

FAIR CLASSIFICATION

Explosion of fairness research over last five years

Fair classification is the most common setup, involving:

- X , some data
- Y , a label to predict
- \hat{Y} , the model prediction
- A , a sensitive attribute (race, gender, age, socio-economic status)

We want to learn a classifier that is:

- accurate
- fair with respect to A

REPRESENTATIONS BEYOND CLUSTERS

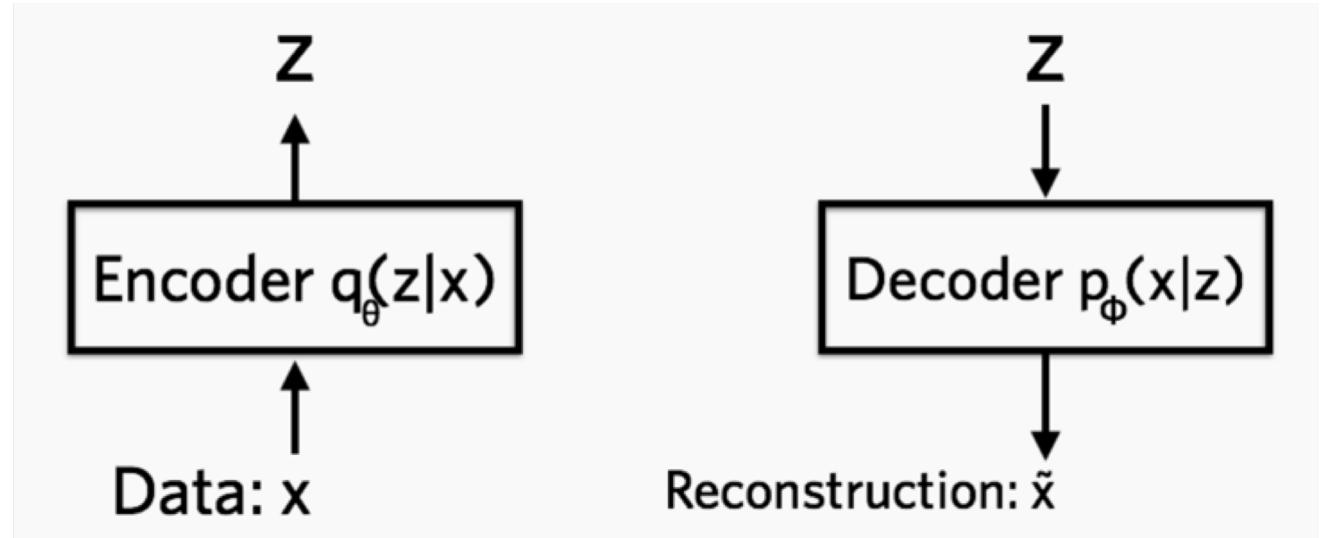
Aim: Replace discrete representation with continuous, multi-dimensional Z

Allow more flexible, nuanced representations

Bring ML arsenal to bear: powerful methods for mapping, embedding in vector spaces: Variational Auto Encoders (VAE)

How to maintain statistical parity in learned representations?

VAE



Re-formulation of autoencoders:

- Each input encoded into a distribution in latent space
- Output prediction obtained by sampling from distribution, mapping through decoder

Allows maximum-likelihood based density modelling:

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - D_{KL} \left(q_{\phi}(\mathbf{z} | \mathbf{x}) \parallel p(\mathbf{z}) \right)$$

MMD

- Suppose we have access to samples from two probability distributions $X \sim P_A$ and $Y \sim P_B$, how can we tell if $P_A = P_B$?
- Maximum Mean Discrepancy (MMD) is a measure of distance between two distributions given only samples from each. [Gretton 2010]

$$\begin{aligned} & \left\| \frac{1}{N} \sum_{n=1}^N \phi(X_n) - \frac{1}{M} \sum_{m=1}^M \phi(Y_m) \right\|^2 \\ &= \frac{1}{N^2} \sum_{n=1}^N \sum_{n'=1}^N \phi(X_n)^\top \phi(X_{n'}) + \frac{1}{M^2} \sum_{m=1}^M \sum_{m'=1}^M \phi(Y_m)^\top \phi(Y_{m'}) - \frac{2}{NM} \sum_{n=1}^N \sum_{m=1}^M \phi(X_n)^\top \phi(Y_m) \\ &= \frac{1}{N^2} \sum_{n=1}^N \sum_{n'=1}^N k(X_n, X_{n'}) + \frac{1}{M^2} \sum_{m=1}^M \sum_{m'=1}^M k(Y_m, Y_{m'}) - \frac{2}{MN} \sum_{n=1}^N \sum_{m=1}^M k(X_n, Y_m) \end{aligned}$$

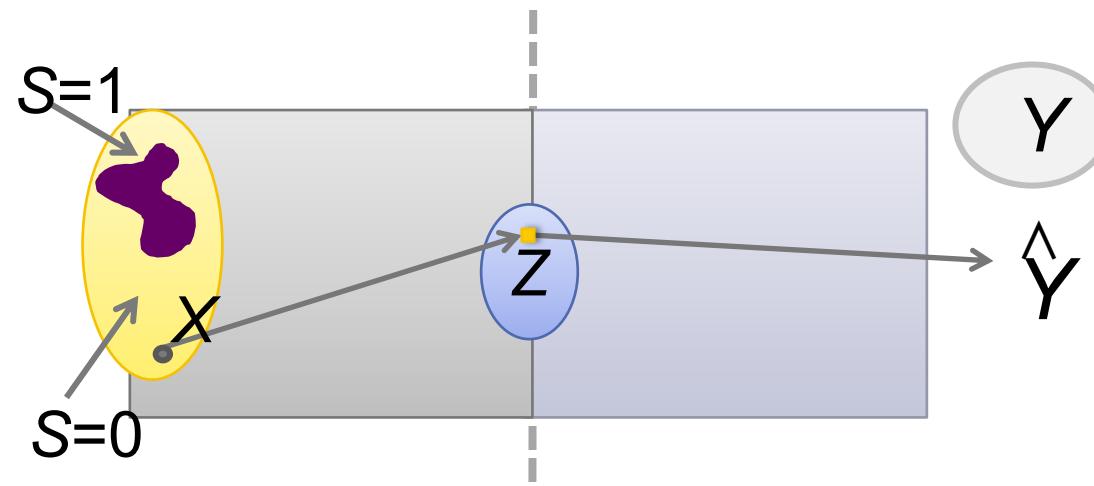
- Our idea: learn to make two distributions indistinguishable
→ small MMD!

VARIATIONAL FAIR AUTOENCODER

VAE with regularizer on latent representations

Match higher-order moments, continuous Z:

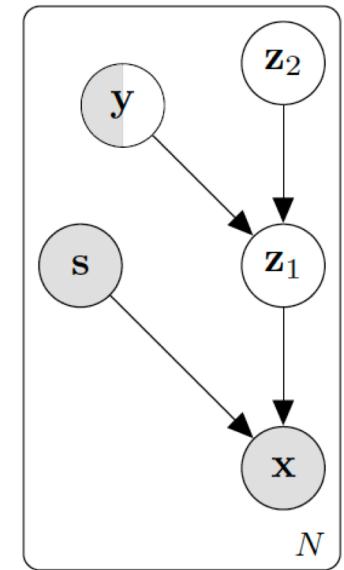
$$\ell_{\text{MMD}}(\mathbf{z}_{1s=0}, \mathbf{z}_{1s=1}) = \| \mathbb{E}_{\tilde{p}(\mathbf{x}|s=0)}[\mathbb{E}_{q(\mathbf{z}_1|\mathbf{x}, s=0)}[\psi(\mathbf{z}_1)]] - E_{\tilde{p}(\mathbf{x}|s=1)}[\mathbb{E}_{q(\mathbf{z}_1|\mathbf{x}, s=1)}[\psi(\mathbf{z}_1)]] \|^2$$



VARIATIONAL FAIR AUTOENCODER

Extend VAE to include some labels y (semi-supervised VAE [Kingma & Welling, 2014]) and “nuisance variable” s

Objective -- maximize:



$$\sum_{n=1}^{N_s} \mathbb{E}_{q_\phi(\mathbf{z}_{1n}|\mathbf{x}_n, \mathbf{s}_n)} [-KL(q_\phi(\mathbf{z}_{2n}|\mathbf{z}_{1n}, \mathbf{y}_n) || p(\mathbf{z}_2)) + \log p_\theta(\mathbf{x}_n|\mathbf{z}_{1n}, \mathbf{s}_n)] + \\ + \mathbb{E}_{q_\phi(\mathbf{z}_{2n}|\mathbf{z}_{1n}, \mathbf{y}_n)} [-KL(q_\phi(\mathbf{z}_{1n}|\mathbf{x}_n, \mathbf{s}_n) || p_\theta(\mathbf{z}_{1n}|\mathbf{z}_{2n}, \mathbf{y}_n))]$$

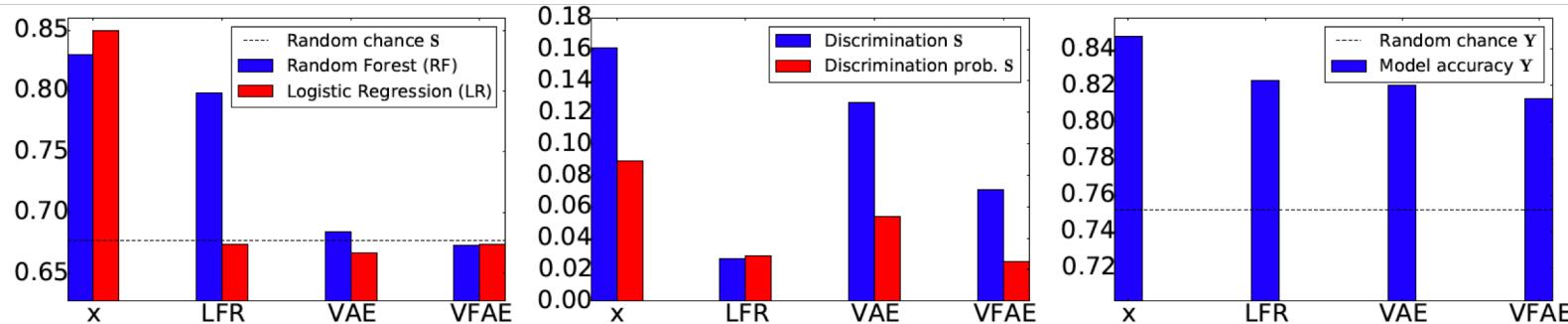
Add for labeled set:

$$\sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_{1n}|\mathbf{x}_n, \mathbf{s}_n)} [-\log q_\phi(\mathbf{y}_n|\mathbf{z}_{1n})]$$

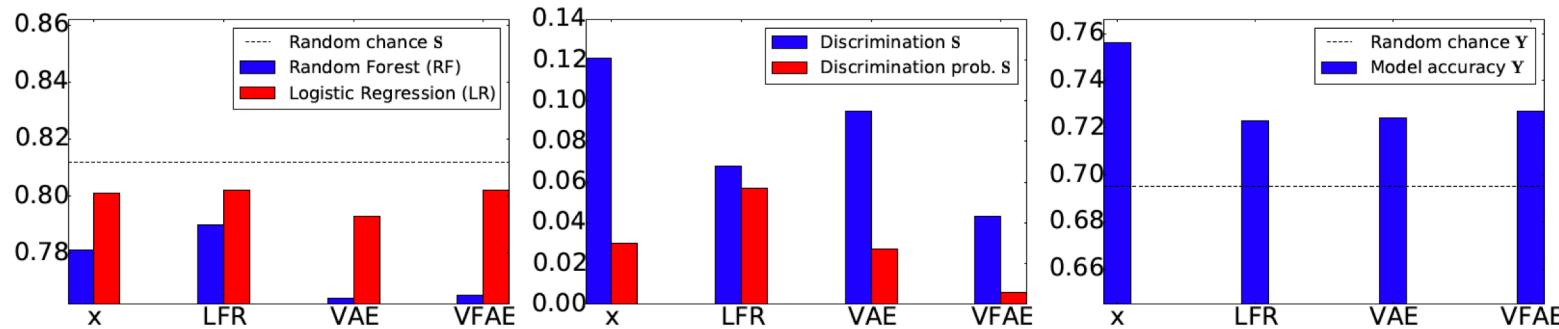
unlabeled set:

$$\sum_{m=1}^M \mathbb{E}_{q_\phi(\mathbf{z}_{1m}|\mathbf{x}_m, \mathbf{s}_m)} [-KL(q(\mathbf{y}_m|\mathbf{z}_{1m}) || p(\mathbf{y}_m))]$$

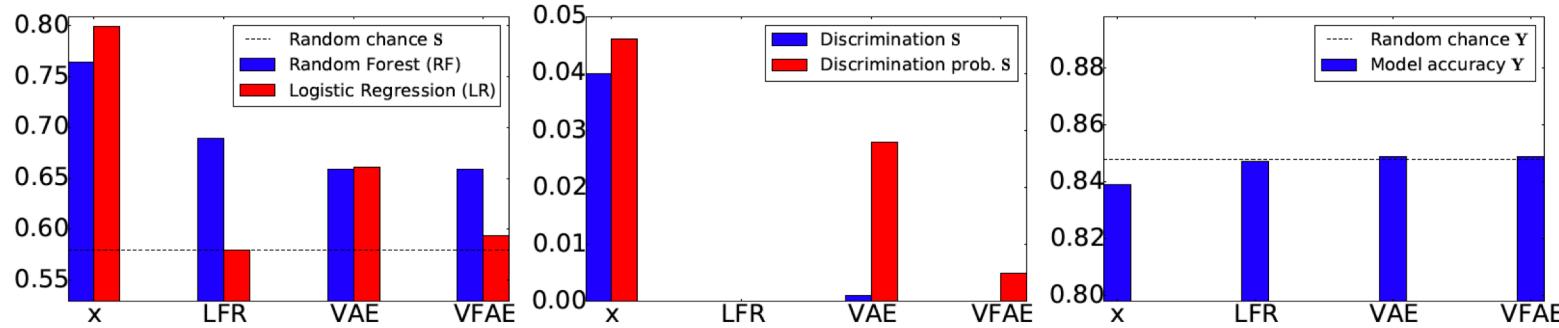
RESULTS



(a) Adult dataset



(b) German dataset



(c) Health dataset

RESULTS

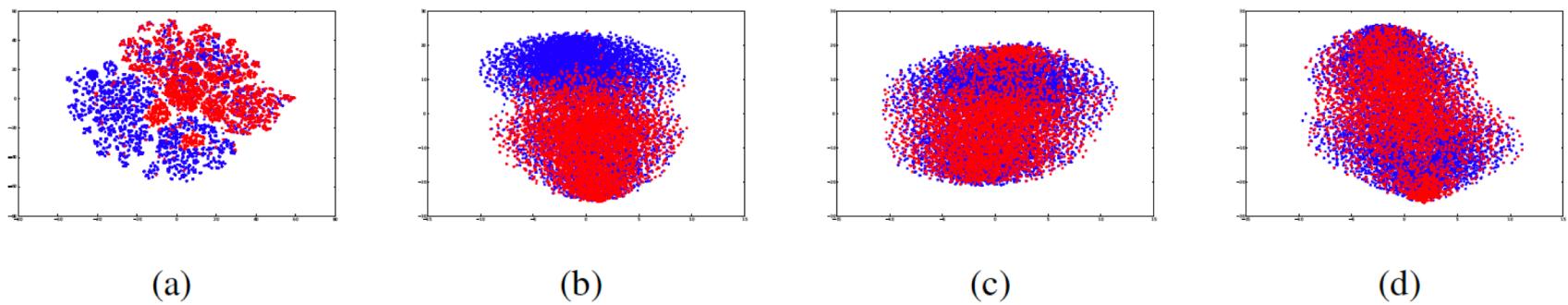


Figure 4: t-SNE (van der Maaten, 2013) visualizations from the Adult dataset on: (a): original \mathbf{x} , (b): latent \mathbf{z}_1 without \mathbf{s} and MMD, (c): latent \mathbf{z}_1 with \mathbf{s} and without MMD, (d): latent \mathbf{z}_1 with \mathbf{s} and MMD. Blue colour corresponds to males whereas red colour corresponds to females.

ADAPTING THE FRAMEWORK

The same idea has many other useful applications, e.g.,

- Eliminating demographic discrimination in deciding who should get transplant surgery
- Removing confounds, such as which scanner produced a medical image

Key: Learning to make two (or more) distributions indistinguishable

DOMAIN ADAPTATION

Natural fit: **domain adaptation**

Make feature representations for source and target domain data indistinguishable

Sentiment classification

- Product reviews (text, tf-idf on words & bigrams)
- Labeled data from source domain, unlabeled data from target domain

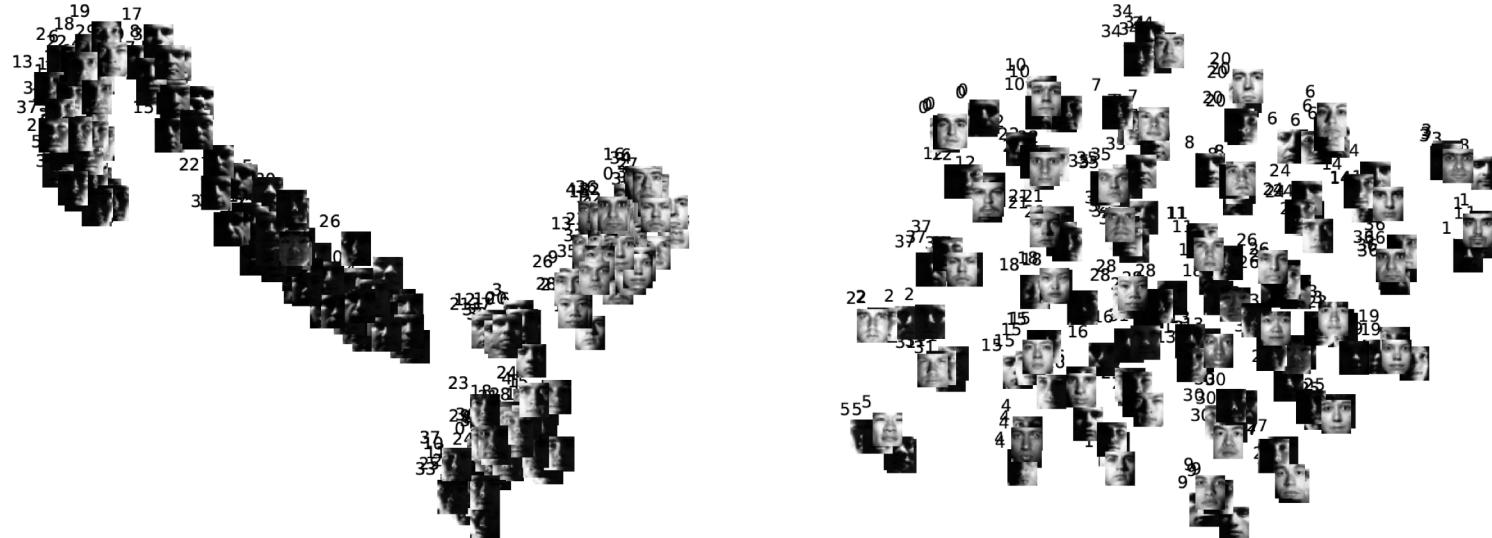
Source - Target	S		Y	
	RF	LR	VFAE	DANN
books - dvd	0.535	0.564	0.799	0.784
books - electronics	0.541	0.562	0.792	0.733
books - kitchen	0.537	0.583	0.816	0.779
dvd - books	0.537	0.563	0.755	0.723
dvd - electronics	0.538	0.566	0.786	0.754
dvd - kitchen	0.543	0.589	0.822	0.783
electronics - books	0.562	0.590	0.727	0.713
electronics - dvd	0.556	0.586	0.765	0.738
electronics - kitchen	0.536	0.570	0.850	0.854
kitchen - books	0.560	0.593	0.720	0.709
kitchen - dvd	0.561	0.599	0.733	0.740
kitchen - electronics	0.533	0.565	0.838	0.843

LEARNING INVARIANT FEATURES

If we have labeled data from all domains, factoring out unwanted domain bias still leads to better generalization.

Make the learned representations invariant to unwanted transformation / variation / bias.

Example: Face identification under different lighting conditions



ADVERSARIAL FAIR LEARNING

Rather than using MMD to ensure learned representation is fair, can use adversarial approach

Adversary takes latent representation (here R) as input and attempts to predict S , then model minimizes:

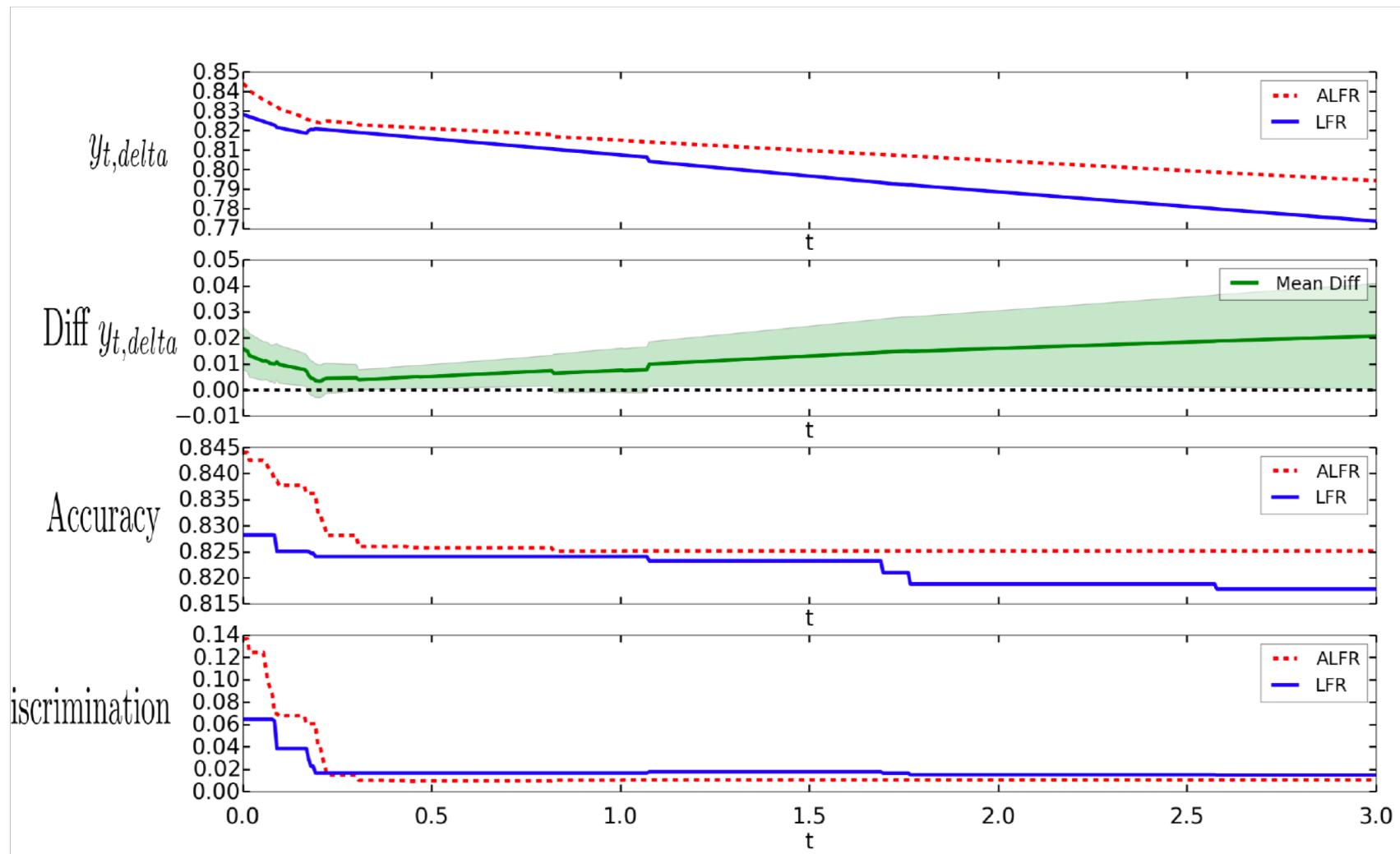
$$D_{\theta, \phi}(R, S) = \mathbb{E}_{X, S} S \cdot \log (\text{Adv}(R)) + (1 - S) \cdot \log (1 - \text{Adv}(R))$$

Combine with reconstruction and classification losses to ensure representation retains info about X, Y

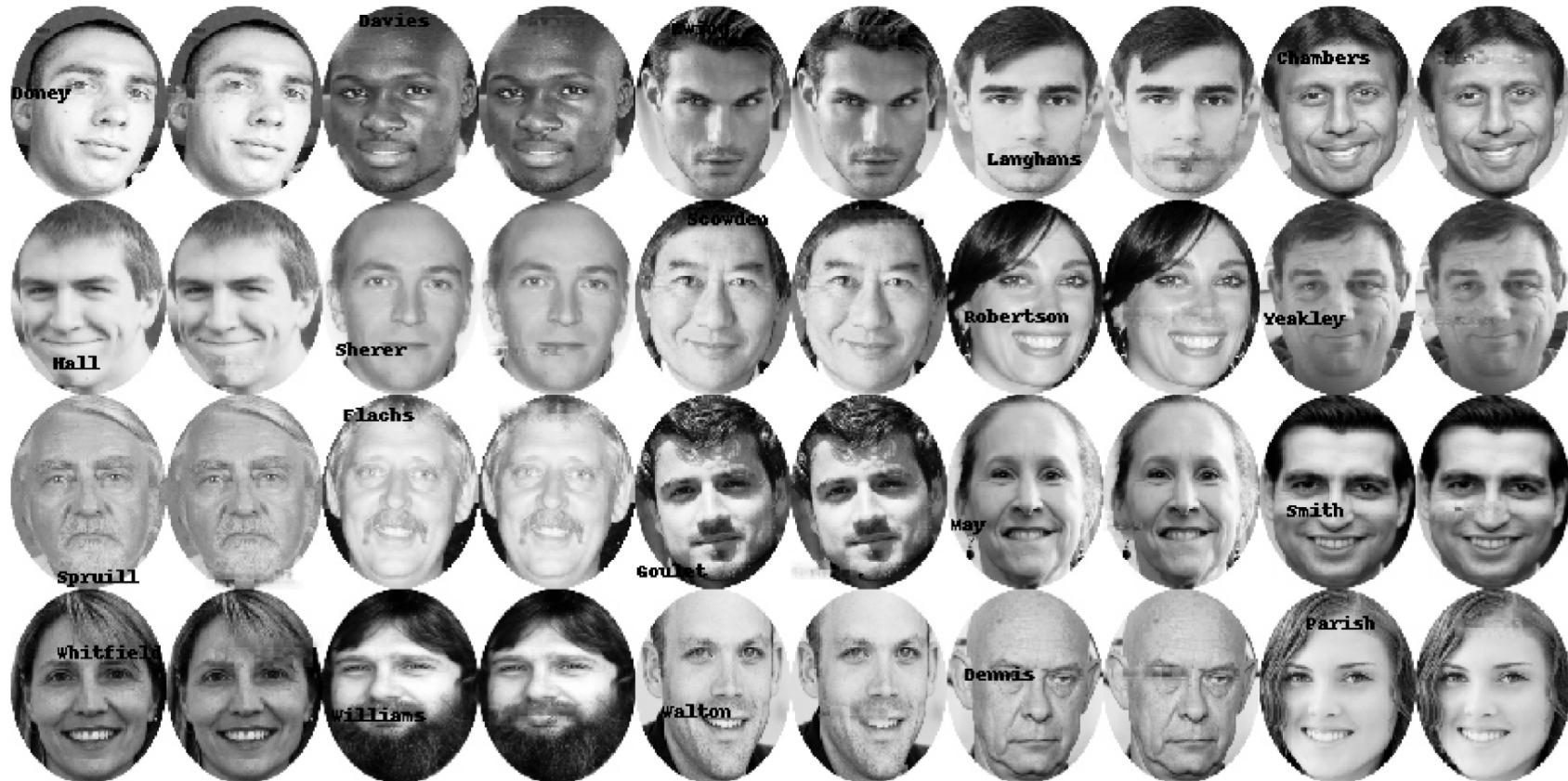
$$C_{\theta}(X, R) = \mathbb{E}_X \|X - \text{Dec}(R)\|_2^2$$

$$E_{\theta}(R, S) = - \mathbb{E}_{X, Y} Y \cdot \log (\text{Pred}(R)) + (1 - Y) \cdot \log (1 - \text{Pred}(R))$$

RESULTS



RESULTS



EQUALIZED ODDS / OPPORTUNITY

Both VFAE and AFLR define fairness as statistical parity

Problems with demographic/statistical parity:

- Coarse measure, not about individuals
- May entail large loss in accuracy

Alternative definition: **equal opportunity** [Hardt, Price, Srebro, 2016]

- Encourage perfect prediction
- But ensure that the prediction errors are balanced between the groups

$$\Pr\{\widehat{Y} = 1 \mid A = 0, Y = y\} = \Pr\{\widehat{Y} = 1 \mid A = 1, Y = y\}, \quad y \in \{0, 1\}$$

FAIR CLASSIFICATION: DEFINITIONS

Definitions based on predicted outcomes:

- Demographic / statistical parity
- Conditional statistical parity (loan conditioned on credit history, amount, employment)

Definitions based on predicted and actual outcomes:

- Balanced PPV (FDR) – predictive equality
- Balanced FNR (TPR) – equal opportunity
- Balanced FNR and FPR – equalized odds

	Actual – Positive	Actual – Negative
Predicted – Positive	True Positive (TP) $PPV = \frac{TP}{TP+FP}$ $TPR = \frac{TP}{TP+FN}$	False Positive (FP) $FDR = \frac{FP}{TP+FP}$ $FPR = \frac{FP}{FP+TN}$
Predicted – Negative	False Negative (FN) $FOR = \frac{FN}{TN+FN}$ $FNR = \frac{FN}{TP+FN}$	True Negative (TN) $NPV = \frac{TN}{TN+FN}$ $TNR = \frac{TN}{TN+FP}$

FAIR CLASSIFICATION: MORE DEFINITIONS

Definitions based on predicted probability and actual outcomes:

- Balanced for positive class
- Balance for negative class
- Calibration

Definitions based on similarity

- Individual fairness

FAIR CLASSIFICATION: DEFINITIONS

Most common way to define fair classification is to require some invariance with respect to the sensitive attribute

- Demographic parity: $\hat{Y} \perp A$
- Equalized Odds: $\hat{Y} \perp A|Y$
- Equal Opportunity: $\hat{Y} \perp A|Y = y$, for some y
- Equal (Weak) Calibration: $Y \perp A|\hat{Y}$
- Equal (Strong) Calibration: $Y \perp A|\hat{Y}$ and $\hat{Y} = P(Y = 1)$
- Fair Subgroup Accuracy: $\mathbb{1}[Y = \hat{Y}] \perp A$

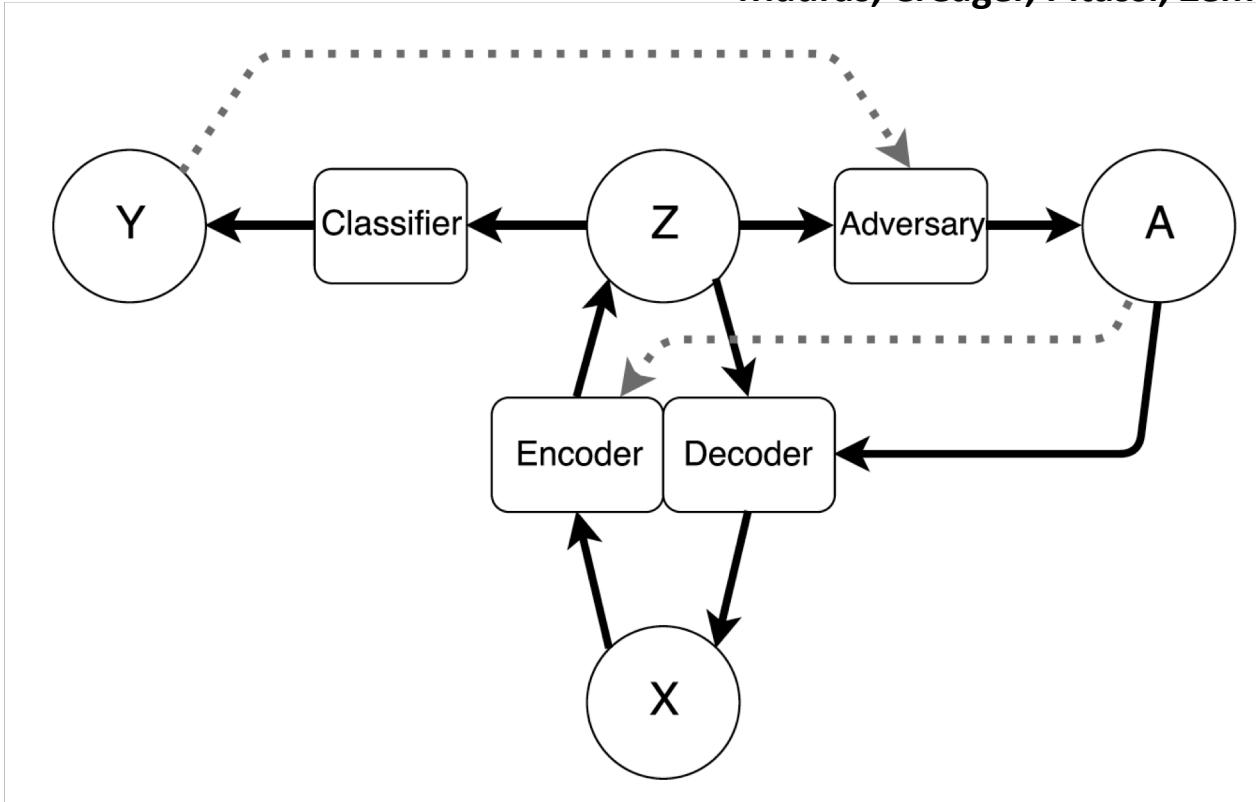
Note: Many of these definitions are incompatible!

WHY FAIR REPRESENTATIONS

- Minimize unfair targeting of disadvantaged groups by vendors (worse lines of credit, lower paying jobs)
- Aim: form a data representation that ensures fair classifications downstream
- Consider two types of unfair vendors:
 1. The **indifferent** vendor: does not care about fairness, only maximizes utility
 2. The **malicious** vendor: doesn't care about utility, discriminates unfairly
- Suggests an adversarial learning scheme

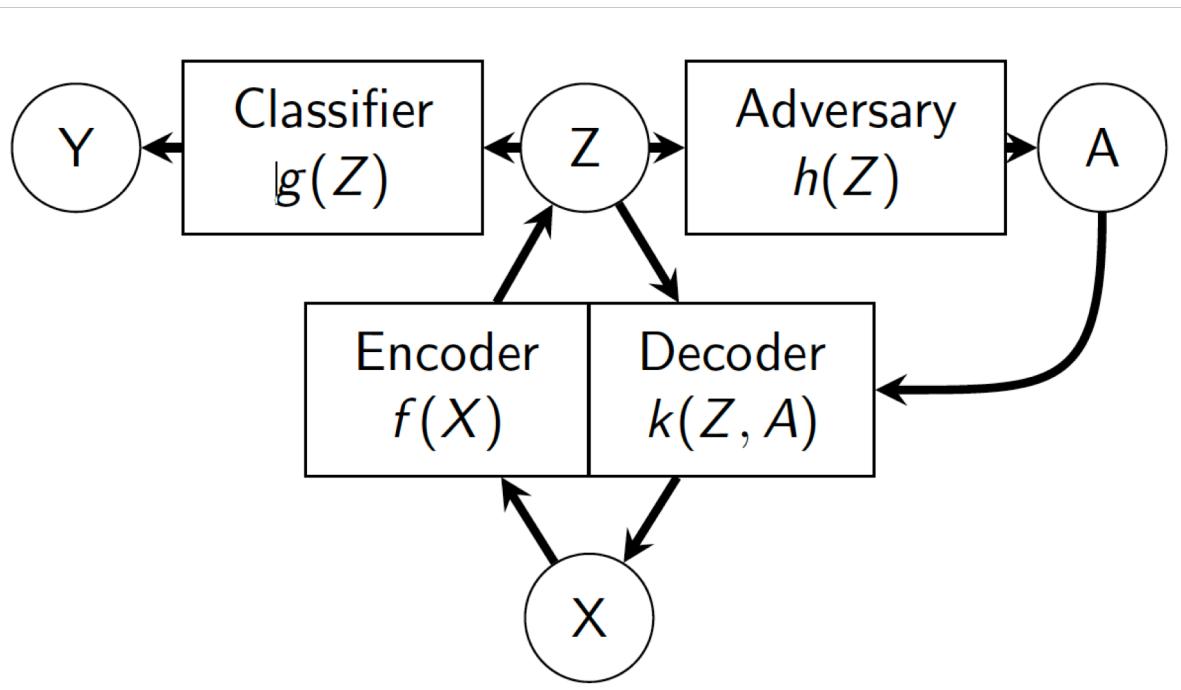
LEARNING ADVERSARILY FAIR TRANSFERABLE REPRESENTATIONS

Madras, Creager, Pitassi, Zemel, 2018



- The classifier is indifferent vendor, forcing the encoder to make the representations useful
- The adversary is the malicious vendor, forcing the encoder to hide the sensitive attributes in the representations

ADVERSARIAL LEARNING IN LAFTR

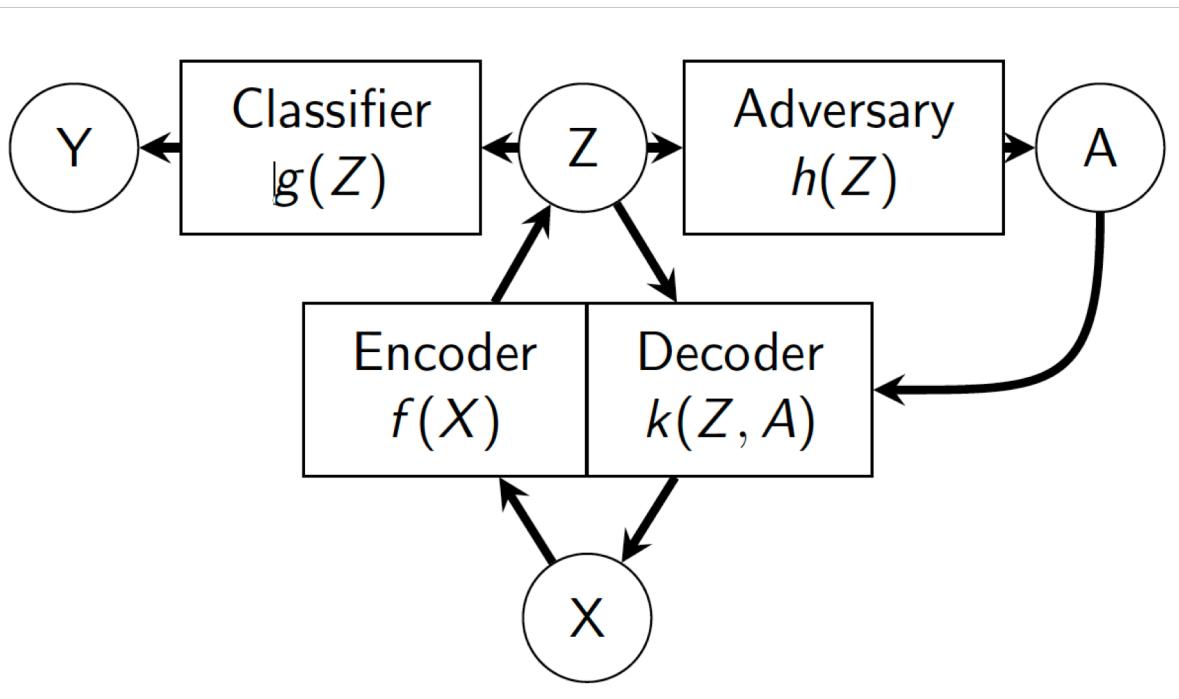


- Our game: encoder-decoder-classifier vs. adversary
- Aim: Learn fair encoder

$$\underset{f,g,k}{\text{minimize}} \underset{h}{\text{maximize}} \mathbb{E}_{X,Y,A} [\mathcal{L}(f, g, h, k)]$$

$$\mathcal{L}(f, g, h, k) = \alpha \mathcal{L}_{\text{Class}} + \beta \mathcal{L}_{\text{Dec}} - \gamma \mathcal{L}_{\text{Adv}}$$

ADVERSARIAL OBJECTIVES



Choice of adversarial objective depends on fairness desideratum

- Demographic parity: $\mathcal{L}_{DP}(h) = \sum_{i \in \{0,1\}} \frac{1}{|\mathcal{D}_i|} \sum_{(x,a) \in \mathcal{D}_i} |h(f(x)) - a|$
- Equalized odds: $\mathcal{L}_{EO}(h) = \sum_{i,j \in \{0,1\}^2} \frac{1}{|\mathcal{D}_i^j|} \sum_{(x,a,y) \in \mathcal{D}_i^j} |h(f(x), y) - a|$
- Equal Opportunity: $\mathcal{L}_{EOpp}(h) = \sum_{i \in \{0,1\}} \frac{1}{|\mathcal{D}_i^1|} \sum_{(x,a) \in \mathcal{D}_i^1} |h(f(x)) - a|$

FROM ADVERSARIAL OBJECTIVES TO FAIRNESS DEFINITIONS

In general: pick the right adversarial loss, encourage the right conditional independencies

- Demographic parity encourages $Z \perp A$ to fool adversary
- Equalized odds encourages $Z \perp A \mid Y$ to fool adversary
- Equal opportunity encourages $Z \perp A \mid Y = 1$ to fool adversary

Note that independencies of $Z = f(x)$ also hold for predictions $\hat{Y} = g(Z)$

We show: In the adversarial limit, these objectives guarantee these fairness metrics!

- The key is to connect predictability of A by the adversary $h(Z)$ to unfairness in the classifier $g(Z)$

EXPERIMENTS

Datasets

1. Adult Income

Size: 45,222 instances, 14 attributes

Task: predict whether or not annual income > 50K

Sensitive feature: Gender

2. Heritage Health

Size: 147,473 instances, 139 attributes

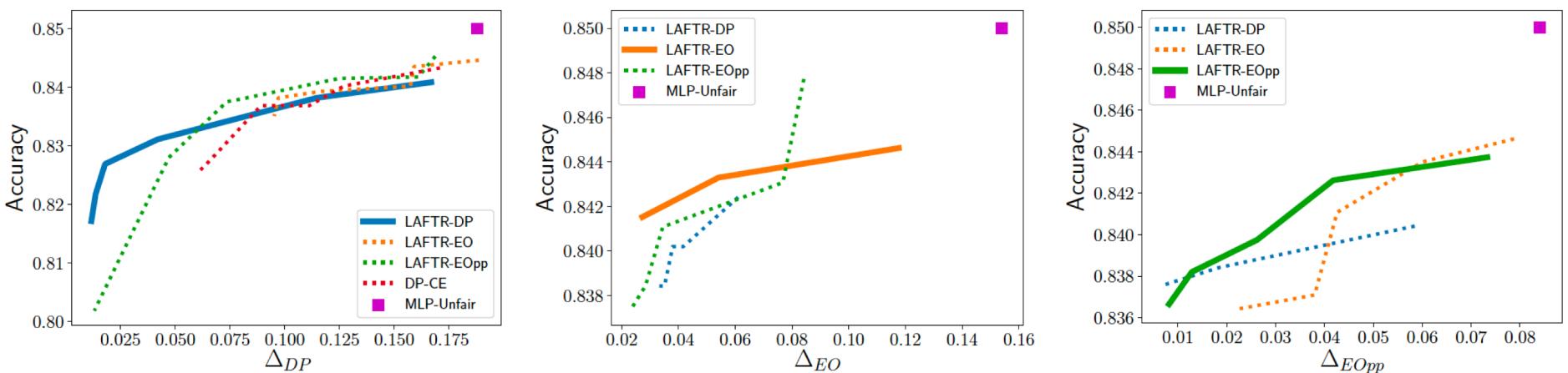
Task: predict patient's Charlson Index (co-morbidity)

Sensitive feature: Age

Models

Encoder, classifier, adversary: each single hidden-layer MLP (8; 20 hidden units)

RESULTS: FAIR CLASSIFICATION



- Train with 2-step process to simulate owner → vendor framework
- Tradeoffs between accuracy and fairness metrics produced by different LAFTR loss functions
- Achieves best solutions, wrt fairness-accuracy tradeoff

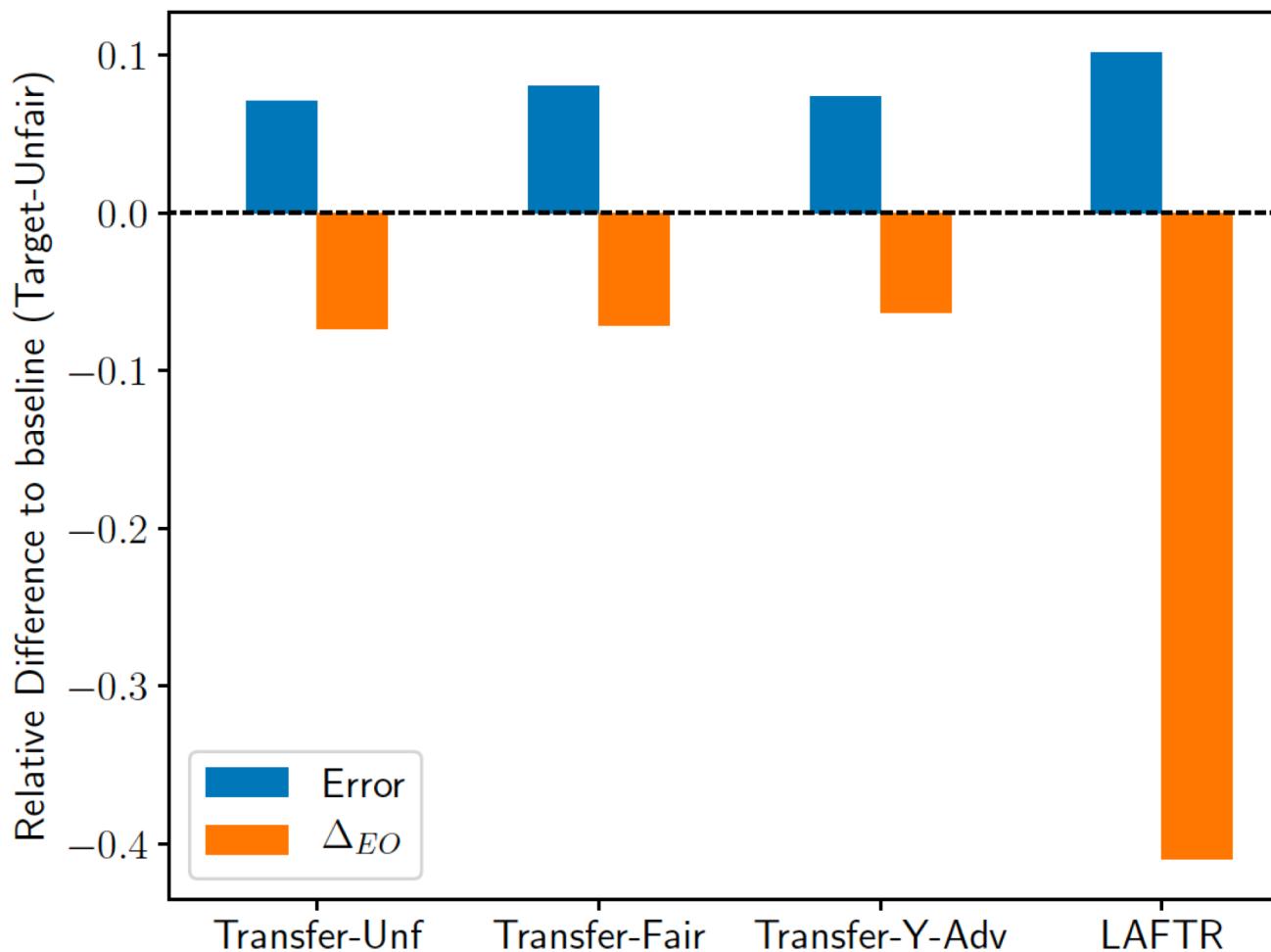
RESULTS: FAIRNESS METRICS

METHOD	Δ_{DP}	Δ_{EO}	Δ_{EOpp}	ACC.
MLP (UNFAIR)	0.381	0.476	0.231	0.785
LAFTR-EO	0.152	0.050	0.036	0.763
	0.143	0.052	0.032	0.752
LAFTR-EOPP	0.087	0.092	0.010	0.742
	0.113	0.063	0.024	0.735
LAFTR-DP	0.041	0.140	0.025	0.731
	0.002	0.196	0.031	0.728

SETUP: FAIR TRANSFER LEARNING

- Downstream vendors will have unknown prediction tasks
- Does fairness transfer?
- We test this as follows:
 - ① Train encoder f on data X , with label Y
 - ② Freeze encoder f
 - ③ On new data X' , train classifier on top of $f(X')$, with new task label Y'
 - ④ Observe fairness and accuracy of this new classifier on new task Y'
- Compare LAFTR encoder f to other encoders
- We use Heritage Health dataset
 - Y is Charlson comorbidity index > 0
 - Y' is whether or not a certain type of insurance claim was made
 - Check for fairness w.r.t. age

RESULTS : FAIR TRANSFER LEARNING



Fair transfer learning on Health dataset. Down is better in both metrics.