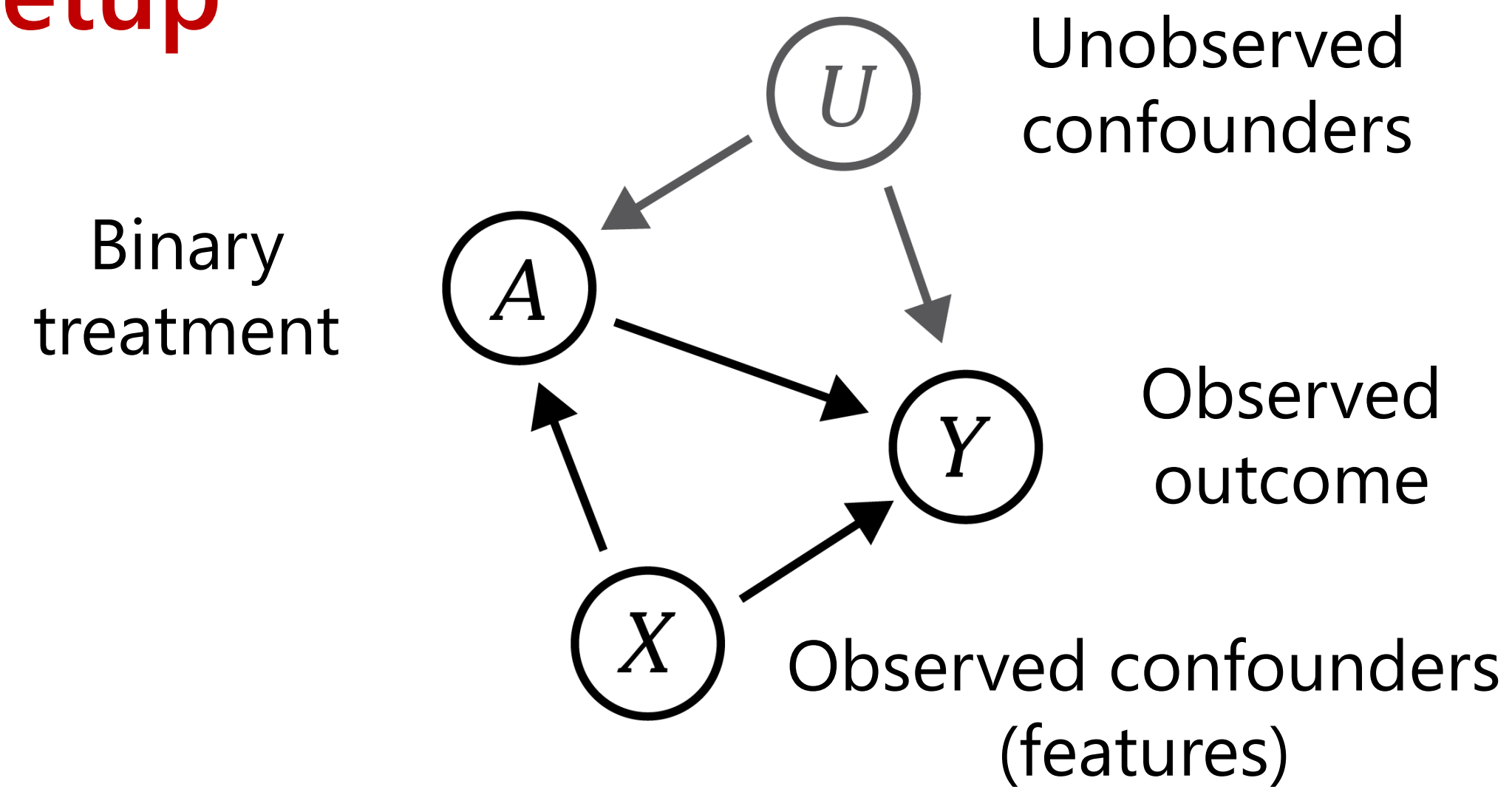


Setup



We want the **bounds** $\tau^+(x), \tau^-(x)$ on the CATE function:
 $\tau(x) = \mathbb{E}[Y(1) - Y(0) | X = x]$

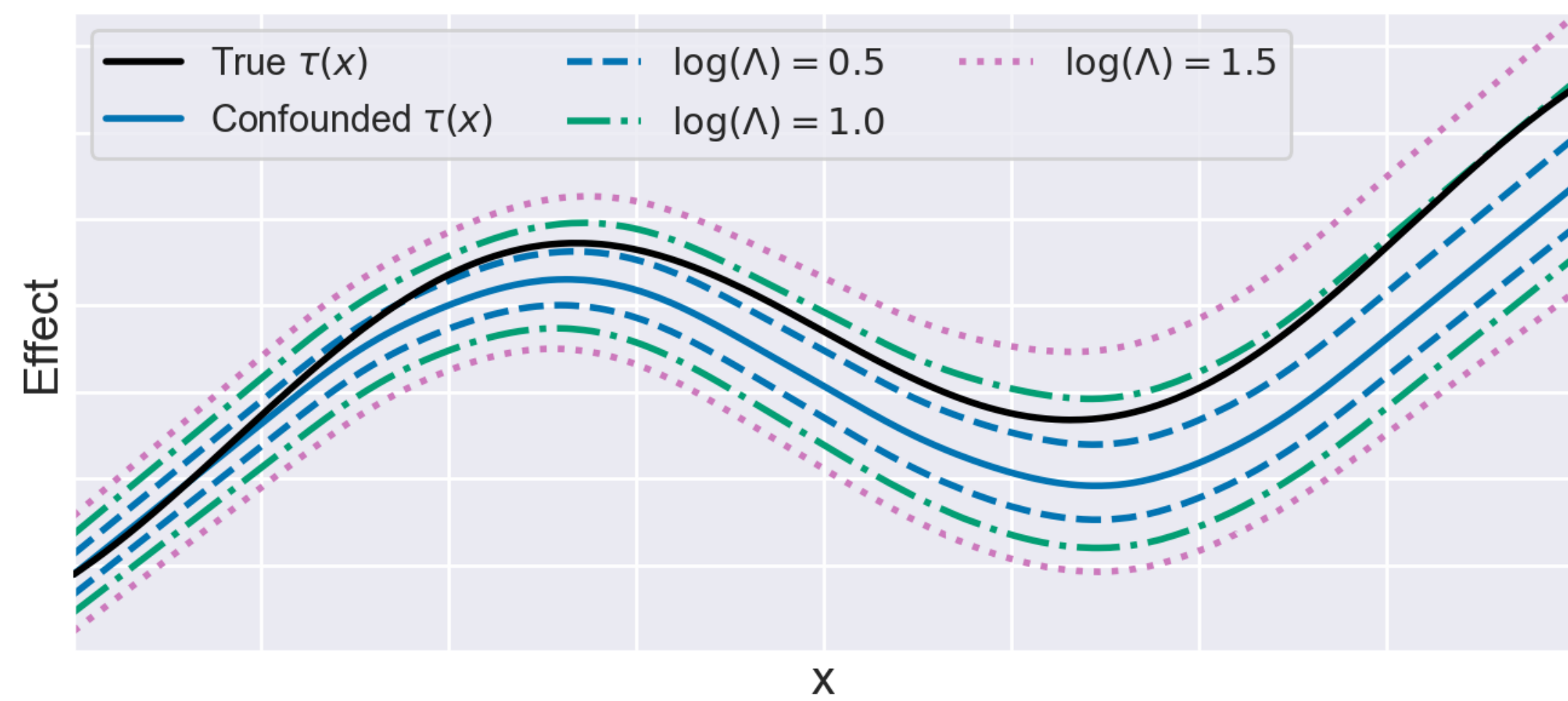
Sensitivity Analysis

Assumption: Marginal Sensitivity Model (MSM).

$$\Lambda^{-1} \leq \frac{e(x, u)}{1 - e(x, u)} \bigg/ \frac{e(x)}{1 - e(x)} \leq \Lambda$$

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

$$e(x, u) = P(A = 1 | X = x, U = u)$$



Identification of CATE Bounds

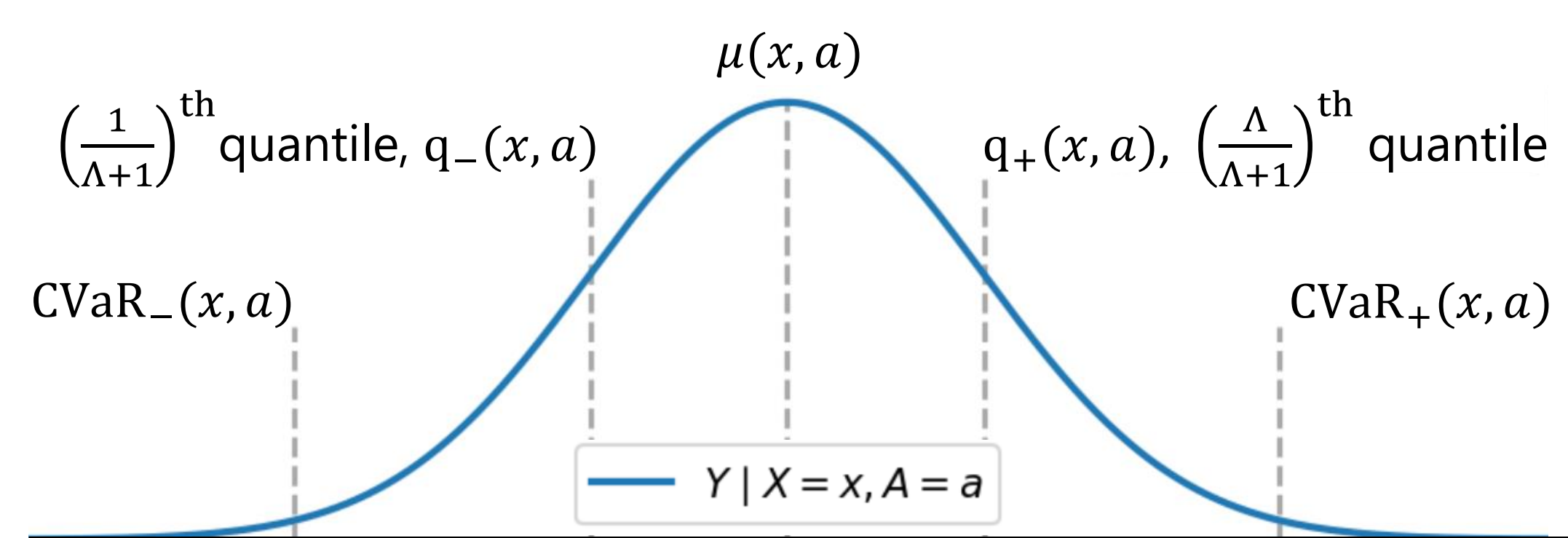
Result 1 (Dorn et al., 2021). $\mu(x, a) = \mathbb{E}[Y | X = x, A = a]$ and $Y^\pm(x, a)$ is the upper (+) / lower (-) sharp bound of $\mathbb{E}[Y(a) | X = x]$ (not identifiable). Then:

$$Y^+(x, 1) = e(x)\mu(x, 1) + (1 - e(x))\rho_+(x, 1)$$

$$Y^-(x, 0) = (1 - e(x))\mu(x, 0) + e(x)\rho_-(x, 0)$$

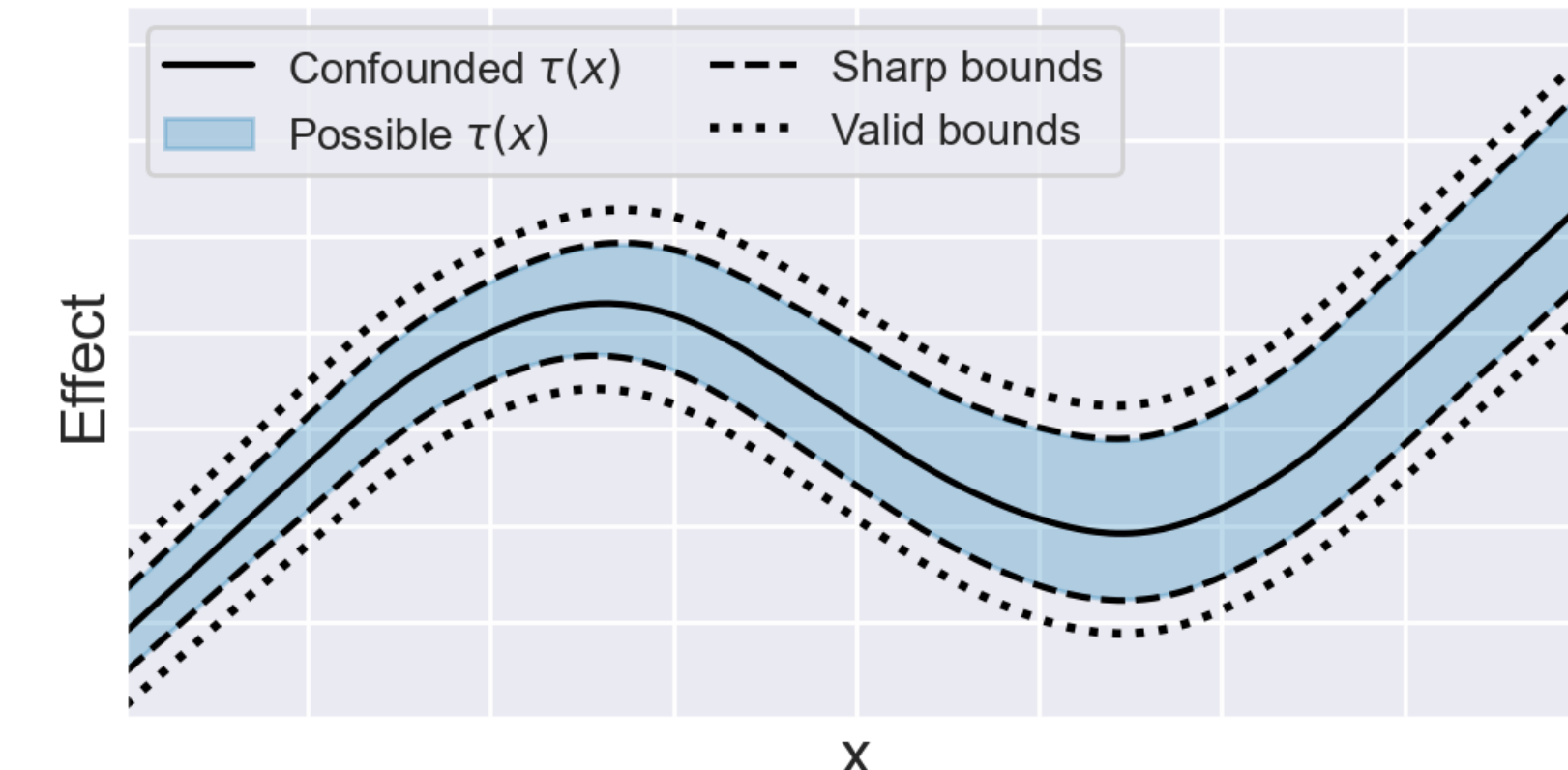
$$\tau^+(x) = Y^+(x, 1) - Y^-(x, 0)$$

where $\rho_\pm(x, a) = \Lambda^{-1}\mu(x, a) + (1 - \Lambda^{-1})CVaR_\pm(x, a)$.



Bound Estimates Should Be...

- **Valid:** $\hat{\tau}(x) \geq \tau(x) + o_p(1)$.
 - **Sharp:** $\hat{\tau}(x) = \tau(x) + o_p(1)$.
 - **Efficient and Robust:** Bounds should converge at desirable rates and have multiple chances at sharp or valid limits.
- Previous works do not achieve all!**



Estimation of CATE Bounds

Naïve "Plug-in" Estimator

Estimate $e(x), \mu(x, a), \rho_\pm(x, a)$ and "plug" them into $Y^\pm(x, a)$ to obtain:

$$\hat{\tau}_{\text{plugin}}^+(x) = \hat{Y}^+(x, 1) - \hat{Y}^-(x, 0)$$

- Inherits bias from the estimated nuisances $\hat{e}(x), \hat{\mu}(x, a), \hat{\rho}_\pm(x, a)$.
- Especially biased when the nuisances are more complex than the CATE bounds.
- Does not yield efficient or robust bounds!

B-Learner

1. Estimate nuisance set $\hat{\eta} = (\hat{e}(x), \hat{q}_\pm(x, a), \hat{\rho}_\pm(x, a))$ in one sample.
2. Get pseudo-outcomes based on the efficient influence function (EIF):

$$Y^+(x, 1) \rightarrow \phi_1^+(Z, \hat{\eta}) = AY + (1 - A)\hat{\rho}_+(X, 1) + \frac{(1 - \hat{e}(X))A}{\hat{e}(X)}(R_+(Z, \hat{q}_+(X, 1)) - \hat{\rho}_+(X, 1))$$

$$Y^-(x, 0) \rightarrow \phi_0^-(Z, \hat{\eta}) = (1 - A)Y + A\hat{\rho}_-(X, 0) + \frac{\hat{e}(X)(1 - A)}{1 - \hat{e}(X)}(R_-(Z, \hat{q}_-(X, 0)) - \hat{\rho}_-(X, 0))$$

$$\tau^+(x) \rightarrow \phi_\tau^+(Z, \hat{\eta}) = \phi_1^+(Z, \hat{\eta}) - \phi_0^-(Z, \hat{\eta})$$

where $\mathbb{E}[R_\pm(Z, q_\pm) | X = x, A = a] = \rho_\pm(x, a)$.

3. Regress pseudo-outcome $\phi_\tau^+(Z, \hat{\eta})$ on features $X \in \mathcal{X}$ in another sample.

Algorithm 1 The B-Learner

input Data $\{(X_i, A_i, Y_i) : i \in \{1, \dots, n\}\}$, folds $K \geq 2$, nuisance estimators, regression learner $\hat{\mathbb{E}}_n$

- 1: **for** $k \in \{1, \dots, K\}$ **do**
- 2: Use data $\{(X_i, A_i, Y_i) : i \neq k - 1 \pmod{K}\}$ to construct nuisance estimates $\hat{\eta}^{(k)} = (\hat{e}^{(k)}, \hat{q}^{(k)}, \hat{\rho}^{(k)})$
- 3: **for** $i = k - 1 \pmod{K}$ **do**
- 4: Set $\hat{\phi}_{\tau, i}^+ = \phi_\tau^+(Z_i, \hat{\eta}^{(k)})$
- 5: **end for**
- 6: **end for**

output $\hat{\tau}^+(x) = \hat{\mathbb{E}}_n[\hat{\phi}_\tau^+ | X = x]$

The B-Learner algorithm with K-fold sample splitting.

Theoretical Guarantees

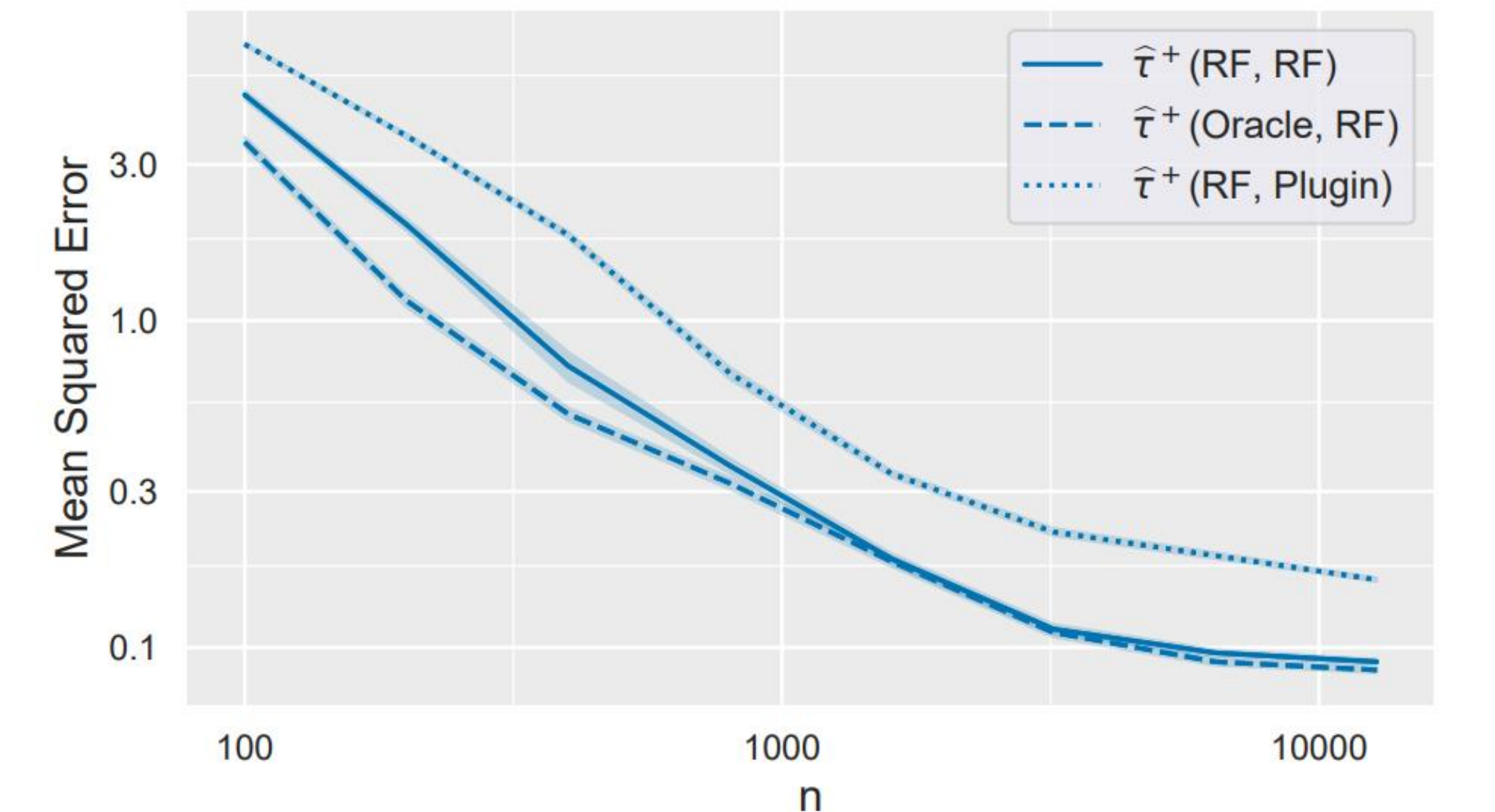
- L_2 validity, sharpness and robustness guarantees for **ERM** second stage estimator $\hat{\mathbb{E}}_n$:
 1. L_2 bias on the order of $\mathcal{E} = \|\hat{e} - e\| \|\hat{\rho} - \rho\| + \|\hat{q} - q\|^2$.
 2. If \hat{q} and either \hat{e} or $\hat{\rho}$ are consistent, the bounds are **sharp** on average.
 3. If \hat{q} is inconsistent, the bounds are still **valid** in expectation.
 4. The B-Learner has **quasi-oracle efficiency**, i.e. the bounds can be learned at a rate dominated by the complexity of the target class.
- **Pointwise** validity, sharpness and robustness guarantees for **linear smoother** second stage estimator $\hat{\mathbb{E}}_n$.

Experiments

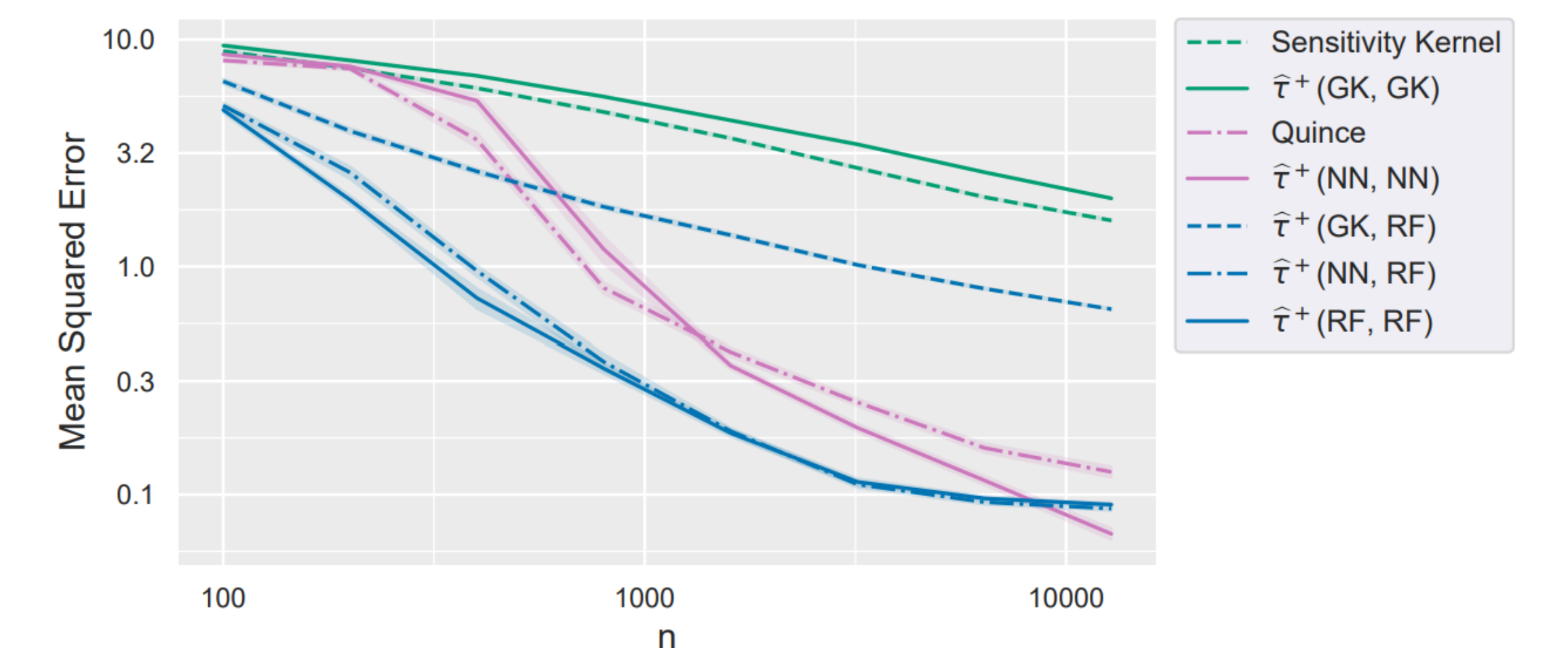
Simulations

$$A \sim \text{Bernoulli}(\text{logit}(0.75X_0 + 0.5))$$

$$Y \sim \mathcal{N}((2A - 1)(X_0 + 1) - 2 \sin((4A - 2)X_0), 1)$$

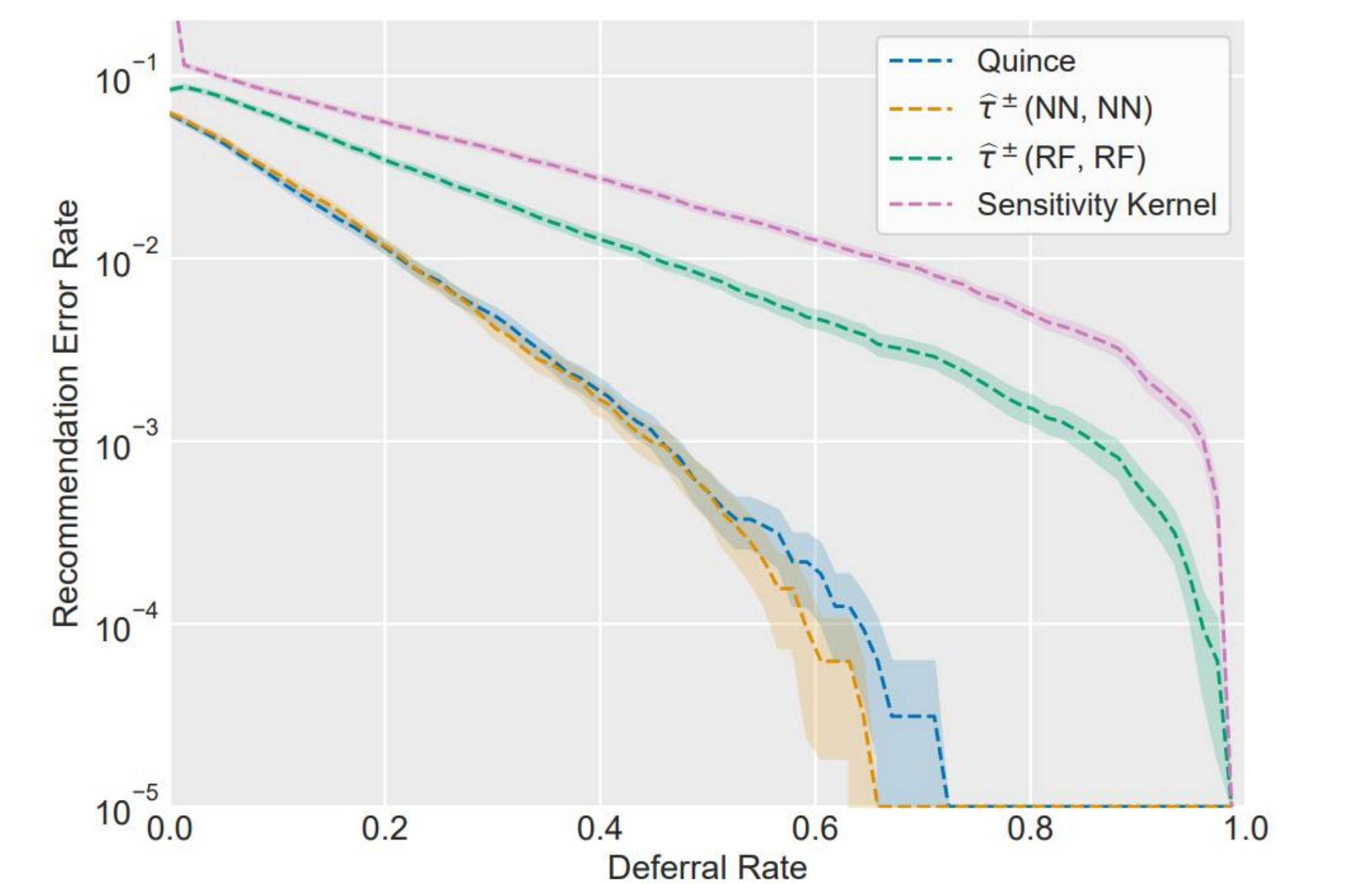


Performance of the B-Learner compared with the *Sensitivity Kernel* (Kallus et al. 2019) and *Quince* (Jesson et al., 2021). GK=Gaussian Kernel, NN=Neural Network, RF=Random Forest.



Performance of the B-Learner compared with the *Sensitivity Kernel* (Kallus et al. 2019) and *Quince* (Jesson et al., 2021). GK=Gaussian Kernel, NN=Neural Network, RF=Random Forest.

IHDP Hidden Confounding



Error recommendation rate for different values of the percentage of deferred points. The x-axis represents different levels of practitioner caution by varying the percentage of recommendations deferred.

¹Cornell University and Cornell Tech

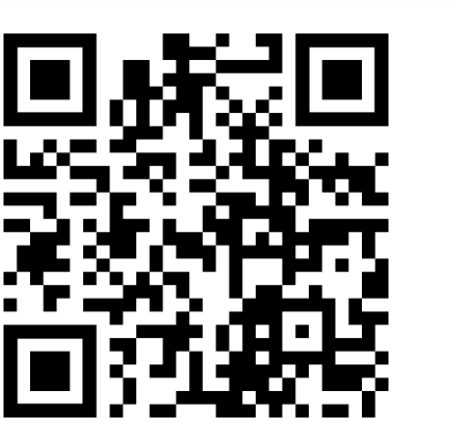
²Princeton University

³Technion, Israel Institute of Technology

⁴OATML, University of Oxford

Correspondence to: Miruna Oprescu

miruna@cs.cornell.edu



Paper



Code