

ECON207 Session 5

OLS using Matrix Algebra

Anthony Tay

This Version: 31 Jul 2024

Session 5

- A Bit More Matrix Algebra
- OLS estimation of the Simple Linear Regression Model
 - Geometric Intuition: Projections
 - Statistic Properties
 - Completion of Basic Theory
 - Some Applications

Session 5.1

Session 5.1 A Bit More Matrix Algebra

- Differentiation of Matrix Forms
- Application to Optimization

Differentiation of Matrix Forms

Let $f(x)$ be some function of the $n \times 1$ vector x

e.g., $f(x) = x^T Ax$ where A is some $n \times n$ matrix of constants

This is just a multivariable function $f(x_1, \dots, x_n)$

We *define*

$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \quad \text{or} \quad \frac{df}{dx^T} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix}$$

Differentiation of Matrix Forms

Example: If $y = x^T A x$ where A is $n \times n$ and x is $n \times 1$, then

$$\frac{dy}{dx} = \frac{d}{dx} (x^T A x) = (A + A^T)x$$

$n = 3$ example:

$$x^T A x = x_1^2 a_{11} + x_2^2 a_{22} + x_3^2 a_{33} + x_1 x_2 (a_{12} + a_{21}) + x_1 x_3 (a_{13} + a_{31}) + x_2 x_3 (a_{23} + a_{32})$$

$$\left. \begin{array}{l} \frac{dx^T A x}{dx_1} = 2x_1 a_{11} + x_2(a_{12} + a_{21}) + x_3(a_{13} + a_{31}) \\ \frac{dx^T A x}{dx_2} = x_1(a_{12} + a_{21}) + 2x_2 a_{22} + x_3(a_{23} + a_{32}) \\ \frac{dx^T A x}{dx_3} = x_1(a_{13} + a_{31}) + x_2(a_{23} + a_{32}) + 2x_3 a_{33} \end{array} \right\} \Rightarrow \frac{dx^T A x}{dx} = (A + A^T)x$$

If A is symmetric, then $\frac{dx^T A x}{dx} = (A + A^T)x = 2Ax$ (compare $f(x) = ax^2 \Rightarrow f'(x) = 2ax$)

Differentiation of Matrix Forms

Suppose x is $n \times 1$ and $y = f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_m(x) \end{bmatrix}$ (i.e., y is a vector of m different functions)

Then we define

$$\frac{dy}{dx^T} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad \text{or} \quad \frac{dy^T}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_2}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_1} \\ \frac{\partial f_1}{\partial x_2} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_n} & \frac{\partial f_2}{\partial x_n} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Differentiation of Matrix Forms

E.g.: if $y = f(x) = Ax$ where $A = (a_{ij})_{m \times n}$ and $x^T = [x_1 \ x_2 \ \dots \ x_n]$, then

$$\frac{dy}{dx^T} = \frac{d(Ax)}{dx^T} = A \quad \text{or} \quad \frac{dy^T}{dx} = \frac{d(x^T A^T)}{dx} = A^T$$

Matrix analogue of the univariate differentiation rule $f(x) = ax \Rightarrow f'(x) = a$

Proof:

- The product Ax is an $m \times 1$ vector whose i -th element is $\sum_{k=1}^n a_{ik}x_k$
- Therefore the (i, j) th element of dy/dx^T is $(\partial/\partial x_j) \sum_{k=1}^n a_{ik}x_k = a_{ij}$
- This says that $dy/dx^T = A$

Differentiation of Matrix Forms

E.g.: If $y = f(x)$ is a scalar-valued function of an $n \times 1$ vector of variables x , then

$$\frac{d}{dx^T} \left(\frac{dy}{dx} \right) = \frac{d}{dx^T} \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial^2 y}{\partial x_1^2} & \frac{\partial^2 y}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 y}{\partial x_1 \partial x_n} \\ \frac{\partial^2 y}{\partial x_2 \partial x_1} & \frac{\partial^2 y}{\partial x_2^2} & \cdots & \frac{\partial^2 y}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 y}{\partial x_n \partial x_1} & \frac{\partial^2 y}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 y}{\partial x_n^2} \end{bmatrix}$$

- This is the Hessian matrix of $y = f(x)$. We usually write $\frac{d}{dx^T} \left(\frac{dy}{dx} \right)$ as $\frac{d^2 y}{dxdx^T}$

Matrix Differentiation (Optimization)

Recall when finding min/max of $f(x, y)$:

- Stationary point: (x^*, y^*) such that $f'_x(x^*, y^*) = 0$ and $f'_y(x^*, y^*) = 0$
- If $f(x, y)$ is concave; stationary point is maximum point
- If $v_1^2 \frac{\partial^2 f(x, y)}{\partial x^2} + 2v_1 v_2 \frac{\partial^2 f(x, y)}{\partial x \partial y} + v_2^2 \frac{\partial^2 f(x, y)}{\partial y^2} < 0$ for all v_1, v_2 not both zero, then $f(x, y)$ is concave
- If $f(x, y)$ is convex; stationary point is minimum point
- If $v_1^2 \frac{\partial^2 f(x, y)}{\partial x^2} + 2v_1 v_2 \frac{\partial^2 f(x, y)}{\partial x \partial y} + v_2^2 \frac{\partial^2 f(x, y)}{\partial y^2} > 0$ for all v_1, v_2 not both zero, then $f(x, y)$ is convex

Matrix Differentiation (Optimization)

The expression

$$v_1^2 \frac{\partial^2 f(x, y)}{\partial x^2} + 2v_1 v_2 \frac{\partial^2 f(x, y)}{\partial x \partial y} + v_2^2 \frac{\partial^2 f(x, y)}{\partial y^2}$$

is the quadratic form involving the Hessian

$$\begin{bmatrix} v_1 & v_2 \end{bmatrix} \begin{bmatrix} \frac{\partial^2 f(x, y)}{\partial x^2} & \frac{\partial^2 f(x, y)}{\partial x \partial y} \\ \frac{\partial^2 f(x, y)}{\partial x \partial y} & \frac{\partial^2 f(x, y)}{\partial y^2} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \quad \text{i.e.,} \quad v^T H_f v$$

- If Hessian is neg. def., i.e., if $v^T H_f v < 0$ for all $v \neq 0_{2 \times 1}$, then $f(x, y)$ is concave
- If Hessian is pos. def., i.e., if $v^T H_f v > 0$ for all $v \neq 0_{2 \times 1}$, then $f(x, y)$ is convex

Matrix Differentiation (Optimization)

Extends to functions of many variables $f(x)$ where x is $n \times 1$

- f concave implies stationary pt is max, f convex implies stationary pt is min
- Stationary point: x^* such that $\frac{df(x^*)}{dx} = 0$
 - If $f(x)$ is convex; stationary point is minimum point
 - If $f(x)$ is concave; stationary point is maximum point
- If Hessian is positive definite, then $f(x)$ is convex
- If Hessian is negative definite, then $f(x)$ is concave

How to check for pos./neg. def'ness? We'll consider only a simple case in the next section

Session 5.2

Session 5.2 OLS Estimation of the MLR

- Setting Up the Multiple Linear Regression Model in Matrix Form
- OLS Estimation of the MLR
 - Estimator Formulas
 - Estimator Variance Formulas
 - R^2 and Adjusted- R^2 (same as before)

Setting Up The MLR Model

Population satisfies $E(Y | X_1, X_2, \dots, X_k) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$

i.e.,

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon, E(\epsilon | X_1, \dots, X_k) = 0$$

X_1, \dots, X_k refer to k different regressors, not k obs. of a variable X

Representative iid sample from the population:

$$\{Y_i, X_{i1}, X_{i2}, \dots, X_{ik}\}_{i=1}^n$$

No perfect collinearity in the sample:

$$c_0 + c_1 X_{i1} + \dots + c_k X_{ik} = 0 \text{ for all } i = 1, \dots, n \iff c_0 = c_1 = \dots = c_k = 0$$

Setting Up The MLR Model

Sample satisfies

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_k X_{ik} + \epsilon_i ,$$

$$E(\epsilon_i | X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1k}, \dots, X_{nk}) = 0$$

$$E(\epsilon_i^2 | X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1k}, \dots, X_{nk}) = \sigma_i^2$$

$$E(\epsilon_i \epsilon_j | X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1k}, \dots, X_{nk}) = 0$$

for all $i, j = 1, \dots, n$, $i \neq j$

More compactly expressed in matrix form

Setting Up The MLR Model

The equations

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_k X_{ik} + \epsilon_i \quad \text{for all } i = 1, \dots, n$$

can be written

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1k} \\ 1 & X_{21} & X_{22} & \cdots & X_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{nk} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\text{i.e., } y = X\beta + \epsilon$$

X_{ij} is i th sample of the j th regressor

Setting Up The MLR Model

The assumption

$$c_0 + c_1 X_{i1} + \cdots + c_k X_{ik} = 0 \text{ for all } i = 1, \dots, n \iff c_0 = c_1 = \cdots = c_k = 0$$

can be written $Xc = 0_{n \times 1} \iff c = 0_{n \times 1}$ or $Xc \neq 0_{n \times 1} \iff c \neq 0_{n \times 1}$

The assumption

$$E(\epsilon_i | X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1k}, \dots, X_{nk}) = 0 \text{ for } i = 1, \dots, n$$

can be written

$$\begin{bmatrix} E(\epsilon_1 | X) \\ E(\epsilon_2 | X) \\ \vdots \\ E(\epsilon_n | X) \end{bmatrix} = 0_{n \times 1} \iff E(\epsilon | X) = 0_{n \times 1}$$

Setting Up The MLR Model

The assumptions

$$E(\epsilon_i^2 | X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1k}, \dots, X_{nk}) = \sigma_i^2 \quad \text{for } i = 1, \dots, n,$$

$$E(\epsilon_i \epsilon_j | X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1k}, \dots, X_{nk}) = 0 \quad \text{for } i, j = 1, \dots, n, i \neq j$$

can be written as

$$\begin{bmatrix} E(\epsilon_1^2 | X) & E(\epsilon_1 \epsilon_2 | X) & \cdots & E(\epsilon_1 \epsilon_n | X) \\ E(\epsilon_2 \epsilon_1 | X) & E(\epsilon_2^2 | X) & \cdots & E(\epsilon_2 \epsilon_n | X) \\ \vdots & \vdots & \ddots & \vdots \\ E(\epsilon_n \epsilon_1 | X) & E(\epsilon_n \epsilon_2 | X) & \cdots & E(\epsilon_n^2 | X) \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix}$$

or

$$Var\epsilon | X = E(\epsilon \epsilon^T | X) = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$$

Setting Up The MLR Model

If the noise terms are homoskedastic, i.e.,

$$E(\epsilon_i^2 \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1k}, \dots, X_{nk}) = \sigma^2 \quad \text{for } i = 1, \dots, n,$$

then we have

$$E(\epsilon\epsilon^T \mid X) = \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix} = \sigma^2 I_n$$

which we will assume for now

OLS Estimation of MLR

In matrix form, the MLR is

$$Y = X\beta + \epsilon, \quad E(\epsilon | X) = 0, \quad E(\epsilon\epsilon^T | X) = \sigma^2 I_n$$

For any potential estimator $\hat{\beta}$, define

- Fitted values: $\hat{y} = X\hat{\beta}$
- Residuals: $\hat{\epsilon} = y - \hat{y} = y - X\hat{\beta}$
- Sum of Squared Residuals $SSR(\hat{\beta})$: $\sum_{i=1}^n \hat{\epsilon}_i^2 = \hat{\epsilon}^T \hat{\epsilon}$

OLS Estimation of MLR

We have

$$\begin{aligned}SSR(\hat{\beta}) &= \hat{\epsilon}^T \hat{\epsilon} = (y - X\hat{\beta})^T (y - X\hat{\beta}) \\&= (y^T - \hat{\beta}^T X^T)(y - X\hat{\beta}) \\&= y^T y - \hat{\beta}^T X^T y - y^T X \hat{\beta} + \hat{\beta}^T X^T X \hat{\beta} \\&= y^T y - 2\hat{\beta}^T X^T y + \hat{\beta}^T X^T X \hat{\beta}\end{aligned}$$

OLS:

$$\begin{aligned}\hat{\beta}^{ols} &= \arg \min_{\hat{\beta}} SSR(\hat{\beta}) \\&= \arg \min_{\hat{\beta}} y^T y - 2\hat{\beta}^T X^T y + \hat{\beta}^T X^T X \hat{\beta}\end{aligned}$$

OLS Estimation of MLR

First and second derivatives are

$$\frac{dSSR(\hat{\beta})}{d\beta} = -2X^T y + 2X^T X \hat{\beta} \quad \text{and} \quad \frac{d^2SSR(\hat{\beta})}{d\beta d\beta^T} = 2X^T X$$

FOC: $\left. \frac{dSSR(\hat{\beta})}{d\beta} \right|_{\hat{\beta}^{ols}} = -2X^T y + 2X^T X \hat{\beta}^{ols} = 0_{(k+1) \times 1}$ which implies

$$\hat{\beta}^{ols} = (X^T X)^{-1} X^T y$$

- $Xc \neq 0 \Leftrightarrow c \neq 0$ means that $X^T X$ is non-singular, and
- $c^T X^T X c = (Xc)^T (Xc) > 0$ (Hessian is pos. def.) so $\hat{\beta}^{ols}$ minimizes $SSR(\hat{\beta})$

OLS Estimation of MLR

Another way of expressing $\hat{\beta}^{ols}$

Writing

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1k} \\ 1 & X_{21} & X_{22} & \cdots & X_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{nk} \end{bmatrix} = \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix}$$

which emphasizes observations.

Note X_{i*} is $1 \times (k + 1)$ vector comprising i th observation of all variables (including '1' for the intercept term)

OLS Estimation of MLR

Then

$$\hat{\beta}^{ols} = (X^T X)^{-1} X^T y$$

$$= \left\{ \begin{bmatrix} X_{1*}^T & X_{2*}^T & \dots & X_{n*}^T \end{bmatrix} \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix} \right\}^{-1} \begin{bmatrix} X_{1*}^T & X_{2*}^T & \dots & X_{n*}^T \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$$

$$= \underbrace{\left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}}_{\text{sum of } k \times k \text{ matrices}} \underbrace{\sum_{i=1}^n X_{i*}^T Y_i}_{\text{sum of } k \times 1 \text{ vectors}}$$

OLS Estimation of MLR (Quick Summary)

- MLR: $Y = X\beta + \epsilon$, $E(\epsilon | X) = 0$, $Var(\epsilon | X) = E(\epsilon\epsilon^T | X) = \sigma^2 I_n$
- OLS Estimator for β : $\hat{\beta}^{ols} = (X^T X)^{-1} X^T y$
- OLS Fitted values:

$$\hat{y}^{ols} = X\hat{\beta}^{ols} = X(X^T X)^{-1} X^T y = Py \text{ where } P = X(X^T X)^{-1} X^T$$

- OLS Residuals:

$$\begin{aligned}\hat{\epsilon}^{ols} &= y - \hat{y}^{ols} = y - X\hat{\beta}^{ols} = y - X(X^T X)^{-1} X^T y \\ &= (I_n - X(X^T X)^{-1} X^T)y = My \text{ where } M = I_n - X(X^T X)^{-1} X^T\end{aligned}$$

Where helpful to do so, I will place “ols” marker in subscript instead of superscript

OLS Estimation of MLR (More on the FOC)

The FOC can be written as: $X^T \hat{\epsilon}^{ols} = 0_{(k+1) \times 1}$

$$\text{Since } X = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1k} \\ 1 & X_{21} & X_{22} & \cdots & X_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{nk} \end{bmatrix} = [i_n \quad X_{*1} \quad X_{*2} \quad \cdots \quad X_{*k}]$$

$$\text{the FOC says: } X^T \hat{\epsilon}^{ols} = \begin{bmatrix} i_n^T \\ X_{*1}^T \\ X_{*2}^T \\ \vdots \\ X_{*k}^T \end{bmatrix} \hat{\epsilon}^{ols} = \begin{bmatrix} i_n^T \hat{\epsilon}^{ols} \\ X_{*1}^T \hat{\epsilon}^{ols} \\ X_{*2}^T \hat{\epsilon}^{ols} \\ \vdots \\ X_{*k}^T \hat{\epsilon}^{ols} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n \hat{\epsilon}_i^{ols} \\ \sum_{i=1}^n X_{i1} \hat{\epsilon}_i^{ols} \\ \sum_{i=1}^n X_{i2} \hat{\epsilon}_i^{ols} \\ \vdots \\ \sum_{i=1}^n X_{ik} \hat{\epsilon}_i^{ols} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

OLS Estimation of MLR (More on the FOC)

The first line $\sum_{i=1}^n \hat{\epsilon}_i^{ols} = 0$ means that $\overline{\hat{\epsilon}^{ols}} = 0$ so

$$\sum_{i=1}^n X_{ij} \hat{\epsilon}_i^{ols} = 0 \iff \text{smp. cov.}(X_{ij}, \hat{\epsilon}_i^{ols}) = 0$$

for each variable X_j , $j = 1, \dots, k$

OLS estimator $\hat{\beta}^{ols}$ makes OLS residuals

- mean zero and
- uncorrelated with each regressor

Goodness of Fit

$SST = SSE + SSR$ decomposition continues to hold in the general MLR case

Let $M_0 = I_n - (1/n)i_n i_n^T$ which is symmetric and idempotent

- $M_0^T = M_0$ and $M_0 M_0 = M_0$

We have

- $M_0 y = \begin{bmatrix} Y_1 - \bar{Y} \\ Y_2 - \bar{Y} \\ \vdots \\ Y_n - \bar{Y} \end{bmatrix}$ and $y^T M_0^T M_0 y = y^T M_0 y = \sum_{i=1}^n (Y_i - \bar{Y})^2$
- $M_0 \hat{\epsilon}^{ols} = \hat{\epsilon}^{ols}$

Goodness of Fit

Therefore

$$y = \hat{y}_{ols} + \hat{\epsilon}_{ols}$$

$$M_0 y = M_0 \hat{y}_{ols} + \hat{\epsilon}_{ols}$$

$$y^T M_0^T M_0 y = (M_0 \hat{y}_{ols} + \hat{\epsilon}_{ols})^T (M_0 \hat{y}_{ols} + \hat{\epsilon}_{ols})$$

$$y^T M_0 y = \hat{y}_{ols}^T M_0 \hat{y}_{ols} + \hat{\epsilon}_{ols}^T \hat{y}_{ols} + \hat{y}_{ols}^T \hat{\epsilon}_{ols} + \hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols}$$

$$y^T M_0 y = \hat{y}_{ols}^T M_0 \hat{y}_{ols} + \hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols}$$

$$SST = SSE + SSR$$

Goodness of Fit

We can use this to define the goodness-of-fit measure:

$$R^2 = 1 - \frac{SSR}{SST}$$

and “Adjusted R^2 ”

$$\text{Adj.-}R^2 = 1 - \frac{\frac{1}{n-k-1}SSR}{\frac{1}{n-1}SST} = 1 - \frac{SSR}{SST} \frac{n-1}{n-k-1}$$

where k is no. of slope coefficients

Session 5.3

Session 5.3 Geometric Perspectives on OLS Estimation of MLR

Before discussing properties, we take a digression and consider OLS estimation of the MLR from a geometric perspective

- A vector space view of OLS estimation
 - Builds intuition
 - Help you understand some of the language used in OLS discussions

Geometric Perspective

Consider one regressor and $n = 3$, i.e.,

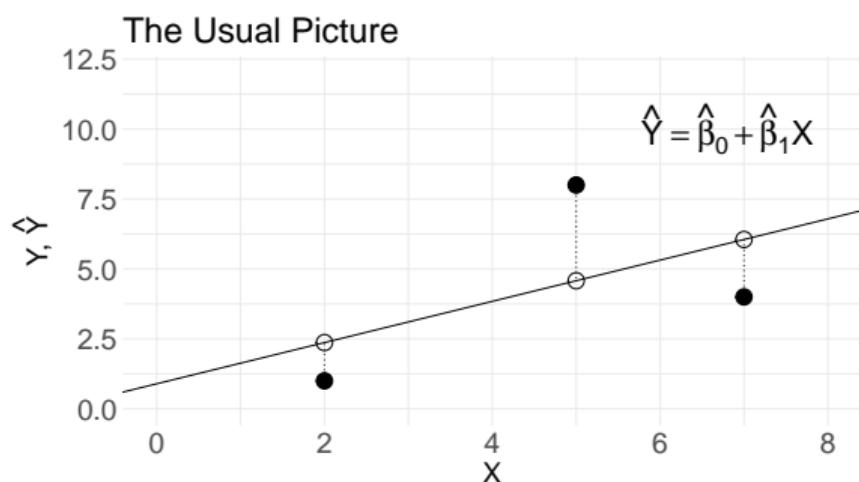
$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} = \begin{bmatrix} 1 & X_{11} \\ 1 & X_{21} \\ 1 & X_{31} \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} + \begin{bmatrix} \hat{\epsilon}_1 \\ \hat{\epsilon}_2 \\ \hat{\epsilon}_3 \end{bmatrix}$$

where

$$(X_{i1}, Y_i) = (2, 1), (5, 8), (7, 4)$$

OLS: find $\hat{\beta}_0$ and $\hat{\beta}_1$ to minimize

$$\hat{\epsilon}_1^2 + \hat{\epsilon}_2^2 + \hat{\epsilon}_3^2$$



Geometric Perspective

We consider now a different perspective

Write regression as $y = i_3 \hat{\beta}_0 + X_{*1} \hat{\beta}_1 + \hat{\epsilon}$, i.e.,

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \hat{\beta}_0 + \begin{bmatrix} X_{11} \\ X_{21} \\ X_{31} \end{bmatrix} \hat{\beta}_1 + \begin{bmatrix} \hat{\epsilon}_1 \\ \hat{\epsilon}_2 \\ \hat{\epsilon}_3 \end{bmatrix}$$

For our data,

$$\begin{bmatrix} 1 \\ 8 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \hat{\beta}_0 + \begin{bmatrix} 2 \\ 5 \\ 7 \end{bmatrix} \hat{\beta}_1 + \begin{bmatrix} \hat{\epsilon}_1 \\ \hat{\epsilon}_2 \\ \hat{\epsilon}_3 \end{bmatrix}$$

The vectors $[1 \ 1 \ 1]^T$ and $[2 \ 5 \ 7]^T$ are vectors in \mathbb{R}^3

Geometric Perspective

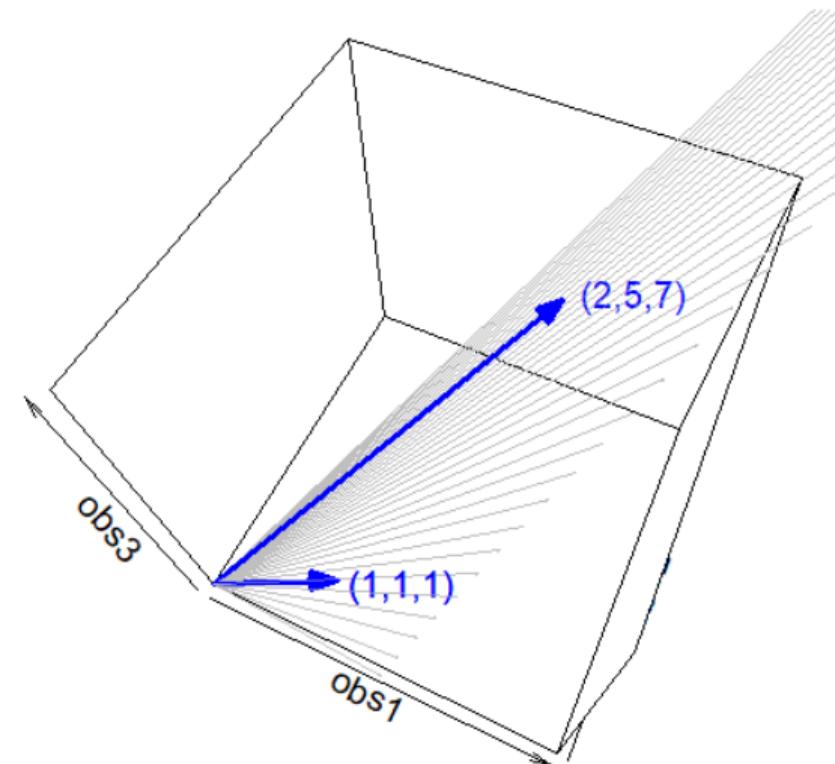
The vectors $[1 \ 1 \ 1]^T$ and $[2 \ 5 \ 7]^T$ are vectors in \mathbb{R}^3 (drawn in blue)

The set of all vectors

$$\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \hat{\beta}_0 + \begin{bmatrix} 2 \\ 5 \\ 7 \end{bmatrix} \hat{\beta}_1$$

(using all possible values of $\hat{\beta}_0$ and $\hat{\beta}_1$) forms a plane in \mathbb{R}^3

Part of plane shown in grey



Geometric Perspective

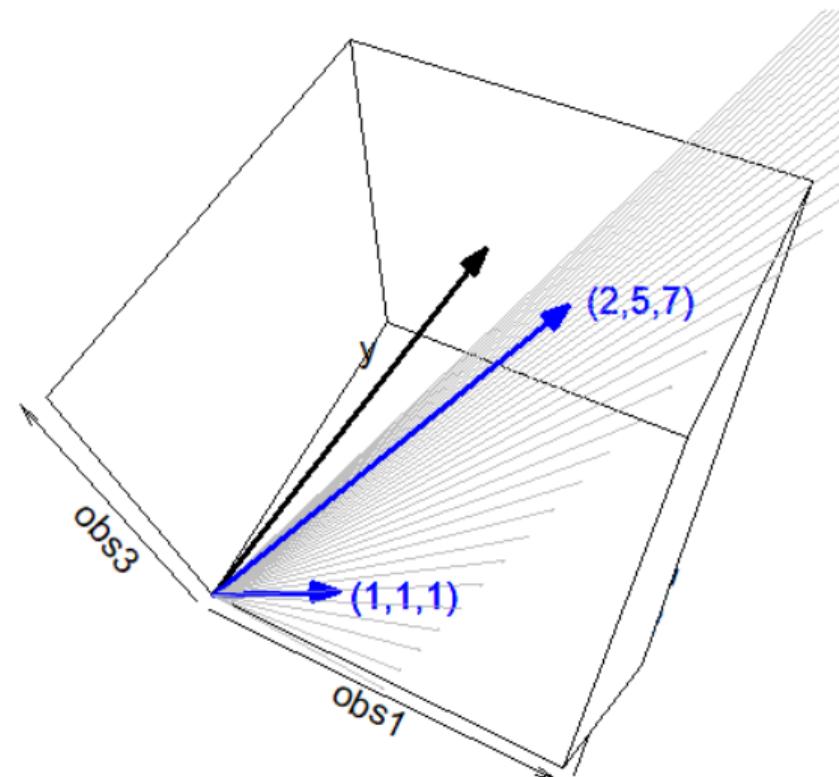
The vector $y = [1 \ 8 \ 4]^T$ is also a vector in \mathbb{R}^3 but it does not lie on the grey plane

y cannot be written as a linear combination of the blue vectors

If it did, the points

$$(X_{i1}, Y_i) = (2, 1), (5, 8), (7, 4)$$

would lie in a straight line when plotted on x - y plane



Geometric Perspective

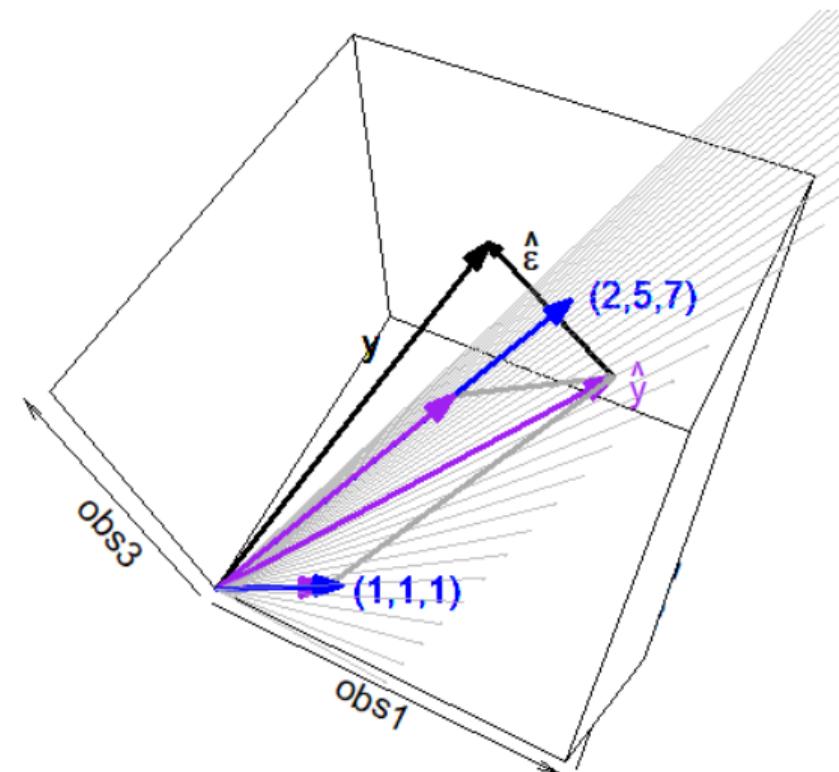
Let \hat{y} be any linear combination of the blue vectors, and $\hat{\epsilon}$ be such that

$$y = \hat{y} + \hat{\epsilon}$$

Qn: Which linear combination of the blue vectors makes \hat{y} “closest” to y (in the sense that $\hat{\epsilon}$ is as “short” as possible)?

Ans: The linear combination such that $\hat{\epsilon}$ is perpendicular to grey plane

- $\hat{\epsilon}$ is perpendicular to $[1 \ 1 \ 1]^T$ and to $[2 \ 5 \ 7]^T$



Geometric Perspective

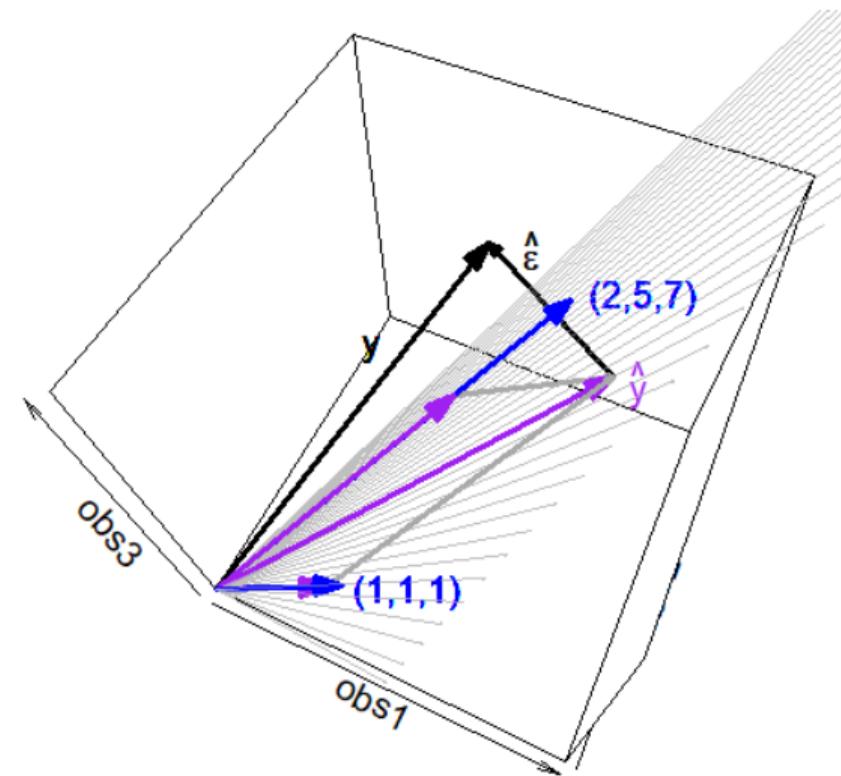
OLS chooses $\hat{\beta}_0$ and $\hat{\beta}_1$ such that

$$\hat{y} = \underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \hat{\beta}_0 + \begin{bmatrix} 2 \\ 5 \\ 7 \end{bmatrix} \hat{\beta}_1}_{\text{in Purple}} \quad \text{and} \quad \hat{\epsilon} = \underbrace{\begin{bmatrix} \hat{\epsilon}_1 \\ \hat{\epsilon}_2 \\ \hat{\epsilon}_3 \end{bmatrix}}_{\text{in Black}}$$

are orthogonal (i.e., perpendicular), thereby minimizing

$$\|\hat{\epsilon}\|^2 = \hat{\epsilon}_1^2 + \hat{\epsilon}_2^2 + \hat{\epsilon}_3^2$$

\hat{y} is called the “orthogonal” projection of Y on to the column space of X



Geometric Perspective

Intuition holds for general MLR case

$$y = X\hat{\beta} + \hat{\epsilon}$$

where y is $n \times 1$ vector, X is $n \times (k + 1)$ matrix whose columns space a $k + 1$ dimensional subspace of \mathbb{R}^n

Cannot draw pictures in n -dimensions!

Nonetheless, OLS minimization solution shows that $\|\hat{\epsilon}\|$ is minimized when $\hat{\beta} = \hat{\beta}^{ols}$, which makes \hat{y} and $\hat{\epsilon}$ orthogonal

$$\hat{y} = X\hat{\beta}^{ols} = \underbrace{X(X^T X)^{-1} X^T y}_{\text{Projection Matrix}}$$

Session 5.4

Session 5.4 Properties of OLS Estimators

- Unbiasedness under basic assumptions
- Standard errors under homoskedasticity and heteroskedasticity
- Best linear unbiasedness under homoskedasticity
- t-tests, F-tests

OLS Estimator Properties: Unbiasedness

- OLS estimator $\hat{\beta}^{ols}$ is unbiased:

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

$$E(\hat{\beta}^{ols} | X) = \beta + E((X^T X)^{-1} X^T \epsilon | X) = \beta + (X^T X)^{-1} X^T E(\epsilon | X) = \beta$$

$$E(\hat{\beta}^{ols}) = \beta$$

But remember points about causal interpretation

- you have estimated $E(Y | X) = X\beta$
- causal interpretation of coefficient only if you have included all relevant determinants of Y
- assuming no issues such as simultaneity, sampling problems, misspecification

OLS Estimator Variance (Homoskedasticity Case)

Under homoskedasticity $\text{Var}(\epsilon\epsilon^T \mid X) = \sigma^2 I_n$, we have

$$\begin{aligned}\text{Var}(\hat{\beta}^{ols} \mid X) &= \text{Var}(\beta + (X^T X)^{-1} X^T \epsilon \mid X) \\ &= (X^T X)^{-1} X^T \text{Var}(\epsilon \mid X) X (X^T X)^{-1} \\ &= (X^T X)^{-1} X^T (\sigma^2 I_n) X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} X^T I_n X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1}\end{aligned}$$

but must estimate σ^2

OLS Estimator Variance (Homoskedasticity Case)

Unbiased estimator of σ^2 is $\widehat{\sigma^2} = \frac{1}{n - k - 1} \sum_{i=1}^n \hat{\epsilon}_{i,ols}^2 = \frac{\hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols}}{n - k - 1}$

Proof: $\hat{\epsilon}_{ols} = My$ where $M = I_n - X(X^T X)^{-1}X^T$

note that

- $MX = (I_n - X(X^T X)^{-1}X^T)X = X - X = 0_{n \times (k+1)}$
- M is symmetric (Exercise!)
- M is idempotent, i.e., $MM = M$ (Exercise!)

Therefore $\hat{\epsilon}_{ols} = My = M(X\beta + \epsilon) = M\epsilon$ and

$$\hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols} = (M\epsilon)^T M\epsilon = \epsilon^T M^T M\epsilon = \epsilon^T M M\epsilon = \epsilon^T M\epsilon$$

OLS Estimator Variance (Homoskedasticity Case)

$$\begin{aligned} E(\hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols} | X) &= E(\epsilon^T M \epsilon | X) \\ &= E(\text{trace}(\epsilon^T M \epsilon) | X) \quad \text{because } \epsilon^T M \epsilon \text{ is a scalar} \\ &= E(\text{trace}(\epsilon \epsilon^T M) | X) \\ &= \text{trace}(E(\epsilon \epsilon^T M) | X) \quad \text{since trace is a sum} \\ &= \text{trace}(E(\epsilon \epsilon^T | X) M) = \text{trace}(\sigma^2 I_n M) = \sigma^2 \text{trace}(M) \\ &= \sigma^2 \text{trace}(I_n - X(X^T X)^{-1} X^T) = \sigma^2 (\text{trace}(I_n) - \text{trace}(X(X^T X)^{-1} X^T)) \\ &= \sigma^2 (n - \text{trace}((X^T X)^{-1} X^T X)) = \sigma^2 (n - \text{trace}(I_{k+1})) \\ &= \sigma^2 (n - k - 1) \end{aligned}$$

Therefore $E\left(\frac{\hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols}}{n - k - 1} \mid X\right) = \sigma^2$, which implies $E\left(\frac{\hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols}}{n - k - 1}\right) = \sigma^2$

OLS Estimator Variance (Heteroskedasticity Case)

$$\text{Var}(\hat{\beta}^{ols} | X) = (X^T X)^{-1} X^T \text{Var}(\epsilon | X) X (X^T X)^{-1}$$

$$= \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} [X_{1*}^T \quad X_{2*}^T \quad \dots \quad X_{n*}^T] \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix} \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix} \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}$$

$$= \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\sum_{i=1}^n \sigma_i^2 X_{i*}^T X_{i*} \right) \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}$$

Heteroskedasticity-Robust Estimator Variance-Covariance Matrix: replace σ_i^2 with $\hat{\epsilon}_{i,ols}^2$

$$\widehat{\text{Var}}(\hat{\beta}^{ols}) = \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\sum_{i=1}^n \hat{\epsilon}_{i,ols}^2 X_{i*}^T X_{i*} \right) \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}$$

OLS (Best Linear Unbiased, under Homoskedasticity)

We have shown that $\hat{\beta}^{ols}$ is unbiased

It is also “linear”

- A linear estimator is one with the form $\tilde{\beta} = Ay$
- Each $\tilde{\beta}_j = \sum_{i=1}^n a_{ji} Y_i$
- OLS estimator is $\hat{\beta}^{ols} = \underbrace{(X^T X)^{-1} X^T}_{A} y$

Now we show: OLS estimators are BLU (they have the smallest variance among all linear unbiased estimators)

OLS (Best Linear Unbiased, under Homoskedasticity)

BLU in the sense that

$$\text{Var}(c^T \hat{\beta}^{ols} | X) \leq \text{Var}(c^T \tilde{\beta} | X)$$

for all $(k + 1) \times 1$ vectors c , and for all unbiased estimators of the form $\tilde{\beta} = By$

- each individual $\hat{\beta}_k$ is BLU
- all linear combinations of $\hat{\beta}$ are BLU

Consider predicting Y at the new observation $X_{0*} = [1 \quad X_{01} \quad \dots \quad X_{0k}]$. OLS predictor is

$$\hat{Y}(X_{0*}) = X_{0*} \hat{\beta}^{ols}$$

which is a linear combination of the parameter estimates in $\hat{\beta}^{ols}$, i.e., OLS prediction rule gives us the most precise linear unbiased prediction of Y at X_{0*}

OLS (Best Linear Unbiased, under Homoskedasticity)

Proof of Efficiency: let $\tilde{\beta} = By$ be an unbiased estimator where $B \neq (X^T X)^{-1} X^T$

Let D be such that $B = D + (X^T X)^{-1} X^T$, so that

$$\begin{aligned}\tilde{\beta} &= By = (D + (X^T X)^{-1} X^T)y \\ &= (D + (X^T X)^{-1} X^T)(X\beta + \epsilon) \\ &= DX\beta + D\epsilon + \beta + (X^T X)^{-1} X^T \epsilon\end{aligned}$$

To ensure unbiasedness of $\tilde{\beta}$, we must assume $DX = 0$

OLS (Best Linear Unbiased, under Homoskedasticity)

Then

$$\begin{aligned}Var(\tilde{\beta} | X) &= E((\tilde{\beta} - \beta)(\tilde{\beta} - \beta)^T | X) \\&= E((D + (X^T X)^{-1} X^T) \varepsilon \varepsilon^T (D + (X^T X)^{-1} X^T)^T | X) \\&= (D + (X^T X)^{-1} X^T) E(\varepsilon \varepsilon^T | X) (D + (X^T X)^{-1} X^T)^T \\&= \sigma^2 (D + (X^T X)^{-1} X^T) (D + (X^T X)^{-1} X^T)^T \\&= \sigma^2 (D + (X^T X)^{-1} X^T) (D^T + X(X^T X)^{-1}) \\&= \sigma^2 [DD^T + (X^T X)^{-1} X^T D^T + DX(X^T X)^{-1} + (X^T X)^{-1}] \\&= \sigma^2 [DD^T + (X^T X)^{-1}] = \sigma^2 DD^T + \sigma^2 (X^T X)^{-1} = \sigma^2 DD^T + Var(\hat{\beta} | X)\end{aligned}$$

OLS (Best Linear Unbiased, under Homoskedasticity)

Therefore

$$\begin{aligned} \text{Var}(c^T \tilde{\beta} | X) &= c^T \text{Var}(\tilde{\beta} | X) c \\ &= c^T (\sigma^2 D D^T + \text{Var}(\hat{\beta} | X)) c \\ &= \sigma^2 c^T D D^T c + c^T \text{Var}(\hat{\beta} | X) c \\ &= \sigma^2 (D^T c)^T D^T c + \text{Var}(c^T \hat{\beta} | X) \geq \text{Var}(c^T \hat{\beta} | X) \end{aligned}$$

NB: $D^T c$ is a vector, therefore $(D^T c)^T D^T c$ is a sum of squares, which cannot be negative

Session 5.5

Session 5.5 Empirical Example

- Estimate Multiple Linear Regression using `earnings.csv` data

$$\begin{aligned} 100 \ln(\text{earn}) = & \beta_0 + \beta_1 \text{educ} + \beta_2 \text{educ}^2 + \beta_3 \text{height} + \beta_4 \text{male} + \\ & \beta_5 \ln(\text{wexp}) + \beta_6 \ln(\text{wexp})^2 + \beta_7 \ln(\text{tenure}) + \beta_8 \ln(\text{tenure})^2 + \\ & \beta_9 \text{age} + \beta_{10} \text{age}^2 + \epsilon \end{aligned}$$

- Everything computed twice
 - Using formulas derived in this session
 - Using existing R functions and packages

```
library(tidyverse)
library(car)      # using linearHypothesis()
library(sandwich) # using vcovHC()
library(lmtest)    # using coeftest()
```

Empirical Example

```
dat <- read_csv("data\\earnings2019.csv", show_col_types=FALSE) %>%
    mutate(learn100 = 100*log(earn), educsq = educ^2, lwexp = log(wexp), lwexp2 = log(wexp)^2,
          ltenure=log(tenure), ltenure2=log(tenure)^2, agesq = age^2)
y <- dat %>% select(learn100) %>% as.matrix()
X <- dat %>% mutate(intercept = 1) %>%
    select(c(intercept, educ, educsq, height, male, lwexp, lwexp2, ltenure, ltenure2, age, agesq)) %>%
    as.matrix()
XTX <- t(X)%*%X                                # (X'X)
XTXinv <- solve(XTX)                            # (X'X)^{-1}
XTy <- t(X)%*%y                                # X'y
betahat <- XTXinv %*% XTy                      # beta estimate
yhat <- X%*%betahat                            # fitted value
ehat <- y - yhat                               # residuals
n <- nrow(X); kplus1 <- ncol(X); df <- n - kplus1 # n - k - 1
SST <- sum((y - mean(y))^2); SSE <- sum((yhat - mean(yhat))^2); SSR <- sum(ehat^2)
s2hat <- SSR/df
betavar <- s2hat*XTXinv
betase <- sqrt(diag(betavar))
tstats <- betahat/betase
pvals <- pt(abs(tstats), df=df, lower.tail=FALSE)*2
R2 <- 1 - SSR/SST; AdjR2 <- 1 - (SSR/df)/(SST/(n-1))
```

Empirical Example

Result using lm() function

```
mdl1 <- lm(learn100 ~ educ + educsq + height + male
            lwexp + lwexp2 + ltenure + ltenure2 + age + agesq,
            sum1 <- summary(mdl1); sum1$coef %>% round(5)
            cat("\nR2:", sum1$r.squared,
                " Adj.R2:", sum1$adj.r.squared)
```

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|-----------|------------|-----------|----------|
| (Intercept) | 81.12348 | 37.53689 | 2.16117 | 0.03073 |
| educ | -15.26088 | 4.81587 | -3.16887 | 0.00154 |
| educsq | 1.00645 | 0.17185 | 5.85653 | 0.00000 |
| height | 1.34326 | 0.26387 | 5.09054 | 0.00000 |
| male | 18.07162 | 2.14105 | 8.44054 | 0.00000 |
| lwexp | -1.58201 | 2.71219 | -0.58330 | 0.55972 |
| lwexp2 | -0.22626 | 0.77976 | -0.29016 | 0.77170 |
| ltenure | 5.22349 | 2.66119 | 1.96284 | 0.04972 |
| ltenure2 | 2.31435 | 0.76945 | 3.00781 | 0.00264 |
| age | 5.82374 | 0.47516 | 12.25627 | 0.00000 |
| agesq | -0.06125 | 0.00517 | -11.84342 | 0.00000 |

R2: 0.303907 Adj.R2: 0.3024965

Result via direct computation

```
results <- cbind(betahat, betase, tstats, pvals)
colnames(results) <- c("Estimate", "Std. Error",
                        "t values", "Pr(>|t|)")
results %>% round(5)
cat("\nR2:", R2, " Adj.R2:", AdjR2)
```

| | Estimate | Std. Error | t values | Pr(> t) |
|-----------|-----------|------------|-----------|----------|
| intercept | 81.12348 | 37.53689 | 2.16117 | 0.03073 |
| educ | -15.26088 | 4.81587 | -3.16887 | 0.00154 |
| educsq | 1.00645 | 0.17185 | 5.85653 | 0.00000 |
| height | 1.34326 | 0.26387 | 5.09054 | 0.00000 |
| male | 18.07162 | 2.14105 | 8.44054 | 0.00000 |
| lwexp | -1.58201 | 2.71219 | -0.58330 | 0.55972 |
| lwexp2 | -0.22626 | 0.77976 | -0.29016 | 0.77170 |
| ltenure | 5.22349 | 2.66119 | 1.96284 | 0.04972 |
| ltenure2 | 2.31435 | 0.76945 | 3.00781 | 0.00264 |
| age | 5.82374 | 0.47516 | 12.25627 | 0.00000 |
| agesq | -0.06125 | 0.00517 | -11.84342 | 0.00000 |

R2: 0.303907 Adj.R2: 0.3024965

Empirical Example: Var-Cov under Homosked. (self-calc.)

```
options(width=400); round(betavar,3) # VAR-COV UNDER HOMOSKEDASTICITY
```

| | intercept | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|-----------|-----------|----------|--------|--------|--------|--------|--------|---------|----------|--------|--------|
| intercept | 1409.018 | -154.411 | 5.479 | -3.632 | 15.606 | 1.019 | 0.115 | -3.887 | 1.708 | -4.886 | 0.050 |
| educ | -154.411 | 23.193 | -0.825 | -0.110 | 1.009 | -0.360 | 0.076 | 0.215 | -0.099 | 0.104 | -0.001 |
| educsq | 5.479 | -0.825 | 0.030 | 0.004 | -0.033 | 0.014 | -0.002 | -0.010 | 0.004 | -0.004 | 0.000 |
| height | -3.632 | -0.110 | 0.004 | 0.070 | -0.385 | -0.016 | 0.003 | 0.009 | -0.002 | -0.004 | 0.000 |
| male | 15.606 | 1.009 | -0.033 | -0.385 | 4.584 | 0.081 | -0.020 | -0.077 | -0.002 | 0.038 | 0.000 |
| lwexp | 1.019 | -0.360 | 0.014 | -0.016 | 0.081 | 7.356 | -1.988 | -0.108 | 0.004 | -0.136 | 0.002 |
| lwexp2 | 0.115 | 0.076 | -0.002 | 0.003 | -0.020 | -1.988 | 0.608 | 0.010 | 0.006 | 0.020 | 0.000 |
| ltenure | -3.887 | 0.215 | -0.010 | 0.009 | -0.077 | -0.108 | 0.010 | 7.082 | -1.935 | -0.120 | 0.002 |
| ltenure2 | 1.708 | -0.099 | 0.004 | -0.002 | -0.002 | 0.004 | 0.006 | -1.935 | 0.592 | 0.012 | 0.000 |
| age | -4.886 | 0.104 | -0.004 | -0.004 | 0.038 | -0.136 | 0.020 | -0.120 | 0.012 | 0.226 | -0.002 |
| agesq | 0.050 | -0.001 | 0.000 | 0.000 | 0.000 | 0.002 | 0.000 | 0.002 | 0.000 | -0.002 | 0.000 |

```
betavar %>% diag %>% as.matrix %>% t %>% sqrt %>% round(3) # STANDARD ERRORS
```

| | intercept | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|------|-----------|-------|--------|--------|-------|-------|--------|---------|----------|-------|-------|
| [1,] | 37.537 | 4.816 | 0.172 | 0.264 | 2.141 | 2.712 | 0.78 | 2.661 | 0.769 | 0.475 | 0.005 |

Empirical Example: Var-Cov under Homoskedasticity with lm()

```
options(width=400); round(vcov(mdl1),3) # VAR-COV UNDER HOMOSKEDASTICITY lm()
```

| | (Intercept) | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|-------------|-------------|----------|--------|--------|--------|--------|--------|---------|----------|--------|--------|
| (Intercept) | 1409.018 | -154.411 | 5.479 | -3.632 | 15.606 | 1.019 | 0.115 | -3.887 | 1.708 | -4.886 | 0.050 |
| educ | -154.411 | 23.193 | -0.825 | -0.110 | 1.009 | -0.360 | 0.076 | 0.215 | -0.099 | 0.104 | -0.001 |
| educsq | 5.479 | -0.825 | 0.030 | 0.004 | -0.033 | 0.014 | -0.002 | -0.010 | 0.004 | -0.004 | 0.000 |
| height | -3.632 | -0.110 | 0.004 | 0.070 | -0.385 | -0.016 | 0.003 | 0.009 | -0.002 | -0.004 | 0.000 |
| male | 15.606 | 1.009 | -0.033 | -0.385 | 4.584 | 0.081 | -0.020 | -0.077 | -0.002 | 0.038 | 0.000 |
| lwexp | 1.019 | -0.360 | 0.014 | -0.016 | 0.081 | 7.356 | -1.988 | -0.108 | 0.004 | -0.136 | 0.002 |
| lwexp2 | 0.115 | 0.076 | -0.002 | 0.003 | -0.020 | -1.988 | 0.608 | 0.010 | 0.006 | 0.020 | 0.000 |
| ltenure | -3.887 | 0.215 | -0.010 | 0.009 | -0.077 | -0.108 | 0.010 | 7.082 | -1.935 | -0.120 | 0.002 |
| ltenure2 | 1.708 | -0.099 | 0.004 | -0.002 | -0.002 | 0.004 | 0.006 | -1.935 | 0.592 | 0.012 | 0.000 |
| age | -4.886 | 0.104 | -0.004 | -0.004 | 0.038 | -0.136 | 0.020 | -0.120 | 0.012 | 0.226 | -0.002 |
| agesq | 0.050 | -0.001 | 0.000 | 0.000 | 0.000 | 0.002 | 0.000 | 0.002 | 0.000 | -0.002 | 0.000 |

```
vcov(mdl1) %>% diag %>% as.matrix %>% t %>% sqrt %>% round(3) # STANDARD ERRORS
```

| | (Intercept) | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|------|-------------|-------|--------|--------|-------|-------|--------|---------|----------|-------|-------|
| [1,] | 37.537 | 4.816 | 0.172 | 0.264 | 2.141 | 2.712 | 0.78 | 2.661 | 0.769 | 0.475 | 0.005 |

Empirical Example: Het-Robust Var-Cov (self-calculated)

```
options(width=400);
varhat_HCO <- XTXinv%*%t(X)%*%diag(as.vector(ehat^2))%*%X%*%XTXinv; round(varhat_HCO,3)
```

| | intercept | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|-----------|-----------|----------|--------|--------|--------|--------|--------|---------|----------|--------|--------|
| intercept | 1345.536 | -144.697 | 5.136 | -3.296 | 14.791 | -2.703 | 1.435 | -10.416 | 3.477 | -5.957 | 0.063 |
| educ | -144.697 | 22.049 | -0.785 | -0.163 | 0.932 | 0.350 | -0.173 | 0.829 | -0.232 | 0.163 | -0.002 |
| educsq | 5.136 | -0.785 | 0.028 | 0.006 | -0.028 | -0.012 | 0.007 | -0.032 | 0.009 | -0.006 | 0.000 |
| height | -3.296 | -0.163 | 0.006 | 0.071 | -0.382 | -0.022 | 0.008 | 0.006 | -0.001 | -0.006 | 0.000 |
| male | 14.791 | 0.932 | -0.028 | -0.382 | 4.563 | 0.076 | -0.048 | 0.098 | -0.053 | 0.075 | -0.001 |
| lwexp | -2.703 | 0.350 | -0.012 | -0.022 | 0.076 | 8.468 | -2.300 | -0.121 | 0.041 | -0.208 | 0.002 |
| lwexp2 | 1.435 | -0.173 | 0.007 | 0.008 | -0.048 | -2.300 | 0.705 | 0.029 | -0.008 | 0.034 | -0.001 |
| ltenure | -10.416 | 0.829 | -0.032 | 0.006 | 0.098 | -0.121 | 0.029 | 7.554 | -2.077 | -0.030 | 0.001 |
| ltenure2 | 3.477 | -0.232 | 0.009 | -0.001 | -0.053 | 0.041 | -0.008 | -2.077 | 0.639 | -0.022 | 0.000 |
| age | -5.957 | 0.163 | -0.006 | -0.006 | 0.075 | -0.208 | 0.034 | -0.030 | -0.022 | 0.266 | -0.003 |
| agesq | 0.063 | -0.002 | 0.000 | 0.000 | -0.001 | 0.002 | -0.001 | 0.001 | 0.000 | -0.003 | 0.000 |

```
varhat_HCO %>% diag %>% as.matrix %>% t %>% sqrt %>% round(3) # STANDARD ERRORS
```

| | intercept | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|------|-----------|-------|--------|--------|-------|-------|--------|---------|----------|-----|-------|
| [1,] | 36.682 | 4.696 | 0.168 | 0.266 | 2.136 | 2.91 | 0.84 | 2.748 | | 0.8 | 0.516 |

Empirical Example : Het-Robust Var-Cov (using packages)

- Can use `hccm()` from `car` package or `vcovHC()` function from `sandwich` package

```
options(width=400); varhat_HC0a = vcovHC(mdl1,type="HC0"); round(varhat_HC0a,3) # HET-ROBUST VAR
```

| | (Intercept) | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|-------------|-------------|----------|--------|--------|--------|--------|--------|---------|----------|--------|--------|
| (Intercept) | 1345.536 | -144.697 | 5.136 | -3.296 | 14.791 | -2.703 | 1.435 | -10.416 | 3.477 | -5.957 | 0.063 |
| educ | -144.697 | 22.049 | -0.785 | -0.163 | 0.932 | 0.350 | -0.173 | 0.829 | -0.232 | 0.163 | -0.002 |
| educsq | 5.136 | -0.785 | 0.028 | 0.006 | -0.028 | -0.012 | 0.007 | -0.032 | 0.009 | -0.006 | 0.000 |
| height | -3.296 | -0.163 | 0.006 | 0.071 | -0.382 | -0.022 | 0.008 | 0.006 | -0.001 | -0.006 | 0.000 |
| male | 14.791 | 0.932 | -0.028 | -0.382 | 4.563 | 0.076 | -0.048 | 0.098 | -0.053 | 0.075 | -0.001 |
| lwexp | -2.703 | 0.350 | -0.012 | -0.022 | 0.076 | 8.468 | -2.300 | -0.121 | 0.041 | -0.208 | 0.002 |
| lwexp2 | 1.435 | -0.173 | 0.007 | 0.008 | -0.048 | -2.300 | 0.705 | 0.029 | -0.008 | 0.034 | -0.001 |
| ltenure | -10.416 | 0.829 | -0.032 | 0.006 | 0.098 | -0.121 | 0.029 | 7.554 | -2.077 | -0.030 | 0.001 |
| ltenure2 | 3.477 | -0.232 | 0.009 | -0.001 | -0.053 | 0.041 | -0.008 | -2.077 | 0.639 | -0.022 | 0.000 |
| age | -5.957 | 0.163 | -0.006 | -0.006 | 0.075 | -0.208 | 0.034 | -0.030 | -0.022 | 0.266 | -0.003 |
| agesq | 0.063 | -0.002 | 0.000 | 0.000 | -0.001 | 0.002 | -0.001 | 0.001 | 0.000 | -0.003 | 0.000 |

```
varhat_HC0a %>% diag %>% as.matrix %>% t %>% sqrt %>% round(3) # STANDARD ERRORS
```

| | (Intercept) | educ | educsq | height | male | lwexp | lwexp2 | ltenure | ltenure2 | age | agesq |
|------|-------------|-------|--------|--------|-------|-------|--------|---------|----------|-------|-------|
| [1,] | 36.682 | 4.696 | 0.168 | 0.266 | 2.136 | 2.91 | 0.84 | 2.748 | 0.8 | 0.516 | 0.006 |

Session 5.6

Session 5.6 Asymptotic Properties of OLS Estimators

- Consistency
- Asymptotic Normality
- Theory behind Heteroskedasticity-Robust Standard Errors

Asymptotic Properties of OLS Estimators

Recall

- Khinchine's LLN: If $\{Z_i\}_{i=1}^n$ iid, $E(Z_i) = \mu < \infty$, then $\bar{Z} \xrightarrow{p} \mu$.
- Lindeberg-Levy CLT: $\{Z_i\}_{i=1}^n$ iid, $E(Z_i) = \mu$ and $Var(Z_i) = \sigma^2 < \infty$, then

$$\sqrt{n}(\bar{Z} - \mu) \xrightarrow{d} \text{Normal}(0, \sigma^2)$$

- if $\mu = 0$, then

$$\sqrt{n}\bar{Z} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i \xrightarrow{d} \text{Normal}(0, \sigma^2).$$

Asymptotic Properties of OLS Estimators

Multivariate versions of these rules

If $\{Z_i\}_{i=1}^n$ iid vectors of random variables with $E(Z_i) = 0$ and $Var(Z_i) = \Omega$, for all i , then

$$\bar{Z} = \frac{1}{n} \sum_{i=1}^n Z_i = \begin{bmatrix} (1/n) \sum_{i=1}^n Z_{1i} \\ (1/n) \sum_{i=1}^n Z_{2i} \\ \vdots \\ (1/n) \sum_{i=1}^n Z_{ki} \end{bmatrix} \xrightarrow{p} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{bmatrix} = \mu$$

$$\sqrt{n} \bar{Z} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i = \begin{bmatrix} (1/\sqrt{n}) \sum_{i=1}^n Z_{1i} \\ (1/\sqrt{n}) \sum_{i=1}^n Z_{2i} \\ \vdots \\ (1/\sqrt{n}) \sum_{i=1}^n Z_{ki} \end{bmatrix} \xrightarrow{d} \text{Normal}_k(0, \Omega)$$

Asymptotic Properties

To talk about limiting distributions, we have to scale $\hat{\beta}$. Use

$$\sqrt{n}(\hat{\beta}^{ols} - \beta) = \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_{i*}^T \epsilon_i \right)$$

Our assumptions and the CLT imply

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n X_{i*} \epsilon_i \xrightarrow{d} \text{Normal}_{k+1}(0, S)$$

therefore

$$\sqrt{n}(\hat{\beta}^{ols} - \beta) = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}}_{\xrightarrow{P} \Sigma_{xx}^{-1}} \underbrace{\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_{i*}^T \epsilon_i \right)}_{\xrightarrow{P} \text{Normal}_{k+1}(0, S)} \xrightarrow{d} \text{Normal}_{k+1}(0, \Sigma_{xx}^{-1} S \Sigma_{xx}^{-1})$$

Asymptotic Properties

That is, $\hat{\beta}^{ols}$ is consistent, with asymptotic variance $Avar(\hat{\beta}^{ols}) = \Sigma_{xx}^{-1} S \Sigma_{xx}^{-1}$. This result justifies the approximation

$$Var(\hat{\beta}^{ols}) \approx (1/n) \Sigma_{xx}^{-1} S \Sigma_{xx}^{-1}.$$

An obvious estimator for Σ_{xx} is

$$\hat{\Sigma}_{xx} = \frac{1}{n} \sum_{i=1}^N X_{i*}^T X_{i*} = \frac{1}{n} X^T X$$

Some additional assumptions (see advanced econometrics textbooks) guarantee

$$\hat{S} = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \xrightarrow{p} S.$$

Asymptotic Properties

This allows us to consistently estimate the asymptotic variance of $\hat{\beta}$ by

$$\widehat{Avar}(\hat{\beta}^{ols}) = \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \right) \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}$$

and justifies the use of

$$\begin{aligned}\widehat{Var}_{HC0}(\hat{\beta}^{ols}) &= \frac{1}{n} \widehat{Avar}(\hat{\beta}^{ols}) \\ &= \frac{1}{n} \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \right) \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \\ &= \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \right) \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}\end{aligned}$$

Roadmap

- (Previous) Session 1: Statistics Review
- (Previous) Session 2: Simple Linear Regression
- (Previous) Session 3: Estimator Standard Errors; Multiple Linear Regression
- (Previous) Session 4: Matrix Algebra
- **This Session 5: OLS using Matrix Algebra**
- *Next Session 6: Hypothesis Testing*
- Session 7: Prediction
- Session 8: Instrumental Variable Regression
- Session 9: Logistic and Other Regressions
- Session 10: Panel Data Regressions
- Session 11: Introduction to Time Series
- Session 12: Time Series Regressions