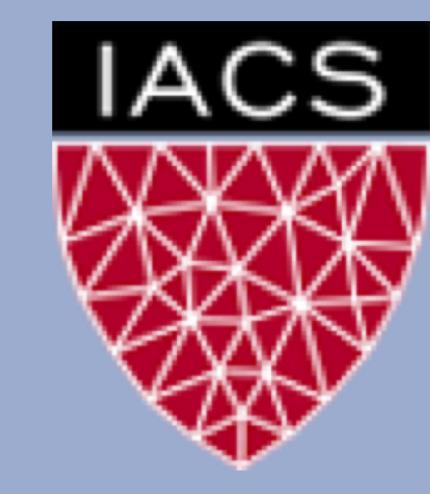




HARVARD

School of Engineering
and Applied Sciences

Movie Genre Multi-label Classification from Plot Descriptions

Andrew Lund, Nicholas Morgan, Amay Umradia, Charles Webb
CS109B - Data Science II - Spring 2018 Final ProjectINSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

Introduction

The overarching goal of this project is to **build a machine learning pipeline from scratch** in order to predict movie genres based on their plot descriptions, with the following goals:

1. Build a dataset of 1000 movies by scraping movie information from The Movie Database (TMDB) and the Internet Movie Database (IMDB).
2. Apply bag-of-words with term frequency inverse document frequency (TFIDF), word2vec (w2v), and doc2vec (d2v) transformations to movie plots.
3. Use conventional machine-learning techniques to classify movies as one or more genres using the transformed plots.

Project Question

“Can we build a robust movie plot and genre dataset through web scraping, then apply bag-of-words and word2vec transformations to plot descriptions to accurately classify movie genres using conventional machine learning multi-label classification techniques?”

Data Scraping & Cleaning

Data scraped using “tmdbsimple” and “imdb” python libraries.



Sources:

1. The Movie Database
2. The Internet Movie Database

Data Description & Design Considerations:

- 1000 movies
- English-language only
- From “top-rated” TMDB section
- Each classified as one or more genres (Figures 2 & 3)
- TMDB genres used as global response variable
- Only first IMDB plot scraped from list of plots
- Genres converted to binary representation (Figure 1)

Title: The Godfather
String Genres: ['Drama', 'Crime']
TMDB Genres: [18, 80]
Binarized Genres: [0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0]

Title: Schindler's List
String Genres: ['Drama', 'History', 'War']
TMDB Genres: [18, 36, 10752]
Binarized Genres: [0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0]

Title: Psycho
String Genres: ['Drama', 'Horror', 'Thriller']
TMDB Genres: [18, 27, 53]
Binarized Genres: [0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0]

Title: The Dark Knight
String Genres: ['Drama', 'Action', 'Crime', 'Thriller']
TMDB Genres: [18, 28, 80, 53]
Binarized Genres: [0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0]

Title: Sing Street
String Genres: ['Comedy', 'Romance', 'Drama', 'Music']
TMDB Genres: [35, 10749, 18, 10402]
Binarized Genres: [0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0]

Figure 1. Binary Genres

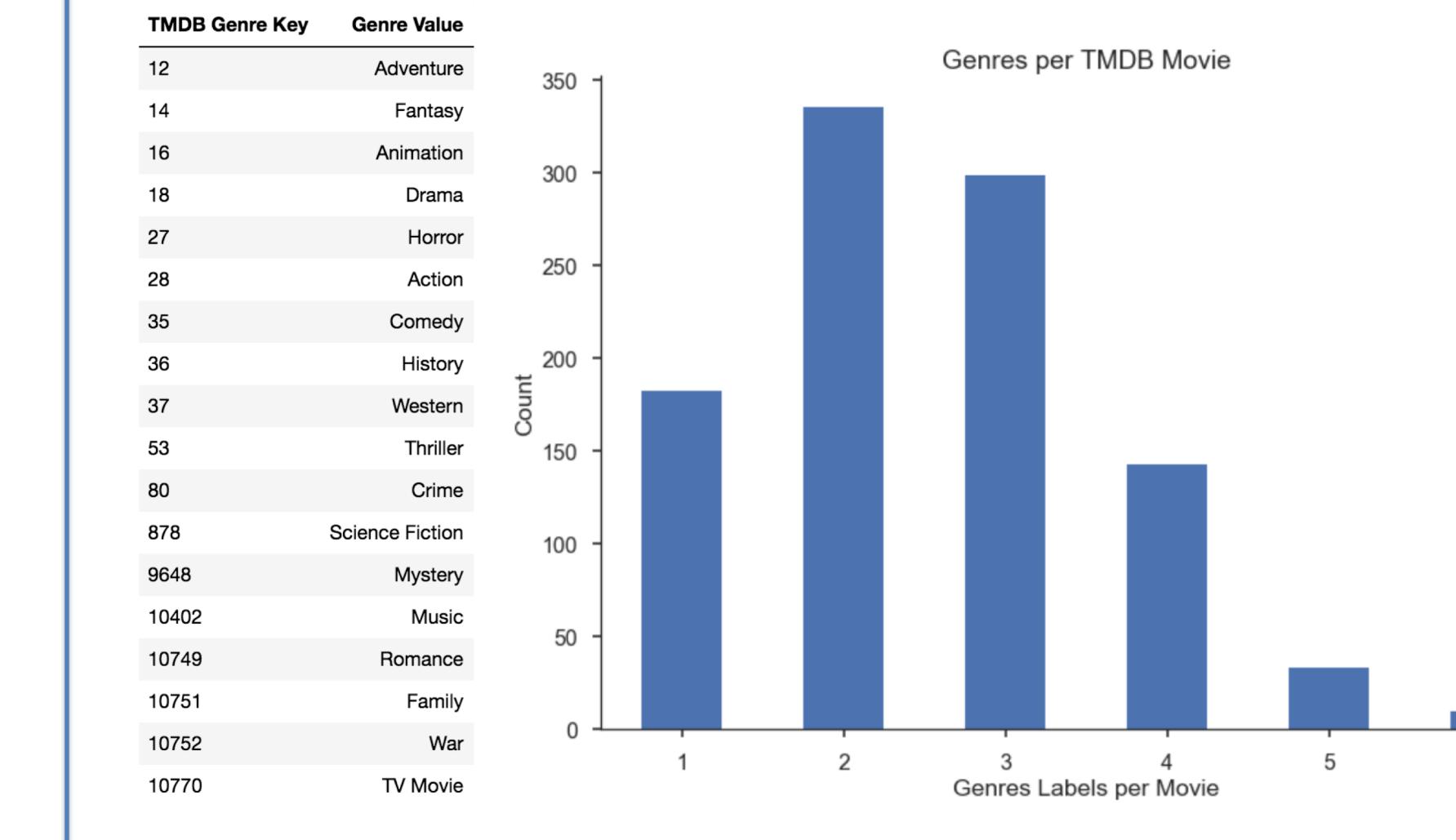


Figure 2. Genre Labels

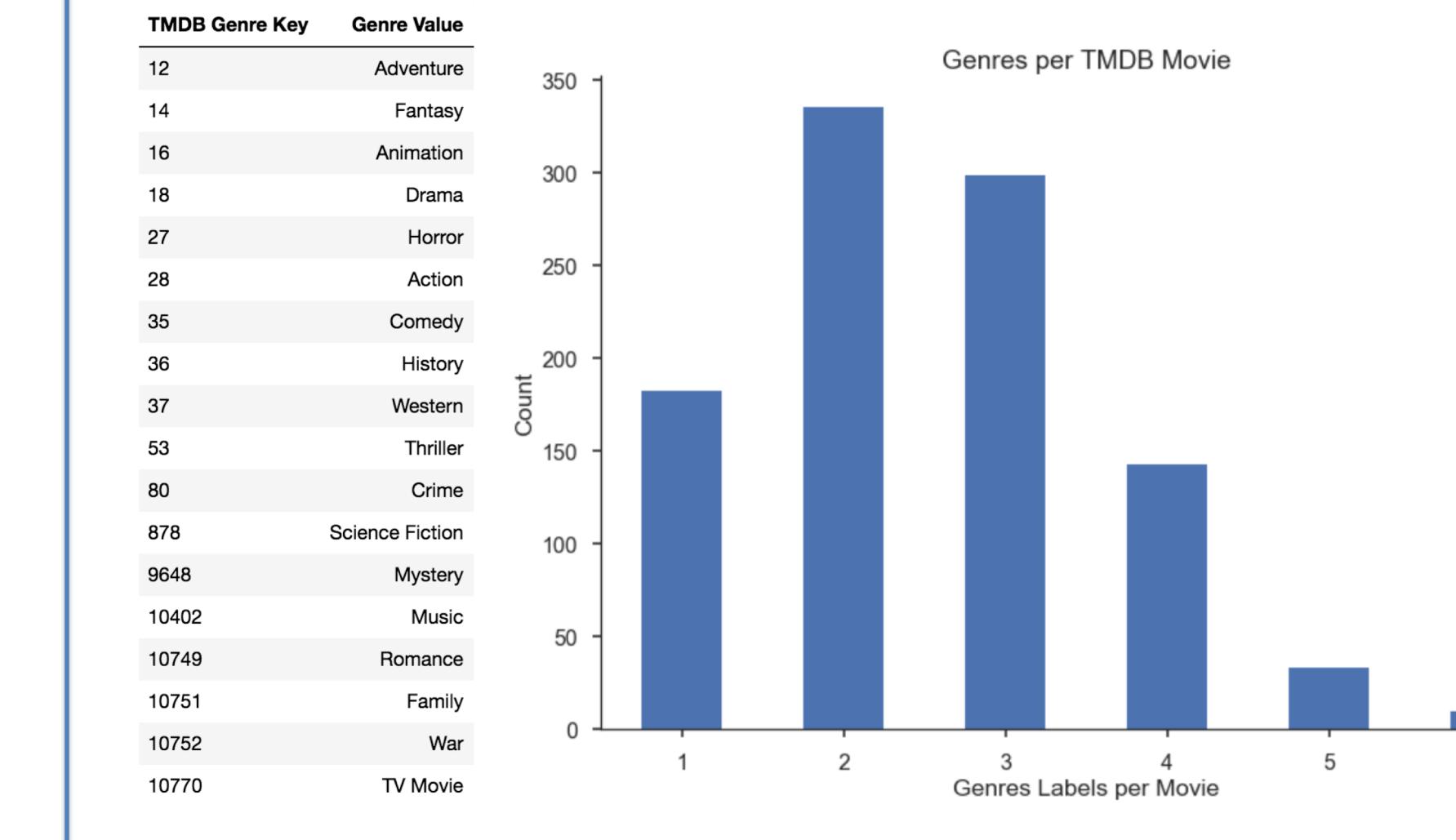


Figure 3. Genres per movie

Pre-modeling Data Analysis

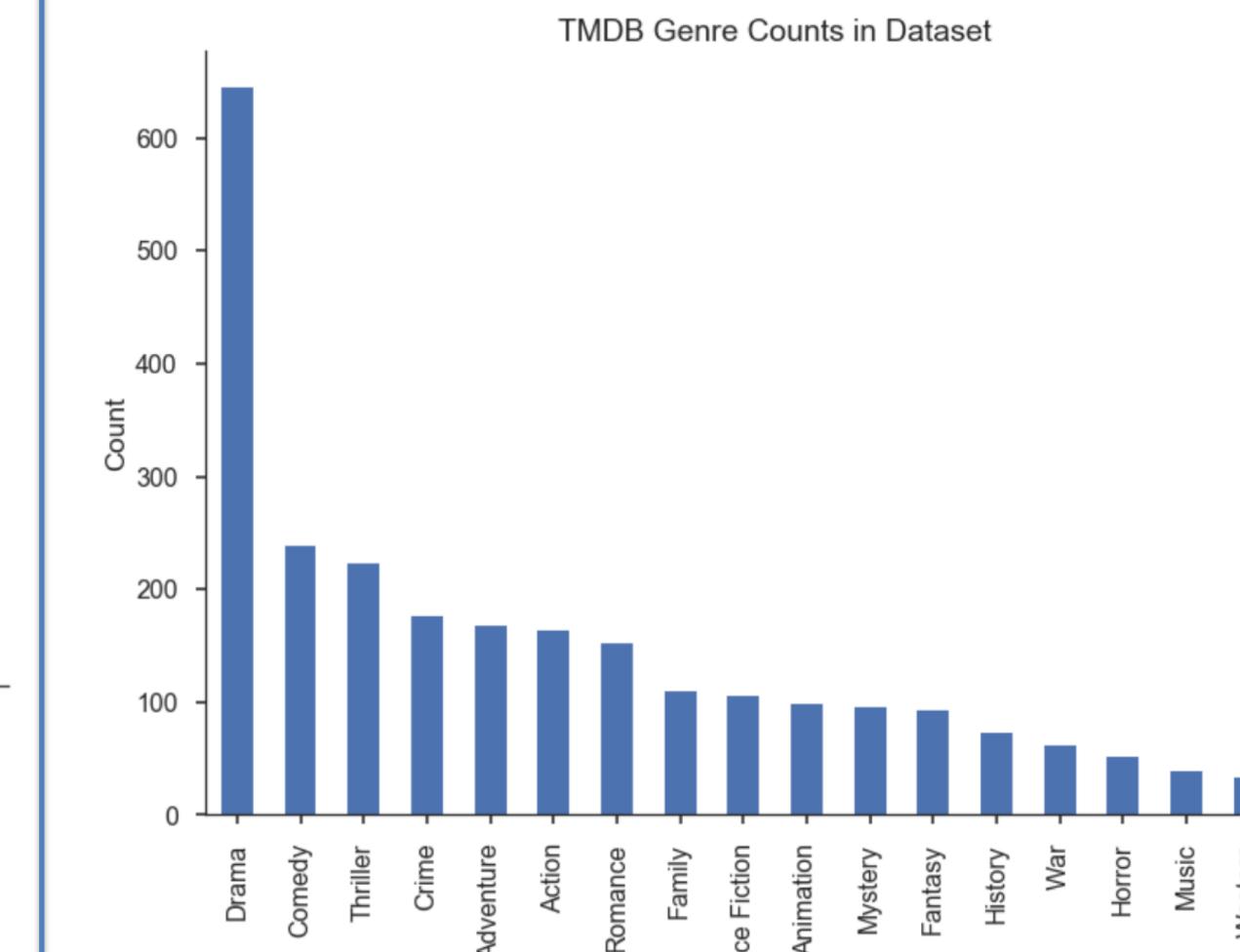


Figure 6. Genre distribution

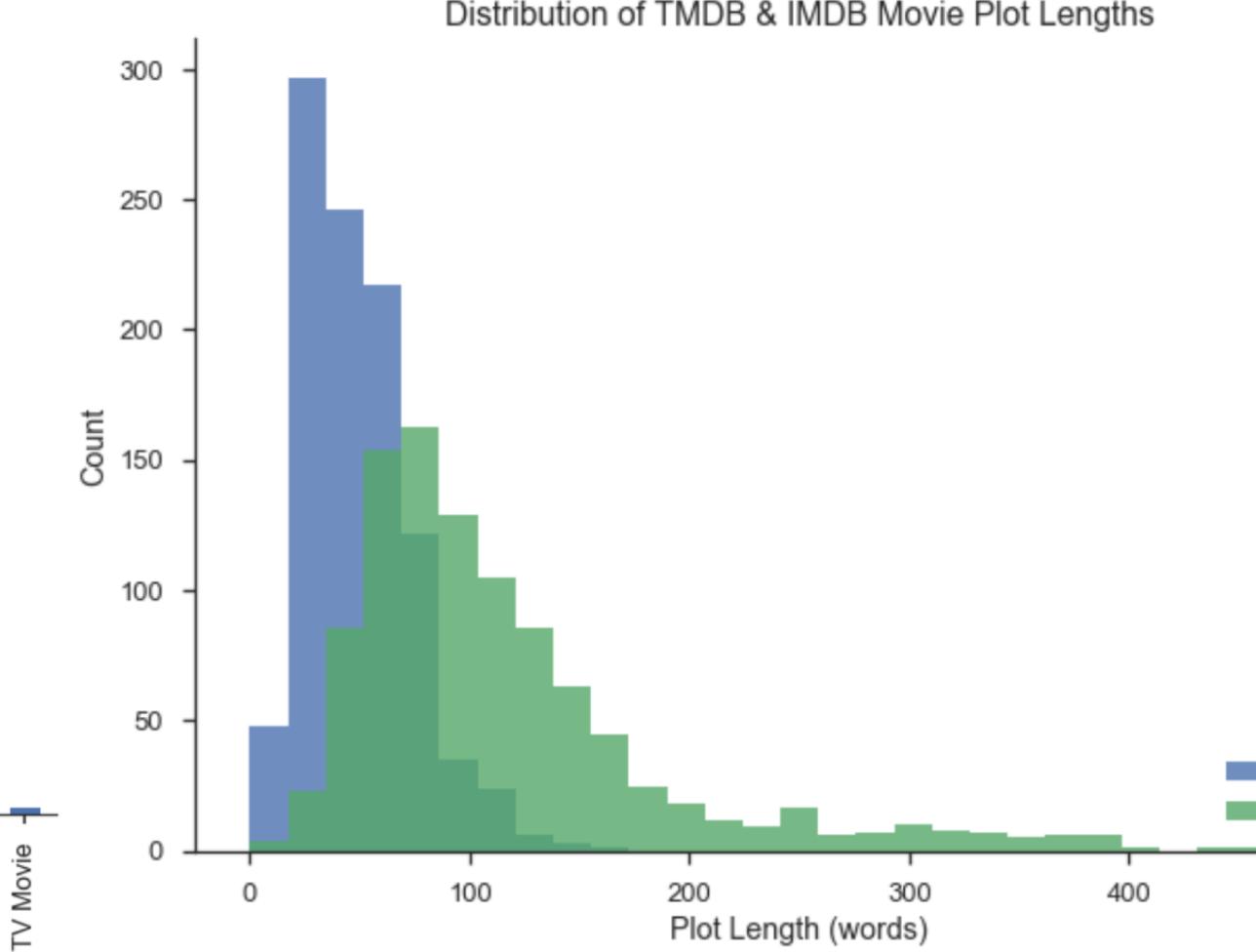


Figure 7. TMDB & IMDB plot lengths distribution

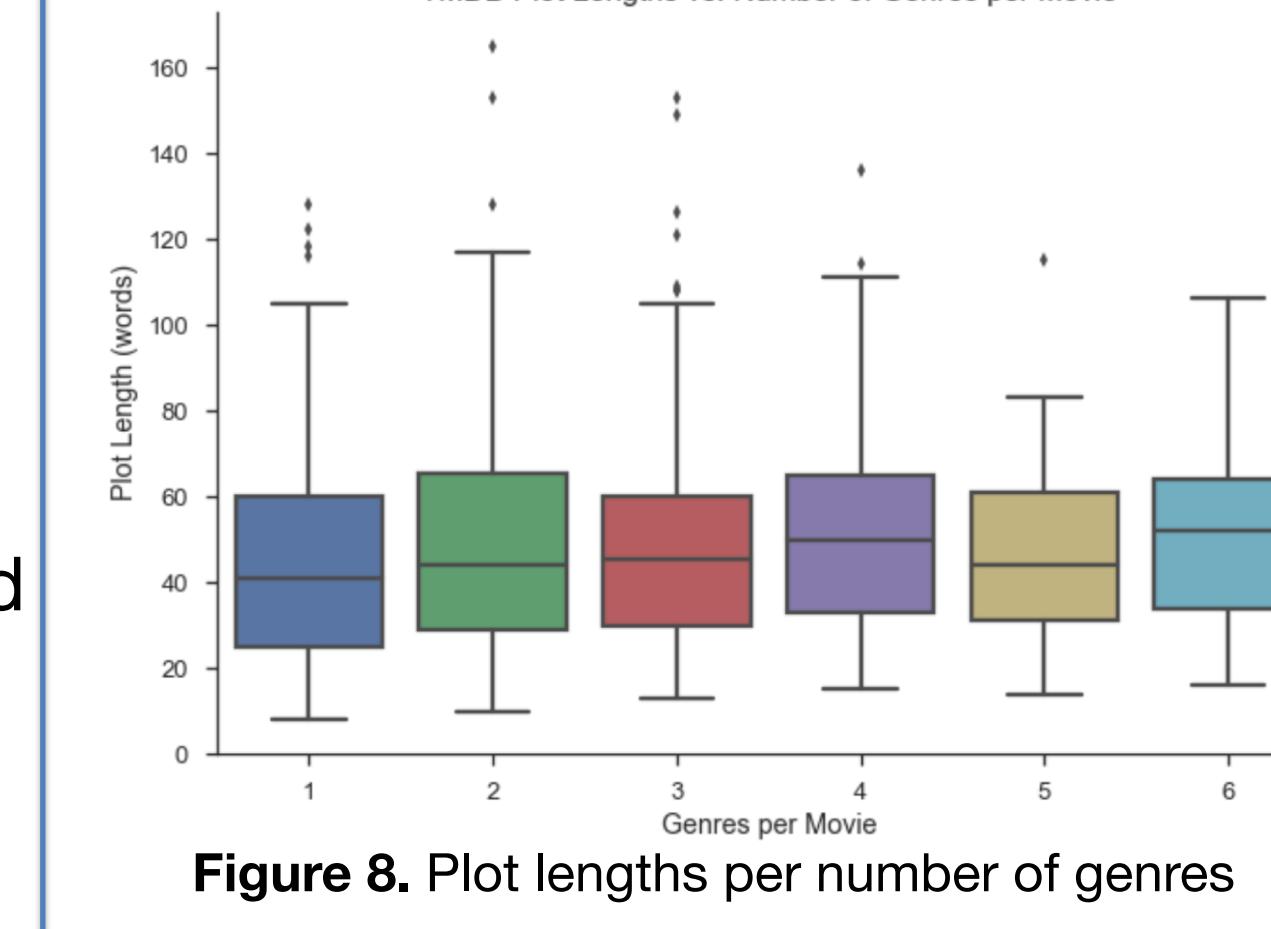


Figure 8. Plot lengths per number of genres

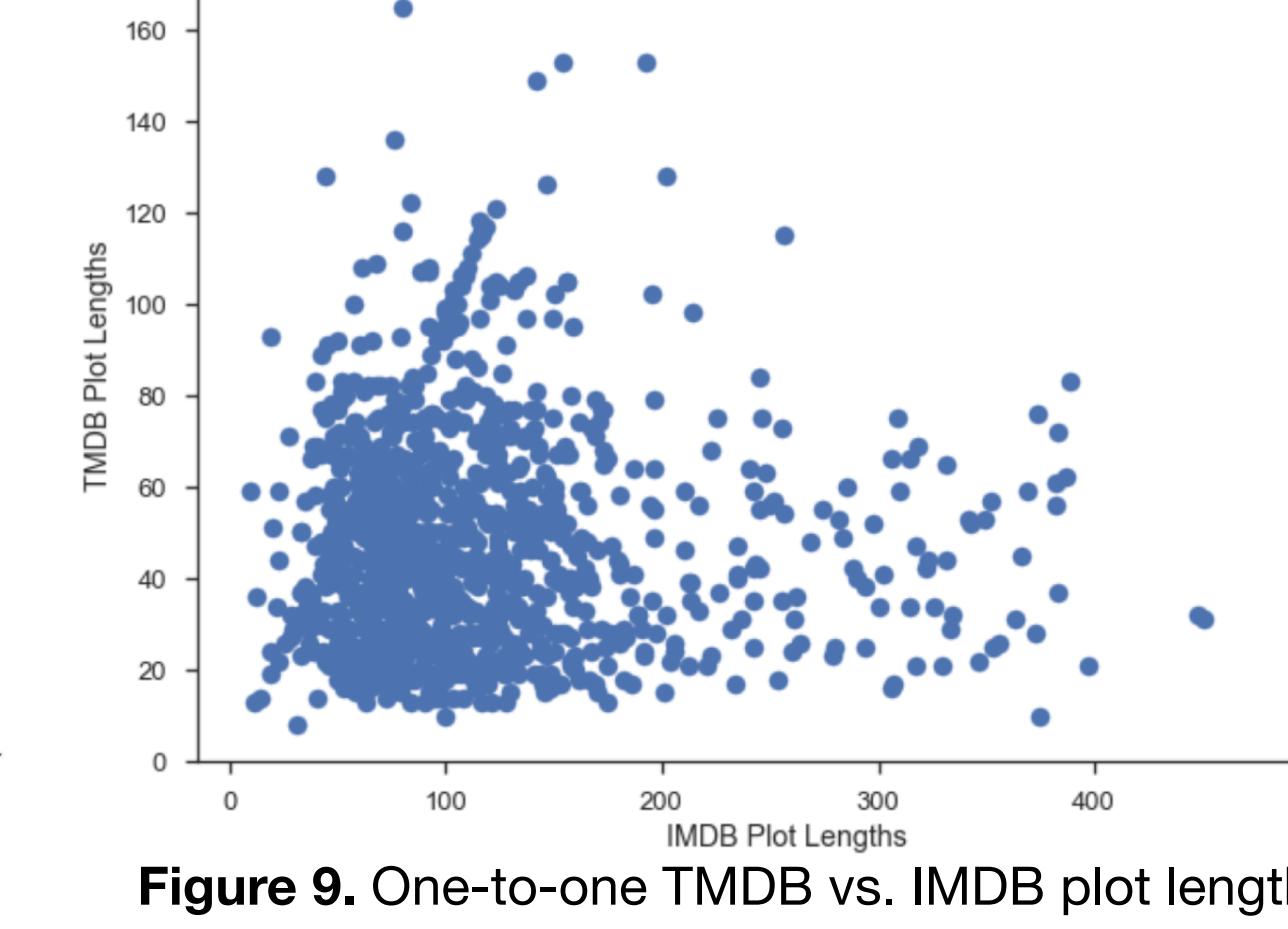


Figure 9. One-to-one TMDB vs. IMDB plot lengths

- Dataset skewed toward dramas (Figure 6)
- 2.5 TMDB genres per film
- IMDB has longer plots (Figure 7)
- One-to-one comparison has most IMDB plots double TMDB length (Figure 9).
- Plot lengths equal across films with multiple genre labels (Figure 8)

Hypothesis

“The longer IMDB or combined IMDB and TMDB plots with potentially richer w2v or d2v representations will lead to improved precision and recall metrics over sparser bag-of-words and shorter TMDB plot vector representations.”

Modeling, Results & Conclusions

Features: 27 Models trained on 9 plot transformations (IMDB, TMDB, combined plots)

Multi-label Classification Models: Naïve-Bayes, Support Vector Machines (SVM), SVM with stochastic gradient descent (SGD)

Result Metrics: precision, recall, F1 (Figure 11)

actual value	P		n		total		Precision = $\frac{T_p}{T_p + F_p}$	Recall = $\frac{T_p}{T_p + F_n}$	$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$
	True Positive	False Positive	False Negative	True Negative	P	n			
p'	T _p	F _p	F _n	T _n	P'	N'			
n'	F _p	T _n	T _p	F _n	N'	P'			
total	P	N							

Figure 11. Result metrics

	actual value	precision	recall	precision	recall
Naive-Bayes-combined_bow	0.725255	0.377953	0.982521	0.926914	
Naive-Bayes-combined_doc_vec	0.166339	0.255906	0.163924	0.254321	
Naive-Bayes-combined_w2v_mean	0.456782	0.267717	0.583491	0.274074	
Naive-Bayes-imdb_bow	0.598621	0.340551	0.948176	0.780741	
Naive-Bayes-imdb_doc_vec	0.166339	0.255906	0.163924	0.254321	
Naive-Bayes-imdb_w2v_mean	0.178442	0.253937	0.605336	0.262716	
Naive-Bayes-tmdb_bow	0.619958	0.301181	0.970122	0.79358	
Naive-Bayes-tmdb_doc_vec	0.166339	0.255906	0.163924	0.254321	
Naive-Bayes-tmdb_w2v_mean	0.316826	0.261811	0.949266	0.265185	
SGD-combined_bow	0.166339	0.255906	0.163924	0.254321	
SGD-combined_doc_vec	0.166339	0.255906	0.163924	0.254321	
SGD-combined_w2v_mean	0.514741	0.409449	0.602544	0.42963	
SGD-tmdb_bow	0.166339	0.255906	0.163924	0.254321	
SGD-tmdb_doc_vec	0.166339	0.255906	0.163924	0.254321	
SGD-tmdb_w2v_mean	0.470563	0.379921	0.619262	0.430617	
SVC-combined_bow	0.678057	0.525591	0.95214	0.985679	
SVC-combined_doc_vec	0.186744	0.525591	0.194323	0.532346	
SVC-combined_w2v_mean	0.558305	0.688976	0.742788	0.893827	
SVC-tmdb_bow	0.807838	0.551181	0.87196	0.974321	
SVC-tmdb_doc_vec	0.153214	0.340551	0.162445	0.37679	
SVC-tmdb_w2v_mean	0.537334	0.685039	0.681017	0.864691	
SVC-tmdb_bow	0.475574	0.496063	0.815313	0.946667	
SVC-tmdb_doc_vec	0.186744	0.525591	0.194323	0.532346	
SVC-tmdb_w2v_mean	0.540201	0.633858	0.748116	0.877531	