

Fitting the Regression Line

- Underlying Model :

model matrix $X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$

$y = \beta_0 + \beta_1 x + \epsilon$, $\epsilon \sim N(0, \sigma^2)$. $\leftrightarrow Y = X\beta + \epsilon$

not regarded as r.v., no measurement error

r.v.: random component

mean structure of y : systematic component

- Coefficients are estimated by minimizing

$\begin{cases} \frac{\partial}{\partial \beta_0} S(\beta_0, \beta_1) = 0 \\ \frac{\partial}{\partial \beta_1} S(\beta_0, \beta_1) = 0 \end{cases} \leftarrow \sum_{i=1}^N (y_i - (\beta_0 + \beta_1 x_i))^2 \equiv S(\beta_0, \beta_1)$

- Least Squares Estimates

Estimated Coefficients :

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

$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$, $var(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$

$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$, $var(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{N} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)$

$\bar{x} = \frac{1}{N} \sum x_i$, $\bar{y} = \frac{1}{N} \sum y_i$.

$Cov(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$

function of x_i 's only irrelevant to β, σ^2

規律: ?%, 隨機: ?%

Explanatory Power of the Model

Q: How well the model explain the data? "goodness of fit" measure

- The total variation in y can be measured by corrected total sum of squares

$CTSS = \sum_{i=1}^N (y_i - \bar{y})^2$

variation in y_i 's

This can be decomposed into two parts (Analysis of Variance (ANOVA)):

source of variation in y_i 's

where $-\hat{y}_i + \hat{y}_i$

$CTSS = \text{RegrSS} + \text{RSS}$

規律 隨機

source of variation: x_i 's \rightarrow $\text{RegrSS} = \text{Regression sum of squares} = \sum_{i=1}^N (\hat{y}_i - \bar{y})^2 = \hat{\beta}_1^2 \sum (x_i - \bar{x})^2$

source of variation: ϵ_i 's \rightarrow $\text{RSS} = \text{Residual sum of squares} = \sum_{i=1}^N (y_i - \hat{y}_i)^2$

$\hat{y}: \text{mean of } \hat{y}_i\text{'s}$

$\hat{\sigma}^2 = \frac{\text{RSS}}{n-p}$

of β

$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ is called the predicted value of y_i at x_i . $\hat{\epsilon}_i$: residuals

- $R^2 = \frac{\text{RegrSS}}{\text{CTSS}} = 1 - \frac{\text{RSS}}{\text{CTSS}}$ measures the proportion of variation in y explained by the fitted model. \leftarrow check the graph in LNp.2-2

ANOVA Table for Simple Linear Regression

↳ decomposition of sum of squares (sources of variation in \hat{y})

ANOVA Table for Simple Linear Regression

Source	Degrees of Freedom	Sum of Squares	Mean Squares
regression	1	$\hat{\beta}_1^2 \sum (x_i - \bar{x})^2$	$\hat{\beta}_1^2 \sum (x_i - \bar{x})^2$
residual	$N - 2$	$\sum_{i=1}^N (y_i - \hat{y}_i)^2$	$\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{(N-2)}$
total (corrected)	$N - 1$	$\sum_{i=1}^N (y_i - \bar{y})^2$	

source of variation: χ_i 's → regression
 source of variation: ϵ_i 's → residual
 RegrSS (Sum of Squares for regression)
 RSS (Residual Sum of Squares)
 CTSS (Corrected Total Sum of Squares) ← RSS under H_0
 RSS under H_A
 R^2 (Coefficient of Determination)
 $F = \frac{\text{MSRegrSS}}{\text{MSRSS}} \sim F_{df_{\text{RegrSS}}, df_{\text{RSS}}}$

overall F-test
 $H_0: y = \beta_0 + \epsilon$
 $H_A: y = \beta_0 + \beta_1 x + \epsilon$

ANOVA Table for Breast Cancer Example

Source	Degrees of Freedom	Sum of Squares	Mean Squares
regression	1	2599.53	2599.53
residual	14	796.91	56.92
total (corrected)	15 (16 obs.)	3396.44	

$F = \frac{\text{MSRegrSS}}{\text{MSRSS}} \sim F_{df_{\text{RegrSS}}, df_{\text{RSS}}}$
 $\sim F_{1, 14}$

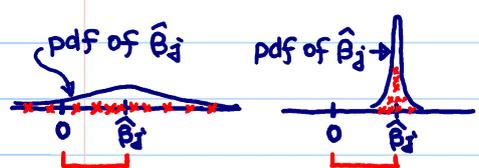
equivalent to the overall F-test in LN p.2-5 for simple regression

t-Statistic

null model = ?
 alternative model = ?

why interested in it?

- To test the null hypothesis $H_0: \beta_j = 0$ against the alternative hypothesis $H_A: \beta_j \neq 0$ under the full model, use the test statistic



$$t_j^2 = F$$

$$t_j = \frac{\hat{\beta}_j - 0}{s.d.(\hat{\beta}_j)}$$

$$s.d.(\hat{\beta}_j) = \hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}}$$

can be changed to other constant a for testing $H_0: \beta_j = a$

i.e., more confident to reject $H_0: \beta_j = 0$

- The higher the value of $|t_j|$, the more significant is the coefficient.
- For 2-sided alternatives, $p\text{-value} = \text{Prob}(|t_{df}| > |t_{obs}|)$, $df = \text{degrees of freedom}$ for the t-statistic, $t_{obs} = \text{observed value of the t-statistic}$. If p-value is very small, then either we have observed something which rarely happens, or H_0 is not true. In practice, if p-value is less than $\alpha = 0.05$ or 0.01 , H_0 is rejected at level α .

under H_0

tests \longleftrightarrow **Confidence Interval**: collection of plausible β 's

$100(1 - \alpha)\%$ confidence interval for β_j is given by

$$\hat{\beta}_j \pm t_{N-2, \frac{\alpha}{2}} \times s.d.(\hat{\beta}_j)$$

critical value

where $t_{N-2, \frac{\alpha}{2}}$ is the upper $\frac{\alpha}{2}$ point of the t distribution with $N - 2$ degrees of freedom.

\therefore 2-sided

df_{RSS}

If the confidence interval for β_j does not contain 0, then H_0 is rejected.

$$\beta_j = 0$$

Test $H_0: \beta_j = \beta_j^*$ vs. $H_A: \beta_j \neq \beta_j^*$

Acceptance region: $|t_j| = \left| \frac{\hat{\beta}_j - \beta_j^*}{s.d.(\hat{\beta}_j)} \right| \leq t_{df_{RSS}, \frac{\alpha}{2}}$

$$\hat{\beta}_j - t \times s.d.(\hat{\beta}_j) \leq \beta_j^* \leq \hat{\beta}_j + t \times s.d.(\hat{\beta}_j)$$

estimate \pm (critical value) \times s.d.(estimate)

Predicted Values and Residuals

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

$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ is the predicted value of y_i at x_i .

$r_i = y_i - \hat{y}_i$ is the corresponding residual.

or denoted by $\hat{\epsilon}$

Standardized residuals are defined as $\frac{r_i}{s.d.(r_i)}$

Studentized

Plots of residuals are extremely useful to judge the "goodness" of fitted model.

under ② & ③, $\hat{\epsilon}$ contains more information

under ① than error distribution.

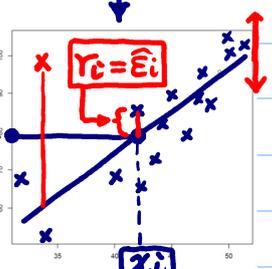
Normal probability plot (will be explained in Unit 3).

Residuals versus predicted values.

Residuals versus covariate x .

under ② & ③

$$y = \underbrace{f(x_1, \dots, x_m)}_{\hat{y}} + \underbrace{\epsilon}_{\hat{\epsilon}}$$



$$\hat{\sigma}^2 = \frac{RSS}{n-p}$$

used in

$$cov(\hat{\epsilon}) = (I - H)\sigma^2$$

hat matrix \hat{H}
 H_{ii} : leverage

future lecture
 LNp.3-22
 ~3-25

- ① If model correct, $\hat{\epsilon}$: carry information about 隨機 ϵ (error dist., error variance)
- \hat{y} : carry information about 規律 f
- ② overfitting: $\epsilon \rightarrow \hat{y}$
- ③ lack of fit: $f \rightarrow \hat{\epsilon}$

Analysis of Breast Cancer Data

The regression equation is

$$M = -21.79 + 2.36 T$$

Predictor	Coef	SE	Coef	T	P
Constant	$\hat{\beta}_0 \rightarrow -21.79$	$s.d.(\hat{\beta}_0) \rightarrow 15.67$		-1.39	0.186
T	$\hat{\beta}_1 \rightarrow 2.3577$	$s.d.(\hat{\beta}_1) \rightarrow 0.3489$		6.76	0.000

S = 7.54466 $\leftarrow \hat{\sigma}$ R-Sq = 76.5% R-Sq(adj) = 74.9%

Analysis of Variance \leftarrow overall F-test

Source	DF	SS	MS	F	P
Regression	1	2599.5	2599.5	45.67	0.000
Residual Error	14	796.9	56.9		
Total	15	3396.4			

$6.76^2 = 45.67$
(only for simple regression)

Unusual Observations

Obs	T	M	Fit	SE Fit	Residual	St Resid
15	31.8	67.30	53.18	4.85	14.12	2.44RX

R denotes an observation with a large standardized residual.

X denotes an observation whose X value gives it large leverage.