

# Group 2 Final Project: Predicting Movies' IMDb Score

...

Adam Kritz, Sahara Ensley, Josh Ting

December 6th, 2021

George Washington University

Introduction to Data Mining - DATS 6103

<https://github.com/Saharae/Final-Project-Group2>

# Introduction

The Internet Movie Database (IMDb) is a popular website for cataloging movies

Over 600,000 movies recorded on IMDb as of September 2021

Users and critics can vote on movies and leave reviews

The IMDb logo is displayed within a yellow rounded rectangle. It consists of the letters "IMDb" in a bold, black, sans-serif font. The "i" is lowercase and has a distinctive dot. The "b" is lowercase and has a thick, rounded stem.

**IMDb**

# Research Question

Movie industry wants to invest in movies that they know will be successful

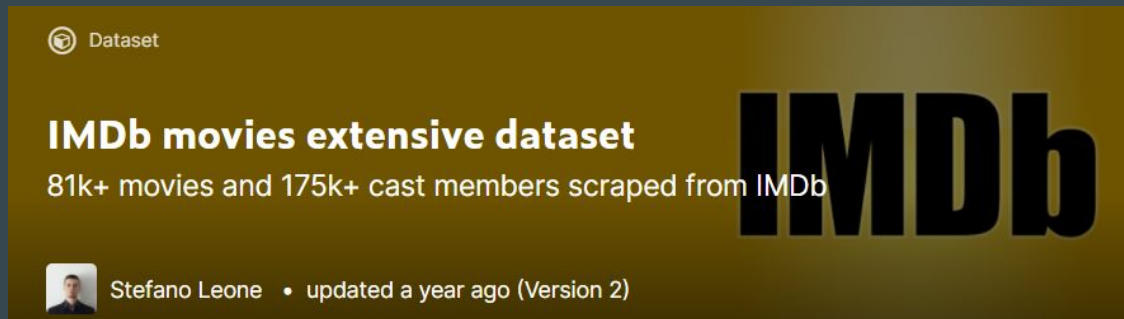
Is there a way to predict if a movie will be enjoyed before it is even released?

Can we build a model to predict the IMDb score of a movie based on its characteristics?

Our research will shed light on whether a model can predict if a movie will be successful solely based on information about the movie

# Data Description

## IMDb Movies Extensive Dataset on Kaggle



Four .csv files initially, narrowed down to Movies and Ratings

Combined the features from Movies and the weighted average vote from Ratings

# Target Variable

Several target variables we were interested in initially: Female average vote, Male average vote, average vote, and weighted average vote

Decided on **weighted average vote**



Weighted average vote appears on IMDb, aka IMDb Rating

Weighted average vote can prevent drastic changes

Weighted average vote ranges from 1-10, to one decimal place

# Variables That Did Not Make The Cut

- Average vote and metascore
  - Both based on the same factors as weighted average vote, could not not be used as predictors
- # of reviews from critics, # of reviews from users, total votes
  - Were not actually features of movies
  - Confusion on definitions
- Language
  - High correlation with country

# Variables We Chose

## Four Numerical Variables

- Duration
- Budget
- USA Gross Income
- Worldwide Gross Income

## Eight Categorical Variables

- Date Published
- Actors
- Writers
- Directors
- Production Company
- Title
- IMDb Description
- Country
- Genre

# Our Modeling Plans

Create a model that could most accurately predict weighted average vote

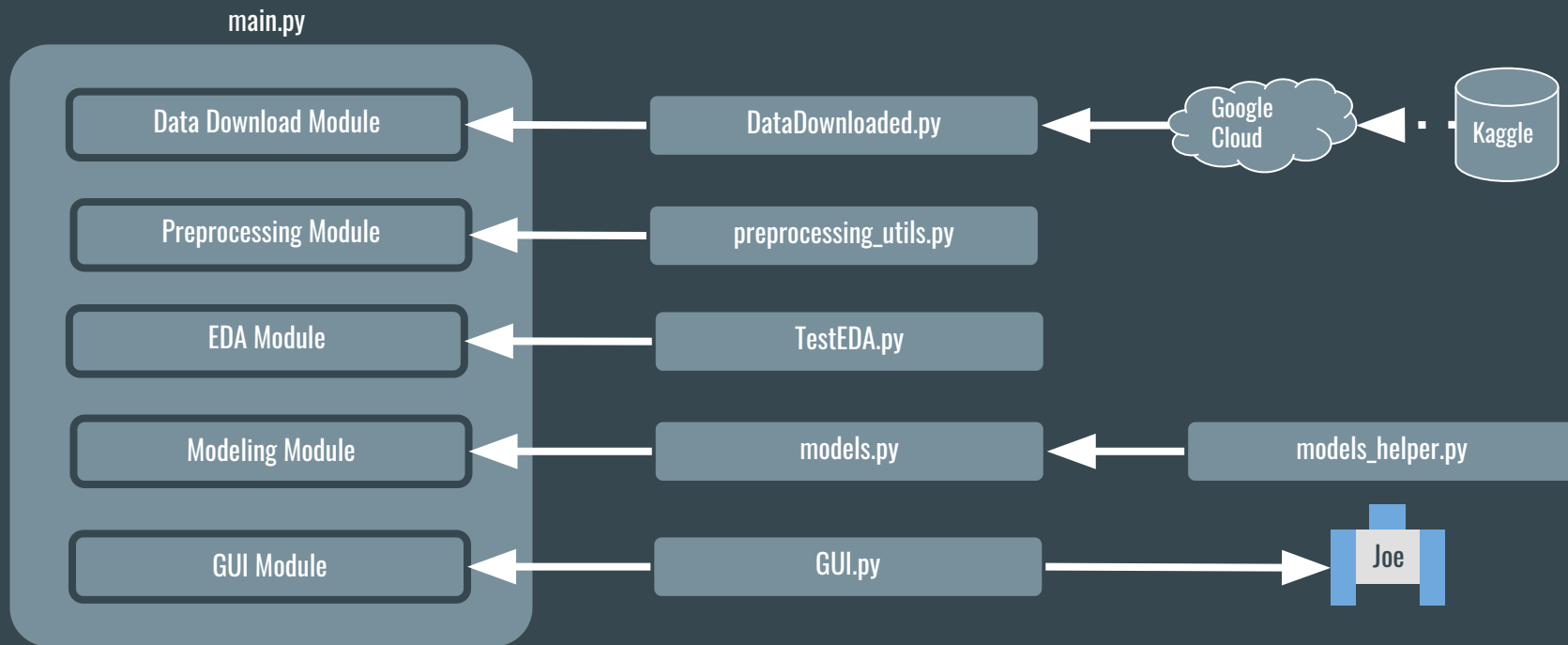
Took a set of models from SKlearn to use for modelling

- Linear Regression
- Random Forest
- Gradient Boosting
- Adaptive Boosting
- K-Nearest Neighbors

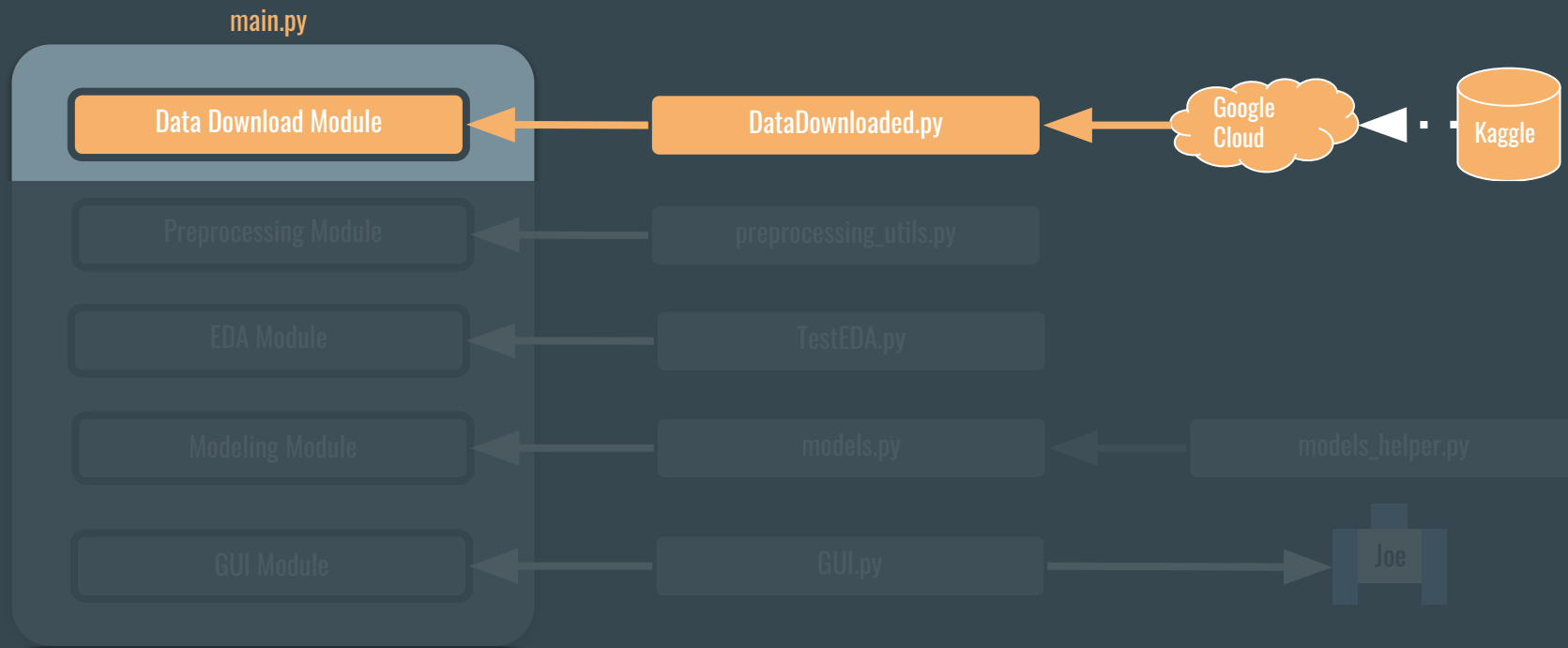
Planned to improve upon the best performers to create the best model



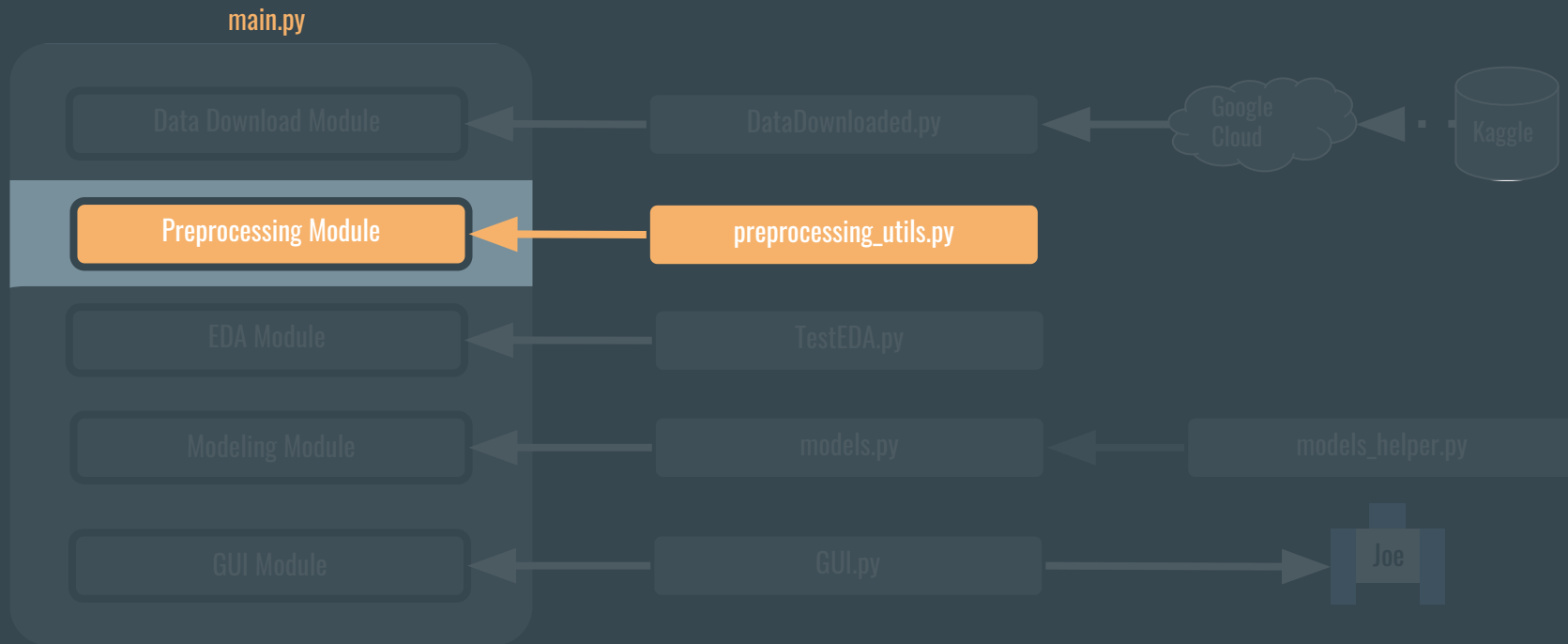
# Code Architecture



# Data Download



# Preprocessing: Code Architecture



# Preprocessing: Initial steps

**Movies**

**Ratings**

**Names**

**Title  
Principals**

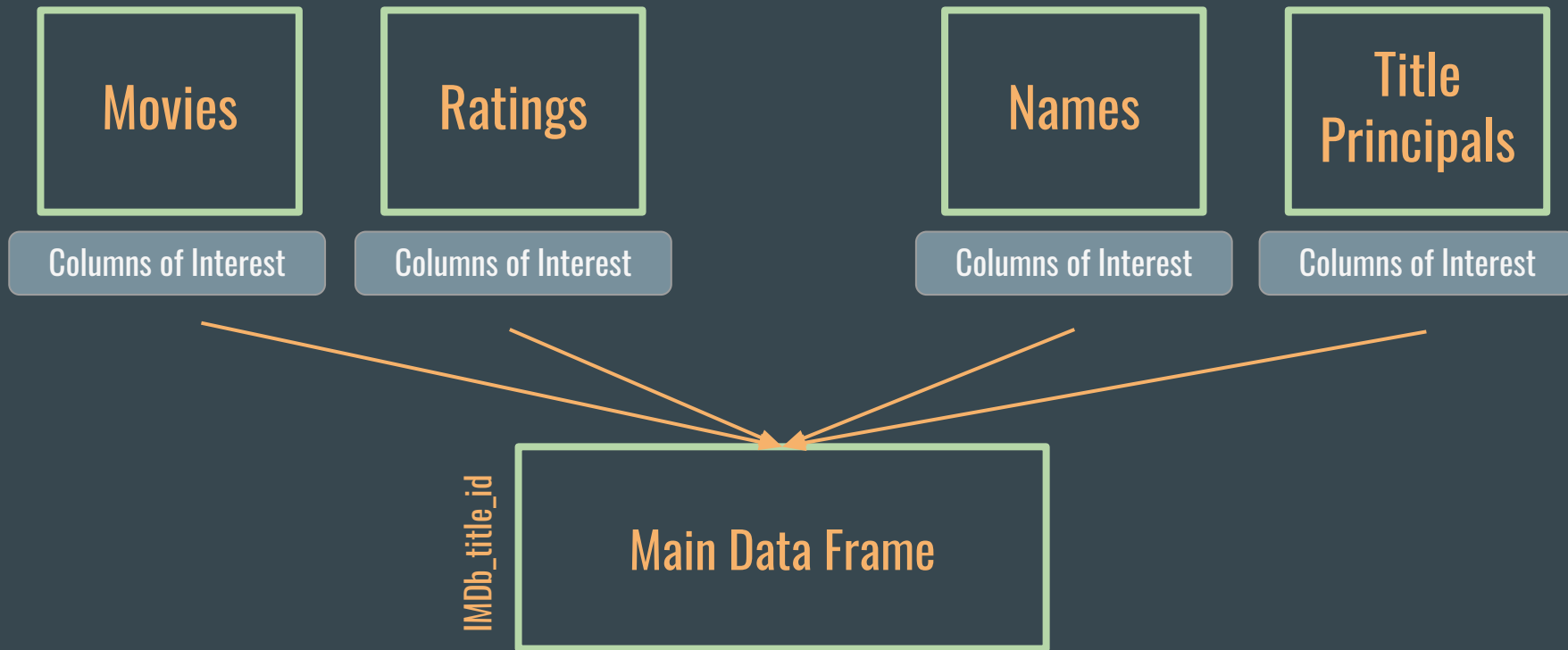
Date Published  
Budget  
Worldwide Income  
Domestic (USA) Income  
Genre  
Country  
Director  
Writer  
Actors  
Production Company  
Title  
Description  
Duration

Weighted Average Vote

Actor ID

Ordering

# Preprocessing: Initial steps



# Preprocessing: Feature Transformation



## Date Published

Budget  
Worldwide Income  
Domestic (USA) Income  
Genre  
Country  
Director  
Writer  
Actors  
Production Company  
Title  
Description



# Preprocessing: Feature Transformation

## Movies

1. All non-US dollar values dropped
2. All remaining values adjusted for inflation based on 2021 CPI

Date Published

Budget

Worldwide Income

Domestic (USA) Income

Genre

Country

Director

Writer

Actors

Production Company

Title

Description

Budget

\$ 500000

(1917)

Budget Adjusted

11722649.57

(2021)

# Preprocessing: Feature Transformation

## Movies

1. Genre Combinations Encoded
2. Transformed into a binary string
3. Binary string expanded into 10 individual columns

Date Published  
Budget  
Worldwide Income  
Domestic (USA) Income  
**Genre**  
Country  
Director  
Writer  
Actors  
Production Company  
Title  
Description

Genre		Genre_1 ... Genre_10
'Action, Comedy, Romance'	→	0 - 0 - 1 - 1 - 0 - 1 - 0 - 0 - 1 - 1



# Preprocessing: Feature Transformation

## Movies

1. First listed country taken as primary country
2. Primary country mapped to region using UN ISO country/region codes

Date Published  
Budget  
Worldwide Income  
Domestic (USA) Income  
Genre  
**Country**  
Director  
Writer  
Actors  
Production Company  
Title  
Description

Country → Region  
'USA, UK, Italy' → Americas

# Preprocessing: Feature Transformation

**Movies**

**Names**

**Title  
Principals**

1. Popularity calculated by frequency
2. Importance calculated by order of mention
3. Weighted Popularity calculated

Date Published  
Budget  
Worldwide Income  
Domestic (USA) Income  
Genre  
Country  
**Director**  
**Writer**  
**Actors**  
Production Company  
Title  
Description

Director → Director Weighted Popularity  
'Steven Spielberg' → 0.00015

# Preprocessing: Feature Transformation

## Movies

1. Popularity calculated by frequency

Date Published  
Budget  
Worldwide Income  
Domestic (USA) Income  
Genre  
Country  
Director  
Writer  
Actors  
Production Company  
Title  
Description

Production Company

‘Warner Bros.’

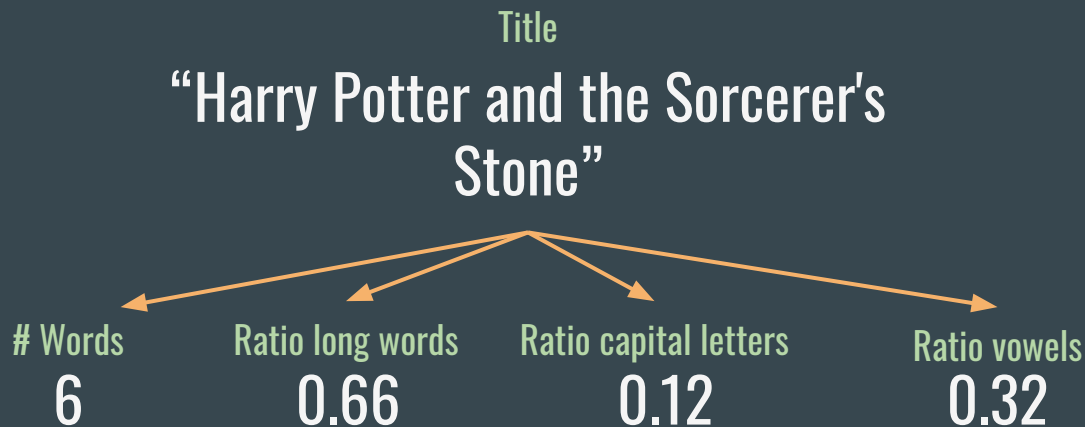
Production Company Popularity

0.115

# Preprocessing: Feature Transformation

**Movies**

Date Published  
Budget  
Worldwide Income  
Domestic (USA) Income  
Genre  
Country  
Director  
Writer  
Actors  
Production Company  
**Title**  
**Description**



# Preprocessing: Feature Transformation

**Movies**

Duration

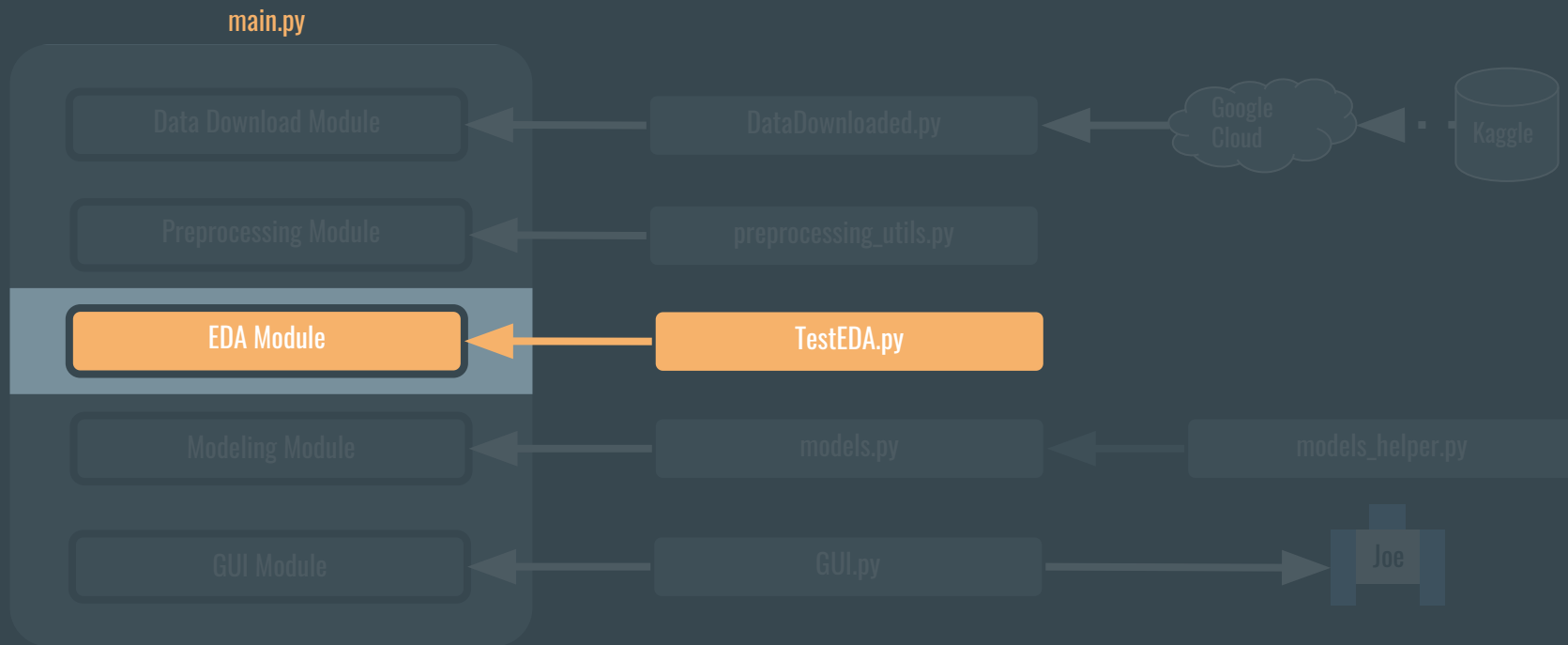
**Ratings**

Weighted Average Vote

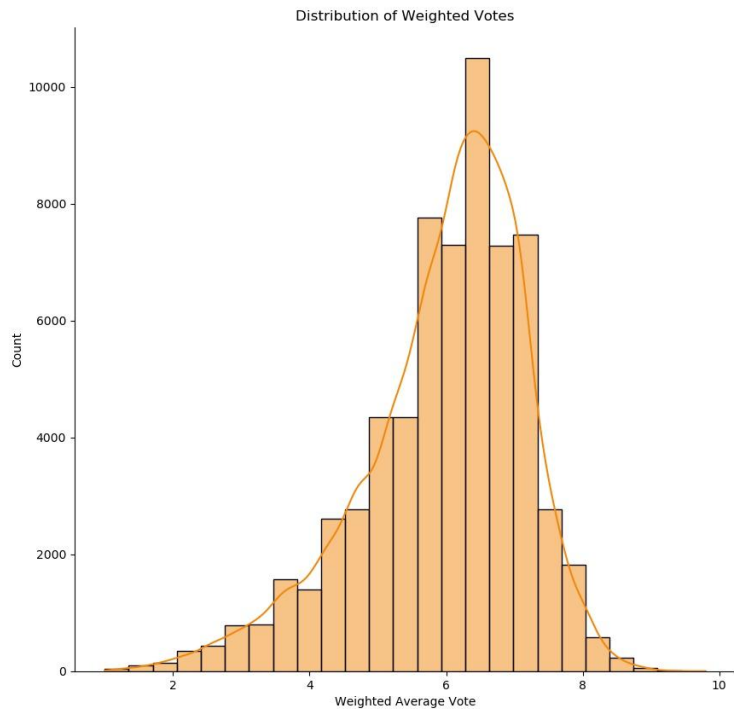
# Preprocessing: Final Steps

1. Data set was split into train, test, and validation
2. All missing values were imputed using the mean
3. All values were scaled

# Preprocessing: Code Architecture



# EDA: Target Variable



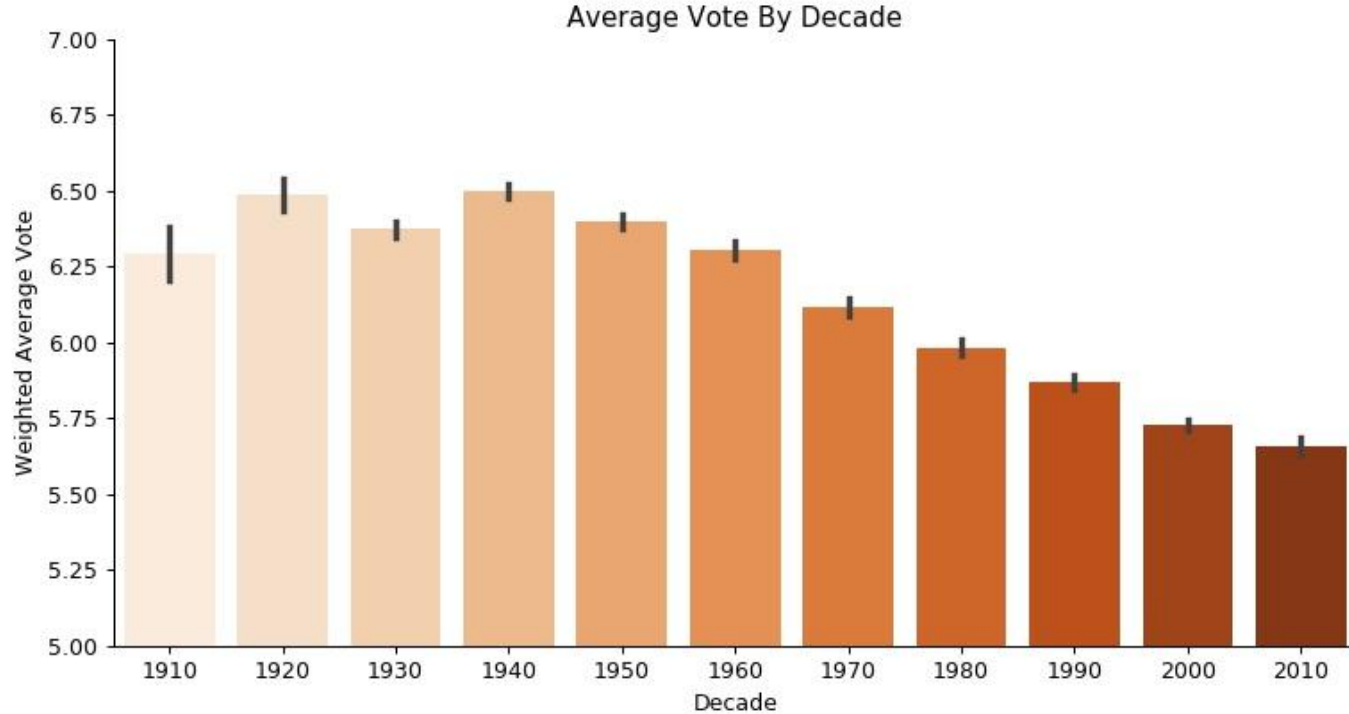
Possible range: 1 - 10

Mean: 5.96

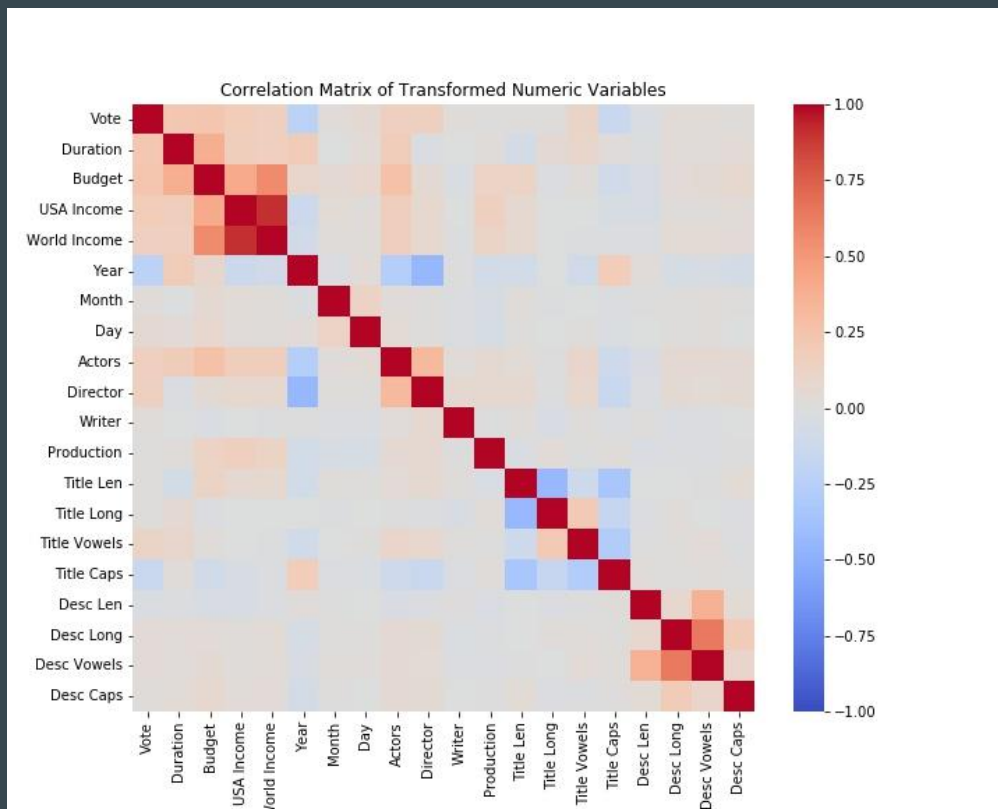
Standard Deviation: 1.19



# EDA: Vote by Decade of Release



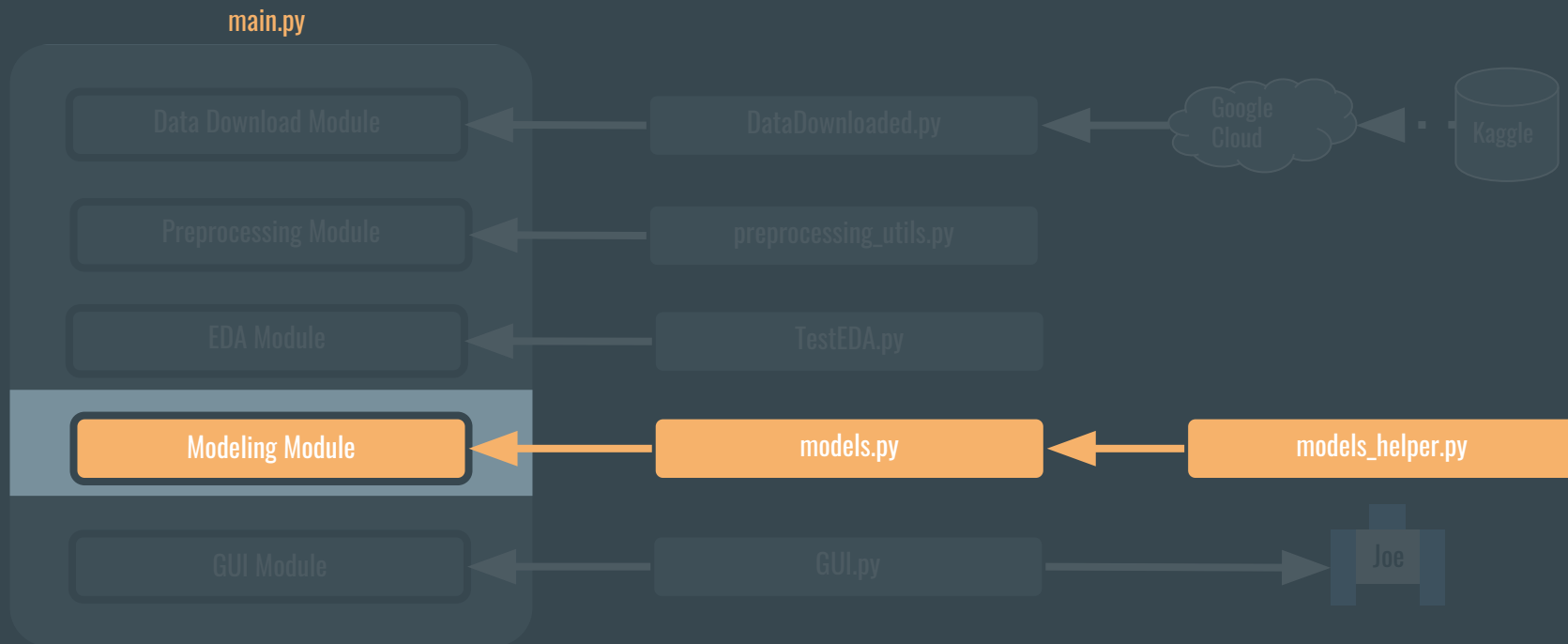
# EDA: Correlation



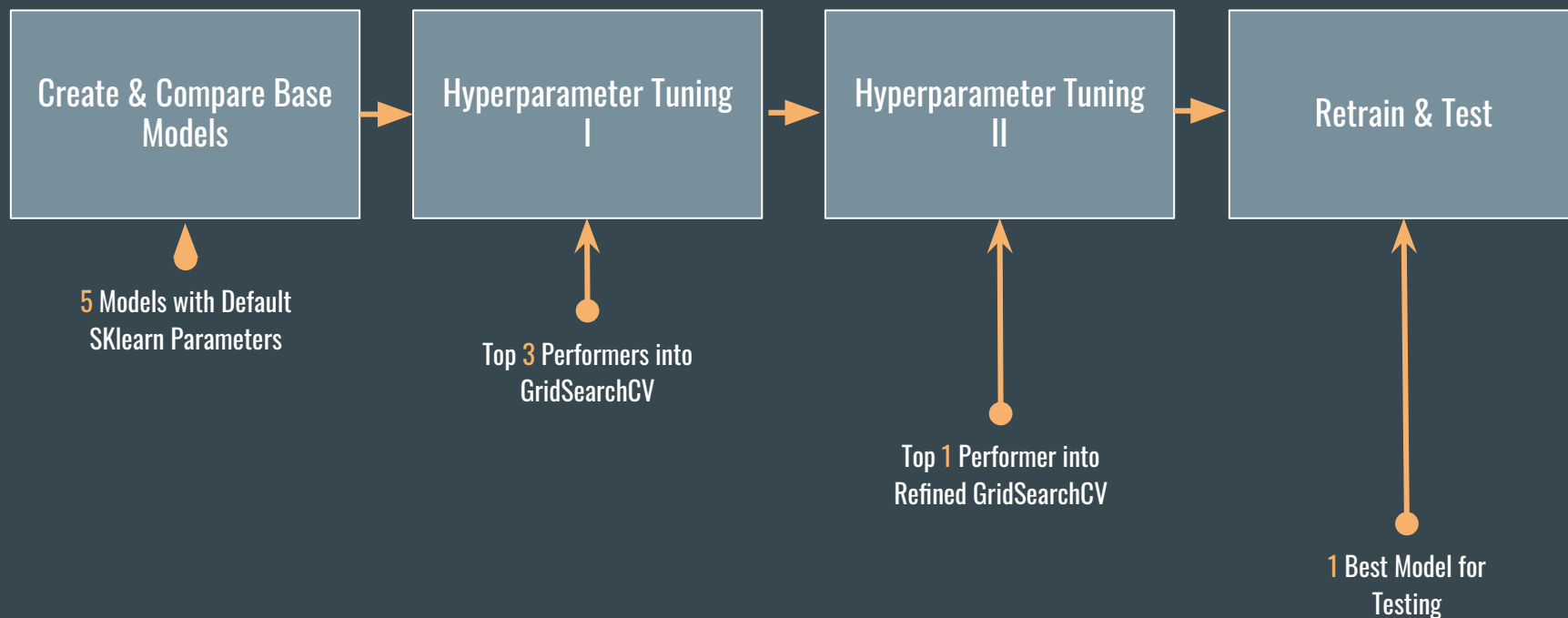
# Modeling: Objective

Create a model that fits well to our dataset and would allow us to predict a satisfiable *weighted average IMDb rating*.

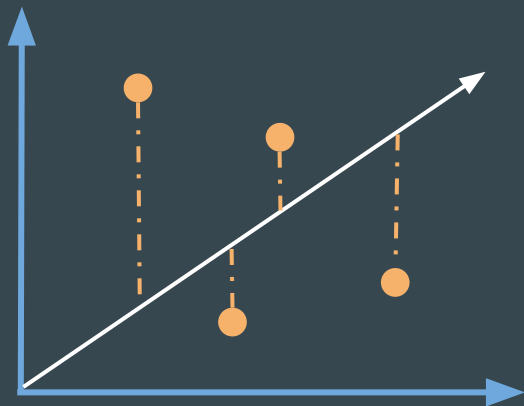
# Modeling: Code Architecture



# Modeling: Methodology



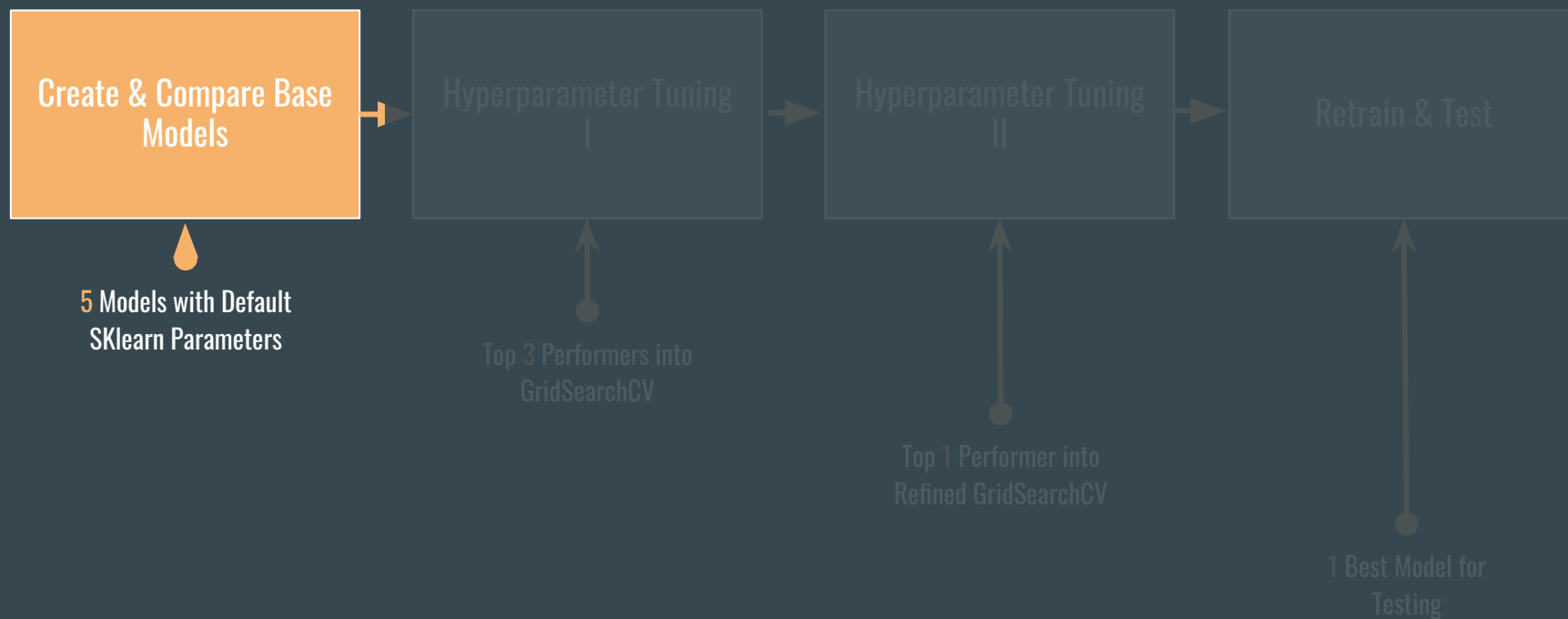
# Modeling: Scoring



$$MSE = \frac{1}{n} \sum_{i=1}^n (y_{actual,i} - y_{predicted,i})^2$$

, where  $n$  is the number of observations in the dataset and  $i$  is the  $i$ th observation in  $n$ .

# Modeling: Base Models

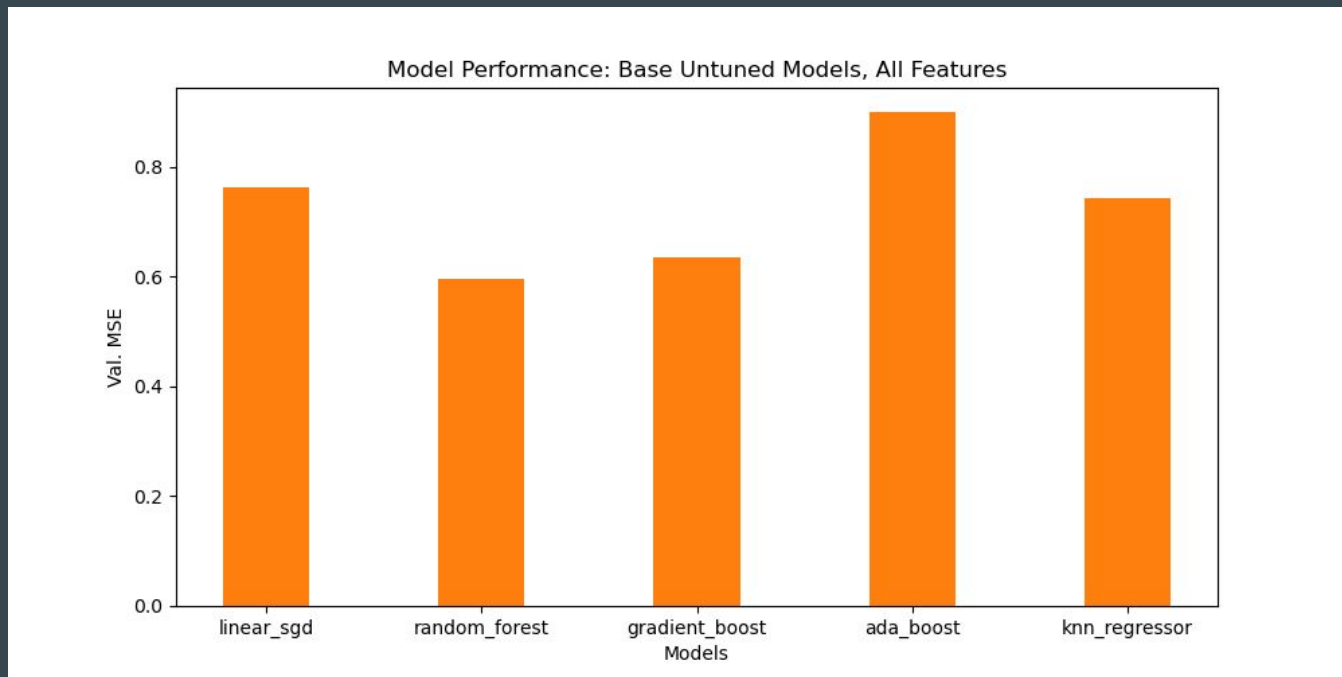


# Modeling: Base Models

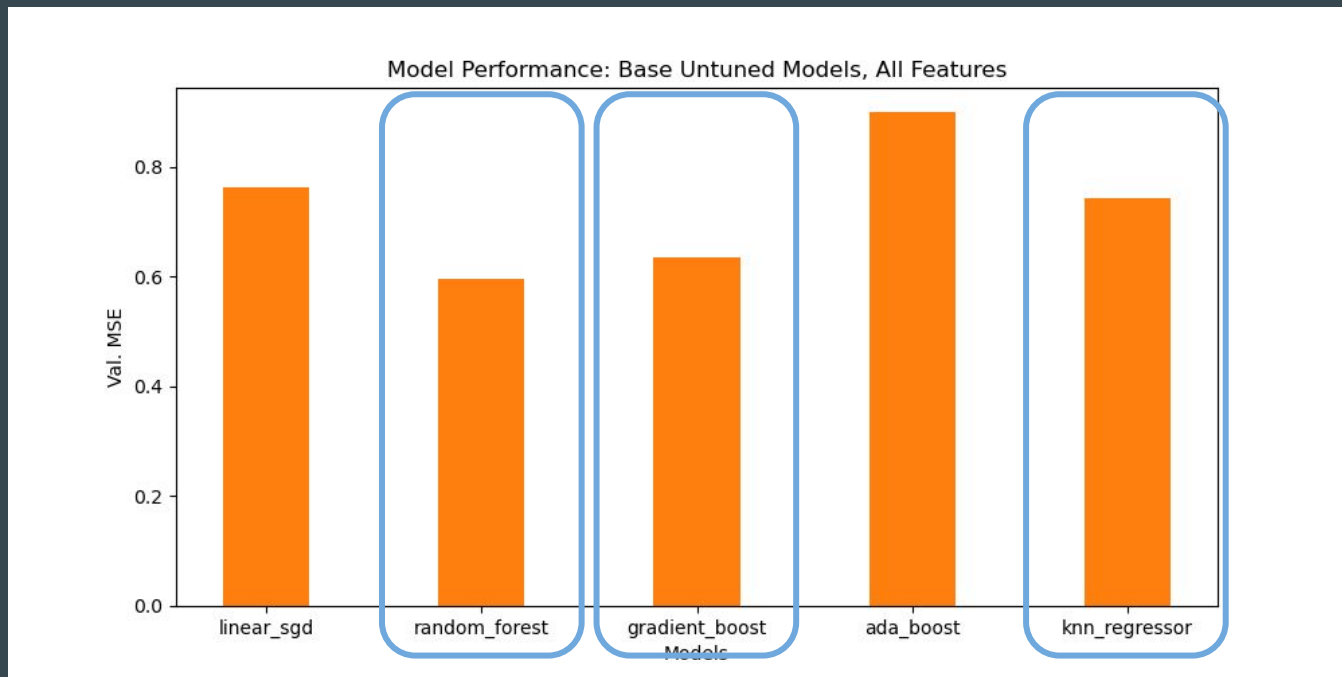
Model Type	SKlearn's Object	Pros	Cons
Linear Regression	SGDRegressor() <sup>6</sup>	Simple to implement and understand.	Doesn't work well for non-linearly separable dataset.
Random Forest	RandomForestRegressor() <sup>7</sup>	Less prone to overfitting.	Large number of trees can be slow to use. Can't extrapolate on data outside range of trained data.
Gradient Boosting	GradientBoostingRegressor() <sup>8</sup>	Fits each subsequent tree on the residuals which can learn very well.	Generally slower to fit than Random Forest and more prone to overfitting.
Adaptive Boosting	AdaBoostRegressor() <sup>9</sup>	Uses many decision stumps and each stump gets a weighted vote.	Not as robust to outliers as it tries to fit to every datapoint.
K-Nearest Neighbors	KNeighborsRegressor() <sup>10</sup>	Simple to understand and low number of hyperparameters.	Not as efficient as dataset grows.



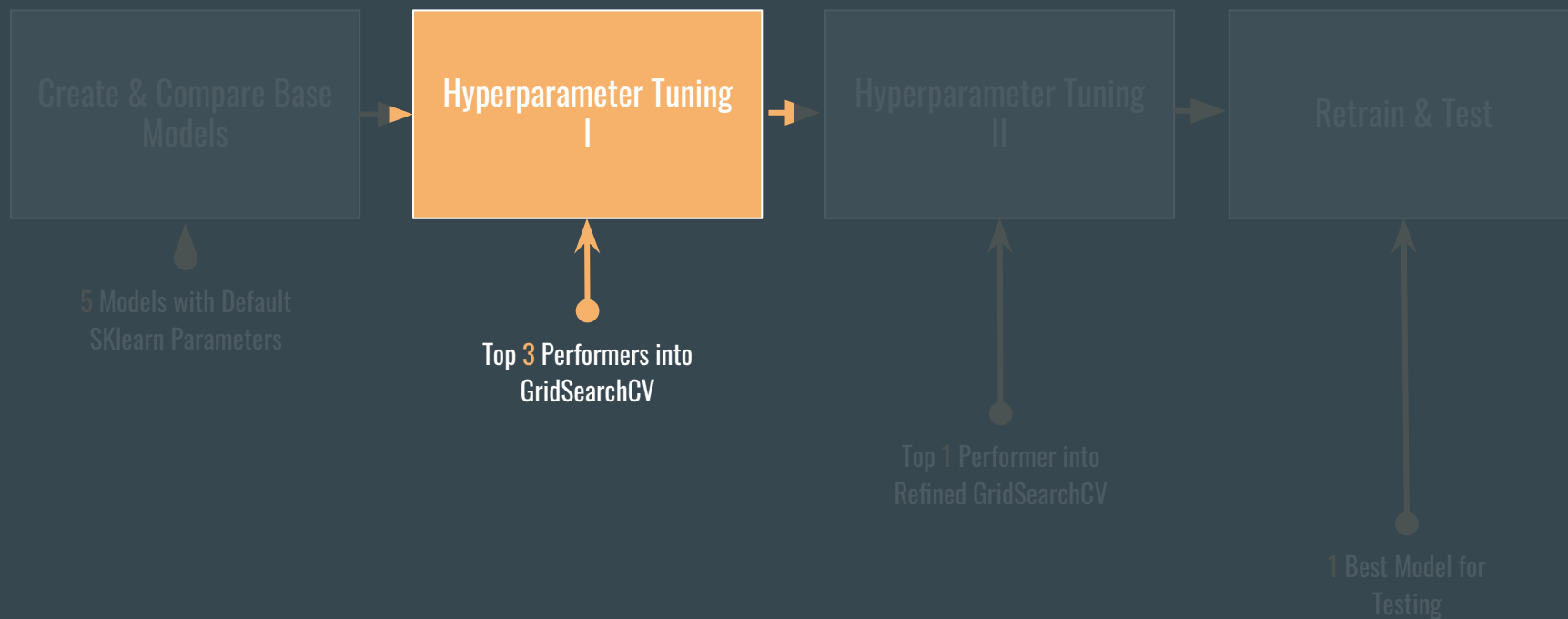
# Modeling: Base Models' Results



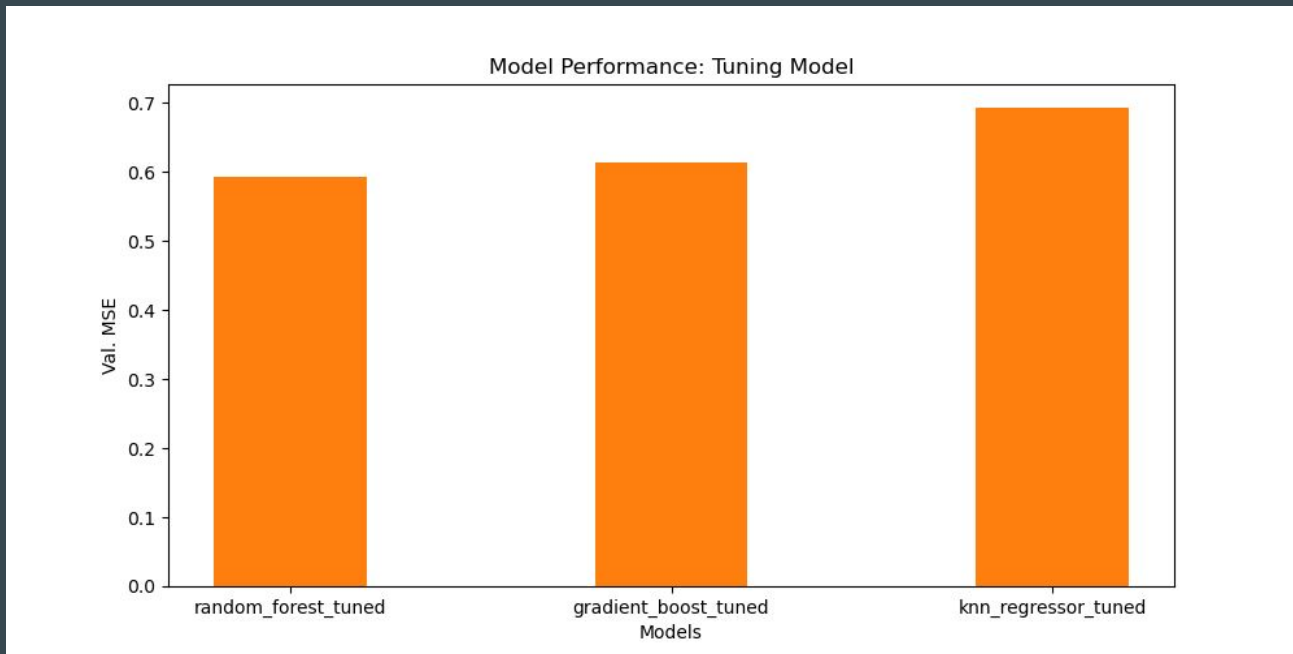
# Modeling: Base Models' Results



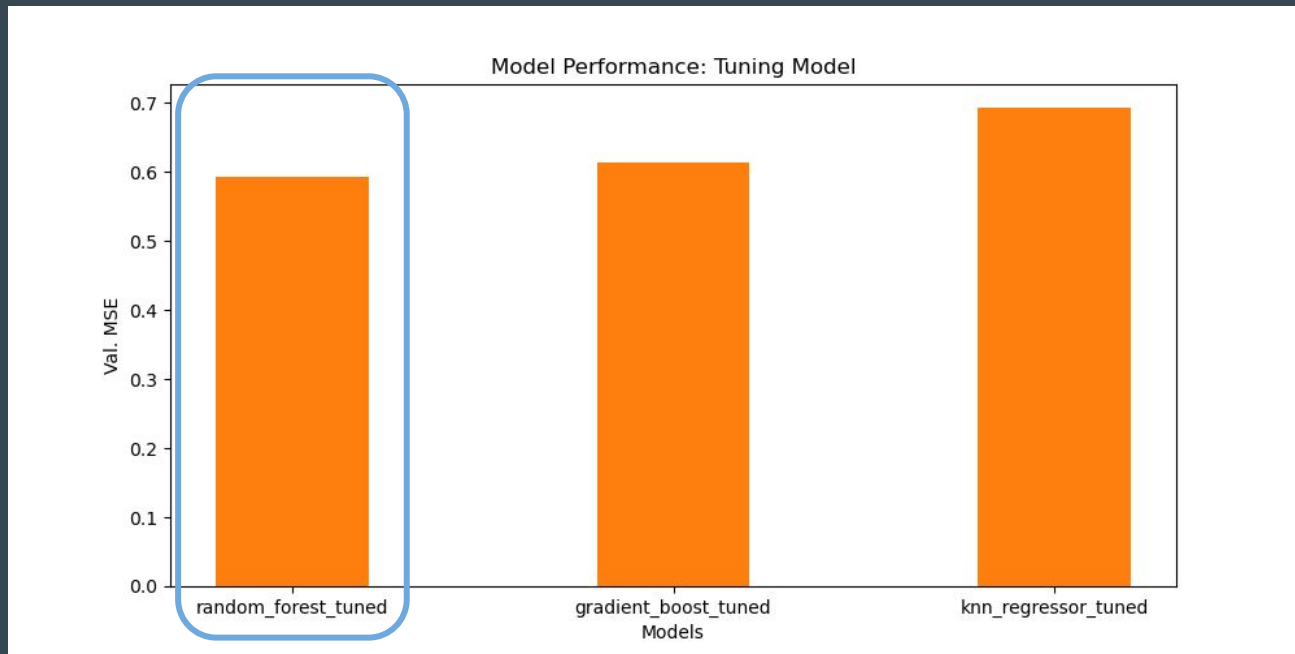
# Modeling: Hyperparameter Tuning I



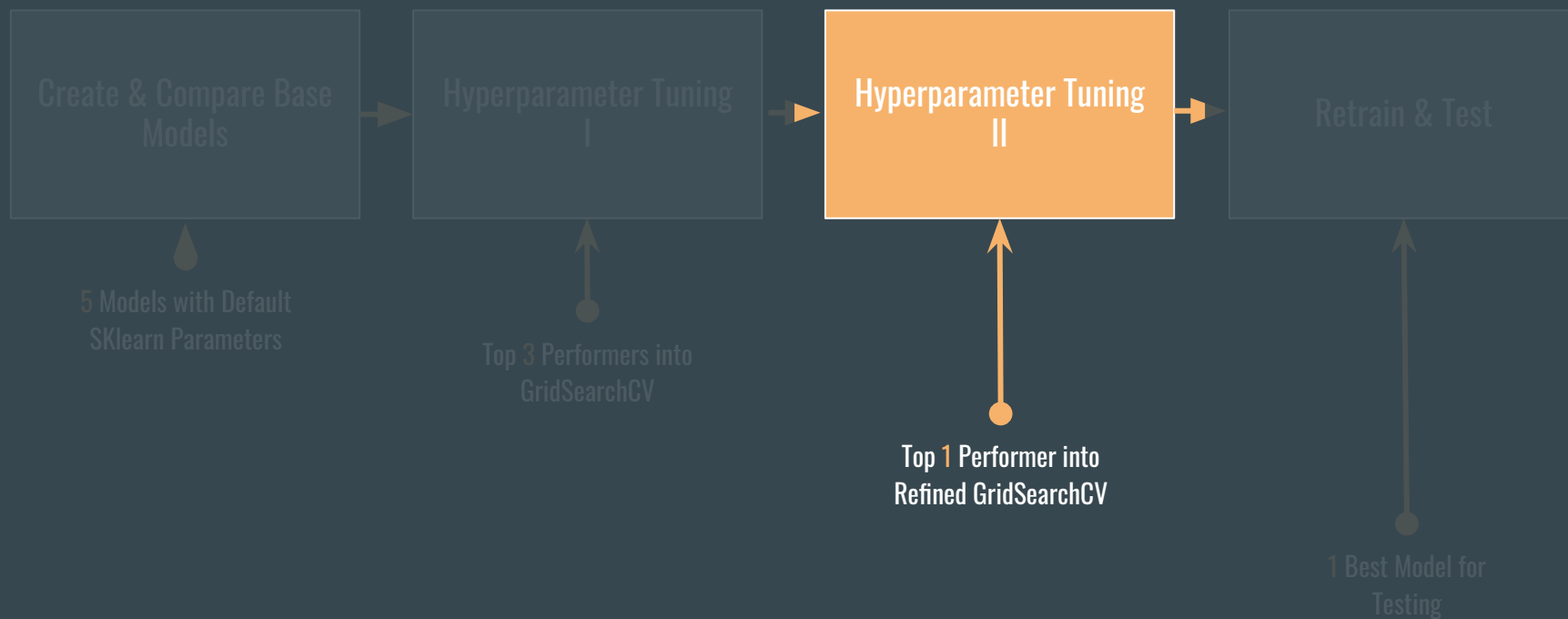
# Modeling: Hyperparameter Tuning I



# Modeling: Hyperparameter Tuning I

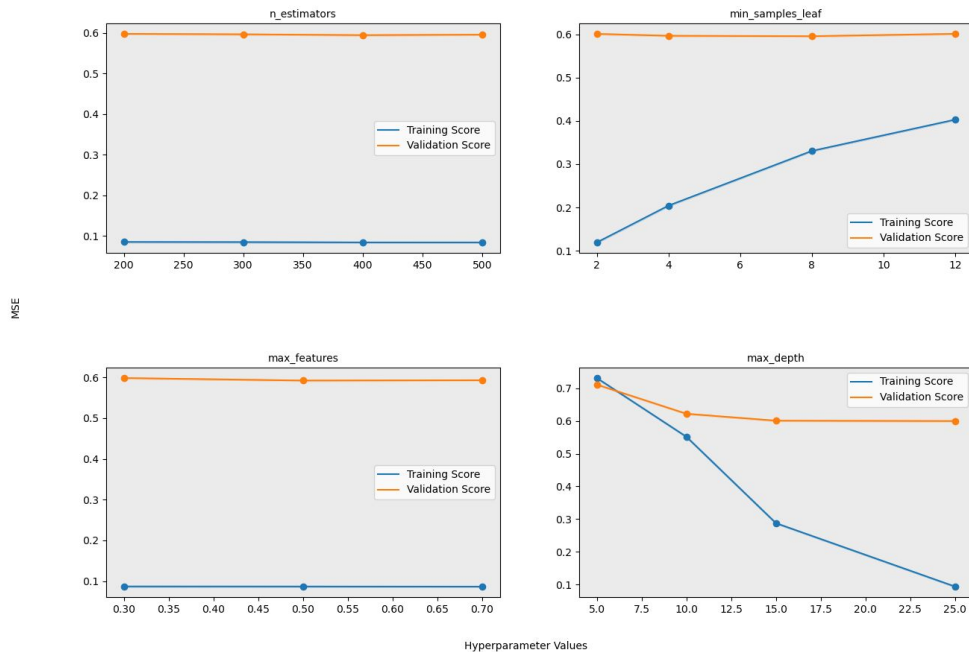


# Modeling: Hyperparameter Tuning II



# Modeling: Hyperparameter Tuning II

Validation Curves for random\_forest\_tuned



## Best Performing Hyperparameters

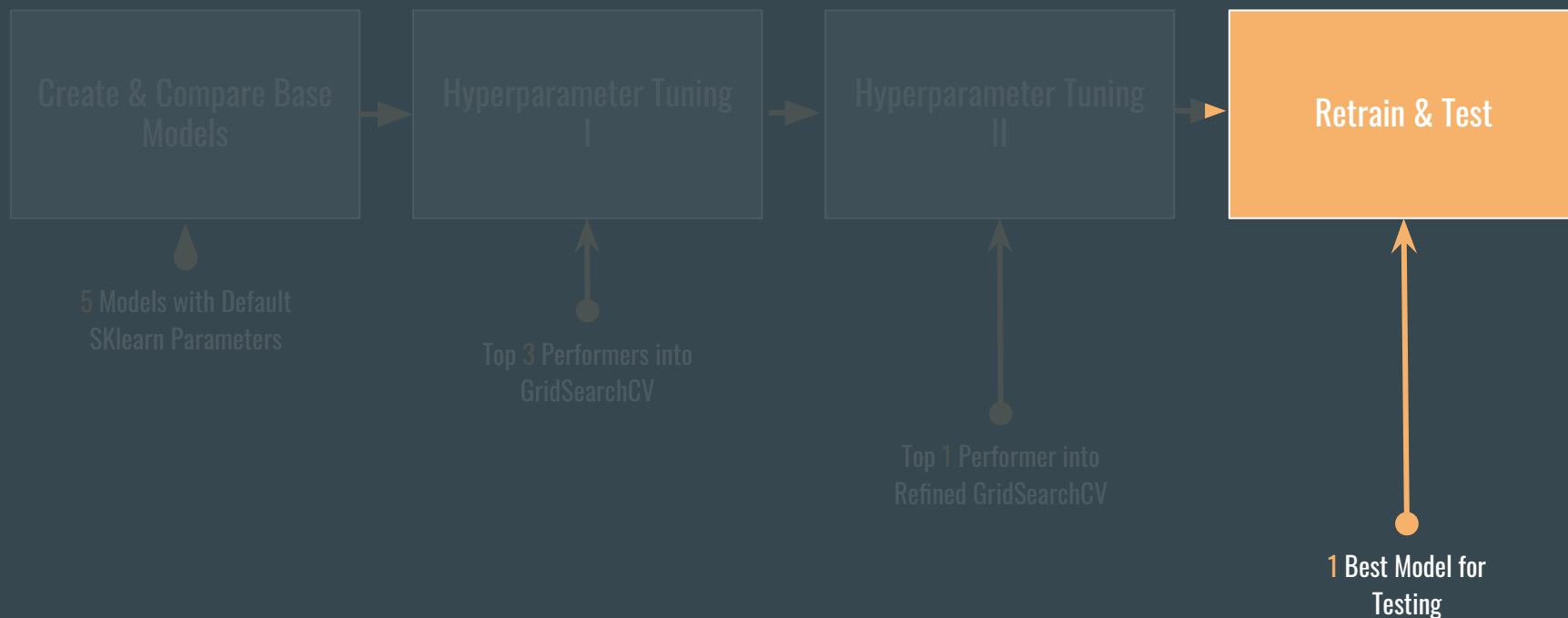
n\_estimators = 500

min\_sampe\_leaf = 2

max\_feature = 0.7

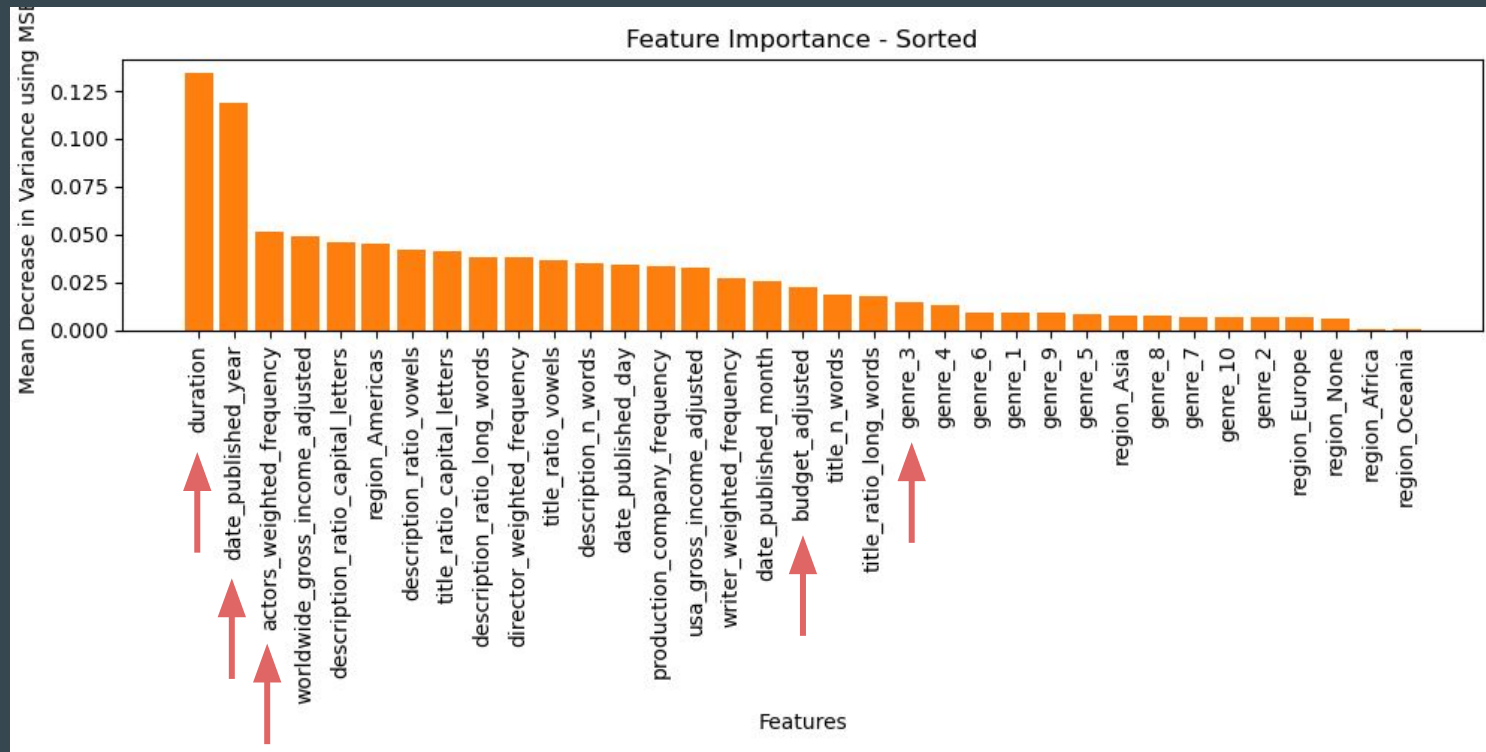
max\_depth = 25

# Modeling: Retrain & Test



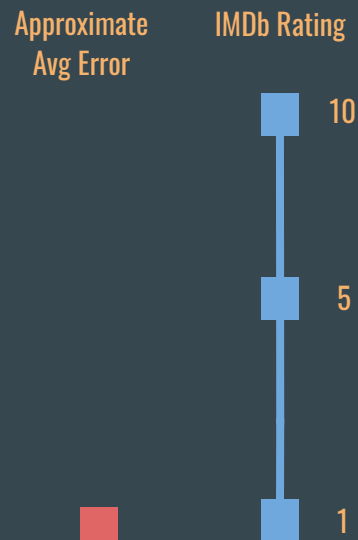


# Modeling: Retrain & Test



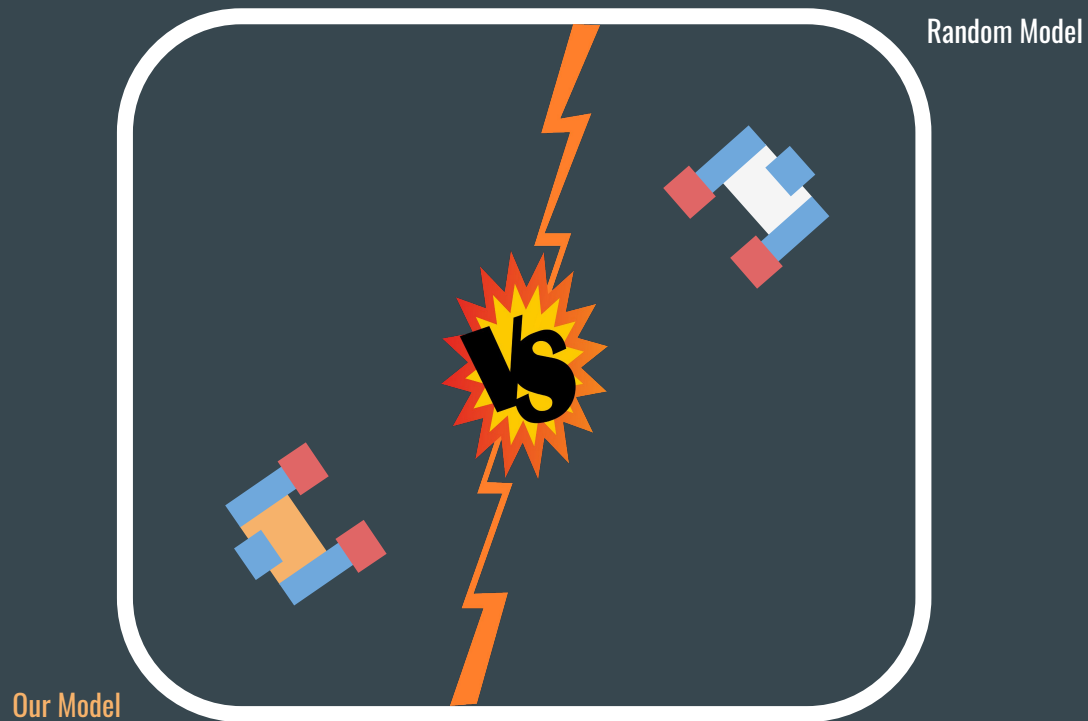
# Results: Test Data

Testing Results
MSE = 0.602
MSE (Inverse Scaling) = 0.869
RMSE (Inverse Scaling) = 0.931



\*Not to scale, just to gain some intuition into how big or small our error is compared against the IMDb Rating range

# Results: Model Evaluation



# Results: Model Evaluation

Our Model



- (4-0): Beat out 4 other models and many more hyperparameter combinations
- Powered by Random Forest
- Highly trained

Random Model

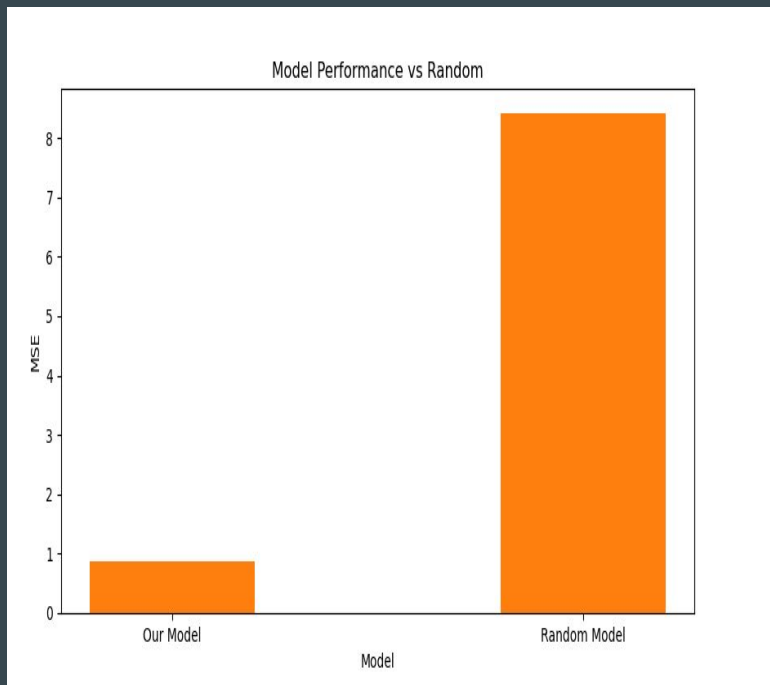


- Is literally just this:  

```
test_Y_random = np.random.uniform(1, 10,  
size=test_Y.shape)
```
- Thinks he's really good at guessing

Place Your Bets!

# Results: Model Evaluation - vs Random Model



Actual Rating	Our Model Predicted Rating	Random Model Predicted Rating	Smaller Error
4.7	5.3	5.4	RM
7.2	6.2	9.7	OM
6.8	6.7	4.5	OM
6.4	6.2	1.7	OM
6.5	5.6	4.2	OM
6.3	6.0	2.5	OM
5.7	6.8	9.8	OM
7.2	6.0	8.9	OM
4.4	6.3	5.4	RM
6.8	6.4	4.6	OM

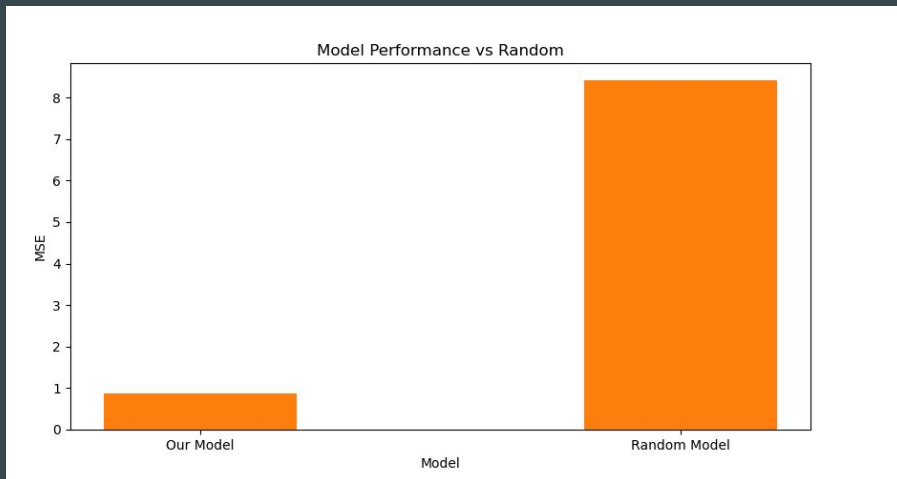
Our model has **smaller error 84%**  
of the time in test data

# Results: Model Evaluation - Distributions Are Different?

## Mann Whitney U Test

Test	H0	p-value	Action
Mann-Whitney U Test	The 2 distributions are equal	3.80e-29	Reject H0

# Results: Model Evaluation - vs Random Model



**Our Model is statistically different from  
Random Model and is a better predictor!**



**Our Model**

# Conclusions

There is a relationship between movies' features and their IMDb score

We were able to successfully create a model that accurately predicts IMDb score based on a movie's features



# Areas of Improvement

Find a dataset on actors w/o missing values

Create a better encoding method for genre

Group movies based on their genres

Reduce the dimensionality of our data

Increase efficiency and speed of modelling



# Future Research

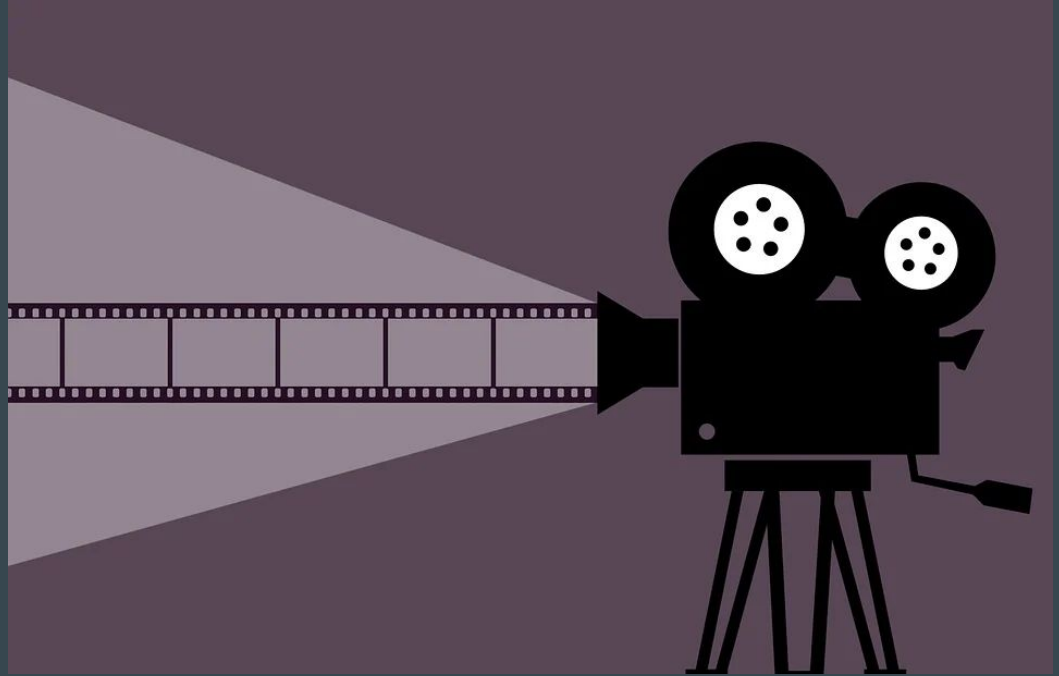
- Different target variables
  - Male versus female votes
  - Gross income
- Titles and description
  - Natural language processing
  - Movie scripts
- Other variables?
  - Screen time
  - Sequel vs non-sequel
  - Release location/method
- Different modeling techniques
  - Voting regressor model
  - Different modeling packages
  - More advanced modeling techniques

# The Bigger Picture

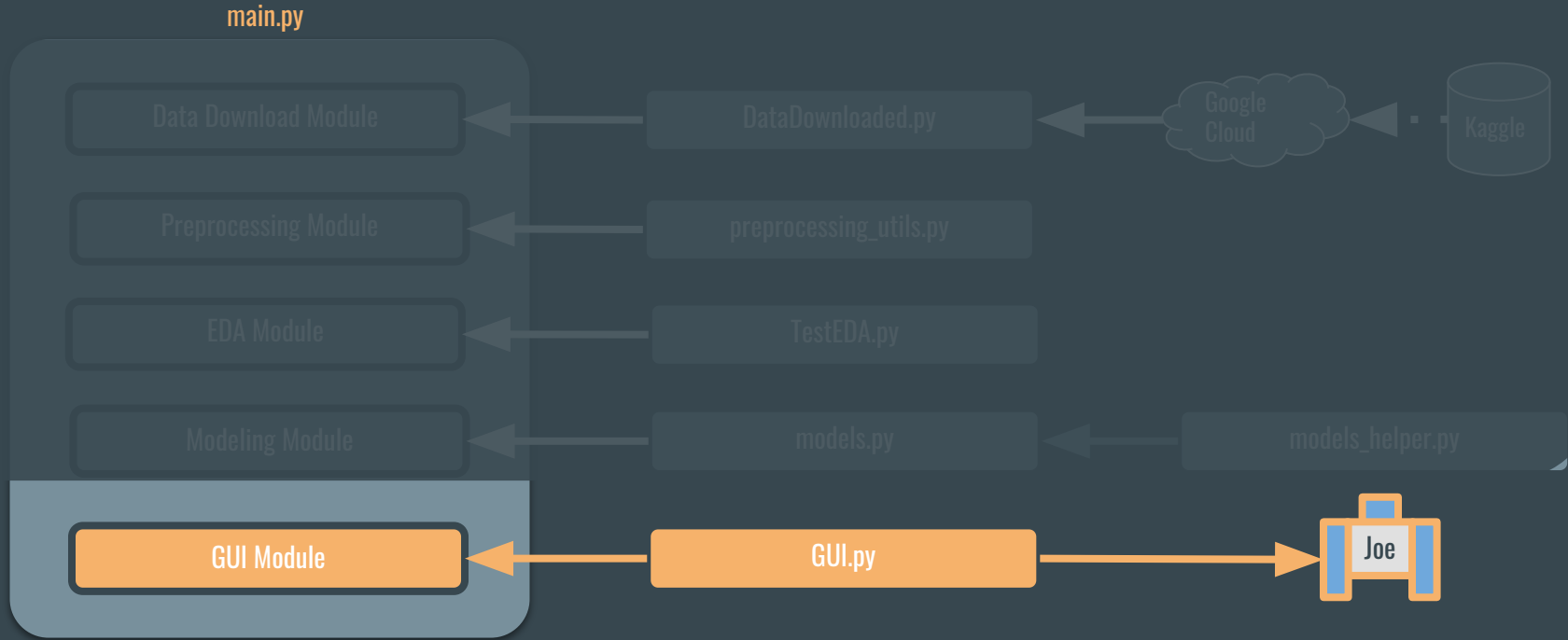
What does this research mean for the movie industry?

How can predictive analytics be used by movie production companies?

Does this pose any problems?

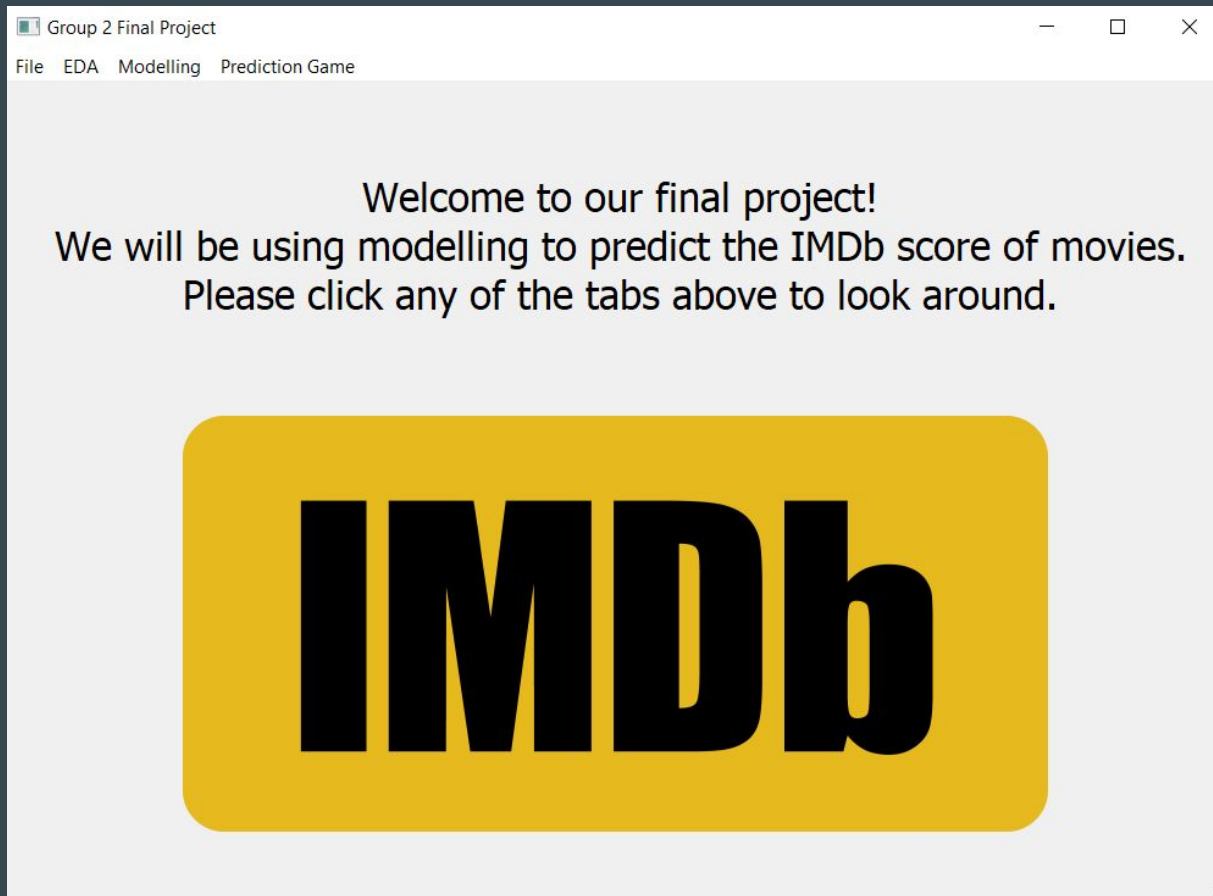


# GUI Demo

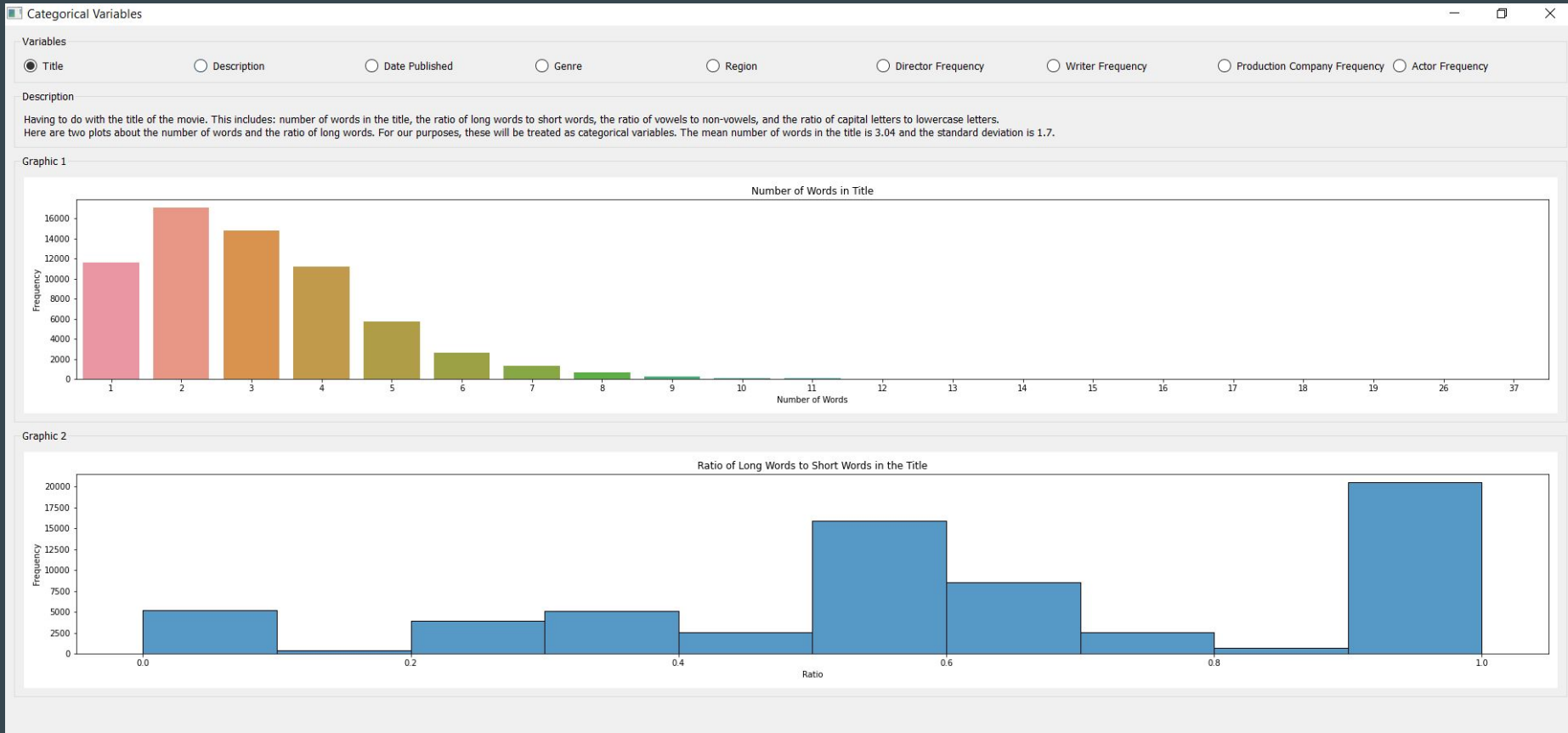


<https://github.com/Saharae/Final-Project-Group2>

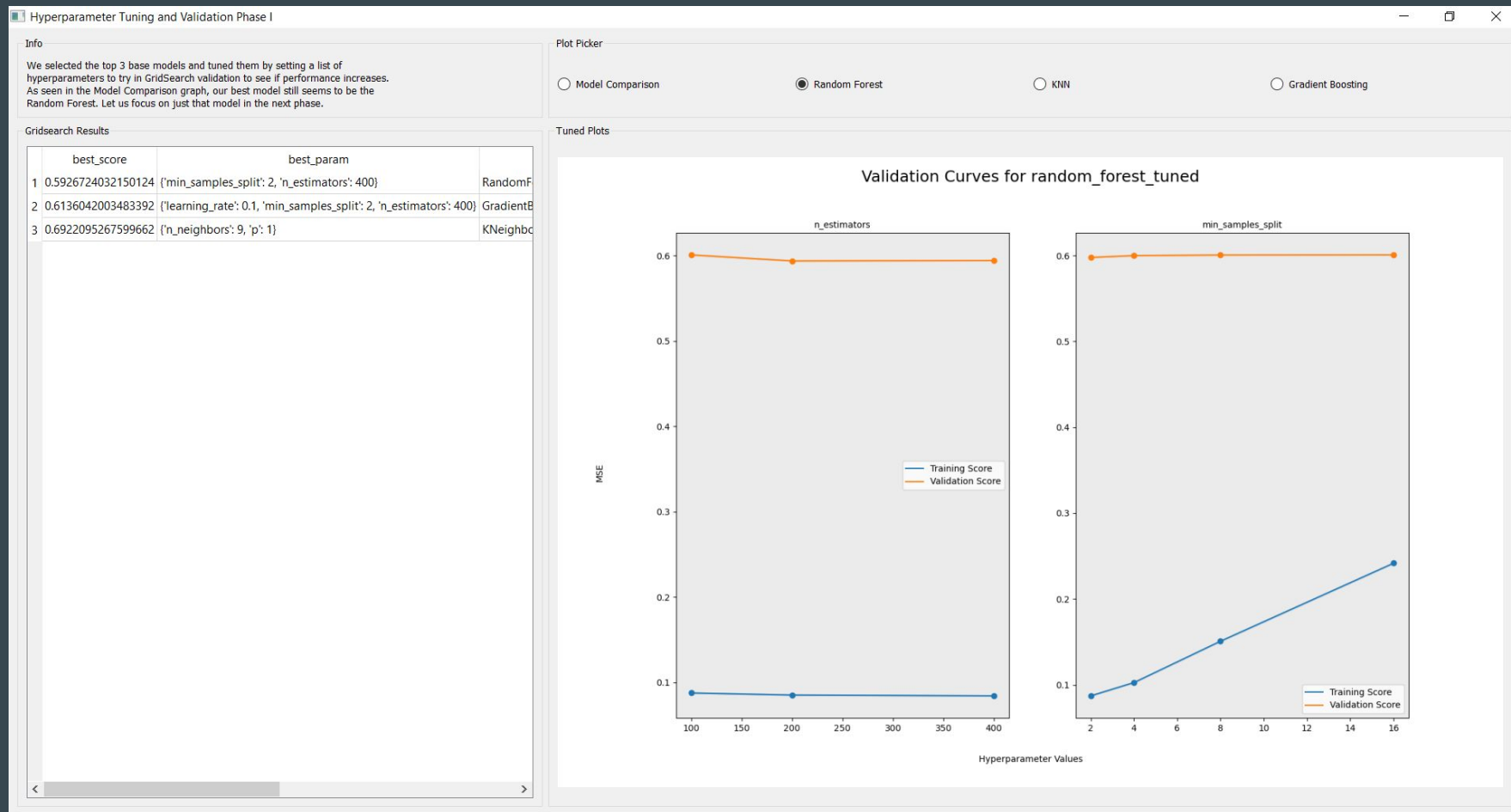
# GUI Images Part 1



# GUI Images Part 2



# GUI Images Part 3



# GUI Images Part 4

Prediction Game

Description

This tool will allow you to make predictions against our best model, to see who can come out on top!

We will give you a list of features of a movie selected randomly from our test set and you will predict the weighted average score.

Our model will also predict the weighted average score, and whoever comes the closest to the real score will win!

Since you are presumably a human, we will give you human readable features for you to make your guess.

Remember, no cheating by looking up the movie online. And if any of the features are missing, it is because they were not in the IMDb dataset, so our model did not get them either.

Random Movie Generator

Generate

Your Movie's Features

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	Duration	Title	Date Published	Director	Writer	Production ...	Actors	Description	Budget	USA Gross ...	Worldwide Gros...	Genre 1	Genre 2	Genre 3	Region
2	87	Un mundo ...	3/15/2013	Gabriel Mariño	Gabriel Mariño	Sobrevivientes ...	Lucía Uribe, ...	A middle-class ...	nan	nan	nan	Drama	None	None	Americas

Input your guess

Lock In!

Results

The results are in...

You Predicted: 6

Our model predicted: 5.815897632

The actual weighted average vote is: 5.7

You lose!



Questions?

# Appendix: Extra Slides

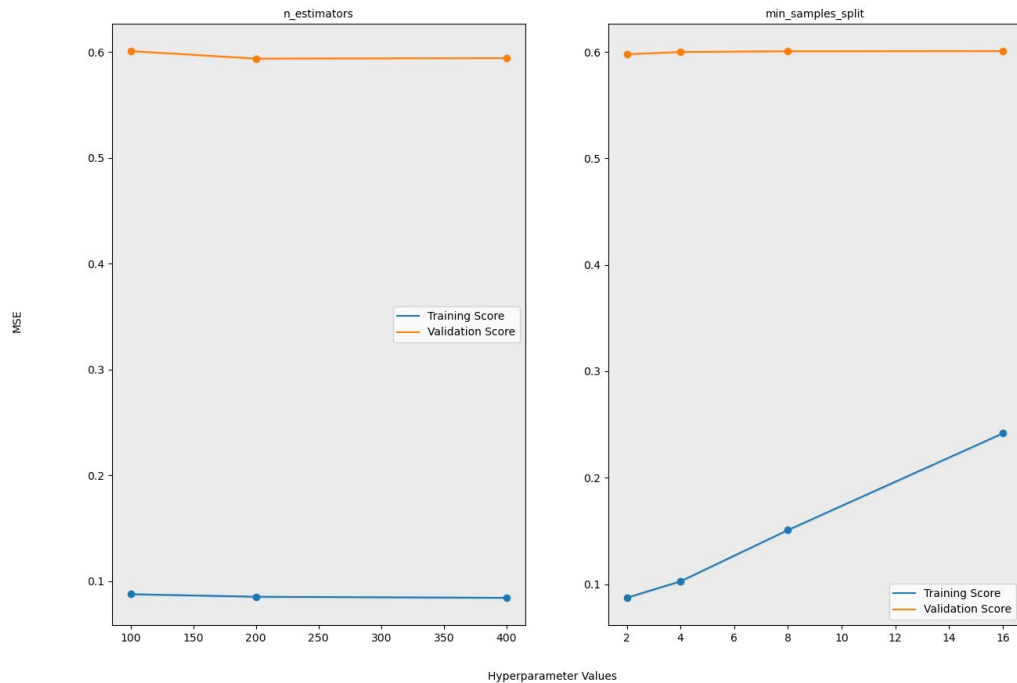
# Modeling: Hyperparameters

Model Type	Hyperparameter	Hyperparameter Meaning	Dtype and Default
Random Forest*	n_estimators	The number of trees in a forest.	int, default = 100
	min_samples_split	Minimum samples required to split an internal node.	int or float, default = 2
	min_samples_leaf	Minimum samples to be at a leaf node.	int or float, default = 1
	max_features	Number of features to consider for each tree.	{'auto', 'sqrt', 'log2'} or int or float, default = 'auto'
	max_depth	Max depth allowed for each tree	int, default = None
Gradient Boosting	learning_rate	Shrinks the contribution of each tree.	float, default = 0.1
	n_estimators	The number of boosting stages to use.	int, default = 100
	min_samples_split	Minimum samples required to split an internal node.	int or float, default = 2
K-Nearest Neighbors	n_neighbors	Number of neighbors to use.	int, default = 5
	p	p = 1 (Manhattan Distance), p = 2 (Euclidean Distance)	int, default = 2

\*Not all listed hyperparameters tested in this first phase  
Source: SKlearn model documentations.

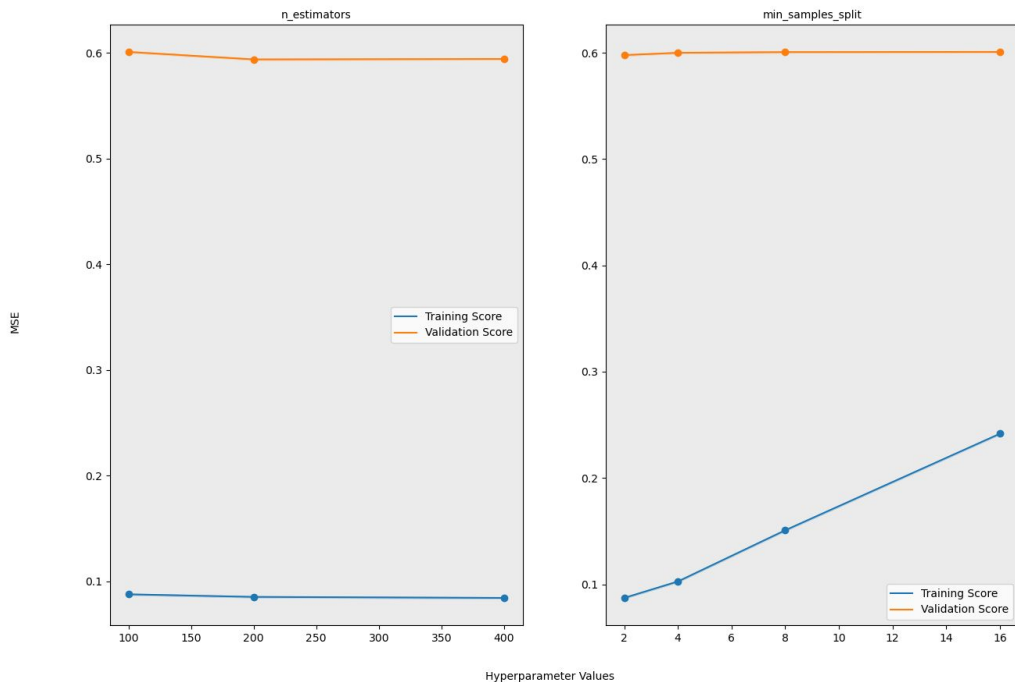
# Modeling: Hyperparameter Tuning I

Validation Curves for random\_forest\_tuned



# Modeling: Hyperparameter Tuning I

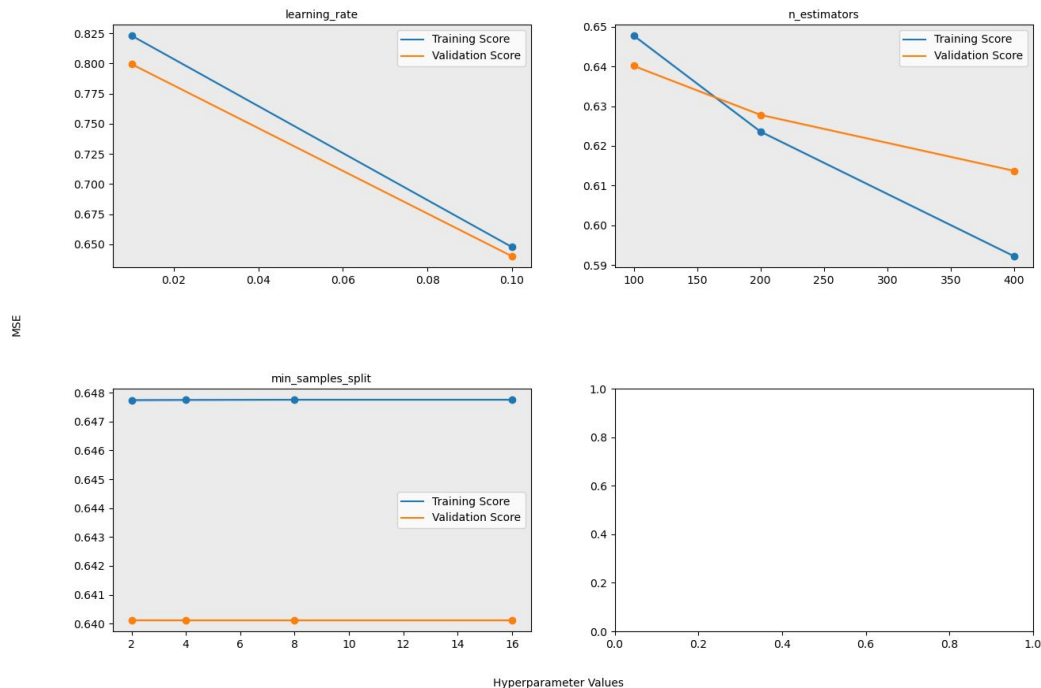
Validation Curves for random\_forest\_tuned



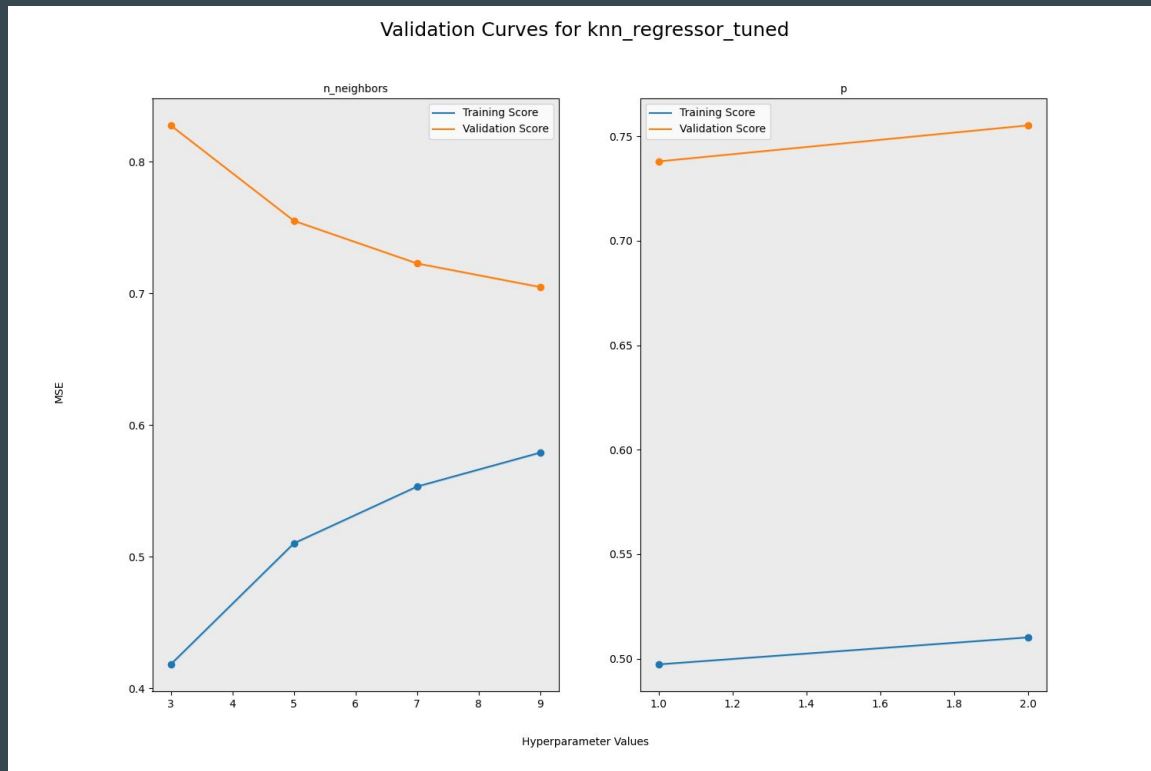
Possible  
overfitting

# Modeling: Hyperparameter Tuning I

Validation Curves for gradient\_boost\_tuned

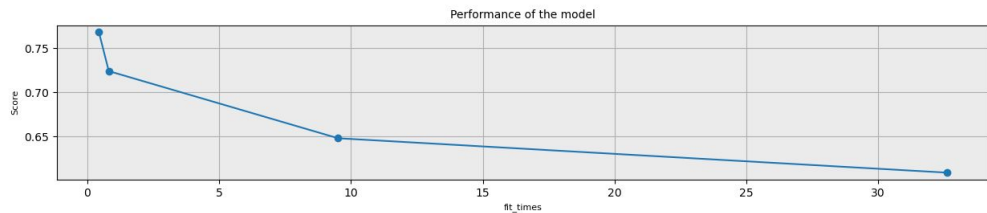
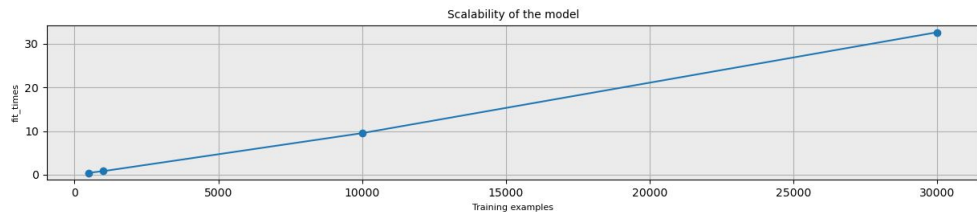
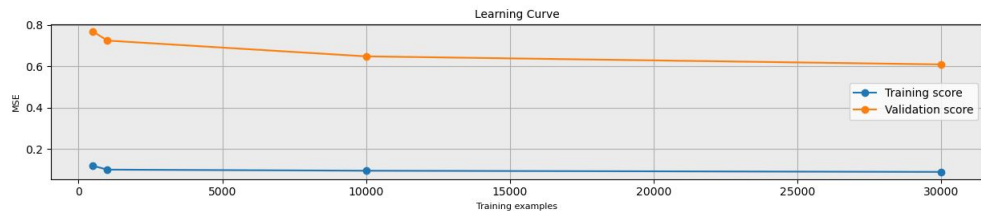


# Modeling: Hyperparameter Tuning I



# Modeling: Hyperparameter Tuning II

Learning Curve for random\_forest\_tuned





# Results: Model Evaluation - Distributions Are Independent?

## 2 Sample T-Test

Test Type	Test	H0	p-value	Action
Test for normality on Our Model Ratings	Shapiro-Wilk Test	Distribution is normal	4.75e-40	Reject H0
Test for normality on Random Model Ratings	Shapiro-Wilk Test	Distribution is normal	0	Reject H0
Test for equal variance	Bartlett Test	Variances are equal	0	Reject H0

# Results: Model Evaluation - Distributions Are Independent?

~~2 Sample T-Test~~

Mann Whitney U  
Test

Test	H0	p-value	Action
Mann-Whitney U Test	The 2 medians of the distributions are the same	3.80e-29	Reject H0

# IMDb License

IMDb, IMDb.COM, and the IMDb logo are trademarks of IMDb.com, Inc. or its affiliates.