

ISL Final Project Phase II

Pooria Assarehha

2025-02-06

```
library(ggplot2)
library(dplyr)
library(tidyr)
library(corrplot)
library(gridExtra)
library(tidyverse)
library(reshape2)
library(moments)

library(caret)

data_path <- "clean_data.csv"
data <- read.csv(data_path, stringsAsFactors = TRUE)
```

Data Preparation

The first two columns assume no role in our estimation, they can be omitted. After reading the file each feature must take its datatype by definition. Then we separate out target variable from predictive features.

```
# Remove the first two columns
data %>% select(!c(URL, Name)) -> data

# Convert binary columns to factors
for (col in names(data)){
  if (all(unique(data[,col]) == c(0,1)) || all(unique(data[,col]) == c(1,0)))
    data[,col] = factor(data[,col])
}

# Define the target variable
target <- data$Amtiaz

# Remove the target variable column from the features
features <- data %>% select(!Amtiaz)
```

Exploratory Data Analysis

Histograms for key numerical variables

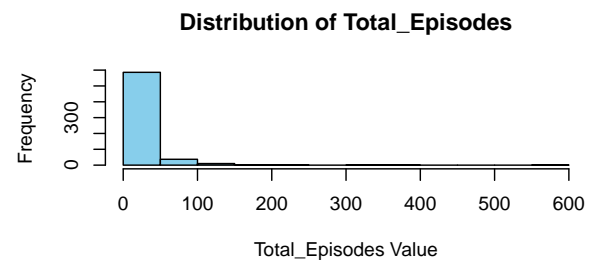
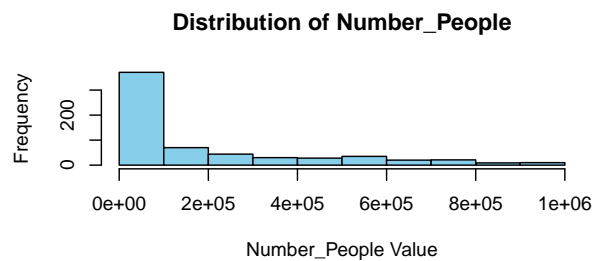
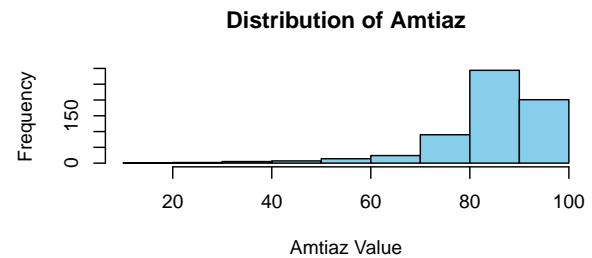
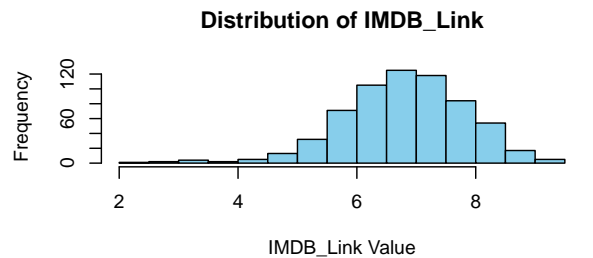
```
num_cols <- c("IMDB_Link", "Amtiaz", "Number_People", "Total_Episodes")

par(mfrow=c(2,2))
for (col in num_cols) {
  hist(
```

```

data[[col]],
main=paste("Distribution of", col),
xlab = paste(col, "Value"),
col="skyblue",
border="black")
}

```



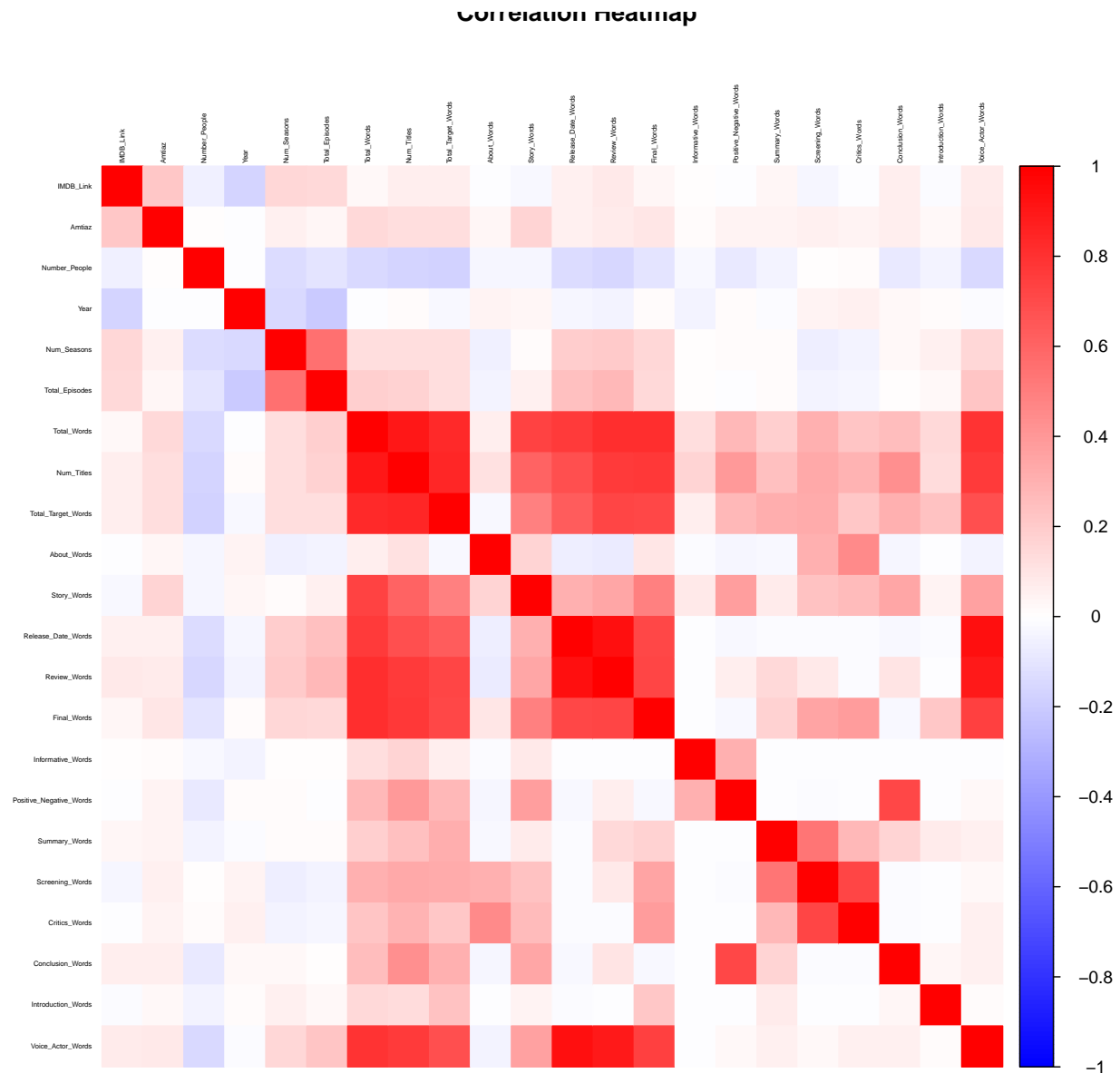
Correlation Heatmap

```

num_data <- data %>% select_if(is.numeric)
corr_matrix <- cor(num_data, use="complete.obs")

corrplot(
  corr_matrix,
  method="color",
  col=colorRampPalette(c("blue", "white", "red"))(200),
  tl.cex=0.35, tl.col="black",
  title="Correlation Heatmap")

```

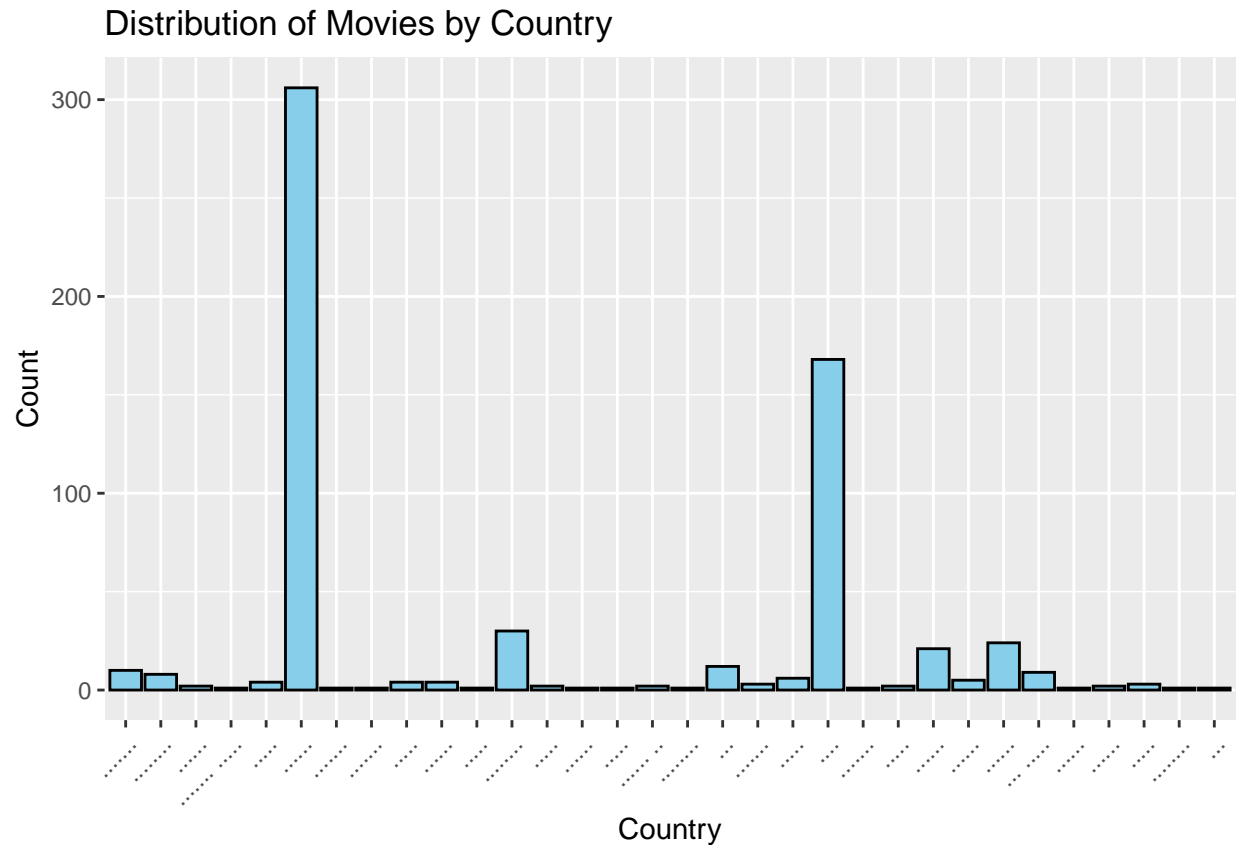


Strong correlations exist between IMDB_Link, Amtiaz, and Number_People, indicating possible relationships worth exploring in modeling.

Most of our features don't show any relation to target variable or any other feature.

Bar chart for Categorical Features, Country distribution

```
data %>%
  ggplot(aes(x=Country)) +
  geom_bar(fill="skyblue", color="black") +
  theme(axis.text.x = element_text(angle=45, hjust=1)) +
  labs(title="Distribution of Movies by Country", x="Country", y="Count")
```



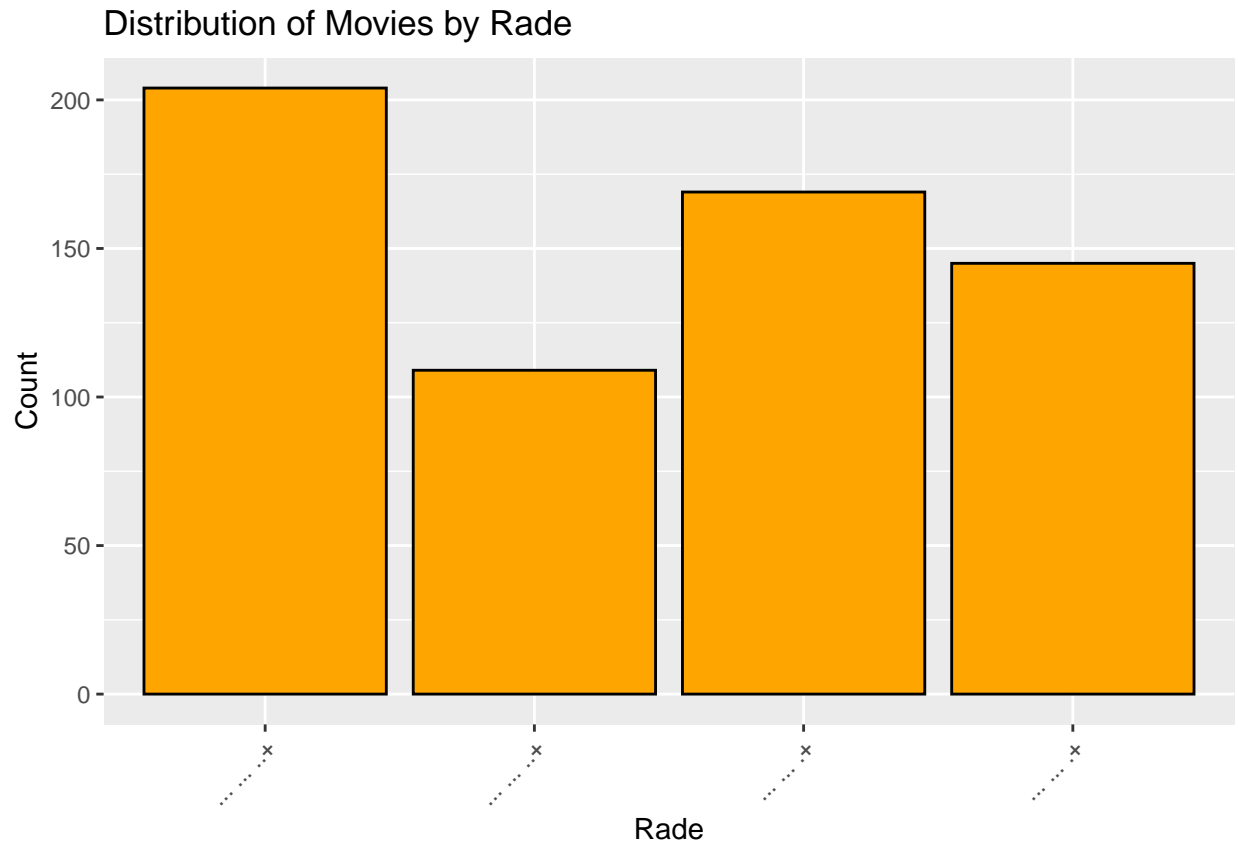
There are Countries that only appear once in our data, hence no inference or estimation can be done with them, let's simply omit those.

```
data %>%
  group_by(Country) %>%
  filter(n() > 1) -> data
nrow(data)
```

```
## [1] 627
```

Bar chart for Rade distribution

```
data %>%
  ggplot(aes(x=Rade)) +
  geom_bar(fill="orange", color="black") +
  theme(axis.text.x = element_text(angle=45, hjust=1)) +
  labs(title="Distribution of Movies by Rade", x="Rade", y="Count")
```



Data Overview: The dataset contains 638 rows and 61 columns. There are both numerical and categorical features. Columns like URL, Name, and Rade are categorical, while others like IMDB_Link, Amtiaz, and Number_People are numerical. The dataset has no missing values after preprocessing.

Numerical Features: IMDB_Link has values ranging from 2.1 to 9.3, with a mean of 6.84. Amtiaz ranges from 17 to 100, with a mean of 84.6. Number_People has high variance, ranging from 1,000 to 998,000. Year values range from 1940 to 2025, with most data points concentrated in recent years. Some numerical columns (like Total_Episodes) have skewed distributions, which might affect modeling.

Categorical Features: Country and Rade should be analyzed further with frequency counts. Many binary genre columns (e.g., Romance, SciFi, Anime) are mostly 0s, meaning most movies don't belong to these genres.

Metric Functions

We choose and define these functions to evaluate out models now on.

```
MAE <- function(model, x_test, y_test) mean(abs(predict(model, x_test) - y_test))
MSE <- function(model, x_test, y_test) mean((predict(model, x_test) - y_test)^2)
Rsqr <- function(model, x_test, y_test) 1 - sum((predict(model, x_test) - y_test)^2) / sum((y_test - mean(y_test))^2)
R2adj <- function(model, x_test, y_test) 1 - ((1 - Rsqr(model, x_test, y_test)) * (nrow(x_test) - 1) / (nrow(x_test) - ncol(x_test)))
```

Predicting Amtiaz

From EDA and corr plot, we know no specific feature that has strong linear correlation with our response/target variable **Amtiaz**. This means simple linear regression won't give us a great prediction.

Simple linear regression

```
set.seed(1)

n <- nrow(data)

train_idx <- sample(1:n, size = 0.9 * n)
test_idx <- setdiff(1:n, train_idx)

train_data <- data[train_idx, ]
test_data <- data[test_idx, ]

lm_model <- lm(Amtiaz~., data = train_data)
res = summary(lm_model)
AIC(lm_model)
```

```
## [1] 4239.751
```

As of Linear Model summary, we see despite having many features, only 5 prove meaningful and there are a lot of features/parameters.

Feature Selection

Stepwise subtest selection

```
#res_step = step(lm_model, direction = 'both')
best_step_lm <- lm(formula = Amtiaz ~ IMDB_Link + Country + Rade + Is_Doblele +
  Story_Words + Series + Adventure + Comedy + Family + Action +
  ShortFilm + Korean, data = train_data)
res <- summary(best_step_lm)
```

```
results <- c(
  mean(abs(best_step_lm$residuals)),
  MAE(best_step_lm, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
  mean(best_step_lm$residuals^2),
  MSE(best_step_lm, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
  res$r.squared,
  Rsq(best_step_lm, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
  res$adj.r.squared,
  R2adj(best_step_lm, test_data %>% select(!'Amtiaz'), test_data$Amtiaz)
)
```

```
results <- round(results, 2)
```

```
cat(paste(results, collapse = " | "))
```

```
## 6.17 | 6.52 | 83.96 | 78.9 | 0.31 | 0.11 | 0.27 | -10.06
```

Using Lasso Selection

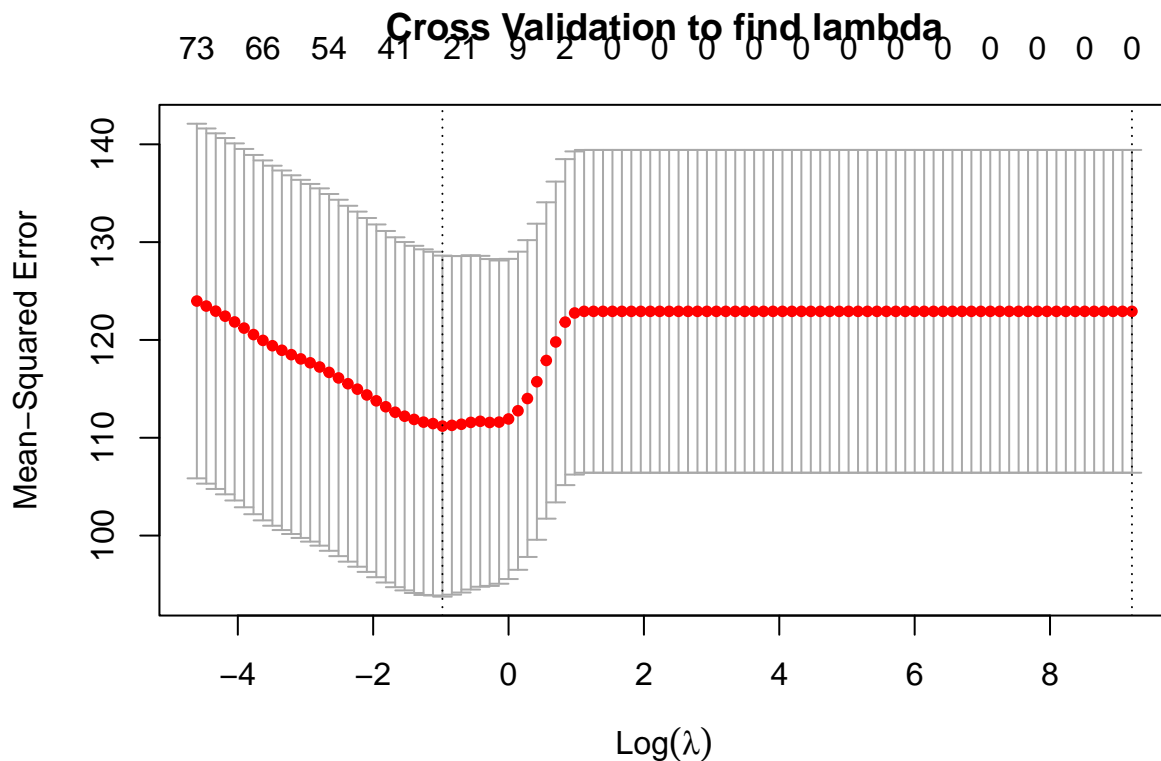
Lasso penalization can be used to select features.

```
library(glmnet)
```

```
## Loading required package: Matrix
##
## Attaching package: 'Matrix'
## The following objects are masked from 'package:tidyr':
##
##     expand, pack, unpack
## Loaded glmnet 4.1-8
x <- model.matrix(Amtiaz ~ ., data)[, -1]
y <- data$Amtiaz
x_train <- x[train_idx, ]
x_test <- x[test_idx, ]
y_train <- y[train_idx]
y_test <- y[test_idx]

lasso_cv <- cv.glmnet(
  x_train, y_train, alpha = 1, # Indicating Lasso
  lambda = 10^seq(4, -2, length = 100)
)

plot(lasso_cv, main = "Cross Validation to find lambda")
```



```
best_lambda_lasso <- lasso_cv$lambda.min

cat("Optimal Lambda for Lasso: ", best_lambda_lasso, "\n")
```

```
## Optimal Lambda for Lasso: 0.3764936
lasso_model <- glmnet(x_train, y_train, alpha = 1, lambda = best_lambda_lasso)

y_pred <- predict(lasso_model, s = best_lambda_lasso, newx = x_test)

results <- c(
  MAE(lasso_model, x_train, y_train),
  MAE(lasso_model, x_test, y_test),
  MSE(lasso_model, x_train, y_train),
  MSE(lasso_model, x_test, y_test),
  Rsq(lasso_model, x_train, y_train),
  Rsq(lasso_model, x_test, y_test),
  R2adj(lasso_model, x_train, y_train),
  R2adj(lasso_model, x_test, y_test)
)

results <- round(results, 2)

cat(paste(results, collapse = " | "))
```

```
## 6.2 | 5.91 | 89.14 | 72.77 | 0.27 | 0.18 | 0.14 | 2.89
```

```
results <- c(
  mean(abs(lm_model$residuals)),
  MAE(lm_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
  mean(lm_model$residuals^2),
  MSE(lm_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
  res$r.squared,
  Rsq(lm_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
  res$adj.r.squared,
  R2adj(lm_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz)
)

results <- round(results, 2)

cat(paste(results, collapse = " | "))
```

```
## 6 | 6.9 | 81.68 | 87.44 | 0.31 | 0.01 | 0.27 | -11.25
```

Model	train MAE	test MAE	train MSE	test MSE	train R^2	test R^2	train Adjusted R^2	test Adjusted R^2
LinReg	6	6.9	81.68	87.44	0.33	0.01	0.23	-11.25
best_step	6.17	6.52	83.96	78.9	0.31	0.11	0.27	-10.06
best_lasso	6.55	6.06	99.04	80.52	0.19	0.09	0.04	3.09

As we see our Linear models (Comparing R^2) are doing no better job than the “Mean predictor” (mean response is the prediction for all). This means features are not predicting the response. So far our models ignored feature interactions, we can turn to models that include interactions well like trees. We know bagging can reduce the variance of trees and boosting can reduce bias.

XG Boost

```
library(xgboost)
```

```
##
```

```
## Attaching package: 'xgboost'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
## slice
```

```
#tuning parameters nrounds(number of repetitions), eta(learning rate), max_depth(trees depth), gamma(mi
```

```
grid <- expand.grid(  
  nrounds = c(50, 100, 150),  
  eta = c(0.01, 0.1, 0.3),  
  max_depth = c(3, 6, 9),  
  gamma = c(0, 1, 5),  
  colsample_bytree = c(0.5, 0.7, 1),  
  min_child_weight = c(1, 3, 5),  
  subsample = c(0.6, 0.8, 1)  
)
```

```
#10-fold cross-validation
```

```
train_control <- trainControl(method = "cv", number = 10)
```

```
#xgb_tuned <- train(x = as.matrix(x_train), y = y_train, method = "xgbTree", trControl = train_control,
```

```
#best_params <- xgb_tuned$bestTune
```

```
#cat("Optimal Parameters for XG Boost : ", paste(best_params, collapse = ", "), "\n")
```

```
cat("Optimal Parameters for XG Boost : ", "50, 3, 0.1, 5, 0.5, 5, 1", "\n")
```

```
## Optimal Parameters for XG Boost : 50, 3, 0.1, 5, 0.5, 5, 1
```

```
xgb_model <- xgboost(data = x_train,  
  label = y_train,  
    nrounds = 50, #best_params$nrounds,  
    eta = 0.1, #best_params$eta,  
    max_depth = 3, #best_params$max_depth,  
  min_child_weight = 5, #best_params$min_child_weight,  
    subsample = 0.5, #best_params$subsample,  
  colsample_bytree = 1, #best_params$colsample_bytree,  
  objective = "reg:squarederror")
```

```
## [1] train-rmse:76.361051
```

```
## [2] train-rmse:68.919525
```

```
## [3] train-rmse:62.207770
```

```
## [4] train-rmse:56.180507
```

```
## [5] train-rmse:50.782654
```

```
## [6] train-rmse:45.963207
```

```
## [7] train-rmse:41.708360
```

```
## [8] train-rmse:37.933850
```

```
## [9] train-rmse:34.520527
```

```
## [10] train-rmse:31.426356
```

```
## [11] train-rmse:28.691359
```

```
## [12] train-rmse:26.154651
```

```
## [13] train-rmse:23.966591
```

```
## [14] train-rmse:22.014074
```

```
## [15] train-rmse:20.272872
```

```
## [16] train-rmse:18.721471
## [17] train-rmse:17.396517
## [18] train-rmse:16.169573
## [19] train-rmse:15.093156
## [20] train-rmse:14.143467
## [21] train-rmse:13.351637
## [22] train-rmse:12.665086
## [23] train-rmse:12.028746
## [24] train-rmse:11.466131
## [25] train-rmse:11.011587
## [26] train-rmse:10.660174
## [27] train-rmse:10.322994
## [28] train-rmse:10.037061
## [29] train-rmse:9.811565
## [30] train-rmse:9.602965
## [31] train-rmse:9.416670
## [32] train-rmse:9.264992
## [33] train-rmse:9.170243
## [34] train-rmse:9.027349
## [35] train-rmse:8.932723
## [36] train-rmse:8.841477
## [37] train-rmse:8.792340
## [38] train-rmse:8.729397
## [39] train-rmse:8.667685
## [40] train-rmse:8.615669
## [41] train-rmse:8.566691
## [42] train-rmse:8.506988
## [43] train-rmse:8.449321
## [44] train-rmse:8.409287
## [45] train-rmse:8.381865
## [46] train-rmse:8.340125
## [47] train-rmse:8.316202
## [48] train-rmse:8.262338
## [49] train-rmse:8.229597
## [50] train-rmse:8.156256
```

```
results <- c(
MAE( xgb_model, x_train, y_train),
MAE( xgb_model, x_test, y_test),
MSE( xgb_model, x_train, y_train),
MSE( xgb_model, x_test, y_test),
Rsqr( xgb_model, x_train, y_train),
Rsqr( xgb_model, x_test, y_test),
R2adj(xgb_model, x_train, y_train),
R2adj(xgb_model, x_test, y_test)
)

results <- round(results, 2)

cat(paste(results, collapse = " | "))
```

```
## 5.42 | 6.47 | 66.52 | 78.38 | 0.46 | 0.11 | 0.35 | 3.03
```

Optimal Parameters for XG Boost : 50, 3, 0.1, 5, 0.5, 5, 1 MSE for XG Boost : 111.038 MAE for XG Boost : 7.076772 R2 for XG Boost : 0.1469957

Model	train MAE	test MAE	train MSE	test MSE	train R^2	test R^2	train Adjusted R^2	test Adjusted R^2
LinReg	6	6.9	81.68	87.44	0.33	0.01	0.23	-11.25
best_step	6.17	6.52	83.96	78.9	0.31	0.11	0.27	-10.06
best_lasso	6.55	6.06	99.04	80.52	0.19	0.09	0.04	3.09
XGBoost	4.47	6.12	46.26	67.16	0.62	0.24	0.55	2.74

Significant Improvement from XGBoost currently the best candidate.

Random forest No Boost

```
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 4.4.2
```

```
## randomForest 4.7-1.2
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##
```

```
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:gridExtra':
```

```
##
```

```
##      combine
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
##      combine
```

```
## The following object is masked from 'package:ggplot2':
```

```
##
```

```
##      margin
```

```
library(tidyverse)
```

```
library(caret)
```

```
rf_model <- randomForest(Amtiaz ~ ., data = train_data, ntree = 100, mtry = 13, importance = TRUE)
```

```
print(rf_model)
```

```
##
```

```
## Call:
```

```
## randomForest(formula = Amtiaz ~ ., data = train_data, ntree = 100, mtry = 13, importance = TRUE)
```

```
##      Type of random forest: regression
```

```
##      Number of trees: 100
```

```
## No. of variables tried at each split: 13
```

```
##
```

```
##      Mean of squared residuals: 103.1444
```

```
##      % Var explained: 15.76
```

```
results <- c(
```

```
MAE( rf_model, train_data %>% select(!'Amtiaz'), train_data$Amtiaz),
```

```
MAE( rf_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
```

```
MSE( rf_model, train_data %>% select(!'Amtiaz'), train_data$Amtiaz),
```

```
MSE( rf_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
```

```
Rsq( rf_model, train_data %>% select(!'Amtiaz'), train_data$Amtiaz),
```

```

Rsq( rf_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz),
R2adj(rf_model, train_data %>% select(!'Amtiaz'), train_data$Amtiaz),
R2adj(rf_model, test_data %>% select(!'Amtiaz'), test_data$Amtiaz)
)

```

```
results <- round(results, 2)
```

```
cat(paste(results, collapse = " | "))
```

```
## 3 | 5.79 | 21.53 | 69.02 | 0.82 | 0.22 | 0.8 | -8.67
```

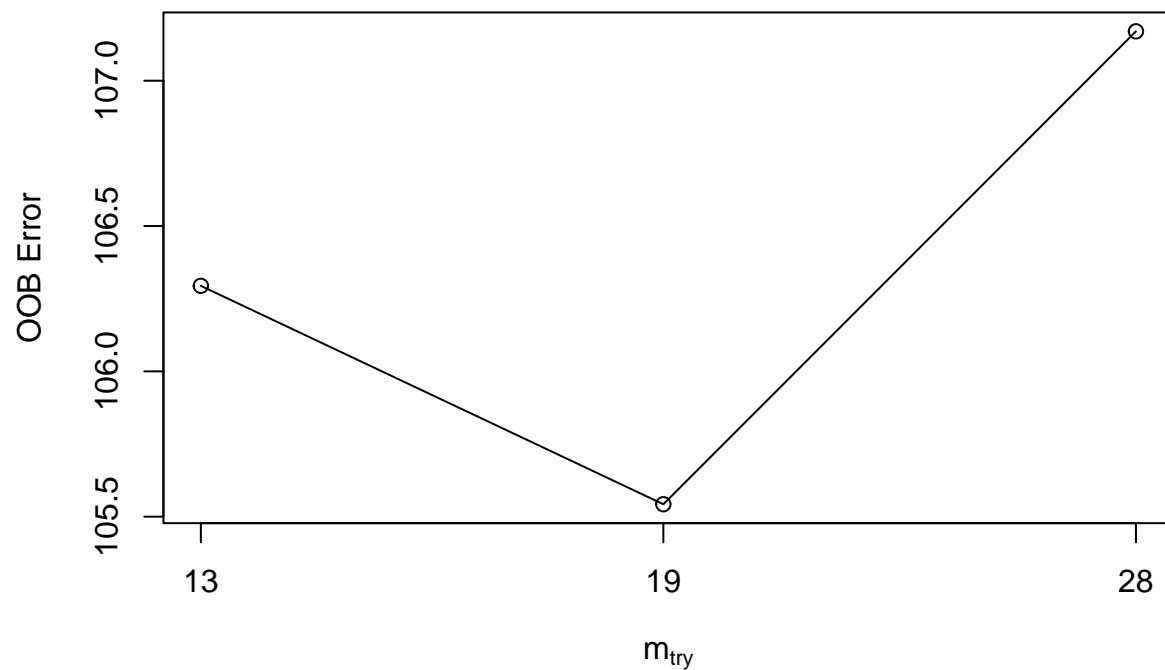
Model	train MAE	test MAE	train MSE	test MSE	train R^2	test R^2	train Adjusted R^2	test Adjusted R^2
LinReg	6	6.9	81.68	87.44	0.33	0.01	0.23	-11.25
best_step	6.17	6.52	83.96	78.9	0.31	0.11	0.27	-10.06
best_lasso	6.55	6.06	99.04	80.52	0.19	0.09	0.04	3.09
XGBoost	4.47	6.12	46.26	67.16	0.62	0.24	0.55	2.74
RandomFrst	5.44	6.02	69.57	76.25	0.43	0.14	0.37	-9.69
tuned_RF	2.97	5.78	21.74	66.29	0.82	0.25	0.8	-8.29

```
tuned_rf <- tuneRF(train_data[-which(names(train_data) == "Amtiaz")], train_data$Amtiaz, stepFactor = 1
```

```

## mtry = 19  OOB error = 105.5428
## Searching left ...
## mtry = 13  OOB error = 106.2943
## -0.007119729 0.01
## Searching right ...
## mtry = 28  OOB error = 107.1698
## -0.01541572 0.01

```



```
print(tuned_rf)
```

```
##      mtry OOBError
## 13      13 106.2943
## 19      19 105.5428
## 28      28 107.1698
```

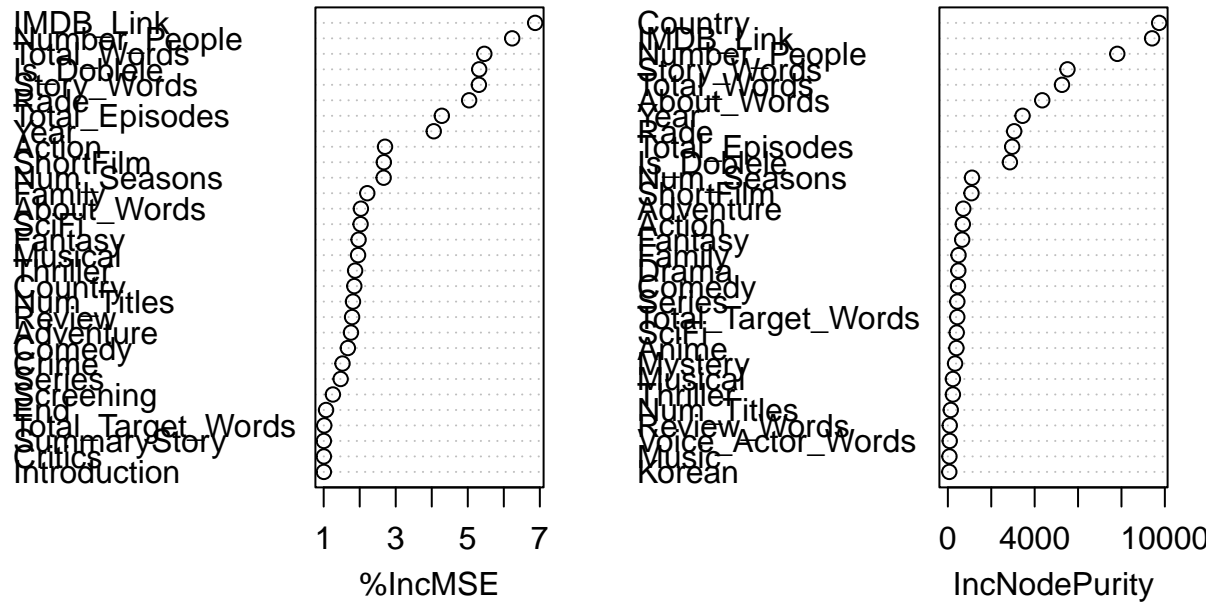
```
importance(rf_model)
```

```
##              %IncMSE IncNodePurity
## IMDB_Link      6.86980923  9411.9984239
## Number_People  6.23264181  7806.7938056
## Country        1.84956446  9736.8278368
## Year           4.05587350  3442.5423781
## Rade           5.03762175  3056.8343337
## Num_Seasons    2.66482108  1113.0230016
## Total_Episodes 4.27983686  2963.2622250
## Publication    -0.98859051   11.9769972
## VoiceActors    -0.79659928   40.5119817
## Review         1.78865762    9.0328034
## Tips          -0.97352916   13.4917334
## End            1.06564350   15.8053974
## Description     0.00000000    4.9207143
## Characters      0.00000000   11.4865656
## InformativeMessages 0.00000000  0.0000000
## PositiveAndNegative -0.29170718  3.3603283
## SummaryStory    1.00503782   1.2293843
```

## Screening	1.25604068	1.7743589
## Critics	1.00503782	9.6693790
## Conclusion	0.40072835	6.8573436
## Introduction	1.00503782	2.3618634
## Total_Words	5.45865211	5259.7008949
## Num_Titles	1.81379278	134.4616073
## Is_Doblele	5.31713832	2858.1818117
## Total_Target_Words	1.01584488	434.3939671
## About_Words	2.03293047	4349.3516623
## Story_Words	5.30784455	5511.9417265
## Release_Date_Words	-0.31715812	41.3675942
## Review_Words	-1.55472228	87.2551790
## Final_Words	-0.37695609	43.8995272
## Informative_Words	0.00000000	1.4948235
## Positive_Negative_Words	-2.74759434	9.4112994
## Summary_Words	0.00000000	0.8829313
## Screening_Words	-0.61840898	3.0183040
## Critics_Words	1.00503782	3.8019608
## Conclusion_Words	0.09299504	19.2463358
## Introduction_Words	0.00000000	2.6948876
## Voice_Actor_Words	0.83594485	85.9195272
## Series	1.47094883	437.0882420
## Animation	0.00000000	3.1083716
## Western	0.00000000	19.3839333
## Adventure	1.75718091	711.1920082
## Comedy	1.67118690	466.0390808
## Family	2.21042966	491.3513558
## Fantasy	1.96883284	658.7675112
## Mystery	0.56229825	331.6981100
## Action	2.70214819	690.2429097
## Romance	0.53561455	48.1307434
## Drama	0.68650845	482.0666584
## SciFi	2.02440705	408.8109623
## ShortFilm	2.67019749	1088.7815991
## Crime	1.52287050	50.9020843
## Musical	1.95431528	233.0723932
## Korean	-0.07271798	66.3105861
## Thriller	1.87300207	233.0573546
## Anime	-0.47305956	395.3543319
## Music	-1.33364117	73.1127324

```
varImpPlot(rf_model)
```

rf_model



SVR

This Model needs separate data prep

```
library(e1071)

## Warning: package 'e1071' was built under R version 4.4.2
##
## Attaching package: 'e1071'
## The following objects are masked from 'package:moments':
##
##      kurtosis, moment, skewness

library(caret)
library(dplyr)

data = 'clean_data.csv'
data <- read.csv(data, stringsAsFactors = TRUE)

data %>%
  group_by(Country) %>%
  filter(n() > 1) -> data
nrow(data)

## [1] 627
```

```
# Remove the first two columns
data %>% select(!c(URL, Name)) -> data
```

```
# Convert columns
data$Country <- as.factor(data$Country)
data$Rade <- as.factor(data$Rade)
data$Amtiaz <- as.numeric(data$Amtiaz)
data$Year <- as.numeric(data$Year)
data$IMDB_Link <- as.numeric(data$IMDB_Link)
# Convert categorical variables to dummy variables
data <- dummyVars(~ ., data = data) %>% predict(data) %>% as.data.frame()
```

Split the data into features and target

```
target <- "Amtiaz"
predictors <- setdiff(names(data), target)
```

train,test

```
set.seed(1)
train_index <- sample(1:nrow(data), size = 0.7 * nrow(data))
svr_train_data <- data[train_index, ]
svr_test_data <- data[-train_index, ]
```

```
# Scale numerical features
preproc <- preProcess(svr_train_data[, predictors], method = c("center", "scale"))
train_data_scaled <- predict(preproc, svr_train_data)
test_data_scaled <- predict(preproc, svr_test_data)
```

```
train_data_scaled$Amtiaz <- svr_train_data$Amtiaz
test_data_scaled$Amtiaz <- svr_test_data$Amtiaz
```

```
svr_model <- svm(Amtiaz ~ ., data = train_data_scaled, kernel = "radial", cost = 1, gamma = 0.1)
```

```
plot(svr_model, train_data_scaled)
summary(svr_model)
```

```
##
## Call:
## svm(formula = Amtiaz ~ ., data = train_data_scaled, kernel = "radial",
##      cost = 1, gamma = 0.1)
##
##
## Parameters:
##   SVM-Type:  eps-regression
##   SVM-Kernel: radial
##      cost:    1
##     gamma:   0.1
##   epsilon:   0.1
##
##
## Number of Support Vectors:  433
```

```
results <- c(
MAE( svr_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
MAE( svr_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz),
```



```

MSE( svr_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
MSE( svr_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz),
Rsqr( svr_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
Rsqr( svr_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz),
R2adj(svr_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
R2adj(svr_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz)
)

```

```
results <- round(results, 2)
```

```
cat(paste(results, collapse = " | "))
```

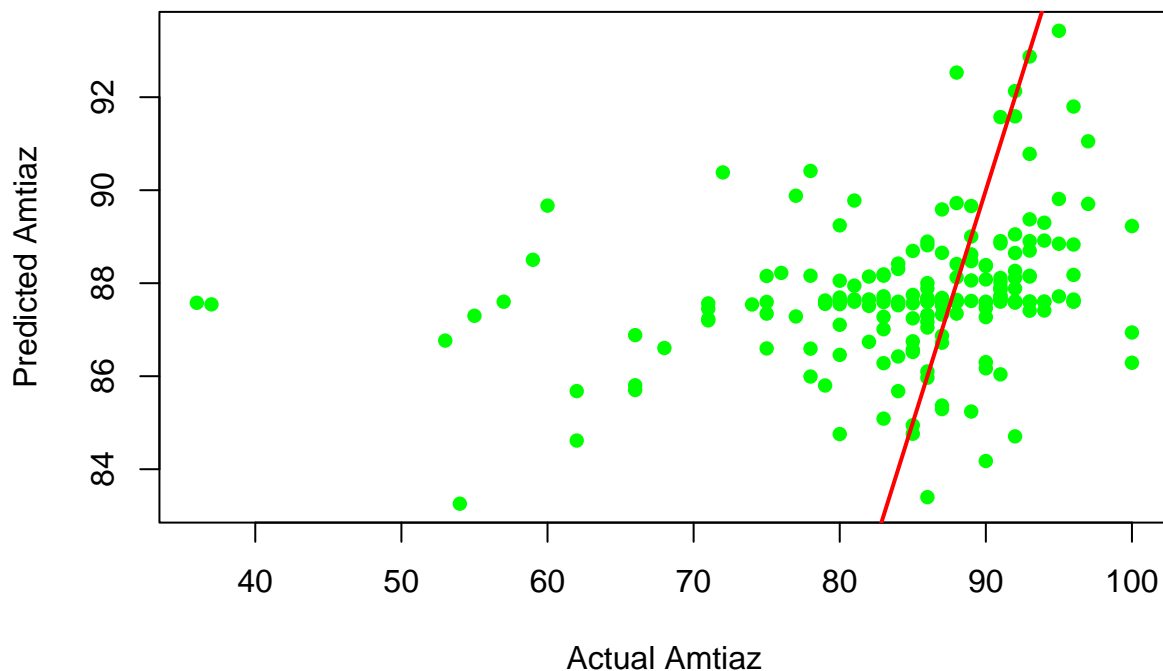
```
## 5.82 | 6.43 | 118.31 | 105.28 | 0.07 | -0.03 | -0.18 | -0.99
```

```
predictions <- predict(svr_model, test_data_scaled)
```

```
# Plot actual vs predicted values
```

```
plot(test_data_scaled$Amtiaz, predictions,
      xlab = "Actual Amtiaz", ylab = "Predicted Amtiaz",
      main = "SVR Predictions vs. Actual Values",
      col = "green", pch = 16)
abline(0, 1, col = "red", lwd = 2)
```

SVR Predictions vs. Actual Values



```

tuned <- tune(svm, Amtiaz ~ ., data = train_data_scaled, kernel = "radial", ranges = list(cost = c(0.1, 1),
best_model <- tuned$best.model
summary(best_model)

```

```
##
```

```
## Call:
## best.tune(METHOD = svm, train.x = Amtiaz ~ ., data = train_data_scaled,
##           ranges = list(cost = c(0.1, 1, 10), gamma = c(0.01, 0.1, 1)),
##           kernel = "radial")
##
##
## Parameters:
##   SVM-Type:  eps-regression
##   SVM-Kernel: radial
##     cost:    10
##     gamma:   0.01
##     epsilon: 0.1
##
##
## Number of Support Vectors: 434

results <- c(
MAE( best_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
MAE( best_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz),
MSE( best_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
MSE( best_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz),
Rsqr( best_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
Rsqr( best_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz),
R2adj(best_model, train_data_scaled %>% select(!'Amtiaz'), train_data_scaled$Amtiaz),
R2adj(best_model, test_data_scaled %>% select(!'Amtiaz'), test_data_scaled$Amtiaz)
)

results <- round(results, 2)

cat(paste(results, collapse = " | "))

## 4.55 | 6.39 | 91.22 | 97.47 | 0.28 | 0.05 | 0.09 | -0.85
```

Model	train MAE	test MAE	train MSE	test MSE	train R^2	test R^2	train Adjusted R^2	test Adjusted R^2
LinReg	6	6.9	81.68	87.44	0.33	0.01	0.23	-11.25
best_step	6.17	6.52	83.96	78.9	0.31	0.11	0.27	-10.06
best_lasso	6.55	6.06	99.04	80.52	0.19	0.09	0.04	3.09
XGBoost	4.47	6.12	46.26	67.16	0.62	0.24	0.55	2.74
RandomFrst	5.44	6.02	69.57	76.25	0.43	0.14	0.37	-9.69
tuned_RF	2.97	5.78	21.74	66.29	0.82	0.25	0.8	-8.29
SVR	5.82	6.43	118.31	105.28	0.07	-0.03	-0.18	-0.99
tuned SVR	4.55	6.39	91.22	97.47	0.28	0.05	0.09	-0.85

Neural Network

An extensive notebook on fitting a neural network is given in python.

```
library(keras)
NN_model = keras_model_sequential() %>%
  layer_dense(units = 128, activation = "relu", input_shape = dim(x_train)[2]) %>%
  layer_dense(units = 60, activation = "relu",) %>%
  layer_dense(units = 15, activation = "relu",) %>%
  layer_dense(units = 1)
```

```

NN_model %>% compile(
  optimizer = "adam",
  loss = "mse"
)

summary(NN_model)

history = NN_model %>% fit(
  x_train, y_train,
  epochs = 50, batch_size = 16,
  validation_data = list(x_test, y_test),
)

```

```

results <- c(
  MAE( NN_model, x_train, y_train),
  MAE( NN_model, x_test, y_test),
  MSE( NN_model, x_train, y_train),
  MSE( NN_model, x_test, y_test),
  Rsq( NN_model, x_train, y_train),
  Rsq( NN_model, x_test, y_test),
  R2adj(NN_model, x_train, y_train),
  R2adj(NN_model, x_test, y_test)
)

```

```

## 18/18 - 0s - 115ms/epoch - 6ms/step
## 2/2 - 0s - 22ms/epoch - 11ms/step
## 18/18 - 0s - 33ms/epoch - 2ms/step
## 2/2 - 0s - 21ms/epoch - 10ms/step
## 18/18 - 0s - 33ms/epoch - 2ms/step
## 2/2 - 0s - 20ms/epoch - 10ms/step
## 18/18 - 0s - 34ms/epoch - 2ms/step
## 2/2 - 0s - 21ms/epoch - 11ms/step

```

```
results <- round(results, 2)
```

```
cat(paste(results, collapse = " | "))
```

```
## 178.09 | 195.45 | 86512.06 | 78718.66 | -705.6 | -888.76 | -838.28 | 2044.14
```

Model	train MAE	test MAE	train MSE	test MSE	train R^2	test R^2	train Adjusted R^2	test Adjusted R^2
LinReg	6	6.9	81.68	87.44	0.33	0.01	0.23	-11.25
best_step	6.17	6.52	83.96	78.9	0.31	0.11	0.27	-10.06
best_lasso	6.55	6.06	99.04	80.52	0.19	0.09	0.04	3.09
XGBoost	4.47	6.12	46.26	67.16	0.62	0.24	0.55	2.74
RandomFrst	5.44	6.02	69.57	76.25	0.43	0.14	0.37	-9.69
tuned_RF	2.97	5.78	21.74	66.29	0.82	0.25	0.8	-8.29
SVR	5.82	6.43	118.31	105.28	0.07	-0.03	-0.18	-0.99
tuned SVR	4.55	6.39	91.22	97.47	0.28	0.05	0.09	-0.85
Neural Net	12.13	13.29	266.01	301.92	-1.17	-2.41	-1.58	8.84

Neural Net when not tuned, performs worse.