

# Singular Propensity Scores: Reducing Variance in Weighted Estimators

Joint Statistical Meeting 2024, Portland, USA

August 5th, 2024

**Kosuke Morikawa**<sup>1</sup> and Keisuke Yano<sup>2</sup>

Graduate School of Engineering Science, Osaka University, Japan<sup>1</sup>

Institute of Statistical Mathematics, Japan<sup>2</sup>

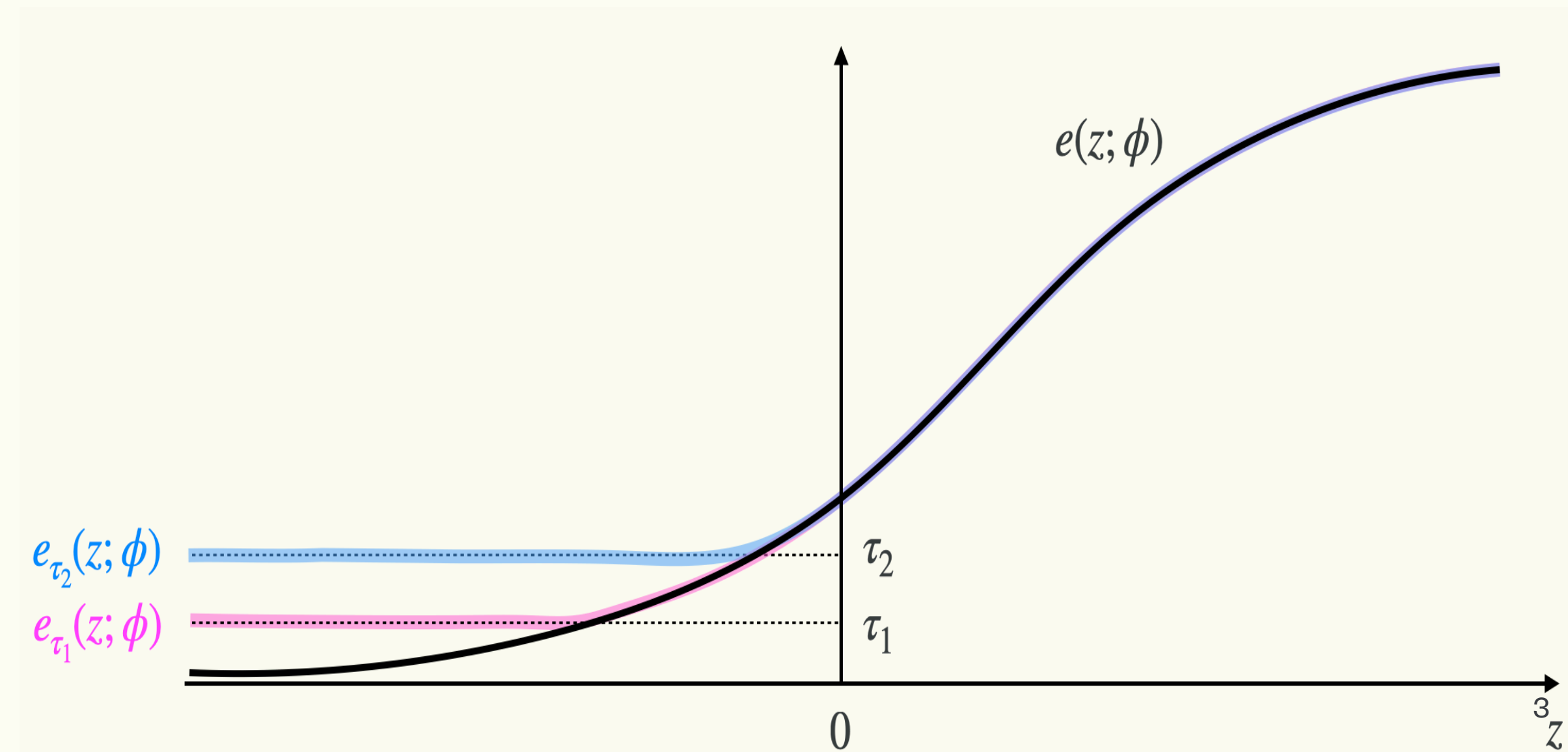
# Introduction

# Brief Summary

- Inverse weighted estimators  $\frac{1}{n} \sum_{i=1}^n \frac{Y_i}{e_i(\phi)}$  have been widely used in many fields, including **causal inference**, **missing data analysis**, and **survey sampling**.

Propensity Score

- However, inverse weights  $1/e_i(\phi)$  can diverge to infinity
- Our contribution is...
  1. to introduce a trimmed propensity score  $e_{\tau,i}(\phi)$  with an appropriate threshold  $\tau$
  2. to propose an information criterion to select the best threshold  $\tau^*$



# Settings: Average Treatment Effect

- $A_i \in \{0,1\}$ : **Treatment**
- $Y_i = A_i Y_i(1) + (1 - A_i) Y_i(0)$ : **Observed outcome**
  - $Y_i(a)$  ( $a = 0,1$ ): **Potential outcomes**
- $X_i$ : **Pre-treatment covariates**
- $Z_i = (X_i, Y_i)$
- $e(z; \phi) = P(A = 1 \mid Z = z; \phi)$ : **Propensity score** (including MNAR case)
- $\theta$ : **Interesting parameter**



✓ **Strongly Ignorability (or MAR).**  
 $A \perp \{Y(0), Y(1)\} \mid X$

• If  $\theta = E\{Y(1)\}$ , an estimator for  $\theta$  is  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \frac{A_i}{e(Z_i)} Y_i$  ← We focus on this estimator

• If  $\theta = E\{Y(0)\}$ ,  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \frac{1 - A_i}{1 - e(Z_i)} Y_i$

# Estimating Functions & Evaluation Metric — Regular Setup —

- $s_\phi(\phi; Z_i, A_i)$ : Scoring function for  $\phi$ 
  - e.g., log-likelihood under MAR:  

$$s_\phi(\phi; Z_i, A_i) = A_i \log e(\phi; X_i) + (1 - A_i) \log \{1 - e(\phi; X_i)\}$$
- $s_\theta(\theta; \phi; Z_i, A_i)$ : Scoring function for  $\theta$ 
  - e.g., IPW estimator  $-\frac{A_i}{e(Z_i; \hat{\phi})}(\theta - Y_i)^2$

$$\sum_{i=1}^n s_\phi(\phi; Z_i, A_i) \xrightarrow{\text{maximizer}} \hat{\phi}$$

$$\sum_{i=1}^n s_\theta(\theta; \hat{\phi}; Z_i, A_i) \xrightarrow{\text{maximizer}} \hat{\theta}$$

- $\nu(\hat{\theta}; \phi^*; Z_i, A_i)$ : Evaluation metric for  $\hat{\theta}$ 
  - e.g.,  $\frac{A_i}{e(Z_i; \phi^*)}(\hat{\theta} - Y_i)^2 \rightarrow$  Objective function  $E_{\hat{\theta}, \tilde{Y}}\{(\hat{\theta} - \tilde{Y})^2\}$ 

True (fixed, for the time being...)

↑

Train

↑

Test

$$E_{\tilde{Z}, \tilde{A}} \left[ E_{\hat{\theta}} \left[ \nu(\hat{\theta}; \phi^*; \tilde{Z}, \tilde{A}) \right] \right] \xrightarrow{\text{minimizer}} \text{Good model (w.r.t. MSE)}$$

Generalization Error (GE)

# Estimating Functions & Evaluation Metric — Regular Setup —

- $s_\phi(\phi; Z_i, A_i)$ : Scoring function for  $\phi$ 
  - e.g., log-likelihood under MAR:  

$$s_\phi(\phi; Z_i, A_i) = A_i \log e(\phi; X_i) + (1 - A_i) \log \{1 - e(\phi; X_i)\}$$
- $s_\theta(\theta; \phi; Z_i, A_i)$ : Scoring function for  $\theta$ 
  - e.g., IPW estimator  $\frac{A_i}{e(Z_i; \hat{\phi})} (\theta - Y_i)^2$   $\longrightarrow$  **unstable**

- $\nu(\hat{\theta}; \phi^*; Z_i, A_i)$ : Evaluation metric for  $\hat{\theta}$ 
  - e.g.,  $\frac{A_i}{e(Z_i; \phi^*)} (\hat{\theta} - Y_i)^2$   

True (fixed, for the time being...)

$$\sum_{i=1}^n s_\phi(\phi; Z_i, A_i) \xrightarrow{\text{maximizer}} \hat{\phi}$$

$$\sum_{i=1}^n s_\theta(\theta; \hat{\phi}; Z_i, A_i) \xrightarrow{\text{maximizer}} \hat{\theta}$$

$$E_{\tilde{Z}, \tilde{A}} \left[ E_{\hat{\theta}} \left[ \nu(\hat{\theta}; \phi^*; \tilde{Z}, \tilde{A}) \right] \right] \xrightarrow{\text{minimizer}} \text{Good model (w.r.t. MSE)}$$

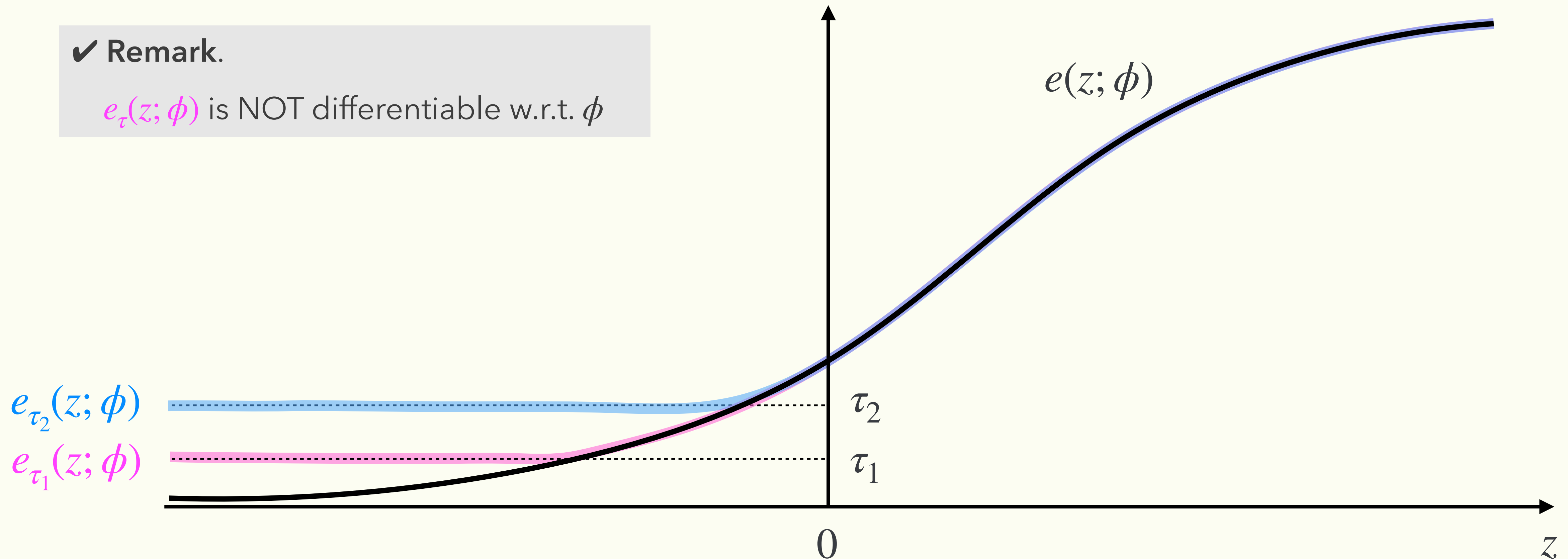
Generalization Error (GE)

# Trimmed propensity score

$$e_{\tau}(Z_i; \phi) = \tau I\{e(Z_i; \phi) \leq \tau\} + e(Z_i; \phi) I\{e(Z_i; \phi) > \tau\}$$

✓ Remark.

$e_{\tau}(z; \phi)$  is NOT differentiable w.r.t.  $\phi$



# Estimating Functions & Evaluation Metric — Singular Setup —

- $s_\phi(\phi; Z_i, A_i)$ : Scoring function for  $\phi$
- $s_\theta(\theta; \phi; Z_i, A_i)$ : Scoring function for  $\theta$

- e.g., IPW estimator  $\frac{A_i}{e_\tau(Z_i; \phi)} (\theta - Y_i)^2$   
② e<sub>τ</sub>(Z<sub>i</sub>; φ) Trimmed propensity

- $\nu(\hat{\theta}[\phi]; \phi^*; Z_i, A_i)$ : Evaluation metric for  $\hat{\theta}[\phi]$

- e.g.,  $\frac{A_i}{e(Z_i; \phi^*)} (\hat{\theta}[\phi] - Y_i)^2$

① **Bayesian Approach**  $\hat{\phi} = \arg \max_{\phi} \exp \left( \sum_{i=1}^n s_\phi(A_i, Z_i) \right) \pi(\phi)$

Pseudo likelihood Prior

$$\phi \sim \exp \left( \sum_{i=1}^n s_\phi(\phi; A_i, Z_i) \right) \pi(\phi)$$

$$\sum_{i=1}^n s_\theta(\theta; \phi; Z_i, A_i) \xrightarrow{\text{maximizer}} \hat{\theta}[\phi]$$

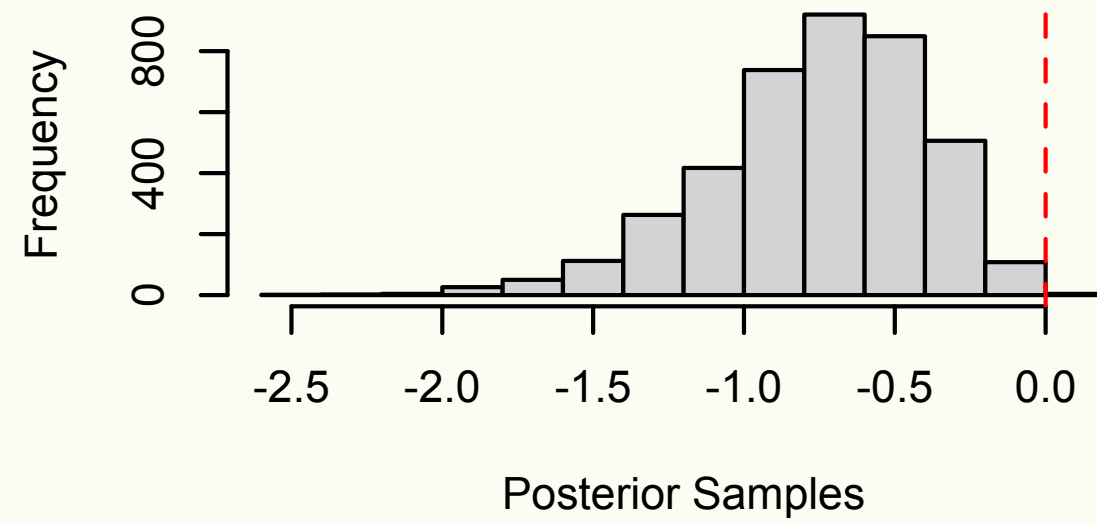
$$E_{\tilde{Z}, \tilde{A}} \left[ E_{\text{pos}, \phi} \left[ \nu(\hat{\theta}[\phi]; \phi^*; \tilde{Z}, \tilde{A}) \mid \underline{\mathcal{S}} \right] \right] \xrightarrow{\text{minimizer}} \text{Good model (w.r.t. MSE)}$$

(Pseudo) Posterior Expectation of  $\phi$        $\{A_i, X_i, Y_i\}_{i=1}^n$  (dataset)

# Posterior Dist. of IPW $\hat{\theta}[\phi]$

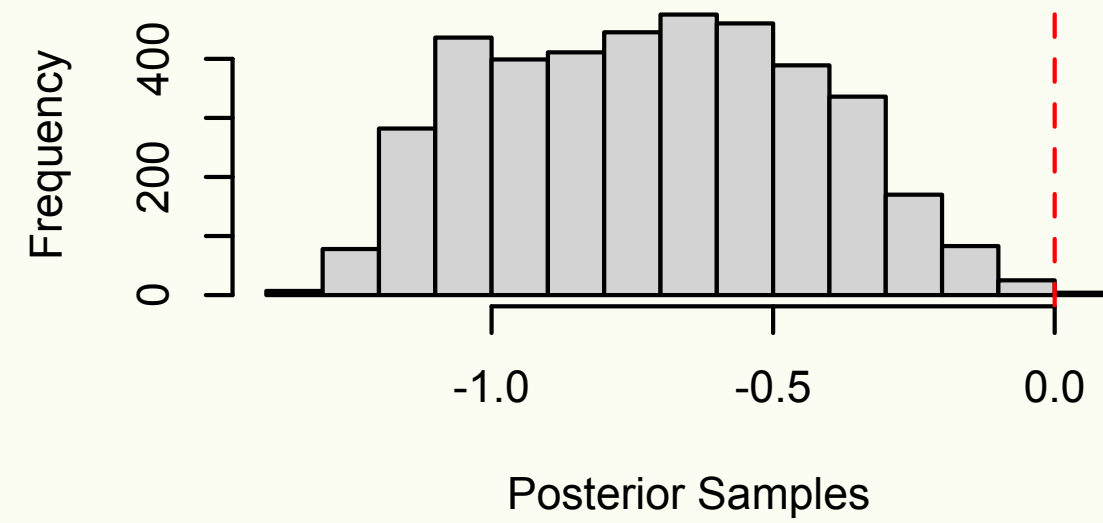
Heavily biased

$\tau = 0$

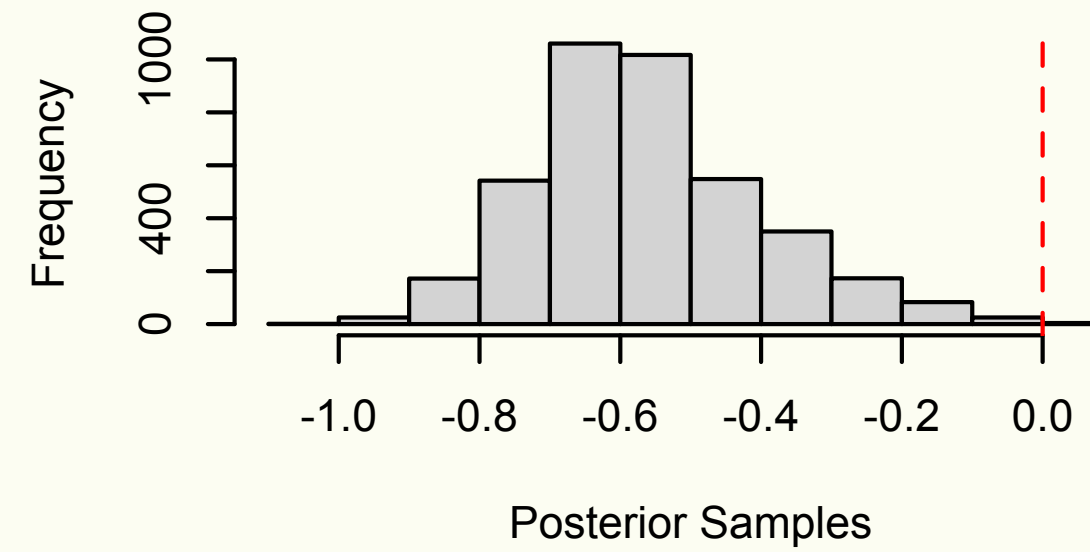


Non-normal distribution

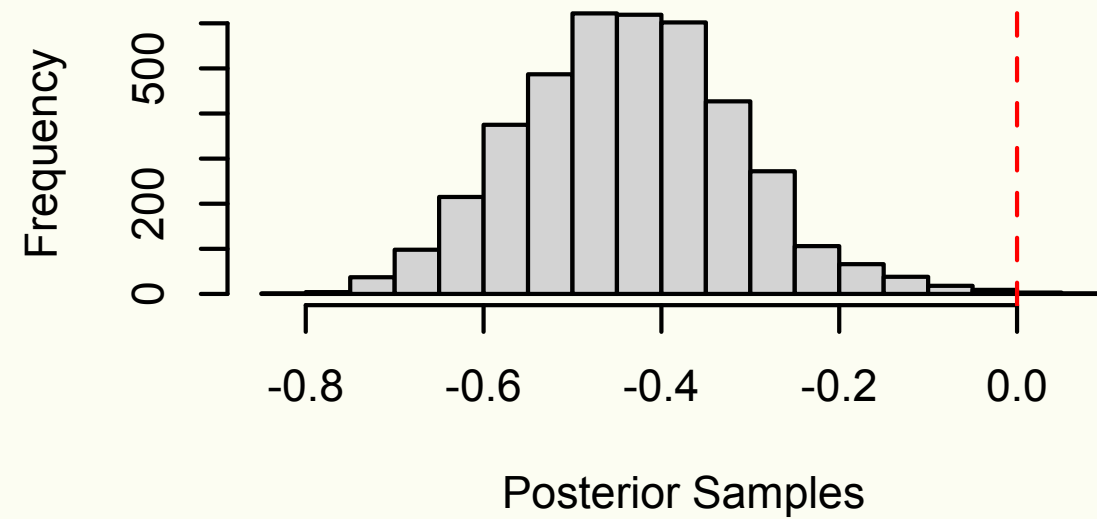
$\tau = 0.01$



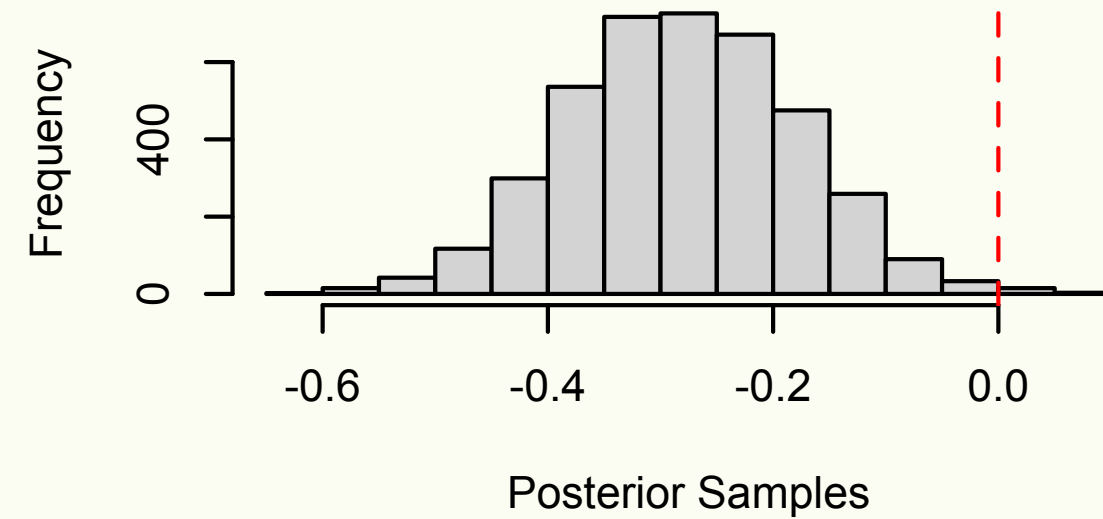
$\tau = 0.02$



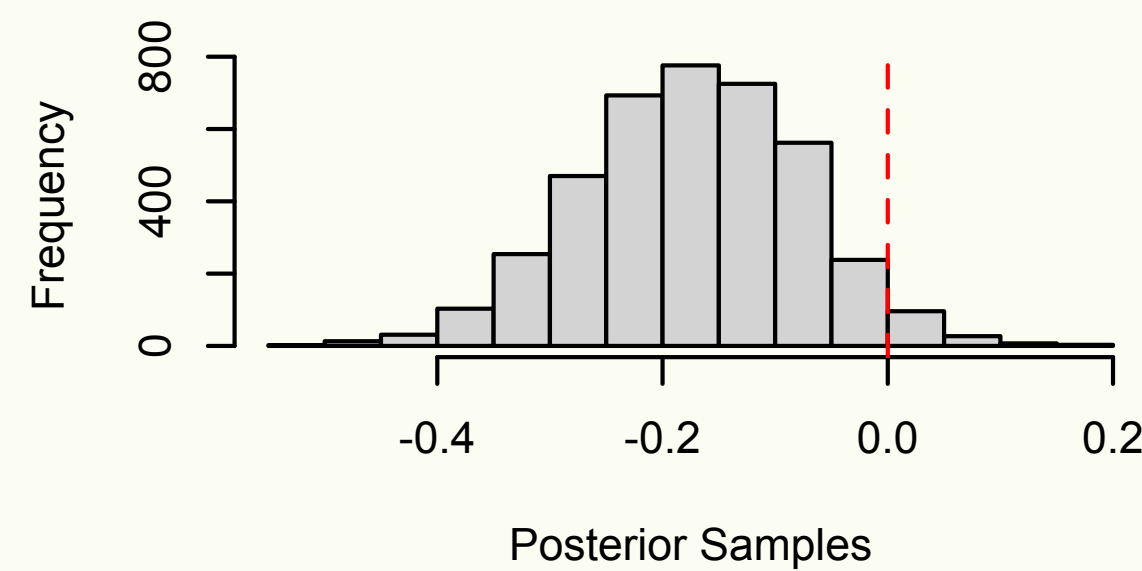
$\tau = 0.03$



$\tau = 0.05$



$\tau = 0.075$

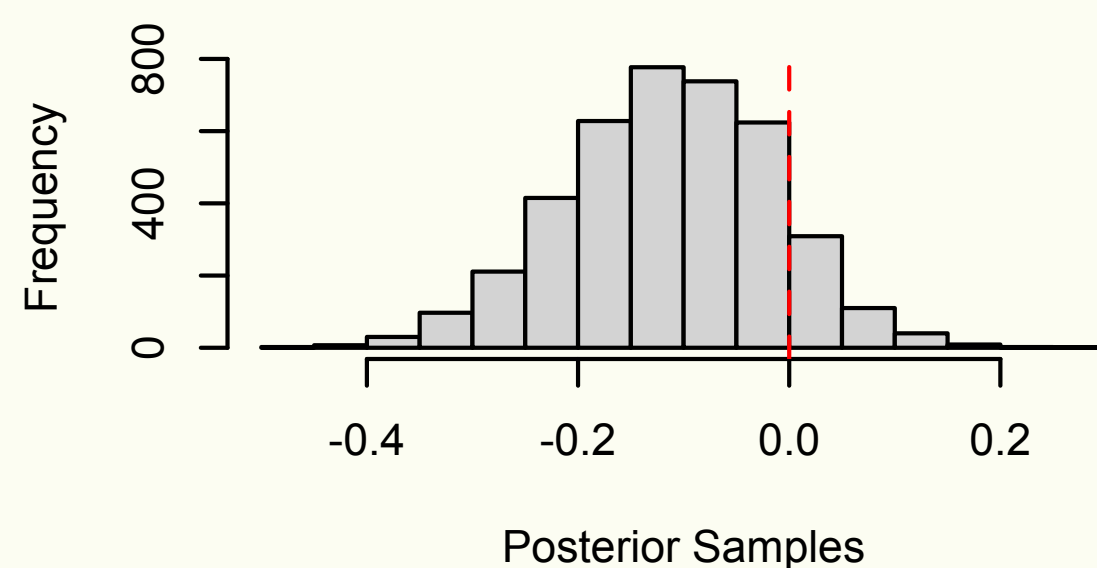


Outcome model

$$Y = X_1 + X_2 - X_3 - X_4 + \varepsilon$$

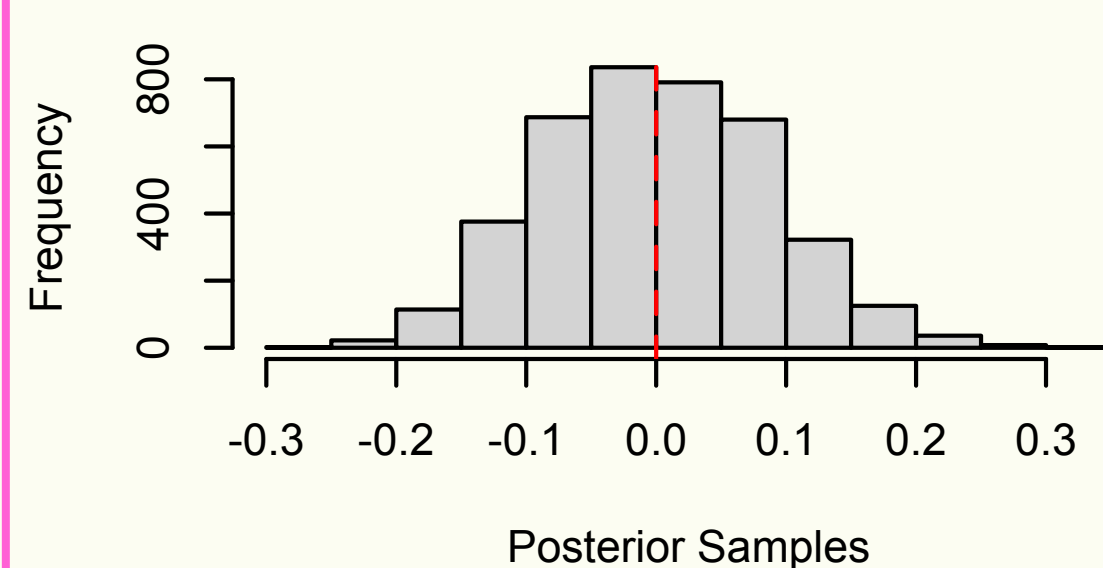
$$X_1, X_2, X_3, X_4, \varepsilon \sim N(0, 1)$$

$\tau = 0.10$

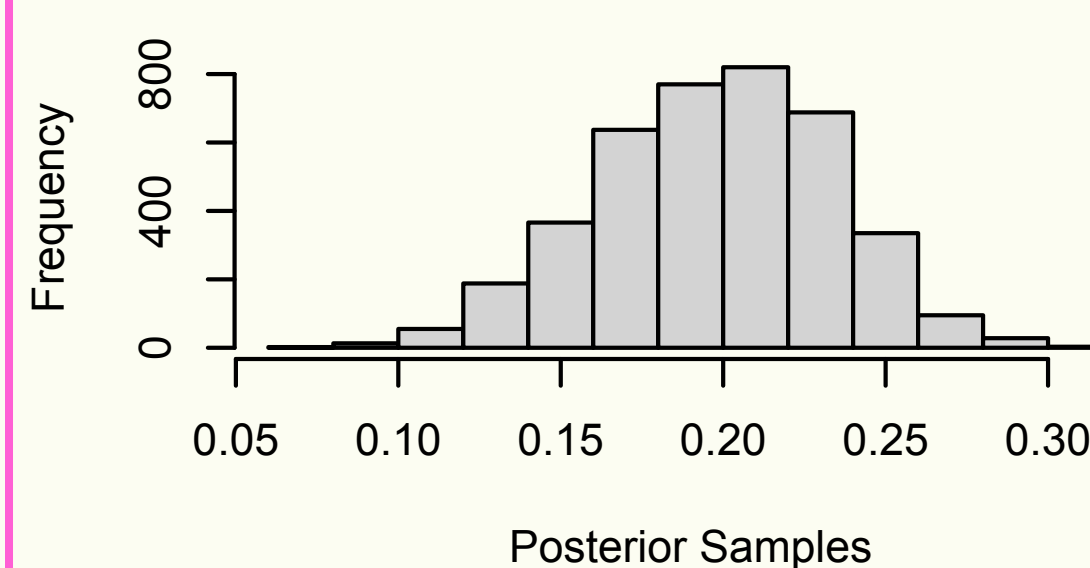


Can we choose this ideal threshold?

$\tau = 0.20$



$\tau = 0.50$



Propensity score

$$e(X) = P(A = 1 | X)$$

$$= \frac{1}{1 + \exp\{(1 + X_1 + X_2 + X_3 - X_4)/2\}}$$

# Estimating Functions & Evaluation Metric — Singular Setup —

- $s_\phi(\phi; Z_i, A_i)$ : Scoring function for  $\phi$
- $s_\theta(\theta; \phi; Z_i, A_i)$ : Scoring function for  $\theta$

- e.g., IPW estimator  $\frac{A_i}{e_\tau(Z_i; \phi)} (\theta - Y_i)^2$ 
  - ②  $e_\tau(Z_i; \phi)$  Trimmed propensity

- $\nu(\hat{\theta}[\phi]; \phi^*; Z_i, A_i)$ : Evaluation metric for  $\hat{\theta}[\phi]$
- e.g.,  $\frac{A_i}{e(Z_i; \phi^*)} (\hat{\theta}[\phi] - Y_i)^2$

① **Bayesian Approach**  $\hat{\phi} = \arg \max_{\phi} \exp \left( \sum_{i=1}^n s_\phi(A_i, Z_i) \right) \pi(\phi)$

Pseudo likelihood      Prior

$$\phi \sim \exp \left( \sum_{i=1}^n s_\phi(\phi; A_i, Z_i) \right) \pi(\phi)$$

$$\sum_{i=1}^n s_\theta(\theta; \phi; Z_i, A_i) \xrightarrow{\text{maximizer}} \hat{\theta}[\phi]$$

$$E_{\tilde{Z}, \tilde{A}} \left[ E_{\text{pos}, \phi} \left[ \nu(\hat{\theta}[\phi]; \phi^*; \tilde{Z}, \tilde{A}) \mid \mathcal{S} \right] \right] \xrightarrow{\text{minimizer}} \text{Good model (w.r.t. MSE)}$$

(Pseudo) Posterior Expectation of  $\phi$        $\{A_i, X_i, Y_i\}_{i=1}^n$  (dataset)

# Problem: estimation of the Generalization Error (GE)

$$E_{\tilde{Z}, \tilde{A}} \left[ \underbrace{E_{\text{pos}, \phi}}_{\text{Pseudo posterior}} \left[ \nu(\hat{\theta}[\phi]; \phi^*; \tilde{Z}, \tilde{A}) \mid \mathcal{S} \right] \right] \approx \frac{1}{n} \sum_{i=1}^n E_{\text{pos}, \phi} \left[ \nu(\hat{\theta}[\phi]; \phi^*; Z_i, A_i) \mid \mathcal{S} \right]$$

Pseudo posterior

Dual use of the same dataset in ESTIMATION and EVALUATION  $\Rightarrow$  Bias

# Leave-One-Out-Cross-Validation Criterion

$$\text{CV} := \frac{1}{n} \sum_{i=1}^n E_{\text{pos}, \phi^{-i}} \left[ \nu(\hat{\theta}[\phi^{-i}]; \phi^*; Z_i, A_i) \mid \underline{\mathcal{S}^{-i}} \right]$$

Dataset excluding  $i$ -th data

- We start from the Cross-validation type criterion
  - 😊 r.v. between "ESTIMATION" and "EVALUATION" is independent
  - 😞 Computationally very intensive (posterior sampling  $\times n$ )

⇒ We derive an information criterion that is asymptotically unbiased as **CV**, by using only "one-set posterior samples".

# List of Relevalent Information Criteria (IC)

IC	Target	Evaluation Metric	Propensity Score (PS)	PS Estimation	Singularity
<b>AIC</b> (Akaike, 1973)	Density function	KL divergence		NA	
<b>WAIC</b> (Watanabe, 2010, JMLR)	Density function	KL divergence		NA	✓
<b>wCp</b> (Baba et al., 2017, Biometrika)	Regression	MSE	✓	△	
<b>PCIC</b> (Iba and Yano, 2023, Neural Comp.)	General	General	✓		✓
<b>Proposed</b>	General	General	✓	✓	✓

# Proposed Information Criterion

# Proposed Information Criterion With Known Outer Propensity

## Theorem 1 (Morikawa and Yano, 2024+)

Under some regularity conditions,

$$E(\text{IC}) = E(\text{CV}) + o(n^{-1})$$

$$\text{IC} = \frac{1}{n} \sum_{i=1}^n E_{\text{pos},\phi} \left[ \nu(\hat{\theta}[\phi]; \phi^*; \mathbf{Z}_i, A_i) \mid \mathcal{S} \right] + \underbrace{I_1 + I_2}_{=\text{Bias corrected term}},$$

$$I_1 = -\frac{1}{n} \sum_{i=1}^n \text{Cov}_{\text{pos},\phi} \left[ \nu(\hat{\theta}[\phi]; \phi^*; \mathbf{Z}_i, A_i), s_{\phi}(\hat{\theta}[\phi]; \phi; \mathbf{Z}_i, A_i) \mid \mathcal{S} \right],$$

$$I_2 = -\frac{1}{n} \sum_{i=1}^n E_{\text{pos},\phi} \left[ \nabla_{\theta} \nu(\hat{\theta}[\phi]; \phi^*; \mathbf{Z}_i, A_i) \left\{ -\sum_{j=1}^n \nabla_{\theta}^2 s_{\theta}(\hat{\theta}[\phi]; \phi; \mathbf{Z}_j, A_j) \right\}^{-1} \nabla_{\theta} s_{\theta}(\hat{\theta}[\phi]; \phi; \mathbf{Z}_i, A_i) \mid \mathcal{S} \right].$$

# Effects of Estimation of Outer Propensity Score

$$\nu(\hat{\theta}[\phi]; \psi; \tilde{Z}, \tilde{A}) = \frac{\tilde{A}}{e(\psi; \tilde{Z})} \left( \tilde{Y} - \frac{\sum_{i=1}^n A_i Y_i / \underbrace{e_{\tau}(\phi; Z_i)}_{\text{Inner propensity}}}{\underbrace{\sum_{i=1}^n A_i / e_{\tau}(\phi; Z_i)}_{=\hat{\theta}[\phi]}} \right)^2$$

Outer propensity
Inner propensity

- We do NOT allow trimming for outer propensity in this study
  - Previous researches have NOT cared about effects of  $\hat{\psi}$  (e.g. MLE)
- In this talk, we also derive the asymptotic bias caused by replacement of  $\phi^*$  with  $\hat{\psi}$



# Key Idea: Data Splitting

$$\frac{1}{n} \sum_{i=1}^n E_{\text{pos}, \phi} \left[ \nabla_{\psi} \nu(\hat{\theta}[\phi]; \psi^*; \boxed{Z_i, A_i}) (\boxed{\hat{\psi}} - \psi^*) \mid \mathcal{S} \right]$$

Dual use of the same dataset  $\Rightarrow$  Data splitting

- Data split:  $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$ 
  - $\mathcal{S}_1$ : Computation of the evaluation metric
  - $\mathcal{S}_2$ : Estimation of  $\psi \rightarrow \hat{\psi}$

- Data-Splitting CV criterion

$$\text{CV}_1 := \frac{1}{n_1} \sum_{i \in \mathcal{S}_1} E_{\text{pos}, \phi^{-i}} \left[ \nu(\hat{\theta}[\phi^{-i}]; \hat{\psi}; Z_i, A_i) \mid \mathcal{S}_1^{-i} \right], \quad \text{where } n_1 = \#\{\mathcal{S}_1\}$$

- We can consider  $\text{CV}_2$  by changing the roles of  $\mathcal{S}_1$  and  $\mathcal{S}_2$
- Our final criterion is  $\text{CV}_1 + \text{CV}_2$

# Proposed Information Criterion With Unknown Outer Propensity

$$\text{IC}_1 = \frac{1}{n_1} \sum_{i \in \mathcal{S}_1} E_{\text{pos}, \phi} \left[ \nu(\hat{\theta}[\phi]; \hat{\psi}; Z_i, A_i) \mid \mathcal{S}_1 \right] + \underbrace{I_1 + I_2 + I_3 + I_4}_{=\text{Bias corrected term}},$$

$$I_1 = -\frac{1}{n_1} \sum_{i \in \mathcal{S}_1} \text{Cov}_{\text{pos}, \phi} \left[ \nu(\hat{\theta}[\phi]; \hat{\psi}; Z_i, A_i), s_\phi(\phi; Z_i, A_i) \mid \mathcal{S}_1 \right],$$

$$I_2 = -\frac{1}{n_1} \sum_{i \in \mathcal{S}_1} E_{\text{pos}, \phi} \left[ \nabla_{\theta} \nu(\hat{\theta}[\phi]; \hat{\psi}; Z_i, A_i) \left\{ -\sum_{j \in \mathcal{S}_1} \nabla_{\theta}^2 s_\theta(\hat{\theta}[\phi]; \phi; Z_j, A_j) \right\}^{-1} \nabla_{\theta} s_\theta(\hat{\theta}[\phi]; \phi; Z_i, A_i) \mid \mathcal{S}_1 \right],$$

$$I_3 = -\frac{1}{n_1} \sum_{i \in \mathcal{S}_1} E_{\text{pos}, \phi} \left[ \nabla_{\psi} \nu(\hat{\theta}[\phi]; \hat{\psi}; Z_i, A_i) \mid \mathcal{S}_1 \right] \underbrace{E(\hat{\psi} - \phi^*)}_{= O_p(n^{-1})},$$

$$I_4 = -\frac{1}{2n_1} \sum_{i \in \mathcal{S}_1} E_{\text{pos}, \phi} \left[ \nabla_{\psi}^2 \nu(\hat{\theta}[\phi]; \hat{\psi}; Z_i, A_i) \mid \mathcal{S}_1 \right] \underbrace{\text{Var}(\hat{\psi})}_{= O_p(n^{-1})},$$

Additional penalty due to the estimation of  $\hat{\psi}$

# Proposed Information Criterion With Unknown Outer Propensity

## Theorem 2 (Morikawa and Yano, 2024+)

Under some regularity conditions,

$$E(\tilde{\text{IC}}) = E(\tilde{\text{CV}}) + o(n^{-1}),$$

where  $\tilde{\text{IC}} = \text{IC}_1 + \text{IC}_2$ ,  $\tilde{\text{CV}} = \text{CV}_1 + \text{CV}_2$

- In practice, to reduce the randomness of data splitting, we can compute the proposed IC multiple times  $\tilde{\text{IC}}^{(b)}$  ( $b = 1, \dots, B$ )
- Then, we can take a vote among minimizer of each  $\tilde{\text{IC}}^{(b)}$  ( $b = 1, \dots, B$ )

# Numerical Study

# Numerical Study (Toy Example)

## 1. Outcome model

$$Y = X_1 + X_2 - X_3 - X_4 + \varepsilon,$$

$$X_1, X_2, X_3, X_4, \varepsilon \sim N(0, 1)$$

## 2. Propensity Score

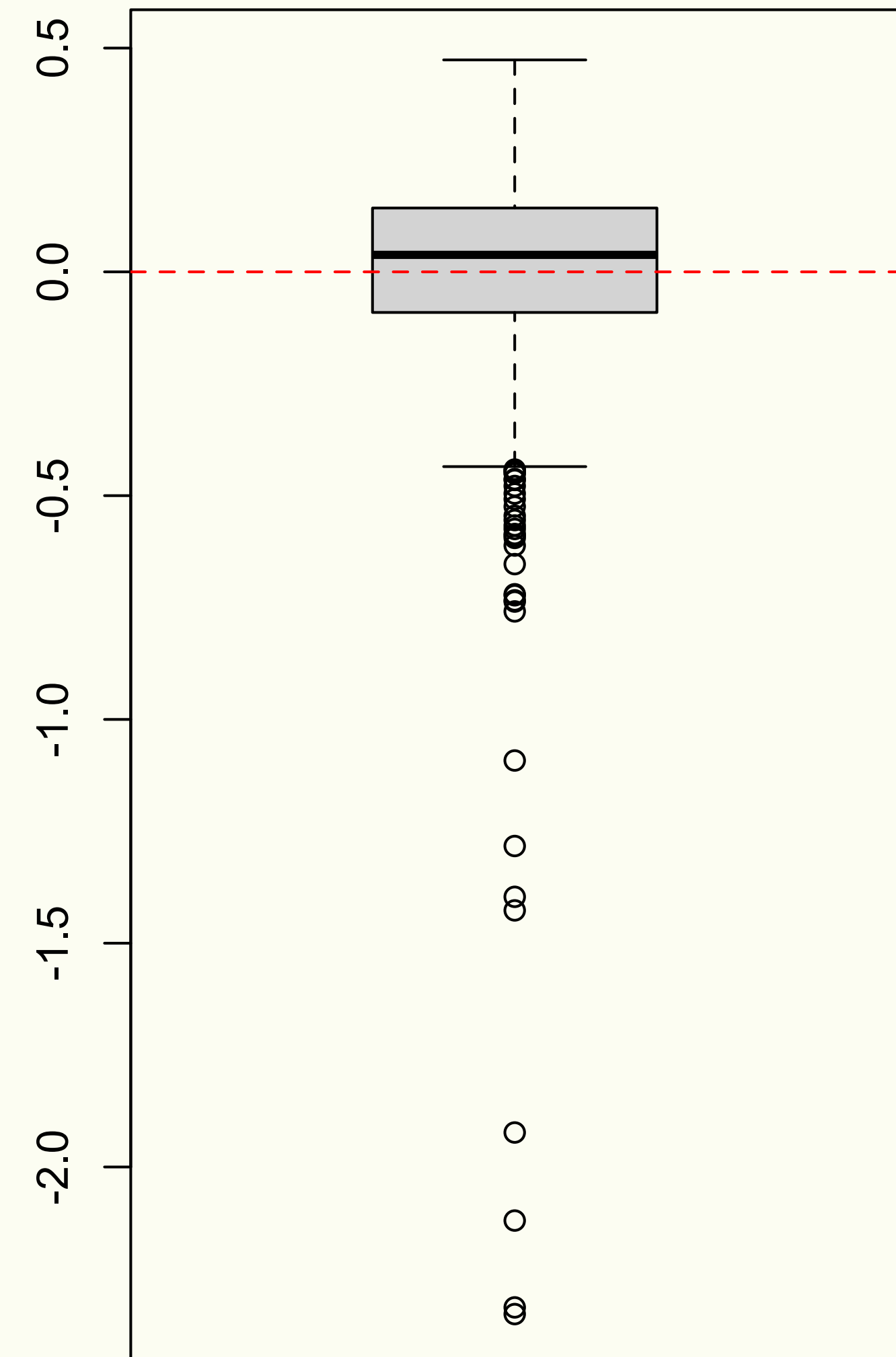
$$e(X) = P(A = 1 | X) = \frac{1}{1 + \exp\{(1 + X_1 + X_2 + X_3 - X_4)/2\}}$$

- Target :  $E(Y)$

- $s_\phi$ : log-likelihood;  $s_\theta$ : IPW,  $\nu(\cdot)$ : IPW  $L_2$ -metric

- $n = 400$ , #iteration = 100

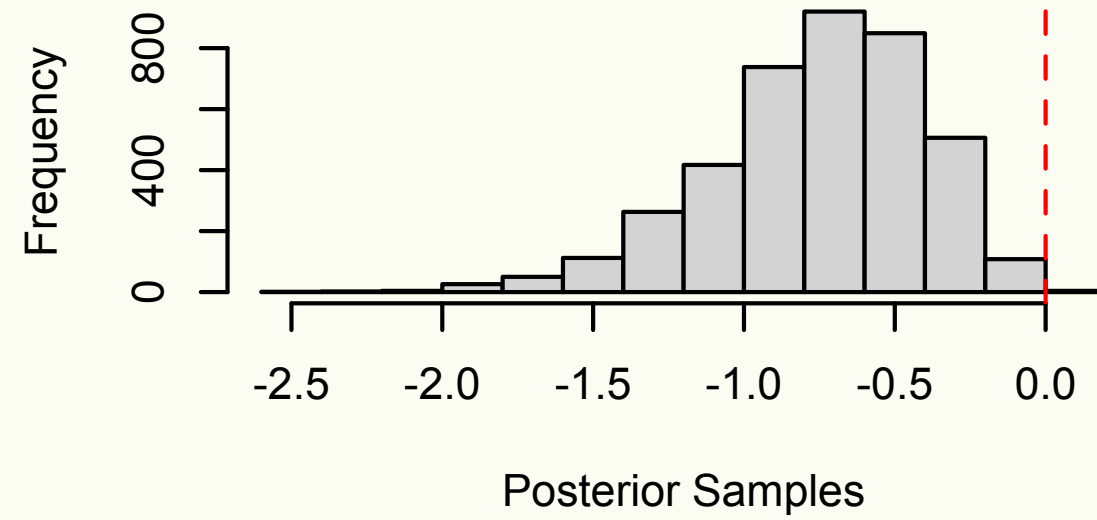
IPW Estimator



# Posterior Dist. of an IPW Estimator (cont'd)

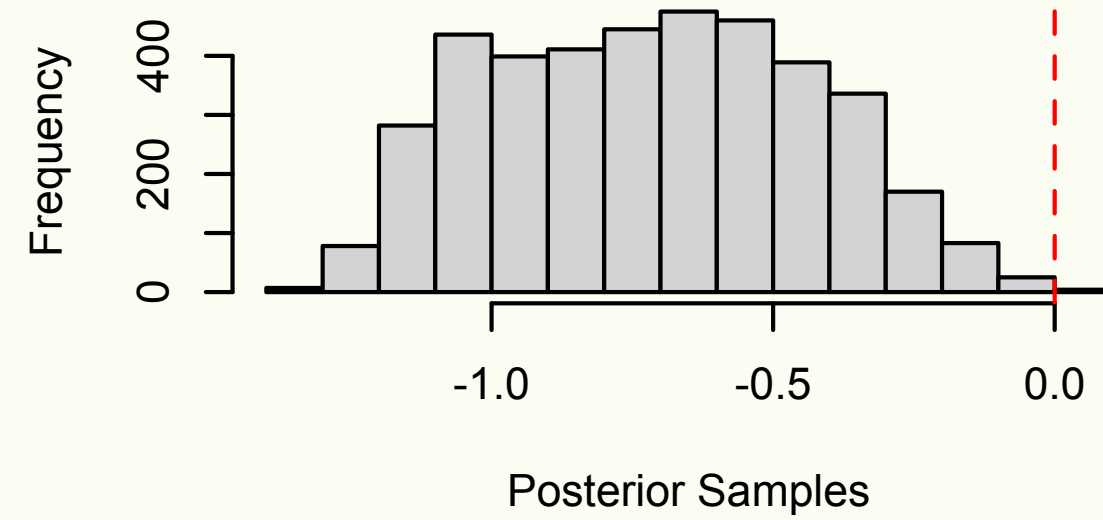
Heavily biased

$\tau = 0$

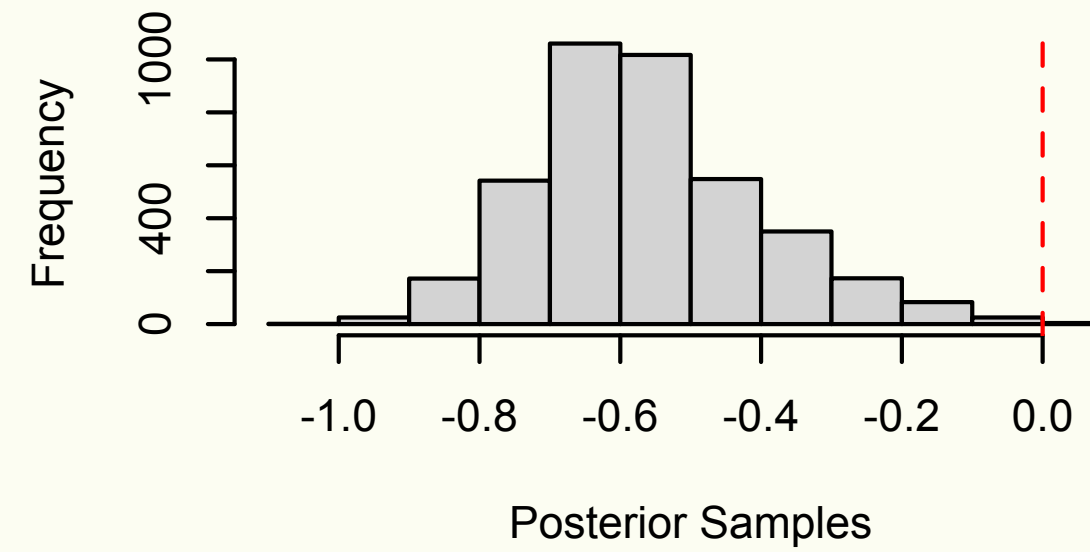


Non-normal distribution

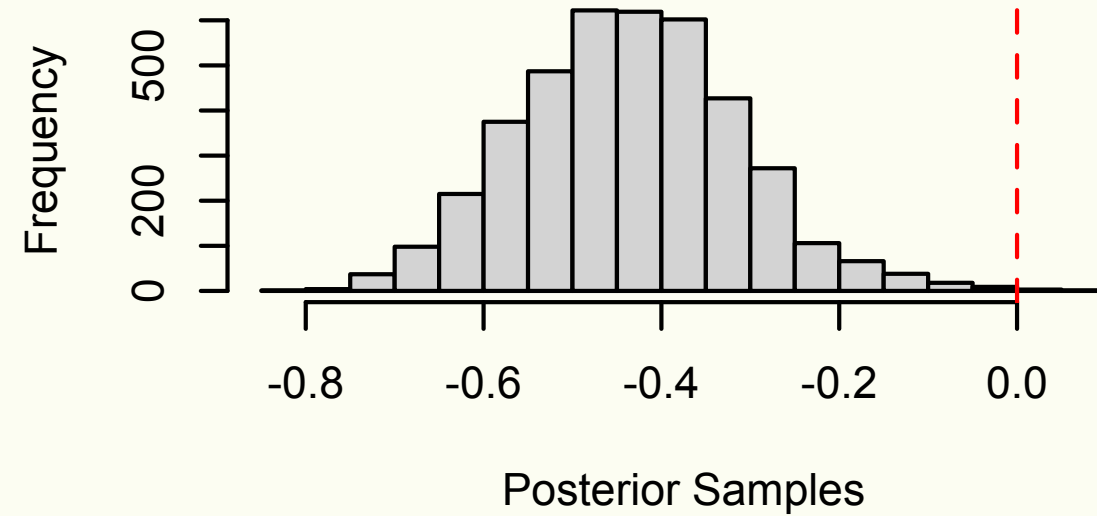
$\tau = 0.01$



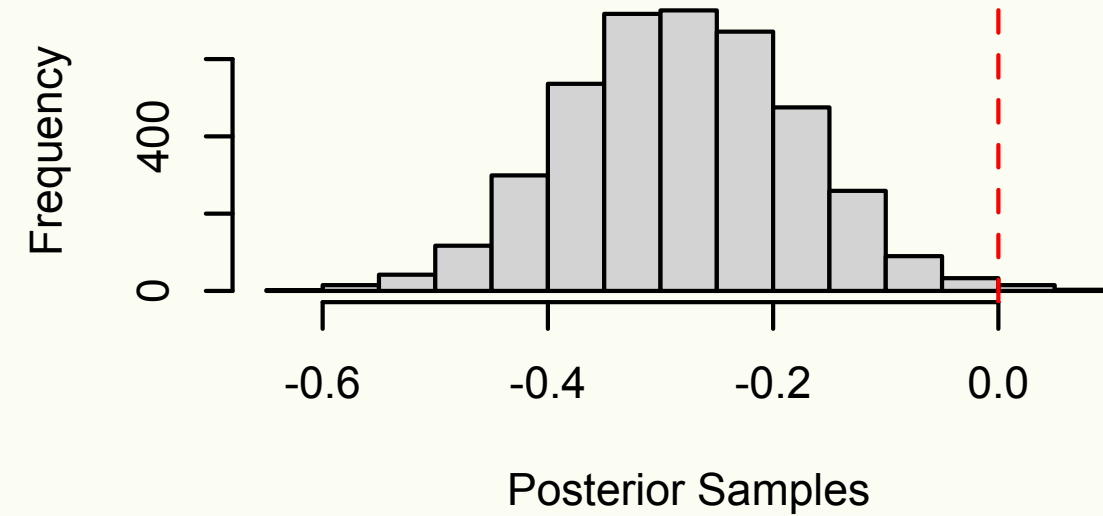
$\tau = 0.02$



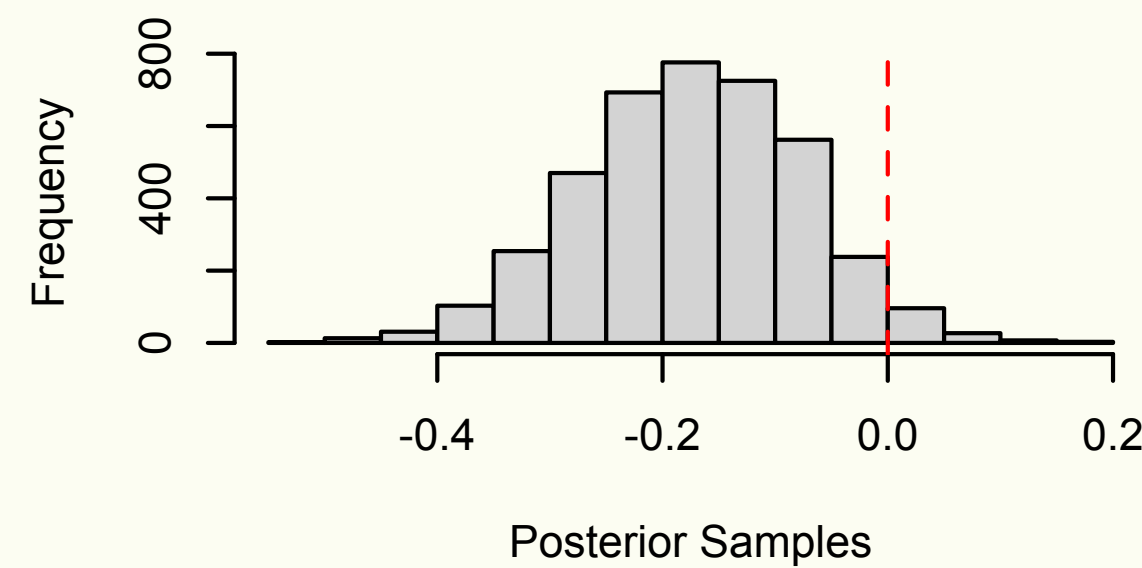
$\tau = 0.03$



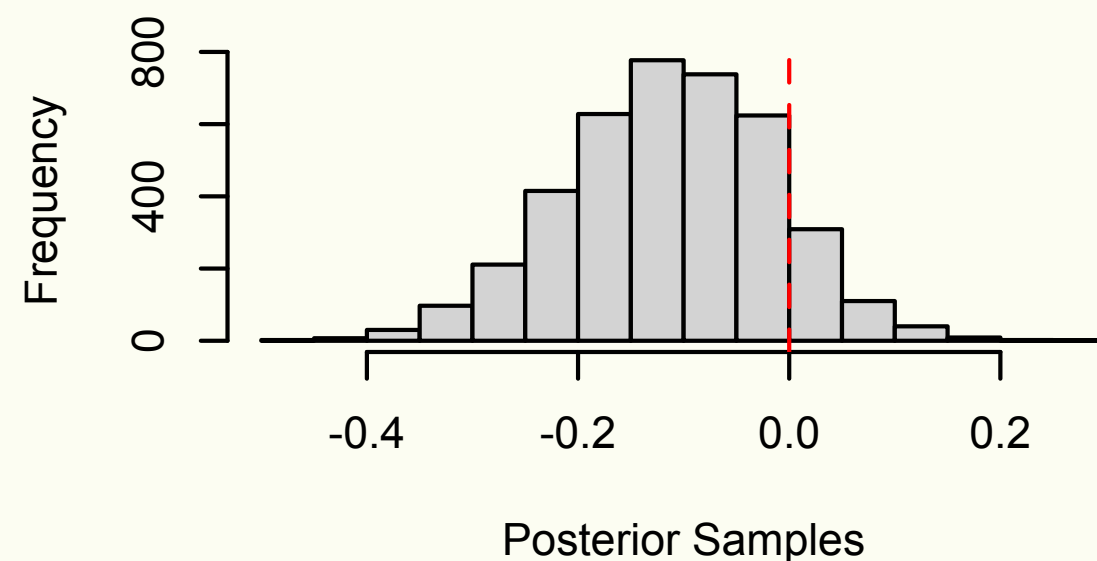
$\tau = 0.05$



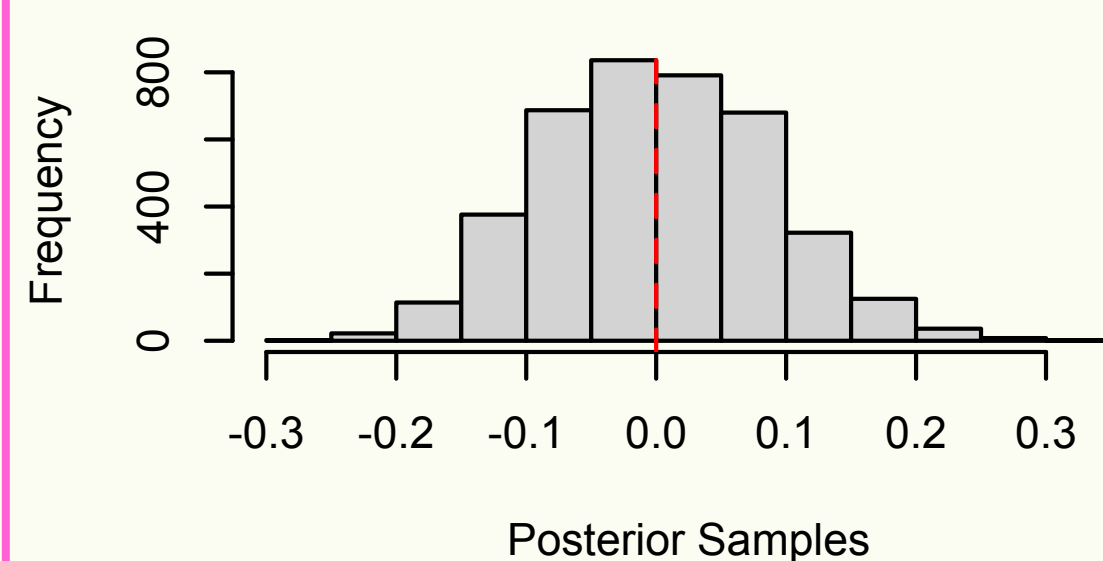
$\tau = 0.075$



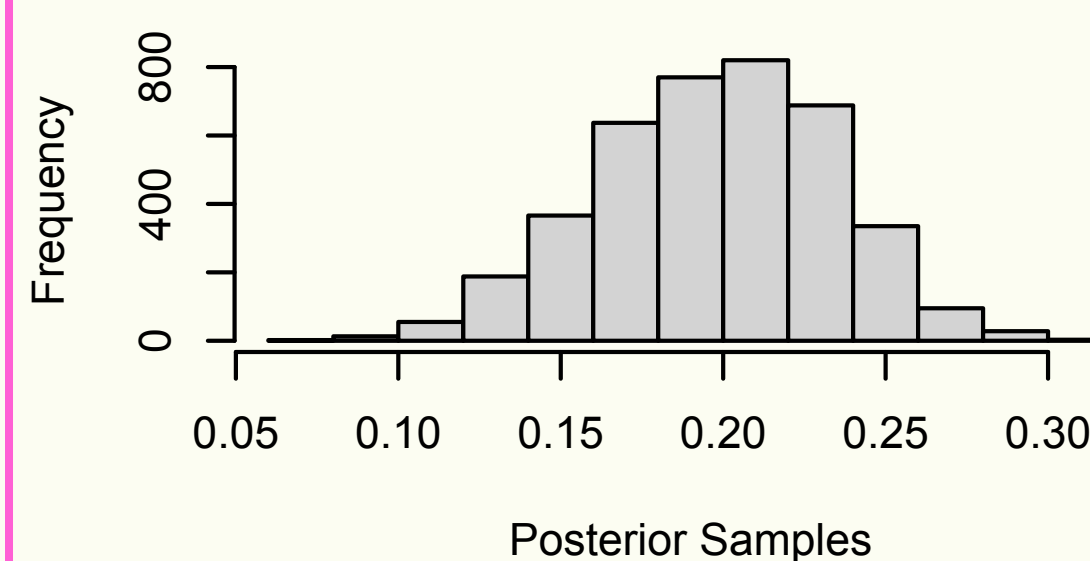
$\tau = 0.10$



$\tau = 0.20$



$\tau = 0.50$



Outcome model

$$Y = X_1 + X_2 - X_3 - X_4 + \varepsilon$$

$$X_1, X_2, X_3, X_4, \varepsilon \sim N(0, 1)$$

Propensity score

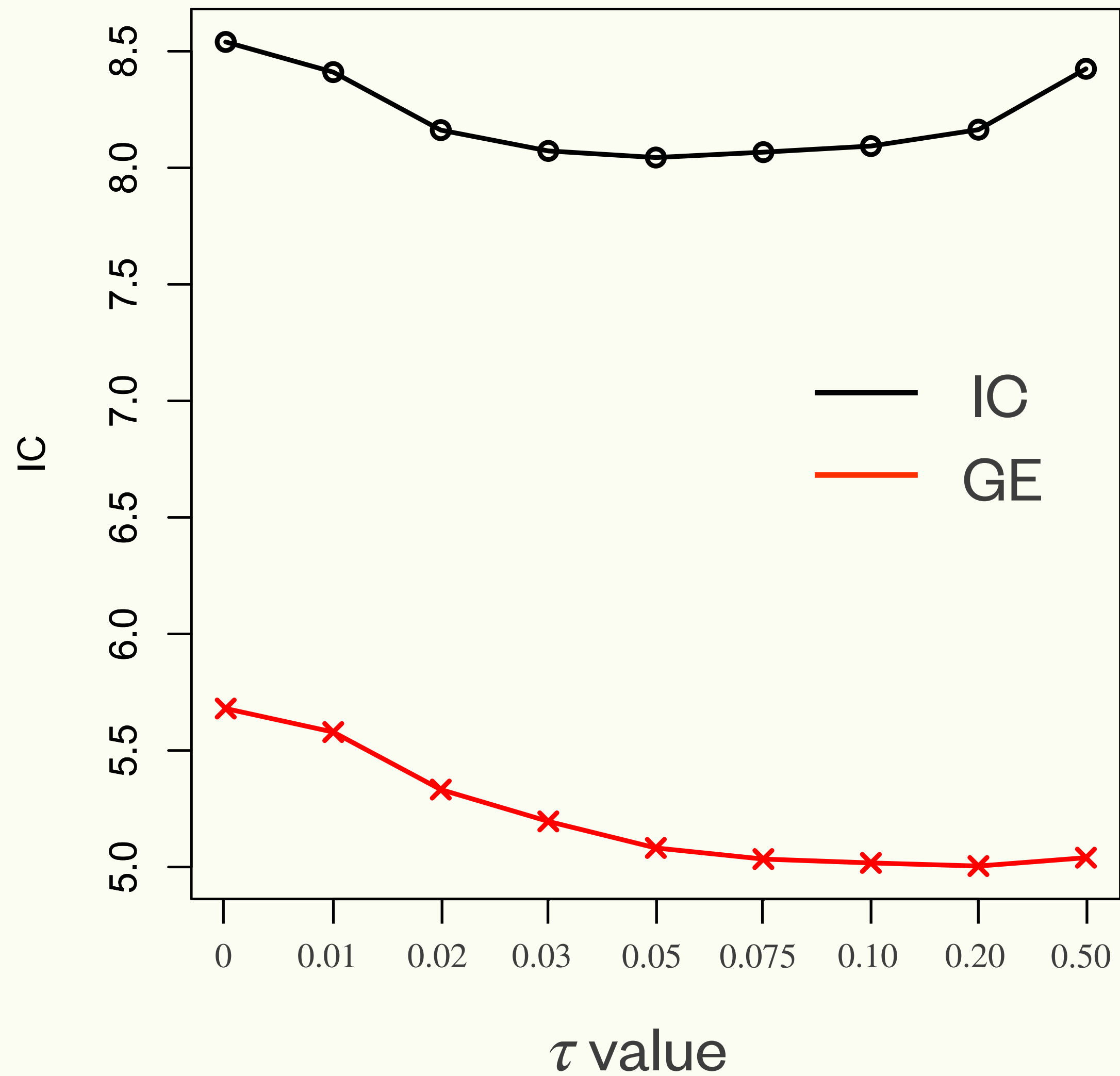
$$e(X) = P(A = 1 | X)$$

$$= \frac{1}{1 + \exp\{(1 + X_1 + X_2 + X_3 - X_4)/2\}}$$

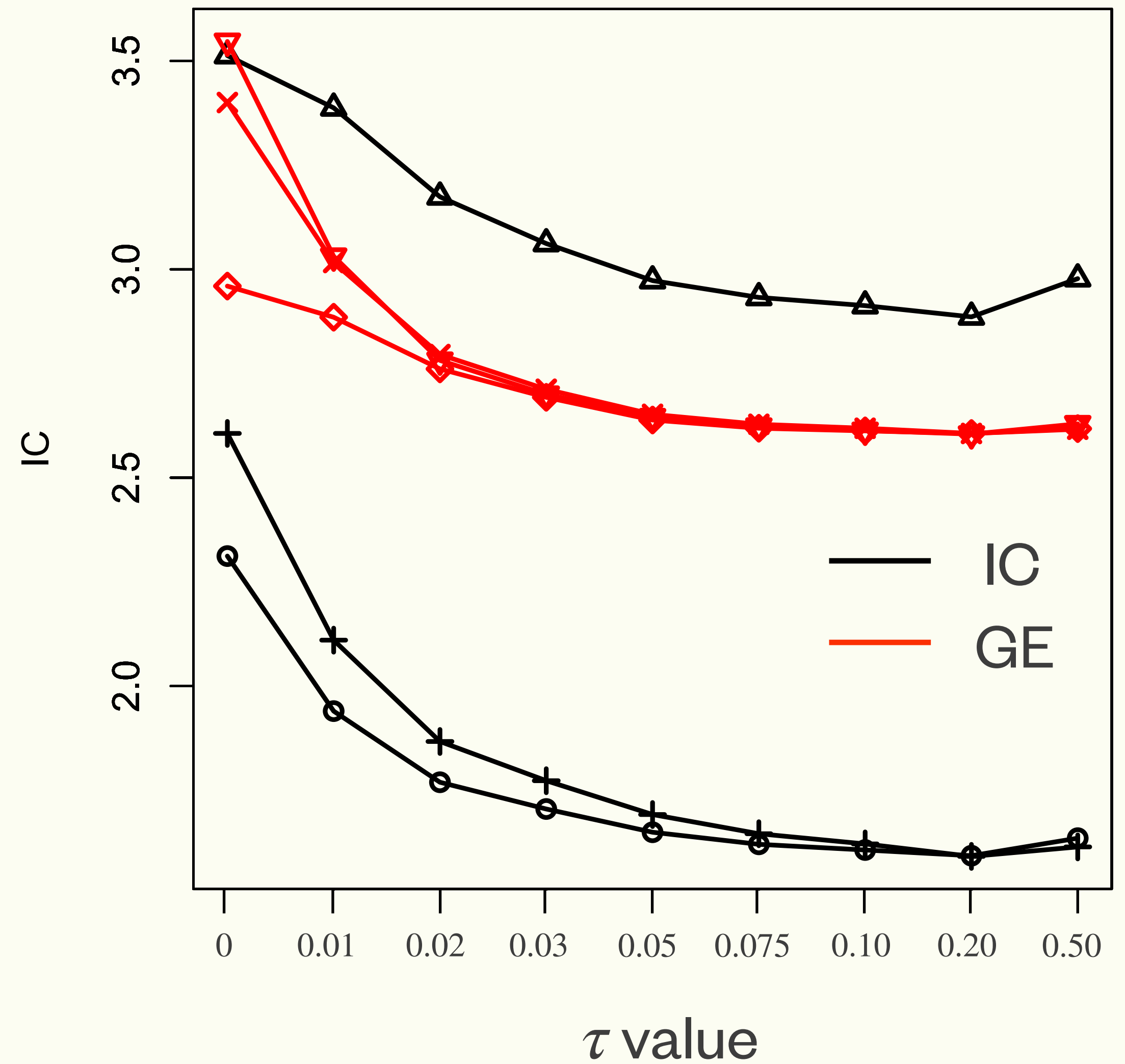
Can we choose this ideal threshold?

# An Application of Proposed IC

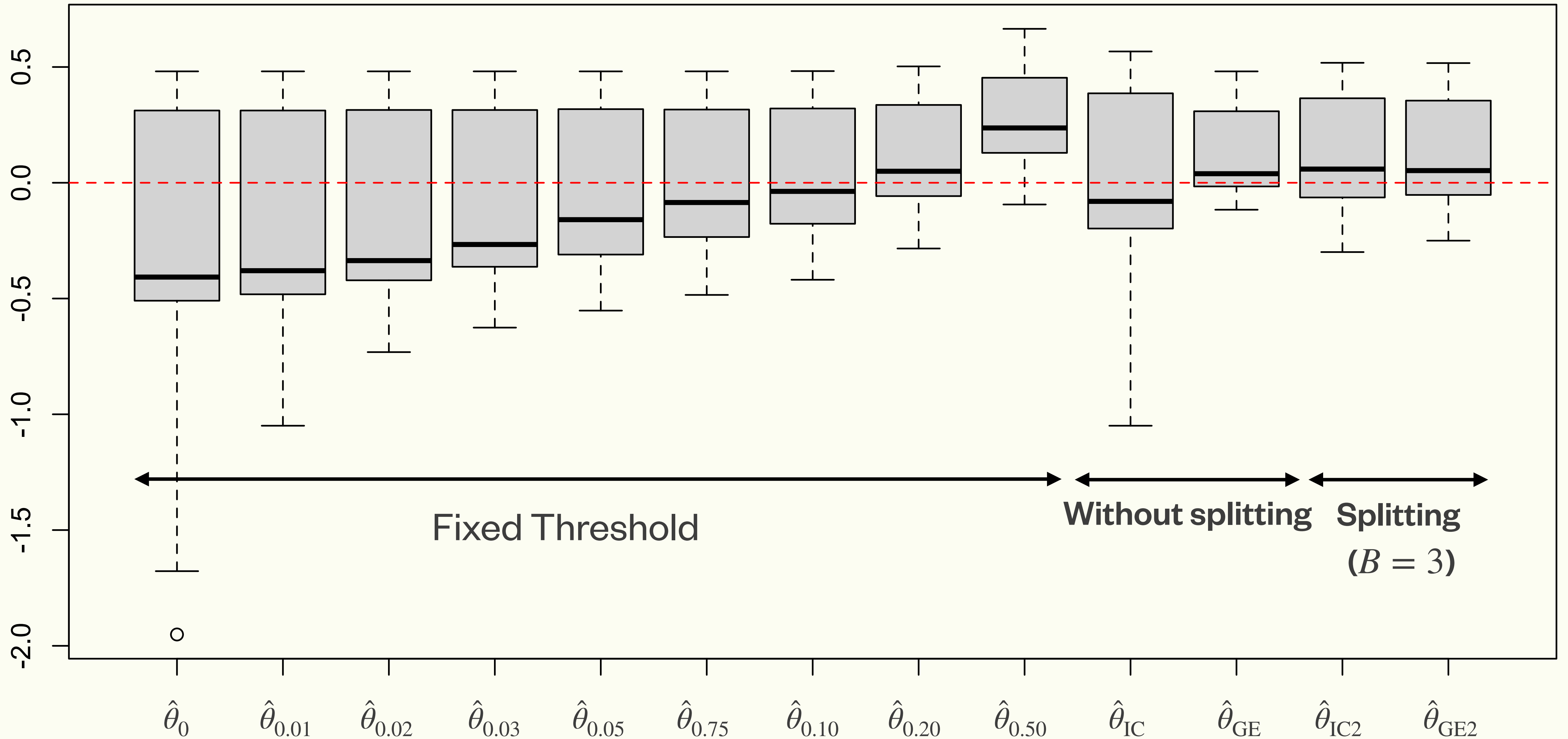
## Without Splitting



## Random Splittings ( $B = 3$ )



# Boxplots of the Estimators



**Conclusion**

# Discussion

- We proposed a statistical framework to use the **trimmed propensity score**

$$\nu(\hat{\theta}[\phi]; \psi; \tilde{Z}, \tilde{A}) = \frac{\tilde{A}}{\underbrace{e(\psi; \tilde{Z})}_{\text{Outer propensity}}} \left( \tilde{Y} - \underbrace{\frac{\sum_{i=1}^n A_i Y_i / \underbrace{e_{\tau}(\phi; Z_i)}_{\text{Inner propensity}}}{\sum_{i=1}^n A_i / \underbrace{e_{\tau}(\phi; Z_i)}_{\text{Inner propensity}}}}_{=\hat{\theta}[\phi]} \right)^2$$

- Crump et al. (2009, Biometrika) and Yang and Ding (2018, Biometrika) considered a statistical inference for trimmed estimand  $E\{Y \cdot I(\alpha_1 \leq e(X) \leq \alpha_2)\}$ 
  - $\alpha_1$  and  $\alpha_2$  are user-specified constants
  - Our estimand is  $E(Y)$
- However, when the propensity is extremely close to zero or one, we also need to consider the trimmed estimand, because the outer propensity also diverges

# Appendix

# Selected Models: #iteration = 100

- Without splitting

- IC:** Proposed Information Criterion

$\varepsilon$	0	0.01	0.02	0.03	0.05	0.075	0.10	0.20	0.50
	1	3	1	1	8	10	35	27	14

- GE:** Generalization error with the use of the true propensity score for the outer propensity

$\varepsilon$	0	0.01	0.02	0.03	0.05	0.075	0.10	0.20	0.50
	15	1	2	1	6	8	13	35	19

- #Random splitting = 3

- IC2:** Proposed Information Criterion

$\varepsilon$	0	0.01	0.02	0.03	0.05	0.075	0.10	0.20	0.50
	0	0	0	1	0	0	21	47	31

- GE2:** Generalization error with the use of the true propensity score for the outer propensity

$\varepsilon$	0	0.01	0.02	0.03	0.05	0.075	0.10	0.20	0.50
	1	3	1	3	8	9	14	40	21

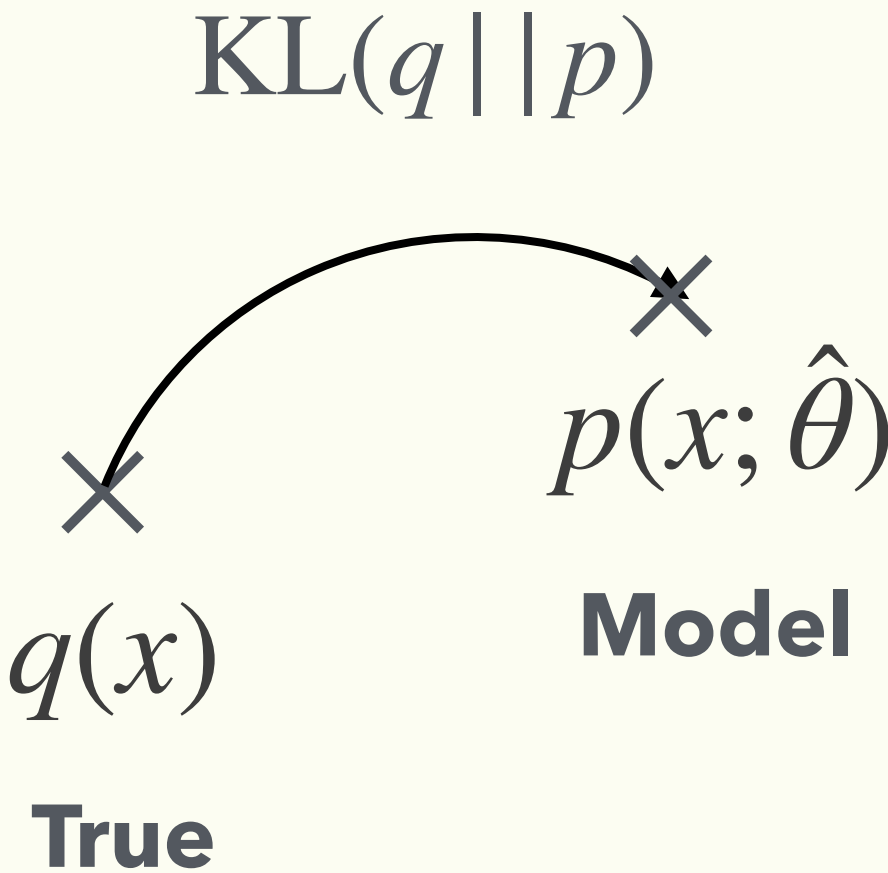
# Classical Information Criterion: AIC $\simeq$ CV-type Criterion

- Estimation of Kullback-Leibler (KL) divergence:

$$\text{KL}(q || p) = \int q(x) \log \frac{q(x)}{p(x; \hat{\theta})} dx \propto - \int q(x) \log p(x; \hat{\theta}) dx =: \tilde{\text{KL}}(q || p) \quad \text{Our Target!!}$$



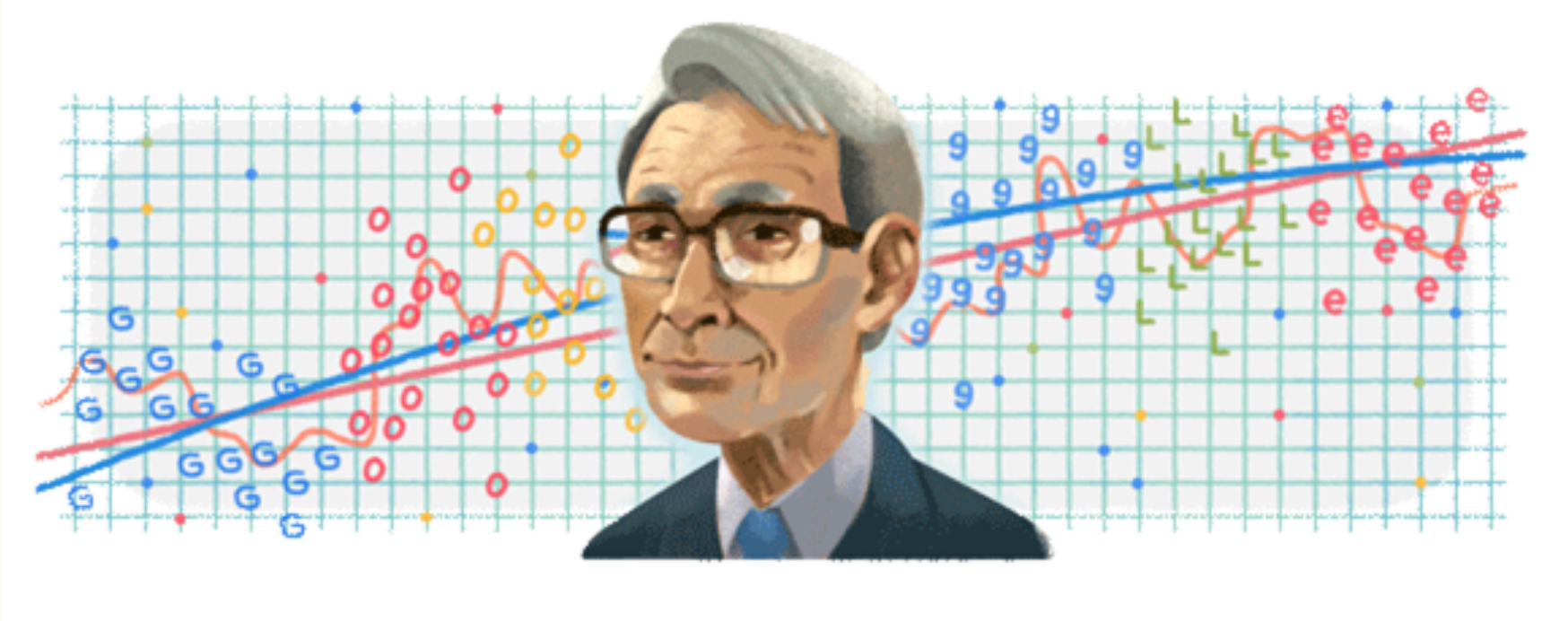
Prof. Akaike with a Google logo



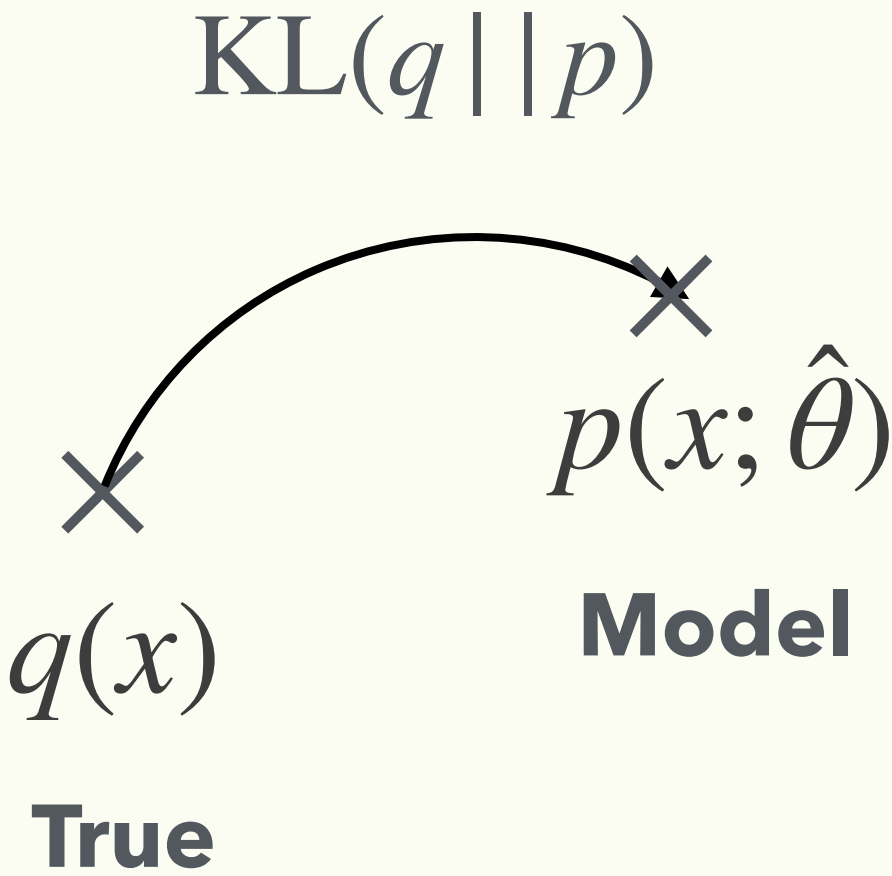
# Classical Information Criterion: AIC $\simeq$ CV-type Criterion

- Estimation of Kullback-Leibler (KL) divergence:

$$\begin{aligned}
 \text{KL}(q || p) &= \int q(x) \log \frac{q(x)}{p(x; \hat{\theta})} dx \propto - \int \boxed{q(x)} \log p(x; \hat{\theta}) dx =: \tilde{\text{KL}}(q || p) && \text{Our Target!!} \\
 &\quad \downarrow \text{Empirical Dist.} \\
 &\approx - \frac{1}{n} \sum_{i=1}^n \log p(X_i; \hat{\theta}(X_1, \dots, X_n))
 \end{aligned}$$



Prof. Akaike (赤池) in Google logo



# Classical Information Criterion: AIC $\simeq$ CV-type Criterion

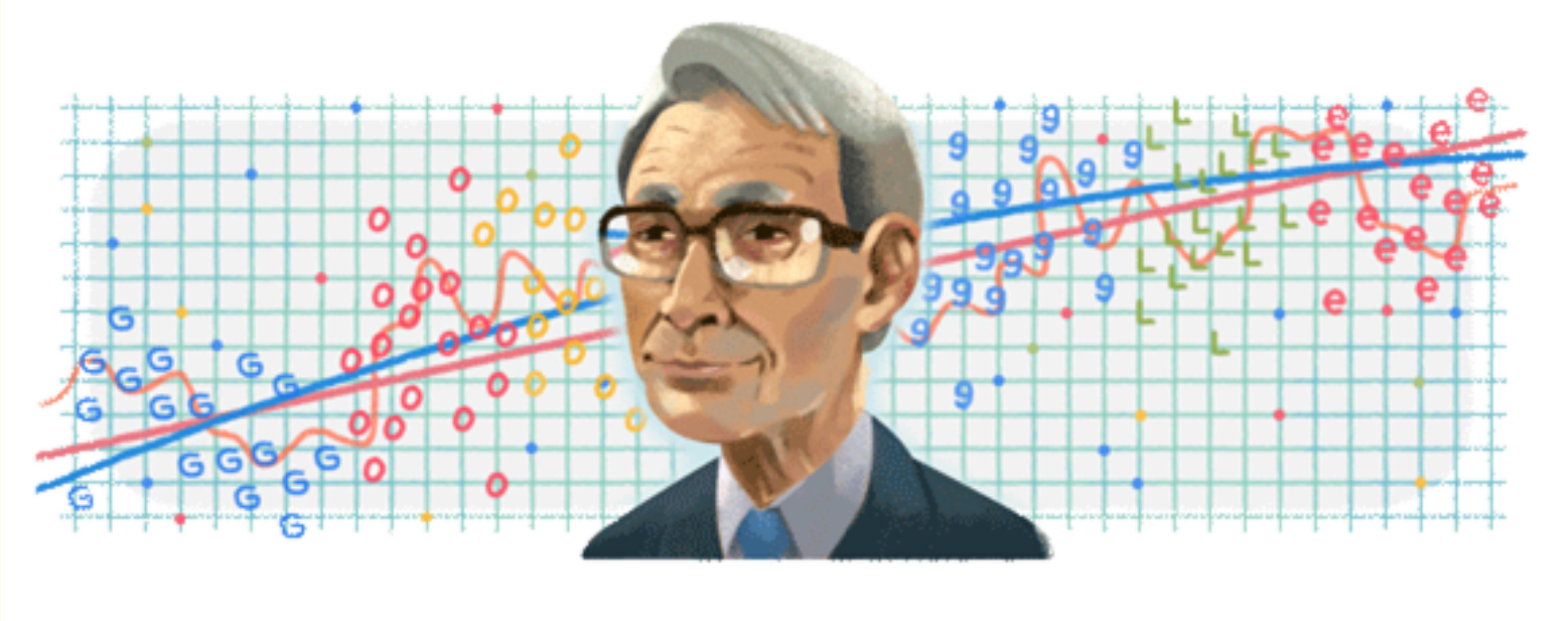
- Estimation of Kullback-Leibler (KL) divergence:

**Our Target!!**

$$KL(q || p) = \int q(x) \log \frac{q(x)}{p(x; \hat{\theta})} dx \propto - \int \boxed{q(x)} \log p(x; \hat{\theta}) dx =: \tilde{KL}(q || p)$$

↓ **Empirical Dist.**

$$\approx -\frac{1}{n} \sum_{i=1}^n \log p(X_i; \hat{\theta}(X_1, \dots, X_n))$$



Prof. Akaike (赤池) in Google logo

- **AIC** (Akaike, 1973)

$$-\frac{1}{n} \sum_{i=1}^n \log p(X_i; \hat{\theta}) + \frac{1}{n} \# \dim(\hat{\theta})$$

**= Asymptotic Bias**

- **CV** (Stone, 1974, JRSSB)

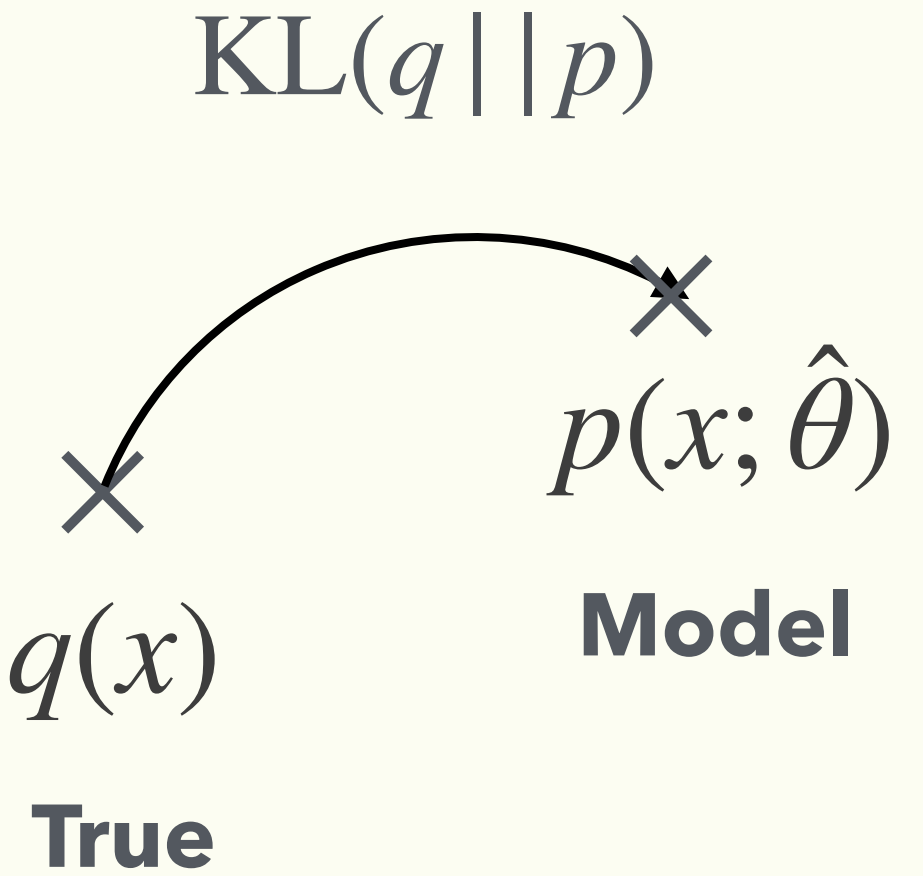
$$-\frac{1}{n} \sum_{i=1}^n \log p(X_i; \hat{\theta}^{(-i)})$$

**Estimator excluding  $i$ -th dataset**

**Remark.**

$$E[CV] = E[AIC] + o(n^{-1})$$

$$= E[\tilde{KL}(q || p)] + o(n^{-1})$$



# Pseudo Posterior & Its Bias Correction

- Pseudo Posterior of  $\phi$

$$\underbrace{\pi_{\text{pos},\phi}(\phi \mid \mathcal{S})}_{\text{Pseudo posterior}} = \frac{\underbrace{\exp\left\{\sum_{i=1}^n s_{\phi}(\phi; Z_i, A_i)\right\}}_{\text{Pseudo likelihood}} \underbrace{\pi(\phi)}_{\text{prior}}}{\int \exp\left\{\sum_{i=1}^n s_{\phi}(\bar{\phi}; Z_i, A_i)\right\} \pi(\bar{\phi}) d\bar{\phi}}$$

- If  $s_{\phi}(\phi; Z_i, A_i) = \log p(Z_i, A_i; \phi)$  (log-likelihood),  $\pi_{\text{pos},\phi}(\phi \mid \mathcal{S})$  is the usual posterior distribution
- Problem: estimation of the Generalization Error (GE)

$$E_{\tilde{Z}, \tilde{A}} \left[ E_{\text{pos},\phi} \left[ \nu(\hat{\theta}[\phi]; \phi^*; \tilde{Z}, \tilde{A}) \mid \mathcal{S} \right] \right] \approx \frac{1}{n} \sum_{i=1}^n E_{\text{pos},\phi} \left[ \nu(\hat{\theta}[\phi]; \phi^*; Z_i, A_i) \mid \mathcal{S} \right]$$

Dual use of the same dataset in ESTIMATION and EVALUATION  $\Rightarrow$  Bias