

1 We thank the reviewers for their consideration of our paper and for their feedback. The consensus appears to be that this
2 is a generally well written paper, exploring a “*potentially impressive*” (R2), “*useful*” (R1, R3) and “*novel*” (R1, R2)
3 connection between fairness and disentanglement with claims backed up by substantial empirical evidence (R1, R2,
4 R3). There appears to be one major concern and 4 minor questions/suggestions, which we kindly address below.

5 **Main concern of R2: use of ground-truth factors to compute disentanglement scores**

6 The main concern of R2 seems to be that the paper relies on “*disentanglement scores, which are computed based on*
7 *ground-truth factors of variations*“ but that these factors “*are generally not available in practice*”. In our view, there
8 appears to be a misunderstanding with regards to the exact goals and contributions of this paper:

- 9 • **Motivation: usefulness of disentanglement notions vs new method.** The main goal of this paper is to *evaluate the*
10 *usefulness* of disentangled representations (in terms of leading to fair predictions), rather than proposing new methods
11 for *learning* disentangled representations. This distinction is critical: While we do train various methods for the
12 purpose of this study, we do not make any assumptions about the *feasibility* of (and *methods* for) learning disentangled
13 representations in the absence of ground-truth factors, and our results have implications for the supervised, semi-
14 supervised, and unsupervised settings. This paper provides evidence that, if one were to find a reliable method for
15 disentanglement (in terms of current metrics), predictions on representations of such an approach would be more fair.
- 16 • **Relevance: validate the motivation of >15 recent ML papers.** Recently, numerous papers have been concerned
17 with learning disentangled representations [7, 8, 12, 17, 18, 20, 31, 32, 33, 34, 43, 50, 57, 61, 66, 74, 76, 82, 84].
18 The key motivation (but also assumption) of these works is that current notions of disentanglement (MIG, DCI, *etc.*)
19 are desirable. Until now there has been little empirical evidence verifying this. In particular, the study of Locatello
20 et al., ICML 2019 was inconclusive in this regard, see their Section 5.5. This paper fills this gap by investigating
21 whether disentanglement is useful for improving fairness of downstream predictions. Our results are novel and
22 provide motivation for further research into (a) disentanglement methods and (b) their application for ML fairness.
- 23 • **We present a heuristic to select fair representations.** As described in Section 4.2, as a by-product of our investiga-
24 tion, we noticed that downstream predictive performance may be used as a way to select fair representations among the
25 representations we trained using SOTA unsupervised disentanglement methods. We argue that this might be an interest-
26 ing heuristic as it only requires labels for the single downstream prediction task but not for the representation learning.

27 **Other concerns:**

28 **R1-R2: Motivation for adjusted metrics in Section 4.2.**

29 We compute the adjusted metrics to answer the question “**Given two representations with the same downstream**
30 **performance, is the more disentangled one also more fair?**”. The key difficulty is that for a given representation
31 there may not be other ones with exactly the same downstream performance. Hence, we compute these adjusted
32 metrics which intuitively measure how much fairer (more disentangled) a representation is compared to an average
33 representation with the same downstream performance. To compute the average fairness (disentanglement) for a
34 given downstream performance, we use a nearest neighbor regressor as a robust non-parametric 1D regression model.

35 **R1: Chain of arguments in “How do we identify fair models?” not fully clear. + Is the argument by R1 correct?**

36 **The understanding of the reviewer is correct.** We will add the suggested graph “fairness <- disentangling ->
37 accuracy” and revise the manuscript such that this argument is explained more concisely.

38 **R2: Theorem 1.**

39 **Entangled representations affect the fairness of downstream classifiers.** The theorem proves a point: the Bayes
40 optimal classifier does not imply fairness on an entangled representation. The proof is by counterexample. We
41 further support this claim with Figure 2, showing that this is a practical issue that arises on trained classifiers and
42 not simply a theoretical corner case. If we can disentangle target and sensitive variables then the classifier will be
43 fair, see lines 145-155. In the particular counterexample of the theorem, there is a trade-off between fairness and
44 accuracy as is common in the fairness literature. This is discussed in lines 135-138 and in the footnote.

45 **R3: Demographic parity.**

46 **We will add more motivation and examples,** in particular relating to previous work (e.g. [9, 92]) and to relevant
47 settings, for example to settings where sensitive variables may not be recorded due to privacy reasons.

48 *Clarity, title, and other comments regarding the presentations of our results.*

49 **We thank the reviewers for their suggestions to improve the presentation.** We will use their feedback to revise the
50 manuscript. In particular, as suggested, we will change the title to “On the Fairness of Disentangled Representations”,
51 improve the discussion on section 4.2 and add a discussion on non-VAE based methods such as ALI. By gradient
52 boosted trees (GBT10000) we mean gradient boosting of decision trees [Friedman, 1999] trained on 10000 examples.