

Or



Using Natural language Processing to distinguish reddit posts about two acclaimed video games

# The Problem:

- Consider two of the most successful and critically acclaimed video game franchises of all time: God of War and The Legend of Zelda
- Both have massive cult followings and have many reddit threads devoted to them.
- There are dozens of subreddits and thousands of posts about two very different games. Let's use NLP to tell which is which.

# Obtaining our Data: Scrape Reddit

url, name, n requests

Our method takes the desired subreddit url, name for outputted .csv, number of API requests. Headers and included in the body and after is initialized as a None type.

API requests

Request JSON and store desired properties in temporary dictionary. Retrieve subreddit name, post title, paragraph body, score, name, link flair text. After each loop delay by one second.

Output to csv

Finally convert post dictionaries into a data frame and give appropriate column names.

# About our data



951 posts collected

post traits scraped	
0	subreddit
1	title
2	score
3	post paragraph
4	link_flair_text



725 posts collected.

# The NLP process

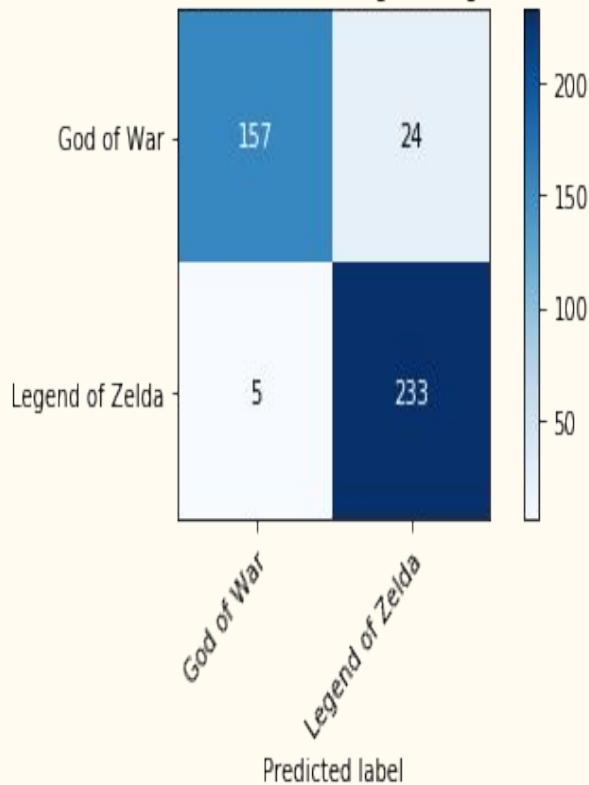
1. Gathered the Data (scrapped reddit's API)
2. Cleaned, Formatted, and Vectorized
3. Selected relevant features and target
4. Train-Test-Split
5. Ran through four models via Gridsearch
6. Analyzed different models performance

# Evaluating Model Performance

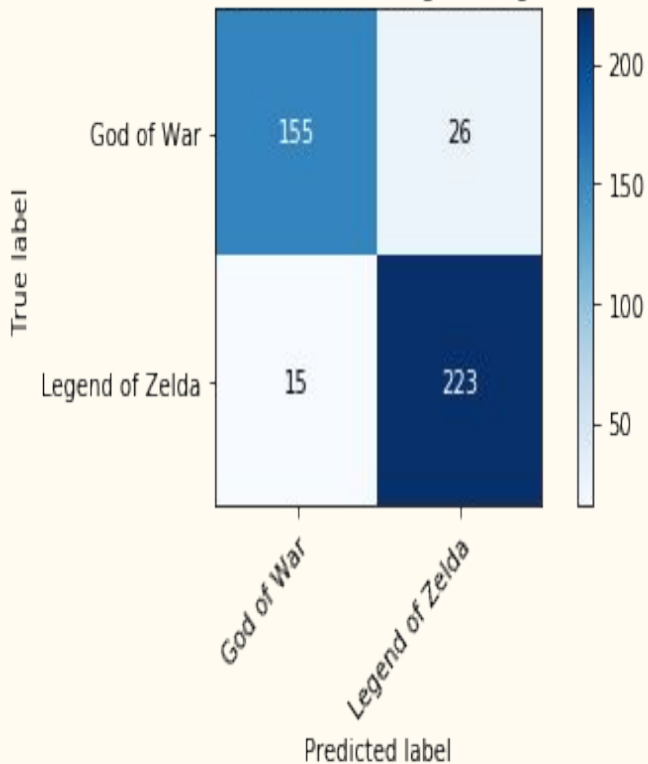
<b>Model:</b>	LogisticRegression	<b>Model:</b>	Random Forrest	<b>Model:</b>	Naive Bayes Gauss
<b>Vectorization:</b>	CountVectorize	<b>Vectorization:</b>	CountVectorize	<b>Vectorization:</b>	CountVectorize
<b>CV Folds</b>	5	<b>CV Folds</b>	5	<b>CV Folds</b>	N/A
<b>Parameter 1</b>	C = 1.3	<b>Parameter 1</b>	n estimators = 106	<b>Parameter 1</b>	N/A
<b>Parameter 2</b>	Penalty = 12	<b>Parameter 2</b>	min sampl split = 2	<b>Parameter 2</b>	N/A
<b>Train Set Score</b>	1.0	<b>Train Set Score</b>	0.9992	<b>Train Set Score</b>	0.9976
<b>Test Set Score</b>	0.93079	<b>Test Set Score</b>	0.9141	<b>Test Set Score</b>	0.9117

# Classification Confusion Matrices

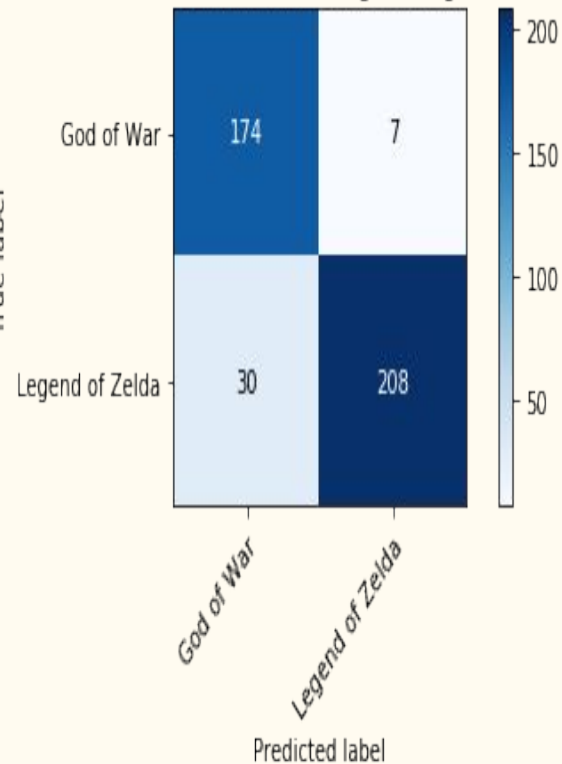
Confusion matrix for LogisticRegression



Confusion matrix for LogisticRegression



Confusion matrix for LogisticRegression



## Recommendations

1. Use the model r/GodofWar and r/truezelda.
2. Try term frequency-inverse document frequency
3. Try the model on different subreddit pairs