

Efficient Inference on High-Dimensional Linear Models With Missing Outcomes

*Yikun Zhang*¹

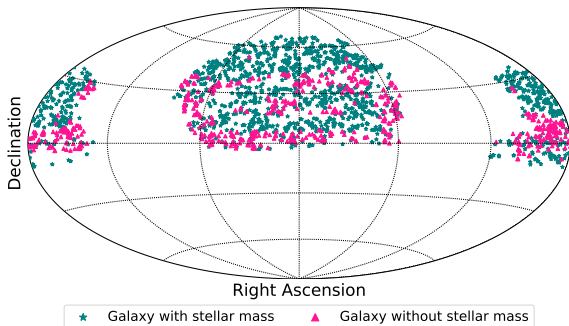
Joint Work with *Alexander Giessing*² and *Yen-Chi Chen*¹

¹Department of Statistics, University of Washington

²Department of Statistics and Data Science, National University of Singapore

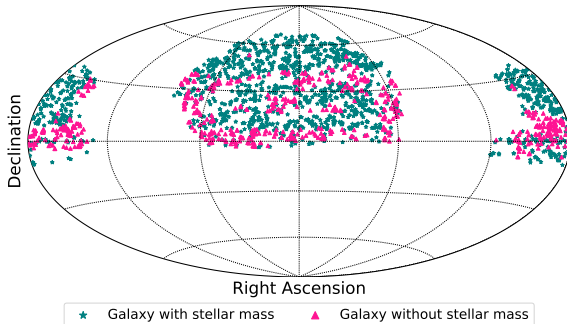
2024 Joint Statistical Meetings

August 6, 2024



Observed galaxies on the high redshift slice $0.4 \sim 0.401$.

► **Notes:** Sloan Digital Sky Survey (SDSS) observes millions of galaxies, but some (estimated) galactic stellar masses are missing in the associated value-added catalog ([Comparat et al., 2017](#)).



Observed galaxies on the high redshift slice $0.4 \sim 0.401$.

► **Scientific Question:**

How can we quantify the uncertainty of the (estimated) stellar mass of a newly observed galaxy based on the spectroscopic and photometric properties?

► **Notes:** Sloan Digital Sky Survey (SDSS) observes millions of galaxies, but some (estimated) galactic stellar masses are missing in the associated value-added catalog ([Comparat et al., 2017](#)).

► **High-dimensional Covariates:**

- Generate nonlinear features to capture complex patterns ([Chang et al., 2015](#); [Belloni et al., 2019](#)).

► **High-dimensional Covariates:**

- Generate nonlinear features to capture complex patterns ([Chang et al., 2015](#); [Belloni et al., 2019](#)).

► **Reasons for Missingness:**

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

► **High-dimensional Covariates:**

- Generate nonlinear features to capture complex patterns (Chang et al., 2015; Belloni et al., 2019).

► **Reasons for Missingness:**

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

► **Statistical Problem:**

How can we conduct valid and efficient inference on the regression function despite missing outcomes?

- ① **Linearity:** The data $\{(Y_i, R_i, X_i)\}_{i=1}^n$ 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 $\|\beta_0\|_0 = s_\beta \ll d$ and $R \in \{0, 1\}$ when Y is missing or not.

- ① **Linearity:** The data $\{(Y_i, R_i, X_i)\}_{i=1}^n$ 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 $\|\beta_0\|_0 = s_\beta \ll d$ and $R \in \{0, 1\}$ when Y is missing or not.

Relaxing the linearity assumption:

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

Our method can be generalized to handle heteroscedastic errors.

- ① **Linearity:** The data $\{(Y_i, R_i, X_i)\}_{i=1}^n$ 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 $\|\beta_0\|_0 = s_\beta \ll d$ and $R \in \{0, 1\}$ when Y is missing or not.

Relaxing the linearity assumption:

- Sparse additive model (Ravikumar et al., 2009);
- Partially linear model (Müller and van de Geer, 2015).

Our method can be generalized to handle heteroscedastic errors.

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

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

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

- ① **Fully Observed Outcomes:** Debiased Lasso (Zhang and Zhang, 2014; van de Geer et al., 2014; Javanmard and Montanari, 2014):

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

- $\hat{\beta}_\lambda$ is a Lasso solution under the regularization parameter $\lambda > 0$;
- $\hat{\Theta} \in \mathbb{R}^{d \times d}$ is an approximation to the matrix inverse $(\frac{1}{n} \sum_{i=1}^n X_i X_i^T)^{-1}$.

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

- 1 **Fully Observed Outcomes:** Debiased Lasso (Zhang and Zhang, 2014; van de Geer et al., 2014; Javanmard and Montanari, 2014):

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

- $\hat{\beta}_\lambda$ is a Lasso solution under the regularization parameter $\lambda > 0$;
 - $\hat{\Theta} \in \mathbb{R}^{d \times d}$ is an approximation to the matrix inverse $(\frac{1}{n} \sum_{i=1}^n X_i X_i^T)^{-1}$.
- 2 **MAR Outcomes:** M-estimation framework with a Lasso-type debiased and doubly robust estimator (Chakraborty et al., 2019).

- ▶ **Drawbacks of Existing Approaches:** Inference on $\beta_0 \in \mathbb{R}^d$.
- ① Need to compute a $d \times d$ debiasing matrix $\hat{\Theta}$.
- ② Require sample splitting or cross fitting for valid inference.

- ▶ **Drawbacks of Existing Approaches:** Inference on $\beta_0 \in \mathbb{R}^d$.
 - ① Need to compute a $d \times d$ debiasing matrix $\hat{\Theta}$.
 - ② Require sample splitting or cross fitting for valid inference.
- ▶ **Our Focus:** Inference on $m_0(x) = x^T \beta_0$.
 - *Computational efficiency:* Our debiasing program is convex and only needs to solve for an n -dimensional weight vector.
 - *Statistical efficiency:* Our estimator is semi-parametrically efficient among all asymptotically linear estimators.

Methodology and Asymptotic Theory

- The debiased Lasso estimator on the complete-case data is given by

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

- The debiased Lasso estimator on the complete-case data is given by

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

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

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

- The debiased Lasso estimator on the complete-case data is given by

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

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

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

- **Issue:** This naive estimator may not be asymptotically normal in general ([van de Geer et al., 2014](#); [Javanmard and Montanari, 2014](#))!

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

► **Idea:** Introduce a weight vector $\mathbf{w} = (w_1, \dots, w_n)^T \in \mathbb{R}^n$ to replace

$$\frac{1}{\sqrt{n}} x^T \widehat{\Theta} X_i \implies w_i \quad \text{for } i = 1, \dots, n$$

and formulate a generic debiased estimator

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

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

► **Idea:** Introduce a weight vector $\mathbf{w} = (w_1, \dots, w_n)^T \in \mathbb{R}^n$ to replace

$$\frac{1}{\sqrt{n}} x^T \widehat{\Theta} X_i \implies w_i \quad \text{for } i = 1, \dots, n$$

and formulate a generic debiased estimator

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

► **Question:** How do we estimate the weight vector $\mathbf{w} = (w_1, \dots, w_n)^T$?

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

$$\mathbb{E} \left[\left(\sqrt{n} m^{\text{debias}}(x; \boldsymbol{w}) - \sqrt{n} m_0(x) \right)^2 \mid X_1, \dots, X_n \right]$$

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

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

► **Notes:** $\pi_i := P(R_i = 1 \mid X_i)$ is the propensity score under the MAR condition.

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

- By Hölder's inequality,

$$\text{"Conditional Bias"} \leq \left[\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \pi_i X_i - x \right\|_\infty \sqrt{n} \|\beta_0 - \beta\|_1 \right]^2.$$

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

- By Hölder's inequality,

$$\text{"Conditional Bias"} \leq \left[\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \pi_i X_i - x \right\|_\infty \sqrt{n} \|\beta_0 - \beta\|_1 \right]^2.$$

- We design our debiasing program as:

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

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

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

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

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

- 2 Obtain consistent propensity score estimates $\widehat{\pi}_i, i = 1, \dots, n$ by *any machine learning method*.

- 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*.
- 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*.

- 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 $m_0(x) = x^T \beta$ 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 How to select the tuning parameter $\gamma > 0$ for our debiasing program?

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

- 2 Is our debiased estimator asymptotically normal?

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

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

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

- ② Is our debiased estimator asymptotically normal?

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

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

► **Primal Program:**

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

► **Primal Program:**

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

► **Dual Program:**

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

► **Primal Program:**

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

► **Dual Program:**

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

► **Primal-Dual Relation:** Under the strong duality,

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

► **Dual Program:**

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

- It is an *unconstrained* optimization problem, and $\gamma > 0$ can be fine-tuned via cross-validation.

► **Dual Program:**

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

- It is an *unconstrained* optimization problem, and $\gamma > 0$ can be fine-tuned via cross-validation.
- Primal-dual relation $\hat{w}_i = -\frac{1}{2\sqrt{n}} \cdot X_i^T \hat{\ell}$, $i = 1, \dots, n$ and dual consistency $\hat{\ell} \xrightarrow{P} \ell_0$ reveal that

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

Theorem (Theorem 7 in [Zhang et al. 2023](#))

Under regularity conditions,

$$\sqrt{n} \left[\widehat{m}^{\text{debias}}(x; \widehat{w}) - m_0(x) \right] \xrightarrow{d} \mathcal{N}(0, \sigma_m^2(x))$$

with $\sigma_m^2(x) = \lim_{n \rightarrow \infty} \sigma_\epsilon^2 \cdot x^T [\mathbb{E}(RXX^T)]^{-1} x$.

Theorem (Theorem 7 in [Zhang et al. 2023](#))

Under regularity conditions,

$$\sqrt{n} \left[\widehat{m}^{\text{debias}}(x; \widehat{w}) - m_0(x) \right] \xrightarrow{d} \mathcal{N} \left(0, \sigma_m^2(x) \right)$$

with $\sigma_m^2(x) = \lim_{n \rightarrow \infty} \sigma_\epsilon^2 \cdot x^T [\mathbb{E}(RXX^T)]^{-1} x$.

- 1 For any fixed dimension $d > 0$, the asymptotic variance

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

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

Theorem (Theorem 7 in [Zhang et al. 2023](#))

Under regularity conditions,

$$\sqrt{n} \left[\widehat{m}^{\text{debias}}(x; \widehat{w}) - m_0(x) \right] \xrightarrow{d} \mathcal{N} \left(0, \sigma_m^2(x) \right)$$

with $\sigma_m^2(x) = \lim_{n \rightarrow \infty} \sigma_\epsilon^2 \cdot x^T \left[\mathbb{E} (RXX^T) \right]^{-1} x$.

- ① For any fixed dimension $d > 0$, the asymptotic variance

$$\sigma_\epsilon^2 \cdot x^T \left[\mathbb{E} (RXX^T) \right]^{-1} x$$

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

- ② Under regularity conditions (Proposition 8 in [Zhang et al. 2023](#)),

$$\widehat{\sigma}_\epsilon^2 \sum_{i=1}^n \widehat{\pi}_i \widehat{w}_i^2 \xrightarrow{P} \sigma_\epsilon^2 \cdot x^T \left[\mathbb{E} (RXX^T) \right]^{-1} x.$$

Sample splitting or cross fitting is often required in debiased inference via machine learning methods ([Chernozhukov et al., 2018](#)).

Sample splitting or cross fitting is often required in debiased inference via machine learning methods ([Chernozhukov et al., 2018](#)).

- Why don't we need sample splitting or cross fitting for estimating the propensity score by any machine learning method?

Sample splitting or cross fitting is often required in debiased inference via machine learning methods ([Chernozhukov et al., 2018](#)).

- Why don't we need sample splitting or cross fitting for estimating the propensity score by any machine learning method?
- ▶ **Answer:** Our asymptotic normality result depends on the *in-sample* estimation error r_π of the propensity score:

$$\max_{1 \leq i \leq n} |\hat{\pi}_i - \pi_i| = O_P(r_\pi) \quad \text{with} \quad \pi_i = \pi(X_i), i = 1, \dots, n.$$

- Our debiased estimator performs even better when the estimated propensity scores on the training data are close to the true ones!!

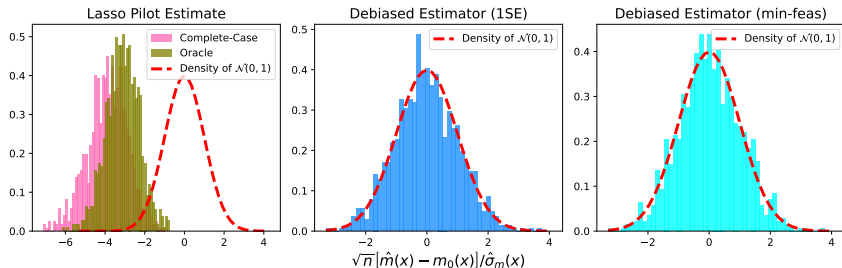
Sample splitting or cross fitting is often required in debiased inference via machine learning methods ([Chernozhukov et al., 2018](#)).

- Why don't we need sample splitting or cross fitting for estimating the propensity score by any machine learning method?
- ▶ **Answer:** Our asymptotic normality result depends on the *in-sample* estimation error r_π of the propensity score:

$$\max_{1 \leq i \leq n} |\hat{\pi}_i - \pi_i| = O_P(r_\pi) \quad \text{with} \quad \pi_i = \pi(X_i), i = 1, \dots, n.$$

- Our debiased estimator performs even better when the estimated propensity scores on the training data are close to the true ones!!
- This permits the use of complex machine learning methods with high learnability ([Steinwart, 2001](#); [Farrell et al., 2021](#); [Gao et al., 2022](#)).

Simulation and Real-World Application



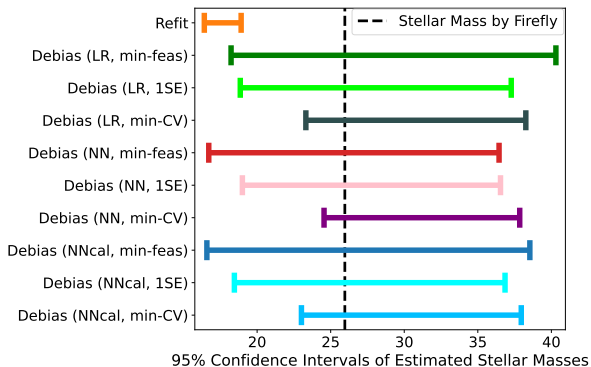
► **Effectiveness of Our Debiased Estimator:**

- Correct the bias of the Lasso pilot estimate.
- Asymptotically normal under a wide range of $\gamma > 0$.

► **Notes:** Our paper contains comprehensive comparisons with other existing methods.

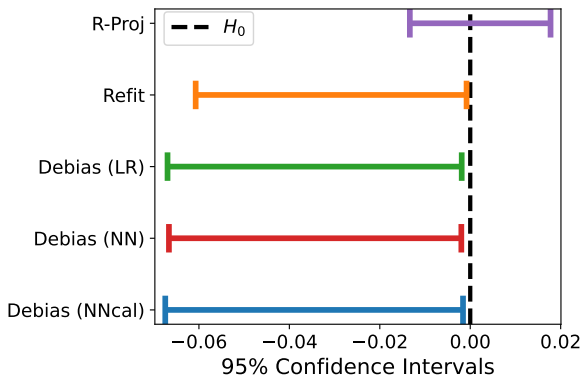
How can we quantify the uncertainty of the (estimated) stellar mass of a galaxy based on the spectroscopic and photometric properties?

How can we quantify the uncertainty of the (estimated) stellar mass of a galaxy based on the spectroscopic and photometric properties?



- The 95% confidence intervals by our debiasing methods cover the true stellar mass of a new galaxy.

Is it statistically significant that the stellar mass of a galaxy is negatively correlated with its distance to the nearby cosmic filament structures?



- 95% confidence intervals by our debiasing methods exclude 0 and are all negative.

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

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

- The dual form explains its computational and statistical efficiencies.

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

- The dual form explains its computational and statistical efficiencies.
- The nuisance propensity score can be nonparametrically estimated without sample splitting or cross fitting.

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

- The dual form explains its computational and statistical efficiencies.
- The nuisance propensity score can be nonparametrically estimated without sample splitting or cross fitting.
- A novel application to the inference on galactic stellar mass.

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

- The dual form explains its computational and statistical efficiencies.
- The nuisance propensity score can be nonparametrically estimated without sample splitting or cross fitting.
- A novel application to the inference on galactic stellar mass.

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](#).

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

- The dual form explains its computational and statistical efficiencies.
- The nuisance propensity score can be nonparametrically estimated without sample splitting or cross fitting.
- A novel application to the inference on galactic stellar mass.

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](#).

Thank you!

- 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.
- 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.
- 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.
- 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.
- M. H. Farrell, T. Liang, and S. Misra. Deep neural networks for estimation and inference. *Econometrica*, 89(1):181–213, 2021.

- 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.
- W. Gao, F. Xu, and Z.-H. Zhou. Towards convergence rate analysis of random forests for classification. *Artificial Intelligence*, 313:103788, 2022.
- J. Jackson. A critique of rees’s theory of primordial gravitational radiation. *Monthly Notices of the Royal Astronomical Society*, 156(1):1P–5P, 1972.
- 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.
- 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.
- I. Steinwart. On the influence of the kernel on the consistency of support vector machines. *Journal of machine learning research*, 2(Nov):67–93, 2001.
- 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.
- 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:

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

- 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

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

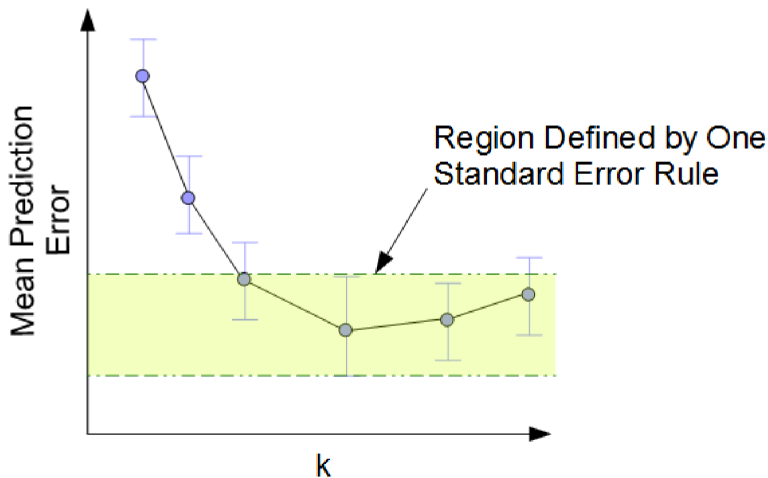


Figure: Illustration of the 1SE rule for selecting the model parameter.

- 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

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

- 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

$$\begin{aligned} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Xi_x(Y_i, R_i, X_i; m_0, \beta_0) \middle| \mathbf{X} \right] &= 0, \\ \sup_{\beta \in \mathcal{T}_n} \left| \left\{ \frac{\partial}{\partial \beta} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Xi_x(Y_i, R_i, X_i; m, \beta) \middle| \mathbf{X} \right] \right\}_{(m_0, \beta_0)} \right|^T (\beta - \beta_0) &\leq \frac{\delta_n}{\sqrt{n}}, \end{aligned} \quad (2)$$

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

- Both conditions in (2) hold true, because for any $\beta \in \mathcal{T}_n$ and some convex set \mathcal{B} containing β_0 , we have that

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

- Our debiasing program optimizes the (estimated) variance among all the estimators satisfying Neyman near-orthogonalization (2).
- (2) also allows our debiasing program to *de-correlate* the Lasso pilot regression from propensity score estimation and weight optimization.

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

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

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

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

- Dual relation $\widehat{w}_i = -\frac{1}{2\sqrt{n}} \cdot X_i^T \widehat{\ell}$ for $i = 1, \dots, n$ and dual consistency $\widehat{\ell} \xrightarrow{P} \ell_0$ reveal that

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

- 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

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

- 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:

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

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.,

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

Methods to be compared:

- “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](#)).

Implementation settings of the above methods:

- 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 over 1000 Monte Carlo experiments:

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

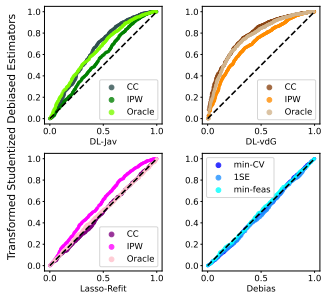
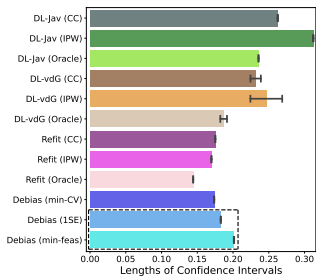
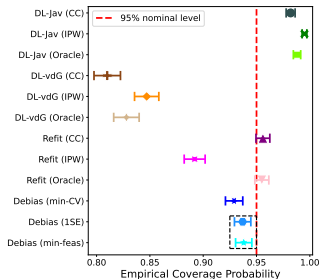
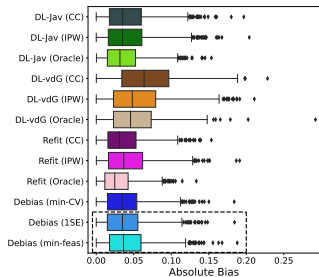


Figure: Sparse β_0^{sp} and sparse $x^{(2)}$ with $X_i \sim \mathcal{N}(\mathbf{0}, \Sigma^{cs})$, $i = 1, \dots, n$.

W Simulation Results Under Gaussian Noises (II)

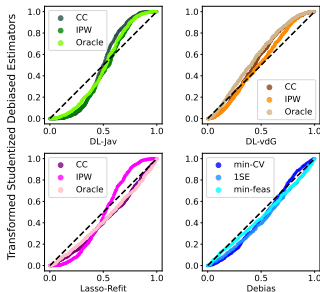
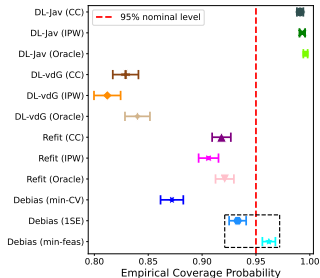
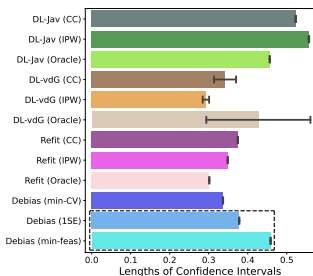
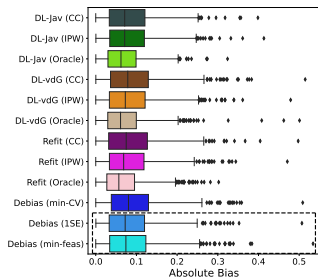


Figure: Pseudo-dense β_0^{pd} and sparse $x^{(2)}$ with $X_i \sim \mathcal{N}_d(\mathbf{0}, \Sigma^{\text{ar}})$, $i = 1, \dots, n$.

Simulation Results Under Laplace(0, 1/√2) Noises

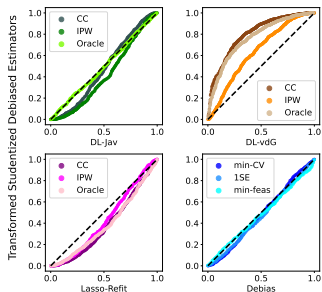
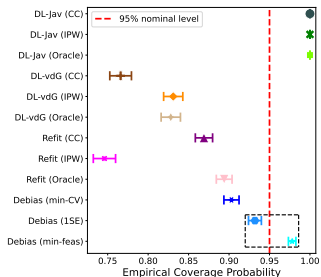
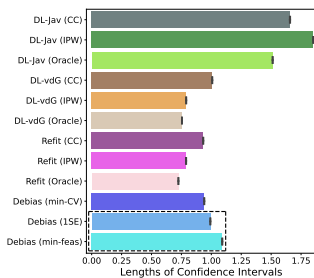
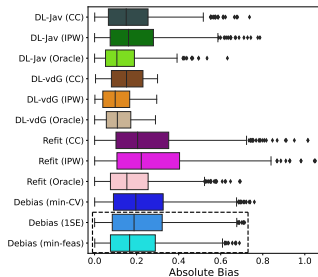


Figure: Dense β_0^{de} and sparse $x^{(2)}$ with $X_i \sim \mathcal{N}(\mathbf{0}, \Sigma^{cs})$, $i = 1, \dots, n$.

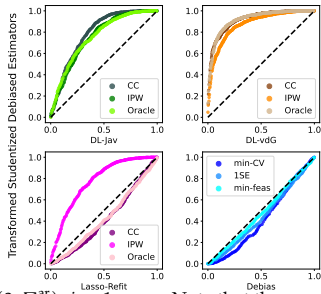
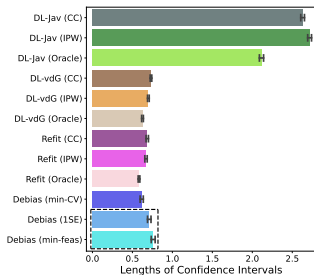
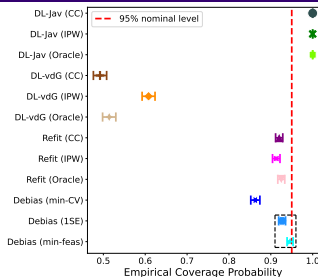
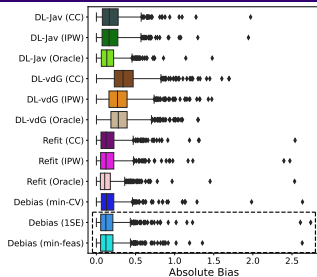
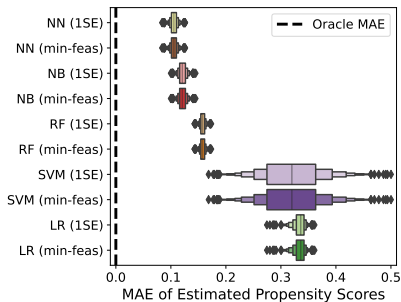
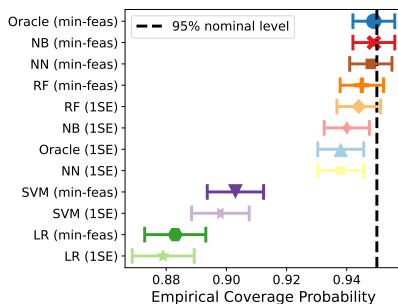
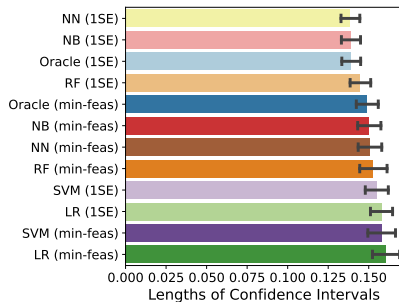
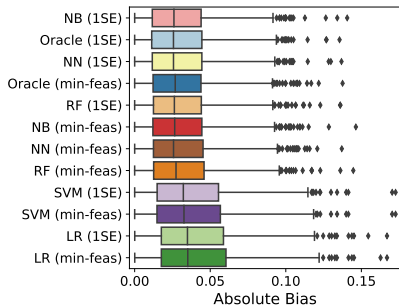
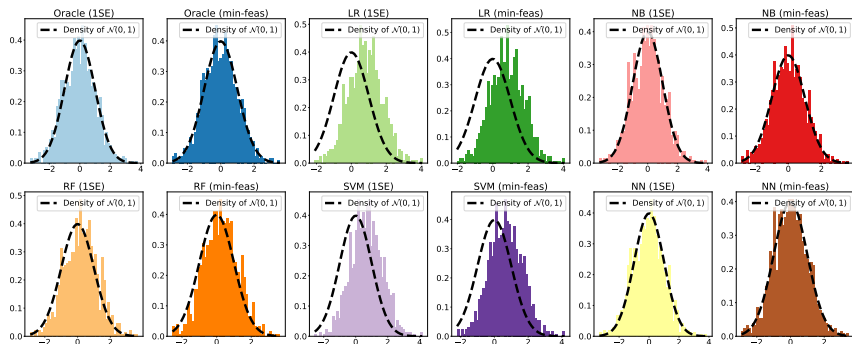


Figure: Pseudo-dense β_0^{pd} and dense $x^{(4)}$ with $X_i \sim \mathcal{N}_d(\mathbf{0}, \Sigma^{\text{ar}})$, $i = 1, \dots, n$. Note that the mean-zero t_2 distribution has *infinite* variance.

- 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 .
- 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.





$$\sqrt{n} [\hat{m}^{debias}(x) - m_0(x)] / \hat{\sigma}_m(x)$$

- 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 $d = 1409$.


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 $E[Y(1)|X, T = 1]$.

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	$Y(0)_n$


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

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

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

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

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

- The extended debiased estimator becomes

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

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


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 $E[Y(1)|X, T = 1]$.

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	$Y(0)_n$


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

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

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

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

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

- The extended debiased estimator becomes

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

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

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.

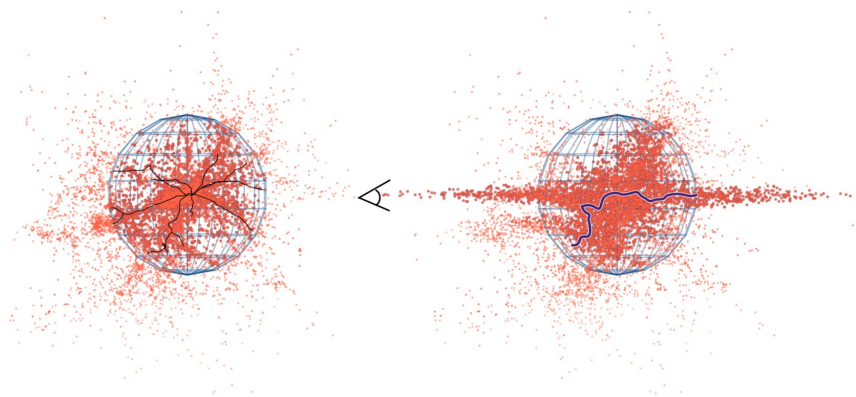


Figure: Redshift distortions along the line of sight (Kuchner et al., 2021).