

Confidence intervals and regions estimation

estimation — point estimation
interval/region estimation

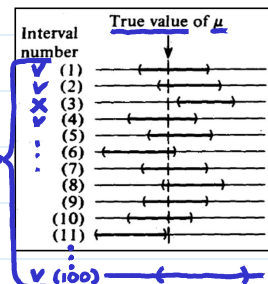
- **Q:** Why need interval/region estimation? What more information can it provide compared to point estimation?

e.g., estimate of $\beta = 3.5$, but accept $H_0: \beta = 0$. How to give such result an explanation? Why point estimation cause such confusing?

$\hat{\beta}$: a random variable

- An interval/region estimation provides
 - plausible values for parameter
 - uncertainty in parameter estimator
 - information about its length and the values it covers may be helpful
 - information related to testing

about 95% of these C.I.'s cover μ



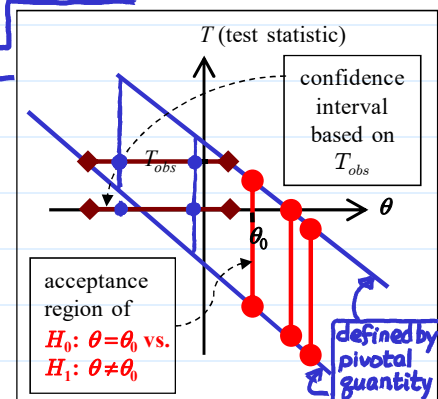
both information are contained/combined in C.I.

A C.I. offers more information than a single test

can use C.I. or C.R. to do test

- meaning of $100(1-\alpha)\%$ confidence interval or region, e.g., 95% confidence interval

- **duality** of interval/region estimation and hypothesis test: For a $100(1-\alpha)\%$ confidence region, any point θ that lies within the region represents a null hypothesis that would not be rejected at the $100\alpha\%$ significance level while every point θ outside represents a null hypothesis that would be rejected.



- Model: $Y = X\beta + \epsilon$, $\epsilon \sim N(0, \sigma^2 I)$; $\hat{\beta}$: OLS estimator $\Rightarrow \hat{\beta} \sim N(\beta, (X^T X)^{-1} \sigma^2)$

- Confidence region for $A\beta$, where A is a full rank $d \times p$ matrix and $d \leq p$

parameters $\rightarrow \zeta = (\zeta_1, \dots, \zeta_d)^T \equiv (N6) \text{ in LNp.4-3}$

$$A\hat{\beta} \sim N(A\beta, A(X^T X)^{-1} A^T \sigma^2) \Rightarrow [(A\hat{\beta} - A\beta)^T [A(X^T X)^{-1} A^T]^{-1} (A\hat{\beta} - A\beta)] / \sigma^2 \sim \chi^2_{d, \lambda}$$

(NI) in LNp.4-3

$(n-p) \hat{\sigma}^2 / \sigma^2 \sim \chi^2_{n-p}$

and they are independent.

$\frac{\chi^2_d / d}{\chi^2_{n-p} / (n-p)} \sim F_{d, n-p}$

$\therefore \hat{\beta}$ & $\hat{\sigma}^2$ are indep. (LNp.4-4)

pivotal quantity $\rightarrow [(A\hat{\beta} - A\beta)^T [A(X^T X)^{-1} A^T]^{-1} (A\hat{\beta} - A\beta)] / (d \hat{\sigma}^2) \sim F_{d, n-p}$

- $100(1-\alpha)\%$ confidence region of $A\beta$: collection of $A\beta$'s (or β) that satisfy

general form $[(A\hat{\beta} - A\beta)^T [A(X^T X)^{-1} A^T]^{-1} (A\hat{\beta} - A\beta)] / (d \hat{\sigma}^2) \leq F_{d, n-p}(\alpha)$

significance level $\rightarrow H_0: A\beta = \zeta_0$
 $H_1: A\beta \neq \zeta_0$

The regions are often ellipsoidally shaped (Q: why?)

Examples: **the overall F-test**

confidence region for β , i.e., $A = I_{p \times p}$

$$(\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \leq (p \hat{\sigma}^2) F_{p, n-p}(\alpha)$$

(Q: What's the confidence region for all effects?)

quadratic form: $M = \Gamma^T \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} \Gamma$

$x^T M x = z^T \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} z = \sum_{i=1}^n \lambda_i z_i^2 = \sum_{i=1}^n z_i^2$ if $\lambda_i \geq 0$ if $\Gamma^T \zeta$

M : positive semi-def.

confidence region of β_{i_2}, β_{j_2} , i.e., $A = \begin{pmatrix} 0, \dots, 0, 1, 0, \dots, 0, 0, 0, \dots, 0 \\ 0, \dots, 0, 0, 0, \dots, 0, 1, 0, \dots, 0 \end{pmatrix}$

test a pair of effects

$$[(\hat{\beta}_{i_2} - \beta_{i_2})^T \quad (\hat{\beta}_{j_2} - \beta_{j_2})^T]^T [A(X^T X)^{-1} A^T]^{-1} (\hat{\beta}_{i_2} - \beta_{i_2} - \beta_{j_2}) \leq (2 \hat{\sigma}^2) F_{2, n-p}(\alpha)$$

$(x^T x)^{-1} = \begin{bmatrix} \ddots & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$

$$\mathbf{Z} \equiv \begin{pmatrix} \hat{\beta}_i \\ \hat{\beta}_j \end{pmatrix} \sim N \left(\begin{pmatrix} \beta_i \\ \beta_j \end{pmatrix}, \sigma^2 \begin{pmatrix} (\mathbf{X}^T \mathbf{X})_{ii}^{-1} & (\mathbf{X}^T \mathbf{X})_{ij}^{-1} \\ (\mathbf{X}^T \mathbf{X})_{ji}^{-1} & (\mathbf{X}^T \mathbf{X})_{jj}^{-1} \end{pmatrix} \right) \equiv N(\mu, \sigma^2 \Sigma) \leftarrow \Sigma = A(\mathbf{X}^T \mathbf{X})^{-1} A^T$$

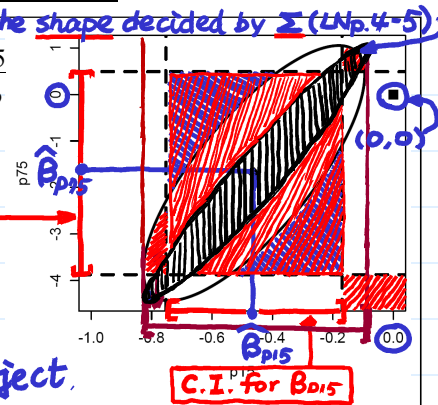
confidence region of β_i and β_j : $\{\mu \mid (\mathbf{Z} - \mu)^T \Sigma^{-1} (\mathbf{Z} - \mu) \leq c\}$ for some c **cf.**

example: confidence region and intervals of β_{p15} and β_{p75} **the shape decided by Σ (LNp.4-5)**

Q1: why the straight lines not tangential to the ellipse?

$$\begin{aligned} 1 - \alpha &= P(\{(\beta_{p15}, \beta_{p75}) \in \text{C. Region}\}) \\ &= P(\{(\beta_{p15}, \beta_{p75}) \in \text{C.I.}_{p15} \times \mathbb{R}\}) \\ &= P(\{(\beta_{p15}, \beta_{p75}) \in \mathbb{R} \times \text{C.I.}_{p75}\}) \end{aligned}$$

C.I. for β_{p75}



Q2: what can you say, based on the plot, about the results of testing $H_0^1: \beta_{p15}=0$, $H_0^2: \beta_{p75}=0$, and $H_0^3: \beta_{p15}=\beta_{p75}=0$?

Ans. H_0^1 : reject, H_0^2 : accept, H_0^3 : reject.

Q3: where will be the point (0,0) located if the data accept H_0^1 , H_0^2 , reject H_0^3 ?

how to explain the result if (0,0) falls in other regions? (exercise)

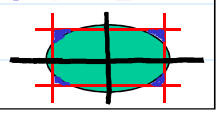
Q4: what is the correlation between $\hat{\beta}_{p15}$ and $\hat{\beta}_{p75}$? how will the shape of ellipse change when the correlation becomes larger or smaller?

$$\frac{(\mathbf{X}^T \mathbf{X})_{ij}^{-1}}{\sqrt{(\mathbf{X}^T \mathbf{X})_{ii}^{-1} (\mathbf{X}^T \mathbf{X})_{jj}^{-1}}}$$

Q5: can you see why the situation in Q3 will happen more frequently when the correlation between $\hat{\beta}_{p15}$ and $\hat{\beta}_{p75}$ gets larger?

LNp.4-13. graph for $\hat{\beta}_j, \hat{\beta}_k$ with strong collinearity

Q6: if $\hat{\beta}_{p15}$ and $\hat{\beta}_{p75}$ are uncorrelated, what would be the shape of the confidence region? why situation in Q3 less possible to occur?



confidence interval for β_i , i.e., $\mathbf{A}=(0, \dots, 0, 1, 0, \dots, 0)$

test just one effect

$$(\hat{\beta}_i - \beta_i)^2 / (\mathbf{X}^T \mathbf{X})_{ii}^{-1} \leq \sigma^2 F_{1, n-p}(\alpha) \Rightarrow |(\hat{\beta}_i - \beta_i) / (\hat{\sigma} \sqrt{(\mathbf{X}^T \mathbf{X})_{ii}^{-1}})| \leq t_{n-p}(\alpha/2)$$

alternative method:

① $\hat{\beta}_i \sim N(\beta_i, \sigma^2 (\mathbf{X}^T \mathbf{X})_{ii}^{-1})$, ② $(n-p) \hat{\sigma}^2 / \sigma^2 \sim \chi^2_{n-p}$, and ③ they are independent

pivotal quantity

$$\Rightarrow \left| \frac{\hat{\beta}_i - \beta_i}{\hat{\sigma} \sqrt{(\mathbf{X}^T \mathbf{X})_{ii}^{-1}}} \right| \leq t_{n-p}(\alpha/2) \Rightarrow \text{C.I.: } \hat{\beta}_i \pm t_{n-p}(\alpha/2) \times \left(\hat{\sigma} \sqrt{(\mathbf{X}^T \mathbf{X})_{ii}^{-1}} \right)$$

$$N(0,1) \sim \frac{\hat{\beta}_i - \beta_i}{\sqrt{(\mathbf{X}^T \mathbf{X})_{ii}^{-1}}} / \sqrt{\frac{(n-p) \hat{\sigma}^2 / \sigma^2}{n-p}} \sim \sqrt{\frac{\chi^2_{n-p}}{n-p}}$$

center of C.I. critical value $se(\hat{\beta}_i)$

confidence interval for prediction of mean response at x_0 (LNp.4-4) $\rightarrow x_0^T \beta \Rightarrow A = x_0^T$

$$x_0^T \hat{\beta} - x_0^T \beta \sim N(0, (x_0^T (\mathbf{X}^T \mathbf{X})^{-1} x_0) \sigma^2) \Rightarrow \left| (x_0^T \hat{\beta} - x_0^T \beta) / (\hat{\sigma} \sqrt{x_0^T (\mathbf{X}^T \mathbf{X})^{-1} x_0}) \right| \leq t_{n-p}(\alpha/2)$$

$$\Rightarrow \text{C.I.: } x_0^T \hat{\beta} \pm t_{n-p}(\alpha/2) \times \left(\hat{\sigma} \sqrt{x_0^T (\mathbf{X}^T \mathbf{X})^{-1} x_0} \right)$$

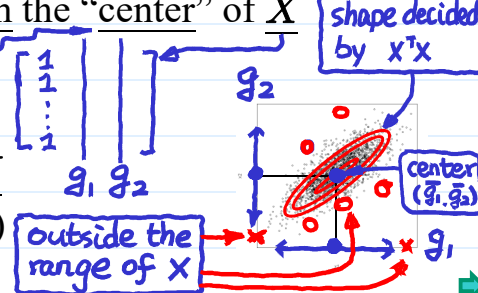
a quadratic form
Mahalanobis distance (future lecture)

Q: for a given dataset and α , the length of the C.I. is related to x_0 only. What x_0 will cause a wider C.I.? Ans: x_0 that is away from the "center" of X

interpolation and extrapolation

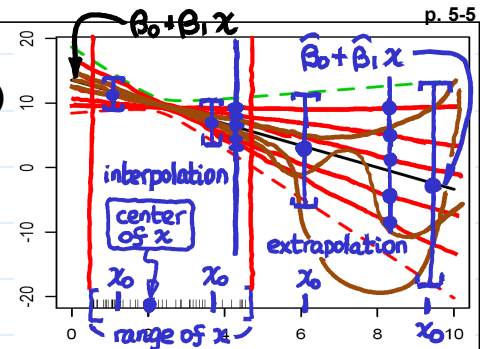
内插 外插

interpolation: x_0 lie "within the range" of X
extrapolation: x_0 lie "outside the range" of X



(Q: fitted model still hold outside the range?)
✓ quantitative x_0 and qualitative x_0

Example: 95% *pointwise* confidence band for prediction of mean responses (model: $y = \beta_0 + \beta_1 x + \varepsilon$)



Q1: why the confidence intervals get wider when we move away from the range of data?

Ans: wider C.I. $\leftarrow \varepsilon$ cause it

Ans: fitted model may approximate the true model very badly or wrongly & you don't have any information about it.

Q2: what's the danger of extrapolation?

Q3: does the widening reflect the possibility that the mean structure of the model may change outside the range?

Q4: does the plot represent a simultaneous confidence band for all prediction of mean response?

wider C.I. \Rightarrow less accurate inference

No. Because the confidence band is constructed under the assumption that the fitted model still hold outside the range of x

No. S.C.B will be wider than P.C.B.

pointwise confidence band

C.I. for prediction of future observation at x_0

$$y_{x_0} \text{ (r.v.)} - (x_0^T \hat{\beta} - (x_0^T \beta + \varepsilon)) \sim N(0, (x_0^T (X^T X)^{-1} x_0 + 1) \sigma^2)$$

estimate \Rightarrow C.I.: $x_0^T \hat{\beta} \pm t_{n-p}(\alpha/2) \times \left(\hat{\sigma} \sqrt{1 + x_0^T (X^T X)^{-1} x_0} \right)$ $se(x_0^T \hat{\beta} - \varepsilon)$

not reduced when $n \uparrow$

a general form for confidence interval: \leftarrow check the C.I. for β_i 's, prediction, ...

estimate \pm (critical value) \times (standard error of estimate)

Reading: Faraway (2005, 1st ed.), 3.4, 3.5

Futher reading: D&S, 5.3, 5.4, 5.5

C.I.: a combination of $\hat{\theta}$ & $se(\hat{\theta})$

with a sampling plan

Sampling model

experimental data vs. observational data

It depends on whether we have control over predictors

examples: yield of crop. \leftarrow response Y

experimental: fertilizer, ...

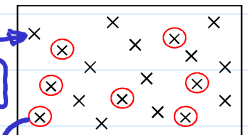
observational: exposure, weather, ...

an object or a unit

(Y, X_1, \dots, X_m)

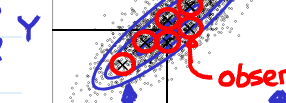
r.v.

population



observational data (no plan or with a sampling plan)

not observed



What's the difference btw no plan & with a sampling plan?

Q: What difference between inferences based on experimental data and observational data?

Ans: experimental data: causation, \leftarrow 因果 \leftrightarrow 因果關係

observational data: often only association (Note. lurking variable)

Q: Is this model description, $Y = X\beta + \varepsilon$, $\varepsilon \sim N(0, \sigma^2 I)$, appropriate for observational data? Note that

(1) observational X are random variables

(2) in LM, X are treated as fixed values, i.e., no distribution assigned for X

Q: difference between "X is random" and "X measured with (random) error"

example: X is random (why?), but measured accurately

measurement error model (future lecture)

$$x^* = x + \delta$$

random error

observed value (random)

true value (fixed or random)

① $X \rightarrow Y$ ② $Y \rightarrow X$
③ $X \xrightarrow{L} Y$ ④ random disturbance

LNp. 5-12

for some data sets, we can regard the data as a sample drawn from a population. In the case, we want to say something about the unknown population values using estimated values that are obtained from the sampld data. (example?)

- better chance to achieve pattern in sample \approx pattern in population \Rightarrow good sample (representative, LNp.1-3)
- the data should be generated using a “(simple) random sample” of the population so that they can be representative \leftarrow each object/unit in the population has the same known probability to be included in the sample
- conditional distribution of multivariate normal: If

Z in LNp.4-3

cf. $(p+q) \times 1$ variables $Z = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix} \right)$

then $Z_1 | Z_2 = z_2 \sim N \left(\mu_1 + \frac{\sigma_{12}}{\sigma_2^2} (z_2 - \mu_2), \sigma_1^2 - \frac{\sigma_{12}^2}{\sigma_2^2} \right)$

Recall simple regression (LNp.3-7)
 $\frac{(y - \mu_y)}{\sigma_y} = \rho \frac{(x - \mu_x)}{\sigma_x} \rightarrow$ cf. regression effect (LNp.2-7)
 $\Rightarrow y = \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x)$
 $\rho = \text{COV}(x, y) \times (\sigma_x^2)^{-1}$

check LNp.3-7? Schur complement

regression line \leftarrow cf. TSS in $R^2 \leftarrow$ cf. RSS in $R^2 \leftarrow$ cf. $\frac{TSS - RSS}{(p+q)\sigma_y^2 / \sigma_x^2} = \rho^2 \times \sigma_y^2$

- an alternative view of regression: data $(y_i, x_i), i=1, \dots, n$, are randomly sampled from a multivariate Normal population, $y_i | x_i = x_i \sim N(\mu_y - \beta^T \mu_x + \beta^T x_i, \sigma^2)$.

vector of p i.v.'s $\begin{bmatrix} Y_i \\ X_i \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_Y \\ \mu_X \end{bmatrix}, \begin{bmatrix} \sigma_Y^2 & \Sigma_{XY}^T \\ \Sigma_{XY} & \Sigma_{XX} \end{bmatrix} \right)$

cf. $y_i | x_i \sim N(x_i^T \beta, \sigma^2)$ $\beta_1 x_1 + \dots + \beta_p x_p$

constant variance σ^2

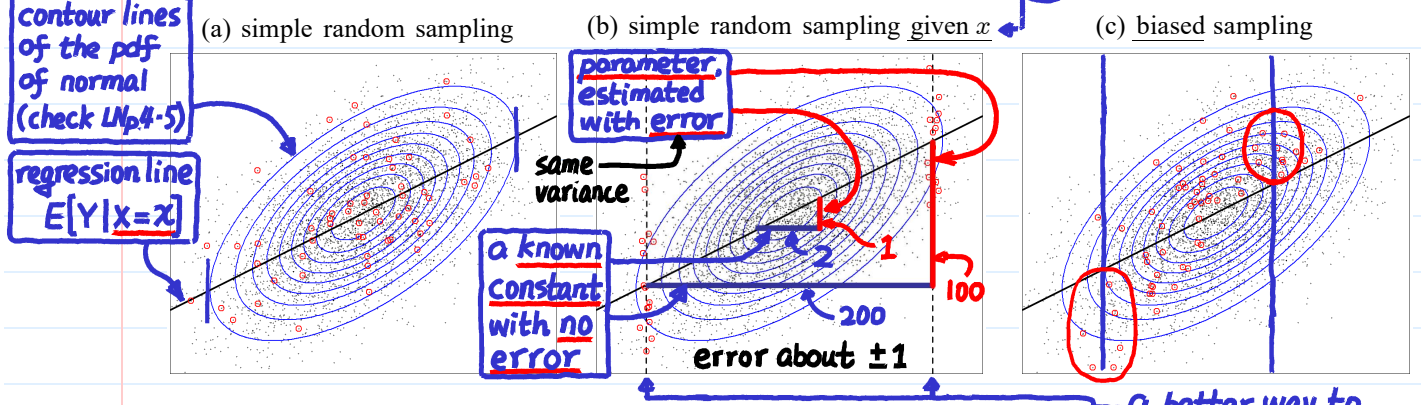
fixed intercept $\mu_y - \beta^T \mu_x$

It is a linear model with $\beta = \Sigma_{XX}^{-1} \Sigma_{XY}$, $\sigma^2 = \sigma_Y^2 - \Sigma_{XY}^T \Sigma_{XX}^{-1} \Sigma_{XY} \equiv \sigma_Y^2 (1 - r^2)$.

$\Sigma_{XY} = 0 \Rightarrow \beta = 0$ \leftarrow $\beta^T = \Sigma_{XY}^T \Sigma_{XX}^{-1}$ \leftarrow R^2 estimate $[r^2 = 1 - (\sigma^2 / \sigma_Y^2)] = \sigma_Y^2 (\rho / \sigma_Y)^2$

When we are interested in the “transformed” parameters, regression can be applied.

Q: what information in these samples is proper? DOE



	(a)	(b)	(c)
Joint distribution (X, Y)	✓	✗	✗
Conditional distribution Y X	✓	✓✓	✗
Marginal distribution X	✓	✗	✓ or ✗

- Q: what will happen if the sample is not random?
 - (i) biased sample \rightarrow sample is not representative
 - (ii) sample of convenience \rightarrow may or may not be a good sample
 - (iii) sample = population \rightarrow (1) estimation is exact (2) testing is questionable.
- these nonrandom samples can cause problems in the inference (e.g., R^2 , LNp.3-18)

❖ **Reading**: Faraway (2005, 1st ed.), 3.8, nonrandom samples
 Weisberg (2005), *Applied Linear Regression*, 3rd Ed., 4.2, 4.3

large population \supset small population \downarrow sampling \leftarrow permutation test

sampling probability of each object/unit could be unknown

Orthogonality

- Q: consider the two models:

model 1: $y = \beta_0 + \beta_1 x_1 + \varepsilon$

model 2: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$

- $\text{span}(x_1) \perp \text{span}(x_2)$
 \Leftrightarrow orthogonality
- $\text{span}(x_1) \supseteq \text{span}(x_2)$
 \Leftrightarrow aliasing (in DOE)

fitted model=model 1: $Y = X_1 \beta_1 + \varepsilon$

true model=model 2: $Y = X_1 \beta_1 + X_2 \beta_2 + \varepsilon$

- $E(\hat{\beta}_1) = \beta_1 + (X_1^T X_1)^{-1} X_1^T X_2 \beta_2$

- $E(X_1 \hat{\beta}_1) = X_1 \beta_1$

Mat matrix of X_1 $+ X_1 (X_1^T X_1)^{-1} X_1^T (X_2 \beta_2)$ LNp.3-11

Note. If fitted model=model 2

- $E(\hat{\beta}_1) = \beta_1$

In general, $\hat{\beta}_1$, in the two models are not identical (Q: why?) \rightarrow check the graph in LNp.4-13

[also, test $H_0: \beta_1=0$ (or c) not identical] 3-13

an exception: when x_1 and x_2 are orthogonal

- $Y = X\beta + \varepsilon = X_1 \beta_1 + X_2 \beta_2 + \varepsilon$, where $\beta = [\beta_1 \beta_2]^T$ and $X = [X_1 \ X_2]$ with the property $X_1^T X_2 = 0 \Rightarrow X_1$ and X_2 are orthogonal (generalization?)

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = (X^T X)^{-1} X^T Y$$

$$= \begin{bmatrix} \square & 0 \\ 0 & \square \end{bmatrix} \begin{bmatrix} X_1^T Y \\ X_2^T Y \end{bmatrix}$$

$$= \begin{bmatrix} (X_1^T X_1)^{-1} X_1^T Y \\ (X_2^T X_2)^{-1} X_2^T Y \end{bmatrix}$$

$$X^T X = \begin{pmatrix} X_1^T X_1 & X_1^T X_2 \\ X_2^T X_1 & X_2^T X_2 \end{pmatrix} = \begin{pmatrix} X_1^T X_1 & 0 \\ 0 & X_2^T X_2 \end{pmatrix}$$

$$\Rightarrow (X^T X)^{-1} = \begin{pmatrix} (X_1^T X_1)^{-1} & 0 \\ 0 & (X_2^T X_2)^{-1} \end{pmatrix}$$

$X = [X_1 \ X_2 \ \dots \ X_k]$

$$X^T X = \begin{bmatrix} X_1^T X_1 & 0 & \dots & 0 \\ 0 & X_2^T X_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & X_k^T X_k \end{bmatrix}$$

By (N4) in LNp.4-3

- Estimation: $\hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T Y$, $\hat{\beta}_2 = (X_2^T X_2)^{-1} X_2^T Y$, and $\hat{\beta}_1, \hat{\beta}_2$ independent

\Rightarrow note that $\hat{\beta}_1$ will be the same regardless of whether X_2 is in the model or not (and vice versa). fitted model: $Y = X_1 \beta_1 + \varepsilon \Rightarrow \hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T Y$

- Q: what if only two predictors, say some x_i in X_1 and some x_j in X_2 , are orthogonal?

$x_i \perp x_j \Leftrightarrow x_i^T x_j = 0 \Leftrightarrow (X^T X)_{ij} = 0$
 But. $X_1^T X_2 \neq 0 \Rightarrow (X^T X)_{ij} \neq 0$
 $\Rightarrow \hat{\beta}_i, \hat{\beta}_j$ may not be indep. \rightarrow