

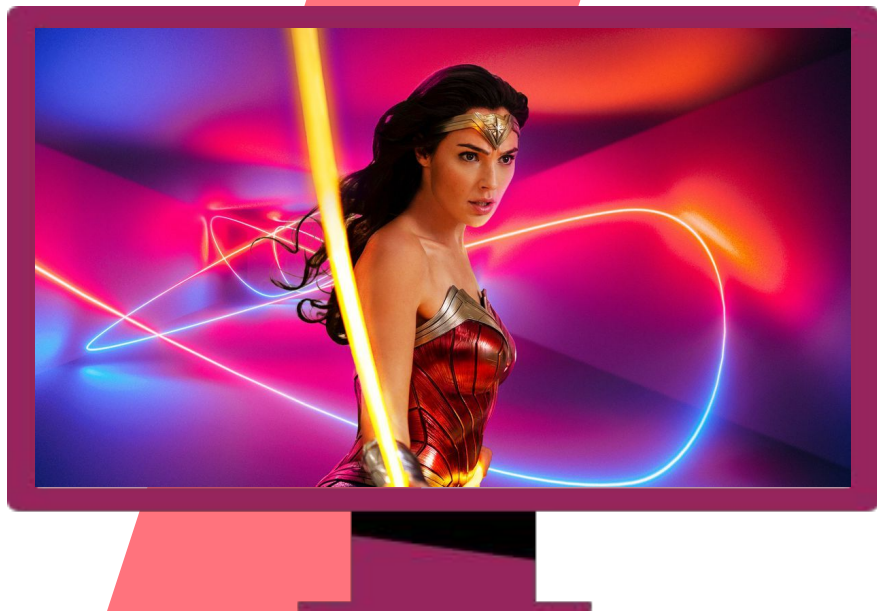


Analyzing Gender Gap In Hollywood Movies

by Okey Ofili



Using Unsupervised and Supervised Models to Understand Gender Roles



Problem Statement

Applying **unsupervised** and **supervised** models on movie roles **scraped** from Wikipedia and IMDB to analyze gender bias in hollywood movies.

6,606

HOLLYWOOD MOVIES

8,266

UNIQUE ACTORS/ACTRESSES

21,250

SCRAPED MOVIE ROLES/CHARACTERS

The Process

1

Scrapping and Merging Data

Methodology for gathering the Data and merging different datasets across IMDB and Wikipedia

2

Exploratory Analysis and Modeling

Inference and Predictive analysis on the Data set.

3

Observations, Issues and Next Steps

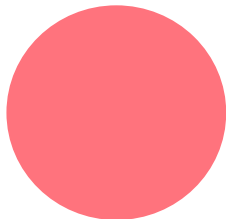
Discuss observations in the data, issues faced during the project execution stage and next steps

“Without data you are just
another person with an opinion.”

– W. Edwards Demmings

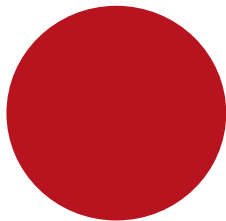


The Scraping Process



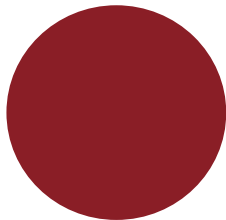
IMDB

This was pre-scraped from kaggle data-set for movies with over 100 votes released 1st Quarter of 2020.



NNDB

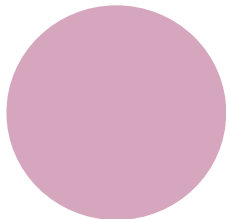
Data was scraped with beautiful-soup to get cast members race, but the data-set was too unbalanced so was not used again.



Wikipedia

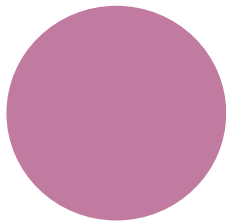
Utilized Wikipedia API to scrape website and get cast roles. This tasks took a whopping 40 hours! Then wrote code to dissect the string returned into usable parts.

The StopWord Process for Names



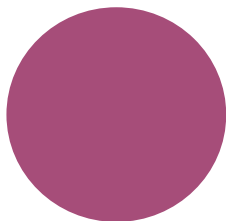
NLTK

Standard NLTK package, I identified all Nouns and then removed names by cross-checking against `wordnet.synsets`



IBM Watson

This was the most accurate prediction of names of Characters and people's, but I hit a API limit in the middle of the process, so had to abandon.



SpacyIO

This was a hybrid between IBM and NLTK, and performed better. It basically is able to extract Named Entity Extraction and identify what is PERSON or CITY. And it was FREE.

1.72 x

RATIO OF MALE ACTORS CHILDREN TO FEMALE ACTORS

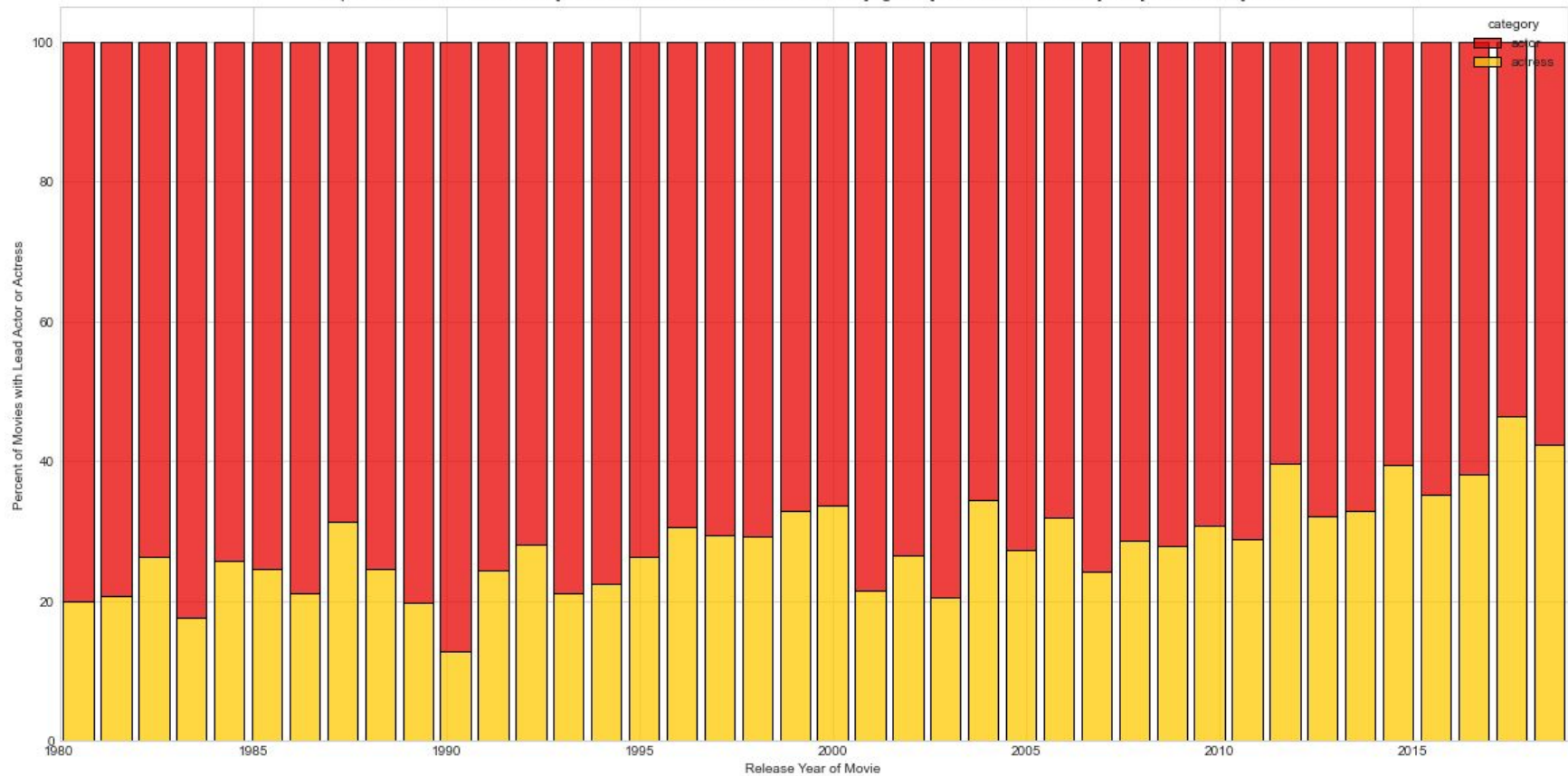
1.57 x

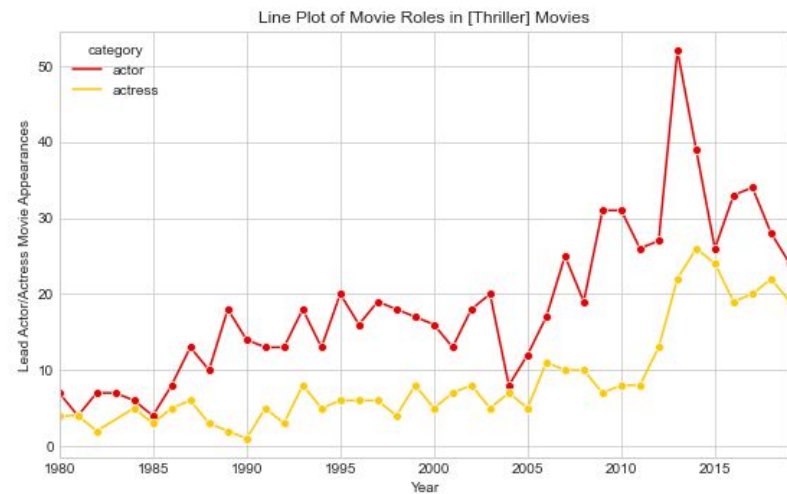
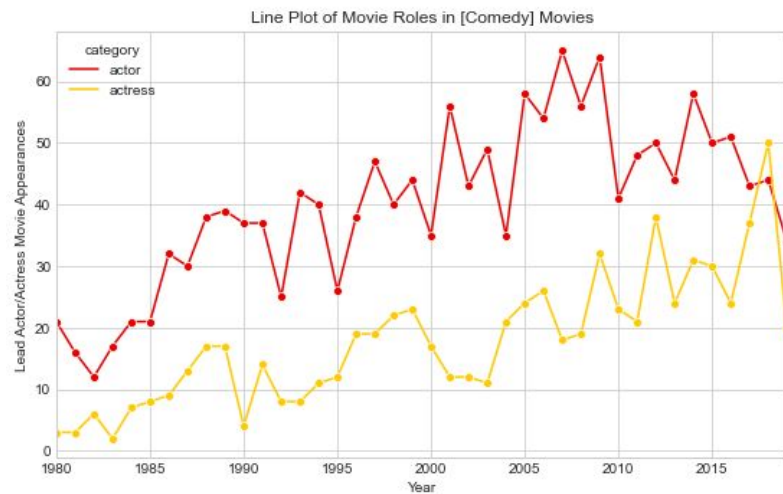
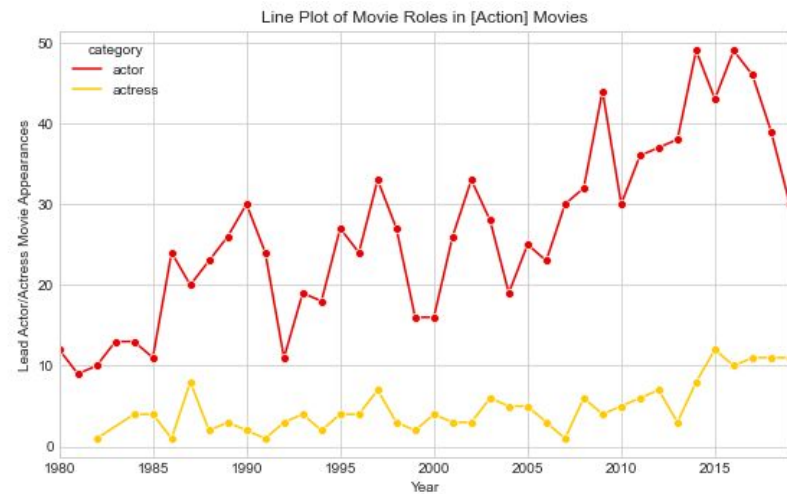
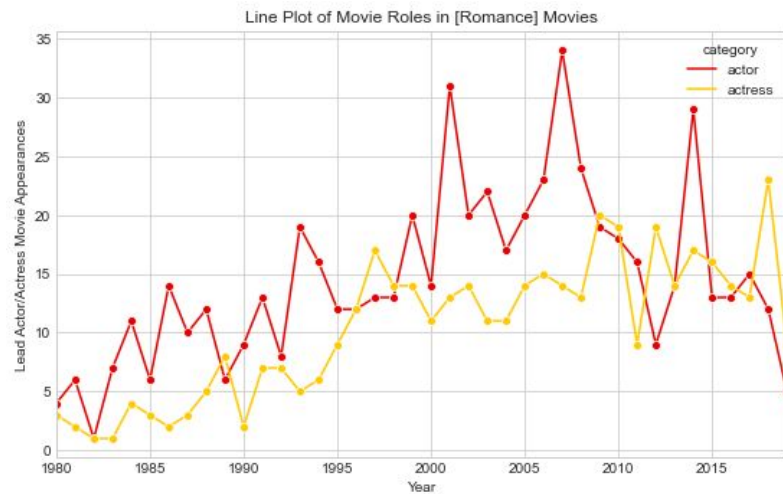
BUDGET RATIO OF MALE LED MOVIES TO FEMALE LED MOVIES

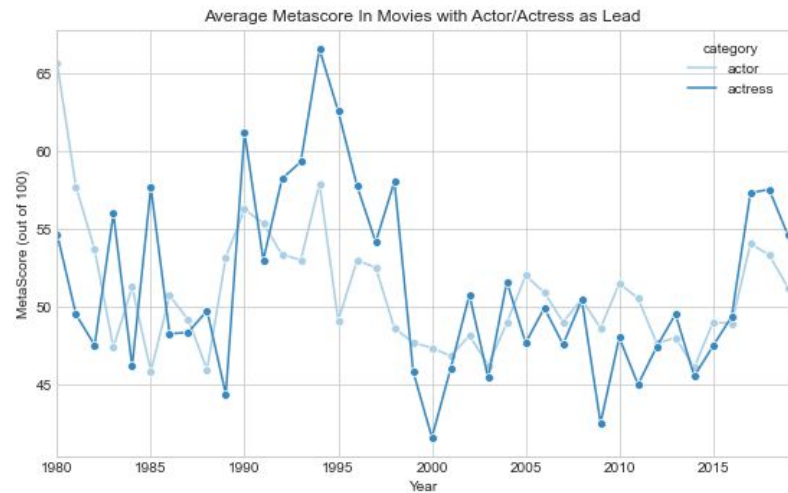
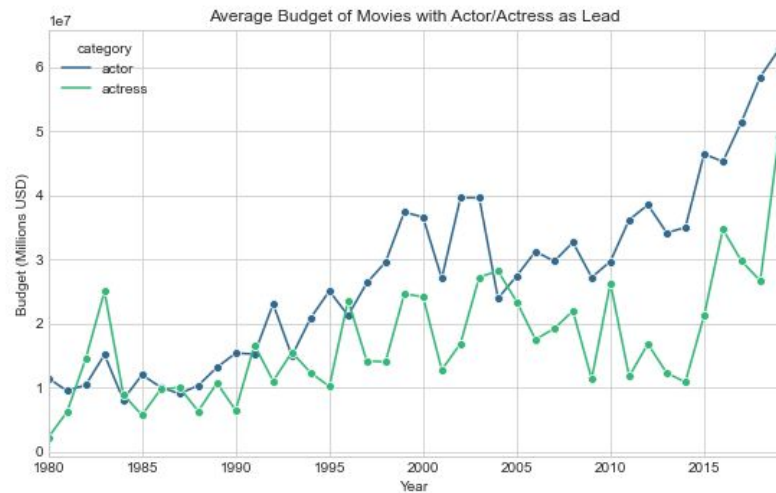
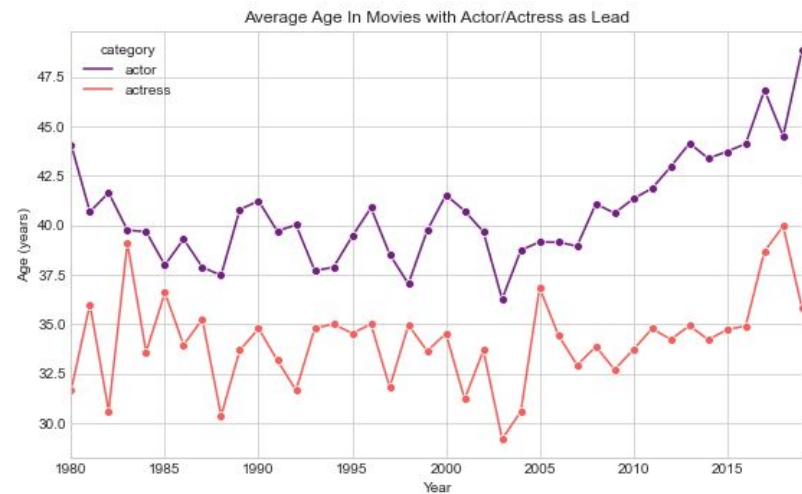
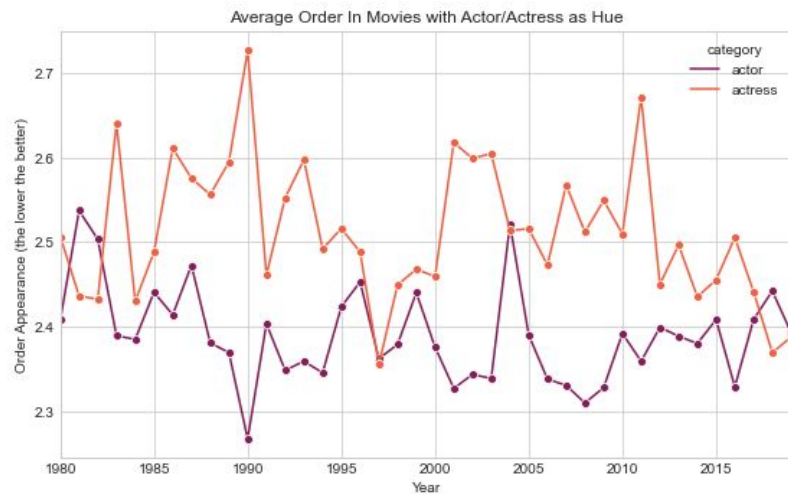
1.00 x

METAScore RATIO OF MOVIES

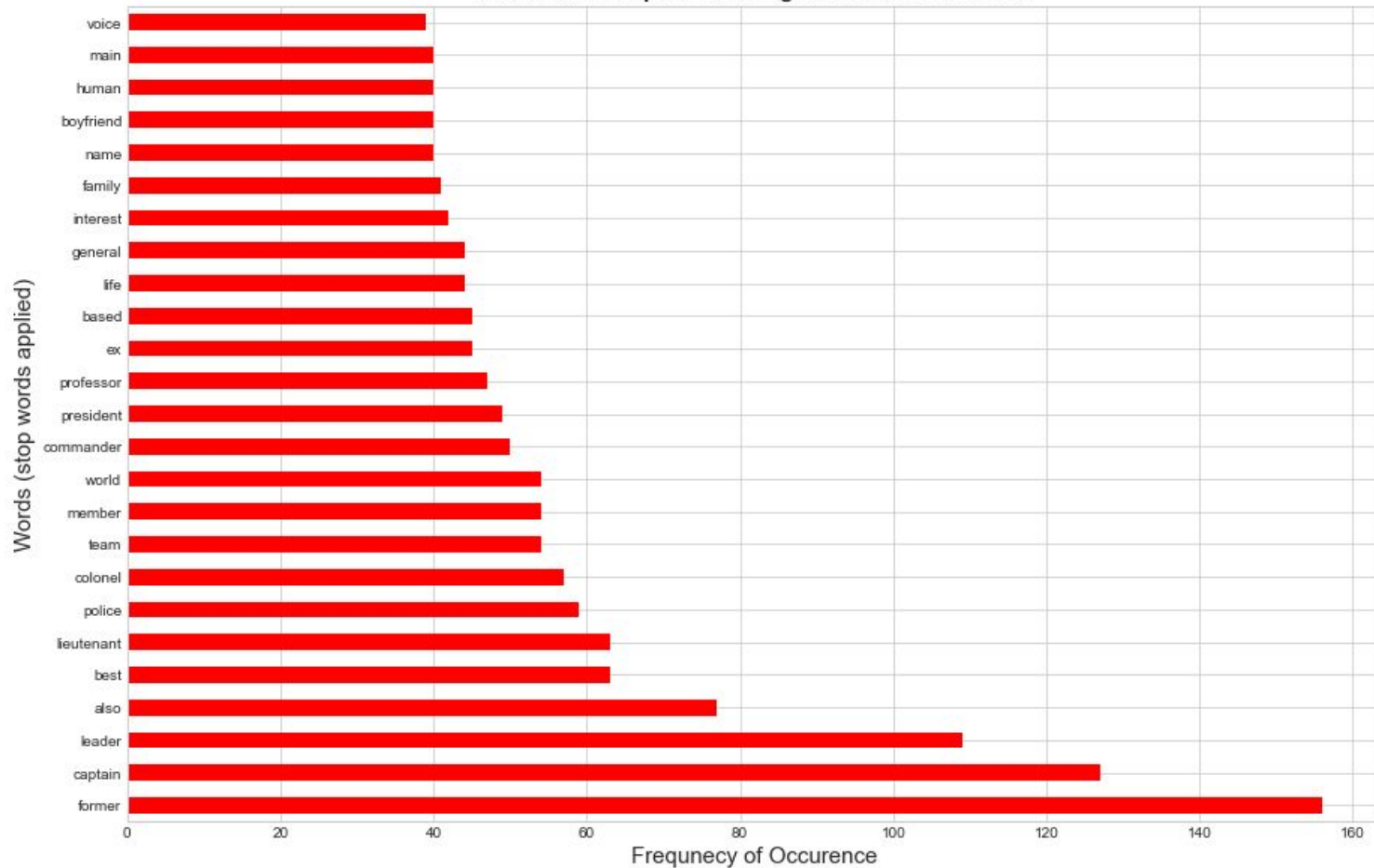
Equalized Stacked Bar Chart of [Percent of Movies with Lead Actor/Actress] against [Year of Movie Release] with [Actor/Actress] as Hue



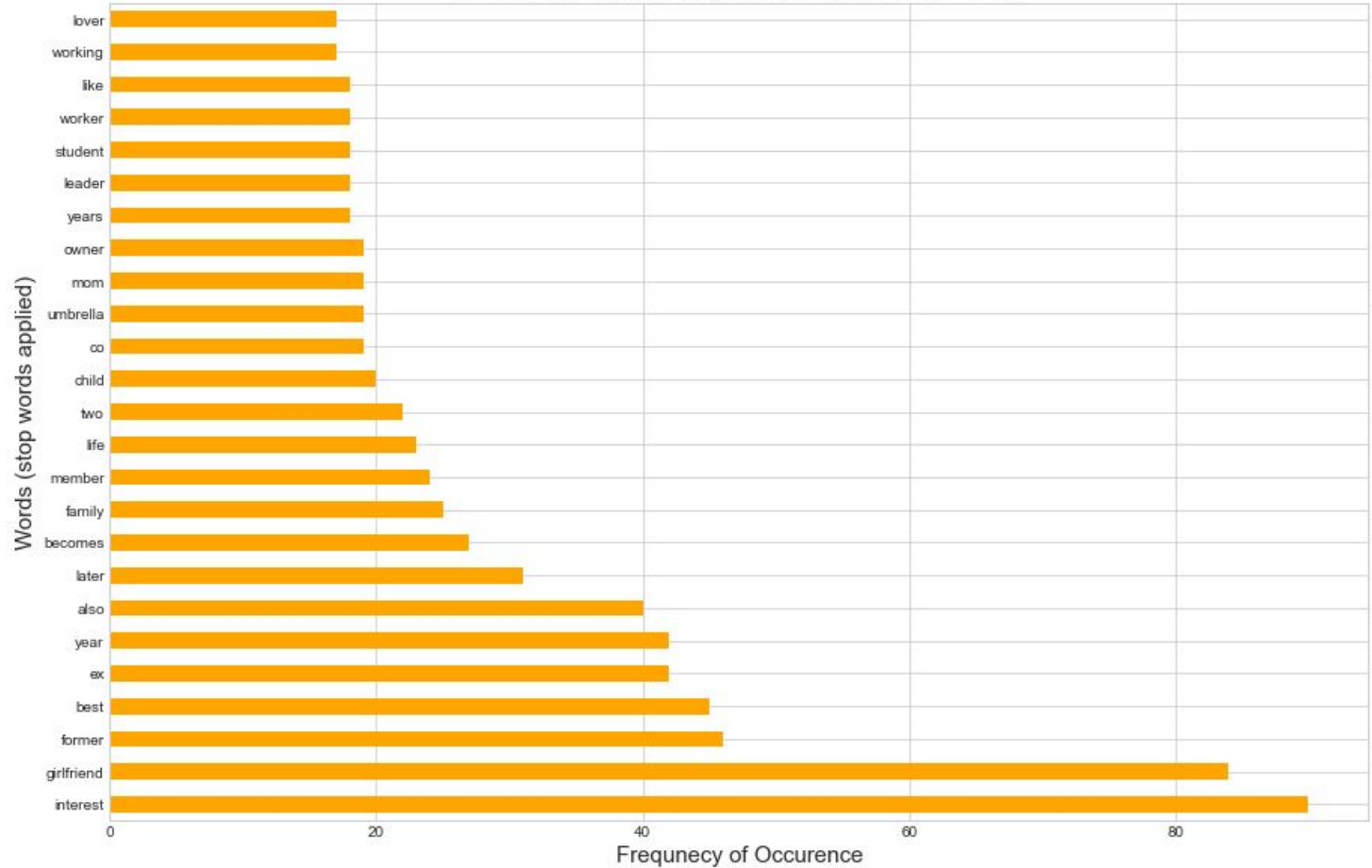




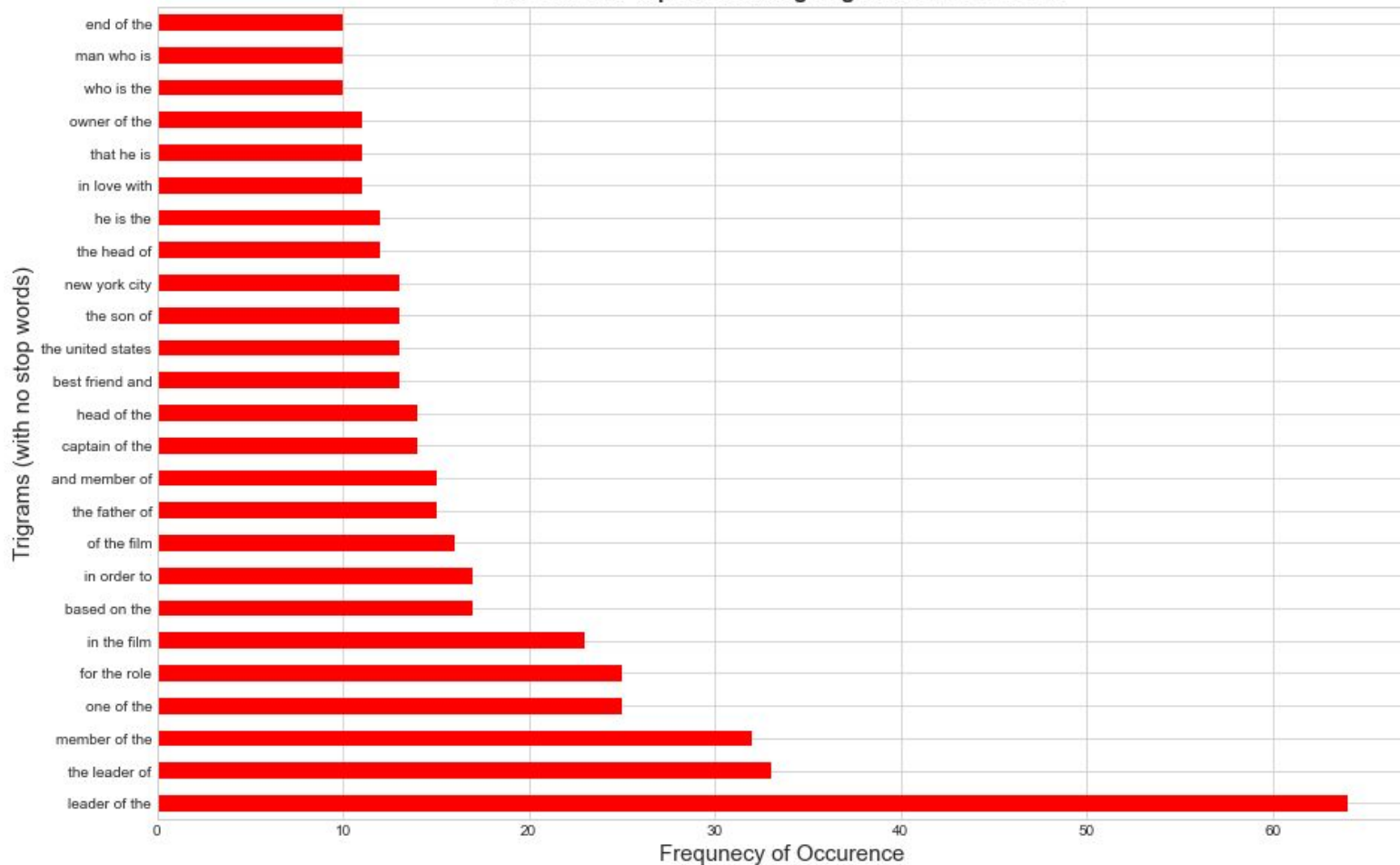
Bar Chart of Top 25 Occuring Words in Actor Roles



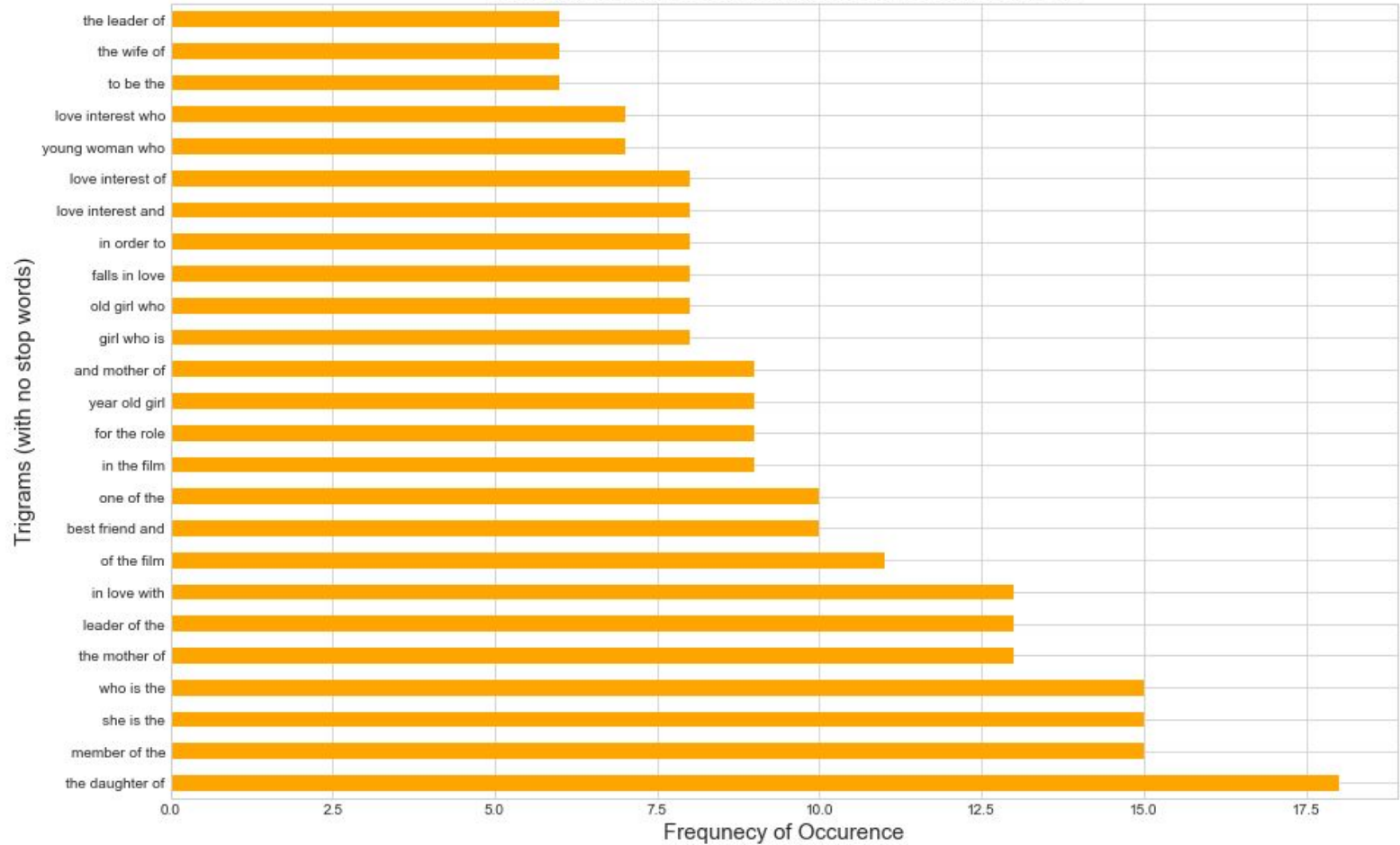
Bar Chart of Top 25 Occuring Words in Actress Roles



Bar Chart of Top 25 Occuring Trigrams in Actor Roles



Bar Chart of Top 25 Occuring Trigrams in Actress Roles



Prediction of Gender Based on Role?

89%

ACCURACY SCORE (NO STOPWORDS)

63.1%

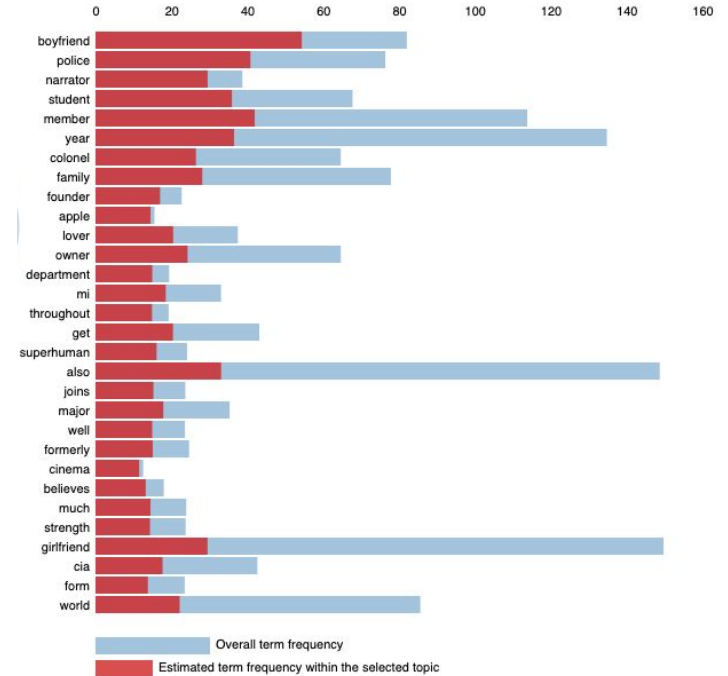
ACCURACY SCORE (WITH STOP WORDS)

Latent Dirichlet Allocation

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 4 (19.8% of tokens)



1. $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$ for topics t ; see Chuang et. al (2012)
2. $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

Issues

Too little text as the ROLES scrapped from wikipedia were mostly stub descriptions.

Syncing scraped data across different datasets leads to data loss. A better approach would be to **focus on one thing. Wikipedia or IMDB.**

Stopwords to remove names ended up removing important words like 'Nurse'



Key Observations

Genre has a clear influence on roles on a Gender basis. For action movies there is a significant gap versus Romance movies.

Movies with male leads tends to have **1.5 times** more budgeted to it, even though average metacritic scores are same across genres.

Next Steps is to get richer Data on the role descriptions or leverage movie transcripts to infer this. Also to get cast race added to dataset and sell it!



Thanks



email **oaofili@gmail.com** for further questions