

MOVIE SUGGESTIONS VS TV SUGGESTIONS

AGENDA

→ Introduction

→ Problem Statement

→ Data Collection

→ Modeling

→ Confusion Matrix

→ Conclusion and
Recommendation

INTRO



PROBLEM STATEMENT

As a part of data science team at film and television production company, I was tasked with creating a model that will be able to take a reddit post and classify which posts were talking about movies and which posts were talking about TV show to find current trends so the company can decide to make movies or television show according to consumer demand.



DATA COLLECTION

- r/MovieSuggestions
- 860 posts

Posted by u/Liftingphilosopher 1 day ago

REQUESTING Another film to watch with my dad? I posted a while ago and got some great suggestions

So during lockdown me and my dad have been trying to watch some great films together, spending time before I go back to college.

Here are a few of the greatest ones we have both enjoyed - uncut gems, good time, Ballad of Buster shruggs, 12 years a slave, moonlight, prisoners, the longest yard, Ted, beasts of no nation, Mowgli, The platform, 28 days later, and 28 weeks later, spirited away, knives out.

And a few we watched separately a while back - gone girl, full metal jacket, kill Bill, platoon, Joker, The revenant, Shawshank redemption, Gran Torino, Superbad, shutter island, little miss sunshine, 300, fantastic mr fox.

So a real range of genres and themes, any suggestions at all are welcome!

117 Comments Award Share Save ...

- r/TelevisionSuggestions
- 996 posts



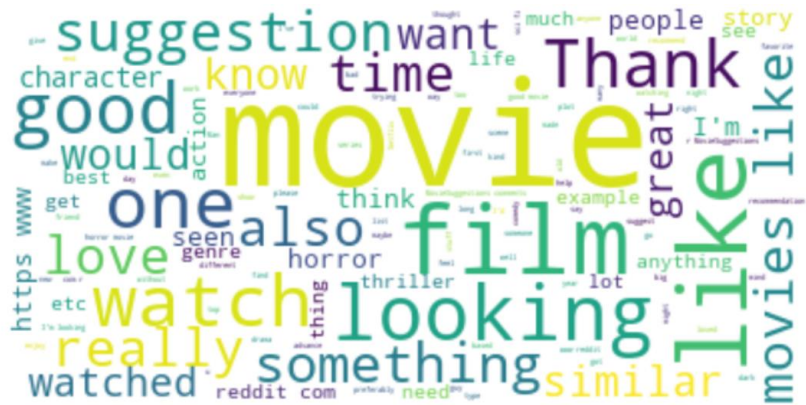
Posted by u/300yearsofexperience 1 day ago

Looking for mystery fantasy like LOST, Lock & Key, October faction, colony, the 100.

As title says, i love those shows that are almost good, but they don't go so slow that they are good, i want something that has some substance and twists with deaths and stuff sure, but most importantly they just move and push forward into mystery. i recently watched the wilds, and it was bad, but close to what I'm looking for, and it got me excited for one of those shows again, just need to find it.

16 Comments Award Share Save ...

Movie Suggestions



Television Suggestions



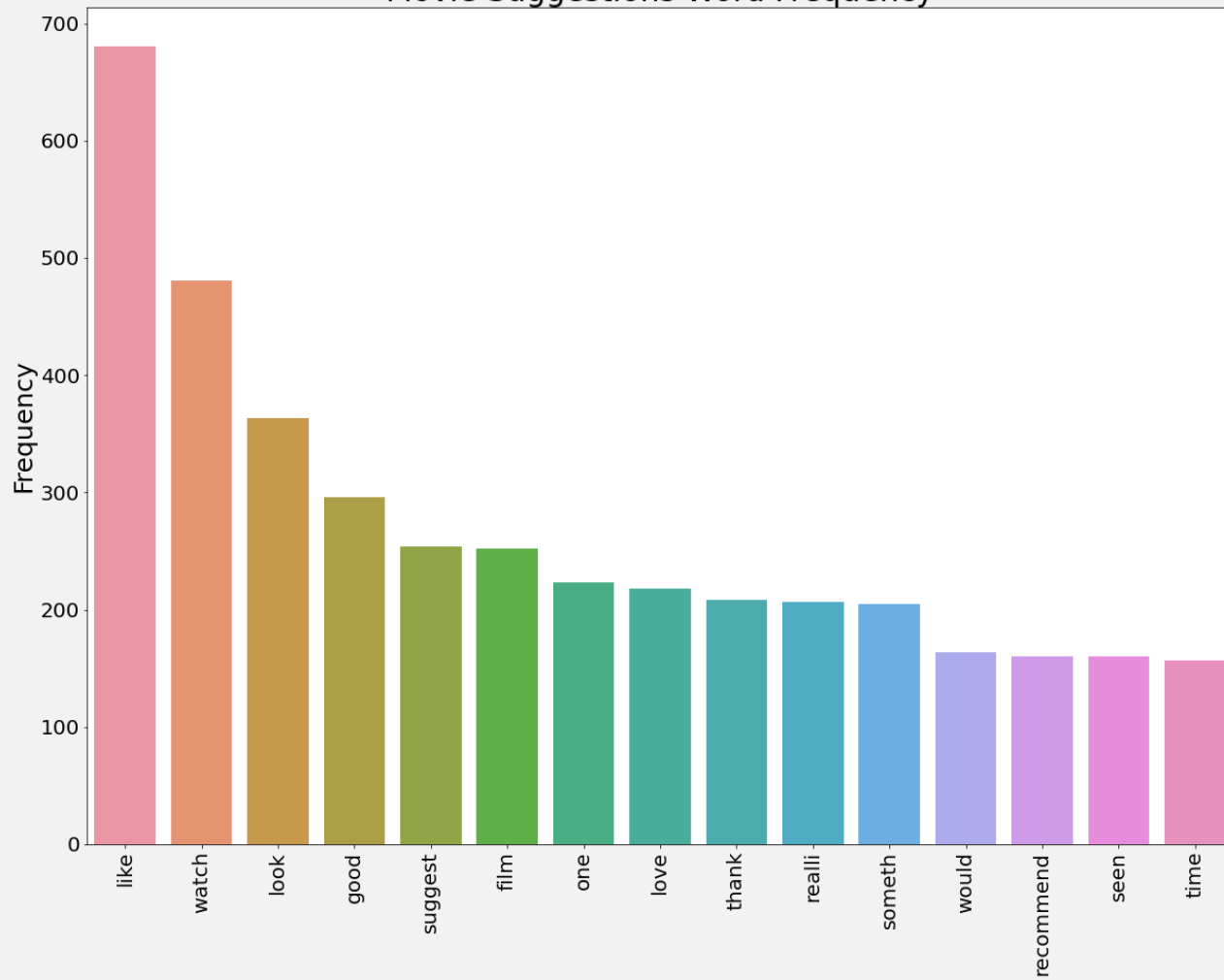
DATA CLEANING

- Remove HTML
- Remove non-letters
- Remove stop words
- Stemming word

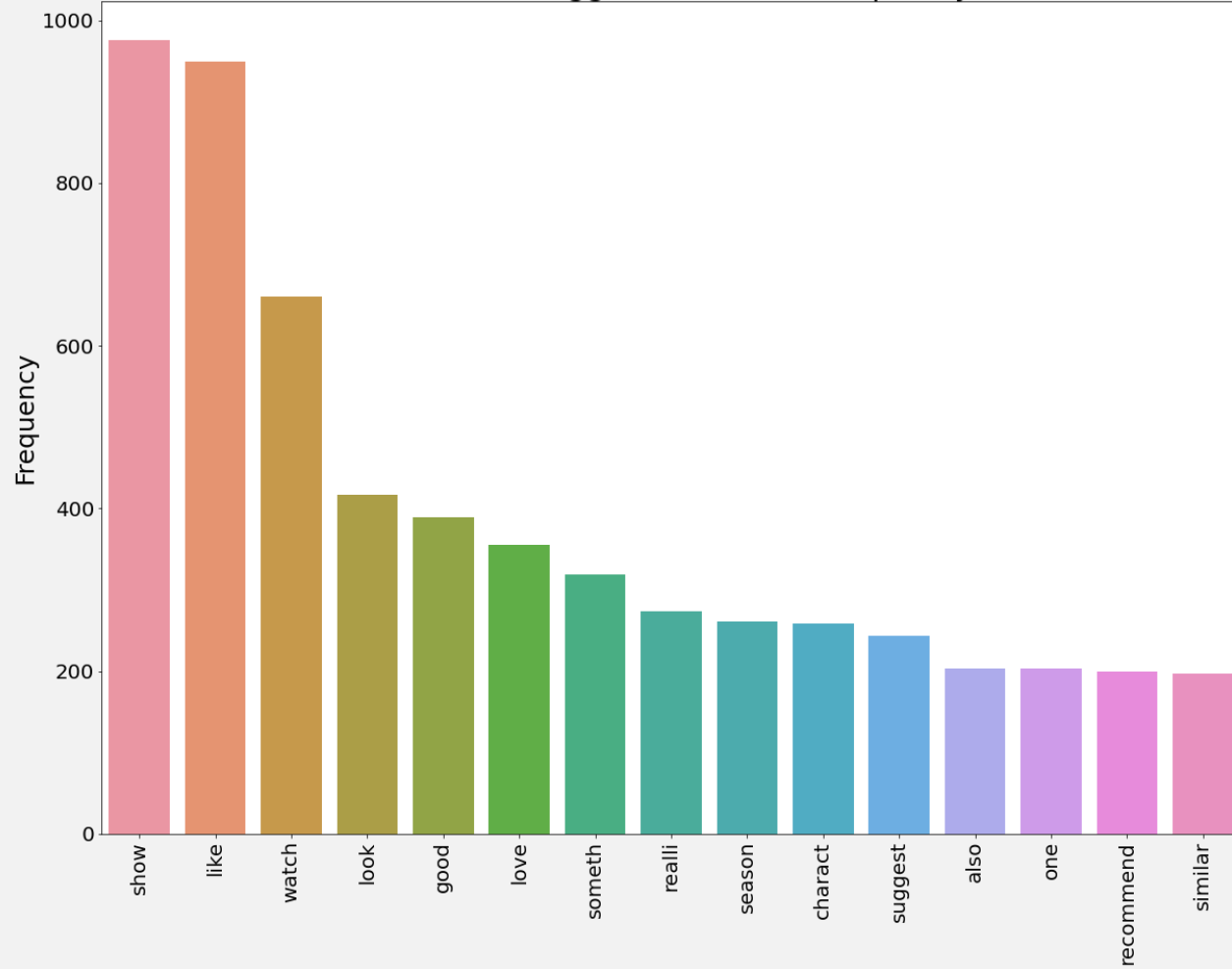
```
In [530]: clean_text(final_df['text'][0])
```

```
Out[530]: 'top februari upi taken comment indic great januari ment ad round thread end januari vote count next month climb climb onto top  
favourit harakiri bruge minari place beyond pine promis young woman solari fall fell top amadeu art self defens brigsbi bear bl  
ood tumbbad wind river name gain biggest gain la hain john wick readi top movement within top rank parasit upon time hollywood  
hereditari midsommar quiet place knife lighthous isl dog joker game night thought top fell climb anyth els'
```

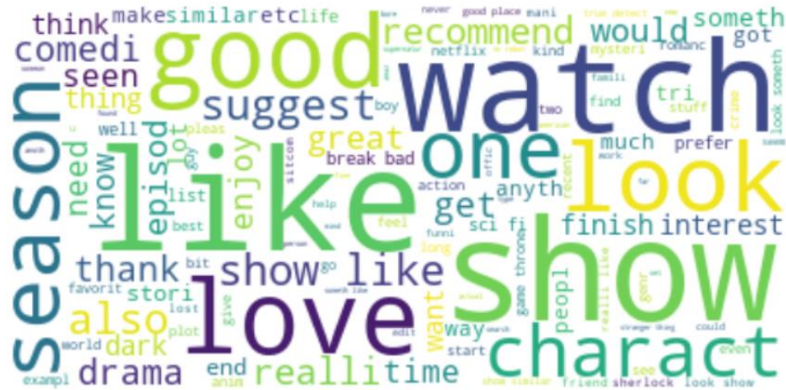
Movie Suggestions Word Frequency



Television Suggestions Word Frequency



Movie Suggestions



Television Suggestions



BASELINE MODEL

```
In [219]: X = final_df['text']  
          y = final_df['y']  
          X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 42, stratify = y)
```

```
In [223]: # Baseline accuracy  
          y_train.value_counts(normalize = True)
```

```
Out[223]: 0    0.536638  
          1    0.463362  
          Name: y, dtype: float64
```

TEXT FEATURE EXTRACTION

```
In [225]: cvec = CountVectorizer()  
          cvec.fit(X_train)  
          X_train = cvec.transform(X_train)  
          X_test = cvec.transform(X_test)
```

```
In [221]: X_train.shape
```

```
Out[221]: (1392,)
```

```
In [222]: X_test.shape
```

```
Out[222]: (464,)
```

```
In [337]: X_train.shape
```

```
Out[337]: (1392, 6494)
```

```
In [338]: X_test.shape
```

```
Out[338]: (464, 6494)
```

MODELING

Model	Train Accuracy	Validation Accuracy	Test Accuracy
Logistic Regression	99%	82%	82%
Multinomial Naïve Bayes	96%	85%	83%
Support Vector Machine	96%	83%	85%

GRID SEARCH

Model	Train Accuracy	Best Score	Test Accuracy
Logistic Regression	96%	83%	83%
Multinomial Naïve Bayes	94%	85%	85%
Support Vector Machine	95%	85%	83%

BEST PREDICTION MODEL

```
In [339]: X = final_df['text']  
y = final_df['y']  
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 42, stratify = y)
```

```
In [340]: cvec = CountVectorizer(max_df = 0.9, max_features = 3000, ngram_range = (1,2))  
cvec.fit(X_train, y_train)  
X_train = cvec.transform(X_train)  
X_test = cvec.transform(X_test)
```

```
In [341]: nb = MultinomialNB(alpha = 1)  
nb.fit(X_train, y_train)  
print(f'Train Accuracy Rate = {nb.score(X_train, y_train)} \nTest Accuracy Rate = {nb.score(X_test, y_test)}')
```

```
Train Accuracy Rate = 0.9425287356321839  
Test Accuracy Rate = 0.8577586206896551
```

```
In [337]: X_train.shape
```

```
Out[337]: (1392, 6494)
```

```
In [338]: X_test.shape
```

```
Out[338]: (464, 6494)
```

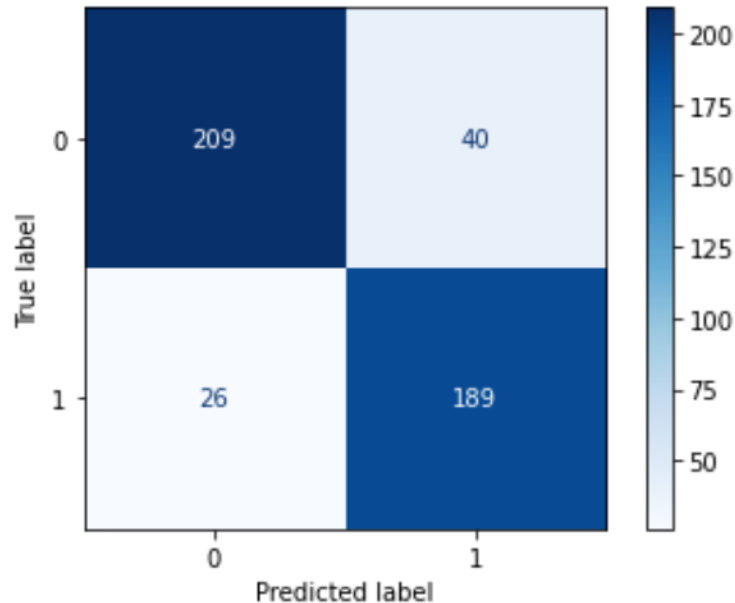
```
In [342]: X_train.shape
```

```
Out[342]: (1392, 3000)
```

```
In [343]: X_test.shape
```

```
Out[343]: (464, 3000)
```

CONFUSION MATRIX



True Positive = 189

True Negative = 209

False Positive: = 40

False Negative = 26

Accuracy = 86%

Misclassification Rate = 14%

Sensitivity = 88%

Specificity = 84%

Precision = 83%

CONCLUSION AND RECOMMENDATION

- Based on train and test accuracy rate multinomial naïve bayes is the best model among the three models.
- This model has accuracy rate at 86%
- Sensitivity or true positive rate at 88%
- Specificity or true negative rate at 84%
- Try another model like Random forest, Adaboost classifier
- Try another feature extraction like TF-IDF or Hashing vectorizer.