

1 Thanks for the reviews! We take your comments to heart and will make all of the small changes suggested. Here we
2 address the more major concerns.

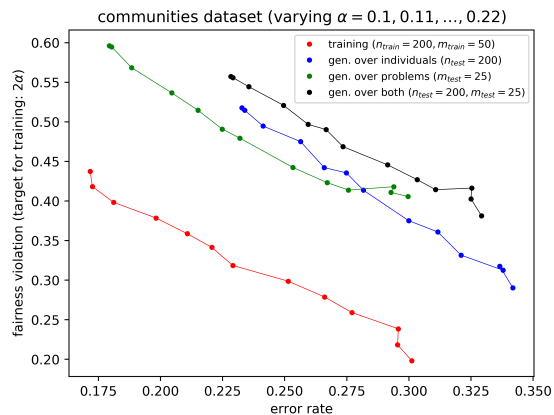
3 **Experiments** First, we agree with reviewer 3 that the main contribution of our paper is conceptual and theoretical: in
4 particular, identifying randomness over the problem distribution as something that can be leveraged to give an individual
5 level guarantee, and then deriving the algorithms and generalization bounds needed to actually realize this guarantee.
6 The purpose of the experiments is modest: to support the theory. Because of this, we chose to illustrate experimentally
7 the one thing that is not actually guaranteed by the theory: namely, the convergence of the algorithm. The theorem we
8 prove only guarantees that the algorithm converges *if we have a cost sensitive classification oracle*. In practice, we
9 do not: we use a regression heuristic — and so we chose the experiment we did to confirm that we get convergence,
10 despite the gap between theory and practice. Reviewer 3 is correct that prior work — in particular, [Kearns et al.] —
11 has shown with similar experiments that *other* oracle efficient algorithms converge in practice with heuristics. But that
12 does not imply that ours will, since our algorithm differs both in the specifics of the dynamics (fictitious play in [Kearns
13 et al.], no-regret vs. best response in our paper) and in the particular problems we are asking the classification oracles to
14 solve. In short, because the framework of “oracle efficiency” leaves a gap between theory and practice, we think of it as
15 good practice to accompany a proof of oracle efficiency with an empirical validation of convergence whenever possible.

16 That being said, Reviewer 3 is of course correct that empirically examining generalization is an interesting thing to do:
17 we didn’t do it in the submission simply because we thought it was *less* interesting than convergence (generalization is
18 guaranteed by our theorems, without any heuristic assumptions) and were mindful of space constraints. But we have
19 now quickly adapted our code to investigate our generalization error, both across data points and across problems. The
20 plot is below:

21 To be consistent with our paper, we trained on exactly
22 the same subset of Communities and Crime that we did
23 in our paper ($n = 200$ datapoints, $m = 50$ problems
24 (selected features from the dataset)). Thus the curve la-
25 belled “training” is the same as the reported in-sample
26 results in our paper. We used a fresh holdout consisting
27 of $n = 200$ datapoints, and $m = 25$ problems (features
28 from the dataset that weren’t previously used) to evaluate
29 our generalization performance over both problems and
30 data points, in terms of both accuracy and fairness viola-
31 tion. Two things stand out:

32 1. As predicted by the theory, our test curves track our
33 training curves, but with higher error and unfairness. In
34 particular, the ordering of the models on the Pareto fron-
35 tier is the same in testing as in training, meaning that the
36 training curve can indeed be used to manage the trade-off
37 out-of-sample as well.

38 2. The gap in error is substantially smaller than would be
39 predicted by our theory: since our training set is so small, our theoretical guarantees are vacuous, but all points plotted
40 in our test Pareto curves are non-trivial in terms of both accuracy and fairness. Presumably the gap in error would
41 narrow on larger training sets. If the paper is accepted and the reviewers feel that generalization experiments should be
42 included, we will make room to do so (and re-run the above experiment on a larger training set, which we can easily do
43 with the luxury of time). Similarly, if the reviewers still feel that the experiments detract from the theory, we are also
44 willing to relegate all experiments to the supplemental material and focus on the main contributions in the body.



45 Other Conceptual Questions

46 **Reviewer 1:** The group-fairness proposal of [Kearns et al.] indeed mitigates the “gerrymandering” concern we cite
47 from their work, but ultimately does not eliminate it because it remains a group-level constraint that only holds for
48 groups that are sufficiently large. With our approach, we are able to reduce these groups to size 1, which makes it an
49 individual level constraint. (We agree this is different than “individual fairness” in Dwork et al. – we use the term more
50 broadly, and will clarify).

51 **Reviewer 3:** We agree that asking for parity of error statistics is not always the right approach, and that our techniques
52 generalize to asking for upper bounds on error rates. We consider equalizing error rates mostly as a canonical example
53 of a popular approach; we will clarify. We will also further emphasize the limitation of needing to be able to observe
54 labels for different problems on the same training set.