



Streamlining the Screenplay Buying Process

By Michael Orlando

Our Problem

Hollywood receives about 50,000 thousand screenplays per year. But continues to make shitty movies. How can we improve this?



Business Understanding

One important element to movies is GENRE.

Genre brings in the target audience.

For example, if Disney was marketing a new Animated Fantasy movie, then the audience would be pretty upset if the film played like a Crime Horror movie.



Business Objective

To build a multi-label classifier, using Natural Language Processing, to classify the genres of a given screenplay.



The Machine -- Simplified

Screenplay

INT. HOUSE - NIGHT

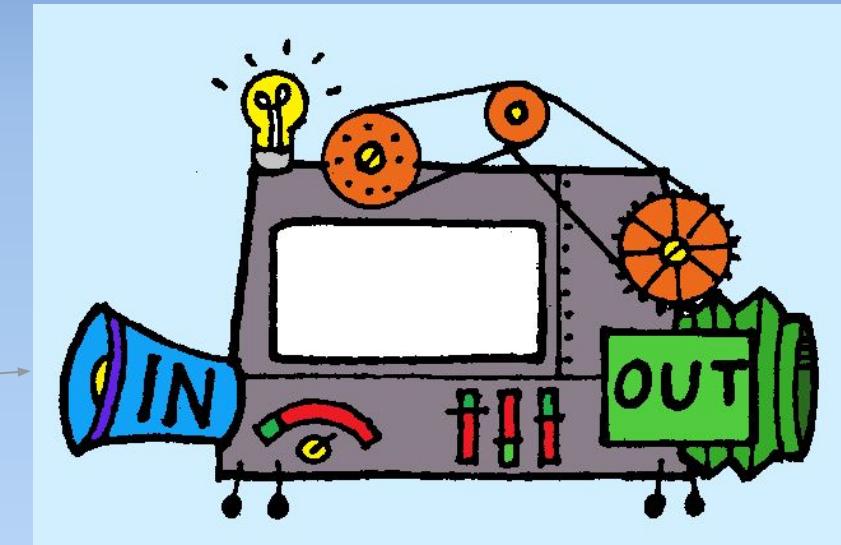
The Phone rings. Jason answers it.

JASON
(into phone)
Hello?

SALLY (V.O.)
Jason, I'm in trouble. I ran out of
gas. Can you bring some?

JASON
(listening)
Yes, I can be there soon.

Jason hangs up and runs out the door.



Genre Classification

[‘Action’, ‘Comedy’]

Data Collection

Screenplays:

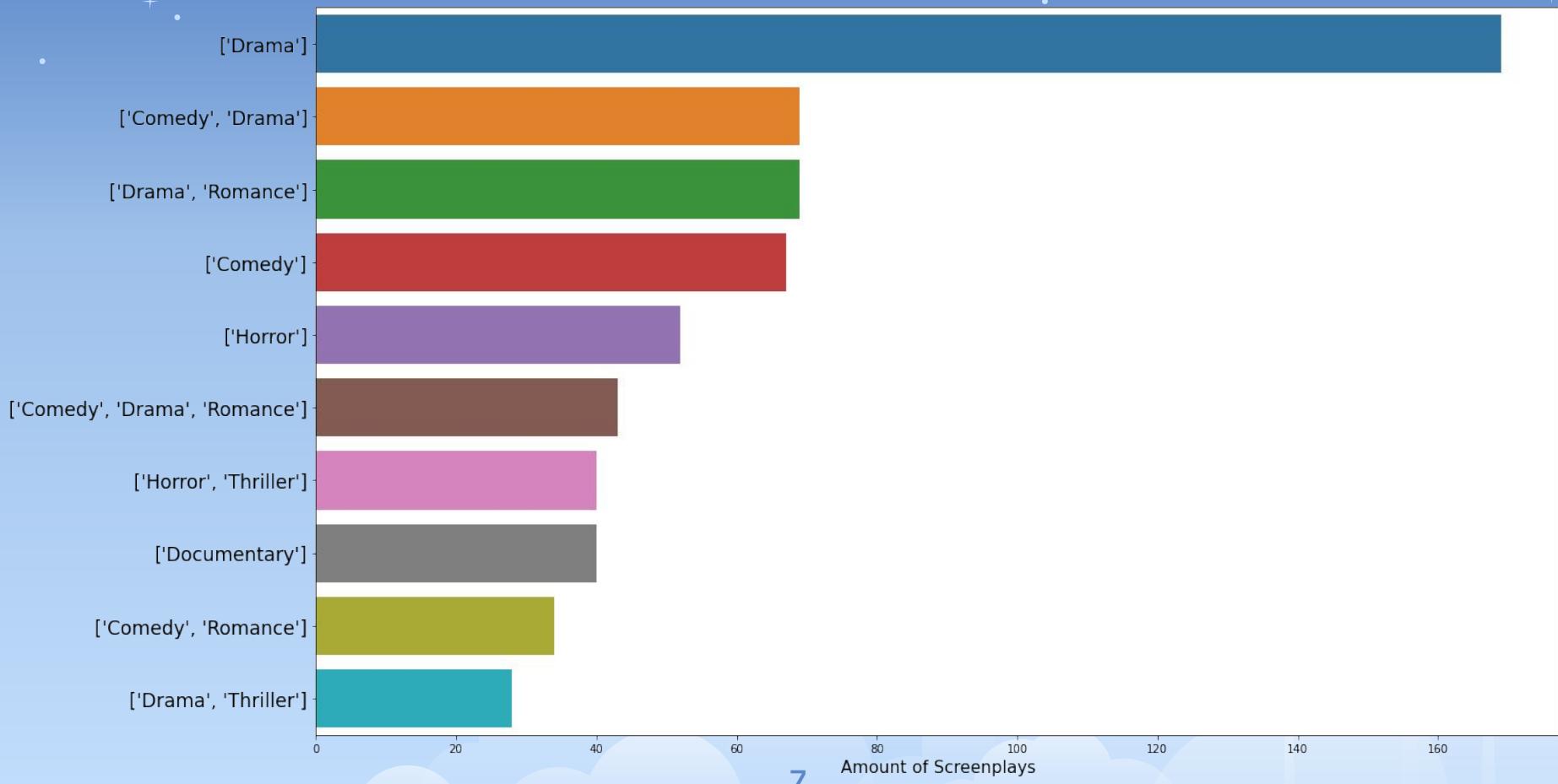
The scripts in pdf format were scrapped from The Script Savant using BeautifulSoup

Genres:

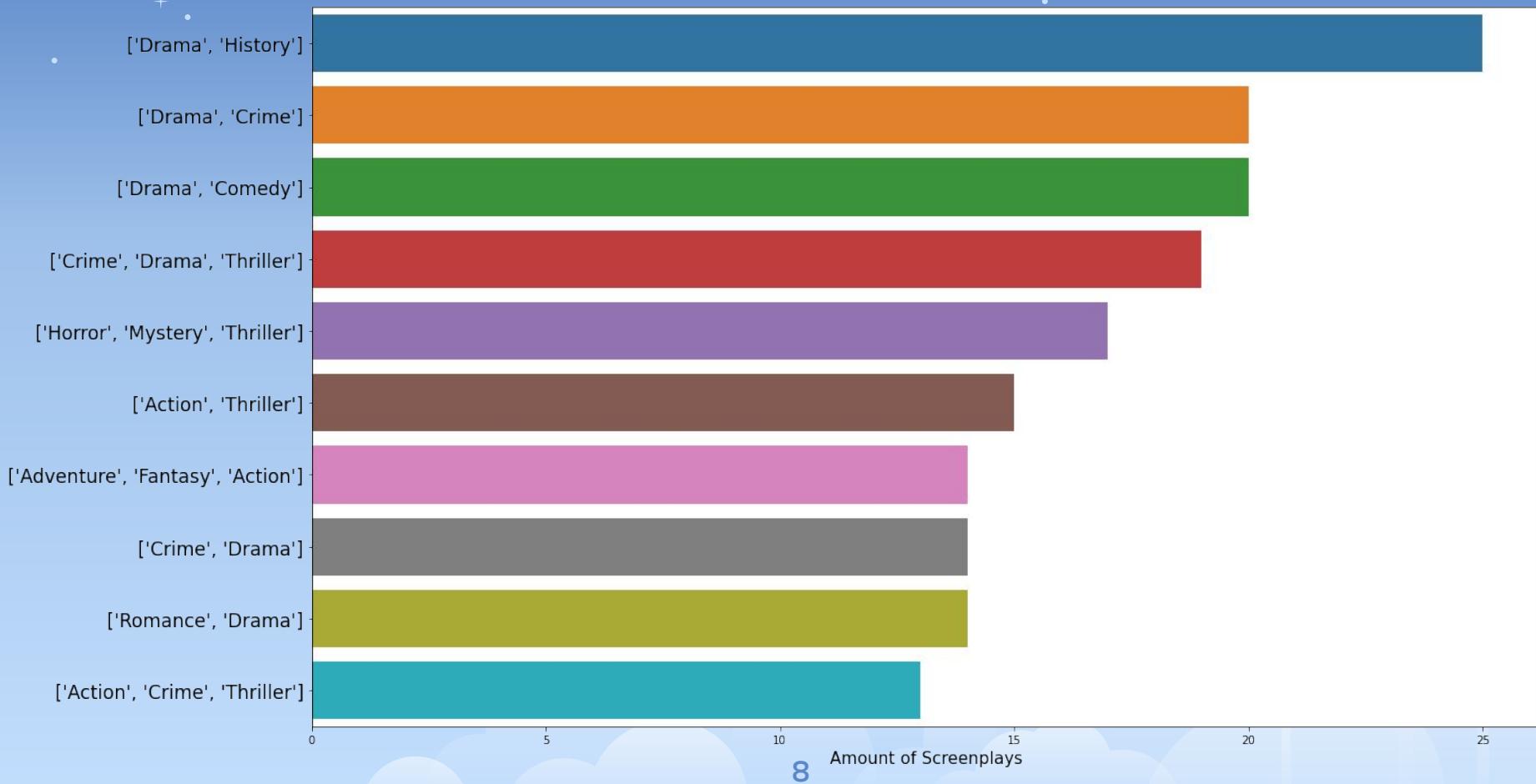
The screenplays were labelled with genres using the Movie Database API



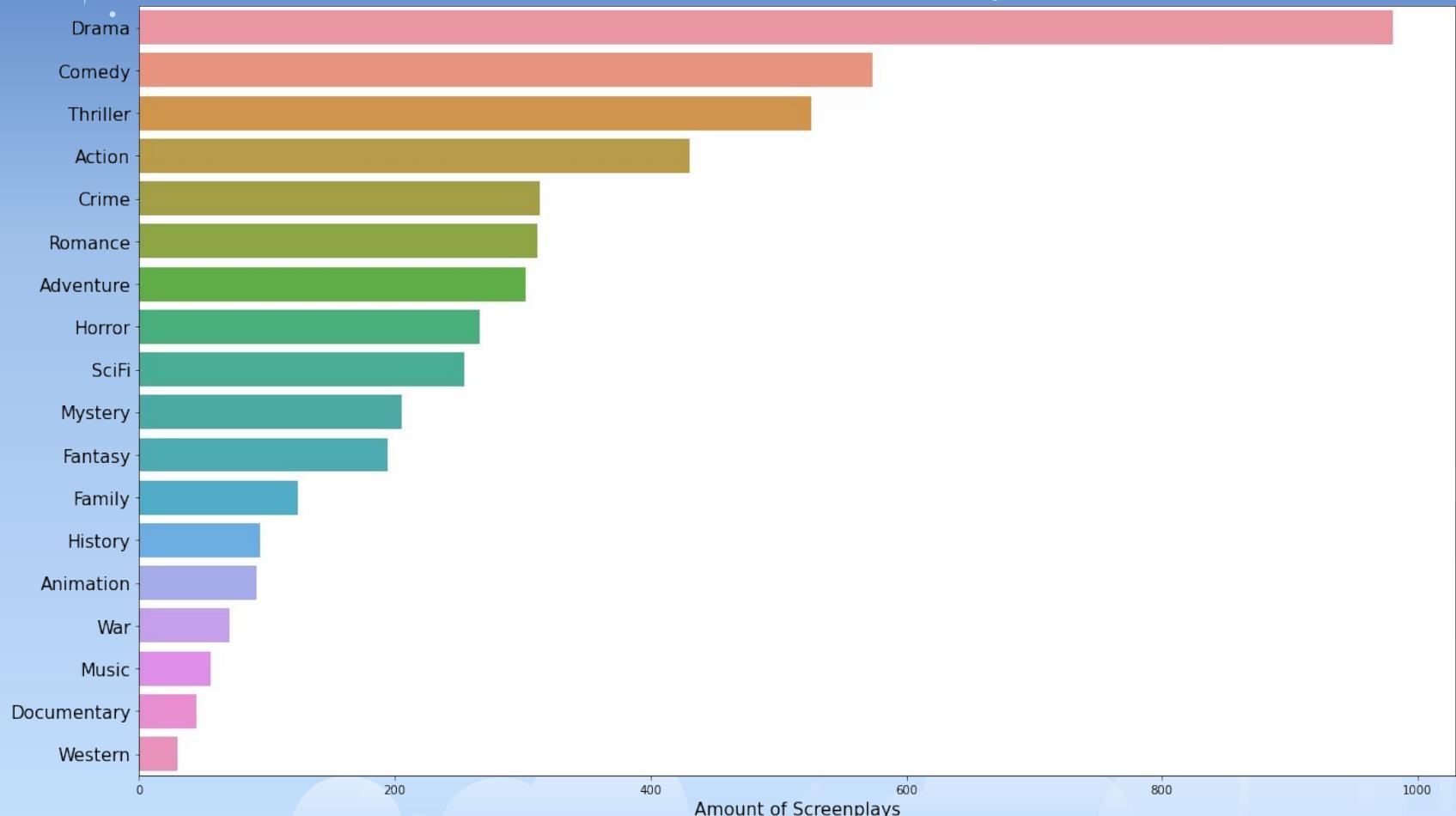
Genre Combination (1-10)



Genre Combination (11-20)



Target Distribution by Single Genre



Overall

- 2000 unique screenplays
- 18 different genres
- 600 unique combinations of genres
- Drama the most frequent in combinations
- Western the least frequent in combinations

Feature Engineering

The feature created is matrix of words
(count or frequency)

I improved the matrix of words by
testing different preprocessing
techniques & eliminating common
words.

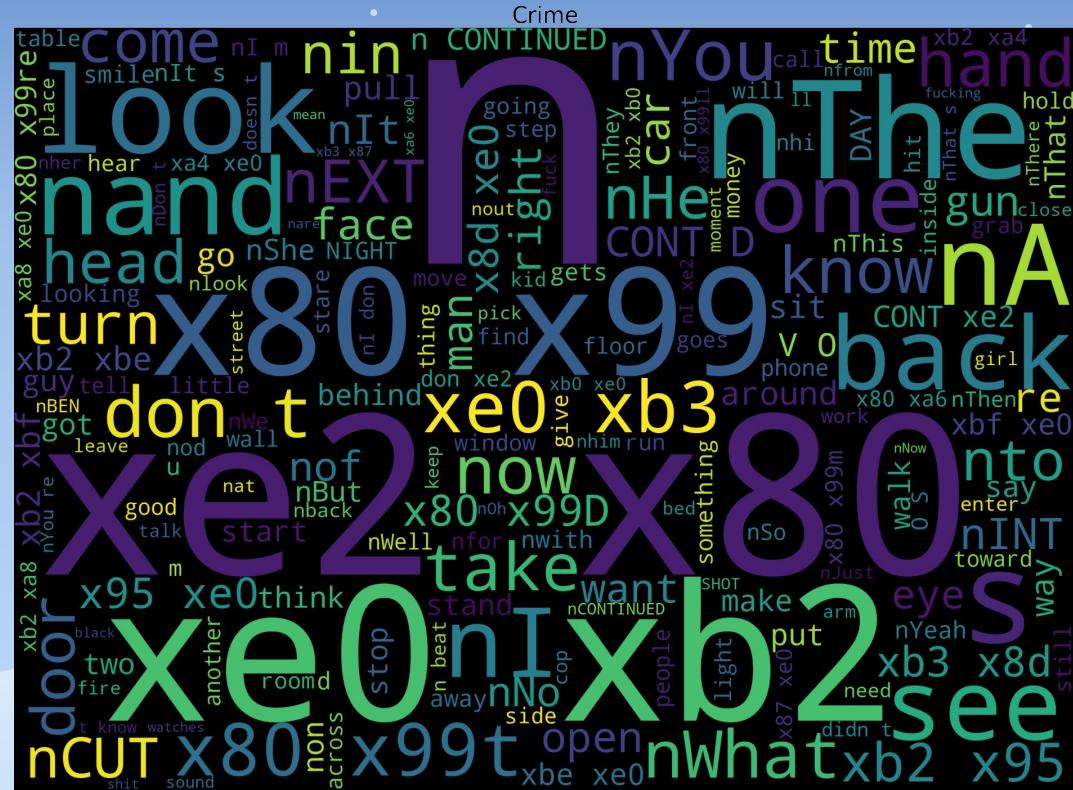
Results were illustrated using word
clouds

For example....

| | Jumps | The | brown | dog | fox | lazy | over | quick | the |
|------|-------|-----|-------|-----|-----|------|------|-------|-----|
| Doc1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 |
| Doc2 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 |

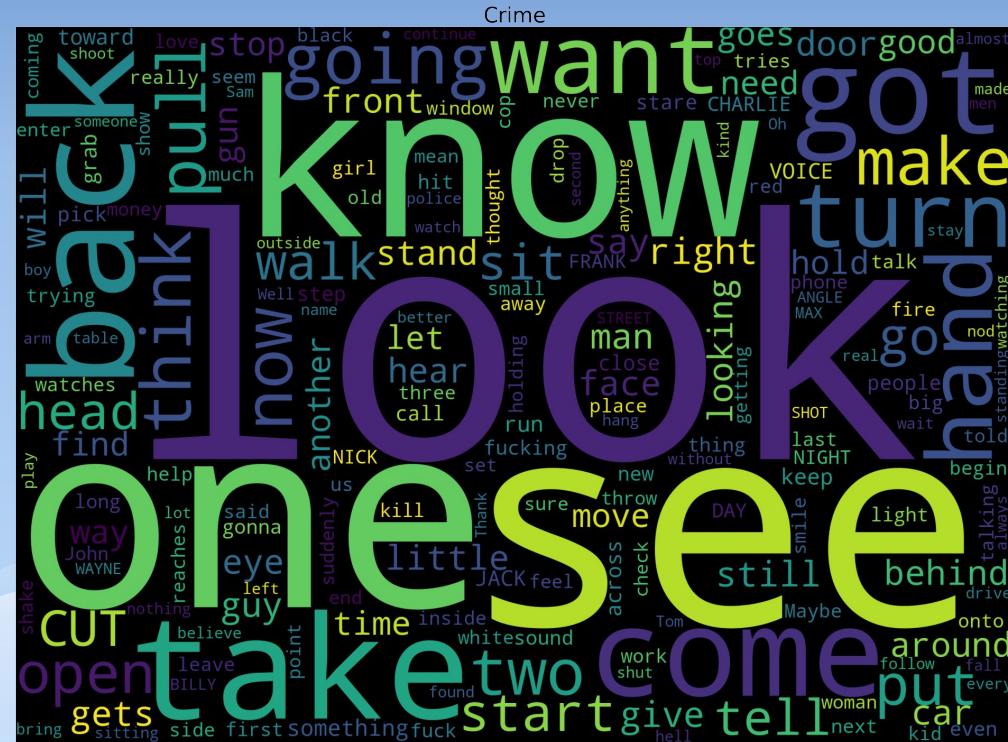
Before Cleaning

Text: Crime Genre

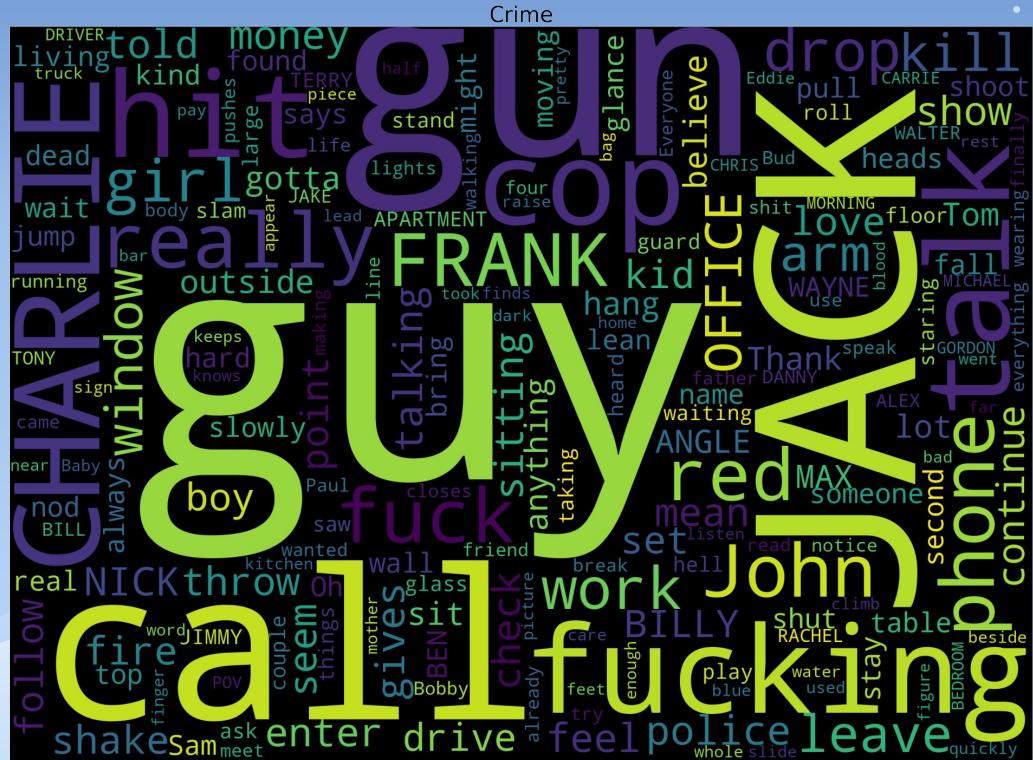


After Cleaning-- Crime Genre

Much better, but still basic



After Cleaning, Removal of Stop Words-- Crime Genre



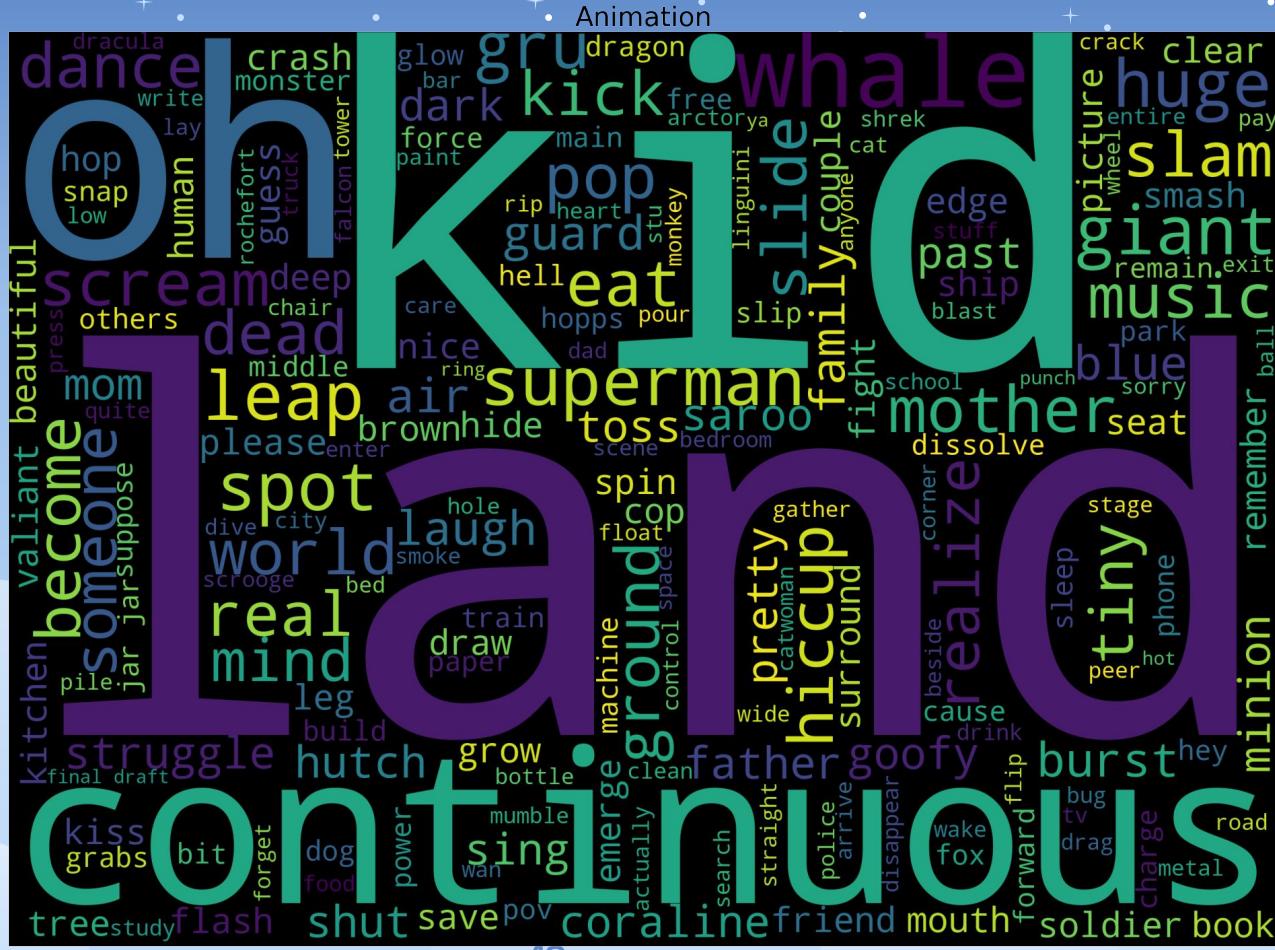
After cleaning,
lemmatizing, and
removal of stop
words: Crime
Genre



Action

Adventure

Animation



Comedy

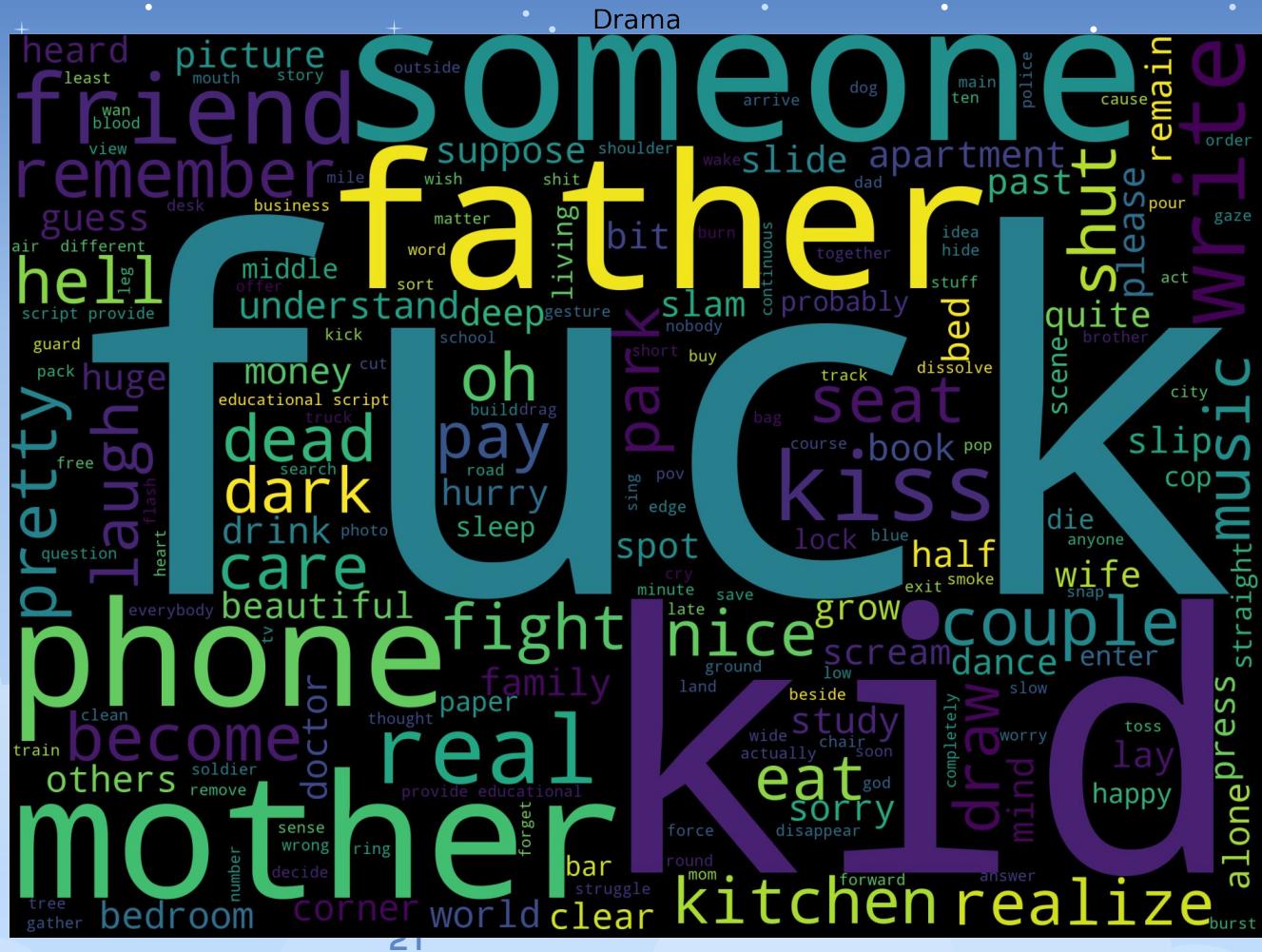
Comedy

Documentary

Documentary

A word cloud visualization titled "Documentary" centered on the word "press". The word "press" appears twice in the center, once in blue and once in green. Other prominent words include "execute" (blue), "continuous" (yellow), "colonel" (yellow), "kittie" (teal), "superman" (purple), "priest" (teal), "president" (blue), "helmut" (purple), "guard" (purple), "couple" (purple), "city" (yellow), "someone" (yellow), "wario" (green), "pay" (purple), "deep" (green), "shut" (black), "slam" (black), "anakin" (purple), "edge" (green), "warrio" (green), "blast" (purple), "happy" (purple), "kitchen" (purple), "khurram" (black), "heard" (black), "eat" (black), "pretty" (black), "others" (black), "train" (black), "faint" (black), "shout" (black), "burst" (black), "remember" (black), "surround" (black), "sell" (black), "remain" (black), "straight" (black), "force" (black), "bhairav" (black), "kikuchiyo" (black), "scene" (black), "fuck" (black), "care" (black), "toad" (black), "money" (black), "slow" (black), "leap" (black), "land" (black), "hell" (black), "round" (black), "paper" (black), "sense" (black), "kick" (black), "control" (black), "nad" (black), "write" (black), "tree" (black), "track" (black), "ghazala" (black), "nice" (black), "order" (black), "entire" (black), "shosanna" (black), "arrive" (black), "answer" (black), "lock" (black), "ten" (black), "bedroom" (black), "half" (black), "snap" (black), "jedi" (black), "breen" (black), "spin" (black), "park" (black), "word" (black), "bank" (black). The background features a faint illustration of a person's head and shoulders.

Drama



Family

Fantasy

mind flash leap cry snap bedroom crack kitchen
death together build building arrive helboy
pov emerge massive remove pay slow spot
afraid hears smash hurry soon arrive wide
smash hurry soon guess music mouth nice
edge stone shoulder strike force wide
music guess surround bag word human answer suppose
nice surround word kiss human write recharge
city blast eat power bed alone shut cut press cause
blast power bed alone someone straight gaze heard
offer glow search track whale sword lock hide
glove glow beneath others remain beautiful
tiny blue remain knife remain
huge enter past couple laugh save toss continuous
remain enter past couple laugh save toss continuous
tiny blue remain knife remain
slam realize peer image grow dance pop
tire circle enter past couple laugh save toss continuous
remain enter past couple laugh save toss continuous
tiny blue remain knife remain
dead creature burst phone

History

Horror

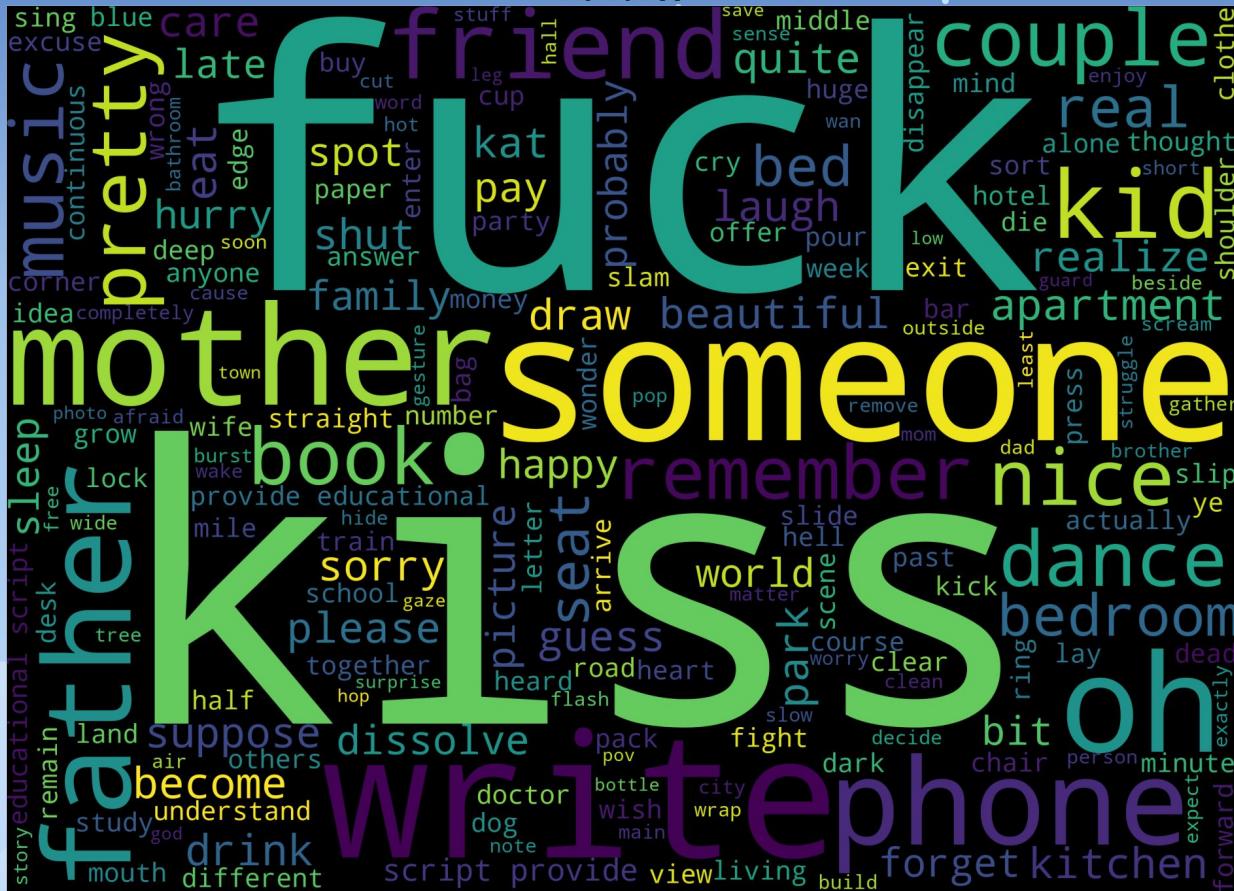
Music

Music
couple mom
city mozart
dog hum
shithurry number
beautiful
police
continuous
guess thought
ya hey
word fight alone
dark
family die
slayer heart
suppose act
probably
friend
act
slow
nobody
jam enter
picture
school boys
others
exit
se dead
hallway
half blood
studio
lay
build blue
become
drag
hawk
chair jolson nice
drum ride
seat money
sort huge
forget together
main hell
sleep ye
photo
Kid
herrick
trip quite
offer
oh
record
god
hour
piano
backstage
sugar
past heard
bag
note grandfather
green
pretty
eazy
party
main
corridor
flash
minute
lock
clear
grow
middle
hotbed
smoke
mother superior
brother burst
hop
edge
excuse
scene
cause
aur
slide wrong
worry
real
kick train
matter remove

Mystery

Romance

Romance



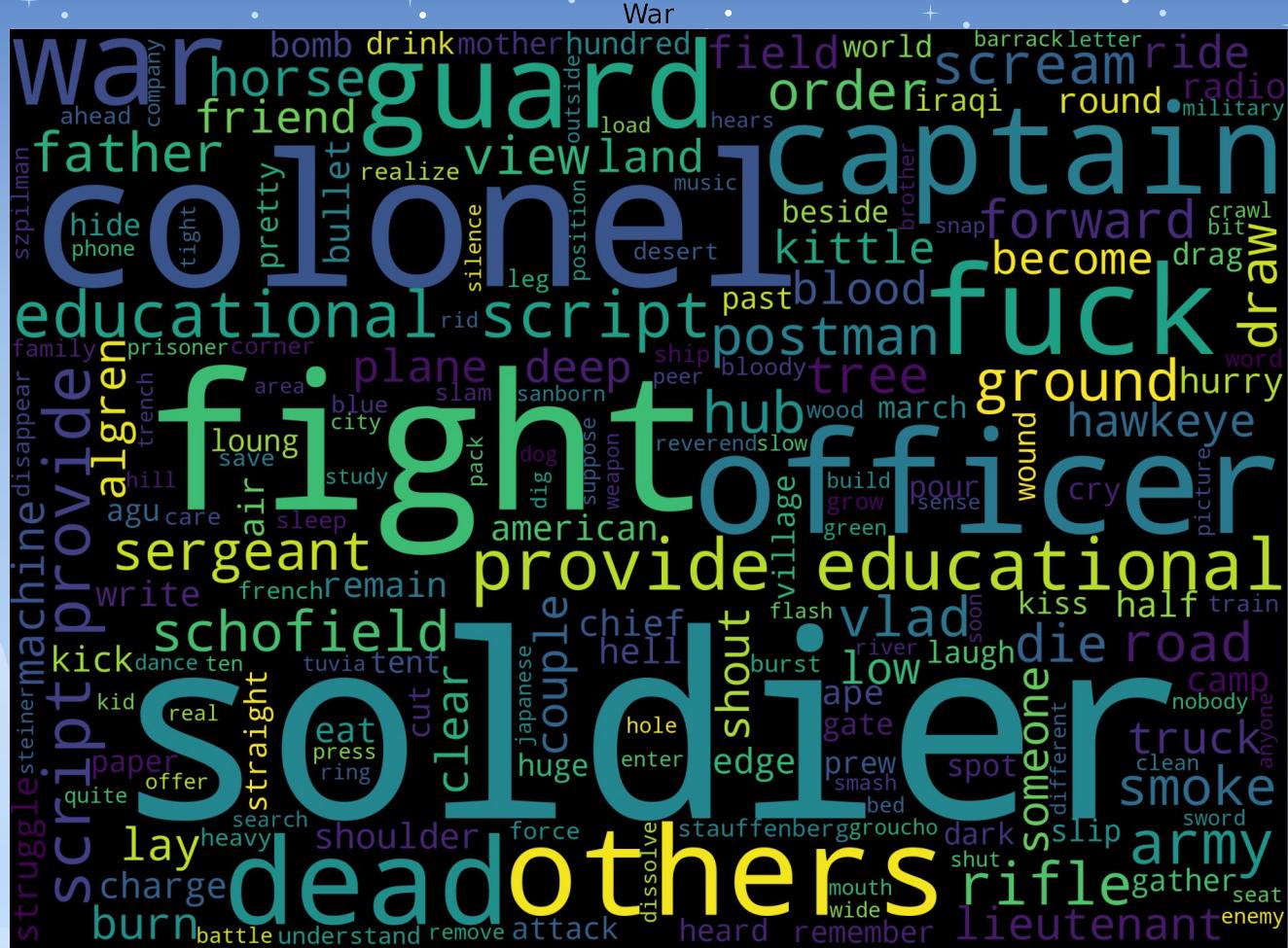
SciFi

SciF

Thriller

Thriller

War



Western

The Machine - Extended Version

Input:

Screenplay
Text



Text Preprocessing
- Lemmatization
- Removal of stopwords

Feature Extraction
- Count
Vectorization
- TFDI

MultiLabel Classifier
- OnevsRestClassifier
- Chain Classifier



Output:

[0,1,0,1,0,0,0,...]

Output

- An array of 18 elements
- Each position of the array is represented by a unique genre
- Elements are either 1 or 0 (True or False)



Example:

Let's say I input a horror, action screenplay into the machine.



The predicted output could be:

[0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]

Action

Thriller

Horror

Metric to Measure Success:

Hamming Loss: The fraction of the wrong labels to the number of labels

Example:

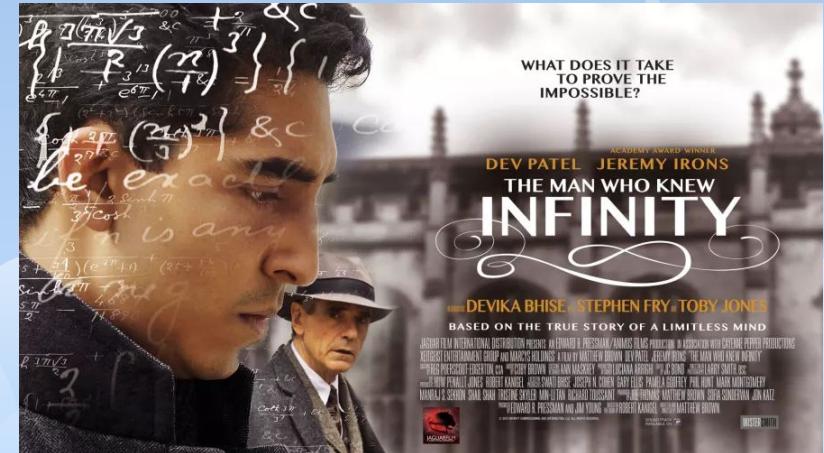
Predicted: [0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0]

Actual: [0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0]

Hamming Loss = 0.055

Goal of the Machine:

To minimize hamming loss; the closer to zero the better



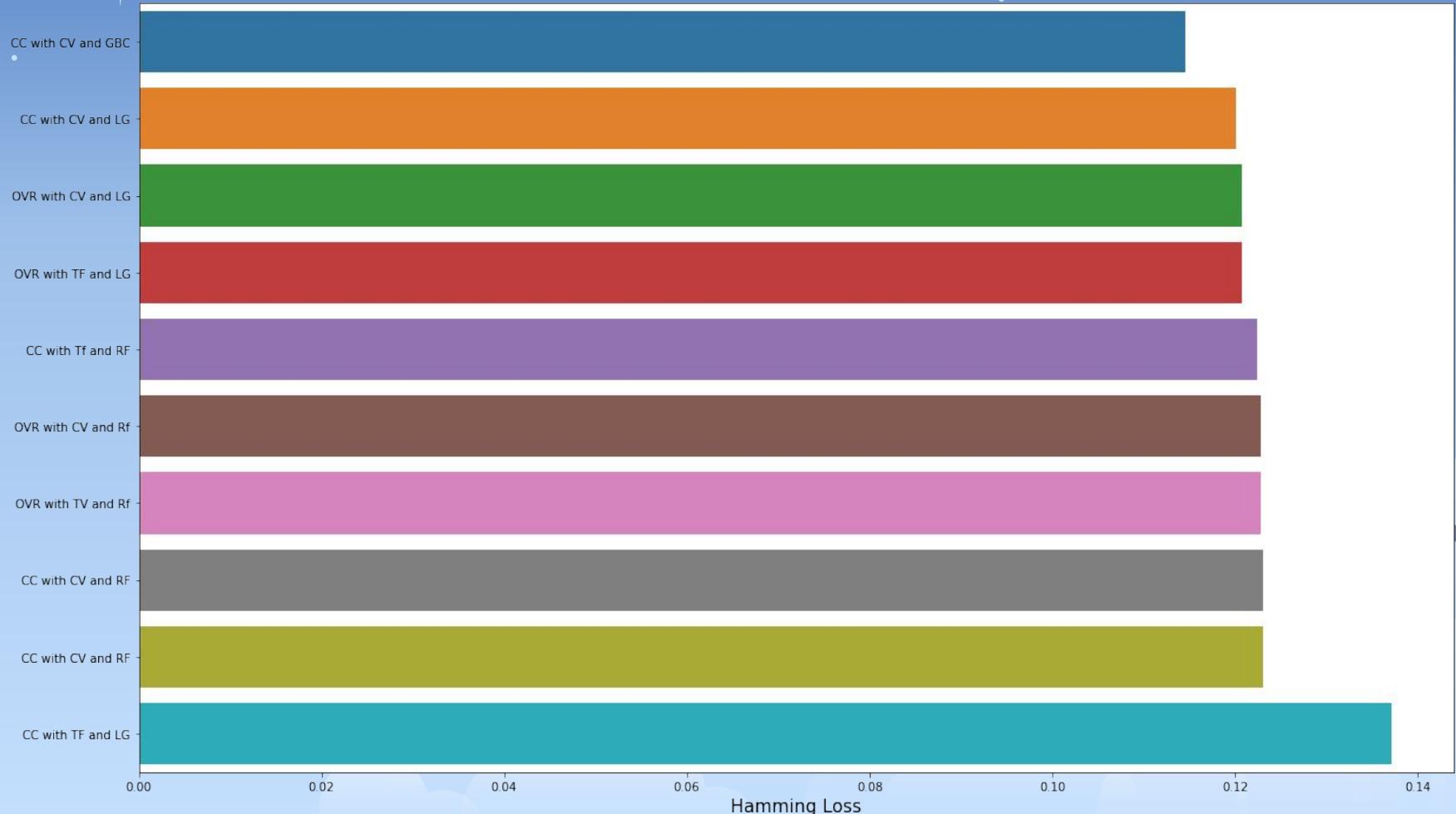
Classification Models Tested:

- Classifier Chain
- OneVsRestClassifier

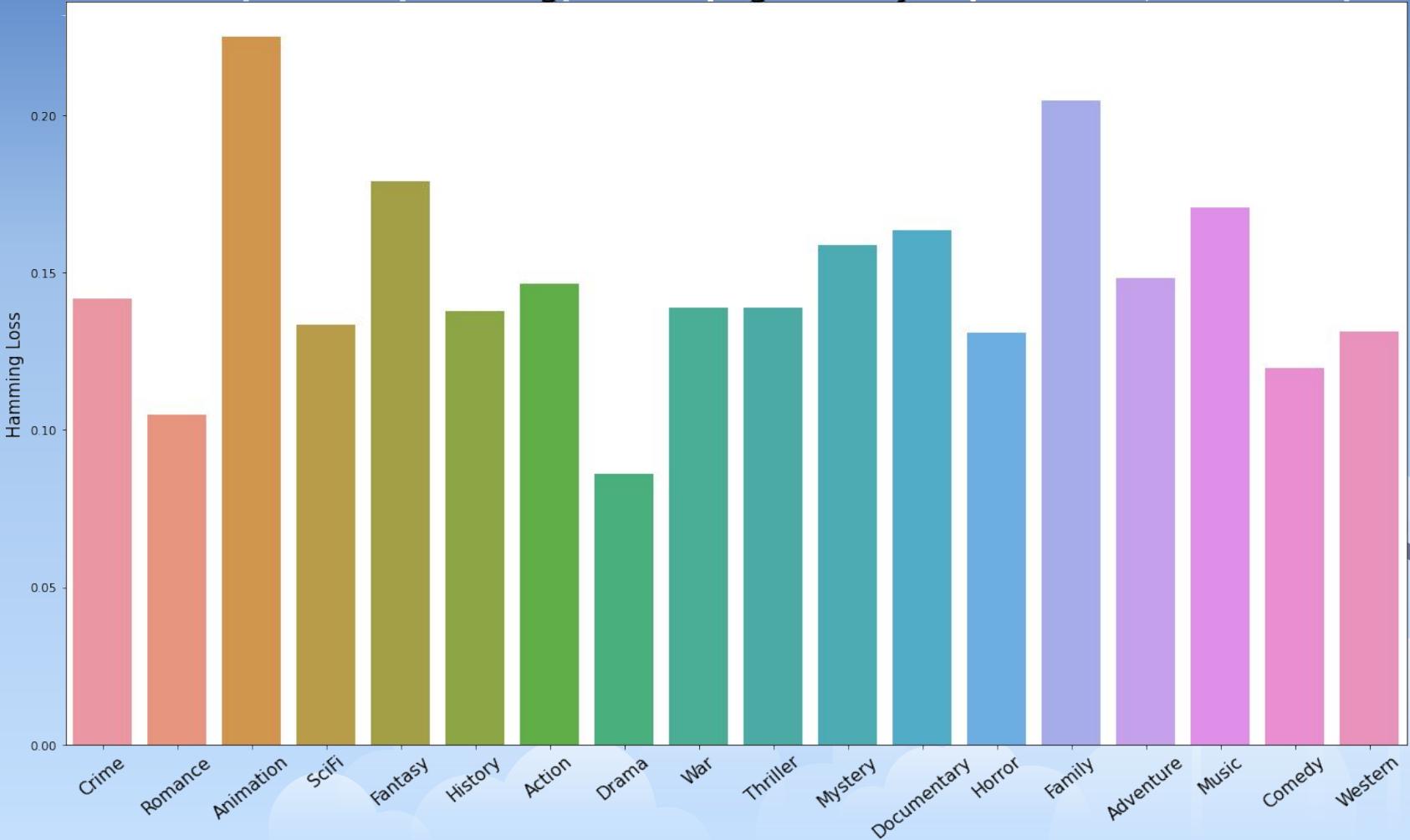
Both **classification models** takes in a classifier as a base estimator (Logistic Regression, RandomForestClassifier, etc) and arranges classifiers into a chain.

Predictions from previous classifiers are used as features for the next classifier

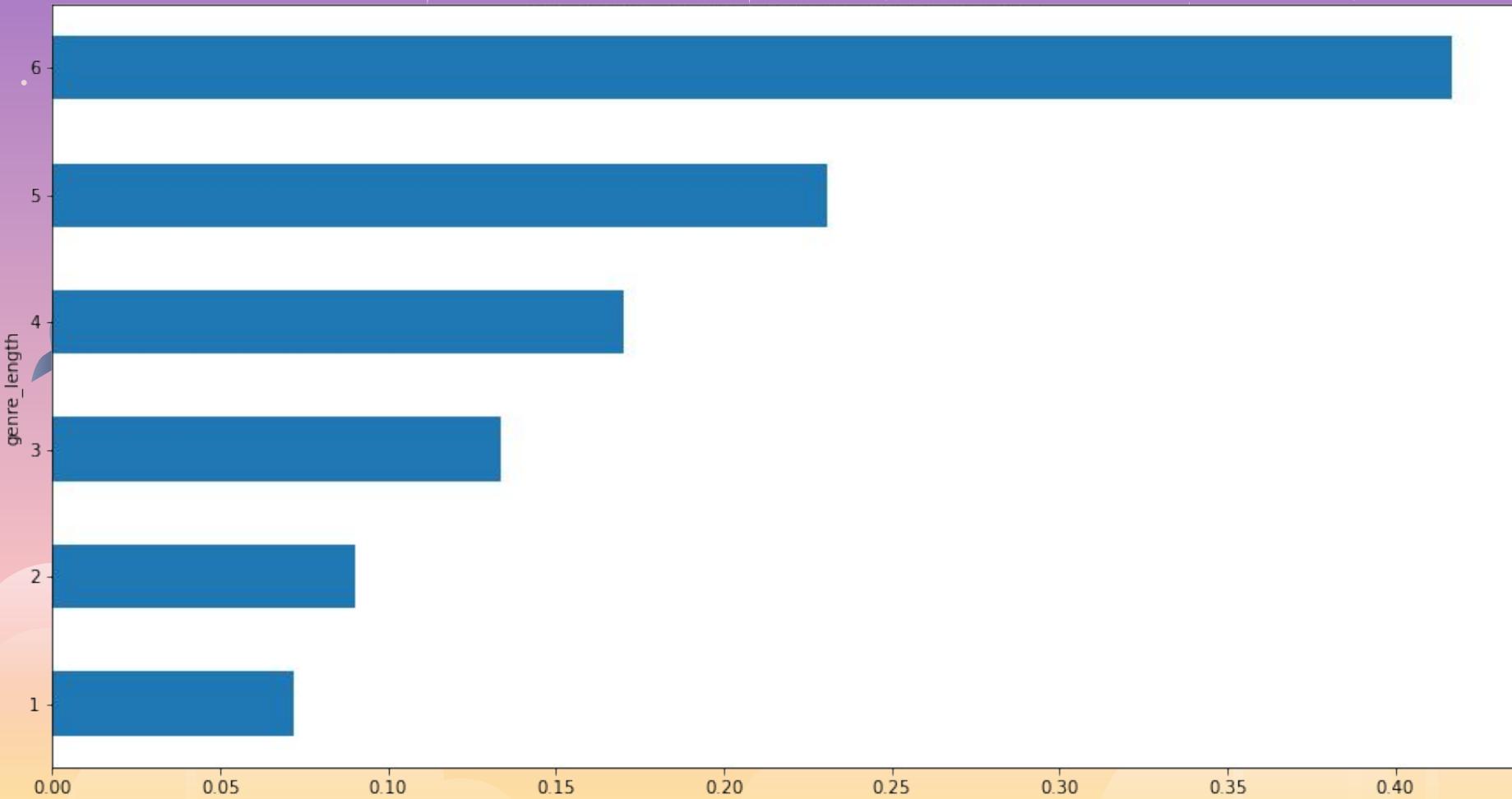
Average Hamming Loss between Model Types



Average Hamming Loss by Genre



Average Hamming Loss by Genre Combination





Movies from the test set ->



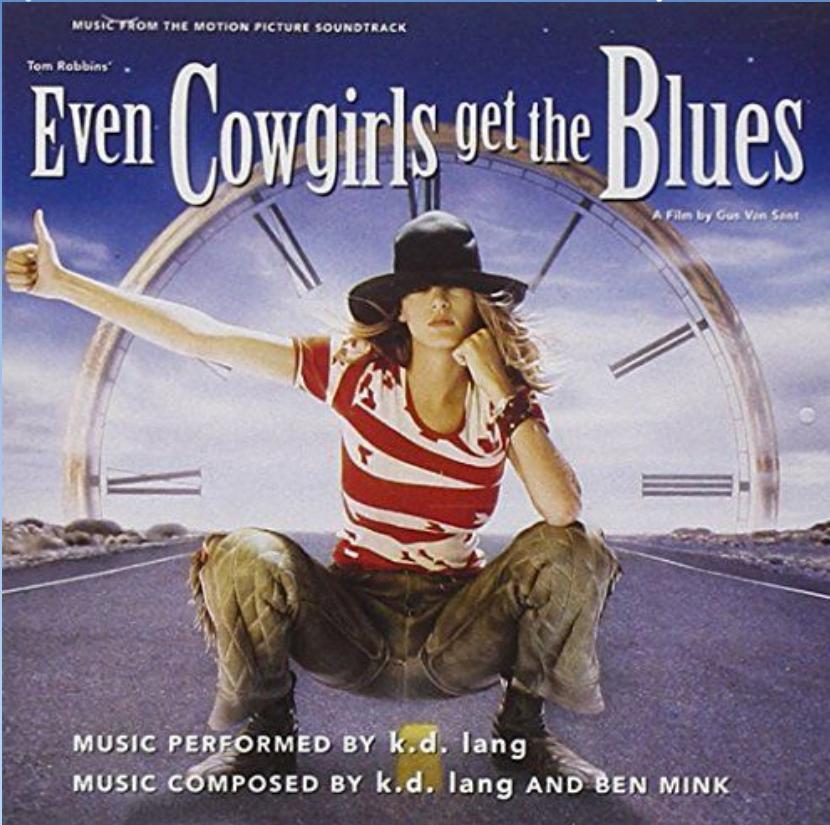
Predicted: SciFi, Action,
Adventure

Actual: SciFi, Adventure

Hamming Loss: 0.055

Tomato Rating: 83%

Box Office: 400 million



Predicted: Drama

Actual: Romance, Drama, Comedy, Western

Hamming Loss: 0.167

Tomato Rating: 18%

Box Office: 1.7 million

A MASTERPIECE
OF MODERN HORROR



A STANLEY KUBRICK FILM

STARRING
JACK NICHOLSON SHELLEY DUVALL "THE SHINING"

WITH
SCATMAN CROTHERS, DANNY LLOYD STEPHEN KING

SCREENPLAY BY
STANLEY KUBRICK & DIANE JOHNSON STANLEY KUBRICK

EXECUTIVE PRODUCER
JAN HARLAN PRODUCED IN ASSOCIATION WITH
THE PRODUCER CIRCLE CO.



From Warner Bros. A Warner Communications Company © MCMXXXI Warner Bros. Inc. All Rights Reserved

Predicted: Comedy

Actual: Thriller, Horror

Hamming Loss: 0.167

Tomato Rating: 85%

Box Office: 47.3 million

Suggestions & Takeaways

- Genres can be identified given a screenplay
- Some screenplays that fit the genre do well with ratings while others don't
 - However, having the right genre doesn't mean high ratings
- Studios can use a genre classifier to pick screenplays

Improvements & Future

- Use other Natural Language Processing Techniques to create new features
 - Topic modeling, etc
- Classify and predict other film elements
 - Theme, Act Structure, etc
- A.I generated screenplays
 - Screenwriters cost a lot of money



Any Questions?

Thank you!



<https://github.com/mikeyo4800>

<https://www.linkedin.com/in/michael-orlando-2835a3149>

