

Moderation Suite for Reddit

Himanshu Singh

May 24, 2023

Contents

List of Figures	2
Listings	2
Abstract	3
1 Introduction	3
2 Methodology	5
3 Implementation	5
4 Deployment	8
5 Conclusion	10
References	11

List of Figures

1	The schematic structure of Reddit	3
2	u/eternalkerri describes the targetting of r/AskHistorians	4
3	r/SubredditDrama members discuss the ban of r/NoNewNormal	4
4	u/NomaiTraveler's suggested notion for "brigading"	5
5	User Overlap Page	9
6	Number of Submissions Chart	9
7	Submission Scores Chart	9
8	Content Type Chart	9
9	Prune Members Page	10
10	Purge Content Page	10

Listings

1	Implementation of User Overlap	6
2	Implementation of Prune Members	6
3	Implementation of Purge Content	7
4	Implementation of Subreddit Trends	8

Abstract

This project aims to reduce the workload of Reddit moderators, who have the responsibility of controlling public discourse in their respective subreddits. The problems faced by them were identified by checking various meta subreddits, and relating them with the established results. We realised most of these could be automatically detected and resolved by performing simple computations and communicating with Reddit API. We combined these protocols with an easy-to-use GUI to make a reliable moderation suite.

1 Introduction

Reddit is a pseudonymous website comprising of over 3.4 million user-run communities, called “subreddits” [4]. Registered users can make submissions and comments to any public subreddit, thereby making it accessible to all the subreddit members, unless it is stuck under a filter. [12]

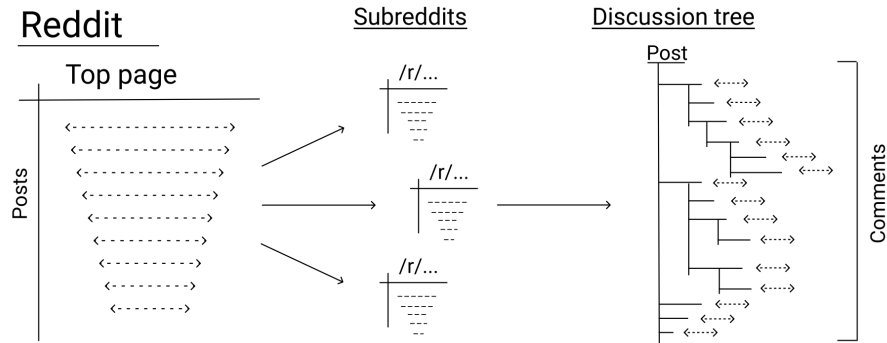


Figure 1: The schematic structure of Reddit [10]

Reddit takes a distributed approach for moderation, with the admins enforcing platform-wide rules, and the “voluntary” moderators deciding on the subreddit specific rules. [7] [8]. Subreddit moderation in itself is a highly subjective task, that often requires discussion amongst the various moderators prior to taking any sort of action. [9] It is thus essential to take into account input from users before automating any action.

Reddit allows “crosslinking” of posts in other communities, allowing an easy influx of users from one subreddit to another. Trogu et al. found that 27.8% crosslinks led to negative interactions, 78% of which were caused by top 1% communities. [17] This can be detrimental for the survival of smaller communities, making them vulnerable to manipulation of submission scores, which directly affects their visibility. [2] Another study, by Datta and Adar noted that 82% of controversial authors have only a single

“anti-social home”. [3] This provides an easy tracking mechanism for users involved in negative interactions.

This suggests that identification of mobilizing communities and subsequent blocking of cross-community interaction can curb majority of conflicts. However, inbuilt moderation tools like AutoModerator and Crowd Control go for simple blocking mechanisms, that does not take into account the nature of involved subreddits. [5] [6] Hence, a new blocking mechanism needs to be explored.

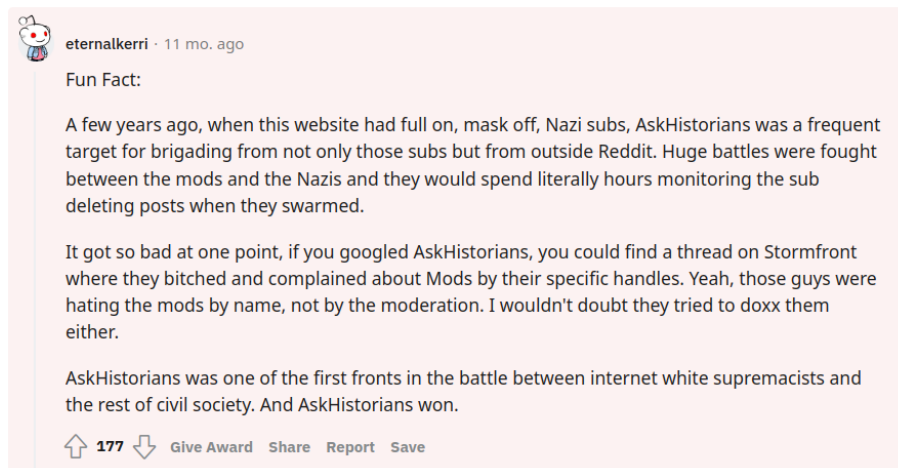


Figure 2: u/eternalkerri describes the targetting of r/AskHistorians [13]

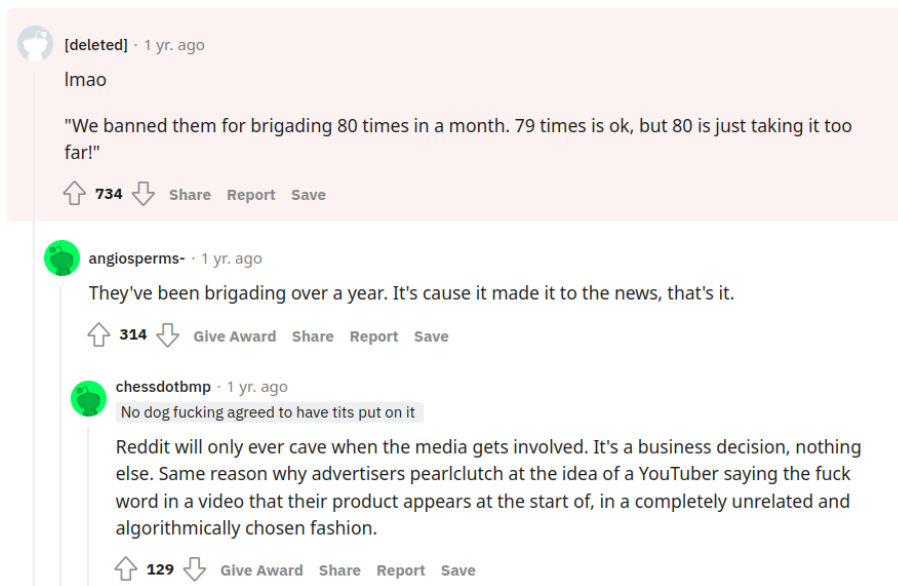


Figure 3: r/SubredditDrama members discuss the ban of r/NoNewNormal [14]



Figure 4: u/NomaiTraveler’s suggested notion for “brigading” [14]

2 Methodology

We start by exploring the subreddits with the highest overlap. From here on, we will be calling the number of shared users between two subreddits as simply “overlap”, and the probability of our subreddit’s member commenting in another subreddit as “participation probability”. [Jump to ‘Implementation of User Overlap’ on page 6] The latter having the advantage of detecting smaller but better overlapping subreddits more accurately.

Based on the results computed, the moderator is asked to flag the subreddits that she deems to be hostile. Additionally, the moderator can flag a list of keywords that she deems to be recurrent in negative interactions, and characterize what she deems to be a “throwaway” account. [18] Then, the list of users to be banned is generated with respect to various tolerance limits, to avoid false positives. [Jump to ‘Implementation of Prune Members’ on page 6]

A complementary feature is provided to purge all the submissions and comments left by users who were either banned earlier or were allowed to pass through by our tolerance limits. Moreover, this protocol allows the moderator to flag a list of domains on similar lines to keywords. [Jump to ‘Implementation of Purge Content’ on page 7]

Lastly, the moderator is provided with a set of charts to help her understand the changes in her subreddit over a period of time. This would help her determine the area and extent of action needed for her subreddit. This includes the number of active users, submissions, comments, scores, type of content, usage of flairs and frequency of edits over a period of time. [Jump to ‘Implementation of Subreddit Trends’ on page 8]

3 Implementation

Around December 25th 2021, the creator of u/Flair_Helper bot, u/Blank-Cheque reported that bots moderating more than 500 subreddits were being limited. [11] Hence, to implement the features discussed in the previous section, we will be looking at an application that can be run locally by any number of users. The application is designed

using GTK widgets for cross-platform compatibility, and written in Python to leverage the benefits offered by PRAW library. [1] [16] To avoid theming inconsistency, we will use libadwaita's `ApplicationWindow` as the top-level container. [15]

Listing 1: Computation of User Overlap ¹

```

1  userlist = []
2  for item in reddit.subreddit(home).new(limit = lim/10):
3      try:
4          author = item.author.name
5          if author not in userlist:
6              if author != "[deleted]":
7                  userlist.append(author)
8      except:
9          continue
10 for item in reddit.subreddit(home).comments(limit = 9*lim/10):
11     try:
12         author = item.author.name
13         if author not in userlist:
14             if author != "[deleted]":
15                 userlist.append(author)
16     except:
17         continue
18
19 overlap = {}
20 for user in userlist:
21     try:
22         profile = {}
23         for item in reddit.redditor(user).new(limit = ulim):
24             if item.subreddit.display_name.startswith("u_") or (item.subreddit.display_name
25                 == home):
26                 continue
27             profile[item.subreddit.display_name] = 1
28         for sub in profile:
29             if sub in overlap:
30                 overlap[sub] += 1
31             else:
32                 overlap[sub] = 1
33
34 prob = {}
35 scale = 100000000/(ulim*len(userlist))
36 for sub in overlap:
37     subcount = reddit.subreddit(sub).subscribers
38     prob[sub] = round(overlap[sub]*scale/subcount, 2)

```

Listing 2: Implementation of Prune Members

```

1  banlist = []
2  currtime = time.time()
3  for user in userlist:
4      profile = reddit.redditor(user)
5      if profile.link_karma < lklim:
6          banlist.append(user)

```

¹For the sake of keeping the listing concise, we have omitted the initialization of several objects, and passing of parameters. The complete implementation is publicly available at <https://github.com/ihsingh2/ms4r>.

```

7         continue
8     if profile.comment_karma < cklim:
9         banlist.append(user)
10        continue
11    if (currttime - profile.created_utc) < agelim:
12        banlist.append(user)
13        continue
14
15    count = 0
16    history = list(profile.new(limit = 100))
17    threshold = min(1, int(len(history)/10))
18    for item in history:
19        try:
20            if count >= threshold:
21                banlist.append(user)
22                break
23            if item.subreddit.display_name in sblacklist:
24                count += 1
25            if isinstance(item, praw.models.Comment):
26                if any(word in item.body for word in wblacklist):
27                    count += 1
28            elif isinstance(item, praw.models.Submission):
29                if any(word in item.title for word in wblacklist):
30                    count += 1
31            if item.is_self:
32                if any(word in item.selftext for word in wblacklist):
33                    count += 1
34        except:
35            continue

```

Listing 3: Implementation of Purge Content

```

1  for item in reddit.subreddit(home).comments():
2      if days != 0:
3          if item.created_utc > rmvafter:
4              item.mod.remove()
5              continue
6      if rmvban:
7          if item.author != None:
8              if item.author.name != "[deleted]":
9                  if any(reddit.subreddit(home).banned(item.author.name)):
10                     item.mod.remove()
11                     continue
12      if any(word in item.body for word in wblacklist):
13          item.mod.remove()
14
15  for item in reddit.subreddit(home).new():
16      if days != 0:
17          if item.created_utc > rmvafter:
18              item.mod.remove()
19              continue
20      if rmvban:
21          if item.author != None:
22              if item.author.name != "[deleted]":
23                  if any(reddit.subreddit(home).banned(item.author.name)):

```

```

24         item.mod.remove()
25         continue
26     if any(domain in item.url for domain in dblacklist):
27         item.mod.remove()
28         continue
29     if any(word in item.title for word in wblacklist):
30         item.mod.remove()
31         continue
32     if item.is_self:
33         if any(word in item.selftext for word in wblacklist):
34             item.mod.remove()

```

Listing 4: Implementation of Subreddit Trends

```

1  # this listing demonstrates the generation of chart for submission scores
2  count = {}
3  for item in reddit.subreddit(home).new():
4      try:
5          date = datetime.datetime.fromtimestamp(item.created_utc).date()
6          key = str(date.day) + '-' + str(date.month)
7          if key in count:
8              count[key][0] = max(count[key][0], item.score)
9              count[key][1] = min(count[key][1], item.score)
10             count[key][2] += item.score
11             count[key][3] += 1
12         else:
13             if len(count) == days:
14                 break
15             count[key] = [item.score, item.score, item.score, 1]
16     except:
17         continue
18
19  vals = np.array(list(count.values()))
20  fig = Figure(figsize = (5, 4), dpi = 100)
21  ax = fig.add_subplot()
22  ax.plot(list(count.keys()), vals[:, 0], label='max_score')
23  ax.plot(list(count.keys()), vals[:, 1], label='min_score')
24  ax.plot(list(count.keys()), vals[:, 2] / vals[:, 3], label='avg_score')
25  ax.legend()
26  canvas = FigureCanvas(fig)

```

4 Deployment

Our application was tested for r/kolkata. 450 comments and 50 recent submissions were scanned to obtain a list of 300 users. 1351 overlapping subreddits were found by scanning the individual user history. The highest overlapping subreddits have been listed in Figure 5.

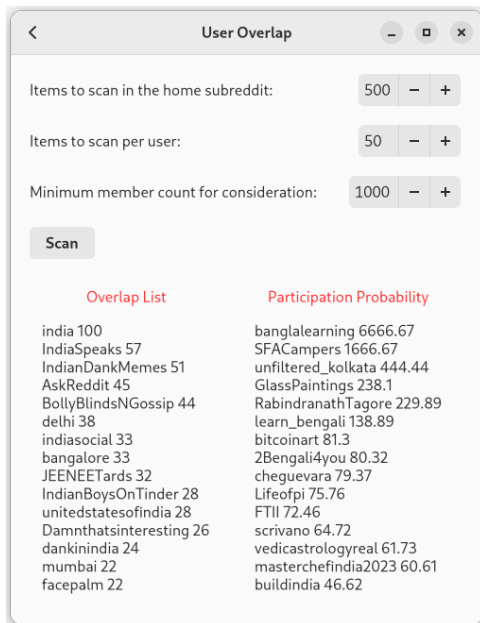


Figure 5: User Overlap Page

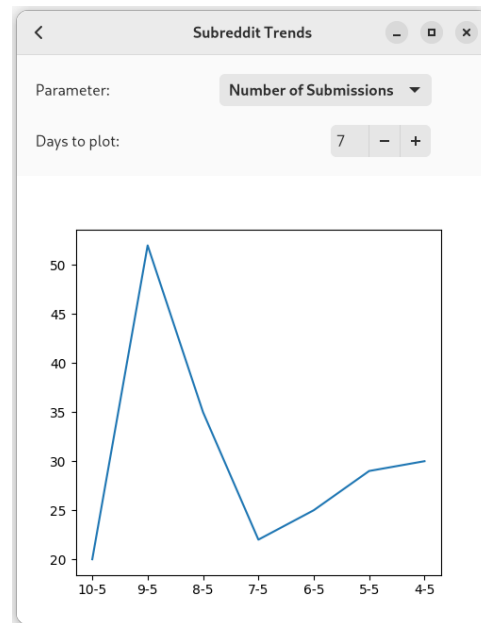


Figure 6: Number of Submissions Chart

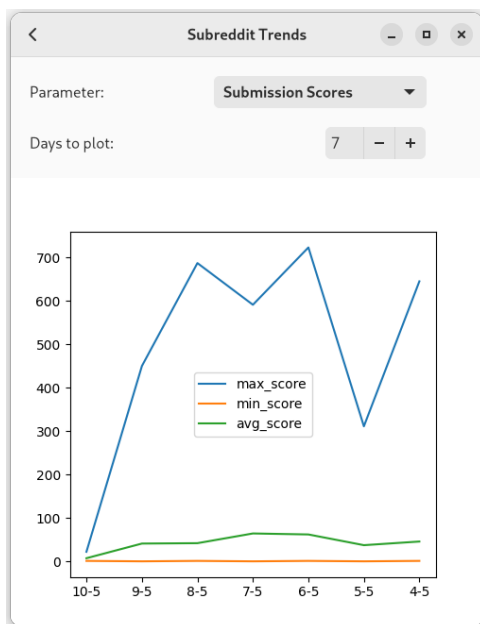


Figure 7: Submission Scores Chart

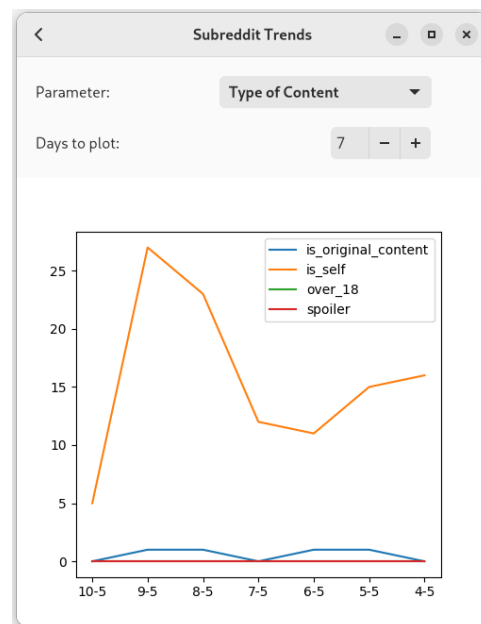


Figure 8: Content Type Chart

Figure 9: Prune Members Page

Figure 10: Purge Content Page

5 Conclusion

In this report, we discussed an easy-to-use application for preventing negative interactions. It is worth mentioning that users may suffer from frequent request limit hits by Reddit API, if their account is relatively new. The possibility of classifying items and user history based on the data available in the moderation logs may be explored to further ease the job. Tracking the changes in upvote ratio over time is a key measure in detecting malpractices. However, its implementation would be very costly and is hence beyond the scope of this project.

References

- [1] Bryce Boe. Praw: The python reddit api wrapper. Retrieved from <https://praw.readthedocs.io/en/latest/>, 2022-11-19.
- [2] Mark Carman, Mark Koerber, Jiuyong Li, Kim-Kwang Raymond Choo, and Helen Ashman. Manipulating visibility of political and apolitical threads on reddit via score boosting. In *2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/ 12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*, pages 184–190, 2018.
- [3] Srayan Datta and Eytan Adar. Extracting inter-community conflicts in reddit. *CoRR*, abs/1808.04405, 2018.

- [4] Metrics for Reddit. New subreddits by month - reddit history. Retrieved from <https://frontpagemetrics.com/history/month>, 2022-05-25.
- [5] Reddit Help. Automoderator. Retrieved from <https://mods.reddithelp.com/hc/en-us/articles/360002561632-AutoModerator>, 2022-11-19.
- [6] Reddit Help. Crowd control. Retrieved from <https://mods.reddithelp.com/hc/en-us/articles/360038129231-What-is-Crowd-Control>, 2022-11-19.
- [7] Reddit Inc. Reddit content policy. Retrieved from <https://www.redditinc.com/policies/content-policy>, 2022-11-18.
- [8] Reddit Inc. Reddit user agreement. Retrieved from <https://www.redditinc.com/policies/user-agreement>, 2022-11-18.
- [9] J. Nathan Matias. The civic labor of volunteer moderators online. *Social Media + Society*, 5(2):2056305119836778, 2019.
- [10] Alexey Medvedev, Renaud Lambiotte, and Jean-Charles Delvenne. *The Anatomy of Reddit: An Overview of Academic Research*, pages 183–204. 05 2019.
- [11] r/Flair_Helper. Flair_helper is currently not accepting new subreddits automatically. Retrieved from https://www.reddit.com/r/Flair_Helper/comments/ro0uqk/flair_helper_is_currently_not_accepting_new/, 2022-10-28.
- [12] r/help. Frequently asked questions. Retrieved from <https://www.reddit.com/r/help/wiki/faq/>, 2022-11-18.
- [13] r/SubredditDrama. Rare skirmish breaks out in r/askhistorians. Retrieved from <https://www.reddit.com/r/SubredditDrama/comments/rd3zr6/>, 2022-11-19.
- [14] r/SubredditDrama. r/nonewnormal has been banned. discuss this dramatic happening here! Retrieved from <https://www.reddit.com/r/SubredditDrama/comments/pfz0d2/>, 2022-11-19.
- [15] Purism SPC. Adw – 1. Retrieved from <https://gnome.pages.gitlab.gnome.org/libadwaita/doc/main/index.html>, 2022-11-19.
- [16] GTK Development Team. Gtk – 4.0. Retrieved from <https://docs.gtk.org/gtk4/>, 2022-11-19.
- [17] Ava Trogus, Daniel Fuchs, Aaron Burtle, and Josh Katz. Community interactions on reddit over time. Retrieved from https://courses.cs.washington.edu/courses/cse481ds/21au/resources/trogusava_3717447_71314956_481DS_Reddit_final_report.pdf.
- [18] Emily van der Nagel. Faceless bodies: Negotiating technological and cultural codes on reddit gonewild. *Scan: Journal of Media Arts Culture*, 10, 01 2013.