

Visualizing fine-grained emotions in Reddit posts through the GoEmotions dataset

Felix Dumont*
MIT EECS
MIT Sloan

Taylor Facen†
MIT EECS
MIT Sloan

ABSTRACT

In this paper, we propose a new interactive visualization site allowing readers to better understand the emotions associated with 483 subreddits. These visualizations are based over more than 58,000 Reddit posts, each classified with one or multiple of 29 distinct emotions in the GoEmotions dataset. Readers can then see the overlap between various Subreddit communities and get immediate insights as to the positivity and dominant emotions of each Subreddit.

In order to achieve this visualization effort, we leverage D3.js [?] to create a mix of an arc diagram, a Sankey diagram, a word cloud, a gauge chart and dynamic text displays. Each display was built to be simple enough to show immediate insights yet allow for various filtering and display clear examples. They allow us to gather powerful insights, such as how certain Subreddits (e.g. r/divorce) can be filled with gratitude and caring emotions and closely link to other similar Subreddits.

We present this site as a contrast to most previous work which in most case focuses on binary emotions (e.g. positive vs negative or hate speech vs normal). We strongly believe that by showing a fine-grained view of the emotions present on social media, we can develop deeper insights for the moderation teams but also possibly improve efforts to identify Subreddits with common negative emotions.

Index Terms: I.7.m [Document and text processing]: Miscellaneous—Life Cycle; K.4.2 [Computers and society]: Social Issues—

1 INTRODUCTION AND MOTIVATION

Hate speech and content moderation has been a recurring topic in both the academia and the mainstream media. Controversies among social media platforms around bullying and cases such as Twitter’s ban of President Donald J Trump’s account have raised awareness about the risk of inappropriate content. However, while most machine learning algorithms and visualizations focus on a rather binary view of a post’s positivity, we aim to take a more holistic view.

Throughout this paper and the associated web visualizations, we analyze over 58k Reddit posts, each classified with one more multiple of 29 distinct emotions. As such, we can not only look at the positivity of each Subreddit, but also what emotions are most dominant. We further analyze the connections between Subreddits, determining what Subreddits share the most common users and how their dominant emotions compare.

Ultimately, we want to show a fine-grained perspective on emotions and hope to show the readers the good, the bad and the ugly, and let them make their own mind by looking at intuitive visualizations and interactivity on their favourite Subreddit.

*e-mail: fdumont@mit.edu

†e-mail: tfacen@mit.edu

2 DATASET

The basis for this project is the GoEmotions dataset, which regroups 58k manually labelled English Reddit posts across 29 emotions. Each post receives at least one label such as neutral, anger, curiosity or admiration. The dataset includes 483 different Subreddits and 49,188 users from 2005 (the start of Reddit) to January 2019.

The GoEmotions dataset was put together by Demszky et al [?] in an effort to provide a richer view of human emotions in online social media messages and create the largest available dataset of emotions. Each post receives multiple annotations from different annotators, unlike other datasets which often focus on less accurate automated methods.

The authors justify the creation of the dataset as follows: “Understanding emotion expressed in language has a wide range of applications, from building empathetic chatbots to detecting harmful online behavior. Advancement in this area can be improved using large-scale datasets with a fine-grained typology, adaptable to multiple downstream tasks.”

3 RELATED WORK

3.1 Emotions Datasets

As mentioned in the first section, there has been several studies on the topic of emotions in social media. However, most studies are focused around concepts such as hate speech. Davidson et al (2017) [1] analyze hate speech in tweets and highlight the challenge in classifying messages considering the various nuances of each message. They further demonstrate that racist and homophobic tweets are more likely to be classified as hate speech. Meanwhile, sexist tweets are commonly reported as offensive instead. Li et al (2017) [?] also create a dataset called DailyDialog and instead analyze emotions in the context of online dialogs, moving away from more traditional binary classifications.

The GoEmotions dataset is relatively recent, having been released in 2020. Little academic work has yet to be done using this dataset outside of its original release paper. However, we will closely follow the use of this new, richer dataset. Until its release, the largest manually labelled dataset of emotions was CrowdFlower (2016) with 39k labeled examples although Bostan and Klinger (2018) [?] qualified it as noisy compared to other similar datasets.

3.2 Hate speech Visualization

Significant work has also been done on the topic of hate speech visualization. Capozzi et al (2018) [?] build a data visualization platform “as a Support to Study, Analyze and Understand the Hate Speech Phenomenon”. They build thorough dashboards showing the target of hate speech as well as the intensity of the tweets and the geolocation (when available). Piazza et al [?] also develop a set of visualizations in which they analyze hate speech and their target across multiple major social media. They show the links between the Facebook actors through an interactive network and leverage a Sankey diagram to break down the messages according to their most common keywords.

As of the submission of this paper, we have not found any interactive visualization on the GoEmotions dataset.

Readers will quickly notice that the r/divorce Subreddit is overwhelmingly positive at 74%, with gratitude, caring and optimism being the dominant emotions. Posts seem to be mostly supportive despite the usually difficult circumstances around posts, with an example given being: "Feel you pain. Stay strong".

Deep-Dive of Emotions in Reddit Posts

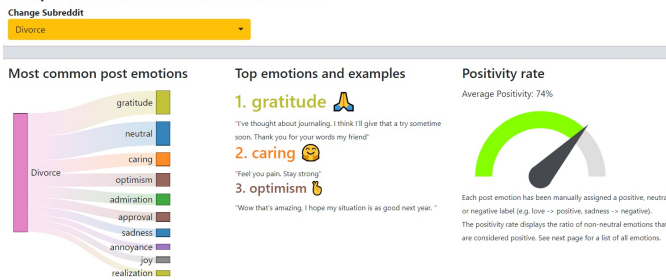


Figure 2: The deep-dive page performs a deep-dive of the emotions of any Subreddit.

5.4 Overall Insights

The advantage of such a broad visualization is there are numerous insights available and that we could elaborate on the results for each Subreddit in details. Still, there are some aggregated insights from the data we would want to emphasize for the reader:

1. Toxicity is present on most Subreddits and can cause great harm. However, the majority of posts in most Subreddits are rather positive (63% of non-neutral emotions are positive). The most common emotions observed are approval and admiration with negative emotions only representing 25% of the observed emotions across all Subreddits.
2. Communities extend beyond individual Subreddits. We observed significant overlap between the communities of many Subreddits. For instance, users in r/Advice are very likely to be active in other question-based Subreddits such as r/AskMenOver30 or r/askwomenadvice. This could lead to several opportunities including crossovers between Subreddits, shared moderation and opportunities for finding new members for a Subreddit.
3. Some Subreddits can have an abnormally high rate of some emotions such as sadness. While sadness is a normal emotion and can be good under many circumstances, it is also important that an individual expressing such an emotion in a Subreddit can also be a sign of deeper problems. As such, understanding which Subreddits are filled with such depressive emotions could help doing more monitoring or interventions towards certain communities or users to reduce potential consequences such as self-harm or suicide.

6 DISCUSSION AND BIASES

[MORE DISCUSSION]

6.1 Biases

First and foremost, there are potential biases in the collection and labelling of the dataset. Reddit posts may not be representative of human interactions as a whole, so insights should keep that in mind and be limited to the scope of the data. Furthermore, the label annotators were all native English speakers from India and may have had their own biases while labelling the data according to the emotions. Without context around some of these posts, the labels can also be inaccurate under certain scenarios.

There is also potential biases in our presentation of the data, although we attempt to limit it as much as possible. Our introduction of a positive or negative category does not take into account the intricacies of each post and could have challenges. The aggregation of the data can also introduce biases, for instance by hiding that a small group of users could account for most posts of a Subreddit.

ACKNOWLEDGMENTS

This work serves as a final project for MIT 6.859 Interactive Data Visualization. We went through several versions before ending up to the one presented here today and have to thank the teaching faculty and the teaching assistant team for helping us learn about D3 and for the constant help and feedback. We also ought to thank our classmates who thoroughly reviewed the first version of this project and provided valuable feedback.

REFERENCES

- [1] T. Davidson, D. Warmesley, M. Macy, and I. Weber. Automated Hate Speech Detection and the Problem of Offensive Language. p. 4.