

Doubly robust inference when combining probability and non-probability samples with high dimensional data

Ma Ting Tin

Department of Statistics, Chinese University of Hong Kong

2021 Summer Lab Meeting

August 17, 2021

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Introduction

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Definition 1.1 (Probability and Non-probability sampling)

In general, we want to understand the DGP or some statistical quantities based on the available sample. There are two type of sampling

- ▶ **[Probability sampling]** The probability samples are selected under known sampling designs and thus are representative of target distribution.
- ▶ **[Non-probability sampling]** Involves non-random selection based on convenience or other criteria, allowing you to easily collect data.
Examples: remote sensing data and web-based volunteer samples

1. Non-probability sample provide rich information about the DGP and might help finite population inference.
2. **AIM:** Propose data integration methods that take advantages of both probability and non-probability sampling.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Definition 1.2 (Auxiliary variable)

Auxiliary variable is a variable of the Sampling frame which is neither the contact variable, nor the identifier.

💡: An Auxiliary variable is any Variable about which information is available prior to data collection and this information is known for all units of the population.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

There are three types of existing methods for data integration.

1. Propensity score adjustment (Rosenbaum and Rubin, 1983)

Model the probability of a unit being selected into non-probability sample followed by estimation and selection biases adjustment (EG: by propensity score weighting or Stratification).

💡: Highly biased and variable under model mis-specification.

2. Calibration weighting (Deville and Sarndal, 1992)

IDEA: Forces the moments or the empirical distribution of auxiliary variable between probability and non-probability sample to be the same

3. Mass imputation: Regard non-probability sample as training data set, and apply imputation to all units in probability sample.

Basic Set-up

1. Introduction
- 2. Basic set-up**
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Notation: Two Samples

Consider a population with finite size N (N is known)

- ▶ Index set of N units $\mathcal{U} := \{1, \dots, N\}$
- ▶ Population consist of $\mathcal{F}_N := \{(X_i, Y_i) : i \in \mathcal{U}\}$
- ▶ **Parameter of interest:** $\mu := N^{-1} \sum_{i=1}^N Y_i$.

Denote the two group of sample as follow:

1. Sample A: probability sample with size n_A .
2. Sample B: non-probability sample with size n_B .

There are some more notation:

- ▶ $\pi_{A,i} := P(i \in A)$, i.e. probability of the i -th individual being included in the probability sample.
- ▶ $\mathbb{1}_{\#,i} := \mathbb{1}(i\text{-th individual is included in sample } \#)$.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Sample	Size	Sampling Weigting π^{-1}	Covariate X	Study Variable Y
A	n_A	✓	✓	✗
B	n_B	✗	✓	✓

To summarize, our dataset is as follow

- ▶ Sample A: $\mathcal{O}_A = \{(d_{A,i}, X_i) : i \in A\}$, where $d_{A,i} := \pi_{A,i}^{-1}$.
- ▶ Sample B: $\mathcal{O}_B = \{(X_i, Y_i) : i \in B\}$.
- ▶ $\mathbb{1}_{A,i}$ and $\mathbb{1}_{B,i}$ for $i = 1, \dots, N$

Remark 0.1

As sampling mechanism for sample B is unknown, Horvitz-Thompson estimation cannot be applied.

The following identification assumption is made:

1. $\pi_B(X) := P(\mathbb{1}_B = 1|X, Y) = P(\mathbb{1}_B = 1|X)$.
2. $\exists \delta_B > 0$ and $\gamma \in (2/3, 1]$ s.t. $\forall X, \pi_B(X) > N^{\gamma-1} \delta_B > 0$.

We then interpret the meaning behind the assumption

1. (1) implies that the selection indicator $\mathbb{1}_B$ and Y are
 $m(X) := E(Y|X) = E(Y|X, \mathbb{1}_B = 1)$ can be estimated solely on basis of sample B .
2. (2) implies $n_B^{-1} = O(N^{-\gamma})$.

Remark 0.2

Notice that assumption (1) is NOT verifiable from the observed data. One should consider more possible predictors for the selection indicator $\mathbb{1}_B$ or outcome Y , resulting in a rich set of variables in X .

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Definition 2.1 (Sampling score)

Sampling (or propensity) score is probability of a unit being selected into the non-probability sample B. Let $\pi_B(X)$ be the sampling score function given X .

Let $\pi_B(X^T\alpha)$ and $m(X^T\beta)$ be the postulated models for $\pi_B(X)$ and $m(X)$ respectively, where α, β are unknown parameters.

- ▶ $\pi_B(X)$ and $m(X)$ are unknown in practice and need to be estimated.
- ▶ Various estimator for μ have been proposed, each requiring different model assumption and estimation strategies.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Definition 2.2 (Inverse Probability Weighting Estimator (IPW))

Given estimator of α , $\hat{\alpha}$, IPW is defined as

$$\hat{\mu}_{IPW} = \hat{\mu}_{IPW}(\hat{\alpha}) = \frac{1}{N} \sum_{i=1}^N \frac{\mathbb{1}_{B,i}}{\pi_B(X_i^T \hat{\alpha})} Y_i.$$

Remark 0.3

Justification for $\hat{\mu}$ requires $\pi_B(X) = \pi_B(X_i^T \alpha)$ and $\hat{\alpha} \xrightarrow{pr} \alpha$.

The estimator $\hat{\alpha}$ can be obtained through following estimating equation:

$$\sum_{i=1}^N \left\{ \frac{\mathbb{1}_{B,i}}{\pi(X_i^T \alpha)} - \frac{\mathbb{1}_{A,i}}{\pi_{A,i}} \right\} h(X_i; \alpha) = 0,$$

where possible choices of $h(X_i; \alpha)$ could be $h(X; \alpha) = X$ or $h(X; \alpha) = \pi(X^T; \alpha)X$.

Definition 2.3 (Outcome regression based on sample A)

The outcome regression estimator is given by

$$\hat{\mu}_{reg} = \hat{\mu}_{reg}(\hat{\beta}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{A,i} d_{A,i} m(X_i^T \hat{\beta}),$$

where $\hat{\beta}$ is obtained by fitting the outcome model based solely on \mathcal{O}_B .

Remark 0.4

- ▶ Justification for $\hat{\mu}_{reg}$ requires $m(X) = m(X_i^T \beta)$ and $\hat{\alpha} \xrightarrow{pr} \alpha$. If one of them are not satisfied, the estimator could be biased.
- ▶ Similar form with IPW, but in IPW, the sampling score has to be estimated and the response value is known, which is exactly opposite with this estimator.

Definition 2.4 (Calibration weighting)

The calibration weighting estimator is given by

$$\hat{\mu}_{cal} = \frac{1}{N} \sum_{i=1}^N \omega_i \mathbb{1}_{B,i} Y_i,$$

where $\{\omega_i : i \in \mathcal{U}\}$ satisfies either

- ▶ $\sum_{i \in A} d_{A,i} X_i = \sum_{i \in B} \omega_i X_i$.
- ▶ $\sum_{i \in A} d_{A,i} m(X_i; \hat{\beta}) = \sum_{i \in B} \omega_i m(X_i; \hat{\beta})$.

Remark 0.5

Justification for $\hat{\mu}_{cal}$ by constraint (1) requires linearity of

- ▶ *Outcome model, i.e. $m(X) = X^t \beta^*$ for some β^* ; or*
- ▶ *Inverse probability of sampling weight, i.e. $\pi_B(X)^{-1} = X^T \alpha^*$ for some α^**

while it requires well-specification of $m(X; \beta)$ for constraint (2).

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Definition 2.5 (Doubly Robust Estimator)

The doubly robust estimator is given by

$$\hat{\mu}_{dr} = \hat{\mu}_{dr}(\hat{\alpha}, \hat{\beta}) = \frac{1}{N} \sum_{i=1}^N \left[\frac{\mathbb{1}_{B,i}}{\pi_B(X_i^T \hat{\alpha})} \{Y_i - m(X_i^T \hat{\beta})\} + \mathbb{1}_{A,i} d_{A,i} m(X_i^T \hat{\beta}) \right],$$

Remark 0.6

The estimator above is doubly robust in the following sense: $\hat{\mu}_{dr} \xrightarrow{pr} \mu$ if **either**

- ▶ $\pi_B(X^T \alpha)$ is correctly specified; or
- ▶ $m(X^T \beta)$ is correctly specified.

Note that the consistency hold without requiring both of them to be well-specified.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Methodology in high dimensional data

Including unnecessary covariates in model increases the bias and variance. Variable selection is important to handle high dimensional covariates.

Let $\alpha \in \mathbb{R}^p$ and $\mathcal{J} \subseteq \{1, \dots, p\}$. Then we define the following

- ▶ $\|\alpha\|_0 := \sum_{j=1}^p \mathbb{1}(\alpha_j \neq 0)$
- ▶ L_1 -norm: $\|\alpha\|_1 := \sum_{j=1}^p |\alpha_j|$
- ▶ L_2 -norm: $\|\alpha\|_2 := \sum_{j=1}^p \sqrt{\alpha_j^2}$
- ▶ L_∞ -norm: $\|\alpha\|_\infty := \max_{1 \leq j \leq p} |\alpha_j|$

Define $\alpha_{\mathcal{J}}$ as subvector of α with indices are in \mathcal{J} and $\mathcal{J}^c = \{1, \dots, p\} \setminus \mathcal{J}$.

Let $\Sigma \in \mathbb{R}^{p \times p}$ and $\Sigma_{\mathcal{J}_1, \mathcal{J}_2}$ be submatrix of Σ formed by rows in \mathcal{J}_1 and columns in \mathcal{J}_2 .

Remark 0.7

Covariates are standardized in prior to the inference so that they have variances approximately equal 1, which stabilize the variable selection procedure.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Assumption 2 (sampling score model):

Assume that $\pi_B(X)$ follows logistic regression model, i.e. $\text{logit}\{\pi_B(X^T\alpha)\} = X^T\alpha$.

The true parameter α , α^* is the parameter that minimize Kullback-Leibler distance between $\pi_B(X)$ and $\pi_B(X^T\alpha)$, i.e.

$$\alpha^* := \arg \min_{\alpha \in \mathbb{R}^p} \mathbb{E} \left[\pi_B(X) \log \left\{ \frac{\pi_B(X)}{\pi_B(X^T\alpha)} \right\} + \{1 - \pi_B(X)\} \log \left\{ \frac{1 - \pi_B(X)}{1 - \pi_B(X^T\alpha)} \right\} \right].$$

Remark 0.8

One can extend the proposed framework to other models such as probit model.

Assumption 3 (outcome model): Assume that $m(X)$ follows a generalized linear regression model, i.e. $m(X) = m(X^T\beta)$ and the true parameter β^* is defined as

$$\beta^* := \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E} \left[\{Y - m(X^T\beta)\}^2 \right].$$

Step 1: Variable Selection

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

The following are estimating functions for α and β

$$U_1(\alpha) = \sum_{i=1}^N \left\{ \frac{\mathbb{1}_{B,i}}{\pi(X_i^T \alpha)} - \frac{\mathbb{1}_{A,i}}{\pi_{A,i}} \right\} X_i \quad \text{and} \quad U_2(\beta) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{B,i} \{Y_i - m(X_i^T \beta)\} X_i,$$

denote $U(\theta) = (U_1(\alpha)^T, U_2(\beta^T))^T$ as the joint estimating function for $\theta := (\alpha^T, \beta^T)^T$

For large p , consider penalized estimating function as

$$U^p(\alpha, \beta) = U(\alpha, \beta) - \begin{pmatrix} q_{\lambda_\alpha}(|\alpha|) \text{sgn}(\alpha) \\ q_{\lambda_\beta}(|\beta|) \text{sgn}(\beta) \end{pmatrix},$$

where

$$q_\lambda(|\theta|) := \lambda \left\{ \mathbb{1}(|\theta| < \lambda) + \frac{(a\lambda - |\theta|)_+}{(a-1)\lambda} \mathbb{1}(|\theta| \geq \lambda) \right\}$$

Step 1: Variable Selection

In order to have better understanding on the penalty function, we rewrite

$$q_{\lambda}(|\theta|) = \begin{cases} \lambda & , \text{ if } |\theta| < \lambda \\ \frac{a\lambda - |\theta|}{(a-1)\lambda} & , \text{ if } \lambda \leq |\theta| < a\lambda \\ 0 & , \text{ if } |\theta| \geq a\lambda \end{cases}$$

Therefore, it is not hard to see that

- ▶ λ represent the penalty size for small $|\theta|$ and it defines the term "small".
- ▶ a defines the term "large" so that no penalty given to $|\theta| \geq a\lambda$.

The penalty of this model is different from lasso/ridge regression. They add penalty s.t. large value of coefficient is not favourable, while it is the opposite in our case.

- ▶ For large $|\alpha_j|$, $q_{\lambda_{\alpha}}(|\alpha_j|) = 0$ and no penalty added.
- ▶ For small $|\alpha_j|$, the penalty forces $\hat{\alpha}_j$ to be zero and excludes j -th element of X from final selected set of variables.

The similar discussion applies to that of $U_2(\beta)$ and $q_{\lambda_{\beta}}$.

Step 2: Doubly Robust Estimation

After obtaining the coefficient estimate in Step 1, denote \mathcal{C} as the index set s.t. those **variable are associated with non-zero coefficient estimate**.

AIM: Re-estimate parameters $\hat{\alpha}, \hat{\beta}$ in the doubly robust estimator $\hat{\mu}_{dr} = \hat{\mu}_{dr}(\hat{\alpha}, \hat{\beta})$ by using only $X_{\mathcal{C}}$ to minimize the asymptotic squared bias of $\hat{\mu}_{dr}$, denoted as $a.bias(\alpha, \beta)$.

The $a.bias(\alpha^*, \beta^*)$ is in the following form

$$E \left[\frac{1}{N} \sum_{i=1}^N \left\{ \frac{\mathbb{1}_{B,i}}{\pi_B(X_i^T \alpha^*)} - 1 \right\} \{Y_i - m(X_i^T \beta^*)\} \right] + E \left\{ \frac{1}{N} \sum_{i=1}^N (\mathbb{1}_{A,i} d_{A,i} - 1) m(X_i^T \beta^*) \right\}$$

So we can compute $\nabla a.bias(\alpha^*, \beta^*)$ and regard its empirical version as the estimating function, i.e.

$$J(\alpha, \beta) := \begin{pmatrix} \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{B,i} \left\{ \frac{1}{\pi_B(X_i^T \alpha)} - 1 \right\} \{Y_i - m(X_i^T \beta)\} X_{i\mathcal{C}} \\ \frac{1}{N} \sum_{i=1}^N \left\{ \frac{\mathbb{1}_{B,i}}{\pi_B(X_i^T \alpha)} - d_{A,i} \mathbb{1}_{A,i} \right\} \frac{\partial}{\partial \beta_{\mathcal{C}}} m(X_i^T \beta) \end{pmatrix}$$

constrained on $\{(\alpha^T, \beta^T)^T \in \mathbb{R}^{2p} : \alpha_{\mathcal{C}^c} = \beta_{\mathcal{C}^c} = 0\}$.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Remark 0.9

1. *Estimating function in Step (1) separates selection for α and β , thus it stabilizes the selection procedure if either one of model is misspecified.*
2. *Estimating function in Step (2) leads to the consistency of $\hat{\mu}_{dr}$ when either one of the model is well-specified.*

In step (2), we consider union of covariates $X_{\mathcal{C}}$, where $\mathcal{C} = \hat{\mathcal{M}}_{\alpha} \cup \hat{\mathcal{M}}_{\beta}$. While there are two other common choices in the literature, which may NOT be robust to model misspecification.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
- 4. Computation**
5. Asymptotic results for variable selection and estimation
6. Simulation

Computation

Computation: Step 1 (MM Algorithm)

Definition 4.1 (Marjorize and Minorize)

A function $g(\theta|\theta^{(t)})$ is said to be majorize (resp. minorize) the function $f(\theta)$ at $\theta^{(t)}$ if

$$g(\theta^{(t)}|\theta^{(t)}) = f(\theta^{(t)}) \quad \text{and} \quad g(\theta|\theta^{(t)}) \geq f(\theta) \quad (\text{resp. } g(\theta|\theta^{(t)}) \leq f(\theta)) \quad \forall \theta$$

Suppose we want to minimize $f(\theta)$. Define $\theta^{(t+1)} = \arg \min_{\theta} g(\theta|\theta^{(t)})$. Then we have

$$f(\theta^{(t+1)}) \leq g(\theta^{(t+1)}|\theta^{(t)}) \leq g(\theta^{(t)}|\theta^{(t)}) = f(\theta^{(t)}).$$

Remark 0.10

- ▶ g is commonly chosen as lower/upper bound of f which is easier to be handled.
- ▶ EM algorithm is a particular case of MM algorithm by setting

$$g(\theta|\theta^{(t)}) := Q(\theta|\theta^{(t)}) - \sum_{i=1}^N \int_{\mathcal{X}} f(z|y_n; \theta^{(t)}) \ln f(z|y_n; \theta^{(t)}) dz.$$

Computation: Step 1 (MM Algorithm)

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

AIM: Solving the penalized estimating equation.

Step 1: By Minorization-maximization (MM) Algorithm, the estimator $\tilde{\theta} = (\tilde{\alpha}, \tilde{\beta})$ satisfies

$$U^p(\tilde{\theta}) = U(\tilde{\theta}) - \begin{pmatrix} q_{\lambda_{\tilde{\alpha}}}(|\tilde{\alpha}|)\text{sgn}(\tilde{\alpha})\frac{|\tilde{\alpha}|}{\epsilon+|\tilde{\alpha}|} \\ q_{\lambda_{\tilde{\beta}}}(|\tilde{\beta}|)\text{sgn}(\tilde{\beta})\frac{|\tilde{\beta}|}{\epsilon+|\tilde{\beta}|} \end{pmatrix} = 0$$

Computation: Step 2 (Coordinate Descent Algorithm)

Definition 4.2 (Coordinate Descent Algorithm)

AIM: Solve $F(x) = 0$ for multivariate function F , where $x \in \mathbb{R}^n$.

1. Initialize $x^0 = (x_1^{(0)}, \dots, x_n^{(0)})$ and step size α (could be non-fixed).
2. There are two possible choice in this step
 - ▶ Update $x_i^{(k+1)} = x_i^{(k)} - \alpha \frac{\partial F}{\partial x_i^{(k)}}(x^{(k)})$. (Similar to the idea of gradient descent).
 - ▶ Set $x_i^{(k+1)} = \arg \min_{y \in \mathbb{R}} f(x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, y, x_{i+1}^{(k)}, \dots, x_n^{(k)})$.

Remark 0.11

- ▶ *Coordinate Descent algorithm prevent the evaluation of large matrix inverse in the Newton-Raphson Algorithm.*
- ▶ *Within the two possible choice, the first one is only available for continuously differentiable function F , while the second one work in general scenario*

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. **Computation**
5. Asymptotic results for variable selection and estimation
6. Simulation

Computation: Step 2 (Coordinate Descent Algorithm)

Step 2: Solve $U^P(\tilde{\theta}) = 0$ by co-ordinate decent algorithm (instead of Newton-Raphson algorithm to lower computational burden): For $k \geq 1$, we write

$$\nabla U(\theta) = \text{diag} \left\{ \frac{\partial U_1(\alpha)}{\partial \alpha^T}, \frac{\partial U_2(\beta)}{\partial \beta^T} \right\}$$

where

$$\frac{\partial U_1(\alpha)}{\partial \alpha^T} = -\frac{1}{N} \sum_{i=1}^N \mathbb{1}_{B,i} \frac{1 - \pi_B(X_i^T \alpha)}{\pi_B(X_i^T \alpha)} X_i X_i^T \text{ and } \frac{\partial U_2(\beta)}{\partial \beta^T} = -\frac{1}{N} \sum_{i=1}^N \mathbb{1}_{B,i} m^{(1)}(X_i^T \beta)^2 X_i X_i^T.$$

With the other coordinates fixed, the coordinate decent algorithm (k -th Newton-Raphson update) for θ_j is

$$\tilde{\theta}_j^{[k]} = \tilde{\theta}_j^{[k-1]} + \frac{U_j(\tilde{\theta}_j^{[k-1]}) - Nq_{\lambda_j}(|\theta_j^{[k-1]}|)\tilde{\theta}_j^{[k-1]}}{\nabla_{jj}(\tilde{\theta}_j^{[k-1]}) + Nq_{\lambda_j}(|\theta_j^{[k-1]}|)}$$

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Computation: Step 3 (K -fold Cross Validation)

Step 3: Apply K -fold CV to select the hyperparameter $(\lambda_\alpha, \lambda_\beta)$.

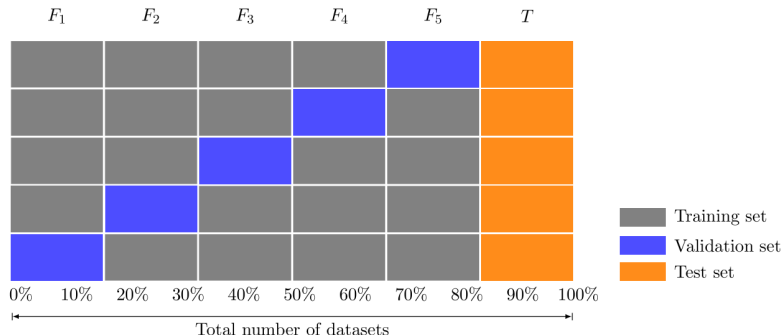


Figure: Idea of K -fold cross validation (Picture extracted from course of STAT4012)

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Computation: Step 3 (K -fold Cross Validation)

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

In this paper, the author suggest to tune λ_α and λ_β one by one. First, as we want to search for the λ_α s.t. the "behaviour" of two sample are similar, so the loss function for λ_α is chosen as

$$L(\lambda_\alpha) := \sum_{j=1}^p \left(\sum_{i=1}^N \left[\frac{\mathbb{1}_{B,i}}{\pi_B \{X_i^T \tilde{\alpha}(\lambda_\alpha)\}} - \frac{\mathbb{1}_{A,i}}{\pi_{A,i}} \right] X_{i,j} \right)^2$$

and we use thre prediction error loss function for selection of λ_β , i.e.

$$L(\lambda_\beta) := \sum_{i=1}^N \mathbb{1}_{B,i} \left[Y_i - m \left\{ X_i^T \tilde{\beta}(\lambda_\beta) \right\} \right]^2$$

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Asymptotic results for variable selection and estimation

Consider sample A with high entropy sampling designs (Berger, 1998). Let $\mathcal{M}_\alpha = \{1 \leq j \leq p : \alpha_j^* \neq 0\}$, $\mathcal{M}_\beta = \{1 \leq j \leq p : \beta_j^* \neq 0\}$ and $\mathcal{M}_\theta := \mathcal{M}_\alpha \cup \{p + \mathcal{M}_\beta\}$. Define $s_\alpha = \|\alpha^*\|_0$, $s_\beta = \|\beta^*\|_0$, $s_\theta = s_\alpha + s_\beta$ and $\lambda_\theta = \min(\lambda_\alpha, \lambda_\beta)$. Assume the following:

1. $\exists \delta_A$ and $\gamma \in (2/3, 1]$ s.t. for all $1 \leq i \leq N$, $\pi_{A,i} \geq N^{\gamma-1} \delta_A > 0$.
2. θ belongs to a compact subset in \mathbb{R}^{2p} and θ^* lies in the interior of that set.
3. $\{X_i : i \in \mathcal{U}\}$ are fixed and uniformly bounded.
4. $\exists c_1, c_2$ s.t. $0 < c_1 \leq \lambda_{\min}(\frac{1}{N} \sum_{i=1}^N X_i^T X_i) \leq \lambda_{\max}(\frac{1}{N} \sum_{i=1}^N X_i^T X_i) \leq c_2 < \infty$, where $\lambda_{\min}(\cdot)$ and $\lambda_{\max}(\cdot)$ refers to the min/max eigenvalue of a matrix (\cdot) respectively.
5. Let $\epsilon_i(\beta) := Y_i - m(X_i^T \beta)$. $\exists c_3, c_4, c_5$ and $\delta > 0$ s.t. $E|\epsilon_i(\beta^{2+\delta})| \leq c_3$ and $E[\exp\{c_4 |\epsilon_i(\beta^*)|\} | X_i] \leq c_5$ for $1 \leq i \leq N$.
6. Let $n := \min(n_A, n_B)$. Assume $m^{(k)}(X_i^T \beta)$ being uniformly bounded away from ∞ on $\mathcal{N}_{\theta, \tau} = \{\theta \in \mathbb{R}^{2p} : \|\theta_\theta - \theta_\theta^*\| \leq \tau \sqrt{s_\theta/n}, \theta_{\theta^c} = 0\}$ for some $\tau > 0$ and $k = 1, 2, 3$.
7. $\min_{j \in \mathcal{M}_\alpha} |\alpha_j^*|/\lambda_\alpha \rightarrow \infty$ and $\min_{k \in \mathcal{M}_\beta} |\beta_k^*|/\lambda_\beta \rightarrow \infty$ as $n \rightarrow \infty$.
8. $s_\theta = o(n^{1/3})$, $\lambda_\alpha, \lambda_\beta \rightarrow 0$, $\log(n)^2 = o(n\lambda_\theta^2)$, $\log(p) = o(n\lambda_\theta^2/\log(n)^2)$, $ps_\theta^4 \log(n)^6 = o(n^3\lambda_\theta^2)$ and $ps_\theta^4 \log(n)^8 = o(n^4\lambda_\theta^4)$ as $n \rightarrow \infty$.

Remark 0.12

1. *If sample A is collected through simple random sampling or poisson sampling, only condition (1) is required.*
2. *Condition (3) specifies a fixed design which is well suited under the finite population inference framework.*
3. *Condition (5) holds for Gaussian distribution, sub-Gaussian distribution (strong tail decay).*
4. *Condition (6) holds for common model.*
5. *Condition (8) specifies the restriction on dimension of covariates p and dimension of the r true non-zero coefficients s_θ .*

Theorem 5.1

Under the assumption stated above, $\exists \tilde{\theta}$ s.t.

▶ $P(|U_j^p(\tilde{\theta})| = 0, j \in \mathcal{M}_\theta) \rightarrow 1.$

$$P\left(\left|U_j^p(\tilde{\theta})\right| \leq \frac{\lambda_\theta}{\log(n)}, j \in \mathcal{M}_\theta^c\right) \rightarrow 1.$$

which implies that $U(\tilde{\theta}) = o_p\{\lambda_\theta / \log(n)\}$. Also,

▶ $P(\tilde{\theta}_{\mathcal{M}_\theta^c} = 0) \rightarrow 1.$

▶ $\tilde{\theta}_{\mathcal{M}_\theta} - \theta_{\mathcal{M}_\theta}^* = O_p\{\sqrt{s_\theta/n}\}.$

implying that for large n , the penalized estimating equation procedure would NOT overselect irrelevant variables and estimate the true non-zero coefficients at the $\sqrt{s_\theta/n}$ convergence rate, which is the so-called oracle property of variable selection.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Under the assumption stated before, if either $\pi_B(X^T \alpha)$ or $m(X^T \beta)$ is correctly specified, then we have

$$\sqrt{n}(\hat{\mu}_{p-dr} - \mu) \xrightarrow{d} N(0, V),$$

where $V := \lim_{n \rightarrow \infty} (V_1 + V_2)$, and

$$V_1 := E \left\{ \frac{n}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\pi_{A,ij} - \pi_{A,i} \pi_{A,j}) \frac{m(X_i^T \beta^*) m(X_j^T \beta^*)}{\pi_{A,i} \pi_{A,j}} \right\}$$

$$V_2 := E \left\{ \frac{n}{N^2} \sum_{i=1}^N \left(\frac{\mathbb{1}_{B,i}}{\pi_{B,i}(X_i^T \alpha^*)} - 1 \right)^2 \left(Y_i - m(X_i^T \beta^*) \right)^2 \right\}$$

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

We can estimate V_1 by its empirical version \hat{V}_1 . For V_2 , we estimate it by

$$\hat{V}_2 = \frac{n}{N^2} \sum_{i=1}^N \left[\left\{ \frac{\mathbb{1}_{B,i}}{\pi_B(X_i^T \hat{\alpha})} - \frac{2l_{B,i}}{\pi_B(X_i^T \hat{\alpha})} (Y_i - m(X_i^T \hat{\beta}))^2 + \mathbb{1}_{A,i} d_{A,i} \hat{\sigma}^2(X_i) \right\} \right],$$

where $\hat{\sigma}^2$ is arbitrary consistent estimator of $\sigma^2(X_i^T \beta^*)$. Define $\hat{V} := \hat{V}_1 + \hat{V}_2$.

Under assumption stated before, if either $\pi_B(X^T \alpha)$ or $m(X^T \beta)$ is well-specified, then

$$\hat{V} \xrightarrow{pr} V.$$

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. **Simulation**

Simulation

- ▶ In all the simulation trials, it perform the best in terms of the magnitude of bias.
- ▶ However in some cases, even though it is nearly unbiased, it show larger variability (Wider confidence band. It might be the tradeoff of the doubly robustness.

1. Introduction
2. Basic set-up
3. Methodology in high dimensional data
4. Computation
5. Asymptotic results for variable selection and estimation
6. Simulation

Thank You!