

High Dimensional M -Estimation & Inference from Observational Data with Incomplete Responses

A Semi-Parametric Doubly Robust Framework

Abhishek Chakraborty¹

Department of Statistics
University of Pennsylvania

Group Meeting
April 24, 2019

¹Joint work with Jiarui Lu, T. Tony Cai and Hongzhe Li.

- Current era of 'big data' and data science \rightsquigarrow rapid influx of large **and high dimensional data** (easily available and computationally tractable).
- Rich information on multitudes of variables at the same place \rightsquigarrow many interesting scientific questions and also **unique statistical challenges!**

- Current era of 'big data' and data science \rightsquigarrow rapid influx of large and high dimensional data (easily available and computationally tractable).
- Rich information on multitudes of variables at the same place \rightsquigarrow many interesting scientific questions and also unique statistical challenges!
- One frequently encountered challenge: incompleteness of the data and in particular, (partial) missingness of the response of interest.
 - Reasons could be 'circumstantial' (e.g. practical constraints such as logistics, time, cost issues etc.), or it could be 'by design' (e.g. due to the 'treatment' assignment/non-assignment mechanism).
 - The response corresponding to a 'treatment' of interest could not be observed for a person who is not 'treated' (and vice versa).

- Current era of 'big data' and data science \rightsquigarrow rapid influx of large and high dimensional data (easily available and computationally tractable).
- Rich information on multitudes of variables at the same place \rightsquigarrow many interesting scientific questions and also unique statistical challenges!
- One frequently encountered challenge: incompleteness of the data and in particular, (partial) missingness of the response of interest.
 - Reasons could be 'circumstantial' (e.g. practical constraints such as logistics, time, cost issues etc.), or it could be 'by design' (e.g. due to the 'treatment' assignment/non-assignment mechanism).
 - The response corresponding to a 'treatment' of interest could not be observed for a person who is not 'treated' (and vice versa).
- Another complication in both cases: observational nature of the data. The missingness mechanism could be informative (not randomized)!

- **Observational data** \rightsquigarrow typically **informative missingness** (or treatment assignment) mechanism. Could **depend** on the person's covariates.
- Often termed **selection bias** or **treatment by indication** or **confounding (in causal inference)** in observational studies. **Has** to be factored in!
- **Need to account for the missingness** in a proper principled way **under minimal conditions** to **ensure valid, unbiased (and robust) inference**.

- **Observational data** \rightsquigarrow typically **informative missingness** (or treatment assignment) mechanism. Could **depend** on the person's covariates.
- Often termed **selection bias** or **treatment by indication** or **confounding (in causal inference)** in observational studies. **Has** to be factored in!
- Need to account for the **missingness** in a proper principled way **under minimal conditions** to **ensure valid, unbiased (and robust) inference**.
- **Relevance**: these issues **occur in virtually any** modern day large scale **observational study** arising in various scientific disciplines, including:
 - **Biomedical studies** (e.g. **electronic health records (EHR) data**); and **Integrative genomics** (e.g. **gene expression data and eQTL studies**).
 - Also **econometrics** (policy evaluation), **computer science**, **finance** etc.

- **Variables of interest:** outcome $Y \in \mathcal{Y} \subseteq \mathbb{R}$ and covariates $\mathbf{X} \in \mathcal{X} \subseteq \mathbb{R}^p$ (possibly high dimensional, compared to the sample size).
 - The supports \mathcal{Y} and \mathcal{X} of Y and \mathbf{X} need **not** be continuous.
- **Main issue:** Y may not always be observed. Let $T \in \{0, 1\}$ denote the indicator of the true Y being observed.
- The (partly) unobserved random vector (T, Y, \mathbf{X}) is assumed to be jointly defined on a common probability space with measure $\mathbb{P}(\cdot)$.

- **Variables of interest:** outcome $Y \in \mathcal{Y} \subseteq \mathbb{R}$ and covariates $\mathbf{X} \in \mathcal{X} \subseteq \mathbb{R}^p$ (possibly high dimensional, compared to the sample size).
 - The supports \mathcal{Y} and \mathcal{X} of Y and \mathbf{X} need **not** be continuous.
- **Main issue:** Y may not always be observed. Let $T \in \{0, 1\}$ denote the indicator of the true Y being observed.
- The (partly) unobserved random vector (T, Y, \mathbf{X}) is assumed to be jointly defined on a common probability space with measure $\mathbb{P}(\cdot)$.
- **Observable data:** $\mathcal{D}_n := \{\mathbf{Z}_i := (T_i, T_i Y_i, \mathbf{X}_i) : i = 1, \dots, n\} \stackrel{iid}{\sim} \mathbf{Z}$, where $\mathbf{Z} := (T, TY, \mathbf{X})$ whose distribution is defined via $\mathbb{P}(\cdot)$.
- **High dimensional setting:** p can diverge with n (including $p \gg n$).

- Generally applicable to **any missing data setting** - with **missing outcomes** Y and (possibly) **high dimensional covariates** \mathbf{X} .
- **Causal inference** problems (via 'potential' outcomes framework).

- Generally applicable to **any missing data setting** - with **missing outcomes** Y and (possibly) **high dimensional covariates** \mathbf{X} .
- **Causal inference** problems (via 'potential' outcomes framework).
 - Here, \mathbf{X} is often called '**confounders**' (for observational studies) or '**adjustment**' variables/features (for randomized trials).
 - **Usual set-up:** **binary 'treatment'** (a.k.a. exposure/intervention) assignment: $\mathcal{T} \in \{0, 1\}$, and **potential outcomes:** $\{Y_{(0)}, Y_{(1)}\}$.
 - **Observed outcome:** $\mathbb{Y} := Y_{(0)}1(\mathcal{T} = 0) + Y_{(1)}1(\mathcal{T} = 1)$, i.e. depending on \mathcal{T} , we **observe only one** of $\{Y_{(0)}, Y_{(1)}\}$.

- Generally applicable to **any missing data setting** - with **missing outcomes** Y and (possibly) **high dimensional covariates** \mathbf{X} .
- **Causal inference** problems (via 'potential' outcomes framework).
 - Here, \mathbf{X} is often called '**confounders**' (for observational studies) or '**adjustment**' variables/features (for randomized trials).
 - **Usual set-up:** **binary 'treatment'** (a.k.a. exposure/intervention) assignment: $\mathcal{T} \in \{0, 1\}$, and **potential outcomes:** $\{Y_{(0)}, Y_{(1)}\}$.
 - **Observed outcome:** $\mathbb{Y} := Y_{(0)}1(\mathcal{T} = 0) + Y_{(1)}1(\mathcal{T} = 1)$, i.e. depending on \mathcal{T} , we **observe only one** of $\{Y_{(0)}, Y_{(1)}\}$.
 - For **each** $j \in \{0, 1\}$, this set-up is **included** based on the '**map**':

$$(T, Y, \mathbf{X}) \leftarrow (T_j, Y_{(j)}, \mathbf{X}) \text{ with } T_j := 1(\mathcal{T} = j) \quad \forall j \in \{0, 1\}.$$

- Generally applicable to **any missing data setting** - with **missing outcomes** Y and (possibly) **high dimensional covariates** \mathbf{X} .
- **Causal inference** problems (via 'potential' outcomes framework).
 - Here, \mathbf{X} is often called '**confounders**' (for observational studies) or '**adjustment**' variables/features (for randomized trials).
 - **Usual set-up:** **binary 'treatment'** (a.k.a. exposure/intervention) assignment: $\mathcal{T} \in \{0, 1\}$, and **potential outcomes:** $\{Y_{(0)}, Y_{(1)}\}$.
 - **Observed outcome:** $\mathbb{Y} := Y_{(0)}1(\mathcal{T} = 0) + Y_{(1)}1(\mathcal{T} = 1)$, i.e. depending on \mathcal{T} , we **observe only one** of $\{Y_{(0)}, Y_{(1)}\}$.
 - For **each** $j \in \{0, 1\}$, this set-up is **included** based on the 'map':

$$(T, Y, \mathbf{X}) \leftarrow (T_j, Y_{(j)}, \mathbf{X}) \text{ with } T_j := 1(\mathcal{T} = j) \quad \forall j \in \{0, 1\}.$$

The case of any multi-category treatment also similarly included.

1 Ignorability assumption: $T \perp\!\!\!\perp Y \mid \mathbf{X}$.

- A.k.a. 'missing at random' (MAR) in the missing data literature.
- A.k.a. 'no unmeasured confounding' (NUC) in causal inference.
- Special case: $T \perp\!\!\!\perp (Y, \mathbf{X})$. A.k.a. missing completely at random (MCAR) in missing data literature, and complete randomization (e.g. randomized trials) in causal inference (CI) literature.

1 Ignorability assumption: $T \perp\!\!\!\perp Y \mid \mathbf{X}$.

- A.k.a. 'missing at random' (MAR) in the missing data literature.
- A.k.a. 'no unmeasured confounding' (NUC) in causal inference.
- Special case: $T \perp\!\!\!\perp (Y, \mathbf{X})$. A.k.a. missing completely at random (MCAR) in missing data literature, and complete randomization (e.g. randomized trials) in causal inference (CI) literature.

2 Positivity assumption (a.k.a. 'sufficient overlap' in CI literature):

- Let $\pi(\mathbf{X}) := \mathbb{P}(T = 1 \mid \mathbf{X})$ be the propensity score (PS), and let $\pi_0 := \mathbb{P}(T = 1)$. Then, $\pi(\cdot)$ is uniformly bounded away from 0:

$$1 \geq \pi(\mathbf{x}) \geq \delta_\pi > 0 \quad \forall \mathbf{x} \in \mathcal{X}, \text{ for some constant } \delta_\pi > 0.$$

- **Rich resources** of data for **discovery research**; fast **growing literature**.



Review Article | Published: 18 May 2011

Using electronic health records to drive discovery in disease genomics

Isaac S. Kohane

Nature Reviews Genetics **12**, 417–428 (2011) |



Review Article | Published: 02 May 2012

Mining electronic health records: towards better research applications and clinical care

Peter B. Jensen, Lars J. Jensen & Søren Brunak

Nature Reviews Genetics **13**, 395–405 (2012) |

- **Rich resources** of data for **discovery research**; fast **growing literature**.



Review Article | Published: 18 May 2011

Using electronic health records to drive discovery in disease genomics

Isaac S. Kohane

Nature Reviews Genetics 12, 417–428 (2011) |

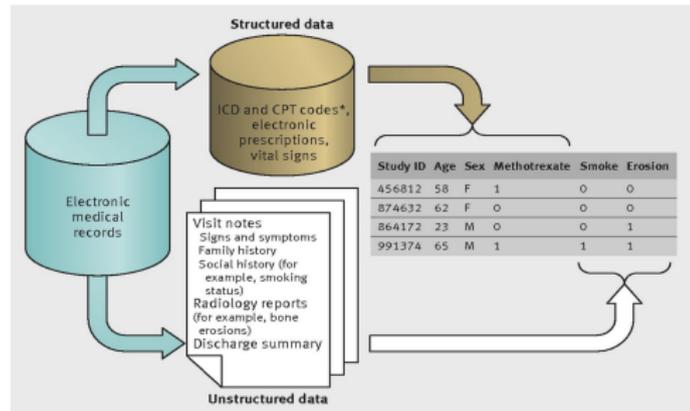


Review Article | Published: 02 May 2012

Mining electronic health records: towards better research applications and clinical care

Peter B. Jensen, Lars J. Jensen & Søren Brunak

Nature Reviews Genetics 13, 395–405 (2012) |



- **Detailed** clinical and phenotypic **data collected electronically** for large **patient cohorts**, as part of routine health care delivery.

- **Rich resources** of data for **discovery research**; fast **growing literature**.



Review Article | Published: 18 May 2011

Using electronic health records to drive discovery in disease genomics

Isaac S. Kohane

Nature Reviews Genetics 12, 417–428 (2011) |

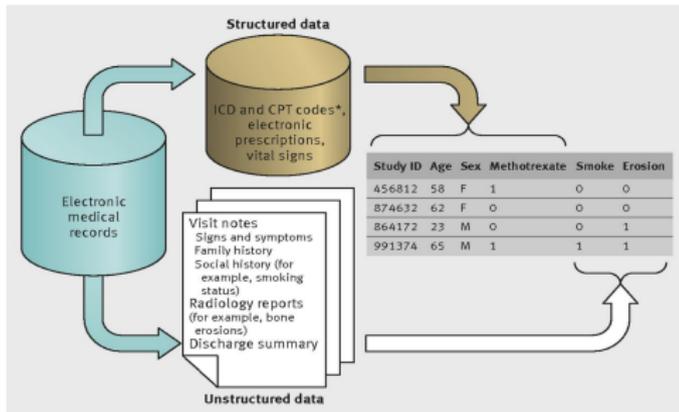


Review Article | Published: 02 May 2012

Mining electronic health records: towards better research applications and clinical care

Peter B. Jensen, Lars J. Jensen & Søren Brunak

Nature Reviews Genetics 13, 395–405 (2012) |



- **Detailed** clinical and phenotypic **data collected electronically for large patient cohorts**, as part of routine health care delivery.
- **Structured data**: ICD codes, medications, lab tests, demographics etc.
- **Unstructured text data** (extracted from clinician notes via NLP): signs and symptoms, family history, social history, radiology reports etc.

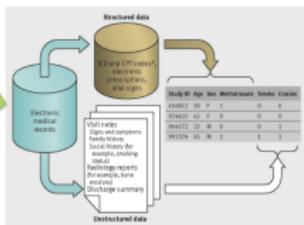
- Information on a **variety** of phenotypes (unlike usual cohort studies).
- Opens up **unique opportunities for** novel **integrative analyses**.

EHR Data: The Promises and the Challenges

- Information on a **variety** of phenotypes (unlike usual cohort studies).
- Opens up **unique opportunities** for novel **integrative analyses**.

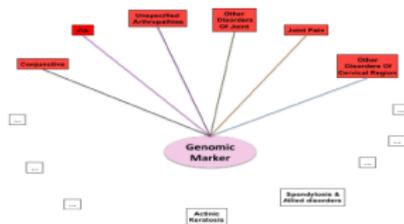


Bio-repository



EMR

- **EHR + Bio-repositories** \rightsquigarrow genome-phenome association networks, **PheWAS studies** and genomic risk prediction of diseases.

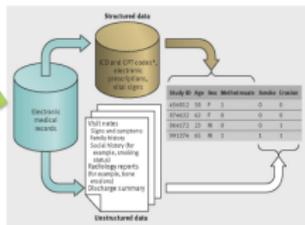


EHR Data: The Promises and the Challenges

- Information on a **variety** of phenotypes (unlike usual cohort studies).
- Opens up **unique opportunities** for novel **integrative analyses**.

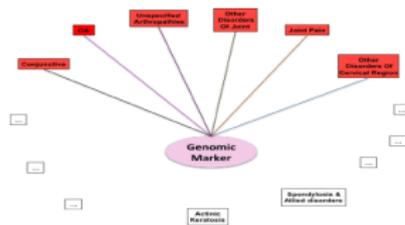


Bio-repository



EMR

- **EHR + Bio-repositories** \rightsquigarrow genome-phenome association networks, **PheWAS studies** and genomic risk prediction of diseases.



- The key **challenges and bottlenecks** for EHR driven research:
 - **Logistic difficulty** in obtaining **validated phenotype (Y)** information.
 - Often **time/labor/cost intensive** (and the ICD codes are imprecise).

- Some examples of missing Y in EHRs and the reason for missingness:

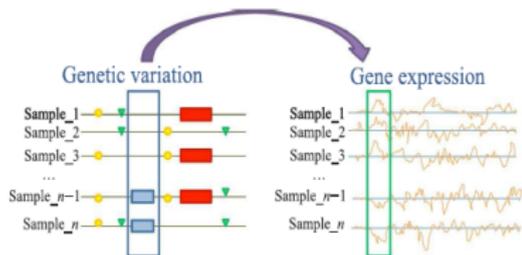
- Some examples of missing Y in EHRs and the reason for missingness:
 - ① $Y \rightsquigarrow$ some (binary) disease phenotype (e.g. Rheumatoid Arthritis). Requires manual chart review by physicians (logistic constraints).
 - ② $Y \rightsquigarrow$ some biomarker (e.g. anti-CCP, an important RA biomarker). Requires lab tests (cost constraints). Similarly, any Y requiring genomic measurements may also have cost/logistics constraints.

- Some examples of missing Y in EHRs and the reason for missingness:
 - ① $Y \rightsquigarrow$ some (binary) disease phenotype (e.g. Rheumatoid Arthritis). Requires manual chart review by physicians (logistic constraints).
 - ② $Y \rightsquigarrow$ some biomarker (e.g. anti-CCP, an important RA biomarker). Requires lab tests (cost constraints). Similarly, any Y requiring genomic measurements may also have cost/logistics constraints.
- Verified phenotypes/treatment response/biomarkers/genomic vars (Y) available **only** for a subset. Clinical features (X) available for **all**.
- Further issues: selection bias/treatment by indication/preferential labeling (e.g. sicker patients get labeled/treated/tested more often).

- Some examples of missing Y in EHRs and the reason for missingness:
 - ① $Y \rightsquigarrow$ some (binary) disease phenotype (e.g. Rheumatoid Arthritis). Requires manual chart review by physicians (logistic constraints).
 - ② $Y \rightsquigarrow$ some biomarker (e.g. anti-CCP, an important RA biomarker). Requires lab tests (cost constraints). Similarly, any Y requiring genomic measurements may also have cost/logistics constraints.
- Verified phenotypes/treatment response/biomarkers/genomic vars (Y) available **only** for a subset. Clinical features (X) available for **all**.
- Further issues: selection bias/treatment by indication/preferential labeling (e.g. sicker patients get labeled/treated/tested more often).
- Causal inference problems (treatment effects estimation): EHRs also facilitate comparative effectiveness research on a large scale.
 - Many treatments/medications (and responses) being observed. All other clinical features (X) serve as potential confounders.

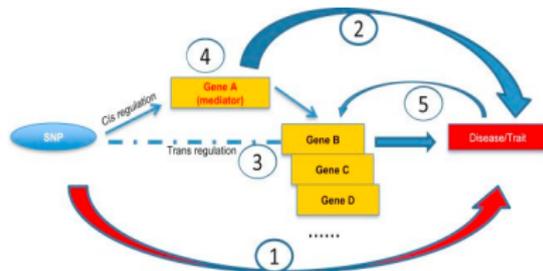
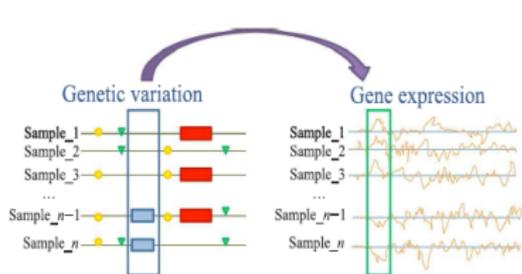
Another Example: eQTL Studies (Integrative Genomics)

- Association studies for gene expression (Y) vs. genetic variants (X).



Another Example: eQTL Studies (Integrative Genomics)

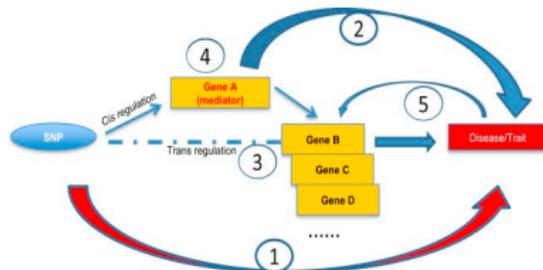
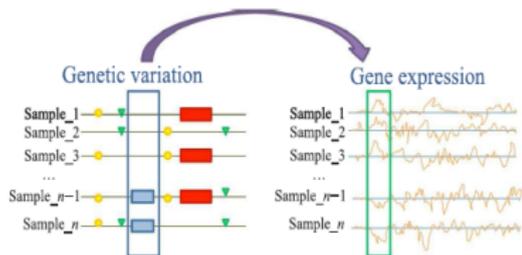
- Association studies for gene expression (Y) vs. genetic variants (X).



- Popular tools in **integrative genomics** (genetic association studies + gene expression profiling) for understanding **gene regulatory networks**.

Another Example: eQTL Studies (Integrative Genomics)

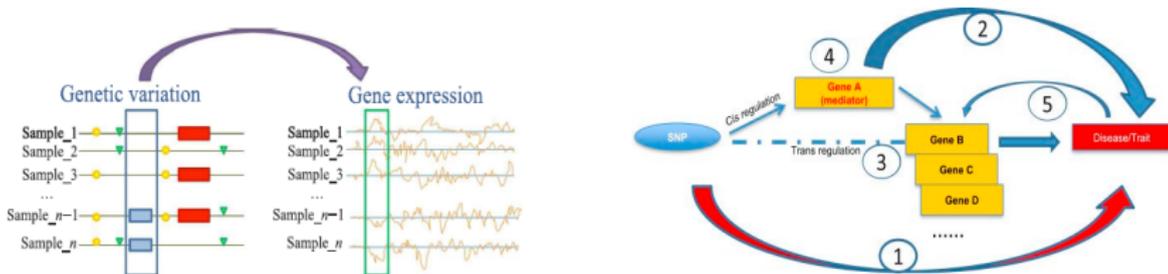
- Association studies for gene expression (Y) vs. genetic variants (X).



- Popular tools in **integrative genomics** (genetic association studies + gene expression profiling) for understanding **gene regulatory networks**.
- **Missing data issue:** gene expression data often missing (loss of power), while genetic variants data often available for a much larger group.

Another Example: eQTL Studies (Integrative Genomics)

- Association studies for gene expression (Y) vs. genetic variants (X).



- Popular tools in **integrative genomics** (genetic association studies + gene expression profiling) for understanding **gene regulatory networks**.
- **Missing data issue**: gene expression data often missing (loss of power), while genetic variants data often available for a much larger group.
- **Causal inference**: estimate the causal effect of any one variant (the 'treatment') on Y while all other variants are potential confounders.

- **Goal for M -estimation:** estimation **and** inference, based on \mathcal{D}_n , of $\theta_0 \in \mathbb{R}^d$ (possibly high dimensional), defined as the risk minimizer:

$$\theta_0 \equiv \theta_0(\mathbb{P}) := \underset{\theta \in \mathbb{R}^d}{\operatorname{arg\,min}} R(\theta), \text{ where } R(\theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)\} \text{ and}$$

$L(\cdot) \in \mathbb{R}^+$ is any 'loss' function that is convex and differentiable in θ . Existence of θ_0 implicitly assumed (guaranteed for most usual probs).

- d can diverge with n (including $d \gg n$). Also, $\theta_0(\mathbb{P})$ is 'model free' (no restrictions on \mathbb{P}). In particular, **no** model assumptions on $Y|\mathbf{X}$.

- **Goal for M -estimation:** estimation **and** inference, based on \mathcal{D}_n , of $\theta_0 \in \mathbb{R}^d$ (possibly high dimensional), defined as the risk minimizer:

$$\theta_0 \equiv \theta_0(\mathbb{P}) := \arg \min_{\theta \in \mathbb{R}^d} R(\theta), \text{ where } R(\theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)\} \text{ and}$$

$L(\cdot) \in \mathbb{R}^+$ is any 'loss' function that is convex and differentiable in θ . Existence of θ_0 implicitly assumed (guaranteed for most usual probs).

- d can diverge with n (including $d \gg n$). Also, $\theta_0(\mathbb{P})$ is 'model free' (no restrictions on \mathbb{P}). In particular, **no** model assumptions on $Y|\mathbf{X}$.
- The **key challenges:** the **missingness** via T (if not accounted for, the estimator **will be inconsistent!**) and the **high dimensional setting**.
- **Need suitable methods** - involves estimation of nuisance functions and careful analyses (due to error terms with complex dependencies).

- **Goal for M -estimation:** estimation **and** inference, based on \mathcal{D}_n , of $\theta_0 \in \mathbb{R}^d$ (possibly high dimensional), defined as the risk minimizer:

$$\theta_0 \equiv \theta_0(\mathbb{P}) := \arg \min_{\theta \in \mathbb{R}^d} R(\theta), \text{ where } R(\theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)\} \text{ and}$$

$L(\cdot) \in \mathbb{R}^+$ is any 'loss' function that is convex and differentiable in θ . Existence of θ_0 implicitly assumed (guaranteed for most usual probs).

- d can diverge with n (including $d \gg n$). Also, $\theta_0(\mathbb{P})$ is 'model free' (no restrictions on \mathbb{P}). In particular, **no** model assumptions on $Y|\mathbf{X}$.
- The **key challenges:** the **missingness** via T (if not accounted for, the estimator **will be inconsistent!**) and the **high dimensional setting**.
- **Need suitable methods** - involves estimation of nuisance functions and careful analyses (due to error terms with complex dependencies).
- **Special (but low- d) case:** $\theta_0 = \mathbb{E}(Y)$ and $L(Y, \mathbf{X}, \theta) = (Y - \theta)^2$. Leads to the **average treatment effect (ATE) estimation** prob in CI.

- The framework includes a broad class of M/Z -estimation problems.
- M -estimation for fully observed data: well studied with rich literature. Classical settings: Van der Vaart (2000); High dimensional settings: Negahban et al. (2012), Loh and Wainwright (2012, 2015) etc.

- The framework includes a broad class of M/Z -estimation problems.
- M -estimation for fully observed data: well studied with rich literature. Classical settings: Van der Vaart (2000); High dimensional settings: Negahban et al. (2012), Loh and Wainwright (2012, 2015) etc.
- Missing data/causal inference problems: semi-parametric inference.
 - Classical settings: vast literature (typically for mean estimation). Tsiatis (2007); Bang and Robins (2005); Robins et al. (1994) etc.

- The framework includes a broad class of M/Z -estimation problems.
- M -estimation for fully observed data: well studied with rich literature. Classical settings: Van der Vaart (2000); High dimensional settings: Negahban et al. (2012), Loh and Wainwright (2012, 2015) etc.
- Missing data/causal inference problems: semi-parametric inference.
 - Classical settings: vast literature (typically for mean estimation). Tsiatis (2007); Bang and Robins (2005); Robins et al. (1994) etc.
 - High dimensional settings (but low dimensional parameters): lot of attention in recent times on mean (or ATE) estimation. Belloni et al. (2014, 2017); Farrell (2015); Chernozhukov et al. (2018).

- The framework includes a broad class of M/Z -estimation problems.
- M -estimation for fully observed data: well studied with rich literature. Classical settings: Van der Vaart (2000); High dimensional settings: Negahban et al. (2012), Loh and Wainwright (2012, 2015) etc.
- Missing data/causal inference problems: semi-parametric inference.
 - Classical settings: vast literature (typically for mean estimation). Tsiatis (2007); Bang and Robins (2005); Robins et al. (1994) etc.
 - High dimensional settings (but low dimensional parameters): lot of attention in recent times on mean (or ATE) estimation. Belloni et al. (2014, 2017); Farrell (2015); Chernozhukov et al. (2018).
- Much less attention when the parameter itself is high dimensional.

- The framework includes a broad class of M/Z -estimation problems.
- M -estimation for fully observed data: well studied with rich literature. Classical settings: Van der Vaart (2000); High dimensional settings: Negahban et al. (2012), Loh and Wainwright (2012, 2015) etc.
- Missing data/causal inference problems: semi-parametric inference.
 - Classical settings: vast literature (typically for mean estimation). Tsiatis (2007); Bang and Robins (2005); Robins et al. (1994) etc.
 - High dimensional settings (but low dimensional parameters): lot of attention in recent times on mean (or ATE) estimation. Belloni et al. (2014, 2017); Farrell (2015); Chernozhukov et al. (2018).
- Much less attention when the parameter itself is high dimensional.
- This work contributes to both literature above: M -estimation + missing data + high dimensional setting and parameter. (Also has applications in heterogeneous treatment effects estimation in CI).

- 1 All standard **high dimensional (HD) regression** problems **with: (a) missing outcomes** and **(b) potentially misspecified (working) models**.

- ① All standard **high dimensional (HD) regression** problems with: (a) **missing outcomes** and (b) **potentially misspecified (working) models**.
- E.g. **squared loss**: $L(Y, \mathbf{X}, \boldsymbol{\theta}) := (Y - \mathbf{X}'\boldsymbol{\theta})^2 \rightsquigarrow$ **linear regression**;
logistic loss: $L(Y, \mathbf{X}, \boldsymbol{\theta}) := \log\{1 + \exp(\mathbf{X}'\boldsymbol{\theta})\} - Y(\mathbf{X}'\boldsymbol{\theta}) \rightsquigarrow$ **logistic regression** (for binary Y), exponential loss (Poisson reg.) so on
 - Note: throughout, **regardless** of any motivating 'working model' being true or not, the **definition of $\boldsymbol{\theta}_0$ is completely 'model free'**.

- 1 All standard **high dimensional (HD) regression** problems with: (a) missing outcomes and (b) potentially misspecified (working) models.
 - E.g. **squared loss**: $L(Y, \mathbf{X}, \boldsymbol{\theta}) := (Y - \mathbf{X}'\boldsymbol{\theta})^2 \rightsquigarrow$ linear regression;
logistic loss: $L(Y, \mathbf{X}, \boldsymbol{\theta}) := \log\{1 + \exp(\mathbf{X}'\boldsymbol{\theta})\} - Y(\mathbf{X}'\boldsymbol{\theta}) \rightsquigarrow$ logistic regression (for binary Y), exponential loss (Poisson reg.) so on . . .
 - Note: throughout, **regardless** of any motivating 'working model' being true or not, the **definition of $\boldsymbol{\theta}_0$ is completely 'model free'**.
- 2 **Series estimation** problems (model free) with missing Y and HD basis functions (instead of \mathbf{X} in Example 1 above). E.g. spline bases.
 - Use the **same choices of $L(\cdot)$** as in Example 1 above with \mathbf{X} replaced by any set of d (possibly HD) **basis functions $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_j(\mathbf{X})\}_{j=1}^d$** .
 - E.g. **polynomial bases**: $\boldsymbol{\Psi}(\mathbf{X}) := \{1, \mathbf{x}_j^k : 1 \leq j \leq p, 1 \leq k \leq d_0\}$. ($d_0 = 1 \rightsquigarrow$ linear bases as in Example 1; $d_0 = 3 \rightsquigarrow$ cubic splines).

Another Application: HD Single Index Models (SIMs)

- Signal recovery in high dimensional single index models (SIMs) with elliptically symmetric design distribution (e.g. \mathbf{X} is Gaussian).
- Let $Y = f(\beta_0' \mathbf{X}, \epsilon)$ with $f : \mathbb{R}^2 \rightarrow \mathcal{Y}$ **unknown** (i.e. β_0 identifiable **only** upto scalar multiples) and $\epsilon \perp\!\!\!\perp \mathbf{X}$ (i.e., $Y \perp\!\!\!\perp \mathbf{X} \mid \beta_0' \mathbf{X}$).

- Signal recovery in high dimensional single index models (SIMs) with elliptically symmetric design distribution (e.g. \mathbf{X} is Gaussian).
- Let $Y = f(\beta_0' \mathbf{X}, \epsilon)$ with $f : \mathbb{R}^2 \rightarrow \mathcal{Y}$ **unknown** (i.e. β_0 identifiable **only** upto scalar multiples) and $\epsilon \perp\!\!\!\perp \mathbf{X}$ (i.e., $Y \perp\!\!\!\perp \mathbf{X} \mid \beta_0' \mathbf{X}$).
- Consider **any** of the regression problems introduced in Example 1.
 - Let $\theta_0 := \arg \min_{\theta \in \mathbb{R}^p} \mathbb{E}\{L(Y, \mathbf{X}'\theta)\}$ for any convex loss function $L(\cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}$ (convex in the second argument). Then, $\theta_0 \propto \beta_0!$
 - A remarkable result due to Li and Duan (1989).

Another Application: HD Single Index Models (SIMs)

- **Signal recovery in high dimensional single index models (SIMs)** with elliptically symmetric design distribution (e.g. \mathbf{X} is Gaussian).
- Let $Y = f(\beta_0' \mathbf{X}, \epsilon)$ with $f : \mathbb{R}^2 \rightarrow \mathcal{Y}$ **unknown** (i.e. β_0 identifiable **only** upto scalar multiples) and $\epsilon \perp\!\!\!\perp \mathbf{X}$ (i.e., $Y \perp\!\!\!\perp \mathbf{X} \mid \beta_0' \mathbf{X}$).
- Consider **any** of the regression problems introduced in Example 1.
 - Let $\theta_0 := \arg \min_{\theta \in \mathbb{R}^p} \mathbb{E}\{L(Y, \mathbf{X}'\theta)\}$ for any convex loss function $L(\cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}$ (convex in the second argument). Then, $\theta_0 \propto \beta_0!$
 - A remarkable result due to Li and Duan (1989).
- **Classic example** of a misspecified parametric model defining θ_0 , yet θ_0 directly relates to an actual (interpretable) semi-parametric model!
 - The proportionality result also preserves any sparsity assumptions.

- **Applications** of all these problems **in causal inference** (estimation of treatment effects with useful applications in **precision medicine**):

- Applications of all these problems in causal inference (estimation of treatment effects with useful applications in precision medicine):
 - 1 Linear heterogeneous treatment effects estimation: application of the linear regression example (twice). Write $\{Y_{(0)}, Y_{(1)}\}$ linearly as:

$$Y_{(j)} = \mathbf{X}'\boldsymbol{\beta}_{(j)} + \epsilon_{(j)}, \quad \mathbb{E}(\epsilon_{(j)}\mathbf{X}) = \mathbf{0} \quad \forall j = 0, 1, \quad \text{so that}$$

$$Y_{(1)} - Y_{(0)} = \mathbf{X}'\boldsymbol{\beta}^* + \epsilon^*, \quad \boldsymbol{\beta}^* := \boldsymbol{\beta}_{(1)} - \boldsymbol{\beta}_{(0)}, \quad \epsilon^* := \epsilon_{(1)} - \epsilon_{(0)}.$$

- Applications of all these problems in causal inference (estimation of treatment effects with useful applications in precision medicine):
 - 1 Linear heterogeneous treatment effects estimation: application of the linear regression example (twice). Write $\{Y_{(0)}, Y_{(1)}\}$ linearly as:

$$Y_{(j)} = \mathbf{X}'\beta_{(j)} + \epsilon_{(j)}, \quad \mathbb{E}(\epsilon_{(j)}\mathbf{X}) = \mathbf{0} \quad \forall j = 0, 1, \quad \text{so that}$$

$$Y_{(1)} - Y_{(0)} = \mathbf{X}'\beta^* + \epsilon^*, \quad \beta^* := \beta_{(1)} - \beta_{(0)}, \quad \epsilon^* := \epsilon_{(1)} - \epsilon_{(0)}.$$

β^* denotes the (model free) linear projection of $Y_{(1)} - Y_{(0)} | \mathbf{X}$. Of interest in HD settings when $\mathbb{E}\{Y_{(1)} - Y_{(0)} | \mathbf{X}\}$ is difficult to model (Chernozhukov et al., 2017; Chernozhukov and Semenova, 2017).

- **Applications** of all these problems **in causal inference** (estimation of treatment effects with useful applications in **precision medicine**):
 - 1 **Linear heterogeneous treatment effects estimation: application of the linear regression example (twice).** Write $\{Y_{(0)}, Y_{(1)}\}$ linearly as:

$$Y_{(j)} = \mathbf{X}'\beta_{(j)} + \epsilon_{(j)}, \quad \mathbb{E}(\epsilon_{(j)}\mathbf{X}) = \mathbf{0} \quad \forall j = 0, 1, \quad \text{so that}$$

$$Y_{(1)} - Y_{(0)} = \mathbf{X}'\beta^* + \epsilon^*, \quad \beta^* := \beta_{(1)} - \beta_{(0)}, \quad \epsilon^* := \epsilon_{(1)} - \epsilon_{(0)}.$$

β^* denotes the (model free) **linear projection of $Y_{(1)} - Y_{(0)} | \mathbf{X}$** . **Of interest in HD settings** when $\mathbb{E}\{Y_{(1)} - Y_{(0)} | \mathbf{X}\}$ is difficult to model (Chernozhukov et al., 2017; Chernozhukov and Semenova, 2017).

- 2 **Average conditional treatment effects (ACTE) estimation via series estimators: application of the series estimation example (twice).**
- 3 **Causal inference via SIMs** (signal recovery, ACTE estimation and ATE estimation): **application of the SIM example (twice).**

- **Some notations:** $m(\mathbf{X}) := \mathbb{E}(Y|\mathbf{X})$ and $\phi(\mathbf{X}, \theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)|\mathbf{X}\}$.

- Some notations: $m(\mathbf{X}) := \mathbb{E}(Y|\mathbf{X})$ and $\phi(\mathbf{X}, \theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)|\mathbf{X}\}$.
- It is generally **necessary** to 'account' for the missingness in Y . The 'complete case' estimator of θ_0 in general will be **inconsistent!**

- **Some notations:** $m(\mathbf{X}) := \mathbb{E}(Y|\mathbf{X})$ and $\phi(\mathbf{X}, \theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)|\mathbf{X}\}$.
- It is generally **necessary** to 'account' for the missingness in Y . The 'complete case' estimator of θ_0 in general will be **inconsistent!**
 - That estimator may be consistent **only if:** (1) $\nabla\phi(\mathbf{X}, \theta_0) = \mathbf{0}$ a.s. for **every** \mathbf{X} (for regression problems, this indicates the 'correct model' case), and/or (2) $T \perp\!\!\!\perp (Y, X)$ (i.e. the MCAR case).
 - Illustration of (1) for sq. loss: $\nabla\phi(\mathbf{X}, \theta_0) = \mathbb{E}\{\mathbf{X}(Y - \mathbf{X}'\theta_0)|\mathbf{X}\} = \mathbf{0}$. Hence, $\mathbb{E}(Y|\mathbf{X}) = \mathbf{X}'\theta_0$ (i.e. a 'linear model' holds for $Y|\mathbf{X}$).

- **Some notations:** $m(\mathbf{X}) := \mathbb{E}(Y|\mathbf{X})$ and $\phi(\mathbf{X}, \theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)|\mathbf{X}\}$.
- It is generally **necessary** to 'account' for the missingness in Y . The 'complete case' estimator of θ_0 in general will be **inconsistent!**
 - That estimator may be consistent **only if:** (1) $\nabla\phi(\mathbf{X}, \theta_0) = \mathbf{0}$ a.s. for **every** \mathbf{X} (for regression problems, this indicates the 'correct model' case), and/or (2) $T \perp\!\!\!\perp (Y, X)$ (i.e. the **MCAR** case).
 - **Illustration of (1) for sq. loss:** $\nabla\phi(\mathbf{X}, \theta_0) = \mathbb{E}\{\mathbf{X}(Y - \mathbf{X}'\theta_0)|\mathbf{X}\} = \mathbf{0}$. Hence, $\mathbb{E}(Y|\mathbf{X}) = \mathbf{X}'\theta_0$ (i.e. a 'linear model' holds for $Y|\mathbf{X}$).
- With θ_0 (and \mathbf{X}) being **high dimensional** (compared to n), we **need** some further **structural constraints** on θ_0 to estimate it using \mathcal{D}_n .
 - We **assume** that θ_0 is **s-sparse**: $\|\theta_0\|_0 := s$ and $s \leq \min(n, d)$.
 - Note: the **sparsity requirement** has attractive (and fairly **intuitive**) **geometric justification** for **all the examples** we have given here.

- Under MAR assmpn., $R(\theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)\} \equiv \mathbb{E}_{\mathbf{X}}\{\phi(\mathbf{X}, \theta)\}$ admits the following **debiased and doubly robust (DDR)** representation:

- Under MAR assmpn., $R(\theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)\} \equiv \mathbb{E}_{\mathbf{X}}\{\phi(\mathbf{X}, \theta)\}$ admits the following **debiased and doubly robust (DDR)** representation:

$$R(\theta) = \mathbb{E}_{\mathbf{X}}\{\phi(\mathbf{X}, \theta)\} + \mathbb{E} \left[\frac{T}{\pi(\mathbf{X})} \{L(Y, \mathbf{X}, \theta) - \phi(\mathbf{X}, \theta)\} \right]. \quad (1)$$

Purely **non-parametric identification** based on the **observable \mathbf{Z}** and the **nuisance functions: $\pi(\mathbf{X})$ and $\phi(\mathbf{X}, \theta)$** (unknown but **estimable**).

- Under MAR assmpn., $R(\theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)\} \equiv \mathbb{E}_{\mathbf{X}}\{\phi(\mathbf{X}, \theta)\}$ admits the following **debiased and doubly robust (DDR)** representation:

$$R(\theta) = \mathbb{E}_{\mathbf{X}}\{\phi(\mathbf{X}, \theta)\} + \mathbb{E} \left[\frac{T}{\pi(\mathbf{X})} \{L(Y, \mathbf{X}, \theta) - \phi(\mathbf{X}, \theta)\} \right]. \quad (1)$$

Purely **non-parametric identification** based on the observable \mathbf{Z} and the **nuisance functions**: $\pi(\mathbf{X})$ and $\phi(\mathbf{X}, \theta)$ (unknown but **estimable**).

- 2^{nd} term is simply 0, can be seen as a 'debiasing' term (of sorts).
 - Plays a **crucial** role in analyzing the empirical version of (1). **Ensures first order insensitivity** to any estimation errors of $\pi(\cdot)$ and $\phi(\cdot)$.

- Under MAR assmpn., $R(\theta) := \mathbb{E}\{L(Y, \mathbf{X}, \theta)\} \equiv \mathbb{E}_{\mathbf{X}}\{\phi(\mathbf{X}, \theta)\}$ admits the following **debiased and doubly robust (DDR)** representation:

$$R(\theta) = \mathbb{E}_{\mathbf{X}}\{\phi(\mathbf{X}, \theta)\} + \mathbb{E} \left[\frac{T}{\pi(\mathbf{X})} \{L(Y, \mathbf{X}, \theta) - \phi(\mathbf{X}, \theta)\} \right]. \quad (1)$$

Purely **non-parametric identification** based on the observable \mathbf{Z} and the **nuisance functions**: $\pi(\mathbf{X})$ and $\phi(\mathbf{X}, \theta)$ (unknown but **estimable**).

- 2^{nd} term is simply 0, can be seen as a 'debiasing' term (of sorts).
 - Plays a **crucial** role in analyzing the empirical version of (1). **Ensures first order insensitivity** to any estimation errors of $\pi(\cdot)$ and $\phi(\cdot)$.
- **Double robustness (DR) aspect**: replace $\{\phi(\mathbf{X}, \theta), \pi(\mathbf{X})\}$ by **any** $\{\phi^*(\mathbf{X}, \theta), \pi^*(\mathbf{X})\}$ and (1) **continues to hold** as long as **one but not necessarily both** of $\phi^*(\cdot) = \phi(\cdot)$ or $\pi^*(\cdot) = \pi(\cdot)$ hold.

- Given **any** estimators $\{\hat{\pi}(\cdot), \hat{\phi}(\cdot)\}$ be of the nuisance fns. $\{\pi(\cdot), \phi(\cdot)\}$, we define our **L_1 -penalized DDR estimator $\hat{\theta}_{\text{DDR}}$** of θ_0 as:

$$\hat{\theta}_{\text{DDR}} \equiv \hat{\theta}_{\text{DDR}}(\lambda_n) := \arg \min_{\theta \in \mathbb{R}^d} \{ \mathcal{L}_n^{\text{DDR}}(\theta) + \lambda_n \|\theta\|_1 \}, \text{ where}$$

$$\mathcal{L}_n^{\text{DDR}}(\theta) := \frac{1}{n} \sum_{i=1}^n \hat{\phi}(\mathbf{X}_i, \theta) + \frac{T_i}{\hat{\pi}(\mathbf{X}_i)} \left\{ L(Y_i, \mathbf{X}_i, \theta) - \hat{\phi}(\mathbf{X}_i, \theta) \right\},$$

$\lambda_n \geq 0$ is the tuning parameter and $\{\hat{\pi}(\cdot), \hat{\phi}(\cdot)\}$ are arbitrary except for **satisfying two basic conditions** regarding their construction:

- Given **any** estimators $\{\widehat{\pi}(\cdot), \widehat{\phi}(\cdot)\}$ be of the nuisance fns. $\{\pi(\cdot), \phi(\cdot)\}$, we define our **L_1 -penalized DDR estimator $\widehat{\theta}_{\text{DDR}}$** of θ_0 as:

$$\widehat{\theta}_{\text{DDR}} \equiv \widehat{\theta}_{\text{DDR}}(\lambda_n) := \arg \min_{\theta \in \mathbb{R}^d} \{ \mathcal{L}_n^{\text{DDR}}(\theta) + \lambda_n \|\theta\|_1 \}, \text{ where}$$

$$\mathcal{L}_n^{\text{DDR}}(\theta) := \frac{1}{n} \sum_{i=1}^n \widehat{\phi}(\mathbf{X}_i, \theta) + \frac{T_i}{\widehat{\pi}(\mathbf{X}_i)} \left\{ L(Y_i, \mathbf{X}_i, \theta) - \widehat{\phi}(\mathbf{X}_i, \theta) \right\},$$

$\lambda_n \geq 0$ is the tuning parameter and $\{\widehat{\pi}(\cdot), \widehat{\phi}(\cdot)\}$ are arbitrary except for satisfying two **basic conditions** regarding their construction:

- $\widehat{\pi}(\cdot)$ obtained from the data $\mathcal{T}_n := \{T_i, \mathbf{X}_i\}_{i=1}^n$ **only**; $\{\widehat{\phi}(\mathbf{X}_i, \theta)\}_{i=1}^n$ obtained in a **'cross-fitted' manner** (via sample splitting).
- Assume (temporarily) $\{\widehat{\pi}(\cdot), \widehat{\phi}(\cdot)\}$ are **both** 'correct'. DR properties (consistency) of $\widehat{\theta}_{\text{DDR}}$ under their misspecifications discussed later.

- **For simplicity, assume** that the gradient $\nabla L(Y, \mathbf{X}, \theta)$ of $L(\cdot)$ satisfies a 'separable form' as follows: for some $\mathbf{h}(\mathbf{X}) \in \mathbb{R}^d$ and $g(\mathbf{X}, \theta) \in \mathbb{R}$,

- **For simplicity**, assume that the gradient $\nabla L(Y, \mathbf{X}, \theta)$ of $L(\cdot)$ satisfies a 'separable form' as follows: for some $\mathbf{h}(\mathbf{X}) \in \mathbb{R}^d$ and $g(\mathbf{X}, \theta) \in \mathbb{R}$,

$$\nabla L(Y, \mathbf{X}, \theta) = \mathbf{h}(\mathbf{X})\{Y - g(\mathbf{X}, \theta)\}, \text{ and hence,}$$

$$\nabla \hat{\phi}(\mathbf{X}, \theta) = \mathbf{h}(\mathbf{X})\{\hat{m}(\mathbf{X}) - g(\mathbf{X}, \theta)\}, \text{ where}$$

$\hat{m}(\mathbf{X})$ denotes the corresponding (cross-fitted) estimator of $m(\mathbf{X})$.

This simplifying assumption holds for all examples given before.

- Assumed form \Rightarrow only need to obtain $\hat{m}(\mathbf{X}_i)$ and not $\hat{\phi}(\mathbf{X}_i, \theta)$.

- **For simplicity**, assume that the gradient $\nabla L(Y, \mathbf{X}, \theta)$ of $L(\cdot)$ satisfies a 'separable form' as follows: for some $\mathbf{h}(\mathbf{X}) \in \mathbb{R}^d$ and $g(\mathbf{X}, \theta) \in \mathbb{R}$,

$$\nabla L(Y, \mathbf{X}, \theta) = \mathbf{h}(\mathbf{X})\{Y - g(\mathbf{X}, \theta)\}, \quad \text{and hence,}$$

$$\nabla \hat{\phi}(\mathbf{X}, \theta) = \mathbf{h}(\mathbf{X})\{\hat{m}(\mathbf{X}) - g(\mathbf{X}, \theta)\}, \quad \text{where}$$

$\hat{m}(\mathbf{X})$ denotes the corresponding (cross-fitted) estimator of $m(\mathbf{X})$.

This simplifying assumption holds for all examples given before.

- Assumed form \Rightarrow only need to obtain $\hat{m}(\mathbf{X}_i)$ and not $\hat{\phi}(\mathbf{X}_i, \theta)$.
- **Implementation algorithm.** $\hat{\theta}_{\text{DDR}}$ can be obtained simply as:

$$\hat{\theta}_{\text{DDR}} \equiv \hat{\theta}_{\text{DDR}}(\lambda_n) := \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n L(\tilde{Y}_i, \mathbf{X}_i, \theta) + \lambda_n \|\theta\|_1 \right\},$$

where $\tilde{Y}_i := \hat{m}(\mathbf{X}_i) + \frac{T_i}{\pi(\mathbf{X}_i)}\{Y_i - \hat{m}(\mathbf{X}_i)\}$, $\forall i$, is a 'pseudo' outcome.

Can use 'glmnet' in R. Pretend to have a 'full' data: $\{\tilde{Y}_i, \mathbf{X}_i\}_{i=1}^n$.

- Assume $L(\cdot)$ is convex and differentiable in θ and $\mathcal{L}_n^{\text{DDR}}(\theta)$ satisfies the Restricted Strong Convexity (RSC) condition (Negahban et al., 2012) at $\theta = \theta_0$. Then, for any choice of $\lambda_n \geq 2 \|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$,

- Assume $L(\cdot)$ is convex and differentiable in $\boldsymbol{\theta}$ and $\mathcal{L}_n^{\text{DDR}}(\boldsymbol{\theta})$ satisfies the Restricted Strong Convexity (RSC) condition (Negahban et al., 2012) at $\boldsymbol{\theta} = \boldsymbol{\theta}_0$. Then, for any choice of $\lambda_n \geq 2 \|\nabla \mathcal{L}_n^{\text{DDR}}(\boldsymbol{\theta}_0)\|_\infty$,

$$\left\| \widehat{\boldsymbol{\theta}}_{\text{DDR}}(\lambda_n) - \boldsymbol{\theta}_0 \right\|_2 \lesssim \lambda_n \sqrt{s}, \text{ and } \left\| \widehat{\boldsymbol{\theta}}_{\text{DDR}}(\lambda_n) - \boldsymbol{\theta}_0 \right\|_1 \lesssim \lambda_n s.$$

where $s := \|\boldsymbol{\theta}_0\|_0$. This is a **deterministic** deviation bound. Holds for any choices of $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ and for any realization of \mathcal{D}_n .

- Assume $L(\cdot)$ is convex and differentiable in θ and $\mathcal{L}_n^{\text{DDR}}(\theta)$ satisfies the Restricted Strong Convexity (RSC) condition (Negahban et al., 2012) at $\theta = \theta_0$. Then, for any choice of $\lambda_n \geq 2 \|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$,

$$\left\| \hat{\theta}_{\text{DDR}}(\lambda_n) - \theta_0 \right\|_2 \lesssim \lambda_n \sqrt{s}, \text{ and } \left\| \hat{\theta}_{\text{DDR}}(\lambda_n) - \theta_0 \right\|_1 \lesssim \lambda_n s.$$

where $s := \|\theta_0\|_0$. This is a **deterministic** deviation bound. Holds for any choices of $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$ and for any realization of \mathcal{D}_n .

- The RSC (or ‘cone’) condition for $\mathcal{L}_n^{\text{DDR}}(\theta)$ is exactly the **same** as the usual RSC condition required under a fully observed data! The fully observed data RSC condition’s validity is well studied.

- Assume $L(\cdot)$ is convex and differentiable in $\boldsymbol{\theta}$ and $\mathcal{L}_n^{\text{DDR}}(\boldsymbol{\theta})$ satisfies the Restricted Strong Convexity (RSC) condition (Negahban et al., 2012) at $\boldsymbol{\theta} = \boldsymbol{\theta}_0$. Then, for any choice of $\lambda_n \geq 2 \|\nabla \mathcal{L}_n^{\text{DDR}}(\boldsymbol{\theta}_0)\|_\infty$,

$$\left\| \widehat{\boldsymbol{\theta}}_{\text{DDR}}(\lambda_n) - \boldsymbol{\theta}_0 \right\|_2 \lesssim \lambda_n \sqrt{s}, \text{ and } \left\| \widehat{\boldsymbol{\theta}}_{\text{DDR}}(\lambda_n) - \boldsymbol{\theta}_0 \right\|_1 \lesssim \lambda_n s.$$

where $s := \|\boldsymbol{\theta}_0\|_0$. This is a **deterministic** deviation bound. Holds for any choices of $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ and for any realization of \mathcal{D}_n .

- The RSC (or ‘cone’) condition for $\mathcal{L}_n^{\text{DDR}}(\boldsymbol{\theta})$ is exactly the **same** as the usual RSC condition required under a fully observed data! The fully observed data RSC condition’s validity is well studied.
- Key quantity of interest:** the random lower bound $\|\nabla \mathcal{L}_n^{\text{DDR}}(\boldsymbol{\theta}_0)\|_\infty$ for λ_n . **Need probabilistic bounds** to determine convergence rate of $\widehat{\boldsymbol{\theta}}_{\text{DDR}}$.

- Bounds on $\|\nabla\mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$ determines the rate of choice of λ_n and hence the convergence rate of $\hat{\theta}_{\text{DDR}}$ (using the deviation bound).
- **Probabilistic** bounds for $\|\nabla\mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$: the basic decomposition

$$\|\nabla\mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty \leq \|\mathbf{T}_{0,n}\|_\infty + \|\mathbf{T}_{\pi,n}\|_\infty + \|\mathbf{T}_{m,n}\|_\infty + \|\mathbf{R}_{\pi,m,n}\|_\infty,$$

- Bounds on $\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$ determines the rate of choice of λ_n and hence the convergence rate of $\hat{\theta}_{\text{DDR}}$ (using the deviation bound).
- **Probabilistic** bounds for $\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$: the basic decomposition

$$\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty \leq \|\mathbf{T}_{0,n}\|_\infty + \|\mathbf{T}_{\pi,n}\|_\infty + \|\mathbf{T}_{m,n}\|_\infty + \|\mathbf{R}_{\pi,m,n}\|_\infty,$$

where $\mathbf{T}_{0,n}$ is the 'main' term (a centered iid average), $\mathbf{T}_{\pi,n}$ is the ' π -error' term involving $\hat{\pi}(\cdot) - \pi(\cdot)$ and $\mathbf{T}_{m,n}$ is the ' m -error' term involving $\hat{m}(\cdot) - m(\cdot)$, while $\mathbf{R}_{\pi,m,n}$ is the ' (π, m) -error' term (usually lower order) involving the product of $\hat{\pi}(\cdot) - \pi(\cdot)$ and $\hat{m}(\cdot) - m(\cdot)$.

- Control each term separately. The analyses are all non-asymptotic and nuanced, especially in order to get sharp rates for $\mathbf{T}_{\pi,n}$ and $\mathbf{T}_{m,n}$.

The Main Goal from Hereon: Probabilistic Bounds for $\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$

- Bounds on $\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$ determines the rate of choice of λ_n and hence the convergence rate of $\hat{\theta}_{\text{DDR}}$ (using the deviation bound).

- **Probabilistic** bounds for $\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty$: the basic decomposition

$$\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty \leq \|\mathbf{T}_{0,n}\|_\infty + \|\mathbf{T}_{\pi,n}\|_\infty + \|\mathbf{T}_{m,n}\|_\infty + \|\mathbf{R}_{\pi,m,n}\|_\infty,$$

where $\mathbf{T}_{0,n}$ is the 'main' term (a centered iid average), $\mathbf{T}_{\pi,n}$ is the ' π -error' term involving $\hat{\pi}(\cdot) - \pi(\cdot)$ and $\mathbf{T}_{m,n}$ is the ' m -error' term involving $\hat{m}(\cdot) - m(\cdot)$, while $\mathbf{R}_{\pi,m,n}$ is the ' (π, m) -error' term (usually lower order) involving the product of $\hat{\pi}(\cdot) - \pi(\cdot)$ and $\hat{m}(\cdot) - m(\cdot)$.

- Control each term separately. The analyses are all non-asymptotic and nuanced, especially in order to get sharp rates for $\mathbf{T}_{\pi,n}$ and $\mathbf{T}_{m,n}$.
- **We show:** $\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty \lesssim \sqrt{(\log d)/n}$ with high probability, and hence $\|\hat{\theta}_{\text{DDR}} - \theta_0\|_2 \lesssim \sqrt{s(\log d)/n}$. So, clearly it is **rate optimal**.

- **Basic (high level) consistency conditions on $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$.** Let $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$ be **any** general and 'correct' estimators of $\{\pi(\cdot), m(\cdot)\}$, and assume they **satisfy** the following **pointwise** convergence rates:

- **Basic (high level) consistency conditions on $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$.** Let $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$ be **any** general and 'correct' estimators of $\{\pi(\cdot), m(\cdot)\}$, and assume they **satisfy** the following **pointwise** convergence rates:

$$|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim_{\mathbb{P}} \delta_{n,\pi} \quad \text{and} \quad |\hat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim_{\mathbb{P}} \xi_{n,m} \quad \forall \mathbf{x} \in \mathcal{X}, \quad (2)$$

for **some** sequences $\delta_{n,\pi}, \xi_{n,m} \geq 0$ such that $(\delta_{n,\pi} + \xi_{n,m})\sqrt{\log(nd)} = o(1)$ and the product $\delta_{n,\pi}\xi_{n,m}(\log n) = o(\sqrt{(\log d)/n})$.

- **Basic (high level) consistency conditions on $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$.** Let $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$ be **any** general and 'correct' estimators of $\{\pi(\cdot), m(\cdot)\}$, and assume they **satisfy** the following **pointwise** convergence rates:

$$|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim_{\mathbb{P}} \delta_{n,\pi} \quad \text{and} \quad |\hat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim_{\mathbb{P}} \xi_{n,m} \quad \forall \mathbf{x} \in \mathcal{X}, \quad (2)$$

for **some** sequences $\delta_{n,\pi}, \xi_{n,m} \geq 0$ such that $(\delta_{n,\pi} + \xi_{n,m})\sqrt{\log(nd)} = o(1)$ and the product $\delta_{n,\pi}\xi_{n,m}(\log n) = o(\sqrt{(\log d)/n})$.

- Under condition (2), along with some more 'suitable' tail assumptions (sub-Gaussian tails etc.), we have: **with high probability**,

$$\|\mathbf{T}_{0,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}}, \quad \|\mathbf{T}_{\pi,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}} \left\{ \delta_{n,\pi} \sqrt{\log(nd)} \right\}, \quad \text{and}$$

- **Basic (high level) consistency conditions on $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$.** Let $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ be **any** general and ‘correct’ estimators of $\{\pi(\cdot), m(\cdot)\}$, and assume they **satisfy** the following **pointwise** convergence rates:

$$|\widehat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim_{\mathbb{P}} \delta_{n,\pi} \quad \text{and} \quad |\widehat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim_{\mathbb{P}} \xi_{n,m} \quad \forall \mathbf{x} \in \mathcal{X}, \quad (2)$$

for **some** sequences $\delta_{n,\pi}, \xi_{n,m} \geq 0$ such that $(\delta_{n,\pi} + \xi_{n,m})\sqrt{\log(nd)} = o(1)$ and the product $\delta_{n,\pi}\xi_{n,m}(\log n) = o(\sqrt{(\log d)/n})$.

- Under condition (2), along with some more ‘suitable’ tail assumptions (sub-Gaussian tails etc.), we have: **with high probability**,

$$\|\mathbf{T}_{0,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}}, \quad \|\mathbf{T}_{\pi,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}} \left\{ \delta_{n,\pi} \sqrt{\log(nd)} \right\}, \quad \text{and}$$

$$\|\mathbf{T}_{m,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}} \left\{ \xi_{n,m} \sqrt{\log(nd)} \right\}, \quad \|\mathbf{R}_{\pi,m,n}\|_\infty \lesssim \delta_{n,\pi} \xi_{n,m} (\log n).$$

- **Basic (high level) consistency conditions on $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$.** Let $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$ be **any** general and ‘correct’ estimators of $\{\pi(\cdot), m(\cdot)\}$, and assume they **satisfy** the following **pointwise** convergence rates:

$$|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim_{\mathbb{P}} \delta_{n,\pi} \quad \text{and} \quad |\hat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim_{\mathbb{P}} \xi_{n,m} \quad \forall \mathbf{x} \in \mathcal{X}, \quad (2)$$

for **some** sequences $\delta_{n,\pi}, \xi_{n,m} \geq 0$ such that $(\delta_{n,\pi} + \xi_{n,m})\sqrt{\log(nd)} = o(1)$ and the product $\delta_{n,\pi}\xi_{n,m}(\log n) = o(\sqrt{(\log d)/n})$.

- Under condition (2), along with some more ‘suitable’ tail assumptions (sub-Gaussian tails etc.), we have: **with high probability**,

$$\|\mathbf{T}_{0,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}}, \quad \|\mathbf{T}_{\pi,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}} \left\{ \delta_{n,\pi} \sqrt{\log(nd)} \right\}, \quad \text{and}$$

$$\|\mathbf{T}_{m,n}\|_\infty \lesssim \sqrt{\frac{\log d}{n}} \left\{ \xi_{n,m} \sqrt{\log(nd)} \right\}, \quad \|\mathbf{R}_{\pi,m,n}\|_\infty \lesssim \delta_{n,\pi} \xi_{n,m} (\log n).$$

- Hence, $\|\nabla \mathcal{L}_n^{\text{DDR}}(\theta_0)\|_\infty \lesssim \sqrt{\frac{\log d}{n}} \{1 + o(1)\}$ with high probability.

- Consider $\hat{\theta}_{\text{DDR}}$ for the squared loss: $L(Y, \mathbf{X}, \theta) := \{Y - \Psi(\mathbf{X})'\theta\}^2$, where $\Psi(\mathbf{X}) \in \mathbb{R}^d$ denotes any HD vector of basis functions of \mathbf{X} .
- Define $\Sigma := \mathbb{E}\{\Psi(\mathbf{X})\Psi(\mathbf{X})'\}$, $\Omega := \Sigma^{-1}$, and let $\hat{\Omega}$ be any reasonable estimator of Ω (and assume Ω is sparse if required).
- We then define the **desparsified** DDR estimator $\tilde{\theta}_{\text{DDR}}$ as follows.

- Consider $\hat{\theta}_{\text{DDR}}$ for the squared loss: $L(Y, \mathbf{X}, \theta) := \{Y - \Psi(\mathbf{X})'\theta\}^2$, where $\Psi(\mathbf{X}) \in \mathbb{R}^d$ denotes any HD vector of basis functions of \mathbf{X} .
- Define $\Sigma := \mathbb{E}\{\Psi(\mathbf{X})\Psi(\mathbf{X})'\}$, $\Omega := \Sigma^{-1}$, and let $\hat{\Omega}$ be any reasonable estimator of Ω (and assume Ω is sparse if required).
- We then define the **desparsified** DDR estimator $\tilde{\theta}_{\text{DDR}}$ as follows.

$$\tilde{\theta}_{\text{DDR}} := \hat{\theta}_{\text{DDR}} + \hat{\Omega} \underbrace{\frac{1}{n} \sum_{i=1}^n \{\tilde{Y}_i - \Psi(\mathbf{X}_i)'\hat{\theta}_{\text{DDR}}\} \Psi(\mathbf{X}_i)}_{\text{Desparsification/Debiasing term}}, \text{ where}$$

$$\tilde{Y}_i := \hat{m}(\mathbf{X}_i) + \frac{T_i}{\hat{\pi}(\mathbf{X}_i)} \{Y_i - \hat{m}(\mathbf{X}_i)\} \text{ are the pseudo outcomes.}$$

- Consider $\hat{\theta}_{\text{DDR}}$ for the squared loss: $L(Y, \mathbf{X}, \theta) := \{Y - \Psi(\mathbf{X})'\theta\}^2$, where $\Psi(\mathbf{X}) \in \mathbb{R}^d$ denotes any HD vector of basis functions of \mathbf{X} .
- Define $\Sigma := \mathbb{E}\{\Psi(\mathbf{X})\Psi(\mathbf{X})'\}$, $\Omega := \Sigma^{-1}$, and let $\hat{\Omega}$ be any reasonable estimator of Ω (and assume Ω is sparse if required).
- We then define the **desparsified** DDR estimator $\tilde{\theta}_{\text{DDR}}$ as follows.

$$\tilde{\theta}_{\text{DDR}} := \hat{\theta}_{\text{DDR}} + \hat{\Omega} \underbrace{\frac{1}{n} \sum_{i=1}^n \{\tilde{Y}_i - \Psi(\mathbf{X}_i)'\hat{\theta}_{\text{DDR}}\} \Psi(\mathbf{X}_i)}_{\text{Desparsification/Debiasing term}}, \text{ where}$$

$$\tilde{Y}_i := \hat{m}(\mathbf{X}_i) + \frac{T_i}{\hat{\pi}(\mathbf{X}_i)} \{Y_i - \hat{m}(\mathbf{X}_i)\} \text{ are the pseudo outcomes.}$$

Debiasing similar (in spirit) to van de Geer et al. (2014), **except its the 'right' one for this problem** (using pseudo outcomes in the full data).

- **Assume:** the basic convergence conditions (2) for $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$, $\Omega\mathbf{X}$ is sub-Gaussian and that $\|\widehat{\Omega} - \Omega\|_1 = O_{\mathbb{P}}(a_n)$, $\|I - \widehat{\Omega}\widehat{\Sigma}\|_{\max} = O_{\mathbb{P}}(b_n)$, with $a_n\sqrt{\log d} = o(1)$ and $b_ns\sqrt{\log d} = o(1)$, where $s := \|\theta_0\|_0$.
- Then, $\widetilde{\theta}_{\text{DDR}}$ satisfies the asymptotic linear expansion **(ALE)**:

- **Assume:** the basic convergence conditions (2) for $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$, $\Omega \mathbf{X}$ is sub-Gaussian and that $\|\widehat{\Omega} - \Omega\|_1 = O_{\mathbb{P}}(a_n)$, $\|I - \widehat{\Omega} \widehat{\Sigma}\|_{\max} = O_{\mathbb{P}}(b_n)$, with $a_n \sqrt{\log d} = o(1)$ and $b_n s \sqrt{\log d} = o(1)$, where $s := \|\theta_0\|_0$.
- Then, $\widetilde{\theta}_{\text{DDR}}$ satisfies the asymptotic linear expansion (**ALE**):

$$(\widetilde{\theta}_{\text{DDR}} - \theta_0) = \frac{1}{n} \sum_{i=1}^n \Omega \{\psi_0(\mathbf{Z}_i)\} + \Delta_n, \text{ where } \|\Delta_n\|_{\infty} = o_{\mathbb{P}}(n^{-\frac{1}{2}})$$

$$\text{and } \psi_0(\mathbf{Z}) := \left[\{m(\mathbf{X}) - \Psi(\mathbf{X})' \theta_0\} + \frac{T}{\pi(\mathbf{X})} \{Y - m(\mathbf{X})\} \right] \Psi(\mathbf{X})$$

with $\mathbb{E}\{\psi_0(\mathbf{Z})\} = \mathbf{0}$. The ALE facilitates inference (e.g. confidence intervals etc.) for any low-d component of θ_0 via Gaussian approx.

- **Assume:** the basic convergence conditions (2) for $\{\hat{\pi}(\cdot), \hat{m}(\cdot)\}$, $\Omega\mathbf{X}$ is sub-Gaussian and that $\|\hat{\Omega} - \Omega\|_1 = O_{\mathbb{P}}(a_n)$, $\|I - \hat{\Omega}\hat{\Sigma}\|_{\max} = O_{\mathbb{P}}(b_n)$, with $a_n\sqrt{\log d} = o(1)$ and $b_ns\sqrt{\log d} = o(1)$, where $s := \|\theta_0\|_0$.

- Then, $\tilde{\theta}_{\text{DDR}}$ satisfies the asymptotic linear expansion (ALE):

$$(\tilde{\theta}_{\text{DDR}} - \theta_0) = \frac{1}{n} \sum_{i=1}^n \Omega\{\psi_0(\mathbf{Z}_i)\} + \Delta_n, \text{ where } \|\Delta_n\|_{\infty} = o_{\mathbb{P}}(n^{-\frac{1}{2}})$$

$$\text{and } \psi_0(\mathbf{Z}) := \left[\{m(\mathbf{X}) - \Psi(\mathbf{X})'\theta_0\} + \frac{T}{\pi(\mathbf{X})}\{Y - m(\mathbf{X})\} \right] \Psi(\mathbf{X})$$

with $\mathbb{E}\{\psi_0(\mathbf{Z})\} = \mathbf{0}$. The ALE facilitates inference (e.g. confidence intervals etc.) for any low-d component of θ_0 via Gaussian approx.

- Further, the ALE is also 'optimal'. The function $\Omega\psi_0(\mathbf{Z}) =: \Psi_{\text{eff}}(\mathbf{Z})$ is the 'efficient' influence function for θ_0 (Robins et al., 1994). Thus, in classical settings, $\tilde{\theta}_{\text{DDR}}$ achieves the semi-parametric efficiency bound.

- Coordinate-wise asymptotic normality of $\tilde{\boldsymbol{\theta}}_{\text{DDR}}$: $\forall 1 \leq j \leq d$,

$$\sqrt{n}(\tilde{\boldsymbol{\theta}}_{\text{DDR}} - \boldsymbol{\theta}_0)_j \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma_{0,j}^2), \text{ where } \sigma_{0,j}^2 := \text{Var}\{\boldsymbol{\Omega}'_j \psi_0(\mathbf{Z})\}.$$

Further, $\max_{1 \leq j \leq d} |\hat{\sigma}_{0,j} - \sigma_{0,j}| = o_{\mathbb{P}}(1)$, where $\hat{\sigma}_{0,j}$ is the plug-in estimator obtained by plugging in $\hat{\boldsymbol{\Omega}}$, $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ in $\text{Var}\{\boldsymbol{\Omega}'_j \psi_0(\mathbf{Z})\}$.

- Can choose $\hat{\boldsymbol{\Omega}}$ to be **any** standard (sparse) precision matrix estimator, e.g. the **node-wise Lasso estimator**. Here, $a_n = s_{\boldsymbol{\Omega}} \sqrt{(\log d)/n}$ and $b_n = \sqrt{(\log d)/n}$ under suitable conditions, with $s_{\boldsymbol{\Omega}} := \max_{1 \leq j \leq d} \|\boldsymbol{\Omega}_j\|_0$.

- Coordinate-wise asymptotic normality of $\tilde{\boldsymbol{\theta}}_{\text{DDR}}$: $\forall 1 \leq j \leq d$,

$$\sqrt{n}(\tilde{\boldsymbol{\theta}}_{\text{DDR}} - \boldsymbol{\theta}_0)_j \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma_{0,j}^2), \text{ where } \sigma_{0,j}^2 := \text{Var}\{\boldsymbol{\Omega}'_j \boldsymbol{\psi}_0(\mathbf{Z})\}.$$

Further, $\max_{1 \leq j \leq d} |\hat{\sigma}_{0,j} - \sigma_{0,j}| = o_{\mathbb{P}}(1)$, where $\hat{\sigma}_{0,j}$ is the plug-in estimator obtained by plugging in $\hat{\boldsymbol{\Omega}}$, $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ in $\text{Var}\{\boldsymbol{\Omega}'_j \boldsymbol{\psi}_0(\mathbf{Z})\}$.

- Can choose $\hat{\boldsymbol{\Omega}}$ to be **any** standard (sparse) precision matrix estimator, e.g. the **node-wise Lasso estimator**. Here, $a_n = s_{\Omega} \sqrt{(\log d)/n}$ and $b_n = \sqrt{(\log d)/n}$ under suitable conditions, with $s_{\Omega} := \max_{1 \leq j \leq d} \|\boldsymbol{\Omega}_j\|_0$.
- The error $\boldsymbol{\Delta}_n$ can be decomposed as: $\boldsymbol{\Delta}_n = \boldsymbol{\Delta}_{n,1} + \boldsymbol{\Delta}_{n,2} + \boldsymbol{\Delta}_{n,3}$, where $\boldsymbol{\Delta}_{n,1} := \frac{1}{n}(\hat{\boldsymbol{\Omega}} - \boldsymbol{\Omega}) \sum_{i=1}^n \boldsymbol{\psi}_0(\mathbf{Z}_i)$, $\boldsymbol{\Delta}_{n,2} := (I_d - \hat{\boldsymbol{\Omega}}\hat{\boldsymbol{\Sigma}})(\hat{\boldsymbol{\theta}}_{\text{DDR}} - \boldsymbol{\theta}_0)$ and $\boldsymbol{\Delta}_{n,3} := \hat{\boldsymbol{\Omega}}(\mathbf{T}_{\pi,n} + \mathbf{T}_{m,n} + \mathbf{R}_{\pi,m,n})$, with $\|\boldsymbol{\Delta}_{n,3}\|_{\infty} \lesssim_{\mathbb{P}} n^{-\frac{1}{2}}$ and

$$\|\boldsymbol{\Delta}_{n,1}\|_{\infty} \lesssim a_n \sqrt{\frac{\log d}{n}} \quad \text{and} \quad \|\boldsymbol{\Delta}_{n,2}\|_{\infty} \lesssim b_n s \sqrt{\frac{\log d}{n}}.$$

- Finally, let $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\} \rightarrow \{\pi^*(\cdot), m^*(\cdot)\}$, with either $\pi^*(\cdot) = \pi(\cdot)$ or $m^*(\cdot) = m(\cdot)$ **but not necessarily both**. Assume the same pointwise convergence conditions and rates $(\delta_{n,\pi}, \xi_{n,m})$ for $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ as in (2), **but** now with $\{\pi(\cdot), m(\cdot)\}$ therein **replaced by** $\{\pi^*(\cdot), m^*(\cdot)\}$.
- Under some ‘suitable’ assumptions, we have: **with high probability**,

- Finally, let $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\} \rightarrow \{\pi^*(\cdot), m^*(\cdot)\}$, with either $\pi^*(\cdot) = \pi(\cdot)$ or $m^*(\cdot) = m(\cdot)$ **but not necessarily both**. Assume the same pointwise convergence conditions and rates $(\delta_{n,\pi}, \xi_{n,m})$ for $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ as in (2), **but** now with $\{\pi(\cdot), m(\cdot)\}$ therein **replaced by** $\{\pi^*(\cdot), m^*(\cdot)\}$.
- Under some ‘suitable’ assumptions, we have: **with high probability**,

$$\|\mathbf{T}_{0,n}\|_{\infty} + \|\mathbf{T}_{\pi,n}\|_{\infty} + \|\mathbf{T}_{m,n}\|_{\infty} \lesssim \sqrt{\frac{\log d}{n}} \{1 + \mathbf{1}_{(\pi^*, m^*) \neq (\pi, m)}\}$$

- Finally, let $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\} \rightarrow \{\pi^*(\cdot), m^*(\cdot)\}$, with either $\pi^*(\cdot) = \pi(\cdot)$ or $m^*(\cdot) = m(\cdot)$ **but not necessarily both**. Assume the same pointwise convergence conditions and rates $(\delta_{n,\pi}, \xi_{n,m})$ for $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ as in (2), **but** now with $\{\pi(\cdot), m(\cdot)\}$ therein **replaced by** $\{\pi^*(\cdot), m^*(\cdot)\}$.
- Under some ‘suitable’ assumptions, we have: **with high probability**,

$$\|\mathbf{T}_{0,n}\|_{\infty} + \|\mathbf{T}_{\pi,n}\|_{\infty} + \|\mathbf{T}_{m,n}\|_{\infty} \lesssim \sqrt{\frac{\log d}{n}} \{1 + \mathbf{1}_{(\pi^*, m^*) \neq (\pi, m)}\}$$

$$\text{and } \|\mathbf{R}_{\pi,m,n}\|_{\infty} \lesssim \{\delta_{n,\pi} \mathbf{1}_{(m^* \neq m)} + \xi_{n,m} \mathbf{1}_{(\pi^* \neq \pi)} + \delta_{n,\pi} \xi_{n,m}\} (\log n).$$

- Finally, let $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\} \rightarrow \{\pi^*(\cdot), m^*(\cdot)\}$, with either $\pi^*(\cdot) = \pi(\cdot)$ or $m^*(\cdot) = m(\cdot)$ **but not necessarily both**. Assume the same pointwise convergence conditions and rates $(\delta_{n,\pi}, \xi_{n,m})$ for $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ as in (2), **but** now with $\{\pi(\cdot), m(\cdot)\}$ therein **replaced by** $\{\pi^*(\cdot), m^*(\cdot)\}$.
- Under some 'suitable' assumptions, we have: **with high probability**,

$$\|\mathbf{T}_{0,n}\|_{\infty} + \|\mathbf{T}_{\pi,n}\|_{\infty} + \|\mathbf{T}_{m,n}\|_{\infty} \lesssim \sqrt{\frac{\log d}{n}} \{1 + \mathbf{1}_{(\pi^*, m^*) \neq (\pi, m)}\}$$

$$\text{and } \|\mathbf{R}_{\pi,m,n}\|_{\infty} \lesssim \{\delta_{n,\pi} \mathbf{1}_{(m^* \neq m)} + \xi_{n,m} \mathbf{1}_{(\pi^* \neq \pi)} + \delta_{n,\pi} \xi_{n,m}\} (\log n).$$

The 2nd and/or 3rd terms also contribute now to the rate $\sqrt{(\log d)/n}$.
 The 4th term is $o(1)$ but **no longer ignorable** (and may be slower).

- Finally, let $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\} \rightarrow \{\pi^*(\cdot), m^*(\cdot)\}$, with either $\pi^*(\cdot) = \pi(\cdot)$ or $m^*(\cdot) = m(\cdot)$ **but not necessarily both**. Assume the same pointwise convergence conditions and rates $(\delta_{n,\pi}, \xi_{n,m})$ for $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$ as in (2), **but** now with $\{\pi(\cdot), m(\cdot)\}$ therein **replaced by** $\{\pi^*(\cdot), m^*(\cdot)\}$.

- Under some ‘suitable’ assumptions, we have: **with high probability**,

$$\|\mathbf{T}_{0,n}\|_{\infty} + \|\mathbf{T}_{\pi,n}\|_{\infty} + \|\mathbf{T}_{m,n}\|_{\infty} \lesssim \sqrt{\frac{\log d}{n}} \{1 + \mathbf{1}_{(\pi^*, m^*) \neq (\pi, m)}\}$$

$$\text{and } \|\mathbf{R}_{\pi,m,n}\|_{\infty} \lesssim \{\delta_{n,\pi} \mathbf{1}_{(m^* \neq m)} + \xi_{n,m} \mathbf{1}_{(\pi^* \neq \pi)} + \delta_{n,\pi} \xi_{n,m}\} (\log n).$$

The 2nd and/or 3rd terms also contribute now to the rate $\sqrt{(\log d)/n}$. The 4th term is $o(1)$ but **no longer ignorable** (and may be slower).

- Regardless, this **establishes general** convergence rates and the **DR property** of $\widehat{\theta}_{\text{DDR}}$ under possible misspecification of $\{\widehat{\pi}(\cdot), \widehat{m}(\cdot)\}$. For the 4th term, **sharper rates need a case-by-case analysis**.

- **Note:** our theory holds generally for **any** choices of $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ under mild conditions (provided they are both 'correct' estimators).
- Under misspecifications, consistency & general non-sharp rates are also established. Sharp rates **need** case-by-case analyses.
 - **Even** for mean (or ATE) estimation problem, this can be quite tricky in HD settings. See Smucler et al. (2019) for a detailed analysis.

- **Note:** our theory holds generally for **any** choices of $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ under mild conditions (provided they are both 'correct' estimators).
 - Under misspecifications, consistency & general non-sharp rates are also established. Sharp rates **need** case-by-case analyses.
 - **Even** for mean (or ATE) estimation problem, this can be quite tricky in HD settings. See Smucler et al. (2019) for a detailed analysis.
- Below we provide only **some** choices of $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ that may be used to implement our theory & methods for $\hat{\theta}_{\text{DDR}}$. In general, one can use any reasonable method (including black box ML methods).
- **Choices of $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$:** we consider estimators from two families.

- **Note:** our theory holds generally for **any** choices of $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ under mild conditions (provided they are both 'correct' estimators).
 - Under misspecifications, consistency & general non-sharp rates are also established. Sharp rates **need** case-by-case analyses.
 - **Even** for mean (or ATE) estimation problem, this can be quite tricky in HD settings. See Smucler et al. (2019) for a detailed analysis.
- Below we provide only **some** choices of $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ that may be used to implement our theory & methods for $\hat{\theta}_{\text{DDR}}$. In general, one can use any reasonable method (including black box ML methods).
- **Choices of $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$:** we consider estimators from two families.
 - Parametric and 'extended' parametric families (series estimators).
 - Semi-parametric single index families.

- If $\pi(\cdot)$ is known, we set $\hat{\pi}(\cdot) := \pi(\cdot)$. Otherwise, we estimate $\pi(\cdot)$ via two (class of) choices of $\hat{\pi}(\cdot)$ (each assumed to be 'correct').

- If $\pi(\cdot)$ is known, we set $\hat{\pi}(\cdot) := \pi(\cdot)$. Otherwise, we estimate $\pi(\cdot)$ via two (class of) choices of $\hat{\pi}(\cdot)$ (each assumed to be 'correct').
- 'Extended' parametric family: $\pi(\mathbf{x}) = g\{\boldsymbol{\alpha}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot) \in [0, 1]$ is a **known** function [e.g. $g_{\text{expit}}(u) := \exp(u)/\{1 + \exp(u)\}$], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\alpha} \in \mathbb{R}^K$ is an unknown (sparse) parameter vector.

- If $\pi(\cdot)$ is known, we set $\hat{\pi}(\cdot) := \pi(\cdot)$. Otherwise, we estimate $\pi(\cdot)$ via two (class of) choices of $\hat{\pi}(\cdot)$ (each assumed to be 'correct').
- 'Extended' parametric family: $\pi(\mathbf{x}) = g\{\boldsymbol{\alpha}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot) \in [0, 1]$ is a **known** function [e.g. $g_{\text{expit}}(u) := \exp(u)/\{1 + \exp(u)\}$], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\alpha} \in \mathbb{R}^K$ is an unknown (sparse) parameter vector.
- Example: $\boldsymbol{\Psi}(\mathbf{X})$ may correspond to the **polynomial bases** of \mathbf{X} upto any fixed degree k . Note: the **special case of linear bases** ($k = 1$) **includes** all **standard parametric regression models**. Further, the case of $\pi(\cdot) = \text{constant}$ (but unknown) i.e. MCAR is **also included**.

- If $\pi(\cdot)$ is known, we set $\hat{\pi}(\cdot) := \pi(\cdot)$. Otherwise, we estimate $\pi(\cdot)$ via two (class of) choices of $\hat{\pi}(\cdot)$ (each assumed to be 'correct').
- 'Extended' parametric family: $\pi(\mathbf{x}) = g\{\boldsymbol{\alpha}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot) \in [0, 1]$ is a **known** function [e.g. $g_{\text{expit}}(u) := \exp(u)/\{1 + \exp(u)\}$], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\alpha} \in \mathbb{R}^K$ is an unknown (sparse) parameter vector.
 - Example: $\boldsymbol{\Psi}(\mathbf{X})$ may correspond to the **polynomial bases** of \mathbf{X} upto any fixed degree k . Note: the **special case of linear bases** ($k = 1$) **includes** all **standard parametric regression models**. Further, the case of $\pi(\cdot) = \text{constant}$ (but unknown) i.e. MCAR is **also included**.
 - **Estimator**: we set $\hat{\pi}(\mathbf{X}) = g\{\hat{\boldsymbol{\alpha}}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $\hat{\boldsymbol{\alpha}}$ denotes **any suitable estimator** (possibly penalized) of $\boldsymbol{\alpha}$ based on $\mathcal{T}_n := \{T_i, \mathbf{X}_i\}_{i=1}^n$.

- If $\pi(\cdot)$ is known, we set $\hat{\pi}(\cdot) := \pi(\cdot)$. Otherwise, we estimate $\pi(\cdot)$ via two (class of) choices of $\hat{\pi}(\cdot)$ (each assumed to be 'correct').
- 'Extended' parametric family: $\pi(\mathbf{x}) = g\{\boldsymbol{\alpha}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot) \in [0, 1]$ is a **known** function [e.g. $g_{\text{expit}}(u) := \exp(u)/\{1 + \exp(u)\}$], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\alpha} \in \mathbb{R}^K$ is an unknown (sparse) parameter vector.
 - Example: $\boldsymbol{\Psi}(\mathbf{X})$ may correspond to the **polynomial bases** of \mathbf{X} upto any fixed degree k . Note: the **special case of linear bases** ($k = 1$) **includes** all **standard parametric regression models**. Further, the case of $\pi(\cdot) = \text{constant}$ (but unknown) i.e. MCAR is **also included**.
 - **Estimator**: we set $\hat{\pi}(\mathbf{X}) = g\{\hat{\boldsymbol{\alpha}}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $\hat{\boldsymbol{\alpha}}$ denotes **any suitable estimator** (possibly penalized) of $\boldsymbol{\alpha}$ based on $\mathcal{T}_n := \{T_i, \mathbf{X}_i\}_{i=1}^n$.
 - **Example of $\hat{\boldsymbol{\alpha}}$** : when $g(\cdot) = g_{\text{expit}}(\cdot)$, $\hat{\boldsymbol{\alpha}}$ may be obtained based on a standard L_1 -penalized logistic regression of $\{T_i \text{ vs. } \boldsymbol{\Psi}(\mathbf{X}_i)\}_{i=1}^n$.

- **Semi-parametric single index family**: $\pi(\mathbf{X}) = g(\boldsymbol{\alpha}'\mathbf{X})$, where $g(\cdot) \in (0, 1)$ is **unknown** and $\boldsymbol{\alpha} \in \mathbb{R}^p$ is a (sparse) unknown parameter (identifiable **only** upto scalar multiples, hence set $\|\boldsymbol{\alpha}\|_2 = 1$ wlog).

- **Semi-parametric single index family**: $\pi(\mathbf{X}) = g(\boldsymbol{\alpha}'\mathbf{X})$, where $g(\cdot) \in (0, 1)$ is **unknown** and $\boldsymbol{\alpha} \in \mathbb{R}^p$ is a (sparse) unknown parameter (identifiable **only** upto scalar multiples, hence set $\|\boldsymbol{\alpha}\|_2 = 1$ wlog).
- Given an estimator $\hat{\boldsymbol{\alpha}}$ of $\boldsymbol{\alpha}$, we estimate $\pi(\mathbf{X}) \equiv \mathbb{E}(T \mid \boldsymbol{\alpha}'\mathbf{X})$ as:

$$\hat{\pi}(\mathbf{x}) \equiv \hat{\pi}(\hat{\boldsymbol{\alpha}}, \mathbf{x}) := \frac{\frac{1}{nh} \sum_{i=1}^n T_i K \{ \hat{\boldsymbol{\alpha}}'(\mathbf{X}_i - \mathbf{x})/h \}}{\frac{1}{nh} \sum_{i=1}^n K \{ \hat{\boldsymbol{\alpha}}'(\mathbf{X}_i - \mathbf{x})/h \}},$$

where $K(\cdot)$ denotes any standard (2^{nd} order) kernel function and $h = h_n > 0$ denotes the bandwidth sequence with $h = o(1)$.

- **Semi-parametric single index family**: $\pi(\mathbf{X}) = g(\alpha' \mathbf{X})$, where $g(\cdot) \in (0, 1)$ is **unknown** and $\alpha \in \mathbb{R}^p$ is a (sparse) unknown parameter (identifiable **only** upto scalar multiples, hence set $\|\alpha\|_2 = 1$ wlog).
- Given an estimator $\hat{\alpha}$ of α , we estimate $\pi(\mathbf{X}) \equiv \mathbb{E}(T | \alpha' \mathbf{X})$ as:

$$\hat{\pi}(\mathbf{x}) \equiv \hat{\pi}(\hat{\alpha}, \mathbf{x}) := \frac{\frac{1}{nh} \sum_{i=1}^n T_i K \{ \hat{\alpha}'(\mathbf{X}_i - \mathbf{x})/h \}}{\frac{1}{nh} \sum_{i=1}^n K \{ \hat{\alpha}'(\mathbf{X}_i - \mathbf{x})/h \}},$$

where $K(\cdot)$ denotes any standard (2^{nd} order) kernel function and $h = h_n > 0$ denotes the bandwidth sequence with $h = o(1)$.

- **Obtaining $\hat{\alpha}$** : In general, **any approach** (if available) from (high dimensional) single index model literature can be used. But if \mathbf{X} is elliptically symmetric, then $\hat{\alpha}$ may be obtained as simply as a standard L_1 -penalized logistic regression of $\{T_i$ vs. $\mathbf{X}_i\}_{i=1}^n$.

- 'Extended' parametric family: $m(\mathbf{x}) = g\{\boldsymbol{\gamma}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot)$ is a **known** 'link' function [e.g. 'canonical' links: identity, expit or exp], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\gamma} \in \mathbb{R}^K$ is an unknown (sparse) parameter vector.

- 'Extended' parametric family: $m(\mathbf{x}) = g\{\boldsymbol{\gamma}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot)$ is a **known** 'link' function [e.g. 'canonical' links: identity, expit or exp], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\gamma} \in \mathbb{R}^K$ is an **unknown** (sparse) parameter vector.
- Example: $\boldsymbol{\Psi}(\mathbf{X})$ may correspond to the **polynomial bases** of \mathbf{X} upto any fixed degree k . Note: the **special case of linear bases** ($k = 1$) **includes** all **standard parametric regression models**.

- **'Extended' parametric family:** $m(\mathbf{x}) = g\{\boldsymbol{\gamma}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot)$ is a **known** 'link' function [e.g. 'canonical' links: identity, expit or exp], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\gamma} \in \mathbb{R}^K$ is an **unknown** (sparse) parameter vector.
 - Example: $\boldsymbol{\Psi}(\mathbf{X})$ may correspond to the **polynomial bases** of \mathbf{X} upto any fixed degree k . Note: the **special case** of linear bases ($k = 1$) **includes** all **standard parametric regression models**.
 - **Estimator:** we set $\widehat{m}(\mathbf{X}) = g\{\widehat{\boldsymbol{\gamma}}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $\widehat{\boldsymbol{\gamma}}$ denotes **any suitable estimator** (possibly penalized) of $\boldsymbol{\gamma}$ **based on** the data subset of **'complete cases'**: $\mathcal{D}_n^{(c)} := \{(Y_i, \mathbf{X}_i) \mid T_i = 1\}_{i=1}^n$.

- **'Extended' parametric family:** $m(\mathbf{x}) = g\{\boldsymbol{\gamma}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $g(\cdot)$ is a **known** 'link' function [e.g. 'canonical' links: identity, expit or exp], $\boldsymbol{\Psi}(\mathbf{X}) := \{\psi_k(\mathbf{X})\}_{k=1}^K$ is **any** set of K basis functions (with $K \gg n$ possibly), and $\boldsymbol{\gamma} \in \mathbb{R}^K$ is an unknown (sparse) parameter vector.
 - Example: $\boldsymbol{\Psi}(\mathbf{X})$ may correspond to the **polynomial bases** of \mathbf{X} upto any fixed degree k . Note: the **special case of linear bases** ($k = 1$) **includes** all **standard parametric regression models**.
 - **Estimator:** we set $\widehat{m}(\mathbf{X}) = g\{\widehat{\boldsymbol{\gamma}}'\boldsymbol{\Psi}(\mathbf{X})\}$, where $\widehat{\boldsymbol{\gamma}}$ denotes **any suitable estimator** (possibly penalized) of $\boldsymbol{\gamma}$ **based on** the data subset of **'complete cases'**: $\mathcal{D}_n^{(c)} := \{(Y_i, \mathbf{X}_i) \mid T_i = 1\}_{i=1}^n$.
 - **Example of $\widehat{\boldsymbol{\gamma}}$:** when $g(\cdot) :=$ any 'canonical' link function, $\widehat{\boldsymbol{\gamma}}$ may be simply obtained based on the respective usual **L_1 -penalized 'canonical' link based regression** (e.g. linear, logistic or poisson) of $\{(Y_i \text{ vs. } \mathbf{X}_i) \mid T_i = 1\}_{i=1}^n$ from the **'complete case' data** $\mathcal{D}_n^{(c)}$.

- **Semi-parametric single index family**: $m(\mathbf{X}) = g(\gamma'\mathbf{X})$, where $g(\cdot)$ is an **unknown** 'link' and $\gamma \in \mathbb{R}^p$ is a (sparse) unknown parameter (identifiable **only** upto scalar multiples, hence set $\|\gamma\|_2 = 1$ wlog).

- **Semi-parametric single index family**: $m(\mathbf{X}) = g(\boldsymbol{\gamma}'\mathbf{X})$, where $g(\cdot)$ is an **unknown** 'link' and $\boldsymbol{\gamma} \in \mathbb{R}^p$ is a (sparse) unknown parameter (identifiable **only** upto scalar multiples, hence set $\|\boldsymbol{\gamma}\|_2 = 1$ wlog).
- Given an estimator $\widehat{\boldsymbol{\gamma}}$ of $\boldsymbol{\gamma}$, we **estimate** $m(\mathbf{X}) \equiv \mathbb{E}(Y \mid \boldsymbol{\gamma}'\mathbf{X}, T)$ as:

$$\widehat{m}(\mathbf{x}) \equiv \widehat{m}(\widehat{\boldsymbol{\gamma}}, \mathbf{x}) := \frac{\frac{1}{nh} \sum_{i=1}^n T_i Y_i K \{ \widehat{\boldsymbol{\gamma}}'(\mathbf{X}_i - \mathbf{x})/h \}}{\frac{1}{nh} \sum_{i=1}^n T_i K \{ \widehat{\boldsymbol{\gamma}}'(\mathbf{X}_i - \mathbf{x})/h \}},$$

where $K(\cdot)$ denotes any standard (2^{nd} order) kernel function, and $h = h_n > 0$ denotes the bandwidth sequence with $h = o(1)$.

- **Semi-parametric single index family**: $m(\mathbf{X}) = g(\gamma'\mathbf{X})$, where $g(\cdot)$ is an **unknown** 'link' and $\gamma \in \mathbb{R}^p$ is a (sparse) unknown parameter (identifiable **only** upto scalar multiples, hence set $\|\gamma\|_2 = 1$ wlog).
- Given an estimator $\widehat{\gamma}$ of γ , we **estimate** $m(\mathbf{X}) \equiv \mathbb{E}(Y \mid \gamma'\mathbf{X}, T)$ as:

$$\widehat{m}(\mathbf{x}) \equiv \widehat{m}(\widehat{\gamma}, \mathbf{x}) := \frac{\frac{1}{nh} \sum_{i=1}^n T_i Y_i K\{\widehat{\gamma}'(\mathbf{X}_i - \mathbf{x})/h\}}{\frac{1}{nh} \sum_{i=1}^n T_i K\{\widehat{\gamma}'(\mathbf{X}_i - \mathbf{x})/h\}},$$

where $K(\cdot)$ denotes any standard (2^{nd} order) kernel function, and $h = h_n > 0$ denotes the bandwidth sequence with $h = o(1)$.

- **Obtaining $\widehat{\gamma}$** : In general, **any approach** (if available) from HD SIM literature can be used on the complete case data subset $\mathcal{D}_n^{(c)}$.
 - If \mathbf{X} is elliptically symmetric and $Y = f(\gamma'\mathbf{X}; \epsilon)$ with f **unknown** and $\epsilon \perp\!\!\!\perp (T, \mathbf{X})$, then $\widehat{\gamma}$ may be obtained as **L_1 -penalized IPW estimator $\widehat{\theta}_{IPW}$** for any 'canonical' link based regression problem.

- For either choices of $\hat{\pi}(\cdot)$, assume that the ingredient estimator $\hat{\alpha}$ satisfies: $\|\hat{\alpha} - \alpha\|_1 \lesssim_{\mathbb{P}} a_n$ for some $a_n = o(1)$. Then, under suitable smoothness and tail assumptions, with high probability (w.h.p.),

- For either choices of $\hat{\pi}(\cdot)$, assume that the ingredient estimator $\hat{\alpha}$ satisfies: $\|\hat{\alpha} - \alpha\|_1 \lesssim_{\mathbb{P}} a_n$ for some $a_n = o(1)$. Then, under suitable smoothness and tail assumptions, with high probability (w.h.p.),
 $|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim a_n = o(1)$, for any fixed $\mathbf{x} \in \mathcal{X}$, (for method 1).

- For either choices of $\hat{\pi}(\cdot)$, assume that the ingredient estimator $\hat{\alpha}$ satisfies: $\|\hat{\alpha} - \alpha\|_1 \lesssim_{\mathbb{P}} a_n$ for some $a_n = o(1)$. Then, under suitable smoothness and tail assumptions, with high probability (w.h.p.),

$$|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim a_n = o(1), \text{ for any fixed } \mathbf{x} \in \mathcal{X}, \text{ (for method 1).}$$

- For method 2 (SIM), assume that $h = o(1)$, $\log(np)/(nh) = o(1)$ and $(a_n/h)\sqrt{\log p} = o(1)$. Then, under some suitable smoothness and tail assumptions, we have: with high probability, for any fixed $\mathbf{x} \in \mathcal{X}$,

$$|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim \left(h^2 + \frac{1}{\sqrt{nh}} \right) + \left(a_n + \frac{\log(np)}{nh} + \frac{a_n^2}{h^2} \right) = o(1).$$

- For either choices of $\hat{\pi}(\cdot)$, assume that the ingredient estimator $\hat{\alpha}$ satisfies: $\|\hat{\alpha} - \alpha\|_1 \lesssim_{\mathbb{P}} a_n$ for some $a_n = o(1)$. Then, under suitable smoothness and tail assumptions, with high probability (w.h.p.),
 $|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim a_n = o(1)$, for any fixed $\mathbf{x} \in \mathcal{X}$, (for method 1).

- For method 2 (SIM), assume that $h = o(1)$, $\log(np)/(nh) = o(1)$ and $(a_n/h)\sqrt{\log p} = o(1)$. Then, under some suitable smoothness and tail assumptions, we have: with high probability, for any fixed $\mathbf{x} \in \mathcal{X}$,

$$|\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \lesssim \left(h^2 + \frac{1}{\sqrt{nh}} \right) + \left(a_n + \frac{\log(np)}{nh} + \frac{a_n^2}{h^2} \right) = o(1).$$

- Usually, we expect the L_1 error rate of $\hat{\alpha}$ to be $a_n = s_\alpha \sqrt{(\log d_*)/n}$ where $s_\alpha := \|\alpha\|_0$ and $d_* = K$ or p (depending on the method).

- For either choices of $\widehat{m}(\cdot)$, assume that the ingredient estimator $\widehat{\gamma}$ satisfies: $\|\widehat{\gamma} - \gamma\|_1 \lesssim_{\mathbb{P}} b_n$ for some $b_n = o(1)$. Then, under suitable smoothness and tail assumptions, we have: **with high probability,**

- For either choices of $\widehat{m}(\cdot)$, assume that the ingredient estimator $\widehat{\gamma}$ satisfies: $\|\widehat{\gamma} - \gamma\|_1 \lesssim_{\mathbb{P}} b_n$ for some $b_n = o(1)$. Then, under suitable smoothness and tail assumptions, we have: **with high probability**,
 $|\widehat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim b_n = o(1)$ for any fixed $\mathbf{x} \in \mathcal{X}$ (for method 1).

- For either choices of $\widehat{m}(\cdot)$, assume that the ingredient estimator $\widehat{\gamma}$ satisfies: $\|\widehat{\gamma} - \gamma\|_1 \lesssim_{\mathbb{P}} b_n$ for some $b_n = o(1)$. Then, under suitable smoothness and tail assumptions, we have: **with high probability**,
$$|\widehat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim b_n = o(1) \quad \text{for any fixed } \mathbf{x} \in \mathcal{X} \quad (\text{for method 1}).$$
- For method 2 (SIM), assume that $h = o(1)$, $\log(np)/(nh) = o(1)$ and $(a_n/h)\sqrt{\log p} = o(1)$. Then, under some suitable smoothness and tail assumptions, we have: **with high probability**, for any fixed $\mathbf{x} \in \mathcal{X}$,

$$|\widehat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim \left(h^2 + \frac{1}{\sqrt{nh}} \right) + \left(b_n + \frac{\log(np)}{nh} + \frac{b_n^2}{h^2} \right) = o(1).$$

- For either choices of $\widehat{m}(\cdot)$, assume that the ingredient estimator $\widehat{\gamma}$ satisfies: $\|\widehat{\gamma} - \gamma\|_1 \lesssim_{\mathbb{P}} b_n$ for some $b_n = o(1)$. Then, under suitable smoothness and tail assumptions, we have: **with high probability**,
$$|\widehat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim b_n = o(1) \quad \text{for any fixed } \mathbf{x} \in \mathcal{X} \quad (\text{for method 1}).$$

- For method 2 (SIM), assume that $h = o(1)$, $\log(np)/(nh) = o(1)$ and $(a_n/h)\sqrt{\log p} = o(1)$. Then, under some suitable smoothness and tail assumptions, we have: **with high probability**, for any fixed $\mathbf{x} \in \mathcal{X}$,

$$|\widehat{m}(\mathbf{x}) - m(\mathbf{x})| \lesssim \left(h^2 + \frac{1}{\sqrt{nh}} \right) + \left(b_n + \frac{\log(np)}{nh} + \frac{b_n^2}{h^2} \right) = o(1).$$

- We typically expect the L_1 error rate of $\widehat{\gamma}$ to be $b_n = s_\gamma \sqrt{(\log d_*)/n}$ where $s_\gamma := \|\alpha\|_0$ and $d_* = K$ or p (depending on the method).

Simulation Studies: The Setup

- Basic parameters: $n = 1000$, $p = 50$ or 500 and $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_p)$.
- Three data generating processes (DGPs) for $Y|\mathbf{X}$ and $T|\mathbf{X}$ as follows:

- Basic parameters: $n = 1000$, $p = 50$ or 500 and $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_p)$.
- Three data generating processes (DGPs) for $Y|\mathbf{X}$ and $T|\mathbf{X}$ as follows:
 - ① “Linear-Linear” DGP:

$$Y = \gamma_0 + \boldsymbol{\gamma}'\mathbf{X} + \varepsilon, \quad \varepsilon|\mathbf{X} \sim \mathcal{N}(0, 1).$$
$$\text{logit}\{\pi(\mathbf{X})\} \equiv \text{logit}\{\mathbb{E}(T|\mathbf{X})\} = \alpha_0 + \boldsymbol{\alpha}'\mathbf{X}.$$

- Basic parameters: $n = 1000$, $p = 50$ or 500 and $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_p)$.
- Three data generating processes (DGPs) for $Y|\mathbf{X}$ and $T|\mathbf{X}$ as follows:

① “Linear-Linear” DGP:

$$Y = \gamma_0 + \boldsymbol{\gamma}'\mathbf{X} + \varepsilon, \quad \varepsilon|\mathbf{X} \sim \mathcal{N}(0, 1).$$
$$\text{logit}\{\pi(\mathbf{X})\} \equiv \text{logit}\{\mathbb{E}(T|\mathbf{X})\} = \alpha_0 + \boldsymbol{\alpha}'\mathbf{X}.$$

② “Quad-Quad” DGP:

$$Y = \gamma_0 + \boldsymbol{\gamma}'\mathbf{X} + \sum_{j=1}^p \gamma_j^* \mathbf{X}_j^2 + \varepsilon, \quad \varepsilon|\mathbf{X} \sim \mathcal{N}(0, 1).$$
$$\text{logit}\{\pi(\mathbf{X})\} \equiv \text{logit}\{\mathbb{E}(T|\mathbf{X})\} = \alpha_0 + \boldsymbol{\alpha}'\mathbf{X}_i + \sum_{j=1}^p \alpha_j^* \mathbf{X}_{ij}^2.$$

- Basic parameters: $n = 1000$, $p = 50$ or 500 and $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_p)$.
- Three data generating processes (DGPs) for $Y|\mathbf{X}$ and $T|\mathbf{X}$ as follows:

① “Linear-Linear” DGP:

$$Y = \gamma_0 + \boldsymbol{\gamma}'\mathbf{X} + \varepsilon, \quad \varepsilon|\mathbf{X} \sim \mathcal{N}(0, 1).$$
$$\text{logit}\{\pi(\mathbf{X})\} \equiv \text{logit}\{\mathbb{E}(T|\mathbf{X})\} = \alpha_0 + \boldsymbol{\alpha}'\mathbf{X}.$$

② “Quad-Quad” DGP:

$$Y = \gamma_0 + \boldsymbol{\gamma}'\mathbf{X} + \sum_{j=1}^p \gamma_j^* \mathbf{X}_j^2 + \varepsilon, \quad \varepsilon|\mathbf{X} \sim \mathcal{N}(0, 1).$$
$$\text{logit}\{\pi(\mathbf{X})\} \equiv \text{logit}\{\mathbb{E}(T|\mathbf{X})\} = \alpha_0 + \boldsymbol{\alpha}'\mathbf{X}_i + \sum_{j=1}^p \alpha_j^* \mathbf{X}_{ij}^2.$$

③ “SIM-SIM” DGP:

$$Y = \gamma_0 + \boldsymbol{\gamma}'\mathbf{X} + c_Y(\boldsymbol{\gamma}'\mathbf{X})^2 + \varepsilon, \quad \varepsilon|\mathbf{X} \sim \mathcal{N}(0, 1).$$
$$\text{logit}\{\pi(\mathbf{X})\} \equiv \text{logit}\{\mathbb{E}(T|\mathbf{X})\} = \alpha_0 + \boldsymbol{\alpha}'\mathbf{X} + c_T(\boldsymbol{\alpha}'\mathbf{X})^2.$$

- Choices of the parameters:

- 1 Covariance matrix Σ_p (for today): $\Sigma_p = I_p$ (identity matrix).
- 2 We set $c_T = 0.2$, $c_Y = 0.3$ and $\gamma_0 = 1$, $\alpha_0 = 0.5$.
- 3 When $p = 50$, $\alpha = 1/\sqrt{5}(1, -1, 0.5, -0.5, 0.5, 0, \dots, 0)$ with $\|\alpha\|_0 = 5$,
 $\gamma = (1, 1, 1, -1, -1, 0.5, 0.5, -0.5, -0.5, -0.5, 0, \dots, 0)$ with $\|\gamma\|_0 = 10$,
 $\alpha^* = (0.25, -0.25, 0, \dots, 0)$ and $\gamma^* = (1, -1, 0.5, 0.5, -0.5, 0, \dots, 0)$.

- Choices of the parameters:
 - 1 Covariance matrix Σ_p (for today): $\Sigma_p = I_p$ (identity matrix).
 - 2 We set $c_T = 0.2$, $c_Y = 0.3$ and $\gamma_0 = 1$, $\alpha_0 = 0.5$.
 - 3 When $p = 50$, $\alpha = 1/\sqrt{5}(1, -1, 0.5, -0.5, 0.5, 0, \dots, 0)$ with $\|\alpha\|_0 = 5$, $\gamma = (1, 1, 1, -1, -1, 0.5, 0.5, -0.5, -0.5, -0.5, 0, \dots, 0)$ with $\|\gamma\|_0 = 10$, $\alpha^* = (0.25, -0.25, 0, \dots, 0)$ and $\gamma^* = (1, -1, 0.5, 0.5, -0.5, 0, \dots, 0)$.
 - 4 When $p = 500$, $\|\alpha\|_0 = 10$ and α consists of three 1s, two -1 s, two 0.5s and three -0.5 s normalized by $1/\sqrt{10}$, while $\|\gamma\|_0 = 15$ and γ consists of three 1s, two -1 s, five 0.5s, five -0.5 s, two 0.25s and three -0.25 s. Further, we set $\alpha^* = (0.25, 0.25, -0.25, -0.25, 0, \dots, 0)$ and $\gamma^* = (1, -1, 0.5, 0.5, -0.5, 0, \dots, 0)$.

- Choices of the parameters:
 - 1 Covariance matrix Σ_p (for today): $\Sigma_p = I_p$ (identity matrix).
 - 2 We set $c_T = 0.2$, $c_Y = 0.3$ and $\gamma_0 = 1$, $\alpha_0 = 0.5$.
 - 3 When $p = 50$, $\alpha = 1/\sqrt{5}(1, -1, 0.5, -0.5, 0.5, 0, \dots, 0)$ with $\|\alpha\|_0 = 5$, $\gamma = (1, 1, 1, -1, -1, 0.5, 0.5, -0.5, -0.5, -0.5, 0, \dots, 0)$ with $\|\gamma\|_0 = 10$, $\alpha^* = (0.25, -0.25, 0, \dots, 0)$ and $\gamma^* = (1, -1, 0.5, 0.5, -0.5, 0, \dots, 0)$.
 - 4 When $p = 500$, $\|\alpha\|_0 = 10$ and α consists of three 1s, two -1 s, two 0.5s and three -0.5 s normalized by $1/\sqrt{10}$, while $\|\gamma\|_0 = 15$ and γ consists of three 1s, two -1 s, five 0.5s, five -0.5 s, two 0.25s and three -0.25 s. Further, we set $\alpha^* = (0.25, 0.25, -0.25, -0.25, 0, \dots, 0)$ and $\gamma^* = (1, -1, 0.5, 0.5, -0.5, 0, \dots, 0)$.
- $\mathbb{K} = 2$ fold cross-fitting used; all simulation settings replicated 500 times.
- $\hat{\Omega}$ obtained as $\hat{\Sigma}^{-1}$ for $p = 50$ and using the nodewise Lasso for $p = 500$.

- Obtain the DDR estimator $\hat{\theta}_{\text{DDR}}$ for linear regression: $\theta_0 = \Sigma^{-1}\mathbb{E}(\mathbf{X}Y)$.

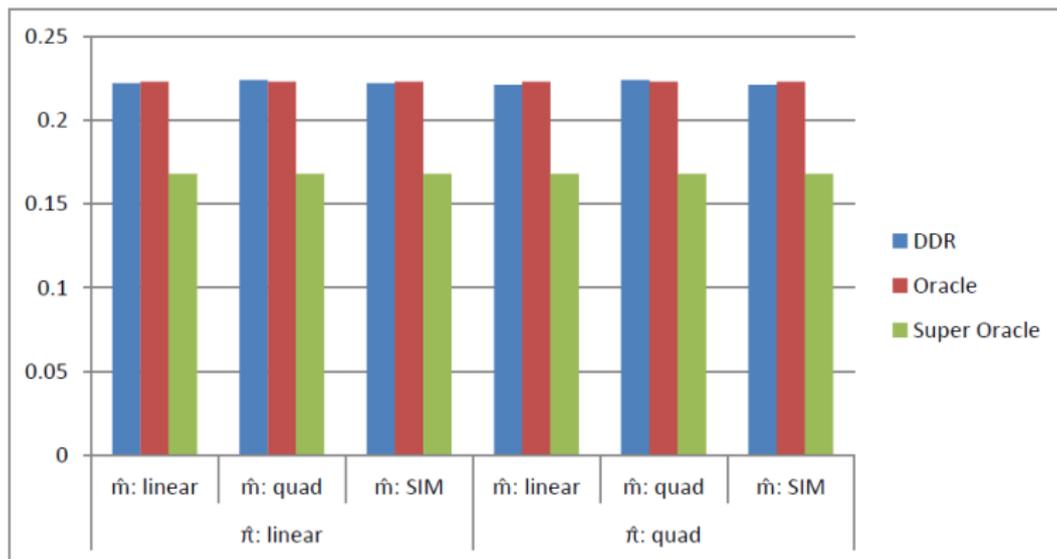
- Obtain the DDR estimator $\hat{\theta}_{\text{DDR}}$ for linear regression: $\theta_0 = \Sigma^{-1}\mathbb{E}(\mathbf{X}Y)$.
- Two choices of the working nuisance models for $\pi(\mathbf{X})$ to obtain $\hat{\pi}(\mathbf{X})$:
 - 1 Linear: L_1 penalized logistic-linear regression.
 - 2 Quad: L_1 penalized logistic-linear regression with quadratic terms.

- Obtain the DDR estimator $\hat{\theta}_{\text{DDR}}$ for linear regression: $\theta_0 = \Sigma^{-1}\mathbb{E}(\mathbf{X}Y)$.
- Two choices of the working nuisance models for $\pi(\mathbf{X})$ to obtain $\hat{\pi}(\mathbf{X})$:
 - 1 Linear: L_1 penalized logistic-linear regression.
 - 2 Quad: L_1 penalized logistic-linear regression with quadratic terms.
- Three choices of the working nuisance models for $m(\mathbf{X})$ to obtain $\hat{m}(\mathbf{X})$:
 - 1 Linear: L_1 penalized linear regression.
 - 2 Quad: L_1 penalized linear regression with quadratic terms.
 - 3 SIM: Single index model (with index parameter estimated via IPW Lasso)

- Obtain the DDR estimator $\hat{\theta}_{\text{DDR}}$ for linear regression: $\theta_0 = \Sigma^{-1}\mathbb{E}(\mathbf{X}Y)$.
- Two choices of the working nuisance models for $\pi(\mathbf{X})$ to obtain $\hat{\pi}(\mathbf{X})$:
 - 1 Linear: L_1 penalized logistic-linear regression.
 - 2 Quad: L_1 penalized logistic-linear regression with quadratic terms.
- Three choices of the working nuisance models for $m(\mathbf{X})$ to obtain $\hat{m}(\mathbf{X})$:
 - 1 Linear: L_1 penalized linear regression.
 - 2 Quad: L_1 penalized linear regression with quadratic terms.
 - 3 SIM: Single index model (with index parameter estimated via IPW Lasso)
- Estimators used for comparison:
 - 1 $\hat{\theta}_{\text{orac}}$ (Oracle): obtained assuming both $\pi(\cdot)$ and $m(\cdot)$ are known.
 - 2 $\hat{\theta}_{\text{full}}$ (Super oracle): obtained assuming a full dataset is observed.
- Criteria: L_2 errors for estimation and coverage probability for inference.

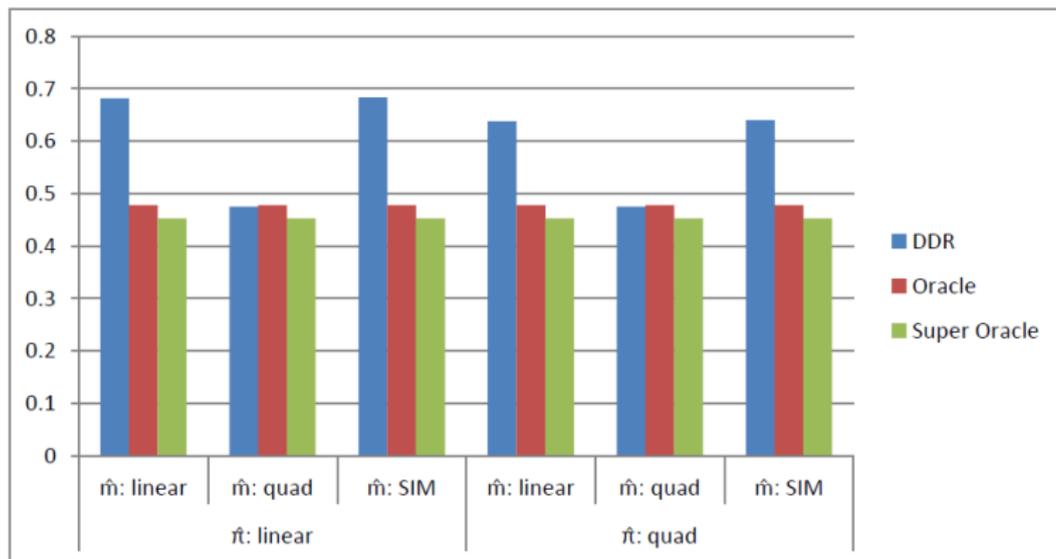
Simulation Results: L_2 Error Comparison ($p = 50$) - I

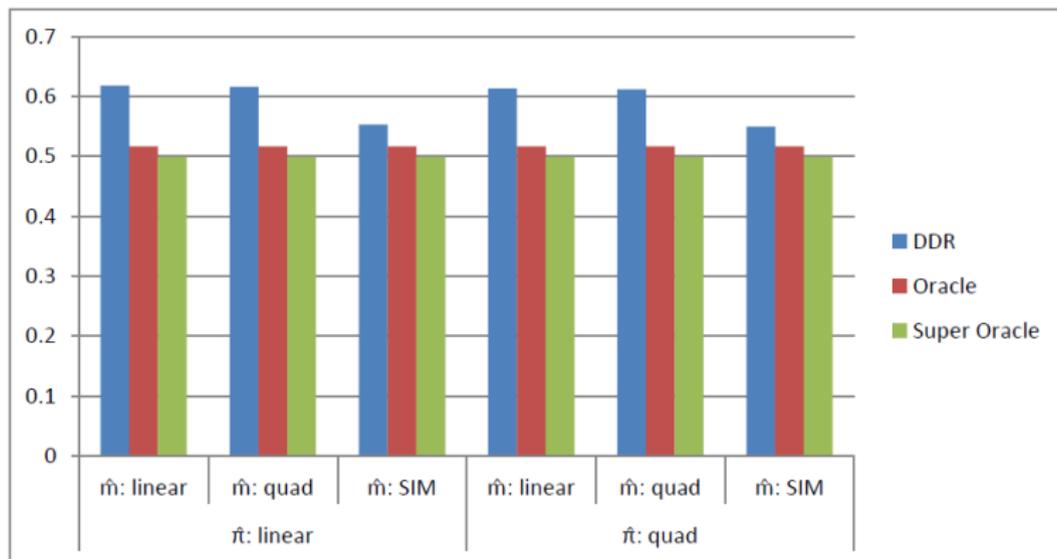
$p = 50$, DGP: Linear-Linear.



Simulation Results: L_2 Error Comparison ($p = 50$) - II

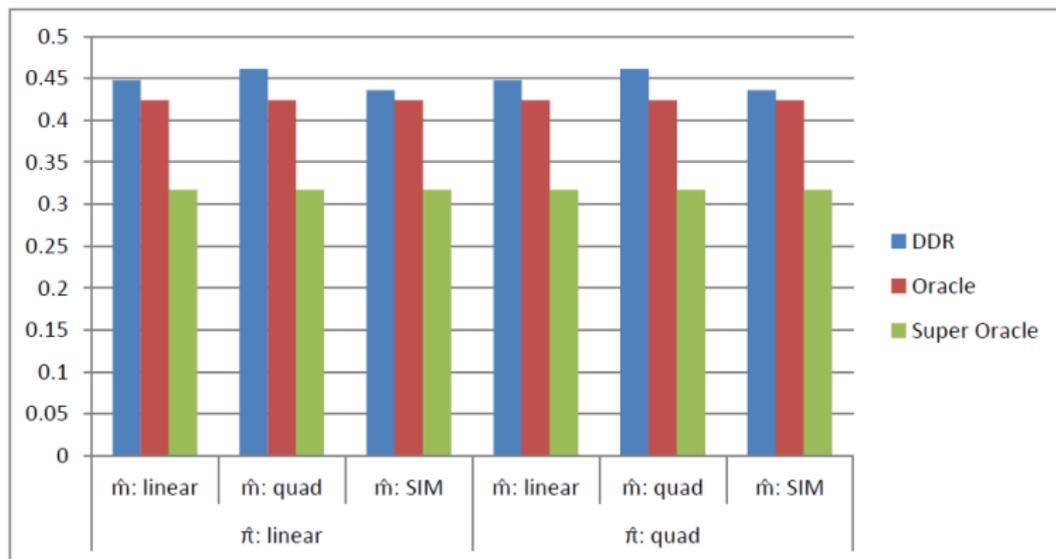
$p = 50$, DGP: Quad-Quad.



$p = 50$, DGP: SIM-SIM.

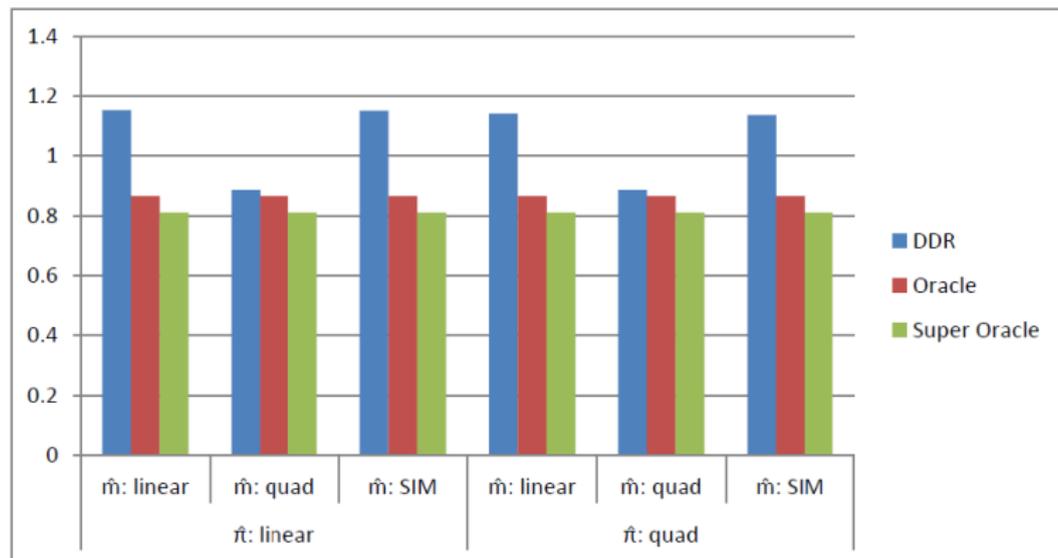
Simulation Results: L_2 Error Comparison ($p = 500$) - I

$p = 500$, DGP: Linear-Linear.



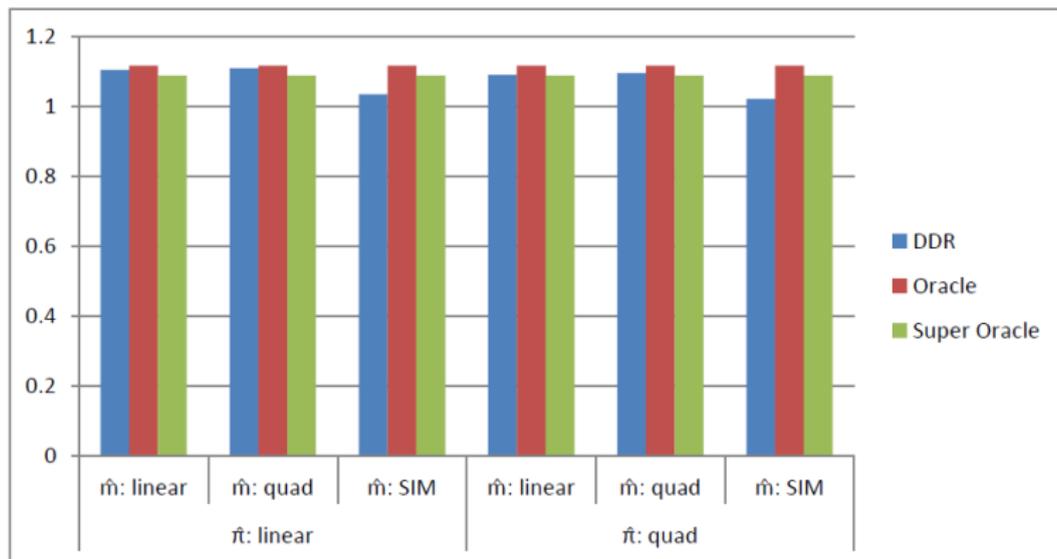
Simulation Results: L_2 Error Comparison ($p = 500$) - II

$p = 500$, DGP: Quad-Quad.



Simulation Results: L_2 Error Comparison ($p = 500$) - III

$p = 500$, DGP: SIM-SIM.



Coverage probability (covg. prob.) of the DDR estimator:

DGP: Linear-Linear.

Coverage probability (covg. prob.) of the DDR estimator:

DGP: Linear-Linear.

① When $p = 50$:

	\hat{m} : linear Average Covg. Prob. (zero coeffs.)	\hat{m} : quad Average Covg. Prob. (zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (zero coeffs.)	\hat{m} : linear Average Covg. Prob. (non-zero coeffs.)	\hat{m} : quad Average Covg. Prob. (non-zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (non-zero coeffs.)
$\hat{\pi}$: logit	0.94 (0.01)	0.94 (0.01)	0.95 (0.01)	0.94 (0.01)	0.94 (0.01)	0.93 (0.01)
$\hat{\pi}$: quad	0.94 (0.01)	0.95 (0.01)	0.95 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)

② When $p = 500$:

	\hat{m} : linear Average Covg. Prob. (zero coeffs.)	\hat{m} : quad Average Covg. Prob. (zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (zero coeffs.)	\hat{m} : linear Average Covg. Prob. (non-zero coeffs.)	\hat{m} : quad Average Covg. Prob. (non-zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (non-zero coeffs.)
$\hat{\pi}$: logit	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	0.92 (0.01)	0.91 (0.02)	0.92 (0.01)
$\hat{\pi}$: quad	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	0.91 (0.02)	0.91 (0.02)	0.92 (0.01)

Coverage probability (covg. prob.) of the DDR estimator:

DGP: Quad-Quad.

Coverage probability (covg. prob.) of the DDR estimator:

DGP: Quad-Quad.

① When $p = 50$:

	\hat{m} : linear Average Covg. Prob. (zero coeffs.)	\hat{m} : quad Average Covg. Prob. (zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (zero coeffs.)	\hat{m} : linear Average Covg. Prob. (non-zero coeffs.)	\hat{m} : quad Average Covg. Prob. (non-zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (non-zero coeffs.)
$\hat{\pi}$: logit	0.94 (0.01)	0.94 (0.01)	0.95 (0.01)	0.88 (0.16)	0.94 (0.01)	0.88 (0.16)
$\hat{\pi}$: quad	0.95 (0.01)	0.94 (0.01)	0.95 (0.01)	0.89 (0.12)	0.94 (0.01)	0.89 (0.12)

② When $p = 500$:

	\hat{m} : linear Average Covg. Prob. (zero coeffs.)	\hat{m} : quad Average Covg. Prob. (zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (zero coeffs.)	\hat{m} : linear Average Covg. Prob. (non-zero coeffs.)	\hat{m} : quad Average Covg. Prob. (non-zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (non-zero coeffs.)
$\hat{\pi}$: logit	0.95 (0.01)	0.94 (0.01)	0.95 (0.01)	0.91 (0.03)	0.92 (0.01)	0.91 (0.05)
$\hat{\pi}$: quad	0.95 (0.01)	0.94 (0.01)	0.95 (0.01)	0.91 (0.03)	0.92 (0.01)	0.91 (0.04)

Coverage probability (covg. prob.) of the DDR estimator:

DGP: SIM-SIM.

Coverage probability (covg. prob.) of the DDR estimator:

DGP: SIM-SIM.

① When $p = 50$:

	\hat{m} : linear Average Covg. Prob. (zero coeffs.)	\hat{m} : quad Average Covg. Prob. (zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (zero coeffs.)	\hat{m} : linear Average Covg. Prob. (non-zero coeffs.)	\hat{m} : quad Average Covg. Prob. (non-zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (non-zero coeffs.)
$\hat{\pi}$: logit	0.94 (0.01)	0.95 (0.01)	0.95 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)
$\hat{\pi}$: quad	0.94 (0.01)	0.95 (0.01)	0.95 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)

② When $p = 500$:

	\hat{m} : linear Average Covg. Prob. (zero coeffs.)	\hat{m} :quad Average Covg. Prob. (zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (zero coeffs.)	\hat{m} : linear Average Covg. Prob. (non-zero coeffs.)	\hat{m} :quad Average Covg. Prob. (non-zero coeffs.)	\hat{m} : SIM Average Covg. Prob. (non-zero coeffs.)
$\hat{\pi}$: logit	0.94 (0.01)	0.95 (0.01)	0.95 (0.01)	0.87 (0.05)	0.88 (0.04)	0.93 (0.02)
$\hat{\pi}$: quad	0.94 (0.01)	0.95 (0.01)	0.95 (0.01)	0.87 (0.05)	0.87 (0.05)	0.93 (0.02)

Consider $n = 50000$ and $p = 50$. In addition, also consider the complete case estimator $\hat{\theta}_{cc}$ (obtained by using only the data with $T_i = 1$).

DGP: Quad-Quad ($p = 50$)

Consider $n = 50000$ and $p = 50$. In addition, also consider the complete case estimator $\hat{\theta}_{cc}$ (obtained by using only the data with $T_i = 1$).

DGP: Quad-Quad ($p = 50$)

L_2 Error Comparison:

model		$\hat{\theta}_{DDR}$	$\hat{\theta}_{orac}$	$\hat{\theta}_{full}$	$\hat{\theta}_{cc}$
\hat{m} : linear	$\hat{\pi}$: logit	0.460 (0.026)	0.072 (0.011)	0.069 (0.01)	0.528 (0.021)
	$\hat{\pi}$: quad	0.204 (0.137)	0.072 (0.011)	0.069 (0.01)	0.528 (0.021)
\hat{m} : quad	$\hat{\pi}$: logit	0.071 (0.010)	0.072 (0.011)	0.069 (0.01)	0.528 (0.021)
	$\hat{\pi}$: quad	0.072 (0.011)	0.072 (0.011)	0.069 (0.01)	0.528 (0.021)
\hat{m} : SIM	$\hat{\pi}$: logit	0.323 (0.019)	0.072 (0.011)	0.069 (0.01)	0.528 (0.021)
	$\hat{\pi}$: quad	0.172 (0.078)	0.072 (0.011)	0.069 (0.01)	0.528 (0.021)

Inference:

	\hat{m} : linear	\hat{m} :quad	\hat{m} : SIM	\hat{m} : linear	\hat{m} :quad	\hat{m} : SIM
	Average Covg. Prob. (zero coeffs.)			Average Covg. Prob. (non-zero coeffs.)		
$\hat{\pi}$: logit	0.94 (0.03)	0.94 (0.03)	0.94 (0.03)	0.68 (0.39)	0.93 (0.03)	0.80 (0.19)
$\hat{\pi}$: quad	0.96 (0.02)	0.94 (0.03)	0.95 (0.02)	0.96 (0.02)	0.94 (0.02)	0.95 (0.02)

Consider $n = 50000$ and $p = 500$. In addition, also consider the complete case estimator $\hat{\theta}_{cc}$ (obtained by using only the data with $T_i = 1$).

DGP: Quad-Quad ($p = 500$)

Consider $n = 50000$ and $p = 500$. In addition, also consider the complete case estimator $\hat{\theta}_{cc}$ (obtained by using only the data with $T_i = 1$).

DGP: Quad-Quad ($p = 500$)

L_2 Error Comparison:

model		$\hat{\theta}_{DDR}$	$\hat{\theta}_{orac}$	$\hat{\theta}_{full}$	$\hat{\theta}_{cc}$
\hat{m} : linear	$\hat{\pi}$: logit	0.297 (0.017)	0.178 (0.009)	0.173 (0.007)	0.325 (0.018)
	$\hat{\pi}$: quad	0.282 (0.113)	0.178 (0.009)	0.173 (0.007)	0.325 (0.018)
\hat{m} : quad	$\hat{\pi}$: logit	0.177 (0.008)	0.178 (0.009)	0.173 (0.007)	0.325 (0.018)
	$\hat{\pi}$: quad	0.180 (0.010)	0.178 (0.009)	0.173 (0.007)	0.325 (0.018)
\hat{m} : SIM	$\hat{\pi}$: logit	0.407 (0.022)	0.178 (0.009)	0.173 (0.007)	0.325 (0.018)
	$\hat{\pi}$: quad	0.294 (0.045)	0.178 (0.009)	0.173 (0.007)	0.325 (0.018)

Inference:

	\hat{m} : linear	\hat{m} :quad	\hat{m} : SIM	\hat{m} : linear	\hat{m} :quad	\hat{m} : SIM
	Average Covg. Prob. (zero coeffs.)			Average Covg. Prob. (non-zero coeffs.)		
$\hat{\pi}$: logit	0.95 (0.02)	0.95 (0.02)	0.95 (0.02)	0.78 (0.32)	0.94 (0.02)	0.75 (0.38)
$\hat{\pi}$: quad	0.95 (0.02)	0.95 (0.02)	0.95 (0.02)	0.94 (0.04)	0.94 (0.02)	0.88 (0.12)

- Bang, H. and Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973.
- Belloni, A., Chernozhukov, V., Fernández-Val, I., and Hansen, C. (2017). Program evaluation and causal inference with high-dimensional data. *Econometrica*, 85(1):233–298.
- Belloni, A., Chernozhukov, V., and Hansen, C. (2014). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68.
- Chernozhukov, V., Demirer, M., Duflo, E., and Fernandez-Val, I. (2017). Generic machine learning inference on heterogenous treatment effects in randomized experiments. *ArXiv preprint arXiv:1712.04802*.
- Chernozhukov, V. and Semenova, V. (2017). Simultaneous inference for best linear predictor of the conditional average treatment effect and other structural functions. *ArXiv preprint arXiv:1702.06240v2*.
- Farrell, M. H. (2015). Robust inference on average treatment effects with possibly more covariates than observations. *Journal of Econometrics*, 189(1):1–23.
- Li, K.-C. and Duan, N. (1989). Regression analysis under link violation. *The Annals of Statistics*, 17(3):1009–1052.

- Loh, P.-L. and Wainwright, M. J. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *The Annals of Statistics*, 40(3):1637.
- Loh, P.-L. and Wainwright, M. J. (2015). Regularized M -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*, 16:559–616.
- Negahban, S. N., Ravikumar, P., Wainwright, M. J., and Yu, B. (2012). A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers. *Statistical Science*, 27(4):538–557.
- Robins, J. M., Rotnitzky, A., and Zhao, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of the American statistical Association*, 89(427):846–866.
- Smucler, E., Rotnitzky, A., and Robins, J. M. (2019). A unifying approach for doubly-robust l_1 regularized estimation of causal contrasts. *ArXiv preprint arXiv:1904.03737v1*.
- Tsiatis, A. (2007). *Semiparametric Theory and Missing Data*. Springer Science & Business Media.
- van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202.
- Van der Vaart, A. W. (2000). *Asymptotic Statistics*, volume 3. Cambridge University Press.

Thank You!