

What if the cure is worse than the disease?

# Policy learning for maximizing a primary outcome while controlling an adverse event

Laura Fuentes, M. Even, J. Josse, A. Chambaz

56es Journées de Statistique de la SFdS, Marseille 2025  
05/06/2025



# I. Context

# I.1-Medical motivations

Classical policy learning: Given patient's characteristics,  
determine the optimal treatment maximizing each patient's outcome

IVF example:

Find the optimal hormone dose to maximize the number of oocyte produced



# I.1-Medical motivations

Classical policy learning: Given patient's characteristics,  
determine the optimal treatment maximizing each patient's outcome

IVF example:

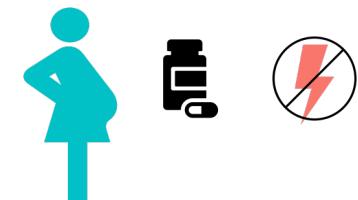
Find the optimal hormone dose to maximize the number of oocyte produced



Our goal: Given patient's characteristics,  
determine the optimal treatment maximizing each patient's outcome  
while controlling an adverse event

IVF example:

Find the optimal hormone dose to maximize the number of oocyte produced  
while avoiding ovarian hyperstimulation



# I.2-Mathematical framework

Set of independent and identically distributed subjects



Covariates:

$$X_i \in \mathcal{X}$$



Binary treatment:

$$A_i \in \{0,1\}$$



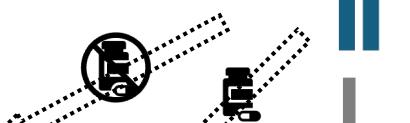
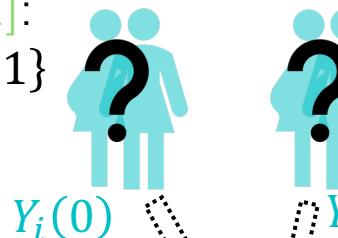
# I.2-Mathematical framework

Set of independent and identically distributed subjects



$$X_i \in \mathcal{X}$$

Potential outcomes [1]:  
 $Y_i(a) \in [0,1], a \in \{0,1\}$



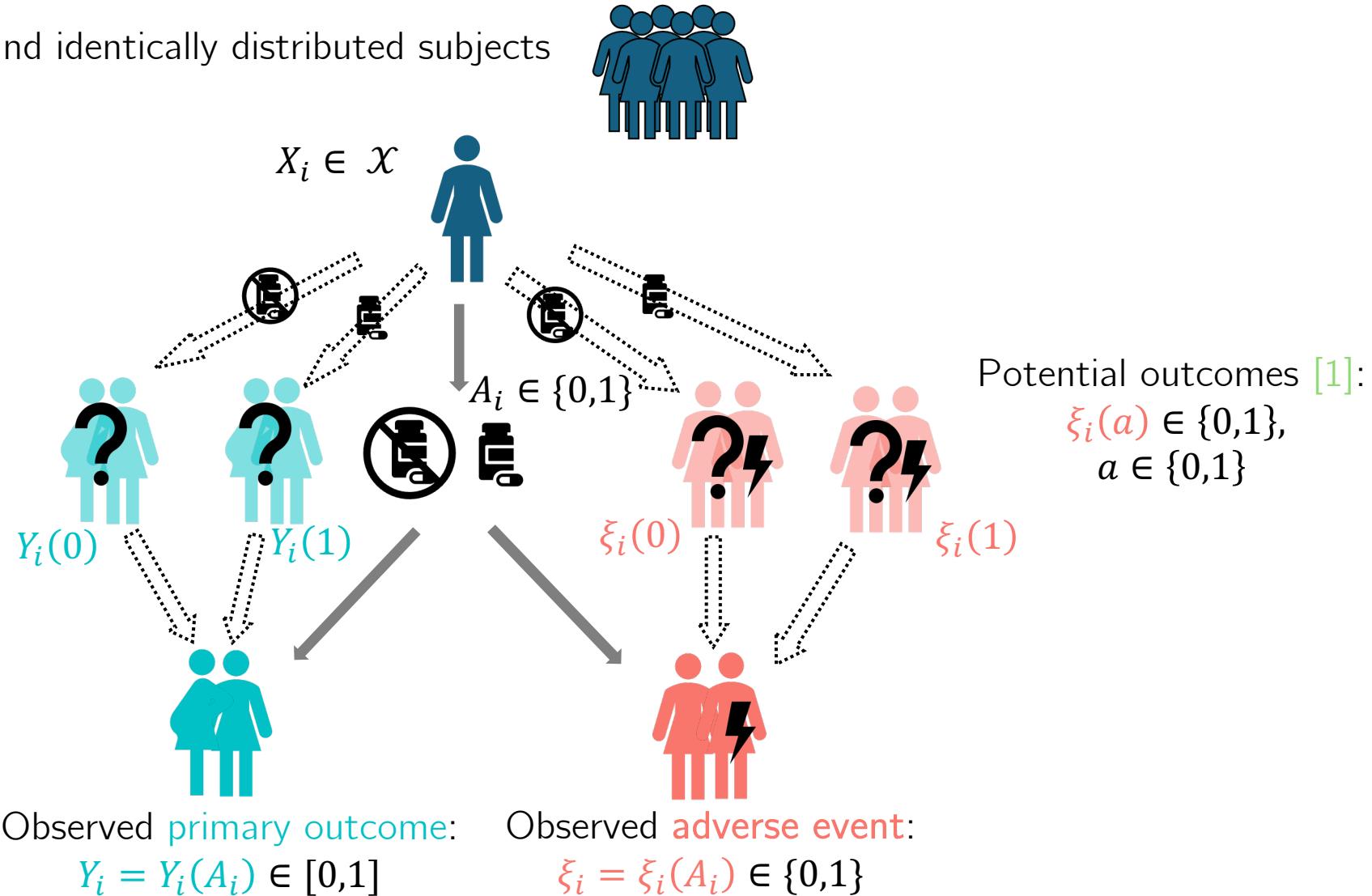
$$A_i \in \{0,1\}$$



Observed primary outcome:  
 $Y_i = Y_i(A_i) \in [0,1]$

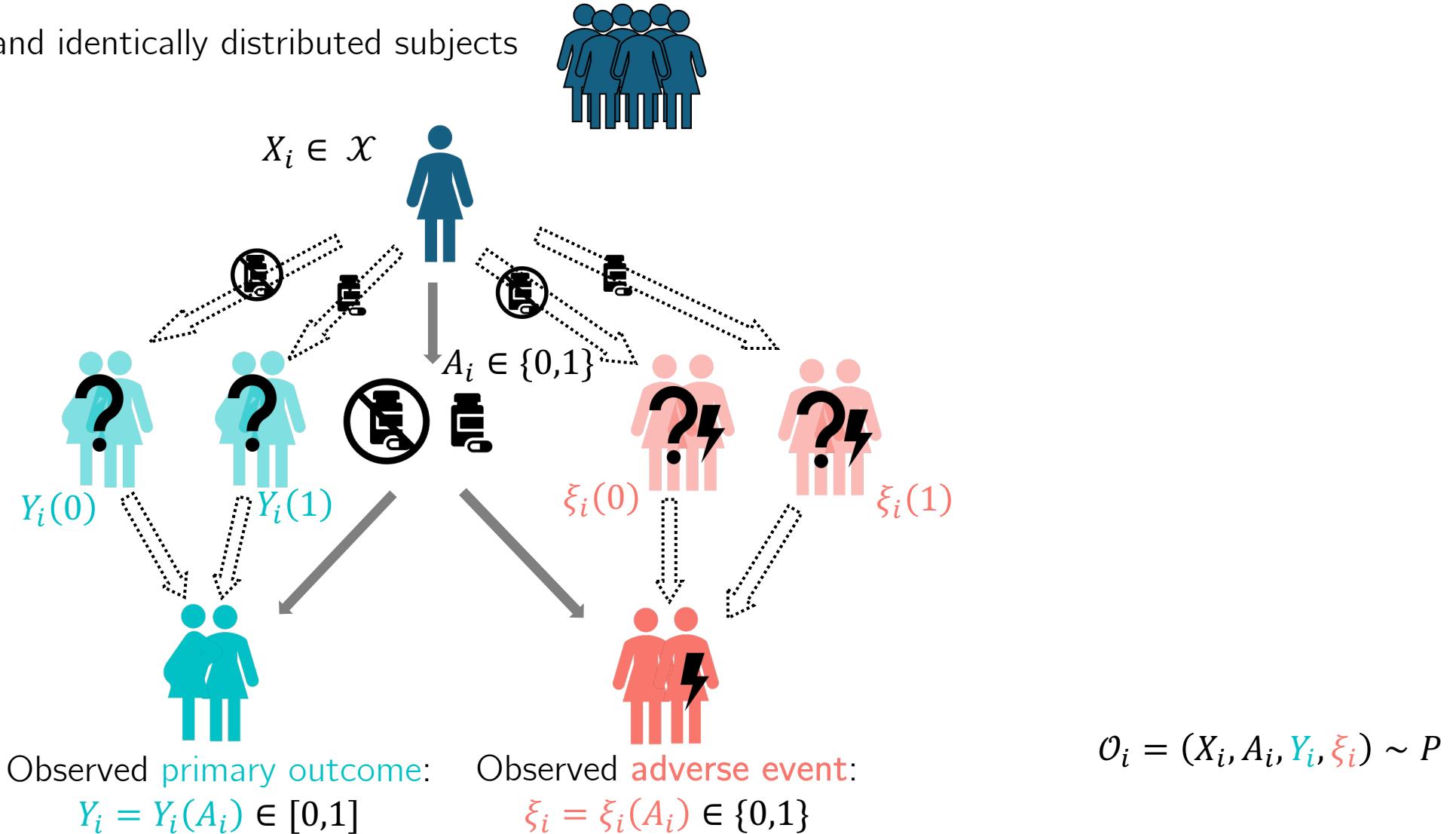
# I.2-Mathematical framework

Set of independent and identically distributed subjects



# I.2-Mathematical framework

Set of independent and identically distributed subjects



# I.2-Policy learning framework

**Policy:**  $\pi \in \Pi$  maps patient's characteristics to a treatment decision

$\pi : X \rightarrow \{0,1\}$  or  $[0,1]$



# I.2-Policy learning framework

**Policy:**  $\pi \in \Pi$  maps patient's characteristics to a treatment decision

$$\pi : X \rightarrow \{0,1\} \text{ or } [0,1]$$



Classical policy learning: Find the value-optimal policy  $\pi^*$ , maximizing the average outcome if  $\pi$  is followed,

$$x \mapsto \pi^*(x) = \mathbf{1}\{\Delta\mu(x) > 0\}$$

$$\left| \begin{array}{l} \mu(a, x) = E[Y|a, x] \\ \Delta\mu(x) = \mu(1, x) - \mu(0, x) \end{array} \right.$$

Covariates	Treatment	$\mu(\mathbf{0}, X)$	$\mu(\mathbf{1}, X)$	$\Delta\mu(X)$	$\pi^*(X)$
4   1   a	1	NA	0.9	NA	?
10   0   b	0	0.6	NA	NA	?
2   1   a	0	0.3	NA	NA	?

# I.2-Policy learning framework

**Policy:**  $\pi \in \Pi$  maps patient's characteristics to a treatment decision

$\pi : X \rightarrow \{0,1\}$  or  $[0,1]$



Classical policy learning: Find the value-optimal policy  $\pi^*$ , maximizing the average outcome if  $\pi$  is followed,

$$x \mapsto \pi^*(x) = \mathbf{1}\{\Delta\mu(x) > 0\}$$

1- Direct: Predict sign of  $\Delta\mu(X)$

$$\left| \begin{array}{l} \mu(a, x) = E[Y|a, x] \\ \Delta\mu(x) = \mu(1, x) - \mu(0, x) \end{array} \right.$$

2- Indirect: Plug-in an estimation of  $\Delta\mu$

$$x \mapsto \hat{\pi}^*(x) = \mathbf{1}\{\hat{\Delta}\mu_n(x) > 0\}$$

Covariates	Treatment	$\hat{\mu}_n(\mathbf{0}, X)$	$\hat{\mu}_n(\mathbf{1}, X)$	$\hat{\Delta}\mu_n(X)$	$\hat{\pi}^*(X)$
4   1   a	1	0.7	< 0.9	0.2	1
10   0   b	0	0.6	> 0.3	-0.3	0
2   1   a	0	0.3	< 0.4	0.1	1

# I.2-Policy learning framework

**Policy:**  $\pi \in \Pi$  maps patient's characteristics to a treatment decision

$\pi : X \rightarrow \{0,1\}$  or  $[0,1]$



Classical policy learning: Find the value-optimal policy  $\pi^*$ , maximizing the average outcome if  $\pi$  is followed,

$$x \mapsto \pi^*(x) = \mathbf{1}\{\Delta\mu(x) > 0\}$$

1- Direct: Predict sign of  $\Delta\mu(X)$

2- Indirect: Plug-in an estimation of  $\Delta\mu$

$$x \mapsto \hat{\pi}^*(x) = \mathbf{1}\{\hat{\Delta}\mu_n(x) > 0\}$$

$$\left| \begin{array}{l} \mu(a, x) = E[Y|a, x] \\ \Delta\mu(x) = \mu(1, x) - \mu(0, x) \end{array} \right.$$



Relies on quality of estimator  $\hat{\mu}_n$

Problematic values with some estimators

Covariates	Treatment	$\hat{\mu}_n(\mathbf{0}, X)$	$\hat{\mu}_n(\mathbf{1}, X)$	$\hat{\Delta}\mu_n(X)$	$\hat{\pi}^*(X)$
4   1   a	1	0.7	0.9	0.2	1
10   0   b	0	0.6	0.3	-0.3	0
2   1   a	0	0.3	0.4	0.1	1

## II. Methods

## II.1- Problem formulation

**EP-learning [2]**: Efficient plug-in risk estimator

1- Plug-in efficient estimator of  $\Delta\mu$

2- Minimize risk (1) for  $\psi \in \Psi$   
(convex space)

$$R: \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X)] \quad (1)$$

 Overlooks adverse events

## II.1- Problem formulation

**EP-learning [2]**: Efficient plug-in risk estimator

- 1- Plug-in efficient estimator of  $\Delta\mu$
- 2- Minimize risk (1) for  $\psi \in \Psi$

$$R: \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X)] \quad (1)$$

 Overlooks adverse events



Add constraint: Guarantee that  $\pi$  does not increase the average probability of adverse event beyond  $\alpha \in [0, \frac{1}{2}]$

$$\begin{aligned} v(A, X) &= E[\xi|A, X] \\ \Delta v(X) &= v(1, X) - v(0, X) \end{aligned}$$

# II.1- Problem formulation

**EP-learning [2]**: Efficient plug-in risk estimator

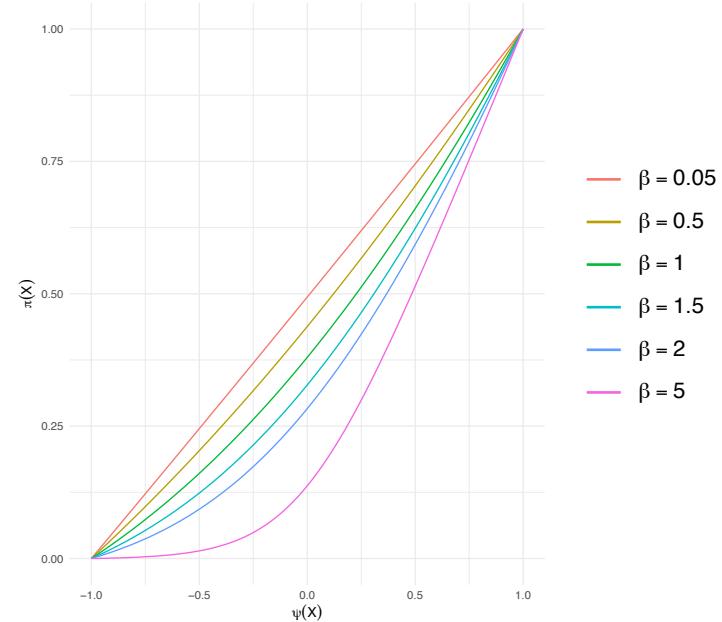
- 1- Plug-in efficient estimator of  $\Delta\mu$
- 2- Minimize risk (1) for  $\psi \in \Psi$

$$R: \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X)] \quad (1)$$

Add constraint: Guarantee that  $\pi$  does **not increase** the average probability of **adverse event** beyond  $\alpha \in \left[0, \frac{1}{2}\right]$

Policy reformulation:

$$\pi(X) = \sigma_\beta \circ \psi(X) \in [0,1]$$



$$\begin{aligned} v(A, X) &= E[\xi|A, X] \\ \Delta v(X) &= v(1, X) - v(0, X) \end{aligned}$$

- ✓ Overcome differentiability issues & strongly convex
- ✓ Confidence measure in treatment recommendation

## II.1- Problem formulation

**EP-learning [2]**: Efficient plug-in risk estimator

- 1- Plug-in efficient estimator of  $\Delta\mu$
- 2- Minimize risk (1) for  $\psi \in \Psi$

$$R: \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X)] \quad (1)$$

Add constraint: Guarantee that  $\pi$  does **not increase** the average probability of **adverse event** beyond  $\alpha \in [0, \frac{1}{2}]$

$$S: \psi \mapsto E\left[\underbrace{\sigma_\beta \circ \psi(X)}_{\pi(X)} \cdot \Delta\nu(X)\right] - \alpha \quad (2)$$

$$\begin{aligned} v(A, X) &= E[\xi|A, X] \\ \Delta\nu(X) &= v(1, X) - v(0, X) \end{aligned}$$

| Assumption: Treated patients are more likely to suffer from an adverse event, that is:  $\Delta\nu(X) \geq 0$

# II.1- Problem formulation

**EP-learning [2]**: Efficient plug-in risk estimator

- 1- Plug-in efficient estimator of  $\Delta\mu$
- 2- Minimize risk (1) for  $\psi \in \Psi$

$$R: \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X)] \quad (1)$$

Add constraint: Guarantee that  $\pi$  does not increase the average probability of **adverse event** beyond  $\alpha \in [0, \frac{1}{2}]$

$$S: \psi \mapsto E[\sigma_\beta \circ \psi(X) \cdot \Delta\nu(X)] - \alpha \quad (2)$$

$$\left| \begin{array}{l} \nu(A, X) = E[\xi|A, X] \\ \Delta\nu(X) = \nu(1, X) - \nu(0, X) \end{array} \right.$$

## Policy optimization objective

Minimize  $R(\psi)$  w.r.t  $\psi \in \Psi$   
s.t.  $S(\psi) \leq 0$

Lagrangian

$$\begin{aligned} \mathcal{L} : \Psi \times \mathbb{R}_+ &\rightarrow \mathbb{R} \\ (\psi, \lambda) &\mapsto R(\psi) + \lambda S(\psi) \end{aligned}$$

# II.1- Problem formulation

**EP-learning [2]**: Efficient plug-in risk estimator

- 1- Plug-in efficient estimator of  $\Delta\mu$
- 2- Minimize risk (1) for  $\psi \in \Psi$

$$R: \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X)] \quad (1)$$

Add constraint: Guarantee that  $\pi$  does not increase the average probability of **adverse event** beyond  $\alpha \in [0, \frac{1}{2}]$

$$S: \psi \mapsto E[\sigma_\beta \circ \psi(X) \cdot \Delta\nu(X)] - \alpha \quad (2)$$

$$\left| \begin{array}{l} \nu(A, X) = E[\xi|A, X] \\ \Delta\nu(X) = \nu(1, X) - \nu(0, X) \end{array} \right.$$

## Policy optimization objective

$$\begin{aligned} & \text{Minimize } R(\psi) \text{ w.r.t } \psi \in \Psi \\ & \text{s.t. } S(\psi) \leq 0 \end{aligned}$$



$$\begin{aligned} \mathcal{L} : \Psi \times \mathbb{R}_+ &\rightarrow \mathbb{R} \\ (\psi, \lambda) &\mapsto R(\psi) + \lambda S(\psi) \end{aligned}$$

## Parameters to fine-tune:

-  $\beta$  policy definition

-  $\lambda$  weight given to constraint

## II.2- Algorithm: EP-learner for constrained policy estimation

1- Estimate  $J$  cross-validated nuisance parameters:  $\hat{\mu}_{n,j}, \hat{v}_{n,j}$

for  $j = 1, \dots, J-1$  do:

    for  $\beta \in B$  do:

        for  $\lambda \in \Lambda$  do:

            2.1- Plug-in  $\hat{\mu}_{n,j}, \hat{v}_{n,j}$  in the objective function:  $\hat{\mathcal{L}}_{n,j}(\psi, \lambda) \mapsto \hat{R}_{n,j}(\psi) + \lambda \hat{S}_{n,j}(\psi)$  (4)

            2.2- Obtain minimizer of (4):  $\hat{\psi}_{n,j} \in \operatorname{argmin}\{\hat{\mathcal{L}}_{n,j}(\psi, \lambda): \psi \in \Psi\}$  [Frank-Wolfe algorithm, 3]

3- Identify the optimal  $\beta^{opt}$  and  $\lambda^{opt}$

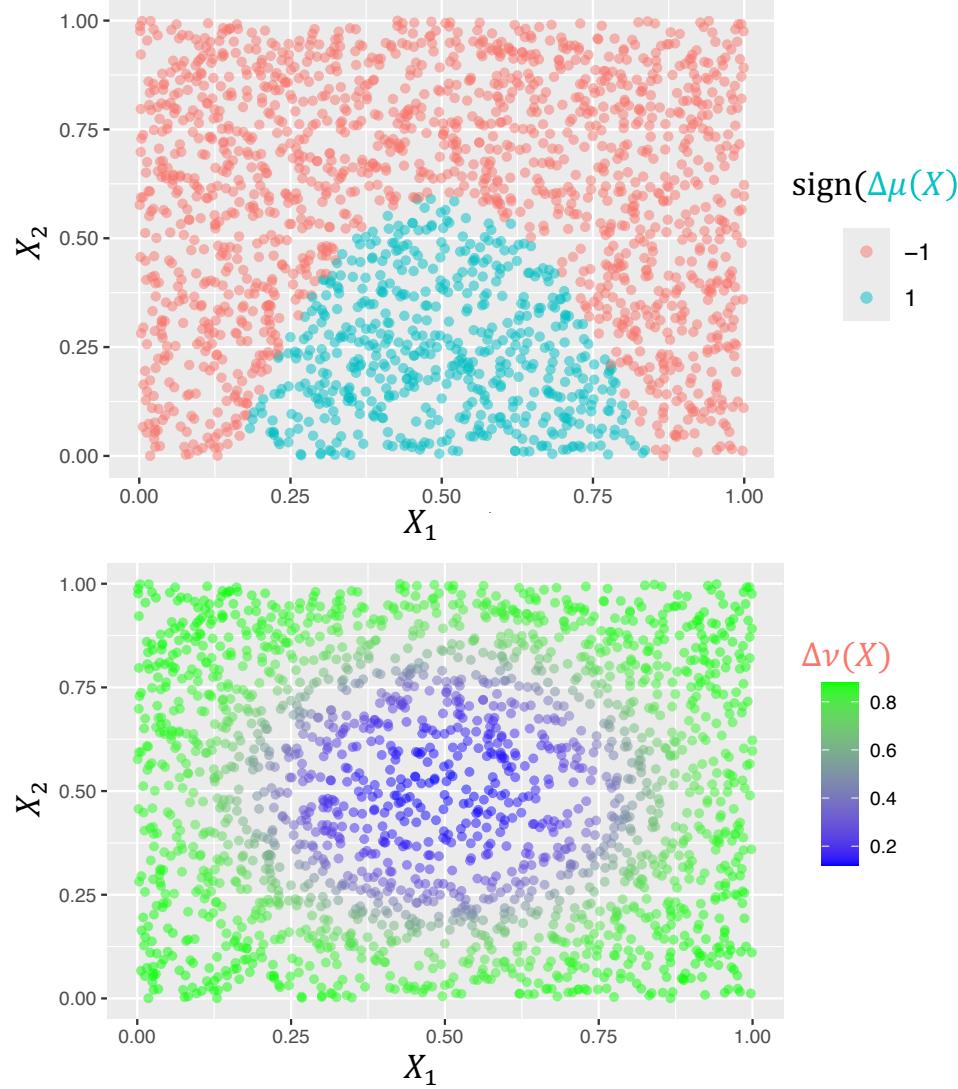
4- Obtain  $\hat{\psi}_n$  minimizer of  $\hat{\mathcal{L}}_{n,J}(\psi, \lambda^{opt})$

return  $\hat{\pi}(X) = \sigma_{\beta^{opt}} \circ \hat{\psi}_n(X)$

# III. Results

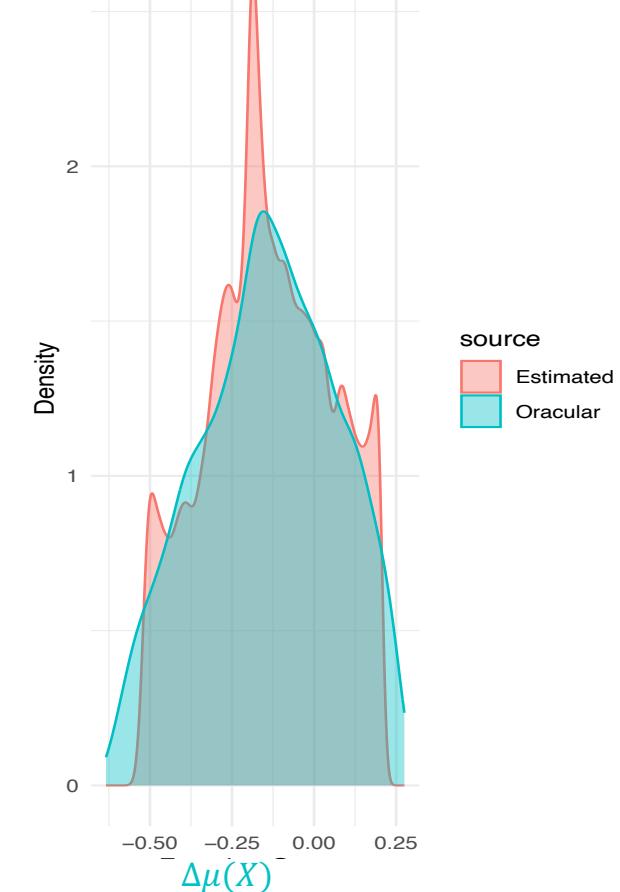
# III.1- Results: synthetic data

Synthetic data

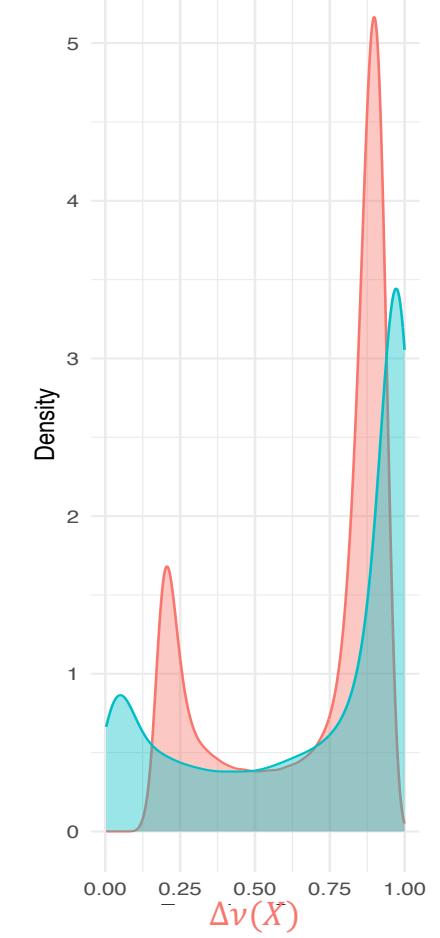


Nuisance parameter estimation

Density of  $\Delta\mu(X)$  vs.  $\hat{\Delta\mu}_n(X)$

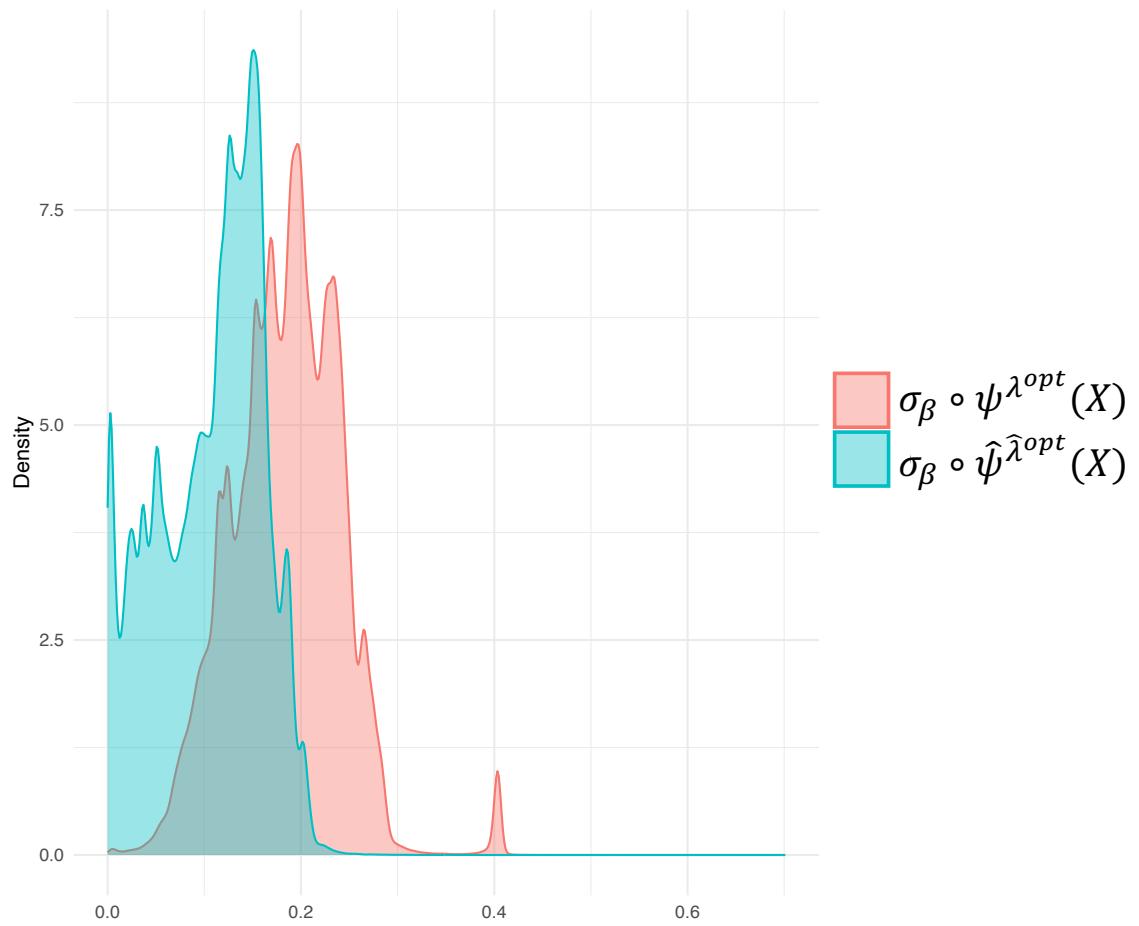


Density of  $\Delta\nu(X)$  vs.  $\hat{\Delta\nu}_n(X)$

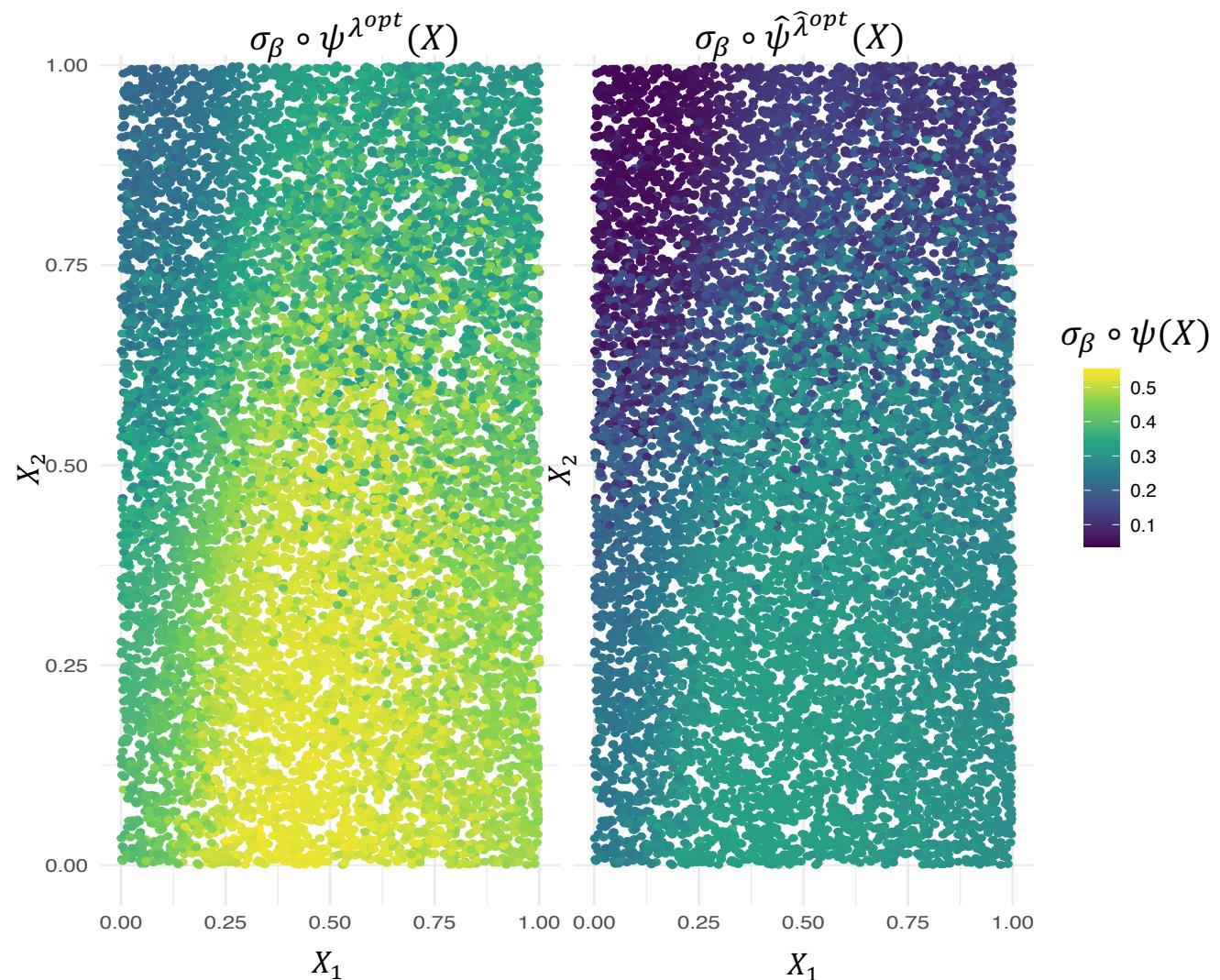


# III.1- Results: synthetic data

Density comparison



Mean treatment probability comparison



# IV. Discussion

# Discussion

- Formulated a policy optimization problem with constraints
- Defined a strongly convex policy in a probabilistic style (in  $[0, 1]$ )
- Adapted the EP-learner for multiple outcomes
- Optimizing policy in a convex function space using the Frank-Wolfe algorithm

## Perspectives:

- Correct plug-in  $\hat{\mathcal{L}}_n(\psi)$ , while **preserving strong-convexity**
- **Alternated optimization**-estimator correction
- Compare to state of the art constrained policy optimization techniques

*Thank you!*

# References

- [1] Rubin, D. B. (2005). Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331.
- [2] Van der Laan, L., Carone, M., and Luedtke, A. (2024). Combining T-learning and DR-learning: a framework for oracle-efficient estimation of causal contrasts arXiv preprint arXiv:2402.01972.
- [3] Jaggi, M. (2013). Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In International conference on machine learning, pages 427–435. PMLR.

# Annexes

# Bias of objective function estimator



Minimize  $\mathcal{L} : \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X) + \lambda\sigma_\beta \circ \psi(X) \cdot \Delta\nu(X)] - \lambda\alpha$

✗  $\mu, \nu$  unknown

Plug-in estimator:  $\hat{\mathcal{L}}_n : \psi \mapsto E_{P_n}[\psi(X)^2 - 2\psi(X) \cdot \hat{\Delta}\mu_n(X) + \lambda\sigma_\beta \circ \psi(X) \cdot \hat{\Delta}\nu_n(X)] - \lambda\alpha$

Solution

# Bias of objective function estimator



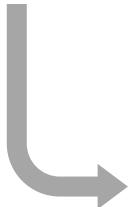
Minimize  $\mathcal{L} : \psi \mapsto E[\psi(X)^2 - 2\psi(X) \cdot \Delta\mu(X) + \lambda\sigma_\beta \circ \psi(X) \cdot \Delta\nu(X)] - \lambda\alpha$

✗  $\mu, \nu$  unknown

Plug-in estimator:  $\hat{\mathcal{L}}_n : \psi \mapsto E_{P_n}[\psi(X)^2 - 2\psi(X) \cdot \hat{\Delta}\mu_n(X) + \lambda\sigma_\beta \circ \psi(X) \cdot \hat{\Delta}\nu_n(X)] - \lambda\alpha$

Solution

Evaluation via  
EIC derivation of  
 $\mathcal{F} : P \mapsto \mathcal{L}_P(\psi)$



$$EIC_{\mathcal{F}}(P)(X, A, Y, \xi) = \nabla_1 \mathcal{F}(P) + \underbrace{\nabla_2 \mathcal{F}(P)}_{E[\nabla_2 \mathcal{F}(P)] = 0}$$

Estimator bias

$$E[\nabla_1 \mathcal{F}(P)(X, A, Y, \xi)] = E \left[ \frac{2A - 1}{P(A|X)} \left( -2\psi(X) \cdot (\textcolor{teal}{Y} - \mu(A, X)) + \lambda\sigma_\beta \circ \psi(X) \cdot (\xi - \nu(A, X)) \right) \right]$$

✗ Non-convex function of  $\psi$

## II.2- Algorithm: Alternated procedure

1- Estimate  $J$  cross-validated nuisance parameters:  $\hat{\mu}_{n,j}, \hat{v}_{n,j}$

for  $j = 1, \dots, J-1$  do:

    for  $\beta \in B$  do:

        for  $\lambda \in \Lambda$  do:

2- Run iteratively alternated procedure for fixed  $\beta, \lambda$ :

- Minimize  $\hat{\mathcal{L}}_{n,j}^k$ :  $\hat{\psi}_{n,j}^k \in \operatorname{argmin}\{\hat{\mathcal{L}}_{n,j}^k(\psi, \lambda): \psi \in \Psi\}$
- Correct  $\hat{\mathcal{L}}_{n,j}^k$  by adjusting nuisance parameters  $\hat{\mu}_{n,j}, \hat{v}_{n,j}$  for fixed  $\hat{\psi}_{n,j}^k$

3- Identify the optimal  $\beta^{opt}$  and  $\lambda^{opt}$

4- Obtain  $\hat{\psi}_n$  minimizer of  $\hat{\mathcal{L}}_{n,J}^k(\psi, \lambda^{opt})$

return  $\hat{\pi}(X) = \sigma_{\beta^{opt}} \circ \hat{\psi}_n(X)$

## II.2- Algorithm: Alternated procedure

**Require:**  $\lambda, \beta, \hat{\mu}_{n,j}, \hat{v}_{n,j}, e_n, \gamma, \text{max\_iter}$

**Initialization:**

**Optimization:**  $\hat{\mu}_{n,j}^0(\epsilon_1) = \hat{\mu}_{n,j}, \hat{v}_{n,j}^0(\epsilon_2) = \hat{v}_{n,j}$

$\psi_{n,j}^0 \in \operatorname{argmin}\{\hat{\mathcal{L}}_{n,j}^0(\psi, \lambda): \psi \in \Psi\}$  [Frank-Wolfe algorithm, 3]

**While**  $(\sum(\hat{\psi}_{n,j}^{k-1} - \hat{\psi}_{n,j}^k)^2) < \gamma \ \& \ k < \text{max\_iter}$ :

$$k = k + 1$$

**Correction:**  $\epsilon \in \operatorname{argmin} \ell_n^k(\epsilon) = \begin{cases} \frac{-1}{n} \sum_{i=1}^n Y_i \log(\hat{\mu}_n^k(\epsilon_1)) + (1 - Y_i) \log(1 - \hat{\mu}_n^k(\epsilon_1)) \\ \frac{-1}{n} \sum_{i=1}^n \xi_i \log(\hat{v}_n^k(\epsilon_2)) + (1 - \xi_i) \log(1 - \hat{v}_n^k(\epsilon_2)) \end{cases}$

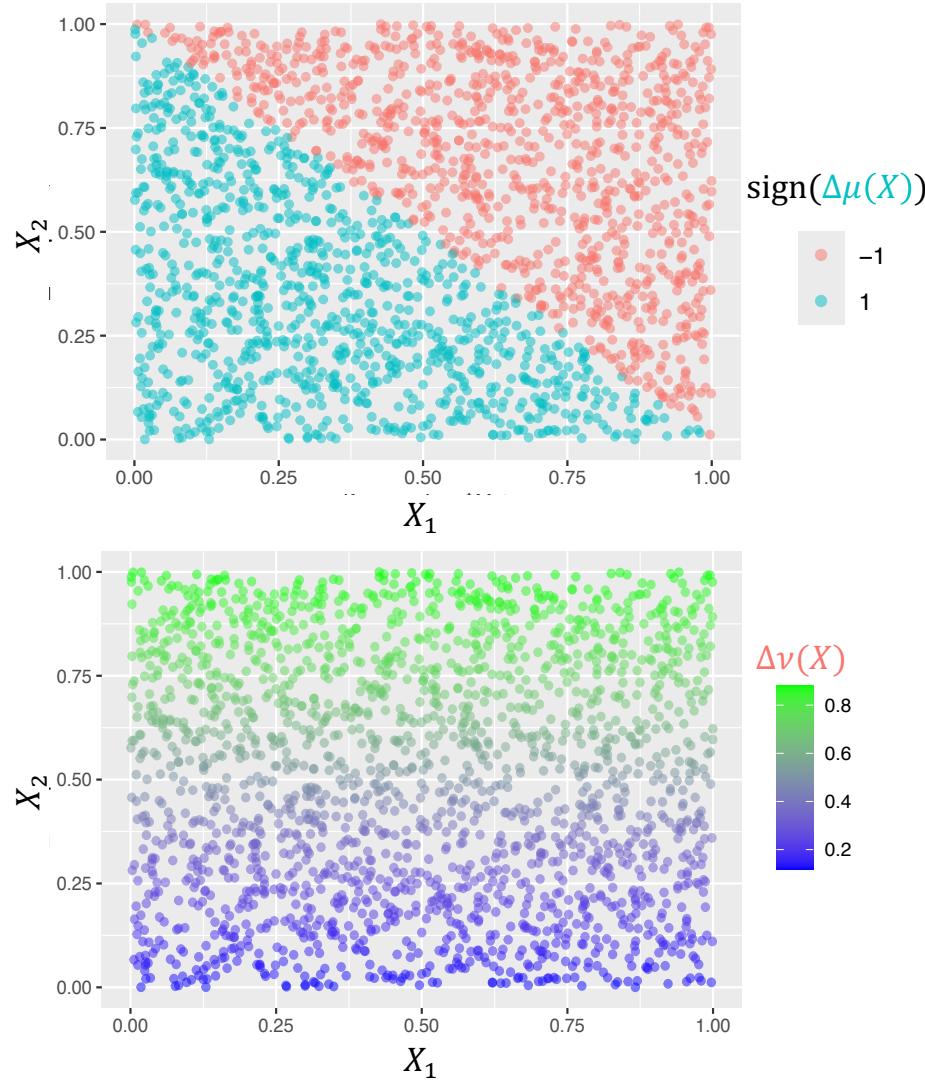
Nuisance parameter update:

$$\hat{\mu}_{n,j}^k(\epsilon_1) = \text{expit}(\text{logit}(\hat{\mu}_{n,j}^0) + \frac{2A-1}{e_n(A,X)} \sum \epsilon_{1,l} \hat{\psi}_l(X)) \quad \hat{v}_{n,j}^k(\epsilon_2) = \text{expit}(\text{logit}(\hat{v}_{n,j}^0) + \frac{2A-1}{e_n(A,X)} \sum \epsilon_{2,l} \sigma_\beta \circ \hat{\psi}_l(X))$$

**Optimization:**  $\hat{\psi}_{n,j}^k \in \operatorname{argmin}\{\hat{\mathcal{L}}_{n,j}^k(\psi, \lambda): \psi \in \Psi\}$  [Frank-Wolfe algorithm, 3]

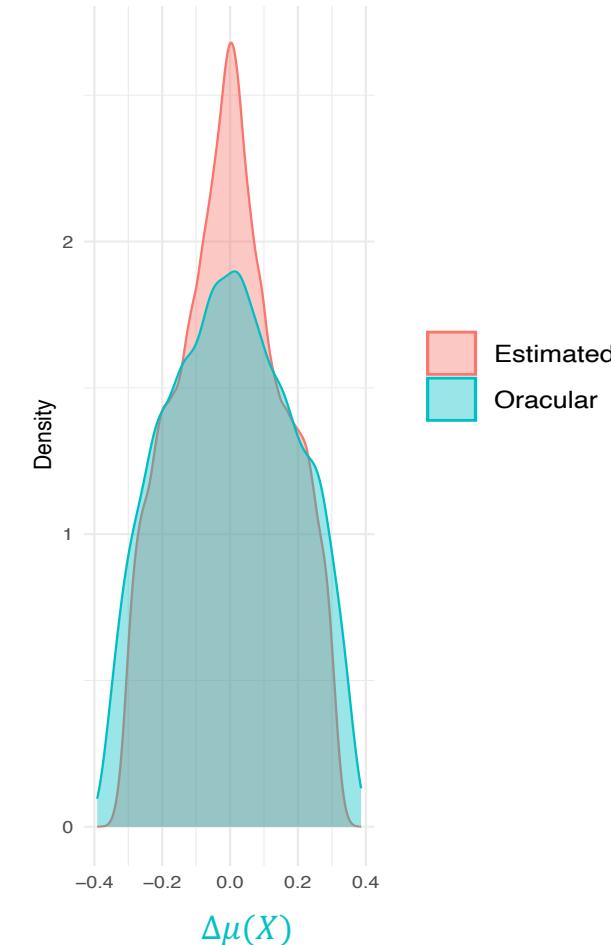
## III.2- Results: synthetic data, second scenario

Synthetic data

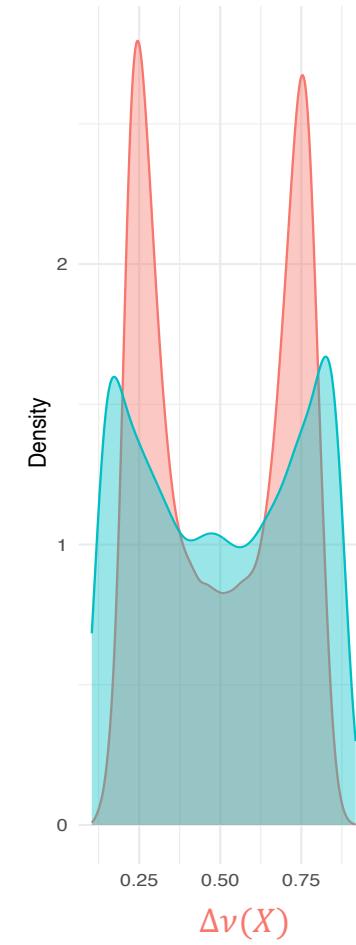


Nuisance parameter estimation

Density of  $\Delta\mu(X)$  vs.  $\widehat{\Delta\mu}_n(X)$

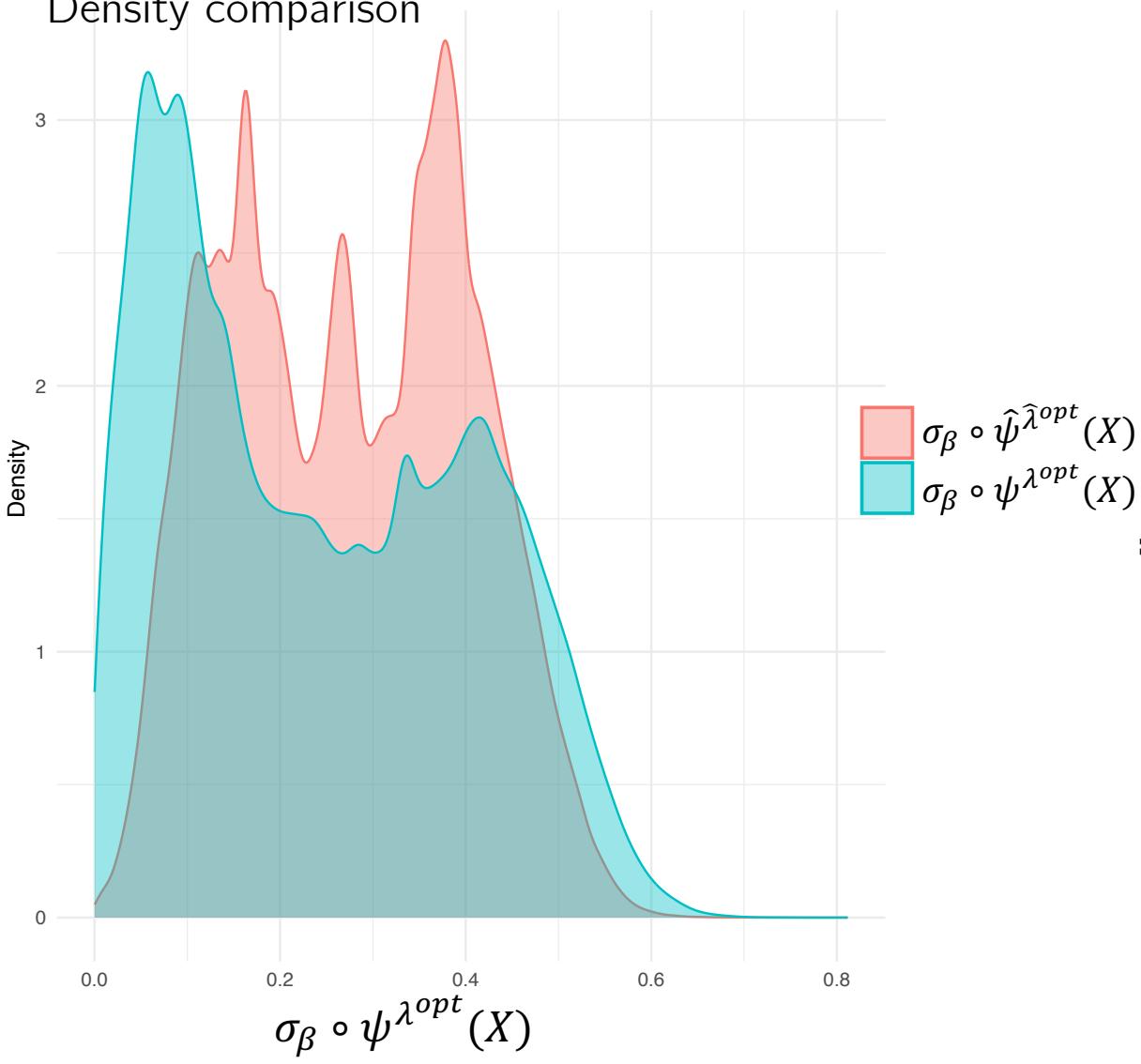


Density of  $\Delta\nu(X)$  vs.  $\widehat{\Delta\nu}_n(X)$



## III.2- Results: synthetic data, second scenario

Density comparison



Mean treatment probability comparison

