

### Problem setup: CATE estimation from high-dimensional observational data

**Input** Observational data  $D \stackrel{i.i.d.}{\sim} P(A, X, Y)$

**Output** Conditional Average Treatment Effect (CATE)

**Difficulty in high-dimensional setup:**

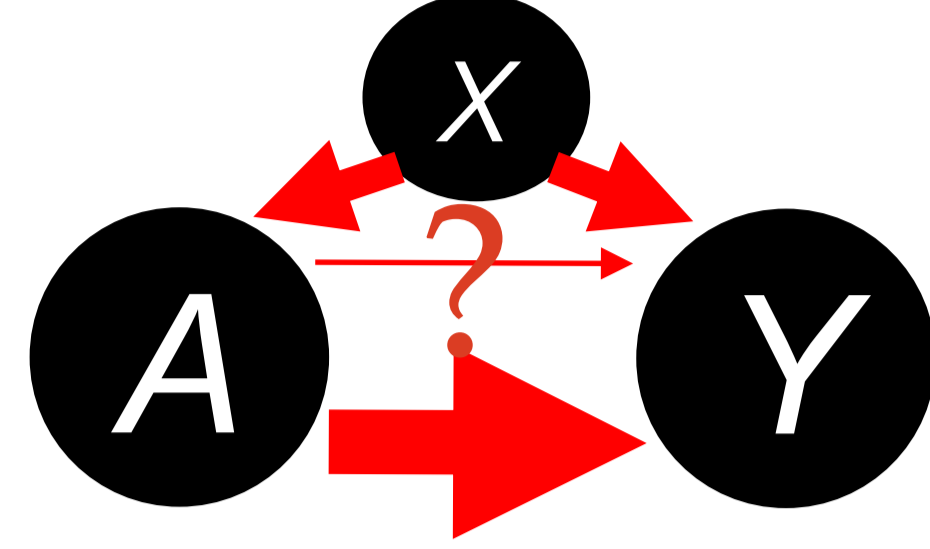
Treatment  $A \in \{0, 1\}$  Features  $X$  Outcome  $Y$  Potential outcomes  $Y^a$

$A$	$X_1$	...	$X_d$	$Y$	$Y^0$	$Y^1$
0	...	...	...	0.2	0.2	?
1	...	...	...	1.1	?	1.1

High-dimensional  $Y = AY^1 + (1 - A)Y^0$  Observed when  $A = a$

$$\text{CATE}(\mathbf{x}) := \mathbb{E}[Y^1 - Y^0 | X = \mathbf{x}]$$

Average treatment effect across individuals with identical feature attributes  $X = \mathbf{x}$



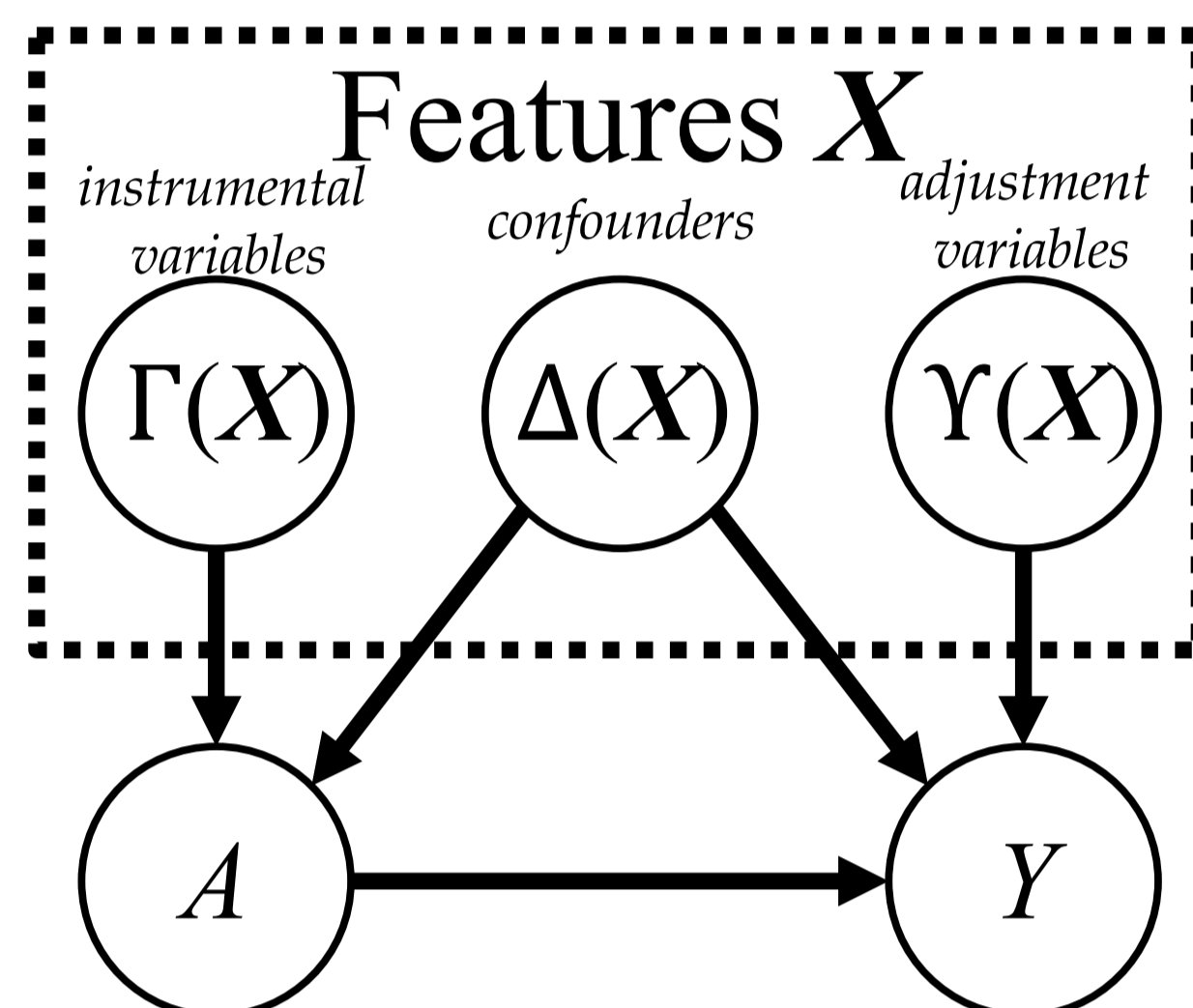
- We often have little prior knowledge about features  $X$
  - Confounders and adjustment variables are often *intertwined*
- Confounder (e.g., age) leads to sample selection bias ⊗
- Adjustment variable (e.g., smoking habit) is predictive of potential outcome ⊗

**Question:**

Can we remove sample selection bias due to confounders while keeping predictive information of adjustment variables?

### Weighted representation learning

**Existing Approach:** Data-driven feature separation



$$A \sim \pi(\Gamma(\mathbf{x}), \Delta(\mathbf{x})) \quad \hat{Y}^a = h^a(\Delta(\mathbf{x}), Y(\mathbf{x}))$$

- Fit propensity score  $\pi(\cdot)$  by

$$\min_{\pi} -\frac{1}{n} \sum_{i=1}^n (a_i \log(\pi(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i))) + (1 - a_i) \log(1 - \pi(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i)))) + \lambda_{\pi} \Omega(\pi)$$

Binary cross entropy loss

- Train  $\Gamma(\cdot), \Delta(\cdot), Y(\cdot), h^0(\cdot), h^1(\cdot)$  by

$$\min_{\Gamma, \Delta, Y, h^0, h^1} \frac{1}{n} \sum_{i=1}^n w_i |y_i - h^a(\Delta(\mathbf{x}_i), Y(\mathbf{x}_i))| + \lambda_{\Gamma} \text{MMD}(\{\Gamma(\mathbf{x}_i)\}_{i:a_i=0}, \{\Gamma(\mathbf{x}_i)\}_{i:a_i=1}) + \lambda_{\Delta} \Omega(\Gamma, \Delta, Y, h^0, h^1)$$

Weighted prediction loss Penalize dependence of  $Y(\mathbf{X})$  on  $A$

$$\text{where } w_i = \frac{P(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i) | A = a_i)}{P(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i) | A = a_i) + P(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i) | A = 1 - a_i)} + \frac{P(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i) | A = 1 - a_i)}{P(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i) | A = a_i) + P(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i) | A = 1 - a_i)}$$

$$\propto \frac{1}{P(A = a_i | \Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i))} := \frac{1}{\pi_{a_i}(\Gamma(\mathbf{x}_i), \Delta(\mathbf{x}_i))}$$

**Weakness:** Inverse probability weight  $w_i$  is numerically unstable: Even slight propensity score estimation error leads to large CATE estimation error

### Weight smoothing with Pareto smoothing

**Advantage:**

- Can obtain a less biased estimator than weight truncation
- Can be combined with self-normalization

**Main idea** Improve CATE estimation stability by Pareto smoothing

- Pareto smoothing [Vehari+; JMLR2024]: Replace the  $M + 1$  largest importance sampling weight values with inverse CDF of generalized Pareto distribution

$$w_{[j]} = \mathbf{I}(j \geq n - M + 1) \hat{F}^{-1} \left( \frac{j - (n - M) - 1/2}{M} \right) + (1 - \mathbf{I}(j \geq n - M + 1)) w_{[j]}$$

$$\text{where } w_{[1]} \leq \dots \leq w_{[n]}, M = \min \left\{ \left\lfloor \frac{n}{3} \right\rfloor, \lfloor 3 \sqrt{n} \rfloor \right\}, \text{ and } F(w) = \begin{cases} 1 - \left(1 + \frac{\xi(w - \mu)}{\sigma}\right)^{-\xi} & (\xi \neq 0) \\ 1 - e^{-\frac{w - \mu}{\sigma}} & (\xi = 0) \end{cases}$$

- Need to compute rank  $\mathbf{r} = \mathbf{r}(\mathbf{w})$ :

- Example:** If  $w_3 \leq w_1 \leq w_2$ , since  $w_1 = w_{[2]}$ ,  $w_2 = w_{[3]}$ ,  $w_3 = w_{[1]}$ ,  $\mathbf{r} = [2, 3, 1]$

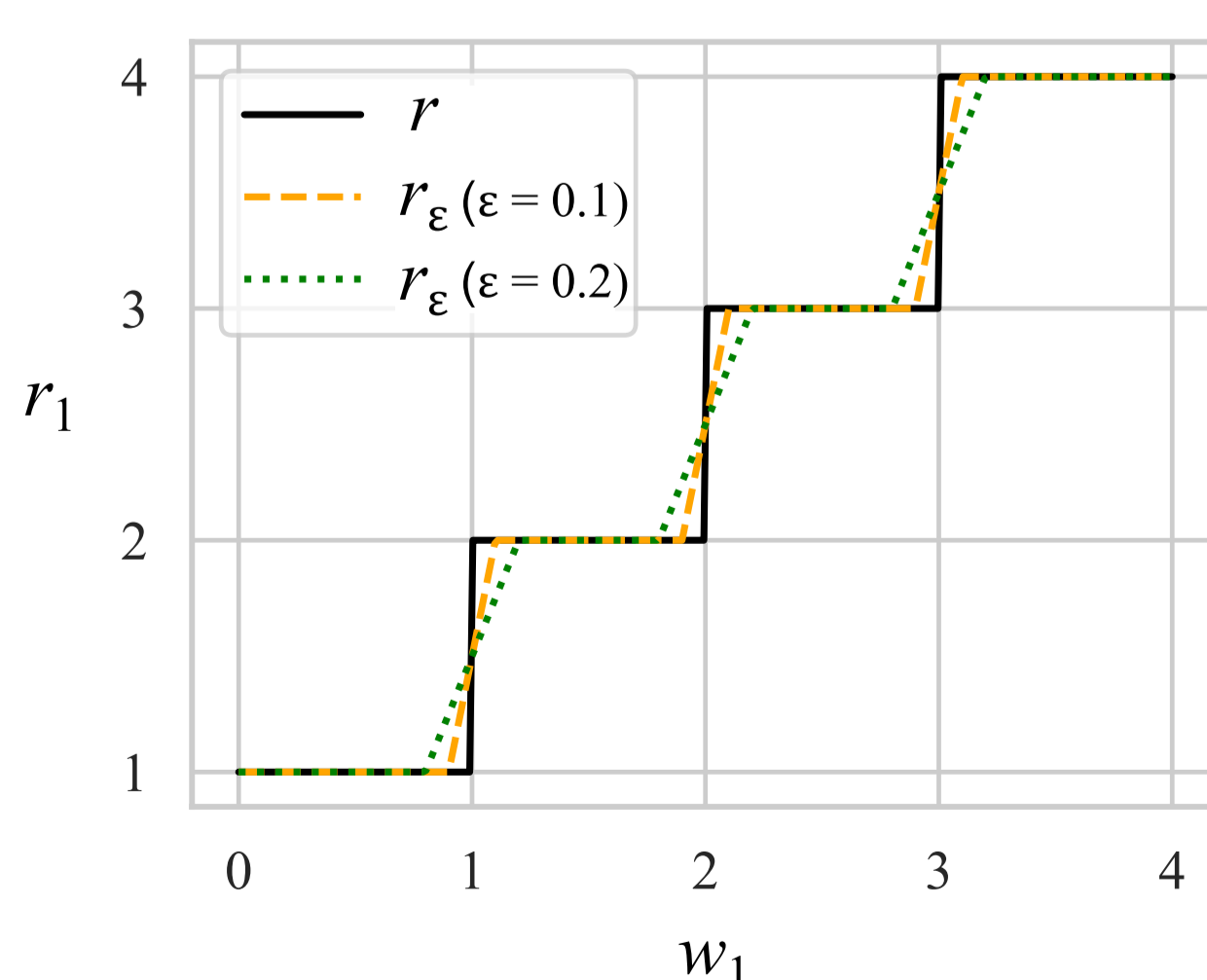
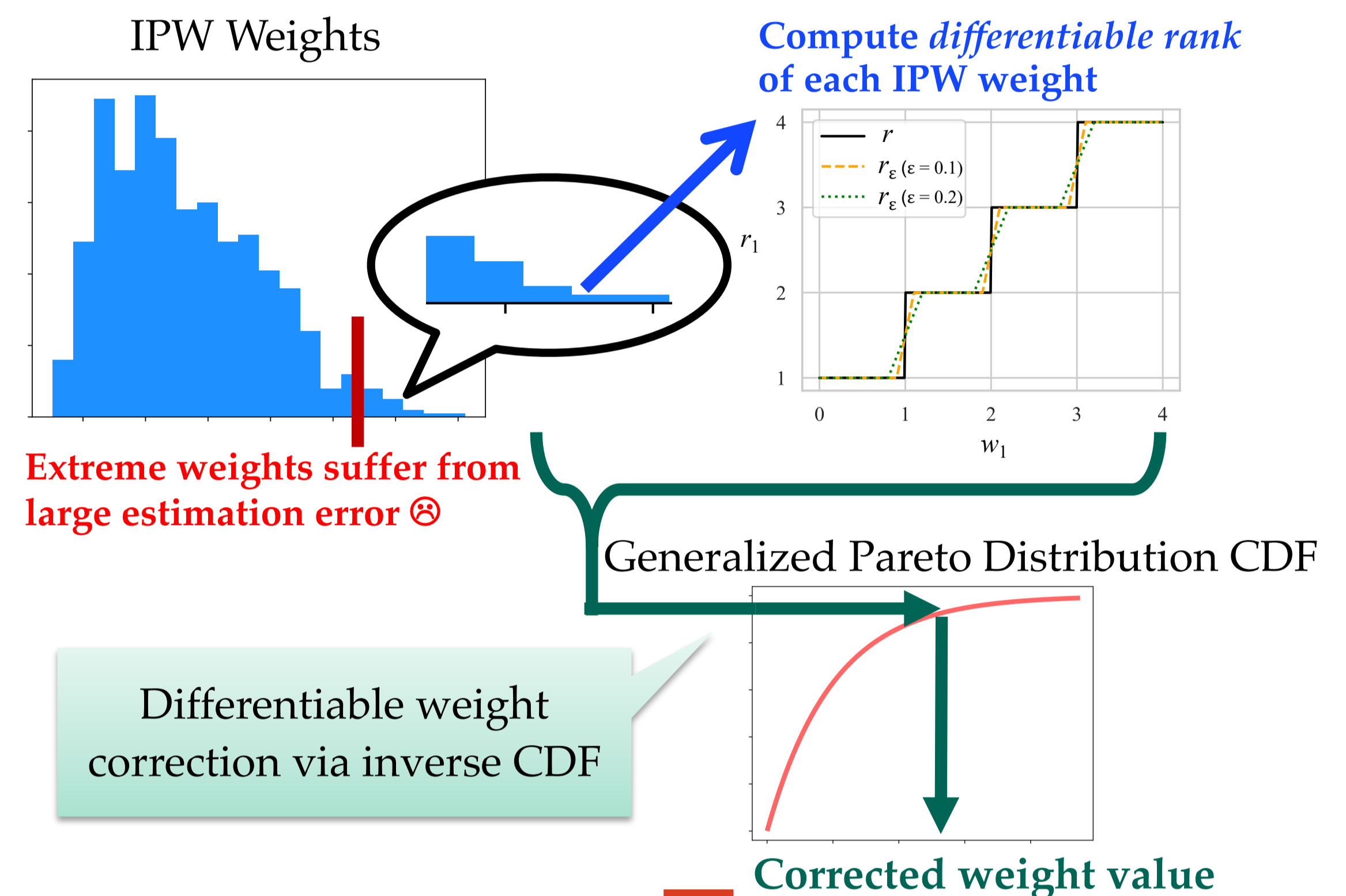


Figure 2: Illustration of rank function  $\mathbf{r} = \mathbf{r}(\mathbf{w})$  (black) and differentiable rank functions  $\mathbf{r} = \mathbf{r}_{\epsilon}(\mathbf{w})$  (orange and green). Here we take input vector  $\mathbf{w} = [w_1, 1, 2, 3]^T$ , vary  $w_1$ 's value, and look at how its rank  $r_1 \in \mathbb{R}$  changes. When regularization parameter  $\epsilon \rightarrow 0$ ,  $\mathbf{r}_{\epsilon}$  converges to  $\mathbf{r}$  [Blondel et al., 2020].

**Difficulty:**  $\mathbf{r}(\mathbf{w})$  is piecewise constant: Gradient is always zero or undefined. We cannot perform gradient back propagation ⊗

### Proposed method

**Proposed method:** Pareto smoothing + Differentiable ranking



End-to-end weighted representation learning for CATE estimation from high-dimensional observational data

- Approximate non-differentiable rank function  $\mathbf{r}$  with differentiable one

- Fast soft rank [Blondel+; ICML2020]: Approximate as a solution to regularized LP

- Approximate indicator function with sigmoid

$$\mathbf{I}(i \geq j) \simeq \zeta(i, j) := \frac{1}{1 + e^{-\kappa(i-j)}}$$

Combining 1. & 2. leads to the following weight replacement formula:

$$\tilde{w}_i = \zeta(r_i, n - M + 1) \tilde{F}^{-1} \left( \zeta \left( \frac{r_i - (n - M) - 1/2}{M} \right) \right) + (1 - \zeta(r_i, n - M + 1)) w_i$$

where  $\zeta(x) := \min \{ \max \{ x, 0 \}, 1 \}$

**Algorithm 1** Differentiable Pareto-Smoothed Weighting (DPSW)

- Initialize the parameters of  $\Gamma, \Delta, Y, \pi, h^0$ , and  $h^1$
- while not converged do
- while not converged do
- Sample mini-batch from  $\mathcal{D} = \{(a_i, \mathbf{x}_i, y_i)\}_{i=1}^n$
- Update  $\pi$  by minimizing cross entropy loss in (2)
- end while
- while not converged do
- Sample mini-batch from  $\mathcal{D} = \{(a_i, \mathbf{x}_i, y_i)\}_{i=1}^n$
- for instance  $i$  in mini-batch do
- Compute weight  $w_i$  by (4)
- end for
- Compute differentiable rank  $\mathbf{r} = \mathbf{r}_{\epsilon}(\mathbf{w})$
- Estimate GPD parameters as  $\hat{\mu}, \hat{\sigma}$ , and  $\hat{\xi}$
- for instance  $i$  in mini-batch do
- Replace each weight  $w_i$  with  $\tilde{w}_i$  in (20)
- end for
- end while
- Update  $\Gamma, \Delta, Y, h^0$ , and  $h^1$  by minimizing prediction loss in (3) with Pareto-smoothed weights  $\{\tilde{w}_i\}$
- end while

**Experimental results:**

**Semi-synthetic data**

Table 1: Mean and standard deviation of test PEHE on semi-synthetic datasets (Lower is better)

Method	News ( $d = 3477$ )	ACIC ( $d = 177$ )
LR-1	3.35 ± 1.28	0.72 ± 0.07
LR-2	5.36 ± 1.75	3.82 ± 0.15
SL	2.83 ± 1.11	1.69 ± 0.52
TL	2.55 ± 0.82	2.23 ± 0.50
XL	2.77 ± 1.01	1.05 ± 0.72
DRL	23.9 ± 5.96	3.77 ± 8.96
CF	3.84 ± 1.67	3.55 ± 0.19
CF DML	2.69 ± 1.06	1.18 ± 0.32
TARNet	4.92 ± 1.80	3.31 ± 0.51
GANITE	2.68 ± 0.66	3.69 ± 0.77
DRCFR	2.38 ± 0.66	0.98 ± 0.07
DRCFR Norm.	2.37 ± 0.94	0.73 ± 0.12
DRCFR Trunc.	2.42 ± 0.79	1.06 ± 0.06
PSW	4.03 ± 1.35	0.71 ± 0.01
DPSW	2.20 ± 0.72	0.57 ± 0.03
DPSW Norm.	2.10 ± 0.66	0.52 ± 0.16

**Synthetic data**

Randomly generate features as

$$[X_{\Gamma}, X_{\Delta}, X_{Y}]^T \in \mathbb{R}^d \quad (d = 15, 18, \dots, 30)$$

Measure the relative difference of average absolute values of the first-layered weight submatrices, e.g.,

$$\frac{|\bar{w}_{\Gamma}^1 - \bar{w}_{\Gamma}^1|}{|\bar{w}_{\Gamma}^1|} \quad \text{where } \mathbf{W}^1 = [\mathbf{W}_{\Gamma}^1, \mathbf{W}_{\Delta}^1, \mathbf{W}_{Y}^1]$$

