

Section: Module 3

OLS

Kentaro Nakamura

GOV 2003

February 13th, 2026

Today's agenda

- Today, we will cover
 - Review of Linear Algebra
 - Matrix derivative
 - OLS
 - Difference between Estimand and Estimator
 - FWL Theorem
 - Gauss Markov Theorem

- Logistics
 - Problem set 1 is due on next Wednesday

- Recommended Textbook:
 - Hansen (2022) Econometrics, Ch2-7, Ch9

Matrix Derivative / Gradients (1)

Definition (The derivative of matrix by a scalar)

Let $A = (a_{ij})_{ij}$ be an $m \times n$ matrix and let each a_{ij} be a function of x . Then, we define the derivative of A with respect to the variable x as follows:

$$\frac{\partial A(x)}{\partial x} = \left(\frac{\partial a_{ij}(x)}{\partial x} \right)_{ij}$$

Example (The derivative of matrix by a scalar)

Let $A(x) = \begin{bmatrix} x^2 + 2x & x^3 - 1 \\ x^3 - x^2 + 1 & x^2 + x \end{bmatrix}$. Then,

$$\frac{\partial A(x)}{\partial x} = \begin{bmatrix} 2x + 2 & 3x^2 \\ 3x^2 - 2x & 2x + 1 \end{bmatrix}$$

Matrix Derivative / Gradients (2)

Definition (The derivative of function by a matrix)

Let $f(X)$ be a function of an $m \times n$ matrix X ($f : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$). Then, the derivative of $f(X)$ with respect to the matrix X is

$$\frac{\partial f(X)}{\partial X} = \left(\frac{\partial f(X)}{\partial x_{ij}} \right)_{ij}$$

Example (The derivative of function by a matrix)

Let $X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix}$ and $f(X) = x_{11}^2 x_{23} + 2x_{12}x_{13} - 5x_{21}x_{22}^2$.

Then,

$$\frac{\partial f(X)}{\partial X} = \begin{bmatrix} 2x_{11}x_{23} & 2x_{13} & 2x_{12} \\ -5x_{22}^2 & -10x_{21}x_{22} & x_{11}^2 \end{bmatrix}$$

Matrix Derivative / Gradients (3)

Definition (The derivative of function by a vector)

Let $f(x)$ be a function of vector $x \in \mathbb{R}^m$ ($f : \mathbb{R}^m \mapsto \mathbb{R}$). Then, we define the first and second derivatives of $f(x)$ as follows:

$$\frac{\partial f(x)}{\partial x} = \left[\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_m} \right]^\top$$
$$\frac{\partial f(x)}{\partial x \partial x^\top} = \left(\frac{\partial^2 f(x)}{\partial x_i \partial x_j} \right)_{ij}$$

Example (The derivative of function by a vector)

Let $x = [x_1, x_2]^\top$ and consider $f(x) = x_1 x_2^3 + x_1^2 x_2^4$. Then,

$$\frac{\partial f(x)}{\partial x} = \begin{bmatrix} x_2^3 + 2x_1 x_2^4 \\ 3x_1 x_2^2 + 4x_1^2 x_2^3 \end{bmatrix}$$
$$\frac{\partial f(x)}{\partial x \partial x^\top} = \begin{bmatrix} 2x_2^4 & 3x_2^2 + 8x_1 x_2^3 \\ 3x_2^2 + 8x_1 x_2^3 & 6x_1 x_2 + 12x_1^2 x_2^2 \end{bmatrix}$$

Some useful formulas

- **Rule 1:** If a is not a function of x ,

$$\frac{\partial a^\top x}{\partial x} = a^\top$$

- **Rule 2:** If A and b are not functions of x ,

$$\frac{\partial b^\top Ax}{\partial x} = A^\top b$$

- **Rule 3:** If A is not a function of x ,

$$\frac{\partial x^\top Ax}{\partial x} = (A + A^\top)^\top x$$

- If A is symmetric, $A + A^\top = 2A$

OLS at the population: Estimand

- Consider $\arg \min_{\beta} \mathbb{E}[(Y_i - \mathbf{X}_i^{\top} \beta)^2]$
- Then, the objective function is written as

$$\mathbb{E}[(Y_i - \mathbf{X}_i^{\top} \beta)^2] = \mathbb{E}[Y_i^2] - 2 \cdot \beta^{\top} \mathbb{E}[\mathbf{X}_i Y_i] + \beta^{\top} \mathbb{E}[\mathbf{X}_i \mathbf{X}_i^{\top}] \beta$$

- By taking the derivative with respect to β , we get

$$\begin{aligned} -2 \cdot \mathbb{E}[\mathbf{X}_i Y_i] + 2 \mathbb{E}[\mathbf{X}_i \mathbf{X}_i^{\top}] \beta &= 0 \\ \Rightarrow \beta &= (\mathbb{E}[\mathbf{X}_i \mathbf{X}_i^{\top}])^{-1} \mathbb{E}[\mathbf{X}_i Y_i] \end{aligned}$$

- BUT we cannot calculate this since we cannot calculate expectation.
 - The quantity of interest, which we cannot calculate from data, is called **estimand**

OLS at the sample: Estimator

- Recall that OLS solves the following optimization problem:

$$\hat{\beta} = \arg \min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2$$

- Notice that

$$\begin{aligned}\|\mathbf{Y} - \mathbf{X}\beta\|_2^2 &= (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) \\ &= \mathbf{Y}^\top \mathbf{Y} - 2\mathbf{Y}^\top \mathbf{X}\beta + \beta^\top \mathbf{X}^\top \mathbf{X}\beta\end{aligned}$$

and thus

$$\nabla_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 = -2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X}\beta$$

- Therefore, $-\mathbf{X}^\top \mathbf{Y} + \mathbf{X}^\top \mathbf{X}\beta = 0$, which gives us

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{Y} \right)$$

- This is calculatable from data \Rightarrow called **estimator**

FWL theorem (1)

- Consider the partition $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$ and $\beta = (\beta_1, \beta_2)$, where

$$\mathbf{Y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \epsilon$$

- Consider the OLS estimator of $\beta = (\beta_1, \beta_2)$, which is defined by

$$(\hat{\beta}_1, \hat{\beta}_2) = \arg \min_{\beta_1, \beta_2} \|\mathbf{Y} - \mathbf{X}_1\beta_1 - \mathbf{X}_2\beta_2\|_2^2$$

Now, $\hat{\beta}_1$ is obtained by nested minimization

$$\hat{\beta}_1 = \arg \min_{\beta_1} \left(\min_{\beta_2} \|\mathbf{Y} - \mathbf{X}_1\beta_1 - \mathbf{X}_2\beta_2\|_2^2 \right)$$

- Now, let's focus on the inner minimization problem. This is simply the OLS of $\mathbf{Y} - \mathbf{X}_1\beta_1$ on \mathbf{X}_2 , which has the solution

$$\hat{\beta}_2 = (\mathbf{X}_2^\top \mathbf{X}_2)^{-1} (\mathbf{X}_2^\top (\mathbf{Y} - \mathbf{X}_1\beta_1))$$

- The residual $\mathbf{Y} - \mathbf{X}_1\beta_1 - \mathbf{X}_2\hat{\beta}_2$ is written as

$$\underbrace{(\mathbf{I}_n - \mathbf{X}_2(\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top)}_{M_2} (\mathbf{Y} - \mathbf{X}_1\beta_1)$$

FWL theorem (2)

- Thus,

$$\begin{aligned} & \arg \min_{\beta_1} \left(\min_{\beta_2} \| \mathbf{Y} - \mathbf{X}_1 \beta_1 - \mathbf{X}_2 \beta_2 \|_2^2 \right) \\ &= \arg \min_{\beta_1} \left((\mathbf{M}_2(\mathbf{Y} - \mathbf{X}_1 \beta_1))^\top (\mathbf{M}_2(\mathbf{Y} - \mathbf{X}_1 \beta_1)) \right) \\ &= \arg \min_{\beta_1} (\mathbf{Y} - \mathbf{X}_1 \beta_1)^\top \mathbf{M}_2 (\mathbf{Y} - \mathbf{X}_1 \beta_1) \end{aligned}$$

- Notice that this is the form of weighted least squares (WLS), which solution is written as

$$\hat{\beta}_1 = (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{Y})$$

- This is equivalently expressed as

$$\begin{aligned} \hat{\beta}_1 &= (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{Y}) \\ &= (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{M}_2 \mathbf{Y}) \quad (\because \mathbf{M}_2^2 = \mathbf{M}_2) \\ &= (\tilde{\mathbf{X}}_1^\top \tilde{\mathbf{X}}_1)^{-1} (\tilde{\mathbf{X}}_1^\top \tilde{\mathbf{e}}) \end{aligned}$$

where $\tilde{\mathbf{X}}_1 = \mathbf{M}_2 \mathbf{X}_1$ and $\tilde{\mathbf{e}} = \mathbf{M}_2 \mathbf{Y}$.

Unbiasedness of OLS Estimator

Theorem

Unbiasedness of OLS Estimator Assume the linearity: $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ with exogeneity $\mathbb{E}[\epsilon | \mathbf{X}] = 0$. Then,

$$\mathbb{E}[\hat{\beta}] = \beta$$

- **Proof**

$$\begin{aligned}\mathbb{E}[\hat{\beta}] &= \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}] \\ &= \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \epsilon)] \\ &= \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon] \\ &= \beta + \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon] \\ &= \beta + \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}[\epsilon | \mathbf{X}]] = \beta\end{aligned}$$

Gauss-Markov Theorem

Theorem (Gauss-Markov Theorem)

Suppose that (i) $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ (linearity), (ii) $\mathbb{E}[\epsilon | \mathbf{X}] = 0$ (exogeneity), and (iii) $\mathbb{V}[\epsilon | \mathbf{X}] = \sigma^2 I_n$ (homoskedasticity). Then, OLS estimator $\hat{\beta}$ is Best Linear Unbiased Estimator in the sense that

$$\mathbb{V}[\tilde{\beta}] \geq \mathbb{V}[\hat{\beta}]$$

for any $\tilde{\beta}$ satisfying (i) $\tilde{\beta} = \mathbf{A}(\mathbf{X})\mathbf{Y}$ (linear) and (ii) $\mathbb{E}[\tilde{\beta} | \mathbf{X}] = \beta$ (unbiasedness).

- But homoskedasticity is rarely true!

Gauss-Markov Theorem: Proof (1)

- Now, $\tilde{\beta}$ satisfies $\tilde{\beta} = \mathbf{A}(\mathbf{X})\mathbf{Y}$. So,

$$\begin{aligned}\tilde{\beta} &= \mathbf{A}(\mathbf{X})\mathbf{Y} = \mathbf{A}(\mathbf{X})\mathbf{Y} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} + \underbrace{(\mathbf{A}(\mathbf{X}) - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)}_{:= \mathbf{B}(\mathbf{X})} \mathbf{Y} \\ &= \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} \mathbf{Y} \\ &= \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} (\mathbf{X}\beta + \epsilon) \\ &= \beta + \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} \epsilon + \mathbf{B}(\mathbf{X})\mathbf{X}\beta\end{aligned}$$

- Now, we assume unbiasedness, which means

$$\mathbb{E}[\tilde{\beta} \mid \mathbf{X}] = \mathbb{E}[\hat{\beta} \mid \mathbf{X}] = \beta. \text{ Thus, for any } \beta,$$

$$\begin{aligned}0 &= \mathbb{E}[\tilde{\beta} \mid \mathbf{X}] - \beta = \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} \mathbb{E}[\epsilon \mid \mathbf{X}] + \mathbf{B}(\mathbf{X})\mathbf{X}\beta \\ &= \mathbf{B}(\mathbf{X})\mathbf{X}\beta\end{aligned}$$

- As a result, we have

$$\tilde{\beta} = \beta + \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} \epsilon$$

Gauss-Markov Theorem: Proof (2)

- Then, recall that for a vector $\mathbf{X} \in \mathbb{R}^k$, the covariance matrix $\mathbb{V}[\mathbf{X}]$ is defined as

$$\mathbb{V}[\mathbf{X}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top]$$

where the (i, j) -th entry is covariance $\text{Cov}(X_i, X_j)$.

- Therefore,

$$\begin{aligned}\mathbb{V}[\tilde{\beta} \mid \mathbf{X}] &= \mathbb{V}[\beta + \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} \epsilon \mid \mathbf{X}] \\ &= \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} \mathbb{V}[\epsilon \mid \mathbf{X}] \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\}^\top \\ &= \sigma^2 \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\} \{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{B}(\mathbf{X})\}^\top \\ &= \sigma^2 \{(\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}(\mathbf{X}) \mathbf{B}(\mathbf{X})^\top\}^\top \\ &\geq \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = \mathbb{V}[\hat{\beta} \mid \mathbf{X}]\end{aligned}$$

- As we consider the variance (which is the diagonal element of covariance matrix), this shows the claim of Gauss-Markov Theorem.
 - This is because the diagonal element of semi-positive definite matrix is non-negative (see Problem Set 0)

Appendix: Partitioned Matrix

- Sometimes, partitioned matrix can help us formulate the proof / computation
 - Let $\mathbf{A}_{11} \in \mathbb{R}^{m \times m}$, $\mathbf{A}_{12} \in \mathbb{R}^{m \times n}$, $\mathbf{A}_{21} \in \mathbb{R}^{n \times m}$, and $\mathbf{A}_{22} \in \mathbb{R}^{n \times n}$.
Then, you can write matrix as $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$
 - You can calculate the addition and multiplication with this partitioned matrix in the usual way

Theorem (Inverse of Partitioned Matrix)

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}_{11}^{-1} + \mathbf{B}_{12} \mathbf{B}_{22}^{-1} \mathbf{B}_{21} & -\mathbf{B}_{12} \mathbf{B}_{22}^{-1} \\ -\mathbf{B}_{22}^{-1} \mathbf{B}_{21} & \mathbf{B}_{22}^{-1} \end{bmatrix}$$

where $\mathbf{B}_{12} := \mathbf{A}_{11}^{-1} \mathbf{A}_{12}$, $\mathbf{B}_{21} := \mathbf{A}_{21} \mathbf{A}_{11}^{-1}$, and $\mathbf{B}_{22} := \mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12}$

- We can later use this to prove partitioned regression in a different way from the last time

Appendix: Partitioned Regression / FWL Theorem

- In this proof, I derive FWL theorem using partitioned matrix.

Theorem (Partitioned Regression)

Consider the partition $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$ and $\beta = (\beta_1, \beta_2)$, and the regression model

$$\mathbf{Y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \epsilon.$$

The least squared estimator $(\hat{\beta}_1, \hat{\beta}_2)$ is written as

$$\hat{\beta}_1 = (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{Y})$$

$$\hat{\beta}_2 = (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{Y})$$

where

$$\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$$

$$\mathbf{M}_2 = \mathbf{I}_n - \mathbf{X}_2(\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top$$

Appendix: Proof: Partitioned Regression / FWL Theorem (1)

- Now, recall that $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$.
- Using partitioned matrix notation,

$$\mathbf{X}^\top \mathbf{Y} = \begin{bmatrix} \mathbf{X}_1^\top \mathbf{Y} \\ \mathbf{X}_2^\top \mathbf{Y} \end{bmatrix}$$

- On the other hand,

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} \mathbf{X}_1^\top \\ \mathbf{X}_2^\top \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1^\top \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{X}_2 \\ \mathbf{X}_2^\top \mathbf{X}_1 & \mathbf{X}_2^\top \mathbf{X}_2 \end{bmatrix}$$

Appendix: Proof: Partitioned Regression / FWL Theorem (2)

- Now, applying the result of inverse for the partitioned matrix, we get

$$\left(\begin{bmatrix} \mathbf{X}_1^\top \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{X}_2 \\ \mathbf{X}_2^\top \mathbf{X}_1 & \mathbf{X}_2^\top \mathbf{X}_2 \end{bmatrix} \right)^{-1} = \begin{bmatrix} \mathbf{X}_{11}^{-1} + \mathbf{Z}_{12} \mathbf{Z}_{22}^{-1} \mathbf{Z}_{21} & -\mathbf{Z}_{12} \mathbf{Z}_{22}^{-1} \\ -\mathbf{Z}_{22}^{-1} \mathbf{Z}_{21} & \mathbf{Z}_{22}^{-1} \end{bmatrix}$$

where $\mathbf{X}_{ij} = \mathbf{X}_i^\top \mathbf{X}_j$, and

$$\mathbf{Z}_{12} = \mathbf{X}_{11}^{-1} \mathbf{X}_{12}$$

$$\mathbf{Z}_{21} = \mathbf{X}_{21} \mathbf{X}_{11}^{-1}$$

$$\mathbf{Z}_{22} = \mathbf{X}_{22} - \mathbf{X}_{21} \mathbf{X}_{11}^{-1} \mathbf{X}_{12}$$

Appendix: Proof: Partitioned Regression / FWL Theorem (3)

- As a result,

$$\begin{aligned} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{X}_{11}^{-1} + \mathbf{Z}_{12}\mathbf{Z}_{22}^{-1}\mathbf{Z}_{21} & -\mathbf{Z}_{12}\mathbf{Z}_{22}^{-1} \\ -\mathbf{Z}_{22}^{-1}\mathbf{Z}_{21} & \mathbf{Z}_{22}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{X}_1^\top \mathbf{Y} \\ \mathbf{X}_2^\top \mathbf{Y} \end{bmatrix} \\ &= \begin{bmatrix} (\mathbf{X}_{11}^{-1} + \mathbf{Z}_{12}\mathbf{Z}_{22}^{-1}\mathbf{Z}_{21}) \mathbf{X}_1^\top \mathbf{Y} - \mathbf{Z}_{12}\mathbf{Z}_{22}^{-1}\mathbf{X}_2^\top \mathbf{Y} \\ -\mathbf{Z}_{22}^{-1}\mathbf{Z}_{21}\mathbf{X}_1^\top \mathbf{Y} + \mathbf{Z}_{22}^{-1}\mathbf{X}_2^\top \mathbf{Y} \end{bmatrix}. \end{aligned}$$

- Therefore,

$$\hat{\beta}_2 = \mathbf{Z}_{22}^{-1}(\mathbf{X}_2^\top \mathbf{Y} - \mathbf{Z}_{21}\mathbf{X}_1^\top \mathbf{Y}) = \mathbf{Z}_{22}^{-1}(\mathbf{X}_2^\top \mathbf{Y} - \mathbf{X}_{21}\mathbf{X}_{11}^{-1}\mathbf{X}_1^\top \mathbf{Y})$$

Appendix: Proof: Partitioned Regression / FWL Theorem (4)

- Recall

$$\begin{aligned}\mathbf{Z}_{22} &= \mathbf{X}_{22} - \mathbf{X}_{21}\mathbf{X}_{11}^{-1}\mathbf{X}_{12} = \mathbf{X}_2^\top \mathbf{X}_2 - \mathbf{X}_2^\top \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \\ &= \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2.\end{aligned}$$

- Also,

$$\mathbf{X}_2^\top \mathbf{Y} - \mathbf{X}_{21}\mathbf{X}_{11}^{-1}\mathbf{X}_1^\top \mathbf{Y} = \mathbf{X}_2^\top \left(\mathbf{I}_n - \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \right) \mathbf{Y} = \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{Y}.$$

- Hence,

$$\hat{\beta}_2 = (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{Y}).$$

- By symmetry (swap the roles of \mathbf{X}_1 and \mathbf{X}_2),

$$\hat{\beta}_1 = (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{Y}).$$