# Domain Adaptation under Missingness Shift

Helen Zhou, Sivaraman Balakrishnan, Zachary C. Lipton Carnegie Mellon University

### Motivation

**Q:** What proportion of individuals  $18 + received \ge 1$  dose of the COVID-19 vaccine? As of October 2021, in southwestern Pennsylvania,



### Motivation

Now, suppose the healthcare provider adopts a new intake form...



Absent any shift in the actual health status of patients, the distribution of observed data would still shift, owing to this sudden change in clerical practices.

#### Domain Adaptation under Missingness Shift (DAMS)

Faced with data from different time periods or locations, each characterized by **different missingness patterns**, how might practitioners use this data to produce the **best possible predictor on a target domain**?

# Definitions





entification

Estima

Experimen<sup>®</sup>



### **Missing Data Definition**

In any environment e with **missing data**, we do not observe underlying clean covariates  $X \in \mathbb{R}^d$  but instead observe corrupted covariates:

 $\widetilde{X} = X \odot \xi$ 

where  $\xi \in \{0,1\}^d$ , and  $(X,Y,\xi) \sim P^e$ .

- $1 \xi$  are missing data indicators, which may or may not be observed.
- $\xi$  could be completely at random, dependent on other covariates, etc.

fication



#### **Missingness Shift Definition**

For a source domain *s* and target domain *t*, assume:

$$P(X,Y) = P^s(X,Y) = P^t(X,Y)$$

**Missingness shift** occurs when the missing data mechanism differs between source s and target t:

 $P^s(\xi \mid \cdot) 
eq P^t(\xi \mid \cdot)$ 

Domain Adaptation under Missingness Shift Definition



Suppose missingness shift occurs between source domain s and target domain t. Labeled source data  $\tilde{X}^s, Y \sim P^s$ , unlabeled target data  $\tilde{X}^t \sim P^t$ .

Goal of **DAMS**: learn an optimal predictor on the corrupted target domain data.

Definitions

tion >

imation 🜔 Experi<u>ments</u>

#### One more detail...



Often, as in our example, missing data indicators are not observed.

(particularly tricky if substantial # of true 0s now indistinguishable from missing)

tification

Estimation

periments



### DAMS with Underreporting Completely at Random

To make this difficult setting tractable, we define the **DAMS with UCAR** setting, which we focus on throughout the rest of this work:

Assume  $\xi$  (unobserved) is parameterized by constant **unknown** missingness rates  $m^s \in [0, 1]^d$  in source and  $m^t \in [0, 1]^d$  in target. Then,

 $egin{aligned} \xi_j^s &\sim ext{Bernoulli}ig(1-m_j^sig) \ \xi_j^t &\sim ext{Bernoulli}ig(1-m_j^tig) \end{aligned}$ 

Definitions

# Cost of Non-Adaptivity





dentification

Estim

Experiment



#### Let's start with a simple example...

**Redundant Features** 



$$egin{aligned} W &= u_W & u_W \sim \mathcal{N}ig(0,\sigma_w^2ig) \ X_1 &= W & u_Y \sim \mathcal{N}ig(0,\sigma_y^2ig) \ X_2 &= W \ Y &= W + u_Y \end{aligned}$$

Suppose *W* is a latent variable,

and we observe 
$$\widetilde{X} = \Big[ \widetilde{X}_1, \widetilde{X}_2 \Big]$$
.

(Failing to adapt
⇒ performance no better
than guess of label mean)

dentification

### Motivating Example

#### **Confounded Features**

 $X_1$  $X_2 \rightarrow Y$ 

$$egin{aligned} X_1 &= u_1 & u_1 \sim \mathcal{N}(0,1) \ X_2 &= a X_1 + u_2 & u_2 \sim \mathcal{N}(0,1) \ Y &= b X_1 + c X_2 + u_Y & u_Y \sim \mathcal{N}(0,1) \end{aligned}$$

for constants a, b, c.

Suppose we observe 
$$\ \ \widetilde{X} = \left[ \widetilde{X}_1, \widetilde{X}_2 
ight]$$

(Failing to adapt
⇒ performance **arbitrarily worse**than guessing label mean)

Definition

Cost adapt

dentification

Estimatic

Experiments

## What if we **observed** missing data indicators $(1 - \xi)$ ?

If  $\xi$  depends on completely observed covariates (or is drawn completely at random), then missingness shift **satisfies the covariate shift assumption**:

$$P^{s}\Big(Y\mid \widetilde{X}'= ilde{x}'\Big)=P^{t}\Big(Y\mid \widetilde{X}'= ilde{x}'\Big),$$

where  ${\widetilde X}' = \left( {\widetilde X}, \xi 
ight).$ 

- There exists an optimal predictor that does not change across domains
- Covariate shift problems are relatively well-studied (Gretton et. al. 2009, Huang et. al. 2006, Shimodaira 2000, Sugiyama et. al. 2007, Tsymbal 2004)
- Extension: leveraging DAMS structure to estimate importance weights more efficiently

cation

### UCAR as L2 Regularization

- Observe that applying mask  $\xi$  (zeroing out covariates with some probability) resembles the mechanism of dropout in neural networks
- We show that for linear models, applying missingness rates to data scaled by  $\frac{1}{1-m}$  is approximately equivalent to **L2 regularization** of  $\beta$  scaled by  $\widetilde{\Delta}_{\text{diag}}$ , where  $\widetilde{\Delta}_{\text{diag}} = \text{diag}\left(\sqrt{\frac{m}{1-m}}\right) \text{diag}(\sqrt{I})$  and  $\text{diag}(\sqrt{I})$  is the Fisher information matrix.

ntification



# Identification



🔹 🔪 Cost <del>c</del>

Identification

Estim

Experimen



### First, let us define *m*-reachability.

For data points  $a \in \mathbb{R}^d$  and  $b \in \mathbb{R}^d$ , we say *b* is **m-reachable** from *a* (denoted  $a \rightsquigarrow b$ ) if there exists some mask  $\xi$  such that  $b = a \odot \xi$ . If  $a \rightsquigarrow b$ , then:

- Dimensions of *a* that are 0 must be a subset of the ones that are 0 in *b*
- Any dimensions that are nonzero in both *a* and *b* must match in value

Identification



#### Identifying corrupted from clean distribution (& vice versa!)

Consider any covariates  $x \in \mathbb{R}^d$  and label  $y \in \mathbb{R}$  .

Identification

Given m, it is straightforward to identify  $\tilde{p}$  from p:

Definitions

Note: we require knowledge of **missingness rates** *m* 

$$ilde{p}_{x,y} = \sum_{z:z \rightsquigarrow x, \ z \in \mathbb{R}^d} p_{z,y} \cdot \prod_{j=1}^d (1 - m_j)^{x_j} m_j^{z_j - x_j}$$
prob. of x, y in
prob. of z, y
in clean dist.
prob. of z, y
in clean dist.
dist.

Interestingly, given  $m \prec 1$ , we can show that p is identifiable from  $\tilde{p}$  (induction on # zeros).

7

### Unfortunately, *m* is not in general identifiable...

Consider two possible source distributions:

 $A \sim ext{Bernoulli}(0.5) \ B \sim ext{Bernoulli}(0.25)$ 

Applying missingness rates  $m_A = 0.5$  to A and  $m_B = 0$  to B yields identical observed corrupted distributions:

$$\widetilde{A} \stackrel{ ext{iid}}{\sim} \widetilde{B} \sim ext{ Bernoulli}(0.25)$$

Thus, the rates are not identifiable.

Identification

### But fortunately, we can identify...

#### 1) Relative missingness rates:

The ratio between non-missingness rates  $1 - m^t$  and  $1 - m^s$  is given by:

Identification

$$rac{P^t( ilde{x}
eq 0)}{P^s( ilde{x}
eq 0)} = rac{(1-m^t)\odot q}{(1-m^s)\odot q} = rac{1-m^t}{1-m^s} \stackrel{\Delta}{=} 1-r^{s
ightarrow t},$$

2) And thus, the labeled target distribution from the labeled source distribution

$${ ilde p}^t_{x,y} = \sum_{z:z \rightsquigarrow x, z \in \mathbb{R}^d} { ilde p}^s_{z,y} \cdot \prod_{j=1}^d ig(1-r_j^{s 
ightarrow t}ig)^{x_j}ig(r_j^{s 
ightarrow t}ig)^{z_j-x_j}$$

# **Estimation**



🔹 🔪 Cost <del>c</del>

Identifi

Estimation

Experiment



21

## Estimating relative missingness rates $r^{s ightarrow t}$

$$\text{Simply compute} \;\; \hat{q}_{j}^{s} = \frac{\operatorname{count} \left( \tilde{x}_{j}^{s} \neq 0 \right)}{n_{s}} \text{,} \quad \hat{q}_{j}^{t} = \frac{\operatorname{count} \left( \tilde{x}_{j}^{t} \neq 0 \right)}{n_{t}} \text{, and} \quad \hat{r}^{s \to t} = 1 - \frac{\hat{q}^{t}}{\hat{q}^{s}}.$$

Using Hoeffding's inequality, we show that with probability  $1 - \delta$ ,

$$\left| \hat{r}^{s o t} - r^{s o t} 
ight| \leq rac{1}{\hat{q}^s} \left( \sqrt{rac{\log\left(4/\delta
ight)}{2n_t}} + \left(1 - r^{s o t}
ight) \sqrt{rac{\log\left(4/\delta
ight)}{2n_s}} 
ight)$$

**Estimation** 

22

#### Non-parametric adjustment procedure

1. Compute 
$$\hat{r}^{s 
ightarrow t} = 1 - rac{\hat{q}^t}{\hat{q}^s}$$

Definitions

2. Check whether  $\hat{r}^{s \to t} \succeq 0$ . If not, this procedure is not a "proper" adjustment.

3. Compute 
$$ilde{r}^{s
ightarrow t}=\maxig(\hat{r}^{s
ightarrow t},0ig)$$
 , where max is elementwise.

4. Apply missingness with rate  $r^{s \rightarrow t}$  to source data to get data distributed identically to target data.

Estimation

Experiments

5. Fit a predictor on the (further) corrupted labeled source data.

#### **Closed-Form Solution for Optimal Linear Predictor**

The optimal linear target predictor is given by:

$$\beta_t^* = \mathbb{E}[\widetilde{X}^{t\top} \widetilde{X}^t]^{-1} \left( r^{s \to t} \odot \mathbb{E}[\widetilde{X}^{s\top} Y^s] \right).$$

We can estimate  $\mathbb{E}[\widetilde{X}^{t\top}\widetilde{X}^{t}]$  in two ways:

- 1. Directly from unlabeled target data
- 2. From source data: for  $i \neq j$ ,

$$\mathbb{E}[\widetilde{X}_t^T \widetilde{X}_t]_{ij} = (1 - r_i^{s \to t})(1 - r_j^{s \to t})\mathbb{E}[\widetilde{X}_s^T \widetilde{X}_s]_{ij}$$
$$\mathbb{E}[\widetilde{X}_t^T \widetilde{X}_t]_{ii} = (1 - r_i^{s \to t})\mathbb{E}[\widetilde{X}_s^T \widetilde{X}_s]_{ii}.$$

# Experiments



🔪 Cost <del>e</del>

entification

Estimati

Experiments



#### Synthetic Experiments

Definitions

We draw 10,000 samples from two simple data-generating processes:

Scenario 1: "Redundant Features"Scenario 2: "Confounded Features" $u_y \sim \mathcal{N}(0,1)$  $u_{x_2} \sim \mathcal{N}(0,1)$  $Z \sim \text{Bernoulli}(0.5)$  $u_y \sim \mathcal{N}(0,1)$  $X_1 = Z$  $X_1 \sim \text{Bernoulli}(0.5)$  $X_2 = Z$  $X_2 = \expit(2X_1 + u_{x_2})$  $Y = Z + u_y$  $Y = X_1 - X_2 + u_y$ 

Experiments

in both, we apply missingness with rates  $m^s = [1 - \epsilon, \epsilon]$  and  $m^t = [\epsilon, 1 - \epsilon]$ .

#### Synthetic Experiments

Scenario 1: "Redundant Features"

$$u_{y} \sim \mathcal{N}(0, 1)$$

$$Z \sim \text{Bernoulli}(0.5)$$

$$X_{1} = Z$$

$$X_{2} = Z$$

$$Y = Z + u_{y}$$

Scenario 2: "Confounded Features"

$$u_{x_2} \sim \mathcal{N}(0, 1)$$
  

$$u_y \sim \mathcal{N}(0, 1)$$
  

$$X_1 \sim \text{Bernoulli}(0.5)$$
  

$$X_2 = \text{expit}(2X_1 + u_{x_2})$$
  

$$Y = X_1 - X_2 + u_y$$



27

#### Semi-synthetic Experiments (real data, synthetic labels)

Average target domain error, given by MSE/Var(Y), on synthetic and semi-synthetic data:

	Rednd.	Confnd.	Confnd. Adult		Bank		Thyroid	
	$\overline{m^s?m^t}$	$m^s ? m^t$	$\overline{m^s \preceq m^t}$	$m^s ? m^t$	$m^s \preceq m^t$	$m^s ? m^t$	$m^s \preceq m^t$	$m^s ? m^t$
Linear Regression Models								
Oracle	0.178	0.206	0.420	0.362	0.338	0.433	0.298	0.251
Source	1.259	1.103	0.437	0.380	0.371	0.480	0.350	0.320
Imputed	1.002	0.918	0.490	0.483	0.501	0.592	0.306	0.358
Closed-form	0.186	0.209	0.422	0.363	0.339	0.442	0.316	0.291
Non-param.	0.473	0.492	0.420	0.373	0.338	0.459	0.293	0.291
XGBoost Models								
Oracle	0.166	0.200	0.398	0.354	0.287	0.453	0.316	0.274
Source	0.166	0.475	0.399	0.379	0.305	0.500	0.310	0.352
Imputed	1.002	1.157	0.512	0.521	0.492	0.708	0.355	0.441
Non-param.	0.425	0.473	0.399	0.392	0.287	0.503	0.310	0.381
MLP Models								
Oracle	0.166	0.201	0.389	0.343	0.295	0.458	0.279	0.230
Source	0.184	0.321	0.399	0.357	0.322	0.499	0.320	0.303
Imputed	1.003	0.924	0.480	0.468	0.484	0.668	0.304	0.345
Non-param.	0.436	0.470	0.389	0.355	0.294	0.487	0.278	0.272

Definition

Estimatio

#### Discussion

- When indicators are provided and missingness depends on completely observed covariates, DAMS can be viewed as a form of covariate shift
- When missing data indicators are not provided, we provide identification and estimation results for the DAMS with UCAR setting
- Experiments validate our findings when assumptions hold



### Extensions

- Seek real-world data well-suited for DAMS with UCAR
- Explore generalizations or relaxations of the DAMS with UCAR assumptions:
  - Allowing dependence on covariates
  - Other noise models



# Thank you!

Feel free to reach out: hlzhou@andrew.cmu.edu