

# NLP Movie Rater

Using Neural Networks and Reddit Comments to Predict Movie Ratings

**Flatiron School // Data Science //**  
**Capstone Project // Zaid Shoorbajee**

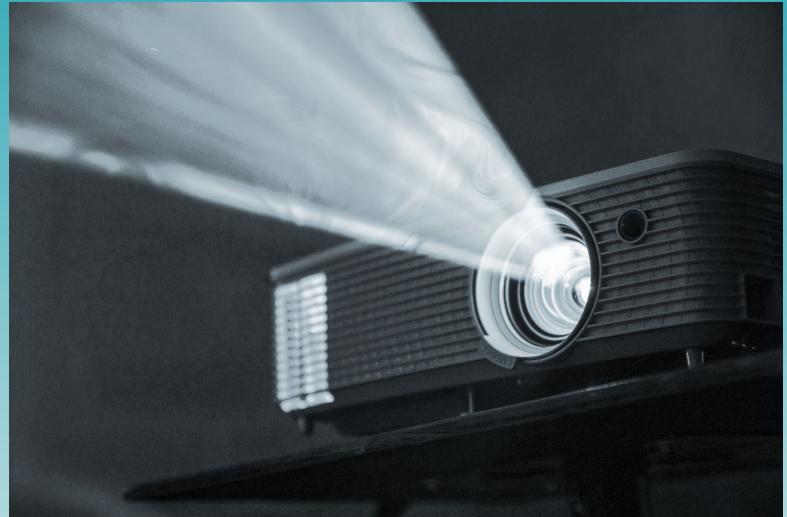
# Overview

- Business Understanding
- Data Understanding
  - Dataset
  - NLP
- Modeling
- Evaluation
- Recommendations
- Limitations & Future Work

# Business Understanding

A movie studio, *Cinedistance*, wants to identify areas of improvement for its movies before they are released.

- Releasing movies early to a focus group of 100 people.
- Prompt: What would you say about this movie in an internet comment section?
- Data science task:
  - Predict the movie's IMDb score based on the focus group's feedback.
  - Draw insights from viewers complaints.



# Business Understanding

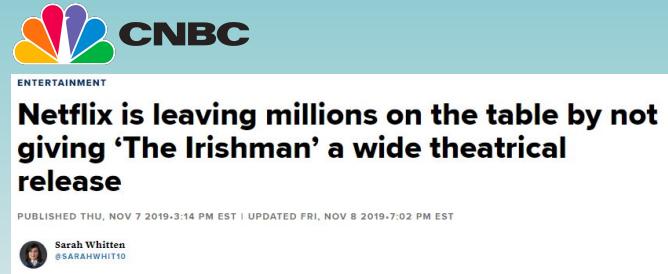
## Why use text to predict movie ratings?

This can inform decisions around marketing, distribution, re-shoots, edits, etc.

If score is low, qualitative feedback is available.



The image shows a news article from NPR's MOVIES section. The title reads "Warner Bros. kills off 'Batgirl' movie, \$90 million in". Below the title, it says "Updated August 3, 2022 - 4:49 PM ET" and "Heard on All Things Considered". The author is listed as ANASTASIA TSIOLCAS, with social media links for Twitter, Facebook, Instagram, and LinkedIn. The background of the slide features a teal gradient.



The image shows a news article from CNBC's ENTERTAINMENT section. The title reads "Netflix is leaving millions on the table by not giving 'The Irishman' a wide theatrical release". Below the title, it says "PUBLISHED THU, NOV 7 2019-3:14 PM EST | UPDATED FRI, NOV 8 2019-7:02 PM EST". The author is listed as Sarah Whitten, with a small profile picture. The background of the slide features a teal gradient.

# Data Understanding: Dataset

Reddit's [r/movies](#) community

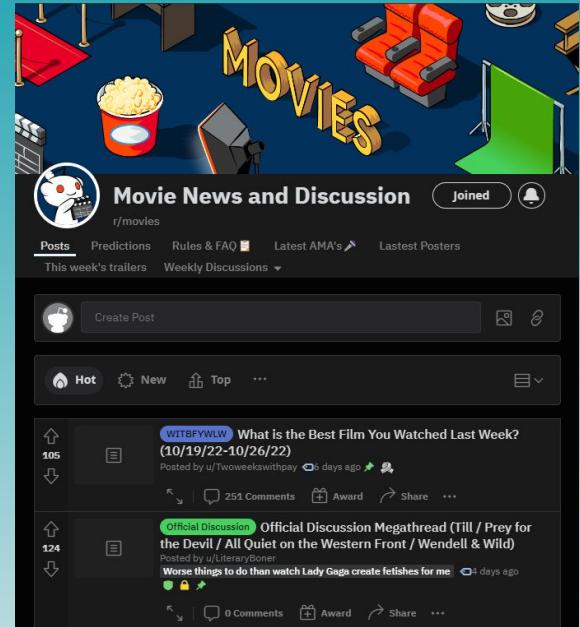
- ~30 million members
- Official discussions of most major releases
- Scraped top 100 comments from as many movie discussions as possible

Features:

- 70,693 text comments from 922 movies
- Years: 2016-2022

Target:

- IMDb movie ratings (1-10)
- Downloaded from IMDb website
- Mean rating: 6.47

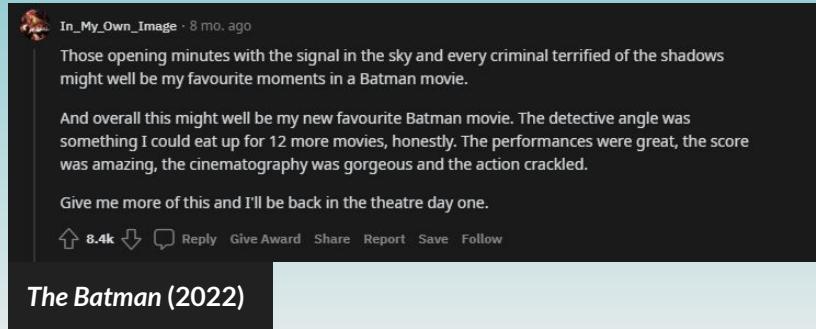
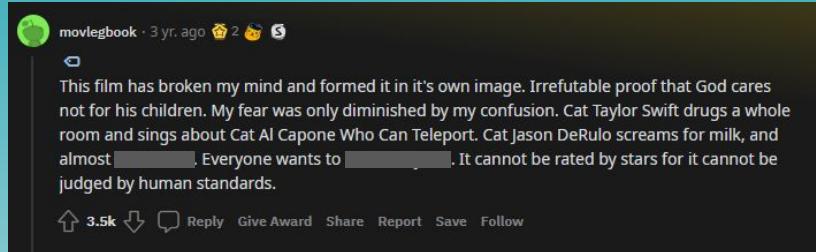


# Data Understanding: NLP

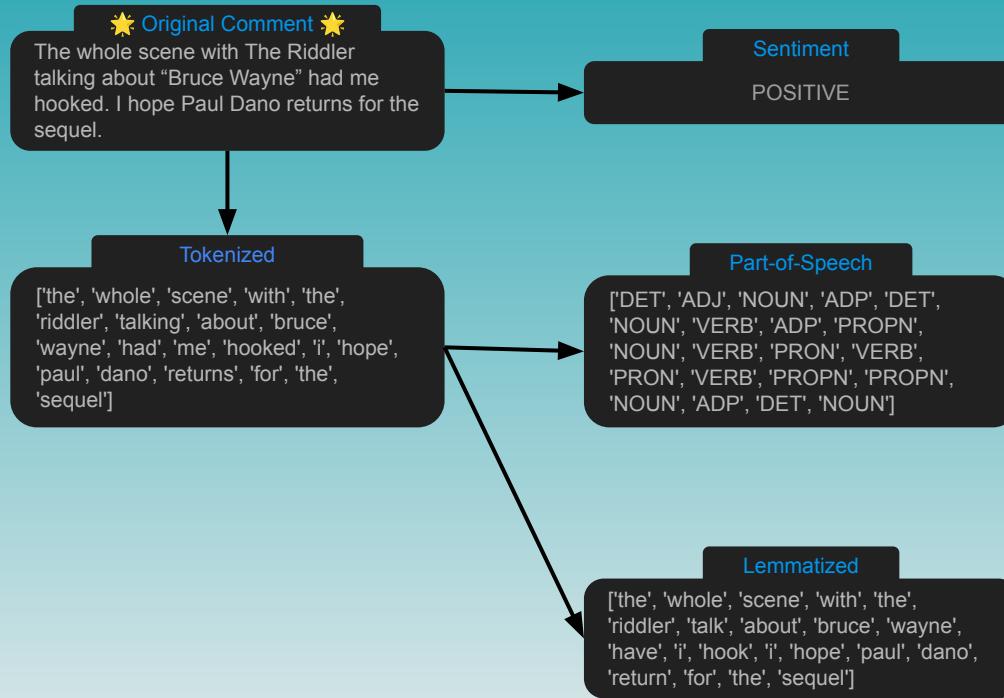
## Natural Language Processing

- Language is messy
- Easy for humans to understand; not for machines
- Unstructured data
- Must be heavily processed

Example r/movies comments



# Data Understanding: NLP



TF-IDF Vectorized (all comments)	
The Batman (TF-IDF)	
city	0.197992
dark	0.180793
end	0.150690
feel	0.279118
know	0.122270
love	0.216030
people	0.123268
scene	0.275356
think	0.231655
way	0.124107

Plus, values of "0" for all other words in the corpus

# Data Understanding: NLP

## Original Comment

The whole scene with The Riddler talking about “Bruce Wayne” had me hooked. I hope Paul Dano returns for the sequel.

## Example NLP Pipeline

- Cleaning; standardizing text
- Tokenization
- Lemmatization
- Removing stop words
- Vectorization

### TF-IDF Vectorized (all comments)

The Batman (TF-IDF)	
city	0.197992
dark	0.180793
end	0.150690
feel	0.279118
know	0.122270
love	0.216030
people	0.123268
scene	0.275356
think	0.231655
way	0.124107

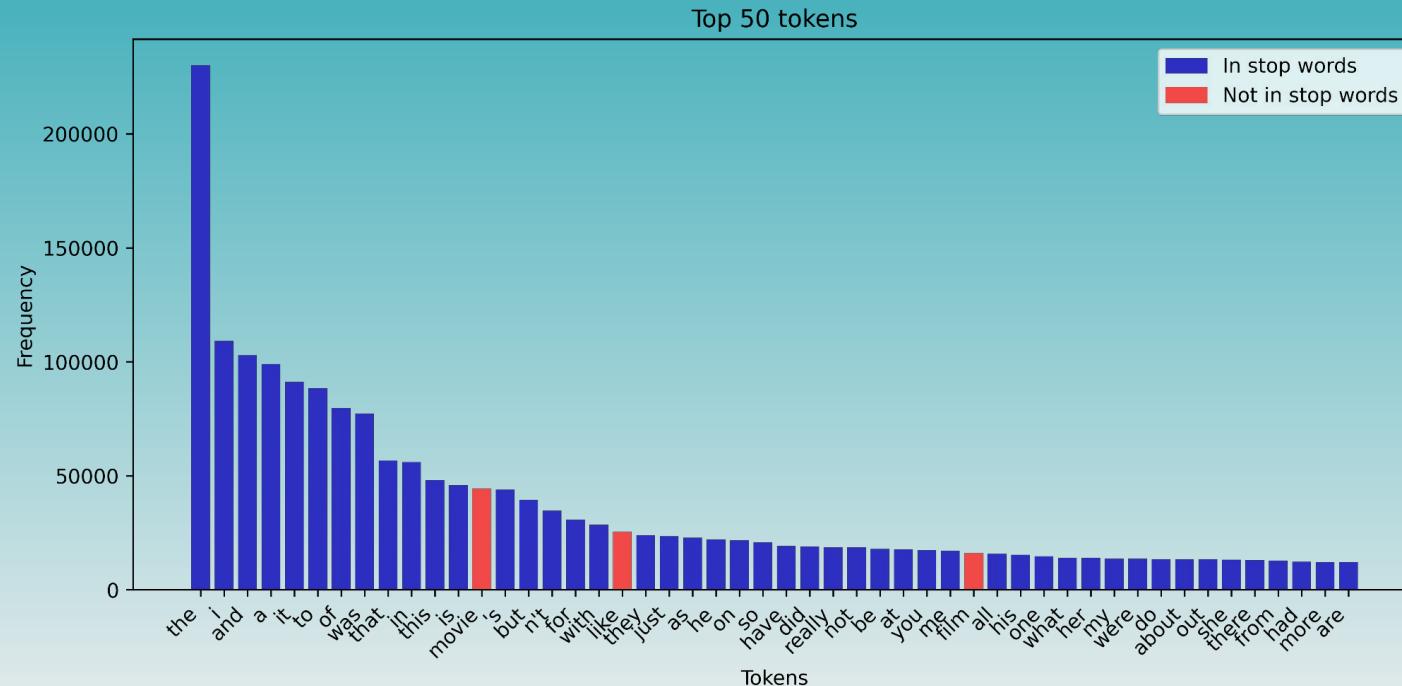
Plus, values of “0”  
for all other words in  
the corpus

## Other features (per comment section)

- Proportion of positive comments
- Proportion of negative comments
- Average character count of comments
- Mean length of words
- Proportion of unique words
- Proportion of stop words

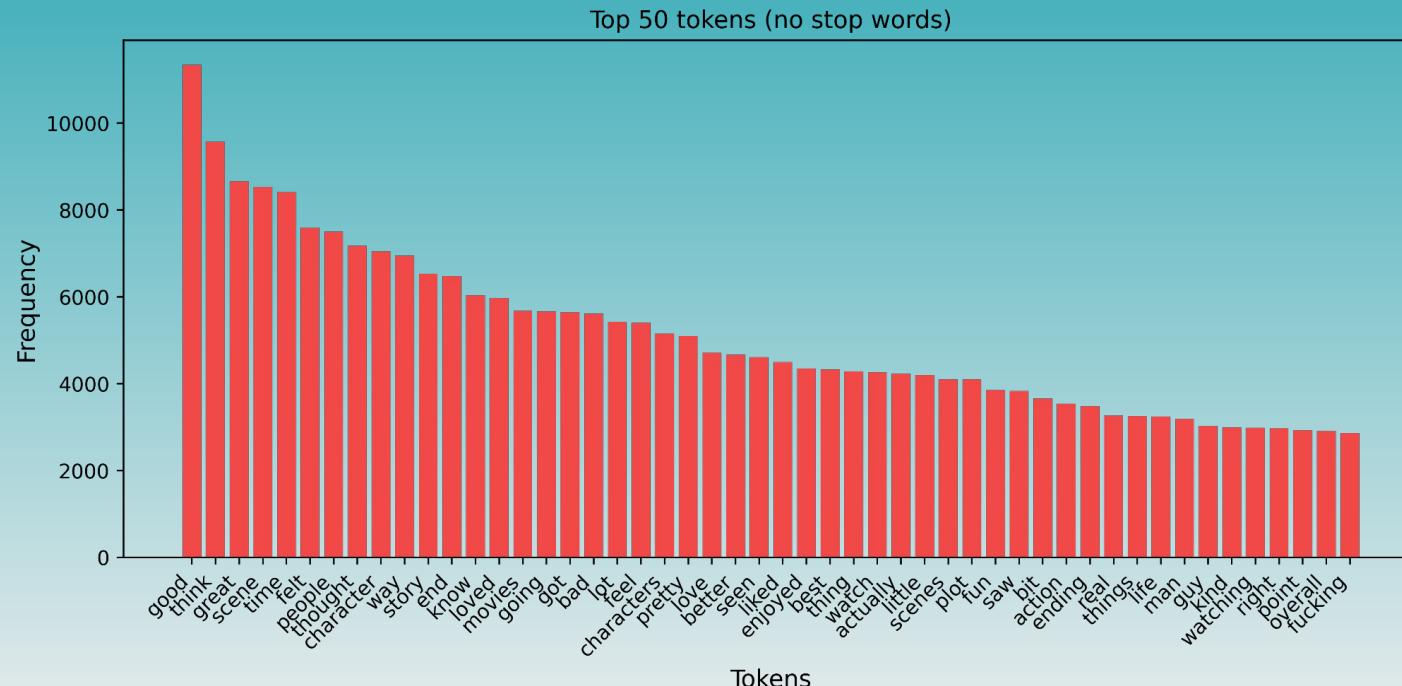
# Data Understanding: NLP

## Before removing stop words



# Data Understanding: NLP

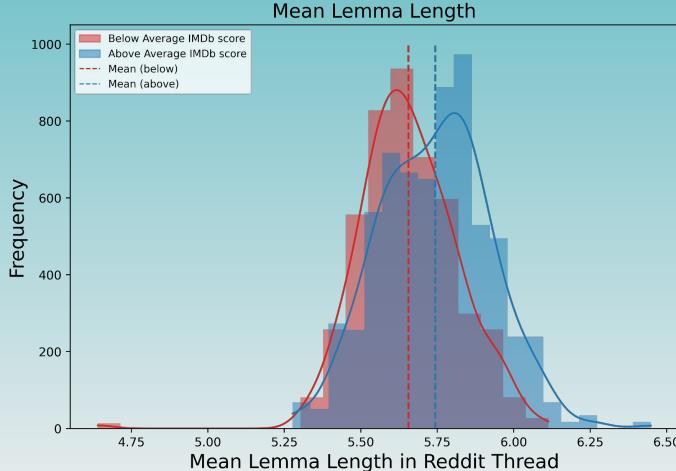
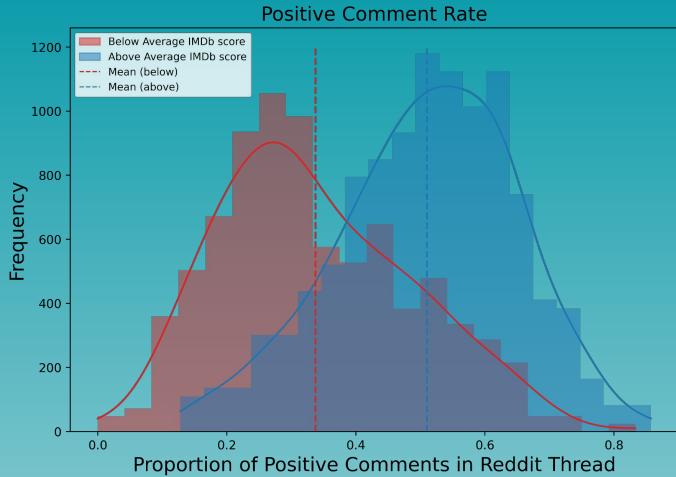
After removing stop words



# Data Understanding: NLP

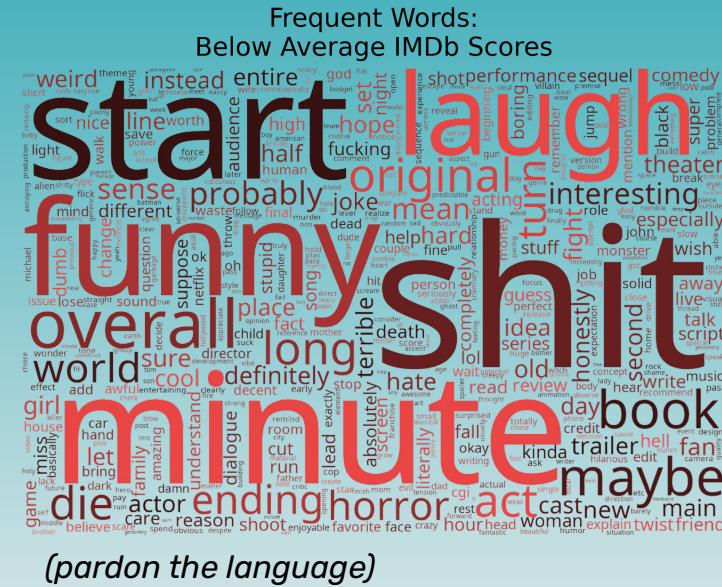
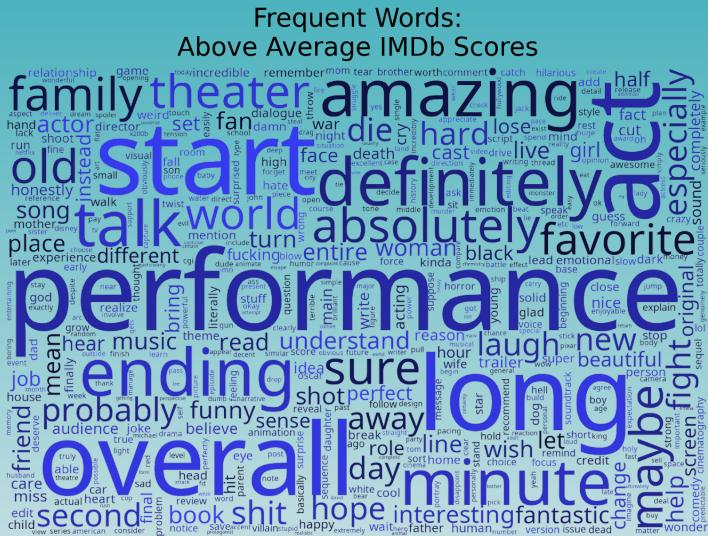
## Other features (per comment section)

- Proportion of positive comments
- Proportion of negative comments
- Average character count of comments
- Mean length of lemmas
- Proportion of unique lemmas
- Proportion of stop words



# Good vs. Bad movies: NLP Approach

Commenters use different vocabulary with movies they liked vs. disliked



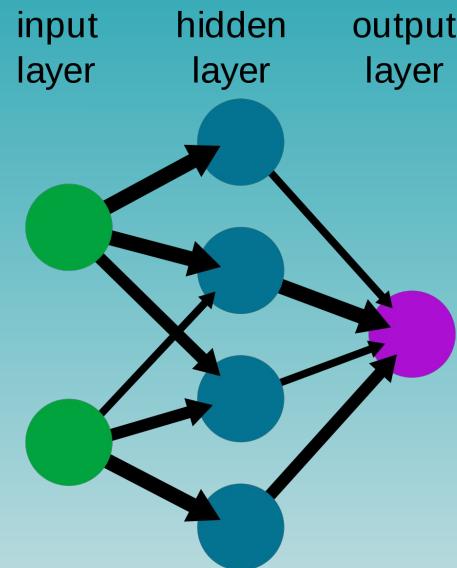
Q: Why is the word “minute” used so often in below average movies?



# Modeling

## Neural Networks:

- Machine learning practice mimicking the neurons of a brain. They can be tuned:
  - Inputs, layers, nodes, activation algorithms
  - Simple baseline model; 13 additional models
- Task: Regression
- Training, validation, and test datasets



# Evaluation

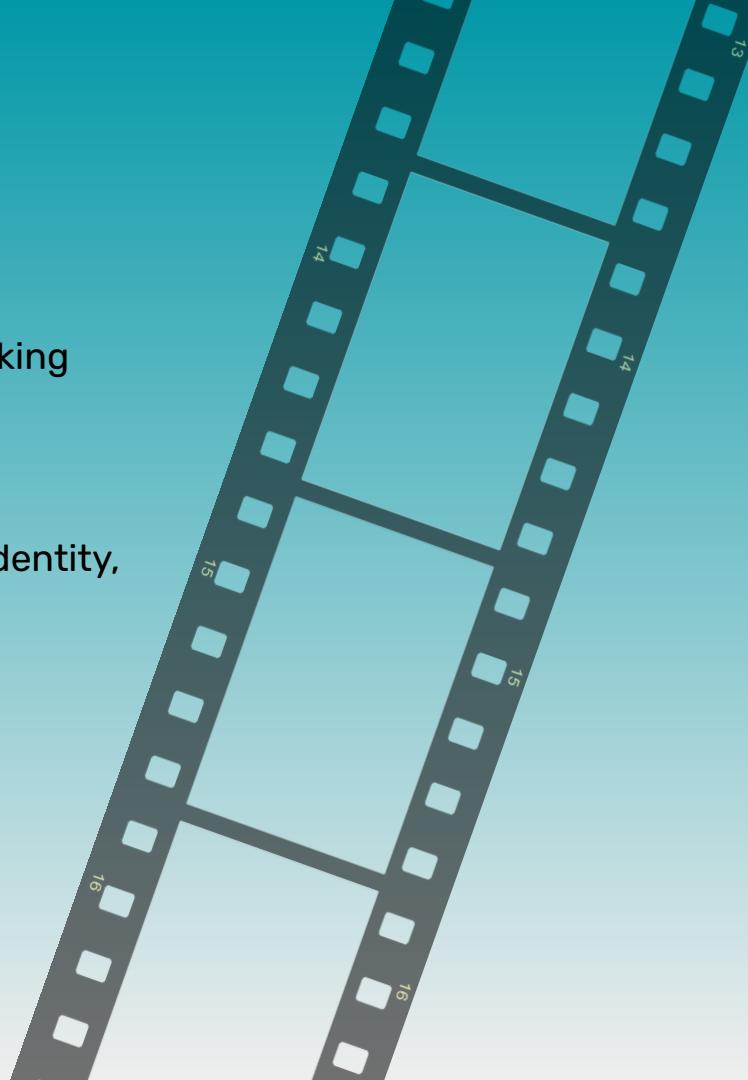
		Baseline	Final
<b>RMSE</b>	On average, how far are predicted ratings from the true value?	0.81	<b>0.57</b>
<b>R<sup>2</sup></b>	How much of the variability in IMDb ratings does the model explain?	29.3%	<b>66.2%</b>

## Example Predictions

Title	Pred. IMDb Score	True IMDb Score
Onward	7.587624	7.4
Birds Of Prey	6.496331	6.0
Hacksaw Ridge	7.453154	8.1
The Irishman	7.711154	7.8
Big Bug	6.382244	5.5
The Legend Of Tarzan	6.075438	6.2
The Nun	5.296925	5.3
365 Days	4.796644	3.3
Hillbilly Elegy	6.538070	6.7
Lion	7.606026	8.0

# Recommendations

- Use the model to predict IMDb scores
- Incorporate predicted score into decision-making
  - Bear in mind: Mean IMDb score is ~6.47
  - Prediction could be ~0.57 points off
- Split feedback by demographic (age, gender identity, etc.) to get more granular results
- If score is low, use qualitative feedback from comments
  - Scenes to reshoot, edit, or cut
  - Where to distribute the movie
  - How to market the movie



# Limitations

- Only had discussions still available on Reddit (922 movies)
- Bias in movie selection
  - Major releases

# Future work

- Experiment with different data
  - Letterboxd vs. Reddit
- Deployment: Python app & web app
- Apply model to YouTube comments on trailers



vs.



# Thank you



[zshoorbajee@gmail.com](mailto:zshoorbajee@gmail.com)



[@zshoorbajee](https://github.com/zshoorbajee)



[linkedin.com/in/zshoorbajee/](https://www.linkedin.com/in/zshoorbajee/)



[zaid.fyi](http://zaid.fyi)