

# LM, LNp.7-1~2 → Residual Analysis: Theory

- Theory: define the **residual** for the  $i^{th}$  observation  $(x_i, y_i)$  as

$$\hat{\epsilon}_i = r_i = y_i - \hat{y}_i, \quad \hat{y}_i = \mathbf{x}_i^T \hat{\beta}$$

$\hat{\epsilon}_i$  often used to check assumptions:

$\hat{y}_i$  contains information given by the model;  $r_i$  is the "difference" between  $y_i$  (observed) and  $\hat{y}_i$  (fitted) and contains information on possible model inadequacy.

①  $\epsilon_i$ 's iid  $N(0, \sigma^2)$   
 ②  $E(Y) = X\beta$  is a correct mean structure

- Vector of residuals  $\mathbf{r} = (r_1, \dots, r_N)^T = \mathbf{y} - \mathbf{X}\hat{\beta} = (\mathbf{I} - \mathbf{H})\mathbf{y}$
- Under the model assumption  $E(\mathbf{y}) = \mathbf{X}\beta$ , it can be shown that

(a)  $E(\mathbf{r}) = \mathbf{0}$ ,  $\mathbf{r} \perp \hat{\mathbf{y}}$  i.e.,  $X\beta$  is correct model

(b)  $\mathbf{r}$  and  $\hat{\mathbf{y}}$  are independent,

Note: variance of  $r_i$   
 $\propto 1 - h_i$  (leverage) =  $1 - H_{ii}$   
 $\rightarrow \propto$  Mahalanobis dist. btwn design pts &  $\bar{X}$

overall pattern  
 individual observation (outlier)

(c) variances of  $r_i$  are nearly constant for "nearly balanced" designs.  
 $\text{cov}(\mathbf{r}) = \sigma^2(\mathbf{I} - \mathbf{H})$

①  $\mathbf{Y} = \mathbf{X}\beta + \epsilon = \hat{\mathbf{Y}} + \hat{\epsilon}$

②  $\mathbf{Y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \epsilon = (\mathbf{X}_1\beta_1 + \mathbf{H}_1\mathbf{X}_2\beta_2) + ((\mathbf{I} - \mathbf{H}_1)\mathbf{X}_2\beta_2 + \epsilon) = \hat{\mathbf{Y}}_{X_1} + \hat{\epsilon}_{X_1}$

Under fitted model:  $\mathbf{Y} = \mathbf{X}_1\beta_1 + \epsilon^*$  ( $\mathbf{H}_1 = \mathbf{X}_1(\mathbf{X}_1^T\mathbf{X}_1)^{-1}\mathbf{X}_1^T$ )



$\sigma_x \leftrightarrow \mu_x$

# Residual Plots → LM, LNp.7-7~9

- Plot  $r_i$  vs.  $\hat{y}_i$  (see Figure 1): It should appear as a parallel band around 0. Otherwise, it would suggest model violation. If spread of  $r_i$  increases as  $\hat{y}_i$  increases, error variance of  $y$  increases with mean of  $y$ . May need a transformation of  $y$ . (Will be explained in future lecture.)

- Plot  $r_i$  from replicates per treatment (see Figure 2): to see if error variance depends on treatment.

qualitative factor  
 Box plots

Note: saturated model

$\hat{r}$  vs. predictor

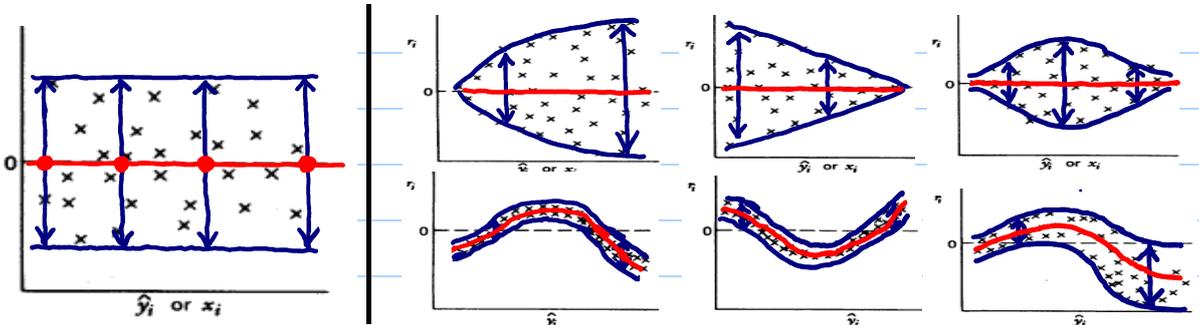
- Plot  $r_i$  vs.  $x_i$ : If not a parallel band around 0, relationship between  $y_i$  and  $x_i$  not fully captured, revise the  $\mathbf{X}\beta$  part of the model.

quantitative factor

- Plot  $r_i$  vs. time sequence: to see if there is a time trend or autocorrelation over time.

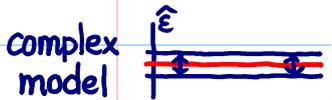
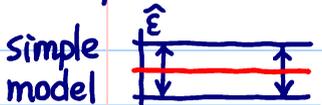
if available (or run order, measure order, ...)

null plot





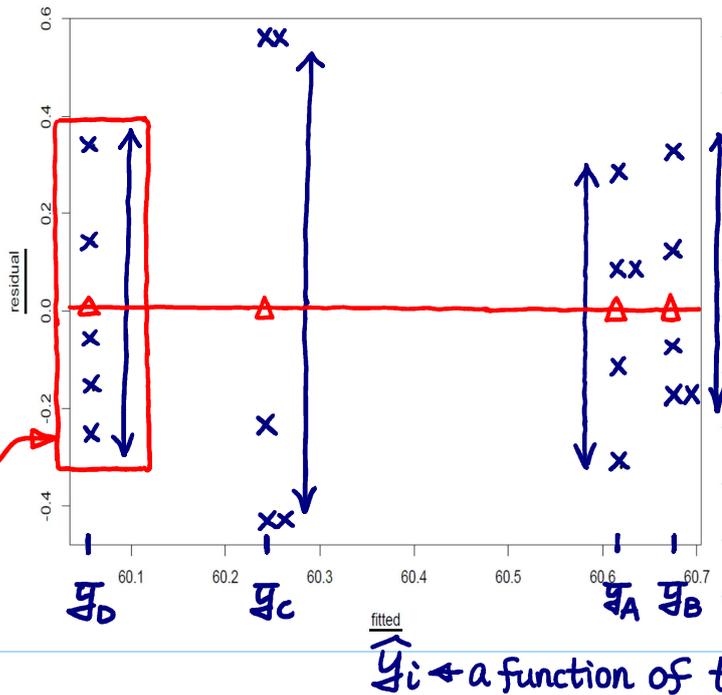
no replicate data



### Plot of $r_i$ vs. $\hat{y}_i$

replicates

allow "pure error" variance estimation



$\hat{y}_i \leftarrow$  a function of treatments A,B,C,D

Figure 1:  $r_i$  vs.  $\hat{y}_i$ , Pulp Experiment

LNp.1

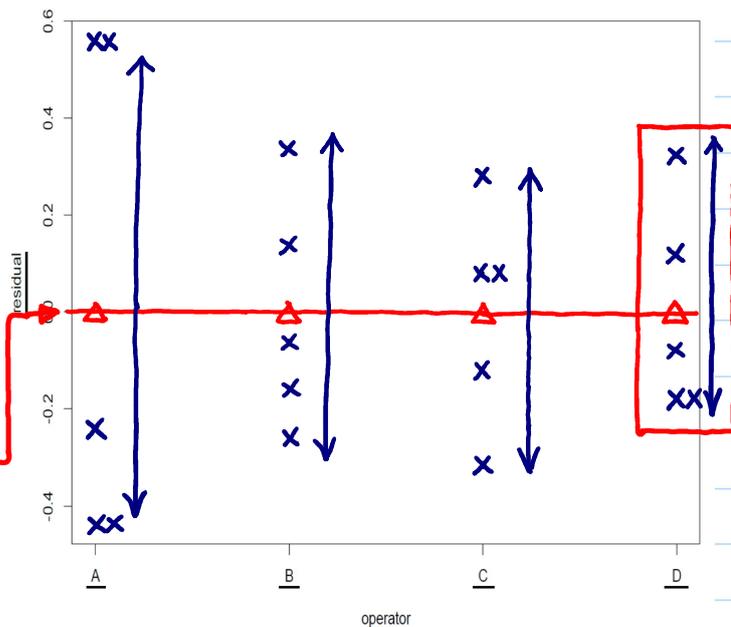


### Plot of $r_i$ vs. treatment

fitted model: saturated

Same 4 groups of residuals as in the residual plot of  $r_i$  vs.  $\hat{y}_i$  (LNp.3-24)

no need to check the trend in the mean of residuals



$$\sum_j \hat{\epsilon}_{ij} = 0 \text{ for any } i$$

$$X = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Figure 2:  $r_i$  vs. treatment, Pulp Experiment

LNp.1

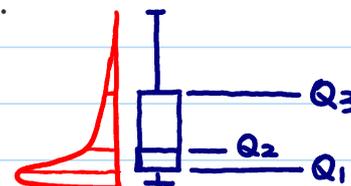
# $Y_1, \dots, Y_n \text{ iid } \sim \text{cdf } F$ ← **Box-(Whisker) Plot**

- A powerful graphical display (due to Tukey) to capture the location, dispersion, skewness and extremity of a distribution. See Figure 3. ↖ LNP.27
- $Q_1 =$  lower quartile (25% quantile),  $Q_3 =$  upper quartile (75% quantile),  $Q_2 =$  median (50% quantile, estimate of location parameter) is the white line in the box.  $Q_1$  and  $Q_3$  are boundaries of the black box.
- $IQR =$  interquartile range (length of box) =  $Q_3 - Q_1$ : measure of dispersion.
- Minimum and maximum of observed values within

$$[Q_1 - 1.5 \times IQR, Q_3 + 1.5 \times IQR]$$

are denoted by two whiskers. Any values outside the whiskers are regarded as extreme values and displayed (possible outliers).

- If  $Q_1$  and  $Q_3$  are not symmetric around the median, it indicates skewness.
- Side-by-side box plots (LNP. 3-2~3) are useful to compare the difference between the distributions of several groups of data.



## Box-(Whisker) Plot

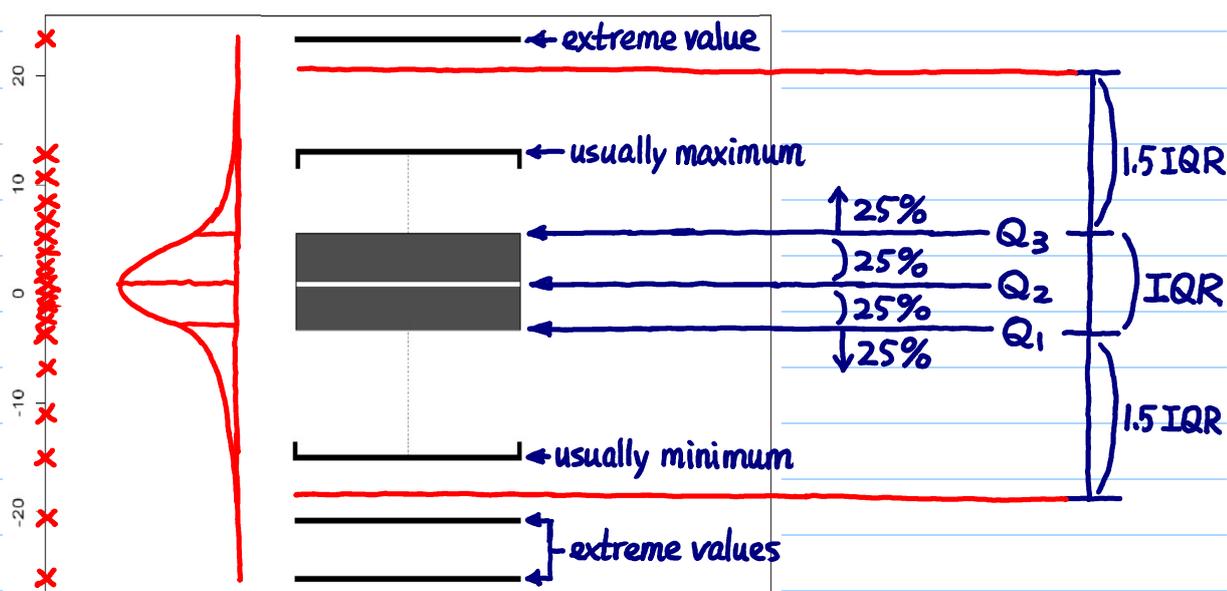


Figure 3: Box-Whisker Plot

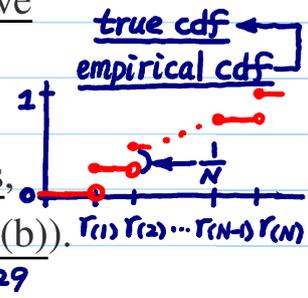
# Normal Probability Plot ← Q-Q plot (LM, LNp.7-15~16)

Original purpose : To test if a distribution is normal, e.g., if the residuals follow a normal distribution (see Figure 5). Q: Why need normality for error?  
 can be used to identify outlier ↗ LNp.30

More powerful use in factorial experiments (discussed in Units 5 and 6).

used to identify significant effects  $\beta_i$ 's ← cf → replace t-tests

- Let  $r_{(1)} \leq \dots \leq r_{(N)}$  be the ordered residuals. The cumulative probability for  $r_{(i)}$  is  $p_i = (i - 0.5)/N$ . Thus the plot of  $p_i$  vs.  $r_{(i)}$  should be S-shaped as in Figure 4(a) if the errors are normal. By transforming the scale of the horizontal axis, the S-shaped curve is straightened to be a line (see Figure 4(b)).



- Normal probability plot of residuals :

$$\left( \Phi^{-1} \left( \frac{i - 0.5}{N} \right), r_{(i)} \right), \quad i = 1, \dots, N, \quad \Phi = \text{normal cdf.}$$

If the errors are normal, it should plot roughly as a straight line. See Figure 5.

↖ LNp.30



# Regular and Normal Probability Plots of Normal

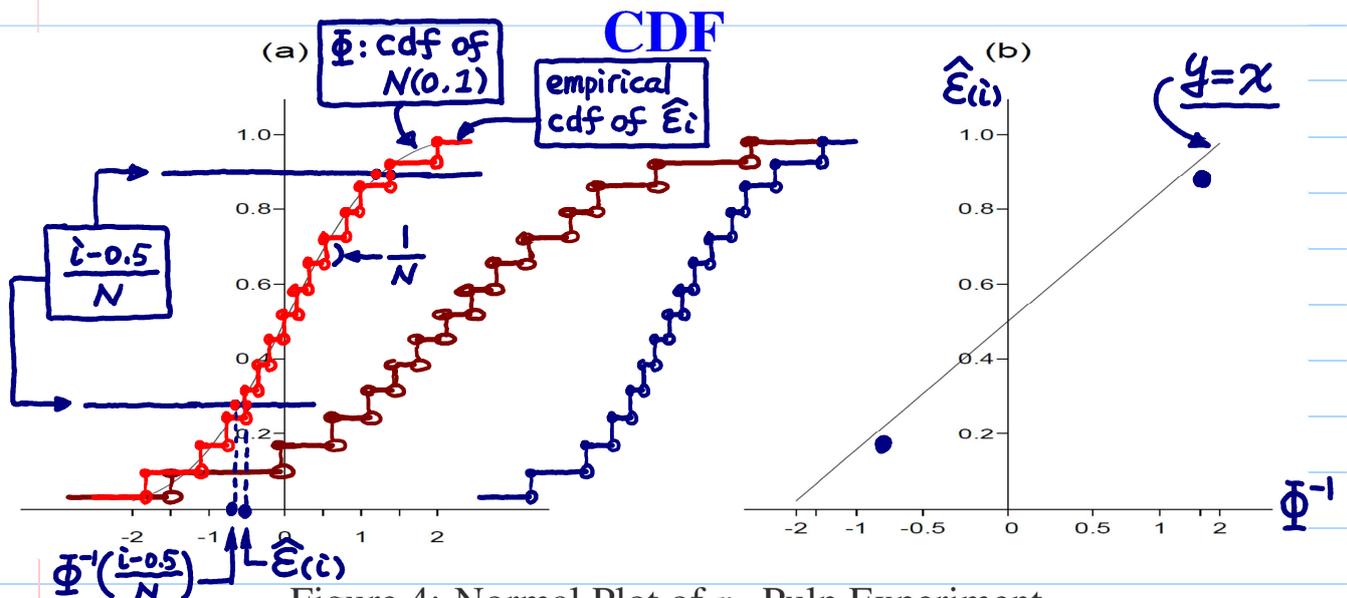


Figure 4: Normal Plot of  $r_i$ , Pulp Experiment

$$Z_1, \dots, Z_N \text{ iid } \sim N(\mu, \sigma^2)$$

$$Z_i = \sigma W_i + \mu \leftarrow W_i = (Z_i - \mu) / \sigma \text{ iid } \sim N(0, 1)$$

Normal probability plot of  $W_i$ 's  $\Rightarrow W_{(i)}$  vs.  $\Phi^{-1}$ :  $y = x$

Normal probability plot of  $Z_i$ 's  $\Rightarrow Z_{(i)}$  vs.  $\Phi^{-1}$ :  $y' = \sigma x + \mu$



# Normal Probability Plot : Pulp Experiment

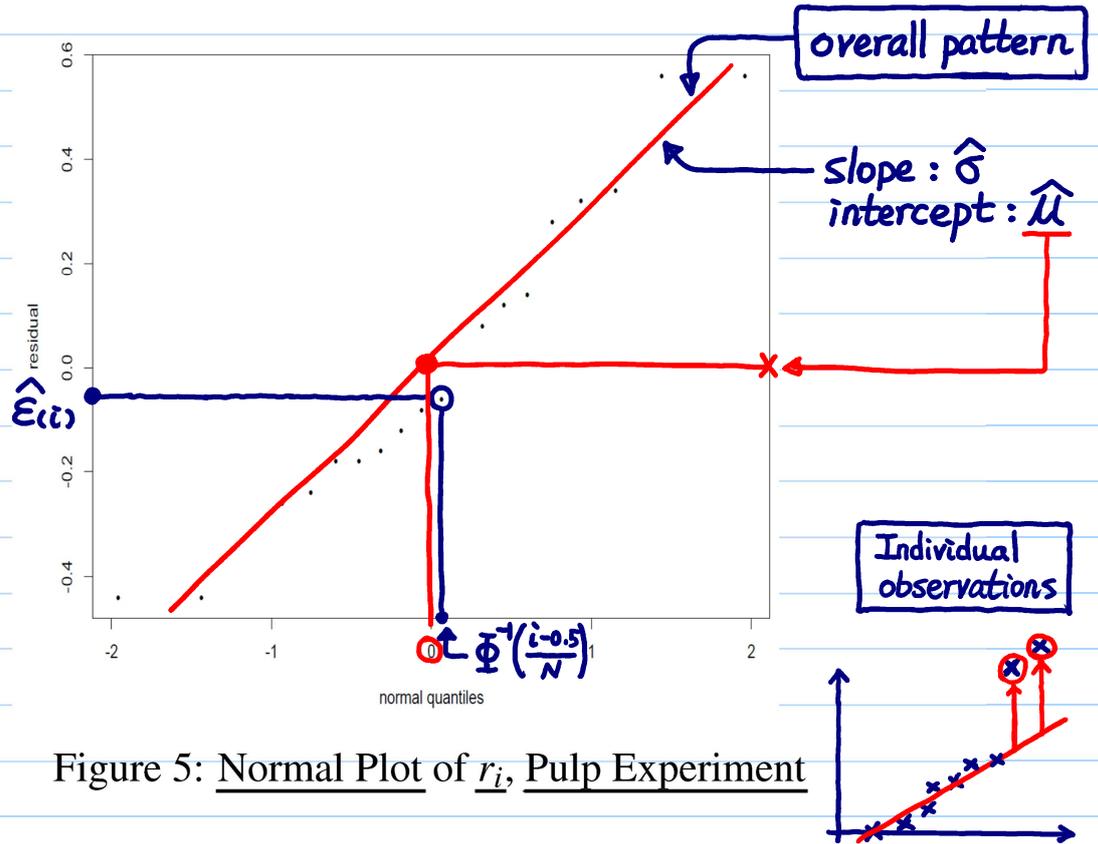


Figure 5: Normal Plot of  $r_i$ , Pulp Experiment

❖ Reading: textbook, 2.6

## Pulp Experiment Revisited

LNp.1

- Compare the 2 scenarios
- (S1) plant has only 4 operators (or only interested in these 4 operators)

LNp.4

effects of operators

$\tau_i$ 's: parameters (unknown fixed values)

after sampling & conditioning

interest: difference btwn the 4 specific  $\tau_i$ 's

cf.

(S2) 4 operators randomly sampled from a large population of operators

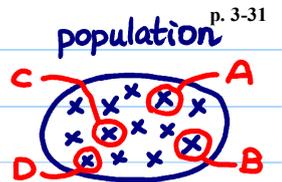
before sampling

$\tau_i$ 's: random variables

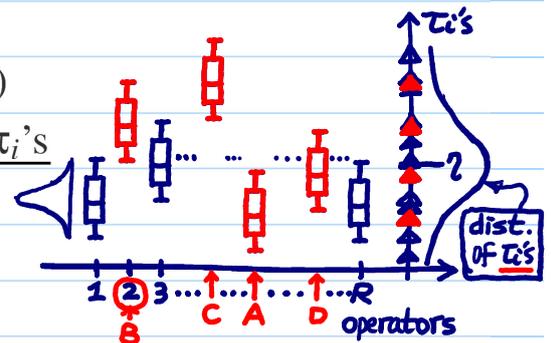
$\mu_i = E(y_{ij})$

interest: difference btwn all operators in this population

Q: Are A,B,C,D a representative sample of the population?



p. 3-31



- In the pulp experiment the effects  $\tau_i$  are called fixed effects because the interest was in comparing the four specific operators in the study. If these four operators were chosen randomly from the population of operators in the plant, the interest would usually be in the variation among all operators in the population. Because the observed data are from operators randomly selected from the population, the variation among operators in the population is referred to as random effects. conditioned on these 4 operators



fixed effect model p. 3-32  
 • One-way random effects model (REM)  $\leftarrow$  cf.  $\rightarrow$  FEM :

FEM in LNp.3-4&3-8  $\leftarrow$  cf.  $\rightarrow$   $y_{ij} = \eta + \tau_i + \epsilon_{ij}$   $\leftarrow$  cf.  $\rightarrow$  whole- & sub-plot errors in LNp.4-47~48

intercept, parameter  $\rightarrow$   $\eta$   $\leftarrow$   $\tau_i$  (r.v.)  $\leftarrow$   $\epsilon_{ij}$

$y_{ij} \sim N(\eta, \sigma_\tau^2 + \sigma^2)$

where  $\epsilon_{ij}$ 's: independent error terms with  $N(0, \sigma^2)$ , Q:  $y_{ij}$ 's indep.?

$\tau_i$ 's: independent  $N(0, \sigma_\tau^2)$ ,  
 and  $\tau_i$  and  $\epsilon_{ij}$  are independent (Why? Give an example.);

REM:  
 parameters  
 $\eta, \sigma_\tau^2, \sigma^2$

$\sigma^2$  and  $\sigma_\tau^2$  are the two variance components of the model. check  $\Sigma^*$

The variance among operators in the population is measured by  $\sigma_\tau^2$ .

$Y = \begin{bmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ \vdots \\ y_{2n_2} \\ \vdots \\ y_{k1} \\ \vdots \\ y_{kn_k} \end{bmatrix}$

REM. (1)  $E(Y) = \eta \mathbf{1}$   
 (2)  $\text{cov}(y_{11}, y_{12}) = \text{cov}(\eta + \tau_1 + \epsilon_{11}, \eta + \tau_1 + \epsilon_{12}) = \sigma_\tau^2$   
 $\text{cov}(y_{11}, y_{21}) = \text{cov}(\eta + \tau_1 + \epsilon_{11}, \eta + \tau_2 + \epsilon_{21}) = 0$

FEM.  $y_{ij}$  indep.  $N(\eta + \tau_i, \sigma^2)$   
 (1)  $E(Y) = X\beta \leftarrow \eta + \tau_i$   
 (2)  $\text{cov}(Y) = \sigma^2 I$   $\leftarrow \chi$

$\text{cov}(Y) = \begin{bmatrix} \square_{n_1 \times n_1} & \circ & \dots \\ \circ & \square_{n_2 \times n_2} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} = \sigma_\tau^2 \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} + \sigma^2 \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}$

$\Sigma^*$

★ relationship btwn  $y_x$  &  $\chi$   
 -  $I_n$  FEM,  $E(Y)$   
 -  $I_n$  REM,  $\text{cov}(Y)$

