

# Impact of heterogenous terminologies in study of sexism across disciplines

*Keywords: sexism; natural language processing; word embedding; bibliometrics; scientific terminologies*

## Extended Abstract

Sexism has long since prevailed in our society. But over the last few years, sexism is seen to be manifested in online social media platforms as well. Women who participate in internet, especially in online social media platforms have to face various forms of violence, predominantly text-based violence.[1] Time and again, online digital platforms have been seen to cultivate anti-women spaces and discourses which are propagated by the technological affordances in these platforms, and ultimately allowing further spread of sexism. Studies have found that “although the category of hate speech online covers both gender and the online environment, it is far too limited to include the diversity of experiences encountered by women.”[2] Furthermore, these online abuses have offline consequences, especially its psychological impact on the victim. As per the findings by Lewis et al.[3], they argue that “online abuse is most usefully conceived as a form of abuse or violence against women and girls, rather than as a form of communication.” To stop propagation of sexism in online spaces, it is essential that more academic studies are promoted across disciplines to understand the vast depths of sexism, both offline and online. Which further adds to the necessity of promoting interdisciplinary research, to culminate the expertise and knowledge from different disciplines in talking the global issue.

This project aims to analyze the study of sexism across disciplines, by looking at the main concepts and vocabulary used with discussing sexism, to underline the key differences in terms of approach and focus across fields. The manuscripts are scraped using the search query for the word “*sexism*” through the Semantic Scholar API, ranging through the years from 2013-2022. After collection, inclusion and exclusion criteria were used to narrow down the search results to include only the relevant papers for this study. For the sake of analysis, only the manuscripts which were published in English were taken. Thereafter, BERTopic model was used on the abstracts of all the paper results, for topic modeling to capture the contextual and semantic structure of the manuscripts and analyze the differences based on both the disciplines. Alongside using its core components: Bag-of-Words(BoW) representation and weighting with c-TF-IDF, additional representation models were used for further fine-tuning of the topic representations – a) by capturing the semantic relationship between keywords and topics using KeyBERTInspired[4]; b) leveraging better topic representation using PartOfSpeech(POS) tagging; and c) decreasing redundancy of keywords/keyphrases providing similar information and maximizing diversity within each abstract using Maximal Marginal Relevance(MMR). To capture the co-occurrence of most-used keywords/key phrases within each discipline, KeyBERT model with parameter ‘keyphrase\_ngram\_range’ set in the range of 1 to 3, was used (Figure 1a). For the purpose of initial analysis, the two distinctive disciplines, namely **Social Science (SS)** and **Computer Science (CS)** has been taken as the starting point. Fields like Engineering and Computer Science were taken under CS studies, while fields like Sociology, Political Science, Philosophy and Psychology were considered under SS studies. We find several differences between papers on sexism in SS and CS, as follows:

- a) **Terminologies:** The keywords associated to sexism differ a lot across disciplines. The lack of standardized terminologies is an issue for cross-disciplinary research, where different

words are used to describe similar concepts. While consistency of terms exists within disciplines, not so much between different disciplines.

- b) **Level of diversity in defining sexism:** Studies in SS have explored to the depths and subjective levels of sexism, alongside observing its intersectionality with other group characteristics, such as profession, age, region of origin and religion (Figure 2b) - building up from the theories of feminist and gender studies. With CS, there have been some studies considering sexism with the other group characters, such as identifying sexism and racism separately from the same data, but not so much in their co-occurrence (Figure 2a).

Owing to the fact that traditionally most studies on sexism stem from SS, a large imbalance in the volume of data between the disciplines of focus was expected, and indeed was reflected in the data gathered. This could account for an unavoidable data bias that impacted in generation of relevant topics for the disciplines. But to maintain a realistic analysis of studies within the last 10 years, creating a 1:1 ratio of the data between the two disciplines, was avoided. Interestingly, an unexpected observation which came up during the data fetching stage was the huge dominance of the field of Medicine in the search results generated. It was seen that the papers were focusing on occupational sexism in the field of medicine, alongside the topic of gender inequity in receiving health care access, services, etc. In fact, additional qualitative analysis of the papers between SS and CS portrayed a notable gap in the study of sexism between both – while CS focuses mostly on the online sexism, SS studies are seldom conducted on online sexism. In SS, most of the “scales” or measures of different forms of sexism tend to look at the context and analyze motivations behind sexist behavior – which mostly, if not always, refer to offline attitudes and behavior. While understanding contexts (and arguably intents) alongside content is immensely useful for developing in-depth inferences about sexism and its different forms, it becomes a hurdle to computationally implement or capture that in CS studies. Especially when considering perpetration of sexism in online digital platforms. Many of the CS studies struggle to capture the context of the text, based on the content. But when it comes to developing the intent of the same, it could often lead to inaccurate interpretation of results, due to an additional impact of the circumstance of the said situation and can be hard to determine in the online space. This is why detecting or measuring sexism can be termed as a subjective NLP task.

This project contributes to adding value in understanding *how* sexism has been studied across the disciplines over the past 10 years, to answer *why*, or rather what are the reasons that could be contributing to the existing separation between the different approaches to the subject. The findings from this project through methodological contributions of both quantitative and qualitative analysis, are important for researchers and non-academic organizations interested to counter sexism, both online and offline, in acknowledging the differences and promoting more interdisciplinary study. By doing so, this project hopes to contribute to both gender and digital equity, and potentially shape the future of interdisciplinary studies of sexism in significant ways.

## References

- [1] K. Barker and O. Jurasz, “Online Misogyny: A Challenge for Digital Feminism?,” *Journal of International Affairs Editorial Board*, vol. 72, no. 2, pp. 95–114, 2019, [Online]. Available: <https://www.jstor.org/stable/26760834>
- [2] D. Ging and E. Siapera, “Special issue on online misogyny,” *Feminist Media Studies*, vol. 18, no. 4, pp. 515–524, Jul. 2018, doi: 10.1080/14680777.2018.1447345.
- [3] R. Lewis, M. Rowe, and C. Wiper, “Online Abuse of Feminists as An Emerging form of Violence Against Women and Girls,” *CRIMIN*, p. azw073, Sep. 2016, doi: 10.1093/bjc/azw073.

- [4] “KeyBERT - BERTopic.”  
<https://maartengr.github.io/BERTopic/api/representation/keybert.html> (accessed Mar. 02, 2023).

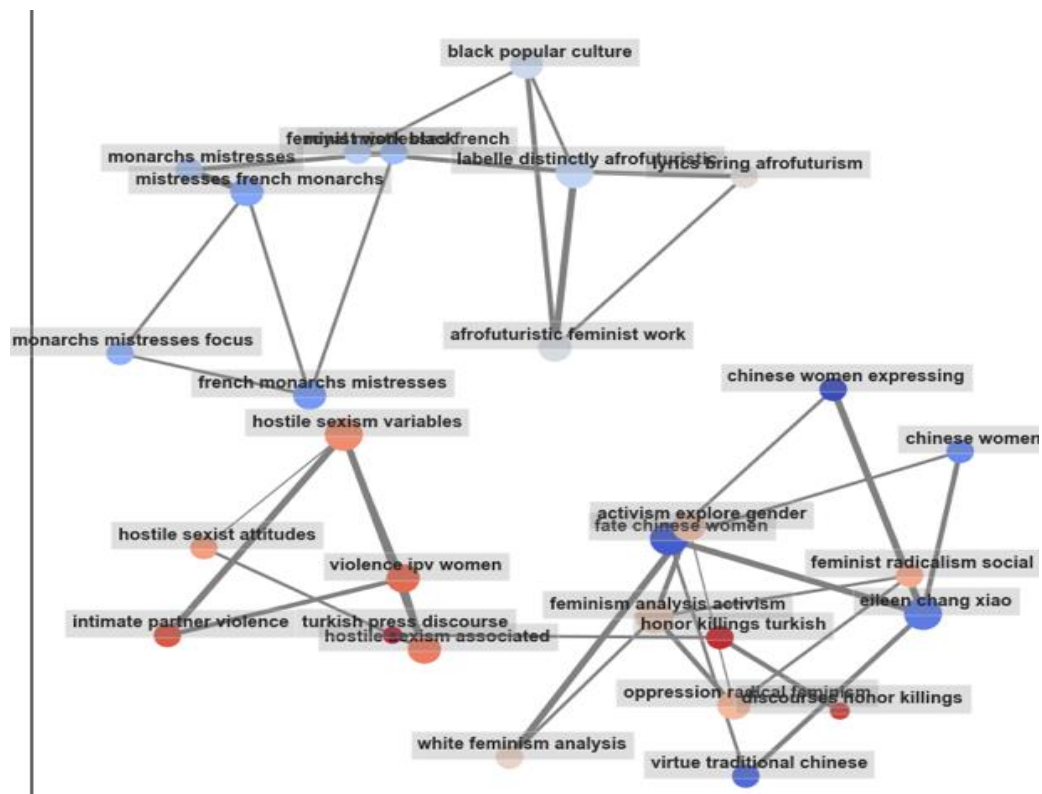


Figure 1(a) – Keyword co-occurrence network in Computer Science papers about sexism, generated using KeyBERT. Nodes represent common (1-3)-grams in CS papers on sexism, links represent co-occurrence of the same terms on the same paper. The width of the links represents the number of times the terms have co-occurred, while the size of the nodes represents how common those terms are in CS papers.

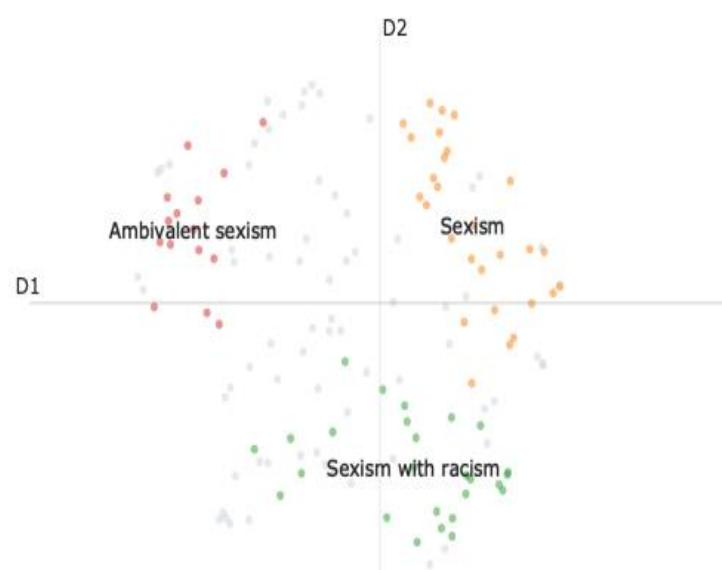


Figure 2(a) – Scatterplot representing Computer Science papers on sexism. Different colors represent topics identified using BERTopic.

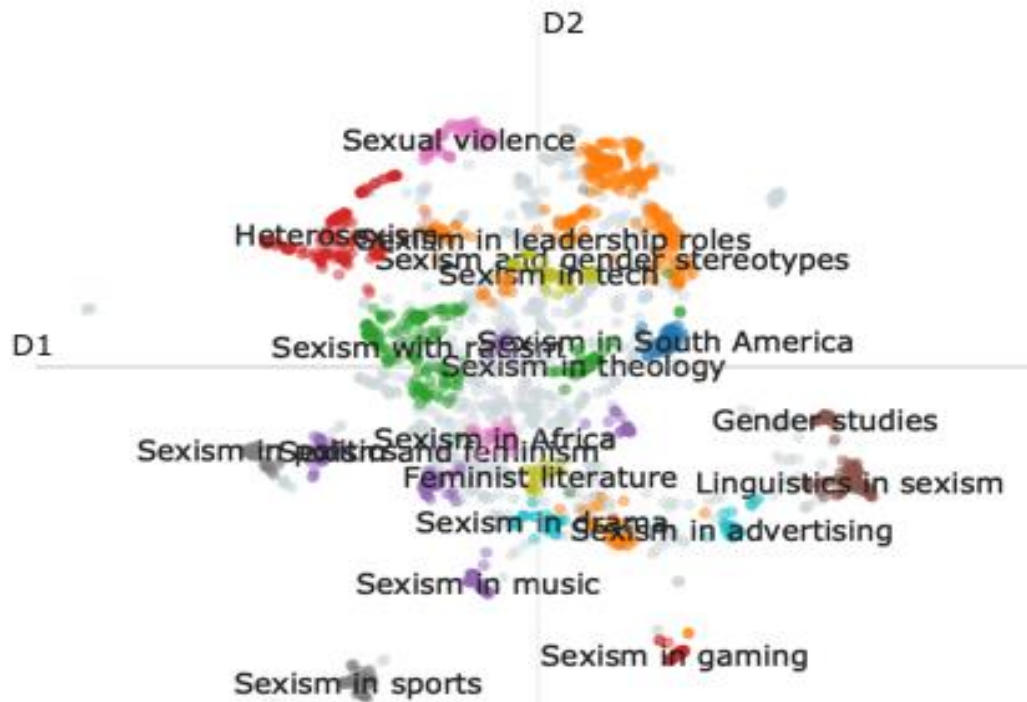


Figure 2(b) – Scatterplot representing Social Science papers on sexism. Different colors represent topics identified using BERTopic.