

## 1. Causal Inference from Observational Data

**Ultimate Goal:** Finding a model that estimates the **Individual Treatment Effect**  $ITE(x) = y^1(x) - y^0(x)$

from an observational dataset in the form of  $\{[x_i, t_i, y_i]\}_{i=1..n}$

with:  $x$ : personal features

→ e.g., values of age, blood work, etc.

$t$ : received treatment chosen from a set of options

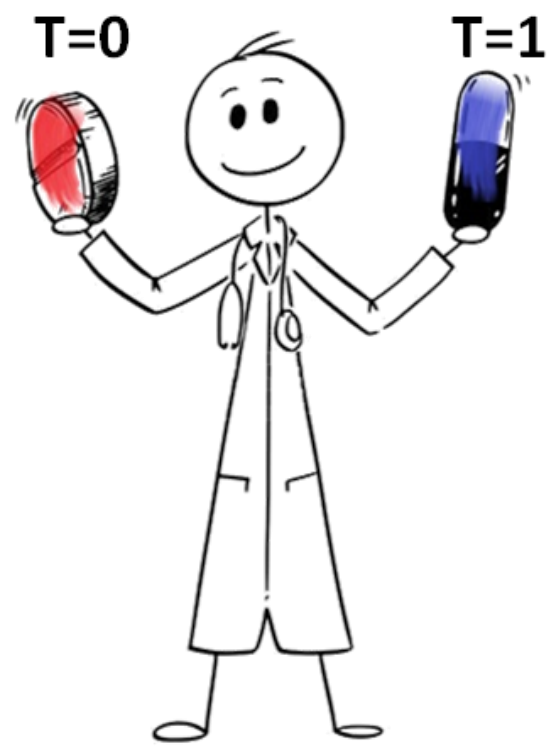
→ e.g., {0: surgery, 1: medication}

$y$ : the observed outcome after receiving the corresponding treatment

→ e.g., survival time



ID	X			T	Y <sup>0</sup>	Y <sup>1</sup>
	Gender	Age	BMI			
Mr. Smith	Male	35	20	0	15	
Mr. Green	Male	22	32	0	22	
Ms. Jones	Female	20	23	1		31
...						

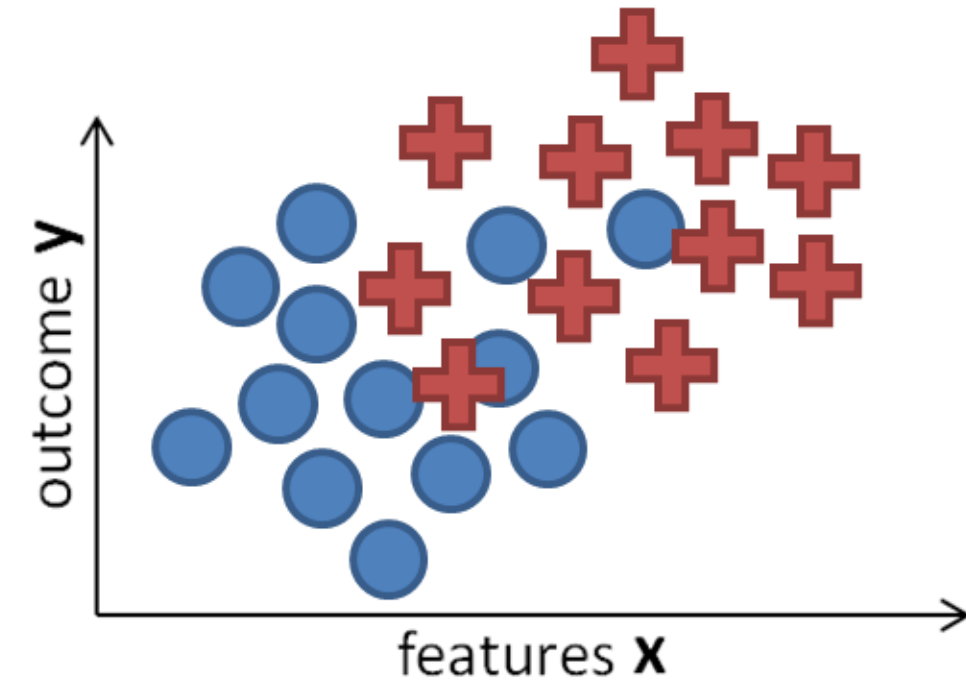


**Challenges:**

1. **Partial information data.** i.e., depending on the received treatment  $t$ , we observe (factual outcome  $y^t$ ) either  $y^0$  or  $y^1$ , but never both. The other outcome (counterfactual outcome  $y^{t'}$ ) is **unobservable**.

2. **Sample selection bias.** i.e., both outcome  $y$  and the treatment  $t$  assignment are dependent on (some) context information  $x$ .

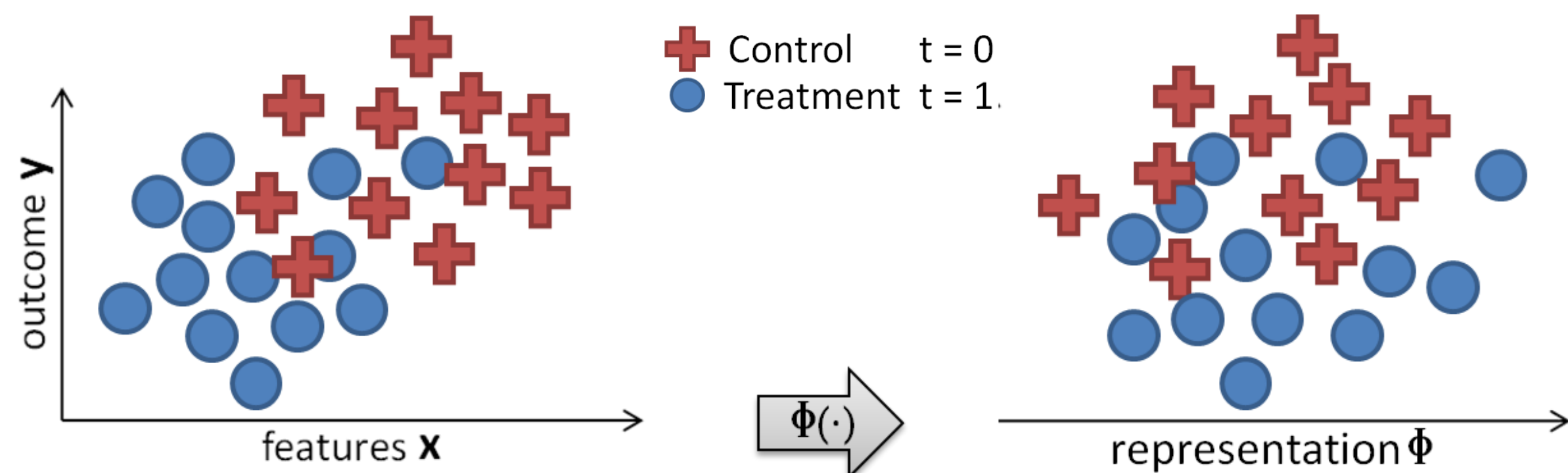
→ e.g., younger {older} patients (part of  $x$ ) are more likely to receive treatment  $t$ : surgery {medication} because they tend to have a faster {complicated} recovery (outcome  $y$ ).



## 2. Shalit et al. (2017) Model Overview (CFR)

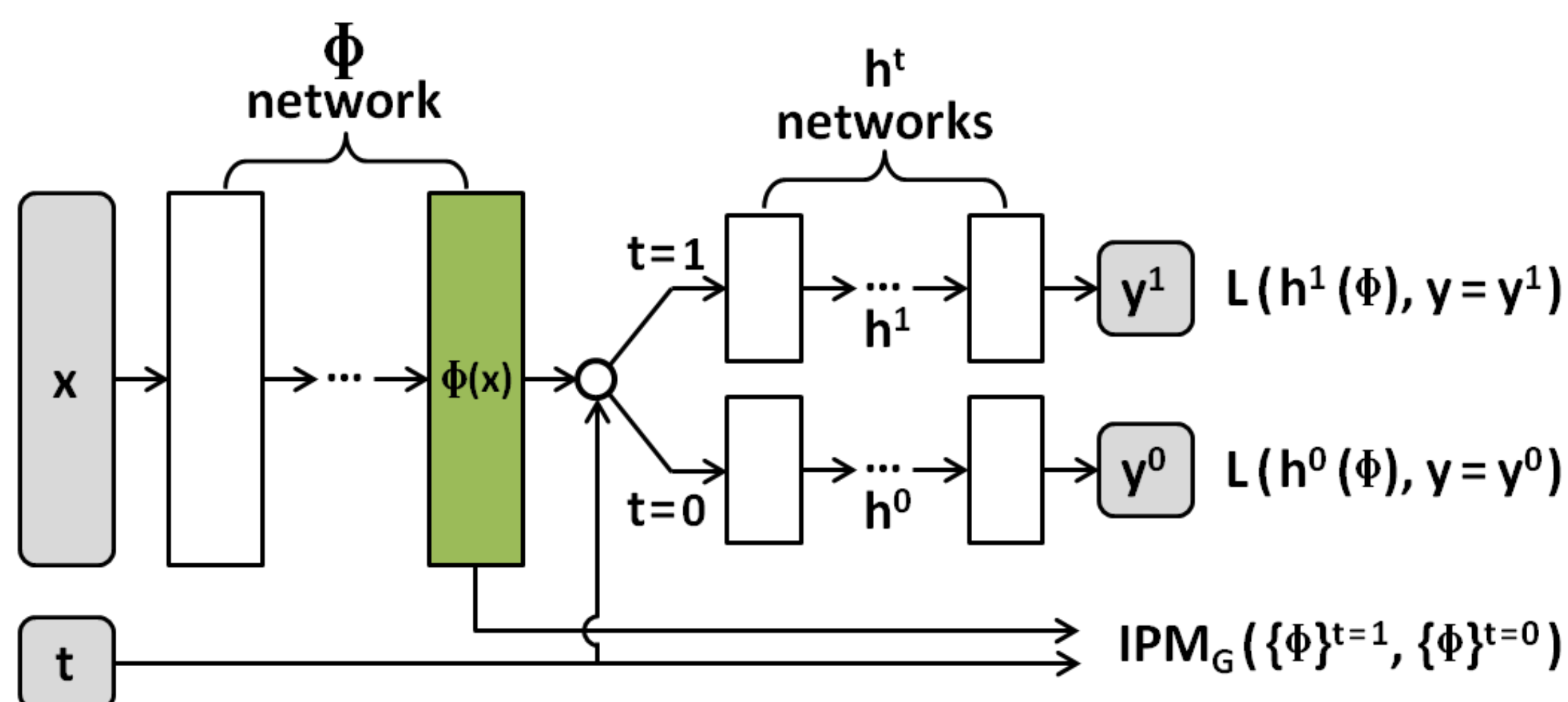
**Their Goal:** Reducing the sample selection bias by learning a common representation space  $\Phi(x)$  such that:

- $\Pr(\Phi(x) | t=0)$  and  $\Pr(\Phi(x) | t=1)$  are as close as possible to each other
- provided that  $\Phi(x)$  retains enough information to accurately predict factual outcomes
- by a learned hypothesis network for each treatment arm (i.e.,  $h^t(x)$ ) that estimate the corresponding outcomes



**Model**

**Structure:**



**Objective:** 
$$\arg \min_{h, \Phi} J(h, \Phi) = \arg \min_{h, \Phi} \left[ \frac{1}{n} \sum_{i=1}^n \omega_i \cdot L[h^{t_i}(\Phi(x_i)), y_i] + \alpha \cdot \text{IPM}_G(\{\Phi(x_i)\}_{i: t_i=0}, \{\Phi(x_i)\}_{i: t_i=1}) + \lambda \cdot \mathcal{R}(h) \right]$$
 → regularization term

where  $L[h^{t_i}(\Phi(x_i)), y_i] = [h^{t_i}(\Phi(x_i)) - y_i]^2$  → **factual loss**

$$\omega_i = \frac{t_i}{u} + \frac{1-t_i}{(1-u)}, \quad \text{with } u = \frac{1}{n} \sum_{i=1}^n t_i = \Pr(t=1)$$

$$\Rightarrow \frac{1}{\Pr(t_i)} = \frac{\Pr(t_i)}{\Pr(t_i)} + \frac{1-\Pr(t_i)}{\Pr(t_i)} = 1 + \frac{\Pr(\neg t_i)}{\Pr(t_i)}$$

$\text{IPM}_G(\{\Phi(x_i)\}_{i: t_i=0}, \{\Phi(x_i)\}_{i: t_i=1})$  → Integral Probability Metric (IPM) is a measure of closeness between two probability distributions; e.g., Maximum Mean Discrepancy (MMD) (Gretton et al., 2012) or Wasserstein distance (Attouch et al., 2014; Cuturi & Doucet, 2014).

Here, IPM measures the discrepancy between empirical  $\Pr(\Phi(x) | t=0)$  and  $\Pr(\Phi(x) | t=1)$  distributions

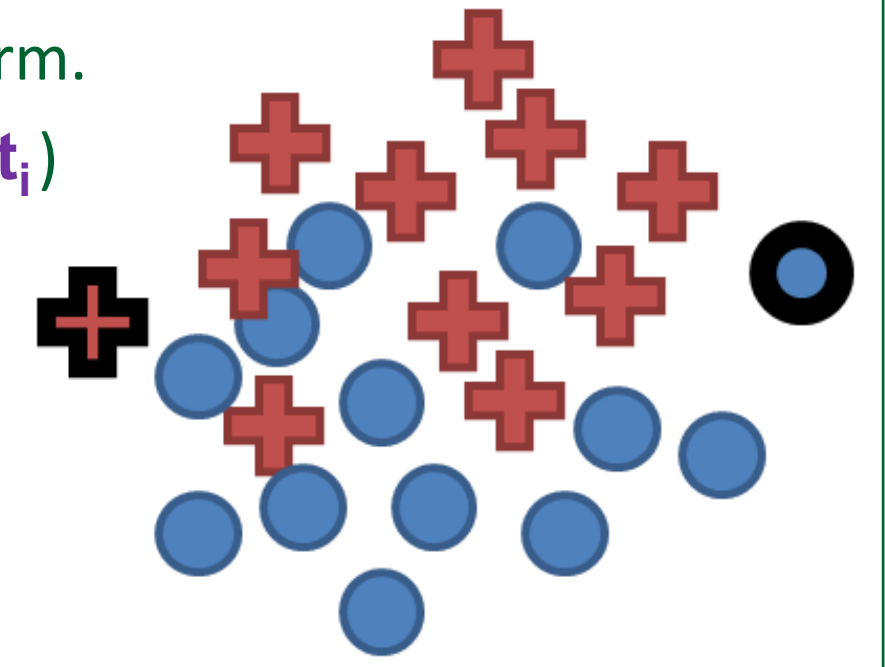
Once the model is trained, we can use it to predict  $y^1$  and  $y^0$ , given as input a feature vector  $x$

This will give us the individual treatment effect  $ITE(x) = y^1(x) - y^0(x)$  for any (novel)  $x$

## 3. Proposed Weighting Scheme (CFR-ISW)

**Our Goal:** Improving the accuracy of estimating **ITE** by incorporating the information extracted from the context of each instance  $\Phi(x)$ , in addition to its respective treatment  $t$ , to assign **sample-specific weights** in the **factual loss** term.

→ For example, if an instance  $x_i$  (with assigned treatment  $t_i$ ) is **far** from other instances with the same assigned treatment (e.g., **samples** in figure) then we force our outcome prediction network to learn this instance well.



**Proposed weights:**

$$\frac{\Pr(\Phi(x_i) | \neg t_i)}{\Pr(\Phi(x_i) | t_i)} = \frac{\frac{\pi(\neg t_i | \Phi(x_i)) \cdot \Pr(\Phi(x_i))}{\Pr(\neg t_i)}}{\frac{\pi(t_i | \Phi(x_i)) \cdot \Pr(\Phi(x_i))}{\Pr(t_i)}} = \frac{\Pr(t_i)}{1 - \Pr(t_i)} \cdot \frac{1 - \pi(t_i | \Phi(x_i))}{\pi(t_i | \Phi(x_i))}$$

where  $\pi(t | \Phi(x))$  is the probability of assigning treatment  $t$  given the context in  $\Phi$  representation space (a.k.a., **propensity score**).

→ We use Logistic Regression (LR) with parameters  $[W, b]$  to fit the propensity score function:

$$\pi(t | \Phi(x)) = \frac{1}{1 + e^{-(2t-1)(\Phi(x) \cdot W + b)}}$$

and learn the parameters by minimizing:  $\min_{W, b} \frac{1}{n} \sum_{i=1}^n C[W, b, \Phi(x_i), t_i]$

where  $C[W, b, \Phi(x), t] = -\log[\pi(t_i | \Phi(x_i))]$

We try to solve this multi-objective optimization problem that iteratively in two steps:

- Minimize  $J(h, \Phi)$  to update the parameters of the representation  $\Phi$  and hypothesis  $h$  networks
- Minimize  $C[W, b, \Phi, t]$  with fixed  $h$  and  $\Phi$  parameters to update parameters of the propensity score function (i.e.,  $W$  and  $b$ ).

## 4. Experiments

**Evaluation Criteria:**  $\text{ENoRMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(1 - \frac{\hat{e}_i}{e_i}\right)^2}$  with  $\begin{cases} \hat{e}_i = \hat{y}_i^1 - \hat{y}_i^0 \\ e_i = y_i^1 - y_i^0 \end{cases}$

where  $\hat{y}$  indicates an outcome predicted by the trained model

**Hyperparameter Selection:** As counterfactual outcomes are **inherently unobservable**, it is not possible to use standard internal cross-validation to select hyperparameters (e.g.,  $\alpha, \lambda$ , etc.).

→ An estimation of the true effect is needed as a **surrogate** for the  $e$  term.

- ❖ Shalit et al. (2017) used the observed outcome  $y_{j(i)}$  of the nearest neighbor in the  $x$  space (referred to as **1-nn**) in the alternative treatment group  $t_{j(i)} = \neg t_i = 1 - t_i$
- ❖ In addition to **1-nn**, we explored two alternatives:
  - 1-nn in the  $\Phi$  space; i.e., **1-nn<sub>Φ</sub>**
  - outcome predicted by the Bayesian Additive Regression Trees (**BART**)

**Synthetic Datasets:** From the 2018 Atlantic Causal Inference Data Challenge

→ The  $x$  matrix is sampled from the Linked Birth and Infant Death Data (**LBIDD**)

❖ **100,000** instances, each with **177** features

→ **24** synthetic datasets were generated from LBIDD;

categorized into **6** groups in terms of the number of instances  $n \in \{1, 2.5, 5, 10, 25, 50\} \times 10^3$

## 5. Results

We compare performance of the following four different methods in terms of ENoRMSE:

- **1-nn:** One nearest neighbor method for finding the counterfactual outcomes
- **BART:** Bayesian Additive Regression Trees method (Chipman et al., 2010) for finding the ITE
- **CFR:** Counterfactual Regression method proposed in (Shalit et al., 2017) for which the best set of hyperparameters is determined based on  $\text{ENoRMSE}_{\text{BART}}$
- **CFR-ISW:** Counterfactual Regression with Importance Sampling Weights (i.e., the proposed method)

Tables report the aggregated ENoRMSE (lower is better). The entry in **bold** is the best for each row.

Comparison of various ITE estimation methods against the proposed CFR-ISW					Hyperparameter selection methods: ENoRMSE <sub>1-nn</sub> vs. ENoRMSE <sub>BART</sub>				
DATASETS	1-NN	BART	CFR	CFR-ISW	CFR		CFR-ISW		
					1-NN	BART	1-NN	BART	
ALL	75.32	20.03	8.92	<b>1.07</b>	15.04	<b>8.92</b>	5.23	<b>1.07</b>	
# INSTANCES	1 k	86.44	141.58	10.51	<b>1.72</b>	<b>8.08</b>	10.51	2.21	<b>1.72</b>
	2.5 k	47.45	23.35	15.27	<b>0.73</b>	36.33	<b>15.27</b>	0.82	<b>0.73</b>
	5 k	38.04	9.51	2.81	<b>0.93</b>	5.79	<b>2.81</b>	1.05	<b>0.93</b>
	10 k	40.25	2.96	1.22	<b>0.81</b>	1.45	<b>1.22</b>	15.11	<b>0.81</b>
	25 k	25.69	1.52	<b>0.89</b>	1.03	1.01	<b>0.89</b>	1.28	<b>1.03</b>
#	50 k	94.45	16.13	11.12	<b>1.14</b>	17.98	<b>11.12</b>	1.14	1.14

**Selected References:**

- Chipman et al., 2010) Chipman, Hugh A, George, Edward I, and McCulloch, Robert E. "BART: Bayesian additive regression trees." *The Annals of Applied Statistics*, 2010.
- (Shalit et al., 2017) Shalit, Uri, Johansson, Fredrik D., and Sontag, David. "Estimating individual treatment effect: generalization bounds and algorithms." *In Proceedings of the 34th International Conference on Machine Learning (ICML)*, pp. 3076–3085, 2017.

Get the paper @

