

# Review Section (Midterm)

## OLS / MLE

Kentaro Nakamura

GOV 2003

March 23th, 2026

# Today's Agenda: Review

- OLS
  - Projection Matrix
  - Annihilator Matrix
  - FWL Theorem
  - Variance of OLS
  - Asymptotics of OLS
  
- MLE
  - Review of Theories
  - Logistic regression
  - Bradley-Terry models

# OLS and Orthogonal Projection

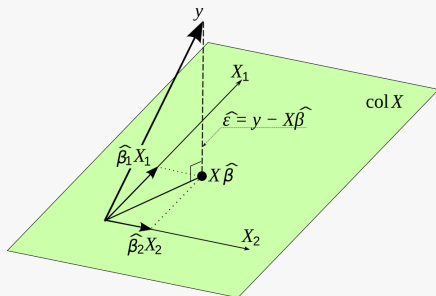
- Consider the linear model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ . Recall that OLS estimator of  $\boldsymbol{\beta}$  is given by

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

- As  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ ,

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \underbrace{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{:=\mathbf{P}_X} \mathbf{Y}$$

- Now, notice that  $\mathbf{P}_X^2 = \mathbf{P}_X$  and  $\mathbf{P}^\top = \mathbf{P}$  (Prove this!), and thus  $\mathbf{P}_X$  is an orthogonal projection matrix



# Annihilator Matrix

## Definition (Annihilator Matrix)

$M = I_n - X(X^T X)^{-1} X^T$  is called **annihilator matrix**.

- Notice that

$$MX = (I_n - X(X^T X)^{-1} X^T)X = X - X(X^T X)^{-1} X^T X = X - X = 0$$

- Annihilator matrix gives you residual since

$$MY = (I_n - X(X^T X)^{-1} X^T)Y = Y - X(X^T X)^{-1} X^T Y = Y - \hat{Y}$$

- Annihilator matrix is a projection onto orthogonal complement of the column space of  $X$
- We can also show that annihilator matrix is symmetric ( $M^T = M$ ) and idempotent  $M^2 = M$ 
  - That is, annihilator matrix is a projection onto orthogonal complement of the column space

## Residual and Annihilator Matrix

- Residual is a projection of  $\mathbf{Y}$  onto  $S^\perp(\mathbf{X})$  (i.e., orthogonal complement of the column space)

- Orthogonality 1:

$$\hat{\mathbf{Y}}^\top \hat{\boldsymbol{\epsilon}} = (\mathbf{P}_X \mathbf{Y})^\top \mathbf{M} \mathbf{Y} = \mathbf{Y}^\top \mathbf{P}_X^\top (\mathbf{I} - \mathbf{P}_X) \mathbf{Y} = 0$$

- Orthogonality 2:

$$\mathbf{X}_k^\top \hat{\boldsymbol{\epsilon}} = (\mathbf{P}_X \mathbf{X}_k)^\top \mathbf{M} \mathbf{Y} = \mathbf{X}_k^\top \mathbf{P}_X^\top (\mathbf{I} - \mathbf{P}_X) \mathbf{Y} = 0$$

- $\mathbf{P}_X \mathbf{X}_k = \mathbf{X}_k$  because

$$\mathbf{P}_X \mathbf{X}_k = \mathbf{P}_X \mathbf{X} \mathbf{e}_k = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{e}_k = \mathbf{X} \mathbf{e}_k = \mathbf{X}_k$$

- Orthogonality 3:  $\mathbf{x}^\top \hat{\boldsymbol{\epsilon}} = 0$  for any  $\mathbf{x} \in S(\mathbf{X})$

- Proof: Recall the definition of column space:  $\mathbf{x} = \mathbf{X} \mathbf{a}$  for some  $\mathbf{a}$ .

$$\mathbf{x}^\top \hat{\boldsymbol{\epsilon}} = (\mathbf{X} \mathbf{a})^\top \hat{\boldsymbol{\epsilon}} = \mathbf{a}^\top \mathbf{X}^\top \mathbf{M} \mathbf{Y} = \mathbf{a}^\top \mathbf{X}^\top (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = 0$$

- Zero mean:  $\sum_{i=1}^n \hat{\epsilon}_i = 0$ :

- As  $\mathbf{X}$  includes intercept, there exists  $\mathbf{a}$  such that  $\mathbf{1} = \mathbf{X} \mathbf{a}$ . As a result, by the previous orthogonality 3, we have zero mean.

# Unbiasedness of OLS

## 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**

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

## Variance of OLS

- Now,

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

and thus

$$\mathbb{V}[\hat{\beta} \mid \mathbf{X}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{V}[\epsilon \mid \mathbf{X}] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

- Under homoskedasticity,  $\mathbb{V}[\epsilon \mid \mathbf{X}] = \sigma^2 I$ .

- Thus,

$$\mathbb{V}[\hat{\beta} \mid \mathbf{X}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

## Variance estimator

- We do not know the true error variance, so we need to propose the estimator  $\hat{\sigma}^2$
- The natural estimator is to replace expectation with average; i.e.,

$$\hat{\sigma}^2 = \frac{1}{n} \hat{\epsilon}^\top \hat{\epsilon} = \frac{1}{n} (\mathbf{M}\epsilon)^\top \mathbf{M}\epsilon = \frac{1}{n} \epsilon^\top \mathbf{M}\epsilon$$

- This is actually biased.
  - To see that, use  $\text{tr}(\epsilon^\top \mathbf{M}\epsilon) = \text{tr}(\mathbf{M}\epsilon\epsilon^\top)$

$$\mathbb{E}[\hat{\sigma}^2 \mid \mathbf{X}] = \frac{1}{n} \text{tr}(\mathbb{E}[\mathbf{M}\epsilon\epsilon^\top \mid \mathbf{X}]) = \frac{1}{n} \text{tr}(\mathbf{M} \underbrace{\mathbb{E}[\epsilon\epsilon^\top \mid \mathbf{X}]}_{\sigma^2 I_n}) = \frac{n-p}{n} \sigma^2$$

where the last line is by  $\text{tr}(\mathbf{M}) = n - p$

- Thus, the unbiased estimator of variance is given by

$$s^2 = \frac{1}{n-p} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{1}{n-p} \epsilon^\top \epsilon$$

## 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

# 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

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

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

# Tools in Asymptotic Statistics

- **Law of Large Numbers (LLN):** If  $X_1, \dots, X_n$  are i.i.d.,

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mathbb{E}[X_i]$$

- **Central Limit Theorem (CLT):** If  $X_1, \dots, X_n$  are i.i.d.,

$$\sqrt{n}(\bar{X} - \mathbb{E}[X_i]) \xrightarrow{d} \mathcal{N}(0, \mathbb{V}[X_i])$$

- **Slutsky's Lemma:** If  $X_n \xrightarrow{d} X$  for some random variable  $X$  and  $Y_n \xrightarrow{P} c$  for some constant  $c$ ,

$$X_n Y_n \xrightarrow{d} cX$$

## Asymptotics of OLS: Consistency

- Now, notice that

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

- By law of large numbers,

$$\frac{1}{n} \mathbf{X}^T \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \xrightarrow{p} \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^T]$$

- Similarly,

$$\frac{1}{n} \mathbf{X}^T \epsilon = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \xrightarrow{p} \mathbb{E}[\mathbf{x}_i^T \epsilon_i] = 0$$

- Thus, by Slutsky,

$$\hat{\beta} - \beta = \left( \frac{1}{n} \mathbf{X}^T \mathbf{X} \right)^{-1} \frac{1}{n} \mathbf{X}^T \epsilon \xrightarrow{p} \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^T]^{-1} \cdot 0 = 0$$

which means that  $\hat{\beta}$  is consistent.

- If we use CLT for the second term, then we can derive the asymptotic normality.

## MLE: Review of Basic Principle

- **Setting:** Observe data  $\mathbf{Y} = (Y_1, \dots, Y_n)$  from a distribution with density (or pmf)  $f_\theta(\mathbf{y})$ , where  $\theta \in \Theta$  is an unknown parameter.
- **Likelihood function:** Given observed data  $\mathbf{y}$ ,

$$L(\theta) = L(\theta | \mathbf{y}) = f_\theta(\mathbf{y}) = \prod_{i=1}^n f_\theta(y_i) \quad (\text{under i.i.d. assumption})$$

- **Likelihood principle:** The likelihood function contains all the information about the unknown parameter of interest.
- **MLE:** The maximum likelihood estimator is

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} L(\theta)$$

# Score Function

## Definition (Score Function)

The score function  $S_n(\mathbf{Y}, \theta)$  (or simply  $S_n(\theta)$ ) for the data  $\mathbf{Y} = (Y_1, \dots, Y_n)$  is defined as

$$S_n(\mathbf{Y}, \theta) = \frac{\partial \log f_\theta(\mathbf{Y})}{\partial \theta} = \frac{\partial \ell(\theta)}{\partial \theta}$$

where  $\ell(\theta)$  is the log-likelihood function.

- **Property:** The expected score is zero. That is,  $\mathbb{E}_\theta[S_n(\mathbf{Y}, \theta)] = 0$ 
  - $\mathbb{E}_\theta[S_n(\mathbf{Y}, \theta)] = \int \frac{f'_\theta(\mathbf{y})}{f_\theta(\mathbf{y})} f_\theta(\mathbf{y}) d\mathbf{y} = \frac{\partial}{\partial \theta} \int f_\theta(\mathbf{y}) d\mathbf{y} = \frac{\partial}{\partial \theta} 1 = 0$

# Fisher Information

## Definition (Fisher Information)

The Fisher information for data  $\mathbf{Y} = (Y_1, \dots, Y_n)$  is defined as

$$I(\theta) := \mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \log f_{\theta}(\mathbf{Y}) \right)^2 \right] = -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta^2} \log f_{\theta}(\mathbf{Y}) \right]$$

- Fisher information gives us the asymptotic variance under the correctly specified models

## Theorem (Fisher's Theorem / MLE Asymptotic Normality)

Let  $Y_1, \dots, Y_n \sim f_{\theta_0}$ . Under the regularity conditions, the MLE estimator  $\hat{\theta}_n$  achieves the asymptotic normality

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} \mathcal{N}(0, I(\theta_0)^{-1})$$

where  $I(\theta_0)$  is the fisher information.

## Example: Logistic Regression (1)

- The likelihood function is

$$L(\beta) = \prod_{i=1}^n \pi(\mathbf{x}_i)^{Y_i} (1 - \pi(\mathbf{x}_i))^{1 - Y_i}$$

and thus the logistic regression is

$$\ell(\beta) = \sum_{i=1}^n \left( Y_i \log \pi(\mathbf{x}_i) + (1 - Y_i) \log(1 - \pi(\mathbf{x}_i)) \right)$$

- The score function is

$$\frac{\partial}{\partial \beta} \ell(\beta) = \sum_{i=1}^n \mathbf{x}_i (Y_i - \pi(\mathbf{x}_i))$$

## Example: Logistic Regression (1)

- The Hessian is

$$\frac{\partial^2}{\partial \boldsymbol{\beta}^2} \ell(\boldsymbol{\beta}) = - \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \pi(\mathbf{x}_i) (1 - \pi(\mathbf{x}_i))$$

- By information equality,

$$I_n(\theta) = -\mathbb{E} \left[ \frac{\partial^2}{\partial \boldsymbol{\beta}^2} \ell(\boldsymbol{\beta}) \right] = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \pi(\mathbf{x}_i) (1 - \pi(\mathbf{x}_i))$$

- Note that this assumes fixed design (i.e.,  $\mathbf{x}_i$  is non-random.)

## Example 2: Bradley Terry Model (1)

- The model is

$$\Pr(j \text{ beats } j') = \frac{e^{\beta_j}}{e^{\beta_j} + e^{\beta_{j'}}},$$

- The likelihood is thus given by

$$L(\beta) = \prod_{1 \leq j < j' \leq J} \left( \frac{e^{\beta_j}}{e^{\beta_j} + e^{\beta_{j'}}} \right)^{W_{jj'}} \left( \frac{e^{\beta_{j'}}}{e^{\beta_j} + e^{\beta_{j'}}} \right)^{W_{j'j}}$$

- And thus the log-likelihood is

$$\ell(\beta) = \sum_{1 \leq j < j' \leq J} \left[ W_{jj'} \beta_j + W_{j'j} \beta_{j'} - (W_{jj'} + W_{j'j}) \log(\theta_j + \theta_{j'}) \right].$$

where  $\theta_j := e^{\beta_j}$

## Example 2: Bradley-Terry Model (2)

- The score is

$$\frac{\partial \ell}{\partial \beta_j} = \sum_{j' \neq j} \left[ W_{jj'} - (W_{jj'} + W_{j'j}) \frac{\theta_j}{\theta_j + \theta_{j'}} \right].$$

- The Hessian is, for  $j \neq k$ ,

$$\frac{\partial^2 \ell}{\partial \beta_j \partial \beta_k} = (W_{jk} + W_{kj}) \frac{\theta_j \theta_k}{(\theta_j + \theta_k)^2},$$

and for  $j = k$ ,

$$\frac{\partial^2 \ell}{\partial \beta_j^2} = - \sum_{j' \neq j} (W_{jj'} + W_{j'j}) \frac{\theta_j \theta_{j'}}{(\theta_j + \theta_{j'})^2}.$$

## Example 2: Bradley-Terry Model (3)

- By information equality,

$$I_n(\beta) = -\mathbb{E}\left[\frac{\partial^2}{\partial\beta\partial\beta^\top}\ell(\beta)\right].$$

- Thus, for  $j \neq k$ ,

$$I_n(\beta)_{jk} = -(W_{jk} + W_{kj})\frac{\theta_j\theta_k}{(\theta_j + \theta_k)^2},$$

and

$$I_n(\beta)_{jj} = \sum_{j' \neq j} (W_{jj'} + W_{j'j})\frac{\theta_j\theta_{j'}}{(\theta_j + \theta_{j'})^2}.$$