

# Subreddit Classification



By: Royce Phillips Jr

# Problem Statement

Reddit is a large platform with millions of submissions and hundreds if not thousands of subreddits. Each subreddit has its own topics and personality.

Is it possible to create a classification model that can decipher between two separate subreddits.

# Which subreddits?

- 'CasualConversation' 1,182,804 users
- 'SeriousConversation'. 34,860 users

Why?

- Sister subreddits?
- Any similar topics?
- Are there any major differences



Image: <https://www.flickr.com/photos/djandywdotcom/31437486496>

# What data??

Used Pushshift's API to pull data

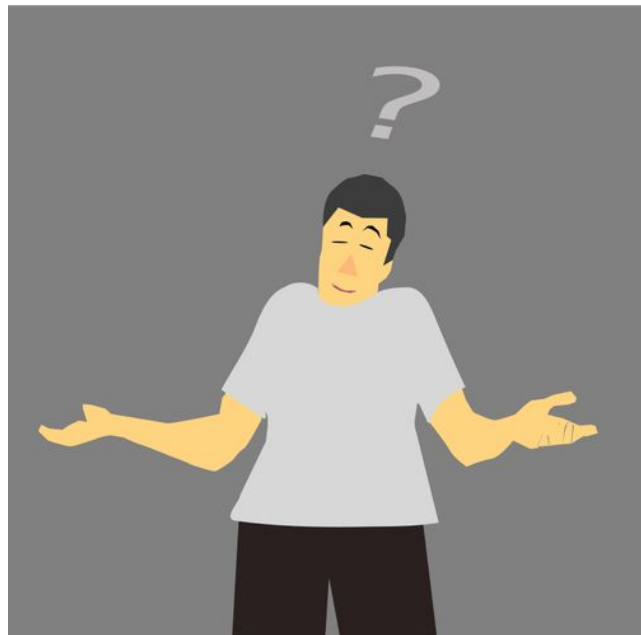


Image: <https://freesvg.org/shrug-gesture>

# Initial cleaning and EDA

- Many null values
- Many submissions had been 'removed'
  - Mostly in 'CasualConversation'
- Unbalanced subreddits

## Initial thoughts

- 'SeriousConversation' more quality?
- 'CasualConversation' too many people posting?

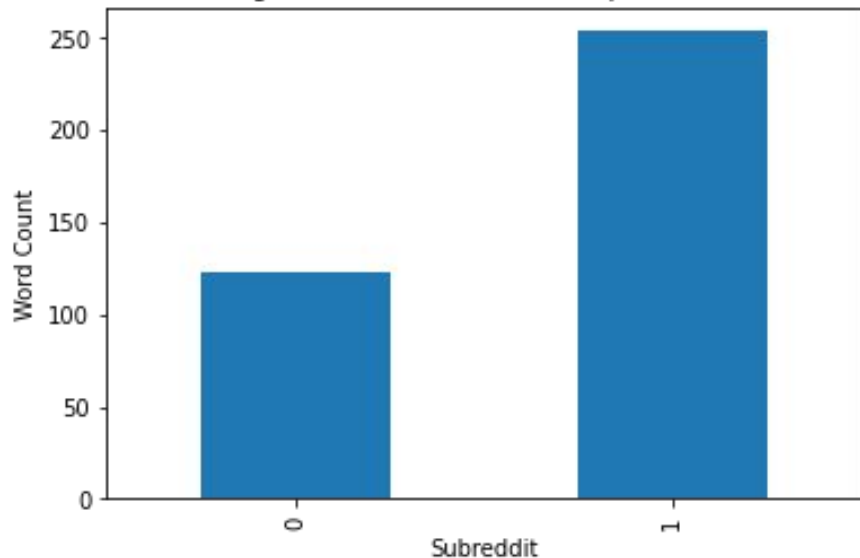
## Some quick facts

	Number of Comments	Score	Title Word Count	SelfText Word Count
Casual Conversation	13	14	10	123
Serious Conversation	10	9	12	254

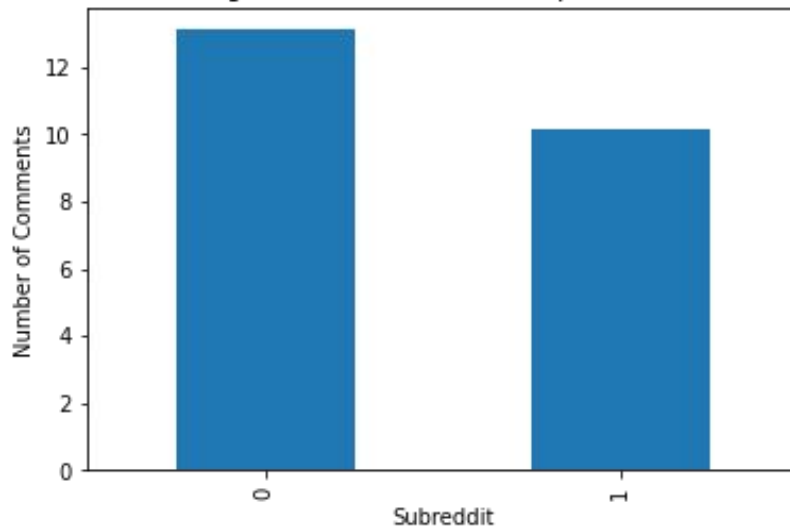
# Comparing Subreddits

0 = 'CasualConversation'  
1 = 'SeriousConversation'

Average Word Count in Selftext per Subreddit



Average Number of Comments per Subreddit



# Which gets more interaction?

	More than 100 Comments	More than 250 Comments
CasualConversation	135 ~ 76.3%	61 ~ 98.4%
SeriousConversation	42 ~ 23.7%	1 ~ 1.6



# Most Commented Post

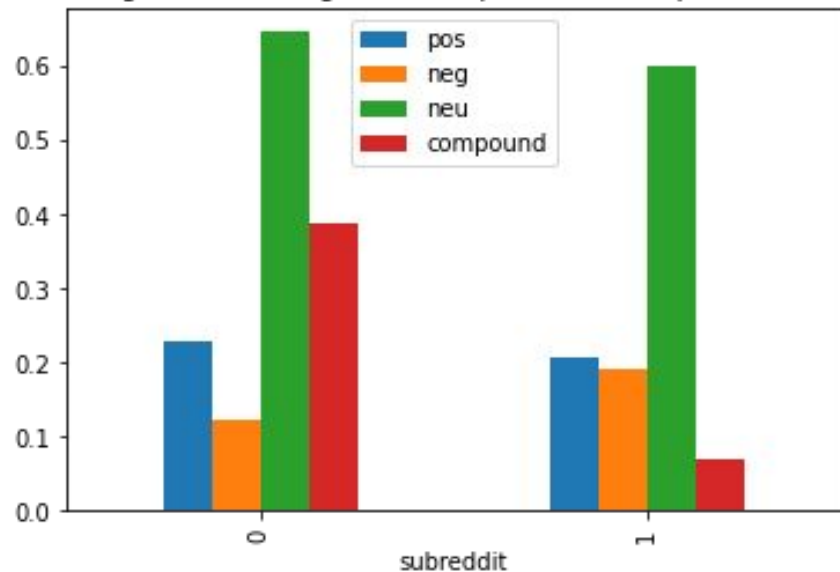
## “CasualConversation”

- 3001 Comments

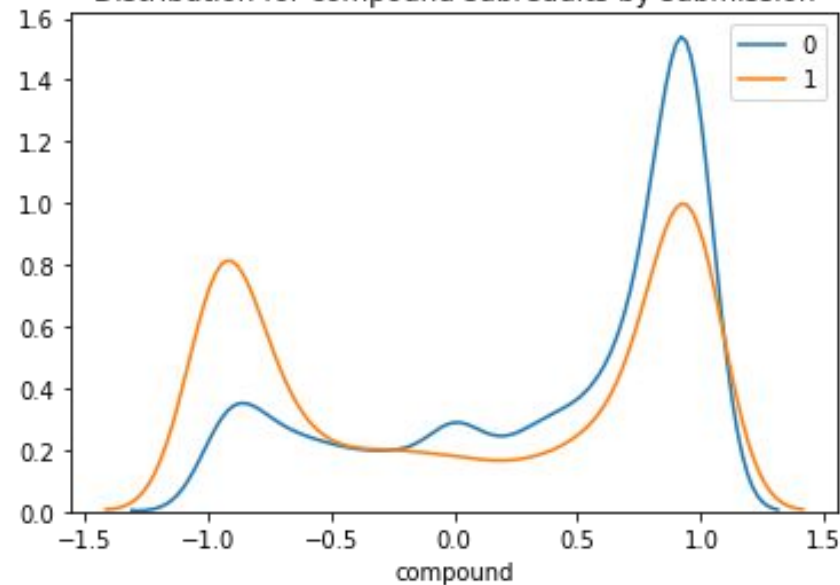
'This coronavirus things has made me realize people would be a lot happier and explore their passions and interests if they didn't have to work so much.'

# Sentiment Analysis

Average Positive/Negative/Compound Scores per Subreddit



Distribution for compound subreddits by submission



# How about models?

- Gridsearch
- Pipelines
- Linear Regression
- Naive Bayes Multinomial

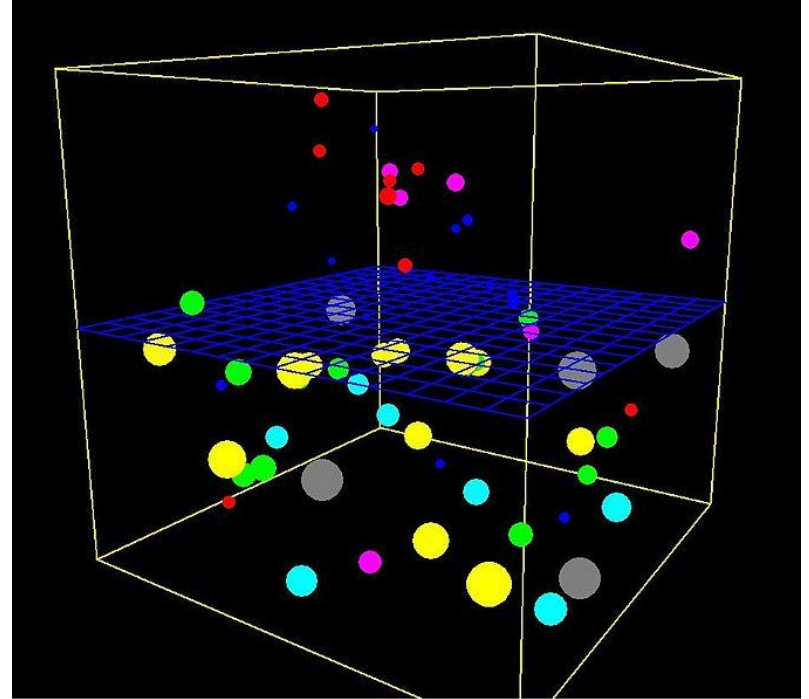


Image: [https://commons.wikimedia.org/wiki/File:Osmosis\\_computer\\_simulation.jpg](https://commons.wikimedia.org/wiki/File:Osmosis_computer_simulation.jpg)

# Top 20 words

## SeriousConversation

	0
<b>religious</b>	2.371477
<b>sexual</b>	2.197139
<b>death</b>	1.962149
<b>died</b>	1.853095
<b>sexually</b>	1.686206
<b>politics</b>	1.667395
<b>trauma</b>	1.660829
<b>cancer</b>	1.633291
<b>drugs</b>	1.562543
<b>dying</b>	1.441792
<b>political</b>	1.379804
<b>sex</b>	1.366498
<b>dead</b>	1.270756
<b>religion</b>	1.240917
<b>discuss</b>	1.192002
<b>much time</b>	1.135842
<b>climate</b>	1.121131
<b>illness</b>	1.077652
<b>passed away</b>	1.058376
<b>porn</b>	1.032302

## CasualConversation

	0
<b>quarantine</b>	-2.026719
<b>wanted share</b>	-1.709266
<b>corona</b>	-1.599827
<b>isolation</b>	-1.434727
<b>lockdown</b>	-1.287376
<b>caring</b>	-1.267222
<b>cake</b>	-1.224584
<b>stay home</b>	-1.188175
<b>social distancing</b>	-1.171571
<b>bought</b>	-1.133384
<b>haha</b>	-1.109965
<b>applied</b>	-1.078960
<b>tldr</b>	-1.061328
<b>series</b>	-1.052097
<b>working home</b>	-1.051667
<b>chatting</b>	-1.022924
<b>coronavirus</b>	-0.991696
<b>virus</b>	-0.977520
<b>holding</b>	-0.971848
<b>covid</b>	-0.961733
<b>thankful</b>	-0.950105

# Ideas to improve model

- Gather more data
- Try different models
  - SVM, KNN, etc.
- See how TFIDF could work



Image:

<https://commons.wikimedia.org/wiki/File:Performance-Evaluation-Process-z.jpg>

# Conclusion

- 80% success rate.
- Amount of words in selftext
- Sentiment analysis
- Number of comments
- Topics
  - 'SeriousConversations' seemed to be more personal.
  - 'CasualConversations' were more revolved around the topic of COVID-19.

## 2nd Most commented : 1814 Comments

'I just ate rich people pasta for the first time'

"I'm about to graduate and get my first full time job in a couple of months. Meanwhile, I'm still living the poor plebejian life. I always get the cheapest 39 cents no name pasta. However, due to all the panic buying, almost all pasta was sold out. Except for the expensive brand pasta, which usually costs 1.43€, but it was reduced to 77 cents for some reason. It's in a fancy carton package with a little see through window. The cheap pasta I get is just a plastic bag. So I got the expensive pasta, but I still bought the cheapest pasta sauce for 79 cents. I cooked the pasta 15 minutes ago. What the actual f@%\$. I cooked it, and it's still thin. But it's not hard. It's stretchy, but it doesn't rip. You don't have to drown it in sauce until it's edible, no, it tastes good even without the sauce. I mixed it with the sauce. What is this. I didn't know pasta could have physical property like that. I can't go back to 39 cents pasta. Now imagine you're rich and this x100 is your average meal."