

Learning Fair Representation using a Parametric Integral Probability Metric

ICML '22

X : Non Sensitive random input vector.

S : Binary Sensitive random input vector.

Y : Output variable.

$h : X \times \{0,1\} \rightarrow \mathbb{Z}$ encoding function

$z := h(X, S)$ representation.

$f_D : \mathbb{Z} \rightarrow X \times \{0,1\}$ decoding function.

Focus: Unsupervised Fair Representation Learning.

Loss to min: $L_{\text{reconstruction}}(f_D \circ h) + \lambda d(P_0^h, P_1^h)$ where $h(X, S) | S=s \sim P_s^h$

Prior: $d(P_0^h, P_1^h) = \sup_{v \in \mathcal{V}} E[v \circ h(X, S)]$

Computation issues.

Can't say that $d(P_0^h, P_1^h) \leq \epsilon \Rightarrow$
bound on predictor's fairness.

\emptyset -fairness of prediction model:

$$DP_{\emptyset}(f) = |E[\emptyset \circ f(X, S) | S=0] - E[\emptyset \circ f(X, S) | S=1]|$$

Proposed method: IPM-based deviance.

$$d_{\mathcal{V}}(P_0, P_1) = \sup_{v \in \mathcal{V}} |E_{P_0}[v(z)] - E_{P_1}[v(z)]|$$

If $d_{\mathcal{V}}(P_0, P_1) \leq \epsilon$,

$$DP_{\emptyset}(f \circ h) = |E_{P_0}[\emptyset \circ f \circ h(X, S)] - E_{P_1}[\emptyset \circ f \circ h(X, S)]|$$

$$\leq \epsilon \text{ if } \emptyset \circ f \in \mathcal{V}$$

Proposed parametric family for \mathcal{V}

$$\mathcal{V}_{\text{sig}} = \{ \sigma(\theta^T x + u) : \theta \in \mathbb{R}^m, u \in \mathbb{R} \}$$

where σ is the sigmoid function.

Lemma 1. $d_{\mathcal{V}_{\text{sig}}}(P_0, P_1) < \epsilon \Rightarrow \sup_{a \in \mathbb{R}^m} \sup_{t \in \mathbb{R}} |P(a^T U_0 \leq t) - P(a^T U_1 \leq t)| < \epsilon$

$$\text{Let } \eta = \frac{1}{1 + e^{-1/\delta}}$$

$$\eta |P(a^T U_0 \leq t) - P(a^T U_1 \leq t)| = |1 - \eta P(a^T U_0 \leq t) - (1 - \eta P(a^T U_1 \leq t))|$$

$$= \left| 1 - \eta P(a^T U_0 \in [t - \delta, t]) - \eta P(a^T U_0 \leq t - \delta) - \right. \\ \left. | (1 - \eta P(a^T U_1 \in [t - \delta, t]) - \eta P(a^T U_1 \leq t - \delta)) \right|$$

$$\leq \eta \sum_{\delta \in \{0,1\}} P(a^T U_{\delta} \leq t - \delta) + \\ |1 - \eta P(a^T U_0 \leq t - \delta) - (1 - \eta P(a^T U_1 \leq t - \delta))|$$

$$= \eta \sum_{\delta \in \{0,1\}} P(a^T U_{\delta} \leq t - \delta) + \\ \left| \sum_{\delta \in \{0,1\}} \left(1 - \eta P(a^T U_{\delta} \leq t - \delta) - E \left[\sigma \left(\frac{a^T U_{\delta} - t}{\delta^2} \right) \right] \right) \right. \\ \left. - E \left[\sigma \left(\frac{a^T U_0 - t}{\delta^2} \right) \right] + E \left[\sigma \left(\frac{a^T U_0 - t}{\delta^2} \right) \right] \right|$$

$$= \eta \sum_{\delta \in \{0,1\}} P(a^T U_{\delta} \leq t - \delta) + \\ \left| \sum_{\delta \in \{0,1\}} \left[E \left[\sigma \left(\frac{a^T U_{\delta} - t}{\delta^2} \right) \right] - (1 - \eta P(a^T U_{\delta} \leq t - \delta)) \right] \right. \\ \left. - E \left[\sigma \left(\frac{a^T U_0 - t}{\delta^2} \right) \right] \right|$$

$$\leq \eta \sum_{\delta \in \{0,1\}} P(a^T U_{\delta} \leq t - \delta) + \\ \left| \sum_{\delta \in \{0,1\}} \left[E \left[\sigma \left(\frac{a^T U_{\delta} - t}{\delta^2} \right) \right] - (1 - \eta P(a^T U_{\delta} \leq t - \delta)) \right] \right| \\ \left| E \left[\sigma \left(\frac{a^T U_1 - t}{\delta^2} \right) \right] - E \left[\sigma \left(\frac{a^T U_0 - t}{\delta^2} \right) \right] \right|$$

$$\leq \eta 2\epsilon + 4\epsilon + \epsilon$$

$$= \frac{2\epsilon}{1 + e^{-1/\delta}} + 5\epsilon$$