# Video Game Sales Prediction

**By:**

**Hari Krishna Rangineni**

**Likith Kumar Miryala**

**Sai Kiran Putta**

# Scope of the Project

➢ **Models performance on the dataset**

➢ **Prove that feature generation from "Neural Networks" are superior!**

# Feature Set

- Name
- Platform
- Year of Release
- Genre
- Critic Score
- User Score
- Critic Count
- User Count
- Global Sales
- Rating
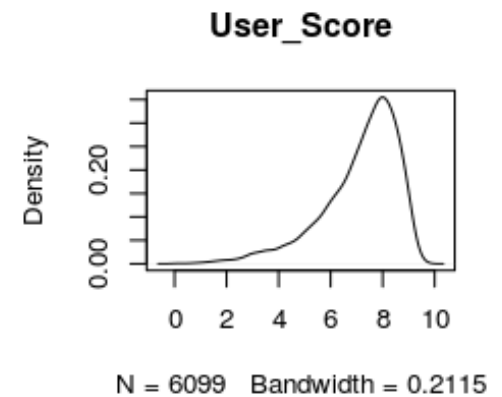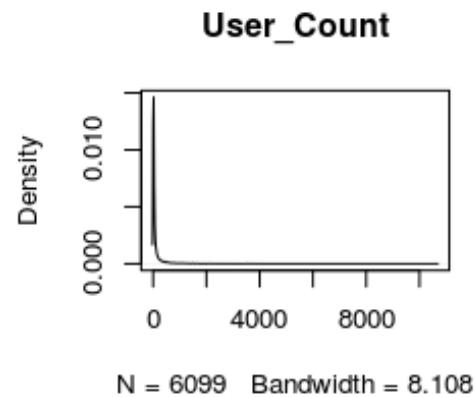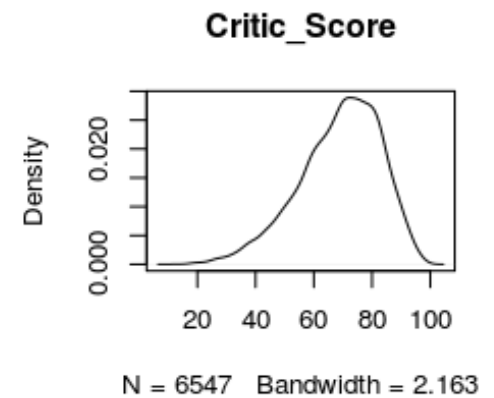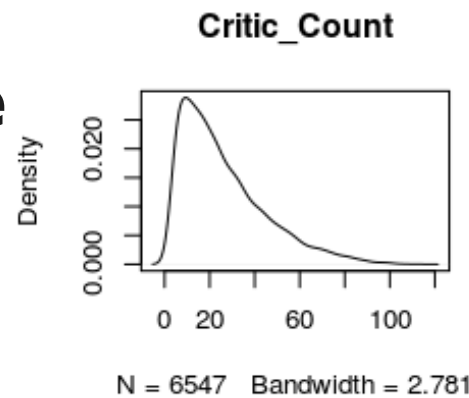
- JP sales
- Other sales

## Deleted Features:

- EU Sales
- Developer
- Publisher

# Mean Vs Median for Imputation

➢**Taking a Median is better than Mean since the distributions are skewed.**



Critic_Count
N = 6547   Bandwidth = 2.781

Critic_Score
N = 6547   Bandwidth = 2.163

User_Count
N = 6099   Bandwidth = 8.108

User_Score
N = 6099   Bandwidth = 0.2115

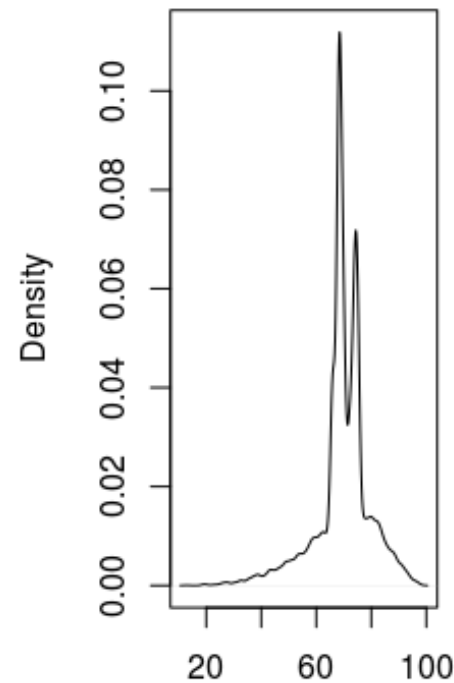# NA Imputation Strategy

➢ **Imputation by Grouping:**

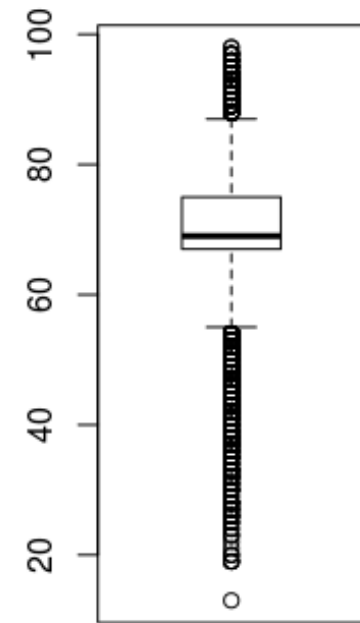**Grouping by Genre and taking a median of feature.**

# Outliers In Data

➢ **If we go by boxplots, most of the data will get deleted.**


Density of Critic_Score
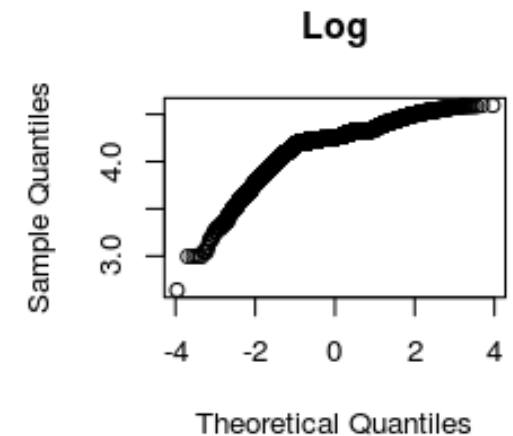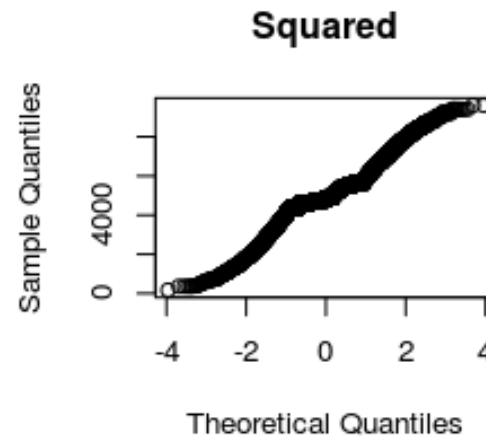

Boxplot of Critic_Score

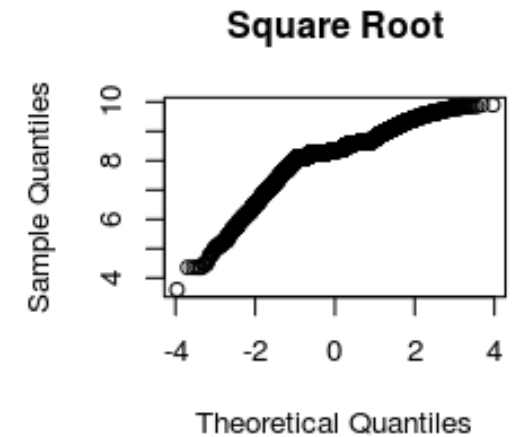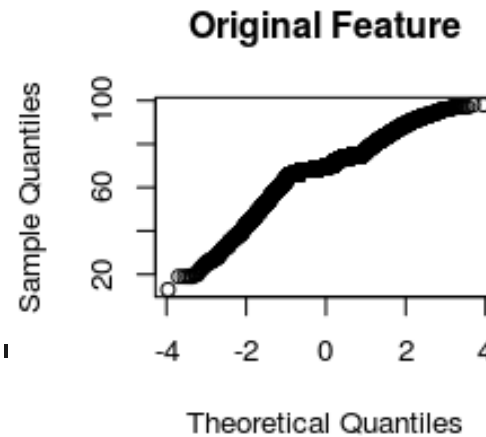# Outliers In Data - 2

## Our approach:

➤ **Do feature transformation first, then remove outliers.**

# Feature Engineering.

➢ **Length of Game Name**

➢ **Length of Publisher Name**

➢ **In Platform: X360 / XB / XOne  -> Xbox**

**PS / PS2 / PS3 / PS4 / PSP / PSV -> PS etc.**

➢**Year by different periods**

**YearI  -> 1980 – 2015 by step of 5 years**

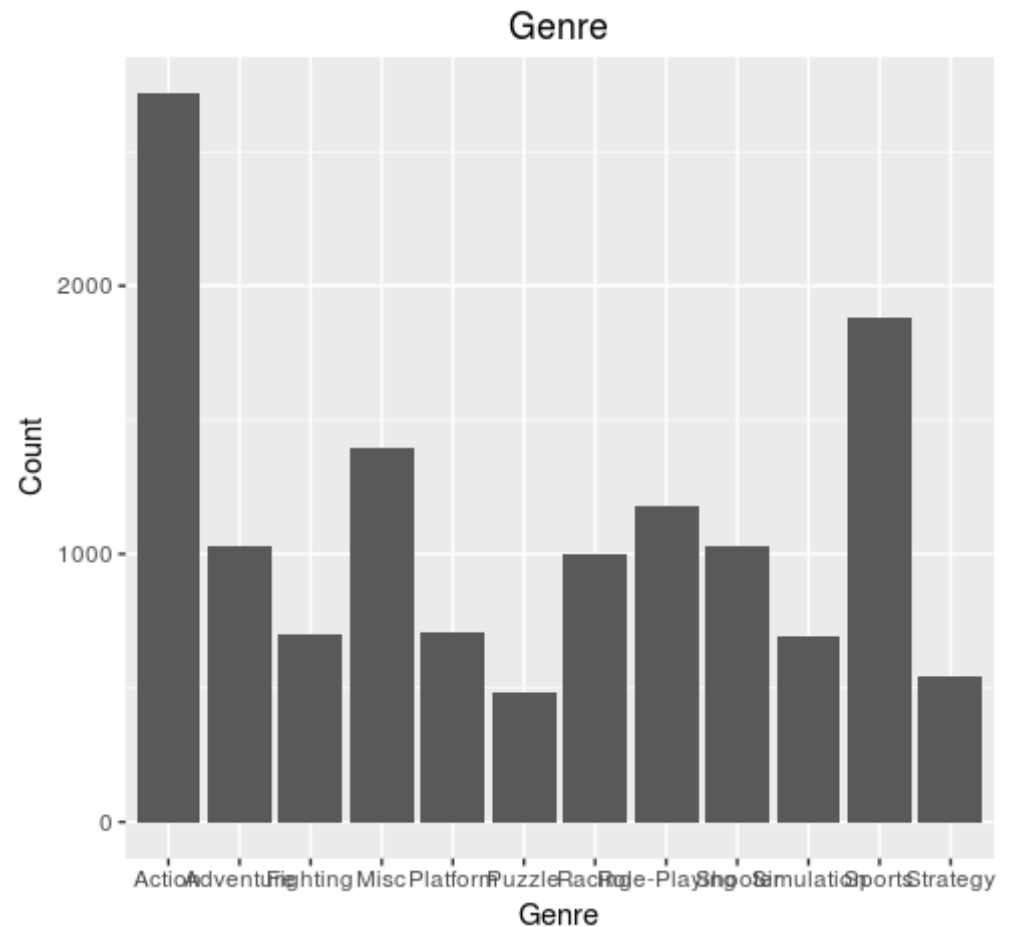**YearII -> 1980 – 2015 by step of 10 years**

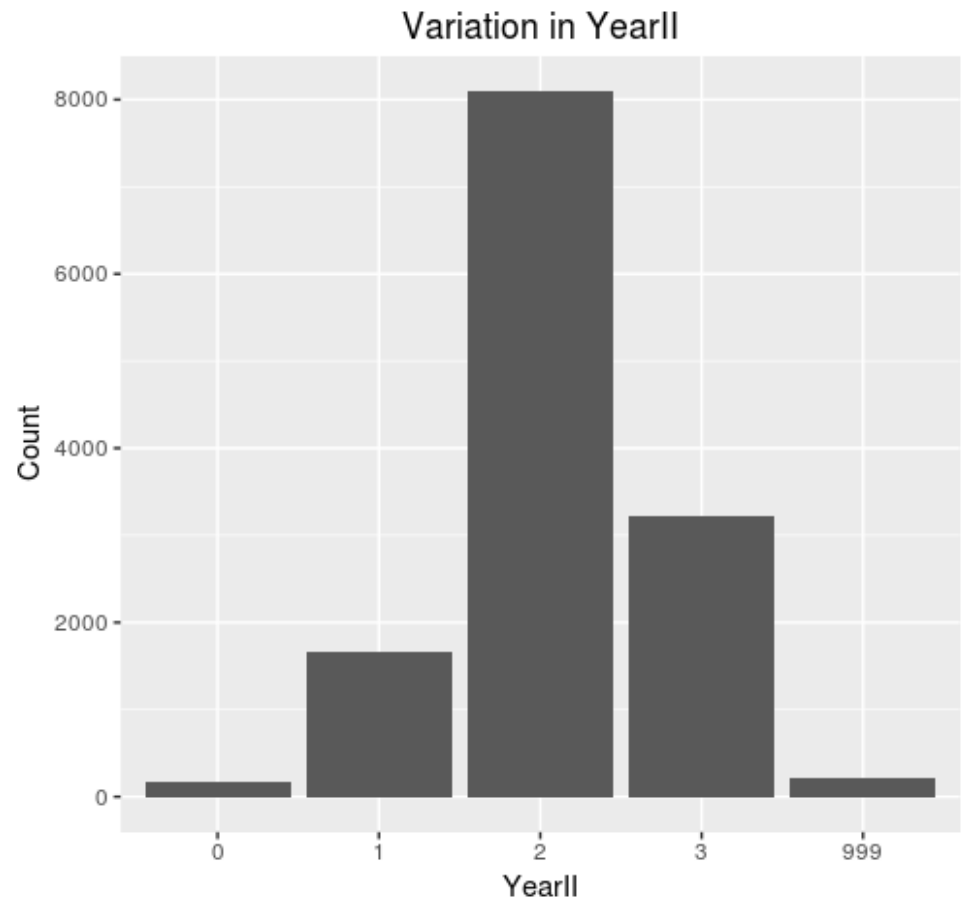# Exploratory Data Analysis.

# Some Interesting Plots

**Inference: More number of Action, Sports related games.**

# Some Interesting Plots - 2

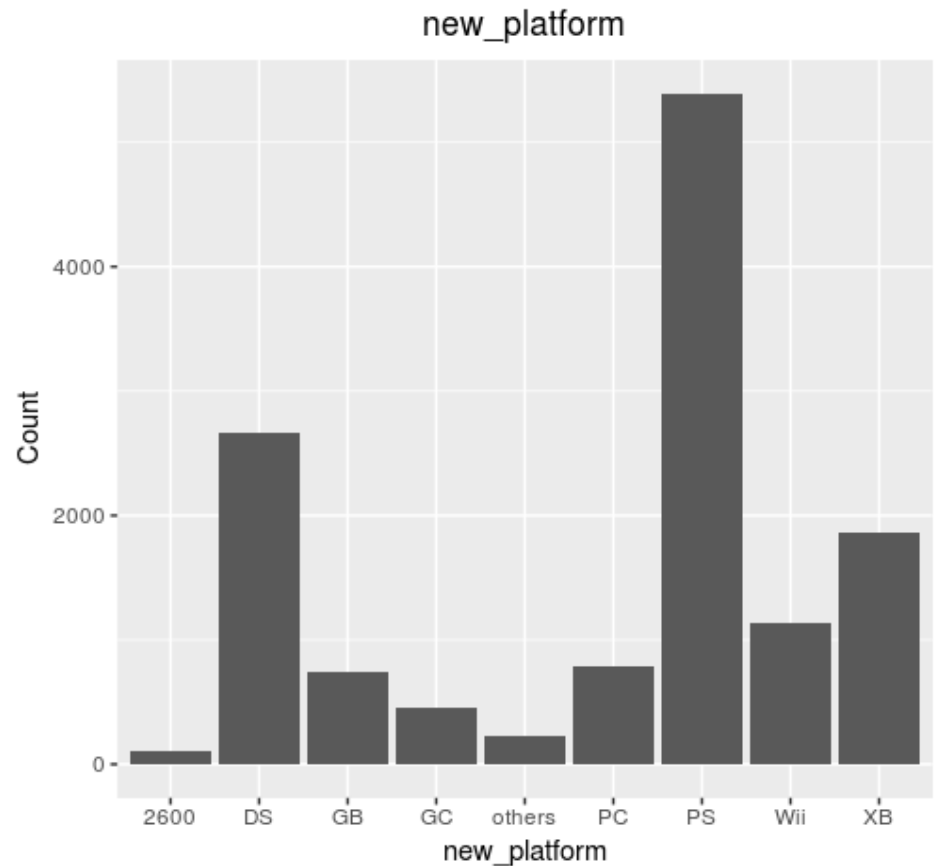**Inference:** **More number of games were released in 2000 – 2010**

**\* 999 is category where we don't know the year.**
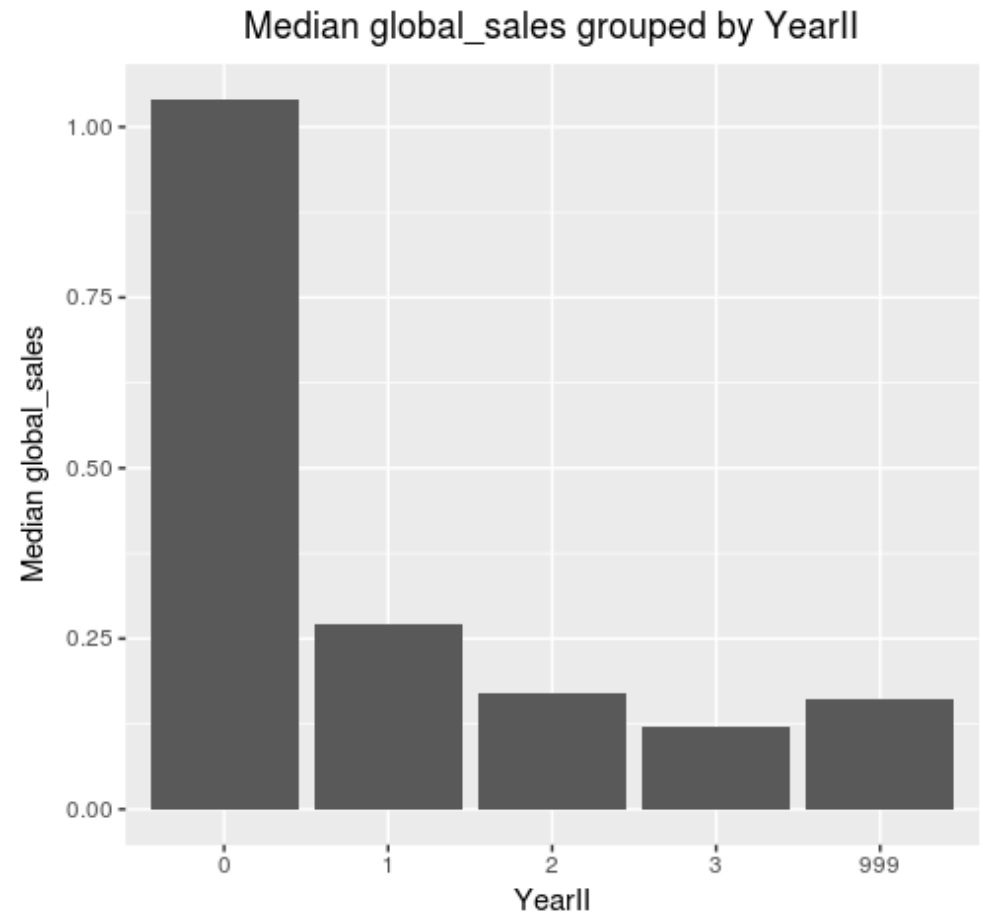

Variation in YearII

# Some Interesting Plots - 3

**Inference: More number of games were released in PlayStation followed by Nintendo DS.**
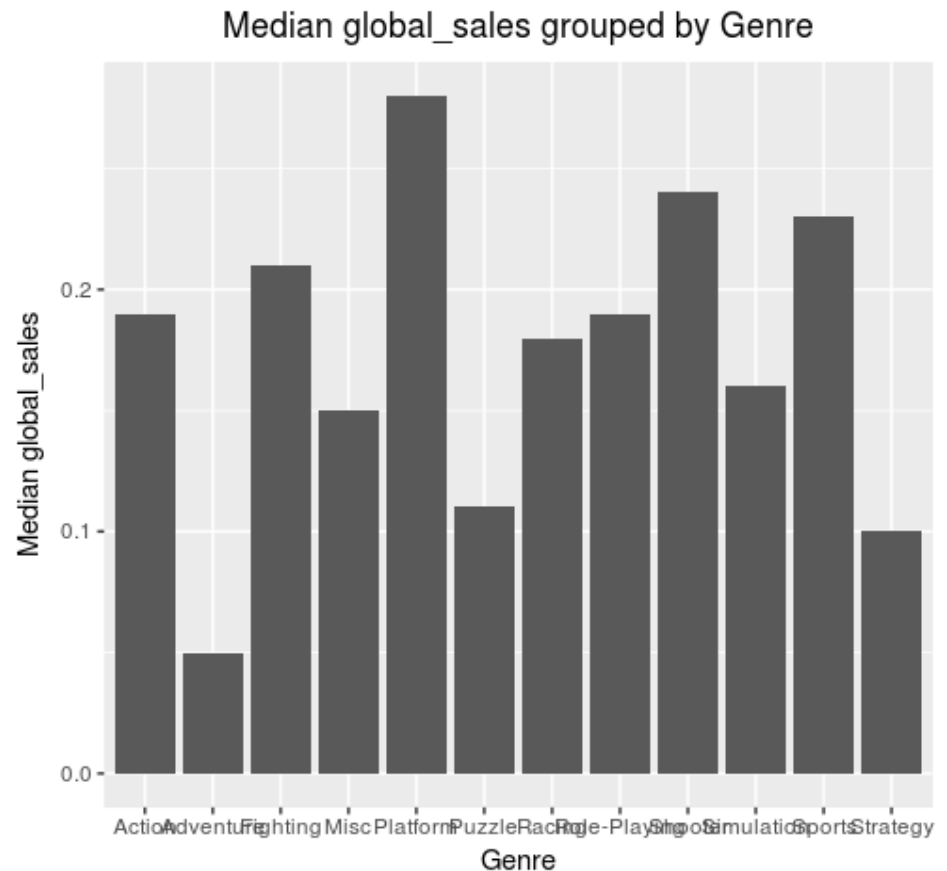
# Some Interesting Plots - 4

**Inference**: Interestingly Global sales in 1980's is very high despite number of games in 2000 – 2010 are more.


Median global_sales grouped by YearII

# Some Interesting Plots - 5

**Inference: Global Sales are more in Platform, Shooter and Sports Genre in the increasing order.**



Median global_sales grouped by Genre

# Some Interesting Plots - 6

**Inference: Global Sales of "Adult Only" rated games is significantly high.**

**New Feature: If a game is AO or K-A, 1 else 0.**



Median global_sales grouped by Rating

# Some Interesting Plots - 7

## Inference:

➤ **Global Sales of 2600 is high since the number of games were more during that era.**

➤ **In current competitors Xbox has high sales!**

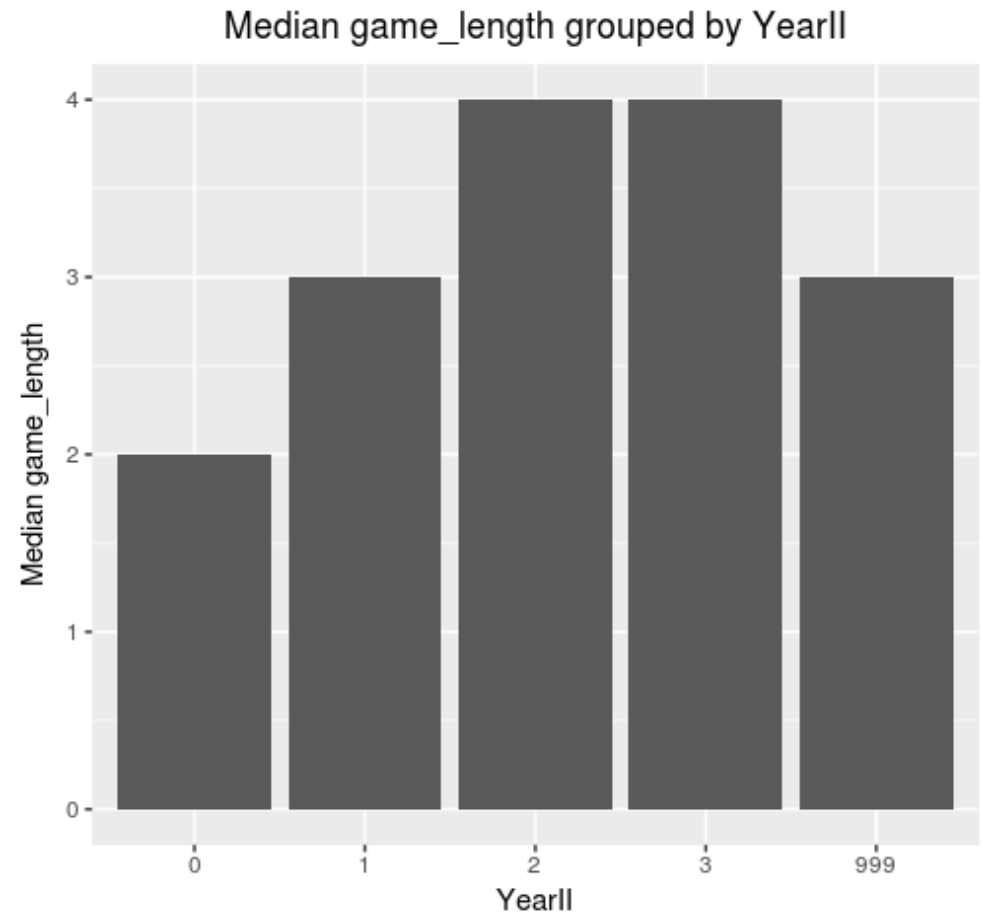➤**New Feature: If 2600/XB/PS 1 else 0.**

➤ **P.S, Any PC fans?!  :P**
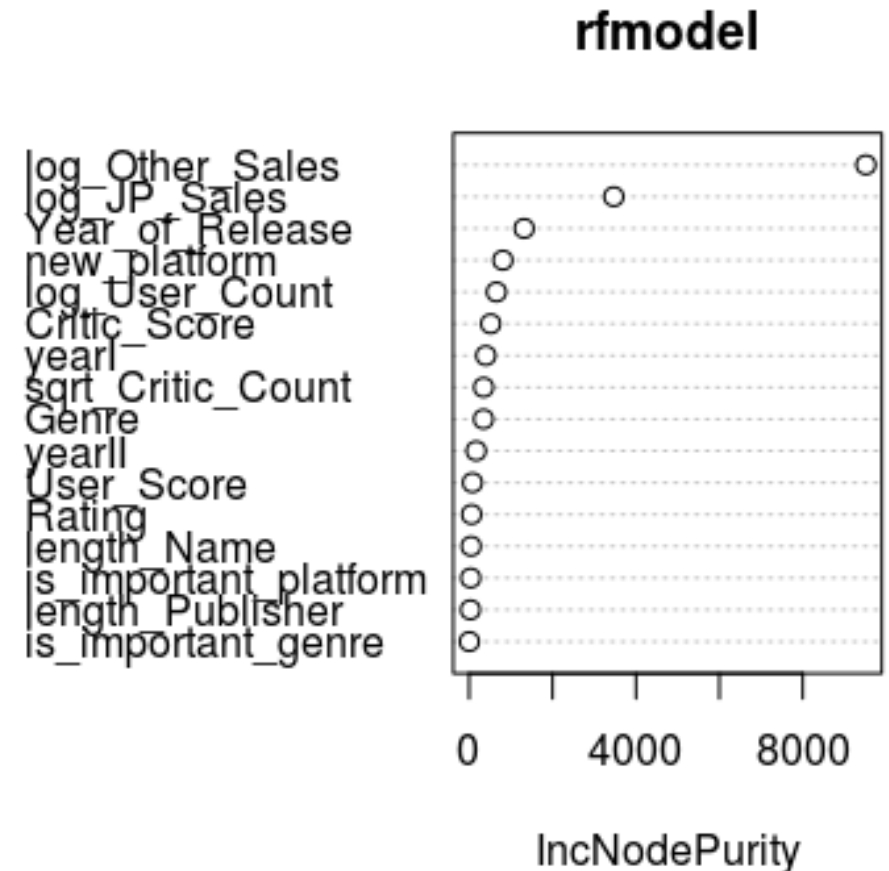
Median global_sales grouped by Platform

**Inference:**

**Interestingly, the length of games have increased with time!**

# Feature Importance – Random Forest

**Inference:**

**Log_Other_Sales, Log_JP_Sales, Year_of_Release are important**

### rfmodel



log_Other_Sales
log_JP_Sales
Year_of_Release
new_platform
log_User_Count
Critic_Score
yearI
sqrt_Critic_Count
Genre
yearII
User_Score
Rating
length_Name
is_important_platform
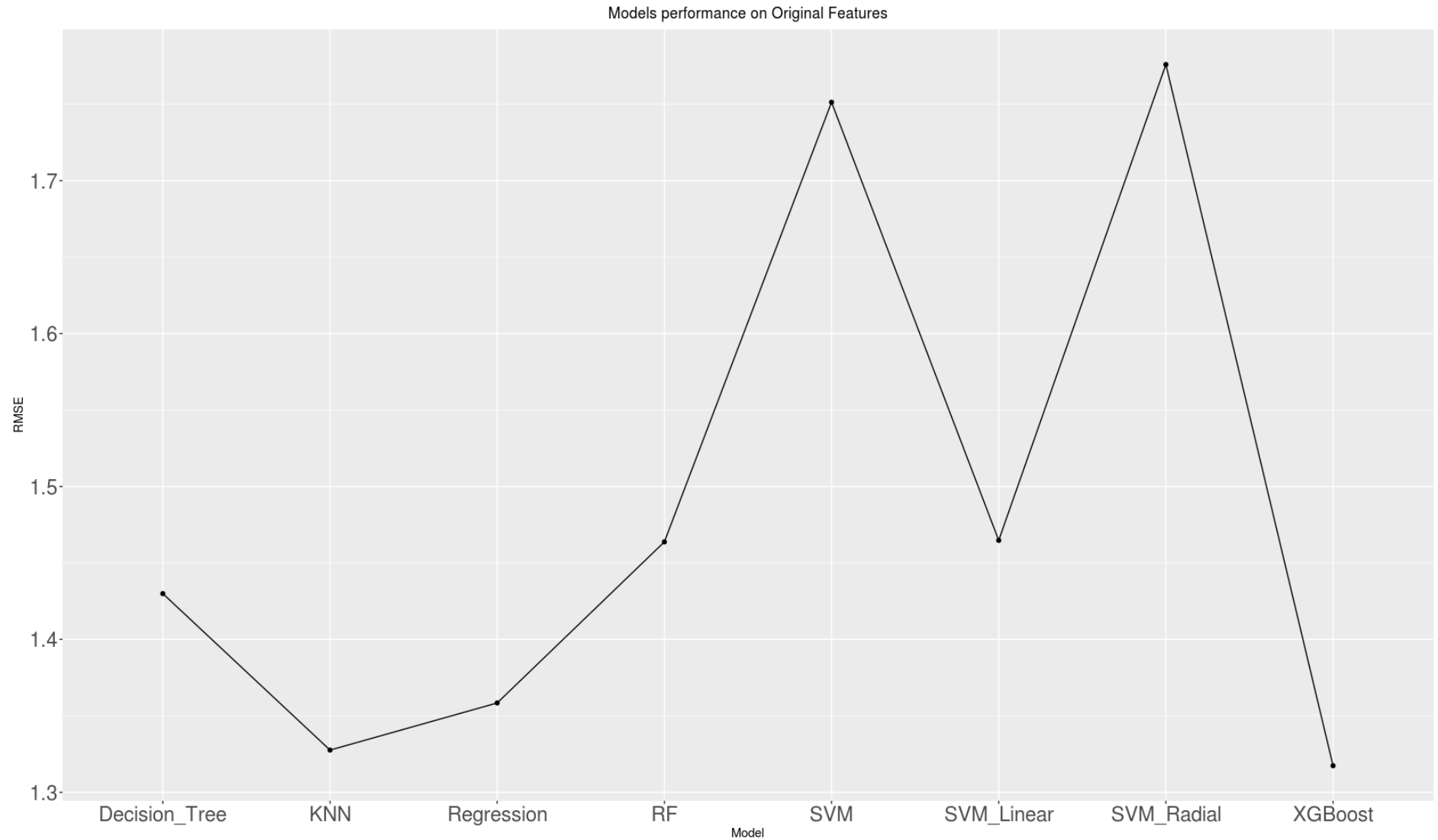length_Publisher
is_important_genre

IncNodePurity

# Feature Importance – XGBoost

## Inference:

**Log_Other_Sales, Log_JP_Sales, If Platform is PS are important features.**

# Models performance on Original Feature



Models performance on Original Features

# Features from Neural Nets – Main Idea

➢ **Create new features from Neural Nets and run other models on them.**

# Features from Neural Nets - How we do it?

➢ **Extracting the activation function's output when we pass an observation.**
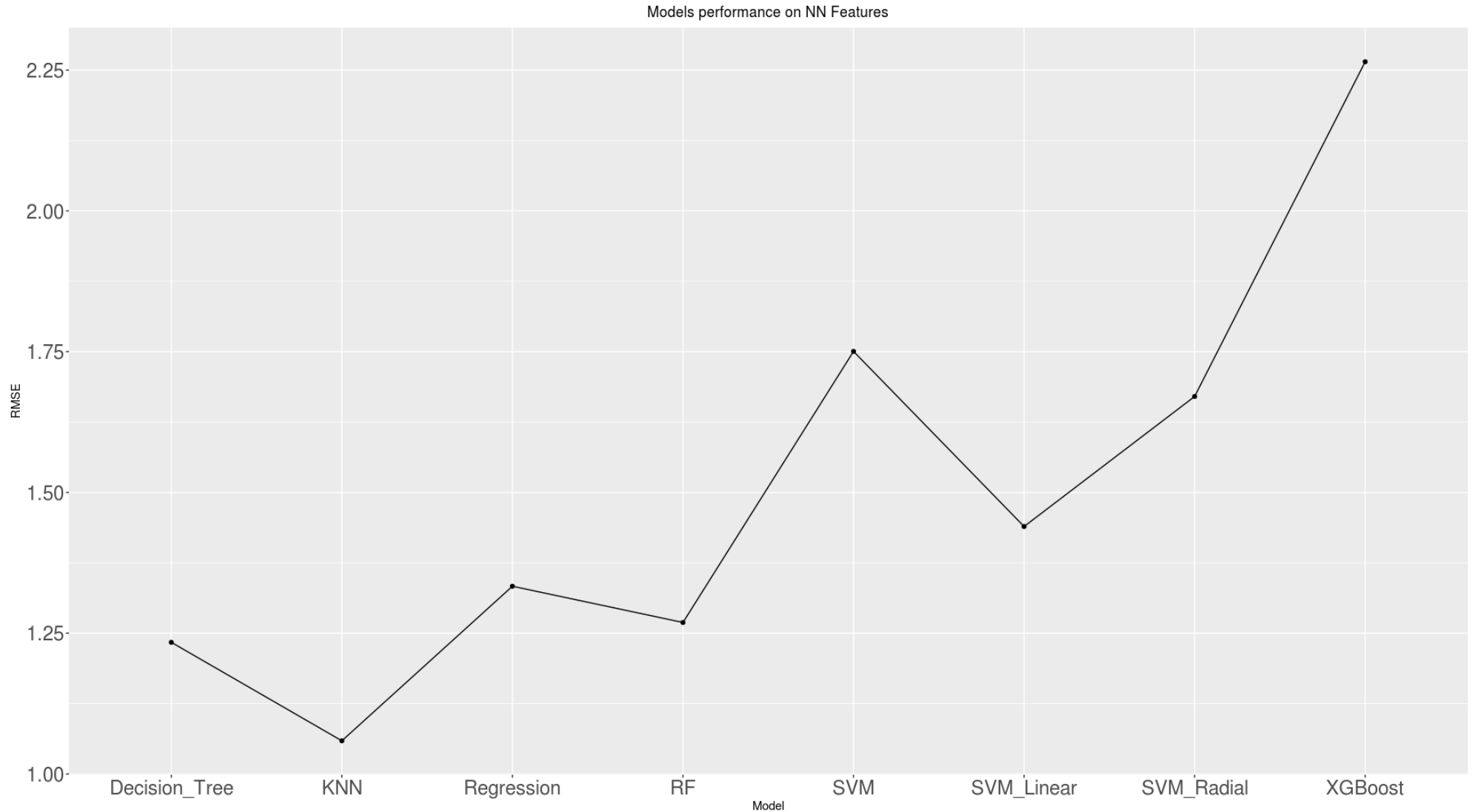
➢ **Neural Network:**

**ActivationFunction : ReLUWithDropout**
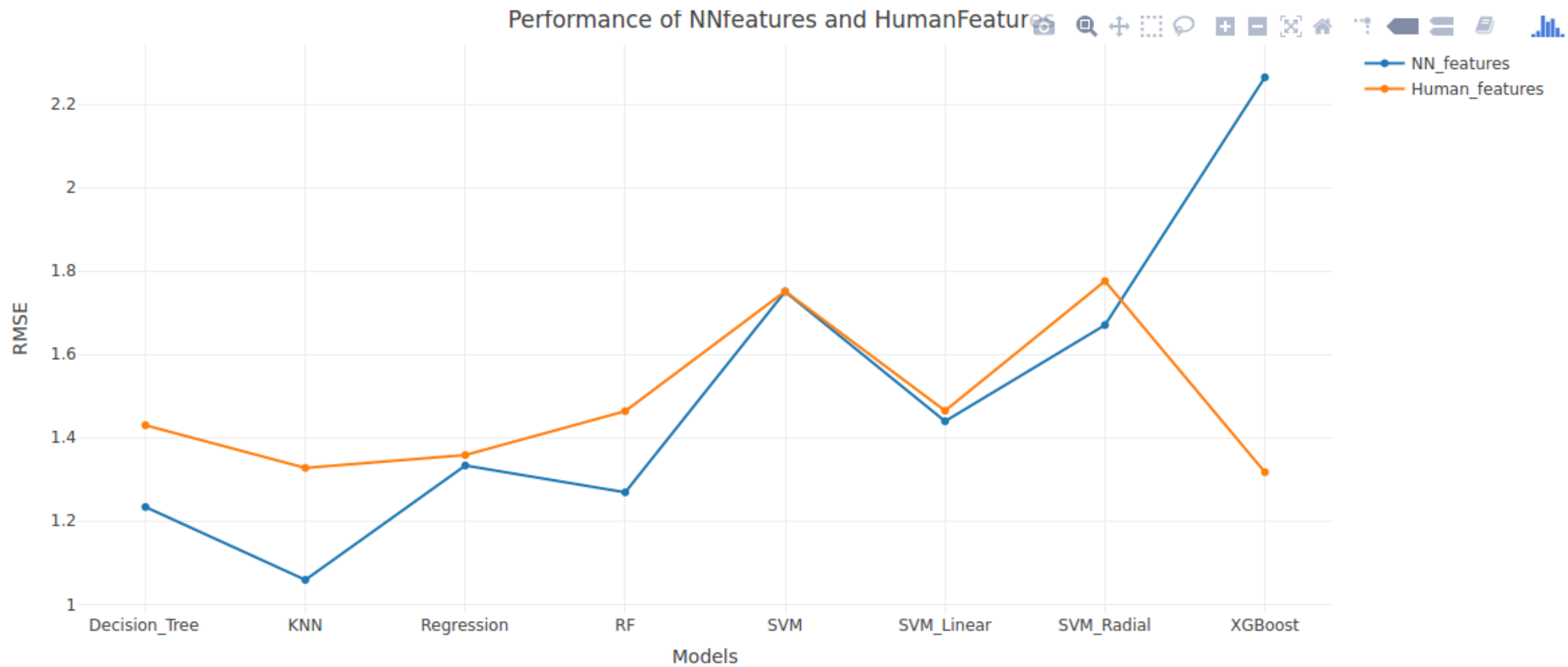
**Hidden Layers : 3 (25, 20, 10)**

**Number of Features: 55**

# Models performance on NN features



Models performance on NN Features

# FaceOff!



Performance of NNfeatures and HumanFeatures

# What does it tell?

> **Features from Neural Networks are clearly superior!**

# Any Questions?