# A Fair Generative Model Using LeCam Divergence

March 21, 2024

Seoul National University

# Table of Contents

## Contribution

- Generate fair synthetic data via LeCam divergence and unlabelled reference dataset.

# Table of Contents

- $\mathcal{X}$ : Data in $\mathbb{R}^D$
- $\mathcal{S}$ : Set of sensitive attributes

# Table of Contents

## Fairness

- Fair synthetic data satisfies

$$P_G(S = 0) = P_G(S = 1)$$

  where $P_{syn}$ is distribution of synthetic data.

- However train data does not satisfy above condition.

$$P_{bias}(S = 0) \neq P_{bias}(S = 1)$$

# Table of Contents

## Setup

- Suppose that the information of $S$ is not available.

- To train fair synthetic data, adopt reference dataset $\mathcal{D}^{\text{ref}}$ which may satisfy

$$P_{ref}(S = 0) \approx P_{ref}(S = 1)$$

- Let $\mathcal{D}^{\text{bias}}$ be train data.

- When the number of train data and reference data are $m_{bias}$ and $m_{ref}$ respectively, we assume that

$$m_{bias} \gg m_{ref}$$

## Original Method

- To train fair synthetic data, we optimize generator $G$ by minimizing

$$\min_{G}(1 - \lambda) \cdot D_f \left( \mathbb{P}_{\mathsf{bias}} \| \mathbb{P}_G \right) + \lambda \cdot D_{\mathsf{fair}} \left( \mathbb{P}_{\mathsf{ref}} \| \mathbb{P}_G \right)$$

where $D_f$ is $f$-divergence and $D_{\mathsf{fair}}$ is fair discrepancy.

# Table of Contents

- Propose Lecam divergence based fairness discrepancy

$$\min_G (1 - \lambda) \cdot D_f \left( \mathbb{P}_{\mathsf{bias}} \| \mathbb{P}_G \right) + \lambda \cdot \mu D_\Delta \left( \mathbb{P}_{\mathsf{ref}} \| \mathbb{P}_G \right)$$

where $\mu$ denotes a non-negative weight, and $D_\Delta \left( \mathbb{P}_{\mathsf{ref}} \| \mathbb{P}_G \right)$ indicates the LC-divergence between $\mathbb{P}_{\mathsf{ref}}$ and $\mathbb{P}_G$ :

$$D_\Delta \left( \mathbb{P}_{\mathsf{ref}} \| \mathbb{P}_G \right) := \sum_{x \in \mathcal{X}} \frac{\left( \mathbb{P}_{\mathsf{ref}} \left( x \right) - \mathbb{P}_G(x) \right)^2}{\mathbb{P}_{\mathsf{ref}} \left( x \right) + \mathbb{P}_G(x)}$$

- $f$-GAN

$$\min_G \max_D \mathbb{E}_{\mathbb{P}_{\text{bias}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [f^*(D(X))]$$

  where $D$ is discriminator and $f^*$ is conjugate of $f$.

- fair $f$-GAN

$$\max_{D} \mathbb{E}_{\mathbb{P}_{\text{bias}}} [D(X)] - \mathbb{E}_{\mathbb{P}_G} [f^*(D(X))]$$

$$\max_{D_{\text{ref}}} \mathbb{E}_{\mathbb{P}_{\text{ref}}} [D_{\text{ref}}(X)] - \mathbb{E}_{\mathbb{P}_G} [D_{\text{ref}}(X)] - \frac{1}{2(\mu + \alpha)} R_\Delta$$

$$\min_{G} -(1 - \lambda)\mathbb{E}_{\mathbb{P}_G} [f^*(D(X))] - \lambda \mathbb{E}_{\mathbb{P}_G} [D_{\text{ref}}(X)]$$

where $\alpha$ denotes an exponential moving average of $D_{\text{ref}}$ v.r.t. reference samples and $R_\Delta$ indicates a regularization term for $D_{\text{ref}}$ defined as:

$$R_\Delta := \mathbb{E}_{\mathbb{P}_{\text{ref}}} \left[ \|D_{\text{ref}}(X) + \alpha\|^2 \right] + \mathbb{E}_{\mathbb{P}_G} \left[ \|D_{\text{ref}}(X) - \alpha\|^2 \right]$$

- Baseline 1 : Unfair method with train and reference data
- Baseline 2 : Unfair method with reference data
- Fairness measure : $\sqrt{\sum_{s=1}^{S}(P_{ref}(S=s) - P_G(S=s))^2}$

| Reference set size | | 25% | 10% | 5% | 2.5% | 1% |
|---|---|---|---|---|---|---|
| Baseline I | Intra FID | $12.00 \pm 0.069$ | $\mathbf{12.73 \pm 0.053}$ | $\mathbf{13.54 \pm 0.074}$ | $\mathbf{13.79 \pm 0.072}$ | $\mathbf{15.89 \pm 0.094}$ |
| | Fairness | $0.495 \pm 0.001$ | $0.554 \pm 0.002$ | $0.559 \pm 0.001$ | $0.566 \pm 0.002$ | $0.576 \pm 0.002$ |
| Baseline II | Intra FID | $23.81 \pm 0.118$ | $32.31 \pm 0.109$ | $40.07 \pm 0.062$ | $67.70 \pm 0.112$ | $92.34 \pm 0.131$ |
| | Fairness | $0.093 \pm 0.002$ | $0.115 \pm 0.002$ | $0.120 \pm 0.003$ | $0.150 \pm 0.003$ | $0.455 \pm 0.002$ |
| Choi et al. (2020) | Intra FID | $20.68 \pm 0.076$ | $25.74 \pm 0.079$ | $30.15 \pm 0.037$ | $30.40 \pm 0.041$ | $31.49 \pm 0.074$ |
| | Fairness | $0.065 \pm 0.002$ | $0.104 \pm 0.002$ | $0.126 \pm 0.001$ | $0.237 \pm 0.003$ | $0.344 \pm 0.002$ |
| Proposed | Intra FID | $\mathbf{11.48 \pm 0.814}$ | $14.50 \pm 0.996$ | $14.64 \pm 0.626$ | $17.16 \pm 1.607$ | $23.11 \pm 0.797$ |
| | Fairness | $\mathbf{0.037 \pm 0.007}$ | $\mathbf{0.039 \pm 0.013}$ | $\mathbf{0.118 \pm 0.007}$ | $\mathbf{0.129 \pm 0.010}$ | $\mathbf{0.146 \pm 0.022}$ |