

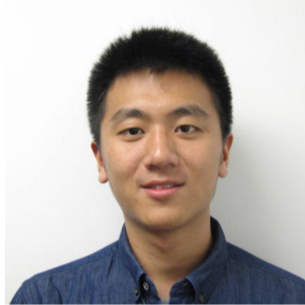
Maxway CRT: Improving the Robustness of Model-X Inference

Shuangning Li

Stanford University

Nov 2021

Joint work with



Molei Liu

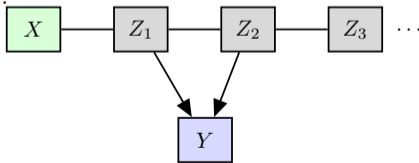
Li & Liu (2021) “Maxway CRT: Improving the Robustness of Model-X Inference” *Draft in preparation.*

Motivation

- ▶ Imagine researchers are interested in whether a particular genetic variant influence a trait.



- ▶ Let X denote the genetic variant. Let Y be the trait.
- ▶ First idea: test for whether $X \perp\!\!\!\perp Y$. Not working because X can be correlated with other variables that influence Y .



Conditional Independence

- ▶ Idea: test for conditional independence!
- ▶ $\mathcal{H}_0 : X \perp\!\!\!\perp Y \mid Z$.
- ▶ Here, Y is the response variable of interest, X is an explanatory variable and Z are confounding variables (potentially high dimensional).
- ▶ There are n i.i.d. samples of (Y, X, Z) denoted as (Y_i, X_i, Z_i) , and let $\mathbf{y} = (Y_1, Y_2, \dots, Y_n)^\top \in \mathbb{R}^n$, $\mathbf{x} = (X_1, X_2, \dots, X_n)^\top \in \mathbb{R}^n$, and $\mathbf{Z} = (Z_{1\cdot}, Z_{2\cdot}, \dots, Z_{n\cdot})^\top \in \mathbb{R}^{n \times p}$.

Conditional Randomization Test

- ▶ Introduced by Candès et al. (2018).

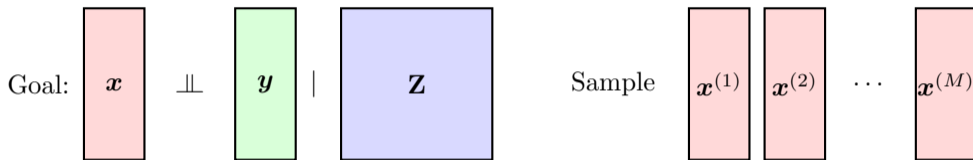
For $m \in \{1, \dots, M\}$

Sample $\mathbf{x}^{(m)}$ from the distribution of $\mathbf{x}|\mathbf{Z}$, independently of \mathbf{x} and \mathbf{y} .

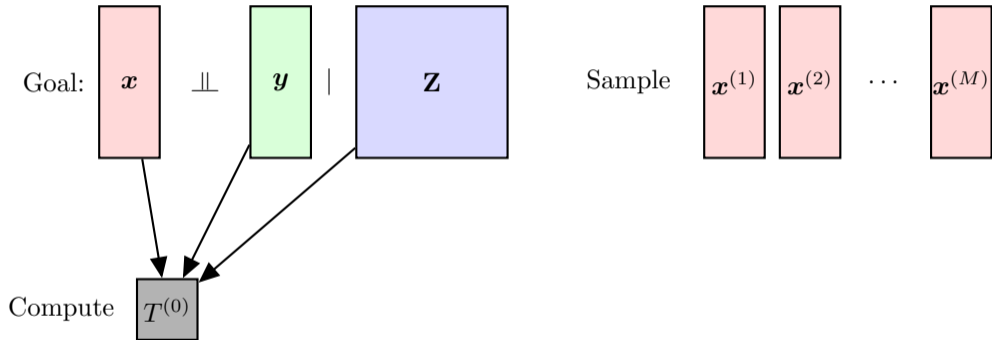
Output The p -value

$$p_{\text{CRT}} = \frac{1}{M+1} \left(1 + \sum_{m=1}^M \mathbb{1} \left\{ T(\mathbf{y}, \mathbf{x}, \mathbf{Z}) \leq T(\mathbf{y}, \mathbf{x}^{(m)}, \mathbf{Z}) \right\} \right).$$

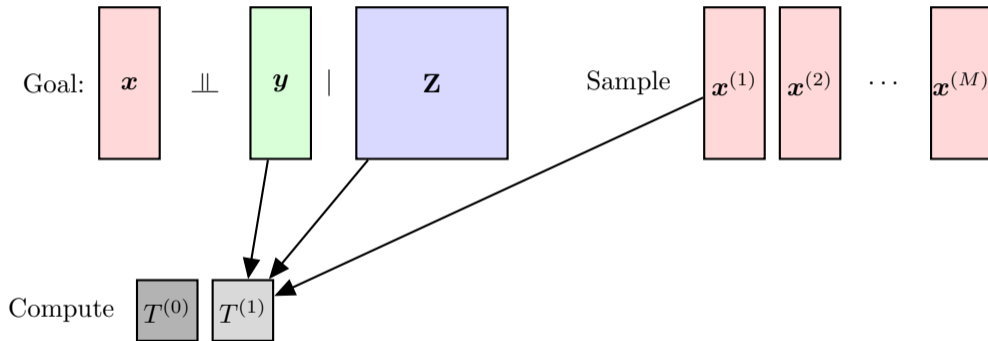
Conditional Randomization Test



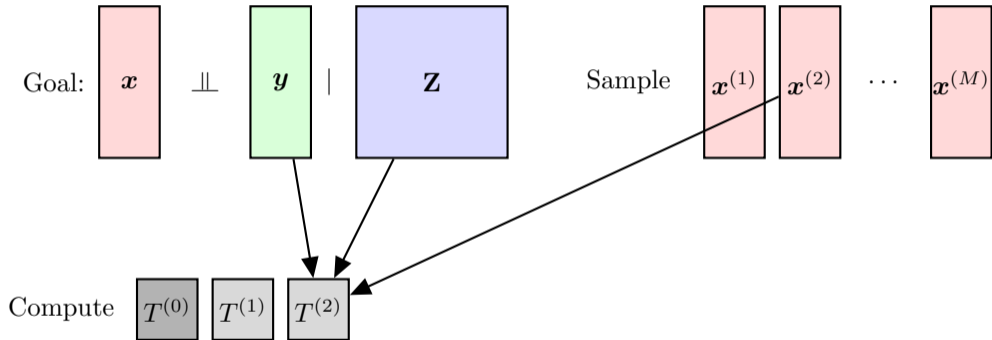
Conditional Randomization Test



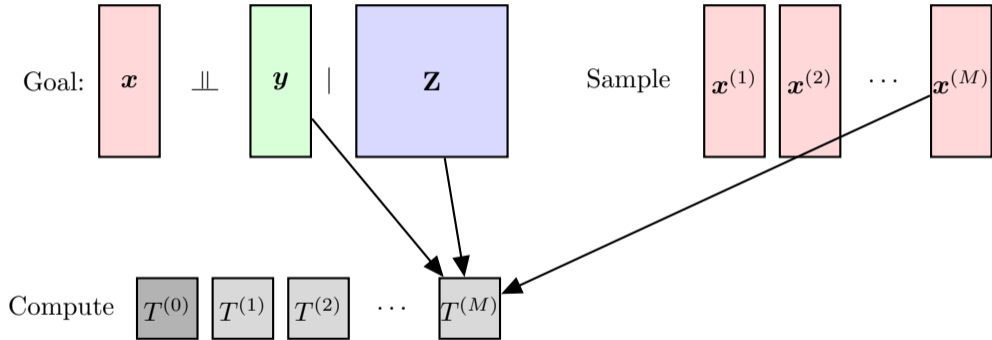
Conditional Randomization Test



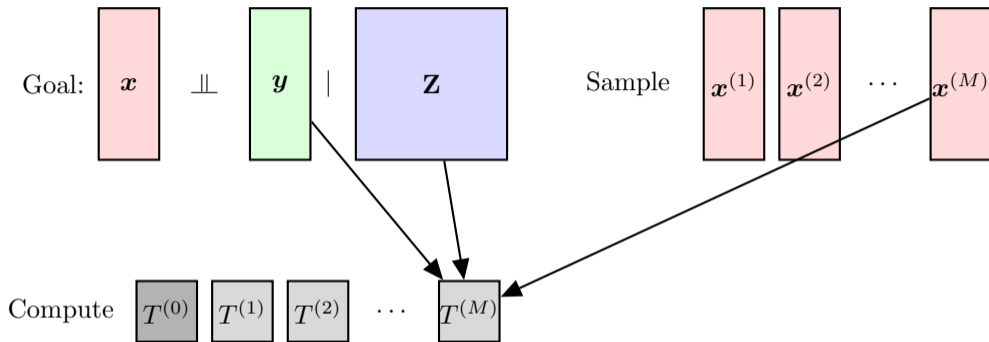
Conditional Randomization Test



Conditional Randomization Test



Conditional Randomization Test



$$p\text{-value } p_j = \frac{1}{M+1} \left(1 + \# \left\{ T^{(m)} : T^{(m)} \geq T^{(0)} \right\} \right)$$

Conditional Randomization Test

Theorem (Candès et al. (2018))

If $X \perp\!\!\!\perp Y \mid Z$, then the p -values from CRT satisfy $\mathbb{P}[p_j \leq \alpha] \leq \alpha$, for any $\alpha \in [0, 1]$. This holds regardless of the test statistic $T(\cdot)$.

- ▶ Requires perfect knowledge of the distribution of $\mathbf{x} \mid \mathbf{Z}$.
- ▶ Let ρ^{*n} be the distribution of $\mathbf{x} \mid \mathbf{Z}$. Assume in CRT, $\mathbf{x}^{(b)}$ are generated instead from ρ^n . Berrett et al. (2020) showed that

$$P(p_{\text{CRT}} \leq \alpha) \leq \alpha + \underbrace{d_{\text{TV}}(\rho^{*n}, \rho^n)}_{\text{Model-X error}}.$$

The bound is tight when M is large.

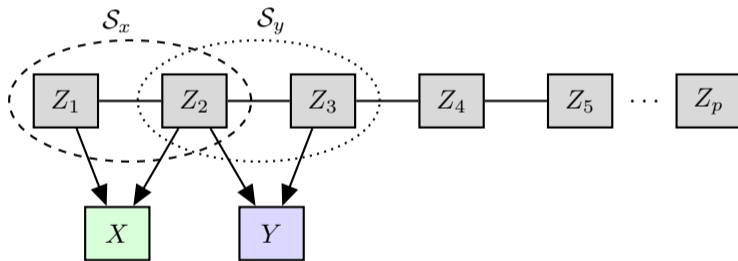
This work

- ▶ In addition to the knowledge of the distribution of $X | Z$, we also have knowledge of the distribution of $Y | Z$. Can we make use of this additional knowledge and make CRT more robust?
- ▶ We propose a Maxway (**M**odel and **A**adjust **X** **W**ith the **A**ssistance of **Y**) CRT.
- ▶ Type-I error inflation of Model-X CRT is Δ_x .
Type-I error inflation of Maxway CRT is $\Delta_x \Delta_y + \Delta_{x|g}$.
Here Δ_x , Δ_y and $\Delta_{x|g}$ are the estimation errors for the distributions of $x | \mathbf{Z}$, $y | \mathbf{Z}$ and $x | g(\mathbf{Z})$ respectively.
- ▶ "double robustness" in type-I error control.

A simple model

Suppose Z is a high dimensional random variable, but X and Y only depend on a small subset of the Z_j 's. Assume that $Z = (Z_1, \dots, Z_p)$, and

$$X = \phi(Z_1, Z_2) + \varepsilon, \quad Y = \psi(Z_2, Z_3) + \eta.$$



- ▶ To implement the original Model-X CRT, need to know S_x and the distribution of $X | Z_{S_x}$.
- ▶ Assume for now that given a set S whose cardinality is not huge, we are able to learn the distribution of $X | Z_S$ accurately.

A simple model

- ▶ A guess of the set: \mathcal{S} . Then (a special version of) the CRT becomes
 1. For $m = 1, 2, \dots, M$:
Sample $\mathbf{x}^{(m)}$ from the distribution of $\mathbf{x} \mid \mathbf{Z}_{\cdot\mathcal{S}}$ independently of (\mathbf{x}, \mathbf{y}) .
 2. Output CRT p -value

$$p_{\text{CRT}} = \frac{1}{M+1} \left(1 + \sum_{m=1}^M \mathbf{1}\{T(\mathbf{y}, \mathbf{x}^{(m)}, \mathbf{Z}_{\cdot\mathcal{S}}) \geq T(\mathbf{y}, \mathbf{x}, \mathbf{Z}_{\cdot\mathcal{S}})\} \right).$$

- ▶ If \mathcal{S} contains $\{1, 2\}$, the p -value is valid.
- ▶ If \mathcal{S} contains $\{2, 3\}$, the p -value is valid as well. This is because $X \perp\!\!\!\perp Y \mid Z_{\mathcal{S}}$. The above procedure can be treated alternatively as a CRT for $(X, Y, Z_{\mathcal{S}})$.
- ▶ Some knowledge of how Y depends on Z can be useful in enhancing robustness in CRT. Without any information or prior knowledge on Y , to achieve validity, the best we can hope for is the set \mathcal{S} to contain \mathcal{S}_x . Extra information on the distribution of Y **relaxes** the condition of the validity of the p -value.

Maxway CRT

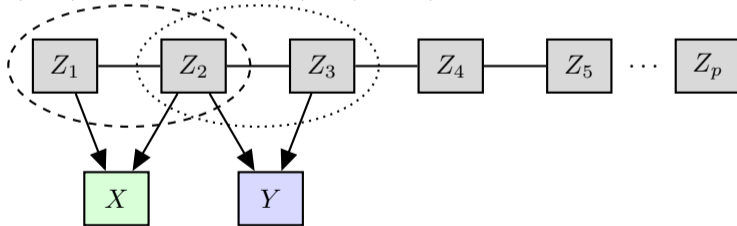
- ▶ Assume that

$$X \stackrel{\text{approximately}}{\perp\!\!\!\perp} Z \mid h(Z), \quad Y \stackrel{\text{approximately}}{\perp\!\!\!\perp} Z \mid g(Z),$$

for some low dimensional functions of g and h .

$$h(Z) = (Z_1, Z_2)$$

$$g(Z) = (Z_2, Z_3)$$



Maxway CRT

- ▶ Assume that

$$X \overset{\text{approximately}}{\perp\!\!\!\perp} Z \mid h(Z), \quad Y \overset{\text{approximately}}{\perp\!\!\!\perp} Z \mid g(Z),$$

for some low dimensional functions of g and h .

- ▶ Sparse linear model: $X = \alpha^\top Z_{\mathcal{S}_x} + \varepsilon$, where $\mathcal{S}_x \subset \{1, \dots, p\}$. $h(Z) = Z_{\mathcal{S}_x}$ or $h(Z) = \alpha^\top Z_{\mathcal{S}_x}$.
- ▶ More general setting: the distribution of X given Z is very complicated.
 - Do a transformation: $R(X, Z) = F_X(X \mid Z)$. If X has a continuous distribution conditional on Z , then $R(X, Z) \sim \text{Unif}[0, 1]$ both marginally and conditionally on Z . Thus $R(X, Z) \perp\!\!\!\perp Z$. We can take h to be a null set.
 - Test whether $R(X, Z)$ is independent of Y conditional on Z .
 - Extracting the residual of X after removing the influence of Z .

Maxway CRT

Let $\rho^*(\cdot | g(Z), h(Z))$ be the conditional distribution of X given $g(Z)$ and $h(Z)$, and let ρ be an estimate of ρ^* .

Model and adjust X with the assistance of Y (Maxway) CRT

1. Sample $\mathbf{x}^{(m)}$ from the distribution of $\rho^n(\cdot | g(\mathbf{Z}), h(\mathbf{Z}))$ independently of (\mathbf{x}, \mathbf{y}) .
2. Output Maxway CRT p -value

$$p_{\text{maxway}} = \frac{1}{M+1} \left(1 + \sum_{m=1}^M \mathbf{1} \left\{ T(\mathbf{x}^{(m)}, \mathbf{y}, g(\mathbf{Z}), h(\mathbf{Z})) \geq T(\mathbf{x}, \mathbf{y}, g(\mathbf{Z}), h(\mathbf{Z})) \right\} \right).$$

Exact Inference

Theorem

Suppose that either of the following conditions holds: (i) each $\mathbf{x}^{(m)}$ is exchangeable with \mathbf{x} given \mathbf{Z} ; (ii) $\mathbf{x}^{(m)}$ is exchangeable with \mathbf{x} given $\{h(\mathbf{Z}), g(\mathbf{Z})\}$, and $\mathbf{Z} \perp\!\!\!\perp \mathbf{y} \mid g(\mathbf{Z})$. Then the Maxway CRT p -value defined is valid, i.e., $\mathbb{P}[p_{\text{maxway}} \leq \alpha] \leq \alpha$ for any $\alpha \in [0, 1]$ under \mathcal{H}_0 .

- ▶ For (i), proof from standard CRT. $T(\mathbf{x}^{(m)}, \mathbf{y}, g(\mathbf{Z}), h(\mathbf{Z}))$ is exchangeable with $T(\mathbf{x}, \mathbf{y}, g(\mathbf{Z}), h(\mathbf{Z}))$.
- ▶ For (ii), $\mathbf{x} \perp\!\!\!\perp \mathbf{y} \mid (g(\mathbf{Z}), h(\mathbf{Z})) \Rightarrow T(\mathbf{x}^{(m)}, \mathbf{y}, g(\mathbf{Z}), h(\mathbf{Z}))$ is exchangeable with $T(\mathbf{x}, \mathbf{y}, g(\mathbf{Z}), h(\mathbf{Z}))$.

Exact Inference

Theorem

Suppose that either of the following conditions holds: (i) each $\mathbf{x}^{(m)}$ is exchangeable with \mathbf{x} given \mathbf{Z} ; (ii) $\mathbf{x}^{(m)}$ is exchangeable with \mathbf{x} given $\{h(\mathbf{Z}), g(\mathbf{Z})\}$, and $\mathbf{Z} \perp\!\!\!\perp \mathbf{y} \mid g(\mathbf{Z})$. Then the Maxway CRT p -value defined is valid, i.e., $\mathbb{P} [p_{\text{maxway}} \leq \alpha] \leq \alpha$ for any $\alpha \in [0, 1]$ under \mathcal{H}_0 .

- ▶ Compared to the conditions for the Model-X CRT p -value to be valid, the conditions stated in this theorem is strictly weaker.
- ▶ Condition (i) requires that the knowledge of the distribution of X given Z is perfect. Condition (ii) requires that $g(Z)$ contains all the information about Y that Z can possibly provide, and that the distribution of X given the low dimensional $g(Z)$ and $h(Z)$ is known.

Robustness of the Maxway CRT

- ▶ When g , h and ρ are not perfect anymore...

Theorem

For any $\alpha \in (0, 1)$,

$$\mathbb{P}[\rho_{\max\text{way}} \leq \alpha] \leq \alpha + 2\mathbb{E}[d_x(\mathbf{Z})d_y(\mathbf{Z})] + \mathbb{E}[d_\rho(g(\mathbf{Z}), h(\mathbf{Z}))],$$

where

$$d_\rho(g(\mathbf{Z}), h(\mathbf{Z})) = d_{\text{TV}}(\rho^{*n}(\cdot | g(\mathbf{Z}), h(\mathbf{Z})), \rho^n(\cdot | g(\mathbf{Z}), h(\mathbf{Z}))),$$

$$d_x(\mathbf{Z}) = d_{\text{TV}}(\rho^{*n}(\cdot | g(\mathbf{Z}), h(\mathbf{Z})), f_{x|\mathbf{Z}}(\cdot | \mathbf{Z})), \text{ and}$$

$$d_y(\mathbf{Z}) = d_{\text{TV}}(f_{y|\mathbf{Z}}(\cdot | \mathbf{Z}), f_{y|g(\mathbf{Z})}(\cdot | g(\mathbf{Z}))).$$

Robustness of the Maxway CRT

Theorem

For any $\alpha \in (0, 1)$,

$$\mathbb{P} [\rho_{\text{maxway}} \leq \alpha] \leq \alpha + 2\mathbb{E} \left[d_x(\mathbf{Z}) d_y(\mathbf{Z}) \right] + \mathbb{E} [d_\rho(g(\mathbf{Z}), h(\mathbf{Z}))],$$

where $d_x(\mathbf{Z}) = d_{\text{TV}} \left(\rho^{*n}(\cdot | g(\mathbf{Z}), h(\mathbf{Z})), f_{x|z}(\cdot | \mathbf{Z}) \right)$.

- ▶ Recall $\rho^{*n}(\cdot | g(\mathbf{Z}), h(\mathbf{Z})) = f_{x|h(Z),g(Z)}(\cdot | g(\mathbf{Z}), h(\mathbf{Z}))$.
- ▶ Captures how independent X is to Z conditional on $h(Z)$.
- ▶ When $X \perp\!\!\!\perp Z | h(Z)$, then the conditional distribution of $X | Z$ would be the same as $X | h(Z)$, thus the term is zero.

Robustness of the Maxway CRT

Theorem

For any $\alpha \in (0, 1)$,

$$\mathbb{P} [p_{\text{maxway}} \leq \alpha] \leq \alpha + 2\mathbb{E} \left[d_x(\mathbf{Z}) d_y(\mathbf{Z}) \right] + \mathbb{E} [d_\rho(g(\mathbf{Z}), h(\mathbf{Z}))],$$

where $d_y(\mathbf{Z}) = d_{\text{TV}}(f_{y|z}(\cdot | \mathbf{Z}), f_{y|g(\mathbf{Z})}(\cdot | g(\mathbf{Z})))$.

- ▶ Captures how independent Y is to Z conditional on $g(Z)$.
- ▶ When $Y \perp\!\!\!\perp Z | g(Z)$, then the conditional distribution of $Y | Z$ would be the same as $Y | g(Z)$, thus the term is zero.

Robustness of the Maxway CRT

Theorem

For any $\alpha \in (0, 1)$,

$$\mathbb{P} [\rho_{\text{maxway}} \leq \alpha] \leq \alpha + 2\mathbb{E} [d_x(\mathbf{Z})d_y(\mathbf{Z})] + \mathbb{E} [d_\rho(g(\mathbf{Z}), h(\mathbf{Z}))],$$

where $d_\rho(g(\mathbf{Z}), h(\mathbf{Z})) = d_{\text{TV}}(\rho^{*n}(\cdot | g(\mathbf{Z}), h(\mathbf{Z})), \rho^n(\cdot | g(\mathbf{Z}), h(\mathbf{Z})))$.

- The d_ρ term is about the accuracy of ρ , i.e., how accurate we can estimate the distribution of X given low dimensional objects $h(Z)$ and $g(Z)$.

Compared to the Model-X CRT

- ▶ Model-X CRT

$$\mathbb{P} [p_{\text{mx}} \leq \alpha] \leq \alpha + d'_x \approx \alpha + d_\rho + d_x.$$

- ▶ Maxway CRT

$$\mathbb{P} [p_{\text{maxway}} \leq \alpha] \leq \alpha + d_\rho + 2d_x d_y.$$

- ▶ For maxway CRT, the test statistic can only be a function of $\mathbf{x}, \mathbf{y}, g(\mathbf{Z}), h(\mathbf{Z})$. Not the most general form. But has a computational advantage (Liu et al., 2020), and it is typically powerful (Katsevich and Ramdas, 2020)

Examples

Example (Gaussian linear example, estimation)

$Y_i = Z_i^T \alpha^* + \varepsilon_i$, $X_i = Z_i^T \beta^* + \eta_i$. Take $g(Z)$ to be an estimate of the mean function $Z_i^T \alpha^*$. Take $h(Z)$ to be an estimate of the mean function $Z_i^T \beta^*$. Estimate the parameters with lasso.

$$\text{Type-I error inflation of the Maxway CRT} \lesssim \sqrt{\frac{n}{N_p}} + \sqrt{\frac{s_\alpha \log(p)n}{N_y}} \sqrt{\frac{s_\beta \log(p)n}{N_x}}.$$

$$\text{Type-I error inflation of the Model-X CRT} \lesssim \sqrt{\frac{s_\beta \log(p)n}{N_x}}.$$

Examples

Example (Gaussian linear example, variable selection)

$Y_i = Z_i^\top \alpha^* + \varepsilon_i$, $X_i = Z_i^\top \beta^* + \eta_i$. Take $g(Z)$ to be an estimate of $Z_{\mathcal{S}^*}$, where \mathcal{S}^* is the support of α^* . Take $h(Z)$ to be an estimate of the mean function $Z_i^\top \beta^*$. Estimate the parameters/support with lasso.

$$\text{Type-I error inflation of the Maxway CRT} \lesssim \sqrt{\frac{ns_\alpha}{N_\rho}} + \delta \sqrt{\frac{s_\beta \log(p)n}{N_x}},$$

where $1 - \delta$ is the probability of exact recovery of the support of α^* .

$$\text{Type-I error inflation of the Model-X CRT} \lesssim \sqrt{\frac{s_\beta \log(p)n}{N_x}}.$$

Examples

Example (Binary X , smooth mean functions)

$Y_i = g^*(Z_{i\cdot}) + \varepsilon_i$, where $\varepsilon_i \sim \mathcal{N}(0, \sigma_y^2)$ independently and $X_i \sim \text{Bern}(h^*(Z_{i\cdot}))$ independently. In the Maxway CRT, we take $\hat{g}(Z_{i\cdot})$ to be an estimate of $g^*(Z_{i\cdot})$ and take $\hat{h}(Z_{i\cdot})$ to be an estimate of $h^*(Z_{i\cdot})$. Assume g^* is α -smooth, h^* is β -smooth. Estimate the mean functions with a kernel method.

$$\text{Type-I error inflation of the Maxway CRT} \lesssim \sqrt{n} N_\rho^{-\frac{\gamma}{2\gamma+2}} + n N_y^{-\frac{\alpha}{2\alpha+p}} N_x^{-\frac{\beta}{2\beta+p}}.$$

$$\text{Type-I error inflation of the Model-X CRT} \lesssim \sqrt{n} N_x^{-\frac{\beta}{2\beta+p}}.$$

Semi-supervised Scenario

- ▶ Labelled data $\mathbf{D} = (\mathbf{y}, \mathbf{x}, \mathbf{Z})$ of sample size n
- ▶ Unlabelled data $\mathbf{D}^u = (\mathbf{x}^u, \mathbf{Z}^u) \rightarrow$ train h and ρ
- ▶ How to train g ? What if we don't have an external dataset of (\mathbf{y}, \mathbf{Z}) ?
 - Train g on \mathbf{D} .
 - Theory cannot be directly applied. Still hope to avoid overfitting.
 - "Cross-fitting". Divide the data \mathbf{D} into K fold. Train g on the $K - 1$ folds and evaluate $g(\mathbf{Z}_i)$ on the other fold.

Simulations

Gaussian Linear Model

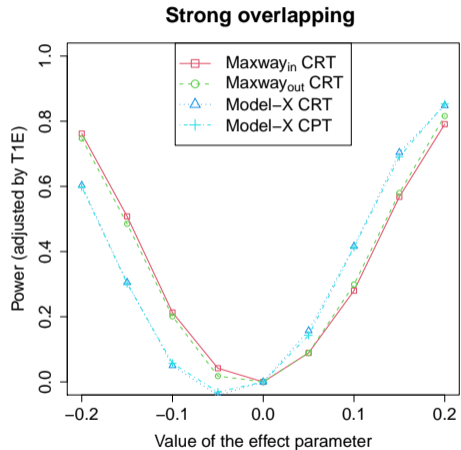
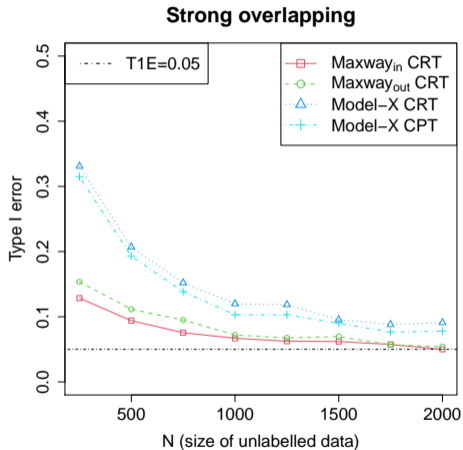
Generate $Z \in \mathbb{R}^p$ from $N(\mathbf{0}, \Sigma)$ where $p = 500$ and $\Sigma = (\sigma_{ij})_{p \times p}$ with $\sigma_{ij} = 0.5^{|i-j|}$. Then generate the conditional gaussian X and Y following:

$$X = 0.3 \sum_{j=1}^5 \nu_j Z_j + \eta \sum_{\ell \in \mathcal{I}_1} \nu_\ell Z_\ell + \epsilon_1; \quad Y = \gamma h(X, Z) + 0.3 \sum_{j=1}^5 \nu_j Z_j + \eta \sum_{\ell \in \mathcal{I}_2} \nu_\ell Z_\ell + \epsilon_2,$$

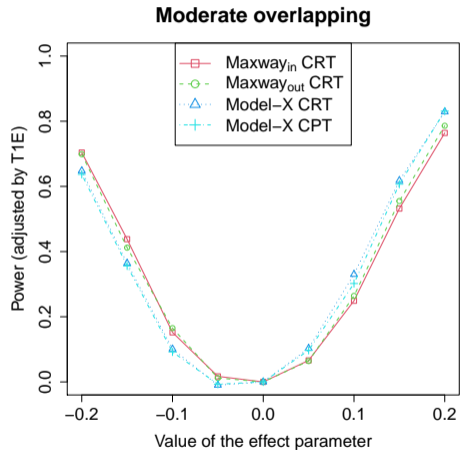
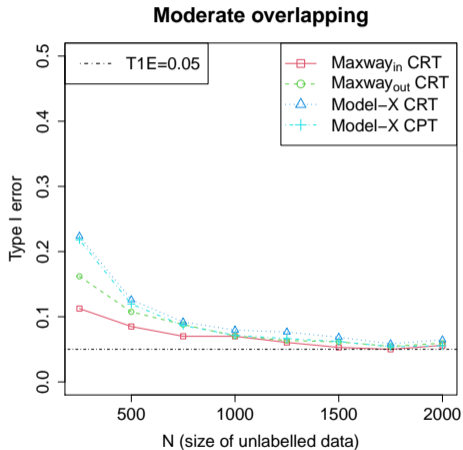
where $\epsilon_1, \epsilon_2 \sim \mathcal{N}(0, 1)$, each ν_j is randomly picked from $\{-1, 1\}$, and $\mathcal{I}_1, \mathcal{I}_2$ are two disjoint sets of indices randomly drawn from $\{6, 7, \dots, p\}$ satisfying $|\mathcal{I}_1| = |\mathcal{I}_2| = 25$.

- ▶ η : how strong is the confounding? $\eta = 0, 0.1, 0.2$: strong, moderate, weak overlapping/confounding.
- ▶ $\gamma = 0$: evaluate type I error.
- ▶ $\gamma \neq 0$: power curve.

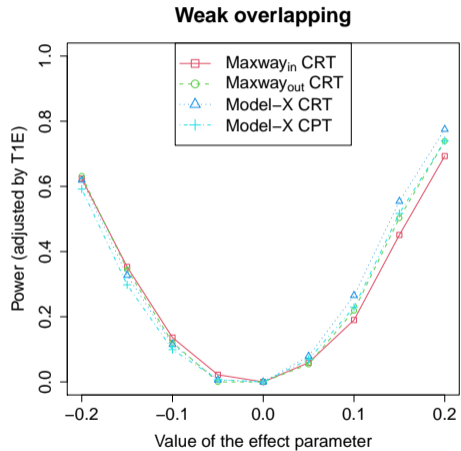
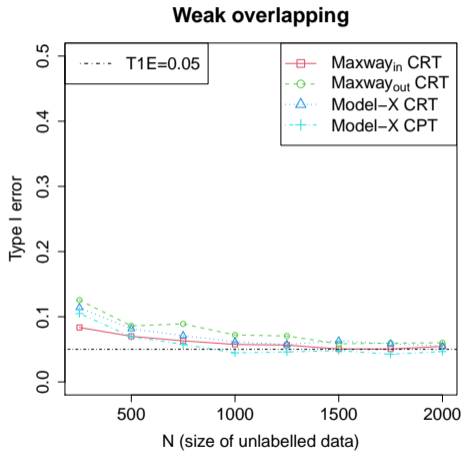
Simulations



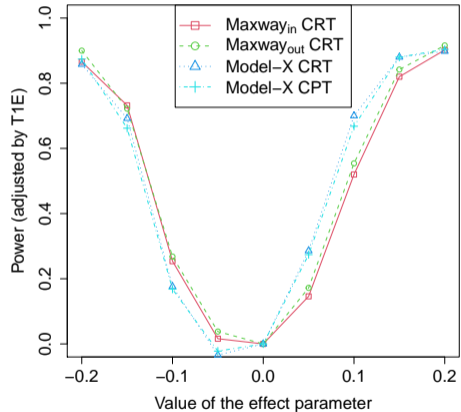
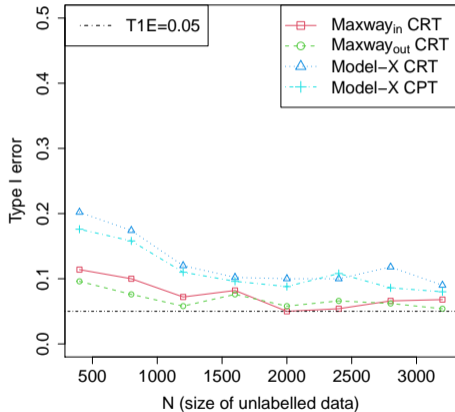
Simulations



Simulations



Simulations: nonlinear model



Real Data Example

- ▶ Statins are one of the most commonly used drug in the United States for lowering the level of low-density lipoprotein (LDL) and the risk of cardiovascular disease (CVD).
- ▶ Working mechanism: HMGCR inhibition
- ▶ Evidences showing that the use of statins could increase the risk for type II diabetes mellitus (DM).



- ▶ Unethical/expensive to conduct randomized control trial.
- ▶ Statins = absence of certain SNP in HMGCR.
- ▶ Test whether SNP $\perp\!\!\!\perp$ diabetes | other variables—gives a biological perspective

Real Data Example

- ▶ UK Biobank
- ▶ Z includes age, gender and genetic variants associated with DM or its related phenotypes including high LDL, high-density lipoprotein (HDL) and BMI.
- ▶ p -values

Statistic	CRT	CPT	Maxway CRT
d_0	0.06	0.06	0.04
d_I	0.16	0.18	0.13

Table: The d_0 and d_I p -values for the dependence of the risk of DM on the treatment of statins functionally represented by the variant rs17238484-G.

- ▶ The Maxway CRT is not generally more conservative than the original Model-X CRT.

Thank you!