

Trade-offs and Guarantees of Adversarial Representation Learning for Information Obfuscation

Jianfeng Chi^{§*}, Han Zhao^{†*}, Yuan Tian[§], Geoffrey J. Gordon[†]

[§]University of Virginia, [†]Carnegie Mellon University

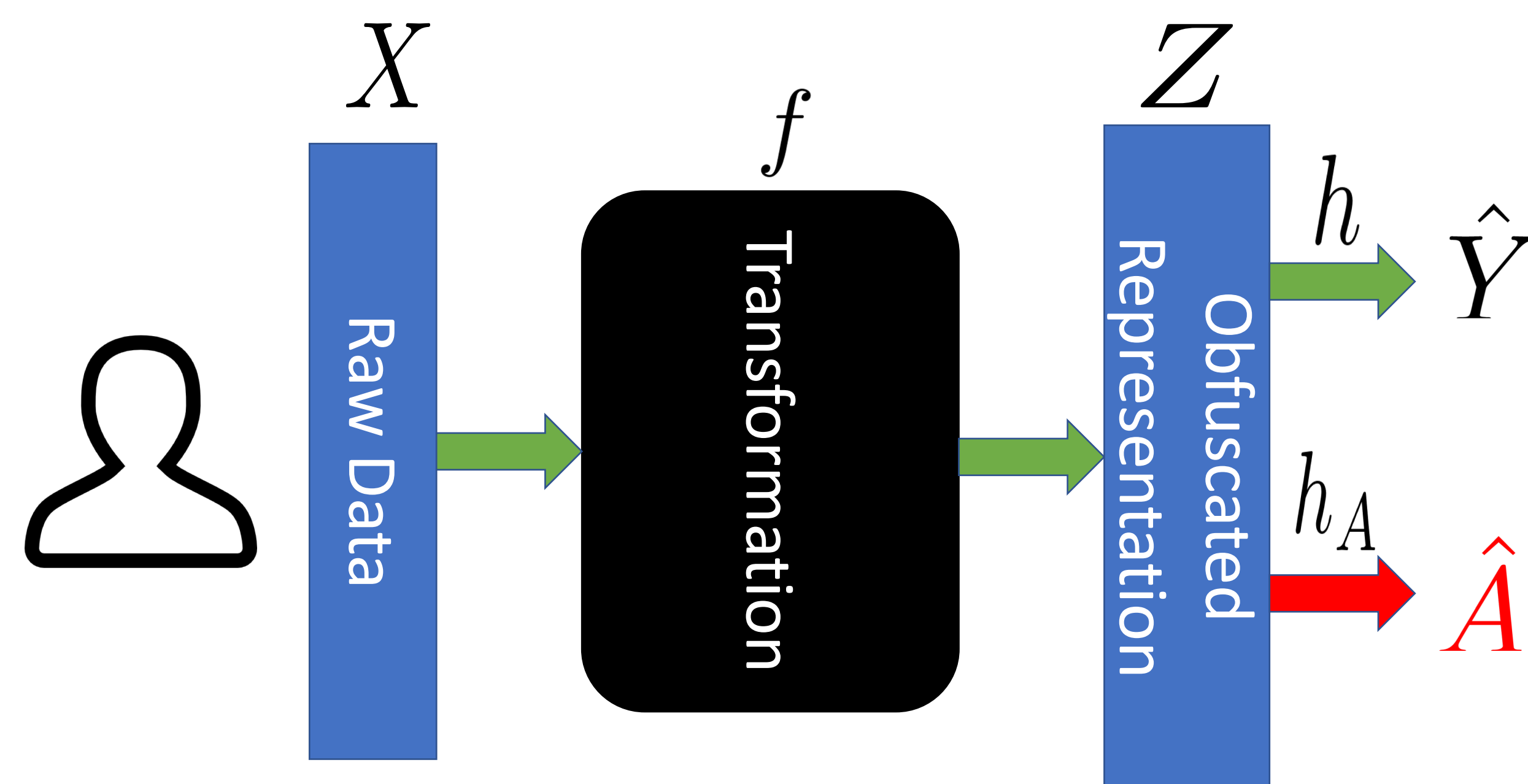
{jfc6ub,yuant}@virginia.edu, {han.zhao,ggordon}@cs.cmu.edu

* equal contribution



Overview

Learning Representations that Obfuscate Sensitive Attributes:



Question:

Can we prevent the information leakage of the sensitive attribute while still maximizing the task accuracy? Furthermore, what is the fundamental trade-off between attribute obfuscation and accuracy maximization in the minimax problem?

Preliminaries

Utility:

$$\text{Acc}(h) := 1 - \mathbb{E}_{\mathcal{D}}[|Y - h(X)|]$$

Attribute Inference Advantage:

$$\text{Adv}(\mathcal{H}_A) := \max_{h_A \in \mathcal{H}_A} \left| \Pr(h_A(X) = 1 \mid A = 1) - \Pr(h_A(X) = 1 \mid A = 0) \right|$$

- $\text{Adv}_A(h) = 0$ iff $I(h(X); A) = 0$ and $\text{Adv}_A(h) = 1$ iff $h(X) = A$ almost surely or $h(X) = 1 - A$
- $\text{Adv}(\mathcal{H}_A) + \min_{h_A \in \mathcal{H}_A} \Pr(h_A(X) = 0 \mid A = 1) + \Pr(h_A(X) = 1 \mid A = 0) = 1$ if \mathcal{H}_A is symmetric: the larger the attribute inference advantage of \mathcal{H}_A , the smaller the minimum sum of Type-I and Type-II error under attacks from \mathcal{H}_A .

Theoretical Analysis

Formal Guarantees against Attribute Inference

$$\min_{h \in \mathcal{H}, f} \max_{h_A \in \mathcal{H}_A} \widehat{\text{Err}}(h \circ f) - \lambda \left(\Pr_{\mathbf{s}}(h_A(f(X)) = 0 \mid A = 1) + \Pr_{\mathbf{s}}(h_A(f(X)) = 1 \mid A = 0) \right) \quad (1)$$

In practice, we have:

$$\min_{h \in \mathcal{H}, f} \max_{h_A \in \mathcal{H}_A} \text{CE}_Y(h \circ f) - \lambda \cdot \text{CE}_A(h_A \circ f) \quad (2)$$

Theorem:

Let f^* be the optimal feature map such that $f^* = \arg \min H(Y \mid Z = f(X)) - \lambda H(A \mid Z = f(X))$ and define $H^* := H(A \mid Z = f^*(X))$. Then for any adversary \hat{A} such that $I(\hat{A}; A \mid Z) = 0$, we have

$$\Pr_{\mathcal{D}^{f^*}}(\hat{A} \neq A) \geq H^* / 2 \lg(6/H^*).$$

Implication: If the obfuscated representation Z contains little information on A , then the inference error made by any adversary has to be large.

Inherent trade-off between Accuracy Maximization and Attribute Obfuscation

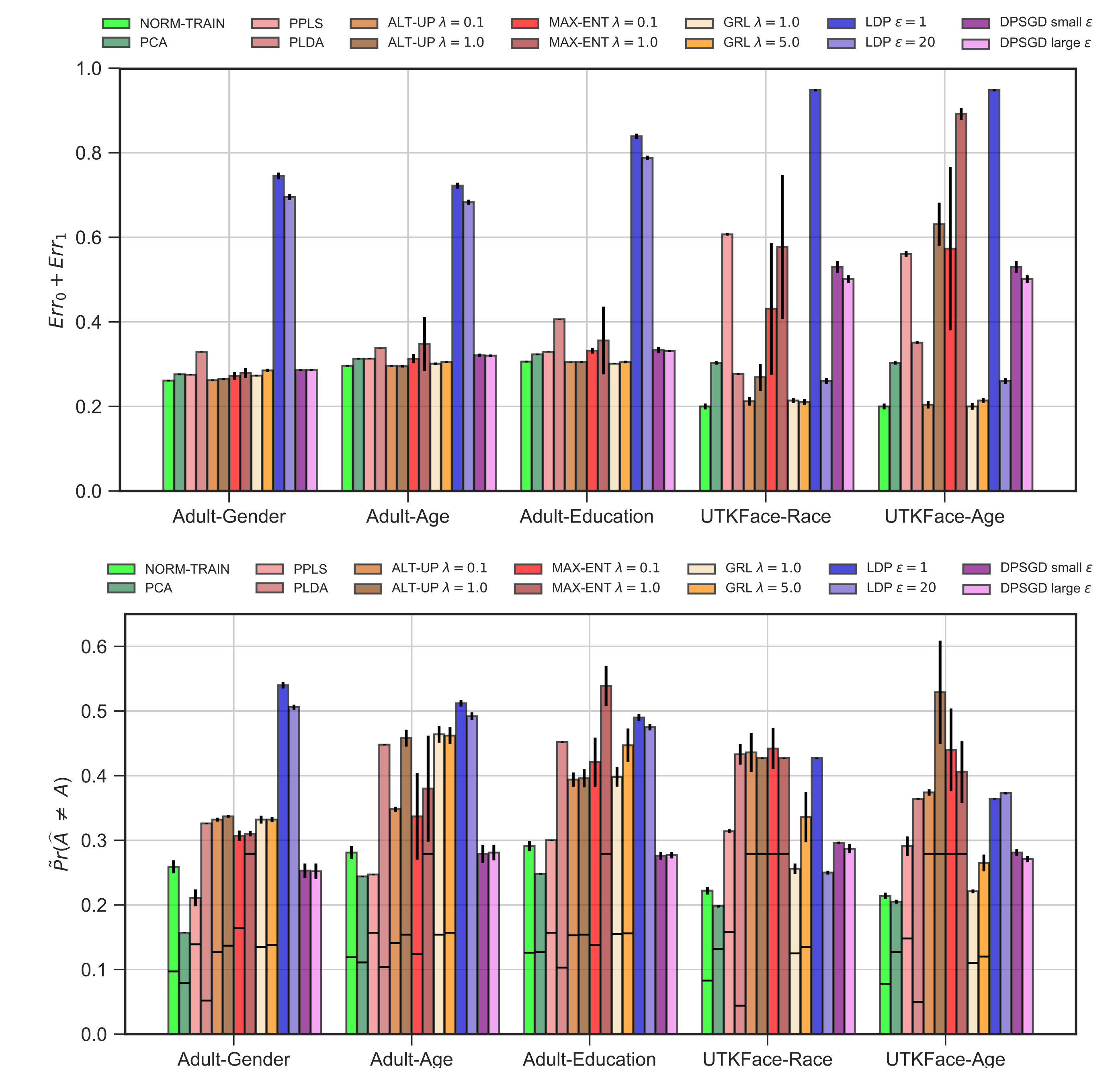
Theorem: Let $\mathcal{H} \subseteq 2^{\mathcal{Z}}$ contains all the measurable functions from \mathcal{Z} to $\{0, 1\}$ and $\mathcal{D}_0^Y, \mathcal{D}_1^Y$ be two distributions over \mathcal{Y} conditioned on $A = 0$ and $A = 1$ respectively. Assume the Markov chain $X \xrightarrow{f} Z \xrightarrow{h} \hat{Y}$ holds, If $\text{Adv}(\mathcal{H}_A \circ f) \leq D_{\text{JS}}(\mathcal{D}_0^Y, \mathcal{D}_1^Y)$, then $\forall h \in \mathcal{H}$, we have

$$\text{Err}_0(h \circ f) + \text{Err}_1(h \circ f) \geq \frac{1}{2} (d_{\text{JS}}(\mathcal{D}_0^Y, \mathcal{D}_1^Y) - \sqrt{\text{Adv}(\mathcal{H}_A \circ f)})^2.$$

Implication: If the label and the sensitive attribute are highly correlated, we cannot obfuscate the sensitive attribute while still maximizing the task accuracy simultaneously.

Empirical Results

(1) Income prediction on the UCI Adult dataset with sensitive attributes: gender, age, and education; (2) Gender estimation on UTKFace dataset with sensitive attributes: age and race.



- The formal guarantees hold for all representation learning based approaches;
- Inherent trade-offs between accuracy maximization and attribute obfuscation exist for all methods;
- Compared to DP-related methods, adversarial representation learning based approaches leads to better trade-offs;

Conclusion: The adversarial representation learning approaches achieve the best trade-offs in terms of attribute obfuscation and accuracy maximization.