# EXPLORING IMDB MOVIES DATASET: EDA AND ML MODELING IN PYTHON

MSDS422B GROUP 9:

**ANAMIKA KUMARI** 

**KARAN RASTOGI** 

**SHRIMATHY SRINIVASAN** 

**GANESH SUBRAMANIAN AUDHIKESAVAN** 



## OUTLINE

- Introduction
- Dataset Walkthrough
- Tools and Libraries Used
- Exploratory Data Analysis
- Data Visualization
- Machine Learning Modeling
- Challenges
- Conclusion & Future Scope



#### INTRODUCTION



The magic of movies transcends borders and languages, touching hearts and sparking imaginations across the globe. From timeless classics to groundbreaking blockbusters, every film has a story to tell, a message to convey, and an impact on its audience.

With IMDb, the world's most comprehensive and trusted source for movie information, we have the power to unlock the secrets behind these cinematic wonders.

With this project, we want to delve deeper into the following:

- 1. Discovering Hidden Gems of Movies
- 2. Genre and Audience Analysis
- 3. Rating and Popularity Patterns
- 4. Directors and the Revenue Budget Analysis
- 5. NLP Analysis on Keywords and Genres



## **DATASET WALKTHROUGH**

**IMDb** 

- Dataset: imdb-movies.csv (Movies information data set from year 1960 to 2015)
- Link: https://www.kaggle.com/datasets/schoolofaitvm/imdbdataset
- This dataset has 10867 records and 21 columns.
- Features/Columns:

id	imdb_id	popularity	budget	revenue	original_ title	cast	homepage	director	tagline	keywords	overview
runtime	genres	production_c ompanies	releas	se_date	vote_count	vote_avera ge	release_ year	budget_a dj	revenue_ adj		

#### Features Type:

S.No	Feature Type	Columns
1.	Int	Id, imdb_id, budget, revenue, runtime, vote_count
2.	Float	Popularity, vote_average, budget_adj, revenue_adj
3.	String	Original_title, cast, homepage, director, tagline, keywords, overview, genres, Production companies,
4.	Date	Release_date, release_year



## **TOOLS AND LIBRARIES USED:**

- Anaconda IDE
- Jupyter Notebook
- Python
- Google Collab

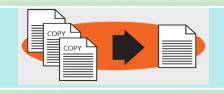
#### Libraries and Packages Used:

- Sklearn
- Matplotlib
- Seaborn
- Numpy
- Pandas
- Tensorflow & Keras

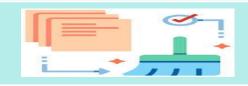












#### **Data Duplication**

- 1. Some entries were duplicated across rows
- 2. Duplicated entries were removed

#### **Null Data Values**

- 1. Columns such as cast(76), homepage (7930), director(44), tagline(2824), keywords (1493), production companies (1030), genres(23) have null values.
- 2. Columns which have low null values have their respective records dropped.

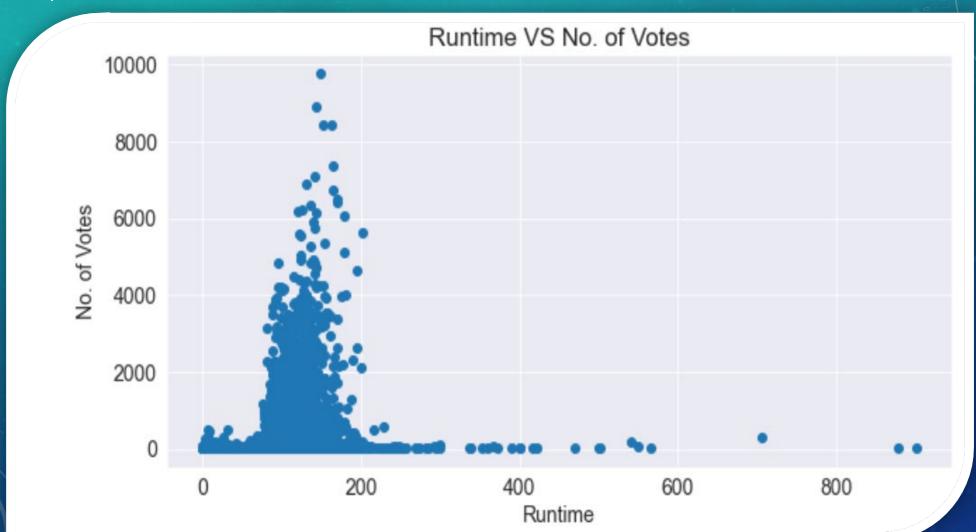
## **Dropping Unnecessary Columns**

- 1. 'homepage',' tagline' were dropped as they contain non relevant information
- 2. NaN values in keywords and production\_companies were replaced with Not Known string for analysis purpose.

## **EXPLORATORY DATA ANALYSIS**

**IMDb** 

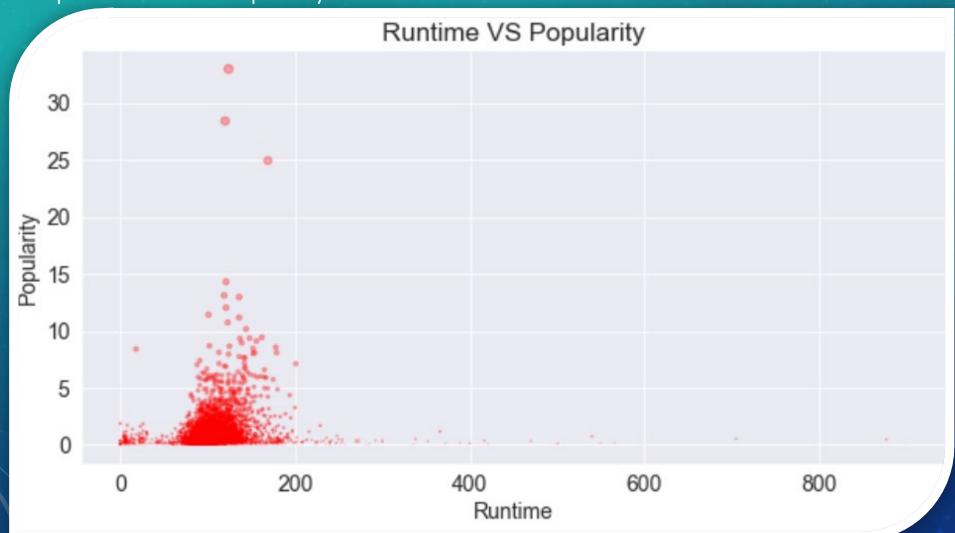
• Graph of Runtime vs Number of Votes



Northwestern
SCHOOL OF
PROFESSIONAL STUDIES

## **EXPLORATORY DATA ANALYSIS**

Graph of Runtime vs Popularity



**IMDb** 



## **CORRELATION HEAT MAP**

#### **IMDb**

#### **Correlation of Movie Features**

/ /										
id	1	-0.0093	-0.14	-0.097	-0.084	-0.033	-0.072	0.51	-0.19	-0.14
popularity	-0.0093	1	0.54	0.66	0.14	0.8	0.22	0.093	0.51	0.61
budget	-0.14	0.54	1	0.73	0.19	0.63	0.088	0.12	0.97	0.62
revenue	-0.097	0.66	0.73	1	0.16	0.79	0.18	0.059	0.71	0.92
runtime	-0.084	0.14	0.19	0.16	1	0.16	0.18	-0.12	0.22	0.18
vote_count	-0.033	0.8	0.63	0.79	0.16	1	0.26	0.11	0.59	0.71
vote_average	-0.072	0.22	0.088	0.18	0.18	0.26	1	-0.13	0.1	0.2
release_year	0.51	0.093	0.12	0.059	-0.12	0.11	-0.13	1	0.02	-0.065
budget_adj	-0.19	0.51	0.97	0.71	0.22	0.59	0.1	0.02	1	0.65
revenue_adj	-0.14	0.61	0.62	0.92	0.18	0.71	0.2	-0.065	0.65	1
	.19	popularity	budget	revenue	runtime	vote_count	vote_average	release_year	budget_adj	revenue_adj

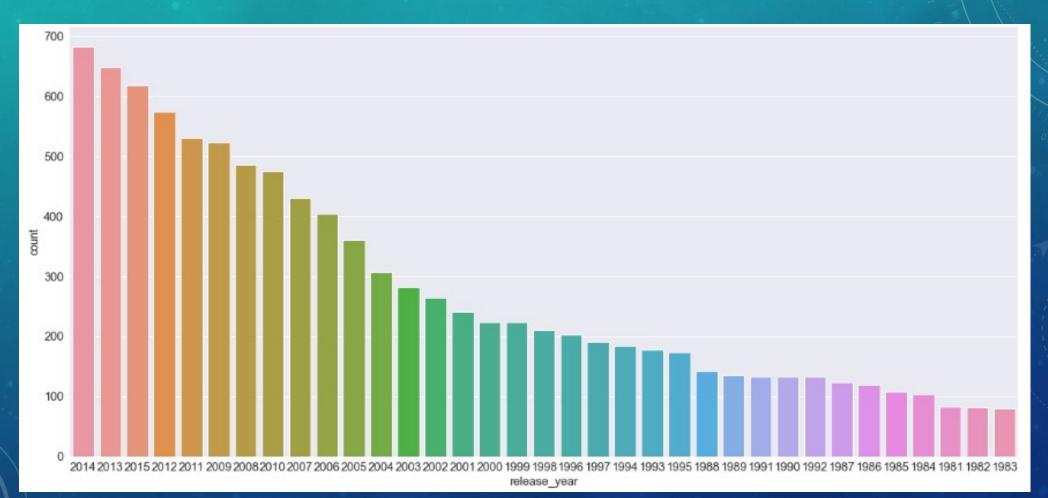
- 1.0 - 0.8 - 0.6 - 0.4 - 0.2 - 0.0

Northwestern SCHOOL OF PROFESSIONAL STUDIES

## **DATA VISUALIZATION**

**IMDb** 

Year wise movie count by graph

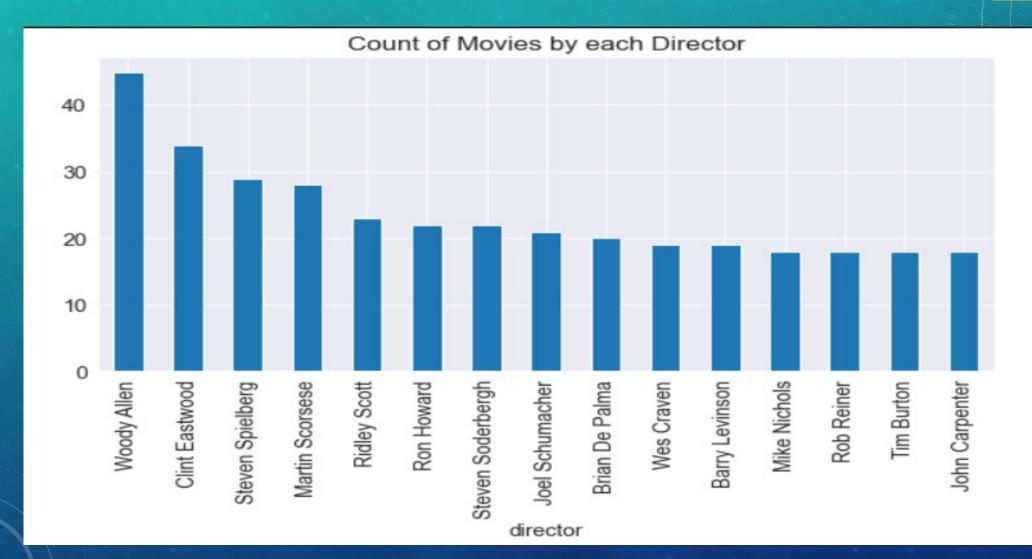


Northwestern SCHOOL OF

SCHOOL OF PROFESSIONAL STUDIES



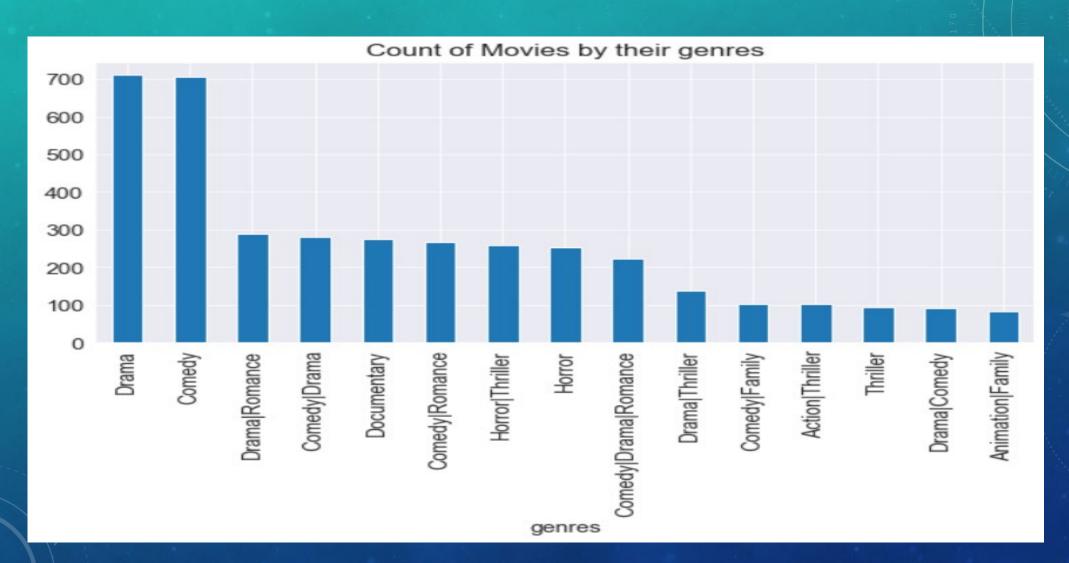








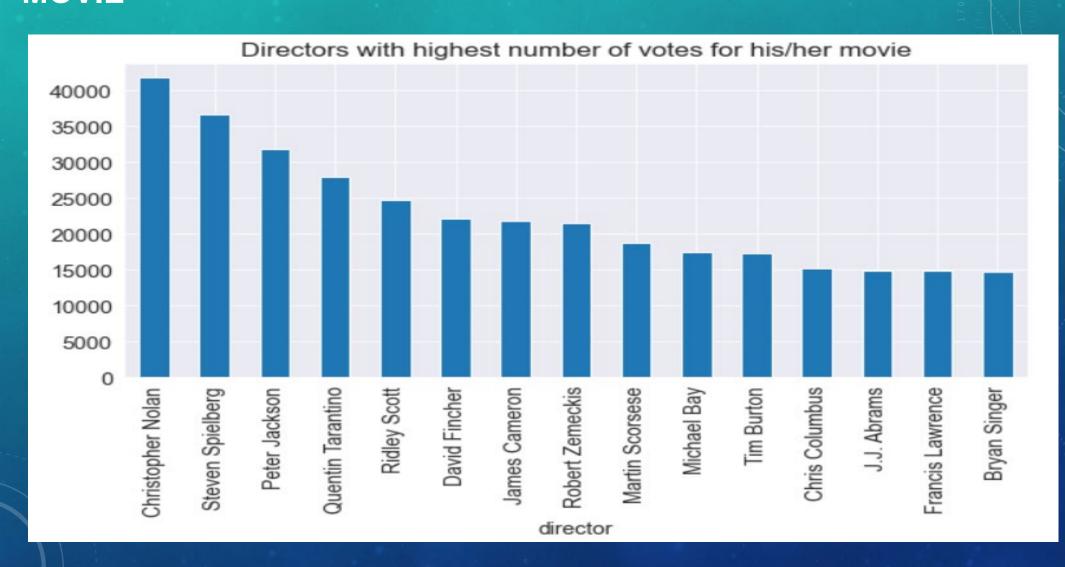






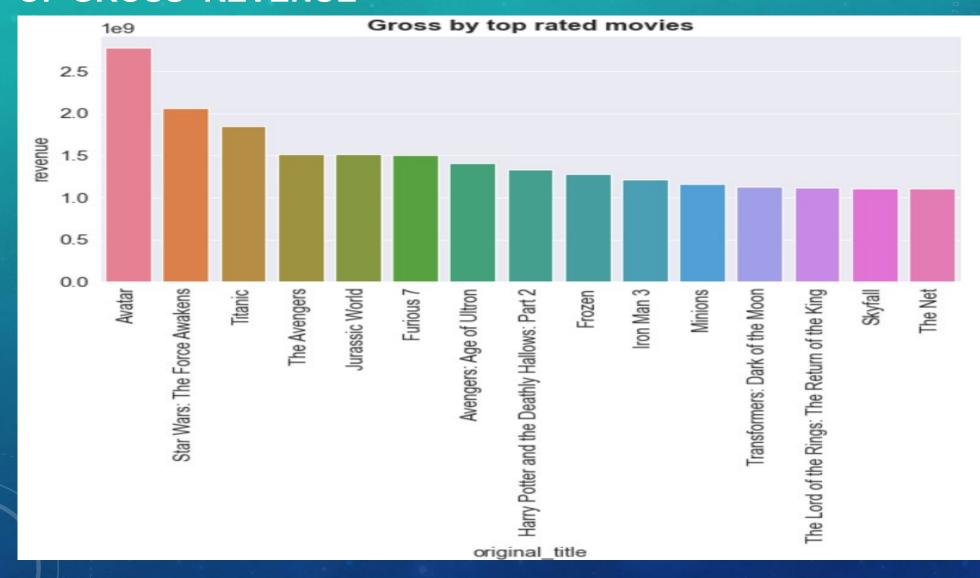
## DIRECTORS WITH HIGHEST NUMBER OF VOTES FOR HIS/HER MOVIE







## BAR PLOT TO SHOW TOP 15 MOVIES BY DESCENDING ORDER OF GROSS REVENUE

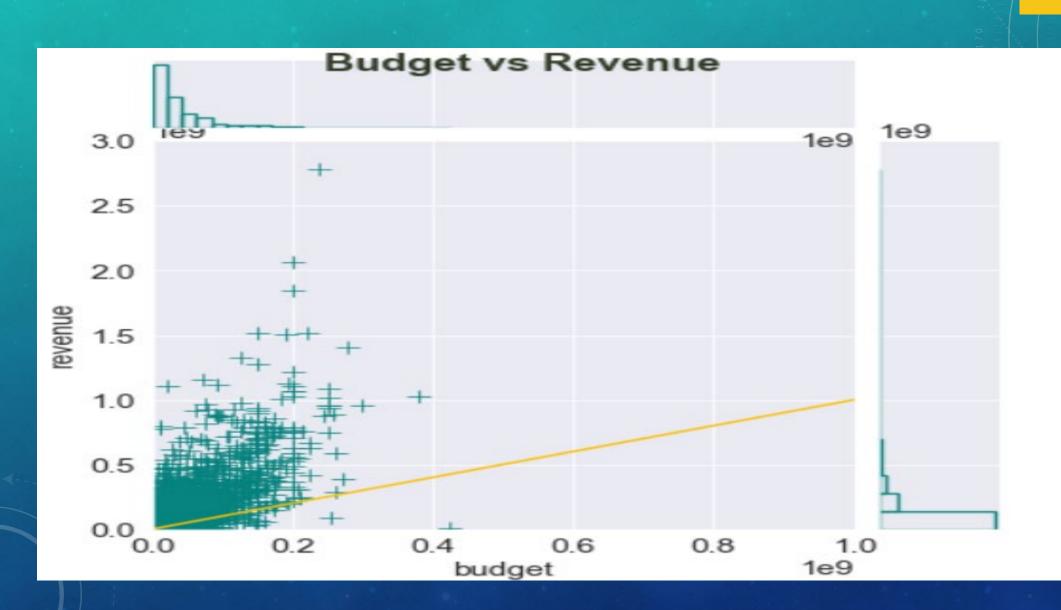






#### **SCATTERPLOT: BUDGET VS REVENUE**

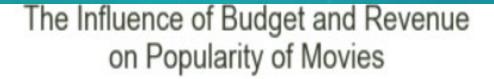


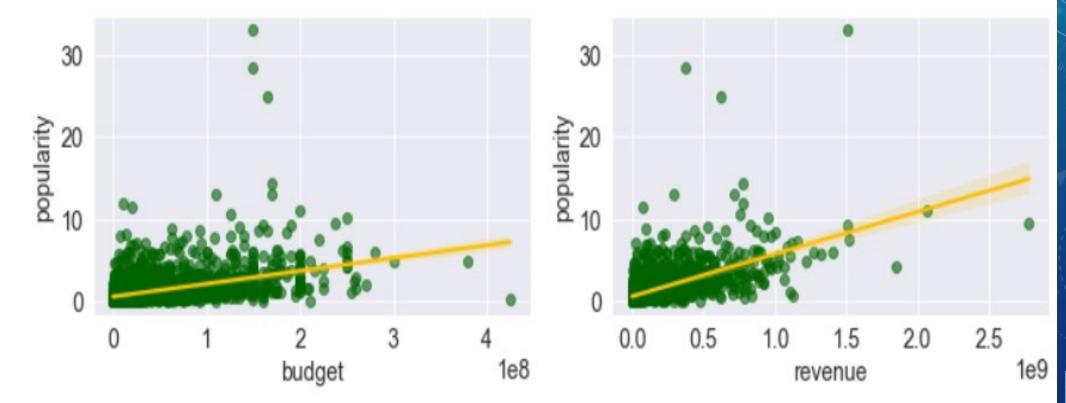


Northwestern
SCHOOL OF
PROFESSIONAL STUDIES

## INFLUENCE OF BUDGET AND REVENUE ON POPULARITY RATING







Northwestern
school of
PROFESSIONAL STUDIES

### ML MODELING: INTENT & USE CASES

Revenue:
Estimation and
Prediction

**Sentiment Analysis** 

**Influence Factors** 

- Linear Regression with Cross Validation
- Linear Regression with Regularization such as Lasso and Ridge
- XGBoost Regressor
- XGBoost Regressor with Cross Validation and Hyperparameter Tuning

- Natural Language
   Processing
   Multinomial NB for
   Keywords, Genres and
   overview features
- LSTM (Long Short-Term Memory)

- Neural Networks
- Gradient Boost



### **ML MODELS OBSERVATIONS AND RESULTS**

Regression Analysis

1. Dependent Variable: Revenue

2. Independent Variable: Popularity, Budget and Vote\_Count

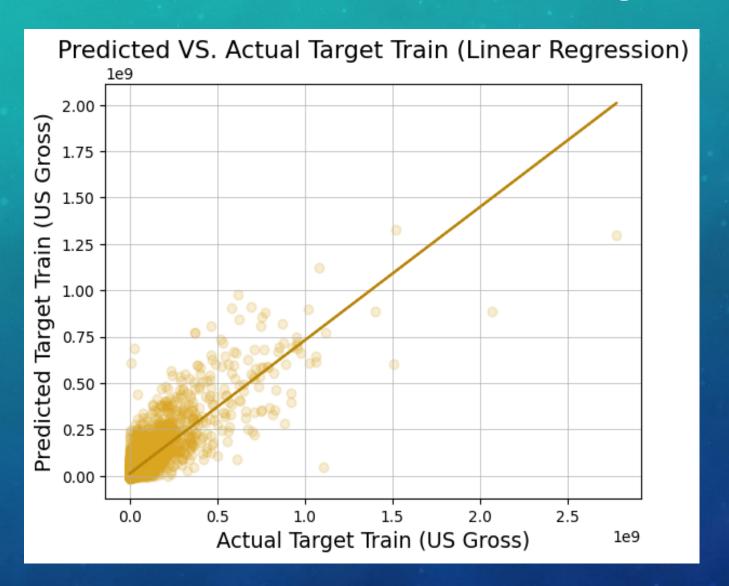
The models and their performance are as follows:

S.No.	Model Name	Training Score	Validation/Test Score
1.	Linear Regression	0.718	0.67
2.	Linear Regression K fold CV	0.718	0.71
3.	Polynomial Regression (degree=2)	0.768	0.69
4.	Ridge Regression	0.718	0.67
5.	Ridge Regression with CV	0.718	0.71
6.	Lasso Regression with CV	0.718	0.713

<sup>\*\*\*</sup> Model highlighted with purple cells are best performing model/s amongst all.



#### **Predicted vs Actual Values with Linear Regression**





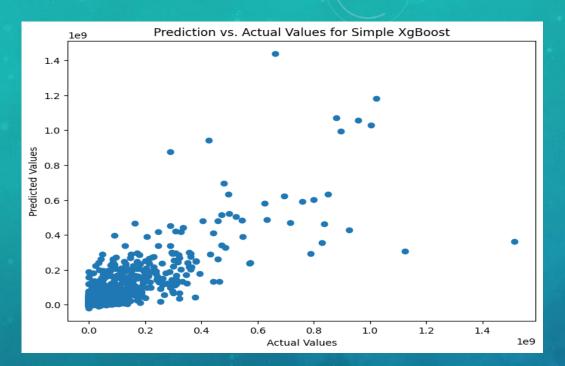
## **REGRESSION ANALYSIS (CONTD)**

 Further we also tried XGBoost Regressor with RFE and cross validation and hyperparameter Tuning.

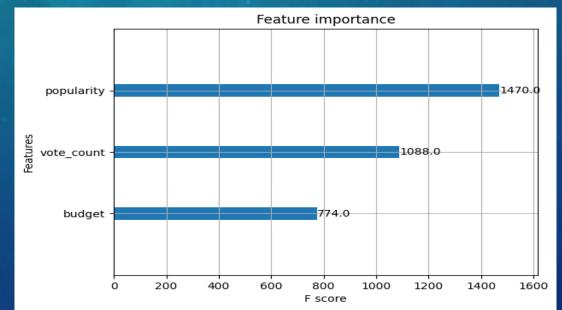
S.No	Regression Model	MAE	RMSE	R2
1.	XGBoost Regressor	26256665.263	75964208.27	0.64
2.	XGBoost Regressor with Cross Validation and Hyper parameter Tuning	22622848.21	71043604.95	0.7139
3.	XGBoost Regressor with CV and Hyperparameter tuning and RFE	26922083.78	67367538.39	0.668

<sup>\*\*\*</sup> Model highlighted with purple cells is best performing model amongst all.









Northwestern
SCHOOL OF
PROFESSIONAL STUDIES

### **SENTIMENT ANALYSIS**

Aim is to use NLP for Sentiment analysis on 'Keywords' and 'Genres' columns to get the Positive,
 Neutral and Negative Sentiment. Metrics Results are as follows:

Sentiment Ana	•	•		
	precision	recall	f1-score	support
negative	0.92	0.47	0.62	285
neutral	0.85	0.98	0.91	1298
positive	0.88	0.68	0.77	291
accuracy			0.86	1874
macro avg	0.88	0.71	0.77	1874
weighted avg	0.87	0.86	0.85	1874

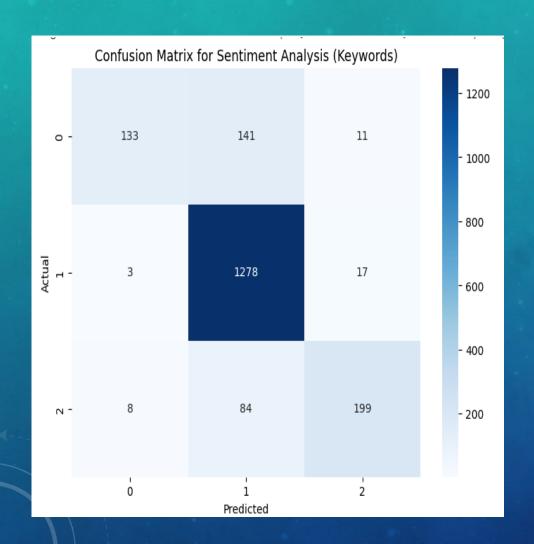
Sentiment Ana	lysis for Gen	res:		
	precision	recall	f1-score	support
negative	1.00	0.97	0.98	29
neutral	1.00	0.99	0.99	1444
positive	0.97	0.99	0.98	401
accuracy			0.99	1874
macro avg	0.99	0.98	0.98	1874
weighted avg	0.99	0.99	0.99	1874

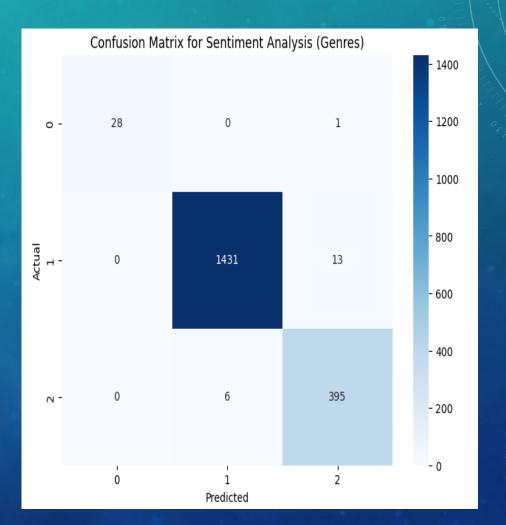
## SENTIMENT ANALYSIS (CONTD...)

#### Sentiment Analysis Models Used Summary

	Keywords	Genres	Overview
Model	MultinomialNB	MultinomialNB	MultinomialNB
Description	Naive Bayes classifier	Naive Bayes classifier	Naive Bayes classifier
Accuracy	0.86	0.99	0.66
Precision	0.92, 0.85, 0.88	1.00, 1.00, 0.97	0.74, 0.00, 0.64
Recall	0.47, 0.98, 0.68	0.97, 0.99, 0.99	0.44, 0.00, 0.94
F1-Score	0.62, 0.91, 0.77	0.98, 0.99, 0.98	0.55, 0.00, 0.76

# CONFUSION MATRIX FOR KEYWORDS & GENRES







# LONG SHORT TERM MODEL(LSTM FOR SENTIMENT ANALYSIS)

- •The LSTM model's performance in sentiment analysis on the combined text data is relatively poor, as indicated by the accuracy of around 54.27%.
- •The consistent validation accuracy throughout the training suggests that the model's architecture or data representation might need further refinement to better capture sentiment nuances from the combined text.
- •Additional feature engineering, hyperparameter tuning, and exploring more sophisticated model architectures could potentially enhance the model's performance.



# COMPARISON BETWEEN DIFFERENT MODELS FOR SENTIMENT ANALYSIS

Model/Scenario	Multinomial NB	Multinomial NB	LSTM (Combined Keywords
	(Keywords)	(Genres)	& Genres)
Accuracy	0.86	0.99	0.54





# INFLUENCE FACTORS ON MOVIE RATINGS

The primary objective of this analysis is to uncover the influence of production companies, directors, genres, and cast members on movie vote ratings. By leveraging machine learning models, we aim to:

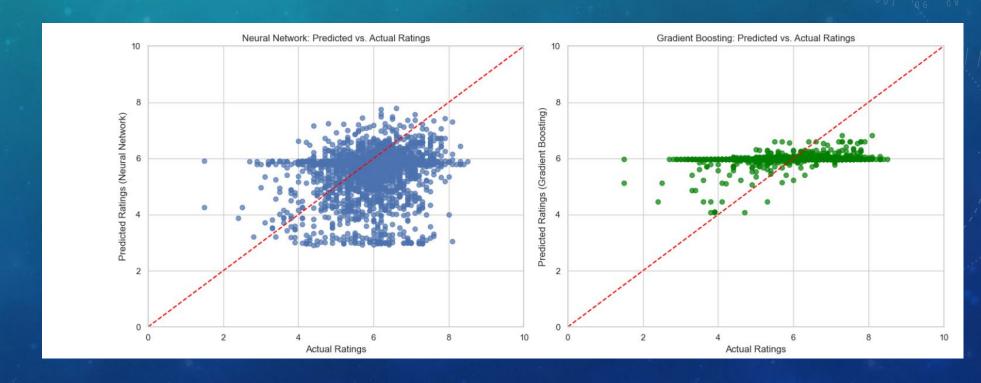
- ✓ Investigate how different factors contribute to movie ratings.
- ✓ Provide actionable insights for the film industry to enhance decision-making processes.
- ✓ Offer a comprehensive analysis of the relationships between key components and ratings.

### **MODEL SELECTION**

Models used for determining influence of various features on the IMDB Ratings:

- Gradient Boosting Regressor: Robust ensemble model
- Neural Network: Complex relationships captured by multi-layer architecture;

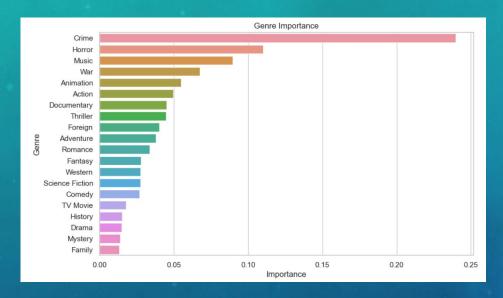
Feature Importance: Both models reveal the impact of production companies, genre, cast, and director collaborations on movie ratings.

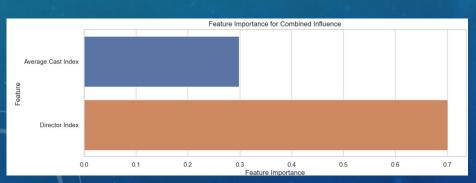


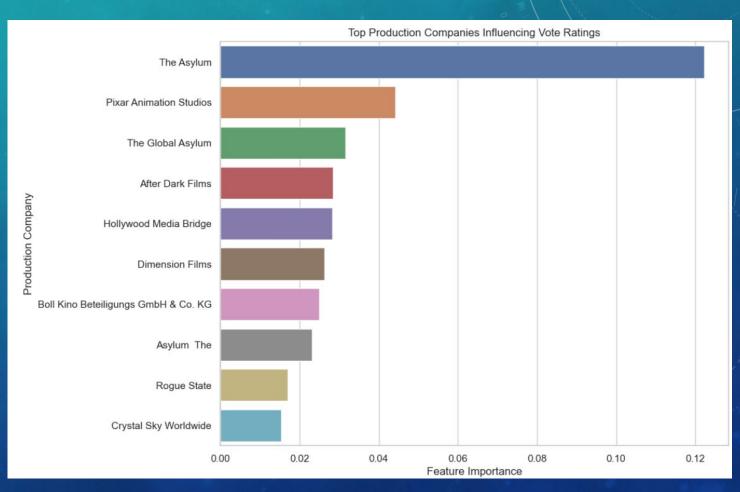
## PERFORMANCE AND INSIGHTS

MODELS ->	GRADIENT BOOSTING REGRESSOR	NEURAL NETWORK	
INFLUENCE FACTORS	MEAN SQU	ARED ERROR (MSE)	0 7
Production Companies	0.801031	1.526187	
Director And Cast Index	0.894547	0.873106	
Genre	0.744598	0.736392	

### FEATURE IMPORTANCE







#### **CHALLENGES**

#### **Elusive Optimal Accuracy**

The XGBoost regressor plateaued at a 71% accuracy, signalling complex feature correlations and multicollinearity. Extensive hyperparameter tuning and cross-validation was needed

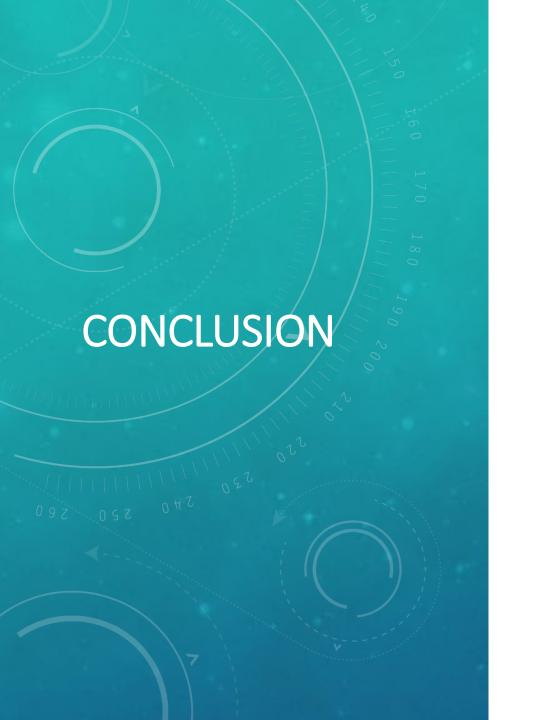
#### **Hurdles in Heteroscedasticity**

The Residual Plot unveiled heteroscedasticity's presence, posing a barrier to effective transformations. The quest for innovative features emerged to enhance model performance.

#### **Navigating Sentiment Analysis**

Sentiment Analysis faced intricacies in capturing positive sentiments accurately.

These challenges unlocked innovative avenues and highlighted the importance and need of additional relevant features for enhanced accuracy.





**Decoding Cinematic Magic**: Our analysis pierced the silver screen's veil, unraveling the enigma of movie influence factors—production companies, directors, genres, and more.



**From Data to Decision**: Armed with insights, the film industry can chart informed paths, leveraging historical data to guide budget decisions, create compelling narratives, and captivate audiences.



**Sentiments Whispered and Heard**: Our sentiment analysis decoded audience sentiments, giving filmmakers the power to comprehend and respond to viewer emotions.



**Challenges as Stepping-Stones**: Through challenges, we found stepping-stones toward innovation, inspiring us to forge ahead with inventive features and techniques.

#### PAVING THE WAY AHEAD...

#### **Explore- Sky Is The Limit:**

Our findings offer a compass for production companies and directors, guiding budget decisions through historical data's lens—votes, popularity scores, revenue, and profit intertwining.

We look forward to discovering more factors that will help us enhance our accuracy and gain better insights

#### **Decoding User Sentiment:**

Merging user reviews with our dataset unveils sentiment nuances, providing filmmakers with a magnifying glass to understand audience perspectives. (IMDB DATASET with 50k records)

#### **Beyond Horizons:**

The voyage doesn't culminate here. The journey forward encompasses a broader horizon, where we expand features and employ advanced techniques to conquer challenges and unveil deeper revelations.



## Thank You!