

# Analyzing Reddit Comment Reception

Aadi Sudan

2022-12-09

## Overview

The goal of this project was to try and determine any connections between sarcasm on the score of a Reddit comment. Included with this report is a .csv file made up of a collection of comments, some sarcastic and some not. Included in the file is the score for each comment and the Subreddit they were posted in, two important factors to take into consideration. The sarcastic comments act as the study group in this experiment, while the nonsarcastic comments act as a baseline for what more typical comment reception looks like. The parent Subreddit is included to keep in mind any possible lurking variables in the analysis, as certain Subreddits may be more prone to a certain type of response.

Together in this report, we will discuss the results that came from analyzing this .csv file to see what the correlation between the nature of a comment and its score is, if any correlation exists at all.

## Introduction/Background Information

I chose to pursue this topic for my project because I use Reddit quite often, and in my age and gender demographic, the site is quite popular. Many of my friends and other people my age that I know don't enjoy using other social media platforms, but Reddit has a certain charm to it for being both simplistic, straight to the point, and lacking the usual pitfalls that come with certain other platforms. Scores on posts have always perplexed me, so I was curious to find a pattern or consistent correlation in the seemingly random formula of users reading a certain type of comment and responding with an upvote or downvote.

For those who may not know, Reddit is a popular social media site where users can post threads to Subreddits - specific groups or fandoms designed to focus in on a particular topic - and other users can leave comments. Comments on Reddit can be upvoted (typically for showing agreement or approval), or downvoted (typically for showing disapproval of the nature or wording of the comment). The 'score', then, is essentially the total number of upvotes minus the total number of downvotes.

Certain posts are more sarcastic in nature, and these posts are typically denoted with the use of a /s at the end of the message. An important thing to note is that not every user uses the /s tag, leading to a few sarcastic posts that become more difficult to track down and categorize. However, the vast majority of the Reddit userbase tends to use this tag to denote when they wish to be sarcastic, so these unfortunate outliers are thankfully few and far between.

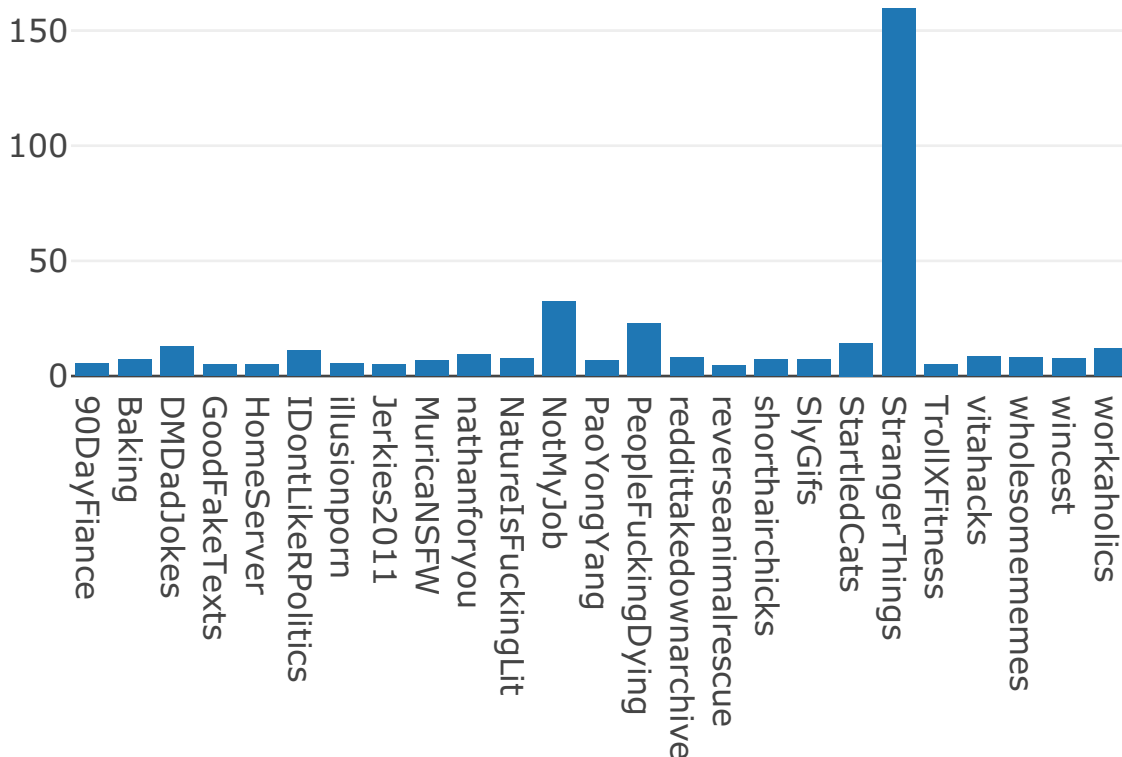
The dataset used in this project was found on Kaggle and consists of several thousands of comments, but for the purposes of this analysis we will only be looking at the top and bottom 25 sarcastic and nonsarcastic comments, as analyzing literal thousands, almost millions, of individual entries is a bit unfeasable for the scope of this project.

The dataset is included in the content portfolio for this project alongside this report and the source code used to manipulate the .csv files and generate the graphs. Now, let's take a look at said graphs and see if we can draw any conclusions from them.

## Sarcastic Comment Analysis

Let's start with sarcastic comments. Below is a graph of the 25 highest scoring Subreddits taken from our dataset. The score here was determined by taking every single entry for a particular Subreddit and finding the weighted average, effectively taking both the number of entries and average score into consideration to craft an accurate comparison to other Subreddits. A reminder that both datasets are simple random samples and were taken with minimal bias, so for the purposes of this project we should be able to rely on them to safely draw our conclusions.

### Top 25 Scores for Sarcastic Reddit Comments

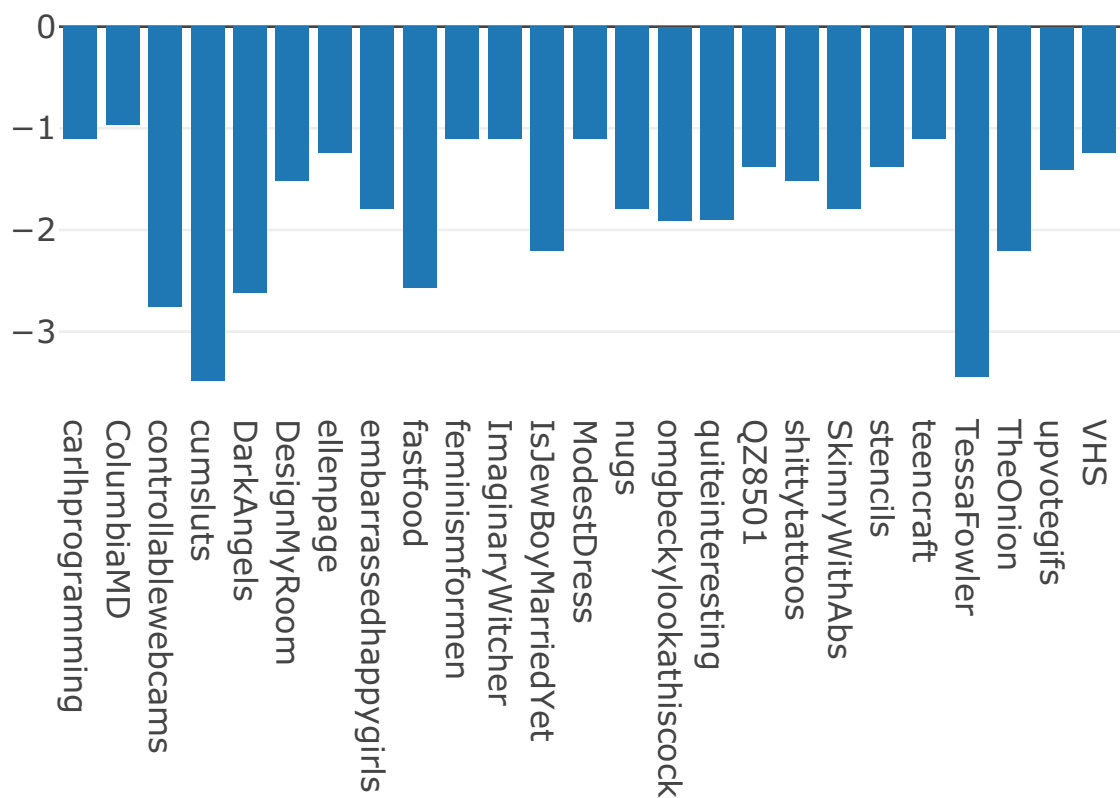


As we can see from this plot, the highest scoring Subreddit by a very large margin was r/StrangerThings with a weighted score of 159.48, an average score of 52.11, and a remarkable 18 entries. While this Subreddit is an extreme outlier, it does tell us a bit more about the nature of r/StrangerThings. Having 18 entries in a simple random sample and still reaching that high means that the Sub is both very active and probably heavier on the sarcastic/potentially comedic side, which isn't all that surprising considering the popularity of the show that it's based on.

Barring this outlier, the majority of the other Subreddits in the top 25 share a relatively consistent weighted score between 4 to 8, with a few local maxima reaching up to 12. Two other highs are r/NotMyJob at 32.25 (average score of 65.25 with 8 entries) and r/PeopleF-ingDying at 23.007 (average score of 55.86 with 7 entries).

Now we're gonna take a quick look at the 25 *lowest* ranked Subreddits.

# Bottom 25 Scores for Sarcastic Reddit Comments

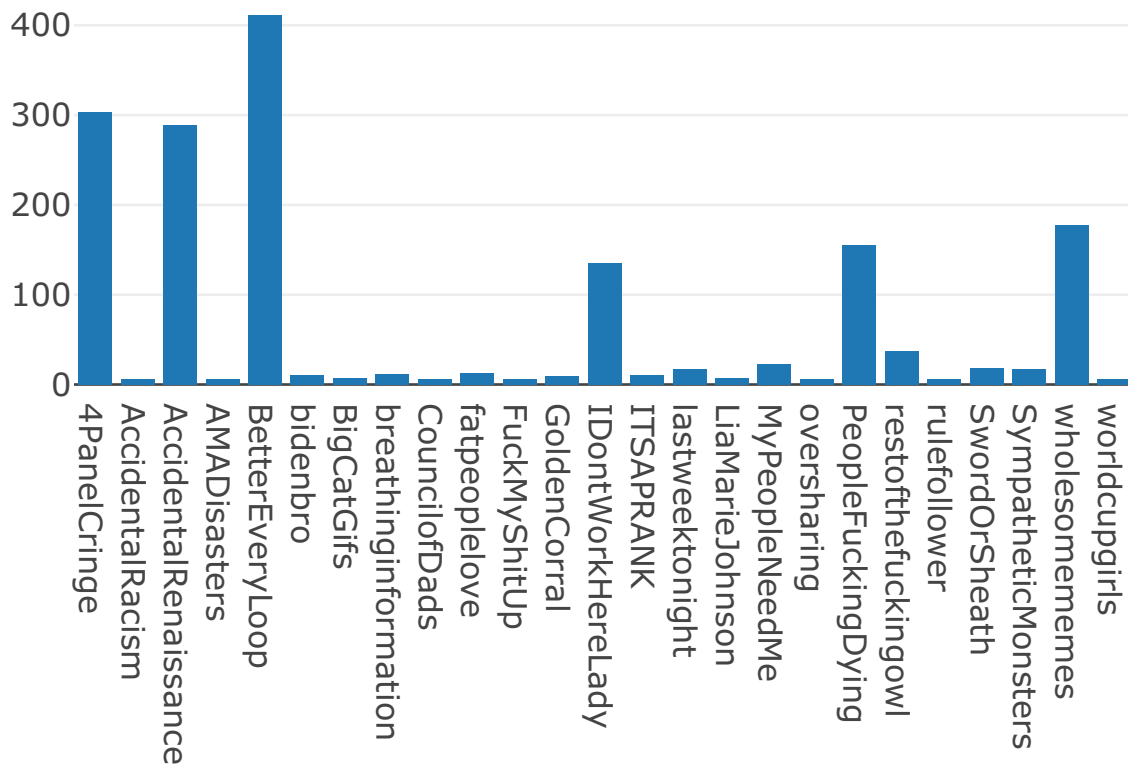


These scores drop into the negatives due to a greater number of downvotes than upvotes, the lowest being r/TessaFowler with a weighted score of -3.45, an average score of -25, and 1 single entry. Apart from that high, which may very well be an outlier due to having only 1 entry, the scores range from -1 to -3.5.

From putting these two datasets together, we can see the range of sarcastic comments in this sample to reach up to 12 (excluding the outlier of r/StrangerThings), and all the way down to -3.5. While a minimum of -3.5 is surprisingly higher than expected, a maximum of 12 is pretty low by Reddit standards, a fact that the next set of plots will highlight in great detail as we shift over to the analysis of nonsarcastic comments.

## Nonsarcastic Comment Analysis

## Top 25 Scores for Nonsarcastic Reddit Comments

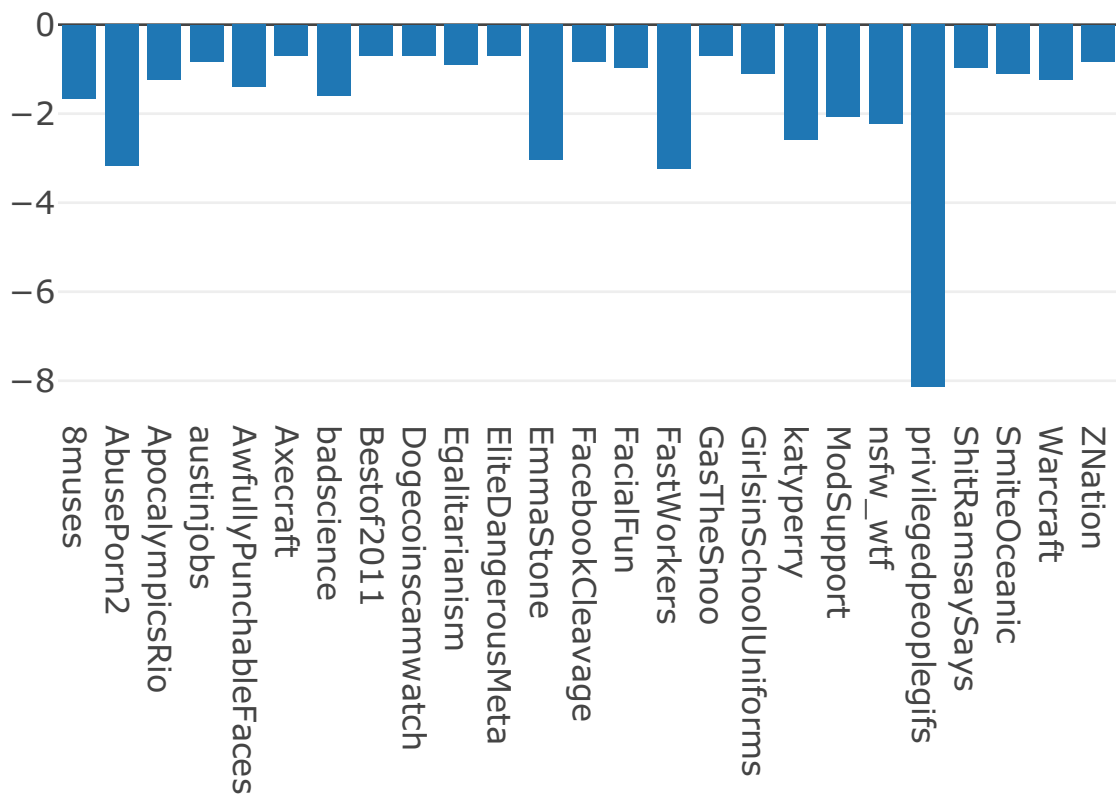


The nonsarcastic comment score distribution is far more diverse, as one might expect. Sarcasm is a variable that we can use to at least somewhat predict the shape of a graph and explain why it looks the way it does, but now that we’ve removed that, we’re in more unpredictable waters.

The three highest Subreddits by far in this sample are r/4PanelCringe with a weighted score of 303.47, average score of 66, and 29 entries, r/AccidentalRenaissance with a weighted score of 289.31, average score of 62.92, and 25 entries, and r/BetterEveryLoop with a weighted score of 411.855, average score of 89.57, and 21 entries. Outside these three local maxima, the rest of the top 25 hold a consistent range of values ranging from 5 to 22, with several scores reaching the 30s and a couple in the 130-150 range.

Despite this rather radical spread of values, on average the nonsarcastic scores remained consistently higher than the sarcastic scores by a wide margin. From these statistics, it would seem that users are more likely to upvote a comment when it is more of an objective statement or serious response to the main thread as opposed to a sarcastic joke. We can investigate that further as we take a look at the bottom 25 nonsarcastic scores.

## Bottom 25 Scores for Nonsarcastic Reddit Comments



The lowest scored Sub ended up being r/privilegedpeoplegifs with a weighted score of -8.13, an average score of -59, and only one entry. It being only one entry means that it's more than likely an outlier, so ignoring that, the rest of the Subs in this bottom 25 list hold scores that range from -0.6 to -3.17, a relatively similar if not slightly higher range than the collection of the lowest sarcastic scores.

By putting this and the highest nonsarcastic score distribution together, we can determine a range for these nonsarcastic scores to be anywhere from a remarkable 150 down to -3.17. Obviously, the high of 150 is exponentially greater than the high of 12 in the sample of sarcastic comments, but what's interesting to note is the number of low scores in both of these collections. The sarcastic sample had 5 scores very close to the -3 threshold, with 2 crossing over to nearly -3.5 and the rest being fairly close behind. In comparison, the nonsarcastic sample had 3 scores that exceeded -3, excluding the outlier of r/privilegedpeoplegifs, and the rest of the values in the dataset remained relatively close together at the approximately -1 mark. While nothing substantial at a glance, it can give us an idea of the larger Reddit landscape, as this sample is pretty representative of the overall population.

The far higher max score indicates, as discussed before, that people are much more likely to upvote a nonsarcastic comment than a sarcastic one. We can interpret this to mean that the Reddit userbase either doesn't enjoy sarcastic comments enough to leave an upvote or they simply find it more compelling to upvote a comment that holds a definitive statement or opinion than one that is sarcastic for the purpose of trying to be funny. The minimum scores being the same, however, indicates that even if users aren't compelled to upvote a sarcastic comment they're not very compelled to downvote it either, at least not any more than they would be for a nonsarcastic comment.

So what does this mean to the average Reddit user? Well, that's harder to tell. There are quite literally a countless number of factors that go into someone's enjoyment or disapproval of a comment on social media, and nailing down a formula for what will get you the most upvotes is nigh impossible. However, with this information, you can at the very least go in with the idea that if you are looking for the highest score possible, then a nonsarcastic comment will probably lean more in your favor than a sarcastic one.

## Post-Project Reflection

All in all, I would say this project was a success. Obviously, as with any statistical analysis on a topic as open-ended and with as many factors to consider as social media, the end conclusion is not true for every instance out there. However, from the use of a representative sample of several thousands of Reddit comments, I think the conclusion reached from the data analysis holds up to respectable scrutiny and can be considered a basic idea to expand upon with further research.

The hardest part of the project probably had to do with trimming the data down to a reasonable level to work with. The initial sample came with literal thousands, almost a million, different entries. Obviously that is *far* too much to be able to illustrate clearly with a barplot, so I narrowed it down to 100. The graphs of the top and bottom 100 scores for each category ended up resembling one another quite a bit, almost like normal models with a leftward skew. However, it was still a bit too much to go any further depth into, so I had to narrow it down once again to the top and bottom 25. I would say in the future, possibly with additional manpower, expanding the analysis to the top 100 would be beneficial to cover a larger proportion of the initial sample and potentially get something that is even more representative of the total population in question.

There is one roadblock in particular that comes to mind during my work on the project. Initially, I wanted to include controversial comments in my analysis of Reddit users' likelihood to upvote or downvote a particular type of comment. Unfortunately, when it came time to interpret the various graphs I crafted, I realized that the score of controversial comments tells one almost nothing. Reddit's algorithm to determine controversy is based off comparing the number of upvotes to the number of downvotes. If the two are close, almost identical, in fact, then the post is considered controversial. For that reason, all controversial comments will have low scores due to the nature of how the site determines what is controversial and what isn't, so looking at the scores of a controversial comment won't explain very much about the average Reddit user's willingness to upvote or downvote it in comparison to a noncontroversial one. For that reason, I was forced to limit myself to sarcastic comments for this project.

In hindsight, one rather important factor I almost completely overlooked was the time of posting. Over time, the score for a comment will obviously change as more viewers leave their opinions on the comment, so filtering to a specific timeslot or date integral could help mitigate the impact of the lurking variable that is time. This would probably be an important change to make to the study in the future.