

1 We thank all the reviewers for the time devoted to provide thoughtful comments.

2 **[Reviewer # 1, Estimating information-theoretic distances]** First, we would like to thank the reviewer for the
3 comprehensive and accurate summary of our work. We are happy that the reviewer found our results to be novel and
4 useful. We agree with the reviewer that estimating the distance between the learned representations are intractable,
5 since the sample complexity of estimating Shannon entropy, mutual information and related concepts are exponential in
6 the dimension of the representation space. That being said, one alternative way to do so is to consider the variational
7 representations of f -divergence and use rich parametrized function class (e.g., neural networks) to approximate these
8 distances. For example, recent work [2] on estimating mutual information has empirically shown that such approach
9 often leads to better estimation result than classic approaches based on nonparametric density estimation. From this
10 perspective, the \mathcal{H} -divergence in Section 3.2 actually serves as a relaxation of the total variation distance and it equals
11 TVD when \mathcal{H} contains all the measurable functions. Hence the lower bound in Proposition 3.1 gives us a practical way
12 to estimate a proxy of TVD in terms of sum of Type-I and Type-II errors in distinguishing group memberships.

13 **[Reviewer # 1, Total variation in Theorem 3.3]** The total variation in Theorem 3.3 is w.r.t the input distributions
14 across groups, i.e., $\mathcal{D}_0(X)$ and $\mathcal{D}_1(X)$. In the remark we use “the distance of representation” to mean “the distance of
15 input distributions”. We will clarify this sentence in our final version to avoid such confusion.

16 **[Reviewer # 1, Comparisons with existing work]** The results in this paper are distinct from the results in [1].
17 Specifically, the trade-off given in [1] (Proposition 8) is in terms of the fairness frontier function under the context
18 of cost-sensitive loss. Roughly speaking, it shows that if the two decision functions are dissimilar to each other, the
19 fairness constraint will not harm too much on the target utility. As a comparison, our results (Theorem 3.1 and 3.2)
20 directly give lower bounds on the sum of errors across groups in terms of the difference in base rates as well as the
21 distance of representations. Our results are also different from those in Madras and Zhang et al.: they gave an upper
22 bound on the demographic parity gap in terms of the loss incurred by an adversary (Theorem 5.1), while ours are about
23 lower bounds on the errors of the target task.

24 **[Reviewer # 1, Other questions]** We use the notation $\mathcal{P} \ll \mathcal{Q}$ to mean that distribution \mathcal{P} is absolutely continuous
25 w.r.t. distribution \mathcal{Q} , i.e., for any measurable event E , if $\mathcal{P}(E) > 0$, then we must have $\mathcal{Q}(E) > 0$ as well. The
26 generator function of KL divergence is indeed $f(t) = t \log t$, and the generator function of the inverse KL divergence
27 is $f(t) = -\log t$. Having identical joint distributions implies that the optimal decision functions are the same across
28 groups, but not the other way around. We also add one more experimental result with $\lambda = 50.0$, and the result is listed
as follows. Compared with the existing results in Table 2, we can see a consistent trend.

	$\text{Err}_{\mathcal{D}}$	$\text{Err}_{\mathcal{D}_0} + \text{Err}_{\mathcal{D}_1}$	$ \text{Err}_{\mathcal{D}_0} - \text{Err}_{\mathcal{D}_1} $	$d_{\text{TV}}(\mathcal{D}_0(\hat{Y}), \mathcal{D}_1(\hat{Y}))$
AdvDebias, $\rho = 50.0$	0.201	0.360	0.112	0.028

29

30 **[Reviewer # 3]** We are happy that the reviewer found our paper to be interesting, theoretically sound and well-written.
31 As stated in the last sentence of the conclusion section, our lower bound naturally implies an algorithm based on
32 instance-reweighting to balance the base rates during fair representation learning. However, the detailed design, analysis
33 and empirical validation of such an algorithm is beyond the scope of current paper. Given that nowadays there are
34 more than tens paper on proposing new algorithms to achieve fairness every year, we believe it would be nice to have a
35 theoretical paper with novel analysis techniques and results to study the fundamental limit of such algorithms. Although
36 it is clear that fairness will compromise utility, before this paper it is still unknown to what extent will it, and how is it
37 related to the difference in terms of base rates across groups. From this perspective, we believe our work is a timely
38 paper that answers the above questions quantitatively. As pointed out by Reviewer 1, our analysis technique using Liese
39 and Vadja lemma is novel and useful. This is of independent interest and we expect its applicability in a broader context.

40 **[Reviewer # 4]** We would like to thank Reviewer 4 for the encouraging comments. As explained in Theorem 2.1.,
41 Chouldechova and Kleinberg et al. mainly proved that positive rate parity and predictive value parity are in general
42 incompatible. This is an impossibility result between two different notions of fairness. As a comparison, we mainly
43 focus on trade-off between utility and fairness. Furthermore, our Theorem 3.1 is a quantitative result in the sense that
44 it not only gives the impossibility statement when base rates are different, but also gives a lower bound on the error
45 that will be incurred by *any* algorithm. Techniques based on instance-reweighting helps to decrease the difference in
46 base rates, and hence we would expect it to help decrease the lower bound as well. This means that we would incur
47 less drop of utility when learning fair representations. Our current technique does not extend to the definition of equal
48 opportunity, and collecting additional data will not help.

49 [1]. The Cost of Fairness in Binary Classification. Menon and Williamson, FAT* 2018.

50 [2]. MINE: Mutual Information Neural Estimation. Belghazi, ICML 2018.