
Reproducing Results for Crossing the Line: Where do Demographic Variables Fit into Humor Detection?

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2

3 **Scope of Reproducibility**

4 Within the original experiment two groups of annotators of size 10 each in age groups 18-25 and 55-70 were relied on
5 to generate the metrics displayed as part of the results. The scope of our study was limited to instances from The Short
6 Jokes dataset on Kaggle with token sizes between 11 and 16 for annotators 21 in number divided into demographic bins
7 by gender: male, female, non-binary, gender being the chief source of demographic diversity under study.

8 **Methodology**

9 The paper was based on gathering direct human response over field values of humorous and/or offensive by presenting
10 short jokes from a varied set of humor genres and subgenres to a diverse audience over different demographic categories
11 segmented by age (18-25, 26-40, 40-55, 56-70), educational qualification as an index of socio-economic status (High
12 School, Undergraduate, Postgraduate), and gender (Male, Female, Non-binary).

13 The same methodology as the original paper of using binary classification ((humorous 1, non-humorous 0), (offensive 1,
14 non-offensive 0)) and adding values demographic bin-wise in studying findings was used.

15 **Results**

16 It was found that inter-annotator agreement was higher, when categorized using demographic data, in this case gender,
17 as exemplified by the instance of a subset of jokes being identified with the keyword sexist were found humorous and
18 offensive by female annotators, with male members ranking jokes in question reporting it as simply offensive.

19 **What was easy**

20 The easy part of the reproduction study was to collect data with regards to short, English jokes representing various
21 genres of humor with a diversity in addressed audiences, expressed sentiment and degrees of simplicity.

22 **What was difficult**

23 The tougher part of recreating the study and determining results was ensuring diversity in representatives of survey
24 response group without insensitive treatment of interviewed people or any prejudice in choosing response.

25 **Communication with original authors**

26 Context was gathered from the completeness of submitted paper and no one to one communication was established
27 with the author or the purpose of this reproducibility study at the time.

28 **Conclusion**

29 The scale of operations could be significantly increased by introducing a more automated form of response collection
30 such as gauging responses to sample inputs on a user forum kept live for engagement and recording responses as per
31 categorized demographic bins having obtained consent to do so in a transparent manner with users agreeing to provide
32 such information for research purposes.

33 **References**

34 Meaney, J. A. "Crossing the Line: Where do Demographic Variables Fit into Humor Detection?" ACL (2020).
35 Paula Cristina Teixeira Fortuna. 2017. Automatic detection of hate speech in text: an overview of the topic and dataset
36 annotation with hierarchical classes.