



We thank all reviewers for their helpful and detailed comments! We have addressed the issue of dual submission in detail in the rebuttal of paper #6290. As R1 notes and we further elucidate, the problem setting, algorithm specifics, and use-case scenarios of the two papers are different and independent – **model bias** of a **pretrained model** for downstream **Monte Carlo evaluation** here vs. **data bias** during **weakly-supervised learning** for **fair data generation** in #6290.

- 6 • **R1: Support of the generative distribution p_θ and p .** Our meta-algorithm takes as input a learned model p_θ and p so
7 satisfying the support assumption is tied to the *training* of p_θ (which we do not consider in this work). Nevertheless for
8 a likelihood-based model, the support assumption can be empirically verified via evaluating p_θ on held-out data. The
9 assumption holds true for most variants of VAEs, flows, and autoregressive which have full support by design. We also
10 consider a more general case where we have only sample access to both p_θ and p , where estimating the support is a
11 computationally hard problem (related to estimating the entropy of arbitrary distribution via samples).
12 To address issues related to the estimation of importance weights via a learned classifier, tricks such as perturbations via
13 small, random Gaussian noise (which has full support), regularization (dropout, early stopping etc.) during training
14 (L306-307), as well as post-processing schemes (L135-143) can be applied. Empirically, we find self-normalization
15 along with early stopping during training (based on validation data) to be sufficient for ensuring good downstream
16 performance for various generative models (GANs, autoregressive models) and modalities considered in this work.
- 17 • **R1: Defining and measuring the bias introduced by p_θ .** In this work, bias is defined w.r.t. any function f defined over
18 the data domain. Given p_{data} , p_θ and f , the bias is defined as the difference in the expected value of f with respect to
19 p_{data} and p_θ (Footnote 1, Page 1). When p_{data} and p_θ are not known directly, the bias can be estimated empirically via
20 Monte Carlo using a sufficiently large number of samples from p_θ and p_{data} e.g., as shown in Table 1 and Appendix B.
- 21 • **R1: High variance in importance reweighting.** As with other applications of importance weighting, the extent of and
22 solutions to the high variance issue are empirically motivated. They could introduce a bias (e.g., clipping) but reduce
23 variance more favorably in the tradeoff. In our setting, the primary limitation was that the estimated importance weights
24 could all be small due to artifacts in the generations that were easy to detect via the binary classifier. While we found
25 self-normalization to be most effective, we note in L142 that schemes for post-processing importance weights could be
26 potentially combined, e.g., self-normalized weights could be clipped when variance is a larger issue.
- 27 • **R1: Choosing the clipping threshold β .** We consider β as a validation hyperparameter with values in $\{0.001, 0.01,$
28 $0.01, 1\}$ chosen to maximally reduce the bias in Monte Carlo evaluation of a downstream function of interest.
- 29 • **R2: Intuition and guidelines for design choices in L135-143.** Self-normalization is applied only for the generated
30 samples (i.e., those that contribute to bias in Monte Carlo evaluation). Like with other applications, the usage is
31 empirically driven. Generative models tend to produce artifacts that are easy to detect via classifiers and hence,
32 the estimated importance weights are very small ($\ll 1$). In all our experiments, self-normalization was essential to
33 circumvent this issue (see expts. in Tables 4, 5 in Appendix where self-normalization leads to a 53% improvement in
34 mean squared error over vanilla importance weighting). It is hyperparameter free and easy to apply. If variance is high,
35 the range of the weights can be restricted via clipping or flattening with hyperparameters β, α tuned on validation set.
- 36 • **R2: Data split for reference scores in L168.** Yes, the split is 50-50.
- 37 • **R2: Running procedure in [45] for long.** Yes, ignoring the high computational requirements of [45] and the fact that
38 the upper bound for rejection sampling is a heuristic estimate, the procedure in [45] could achieve the same effect as the
39 proposed importance weighting approach.
- 40 • **R2, R3: Calibration.** We believe the default calibration behavior is largely due to the fact that our **binary** classifiers
41 distinguishing real and fake data do not require very complex neural networks architectures and training tricks that lead
42 to miscalibration for **multi-class classification**. As shown by Niculescu-Mizil & Caruana (2005), shallow networks are
43 well-calibrated and Guo et al. (2017) further argue that a major reason for miscalibration is the use of a softmax loss
44 typical for multi-class problems. Top-left figure shows example calibration curves for the experiment in 5.1.
- 45 • **R3: Interaction of post-hoc normalization schemes with calibration.** While calibration is necessary for a sound
46 density ratio estimation procedure, the utility of the derived importance weights for downstream tasks depends on the
47 underlying expectation of interest. These expectations are evaluated with finite samples and hence, the asymptotic
48 properties of importance weighting (e.g., unbiasedness) are traded off for improved downstream performance using
49 self-normalization and other post-processing schemes.
- 50 • **R3: Domain adaptation.** We clarify that we are considering the task of multi-class classification and not domain
51 adaptation (L179-181). As we note in L182-183, the Omniglot dataset is a particularly relevant test bed for data
52 augmentation since there are a large number of classes and a few number of training examples per class. We will
53 consider other related scenarios in future version!
- 54 • **R3: $D_g + \text{LFIW}$ vs. D_g .** Note that this experiment does not only involve Monte Carlo evaluation of a supervised loss
55 but also optimization via gradient methods. In the absence of real data D_{cl} , the classifier training is dominated by D_g
56 and correcting the bias in the dataset via LFIW towards an unseen dataset (D_{cl}) can potentially have limited gains.
- 57 • **R3: Modes getting closer in Fig 1.** As modes get closer, the importance weights will approach 1 (and the class
58 probabilities will approach 0.5) since the mismatch in generative model and data distributions will accordingly decrease.