# Fair Generative Modeling via Weak Supervision

Kristy Choi   Aditya Grover   Trisha Singh  Rui Shu  Stefano Ermon

ICML  2020

# Motivation



Figure 1. Samples from a baseline BigGAN that reflect the gender bias underlying the true data distribution in CelebA. All faces above the orange line (67%) are classified as female, while the rest are labeled as male (33%).

biased dataset will result that generative model have biased performance
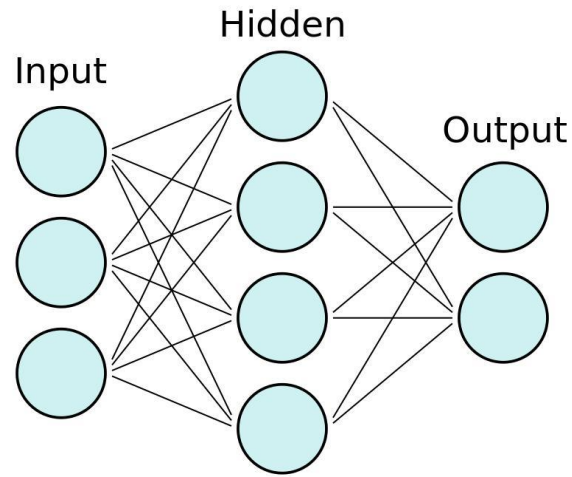
# Setup

two dataset :

$\mathcal{D}_{\mathrm{ref}}$ ⟶ A small unbiased dataset

$\mathcal{D}_{\mathrm{bias}}$ ⟶ A large biased dataset

Goal:  generated data Pdata $\qquad p_{\mathrm{data}} = \bar{p}_{\mathrm{ref}}$

# Method: Importance Reweighting

binary classification

Input
Hidden
Output

$Y = 1 \longrightarrow \mathcal{D}_{\text{ref}}$

$Y = 0 \longrightarrow \mathcal{D}_{\text{bias}}$

Loss:

$$NCE(c) := \frac{1}{\gamma + 1} \mathbb{E}_{p_{\text{ref}}(\mathbf{x})}[\log c(Y = 1|\mathbf{x})]$$
$$+ \frac{\gamma}{\gamma + 1} \mathbb{E}_{p_{\text{bias}}(\mathbf{x})}[\log c(Y = 0|\mathbf{x})].$$

Reweighting:

$$w(\mathbf{x}) = \frac{p_{\text{ref}}(\mathbf{x})}{p_{\text{bias}}(\mathbf{x})} = \gamma \frac{c^*(Y = 1|x)}{1 - c^*(Y = 1|x)} \qquad (5)$$
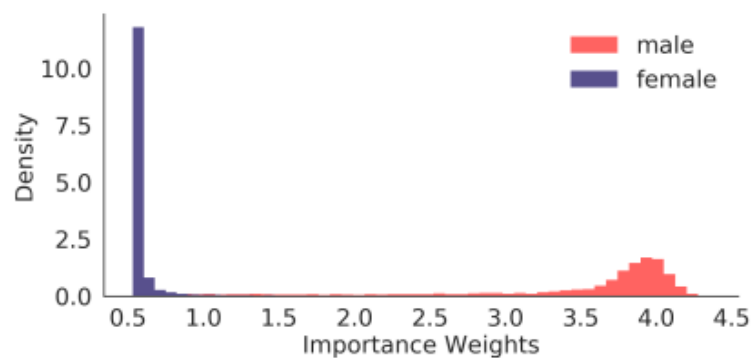
**Algorithm 1** Learning Fair Generative Models

**Input:** $\mathcal{D}_{\text{bias}}, \mathcal{D}_{\text{ref}}$, Classifier and Generative Model Architectures & Hyperparameters
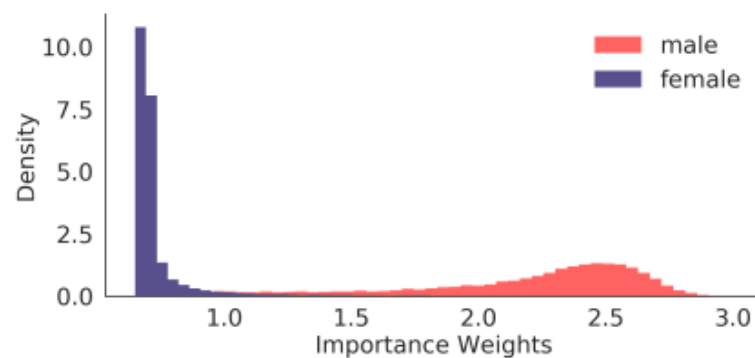
**Output:** Generative Model Parameters $\theta$

1: ▷ Phase 1: Estimate importance weights
2: Learn binary classifier $c$ for distinguishing $(\mathcal{D}_{\text{bias}}, Y = 0)$ vs. $(\mathcal{D}_{\text{ref}}, Y = 1)$
3: Estimate importance weight $\hat{w}(\mathbf{x}) \leftarrow \frac{c(Y=1|\mathbf{x})}{c(Y=0|\mathbf{x})}$ for all $\mathbf{x} \in \mathcal{D}_{\text{bias}}$ (using Eq. 5)
4: Set importance weight $\hat{w}(\mathbf{x}) \leftarrow 1$ for all $\mathbf{x} \in \mathcal{D}_{\text{ref}}$
5:
6: ▷ Phase 2: Minibatch gradient descent on $\theta$ based on weighted loss
7: Initialize model parameters $\theta$ at random
8: Set full dataset $\mathcal{D} \leftarrow \mathcal{D}_{\text{bias}} \cup \mathcal{D}_{\text{ref}}$
9: **while** training **do**
10:     Sample a batch of points $B$ from $\mathcal{D}$ at random
11:     Set loss $\mathcal{L}(\theta; \mathcal{D}) \leftarrow \frac{1}{|B|} \sum_{\mathbf{x}_i \in B} \hat{w}(\mathbf{x}_i) \ell(\mathbf{x}_i, \theta)$
12:     Estimate gradients $\nabla_\theta \mathcal{L}(\theta; \mathcal{D})$ and update parameters $\theta$ based on optimizer update rule
13: **end while**
14: **return** $\theta$

$$w(\mathbf{x}) = \frac{p_{\text{ref}}(\mathbf{x})}{p_{\text{bias}}(\mathbf{x})} = \gamma \frac{c^*(Y = 1|x)}{1 - c^*(Y = 1|x)} \quad (5)$$

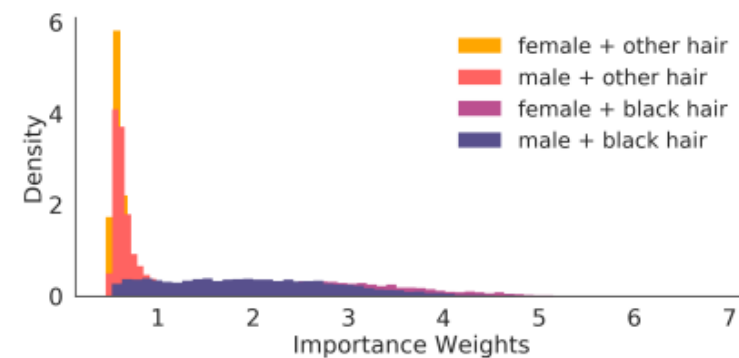# Experiments



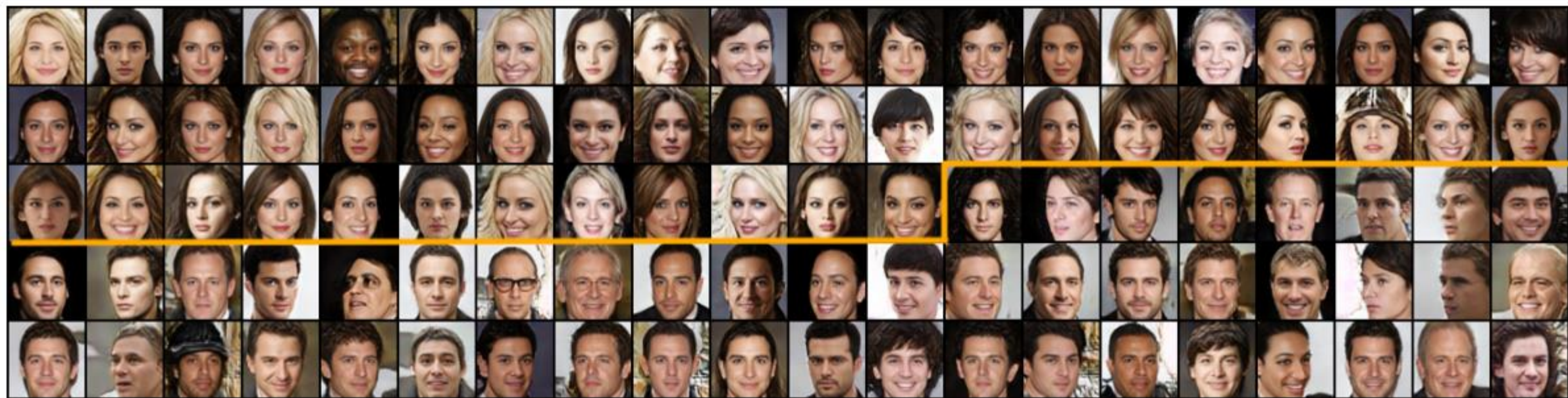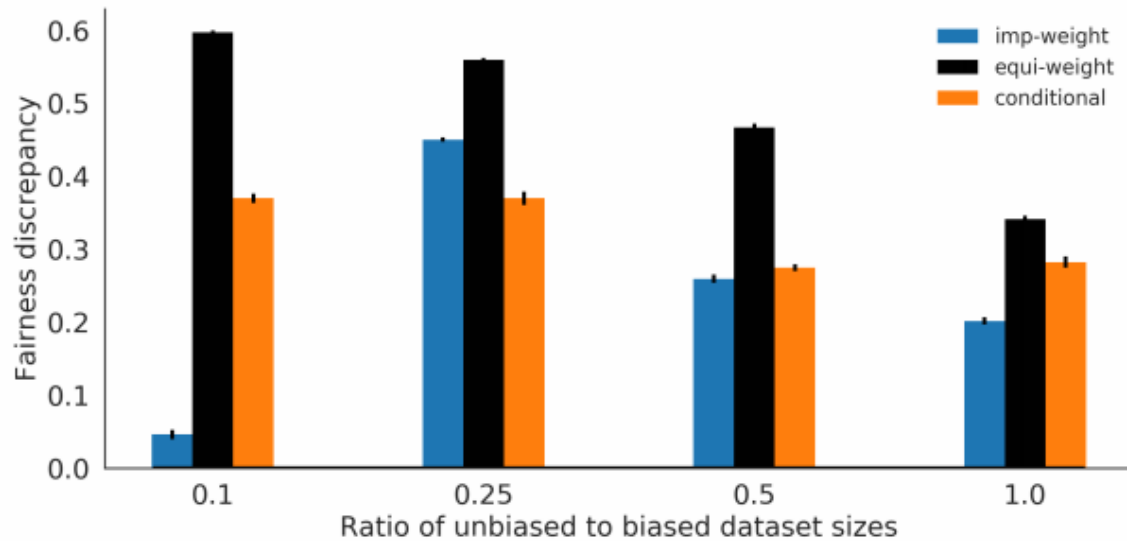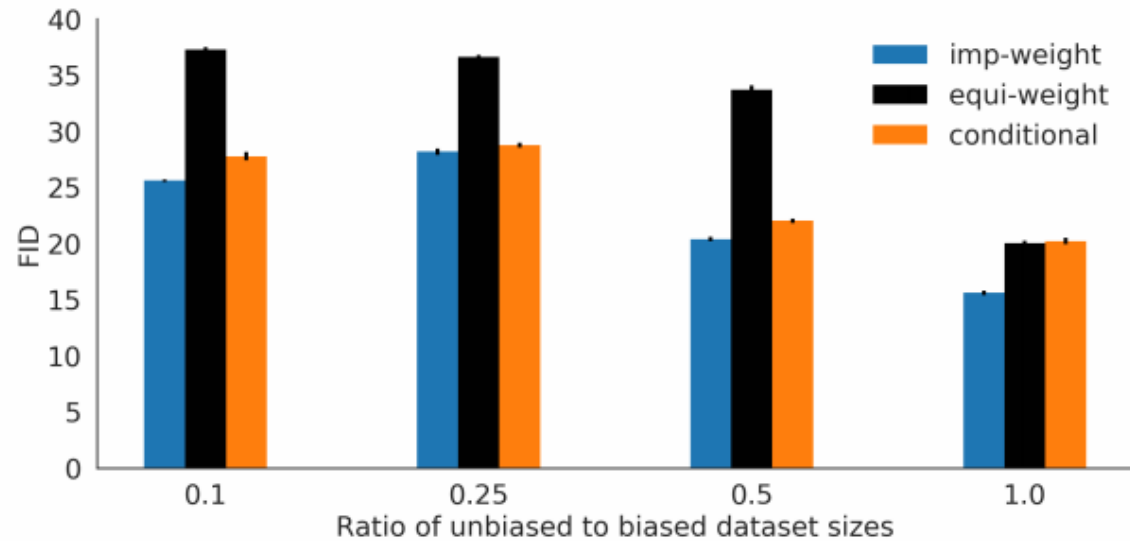(a) single, `bias=0.9`   (b) single, `bias=0.8`   (c) multi

*Figure 2.* Distribution of importance weights for different latent subgroups. On average, The underrepresented subgroups are upweighted while the overrepresented subgroups are downweighted.

(a) Samples generated via importance reweighting with subgroups separated by the orange line. For the 100 samples above, the classifier concludes 52 females and 48 males.
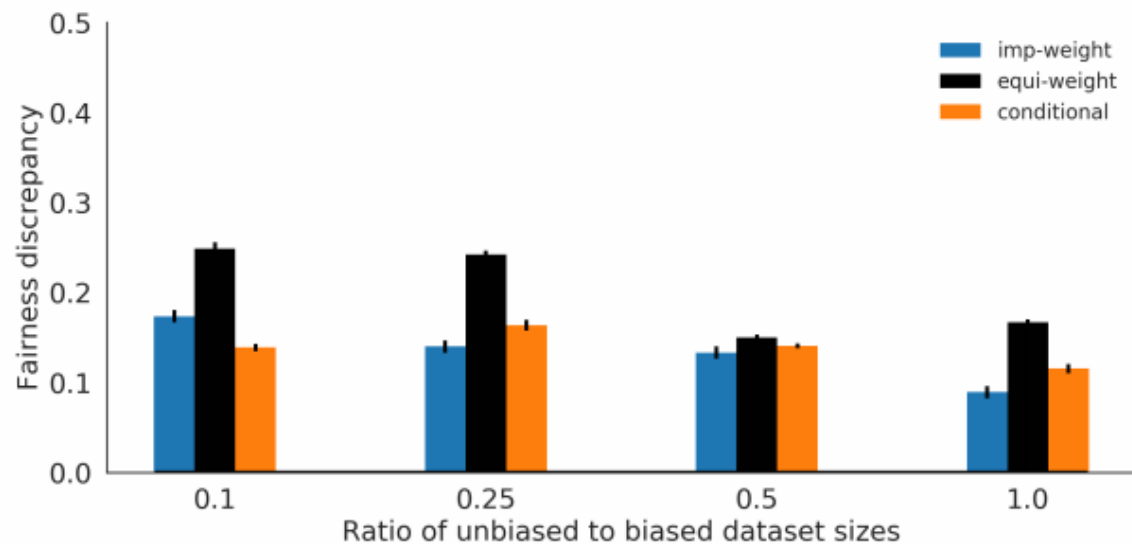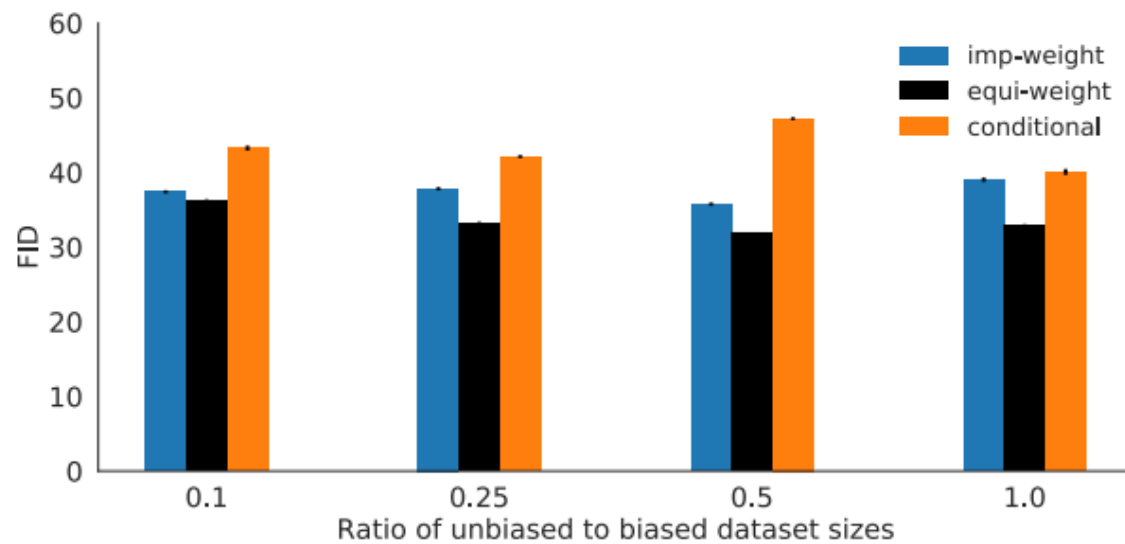
(b) Fairness Discrepancy

(c) FID

*Figure 3.* Single-Attribute Dataset Bias Mitigation for `bias`=0.9. Lower discrepancy and FID is better. Standard error in (b) and (c) over 10 independent evaluation sets of 10,000 samples each drawn from the models. We find that on average, `imp-weight` outperforms the `equi-weight` baseline by 49.3% and the `conditional` baseline by 25.0% across all reference dataset sizes for bias mitigation.

(a) Samples generated via importance reweighting. For the 100 samples above, the classifier concludes 37 females and 20 males without black hair, 22 females and 21 males with black hair.

(b) Fairness Discrepancy

(c) FID

*Figure 4.* Mult-Attribute Dataset Bias Mitigation. Standard error in (b) and (c) over 10 independent evaluation sets of 10,000 samples each drawn from the models. Lower discrepancy and FID is better. We find that on average, `imp-weight` outperforms the `equi-weight` baseline by 32.5% and the `conditional` baseline by 4.4% across all reference dataset sizes for bias mitigation.