating, Average Room Rent are the most influential variables

# Regression Algorithm Results IV

| S.No | Regression - Neural Networks | Input Pre-processing & Other | Train Model Performance (Percentage) | Hyper Parameter &Others | Test Model Accuracy (Percentage) |
|---|---|---|---|---|---|
| 9 | Neural Networks – Keras Package | - Log Transform of Y Parameter<br>- Scaling, Center values of X Parameter | Train<br>loss: 0.0175 –<br>mean_squared_error: 0.0175<br><br>Validation<br>loss: 0.0135<br>mean_squared_error: 0.0135 | Compile function:<br>optimizer = "rmsprop",<br>loss = "mse",<br>metrics = c("mse")<br><br>Training<br>epochs = 15,<br>batch_size = 32,<br>validation_split = 0.2 | RMSE        : 0.1215<br>R-squared : 0.9223<br>MAE          : 0.0954<br>MAPE        : 0.0079 |



# network architecture
layer_dense(units = 50, activation = "relu", input_shape = ncol(train_x)) %>%
layer_dense(units=20,activation="relu")%>%
layer_dense(units = 5, activation = "relu")

```
Layer (type)                        Output Shape                    Param #
================================================================================
dense_19 (Dense)                    (None, 50)                      2000

dense_20 (Dense)                    (None, 20)                      1020

dense_21 (Dense)                    (None, 5)                       105
================================================================================
Total params: 3,125
Trainable params: 3,125
Non-trainable params: 0
```

# Classification Algorithm Results

Target Variable : Amount With Drawn
Class: High, Low

| S.No | Classification Algorithm | Input Pre-processing & Other | Train Model Accuracy (Percentage) | Hyper Parameter Condition & Optimum value | Test Model Accuracy (Percentage) |
|---|---|---|---|---|---|
| 1 | Naïve Bayes Classification algorithm | - 3 Fold Repeated Cross Validation<br>- No scaling | 90.05 | Laplace=1 to 5, usekernel=TRUE,FALSE, adjust=1 to 5<br><br>Optimum Values:<br>Laplace = 1.5,<br>Usekernel=TRUE<br>Adjust=2 | Accuracy : 90.02<br>Sensitivity : 88.84<br>Specificity : 91.20 |
| 2 | K-Nearest Neighbours algorithm. | - Without Scaling | | K = 7 to 17<br>Optimum Value:<br>K = 11 | Accuracy : 98.23<br>Sensitivity : 98.30<br>Specificity : 98.17 |
| | | - With Scaling<br><br>(Note: Scaling execution was little faster than the other) | | K = 7 to 17<br>Optimum Value:<br>K = 11 | Accuracy : 98.48<br>Sensitivity : 98.46<br>Specificity : 98.50 |

# Learnings

Many Zero Probabilities were identified during execution of Naïve Bayes algorithm. Used Laplace smoothing to prevent such errors

# Clustering Algorithm Results

| S.No | Clustering Algorithm | Input Pre-processing & Other | Hyper Parameter Condition & Optimum value | Observations |
|------|---------------------|------------------------------|-------------------------------------------|--------------|
| 1 | K means Clustering | - PCA performed on input parameters to reduce the dimensions<br><br>- Base on the screeplot, Dimension from PC1 to PC 17 are selected and clustering is performed | Kmeans Parameters<br>K = 1 to 30<br>Max.Iter = 100<br>Algorithm = MacQueen | Though, PC1 to PC17 is chosen from Screeplot (eigen value >1), the total variation explained is only 63%<br><br>Optimum no. of clusters couldn't be established using Elbow Method.<br><br>Attempt is made to plot 10 cluster on a 2 dimensional space |
| 2 | K Means Clustering – ClusterR Package | - PCA performed and dimensions are reduced to 2 | max_clusters =15<br>criterion = WCSSE<br>max_iters =100<br>initializer = kmeans++ | PC1 and PC2 explains only 15% variation in the original data<br><br>Sharp decline observed at cluster 3 and 5.<br><br>Attempt is made to plot 5 cluster on a 2 dimensional space |

## Learnings

Optimum cluster based on elbow method couldn't be figured out. The total variation explained by PCA is less

# Clustering Algorithm Results



Screeplot of the all 39 PCs

PC1 – PC17
No elbow identified

5 Cluster

Cluster plot

10 Cluster

PC1 – PC2
Cluster 3,5

# Model Deployment – R Shiny

# Model Deployment – Plumber API



## 1

### swagger

http://127.0.0.1:8000/swagger.json?schemes=http&host=127.0.0.1:8000&path=/    Explore

## Run Insurance predictions,Target Variable: Amount withdrawn 1.0.0

[ Base url: 127.0.0.1:8000/ ]
http://127.0.0.1:8000/swagger.json?schemes=http&host=127.0.0.1:8000&path=/

This API takes various input continuous parameters indicates Amount withdrawn(Target Variable)

Schemes
HTTP

### default

**POST** /predict

**Parameters**    Try it out

| Name | Description |
|---|---|
| Established_years | Please enter numeric values between 9 and 13 |

## 2

| | |
|---|---|
| 11 | |
| (query) | |
| No_Times_Visited_Doctor_in_Last_Month | Please enter numeric values between 0 and 8 |
| number | 2 |
| (query) | |
| Average_Room_Rent | Please enter numeric values between 339 and 4717 |
| number | 1097 |
| No_Days_Stay_Recommended | Please enter numeric values between 1 and 10 |
| number | 5 |
| Premium_Amount | Please enter numeric values between 1401 and 215445 |
| number | 7608 |
| (query) | |
| Doctor_Visited_Before_Days | Please enter numeric values between 22 and 318 |
| number | 53 |
| (query) | |
| Number_of_Months_as_Customer | Please enter numeric values between 22 and 205 |
| number | 193 |
| (query) | |

Execute

## 3

**Curl**

```
curl -X POST "http://127.0.0.1:8000/predict?
Established_years=11&No_Times_Visited_Doctor_in_Last_Month=2&Average_Room_Rent=1097&No_Days_Stay_Recommended=5&Premium_Amount=7608&Doctor_Visited_Before_Days=53&Number_of_Mont
hs_as_Customer=193" -H  "accept: application/json"
```

**Server response**

| Code | Details |
|---|---|
| 200 | **Response body** |

```
----------------
Amount With Drawn 127012.488532534
----------------
```

**Response headers**

```
date: Sat, 13 Jul 2019 02:37:59 GMT, Sat, 13 Jul 2019 2:37:59 AM GMT
content-type: text/html; charset=utf-8
connection: close
content-length: 70
```

Responses

| Code | Description |
|---|---|
| 200 | Returns the Amount Withdrawn prediction from linear model |
| default | Default response. |

THANK YOU