

Efficient Inference on High-Dimensional Linear Models With Missing Outcomes

Yikun Zhang

Joint Work with *Alexander Giessing* and *Yen-Chi Chen*

Department of Statistics,
University of Washington

November 8, 2023 at Casual Inference and Missing Data Reading Group

- 1 Introduction
- 2 Methodology: Efficient Debiasing Method
- 3 Theory: Consistency and Asymptotic Normality
- 4 Comparative Simulations
- 5 Real-World Applications: Stellar Mass Inference Problem
- 6 Conclusions and Future Works

Introduction



Consider a random sample $\{(Y_i, R_i, X_i)\}_{i=1}^n$ drawn from the joint distribution of (Y, R, X) , where

- $Y \in \mathbb{R}$ is the outcome variable that could potentially be missing;
- $R \in \{0, 1\}$ is the indicator of Y being observed;
- $X \in \mathbb{R}^d$ is the high-dimensional covariate vector with $d \gg n$.

Consider a random sample $\{(Y_i, R_i, X_i)\}_{i=1}^n$ drawn from the joint distribution of (Y, R, X) , where

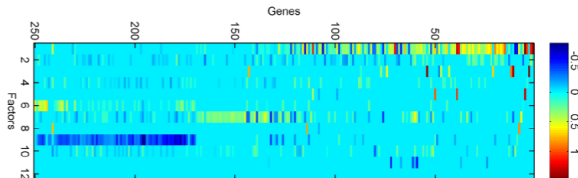
- $Y \in \mathbb{R}$ is the outcome variable that could potentially be missing;
- $R \in \{0, 1\}$ is the indicator of Y being observed;
- $X \in \mathbb{R}^d$ is the high-dimensional covariate vector with $d \gg n$.

► **Central Question of Interest:**

How can we conduct statistically and computationally efficient inference on $m_0(x) = E(Y|X = x)$ despite missing outcomes?

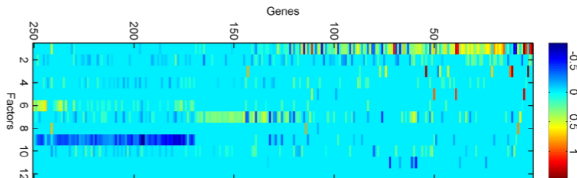
- 1 The covariates are easier to obtain within some population.

- ① The covariates are easier to obtain within some population.



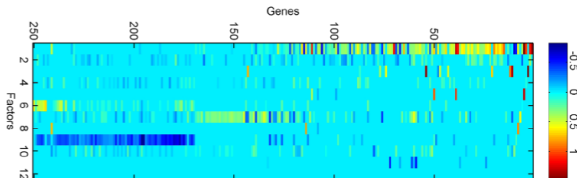
- Micro-array gene expression data in biology ([Carvalho et al., 2008](#)).
- Home-price data with cross-sectional effects ([Fan et al., 2011](#)).

- 1 The covariates are easier to obtain within some population.



- Micro-array gene expression data in biology ([Carvalho et al., 2008](#)).
 - Home-price data with cross-sectional effects ([Fan et al., 2011](#)).
- 2 Incorporating as many covariates as possible can control for potential confounders in causal inference ([Wyss et al., 2022](#)).

- ① The covariates are easier to obtain within some population.



- Micro-array gene expression data in biology ([Carvalho et al., 2008](#)).
 - Home-price data with cross-sectional effects ([Fan et al., 2011](#)).
- ② Incorporating as many covariates as possible can control for potential confounders in causal inference ([Wyss et al., 2022](#)).
- ③ Generating high-dimensional covariates with interaction terms or spline features enables the simple parametric (e.g., linear) model to capture complex patterns ([Belloni et al., 2019](#)).

The response/outcome variable in observational data could be missing.

The response/outcome variable in observational data could be missing.

- 1 Participants may drop out from the study in clinical trials ([Higgins et al., 2008](#)).

The response/outcome variable in observational data could be missing.

- 1 Participants may drop out from the study in clinical trials ([Higgins et al., 2008](#)).
- 2 The semi-supervised learning, where additional samples without labels are provided, is a missing-outcome problem ([Chapelle et al., 2006](#)).

The response/outcome variable in observational data could be missing.

- 1 Participants may drop out from the study in clinical trials ([Higgins et al., 2008](#)).
 - 2 The semi-supervised learning, where additional samples without labels are provided, is a missing-outcome problem ([Chapelle et al., 2006](#)).
- **More Concrete Example:** Some (estimated) stellar masses of the observed galaxies in the Sloan Digital Sky Survey (SDSS-IV) are missing in the Firefly value-added catalog ([Comparat et al., 2017](#)).

W Motivations: Stellar Mass Inference Problem

The missingness of (estimated) stellar masses is due to

- Limiting usage of the observational run in SDSS-IV for galaxy targets;
- Potential data contamination;
- Misclassification of galaxies as stars.

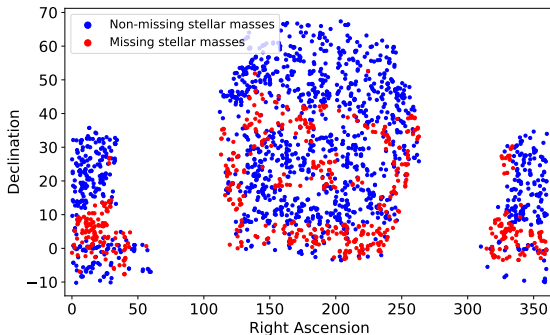


Figure 1: Galaxy distribution at a high redshift slice $0.4 \sim 0.401$.

W Motivations: Stellar Mass Inference Problem

The missingness of (estimated) stellar masses is due to

- Limiting usage of the observational run in SDSS-IV for galaxy targets;
- Potential data contamination;
- Misclassification of galaxies as stars.

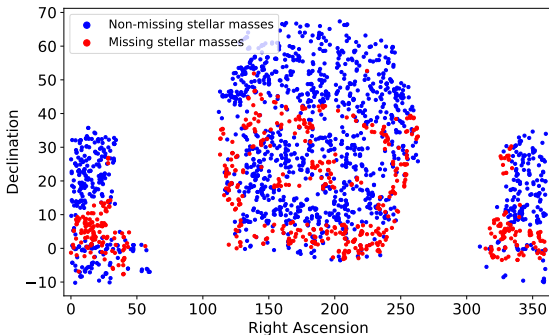


Figure 1: Galaxy distribution at a high redshift slice $0.4 \sim 0.401$.

► **Scientific Question:** *How can we conduct valid inference on the (estimated) stellar mass based on the spectroscopic and photometric properties?*

To tackle the challenges of high-dimensional data with missing outcomes, we impose two basic assumptions.

To tackle the challenges of high-dimensional data with missing outcomes, we impose two basic assumptions.

- ① (*Linearity*) The data $\{(Y_i, R_i, X_i)\}_{i=1}^n \subset \mathbb{R} \times \{0, 1\} \times \mathbb{R}^d$ are i.i.d. observations from a sparse linear model

$$Y = X^T \beta_0 + \epsilon \quad \text{with} \quad \mathbb{E}(\epsilon|X) = 0 \quad \text{and} \quad \mathbb{E}(\epsilon^2|X) = \sigma_\epsilon^2,$$

where $\|\beta_0\|_0 = \sum_{k=1}^d \mathbb{1}_{\{\beta_{0k} \neq 0\}} = s_\beta \ll d$.

To tackle the challenges of high-dimensional data with missing outcomes, we impose two basic assumptions.

- ① (*Linearity*) The data $\{(Y_i, R_i, X_i)\}_{i=1}^n \subset \mathbb{R} \times \{0, 1\} \times \mathbb{R}^d$ are i.i.d. observations from a sparse linear model

$$Y = X^T \beta_0 + \epsilon \quad \text{with} \quad \mathbb{E}(\epsilon|X) = 0 \quad \text{and} \quad \mathbb{E}(\epsilon^2|X) = \sigma_\epsilon^2,$$

where $\|\beta_0\|_0 = \sum_{k=1}^d \mathbb{1}_{\{\beta_{0k} \neq 0\}} = s_\beta \ll d$.

► **Notes:** The linearity assumption can be relaxed to

- Sparse additive model ([Ravikumar et al., 2009](#));
- Partially linear model ([Müller and van de Geer, 2015](#));
- Approximately/weakly sparse linear model ([Belloni et al., 2019](#)).

To tackle the challenges of high-dimensional data with missing outcomes, we impose two basic assumptions.

- ① (*Linearity*) The data $\{(Y_i, R_i, X_i)\}_{i=1}^n \subset \mathbb{R} \times \{0, 1\} \times \mathbb{R}^d$ are i.i.d. observations from a sparse linear model

$$Y = X^T \beta_0 + \epsilon \quad \text{with} \quad E(\epsilon|X) = 0 \quad \text{and} \quad E(\epsilon^2|X) = \sigma_\epsilon^2,$$

where $\|\beta_0\|_0 = \sum_{k=1}^d \mathbb{1}_{\{\beta_{0k} \neq 0\}} = s_\beta \ll d$.

► **Notes:** The linearity assumption can be relaxed to

- Sparse additive model (Ravikumar et al., 2009);
- Partially linear model (Müller and van de Geer, 2015);
- Approximately/weakly sparse linear model (Belloni et al., 2019).

- ② (*Missing At Random; MAR*) $Y_i \perp\!\!\!\perp R_i | X_i$ for $i = 1, \dots, n$.

The existing works focus mainly on the statistical inference on $\beta_0 \in \mathbb{R}^d$.

The existing works focus mainly on the statistical inference on $\beta_0 \in \mathbb{R}^d$.

- ① (*Fully observed outcomes*) Debiased Lasso is applicable (Zhang and Zhang, 2014; van de Geer et al., 2014; Javanmard and Montanari, 2014):

$$\hat{\beta}^{\text{debias}} = \hat{\beta}_\lambda + \frac{1}{n} \hat{\Theta} \sum_{i=1}^n X_i (Y_i - X_i^T \hat{\beta}_\lambda),$$

- $\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right]$ is a Lasso solution with the regularization parameter $\lambda > 0$;
- $\hat{\Theta} \in \mathbb{R}^{d \times d}$ is an approximation to the matrix inverse $\left(\frac{1}{n} \sum_{i=1}^n X_i X_i^T \right)^{-1}$.

The existing works focus mainly on the statistical inference on $\beta_0 \in \mathbb{R}^d$.

- ① (*Fully observed outcomes*) Debiased Lasso is applicable ([Zhang and Zhang, 2014](#); [van de Geer et al., 2014](#); [Javanmard and Montanari, 2014](#)):

$$\hat{\beta}^{\text{debias}} = \hat{\beta}_\lambda + \frac{1}{n} \hat{\Theta} \sum_{i=1}^n X_i (Y_i - X_i^T \hat{\beta}_\lambda),$$

- $\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right]$ is a Lasso solution with the regularization parameter $\lambda > 0$;
 - $\hat{\Theta} \in \mathbb{R}^{d \times d}$ is an approximation to the matrix inverse $\left(\frac{1}{n} \sum_{i=1}^n X_i X_i^T \right)^{-1}$.
- ② (*MAR outcomes*) [Chakraborty et al. \(2019\)](#) proposed an M-estimation framework with a Lasso-type debiased and doubly robust estimator.

► **Drawbacks of the Existing Approaches:**

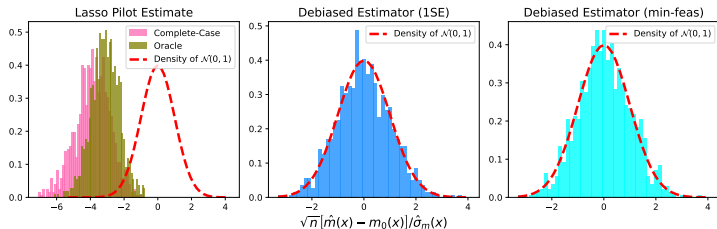
- ① (*Computational issue*) They require a good approximation to the $d \times d$ debiasing matrix $\hat{\Theta}$.
- ② (*Loss of statistical efficiency*) Sample splitting or cross-fitting is necessary for the M-estimation framework.

► **Drawbacks of the Existing Approaches:**

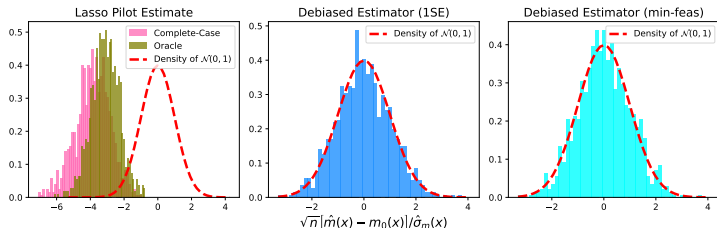
- ① (*Computational issue*) They require a good approximation to the $d \times d$ debiasing matrix $\hat{\Theta}$.
- ② (*Loss of statistical efficiency*) Sample splitting or cross-fitting is necessary for the M-estimation framework.

► **Our Contributions:** Focus on the inference of $m_0(x) = x^T \beta_0$ instead.

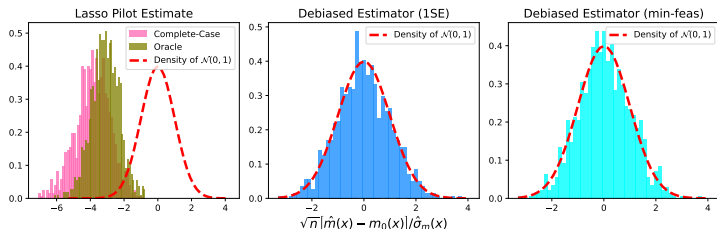
- (*Computational efficiency*) Our core debiasing program is convex and only needs to solve for a n -dimensional weight vector.
- (*Statistical efficiency*) Our debiased estimator is semi-parametrically efficient among all asymptotically linear estimators.



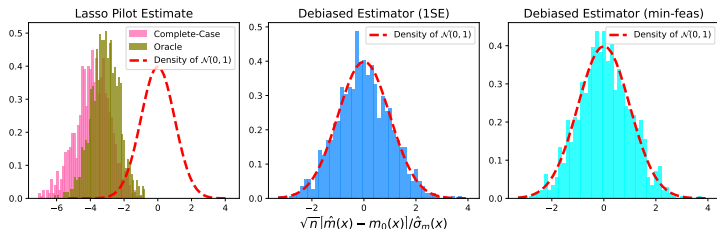
- 1 Introduce our efficient debiasing method for inferring $m_0(x) = x^T \beta_0$.



- 1 Introduce our efficient debiasing method for inferring $m_0(x) = x^T \beta_0$.
 - Estimate $\pi(X) = P(R = 1|X)$ via any machine learning methods.
 - Design our debiasing program based on bias-variance trade-offs.
 - Fine-tune the program from its dual so as to debias the Lasso solution.



- 1 Introduce our efficient debiasing method for inferring $m_0(x) = x^T \beta_0$.
 - Estimate $\pi(X) = P(R = 1|X)$ via any machine learning methods.
 - Design our debiasing program based on bias-variance trade-offs.
 - Fine-tune the program from its dual so as to debias the Lasso solution.
- 2 Discuss the asymptotic normality and semi-parametric efficiency of our final debiased estimator.



- 1 Introduce our efficient debiasing method for inferring $m_0(x) = x^T \beta_0$.
 - Estimate $\pi(X) = P(R = 1|X)$ via any machine learning methods.
 - Design our debiasing program based on bias-variance trade-offs.
 - Fine-tune the program from its dual so as to debias the Lasso solution.
- 2 Discuss the asymptotic normality and semi-parametric efficiency of our final debiased estimator.
- 3 Demonstrate the finite-sample performances via simulations and present an application to the stellar mass inference problem.

Methodology



For any fixed $\lambda > 0$, the Lasso solution (on the complete-case data) is a biased estimator of $\beta_0 \in \mathbb{R}^d$:

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

► **Question:** How can we correct for the bias in $\hat{\beta}_\lambda$ or $\hat{m}(x) = x^T \hat{\beta}_\lambda$?

For any fixed $\lambda > 0$, the Lasso solution (on the complete-case data) is a biased estimator of $\beta_0 \in \mathbb{R}^d$:

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

► **Question:** How can we correct for the bias in $\hat{\beta}_\lambda$ or $\hat{m}(x) = x^T \hat{\beta}_\lambda$?

- Optimality/KKT condition reads

$$\frac{1}{n} \sum_{i=1}^n R_i X_i \left(Y_i - X_i^T \hat{\beta}_\lambda \right) = \lambda \hat{z} \quad \text{with} \quad \hat{z} \in \partial \left\| \hat{\beta}_\lambda \right\|_1 \in \mathbb{R}^d. \quad (1)$$

For any fixed $\lambda > 0$, the Lasso solution (on the complete-case data) is a biased estimator of $\beta_0 \in \mathbb{R}^d$:

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

► **Question:** How can we correct for the bias in $\hat{\beta}_\lambda$ or $\hat{m}(x) = x^T \hat{\beta}_\lambda$?

- Optimality/KKT condition reads

$$\frac{1}{n} \sum_{i=1}^n R_i X_i \left(Y_i - X_i^T \hat{\beta}_\lambda \right) = \lambda \hat{z} \quad \text{with} \quad \hat{z} \in \partial \left\| \hat{\beta}_\lambda \right\|_1 \in \mathbb{R}^d. \quad (1)$$

- Linearity assumption $Y_i = X_i^T \beta_0 + \epsilon_i$ for $i = 1, \dots, n$ implies that

$$\frac{1}{n} \sum_{i=1}^n R_i X_i \epsilon_i + \hat{\Sigma} \left(\beta_0 - \hat{\beta}_\lambda \right) = \lambda \hat{z} \quad \text{with} \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n R_i X_i X_i^T.$$

- Given an approximation $\hat{\Theta} \in \mathbb{R}^{d \times d}$ to $\hat{\Sigma}^{-1}$, it becomes

$$\hat{\beta}_\lambda - \beta_0 + \hat{\Theta} \lambda \hat{z} = \underbrace{\frac{1}{n} \sum_{i=1}^n R_i \hat{\Theta} X_i \epsilon_i}_{\text{Stochastic error } \sim \mathcal{N}_d(0, \tilde{\Sigma})} + \underbrace{\left(\hat{\Theta} \hat{\Sigma} - I_d \right) \left(\beta_0 - \hat{\beta}_\lambda \right)}_{\text{Asymptotically negligible bias}}.$$

- Given an approximation $\hat{\Theta} \in \mathbb{R}^{d \times d}$ to $\hat{\Sigma}^{-1}$, it becomes

$$\hat{\beta}_\lambda - \beta_0 + \hat{\Theta} \lambda \hat{z} = \underbrace{\frac{1}{n} \sum_{i=1}^n R_i \hat{\Theta} X_i \epsilon_i}_{\text{Stochastic error } \sim \mathcal{N}_d(0, \tilde{\Sigma})} + \underbrace{(\hat{\Theta} \hat{\Sigma} - I_d) (\beta_0 - \hat{\beta}_\lambda)}_{\text{Asymptotically negligible bias}}.$$

- By KKT condition (1), the debiased Lasso estimate is thus given by

$$\begin{aligned} \hat{\beta}^{\text{debias}} &= \hat{\beta}_\lambda + \hat{\Theta} \lambda \hat{z} \\ &= \hat{\beta}_\lambda + \frac{1}{n} \sum_{i=1}^n R_i \hat{\Theta} X_i (Y_i - X_i^T \hat{\beta}_\lambda). \end{aligned}$$

- Given an approximation $\hat{\Theta} \in \mathbb{R}^{d \times d}$ to $\hat{\Sigma}^{-1}$, it becomes

$$\hat{\beta}_\lambda - \beta_0 + \hat{\Theta} \lambda \hat{z} = \underbrace{\frac{1}{n} \sum_{i=1}^n R_i \hat{\Theta} X_i \epsilon_i}_{\text{Stochastic error } \sim \mathcal{N}_d(0, \tilde{\Sigma})} + \underbrace{\left(\hat{\Theta} \hat{\Sigma} - I_d \right) \left(\beta_0 - \hat{\beta}_\lambda \right)}_{\text{Asymptotically negligible bias}}.$$

- By KKT condition (1), the debiased Lasso estimate is thus given by

$$\begin{aligned} \hat{\beta}^{\text{debias}} &= \hat{\beta}_\lambda + \hat{\Theta} \lambda \hat{z} \\ &= \hat{\beta}_\lambda + \frac{1}{n} \sum_{i=1}^n R_i \hat{\Theta} X_i \left(Y_i - X_i^T \hat{\beta}_\lambda \right). \end{aligned}$$

- A candidate debiased estimator for $m_0(x) = x^T \beta_0$ is

$$\hat{m}^{\text{debias}}(x) = x^T \hat{\beta}^{\text{debias}} = x^T \hat{\beta}_\lambda + \frac{1}{n} x^T \hat{\Theta} \sum_{i=1}^n R_i X_i \left(Y_i - X_i^T \hat{\beta}_\lambda \right).$$

$$\hat{m}^{\text{debias}}(x) = x^T \hat{\beta}^{\text{debias}} = x^T \hat{\beta}_\lambda + \frac{1}{n} x^T \hat{\Theta} \sum_{i=1}^n R_i X_i \left(Y_i - X_i^T \hat{\beta}_\lambda \right).$$

► **Issue:** Fitting the debiasing matrix $\hat{\Theta} \in \mathbb{R}^{d \times d}$ is computationally inefficient; see, e.g., the nodewise regression ([Meinshausen and Bühlmann, 2006](#); [van de Geer et al., 2014](#)).

$$\hat{m}^{\text{debias}}(x) = x^T \hat{\beta}^{\text{debias}} = x^T \hat{\beta}_\lambda + \frac{1}{n} x^T \hat{\Theta} \sum_{i=1}^n R_i X_i \left(Y_i - X_i^T \hat{\beta}_\lambda \right).$$

► **Issue:** Fitting the debiasing matrix $\hat{\Theta} \in \mathbb{R}^{d \times d}$ is computationally inefficient; see, e.g., the nodewise regression (Meinshausen and Bühlmann, 2006; van de Geer et al., 2014).

► **Solution:** Introduce the weight vector $\hat{w} = (\hat{w}_1, \dots, \hat{w}_n)^T \in \mathbb{R}^n$ with (Giessing and Wang, 2023)

$$\hat{w}_i = \begin{cases} \frac{1}{\sqrt{n}} x^T \hat{\Theta} X_i & R_i = 1, \\ 0 & R_i = 0, \end{cases}$$

for $i = 1, \dots, n$ so that our final debiased estimator becomes

$$\hat{m}^{\text{debias}}(x; \hat{w}) = x^T \hat{\beta} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i \left(Y_i - X_i^T \hat{\beta} \right). \quad (2)$$

$$\hat{m}^{\text{debias}}(x) = x^T \hat{\beta}^{\text{debias}} = x^T \hat{\beta}_\lambda + \frac{1}{n} x^T \hat{\Theta} \sum_{i=1}^n R_i X_i \left(Y_i - X_i^T \hat{\beta}_\lambda \right).$$

► **Issue:** Fitting the debiasing matrix $\hat{\Theta} \in \mathbb{R}^{d \times d}$ is computationally inefficient; see, e.g., the nodewise regression (Meinshausen and Bühlmann, 2006; van de Geer et al., 2014).

► **Solution:** Introduce the weight vector $\hat{w} = (\hat{w}_1, \dots, \hat{w}_n)^T \in \mathbb{R}^n$ with (Giessing and Wang, 2023)

$$\hat{w}_i = \begin{cases} \frac{1}{\sqrt{n}} x^T \hat{\Theta} X_i & R_i = 1, \\ 0 & R_i = 0, \end{cases}$$

for $i = 1, \dots, n$ so that our final debiased estimator becomes

$$\hat{m}^{\text{debias}}(x; \hat{w}) = x^T \hat{\beta} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i \left(Y_i - X_i^T \hat{\beta} \right). \quad (2)$$

► **Question:** How do we estimate the weight vector $\hat{w} = (\hat{w}_1, \dots, \hat{w}_n)^T$?

Consider the generic debiased estimator $m^{\text{debias}}(x; \mathbf{w})$ from (2) as:

$$m^{\text{debias}}(x; \mathbf{w}) = x^T \beta + \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i R_i (Y_i - X_i^T \beta). \quad (3)$$

Consider the generic debiased estimator $m^{\text{debias}}(x; \mathbf{w})$ from (2) as:

$$m^{\text{debias}}(x; \mathbf{w}) = x^T \beta + \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i R_i (Y_i - X_i^T \beta). \quad (3)$$

The conditional mean squared error of $\sqrt{n} m^{\text{debias}}(x; \mathbf{w})$ is given by

$$\begin{aligned} & \mathbb{E} \left[\left(\sqrt{n} m^{\text{debias}}(x; \mathbf{w}) - \sqrt{n} m_0(x) \right)^2 \middle| X_1, \dots, X_n \right] \\ &= \underbrace{\sigma_\epsilon^2 \sum_{i=1}^n w_i^2 \pi(X_i)}_{\text{Main Conditional Variance}} + \underbrace{\left[\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \pi(X_i) X_i - x \right)^T \sqrt{n} (\beta_0 - \beta) \right]^2}_{\text{Conditional Bias}} \\ &+ \underbrace{(\beta_0 - \beta)^T \left[\sum_{i=1}^n w_i^2 \pi(X_i) (1 - \pi(X_i)) X_i X_i^T \right] (\beta_0 - \beta)}_{\text{Asymptotically Negligible Conditional Variance}}, \end{aligned}$$

where $\pi(X) = P(R = 1|X)$ is the propensity score under MAR condition.

$$\begin{aligned}
& \mathbb{E} \left[\left(\sqrt{n} m^{\text{debias}}(x; \mathbf{w}) - \sqrt{n} m_0(x) \right)^2 \middle| X_1, \dots, X_n \right] \\
& \asymp \underbrace{\sigma_\epsilon^2 \sum_{i=1}^n w_i^2 \pi(X_i)}_{\text{Main Conditional Variance}} + \underbrace{\left[\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \pi(X_i) X_i - x \right)^T \sqrt{n} (\beta_0 - \beta) \right]^2}_{\text{Conditional Bias}}.
\end{aligned}$$

- By Hölder's inequality, the "Conditional Bias" is upper bounded by

$$\left[\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \pi(X_i) X_i - x \right\|_\infty \sqrt{n} \|\beta_0 - \beta\|_1 \right]^2.$$

$$\begin{aligned}
& \mathbb{E} \left[\left(\sqrt{n} m^{\text{debias}}(x; \mathbf{w}) - \sqrt{n} m_0(x) \right)^2 \middle| X_1, \dots, X_n \right] \\
& \asymp \underbrace{\sigma_\epsilon^2 \sum_{i=1}^n w_i^2 \pi(X_i)}_{\text{Main Conditional Variance}} + \underbrace{\left[\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \pi(X_i) X_i - x \right)^T \sqrt{n} (\beta_0 - \beta) \right]^2}_{\text{Conditional Bias}}.
\end{aligned}$$

- By Hölder's inequality, the “Conditional Bias” is upper bounded by

$$\left[\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \pi(X_i) X_i - x \right\|_\infty \sqrt{n} \|\beta_0 - \beta\|_1 \right]^2.$$

- We design our core debiasing program as:

$$\min_{\mathbf{w} \in \mathbb{R}^n} \sum_{i=1}^n \hat{\pi}_i w_i^2 \quad \text{subject to} \quad \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_\infty \leq \frac{\gamma}{n},$$

where $\gamma > 0$ is a tuning parameter and $\hat{\pi}_i$ is a consistent estimate of the propensity score $\pi(X_i)$ for $i = 1, \dots, n$.

- 1 Compute the Lasso pilot estimate $\hat{\beta}_\lambda$ on the complete-case data

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

- 1 Compute the Lasso pilot estimate $\hat{\beta}_\lambda$ on the complete-case data

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

- 2 Obtain consistent propensity score estimates $\hat{\pi}_i, i = 1, \dots, n$ by *any machine learning method* based on $\{(X_i, R_i)\}_{i=1}^n \subset \mathbb{R}^d \times \{0, 1\}$.

- 1 Compute the Lasso pilot estimate $\hat{\beta}_\lambda$ on the complete-case data

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

- 2 Obtain consistent propensity score estimates $\hat{\pi}_i, i = 1, \dots, n$ by *any machine learning method* based on $\{(X_i, R_i)\}_{i=1}^n \subset \mathbb{R}^d \times \{0, 1\}$.
- 3 Solve the debiasing program defined as:

$$\min_{w \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \hat{\pi}_i w_i^2 : \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_\infty \leq \frac{\gamma}{n} \right\}.$$

- ① Compute the Lasso pilot estimate $\hat{\beta}_\lambda$ on the complete-case data

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

- ② Obtain consistent propensity score estimates $\hat{\pi}_i, i = 1, \dots, n$ by *any machine learning method* based on $\{(X_i, R_i)\}_{i=1}^n \subset \mathbb{R}^d \times \{0, 1\}$.
- ③ Solve the debiasing program defined as:

$$\min_{w \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \hat{\pi}_i w_i^2 : \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_\infty \leq \frac{\gamma}{n} \right\}.$$

- ④ Define the debiased estimator for $m_0(x)$ as:

$$\hat{m}^{\text{debias}}(x; \hat{w}) = x^T \hat{\beta} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i (Y_i - X_i^T \hat{\beta}).$$

- 1 Compute the Lasso pilot estimate $\hat{\beta}_\lambda$ on the complete-case data

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \left[\frac{1}{2n} \sum_{i=1}^n R_i (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_1 \right].$$

- 2 Obtain consistent propensity score estimates $\hat{\pi}_i, i = 1, \dots, n$ by *any machine learning method* based on $\{(X_i, R_i)\}_{i=1}^n \subset \mathbb{R}^d \times \{0, 1\}$.
- 3 Solve the debiasing program defined as:

$$\min_{w \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \hat{\pi}_i w_i^2 : \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_\infty \leq \frac{\gamma}{n} \right\}.$$

- 4 Define the debiased estimator for $m_0(x)$ as:

$$\hat{m}^{\text{debias}}(x; \hat{w}) = x^T \hat{\beta} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i (Y_i - X_i^T \hat{\beta}).$$

- 5 Construct the asymptotic $(1 - \tau)$ -level confidence interval for $m_0(x)$ as:

$$\left[\hat{m}^{\text{debias}}(x; \hat{w}) \pm \Phi^{-1} \left(1 - \frac{\tau}{2} \right) \cdot \hat{\sigma}_\epsilon \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\pi}_i \hat{w}_i^2} \right] \quad \text{with } \Phi(\cdot) \text{ being the CDF of } \mathcal{N}(0, 1).$$

There are two unanswered questions in our proposed debiasing inference procedure:

- 1 How can we select the tuning parameter $\gamma > 0$ for our debiasing program?

$$\min_{w \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \hat{\pi}_i w_i^2 : \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \leq \frac{\gamma}{n} \right\}.$$

- 2 Why is the asymptotic $(1 - \tau)$ -level confidence interval for $m_0(x)$ valid?

$$\left[\hat{m}^{\text{debias}}(x; \hat{w}) \pm \Phi^{-1} \left(1 - \frac{\tau}{2} \right) \cdot \hat{\sigma}_{\epsilon} \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\pi}_i \hat{w}_i^2} \right] \quad \text{with } \Phi(\cdot) \text{ being the CDF of } \mathcal{N}(0, 1).$$

There are two unanswered questions in our proposed debiasing inference procedure:

- 1 How can we select the tuning parameter $\gamma > 0$ for our debiasing program?

$$\min_{w \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \hat{\pi}_i w_i^2 : \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \leq \frac{\gamma}{n} \right\}.$$

- 2 Why is the asymptotic $(1 - \tau)$ -level confidence interval for $m_0(x)$ valid?

$$\left[\hat{m}^{\text{debias}}(x; \hat{w}) \pm \Phi^{-1} \left(1 - \frac{\tau}{2} \right) \cdot \hat{\sigma}_{\epsilon} \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\pi}_i \hat{w}_i^2} \right] \quad \text{with } \Phi(\cdot) \text{ being the CDF of } \mathcal{N}(0, 1).$$

► **Answer:** The above two questions can be addressed by the *dual formulation/solution* of our debiasing program!

The primal form of our debiasing program is a quadratic programming problem with a box constraint:

$$\min_{w \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \hat{\pi}_i w_i^2 : \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \leq \frac{\gamma}{n} \right\}.$$

The primal form of our debiasing program is a quadratic programming problem with a box constraint:

$$\min_{\mathbf{w} \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \hat{\pi}_i w_i^2 : \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \leq \frac{\gamma}{n} \right\}.$$

Proposition (Proposition 1 in [Zhang et al. 2023](#))

The dual form of our debiasing program is given by

$$\min_{\ell \in \mathbb{R}^d} \left\{ \frac{1}{4n} \sum_{i=1}^n \hat{\pi}_i [X_i^T \ell]^2 + x^T \ell + \frac{\gamma}{n} \|\ell\|_1 \right\}.$$

If the strong duality holds, we further have that

$$\hat{w}_i = -\frac{1}{2\sqrt{n}} \cdot X_i^T \hat{\ell} \quad \text{for } i = 1, \dots, n,$$

where $\hat{\mathbf{w}} \in \mathbb{R}^n$ and $\hat{\ell} \in \mathbb{R}^d$ are the solutions to the primal and dual debiasing program, respectively.

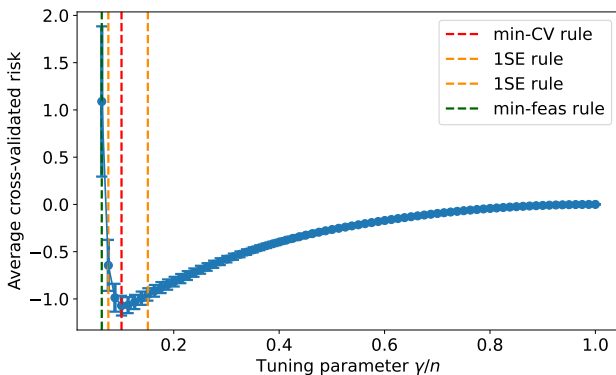
The dual form of our debiasing program is an *unconstrained* quadratic programming problem:

$$\min_{\ell \in \mathbb{R}^d} \left\{ \frac{1}{4n} \sum_{i=1}^n \widehat{\pi}_i [X_i^T \ell]^2 + x^T \ell + \frac{\gamma}{n} \|\ell\|_1 \right\}.$$

The dual form of our debiasing program is an *unconstrained* quadratic programming problem:

$$\min_{\ell \in \mathbb{R}^d} \left\{ \frac{1}{4n} \sum_{i=1}^n \hat{\pi}_i [X_i^T \ell]^2 + x^T \ell + \frac{\gamma}{n} \|\ell\|_1 \right\}.$$

We can fine-tune $\gamma > 0$ by cross-validation.



- Consider the regression function $m \equiv m(x) \in \mathbb{R}$ as the main parameter to be inferred and $\beta \in \mathbb{R}^d$ as the high-dimensional nuisance parameter.
- Our generic debiased estimator $m^{\text{debias}}(x, \mathbf{w})$ solves the sample-based estimating equation

$$\frac{1}{n} \sum_{i=1}^n \Xi_x(Y_i, R_i, X_i; m^{\text{debias}}, \beta) = m^{\text{debias}}(x; \mathbf{w}) - x^T \beta - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot R_i (Y_i - X_i^T \beta) = 0.$$

- Consider the regression function $m \equiv m(x) \in \mathbb{R}$ as the main parameter to be inferred and $\beta \in \mathbb{R}^d$ as the high-dimensional nuisance parameter.
- Our generic debiased estimator $m^{\text{debias}}(x, w)$ solves the sample-based estimating equation

$$\frac{1}{n} \sum_{i=1}^n \Xi_x(Y_i, R_i, X_i; m^{\text{debias}}, \beta) = m^{\text{debias}}(x; w) - x^T \beta - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i R_i (Y_i - X_i^T \beta) = 0.$$

- The Neyman near-orthogonalization condition ([Chernozhukov et al., 2018](#)) given $\mathbf{X} = (X_1, \dots, X_n)^T \in \mathbb{R}^{n \times d}$ at $(m_0, \beta_0) = (x^T \beta_0, \beta_0)$ requires

$$\begin{aligned} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Xi_x(Y_i, R_i, X_i; m_0, \beta_0) \middle| \mathbf{X} \right] &= 0, \\ \sup_{\beta \in \mathcal{T}_n} \left| \left\{ \frac{\partial}{\partial \beta} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Xi_x(Y_i, R_i, X_i; m, \beta) \middle| \mathbf{X} \right] \right\}_{(m_0, \beta_0)}^T (\beta - \beta_0) \right| &\leq \frac{\delta_n}{\sqrt{n}}, \end{aligned} \tag{4}$$

where \mathcal{T}_n is a properly shrinking neighborhood of β_0 and $\delta_n = o(1)$.

- Both conditions in (4) hold true, because for any $\beta \in \mathcal{T}_n$ and some convex set \mathcal{B} containing β_0 , we have that

$$\begin{aligned}
 & \left| \left\{ \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta} \mathbb{E} [\Xi_x(Y_i, R_i, X_i; m, \beta) | X] \Big|_{(m_0, \beta_0)} \right\}^T (\beta - \beta_0) \right| \\
 &= \left| \left[x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \pi(X_i) X_i \right]^T (\beta_0 - \beta) \right| \\
 &\leq \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \|\beta - \beta_0\|_1 \quad \text{by Hölder's inequality} \\
 &\leq \frac{\gamma}{n} \|\beta - \beta_0\|_1 \quad \text{by the box constraint in our debiasing program} \\
 &\leq \frac{\delta_n}{\sqrt{n}} \quad \text{by setting } \mathcal{T}_n = \left\{ \beta \in \mathcal{B} \subset \mathbb{R}^d : \|\beta - \beta_0\|_1 \leq \frac{\sqrt{n}\delta_n}{\gamma} \right\}.
 \end{aligned}$$

- Both conditions in (4) hold true, because for any $\beta \in \mathcal{T}_n$ and some convex set \mathcal{B} containing β_0 , we have that

$$\begin{aligned}
 & \left| \left\{ \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta} \mathbb{E} [\Xi_x(Y_i, R_i, X_i; m, \beta) | X] \Big|_{(m_0, \beta_0)} \right\}^T (\beta - \beta_0) \right| \\
 &= \left| \left[x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \pi(X_i) X_i \right]^T (\beta_0 - \beta) \right| \\
 &\leq \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \|\beta - \beta_0\|_1 \quad \text{by Hölder's inequality} \\
 &\leq \frac{\gamma}{n} \|\beta - \beta_0\|_1 \quad \text{by the box constraint in our debiasing program} \\
 &\leq \frac{\delta_n}{\sqrt{n}} \quad \text{by setting } \mathcal{T}_n = \left\{ \beta \in \mathcal{B} \subset \mathbb{R}^d : \|\beta - \beta_0\|_1 \leq \frac{\sqrt{n}\delta_n}{\gamma} \right\}.
 \end{aligned}$$

- Our debiasing program optimizes the (estimated) variance among all the estimators satisfying Neyman near-orthogonalization (4).

- Both conditions in (4) hold true, because for any $\beta \in \mathcal{T}_n$ and some convex set \mathcal{B} containing β_0 , we have that

$$\begin{aligned}
 & \left| \left\{ \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta} \mathbb{E} [\Xi_x(Y_i, R_i, X_i; m, \beta) | X] \Big|_{(m_0, \beta_0)} \right\}^T (\beta - \beta_0) \right| \\
 &= \left| \left[x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \pi(X_i) X_i \right]^T (\beta_0 - \beta) \right| \\
 &\leq \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \|\beta - \beta_0\|_1 \quad \text{by Hölder's inequality} \\
 &\leq \frac{\gamma}{n} \|\beta - \beta_0\|_1 \quad \text{by the box constraint in our debiasing program} \\
 &\leq \frac{\delta_n}{\sqrt{n}} \quad \text{by setting } \mathcal{T}_n = \left\{ \beta \in \mathcal{B} \subset \mathbb{R}^d : \|\beta - \beta_0\|_1 \leq \frac{\sqrt{n}\delta_n}{\gamma} \right\}.
 \end{aligned}$$

- Our debiasing program optimizes the (estimated) variance among all the estimators satisfying Neyman near-orthogonalization (4).
- (4) also allows our debiasing program to *de-correlate* the Lasso pilot regression from propensity score estimation and weight optimization.

Asymptotic Theory



- **Goal:** Establish the asymptotic normality of our debiased estimator

$$\hat{m}^{\text{debias}}(x; \hat{w}) = x^T \hat{\beta} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i \left(Y_i - X_i^T \hat{\beta} \right).$$

- **Goal:** Establish the asymptotic normality of our debiased estimator

$$\hat{m}^{\text{debias}}(x; \hat{\mathbf{w}}) = x^T \hat{\beta} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i \left(Y_i - X_i^T \hat{\beta} \right).$$

- **Naive Attempt:** Linearity assumption $Y_i = X_i^T \beta_0 + \epsilon_i$ for $i = 1, \dots, n$ implies that

$$\sqrt{n} \left[\hat{m}^{\text{debias}}(x; \hat{\mathbf{w}}) - m_0(x) \right] = \underbrace{\sum_{i=1}^n \hat{w}_i R_i \epsilon_i}_{\text{Not an i.i.d. sum!}} + \left[x - \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i X_i \right]^T \sqrt{n} (\hat{\beta} - \beta_0),$$

- **Goal:** Establish the asymptotic normality of our debiased estimator

$$\hat{m}^{\text{debias}}(x; \hat{\mathbf{w}}) = x^T \hat{\beta} + \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i \left(Y_i - X_i^T \hat{\beta} \right).$$

- **Naive Attempt:** Linearity assumption $Y_i = X_i^T \beta_0 + \epsilon_i$ for $i = 1, \dots, n$ implies that

$$\sqrt{n} \left[\hat{m}^{\text{debias}}(x; \hat{\mathbf{w}}) - m_0(x) \right] = \underbrace{\sum_{i=1}^n \hat{w}_i R_i \epsilon_i}_{\text{Not an i.i.d. sum!}} + \left[x - \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{w}_i R_i X_i \right]^T \sqrt{n} (\hat{\beta} - \beta_0),$$

- **Solution:** With the dual relation $\hat{w}_i = -\frac{1}{2\sqrt{n}} \cdot X_i^T \hat{\ell}$, $i = 1, \dots, n$, we obtain

$$\begin{aligned} \sqrt{n} \left[\hat{m}^{\text{debias}}(x; \hat{\mathbf{w}}) - m_0(x) \right] &= -\frac{1}{2\sqrt{n}} \sum_{i=1}^n R_i \epsilon_i X_i^T \hat{\ell} + \left[x + \frac{1}{2n} \sum_{i=1}^n R_i X_i X_i^T \hat{\ell} \right]^T \sqrt{n} (\beta_0 - \hat{\beta}) \\ &= \underbrace{-\frac{1}{2\sqrt{n}} \sum_{i=1}^n R_i \epsilon_i X_i^T \ell_0(x)}_{\text{i.i.d. sum!}} + \underbrace{\text{“Bias terms”}}_{o_p(1)}. \end{aligned}$$

- 1 The covariate vector $X \in \mathbb{R}^d$ and the noise $\epsilon \in \mathbb{R}$ are sub-Gaussian.

- 1 The covariate vector $X \in \mathbb{R}^d$ and the noise $\epsilon \in \mathbb{R}$ are sub-Gaussian.
- 2 There exists a constant $\kappa_R > 0$ such that

$$\inf_{v \in \mathbb{S}^{d-1}} \mathbb{E} [R(X^T v)^2] \geq \kappa_R^2 \quad \text{with} \quad \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}.$$

- ① The covariate vector $X \in \mathbb{R}^d$ and the noise $\epsilon \in \mathbb{R}$ are sub-Gaussian.
- ② There exists a constant $\kappa_R > 0$ such that

$$\inf_{v \in \mathbb{S}^{d-1}} \mathbb{E} [R(X^T v)^2] \geq \kappa_R^2 \quad \text{with} \quad \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}.$$

- ③ Given any $n \geq 1$ and $\delta \in (0, 1)$, there exists $r_\pi \equiv r_\pi(n, \delta) > 0$ such that

$$\mathbb{P} \left(\max_{1 \leq i \leq n} |\hat{\pi}_i - \pi_i| > r_\pi \right) < \delta \quad \text{with} \quad \pi_i = \pi(X_i), i = 1, \dots, n.$$

- ① The covariate vector $X \in \mathbb{R}^d$ and the noise $\epsilon \in \mathbb{R}$ are sub-Gaussian.
- ② There exists a constant $\kappa_R > 0$ such that

$$\inf_{v \in \mathbb{S}^{d-1}} \mathbb{E} [R(X^T v)^2] \geq \kappa_R^2 \quad \text{with} \quad \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}.$$

- ③ Given any $n \geq 1$ and $\delta \in (0, 1)$, there exists $r_\pi \equiv r_\pi(n, \delta) > 0$ such that

$$\mathbb{P} \left(\max_{1 \leq i \leq n} |\hat{\pi}_i - \pi_i| > r_\pi \right) < \delta \quad \text{with} \quad \pi_i = \pi(X_i), i = 1, \dots, n.$$

- ④ Define the population dual program as:

$$\min_{\ell \in \mathbb{R}^d} \left\{ \frac{1}{4} \mathbb{E} [R(X^T \ell)^2] + x^T \ell \right\},$$

whose exact solution is $\ell_0(x) = -2 [\mathbb{E} (RXX^T)]^{-1} x$. We assume that the r_ℓ -approximation $\tilde{\ell}(x)$ to $\ell_0(x)$ is sparse with $r_\ell \in [0, \frac{1}{2}]$, i.e.,

$$s_\ell(x) = \left\| \tilde{\ell}(x) \right\|_0 \ll \min\{n, d\} \quad \text{with} \quad \tilde{\ell}(x) = \arg \min_{u \in \mathbb{R}^d} \{ \|u\|_0 : \|u - \ell_0(x)\|_2 \leq r_\ell \|\ell_0(x)\|_2 \}.$$

- ① **Consistency of Lasso pilot estimate:** If $\lambda \asymp \sigma_\epsilon \sqrt{\frac{\log d}{n}}$ with $\log d = o(n)$, then $\left\| \hat{\beta} - \beta_0 \right\|_2 = O_P \left(\frac{1}{\kappa_R^2} \sqrt{\frac{s_\beta \log d}{n}} \right)$.
- ② **Consistency of the solution to the dual debiasing program:** If r_ℓ shrinks to 0 in a certain rate and $\frac{\gamma}{n} \asymp \frac{\|x\|_2}{\kappa_R} \sqrt{\frac{\log d}{n}} + \frac{\|x\|_2}{\kappa_R^2} \cdot r_\pi$, then

$$\left\| \hat{\ell}(x) - \ell_0(x) \right\|_2 = O_P \left(\frac{1}{\kappa_R^3} \sqrt{\frac{s_\ell(x) \log d}{n}} + \frac{r_\ell}{\kappa_R^4} + \frac{r_\pi \sqrt{s_\ell(x)}}{\kappa_R^4} \right).$$

Note: Under the same choice of $\gamma > 0$, the strong duality holds.

- ① **Consistency of Lasso pilot estimate:** If $\lambda \asymp \sigma_\epsilon \sqrt{\frac{\log d}{n}}$ with $\log d = o(n)$, then $\left\| \hat{\beta} - \beta_0 \right\|_2 = O_P \left(\frac{1}{\kappa_R^2} \sqrt{\frac{s_\beta \log d}{n}} \right)$.
- ② **Consistency of the solution to the dual debiasing program:** If r_ℓ shrinks to 0 in a certain rate and $\frac{\gamma}{n} \asymp \frac{\|x\|_2}{\kappa_R} \sqrt{\frac{\log d}{n}} + \frac{\|x\|_2}{\kappa_R^2} \cdot r_\pi$, then

$$\left\| \hat{\ell}(x) - \ell_0(x) \right\|_2 = O_P \left(\frac{1}{\kappa_R^3} \sqrt{\frac{s_\ell(x) \log d}{n}} + \frac{r_\ell}{\kappa_R^4} + \frac{r_\pi \sqrt{s_\ell(x)}}{\kappa_R^4} \right).$$

Note: Under the same choice of $\gamma > 0$, the strong duality holds.

Theorem (Theorem 7 in Zhang et al. 2023)

If $\frac{(1+\kappa_R^2)s_{\max} \log(nd)}{\kappa_R^4} = o(\sqrt{n})$, $\frac{(1+\kappa_R^4)\sqrt{s_{\max} \log(nd)}}{\kappa_R^6} (r_\ell + r_\pi) = o(1)$, and $\|x\|_2 = O(1)$ with $s_{\max} = \{s_\beta, s_\ell(x)\}$, then

$$\frac{\sqrt{n} [\hat{m}^{\text{debias}}(x; \hat{w}) - m_0(x)]}{\sigma_m(x)} \xrightarrow{d} \mathcal{N}(0, 1).$$

- ① Our growth requirement $s_{\max} = o\left(\frac{\sqrt{n}}{\log d}\right)$ on the sparsity level is a standard and *essentially necessary* condition for asymptotic normality; see Section 8.6 of [Jankova and van de Geer \(2018\)](#).

- ① Our growth requirement $s_{\max} = o\left(\frac{\sqrt{n}}{\log d}\right)$ on the sparsity level is a standard and *essentially necessary* condition for asymptotic normality; see Section 8.6 of [Jankova and van de Geer \(2018\)](#).

- ② Given any dimension $d > 0$, the asymptotic variance of our debiased estimator

$$\sigma_m^2(x) = \sigma_\epsilon^2 \cdot x^T [E(RXX^T)]^{-1} x$$

attains the *semi-parametric efficiency bound* among all asymptotically linear estimators under MAR outcomes ([Müller and Keilegom, 2012](#)).

- ① Our growth requirement $s_{\max} = o\left(\frac{\sqrt{n}}{\log d}\right)$ on the sparsity level is a standard and *essentially necessary* condition for asymptotic normality; see Section 8.6 of [Jankova and van de Geer \(2018\)](#).

- ② Given any dimension $d > 0$, the asymptotic variance of our debiased estimator

$$\sigma_m^2(x) = \sigma_\epsilon^2 \cdot x^T [E(RXX^T)]^{-1} x$$

attains the *semi-parametric efficiency bound* among all asymptotically linear estimators under MAR outcomes ([Müller and Keilegom, 2012](#)).

Proposition (Proposition 8 in [Zhang et al. 2023](#))

If $\frac{(1+\kappa_R^3)}{\kappa_R^5} \sqrt{\frac{s_\ell(x) \log(nd)}{n}} = o(1)$, $\frac{(1+\kappa_R^4)}{\kappa_R^6} [r_\ell + r_\pi \sqrt{s_\ell(x)}] = o(1)$, and $\|x\|_2 = O(1)$, then

$$\left| \sum_{i=1}^n \hat{\pi}_i \hat{w}_i^2 - x^T [E(RXX^T)]^{-1} x \right| = o_P(1).$$

Our theoretical results also provide insightful answers to the following two questions:

- Why don't we need sample splitting or cross fitting?
- Why can we estimate the propensity score by any machine learning methods without worrying about the overfitting issue?

Our theoretical results also provide insightful answers to the following two questions:

- Why don't we need sample splitting or cross fitting?
- Why can we estimate the propensity score by any machine learning methods without worrying about the overfitting issue?

► **Answer:** Our asymptotic normality result depends on the *in-sample* estimation error r_π of the propensity score; recall that

$$P\left(\max_{1 \leq i \leq n} |\hat{\pi}_i - \pi_i| > r_\pi\right) < \delta \quad \text{with} \quad \pi_i = \pi(X_i), i = 1, \dots, n.$$

- In other words, our debiased estimator performs even better when we overfit the propensity scores $\pi(X_i) = P(R_i = 1|X_i), i = 1, \dots, n$.

Our theoretical results also provide insightful answers to the following two questions:

- Why don't we need sample splitting or cross fitting?
- Why can we estimate the propensity score by any machine learning methods without worrying about the overfitting issue?

► **Answer:** Our asymptotic normality result depends on the *in-sample* estimation error r_π of the propensity score; recall that

$$P\left(\max_{1 \leq i \leq n} |\hat{\pi}_i - \pi_i| > r_\pi\right) < \delta \quad \text{with} \quad \pi_i = \pi(X_i), i = 1, \dots, n.$$

- In other words, our debiased estimator performs even better when we overfit the propensity scores $\pi(X_i) = P(R_i = 1|X_i), i = 1, \dots, n$.
- This coincides with “*benign overfitting*” in linear regression or neural networks (Bartlett et al., 2020; Li et al., 2021; Cao et al., 2022).

Comparative Simulations



We compare our debiasing method with L_1 -penalized logistic regression for the propensity score estimation with several existing methods:

- “DL-Jav”: The debiased Lasso by [Javanmard and Montanari \(2014\)](#).
- “DL-vdG”: The debiased Lasso by [van de Geer et al. \(2014\)](#).
- “Refit”: Run the regular least-square regression on the support set of the Lasso pilot estimate ([Belloni and Chernozhukov, 2013](#)).

We compare our debiasing method with L_1 -penalized logistic regression for the propensity score estimation with several existing methods:

- “DL-Jav”: The debiased Lasso by [Javanmard and Montanari \(2014\)](#).
- “DL-vdG”: The debiased Lasso by [van de Geer et al. \(2014\)](#).
- “Refit”: Run the regular least-square regression on the support set of the Lasso pilot estimate ([Belloni and Chernozhukov, 2013](#)).

These methods to be compared are implemented on

- Complete-case (CC) data $\{(X_i, Y_i, R_i = 1)\}_{i=1}^n$;
- Inverse probability weighted (IPW) data $\left\{ \left(\frac{X_i}{\sqrt{\hat{\pi}_i}}, \frac{Y_i}{\sqrt{\hat{\pi}_i}}, R_i = 1 \right) \right\}_{i=1}^n$;
- Oracle fully observed data (X_i, Y_i) for $i = 1, \dots, n$.

We compare our debiasing method with L_1 -penalized logistic regression for the propensity score estimation with several existing methods:

- “DL-Jav”: The debiased Lasso by [Javanmard and Montanari \(2014\)](#).
- “DL-vdG”: The debiased Lasso by [van de Geer et al. \(2014\)](#).
- “Refit”: Run the regular least-square regression on the support set of the Lasso pilot estimate ([Belloni and Chernozhukov, 2013](#)).

These methods to be compared are implemented on

- Complete-case (CC) data $\{(X_i, Y_i, R_i = 1)\}_{i=1}^n$;
- Inverse probability weighted (IPW) data $\left\{ \left(\frac{X_i}{\sqrt{\hat{\pi}_i}}, \frac{Y_i}{\sqrt{\hat{\pi}_i}}, R_i = 1 \right) \right\}_{i=1}^n$;
- Oracle fully observed data (X_i, Y_i) for $i = 1, \dots, n$.

Evaluation metrics on 1000 Monte Carlo experiments include

- Average absolute bias $|\hat{m}^{\text{debias}}(x) - m_0(x)|$;
- Average coverage of the yielded 95% confidence intervals;
- Average length of the yielded 95% confidence intervals.

W Simulation Results Under Gaussian Noises (I)

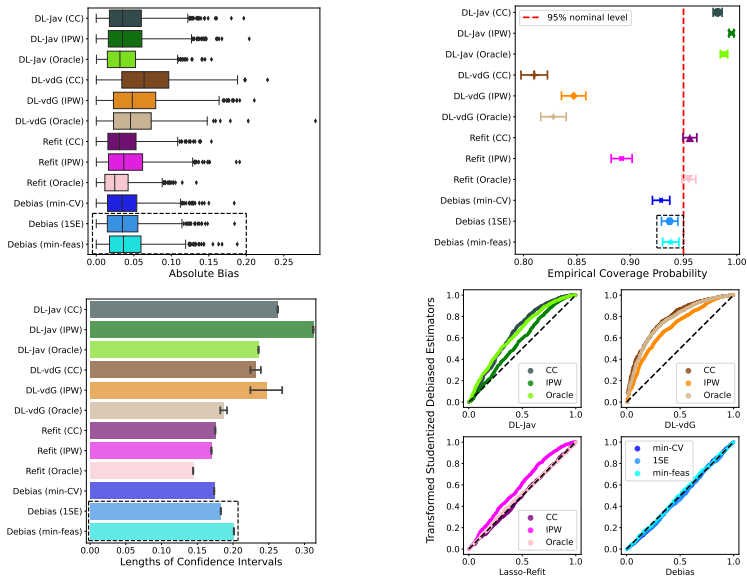


Figure 2: Sparse β_0^{sp} and sparse $x^{(2)}$ with $X_i \sim \mathcal{N}_d(\mathbf{0}, \Sigma^{\text{cs}})$, $i = 1, \dots, n$.

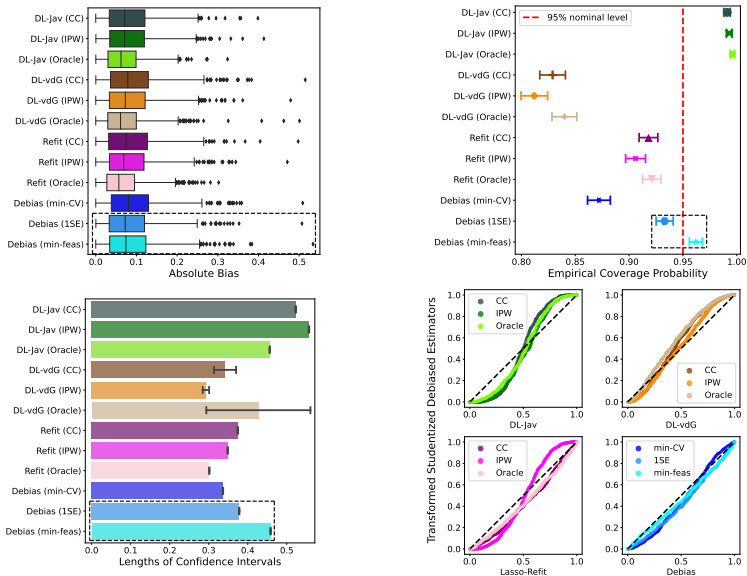


Figure 3: Pseudo-dense β_0^{pd} and sparse $x^{(2)}$ with $X_i \sim \mathcal{N}_d(\mathbf{0}, \Sigma^{\text{ar}})$, $i = 1, \dots, n$.

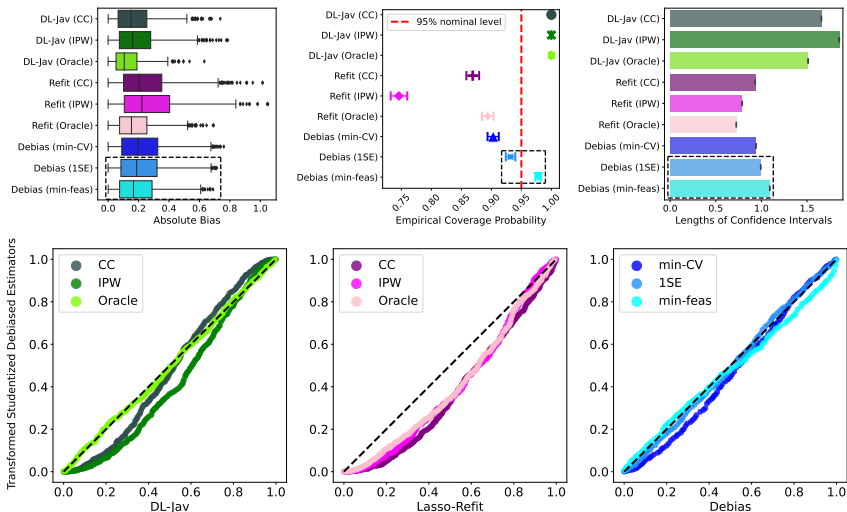


Figure 4: Dense β_0^{de} and sparse $x^{(4)}$ with $X_i \sim \mathcal{N}_d(\mathbf{0}, \Sigma^{cs})$, $i = 1, \dots, n$.

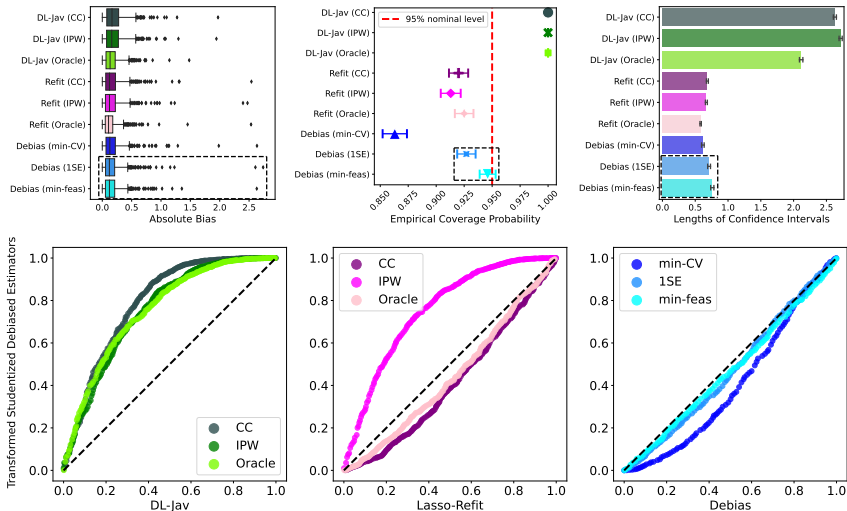


Figure 5: Pseudo-dense β_0^{pd} and dense $x^{(4)}$ with $X_i \sim \mathcal{N}_d(\mathbf{0}, \Sigma^{\text{ar}})$, $i = 1, \dots, n$. Note that the mean-zero t_2 distribution has *infinite* variance.

- ① True propensity score model: $P(R_i = 1|X_i) = \Phi\left(-4 + \sum_{k=1}^K Z_{ik}\right)$, where (Z_{i1}, \dots, Z_{iK}) contains all polynomial combinations of the first eight components X_{i1}, \dots, X_{i8} of $X_i \in \mathbb{R}^{1000}$ with degrees ≤ 2 .

- ① True propensity score model: $P(R_i = 1|X_i) = \Phi\left(-4 + \sum_{k=1}^K Z_{ik}\right)$, where (Z_{i1}, \dots, Z_{iK}) contains all polynomial combinations of the first eight components X_{i1}, \dots, X_{i8} of $X_i \in \mathbb{R}^{1000}$ with degrees ≤ 2 .
- ② Estimate the propensity scores $\pi(X_i), i = 1, \dots, n$ by the following nonlinear/nonparametric machine learning methods:
 - **Gaussian Naive Bayes (“NB”).**
 - **Random Forest (“RF”):** 100 trees, bootstrapping samples, and the Gini impurity.
 - **Support Vector Machine (“SVM”):** Gaussian radial basis function.
 - **Neural Network (“NN”):** Two hidden layers of size 80×50 and ReLU $h(x) = \max\{x, 0\}$ as the activation function.

- ① True propensity score model: $P(R_i = 1|X_i) = \Phi\left(-4 + \sum_{k=1}^K Z_{ik}\right)$, where (Z_{i1}, \dots, Z_{iK}) contains all polynomial combinations of the first eight components X_{i1}, \dots, X_{i8} of $X_i \in \mathbb{R}^{1000}$ with degrees ≤ 2 .
- ② Estimate the propensity scores $\pi(X_i), i = 1, \dots, n$ by the following nonlinear/nonparametric machine learning methods:
 - **Gaussian Naive Bayes (“NB”).**
 - **Random Forest (“RF”):** 100 trees, bootstrapping samples, and the Gini impurity.
 - **Support Vector Machine (“SVM”):** Gaussian radial basis function.
 - **Neural Network (“NN”):** Two hidden layers of size 80×50 and ReLU $h(x) = \max\{x, 0\}$ as the activation function.
- ③ Include an extra evaluation metric as the average mean absolute error (“Avg-MAE”) for the estimated propensity scores.

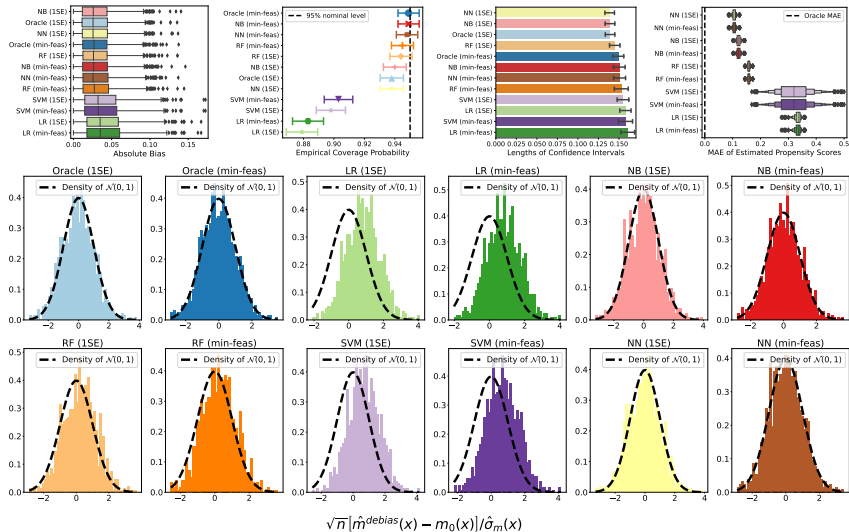
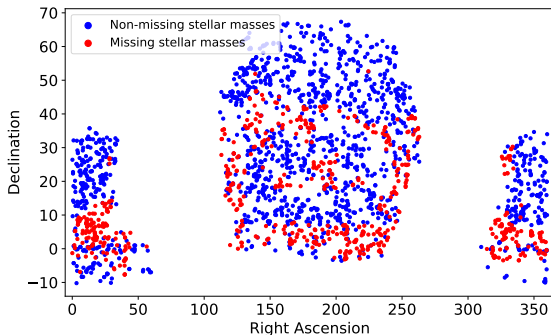


Figure 6: Sparse β_0^{sp} and (weakly) dense $x^{(4)}$.

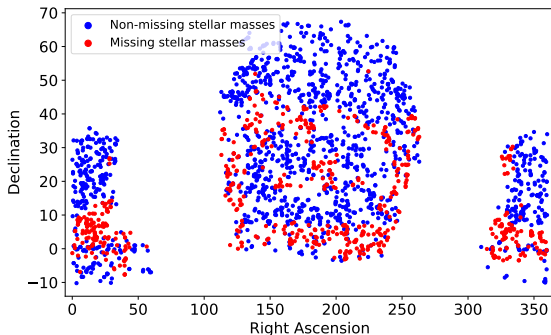
Real-World Applications



Recall that some estimated stellar masses of the observed galaxies in SDSS-IV are missing in the most recent Firefly value-added catalog.



Recall that some estimated stellar masses of the observed galaxies in SDSS-IV are missing in the most recent Firefly value-added catalog.



► Scientific Questions:

- 1 *How can we conduct valid inference on the (estimated) stellar mass based on the spectroscopic and photometric properties?*
- 2 *Is it statistically significant that the stellar mass of a galaxy is negatively correlated with its distance to the nearby cosmic filament structures?*

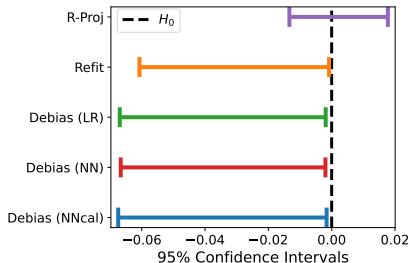
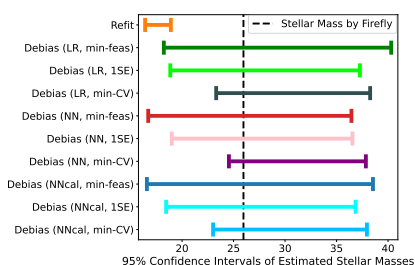
- 1 Consider all the observed galaxies by SDSS-IV within a thin redshift slice $0.4 \sim 0.4005$, among which 30.2% of their stellar masses are missing in the Firefly value-added catalog.

- 1 Consider all the observed galaxies by SDSS-IV within a thin redshift slice $0.4 \sim 0.4005$, among which 30.2% of their stellar masses are missing in the Firefly value-added catalog.
- 2 Fetch their spectroscopic and photometric properties from SDSS-IV DR16 database similar to the input catalog of [Chang et al. \(2015\)](#).

- 1 Consider all the observed galaxies by SDSS-IV within a thin redshift slice $0.4 \sim 0.4005$, among which 30.2% of their stellar masses are missing in the Firefly value-added catalog.
- 2 Fetch their spectroscopic and photometric properties from SDSS-IV DR16 database similar to the input catalog of [Chang et al. \(2015\)](#).
- 3 Apply feature transformation, remove highly linearly correlated covariates, and generate univariate B-spline base covariates of polynomial order 3 with 40 knots.

- 1 Consider all the observed galaxies by SDSS-IV within a thin redshift slice $0.4 \sim 0.4005$, among which 30.2% of their stellar masses are missing in the Firefly value-added catalog.
- 2 Fetch their spectroscopic and photometric properties from SDSS-IV DR16 database similar to the input catalog of [Chang et al. \(2015\)](#).
- 3 Apply feature transformation, remove highly linearly correlated covariates, and generate univariate B-spline base covariates of polynomial order 3 with 40 knots.
- 4 Incorporate RA, DEC, and the angular diameter distances from the galaxies to the two-dimensional spherical cosmic filaments by [Zhang and Chen \(2023\)](#); [Zhang et al. \(2022\)](#).

- 1 Consider all the observed galaxies by SDSS-IV within a thin redshift slice $0.4 \sim 0.4005$, among which 30.2% of their stellar masses are missing in the Firefly value-added catalog.
 - 2 Fetch their spectroscopic and photometric properties from SDSS-IV DR16 database similar to the input catalog of [Chang et al. \(2015\)](#).
 - 3 Apply feature transformation, remove highly linearly correlated covariates, and generate univariate B-spline base covariates of polynomial order 3 with 40 knots.
 - 4 Incorporate RA, DEC, and the angular diameter distances from the galaxies to the two-dimensional spherical cosmic filaments by [Zhang and Chen \(2023\)](#); [Zhang et al. \(2022\)](#).
 - 5 Control for the confounding effects by including the distances from galaxies to candidate galaxy clusters.
- **Final Dataset:** $n = 1185$ and $d = 1409$.



- *Left Panel:* 95% confidence intervals by different debiasing methods for the estimated stellar mass of a new galaxy.
- *Right Panel:* 95% confidence intervals by different debiasing methods for the estimated regression coefficient associated with the distance to nearby cosmic filaments.

Conclusions and Future Works



We develop an efficient debiasing method for conducting valid inference on high-dimensional linear models with MAR outcomes.

We develop an efficient debiasing method for conducting valid inference on high-dimensional linear models with MAR outcomes.

- Its computational and statistical efficiencies follow from the dual formulation.
- Sample splitting and cross fitting are not required, and the nuisance propensity score can be estimated by any machine learning method.
- We provide interpretations to our debiasing method from the viewpoints of bias-variance trade-off and Neyman near-orthogonalization.
- Comprehensive simulation studies and motivating applications demonstrate the potential of our proposed debiasing method.

W Potential Application to Causal Inference (I)

The observable data in causal inference are

$$\{(\mathbb{Y}_i, T_i, X_i)\}_{i=1}^n \subset \mathbb{R} \times \{0, 1\} \times \mathbb{R}^d.$$

- $T_i \in \{0, 1\}$ is a binary treatment assignment indicator;
- $\mathbb{Y}_i = T_i \cdot Y(1)_i + (1 - T_i) \cdot Y(0)_i$ with $Y(0), Y(1)$ as potential outcomes.

► **Objective:** Conduct valid inference on the regression function (or conditional mean outcome) of the treatment group.

Treatment Group	X_1^T	$Y(1)_1$
	\vdots	\vdots
	$X_{\frac{n}{2}}^T$	$Y(1)_{\frac{n}{2}}$
Control Group	$X_{\frac{n}{2}+1}^T$	$Y(0)_{\frac{n}{2}+1}$
	\vdots	\vdots
	X_n^T	$Y(0)_n$


 $E(Y|X, T = 1)$
based on
 $\{(Y(1)_i, X_i)\}_{i=1}^{\frac{n}{2}}$

Figure 8: Traditional approaches for inferring $E(Y|X, T = 1)$.

W Potential Application to Causal Inference (I)

The observable data in causal inference are

$$\{(\mathbb{Y}_i, T_i, X_i)\}_{i=1}^n \subset \mathbb{R} \times \{0, 1\} \times \mathbb{R}^d.$$

- $T_i \in \{0, 1\}$ is a binary treatment assignment indicator;
- $\mathbb{Y}_i = T_i \cdot Y(1)_i + (1 - T_i) \cdot Y(0)_i$ with $Y(0), Y(1)$ as potential outcomes.

► **Objective:** Conduct valid inference on the regression function (or conditional mean outcome) of the treatment group.

Treatment Group	X_1^T	$Y(1)_1$
	\vdots	\vdots
	$X_{\frac{n}{2}}^T$	$Y(1)_{\frac{n}{2}}$
Control Group	$X_{\frac{n}{2}+1}^T$	$Y(0)_{\frac{n}{2}+1}$
	\vdots	\vdots
	X_n^T	$Y(0)_n$

$E(Y|X, T = 1)$
based on
 $\{(Y(1)_i, T_i, X_i)\}_{i=1}^n$

Figure 8: Our approach for inferring $E(Y|X, T = 1)$.

Our debiasing method can be extended to valid inference on the linear average conditional treatment effect (ACTE)

$$E[Y(1) - Y(0)|X]$$

with no unmeasured confounding and high-dimensional covariates.

Our debiasing method can be extended to valid inference on the linear average conditional treatment effect (ACTE)

$$E[Y(1) - Y(0)|X]$$

with no unmeasured confounding and high-dimensional covariates.

- The modified debiasing program with tuning parameters $\gamma_1, \gamma_2 > 0$ is

$$\begin{aligned} & \arg \min_{w_{(0)}, w_{(1)} \in \mathbb{R}^n} \sum_{i=1}^n \left[\hat{\pi}_i w_{i(1)}^2 + (1 - \hat{\pi}_i) w_{i(0)}^2 \right] \\ \text{s.t. } & \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_{i(1)} \cdot \hat{\pi}_i \cdot X_i \right\|_{\infty} \leq \frac{\gamma_1}{n} \quad \text{and} \quad \left\| x - \frac{1}{\sqrt{n}} \sum_{i=1}^n w_{i(0)} (1 - \hat{\pi}_i) X_i \right\|_{\infty} \leq \frac{\gamma_2}{n}. \end{aligned}$$

- The extended debiased estimator becomes

$$\begin{aligned} & \hat{m}^{\text{debias}}(x; \hat{w}_{(1)}, \hat{w}_{(0)}) \\ &= x^T (\hat{\beta}_{(1)} - \hat{\beta}_{(0)}) + \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[\hat{w}_{i(1)} \cdot T_i \left(\mathbb{Y}_i - X_i^T \hat{\beta}_{(1)} \right) - \hat{w}_{i(0)} \cdot (1 - T_i) \left(\mathbb{Y}_i - X_i^T \hat{\beta}_{(0)} \right) \right]. \end{aligned}$$

- The efficiency theory for this modified procedure is worth studying!

Thank you!

More details can be found in

[1] Y. Zhang, A. Giessing, and Y.-C. Chen. Efficient Inference on High-Dimensional Linear Models with Missing Outcomes. *arXiv preprint*, 2023. <https://arxiv.org/abs/2309.06429>.

Python Package: [Debias-Infer](#) and R Package: [DebiasInfer](#).



- A. Agrawal, R. Verschueren, S. Diamond, and S. Boyd. A rewriting system for convex optimization problems. *Journal of Control and Decision*, 5(1):42–60, 2018.
- P. L. Bartlett, P. M. Long, G. Lugosi, and A. Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- A. Belloni and V. Chernozhukov. Least squares after model selection in high-dimensional sparse models. *Bernoulli*, 19(2):521–547, 2013.
- A. Belloni, V. Chernozhukov, and K. Kato. Valid post-selection inference in high-dimensional approximately sparse quantile regression models. *Journal of the American Statistical Association*, 114(526):749–758, 2019.
- L. Breiman, J. Friedman, C. J. Stone, and R. Olshen. *Classification and Regression Trees*. Chapman and Hall/CRC, 1984.
- Y. Cao, Z. Chen, M. Belkin, and Q. Gu. Benign overfitting in two-layer convolutional neural networks. *Advances in neural information processing systems*, 35:25237–25250, 2022.
- C. M. Carvalho, J. Chang, J. E. Lucas, J. R. Nevins, Q. Wang, and M. West. High-dimensional sparse factor modeling: applications in gene expression genomics. *Journal of the American Statistical Association*, 103(484):1438–1456, 2008.
- A. Chakraborty, J. Lu, T. T. Cai, and H. Li. High dimensional m-estimation with missing outcomes: A semi-parametric framework. *arXiv preprint arXiv:1911.11345*, 2019.
- Y.-Y. Chang, A. van der Wel, E. da Cunha, and H.-W. Rix. Stellar masses and star formation rates for 1 m galaxies from sdss+ wise. *The Astrophysical Journal Supplement Series*, 219(1):8, 2015.
- O. Chapelle, B. Schölkopf, and A. Zien. *Semi-Supervised Learning*. The MIT Press, 2006.
- Y. Chen and Y. Yang. The one standard error rule for model selection: does it work? *Stats*, 4(4):868–892, 2021.

- V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 01 2018.
- J. Comparat, C. Maraston, D. Goddard, V. Gonzalez-Perez, J. Lian, S. Meneses-Goytia, D. Thomas, J. R. Brownstein, R. Tojeiro, A. Finoguenov, et al. Stellar population properties for 2 million galaxies from sdss dr14 and deep2 dr4 from full spectral fitting. *arXiv preprint arXiv:1711.06575*, 2017.
- S. Diamond and S. Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- J. Fan, J. Lv, and L. Qi. Sparse high-dimensional models in economics. *Annu. Rev. Econ.*, 3(1):291–317, 2011.
- A. Fu, B. Narasimhan, and S. Boyd. CVXR: An R package for disciplined convex optimization. *Journal of Statistical Software*, 94(14):1–34, 2020. doi: 10.18637/jss.v094.i14.
- A. Giessing and J. Wang. Debiased inference on heterogeneous quantile treatment effects with regression rank scores. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, page qkad075, 08 2023.
- J. P. Higgins, I. R. White, and A. M. Wood. Imputation methods for missing outcome data in meta-analysis of clinical trials. *Clinical trials*, 5(3):225–239, 2008.
- J. Jackson. A critique of rees’s theory of primordial gravitational radiation. *Monthly Notices of the Royal Astronomical Society*, 156(1):1P–5P, 1972.
- J. Jankova and S. van de Geer. Semiparametric efficiency bounds for high-dimensional models. *The Annals of Statistics*, 46(5):2336–2359, 2018.
- A. Javanmard and A. Montanari. Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research*, 15(1):2869–2909, 2014.

- N. Kaiser. Clustering in real space and in redshift space. *Monthly Notices of the Royal Astronomical Society*, 227(1):1–21, 1987.
- U. Kuchner, A. Aragón-Salamanca, A. Rost, F. R. Pearce, M. E. Gray, W. Cui, A. Knebe, E. Rasia, and G. Yepes. Cosmic filaments in galaxy cluster outskirts: quantifying finding filaments in redshift space. *Monthly Notices of the Royal Astronomical Society*, 503(2):2065–2076, 2021.
- Z. Li, Z.-H. Zhou, and A. Gretton. Towards an understanding of benign overfitting in neural networks. *arXiv preprint arXiv:2106.03212*, 2021.
- N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the Lasso. *The Annals of Statistics*, 34(3):1436 – 1462, 2006.
- P. Müller and S. van de Geer. The partial linear model in high dimensions. *Scandinavian Journal of Statistics*, 42(2):580–608, 2015.
- U. U. Müller and I. V. Keilegom. Efficient parameter estimation in regression with missing responses. *Electronic Journal of Statistics*, 6(none):1200 – 1219, 2012.
- P. Ravikumar, J. Lafferty, H. Liu, and L. Wasserman. Sparse additive models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 71(5):1009–1030, 2009.
- T. Sun and C.-H. Zhang. Scaled sparse linear regression. *Biometrika*, 99(4):879–898, 2012.
- S. van de Geer, P. Bühlmann, Y. Ritov, and R. Dezeure. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202, 2014.
- S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
- R. Wyss, C. Yanover, T. El-Hay, D. Bennett, R. W. Platt, A. R. Zullo, G. Sari, X. Wen, Y. Ye, H. Yuan, et al. Machine learning for improving high-dimensional proxy confounder adjustment in healthcare database studies: An overview of the current literature. *Pharmacoepidemiology and Drug Safety*, 31(9): 932–943, 2022.

- C.-H. Zhang and S. S. Zhang. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014.
- Y. Zhang and Y.-C. Chen. Linear convergence of the subspace constrained mean shift algorithm: from euclidean to directional data. *Information and Inference: A Journal of the IMA*, 12(1):210–311, 2023.
- Y. Zhang, R. S. de Souza, and Y.-C. Chen. Sconce: a cosmic web finder for spherical and conic geometries. *Monthly Notices of the Royal Astronomical Society*, 517(1):1197–1217, 2022.
- Y. Zhang, A. Giessing, and Y.-C. Chen. Efficient inference on high-dimensional linear models with missing outcomes. *arXiv preprint arXiv:2309.06429*, 2023.

- ① **Lasso pilot estimate:** We adopt the scaled Lasso (Sun and Zhang, 2012) with its universal regularization parameter $\lambda_0 = \sqrt{\frac{2 \log d}{n}}$ as the initialization. Specifically, it iteratively updates $\hat{\beta}(\tilde{\lambda}), \hat{\sigma}_\epsilon(\tilde{\lambda}), \tilde{\lambda}$ via the jointly convex optimization program:

$$\left(\hat{\beta}(\tilde{\lambda}), \hat{\sigma}_\epsilon(\tilde{\lambda}) \right) = \arg \min_{\beta \in \mathbb{R}^d, \sigma_\epsilon > 0} \left[\frac{1}{2n\sigma_\epsilon} \sum_{i=1}^n R_i \left(Y_i - X_i^T \beta \right)^2 + \frac{\sigma_\epsilon}{2} + \tilde{\lambda} \|\beta\|_1 \right].$$

- ② **Debiasing program:** We solve the primal program by Python package “CVXPY” (Diamond and Boyd, 2016; Agrawal et al., 2018) or R package “CVXR” (Fu et al., 2020). For the dual program, we formulate a coordinate descent algorithm (Wright, 2015) as:

$$\left[\hat{\ell}(x) \right]_j \leftarrow \frac{\mathcal{S}_{\frac{\gamma}{n}} \left(-\frac{1}{2n} \sum_{i=1}^n \hat{\pi}_i \left(\sum_{k \neq j} X_{ik} X_{jk} \left[\hat{\ell}(x) \right]_k \right) - x_j \right)}{\frac{1}{2n} \sum_{i=1}^n \hat{\pi}_i X_{ij}^2} \quad \text{for } j = 1, \dots, d,$$

where $\mathcal{S}_{\frac{\gamma}{n}}(u) = \text{sign}(u) \cdot \left(u - \frac{\gamma}{n} \right)_+$ is the soft-thresholding operator.

- Suppose that we conduct a K -fold cross-validation on a candidate set $\Gamma = \{\gamma_1, \dots, \gamma_m\}$ of the tuning parameter.
- For each $\gamma_i \in \Gamma$, we compute the cross-validated risk or error on each fold of the data as:

$$CV_k(\gamma_i), \quad k = 1, \dots, K.$$

- For each $\gamma_i \in \Gamma$, we calculate the standard error of $CV_1(\gamma_i), \dots, CV_K(\gamma_i)$ as:

$$SD(\gamma_i) = \sqrt{\text{Var}(CV_1(\gamma_i), \dots, CV_K(\gamma_i))}, \quad SE(\gamma_i) = SD(\gamma_i)/\sqrt{K}.$$

- Let

$$CV(\gamma) = \frac{1}{K} \sum_{k=1}^K CV_k(\gamma) \quad \text{and} \quad \hat{\gamma} = \arg \min_{\gamma \in \Gamma} CV(\gamma).$$

The 1SE rule ([Breiman et al., 1984](#); [Chen and Yang, 2021](#)) selects $\gamma_{1SE} \in \Gamma$ with as the one with the smallest $CV(\gamma)$ such that

$$CV(\gamma_{1SE}) \geq CV(\hat{\gamma}) + SE(\hat{\gamma}).$$

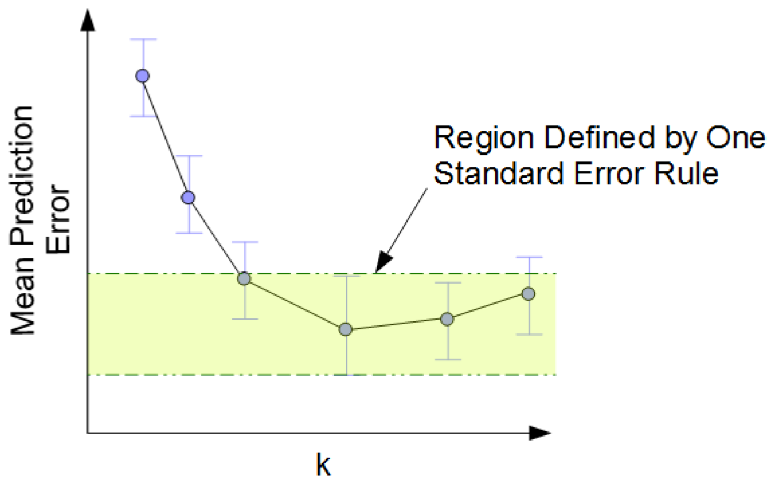


Figure 9: Illustration of the 1SE rule for selecting the model parameter.

The galaxy distribution is distorted along the line of sight due to the peculiar velocities of galaxies, *i.e.*, the so-called *finger-of-god* (Jackson, 1972) and *Kaiser* (Kaiser, 1987) effects.

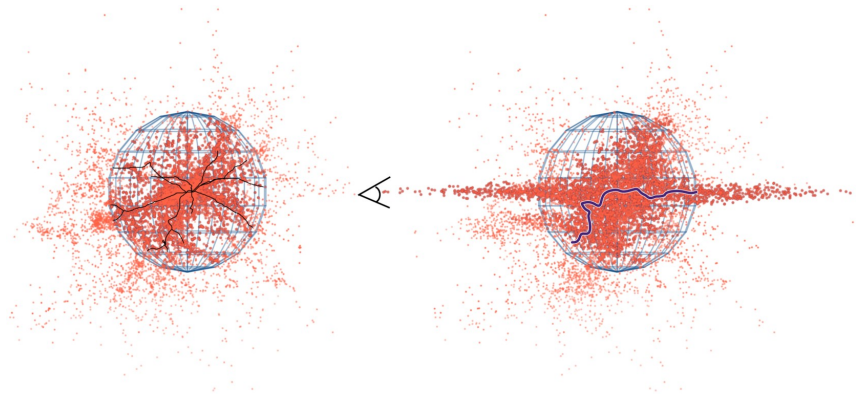


Figure 10: Redshift distortions along the line of sight (Kuchner et al., 2021).