
Appendix For Debiasing Pretrained Generative Models by Uniformly Sampling Semantic Attributes

Anonymous Author(s)

Affiliation

Address

email

A Appendix

A.1 Corrections

We unfortunately had an error in Definition 3, and in Figure 5.

A.1.1 Correction for Definition 3

We had an error in the subscripts of the summations in Definition 3. The statement should have been:

Definition (\mathbb{P}_E^λ). Define $\mathbb{P}_E^\lambda = \sum_{i=1}^{|\mathcal{Y}|} \lambda_i E_{:,i}$ as a distribution over \mathcal{Y} determined by prediction-conditional error matrix E and $\lambda = \{\lambda_1, \lambda_2, \dots, \lambda_{|\mathcal{Y}|}\}$, $\lambda_i \in \mathbb{R}_{\geq 0}$, $\sum_{i=1}^{|\mathcal{Y}|} \lambda_i = 1$.

A.1.2 Correction for Figure 5

We incorrectly transposed the Antimode and Mode Polarity Sampling results in this figure. Corrected figure is shown below.

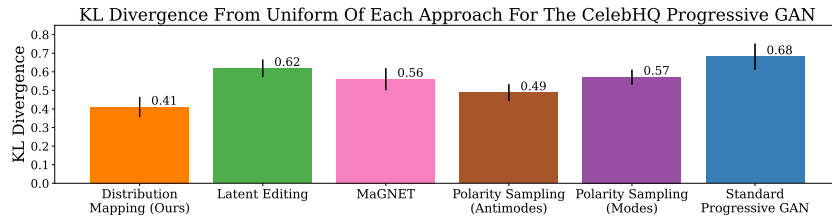


Figure 1: KL Divergence between the distribution over the semantic space for the output of each method (*lower is better*).

A.2 Calculating E for The Generated Distribution

The error rates reported for a classifier C_ϕ are typically reported on the distribution on the distribution of fit’s training data, $\mathbb{P}_{training}$. However, the distribution \mathbb{P}_{G_θ} of the generative model

G_θ may differ from the training distribution. Additionally, rather than reporting $P(\mathbf{y}|\hat{\mathbf{y}})$, often times the error rates are given in a confusion matrix $C_{\hat{\mathbf{y}}|\mathbf{y}}$ where $C_{\hat{\mathbf{y}}|\mathbf{y}}[i, j] = P(\hat{\mathbf{y}}=i|\mathbf{y}=j)$. Thankfully, we can construct the error rate matrix E for the generative distribution \mathbb{P}_{G_θ} under the simplifying assumption that the difference between \mathbb{P}_{G_θ} and $\mathbb{P}_{training}$ can be explained as a label shift [1, 3].

By Bayes' Theorem, we know that

$$P(\mathbf{y}|\hat{\mathbf{y}}) = P(\hat{\mathbf{y}}|\mathbf{y}) \frac{P(\mathbf{y})}{P(\hat{\mathbf{y}})}.$$

Under the label shift assumption, $P(\hat{\mathbf{y}}|\mathbf{y})$ stays the same between $\mathbb{P}_{training}$ and \mathbb{P}_{G_θ} . Additionally, $P(\mathbf{y})$ can be calculated for \mathbb{P}_{G_θ} under label shift [1, 3]. Lastly, $P(\hat{\mathbf{y}})$ can be approximated for \mathbb{P}_{G_θ} by finding the proportion predicted for each class on a large sample from the generative model. Thus, E can be calculated as:

$$E = C_{\hat{\mathbf{y}}|\mathbf{y}} \frac{P_{G_\theta}(\mathbf{y})}{P_{G_\theta}(\hat{\mathbf{y}})}.$$

18 A.3 Distribution of Races Generated By Progressive GAN

19 We show the two best performing methods' distributions on Progressive GAN, along with the
20 distribution of the unmodified ProgressiveGAN, over the Race attribute.

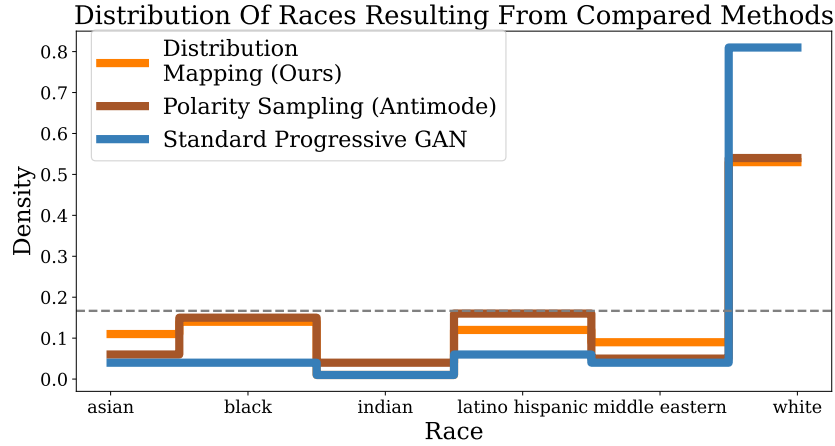


Figure 2: Distribution of our approach, Polarity Antimode Sampling (next best), and the standard generator.

21 A.4 Implementation Details

22 Ground Truth Shape Classifier

Layer (type)		Output Shape
Conv2d-1		[-1, 32, 16, 16]
ReLU-2		[-1, 32, 16, 16]
Conv2d-3		[-1, 64, 8, 8]
ReLU-4		[-1, 64, 8, 8]
Conv2d-5		[-1, 128, 4, 4]
ReLU-6		[-1, 128, 4, 4]
Conv2d-7		[-1, 256, 2, 2]
ReLU-8		[-1, 256, 2, 2]
Conv2d-9		[-1, 2, 1, 1]

37 Encoder for Shapes VAE

Layer (type)	Output Shape
Conv2d-1	[-1, 32, 16, 16]
ReLU-2	[-1, 32, 16, 16]
Conv2d-3	[-1, 64, 8, 8]
ReLU-4	[-1, 64, 8, 8]
Conv2d-5	[-1, 128, 4, 4]
ReLU-6	[-1, 128, 4, 4]
Conv2d-7	[-1, 256, 2, 2]
ReLU-8	[-1, 256, 2, 2]
Conv2d-9	[-1, code_dim, 1, 1]

51 Decoder for Shapes VAE

Layer (type)	Output Shape
ConvTranspose2d-1	[-1, 256, 2, 2]
ReLU-2	[-1, 256, 2, 2]
ConvTranspose2d-3	[-1, 128, 8, 8]
ReLU-4	[-1, 128, 8, 8]
ConvTranspose2d-5	[-1, 64, 16, 16]
ReLU-6	[-1, 64, 16, 16]
ConvTranspose2d-7	[-1, 32, 32, 32]
ReLU-8	[-1, 32, 32, 32]
ConvTranspose2d-9	[-1, 3, 64, 64]
Sigmoid-10	[-1, 3, 64, 64]

66 Biased Age Classifier (Note: Target value was normalized age, made binary after)

Layer (type)	Output Shape
Conv2d-1	[-1, 2, 32, 32]
BatchNorm2d-2	[-1, 2, 32, 32]
LeakyReLU-3	[-1, 2, 32, 32]
Dropout-4	[-1, 2, 32, 32]
Conv2d-5	[-1, 4, 16, 16]
BatchNorm2d-6	[-1, 4, 16, 16]
LeakyReLU-7	[-1, 4, 16, 16]
Dropout-8	[-1, 4, 16, 16]
Conv2d-9	[-1, 8, 8, 8]
BatchNorm2d-10	[-1, 8, 8, 8]
LeakyReLU-11	[-1, 8, 8, 8]
Dropout-12	[-1, 8, 8, 8]
Flatten-13	[-1, 512]
Linear-14	[-1, 64]
LeakyReLU-15	[-1, 64]
Linear-16	[-1, 1]
Sigmoid-17	[-1, 1]

88 Ground Truth Age Classifier (Note: Target value was normalized age; made binary after)

Layer (type)	Output Shape
Conv2d-1	[-1, 8, 32, 32]

93	BatchNorm2d-2	[-1, 8, 32, 32]
94	LeakyReLU-3	[-1, 8, 32, 32]
95	Dropout-4	[-1, 8, 32, 32]
96	Conv2d-5	[-1, 16, 16, 16]
97	BatchNorm2d-6	[-1, 16, 16, 16]
98	LeakyReLU-7	[-1, 16, 16, 16]
99	Dropout-8	[-1, 16, 16, 16]
100	Conv2d-9	[-1, 32, 8, 8]
101	BatchNorm2d-10	[-1, 32, 8, 8]
102	LeakyReLU-11	[-1, 32, 8, 8]
103	Dropout-12	[-1, 32, 8, 8]
104	Flatten-13	[-1, 2048]
105	Linear-14	[-1, 64]
106	LeakyReLU-15	[-1, 64]
107	Linear-16	[-1, 1]
108	Sigmoid-17	[-1, 1]
109	=====	

110 The distribution mapper used default architecture of SDV's CTGAN¹ version 0.6.0, except for in the
111 ProgressiveGAN experiment where embedding_dim =512, generator_dim =(512,512) were
112 passed as arguments.

113 For the networks we trained, we utilized the Adam optimizer [2] with learning rate between 0.002
114 and 0.0001.

115 The linear classifier utilized Scikit-Learn's LinearSVC (for latent editing) and RidgeClassifier for the
116 biased Shapes classifier.

117 A.5 Proof of Lemma 1

118 *Proof.* First, note that if $1^{|\mathcal{Y}|} \in \text{Cone}(E)$, then likewise $\frac{1}{|\mathcal{Y}|}1^{|\mathcal{Y}|} \in \text{Cone}(E)$.

119 Let $\mathbf{z}' \sim \mathbb{P}_{z|C_\phi=i}$; i.e., z is a draw from the distribution of noise such that the classifiers prediction of
120 the generated sample corresponding to z' is group i .

121 Let $(C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=i}$ be the pushforward distribution of the perfect classifier C' 's output when
122 conditioned on the generator's output of draws from $\mathbb{P}_{z|C_\phi=i}$. Then, $(C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=i} =$
123 $[Pr(\mathbf{y} = 1|C_\theta = i), Pr(\mathbf{y} = 2|C_\theta = i), \dots, Pr(\mathbf{y} = N|C_\theta = i)] = E_{:,i}$. Thus, $\text{Cone}(\{(C' \circ$
124 $G_\theta)_* \mathbb{P}_{z|C_\phi=i}, \dots, (C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=|\mathcal{Y}|}\}) = \text{Cone}(E)$. Therefor, following from above, $\frac{1}{|\mathcal{Y}|}1^{|\mathcal{Y}|} \in$
125 $\text{Cone}(\{(C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=i}, \dots, (C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=|\mathcal{Y}|}\})$. This means that $\exists \lambda_1, \lambda_2, \dots, \lambda_{|\mathcal{Y}|}$ s.t.
126 $\lambda_1(C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=i} + \dots + \lambda_{|\mathcal{Y}|}(C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=|\mathcal{Y}|} = [\frac{1}{|\mathcal{Y}|}, \dots, \frac{1}{|\mathcal{Y}|}] = \text{Unif} \mathcal{Y}$. This is equiva-
127 lent to saying that $C'(G_\theta(\mathbf{z})) \sim \text{Unif}(\mathcal{Y})$ for $\mathbf{z} \sim \sum_{i=1}^{|\mathcal{Y}|} \lambda_i \mathbb{P}_{z|C_\phi=i} = \mathbb{Q}^\lambda$. Thus, by definition \mathbb{Q}^λ
128 is a Fair Noise Distribution.

129 □

130 A.6 Proof of Lemma 2

131 *Proof.* Note that the sign of the coefficient of the cross product $E_{:,1} \times E_{:,2}$ is $P(\mathbf{y} = 1|\hat{\mathbf{y}} =$
132 $1)P(\mathbf{y} = 2|\hat{\mathbf{y}} = 2) - P(\mathbf{y} = 1|\hat{\mathbf{y}} = 2)P(\mathbf{y} = 2|\hat{\mathbf{y}} = 1)$. Also note that $E_{:,1} \times [0.5, 0.5]$ is
133 $0.5P(\mathbf{y} = 1|\hat{\mathbf{y}} = 1) - 0.5P(\mathbf{y} = 2|\hat{\mathbf{y}} = 1)$.

134 Additionally, $P(\mathbf{y} = 1|\hat{\mathbf{y}} = 1)P(\mathbf{y} = 2|\hat{\mathbf{y}} = 2) > P(\mathbf{y} = 1|\hat{\mathbf{y}} = 1)0.5 > 0$, and $0 < P(\mathbf{y} = 1|\hat{\mathbf{y}} =$
135 $2)P(\mathbf{y} = 2|\hat{\mathbf{y}} = 1) < 0.5P(\mathbf{y} = 2|\hat{\mathbf{y}} = 1)$. Thus, the coefficient of $E_{:,1} \times E_{:,2}$ is greater than
136 $E_{:,1} \times [0.5, 0.5]$, while there signs are equal. This implies that $[0.5, 0.5]$ is in between $E_{:,1}$ and $E_{:,2}$.
137 Thus, $[0.5, 0.5] \in \text{cone}(E)$. The rest of the proof follows directly from Lemma 1. □

¹https://sdv.dev/SDV/user_guides/single_table/ctgan.html#how-to-modify-the-ctgan-hyperparameters

138 A.7 Proof of Proposition 1

139 *Proof.* Note that \mathbb{P}_E^λ has density $[\sum_i \lambda_i Pr(\mathbf{y} = 1|\hat{\mathbf{y}} = i), \dots, \sum_i \lambda_i Pr(\mathbf{y} = N|\hat{\mathbf{y}} = i)]$. For ease
 140 of notation let us refer to $\sum_i \lambda_i Pr(\mathbf{y} = m|\hat{\mathbf{y}} = i)$ as r_m^λ .

141 Then,

$$\begin{aligned} KL\{\mathbb{P}_E^\lambda || Unif(\mathcal{Y})\} &= \sum_m r_m^\lambda \log\left(\frac{r_m^\lambda}{u}\right) \\ &= \sum_m \left(r_m^\lambda \log(r_m^\lambda) - r_m^\lambda \log\left(\frac{1}{|\mathcal{Y}|}\right)\right) \\ &= \sum_m r_m^\lambda \log(r_m^\lambda) - \sum_m r_m^\lambda \log\left(\frac{1}{|\mathcal{Y}|}\right) \end{aligned}$$

142 Note that $\log\left(\frac{1}{N}\right)$ is constant for each term in the second summation. Thus,

$$\begin{aligned} &= \sum_m r_m^\lambda \log(r_m^\lambda) - \log\left(\frac{1}{N}\right) \sum_m r_m^\lambda \\ &= \sum_m r_m^\lambda \log(r_m^\lambda) - \log\left(\frac{1}{N}\right), \end{aligned}$$

143 As $\log\left(\frac{1}{N}\right)$ does not depend on r_m^λ ,

$$\begin{aligned} \operatorname{argmin}_\lambda KL\{\mathbb{P}_E^\lambda || Unif(\mathcal{Y})\} &= \operatorname{argmin}_\lambda \sum_m r_m^\lambda \log(r_m^\lambda) \\ &= \operatorname{argmin}_\lambda -H(\mathbb{P}_E^\lambda) \\ &= \operatorname{argmax}_\lambda H(\mathbb{P}_E^\lambda) \end{aligned}$$

144 □

145 A.8 Proof of Proposition 2

Proof.

$$\begin{aligned} (C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=i} &= [Pr(\mathbf{y} = 1|C_\theta = i), Pr(\mathbf{y} = 2|C_\theta = i), \dots, Pr(\mathbf{y} = N|C_\theta = i)] \\ \implies \sum_i \lambda_i (C' \circ G_\theta)_* \mathbb{P}_{z|C_\phi=i} &= [\sum_i \lambda_i Pr(\mathbf{y} = 1|C_\theta = i), \sum_i \lambda_i Pr(\mathbf{y} = 2|C_\theta = i), \\ &\quad \dots, \sum_i \lambda_i Pr(\mathbf{y} = N|C_\theta = i)] \\ \implies \mathbb{P}_E &= \mathbb{Q}^\lambda \end{aligned}$$

146 □

147 A.9 Proof of Theorem 1

148 The first statement follows directly from Proposition 1 and Proposition 2.

149 If $C_\phi = C'$, then $\{\frac{E_{:,1}}{|E_{:,1}|}, \dots, \frac{E_{:,|\mathcal{Y}|}}{|E_{:,|\mathcal{Y}|}|}\}$ forms a standard basis of $\mathbb{R}^{|\mathcal{Y}|}$, and therefor $1^{|\mathcal{Y}|}$ is in
 150 $Cone(E)$. Thus, $\mathbb{Q}^{\lambda*}$ is a Fair Noise Distribution by Lemma 1.

151 References

- 152 [1] S. Garg, Y. Wu, S. Balakrishnan, and Z. Lipton. A Unified View of Label Shift Estimation.
 153 In *Advances in Neural Information Processing Systems*, volume 33, pages 3290–3300. Curran
 154 Associates, Inc., 2020.

- 155 [2] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization, Jan. 2017.
156 arXiv:1412.6980 [cs].
- 157 [3] T. Sipka, M. Sulc, and J. Matas. The Hitchhiker’s Guide to Prior-Shift Adaptation. In 2022
158 *IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2031–2039,
159 Waikoloa, HI, USA, Jan. 2022. IEEE.