

Nesta attacks study

Patrycja Tempska and Rafal Urbaniak

Contents

1 Section Abstract	1
2 Introduction	1
3 Automated versus Human-based Moderation	2
4 Pro-active and reactive moderation	3
5 Collective Intelligence Approach to Counter-speech	4
6 Experimental Design	4
7 Results	5
References	5

1 Section Abstract

This article describes an experimental intervention study based in a naturalistic, digital setting (Q&A forum - Reddit), utilizing a collective intelligence approach to content moderation and reduction of the level of verbal aggression among a selected group of Reddit users who regularly attack other community members. Collective Intelligence in this sense means exploring the collaboration between human and machine intelligence to develop solutions to social challenges. Artificial Intelligence was used to detect verbal aggression (personal attacks) and notify human volunteers about attacks. Volunteers after receiving notifications employed interventions based on norm or empathy promotion. We find that only those who were sanctioned with norms-inducing interventions had their personal attacks' level significantly decreased.

2+2 #use this formatting for chunks

[1] 4

2 Introduction

Although much effort has been made in order to tackle the problem of verbal aggression and harassment online, looking at various reports and surveys, it remains a common hindrance for people engaging with social media in their everyday lives. The situation got exacerbated amidst the COVID19 pandemic, during which a majority of our social life moved to cyberspace. During this shift, there was an increase in cyberbullying attitudes and perpetration (Barlett, Simmers, Roth, & Gentile (2021)), 90% increase in public reports of illegal online content¹, including 114% increase in non-consensual sharing of intimate images, 30% increase in cyberbullying, as well as 40%

¹<https://www.aspistrategist.org.au/australias-esafety-commissioner-targets-abuse-online-as-covid-19-supercharges-cyberbullying/>

of increase in adults reporting online harassment. According to a report conducted by company Light \footnote{\href {https://light.com/Toxicity_during_coronavirus_Report-Light.pdf}}{https://light.com/Toxicity_during_coronavirus_Report-Light.pdf}}, hate speech directed towards China and the Chinese went up by 900% on Twitter. Gaming platforms were in the spotlight as well, with a 40% increase in toxicity on Discord.

But alongside the growing need for even more efficient and proactive moderation, the capacity to execute it did not go hand in hand, forcing companies and policymakers to rethink the current model of moderation processes and workforce. Due to the COVID19 restrictions including social distancing, a lot of those serving the role of moderators had to be sent home² without the ability to work remotely because of the constraints affiliated with restrictive non-disclosure agreements (NDA) among others. Curtailing the moderators' workforce was accompanied by more agency given to algorithms and AI-based moderation. Those changes, as argued by Gerrard (2020), can be seen as a serious red flag in terms of safety for all users on online platforms.

3 Automated versus Human-based Moderation

The hindrances and threats that go along with the Artificial Intelligence-based methods for moderation have been widely debated, with the most critical discussions revolving around technology performance (MacAvaney et al. (2019), Schmidt & Wiegand (2017)). State-of-the-art solutions are mostly governed by statistical methods including deep learning and machine learning (LeCun, Bengio, & Hinton (2015), Sejnowski (2020), Jordan & Mitchell (2015)). Their performance is inherently tied to the amount of data being fed to the system and the quality of its annotation. At different stages of the process, from datasets gathering and preparation, annotation to the training or algorithms themselves, biases seem to be omnipresent (Binns, Veale, Van Kleek, & Shadbolt (2017), Geva, Goldberg, & Berant (2019) Mehrabi, Morstatter, Saxena, Lerman, & Galstyan (2021)). Users' of online services are also creative in their strategies to circumvent automated content moderation systems and as shown by Gröndahl, Pajola, Juuti, Conti, & Asokan (2018), current techniques are vulnerable to the most common evasion attacks like word changes (inserting typos and leetspeak), word-boundary changes (inserting or removing whitespace), or word appending (appending common or non-hateful words like "love"). Lack of generalisability of the models - an ability to perform well on datasets coming from sources other than the one used for training is a serious shortcoming as well (Yin & Zubiaga (2021), Swamy, Jamatia, & Gambäck (2019), Rosa et al. (2019)). As shown by Wu, Ribeiro, Heer, & Weld (2019), Lipton & Steinhardt (2019), and Musgrave, Belongie, & Lim (2020) in practice, creations of models often lacks thorough error analysis and legitimate experimental methodology, which can result in non-reproducibility. This is also connected with a potential lack of thorough understanding of the limitations of the models and spurious conclusions being made to a wider public. Specifically, @Lipton & Steinhardt (2019) distinguishes four dysfunctional patterns occurring in the current research paradigm in the industry and academia alike. First, the inability to draw a clear distinction between speculation and explanation, with the first one often being disguised as the second. Second, inability for successful identification of the sources of empirical gains (whether it was problem formulation, optimization of the heuristics, data-preprocessing, hyperparameter tuning, or perhaps yet another aspect). Third, "mathiness" - the use obscure language and often covering weak argumentation with the alluring but often apparent depth of technical jargon. Last but not least - misuse of language. This includes suggestive definitions without proper explanation of what they mean in the context (e.g. inflating good performance in simple NLP tasks to human-level understanding), overloading the papers with technical terminology, or suitcase words (those words that can encompass a variety of meanings, e.g. consciousness).

Yet another obstacle in the process is the lack of gold standard in dataset creation and taxonomies of abusive language being used for instance in the process of annotating different datasets. Frequently people obtain data from various sources and do not follow any universally used instructions when it comes to its annotation, leading to discrepancies between various datasets being tagged within one domain (e.g. hate speech). Lack of expert annotators and proper annotation criteria and instructions are also widespread, with the common practice hiring untrained workers from Mechanical Turk or other crowdsourcing platforms.

²<https://qz.com/india/1976450/facebook-covid-19-lockdowns-hurt-content-moderation-algorithms/>

Although there are some initiatives developed in response, most notably, functional tests for Hate Speech Detection Models created by Röttger et al. (2020), or the Online Safety Data Initiative (OSDI) LINK <https://onlinesafetydata.blog.gov.uk/about-us/>, focused on projects related to improving access to data, standardizing the description of online harms, as well as creating tools and benchmarks for evaluation of technologies focused on safety, much effort must be made before wider adoption of such solutions comes into force.

At the same time, only automated methods can scan through the massive amount of content being generated every day on different platforms. On Facebook, there are more than 3B comments and likes daily (<https://martech.org/facebook-3-2-billion-likes-comments-every-day/>), 500M tweets are sent daily on Twitter (<https://www.oberlo.com/blog/twitter-statistics>), and over 2B comments made by users of Reddit in 2020 (https://old.reddit.com/r/blog/comments/k967mm/reddit_in_2020/) which is almost 3M comments made daily. With this amount of content, it's either impossible or extremely costly to scale the moderation workforce. One can also have doubts about the ethical aspects of hiring workers who are often unaware of how this kind of task will affect their well-being. Being submersed in the cyber-Augean stables takes a toll on many - as examined by Roberts (2014) & Roberts (2016), workers who are hired for such tasks are often low-status and low-wage, isolated, and asked to keep what they've seen in secret under restrictive NDAs. Screening through the reported user-generated content is connected with exposure to violent and deeply disturbing materials, with child pornography, murders, or suicides as examples of the most extreme cases. This can lead to serious psychological damage, like depression, or PTSD (Roberts (2014)). Some of the employees filed a lawsuit against Facebook and as a result, the company agreed to pay \$52M in compensation for mental health issues developed during the job (<https://www.theverge.com/2020/5/12/21255870/facebook-content-moderator-settlement-scola-ptsd-mental-health>). Taking into consideration that Facebook employs 15K moderators (<https://www.forbes.com/sites/johnkoetsier/2020/06/09/300000-facebook-content-moderation-mistakes-daily-report-says/?sh=6c6bdbbc54d0>) and most likely more are needed to keep up with the growing amount of content, with the parallel considerations about the negative effects of content moderation on mental health, a collaboration between humans and machines in this area seems inevitable.

4 Pro-active and reactive moderation

There are different approaches when it comes to the moderation of online content. One can follow the workflow of reactive moderation, which happens once the content is published. Harmful messages can be either reported by the users of the platforms or automated methods and then sent for review to moderators. A set of actions can be then taken depending on the platform and their community guidelines - on the content or user level. A harmful message can be deleted, made invisible to other users, or certain profanities can be altered with special signs to censor them. Depending on the type and amount of infraction, a particular user can be warned, muted or shadowbanned, or banned from further participation in the community for a period of time. The weakness of the reactive method is that the damage is done. Whoever is the recipient of the abusive message has the chance to see it and potentially suffer (Hoff & Mitchell (2009), Keipi, Näsi, Oksanen, & Räsänen (2016), Wwachs2019associations). Yet another weakness connected with relying solely on human reports is the content that is harmful but unreported by a recipient or any bystander. Although the exact scale of unreported content is not known, various self-report studies show that a lot of children, teens, or even adults do not report cyberbullying or harassment online (LINKS: <https://www.ctvnews.ca/canada/as-the-pandemic-forces-us-online-lgbtq2s-teens-deal-with-cyberbullying-1.5430945> ; <https://www.adl.org/free-to-play-2020> ; <https://www.ditchthelabel.org/wp-content/uploads/2017/05/InGameAbuse.pdf>).

Yet another type of moderation can be distinguished as pro-active or pre-moderation. In pre-moderation, automated methods are either based on Artificial Intelligence or other less sophisticated tools (e.g. blacklists) and can screen the content before it gets published. If a type of harmful message gets detected, it can be removed before reaching the recipient. Due to the aforementioned dubious performance of state-of-the-art statistical methods, particularly low precision, they are rarely used autonomously.

Pro-active moderation can be utilized using AI or other methods to promote socially positive engagement. Instead of or in the combination with punitive solutions like privileges restriction, one can induce empathy or community norms with counter-speech. Counter speech as described by Danger-

ous Speech Project (<https://dangerousspeech.org/counterspeech/>) is “any direct response to hateful or harmful speech which seeks to undermine it.” As examined by Munger (2017), counter-speech can be effective in the reduction of racist tweets (although only in the condition in which a white male with high followers was approaching another white male). In a study conducted by Bilewicz et al. (2021), a bot ...

5 Collective Intelligence Approach to Counter-speech

6 Experimental Design

This was a 6-months field experiment in a digital setting conducted on a popular Q&A and news forum, Reddit (www.reddit.com). We formed treatment and control groups based on three main criteria: 1. During the intervention period, we have expected to have 20 active volunteers at any given time, each willing to conduct 10 interventions daily. Thus, we needed approximately 200 attacks daily generated by the treatment groups. 2. Recruitment of users who regularly attack others was necessary to measure the effect of interventions at the end. 3. Recruitment of users who were active during the whole preliminary monitoring period was necessary to minimize the risk of attrition during the study. First, we spent 9 weeks for preliminary Reddit’s monitoring period to recruit only those users who generate at least one personal attack per week and were active during the whole period of the monitoring. As mentioned, sustained activity was crucial to minimize the risk users becoming inactive during the course of the study. Next, we have calculated the daily average number of personal attacks generated by the group who met the above criteria (which resulted in 357 attacks per day - 1.94 attacks daily per person on average). Knowing that we need around 200 attacks daily (just enough so our volunteers can keep up with the volume) we have randomly selected 390 people for our study groups (195 per group). The rest (304 people) were selected as our control group.

Three groups were formed - two treatment groups (195 people in each) and one control group (304 people). The first treatment group received counter-speech interventions based on normative influence, while the second one received interventions based on empathetic influence. Users in the control group did not receive any intervention.

The duration of the experiment, 6 months, was divided into three 2-months periods. The first two months served as a monitoring period to properly select groups and establish baselines. The next 2 months served as treatment period, during which groups received counter-speech comments from volunteers, in response to personal attacks detected by the Artificial Intelligence-based system. The last 2-months served as the post-treatment monitoring period to gather the data needed to evaluate the effectiveness of interventions.

Normative interventions can refer either to general social norms of civility and respect, community standards, a particular subreddit rule, or a descriptive norm among others. E.g. Insulting others is against Reddit’s policy. Hey there, we do not call each other like that here. Empathetic interventions can refer to the emotional state of the recipient or the sender of the attack or both. They are designed to evoke an empathetic response. E.g. There is a human being on the other side who might be hurt by your words. I see you are frustrated but remember the human.

Cyberviolence was defined in this experiment as a personal attack - any kind of verbal harassment, insult, or threat directed against the interlocutor in a text-based conversation online. Those were detected using Samurai Labs’ cyberviolence detection system.

The following hypotheses were formulated:

H1: If a group of human volunteers notified by an AI-based cyberviolence detection system about cyberviolence generated by the treatment group users (cyberviolence will be defined as a personal attack, harassment, or a threat targeted against an interlocutor) responds with counter-speech interventions, this will result in a decreased cyberviolence level for the whole group after the intervention period.

H2: If two groups receive different types of interventions (empathy-based or normative), then the decrease in cyberviolence will be larger in the case of normative interventions in comparison to the empathy-based ones.

7 Results

References

- Barlett, C. P., Simmers, M. M., Roth, B., & Gentile, D. (2021). Comparing cyberbullying prevalence and process before and during the COVID-19 pandemic. *The Journal of Social Psychology*, 1–11.
- Bilewicz, M., Tempa, P., Leliwa, G., Dowgiało, M., Tańska, M., Urbaniak, R., & Wroczyński, M. (2021). Artificial intelligence against hate: Intervention reducing verbal aggression in the social network environment. *Aggressive Behavior*, 47(3), 260–266.
- Binns, R., Veale, M., Van Kleek, M., & Shadbolt, N. (2017). Like trainer, like bot? Inheritance of bias in algorithmic content moderation. *International Conference on Social Informatics*, 405–415. Springer.
- Gerrard, Y. (2020). <? covid19?> The COVID-19 mental health content moderation conundrum. *Social Media+ Society*, 6(3), 2056305120948186.
- Geva, M., Goldberg, Y., & Berant, J. (2019). Are we modeling the task or the annotator? An investigation of annotator bias in natural language understanding datasets. *arXiv Preprint arXiv:1908.07898*.
- Gröndahl, T., Pajola, L., Juuti, M., Conti, M., & Asokan, N. (2018). All you need is" love" evading hate speech detection. *Proceedings of the 11th ACM Workshop on Artificial Intelligence and Security*, 2–12.
- Hoff, D. L., & Mitchell, S. N. (2009). Cyberbullying: Causes, effects, and remedies. *Journal of Educational Administration*.
- Jordan, M. I., & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245), 255–260.
- Keipi, T., Näsi, M., Oksanen, A., & Räsänen, P. (2016). *Online hate and harmful content: Cross-national perspectives*. Taylor & Francis.
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444.
- Lipton, Z. C., & Steinhardt, J. (2019). Troubling trends in machine learning scholarship: Some ML papers suffer from flaws that could mislead the public and stymie future research. *Queue*, 17(1), 45–77.
- MacAvaney, S., Yao, H.-R., Yang, E., Russell, K., Goharian, N., & Frieder, O. (2019). Hate speech detection: Challenges and solutions. *PloS One*, 14(8), e0221152.
- Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6), 1–35.
- Munger, K. (2017). Tweetment effects on the tweeted: Experimentally reducing racist harassment. *Political Behavior*, 39(3), 629–649.
- Musgrave, K., Belongie, S., & Lim, S.-N. (2020). A metric learning reality check. *European Conference on Computer Vision*, 681–699. Springer.
- Roberts, S. T. (2014). *Behind the screen: The hidden digital labor of commercial content moderation*. University of Illinois at Urbana-Champaign.
- Roberts, S. T. (2016). *Commercial content moderation: Digital laborers' dirty work*.
- Rosa, H., Pereira, N., Ribeiro, R., Ferreira, P. C., Carvalho, J. P., Oliveira, S., ... Trancoso, I. (2019). Automatic cyberbullying detection: A systematic review. *Computers in Human Behavior*, 93, 333–345.
- Röttger, P., Vidgen, B., Nguyen, D., Waseem, Z., Margetts, H., & Pierrehumbert, J. (2020). Hatecheck: Functional tests for hate speech detection models. *arXiv Preprint arXiv:2012.15606*.
- Schmidt, A., & Wiegand, M. (2017). A survey on hate speech detection using natural language processing. *Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media*, 1–10.
- Sejnowski, T. J. (2020). The unreasonable effectiveness of deep learning in artificial intelligence. *Proceedings of the National Academy of Sciences*, 117(48), 30033–30038.
- Swamy, S. D., Jamatia, A., & Gambäck, B. (2019). Studying generalisability across abusive language detection datasets. *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, 940–950.
- Wu, T., Ribeiro, M. T., Heer, J., & Weld, D. S. (2019). Errudite: Scalable, reproducible, and testable

error analysis. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 747–763.

Yin, W., & Zubiaga, A. (2021). Towards generalisable hate speech detection: A review on obstacles and solutions. *PeerJ Computer Science*, 7, e598.