

Movie Genre Classification



X



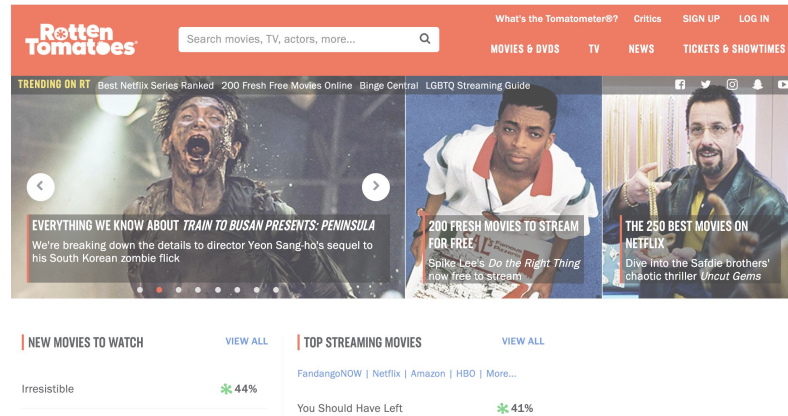
By: Pallavi Chhaparwal, Wong Yong Tian, Desmond Poo, Dominic Teo

Rotten Tomatoes®

Rotten Tomatoes and the Tomatometer score are the world's most trusted recommendation resources for quality entertainment. As the leading online aggregator of movie and TV show reviews from critics, they provide fans with a comprehensive guide to what's Fresh – and what's Rotten – in theaters and at home.

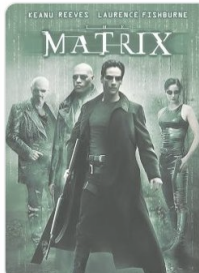
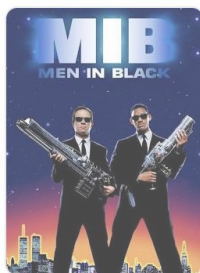
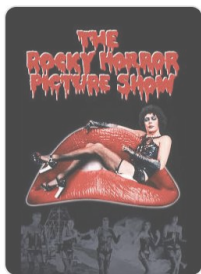
Lately, feedback from customers on not being able to find or match certain titles to respective genres led the Rotten tomato to discover some inaccuracy in existing movie genre classification process.

To improve on genre accuracy, the Management at Rotten Tomatoes has tasked an exploratory team to build a pilot model focusing on two popular genres: **Science Fiction** and **Horror** using Reddit as a data source.

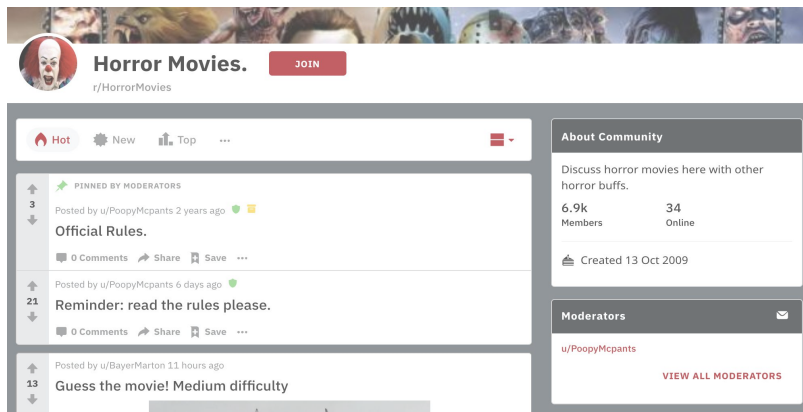


Problem Statement

How can we predict a Genre using keywords from subreddit posts?

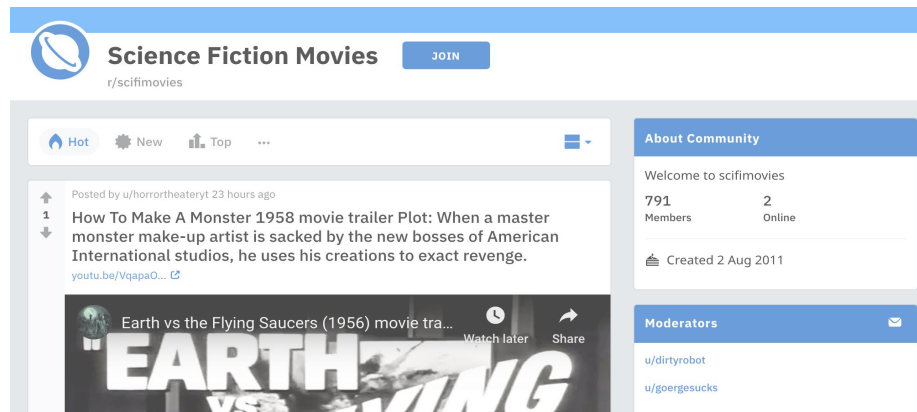


Subreddits



/r/HorrorMovies

- 892 post collected
- < 10 new post a day
- No advertisements
- Text and media mix
- 6.9k members
- Has Official Rules
- High member interaction



/r/scifimovies

- 345 post collected
- < 10 new post a day
- No advertisements
- Text and media mix
- 791 members
- No Official Rules
- Low member interaction

Subreddits



↑
1
↓
Posted by u/zmobiegirl 3 days ago
Help Finding An Elusive Zombie Movie!

Okay. I watched this movie over a decade ago on Netflix, back when streaming was really just starting to take off. It was about a guy traveling alone across a desert during a zombie apocalypse type scenario. If he encountered a zombie, he'd punch them to death.

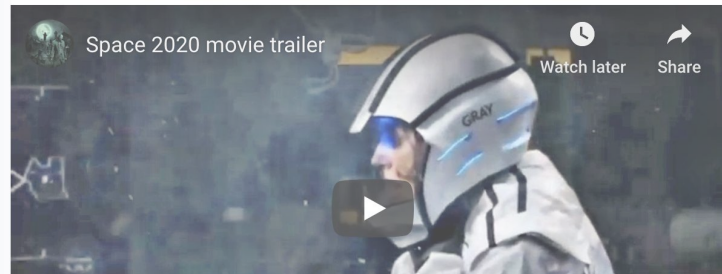
He found an isolated house with survivors in it - a man and woman who were in a romantic relationship and the woman's sister. The man had a "sexy" zombie secretly tied up in a shed out back, and I believe he engaged in adult type activities with said zombie. The woman's sister had a face on her abdomen that could talk... I think it was explained that it was a conjoined twin type situation.

I don't remember much more than that, but geez, do I want to find it again. Any ideas?
Help!

2 Comments Share Save ...



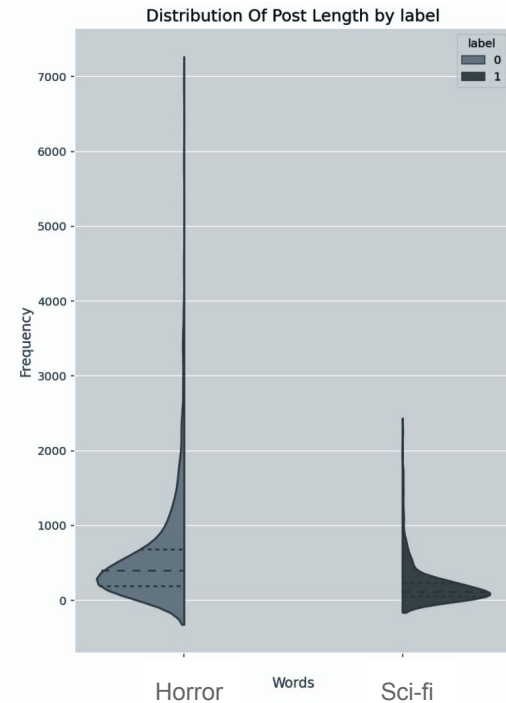
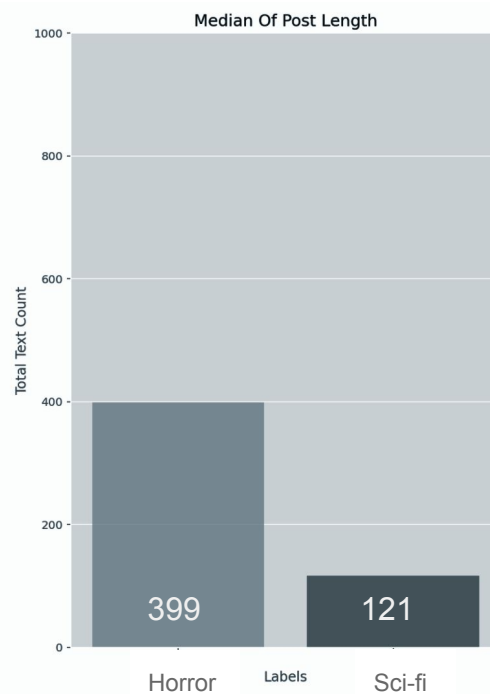
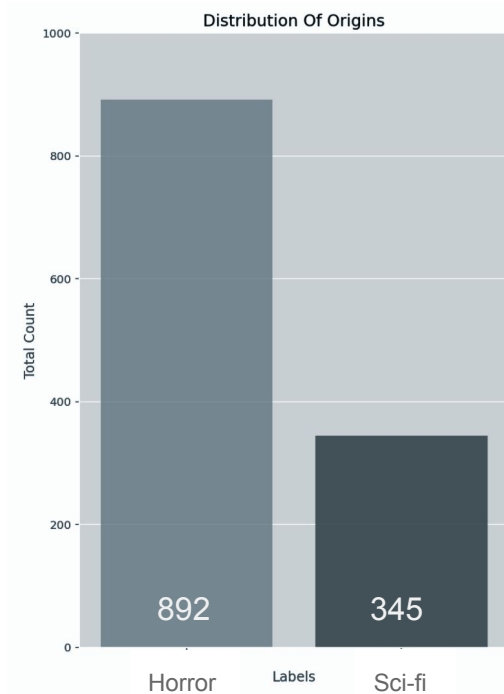
↑
1
↓
Posted by u/horrortheateryt 7 days ago
Space 2020 movie trailer Plot: In the year 2050, Dr. Ada Gray and her fellow astronauts aboard The Udo fight for survival after an accident leaves them stranded in deep space.
youtu.be/vlo4Dq...



Visually inspecting our subreddits:

- High members and high interactions can explain the number of post collected from Horror
- We also found that Sci-fi post are more condensed, plot related and have fewer noise

Subreddits



Our Data

There is a total of 1238 post.

893 belonging to Horror and 345 belonging to Sci-fi.

Issue: Imbalanced data



How we can overcome this:

- We can mediate this by picking the right models : Logistic Regression, Naïve Bayes Classifiers
- Scoring methods which are favourable towards imbalanced data :
Sensitivity-Specificity Metrics, Precision-Recall Metrics, ROC_AUC, F1
- Others : Stratified Cross Validation

Our Baseline

1	0.720032
0	0.279968

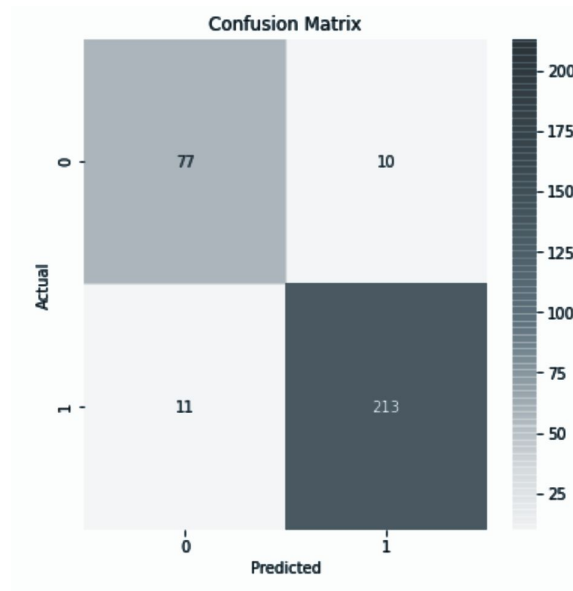
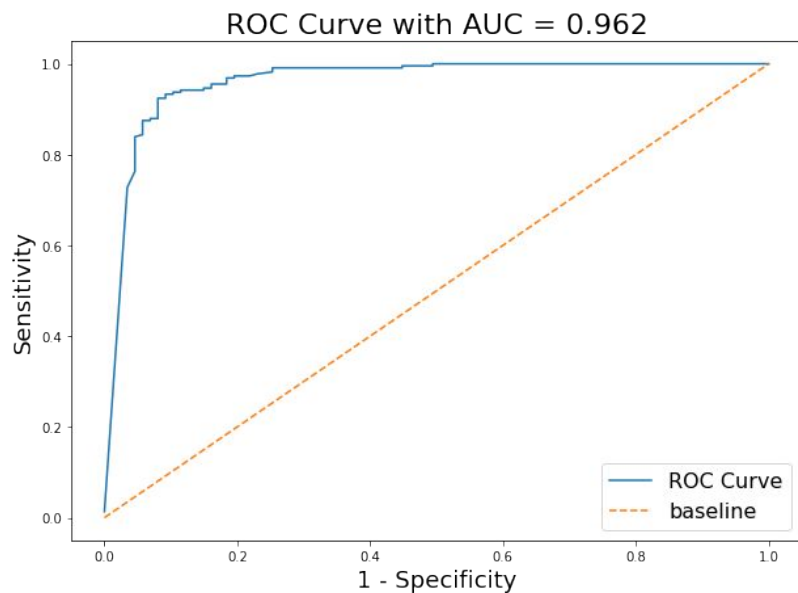
- Baseline model
 - 1 represent Horror movie
 - 0 represent Sci-fi movie
- Our baseline model accuracy is **72.0%** for predicting Horror movie posts.

Our Models

- Tested 4 models:
 - **Logistic Regression vs Multinomial Naive Bayes**
 - **Count Vectorizer vs TFIDF Vectorizer**
- Evaluated Sensitivity, Specificity, Precision and ROC AUC

Classifier	Logistic Regression with CountVectorizer	Logistic Regression with TFIDFVectorizer	Multinomial NB with CountVectorizer	Multinomial NB with TFIDFVectorizer
Sensitivity	93.8	98.2	99.1	99.1
Specificity	88.5	77.0	67.8	55.2
Precision	95.5	91.7	88.8	85.1
ROC AUC	96.2	97.4	97.3	95.5

Our Final Model



- **Logistic Regression model with Count Vectorizer** performed best overall
- Best Specificity and Precision scores
- High Sensitivity and ROC AUC

Our Words

- Best predictive words for Sci-fi movies:
 - **"Sci", "fi", "robot", "alien", "war", "ship", "earth", "scifi", "star", "arrival", "planet", "space" and "clone"**
- Best predictive words for Horror movies:
 - **"Horror", "scary", and "haunting"**

Word	Coefficient
fi	-1.609759
sci	-1.609759
trailer	-1.402131
robot	-0.942322
alien	-0.931292
war	-0.917861
plot	-0.905834
review	-0.874224
ship	-0.844599
earth	-0.834473
scifi	-0.809075
else	-0.771269
star	-0.732366
arrival	-0.721598
planet	-0.683273
short	-0.682440
space	-0.597077
clone	-0.568945
two	-0.553655
possibly	-0.544175

Word	Coefficient
day	0.622345
worth	0.626837
looking	0.630703
next	0.645085
little	0.653564
point	0.701007
best	0.725860
haunting	0.739609
opinion	0.790999
favorite	0.811128
guess	0.841198
thought	0.944435
scary	0.967569
good	0.987151
guy	0.991624
watch	1.007974
anyone	1.054510
watched	1.103585
seen	1.321873
horror	2.369752

Our Words

Horror Movies



SciFi Movies



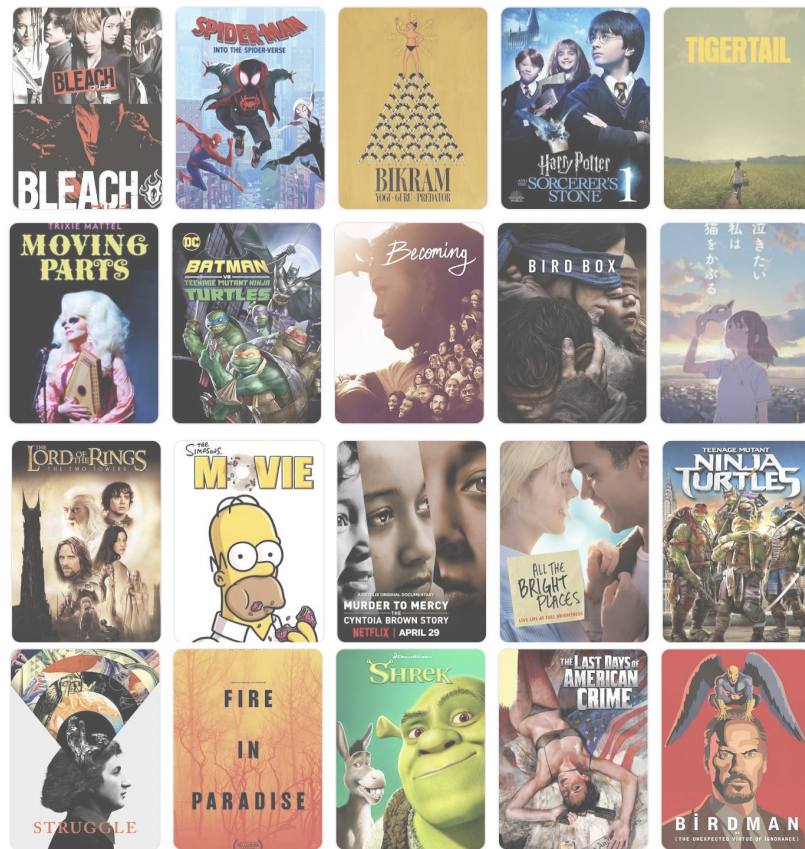
Limitations

- Insufficient data
 - small sample
 - unaccounted words
- Only classifies two genres



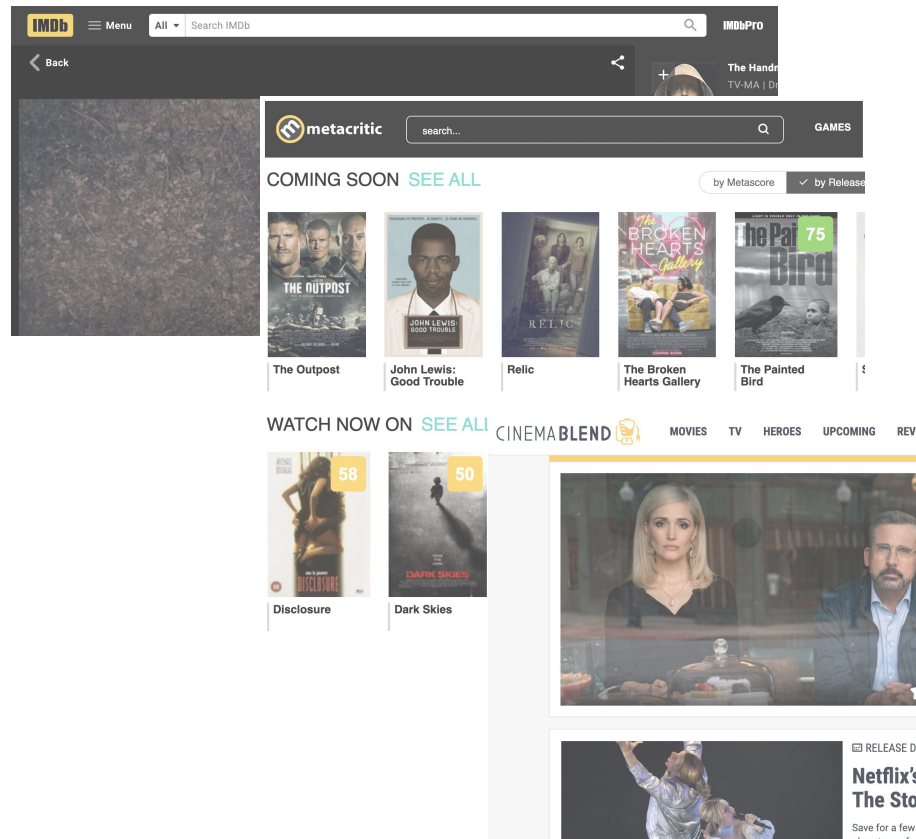
Recommendations & next step

1. Obtain more data: Sci-fi and Horror related, and others
2. Train a model with multiple target variables
3. Expand the data acquisition to include all existing genres listed in Rotten Tomatoes



Takeaways & Improvements

- Model accurately predicts subreddit source
 - Learned associative words
 - Horror subreddit is less descriptive than Sci-fi (dataset is key)
- Change source of data (reviews/websites)





Thank you! Questions?