[1]     Does Double-Hard Debias Keep its Promises? - A Replication Study

[2]     Kristina Kobrock[1], Meike Korsten[1], & Sonja Börgerding[1]

[3]     [1] University of Osnabrück

[4]                         Author Note

[5]     The authors made the following contributions. Kristina Kobrock: müssen wir

[6] natürlich nicht ausfüllen; Meike Korsten: . . . ; Sonja Börgerding: . . . .

Does Double-Hard Debias Keep its Promises? - A Replication Study

Just edit the text! There are several websites that help with (r)markdown, e.g. https://rmarkdown.rstudio.com/. If you don't want to edit directly, but rather comment use: Cite using: Wang et al. (2020) said that... or simply blablabla (Wang et al., 2020). References must be added to r-references.bib file.

## Introduction

- shall we include an abstract?

### Bias in Embeddings.

- some text on bias in embeddings and why it's useful to do research in that field
- summary of related approaches (or put it in the discussion ?)

### Motivation.

- rough description of the specific study (Wang et al.) & motivation of replication attempt (maybe some general remarks on the need for replication studies)

## Implementation

- task description
- explanation of model and training choices

### Datasets and Preliminaries.

### Hard Debias.

### Double-Hard Debias.

26 **Evaluation**

27 - stick to the paper

28   **Baselines.**

29   **Evaluation of Debiasing Performance.**

30   **Debiasing in Downstream Applications.**

31   ***Coreference Resolution.***

32 - no replication possible, no code provided

33   **Debiasing at Embedding Level.**

34   ***The Word Embeddings Association Test (WEAT).***

35   ***Neighborhood Metric.***

36 - discuss that this is shady in the discussion part

37 **Analysis of Retaining Word Semantics**

38   ***Word Analogy.***

39   ***Concept Categorization.***

40 **Discussion**

41 - analysis of results and evaluation of performance evaluation

42 - ablation studies (not applicable)

43 - discuss the results and what could be (partly) replicated and what not

44 **Conclusion**

# References

Wang, T., Lin, X. V., Rajani, N. F., McCann, B., Ordonez, V., & Xiong, C. (2020). Double-hard debias: Tailoring word embeddings for gender bias mitigation. *Association for computational linguistics (acl).*