

Subjective *Isms*? On the Danger of Conflating Hate and Offense in Abusive Language Detection

Anonymous ACL submission

Abstract

Natural language processing (NLP) research has begun to embrace the notion of annotator *subjectivity*, motivated by variations in labelling. This approach understands each annotator's view as valid, which can be highly suitable for tasks that embed subjectivity, e.g., sentiment analysis. However, this construction may be inappropriate for tasks such as hate speech detection, as it affords equal validity to all positions on e.g., sexism or racism. We argue that the conflation of hate and offence can invalidate findings on hate speech, and call for future work to be situated in theory, disentangling hate from its orthogonal concept, offense.

1 Introduction

Recently, Natural Language Processing (NLP) researchers have dedicated significant efforts towards tasks under the umbrella of online abuse detection. For example, racism (e.g. Talat, 2016; Talat and Hovy, 2016), sexism and misogyny (e.g. Jiang et al., 2022; Zeinert et al., 2021), xenophobia (e.g. Ross et al., 2016), homophobia (Dias Oliva et al., 2021), and transphobia (e.g. Chakravarthi et al., 2022) have been all been proposed as suitable for automated identification using NLP methods. Collectively these can be referred to as *isms*. *Isms* are prejudices, stereotyping, or discrimination on the basis on some personal characteristic. For example, sexism is defined as prejudice, stereotyping, or discrimination, typically against women, on the basis of sex or gender (Masequesmay, 2008).

This line of research has been faced with high annotator disagreement (e.g. Leonardelli et al., 2021), and as a result has conceptualised this as an indication that the concepts themselves are subjective. For example, Rottger et al. (2022) argue that labelling these phenomena is inherently subjective and can either be addressed as *descriptive*, i.e., encouraging annotator subjectivity, or *prescriptive*,

i.e., discouraging it. By constructing abuse as individually subjective, social norms are disregarded in favour of an approach that is blind to existing conditions of marginalisation. This stands in contrast to early work in the field, which sought to tease apart the distinction between offensiveness and hate (Davidson et al., 2017), and sought frameworks to identify the particular vectors which indicated hate (Talat et al., 2017; Wright et al., 2017).

Discrimination is also an area subject to policy and regulatory debates. Policy often distinguishes *hate* from *offence*. For instance, in their definition of sexism, the European Institute for Gender Equality (EIGE) position sexism as the *presence* rather than the offensiveness of a gendered stereotype:

‘Sexism is linked to beliefs around the fundamental nature of women and men and the roles they should play in society. Sexist assumptions about women and men, which manifest themselves as gender stereotypes, can rank one gender as superior to another.’

In this position paper, we consider such *isms* and how offence and hate are orthogonal concepts that can be mutually informative, and argue that their conflation can delegitimise research artefacts and findings. That is, we contend that the hatefulness of a statement is invariant of a reader's position on whether it should be allowed within a particular public forum. Consider for instance the use of gendered slurs: while inappropriate for a general audience (e.g., a public debate) they may be appropriate for others (e.g., academic work exploring the uses of expletives). In particular, we argue that *isms* are culturally defined, whereas offence is a subjective experience. Thus, we argue that it the presence of a stereotype that determines if a statement is hate speech, rather than individual perceptions of its offensiveness. Understanding *isms* as culturally defined, and offence as individually subjective allows us to distinguish any offence

caused to a reader from whether a message contains hate speech. To this end, we call for researchers to adopt new guidelines for annotating online abuse tasks which delineate the degree of offence caused by statements from the phenomenon itself.

2 Understanding Subjectivity

Recent efforts in NLP have constructed annotation as a subjective task, without attending to what other fields have understood subjectivity to be. *Subjectivity* has been given as the cause for “humans (e.g. dataset annotators) [being] sensitive to sensory demands, cognitive fatigue, and external factors that affect judgements made at a particular place and point in time” (Alm, 2011). Philosophy however constructs *subjectivity* as concerned with people’s differing perspectives, formed by factors such as cultural and individual experiences (Solomon, 2005). This construction of subjectivity purports that the only valid knowledge is based on personal experiences, thereby negating the existence of objective or communal truths. In contrast, *relativism* proposes that criteria of judgement are relative to a culture or society (Baghrmian, 2004). For instance, while humour would on one hand be subjective, we can understand concepts such as beauty standards to be culturally defined.

Hate speech detection, in particular, has frequently been argued to be a subjective task (e.g. Almanea and Poesio, 2022; Basile, 2020; Sandri et al., 2023). Under this framing, researchers also collapse data labelled as offensive with that labelled as hate speech (e.g. Leonardelli et al., 2021), thereby further conflating offence and hate. For instance, Akhtar et al. (2021) posit that ‘judging whether a message contains hate speech is quite subjective, given the nature of the phenomenon’. When categories of abuse are described as subjective, we understand that there is no ground truth, and wider cultural norms do not impact what constitutes hate. Within the concept of *isms*, we argue that is the wrong approach and that these are culturally defined. That is, we argue that, for a stereotype or norm, there *is* a ground truth given by the cultural and temporal context a statement is made in.

2.1 Stereotypes as Socially-defined Artefacts

Isms are a term given to various forms of marginalization and concepts such as racism, sexism, transphobia, etc. Such *isms* rely on tropes and stereotypes about a target group (Manne, 2017). They describe beliefs about the way a group is and how

it ought to be (Ellemers, 2018). Although stereotypes are held by individuals, they are formed collectively. For example, stereotypes are observable: we can catalogue the content of gender stereotypes within a culture (Prentice and Carranza, 2002), suggesting these are not solely individual but instead exist in the ‘collective brain’.

Haslam et al. (1997) argue that stereotypes emerge when individuals are acting in terms of a common social identity. Although the belief that stereotypes are simply an inferior representation of an unfamiliar group may be alluring, they serve to represent group-based realities: they represent (and accentuate) perceived differences between then in- and out-group (Haslam et al., 1997). Through the lens of self-categorisation theory, Haslam et al. (1997) argue that stereotypes are a social force—they reassure individuals of their belonging to a group ‘by: (1) enhancing perceived in-group homogeneity; (2) providing associated expectations of mutual agreement; and (3) producing pressure to actively reach consensus through mutual influence’. Uniformity of belief is thus the very essence of a stereotype. The shared nature of stereotypes is what causes their severity, a single individual holding and acting on discriminatory beliefs is less consequential than a group holding and acting on the same beliefs. However, because stereotypes are collective, they are also fuzzy; while individuals in the in-group are at least aware of stereotypes, they do not necessarily believe in them. This is in part why the degree of offence to *isms* may vary. Group memberships and social relations play a key role in shaping cognition, leading to the application and salience of stereotypes to be context-dependent but consensual at the group level nonetheless.

2.2 Acceptability as a Social Norm

Generally speaking, some *isms* are less socially acceptable nowadays than they were a century ago due to the social justice movements of the last century. Such movements have, in some countries, resulted in an increased public awareness of the harms caused by stereotypes, making support for some of them less socially acceptable. That is, the Overton Window, a political theory that describes the spectrum of acceptable policies and discourse, has shifted to make it less socially acceptable to hold particular stereotypical beliefs. The result of such a shift is that people do not wish to label statements they agree with as an *ism* lest they be labelled as **ists* themselves. For instance,

homophobia has become less tolerated in many countries, and individuals do not want their statements, or them, to be labelled as homophobic. Yet while being labelled as homophobic is perceived as undesirable, this does not mean that homophobic comments are not made, and policies not pursued. For example, in the United States of America, the American Civil Liberties Union has currently flagged more than 500 legal bills as anti-LGBTQ (American Civil Liberties Union, 2023). Thus, despite forward progress on some forms of discrimination and isms (Azcona et al., 2023; Menasce Horowitz, 2023), there are still socially acceptable isms that come in two general flavours: the benevolent isms and the scientific isms.

The Benevolent *Ism Some stereotypes may be seen as ‘positive’ and therefore not recognised by some as hateful. The existence of ‘benevolent’ stereotypes (Jha and Mamidi, 2017), such as ‘neosexism’ (Tougas et al., 1995)—those without clear negative connotations—means that annotators may be unlikely to recognise them as harmful. For example, the seemingly positive stereotype in Western nations that Asians are successful, high-achievers leads to their vilification (for being *too* high-achieving) and the perception that they lack interpersonal skills (Wong and Halgin, 2006). These stereotypes may also cause indirect harm to the individuals who may feel they are not living up to what is expected from them (Haslam et al., 1997). We might be tempted to only oppose or target stereotypes that imply or directly state that a certain group is inferior, however this approach would leave many of the issues of stereotyping unaddressed. For example, not addressing claims such as ‘women need to be protected’ or that ‘women’s bodies are more aesthetically pleasing’ suggests that the perception of women as inferior, or inherently sexualised, should remain acceptable.

The Scientific *Ism This *ism* uses evolutionary biology as evidence for stereotypes. In this case, different groups are proposed as differing on the basis of *natural* differences, such as physiology. One such example is the idea that women are naturally more nurturing than men due to imaginations of gender roles of the past. However, investigations of hunter-gatherer societies indicate that this idea may not be an accurate reflection of past societies and social evolution (Hewlett and Macfarlan, 2010). The idea of evolutionary psychology as evidence stems from Social Darwinism (Miller, 2011), which ar-

gues that one cannot accuse nature of being *-ist*, and therefore any generalisation based on biology cannot be labelled as such. Such pseudo-scientific *isms* are commonly used as a rationalisation for the ‘objective’ differences between dominant and marginalised groups (e.g. Browne (2006)).

2.3 Separating Isms and Offensiveness

So far, we have established that *isms* are rooted within socio-cultural contexts, and, while not necessarily factual or objective, exist as normative and therefore stable concepts, given their socio-cultural and temporal situations. Due to being norms rooted in a socio-cultural context, *isms* can be the cause of harms to members of targeted groups, e.g., psychological harms to the targeted group, present barriers to harmonious community relations, or pose threats to law and order (Barendt, 2019).

Offensiveness can generally be understood as moral outrage or disgust (Sneddon, 2020). As the existence of *isms* can be harmful, it is tempting to suggest that they should always be constructed as offensive. However, this would not afford the disagreements that have been observed in annotation of various *isms*, a task that typically has a high level of disagreement. Such disagreement can be accounted for by considering the degree of offence taken as subjective. That is, the degree of offence is knowable only by each annotator. According to Sneddon (2020), we tend to give claims of offensiveness more credence than they deserve. In general, people tend to be more offended about topics that particularly matter to them, and these are impacted by one’s identity: A citizen of the USA is more likely to be offended by the burning of their national flag than a European. That is to say, when we are offended, we take the object of offence as a personal affront. This has material consequences when it comes to modelling *isms* as offensive.

This approach is often motivated by the desire to maintain minority voices within the annotation pool (Abercrombie et al., 2022), and in doing so, argues that disagreements are often the result the subjectivity of the annotation task.

3 Annotator Competency

Dataset labelling in NLP is typically performed by annotators recruited either as crowd-sourced workers (e.g. Abercrombie et al., 2023; Basile et al., 2019; Fersini et al., 2018), academics or students available to the researchers (e.g. Cercas Curry et al., 2021; Fanton et al., 2021; Jiang et al., 2022), or peo-

ple deemed to hold expertise in the phenomena (e.g. Talat, 2016; Vidgen et al., 2021; Zeinert et al., 2021). However, the Standpoint Theory (Harding, 1991) argues that annotators, can largely only be competent within their own lived experiences, regardless of training. Without lived experience, annotators may not be able to gain a full understanding of the *ism* under consideration. For instance, Larimore et al. (2021) found that white annotators were far less competent in identifying anti-Black racism than Black annotators. Annotator guidelines and labelling taxonomies, no matter how thoroughly and carefully constructed are not capable of adjusting for a lifetime of lived experience. It is not, therefore, inherent subjectivity within the task, but rather differences in annotator ability due to their personal standpoint that impact on annotators' ability to recognise whether hate speech or abuse is present. Sometimes even if an individual does recognise the target phenomenon, they may choose to ignore it for political reasons (Marable, 1995).

4 Towards a New Formulation of *Isms* as Cultural Formation of Societal Norms

Given our understanding of *isms* as culturally relative constructions and *offence* as an individually subjective concept, we propose that *isms* can best be understood as cultural formations of societal norms. That is, *isms* encode norms, which are inherently fuzzy at the border (Hall, 1997). When creating data for *isms*, researchers often work at the fuzzy borders of acceptability. In operating at these borders, and developing computational methods to draw them, research delineates what is acceptable from that which is not. While such borders are inherently messy, through an understanding of determining acceptability as cultural norms, we can refocus our attention towards the question of how such norms and borders should be drawn.

For instance, Douglas (1978) argues that determining what is 'dirt' is a cultural process which strengthens communities and builds community cohesion. That is, while encountering an offensive instance, i.e., an instance of sexism, can be destabilising to a community, the process with which the community makes a determination, and the determination itself, allows for the community to reify itself. This is particularly important as we can come to understand that *isms* are culturally defined objects, and identifying the borders of acceptability necessitates an ongoing negotiation with the communities in question (Thylstrup and Talat, 2020).

Within this formulation of *isms*, we can come to understand *isms* as distinct from *offence*. Thus, this formulation of *isms* provides space for both a cultural understanding of *isms* whilst making space for *offence* as an individual and subjective notion.

5 Conclusion: Implications for NLP

If, as we propose, the task of identifying *isms* in language is not subjective, we must conclude that annotator differences are irrelevant at the individual level for such tasks. Rather, they are symptoms of disagreement with the degree to which an *ism* causes offence to individual annotators.

At the group level, we must be careful that conflicting responses are not treated equally. If a minority of people with the necessary lived experience (e.g. to recognise misogyny) disagree with the majority who don't, then that matters. For example, Gordon et al. (2022) attempt to pick out the 'correct' minority perspectives from the wider pool of annotators for each instance to be analysed and Fleisig et al. (2023) specifically assume that the majority of annotators are likely to be 'wrong', that is, they will not recognise the target phenomenon.

Construction of the desired classification schema based on societal norms comes with its challenges. While prescriptivist annotation based on agreed societal norms may be desired, it can be difficult or even impossible to implement comprehensively in practice. One reason for this is that it is probably not possible to recruit annotators with the correct standpoint or competencies to recognise every instance—or indeed to know what those characteristics might be. Another is the nature of building classification schema. While a clearly defined, unambiguous, comprehensive and static *Aristotelian* classification scheme may be desired rather than *prototypical* classification,¹ it can be hard or even impossible to implement, and people generally resort to the latter (Bowker and Star, 2000, p. 61-62).

Despite this, we believe that it is vital that *isms* like misogyny and other hate and abuse not be constructed as individually subjective, but rather as culturally formed societal norms. While there may be much to gain from examining the responses of individual annotators to these tasks, NLP researchers should be careful not to conflate individual differences with inherent subjectivity of tasks.

¹The Aristotelian and prototypical classification paradigms (Bowker and Star, 2000) could be said to mirror Rottger et al.'s (2022) prescriptivist and descriptivist paradigms.

References

- Gavin Abercrombie, Valerio Basile, Sara Tonelli, Verena Rieser, and Alexandra Uma, editors. 2022. *Proceedings of the 1st Workshop on Perspectivist Approaches to NLP @LREC2022*. European Language Resources Association, Marseille, France.
- Gavin Abercrombie, Dirk Hovy, and Vinodkumar Prabhakaran. 2023. *Temporal and second language influence on intra-annotator agreement and stability in hate speech labelling*. In *Proceedings of the 17th Linguistic Annotation Workshop (LAW-XVII)*, pages 96–103, Toronto, Canada. Association for Computational Linguistics.
- Sohail Akhtar, Valerio Basile, and Viviana Patti. 2021. *Whose opinions matter? Perspective-aware models to identify opinions of hate speech victims in abusive language detection*.
- Cecilia Ovesdotter Alm. 2011. Subjective natural language problems: Motivations, applications, characterizations, and implications. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 107–112.
- Dina Almanea and Massimo Poesio. 2022. *ArMIS - the Arabic misogyny and sexism corpus with annotator subjective disagreements*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2282–2291, Marseille, France. European Language Resources Association.
- American Civil Liberties Union. 2023. *Mapping Attacks on LGBTQ Rights in U.S. State Legislatures*.
- Ginette Azcona, Antra Bhatt, Guillem Fortuny Fillo, Yongyi Min, Heather Page, and Sokunpanha You. 2023. *Progress on the Sustainable Development Goals: The gender snapshot 2023*. United Nations Entity for Gender Equality and the Empowerment of Women (UN Women) Department of Economic and Social Affairs (DESA).
- Maria Baghrmian. 2004. *Relativism*. Routledge.
- Eric Barendt. 2019. What is the harm of hate speech? *Ethical Theory and Moral Practice*, 22:539–553.
- Valerio Basile. 2020. It’s the end of the gold standard as we know it. on the impact of pre-aggregation on the evaluation of highly subjective tasks. In *Proceedings of The 19th International Conference of the Italian Association for Artificial Intelligence*, volume 2776, pages 31–40. CEUR-WS.
- Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. *SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in Twitter*. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Geoffrey C. Bowker and Susan Leigh Star. 2000. *Sorting Things Out: Classification and Its Consequences*. The MIT Press.
- Kingsley R Browne. 2006. Sex, power, and dominance: The evolutionary psychology of sexual harassment. *Managerial and Decision Economics*, 27(2-3):145–158.
- Amanda Cercas Curry, Gavin Abercrombie, and Verena Rieser. 2021. *ConvAbuse: Data, analysis, and benchmarks for nuanced abuse detection in conversational AI*. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7388–7403, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Bharathi Raja Chakravarthi, Ruba Priyadarshini, Thenmozhi Durairaj, John McCrae, Paul Buitelaar, Prasanna Kumaresan, and Rahul Ponnusamy. 2022. *Overview of the shared task on homophobia and transphobia detection in social media comments*. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 369–377, Dublin, Ireland. Association for Computational Linguistics.
- Thomas Davidson, Dana Warmesley, Michael Macy, and Ingmar Weber. 2017. *Automated hate speech detection and the problem of offensive language*. *Proceedings of the International AAAI Conference on Web and Social Media*, 11(1):512–515.
- Thiago Dias Oliva, Dennys Marcelo Antonialli, and Alessandra Gomes. 2021. *Fighting Hate Speech, Silencing Drag Queens? Artificial Intelligence in Content Moderation and Risks to LGBTQ Voices Online*. *Sexuality & Culture*, 25(2):700–732.
- Mary Douglas. 1978. *Purity and Danger: An Analysis of the Concepts of Pollution and Taboo*, repr edition. Routledge, London.
- EIGE. Sexism at work handbook. https://eige.europa.eu/publications-resources/toolkits-guides/sexism-at-work-handbook/part-1-understand/what-sexism?language_content_entity=en. Accessed June 20 2023.
- Naomi Ellemers. 2018. Gender stereotypes. *Annual review of psychology*, 69:275–298.
- Margherita Fanton, Helena Bonaldi, Serra Sinem Tekiroğlu, and Marco Guerini. 2021. *Human-in-the-loop for data collection: a multi-target counter narrative dataset to fight online hate speech*. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3226–3240, Online. Association for Computational Linguistics.
- Elisabetta Fersini, Paolo Rosso, and Maria Anzovino. 2018. Overview of the task on automatic misogyny identification at IberEval 2018. In *Iberval@ sepln*, volume 2150, pages 214–228.

- Nanna Thylstrup and Zeerak Talat. 2020. [Detecting ‘dirt’ and ‘toxicity’: Rethinking content moderation as pollution behaviour.](#)
- Francine Tougas, Rupert Brown, Ann M Beaton, and Stéphane Joly. 1995. Neosexism: Plus ça change, plus c’est pareil. *Personality and social psychology bulletin*, 21(8):842–849.
- Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2021. [Learning from the worst: Dynamically generated datasets to improve online hate detection.](#) In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1667–1682, Online. Association for Computational Linguistics.
- Frieda Wong and Richard Halgin. 2006. [The “Model Minority”: Bane or Blessing for Asian Americans?](#) *Journal of Multicultural Counseling and Development*, 34(1):38–49.
- Lucas Wright, Derek Ruths, Kelly P Dillon, Haji Mohammad Saleem, and Susan Benesch. 2017. [Vectors for counterspeech on Twitter.](#) In *Proceedings of the First Workshop on Abusive Language Online*, pages 57–62, Vancouver, BC, Canada. Association for Computational Linguistics.
- Philine Zeinert, Nanna Inie, and Leon Derczynski. 2021. [Annotating online misogyny.](#) In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3181–3197, Online. Association for Computational Linguistics.