# IMDb Movie Rating Prediction

# Team

Halil Kolatan

Muhammed Maral

# Table of contents

# Introduction

## Objectives:

1. Unleashing the Power of Data

2. Understanding the importance of movie Ratings

3. Motivation towards predicting IMDb movie Ratings

# Unleashing the **Power** of Data

$42.5B

200

7.5M

Global Revenue

Update Per min

Titles in Database

# Importance of **Movie Ratings**

| 10,000 | %95 | %20 ↑ |
|---|---|---|
| Movies Analyzed | Decision Making | Box-office Revenue |

# Motivation towards **Rating Prediction**



Audience Preferences

Accurate Planning

Improving the UX

# Data Collection

## Tools Used:



## Scraped Data:

- Title
- Publish Year
- Runtime
- Genre
- IMDb Rating

- Metascore
- Movie Votes
- Budget
- Gross Revenue



## Movie, Sci-Fi (Sorted by Number of Votes Descending)

1-50 of 17,096 titles. | Next »      View Mode: Compact | **Detailed**

Sort by: Popularity | A-Z | User Rating | **Number of Votes**▼ | US Box Office | Runtime | Year | Release Date | Date of Your Rating | Your Rating

**1. Inception** (2010)

12A | 148 min | Action, Adventure, Sci-Fi

⭐ **8.8** | ☆ Rate this      74 Metascore

A thief who steals corporate secrets through the use of dream-sharing technology is given the inverse task of planting an idea into the mind of a C.E.O., but his tragic past may doom the project and his team to disaster.

Director: Christopher Nolan | Stars: Leonardo DiCaprio, Joseph Gordon-Levitt, Elliot Page, Ken Watanabe

Votes: 2,408,282 | Gross: $292.58M

**2. The Matrix** (1999)

15 | 136 min | Action, Sci-Fi

⭐ **8.7** | ☆ Rate this      73 Metascore

When a beautiful stranger leads computer hacker Neo to a forbidding underworld, he discovers the shocking truth--the life he knows is the elaborate deception of an evil cyber-intelligence.

Directors: Lana Wachowski, Lilly Wachowski | Stars: Keanu Reeves, Laurence Fishburne, Carrie-Anne Moss, Hugo Weaving
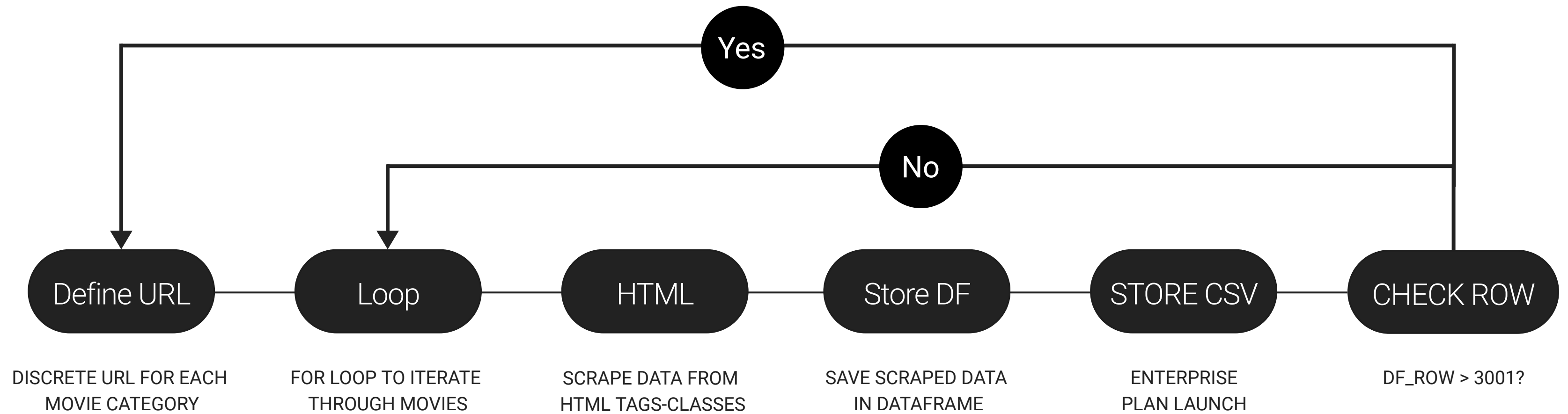
Votes: 1,954,301 | Gross: $171.48M

**Title Type**

Feature Films (17,096)

**Genres**

Sci-Fi (17,096)
Thriller (3,729)
Horror (3,197)
Comedy (2,931)
Mystery (1,412)
Romance (855)
Crime (587)
War (200)
Music (148)
Western (114)
Sport (41)

# Web Scraping Flowchart



Yes

No

Define URL

Loop

HTML

Store DF

STORE CSV

CHECK ROW

DISCRETE URL FOR EACH
MOVIE CATEGORY

FOR LOOP TO ITERATE
THROUGH MOVIES

SCRAPE DATA FROM
HTML TAGS-CLASSES

SAVE SCRAPED DATA
IN DATAFRAME

ENTERPRISE
PLAN LAUNCH

DF_ROW > 3001?

# Web Scraping DataFrame

| | Title | Year | Runtime | Genre | Rating | Budget | Gross US & Canada | Votes | Metascore |
|---|---|---|---|---|---|---|---|---|---|
| **0** | The Dark Knight | (2008) | 152 min | Action - Crime - Drama | 9.0 | $185,000,000 | 534,858,444 | 2710261 | 84.0 |
| **1** | Inception | (2010) | 148 min | Action - Adventure - Sci-Fi | 8.8 | $160,000,000 | 292,576,195 | 2405851 | 74.0 |
| **2** | The Matrix | (1999) | 136 min | Action - Sci-Fi | 8.7 | $63,000,000 | 171,479,930 | 1952536 | 73.0 |
| **3** | The Lord of the Rings: The Fellowship of the Ring | (2001) | 178 min | Action - Adventure - Drama | 8.8 | $93,000,000 | 315,544,750 | 1911195 | 92.0 |
| **4** | The Lord of the Rings: The Return of the King | (2003) | 201 min | Action - Adventure - Drama | 9.0 | $94,000,000 | 377,845,905 | 1882385 | 94.0 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **12975** | World War Z | (2013) | 116 min | Action - Adventure - Horror | 7.0 | $190,000,000 | 202,359,711 | 685861 | 63.0 |
| **12976** | 2001: A Space Odyssey | (1968) | 149 min | Adventure - Sci-Fi | 8.3 | $12,000,000 | 56,954,992 | 683552 | 84.0 |
| **12977** | The Hunger Games: Catching Fire | (2013) | 146 min | Action - Adventure - Sci-Fi | 7.5 | $130,000,000 | 424,668,047 | 682086 | 76.0 |
| **12978** | Spider-Man: Homecoming | (2017) | 133 min | Action - Adventure - Sci-Fi | 7.4 | $175,000,000 | 334,201,140 | 677261 | 73.0 |
| **12979** | Wonder Woman | (2017) | 141 min | Action - Adventure - Fantasy | 7.4 | $149,000,000 | 412,563,408 | 672854 | 76.0 |

12980 rows × 9 columns

# Data Cleaning

## Removed Rows

**Rows with no US Revenue**

**Rows with no Budget**     **Duplicated Movies**

**Rows with no Metascore**     **New Movies**

New Dataframe Rows

12980 → 4852 rows

## Edited Columns

Year
(2018) object → 2018 int64

Runtime
158 min object → 158 int64

Budget
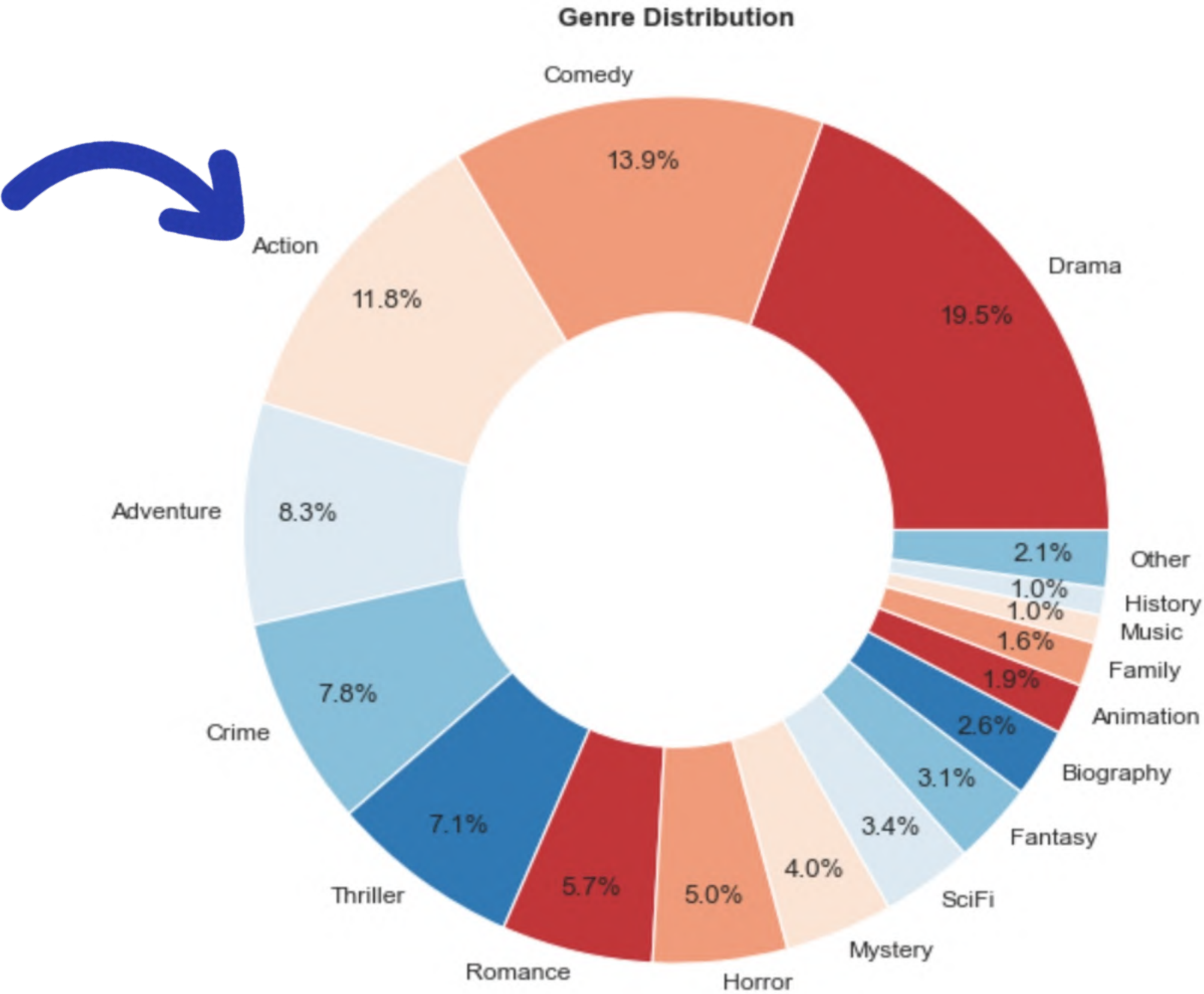$185,000,000 → 185000000 fl64

Gross US Canada
534,858,444 → 534858444 fl64

# Exploratory Data Analysis
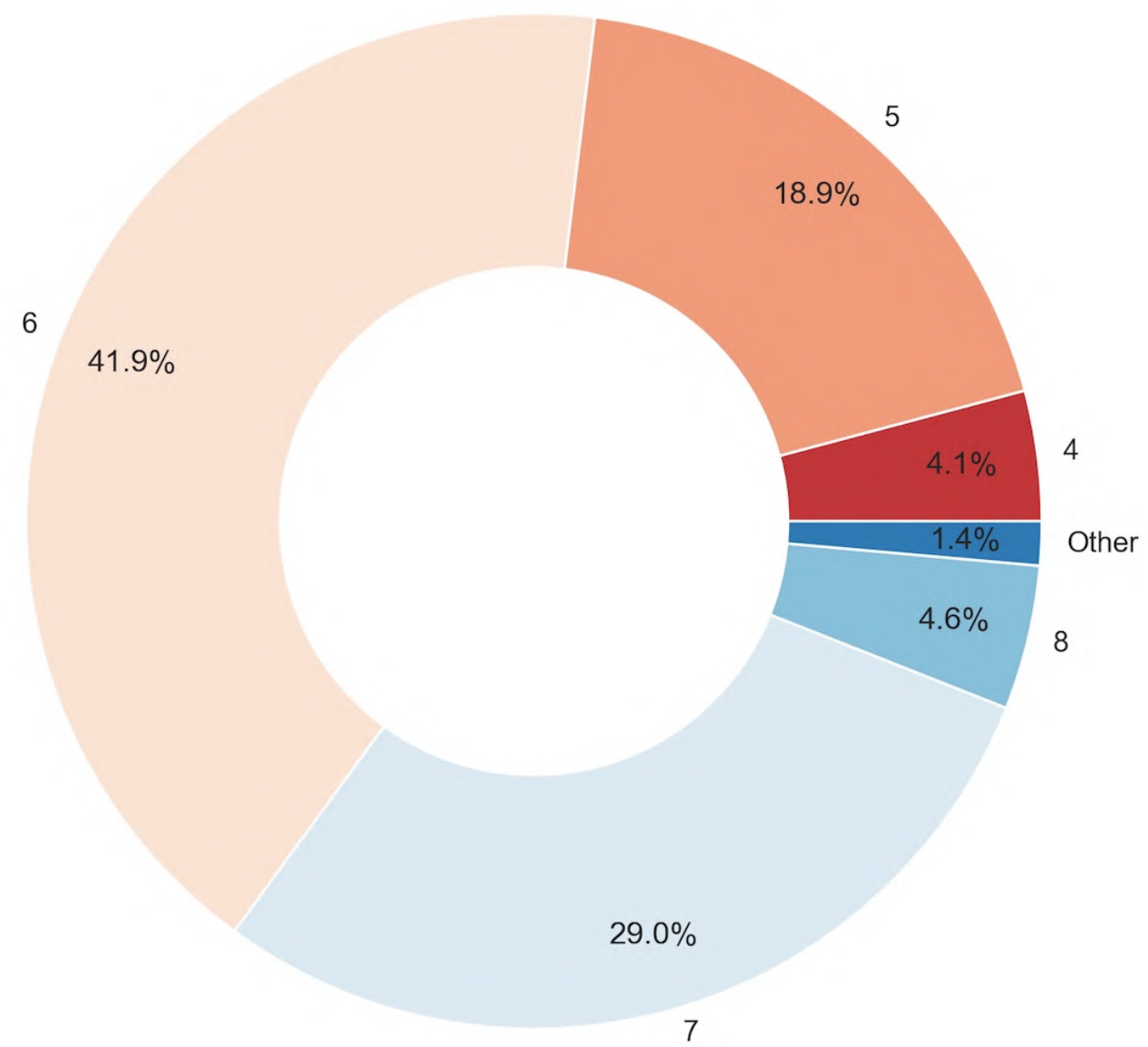
Applying Dummy Variable to the Movie Genre's

| | Action | Adult | Adventure | Animation | Biography | Comedy | Crime | Drama | Fami |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | |
| **1** | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | |
| **2** | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| **3** | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | |
| **4** | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | |
| **4847** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | |
| **4848** | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | |
| **4849** | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | |
| **4850** | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | |
| **4851** | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | |

4852 rows × 22 columns



**Genre Distribution**

Comedy 13.9%
Drama 19.5%
Action 11.8%
Adventure 8.3%
Crime 7.8%
Thriller 7.1%
Romance 5.7%
Horror 5.0%
Mystery 4.0%
SciFi 3.4%
Fantasy 3.1%
Biography 2.6%
Animation 1.9%
Family 1.6%
Music 1.0%
History 1.0%
Other 2.1%

# Visualizations

**Rating Distribution**



5 — 18.9%
6 — 41.9%
4 — 4.1%
Other — 1.4%
8 — 4.6%
7 — 29.0%

**Film Count by Year**



1990, 2000's — 18.0%
1980, 1990's — 11.6%
1970, 1980's — 3.9%
Other — 4.9%
2010, 2020's — 30.6%
2000, 2010's — 31.0%

# Visualizations

Try Pitch

# Feature Engineering

Gross US&Canada %30 - %40

Gross Worldwide %60 - %70

New Column:

Estimated WorldWide Gross

New Column:

Score = (Rating + (Metascore/10)) / 2

Estimated WorldWide Gross - Budget

Estimated Revenue

Try Pitch

# Train-Validation-Test Split

%60 Train

%20 Validation

%20 Test

# Feature Engineering

New Column:

Score = (Rating + (Metascore/10)) / 2

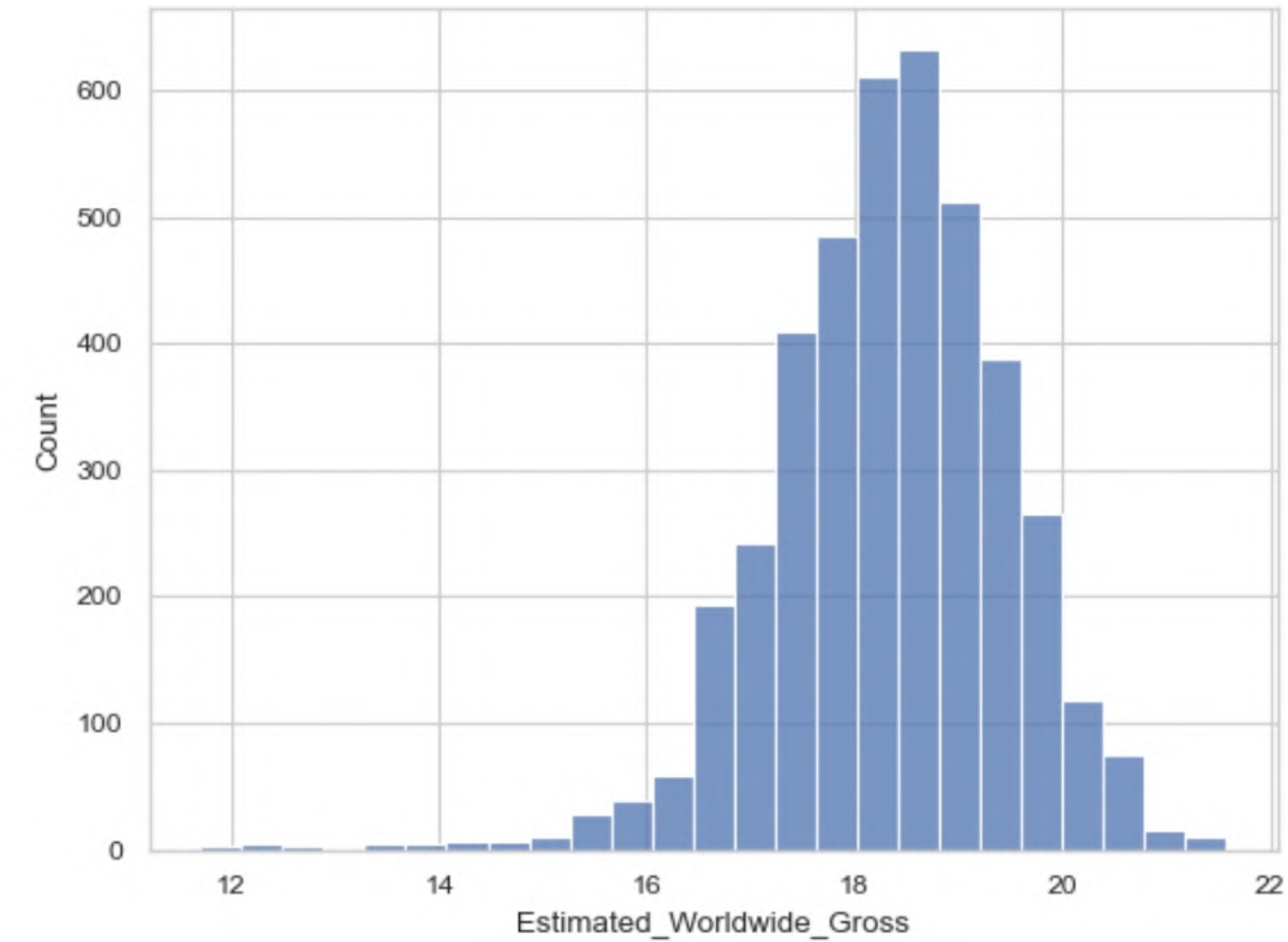Result = Multicollinearity

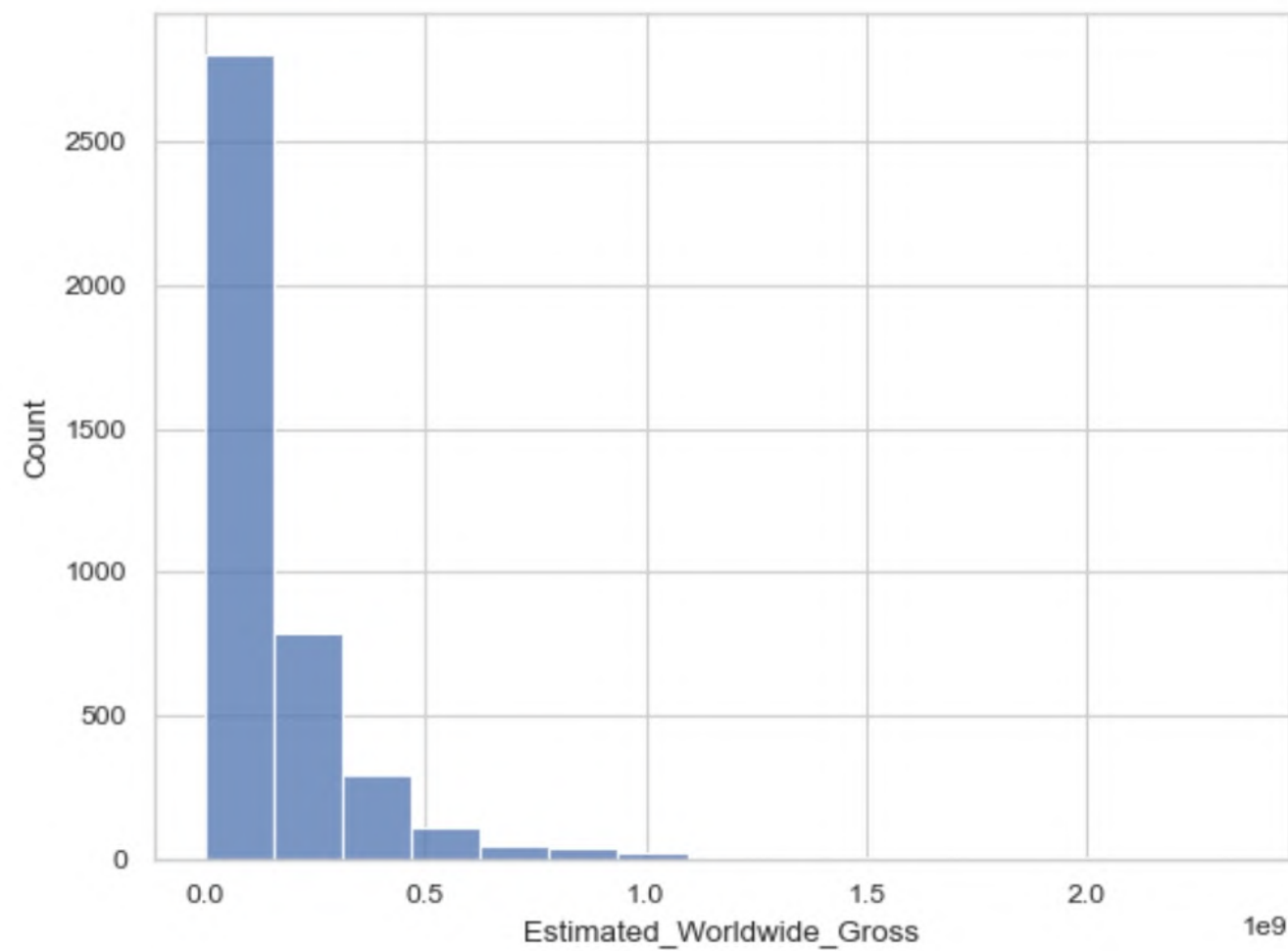# Comparison of Pearson Spearman
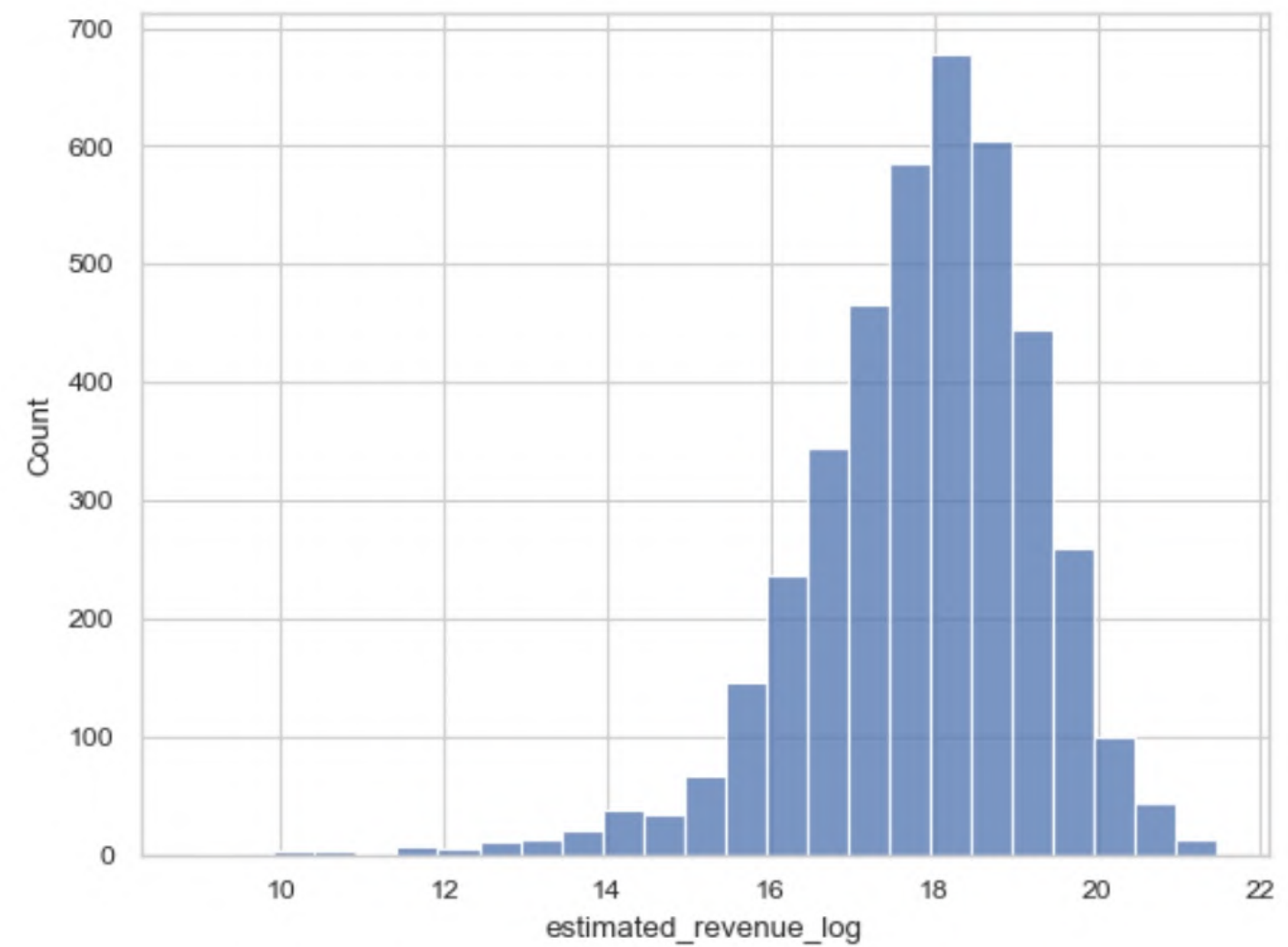


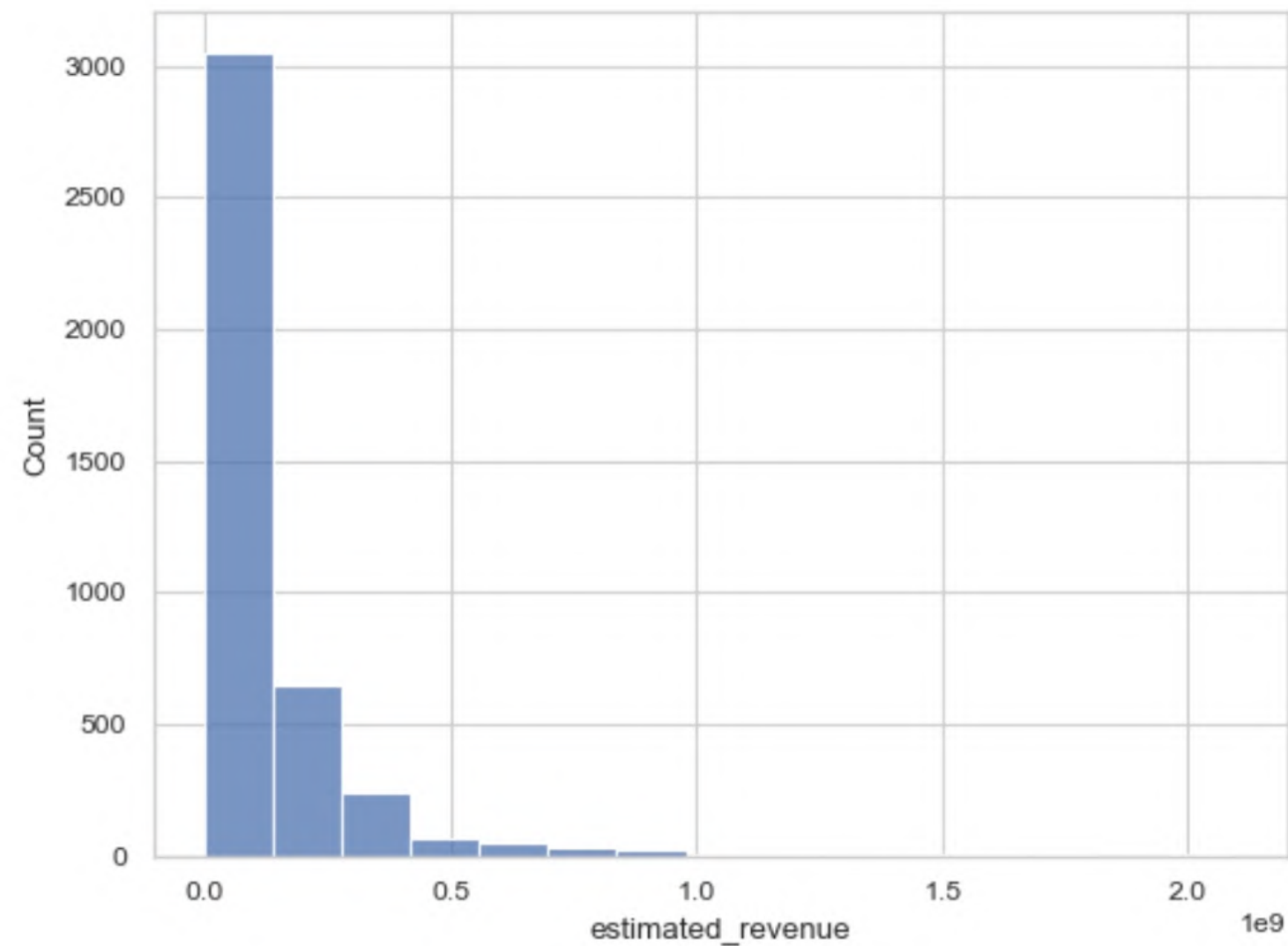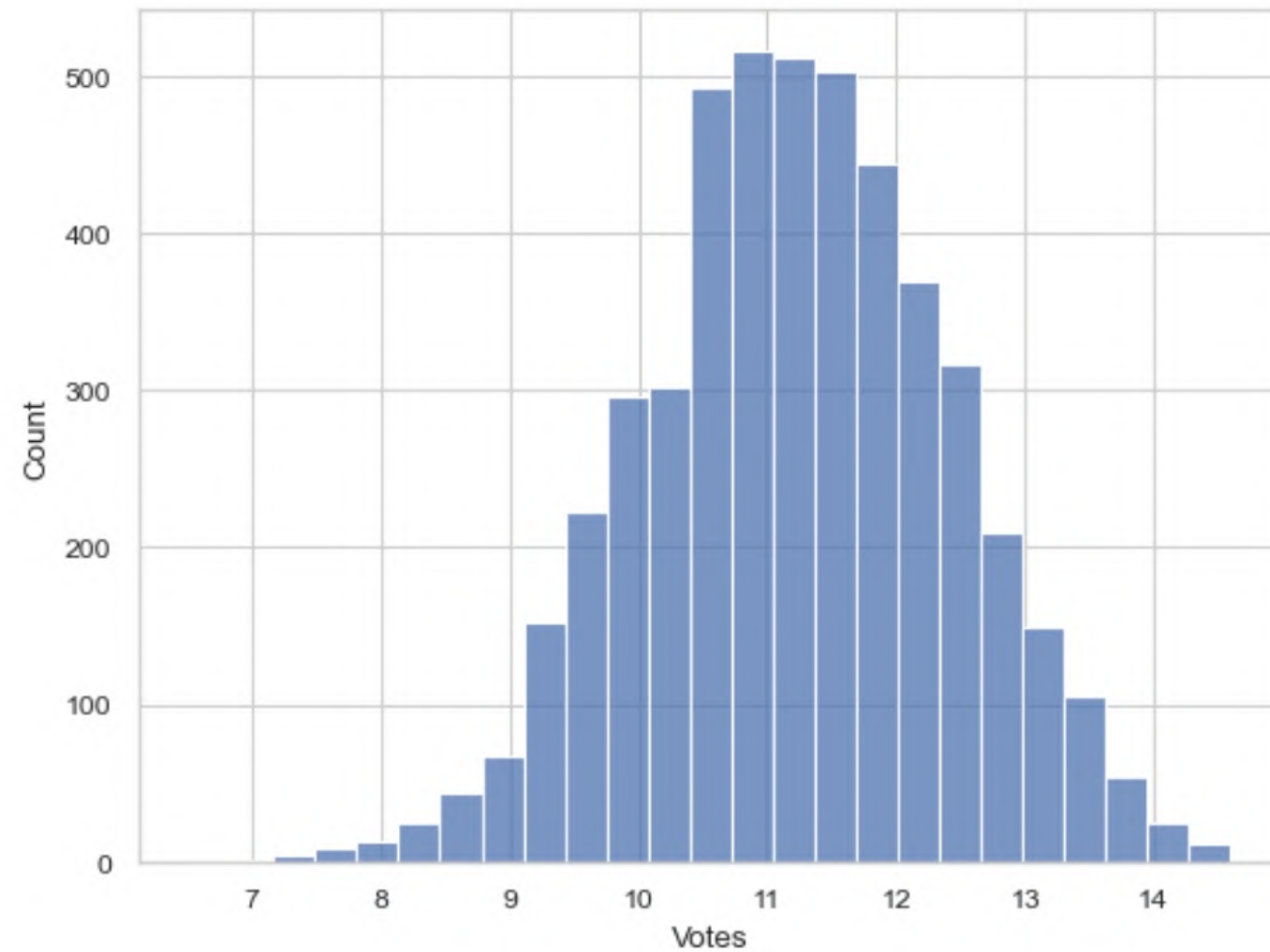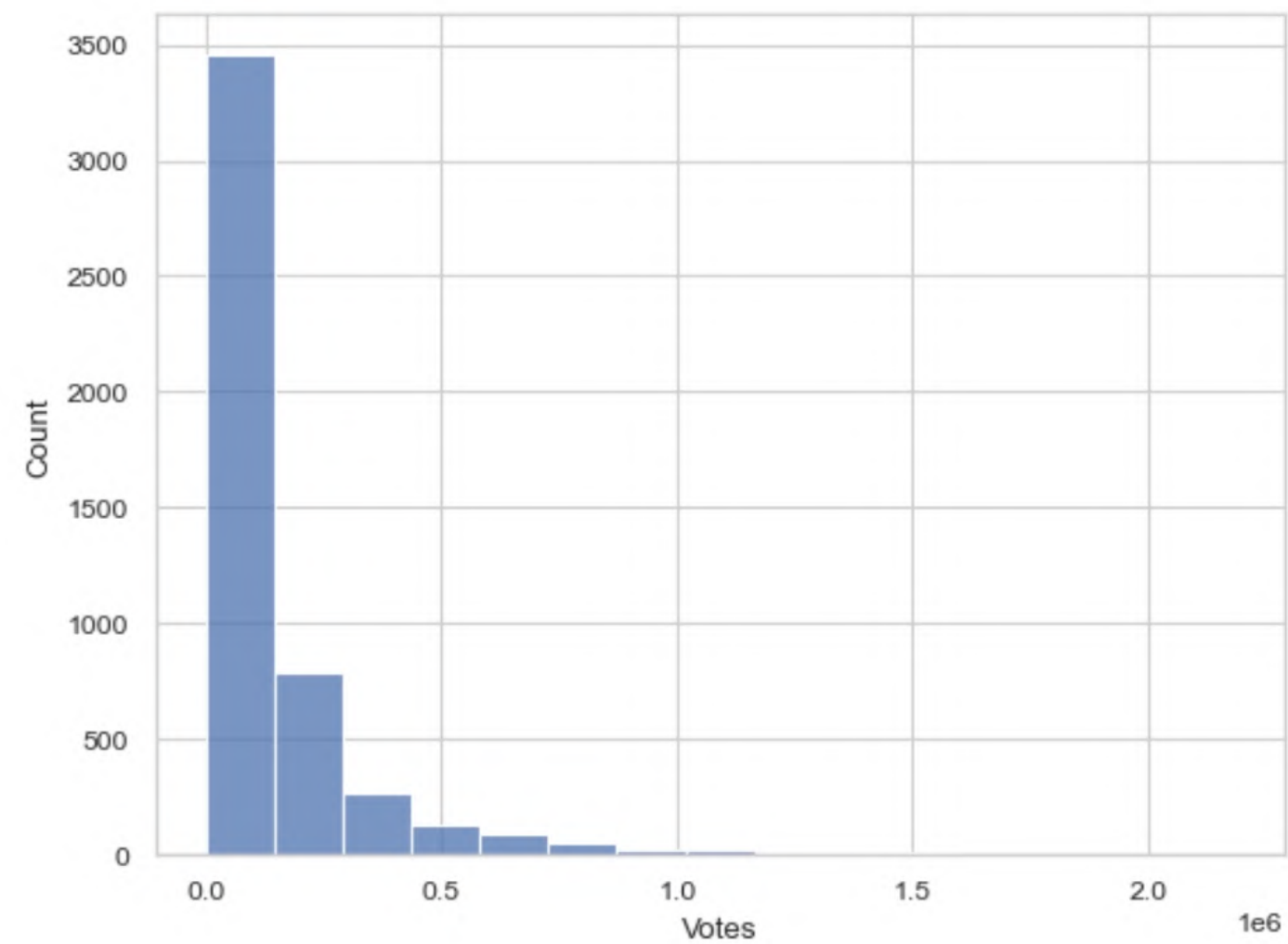Linear Relationship

Rank Correlation

# Model Development

- Analyze the behavior of the dataset. ⟶  - Filter&Manipulate The behavior

# Model Development

- Analyze the behavior of the dataset. ⟶  • Filter&Manipulate The behavior
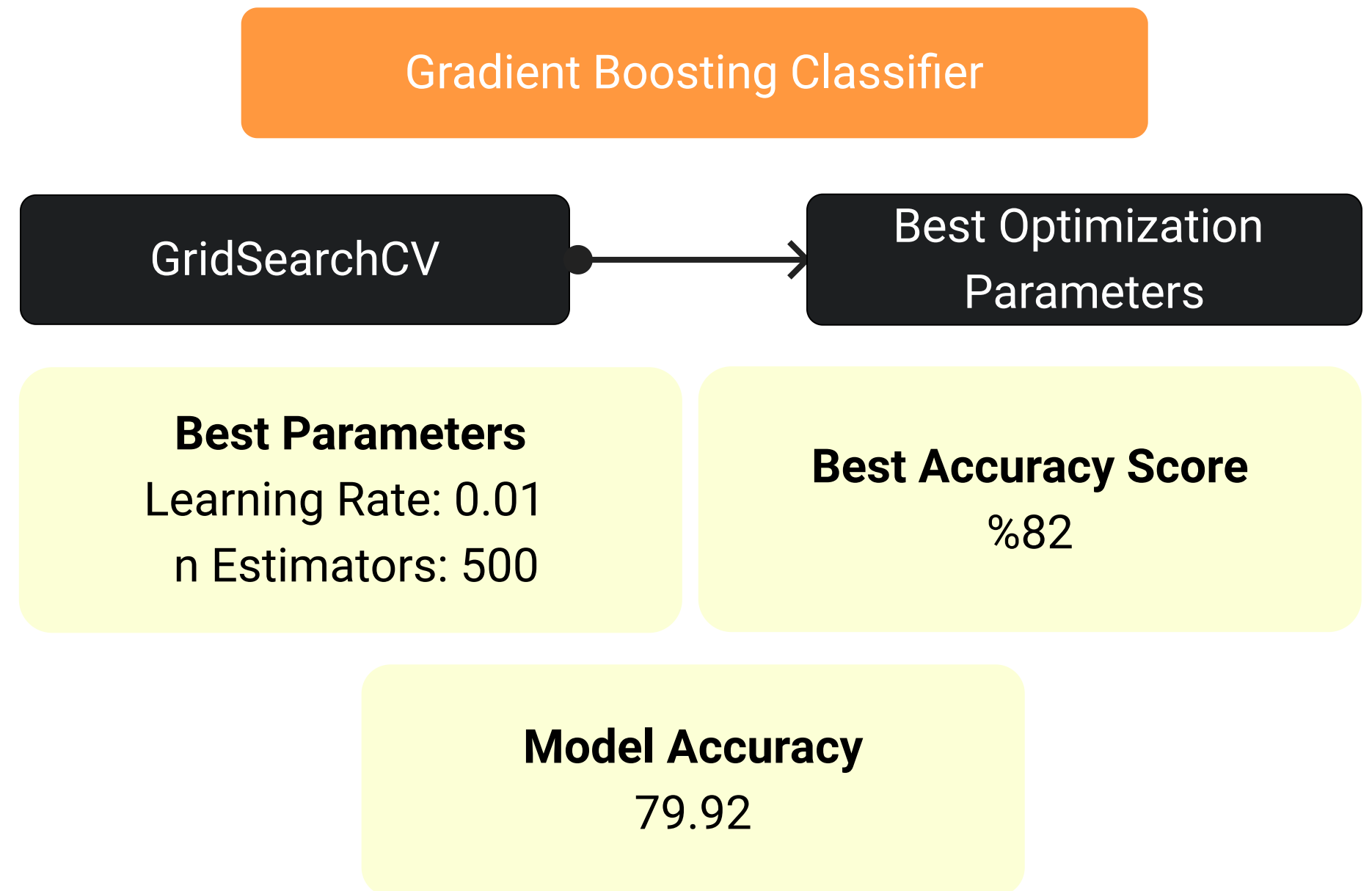
# Model Development

- Analyze the behavior of the dataset. ⟶ - Filter&Manipulate The behavior

# Multiple-Model Techniques

| Model | Accuracy (%) |
|---|---|
| Logistic Regression | 66.53 |
| Naive Bayes | 46.65 |
| Decision Tree (CART) | 73.33 |
| K-NN | 59.53 |
| SVM | 63.95 |
| Gradient Boosting Classifier | 80.02 |
| AdaBoost Classifier | 74.67 |
| Bagging Classifier | 77.75 |
| Random Forest Classifier | 79.71 |
| MLP Classifier | 59.42 |

Classification Accuracies

**Gradient Boosting Classifier**

GridSearchCV → Best Optimization Parameters

**Best Parameters**
Learning Rate: 0.01
n Estimators: 500

**Best Accuracy Score**
%82

**Model Accuracy**
79.92

# Gradient Boosting Classifier

**Model Control**

| Accuracy | → | %79.92 |

| Precision | → | %78.99 |

| Recall | → | %79.92 |

| F1-Score | → | %78.73 |

Confusion Matrix

## Model Evaluation - Train

%60 Train

```
MSE:  0.2931277793963753
R2 Score:  0.64234148086677816
```

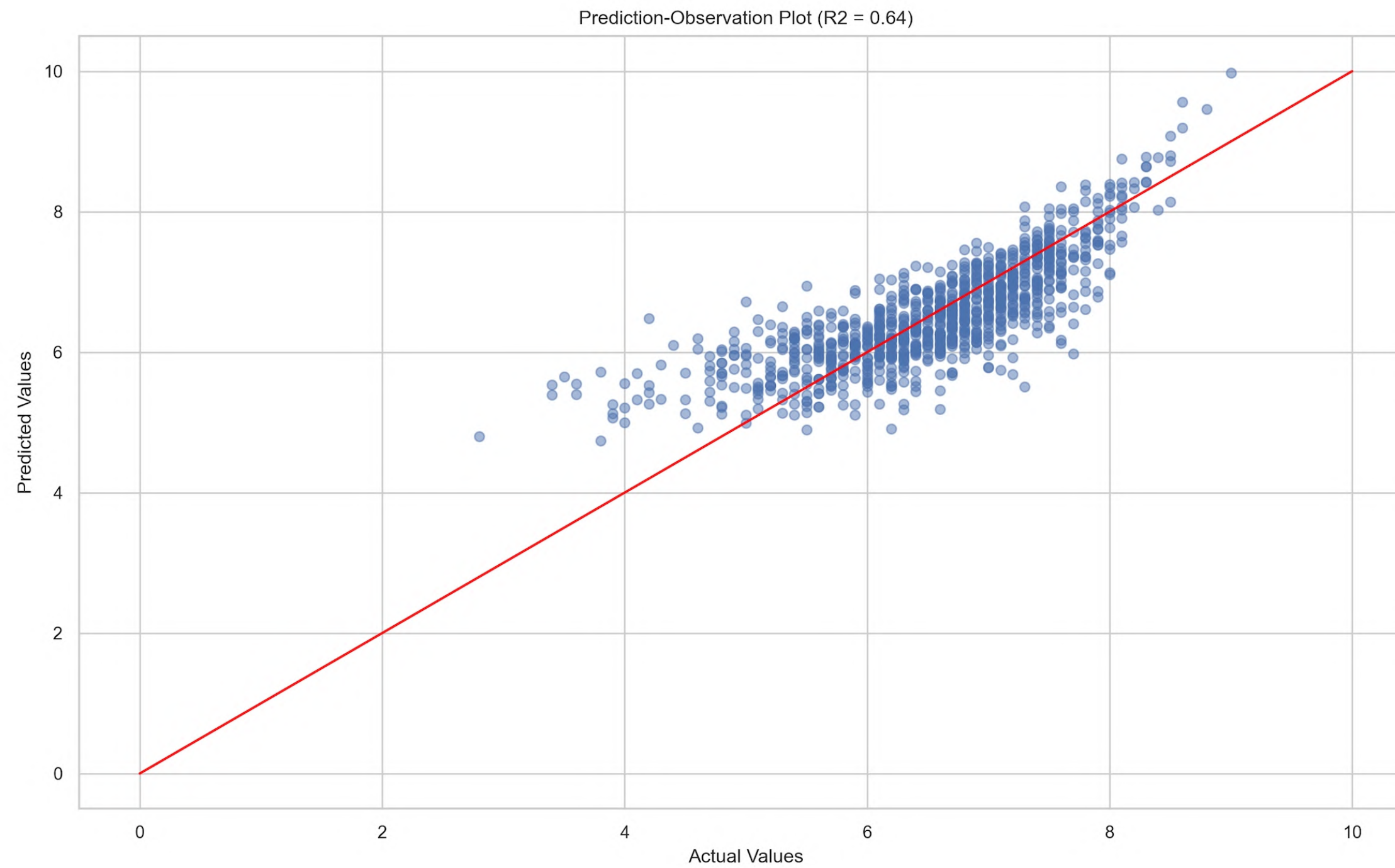## Model Evaluation - Validation

%20 Validation

```
Validation MSE:  0.29312779398516225
Validation R2 Score:  0.654345566892436
```

## Model Evaluation - Test

%20 Test

```
Test MSE:  0.30294564675954655
Test R2 Score:  0.654345566892436
```

# Model Development



Prediction-Observation Plot (R2 = 0.64)

# Conclusion & Future Work

Features: "Year", "Runtime", "Gross US & Canada", "Votes", "Metascore", "Estimated Revenue", "Budget"

Target: Rating

Success Rate: %65.5

## Future Work
Director Data
Actor Data