

# METHOD-OF-MOMENTS INFERENCE FOR GLMs

Xingyu Chen, Lin Liu, Rajarshi Mukherjee



Arxiv Link



Github Link

## Introduction

### Problem Setting and Questions:

- Samples:  $(A_i \in \mathbb{R}, \mathbf{X}_i \in \mathbb{R}^p)_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$
- The distribution  $\mathbb{P}$  is parameterized as:  
 $\mathbf{X} \sim \mathbb{P}_{\mathbf{X}}, \quad \mathbb{E}(A|\mathbf{X}) = \phi(\boldsymbol{\alpha}^\top \mathbf{X}), \quad \text{var}(A|\mathbf{X}) = \sigma^2(\mathbf{X})$   
 where  $\mathbf{X}$  has mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$
- Asymptotic regime:  $\frac{p}{n} \rightarrow \delta \in [0, +\infty)$  as  $n \rightarrow \infty$

### Questions

How can we conduct inference on  $\boldsymbol{\alpha}$ ?  
 How can we conduct inference on  $\|\boldsymbol{\alpha}\|_{\boldsymbol{\Sigma}}^2 := \boldsymbol{\alpha}^\top \boldsymbol{\Sigma} \boldsymbol{\alpha}$ ?  
 How can we conduct inference on other functionals?

We focus on four settings, ranging from simple to complex, to introduce our methodology:

- **Case I (Gaussian, known  $\boldsymbol{\mu} = 0$  & known  $\boldsymbol{\Sigma}$ ):**  
 $\mathbf{X} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma}), \boldsymbol{\mu} = 0, \boldsymbol{\Sigma}$  is known
- **Case II (Gaussian, unknown  $\boldsymbol{\mu}$  & known  $\boldsymbol{\Sigma}$ ):**  
 $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \boldsymbol{\mu}$  is unknown,  $\boldsymbol{\Sigma}$  is known
- **Case III (Gaussian, unknown  $\boldsymbol{\mu}$  & unknown  $\boldsymbol{\Sigma}$ ):**  
 $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$  both  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  are unknown
- **Case IV (Missing Data):**  
 $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \boldsymbol{\mu}$  is unknown,  $\boldsymbol{\Sigma}$  is known.  
 In addition, we observe  $(Y_i)_{i=1}^n$  with  $\mathbb{E}[Y_i|\mathbf{X}_i] = \beta^\top \mathbf{X}_i$   
 and assume  $A \perp Y | \mathbf{X}$ .

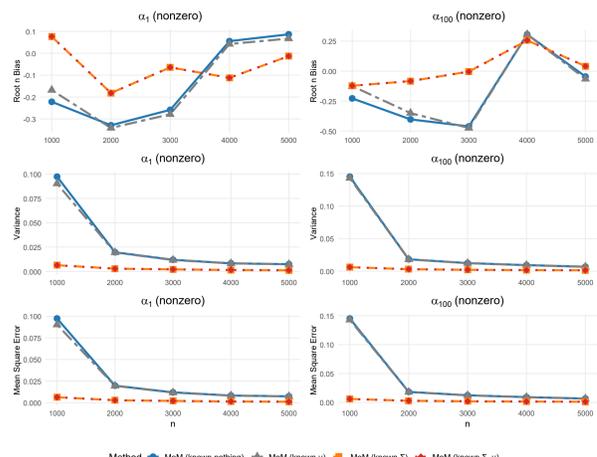


Figure 1: Root- $n$  bias, variance, and mean squared error of estimators for  $\alpha_1$  and  $\alpha_{100}$  under different assumptions on  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , with  $p/n = 0.4$ .

## Estimation in Case I

**Lemma 1 (Stein's Lemma)** Let  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $f: \mathbb{R}^p \rightarrow \mathbb{R}$  be a differentiable function such that all expectations below are finite. Then:

$$\mathbb{E}[\mathbf{X}f(\mathbf{X})] = \boldsymbol{\Sigma} \mathbb{E}[\nabla f(\mathbf{X})] + \mathbb{E}[f(\mathbf{X})]\boldsymbol{\mu} \quad (1)$$

For  $f(\mathbf{X}) = \phi(\boldsymbol{\alpha}^\top \mathbf{X})$ , we have:

$$\mathbb{E}[\mathbf{X}A] = \mathbb{E}[\mathbf{X}\phi(\boldsymbol{\alpha}^\top \mathbf{X})] = \mathbb{E}[\phi'(\mathbf{Z})\boldsymbol{\Sigma}\boldsymbol{\alpha} + \mathbb{E}[f(\mathbf{X})]\boldsymbol{\mu}] \quad (2)$$

Here and below,  $\mathbf{Z} \sim \boldsymbol{\alpha}^\top \mathbf{X} \sim \mathcal{N}(\boldsymbol{\alpha}^\top \boldsymbol{\mu}, \|\boldsymbol{\alpha}\|_{\boldsymbol{\Sigma}}^2)$ , with  $\lambda_{\boldsymbol{\alpha}} := \boldsymbol{\alpha}^\top \boldsymbol{\mu}$  and  $\gamma_{\boldsymbol{\alpha}}^2 := \|\boldsymbol{\alpha}\|_{\boldsymbol{\Sigma}}^2$ . We define  $f_i(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) := \mathbb{E}[\phi^{(i)}(\mathbf{Z})]$ .

Since  $\boldsymbol{\mu} = \mathbf{0}$  in Case I:

### Identification Equations for $\boldsymbol{\mu} = 0$

$$\begin{aligned} m_{\alpha_j} &:= \mathbb{E}[A\mathbf{X}_j^\top] \boldsymbol{\Sigma}^{-1} \mathbf{e}_j = f_1(0, \gamma_{\boldsymbol{\alpha}}^2) \alpha_j, \\ m_{\mathbf{X}_{A,2}} &:= \mathbb{E}[A_1 \mathbf{X}_1^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}_2 A_2] = f_1(0, \gamma_{\boldsymbol{\alpha}}^2)^2 \cdot \gamma_{\boldsymbol{\alpha}}^2 =: \Psi(\gamma_{\boldsymbol{\alpha}}^2) \end{aligned} \quad (3)$$

Since  $\boldsymbol{\Sigma}$  is **known** in Case I, we construct  $U$ -statistics to unbiasedly estimate the required moments:

### Moment Estimators for $\boldsymbol{\mu} = 0$

$$\begin{aligned} \widehat{m}_{\alpha_j} &:= \frac{1}{n} \sum_{1 \leq i \leq n} A_i \mathbf{X}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_j \\ &=: \mathbb{U}_{n,1}[A_1 \mathbf{X}_1^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_j], \\ \widehat{m}_{\mathbf{X}_{A,2}} &:= \frac{1}{n(n-1)} \sum_{1 \leq i_1 \neq i_2 \leq n} A_{i_1} \mathbf{X}_{i_1}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}_{i_2} A_{i_2} \\ &=: \mathbb{U}_{n,2}[A_1 \mathbf{X}_1^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}_2 A_2] \end{aligned} \quad (4)$$

We then plug the moment estimators from (4) into the identification equations (3) and solve the system to obtain the estimator for the parameter of interest.

$$\begin{aligned} \widehat{\gamma}_{\boldsymbol{\alpha}}^2 &= \Psi^{-1}(\widehat{m}_{\mathbf{X}_{A,2}}) \\ \widehat{\boldsymbol{\alpha}}_j &= \frac{\widehat{m}_{\alpha_j}}{f_1(0, \widehat{\gamma}_{\boldsymbol{\alpha}}^2)} \end{aligned} \quad (5)$$

## Estimation in Case II

Since  $\boldsymbol{\mu}$  is unknown in Case II, the identification equations are as follows:

### Identification Equations for Unknown $\boldsymbol{\mu}$

$$\begin{aligned} m_A &:= \mathbb{E}[A] = f_0(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2), \\ m_{\mathbf{X},2} &:= \mathbb{E}[\mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}, \\ m_{\mathbf{X}_A, \mathbf{X}} &:= \mathbb{E}[A\mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \mathbb{E}[\mathbf{X}] \\ &= m_A \cdot m_{\mathbf{X},2} + f_1(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \lambda_{\boldsymbol{\alpha}}, \\ m_{\mathbf{X}_{A,2}} &:= \mathbb{E}[A\mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \mathbb{E}[A\mathbf{X}] \\ &= m_A^2 \cdot m_{\mathbf{X},2} + f_1^2(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \gamma_{\boldsymbol{\alpha}}^2 + \\ &\quad 2 \cdot m_A \cdot f_1(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \lambda_{\boldsymbol{\alpha}}, \\ m_{\nu_j} &:= \mathbb{E}[\mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \mathbf{e}_j = \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_j = \boldsymbol{\nu}^\top \mathbf{e}_j = \nu_j, \\ m_{\alpha_j} &:= \mathbb{E}[A\mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \mathbf{e}_j \\ &= f_0(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \nu_j + f_1(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \boldsymbol{\alpha}_j. \end{aligned} \quad (6)$$

The first four moments in (6) can be reduced to two equations, which form a diffeomorphism map  $\Psi_{GLM}$ :

### Reduced Identification Equations for Unknown $\boldsymbol{\mu}$

$$\begin{aligned} m_1 &:= m_A = f_0(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2), \\ m_2 &:= m_{\mathbf{X}_{A,2}} + m_A^2 \cdot m_{\mathbf{X},2} - 2 \cdot m_A \cdot m_{\mathbf{X}_A, \mathbf{X}} \\ &= f_1^2(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \gamma_{\boldsymbol{\alpha}}^2, \\ \Psi_{GLM} &: (\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \rightarrow (m_1, m_2). \end{aligned} \quad (7)$$

Since  $\boldsymbol{\Sigma}$  is **known** in Case II:

### Moment Estimators for Unknown $\boldsymbol{\mu}$

$$\begin{aligned} \widehat{m}_1 &:= \widehat{m}_A = \mathbb{U}_{n,1}[A], \\ \widehat{m}_2 &:= \widehat{m}_{\mathbf{X}_{A,2}} + \widehat{m}_A^2 \cdot \widehat{m}_{\mathbf{X},2} - 2 \cdot \widehat{m}_A \cdot \widehat{m}_{\mathbf{X}_A, \mathbf{X}}, \\ \text{where } \widehat{m}_{\mathbf{X},2} &:= \mathbb{U}_{n,2}[\mathbf{X}_1^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}_2], \\ \widehat{m}_{\mathbf{X}_A, \mathbf{X}} &:= \mathbb{U}_{n,2}[A_1 \mathbf{X}_1^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}_2], \\ \widehat{m}_{\mathbf{X}_{A,2}} &:= \mathbb{U}_{n,2}[A_1 \mathbf{X}_1^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}_2 A_2], \\ \widehat{m}_{\nu_j} &:= \mathbb{U}_{n,1}[\mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \mathbf{e}_j, \\ \widehat{m}_{\alpha_j} &:= \mathbb{U}_{n,1}[A\mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \mathbf{e}_j. \end{aligned} \quad (8)$$

We then plug the moment estimators from (8) into the reduced identification equations (7):

$$\begin{aligned} (\widehat{\lambda}_{\boldsymbol{\alpha}}, \widehat{\gamma}_{\boldsymbol{\alpha}}^2) &:= \Psi_{GLM}^{-1}(\widehat{m}_1, \widehat{m}_2), \\ \widehat{\boldsymbol{\alpha}}_j &:= \frac{\widehat{m}_{\alpha_j} - f_0(\widehat{\lambda}_{\boldsymbol{\alpha}}, \widehat{\gamma}_{\boldsymbol{\alpha}}^2) \cdot \widehat{m}_{\nu_j}}{f_1(\widehat{\lambda}_{\boldsymbol{\alpha}}, \widehat{\gamma}_{\boldsymbol{\alpha}}^2)}. \end{aligned} \quad (9)$$

## Estimation in Case III

In Case III, the identification equations remain the same as in Case II, i.e.,

### Identification Equations for Unknown $\boldsymbol{\mu}$

Same as (7).

The knowledge of  $\boldsymbol{\Sigma}$  influences the construction of moment estimators. We propose two methods to address this, each with theoretical guarantees.

One method involves using a sample splitting strategy with weighted sample covariance:

### Moment Estimators Unknown $\boldsymbol{\Sigma}$ (Sample Splitting)

$$\begin{aligned} \widetilde{\boldsymbol{\Sigma}} &:= \frac{1}{\frac{n}{2} - p - 1} \sum_{j \in I_2} (\mathbf{X}_j - \bar{\mathbf{X}}_{I_2})(\mathbf{X}_j - \bar{\mathbf{X}}_{I_2})^\top, \\ \text{where } \bar{\mathbf{X}}_{I_2} &:= \frac{1}{n/2} \sum_{j \in I_2} \mathbf{X}_j, \\ \widehat{m}_{\mathbf{X}_{A,2}} &:= \frac{1}{\frac{n}{2}(\frac{n}{2} - 1)} \sum_{i_1 \neq i_2 \in I_1} A_{i_1} \mathbf{X}_{i_1}^\top \widetilde{\boldsymbol{\Sigma}}^{-1} \mathbf{X}_{i_2} A_{i_2}. \end{aligned} \quad (10)$$

Another method uses Chebyshev polynomials to approximate  $\boldsymbol{\Sigma}^{-1}$ :

### Moment Estimators Unknown $\boldsymbol{\Sigma}$ (Chebyshev)

$$\begin{aligned} \boldsymbol{\Sigma}^{-1} &\approx \sum_{l=0}^J c_l \boldsymbol{\Sigma}^l, \\ \widehat{m}_{\mathbf{X}_{A,2}} &:= \sum_{l=0}^J c_l \mathbb{U}_{n,l+2} \left[ A_1 \mathbf{X}_1^\top \left( \prod_{s=3}^{l+2} \mathbf{X}_s \mathbf{X}_s^\top \right) \mathbf{X}_2 A_2 \right]. \end{aligned} \quad (11)$$

## Estimation in Case IV

Now, we apply the methods discussed above to the Missing Data setting, also considered in [1].

Here, the observed data are  $(A_i, X_i, A_i \cdot Y_i)_{i=1}^n$ , and the parameter of interest is  $\boldsymbol{\psi} := \mathbb{E}[Y] = \mathbb{E}[\beta^\top \mathbf{X}] = \beta^\top \boldsymbol{\mu}$ .

### Identification Equations

Equations in (7), plus two additional equations:

$$\begin{aligned} m_{AY} &:= \mathbb{E}[AY] = m_A \cdot \boldsymbol{\psi} + f_1(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \gamma_{\alpha, \beta}, \\ m_{\mathbf{X}_{AY, \mathbf{X}}} &:= \mathbb{E}[Y A \mathbf{X}^\top] \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ &= (m_A + m_{\mathbf{X}_A, \mathbf{X}}) \cdot \boldsymbol{\psi} + \\ &\quad \{m_{\mathbf{X},2} \cdot f_1(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) + f_2(\lambda_{\boldsymbol{\alpha}}, \gamma_{\boldsymbol{\alpha}}^2) \cdot \lambda_{\boldsymbol{\alpha}}\} \cdot \gamma_{\alpha, \beta}, \\ \text{where } \gamma_{\alpha, \beta} &:= \boldsymbol{\alpha}^\top \boldsymbol{\Sigma} \boldsymbol{\beta}. \end{aligned} \quad (12)$$

As in [1], we use the knowledge of  $\boldsymbol{\Sigma}$ , so the additional moment estimators are the same as those in (8), and the final estimator follows similarly. Comparison below:

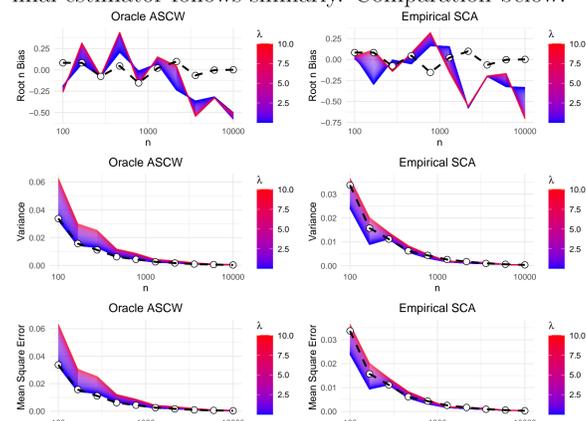


Figure 2: Root- $n$  Bias, Variance, and Mean Squared Error of Estimators for  $\boldsymbol{\psi}$ , comparing our method with two Ridge regression-based methods from [1].

### Theorem (Informal), Gaussian, CAN

When  $\boldsymbol{\Sigma}$  is known, under some mild conditions, the above estimators are all  $\sqrt{n}$ -consistent.

Further assume that  $\sqrt{p}\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\mu}$ ,  $\sqrt{p}\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\alpha}$ , and  $\sqrt{p}\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\beta}$  and their inner products with respect to  $\boldsymbol{\Sigma}^{-1}$ ,  $\boldsymbol{\Sigma}$ ,  $\boldsymbol{\Sigma}^2$ ,  $\boldsymbol{\Sigma}^3$  converge to a nontrivial distribution. Then, the above estimators, after scaling by  $\sqrt{n}$ , converge to a normal distribution.

When  $\boldsymbol{\Sigma}$  is unknown, for the method in (10), which requires  $\frac{n}{2} > p + 3$ , we can show that our method is  $\sqrt{n}$ -consistent. For the method in (11), we can show that our method is consistent.

### Theorem (Informal), Universality

When  $\mathbf{X}$  violates the Gaussian distribution, but  $\boldsymbol{\Sigma}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$  has zero mean and unit variance, the above identification equations (3), (6) hold with error  $\mathcal{O}(n^{-3/4})$ .

Thus, the above consistent results for the Gaussian model will still hold.

### Variance Estimator

We use a bootstrap method to estimate the variance of the  $U$ -statistics and then apply the Delta method to estimate the variance of the parameters of interest.

## References

- [1] Michael Celentano and Martin J Wainwright. Challenges of the inconsistency regime: Novel debiasing methods for missing data models. *arXiv preprint arXiv:2309.01362*, 2023.