

# Sum of Squares Decomposition & the F-test

An alternative way to test whether a predictor  $X$  is a "significant predictor" of  $Y$  is through something called the "sum of squares decomposition".

The decomposition provides a breakdown of the total variation in  $Y$  into 2 pts:

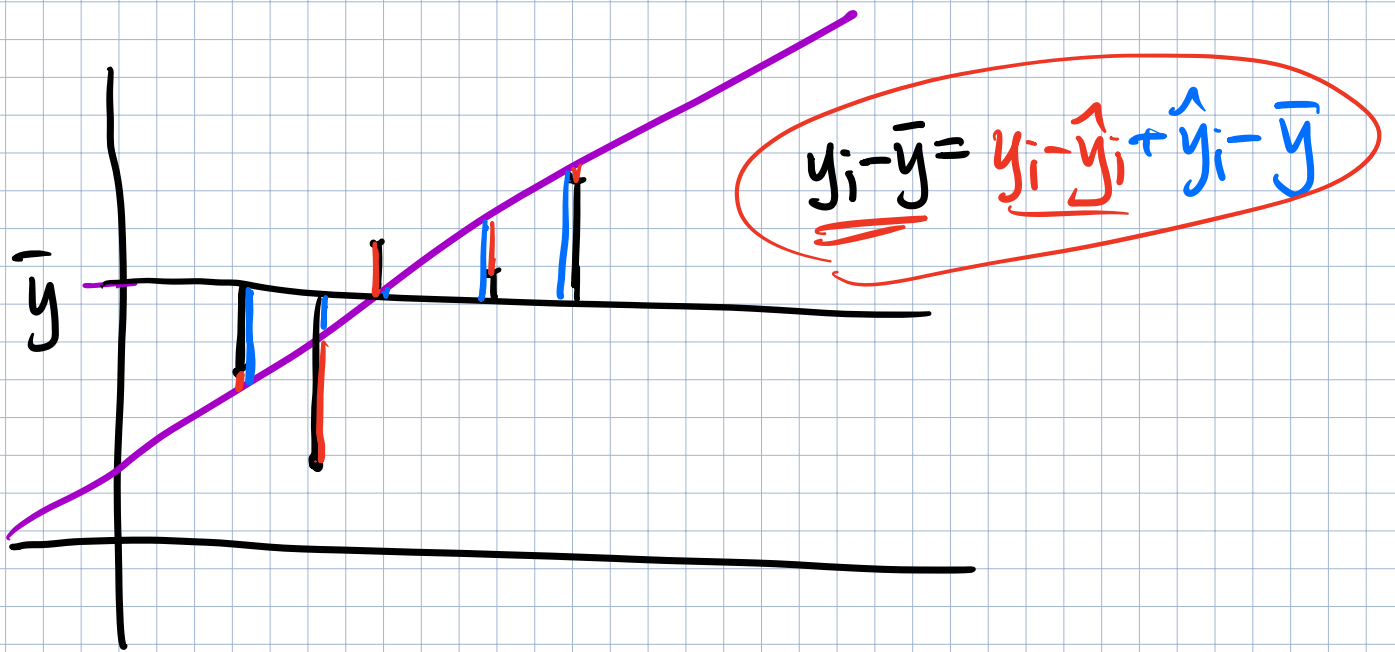
① the sum of sq. errors (SSE)

② the regression sum of squares (SSR)

↓

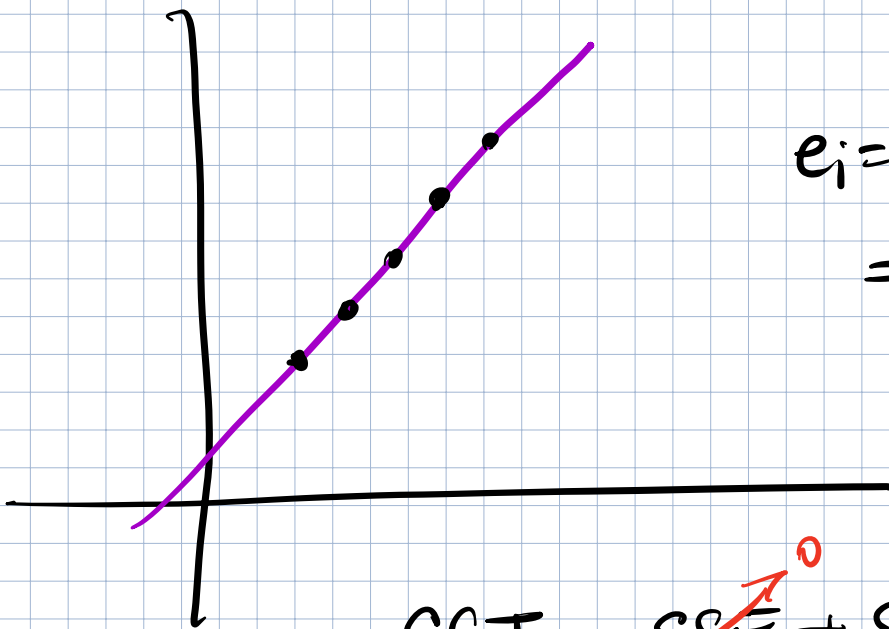
$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{SSE} + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{SSR}$$

"total sum of squares"



To build our intuition let's look @ some extreme cases:

① our predictions are perfect:  $y_i = \hat{y}_i \forall i$



$$e_i = y_i - \hat{y}_i = 0 \quad \forall i$$

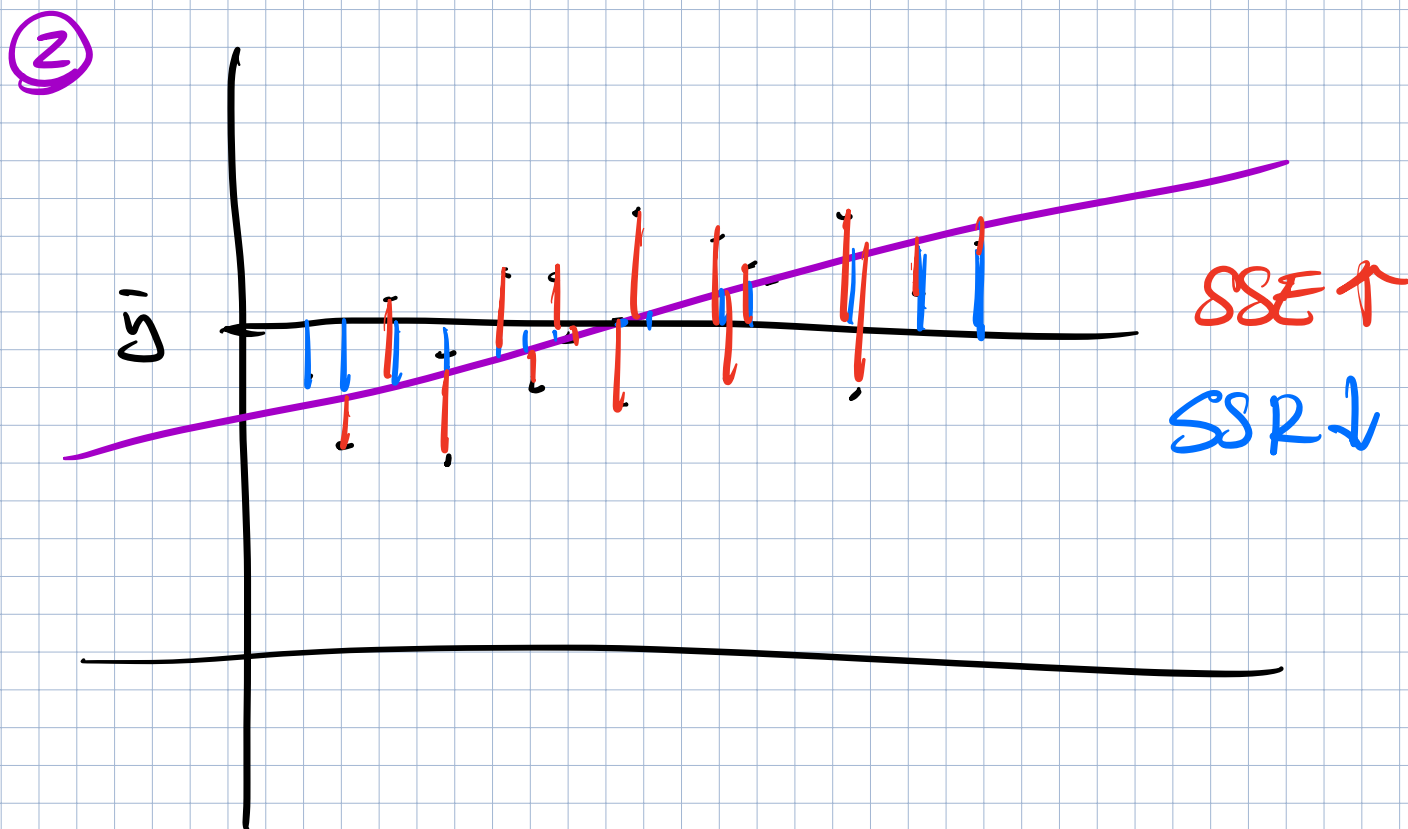
$$\Rightarrow \sum_i e_i^2 = 0$$

$$SSE = 0$$

$$SST = \cancel{SSE}^0 + SSR$$

$$SST = SSR$$

Intuition  $\Rightarrow$  reject  $H_0 \Rightarrow$  the slope is not 0  
 $\Rightarrow$   $X$  is a useful predictor of  $y$ .



If SSR is low & SSE is high

$\Rightarrow$  fail to reject  $H_0$

$\Rightarrow$   $X$  is not useful in predicting  $y$ .

$$e_i = y_i - \hat{y}_i$$

$$y_i = \hat{y}_i + e$$

theoretical version

$$\text{Var}(y) = \text{Var}(\hat{y}) + \text{Var}(e)$$

$$\sum_i (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n e_i^2$$

$$\text{SST} = \text{SSR} + \text{SSE}$$

need to prove...

Now let's explicitly define a statistic:

$$F = \frac{\text{SSR} / 1}{\text{SSE} / \underline{\underline{n-2}}} \sim F_{1, n-2}$$

Define an  $F$  RV:

Suppose ①  $U \sim \chi_n^2$

②  $V \sim \chi_m^2$

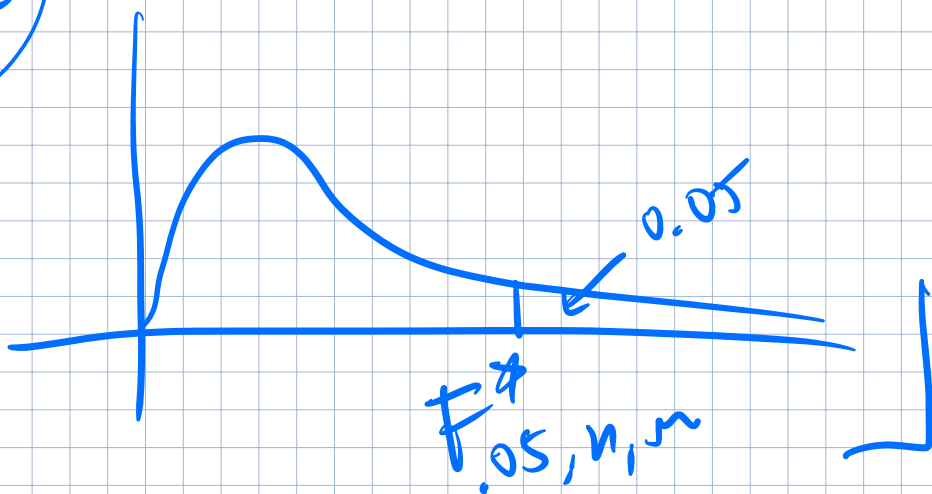
③  $U \perp V$

$\perp$  = "independent"

$\perp$

then:

$$F = \frac{U/n}{V/m} \sim F_{n,m}$$



To show that  $F = \frac{MSR}{MSE}$

$$F = \frac{SSR/1}{SSE/n-2} \sim F_{1, n-2} \text{ we}$$

need to know:

$p = \# \text{ of betas}$

$$\frac{SSR}{\sigma^2} \sim \chi^2_{\underline{p-1}} = p-1$$

$$\frac{SSE}{\sigma^2} \sim \chi^2_{\underline{n-2}} = n-p$$

$$\frac{SST}{\sigma^2} \sim \chi^2_{n-1}$$

$\downarrow$   $SSR \perp SSE$

2 params est  $\downarrow$

$$SSE = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad n-2 \text{ df}$$

$$\sum_{i=1}^n (y_i - \underbrace{\beta_0}_{\theta} - \beta_1 x_i)^2 \sim \chi^2_n$$

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

one parameter  
est  
↓  
n-1 df

$$\frac{SST}{\sigma^2} = \frac{SSE}{\sigma^2} + \frac{SSR}{\sigma^2}$$

$$\chi^2_{n-1} = \chi^2_{n-2} + \chi^2_1$$

Decision Rule:

$$\text{If } F > F_{1-\alpha, 1, n-2}^*$$

⇒ reject  $H_0$

⇒  $X$  is a sig. pred. of  $y$ .

$$\sum_i (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

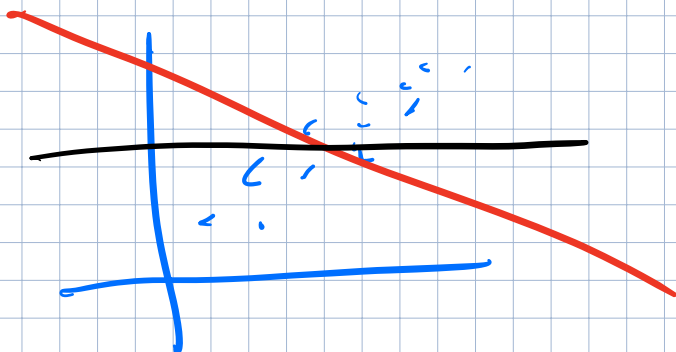
$$= \underbrace{\sum_i (y_i - \hat{y}_i)^2}_{SSE} + \underbrace{2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})}_{=0}$$

$$+ \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{SSR}$$

Coef. of Determination

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

① For the LSRL,  $0 \leq R^2 \leq 1$





② For SLR

it turns out that

$$R^2 = r^2 \quad \text{where } r \text{ is}$$

the sample correlation

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$