

# Data Integration Using Covariate Summaries from External Sources

Facheng Yu

Department of Statistics, University of Washington

Joint work with Yuqian Zhang

Institute of Statistics and Big Data, Renmin University of China

December 2, 2024

- A brief overview of data integration
- The use of summary statistics
- Our contribution
  - Efficient estimation using external covariate summaries
  - Efficient estimation for large external datasets
- Methodology and asymptotic theory
  - Mean estimation under data homogeneity
  - Correcting selection bias under data heterogeneity

## A brief overview of data integration

### Example 1. Semi-supervised learning

- > A labeled dataset (primary data) and an unlabeled dataset (external data).
- > Applications in linear regression, mean estimation, etc.
- > Data is homogeneous (primary and external sources share the same marginal distribution).

### Example 2. Causal inference

- > Integration of randomized controlled trials and observational studies.
- > Data could be heterogeneous (primary and external sources could preserve different marginal distributions).

**Question:** Popular integration methods usually require individualized covariates from external sources, which might be unavailable due to concerns over accessibility, privacy, storage, or cost.

## The use of summary statistics

The use of summary statistics has been explored in various causal inference and data integration studies. For instance,

- > Mendelian randomization in genome-wide association studies (GWAS).
- > Data fusion in meta-analysis.
- > Require access to outcome-related information from external sources, such as  $\hat{\beta}_{\text{least-squared}}$ .

**Question:** How to conduct efficient estimation in scenarios where

- > only summary statistics of the external covariates are available,
- > and the outcome of external sources has not yet been generated?

## Our contribution

Throughout the paper, we focus on mean estimation problems. Let  $(\Gamma, X, Y)$  be an independent copy of the complete data, where

- >  $\Gamma \in \{0, 1\}$ : the group indicator, 1 for primary source and 0 for external source,
- >  $X \in \mathbb{R}^d$ : the covariate vector,
- >  $Y \in \mathbb{R}$ : the outcome.

The goal is to conduct the point estimates for

$$\theta_g = \mathbb{E}[Y] \text{ (Generalizability)} \quad \text{and} \quad \theta_t = \mathbb{E}[Y \mid \Gamma = 0] \text{ (Transportability)}.$$

### Note.

- > For the entire dataset with sample size  $n$ , the primary data can be written as  $\mathcal{D}_n := (\Gamma_i, \Gamma_i X_i, \Gamma_i Y_i)_{i=1}^n$ , where  $(\Gamma_i, X_i, Y_i)_{i=1}^n$  contains i.i.d. samples.
- > Since the outcome of external sources has not yet been generated, estimation on  $\theta_t$  is meaningful.

## Our contribution

With the sample mean  $\bar{X}_0$  and sample gram matrix  $\bar{\Xi}_0$  from the external source:

$$\bar{X}_0 = \frac{\sum_{i=1}^n (1 - \Gamma_i) X_i}{\sum_{i=1}^n (1 - \Gamma_i)}, \quad \bar{\Xi}_0 = \frac{\sum_{i=1}^n (1 - \Gamma_i) X_i X_i^\top}{\sum_{i=1}^n (1 - \Gamma_i)},$$

our results demonstrate that

- > Knowing  $\bar{X}_0$  is sufficient to construct consistent (doubly robust) estimators for  $\theta_g = \mathbb{E}[Y]$  and  $\theta_t = \mathbb{E}[Y \mid \Gamma = 0]$  under homogeneity and heterogeneity.
- > If  $\bar{\Xi}_0$  is accessible additionally, asymptotic inference can be conducted.
- > Our estimations remain efficient for particularly large external data, i.e.,  $\gamma_n = \mathbb{P}(\Gamma = 1) \rightarrow 0$  as  $n \rightarrow \infty$ .

**Note.** Comment on  $\gamma_n = \mathbb{P}(\Gamma = 1) \rightarrow 0$ :

- > Our measure  $\mathbb{P}$  is defined on the space of  $(\Gamma_i, X_i, Y_i)_{i=1}^n$ .
- > As  $n \rightarrow \infty$ , the measure could change while  $(\Gamma_i, X_i, Y_i)_{i=1}^n$  remains independent.

## Inspiration from linear model under MCAR

We start with the scenario of data homogeneity, which leads to the missing completely at random (MCAR) assumption:

Assumption 1 (MCAR)

$$\Gamma \perp (X, Y).$$

Define the population slope  $\beta^* = \arg \min_{\beta \in \mathbb{R}^d} \mathbb{E}[(Y - X^\top \beta)^2]$ .

Observe that under MCAR,

- > By KKT conditions,  $\theta = \mathbb{E}[Y] = \mathbb{E}[X]^\top \beta^*$ .
- > Replace  $\mathbb{E}[X]$  by the sample mean  $\bar{X}_{\text{all}}$  such that

$$\bar{X}_{\text{all}} = n^{-1} \sum_{i=1}^n X_i = (1 - n^{-1} \sum_{i=1}^n \Gamma_i) \bar{X}_0 + n^{-1} \sum_{i=1}^n \Gamma_i X_i.$$

**Note.** When  $\beta^*$  is estimated via least squares  $\hat{\beta}_{\text{least-squares}}$ , this approach aligns with the semi-supervised least squares (SSLS) estimator  $\hat{\theta}_{\text{SSLS}}$  introduced by Zhang et al. (2019).

- > In low dimensions,  $\hat{\theta}_{\text{SSLS}}$  is CAN and at least as efficient as  $\bar{Y}_1$ ,
- > require the covariate dimension  $d$  satisfying  $d = o(\sqrt{n_P})$ , where  $n_P = \sum_{i=1}^n \Gamma_i$  is the primary sample size.

# Inspiration from linear model and AIPW under MCAR

Let

$$m(X) = \mathbb{E}[Y | X] \text{ and } \gamma_n(X) = \mathbb{P}(\Gamma = 1 | X).$$

Let  $m^*(\cdot)$  and  $\gamma_n^*(\cdot)$  represent the working models for the outcome regression and propensity score.

- > Under MCAR,  $\gamma_n(X) = \gamma_n = \mathbb{P}(\Gamma = 1)$  and  $\theta_g = \theta_t =: \theta$ .
- > Use the naive estimator  $\hat{\gamma}_n = n^{-1} \sum_{i=1}^n \Gamma_i$ , then  $\gamma_n^* = \gamma_n$ .
- > Thus, with  $m^*(X) = X^\top \beta^*$ , AIPW always gives the consistent estimator for  $\theta$ :

$$\begin{aligned} \theta &= \mathbb{E} \left[ X^\top \beta^* + \frac{\Gamma}{\gamma_n^*} (Y - X^\top \beta^*) \right] \\ &\approx \bar{X}_{\text{all}}^\top \beta^* + \mathbb{E}_n \left[ \frac{\Gamma}{\gamma_n^*} (Y - X^\top \beta^*) \right]. \end{aligned}$$

**Note.**

- > To apply to high-dimensional setting, we can estimate  $\beta^*$  by Lasso.
- > To relax the sparsity assumption, we can use cross-fitting.



## Procedure under MCAR

Step 1: Divide the index set  $[n]$  into  $K$  disjoint subsets  $\mathcal{I}_1, \dots, \mathcal{I}_K$  with equal sizes such that  $n_k := |\mathcal{I}_k| = n/K$  for  $k \in [K]$ . Let  $\hat{\gamma}_k = n_k^{-1} \sum_{i \in \mathcal{I}_k} \Gamma_i$  and  $\mathcal{I}_{-k} = [n] \setminus \mathcal{I}_k$ .

Step 2: For each  $k \in [K]$ , compute the Lasso estimator  $\hat{\beta}^{(-k)}$ : with some  $\lambda_n \geq 0$ ,

$$\hat{\beta}^{(-k)} = \arg \min_{\beta \in \mathbb{R}^d} \left\{ \frac{\sum_{i \in \mathcal{I}_{-k}} \Gamma_i (Y_i - X_i^\top \beta)^2}{\sum_{i \in \mathcal{I}_{-k}} \Gamma_i} + \lambda_n \|\beta\|_1 \right\}.$$

Step 3: The mean estimator is proposed as:

$$\hat{\theta} = n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} (\Gamma_i X_i + (1 - \Gamma_i) \bar{X}_0)^\top \hat{\beta}^{(-k)} + n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} \frac{\Gamma_i}{\hat{\gamma}_k} (Y_i - X_i^\top \hat{\beta}^{(-k)}). \quad (0.1)$$

When the sample gram matrix  $\Xi_0$  is also observable, the corresponding asymptotic variance estimator is defined as:

$$\begin{aligned} \hat{\sigma}^2 = & n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} (1 - \Gamma_i) \hat{\beta}^{(-k), \top} \Xi_0 \hat{\beta}^{(-k)} \\ & + n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} \left\{ \Gamma_i X_i^\top \hat{\beta}^{(-k)} + \frac{\Gamma_i}{\hat{\gamma}_k} (Y_i - X_i^\top \hat{\beta}^{(-k)}) \right\}^2 - \hat{\theta}^2. \end{aligned} \quad (0.2)$$

# Asymptotic theory

## Assumption 2

Let the following conditions hold with constants  $\kappa_l, \sigma, \sigma_w, \delta_w > 0$ : (a)  $X$  is a sub-Gaussian random vector with  $\|X^\top v\|_{\psi_2} \leq \sigma \|v\|_2, \forall v \in \mathbb{R}^d$ . In addition,  $\|X^\top \beta^*\|_{\psi_2} \leq \sigma$  and  $\inf_{v \in \mathbb{R}^d, \|v\|_2=1} \mathbb{E}[(X^\top v)^2] \geq \kappa_l$ . (b)  $w = Y - X^\top \beta^*$  is a sub-Gaussian random variable with  $\|w\|_{\psi_2} \leq \sigma_w$  and  $\mathbb{E}[w^2] \geq \delta_w$ .

The following theorem characterizes asymptotic properties of the mean estimator  $\hat{\theta}$ .

## Theorem 3

Suppose that Assumptions 1 and 2 hold. Choose  $\lambda_n \asymp \sqrt{\log d / (n\gamma_n)}$ . If  $n\gamma_n \gg (\log n)^2 \log d$  and the sparsity level satisfies  $s = \|\beta^*\|_0 = o(n\gamma_n / \log d)$ , then as  $n, d \rightarrow \infty$ ,  $\hat{\theta} - \theta = O_p((n\gamma_n)^{-1/2})$ ,  $\hat{\sigma}^2 = \sigma_n^2 \{1 + o_p(1)\}$ , and  $\hat{\sigma}^{-1} \sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} \mathcal{N}(0, 1)$ , where  $\sigma_n^2 = \text{Var}\left(X^\top \beta^* + \frac{\gamma}{\gamma_n} (Y - X^\top \beta^*)\right) = \gamma_n^{-1} \text{Var}(Y) + (1 - \gamma_n^{-1}) \text{Var}(X^\top \beta^*)$ .

In the following, we consider the possibility of data heterogeneity between primary and external data and assume the following missing at random (MAR) condition instead

Assumption 4 (MAR)

$\Gamma \perp Y \mid X$  and  $\gamma_n(X) := \mathbb{P}(\Gamma = 1 \mid X) > 0$  *almost surely*.

We are interested in

$$\theta_g = \mathbb{E}[Y] \text{ (generalizability)} \quad \text{and} \quad \theta_t = \mathbb{E}[Y \mid \Gamma = 0] \text{ (transportability)}.$$

For AIPW, as long as either  $m^*(\cdot) = m(\cdot)$  or  $\gamma_n^*(\cdot) = \gamma_n(\cdot)$ ,

$$\begin{aligned}\theta_g &= \mathbb{E} \left[ m^*(X) + \frac{\Gamma}{\gamma_n^*(X)} (Y - m^*(X)) \right], \\ \theta_t &= \mathbb{E} \left[ \frac{1 - \Gamma}{1 - \gamma_n} m^*(X) + \frac{\Gamma(1 - \gamma_n^*(X))}{\gamma_n^*(X)(1 - \gamma_n)} (Y - m^*(X)) \right].\end{aligned}$$

Consider a linear outcome model  $m^*(X) = X^\top \beta_{OR}^*$ . The above representations can also be expressed as:

$$\begin{aligned}\theta_g &= (1 - \gamma_n) \mathbb{E}[X \mid \Gamma = 0]^\top \beta_{OR}^* + \mathbb{E} \left[ \Gamma X^\top \beta_{OR}^* + \frac{\Gamma}{\gamma_n^*(X)} (Y - X^\top \beta_{OR}^*) \right], \\ \theta_t &= \mathbb{E}[X \mid \Gamma = 0]^\top \beta_{OR}^* + \mathbb{E} \left[ \frac{\Gamma(1 - \gamma_n^*(X))}{\gamma_n^*(X)(1 - \gamma_n)} (Y - X^\top \beta_{OR}^*) \right].\end{aligned}$$

**Note.** With linear model, we can still construct consistent estimator for  $\theta_g$  and  $\theta_t$ .

## Estimation for nuisances under MAR

Consider a cross-fitted estimation. For each fold  $\mathcal{I}_k$ , and separate  $\mathcal{I}_{-k} = [n] \setminus \mathcal{I}_k$  into two disjoint subsets  $\mathcal{I}_{-k,\alpha}, \mathcal{I}_{-k,\beta}$  with same sizes  $M$ .

> Construct a logistic propensity estimator for  $\gamma_n^*(X) = \text{expit}(X^\top \alpha_{PS}^*)$  as

$$\hat{\alpha}_{PS}^{(-k)} = \arg \min_{\alpha \in \mathbb{R}^d} \left\{ M^{-1} \sum_{i \in \mathcal{I}_{-k,\alpha}} \left\{ (1 - \Gamma_i) \bar{X}_0^\top \alpha + \Gamma_i \exp(-X_i^\top \alpha) \right\} + \lambda_\alpha \|\alpha\|_1 \right\}.$$

> Construct a linear outcome estimator

$$\hat{\beta}_{OR}^{(-k)} = \arg \min_{\beta \in \mathbb{R}^d} \left\{ M^{-1} \sum_{i \in \mathcal{I}_{-k,\beta}} \Gamma_i \exp(-X_i^\top \hat{\alpha}_{PS}^{(-k)}) (Y_i - X_i^\top \beta)^2 + \lambda_\beta \|\beta\|_1 \right\}.$$

To leveraging conditional independence, we also define a oracle nuisance estimator assuming  $\alpha_{PS}^*$  is known.

$$\tilde{\beta}_{OR}^{(-k)} = \arg \min_{\beta \in \mathbb{R}^d} \left\{ M^{-1} \sum_{i \in \mathcal{I}_{-k,\beta}} \Gamma_i \exp(-X_i^\top \alpha_{PS}^*) (Y_i - X_i^\top \beta)^2 + \lambda_\beta \|\beta\|_1 \right\}.$$

## Asymptotic theory

When  $\bar{X}_0$  is observed,  $\theta_g = \mathbb{E}[Y]$  can be estimated by

$$\hat{\theta}_g = n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} \left\{ (\Gamma_i X_i + (1 - \Gamma_i) \bar{X}_0)^\top \hat{\beta}_{OR}^{(-k)} + \frac{\Gamma_i}{g(X_i^\top \hat{\alpha}_{PS}^{(-k)})} (Y_i - X_i^\top \hat{\beta}_{OR}^{(-k)}) \right\}.$$

When  $\Xi_0$  is further observed, we can construct the asymptotic variance estimators as:

$$\begin{aligned} \hat{\sigma}_g^2 = & n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} (1 - \Gamma_i) \hat{\beta}_{OR}^{(-k), \top} \Xi_0 \hat{\beta}_{OR}^{(-k)} - \hat{\theta}_g^2 \\ & + n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} \left\{ \Gamma_i X_i^\top \hat{\beta}_{OR}^{(-k)} + \frac{\Gamma_i}{g(X_i^\top \hat{\alpha}_{PS}^{(-k)})} (Y_i - X_i^\top \hat{\beta}_{OR}^{(-k)}) \right\}^2 \end{aligned}$$

# Assumptions

Let  $g(x) = \text{expit}(x)$ . For propensity score:

- >  $k_0(1 - \gamma_n)/\gamma_n \leq (1 - g(X^\top \alpha_{PS}^*))/g(X^\top \alpha_{PS}^*) \leq k_0^{-1}(1 - \gamma_n)/\gamma_n$  almost surely.
- >  $\mathbb{P}(\Gamma = 0) \geq c_0$ .
- >  $\mathbb{E}[\gamma_n^q(X)] \leq \nu \gamma_n^q$  for some  $q > 1$  and  $\nu > 0$ .

In addition,

- > For each  $j \in \{0, 1\}$ , and conditional on  $\Gamma = j$ ,  $X$  is a sub-Gaussian random vector and  $X^\top \beta_{OR}^*$  be a sub-Gaussian random variable, both with parameter  $\sigma$ .
- >  $\inf_{v \in \mathbb{R}^d, \|v\|_2=1} \mathbb{E} \left[ (X^\top v)^2 \mid \Gamma = 1 \right] \geq \kappa_l$ .
- > The residual  $w_{OR} = Y - X^\top \beta_{OR}^*$  is a sub-Gaussian random variable with parameter  $\sigma_w$ ,  $\mathbb{E} [w_{OR}^8 \mid \Gamma = 1] \leq \sigma_w^8$ , and  $\mathbb{E} [w_{OR}^2 \mid \Gamma = 1] \geq \delta_w$ .

## Theorem 5

Under MAR and regular assumptions. Choose  $\lambda_\alpha \asymp \lambda_\beta \asymp \sqrt{\log d / (n\gamma_n)}$ . Let  $n\gamma_n \gg (\log n)^2 \log d$ ,  $\|\alpha_{PS}^*\|_0 \|\beta_{OR}^*\|_0 = o(n\gamma_n / (\log n (\log d)^2))$ , and either of the following conditions hold: (1) (Correct OR model)  $\mu(X) = X^\top \beta_{OR}^*$  or (2) (Correct PS model)  $\gamma_n(X) = g(X^\top \alpha_{PS}^*)$  and  $\|\alpha_{PS}^*\|_0 = o(\sqrt{n\gamma_n} / \log d)$ . Then as  $n, d \rightarrow \infty$ ,

(a)  $\hat{\theta}_g - \theta_g = O_p\left((n\gamma_n)^{-1/2}\right)$ ,  $\hat{\sigma}_g^2 = \sigma_g^2 \{1 + o_p(1)\}$ , and  $\hat{\sigma}^{-1} \sqrt{n}(\hat{\theta}_g - \theta_g) \xrightarrow{d} \mathcal{N}(0, 1)$ ,

where  $\sigma_g^2 = \text{Var}\left(X^\top \beta_{OR}^* + \frac{\Gamma}{g(X^\top \alpha_{PS}^*)} (Y - X^\top \beta_{OR}^*)\right)$ .

(b)  $\hat{\theta}_t - \theta_t = O_p\left((n\gamma_n)^{-1/2}\right)$ ,  $\hat{\sigma}_t^2 = \sigma_t^2 \{1 + o_p(1)\}$ , and  $\hat{\sigma}_t^{-1} \sqrt{n}(\hat{\theta}_t - \theta_t) \xrightarrow{d} \mathcal{N}(0, 1)$ ,

where  $\sigma_t^2 = \mathbb{E}\left[\left\{\frac{1-\Gamma}{1-\gamma_n} (X^\top \beta_{OR}^* - \theta_t) + \frac{\Gamma(1-g(X^\top \alpha_{PS}^*))}{(1-\gamma_n)g(X^\top \alpha_{PS}^*)} (Y - X^\top \beta_{OR}^*)\right\}^2\right]$ .



We also study some extensions of our method in

- > Estimate ATE in causal inference
- > Estimate ATE on the external population in causal inference.

Zhang, A., Brown, L. D., and Cai, T. T. (2019). Semi-supervised inference: General theory and estimation of means. The Annals of Statistics, 47(5):2538–2566.