

Estimating β

• some versions of linear model:

➤ observed data $(y_i, x_{i1}, \dots, x_{im}), i=1, \dots, n$. Regard y_i as a realization of a random variable Y_i . **functional form**

$$Y_i = \sum_{j=0}^{p-1} \beta_j \cdot g_j(x_{i1}, \dots, x_{im}) + \epsilon_i, \quad E(\epsilon_i) = 0, \quad \text{Var}(\epsilon_i) = \sigma^2, \quad \text{Cor}(\epsilon_{i_1}, \epsilon_{i_2}) = 0, \quad \text{for } i_1 \neq i_2$$

Annotations: Y_i is a r.v. (random variable), β_j are not r.v. (not random variables), ϵ_i is a random component, $\sum_{j=0}^{p-1} \beta_j \cdot g_j(\dots)$ is a deterministic component.

matrix form

Y: random, X_D : fixed; $Y = X_M \beta + \epsilon, \quad E(\epsilon) = 0, \quad \text{Var}(\epsilon) = \sigma^2 I$

Annotations: Y is a random vector with $E(Y) = X_M \beta$, $\text{Var}(Y) = \sigma^2 I$. X_M is the model matrix (e.g. data from DOE). $\epsilon \equiv Y - X_M \beta$ is the error vector. $\text{Cov}(\epsilon) = \sigma^2 I$ is the variance-covariance matrix.

Y: random, X_D : fixed; Y is a random vector with $E(Y) = X_M \beta$, $\text{Var}(Y) = \sigma^2 I$ ← easier to be extended to GLM (e.g. binomial y_x , Poisson y_x , ...)

Y: random, X_D : random; $E(Y|X_D) = X_M \beta, \quad \text{Var}(Y|X_D) = \sigma^2 I$

(X_0, Y) : joint distribution $\rightarrow Y|X_0 = X_M \beta + \epsilon, \quad E(\epsilon) = 0, \quad \text{Var}(\epsilon) = \sigma^2 I$.

$\rightarrow Y|X_0$: conditional dist. ← Linear model is built on this.

Notes: 1. At this stage, no specific distribution assumption imposed on ϵ ; only assume they (1) have zero mean, (2) have constant variance, and (3) are mutually uncorrelated.

2. parameters in the model: β and σ^2

i.e., ϵ_i 's can have any distributional forms (non-parametric approach)

sampling model (future lecture)

• a geometric view – an example

model: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i, i=1, 2, 3 \Rightarrow$

$Y = X\beta + \epsilon$

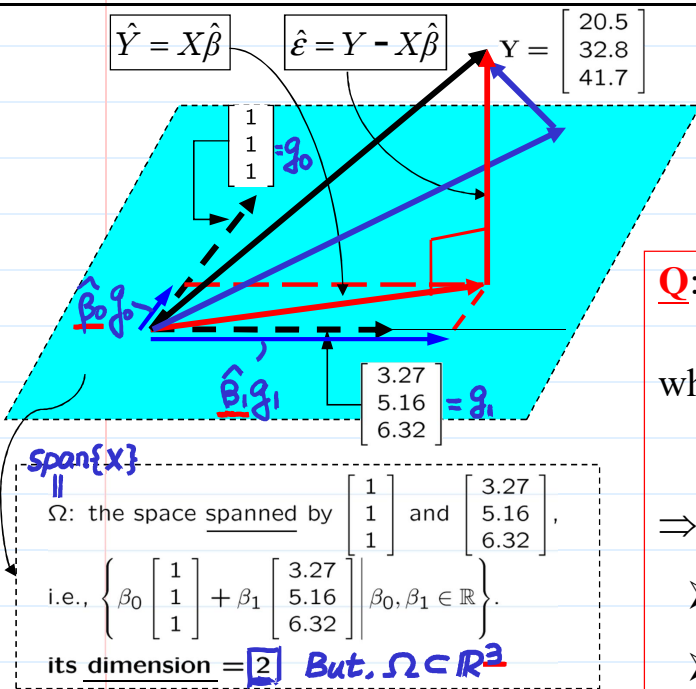
Sum of columns in X multiplied by β



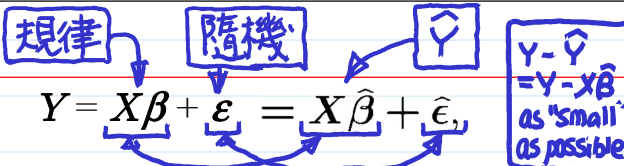
$$\begin{bmatrix} 20.5 \\ 32.8 \\ 41.7 \end{bmatrix} = \begin{bmatrix} 1 & 3.27 \\ 1 & 5.16 \\ 1 & 6.32 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \end{bmatrix} \Rightarrow \begin{bmatrix} 20.5 \\ 32.8 \\ 41.7 \end{bmatrix} = \beta_0 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + \beta_1 \begin{bmatrix} 3.27 \\ 5.16 \\ 6.32 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \end{bmatrix}$$

Annotations: β_0 is g_0 , β_1 is g_1 . $g_0, g_1 \in \mathbb{R}^3$.

of observations = 3, # of parameters in $\beta = 2 \Rightarrow Y \in \mathbb{R}^3, \beta \in \mathbb{R}^2$



estimation of β (i.e., finding $\hat{\beta}$) is equivalent to finding a vector (i.e., $X\hat{\beta}$) on the 2-dim model space $\Omega = \text{span}\{x\}$



Q: $Y = X\beta + \epsilon = X\hat{\beta} + \hat{\epsilon}$

what $\hat{\beta}$ would best "separate" $X\hat{\beta}$ from $\hat{\epsilon}$?

(Q: What's its meaning?)

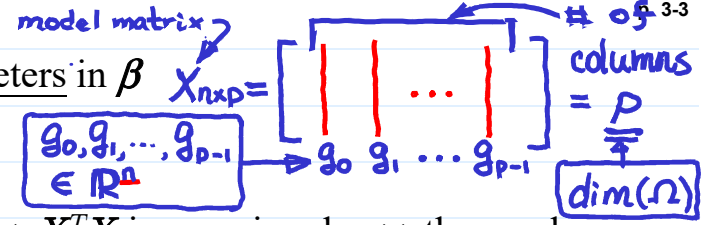
\Rightarrow find $\hat{\beta}$ s.t. $X\hat{\beta}$ is "close" to $Y \leftarrow \hat{Y} \approx Y$

➤ Q: is it reasonable? as similar as possible

➤ Q: what's the measurement of closeness?

• a geometric view – a general description

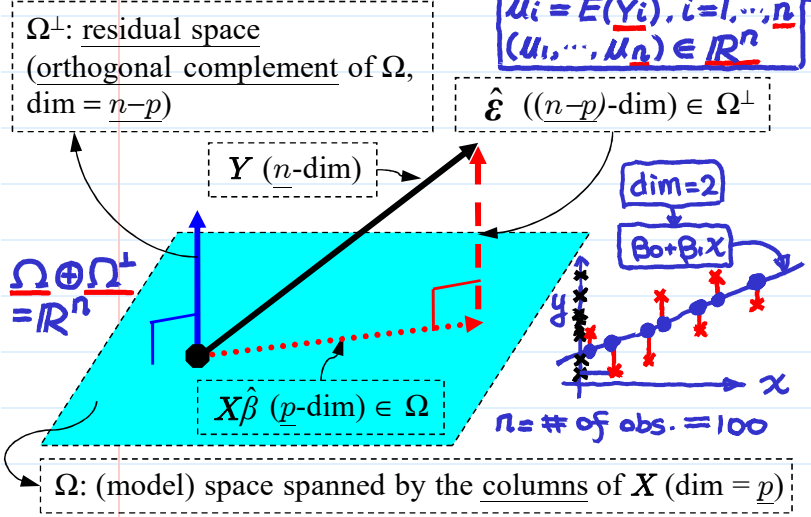
- $n = \#$ of observations; $p = \#$ of parameters in β
- $Y \in \mathbb{R}^n$, $\beta \in \mathbb{R}^p$
- assume $p \leq n$
- assume X is of full rank, $\text{rank}(X) = p$ ($\Leftrightarrow X^T X$ is non-singular \Leftrightarrow the p columns of X are linearly independent \Leftrightarrow the model space Ω spanned by the columns of X is of dimension p)



$$\Omega \subset \mathbb{R}^n, \Omega^\perp \subset \mathbb{R}^n$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} = Y = X\beta + \varepsilon = \beta_0 \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + \beta_1 \begin{bmatrix} g_{11} \\ g_{21} \\ \dots \\ g_{n1} \end{bmatrix} + \dots + \beta_{p-1} \begin{bmatrix} g_{1,p-1} \\ g_{2,p-1} \\ \dots \\ g_{n,p-1} \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_n \end{bmatrix}$$

$X^T X$: a $p \times p$ matrix



$Y \in \mathbb{R}^n$ (data) = $X\beta$ (systematic structure) + $\varepsilon \in \mathbb{R}^n$ (random error)

$n\text{-dim} = p\text{-dim} + (n-p)\text{-dim}$
 $Y = X\hat{\beta} + \hat{\varepsilon}$

Use a simpler structure (p -dim) to describe the complex (n -dim) Y

\Rightarrow if successful (i.e., $Y \approx X\hat{\beta}$), main structure of Y is p -dim

規律 \rightarrow

find $\hat{\beta}$ such that $X\hat{\beta}$ is as "close" as possible to Y , where closeness is measured by Euclidean distance (Q: when, i.e. under what assumption, is this a reasonable measurement of closeness?)

least square criterion: $\sum_{i=1}^n [y_i - (X\hat{\beta})_i]^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \|\hat{\varepsilon}\|^2$

Assumptions: $\text{Var}(\varepsilon) = \sigma^2 I$ (constant variance, uncorrelated)

$X\hat{\beta}$ is the orthogonal projection of Y onto Ω , (the space spanned by the columns of X), i.e.,

Surrogate of $X\beta$ (規律)

$$X\hat{\beta} = \frac{XY}{X(X^T X)^{-1} X^T Y}$$

an $n \times n$ matrix

IF $\Omega_1 = \text{span}\{x_1\} = \text{span}\{x_2\} = \Omega_2$, $x_1(x_1^T x_1)^{-1} x_1^T = H_1 = H_2 = x_2(x_2^T x_2)^{-1} x_2^T$

where $X(X^T X)^{-1} X^T$ is the orthogonal projection matrix of Ω (called hat matrix H).

- $\hat{\varepsilon}$: residuals = difference between Y and $X\hat{\beta}$ (i.e., $Y - X\hat{\beta}$)
 - surrogate of ε (error) OPM of $\Omega^\perp \rightarrow (I - H)Y = \hat{\varepsilon}$
 - $X\hat{\beta} \perp \hat{\varepsilon}$, \perp : geometrically orthogonal $\Rightarrow (X\hat{\beta})^T \hat{\varepsilon} = 0$
- different properties, r.v.'s, $\text{cov}(X\hat{\beta}, \hat{\varepsilon}) = 0_{n \times n}$

(Note: actually, the two random vectors are uncorrelated \Rightarrow when normal distribution assumption is imposed on ε , they become independent)

$n-p$: degree of freedom of $\hat{\varepsilon}$ (Q: what is a geometric interpretation for the degree of freedom?) (Q: what is the constraint on $\hat{\varepsilon}$?)

$X^T \hat{\varepsilon} = 0_{p \times 1}$

$\dim(\varepsilon) = n$ cf. $\hat{\varepsilon} \in \Omega^\perp \subset \mathbb{R}^n$, $\dim(\Omega^\perp) = n-p$

$\hat{\beta}$ is the ordinary least square estimator

If $\mathbb{1} \in X$ or $\mathbb{1} \in \Omega \Rightarrow \mathbb{1}^T \hat{\varepsilon} = 0$
 \uparrow vector of 1. $\uparrow \sum_{i=1}^n \hat{\varepsilon}_i$

- (ordinary) least square estimator **distribution form not specified**

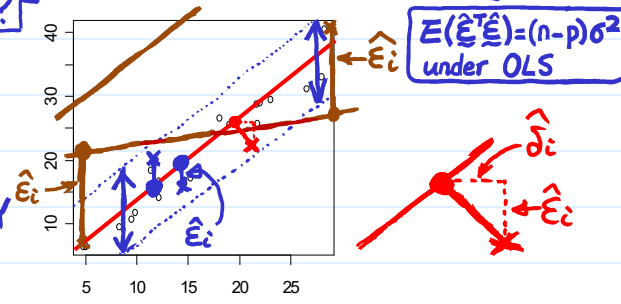
➤ assume ϵ are (i) uncorrelated (ii) equal variance ($Var(\epsilon) = \sigma^2 I$)

➤ define the best $\hat{\beta}$ as that minimizes sum of squared error: $\epsilon^T \epsilon = \sum_{i=1}^n \hat{\epsilon}_i^2$ **length² of $\hat{\epsilon}$**

(Q: why?) **reasonable criterion?**

Why not minimize $\sum_{i=1}^n |\epsilon_i|$? or minimize $\max_{1 \leq i \leq n} |\epsilon_i|$?

- easy to calculate
 - variation due to random error only appears on y-axis
 - variation minimization $\Rightarrow \hat{Y}$ close to Y
 - same scale in y_i 's



➤ $\epsilon^T \epsilon = (Y - X\beta)^T (Y - X\beta) = Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta$ (*)

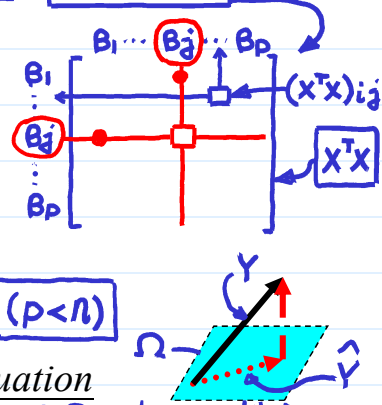
$Y^T - \beta^T X^T =$ \Rightarrow a second-order polynomial of β

➤ One method of finding the minimizer is to differentiate (*) w.r.t. β and set the derivatives equal to zero

$\Rightarrow \frac{\partial}{\partial \beta} \epsilon^T \epsilon = -2X^T Y + 2X^T X \beta = 0$

➤ By calculus, $\hat{\beta}$ is the solution of **X^T : $p \times n$ matrix ($p < n$)**

$X^T \hat{Y} = X^T Y \Leftarrow \hat{Y} \Leftarrow X^T X \beta = X^T Y$ \Leftarrow called normal equation **sufficient statistics of β (under normality)**



➤ assume $X^T X$ is non-singular (Q: when would it be singular?), **check LNp.3-3 Condition for X to be of full rank**

a linear function of y_i 's

unique solution
 $\hat{\beta} = (X^T X)^{-1} X^T Y \Rightarrow X \hat{\beta} = X (X^T X)^{-1} X^T Y \equiv \underline{HY}$

➤ $H_{n \times n} = X(X^T X)^{-1} X^T$ is called hat matrix in statistics. It is the orthogonal projection matrix onto Ω , the space spanned by the columns of X

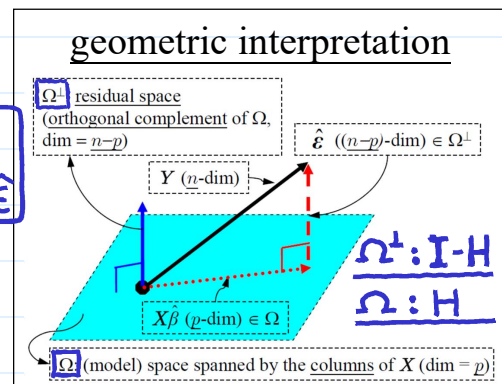
$H^2 = [X(X^T X)^{-1} X^T][X(X^T X)^{-1} X^T] = X(X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H$

- $H^2 = H \Rightarrow$ idempotent **a matrix with the 2 properties must be an orthogonal projection matrix**
- $H^T = H \Rightarrow$ symmetric

(exercise) $(I-H)^2 = (I-H)$ and $(I-H)^T = (I-H)$
 $H(I-H) = (I-H)H = 0$ **$n \times n$ matrix**
 $(I-H)Y = Y - HY = Y - X\hat{\beta} = \hat{\epsilon}$ **their columns orthogonal**

the eigenvalues of H (or $I-H$) are either 0 or 1; # of 1's = p (or $n-p$); # of 0's = $n-p$ (or p)

- $Hx = x$ if and only if x lies in Ω
- $Hx = 0$ if and only if x lies in Ω^\perp , the orthogonal complement of Ω



- predicted values: $\hat{Y} = X \hat{\beta} = \underline{HY}$ \leftarrow c.f. $X\beta$ **unknown** $Y = X\beta + \epsilon$
- residuals: $\hat{\epsilon} = Y - X \hat{\beta} = Y - \hat{Y} = (I-H)Y$ \leftarrow c.f. ϵ $\epsilon^T \epsilon$ **c.f.**
- residual sum of squares (RSS): $\hat{\epsilon}^T \hat{\epsilon} = [Y^T (I-H)^T][(I-H)Y] = Y^T (I-H)Y$
 $= \sum_{i=1}^n \hat{\epsilon}_i^2 =$ **length² of $\hat{\epsilon}$** $= \epsilon^T (I-H) \epsilon$

• examples of calculating ordinary least square estimator $\hat{\beta}$

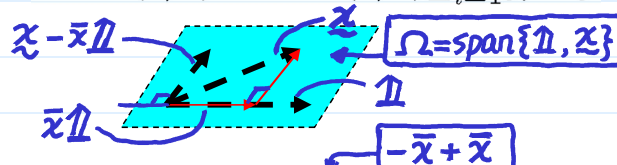
➤ example 1 (one-sample problem). functional form: $y_i = \mu + \epsilon_i, i = 1, \dots, n$.

$Y = X\beta + \epsilon$ ← **matrix form**

$Y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix} \equiv \mathbf{1}, \beta = [\mu]$

$X^T X = \mathbf{1}^T \mathbf{1} = n$

$\hat{\beta} = (X^T X)^{-1} X^T Y$
 $= (1/n) \mathbf{1}^T Y = (1/n) \sum_{i=1}^n y_i = \bar{y}$



➤ example 2 (simple regression). functional form: $y_i = \beta_0 + \beta_1(x_i) + \epsilon_i, i = 1, \dots, n$.

$Y = X\beta + \epsilon$ ← **matrix form** $\Rightarrow y_i = (\beta_0 + \beta_1 \bar{x}) + \beta_1(x_i - \bar{x}) + \epsilon_i \equiv \beta'_0 + \beta_1(x_i - \bar{x}) + \epsilon_i$

$Y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 & x_1 - \bar{x} \\ \dots & \dots \\ 1 & x_n - \bar{x} \end{bmatrix}, \beta = \begin{bmatrix} \beta'_0 \\ \beta_1 \end{bmatrix}$

$X^T X = \begin{bmatrix} 1 & \dots & 1 \\ x_1 - \bar{x} & \dots & x_n - \bar{x} \end{bmatrix} \begin{bmatrix} 1 & x_1 - \bar{x} \\ \dots & \dots \\ 1 & x_n - \bar{x} \end{bmatrix}$
 $= \begin{bmatrix} n & 0 \\ 0 & \sum_{i=1}^n (x_i - \bar{x})^2 \end{bmatrix}$

$\hat{\beta} = (X^T X)^{-1} X^T Y$
 $\hat{\beta}_0 = \hat{\beta}'_0 - \hat{\beta}_1 \bar{x} = \bar{y} - \hat{\rho} \frac{\hat{\sigma}_y}{\hat{\sigma}_x} \bar{x}$

regression effect: $y_i = (\bar{y} - \hat{\rho} \frac{\hat{\sigma}_y}{\hat{\sigma}_x} \bar{x}) + \hat{\rho} \frac{\hat{\sigma}_y}{\hat{\sigma}_x} x_i + \epsilon_i$
 $(y_i - \bar{y}) = \hat{\rho} \left(\frac{x_i - \bar{x}}{\hat{\sigma}_x} \right) + \left(\frac{\epsilon_i}{\hat{\sigma}_y} \right)$

$\hat{\beta}'_0 = \bar{y}$
 $\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \hat{\rho} \frac{\hat{\sigma}_y}{\hat{\sigma}_x}$

• mean and covariance matrix of OLS estimator $\hat{\beta}$

$\hat{\beta} = (X^T X)^{-1} X^T Y$ is a $p \times 1$ vector of random variables, so

➤ mean: $E(\hat{\beta}) = (X^T X)^{-1} X^T E(Y) = (X^T X)^{-1} X^T X \beta = \beta$ (i.e., unbiased)

i.e. we can control the X in data collection → **Recall. LN p. 1-2, 5 steps**
 2nd step: data collection

➤ $Cov(\hat{\beta}) = (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} = (X^T X)^{-1} \sigma^2$ (\Rightarrow irrelevant to Y and β)
Note: if we can control X , can decide the var-cov matrix before observing Y

standard error: $\hat{\theta}$ estimate $\rightarrow \theta$
 $E^*(\hat{\beta} \hat{\beta}^T) = E^*[(X^T X)^{-1} X^T Y Y^T X (X^T X)^{-1}]$
 a convenient notation (- mean) $E^*(Y Y^T) = \sigma^2 I$
FYI If $\epsilon \sim N(0, \sigma^2 I)$, then this is the inverse of Fisher information matrix

Var($\hat{\theta}$) Since $\hat{\beta}$ is a random vectors, $(X^T X)^{-1} \sigma^2$ is a variance-covariance matrix.
 ➤ $se(\hat{\beta}_i) = \sqrt{(X^T X)^{-1}_{ii}} \hat{\sigma}$ ← c.f. → $Var(\hat{\beta}_i) = (X^T X)^{-1}_{ii} \sigma^2$ (symmetric positive semi-definite)

➤ how to calculate the correlation between $\hat{\beta}_i$ and $\hat{\beta}_j$? (**Q:** Is this correlation

influenced by σ ?) (**Q:** What's the difference between this correlation and the correlation between predictors g_i and g_j ? **Ans:** the former is calculated from $(X^T X)^{-1}$ while the latter from $X^T X$)

also irrelevant to β
 $\sqrt{(X^T X)^{-1}_{ii}} \hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}} \hat{\sigma}$
 $X = \begin{bmatrix} | & \dots & | & \dots & | \\ \dots & & \dots & & \dots \\ | & \dots & | & \dots & | \end{bmatrix}$
 assume $g_i \perp g_j$
 If $X^T X = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$, $\Rightarrow (X^T X)^{-1} = \begin{bmatrix} 1/\sigma^2 & 0 \\ 0 & 1/\sigma^2 \end{bmatrix}$
 check LN p. 3-7, EX 2.
 $g_i \perp g_j \Rightarrow cor(\hat{\beta}_i, \hat{\beta}_j) = 0$
 In general cases, $(X^T X)_{ij} = 0 \Rightarrow (X^T X)^{-1}_{ij} = 0$

❖ Reading: Faraway (2005, 1st ed.), 2.3, 2.4, 2.5 LN p. 3-7, $g_i \perp g_j$
 ❖ Further reading: D&S, 4.4, 5.1, 5.2, 20.1, 20.2