

Fairness Transferability Subject to Bounded Distribution Shift

Yatong Chen, Reilly Raab, Jialu Wang, Yang Liu

Presented by Seonghyeon Kim

Formulation

- $(X, Y, G) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{G}$: Random vector of a triplet of output, input and sensitive variable
- \mathcal{S}, \mathcal{T} : Source distribution and target distribution
- $\mathcal{P}(\cdot)$: Space of distributions over a given space
- $\Pi := \{\pi : \mathcal{X} \times \mathcal{G} \rightarrow \mathcal{P}(\mathcal{Y})\}$: Space of nondeterministic policies
- $\hat{Y}|X, G \sim \pi(X, G)$: Predicted label

Disparity Function

- $\Delta^*(\pi, \mathcal{T})$: Disparity of a policy π on a distribution \mathcal{T} .
- Demographic Parity (DP)

$$\Delta^*_{\text{DP}}(\pi, \mathcal{T}) := \sum_{g, h \in \mathcal{G}} \left| \Pr_{\pi, \mathcal{T}}(\hat{Y}=1 \mid G=g) - \Pr_{\pi, \mathcal{T}}(\hat{Y}=1 \mid G=h) \right|$$

- Equal Opportunity (EOP)

$$\Delta^*_{\text{EOP}}(\pi, \mathcal{T}) := \sum_{g, h \in \mathcal{G}} \left| \Pr_{\pi, \mathcal{T}}(\hat{Y}=1 \mid G=g, Y=1) - \Pr_{\pi, \mathcal{T}}(\hat{Y}=1 \mid G=h, Y=1) \right|$$

Vector-Bounded Distribution Shift

- Divergence : $K(P \parallel Q)$ is a *divergence* if and only if K satisfies the followings.
 - $K(P \parallel Q) \geq 0$ for all P and Q
 - $K(P \parallel Q) = 0 \iff P = Q$
- Group-vectorized shift : For given divergences $K_1, \dots, K_{|\mathcal{G}|}$,

$$\mathbf{D}(\mathcal{T} \parallel \mathcal{S}) := \left[K_g \left(\Pr_{\mathcal{T}}(X, Y \mid G=g) \parallel \Pr_{\mathcal{S}}(X, Y \mid G=g) \right) \right]_{g=1, \dots, |\mathcal{G}|}.$$

- Vector-bounded distribution shift : For given a vector $\mathbf{B} \succeq 0$ (element-wise inequality),

$$\mathbf{D}(\mathcal{T} \parallel \mathcal{S}) \preceq \mathbf{B}.$$

General Bound

- Supremum disparity function :

$$v(\Delta^*, \mathcal{D}, \pi, \mathcal{S}, \mathcal{B}) := \sup_{\mathcal{D}(\mathcal{T} \parallel \mathcal{S}) \leq \mathcal{B}} \Delta^*(\pi, \mathcal{T})$$

Theorem (Lipshitz Upper Bound)

If there exists an L such that $\nabla_b v(\Delta^*, \mathcal{D}, \pi, \mathcal{S}, b) \leq L$, everywhere along some curve as b varies from 0 to B , then for $\mathcal{D}(\mathcal{T} \parallel \mathcal{S}) \leq B$

$$\Delta^*(\pi, \mathcal{T}) \leq \Delta^*(\pi, \mathcal{S}) + L \cdot B.$$

Restricted Distribution Shift

- Covariate shift :

$$\Pr_{\mathcal{T}}(Y | X, G) = \Pr_{\mathcal{S}}(Y | X, G)$$

- Label shift :

$$\Pr_{\mathcal{T}}(X | Y, G) = \Pr_{\mathcal{S}}(X | Y, G)$$

Covariate Shift - Demographic Parity

- Covariate shift : $\Pr_{\mathcal{T}}(Y | X, G) = \Pr_{\mathcal{S}}(Y | X, G)$
- Density ratio : $\omega_g(\mathcal{T}, \mathcal{S}, x) := \frac{p_{\mathcal{T}}(X=x | G=g)}{p_{\mathcal{S}}(X=x | G=g)}$
- Note that $\omega_g(\mathcal{T}, \mathcal{S}, x)$ is extremely large, when $p_{\mathcal{T}}(X=x | G=g)$ is sufficiently large and $p_{\mathcal{S}}(X=x | G=g) \approx 0$
- $\beta_g := \Pr_{\pi, \mathcal{S}}(\hat{Y}=1 | G=g)$

Theorem (Bound for DP under Covariate Shift)

For demographic parity between two groups under covariate shift,

$$\Delta^*_{DP}(\pi, \mathcal{T}) \leq \Delta^*_{DP}(\pi, \mathcal{S}) + \sum_g (\beta_g(1 - \beta_g) \cdot \text{Var}_{\mathcal{S}}[\omega_g(\mathcal{T}, \mathcal{S}, x)])^{1/2}$$

Covariate Shift - Equal Opportunity

- Covariate shift : $\Pr_{\mathcal{T}}(Y | X, G) = \Pr_{\mathcal{S}}(Y | X, G)$
- $\beta_g^+(\pi, \mathcal{T}) := \Pr_{\pi, \mathcal{T}}(\hat{Y}=1 | Y=1, G=g)$
- $\Delta^*_{EOp}(\pi, \mathcal{T}) = \sum_{g, g' \in \mathcal{G}} |\beta_g^+(\pi, \mathcal{T}) - \beta_{g'}^+(\pi, \mathcal{T})|$

Theorem (Bound for EOp under Covariate Shift)

Subject to covariate shift and any given D, B , assume extremal values for β_g^+ , i.e.,

$$\forall g, \quad (D_g(\mathcal{T} \| \mathcal{S}) < B_g) \implies (l_g \leq \beta_g^+(\pi, \mathcal{T}) \leq u_g)$$

it follows that

$$v(\Delta^*_{EOp}, D, \pi, \mathcal{S}, B) \leq \max_{r \in \prod_{g=1}^{|\mathcal{G}|} \{l_g, u_g\}} \sum_{h, h' \in \mathcal{G}} |r_h - r_{h'}|.$$

Label Shift - Demographic Parity

- Label shift : $\Pr_{\mathcal{T}}(X \mid Y, G) = \Pr_{\mathcal{S}}(X \mid Y, G)$
- Note that $\Pr_{\pi, \mathcal{T}}(\hat{Y} \mid Y, G) = \Pr_{\pi, \mathcal{S}}(\hat{Y} \mid Y, G)$, so EOp is invariant.
- $\beta_g^+ := \Pr_{\pi}(\hat{Y}=1 \mid Y=1, G=g); \quad \beta_g^- := \Pr_{\pi}(\hat{Y}=1 \mid Y=0, G=g)$
- $Q_g(\mathcal{T}) := \Pr_{\mathcal{T}}(Y = 1 \mid G = g)$

Theorem (Bound for DP under Label Shift)

For DP under the bounded label-shift assumption $\forall g, |Q_g(\mathcal{S}) - Q_g(\mathcal{T})| \leq B_g$,

$$\Delta_{DP}^*(\pi, \mathcal{T}) \leq \Delta_{DP}^*(\pi, \mathcal{S}) + (|\mathcal{G}| - 1) \sum_g B_g \left| \beta_g^+ - \beta_g^- \right|$$