

## 12.1 The simple linear Regression Model

### 1. A linear probabilistic model

For any fixed value of the independent variable  $x$ ,  
the dependent variable

$$Y = \beta_0 + \beta_1 x + \varepsilon$$

$\varepsilon$  is a RV, assume randomly distributed with

$$E(\varepsilon) = 0, V(\varepsilon) = \sigma^2$$

## 12.2. Estimating Model Parameters

### 1. least square

$$\hat{y}_i = b_0 + b_1 x_i$$

$$b_1 = \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

$$S_{xy} = \sum x_i y_i - (\sum x_i)(\sum y_i)/n.$$

$$S_{xx} = \sum x_i^2 - (\sum x_i)^2/n$$

$$b_0 = \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

### 2. Estimating $\sigma^2$ and $\sigma$

#### - Error sum of square

$$SSE = \sum [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

#### Estimate of $\sigma^2$

$$\hat{\sigma}^2 = S^2 = \frac{SSE}{n-2} = \frac{\sum (y_i - \hat{y})^2}{n-2}$$

$$SSE = S_{yy} - \text{slope} \cdot S_{xy}$$

- The coefficient of determination.

$$r^2 = 1 - \frac{SSE}{SST}$$

It's interpreted as the proportion of observed variance that can be explained by the model

### 12.3. Inference About the Slope Parameter $\beta_1$

- Proposition:

1. The mean value of  $\hat{\beta}_1$  is  $E(\hat{\beta}_1) = \mu_{\hat{\beta}_1} = \hat{\beta}_1$

2. Variance and standard deviation of  $\hat{\beta}_1$

$$V(\hat{\beta}_1) = \sigma_{\hat{\beta}_1}^2 = \frac{\sigma^2}{S_{xx}}$$

$$\sigma_{\hat{\beta}_1} = \sqrt{\frac{\sigma^2}{S_{xx}}} \quad S_{xx} = \sum (x_i - \bar{x})^2$$

3. Estimate  $\hat{\beta}_1$  has normal distribution.

- Theorem

$$1. T = \frac{\hat{\beta}_1 - \beta_1}{S/\sqrt{S_{xx}}} = \frac{\hat{\beta}_1 - \beta_1}{S\hat{\beta}_1}$$

has  $t$  distribution with  $df = n-2$

- Confidence Interval for  $\beta_1$

• A  $100(1-\alpha)\%$  CI for  $\beta_1$

$$\hat{\beta}_1 \pm t_{\alpha/2, n-2} \cdot S\hat{\beta}_1$$

slope

of line.

- Model utility test

1. Null hypothesis:  $H_0: \beta_1 = \beta_{10}$

Model utility test

$$H_0: \beta_1 = 0$$

$$H_a: \beta_1 \neq 0$$

## 2. Test statistic

$$t = \frac{\hat{\beta}_1 - \beta_{10}}{S\hat{\beta}_1} \quad df = n-2$$

$$t = \frac{\hat{\beta}_1}{S\hat{\beta}_1}$$

## 3. Alternative Hypothesis

$$H_a: \beta_1 > \beta_{10}$$

P-value

Right of  $t_{n-2}$ .

$$H_a: \beta_1 < \beta_{10}$$

Left of  $t_{n-2}$ .

$$H_a: \beta_1 \neq \beta_{10}$$

2 · right of  $|t_{n-2}|$

## 12.4. Inference Concerning $\mu_{Y|x^*}$ and the

### 1. Prediction of future $Y$ values.

#### - Proposition

- Let  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x^*$ , where  $x^*$  is some fixed value of  $x$ .
- The mean value of  $\hat{Y}$  is.

$$E(\hat{Y}) = E(\hat{\beta}_0 + \hat{\beta}_1 x^*)$$

$$= (\hat{\beta}_0 + \hat{\beta}_1 x^*) = \beta_0 + \beta_1 x^*.$$

Thus,  $\hat{\beta}_0 + \hat{\beta}_1 x^*$  is an unbiased estimate for  $\beta_0 + \beta_1 x^*$  ( $\mu_{Y|x^*}$ ).

#### - The variance of $\hat{Y}$

$$V(Y) = \sigma^2 \left[ \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right]$$

$$S_{\hat{Y}} = S \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$$

#### - $\hat{Y}$ has a normal distribution.

## 2. Inference Concerning $\mu_{Y|x^*}$

- The variable  $T$

$$T = \frac{\bar{Y} - (\beta_0 + \beta_1 x^*)}{S_{\bar{Y}}}$$

has  $t$ -distribution with  $n-2$  df.

- A  $100(1-\alpha)\%$  CI for  $\mu_{y|x^*}$ ,

the expected value of  $Y$  when  $x=x^*$ , is

$$\hat{y} \pm t_{\alpha/2, n-2} \cdot S_{\bar{Y}}$$

### 3. A Prediction Interval for a future Value of $Y$ .

A  $100(1-\alpha)\%$  PI for a future  $Y$  observation to be made when  $x=x^*$  is

$$\begin{aligned} & \hat{\beta}_0 + \hat{\beta}_1 x^* \pm t_{\alpha/2, n-2} \cdot S \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}} \\ &= \hat{y} \pm t_{\alpha/2, n-2} \cdot \sqrt{S^2 + S_{\bar{Y}}^2} \end{aligned}$$

## 12.5 Correlation.

$$\begin{aligned} S_{xy} &= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \cdot \sum_{i=1}^n y_i}{n} \end{aligned}$$

## • Correlation

$$r = \frac{S_{xy}}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

$$= \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}}$$

- Properties of  $r$

- the value of  $r$  is independent of unit
- $-1 \leq r \leq 1$
- $r=1$  iff all points lie on a straight line with positive slope.

## 2. Absence of Correlation.

- Let  $R$  denote sample correlation coefficient.
- When  $H_0: \rho = 0$  is true.

$$T = \frac{R\sqrt{n-2}}{\sqrt{1-R^2}} \text{ with } n-2 \text{ df.}$$

After Hypothesis.

$$H_a: \rho > 0$$

$$H_a: \rho < 0$$

$$H_a: \rho \neq 0$$

P-value

Right of  $T$  curve

left of  $T$  curve

2. (Right of  $|t|$ )

## 3. When $(x_1, y_1), \dots, (x_n, y_n)$ is a sample from a bivariate normal distribution.

$$V = \frac{1}{2} \ln \left( \frac{1+R}{1-R} \right).$$

has approximately a normal distribution.

$$\mu_V = \frac{1}{2} \ln \left( \frac{1+\rho}{1-\rho} \right) \quad \sigma_V^2 = \frac{1}{n-3}.$$

## 4. A $100(1-\alpha)\%$ CI for $\rho$ is

$$\left( \frac{e^{2C_1}-1}{e^{2C_1}+1}, \frac{e^{2C_2}-1}{e^{2C_2}+1} \right).$$

where  $C_1, C_2$  are left, right endpoints.

- The most of softwares provide regression output such as least squares estimates,  $r^2$ , regression sum of squares and other information upon request.

use  $\hat{y}$  to write fitted line.

The regression equation is  
cet num = 75.2 - 0.209 iod val

Predictor	Coef	$\hat{\beta}_0$	$\hat{\beta}_1$	SE Coef	$s_{\hat{\beta}_1}$	T	P
Constant	75.212			2.984		25.21	0.000
iod val	-0.20939		0.03109		6.73		0.000
	$s = 2.56450$	$R-sq = 79.1\%$	$100r^2$	$R-sq(\text{adj}) = 77.3\%$			
Analysis of Variance							
SOURCE	DF	SS	MS	F	P		
Regression	1	298.25	298.25	45.35	0.000		
Error	12	78.92	6.58				
Total	13	377.17					

- The  $F$  test gives exactly the same result as the model utility  $t$  test because  $t^2 = f$  and  $t_{\alpha/2, n-2}^2 = F_{\alpha, 1, n-2}$ .

Source of Variation	df	Sum of Squares	Mean Square	f
Regression	1	SSR	SSR	$\frac{\text{SSR}}{\text{SSE}/(n-2)}$
Error	$n-2$	SSE	$s^2 = \frac{\text{SSE}}{n-2}$	
Total	$n-1$	SST		

$$r^2 = 1 - \frac{SSE}{SST}$$

Assumption of Regression line.

① Linearity between  $x$  and  $y$ .

Check by scatter plot of  $x$  and  $y$ .

②  $y$ 's have normal distribution  $N(\mu_{xi}, \sigma^2)$

Normal Probability Plot.

③ The variance of the residual is the same for all  $x$ 's.  
By scatter residual vs. fit.

