

$E(\hat{\beta}_0) = \beta_0$ $E(\hat{\beta}_1) = \beta_1$
 $\text{Var}(\hat{\beta}_0) = ?$ $\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$

* express
 $\hat{\beta}_1 = \sum a_i y_i$
 depends on x

$$\hat{\sigma}^2 = \frac{\sum \hat{e}_i^2}{n-2}$$

$$E(\hat{\sigma}^2) = \sigma^2$$

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

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$



Analysis of variance (ANOVA)

- Regression is the study of dependence of variables

- $y_i = \beta_0 + \beta_1 x_i + e_i$ ✓ *not*
- $\beta_1 = 0 \rightarrow x$ and y are dependent
- $\beta_1 \neq 0 \rightarrow x$ and y are *not* dependent

- Question:

- Are x and y dependent?

- Answer:

- Method 1) test whether $\beta_1 = 0$
- Method 2) Compare the two models

Model 1
 $\beta_1 = 0$

$$\rightarrow \bullet E(y|x) = \beta_0 \quad \text{i.e. } y_i = \beta_0 + e_i$$

$$\bullet E(y|x) = \beta_0 + \beta_1 x \quad \text{i.e. } y_i = \beta_0 + \beta_1 x_i + e_i$$

← Model 2 : β_1 free



Analysis of variance (ANOVA)

- Analysis of variance (ANOVA) is a method that compares two models of mean functions

- $E(y|x) = \beta_0$
- $E(y|x) = \beta_0 + \beta_1 x$

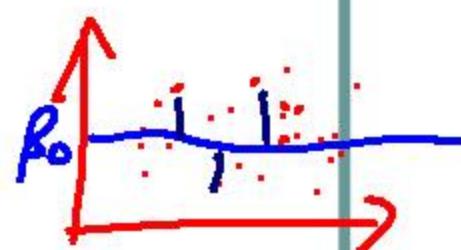
- For the first model: $E(y|x) = \beta_0$

- β_0 can be estimated by minimizing $\sum(y_i - \beta_0)^2$

- differentiate w.r.t. β_0 gives $\hat{\beta}_0 = \bar{y}$

- Residual sum of square RSS is

$$\sum(y_i - \hat{\beta}_0)^2 = \sum(y_i - \bar{y})^2 = SYY$$



$$\leftarrow RSS(\beta_0)$$

$$\frac{\partial RSS}{\partial \beta_0} = -2 \sum (y_i - \beta_0) = 0$$

$$\Rightarrow \hat{\beta}_0 = \frac{\sum y_i}{n} = \bar{y}$$



- Compare

- $E(y|x) = \beta_0$
- $E(y|x) = \beta_0 + \beta_1 x$

$$\text{SXX} = \sum (x_i - \bar{x})^2$$

- For the first model: $E(y|x) = \beta_0$

- Residual sum of square, RSS_1 , is

$$\sum (y_i - \hat{\beta}_0)^2 = \sum (y_i - \bar{y})^2 = \text{SYY}$$

- For the second model: $E(y|x) = \beta_0 + \beta_1 x$

- Residual sum of square, RSS_2 , is

$$\sum_{i=1}^n [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2 = \text{SYY} - \frac{(\text{SXY})^2}{\text{SXX}}$$

- $\text{RSS}_1 > \text{RSS}_2 \dots$ Is the 2nd model always better?



Analysis of variance (ANOVA)

- Difference sum of square due to regression (SSreg)
 - $SS_{\text{reg}} = \underline{\text{RSS}_1} - \underline{\text{RSS}_2} = \frac{(SXY)^2}{SXX} = \left[\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{SXX} \right]^2$
 - Large SSreg → 2nd model explains much more variation
 - How large is large?
- Study the distribution of SSreg under model 1 (idea)
 - After some algebra, $SS_{\text{reg}} = [\sum (\frac{x_i - \bar{x}}{SXX}) y_i]^2$
 - By CLT, $\sum (\frac{x_i - \bar{x}}{SXX}) y_i$ is approximately $N(0, \sigma^2)$
 - $SS_{\text{reg}} \sim \sigma^2 (N(0, 1))^2 = \sigma^2 X_1^2$ (Chi-square with d.f. 1)
 - $\hat{\sigma}^2 = \frac{\sum \hat{e}_i^2}{n-2} \sim \sigma^2 X_{n-2}^2 / (n-2)$ (Chi-square with d.f. n-2)
 - $SS_{\text{reg}} / \hat{\sigma}^2 \sim X_1^2 / (X_{n-2}^2 / (n-2)) = F(1, n-2)$

• ANOVA table: a break-down of squares (variation)

$$SYY = SS_{reg} + RSS$$

Source	df	SS	MS	F	p-value
Regression	1	SS_{reg}	$SS_{reg}/1$	$MS_{reg}/\hat{\sigma}^2$	
Residual	$n - 2$	RSS	$\hat{\sigma}^2 = RSS/(n - 2)$		↑ checking $F(1, n-2)$
Total	$n - 1$	SYY			

$$\sum_{i=1}^n [y_i - \bar{y}]^2 = \sum_{i=1}^n [y_i - \hat{y}_i]^2 + \sum_{i=1}^n [\hat{y}_i - \bar{y}]^2$$

$$TSS = SYY = RSS + SS_{reg}$$

Variation of the data

Variation not explained by regression

Variation explained by regression