

CLICK TO SEE SPOILER

An abstract graphic on the left side of the page, featuring a solid orange background. Overlaid on this are several geometric shapes in black and white. These include a large black circle in the upper right, a white triangle pointing downwards in the upper left, a black curved shape on the left, and a white triangle pointing to the right in the lower left. A thin black vertical line is also present in the lower right of the orange area.

# Detecting Movie Spoilers

---

BRIAN TRACY

# Background

---



Information that reveals important plot details or surprises



Reviews influence consumer attitude and decisions on seeing a movie

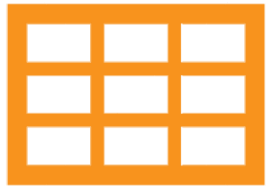


Users will avoid resources that present them with spoilers



# IMDb

## Business Problem



Current model relies on site moderation and user reporting



Automated flagging systems can free resources for other problems



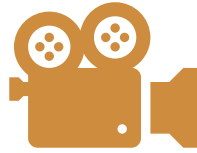
Goal: increase user confidence they will not be presented with spoilers

# Data Understanding

---

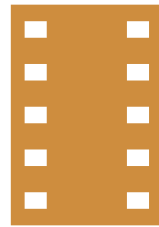


Data sourced from  
Kaggle



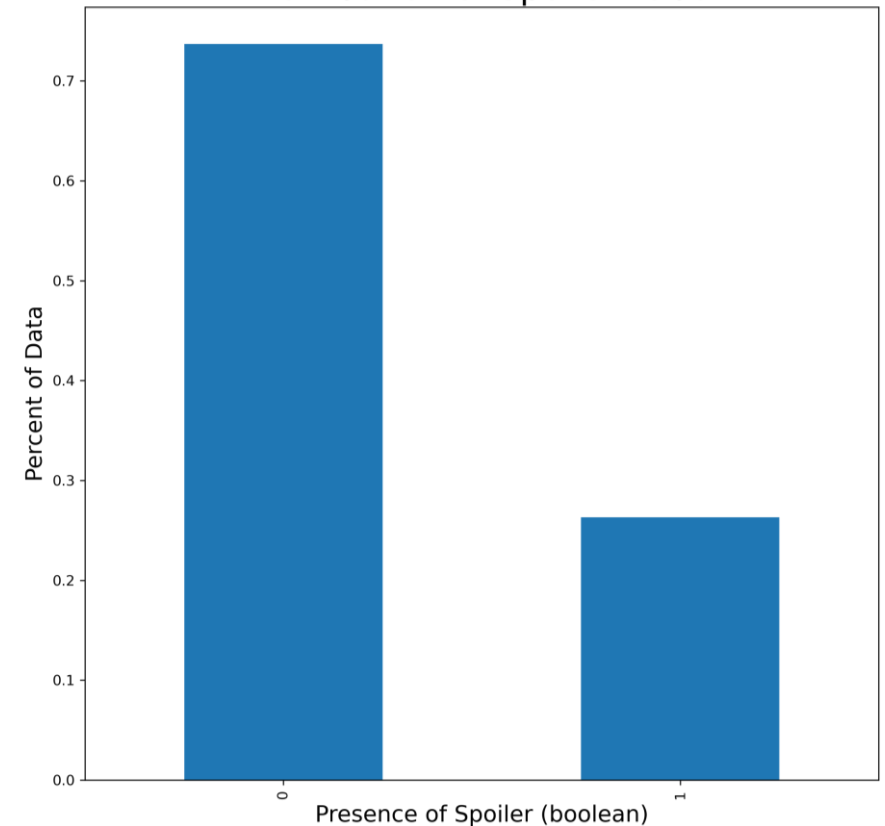
573,913 reviews  
between 1998 and  
2018

1572 unique films being  
reviewed

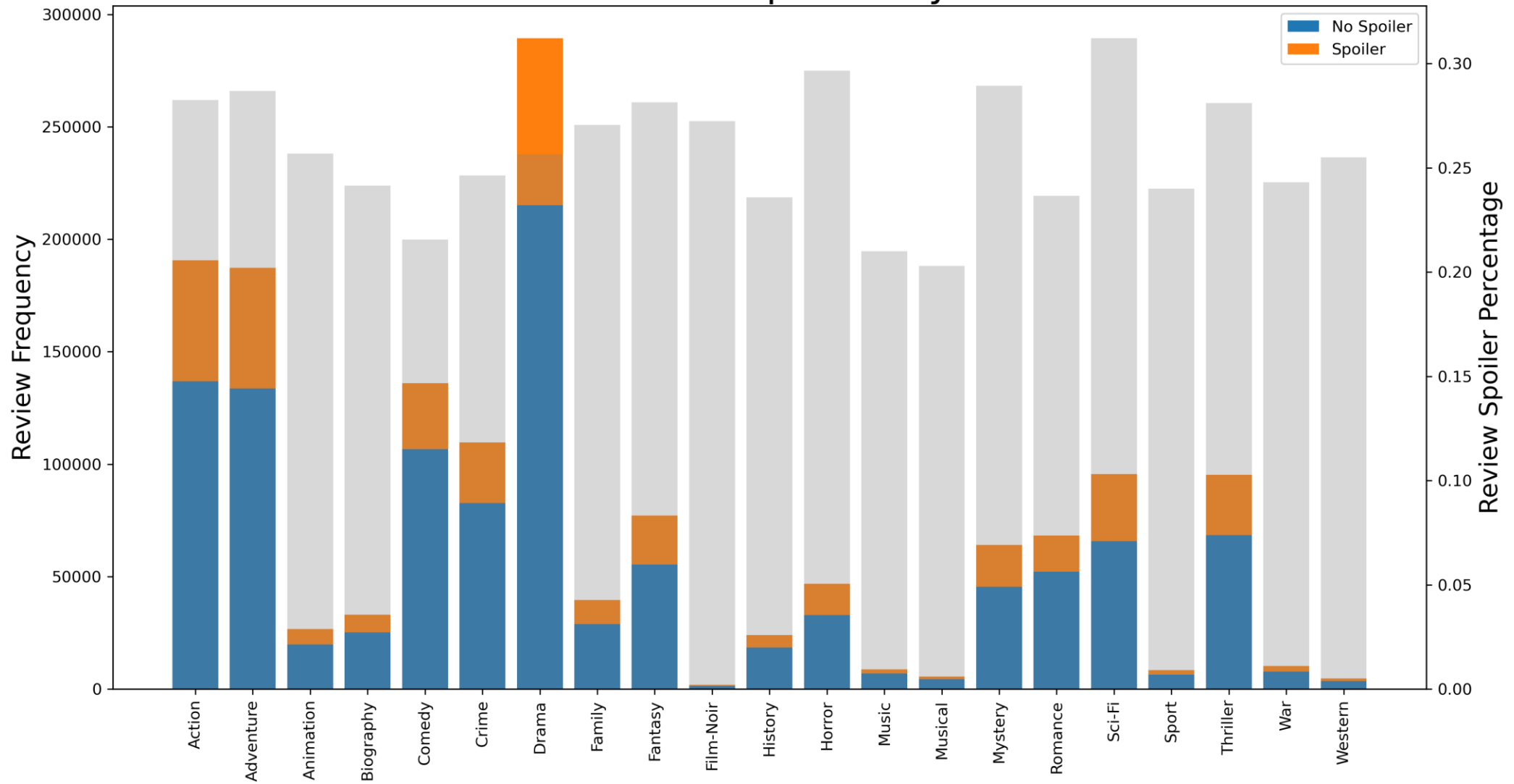


Additional film  
metadata and  
spoiler tagging

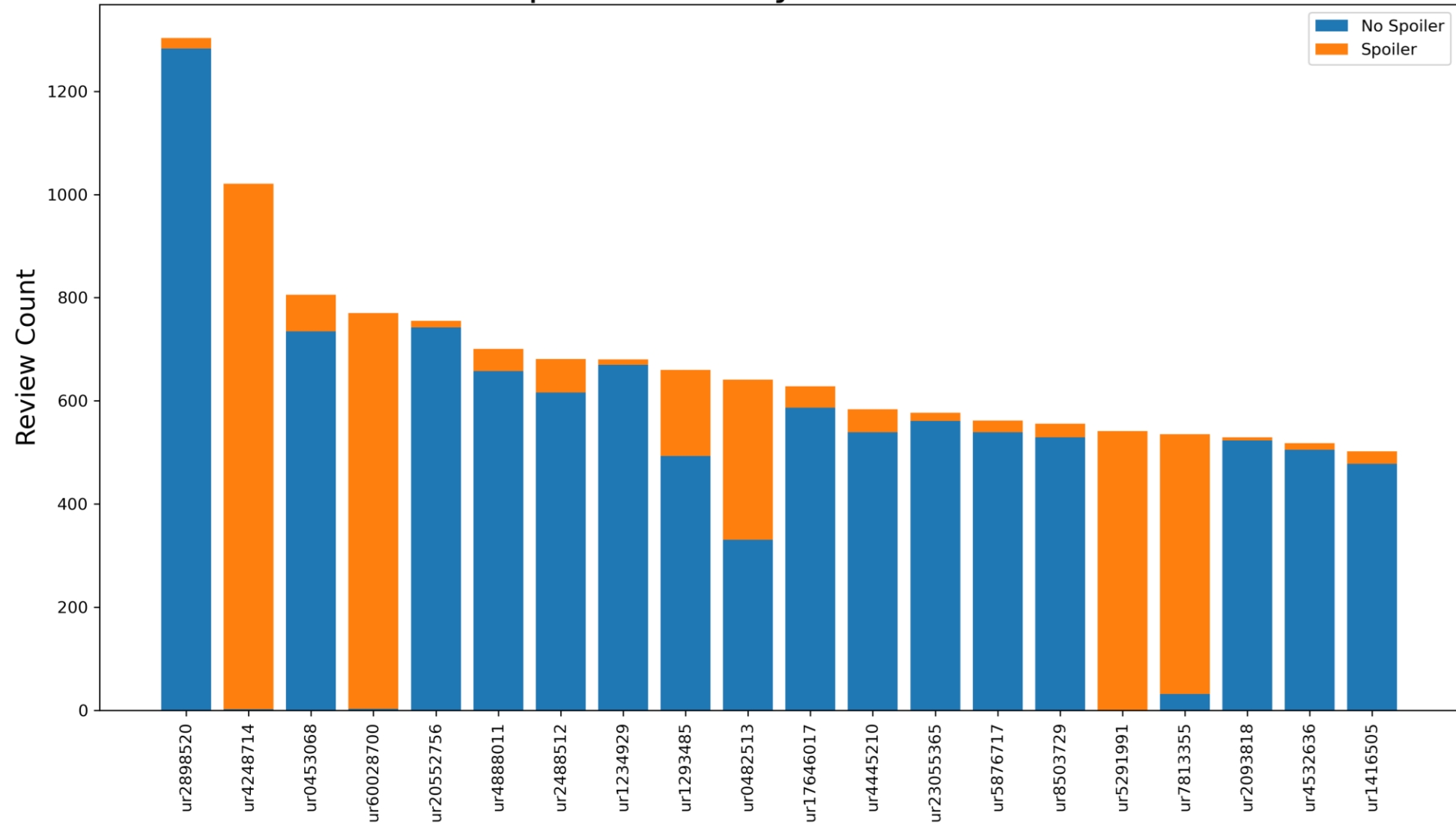
Distribution of Spoiler Label



## Distribution of Spoilers by Genre



# Top 20 Users by Total Reviews



[illegible]

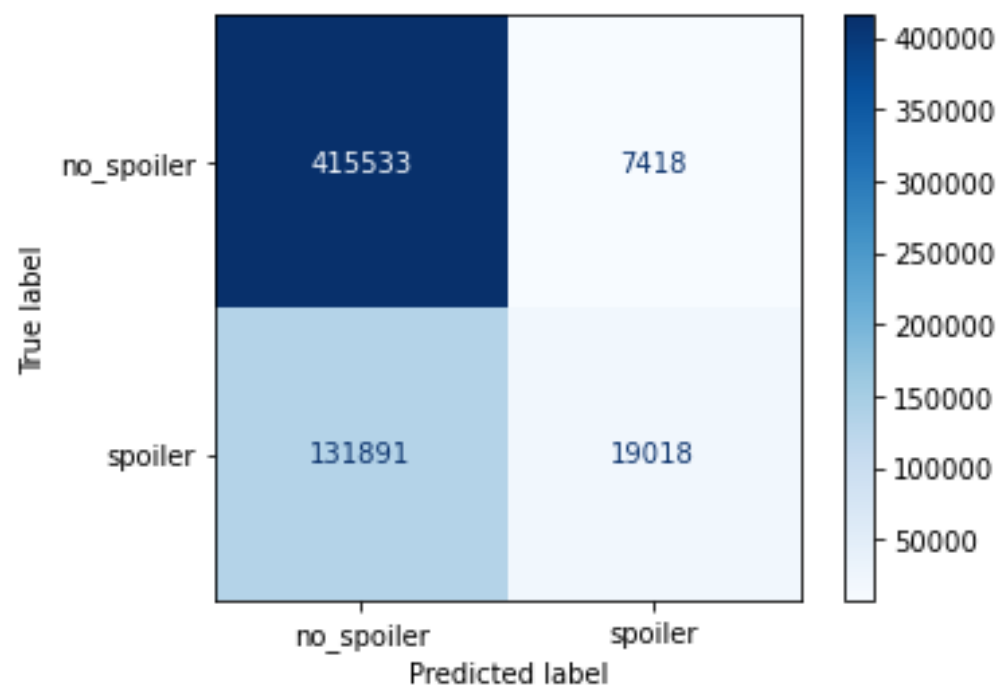


[illegible]

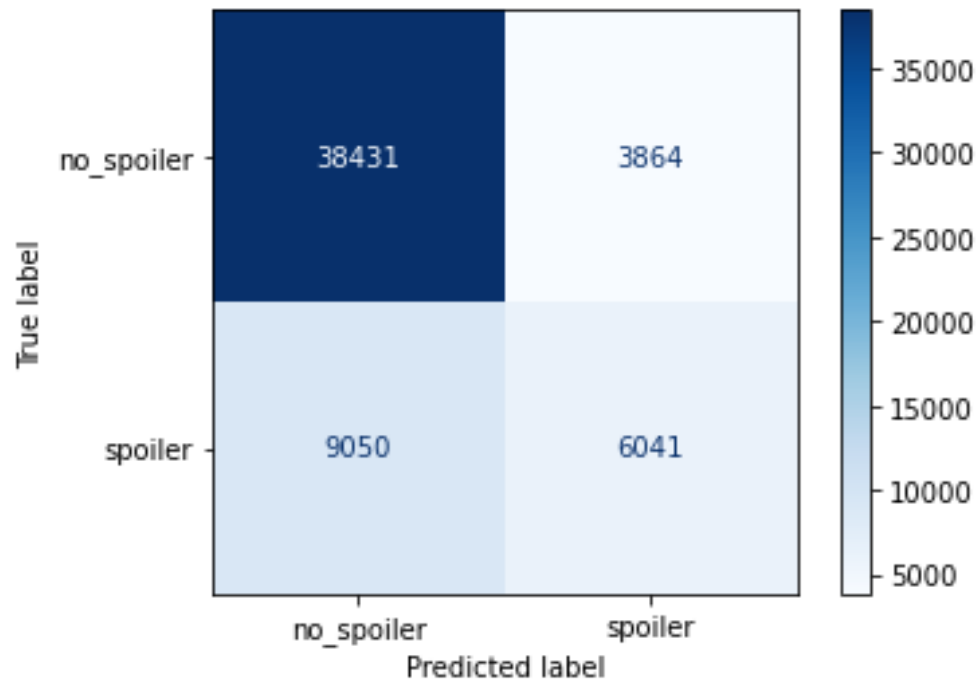
# Baseline modeling

- Premise: users will self-tag spoilers
  - Does not account for negation
- Also modeled non-language features
  - Results in all reviews predicted as no spoiler

	precision	recall	f1-score	support
no_spoiler	0.76	0.98	0.86	422951
spoiler	0.72	0.13	0.21	150909
accuracy			0.76	573860
macro avg	0.74	0.55	0.54	573860
weighted avg	0.75	0.76	0.69	573860



	precision	recall	f1-score	support
no_spoiler	0.81	0.91	0.86	42295
spoiler	0.61	0.40	0.48	15091
accuracy			0.77	57386
macro avg	0.71	0.65	0.67	57386
weighted avg	0.76	0.77	0.76	57386



# Results

- Iterative testing with different encoding tools and parameters
- Evaluating using f1-score and overall accuracy
- Utilized several different encoding methods for bag of words modeling

# Next Steps

---

- Scrape movie-centric subreddits
- Explore cosine similarity
- Use transfer learning



# Questions?

---



brtracy1984@gmail.com



<https://github.com/brtracy>