

# Sum of Squares Decomposition & the F-test

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

The decomposition gives us a breakdown of the total variance of  $\{Y_i\}_{i=1}^n$  into 2 parts:  $\sim SST$

- ① the sum of squared errors (SSE)
- ② the sum of squares regression

(SSR)

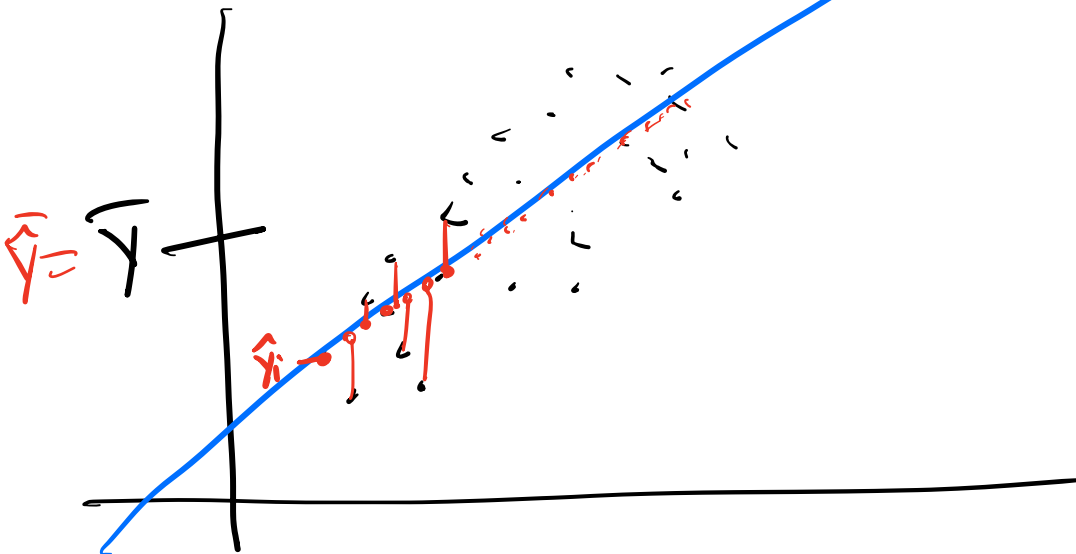
Mathematically:

Total sum of squares in  $Y =$  Regression sum of sq. +  
Sum of squared error

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

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

$\uparrow$   
 $\bar{y}_i$

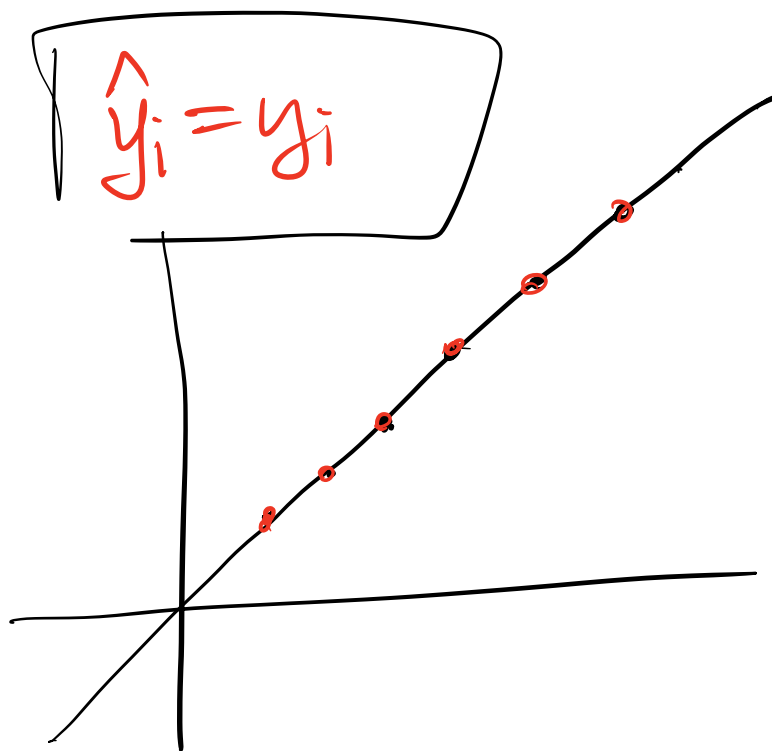


let's construct a statistic which tests whether  $x$  is a significant predictor of  $y$ :

$$H_0: \beta_1 = 0 \quad \text{vs} \quad H_1: \beta_1 \neq 0.$$

To build intuition think about the extremes:

① our predictions are perfect



0

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

$\uparrow$   
 $\hat{y}_i$

If  $\text{SSE} = 0$

$$\text{SSR} = \text{SST} \Rightarrow$$

Intuition:

X predicts  
Y well

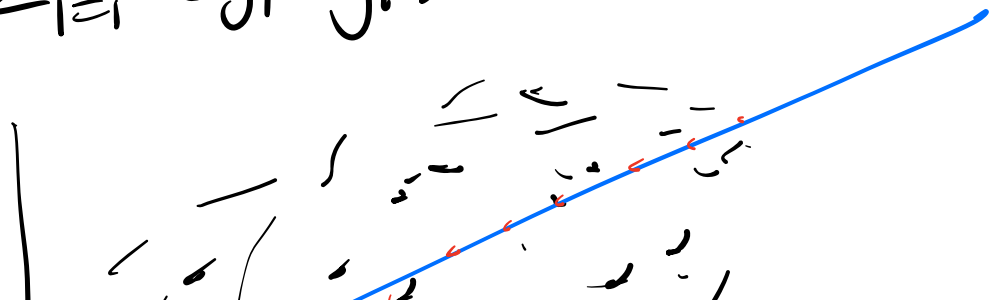


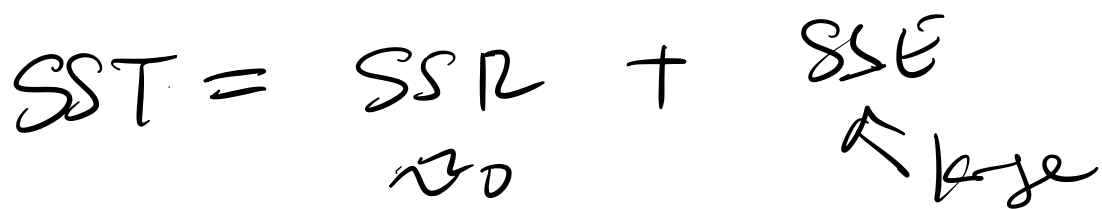
Reject  $H_0$  in  
the t-test.

Other Extreme:

② If X was a really bad  
predictor of Y.....

$$\Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \text{SSE} \uparrow \uparrow \uparrow$$





$\Rightarrow$  t-test would fail to reject

theoretical:

$$\text{Var}(y_i) = \text{Var}(\hat{y}_i) + \text{Var}(e_i)$$

c.f. sample:

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

Exercise!

Agree that  $SST = SSR + SSE$ .

WTS:

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

Step 1

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\underbrace{e_i}_{\text{circled}} + \underbrace{d_i})^2$$

$$= \sum_{i=1}^n (e_i^2 + 2e_i d_i + d_i^2)$$

$$= \underbrace{\sum_{i=1}^n e_i^2} + \underbrace{2 \sum_{i=1}^n e_i d_i} + \underbrace{\sum_{i=1}^n d_i^2}$$

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

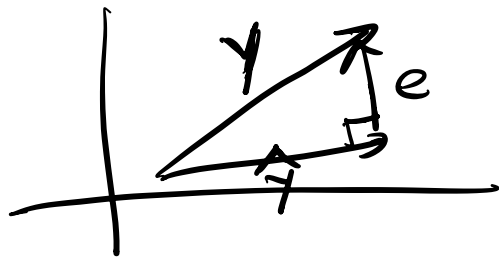
$$= \sum_{i=1}^n (e_i \hat{y}_i - e_i \bar{y})$$

$$\begin{aligned}
 &= \sum_{i=1}^n e_i \hat{y}_i - \sum_{i=1}^n e_i \bar{y} \\
 &\searrow = \underline{\underline{\sum_{i=1}^n e_i \hat{y}_i}} - \bar{y} \sum_{i=1}^n e_i = 0 \\
 &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2
 \end{aligned}$$

Geometric Concept

$Y = (y_1, \dots, y_n)$  ← observed data

$$Y = \underset{\substack{\parallel \\ (\hat{y}_1, \dots, \hat{y}_n)}}{\hat{Y}} + \underset{\substack{\parallel \\ (e_1, \dots, e_n)}}{e} \Rightarrow$$



I will tell you that

$$\textcircled{1} \quad \frac{SSE}{\sigma^2} \sim \chi_{n-2}^2,$$

$$\textcircled{2} \quad \frac{SSR}{\sigma^2} \sim \chi_1^2, \quad \&$$

$$\textcircled{3} \quad \frac{SST}{\sigma^2} \sim \chi_{n-1}^2$$

To construct a statistic which follows an  $F$ -dist, we can take:

$$F = \frac{SSR/1}{SSE/n-2} = \frac{MSR}{MSE}$$

---



## F Table:

SS	notation	df	mean square
Reg.	SSR	1	$MSR = SSR/1$
Err.	SSE	$n-2$	$MSE = SSE/n-2$
Tot.	SST	$n-1$	

$$F = \frac{MSR}{MSE} \overset{tho}{\sim} F_{1, n-2}$$

high F-val  $\Rightarrow$  reject  $H_0$

### Decision Rule:

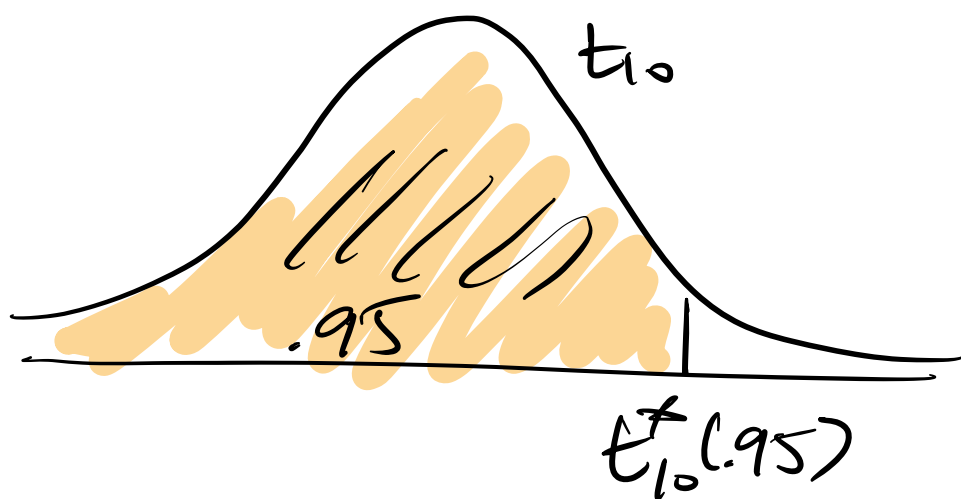
If  $F > F_{1, n-2}^*(1-\alpha)$  then reject  $H_0$ .

$\Rightarrow$  X is a significant predictor of Y.

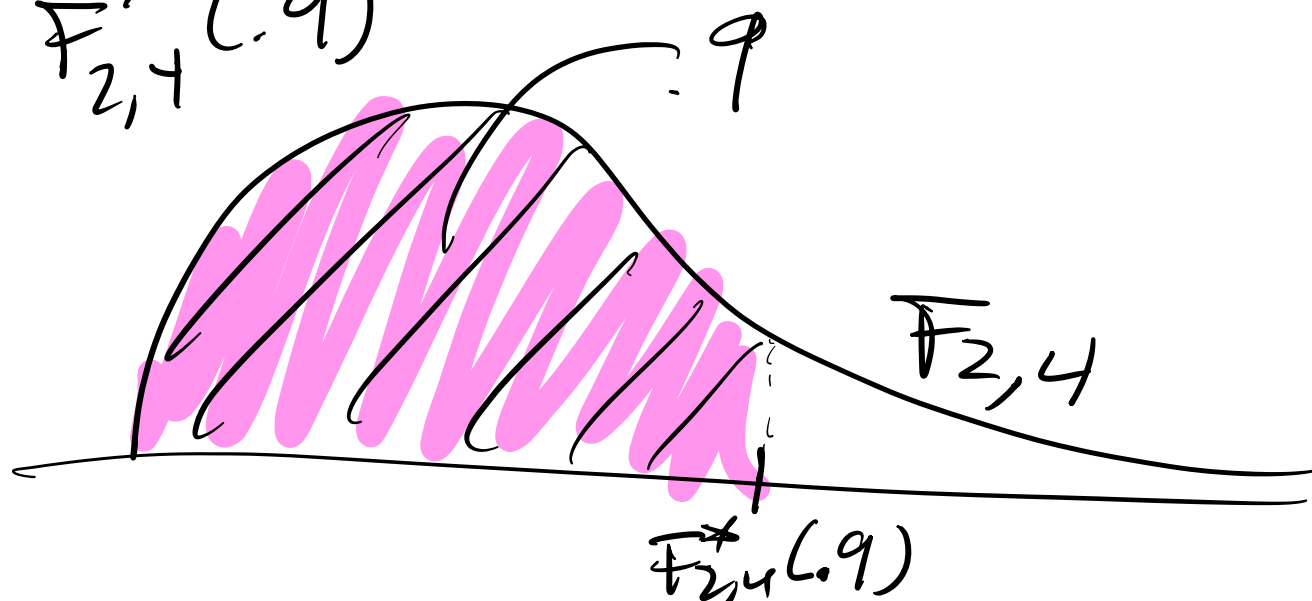
What are the  $\star$ s denoting?  
"critical values" or "quantiles".

EX:

$$t_{10}^{\star}(.95)$$



$$F_{2,4}^{\star}(.9)$$



## Coefficient of Determination

The sum of squares decomposition can also be used to evaluate "goodness-of-fit."

To do this we define the coefficient of determination,  $R^2$ :

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

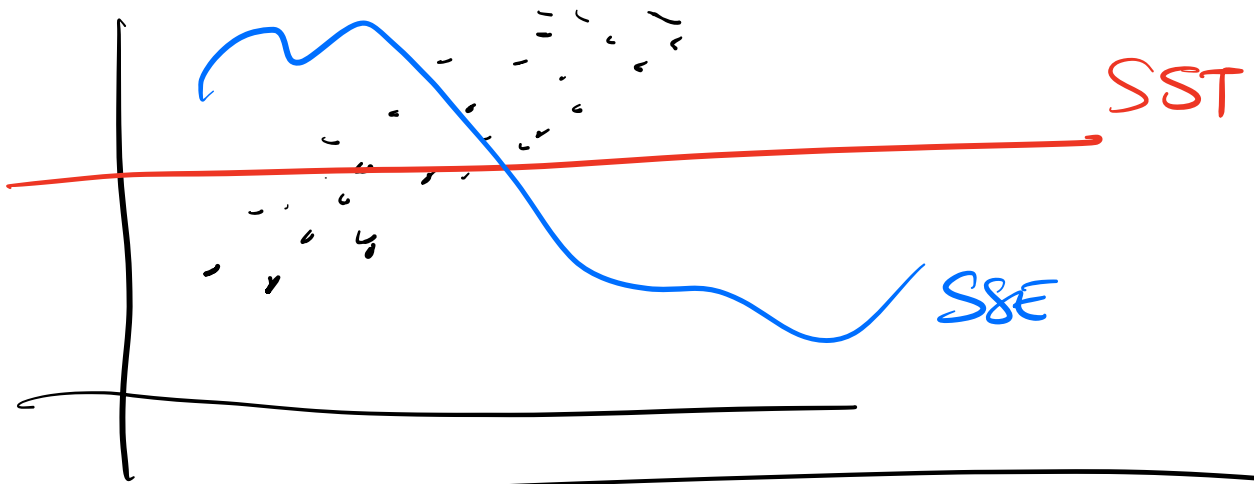
Facts about  $R^2$ :

①  $0 \leq R^2 \leq 1$  for SLR

In general,  $R^2$  doesn't have to be positive if we take

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

When would  $SSE > SST$   
 $(\Rightarrow R^2 < 0)$



② In SLRs we have

$$R^2 = r^2 \quad \text{where} \quad r = \text{corr. coeff.}$$

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

Consequence:

EX:

If I give you  $SSE = 4$ ,  $SST = 10$

&  $\text{sign}(\hat{\beta}_1) < 0$ , you can  
figure out  $r$ .

EX:

$$r = \sqrt{R^2}$$

$$R^2 = 1 - SSE / SST \\ = 1 - 4/10 = .6$$

$$r = -\sqrt{.6}$$

Interpreting  $R^2$ : observed

"  $(R^2)\%$  of variation in  $y$  ↓

is explained by the LSRL <sup>fitted</sup>  
the variation in  $x$