

The background of the slide features the word "NETFLIX" in a large, dark red, sans-serif font. The letters are slightly transparent, allowing a grid pattern and vertical light streaks to be visible through them. The grid is composed of thin, dark lines forming a square pattern. The vertical light streaks are in shades of red and white, running from top to bottom.

Netflix Movie Analysis & Rating Prediction

Managing Big Data Final Project - Team7

Wei (Will) Jiang, Xiang (Maggie) Meng, Paakhi Srivastava, Yulin Gai, Joe Ebby Karuthedath

Background

Netflix: The Beginning

- Launched its SVOD (Streaming Video on Demand) service in 2007
- Media Clearinghouses: gathering of B-rated movies
- Networks were initially reluctant to give up rights to popular content

Netflix: Now

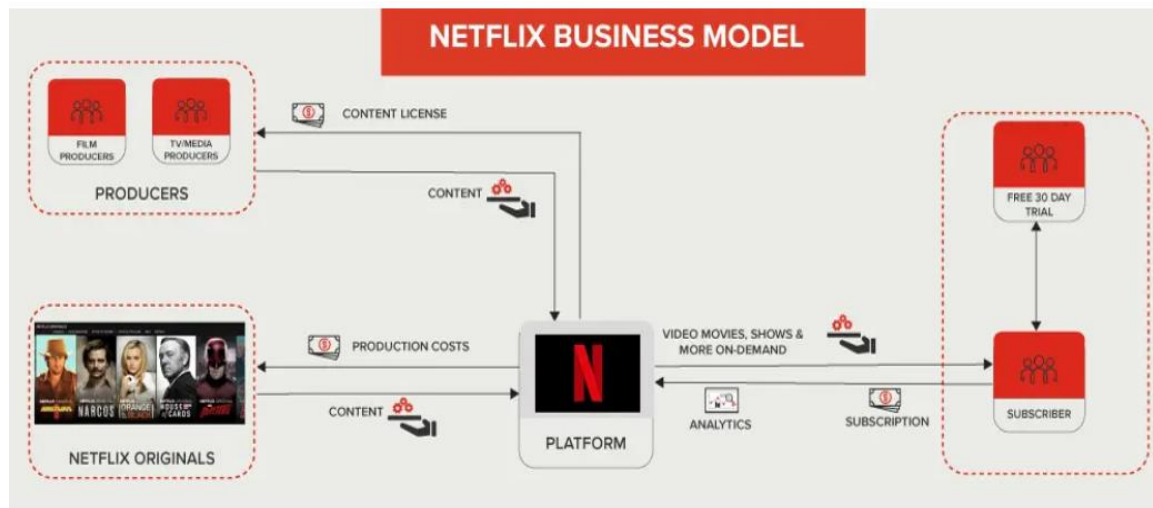
- Departure from cable, many people started to substitute cable with Netflix
- Netflix launched its android and ios app in 2010
- Rise was aided by technological advancements: Internet speeds, smart phones
- Market cap : \$291 billion

Netflix: The End?

- Challenge posed by new competitors: Amazon Prime video, Disney+
- Eroding market share

Business Understanding

- Major Costs: Acquiring content licenses and producing original content
- Optimizing this cost is essential for profits and maintaining market share
- The model attempts to predict, if a movie is **good or bad** before the movie is acquired by the platform.



Data Collection & Storage

kaggle

netflix_titles.csv

Detail Compact Column 12 of 12 columns

About this file

All TV Shows and Movies meta data on Netflix. Updated every month.

#	show_id	#	A type	#	A title	#	A director	#	A cast
Unique ID for every Movie / TV Show	Identifier - A Movie or TV Show	Title of the Movie / TV Show	Director of the Movie	Actors involved in the movie / show					
8807 unique values	Movie TV Show	70% 30%	8807 unique values	[null] Rajiv Chhola David Attenborough Other (17614)	35% 0% 70%	[null] 0% 0%	9% 0% 90%		
x1	Movie		Book Johnson Is Dead	Kirsten Johnson					
x2	TV Show		Good & Water						
x3	TV Show		Bandstand	Julien Leclercq					

Summary

- 1 file
- 12 columns

Feature films [edit]

Title	Genre	Premiere	Runtime	Language
Ghost Stories	Horror anthology	January 1, 2020	2 hours, 24 min.	Hindi
A Fall from Grace	Thriller	January 17, 2020	2 hours	English
Airplane Mode	Comedy	January 23, 2020	1 hour, 36 min.	Portuguese
Horse Girl	Drama	February 7, 2020	1 hour, 44 min.	English
To All the Boys: P.S. I Still Love You	Romantic comedy	February 12, 2020	1 hour, 42 min.	English
Isi & Cesi	Romantic comedy	February 14, 2020	1 hour, 53 min.	German
The Last Thing He Wanted	Political thriller	February 21, 2020	1 hour, 55 min.	English
Yeh Ballet	Drama	February 21, 2020	1 hour, 57 min.	Hindi
All The Bright Places	Romance	February 28, 2020	1 hour, 48 min.	English
Guilty	Thriller	March 6, 2020	1 hour, 59 min.	Hindi
Spenser Confidential	Action comedy	March 6, 2020	1 hour, 51 min.	English
Lost Girls	Crime drama	March 13, 2020	1 hour, 35 min.	English



WIKIPEDIA
The Free Encyclopedia

IMDb

	tconst	averageRating	numVotes	title	id	url	review	rating
0	tt1365050	7.7	78276.0	Beasts of No Nation	tt1365050	tt1365050	[I have seen several documentaries about the i...	6
1	tt1365050	7.7	78276.0	Beasts of No Nation	tt1365050	tt1365050	[Agu is a young boy in an unnamed African nati...	8
	tt1365050	7.7	78276.0	Beasts of No Nation	tt1365050	tt1365050	[Director Cary Fukunaga brings us the turmoil ...	5
	tt1365050	7.7	78276.0	Beasts of No Nation	tt1365050	tt1365050	[Agu (Abraham Attah) is a young lad growing up...	10
	tt1365050	7.7	78276.0	Beasts of No Nation	tt1365050	tt1365050	[You know when you hear other people whining a...	9

	tt15469820	6.4	5442.0	Britney vs Spears	tt15469820	tt15469820	[In some ways, the fascination with Britney Sp...	6
3876	tt15469820	6.4	5442.0	Britney vs Spears	tt15469820	tt15469820	[I gave up half way through..... says it all r...	3
3877	tt15469820	6.4	5442.0	Britney vs Spears	tt15469820	tt15469820	[That documentary was amazing! It revealed so ...	9
3878	tt15469820	6.4	5442.0	Britney vs Spears	tt15469820	tt15469820	[In this documentary, despite the fact that th...	8
3879	tt15469820	6.4	5442.0	Britney vs Spears	tt15469820	tt15469820	[Conservatorships and guardianships are among ...	8

Storage for future analysis



amazon
DynamoDB

NETFLIX

The MySQL logo, featuring the word "MySQL" in blue and orange text with a blue dolphin icon above the "SQL" part, and the Python logo, consisting of two interlocking snakes, one blue and one yellow.



is_US
0/1

↓

1 = general audience
2 = parent guidance needed
3 = adults only

audience_class
1/2/3

audience_class
1/2/3

↓ Threshold: 50% Percentile

is_popular
0/1

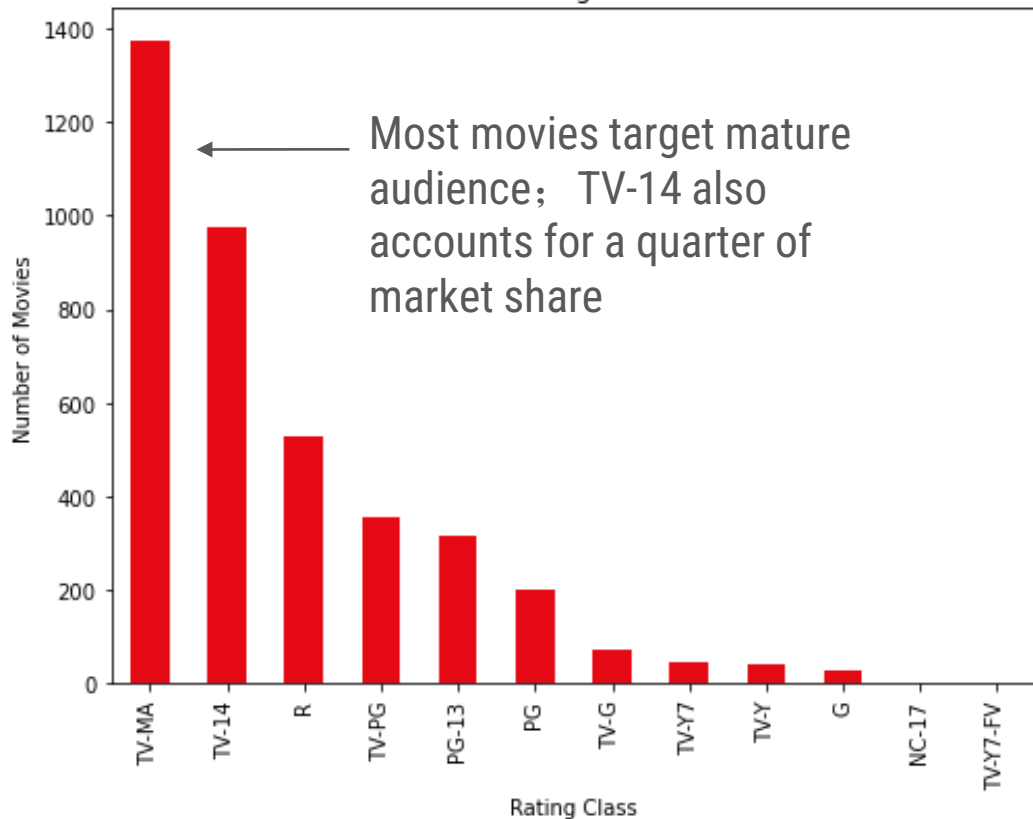
NETFLIX

is_popular
0/1

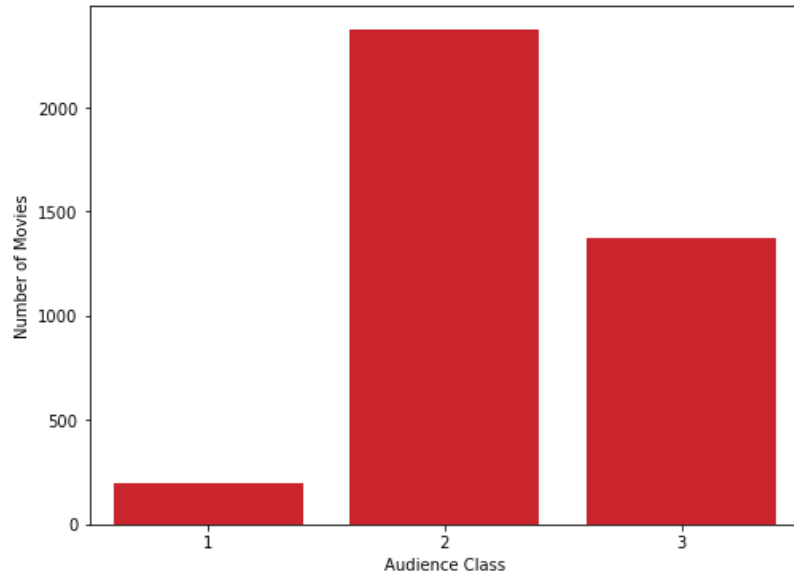
NETFLIX

Exploratory Data Analysis

Netflix Movies Rating Class Distribution

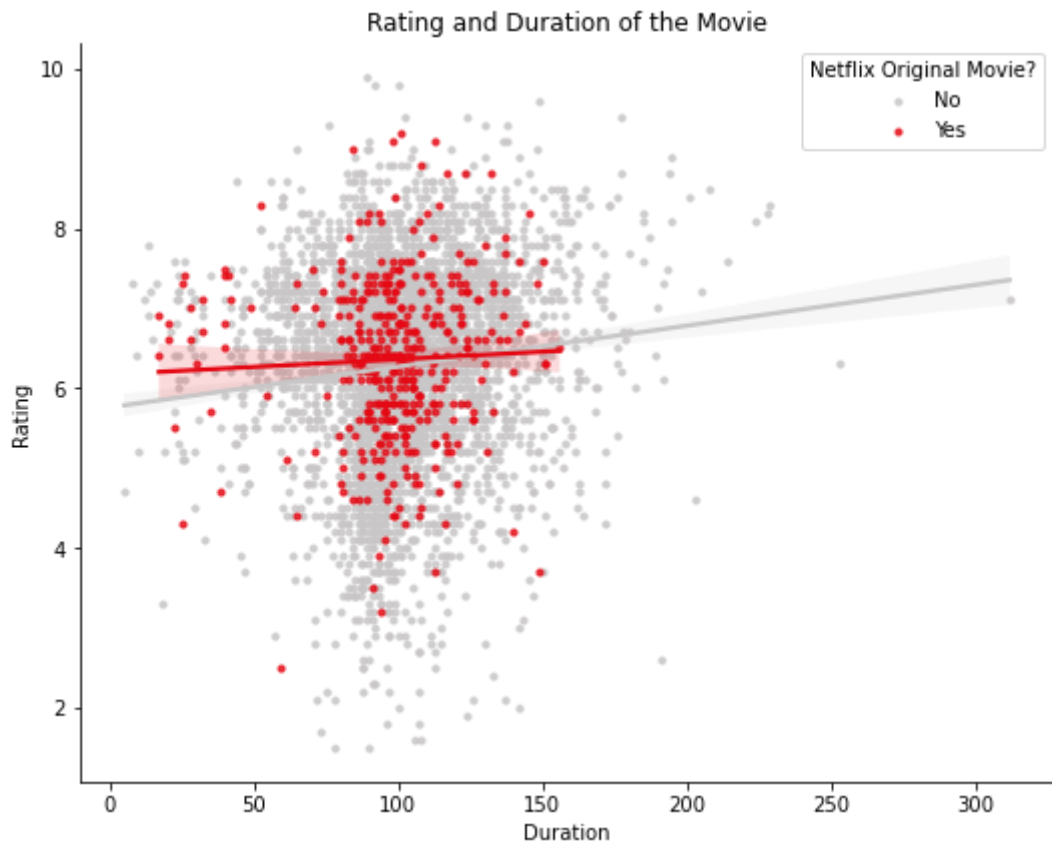


Netflix Movies Audience Class Distribution



Most movies are in the audience class 2 (need parental guidance), pointing out market opportunities in class 2 movies.

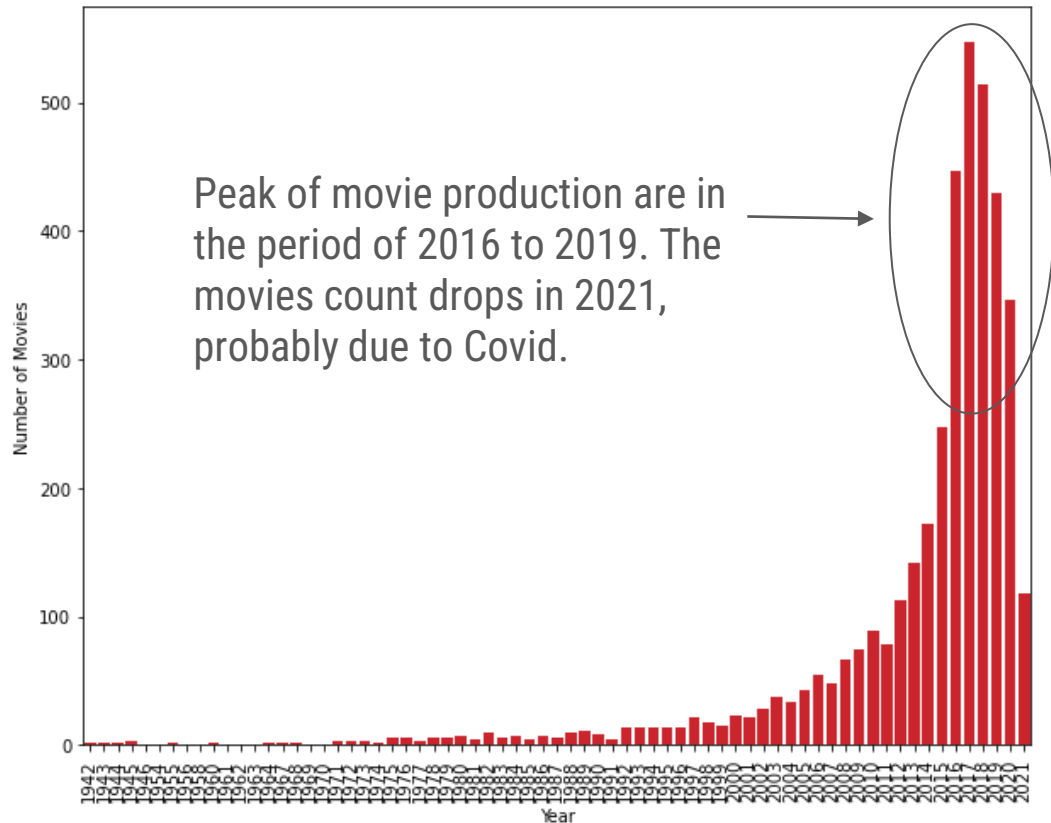
Exploratory Data Analysis



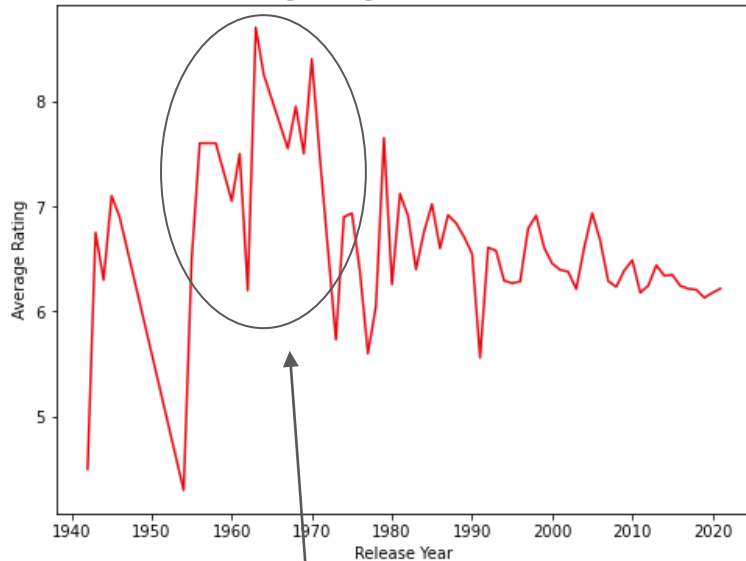
Comparing the regression lines that plot by scatter plots, we can see there is a stronger positive relationship between duration and movie ratings for non-Netflix-original movies, which indicates that we can include movie duration in our model.

Exploratory Data Analysis

Release Year of Movies in Netflix



Average Rating in Each Release Year



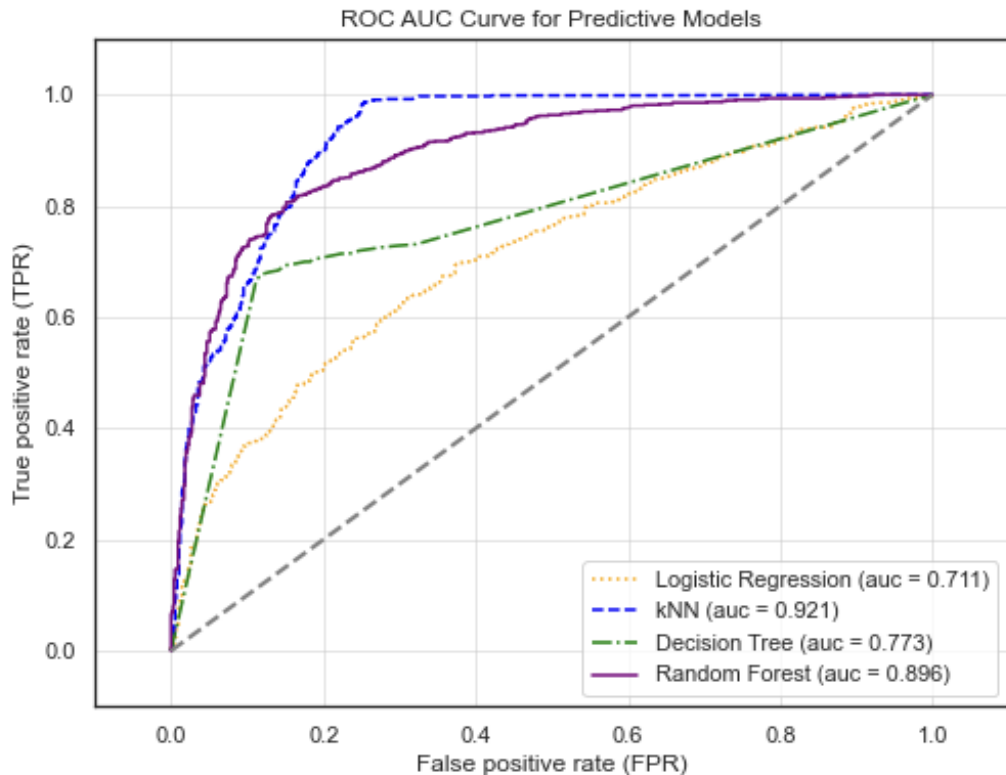
However, highest rating scores appear in 1960s, showing that large movie production does not lead to higher scores.

Modeling

- Split the data into test and training sets
- Ran a Logistic Regression, KNN, Decision Tree, and a Random Forest with parameter tuning
- Use the AUC Score to select the best model

	Logistic Regression	KNN	Decision Tree	Random Forest
Hyperparameter Grid Search	'lr_C': 10.0, 'lr_penalty': 'l2'	'knn_n_neighbors': 49, 'knn_weights': 'distance'	'criterion': 'gini', 'max_depth': 18, 'min_samples_leaf': 1, 'min_samples_split': 2	'max_depth': 16, 'n_estimators': 100
Non-nested CV accuracy	0.6567	0.7558	0.7504	0.7876
AUC Score	0.71	0.90	0.76	0.87

Model Evaluation



To evaluate our model by using test data, we created an ROC curve on the test data where we can see the KNN model surpasses the others with a AUC of 0.92.

Business Value

- Netflix Currently spends between \$100 and \$250 million on blockbuster movies
- The model will aid Netflix executives while making the decision of weather to go ahead producing an in house content.
- And while making a decision to acquire content licenses from a third party producer.

Future Improvements

- For **Opportunities**, it can be said a very friendly period for the development of the streaming platform during the pandemic
- For **Threats**, however, Netflix also has many competitors competing with them for loyal users, and there are also some further improvements in the script that can be brought about, like using NLP for assessing the scripts of movies to categorize them into different clusters, allowing executives to get a clearer picture

NETFLIX

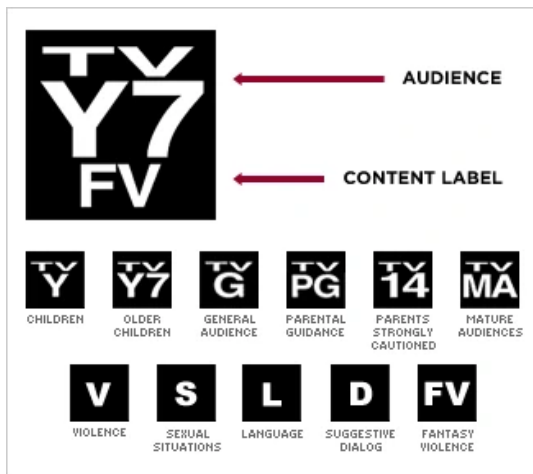
— & —

CHILL

**Thank You
&
Questions!**

U.S. Movie & TV Rating System

TV Rating System



Movie Rating System

