

How to Measure Bias in Automatic Speech Recognition Systems?

A Bias Metric Without a Reference Group

Tereza Ležovičová

EEMCS, Delft University of Technology, The Netherlands

25-06-2024

1 Introduction

Automatic speech recognition (ASR) is concerned with converting spoken language into text. However, as human speech is naturally highly variable, depending on context, region, or age, producing an objective transcription of the speech [2] is challenging. Research shows that ASR models show bias, e.g., between genders and races [4].

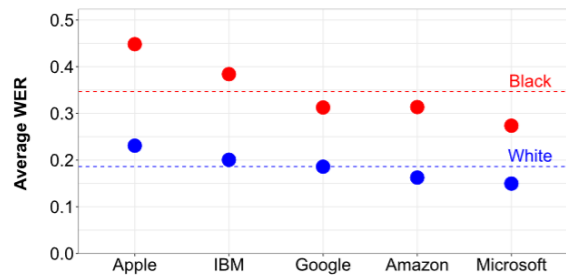


Figure: Average WER across multiple models [4]

Group fairness

Accepting similar outcomes across different speaker (demographic) groups [7]

Bias

A performance difference, specifically an error rate difference, between various speaker groups [6]

1.1 Bias metric

$$G2_{\text{reference group diff}} = b_{\text{spk}_g} - b_{\text{reference group}}$$

$G2_{\text{reference group diff}}$ refers to Bias metric

b refers to Base metric

- in our case Word Error Rate (WER)

2 Problem description

Existing bias metrics typically use a "reference group". The reference group serves as a baseline, i.e., its word error rate is compared to that of other groups to observe potential bias in the system.

Limitations of this approach:

- Binary decision outcomes**
 - The privileged group and the non-privileged group [5]
 - No insight into inter-group biases
- Standardizing reference group characteristics**
 - Ethical concern regarding the standardization of what constitutes a "normal" speech characteristic

Aim

Approach to measuring bias in ASR systems that eliminates the use of a reference group

3 Methodology

I evaluate these models on my bias metric:

- the E2E **Conformer** model trained on CGN dataset (*NoAug*, *SpAug*, *SpSpecAug*)
- Open AI's **Whisper** model (W_s , $W_s Ft_{cgn}$)

I use **Jasmin dataset** as my testing set with two speaking styles (Rd, HMI) and speaker groups:

- Dutch Children (DC)
- Dutch Teenagers (DT)
- Dutch Older Adults (DOA)
- Non-native Teenagers (NnT)
- Non-native Adults (NnA)

4 Final Design

$$b_{\text{avg}(MS)} = \frac{1}{G} \sum_{g \in G} WER_{g(MS)}$$

$$G2_{\text{avg log ratio}}(g(MS)) = -\ln \frac{b_{g(MS)}}{b_{\text{avg}(MS)}}$$

$$\text{Sum}_M = \sum_{g \in G} B(g_{MRd}) + B(g_{MHMI})$$

5 Limitations and Future Work

- Hard to estimate effectiveness of bias metric as there is no ground truth to bias
- Simulated environment where one or multiple groups are artificially biased [1]

6 Results

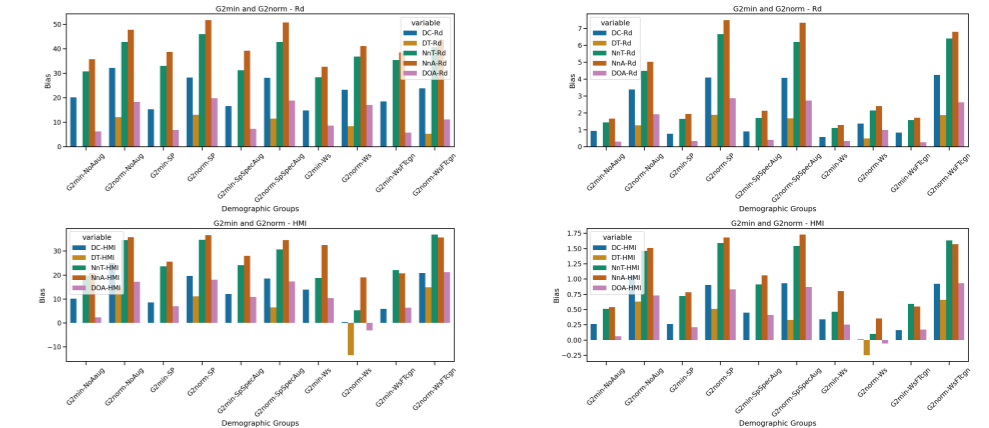


Figure: $G2_{\text{min diff}}$ and $G2_{\text{norm diff}}$ [6]

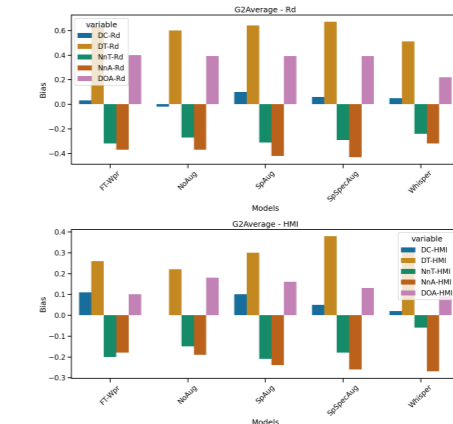


Figure: $G2_{\text{avg log ratio}}$ [3]

Figure: $G2_{\text{min reldiff}}$ and $G2_{\text{norm reldiff}}$ [6]

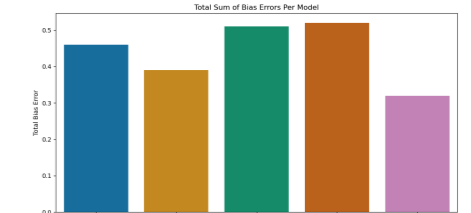


Figure: $\sum G2_{\text{avg log ratio}}$ [3, 1]

References

- A. Elobaid, N. Ramoly, L. Younes, S. Papadopoulos, E. Ntoutsis, and I. Kompatsiaris. "Sum of Group Error Differences: A Critical Examination of Bias Evaluation in Biometric Verification and a Dual-Metric Measure". In: *18th International Conference on Automatic Face and Gesture Recognition (FG)* (Apr. 2024). doi: 10.48550/arXiv.2404.15385.
- S. Feng, O. Kudina, B. M. Halpern, and O. Scharenborg. "Quantifying Bias in Automatic Speech Recognition". In: (Apr. 2021). doi: 10.48550/arXiv.2103.15122.
- W. Hutiri, T. Patel, A. Ding, and O. Scharenborg. "As Biased as You Measure: Methodological Pitfalls of Bias Evaluations in Speaker Verification Research". Submitted to *Interspeech 2024*. 2024.
- A. Koenecke et al. "Racial disparities in automated speech recognition". In: *Proceedings of the National Academy of Sciences* 117.14 (Apr. 2020), pp. 7684–7689.
- K. Lum, Y. Zhang, and A. Bower. "De-biasing "bias" measurement". In: *June 2022*, pp. 379–389. doi: 10.1145/3531146.3533105.
- T. Patel, W. Hutiri, A. Ding, and O. Scharenborg. "How to Evaluate Automatic Speech Recognition: Comparing Different Performance and Bias Measures". *Work in progress*. 2024.
- M.-H. Yeh, B. Metevier, A. Hoag, and P. Thomas. "Analyzing the Relationship Between Difference and Ratio-Based Fairness Metrics". In: (June 2024), pp. 518–528. doi: 10.1145/3630106.3658922.