

## Inference

### ① Sampling distribution of $\hat{\beta}_0$ , $\hat{\beta}_1$ and $\hat{\sigma}^2$

Recall that

$$E(\hat{\beta}_1) = \beta_1, \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{s_{xx}}, \quad \text{and } E(\hat{\sigma}^2) = \sigma^2$$

$$E(\hat{\beta}_0) = \beta_0, \quad \text{Var}(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}_1^2}{s_{xx}} \right)$$

These results rely on the SLR assumption

$$E(\varepsilon | X_1=x_1) = 0, \quad \text{Var}(\varepsilon | X_1=x_1) = \sigma^2$$

$\Rightarrow$  Additional assumption: Gaussian-Noise SLR

$\Rightarrow$  conditional Gaussian distributed  $\gamma$ .

$$\Rightarrow \hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{s_{xx}})$$

$$\hat{\beta}_0 \sim N(\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}_1^2}{s_{xx}} \right))$$

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p} \quad (p=2 \text{ for SLR})$$

### 1.1 Gaussian-Noise Simple Linear Regression Model

Define Gaussian-Noise SLR as

- some
- ① The distribution of  $X$  is arbitrary (and perhaps  $X$  is even non-random)
  - ②  $\gamma = \beta_0 + \beta_1 x_1 + \varepsilon$

If  $x_1=x_1$ , then  $\gamma = \beta_0 + \beta_1 x_1 + \varepsilon$  for some coefficients  $\beta_0, \beta_1$  and random noise  $\varepsilon$

$$E(\varepsilon) = 0 \\ \text{Var}(\varepsilon) = \sigma^2$$

③  $\varepsilon \sim N(0, 1)$

- ④  $\varepsilon$  is independent of  $x_1$  and independent across observations.

$\varepsilon$  is uncorrelated with  $x_1$

linear correlation

Those having 0 correlation might depend on linearly

Stronger assumption

# 1.2 Sampling distribution of $(\hat{\beta}_1 - \beta_1) / se(\hat{\beta}_1)$ and $(\hat{\beta}_0 - \beta_0) / se(\hat{\beta}_0)$

Since  $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$   $\hat{\beta}_0 \sim N(\beta_0, \sigma^2(\frac{1}{n} + \frac{\bar{x}_1^2}{S_{xx}}))$

$$\Rightarrow \frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} \stackrel{\text{unknown}}{\sim} N(0, 1) \quad \frac{\hat{\beta}_0 - \beta_0}{se(\hat{\beta}_0)} \stackrel{\text{unknown}}{\sim} N(0, 1)$$

• unknown                                  • unknown

Replace  $se(\hat{\beta}_1)$  and  $se(\hat{\beta}_0)$  with  $ese(\hat{\beta}_1)$  and  $ese(\hat{\beta}_0)$

$$\Rightarrow T_1 = \frac{\hat{\beta}_1 - \beta_1}{ese(\hat{\beta}_1)} \sim t_{n-2} \quad T_0 = \frac{\hat{\beta}_0 - \beta_0}{ese(\hat{\beta}_0)} \sim t_{n-2}$$

where  $p=2$  for AN-SLR

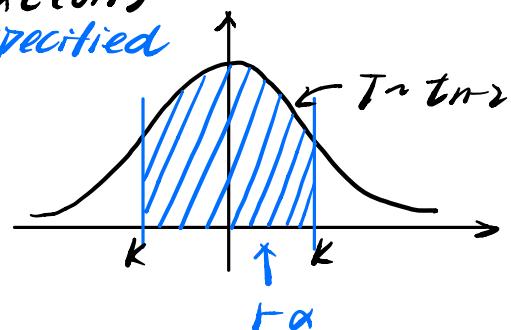
We can compute  $ese(\hat{\beta}_1) = \sqrt{\frac{MS_{RES}}{S_{xx}}}$  and  $ese(\hat{\beta}_0) = \sqrt{MS_{RES}(\frac{1}{n} + \frac{\bar{x}_1^2}{S_{xx}})}$  from the data without knowing  $\sigma^2$ .

## ② Confidence Interval

For  $\beta_1$ . Suppose that  $f_T$  is the density of a  $T \sim t_{n-2}$  distribution, define

$$k \equiv t_{\frac{\alpha}{2}, n-2} > 0$$

such that  $P(-k \leq T \leq k) = 1 - \alpha$   $\alpha \in (0, 1)$   
• user specified



A  $1 - \alpha$  level confidence interval for  $\beta_1$  is

$$CI(\hat{\beta}_1) = [\hat{\beta}_1 - k \cdot ese(\hat{\beta}_1), \hat{\beta}_1 + k \cdot ese(\hat{\beta}_1)]$$

$$P(\beta_1 \in CI(\hat{\beta}_1))$$

$$= P(\hat{\beta}_1 - k \cdot ese(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + k \cdot ese(\hat{\beta}_1))$$

$$= P(-k \leq \frac{\hat{\beta}_1 - \beta_1}{ese(\hat{\beta}_1)} \leq k) \quad T_1$$

$$= 1 - \alpha.$$

SLR

Estimation

$\hat{\beta}_1, \text{esec}(\hat{\beta}_1), \hat{r}^2$

✓

AN-SLR

✓

not necessary

unbiasedness  $E(\hat{\beta}_1)$

Variance  $\text{Var}(\hat{\beta}_1)$

✓

✓

not necessary

confidence interval

✓

hypothesis testing  
(t-test)

✓

F-test

✓

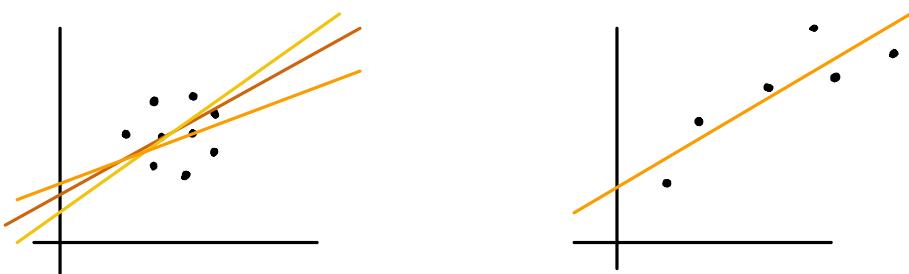
Prediction interval

✓

The upper and lower bound for  $\text{CI}(\hat{\beta}_1)$  are random, since their formula involve the estimator  $\hat{\beta}_1$ , which could be potentially random if  $\hat{\beta}_1$  is computed using a different set of samples.

### [ Interpretation ]

- The interval  $\text{CI}(\hat{\beta}_1)$  traps  $\beta_1$  with probability  $1-\alpha$ .
- $\beta_1$  is non-random
- The width of the CI is  $2 \cdot k \cdot \text{esec}(\hat{\beta}_1)$ 
  - As  $\alpha$  shrinks, the interval widens
  - As  $n$  grows, the interval shrinks.
  - As  $\sigma^2$  increases, the interval widens.
  - As  $S_{xx}$  grows, the interval shrinks.



A  $1-\alpha$  CI( $\hat{\beta}_0$ ) is

$$\text{CI}(\hat{\beta}_0) = [\hat{\beta}_0 - k \cdot \text{esec}(\hat{\beta}_0), \hat{\beta}_0 + k \cdot \text{esec}(\hat{\beta}_0)]$$

$k$  remains the same  
as long as the distribution  
and sample size don't change.

## ③ Hypothesis Testing

two sided test  $\begin{cases} H_0: \beta_1 = c \\ H_1: \beta_1 \neq c \end{cases}$

Rule: We reject  $H_0$  in hypothesis test with significant level  $\alpha$ , if  $100(1-\alpha)\%$  confidence interval  $CI(\hat{\beta}_1)$  doesn't cover  $c$ .

Why: Let's set  $\alpha = 0.01$  for example.

Because if  $H_0$  is true, then highly likely (with probability  $1-\alpha=0.99$ )  $CI(\hat{\beta}_1)$  should trap  $\beta_1 = c$ .

- either in fact  $\beta_1 = c$  ( $H_0$  is true), and therefore you just (by chance) observe a very rare event that  $CI(\hat{\beta}_1)$  doesn't trap  $\beta_1 = c$  (with probability  $\alpha = 0.01$ )
- or just because  $\beta_1 \neq c$ , and  $H_0$  should be rejected.

If we reject  $H_0$  when  $CI(\hat{\beta}_1)$  doesn't cover  $c$  i.e.  $c \notin CI(\hat{\beta}_1)$

$$\Leftrightarrow c \geq \hat{\beta}_1 + k \cdot \text{ese}(\hat{\beta}_1) \text{ or } c \leq \hat{\beta}_1 - k \cdot \text{ese}(\hat{\beta}_1)$$

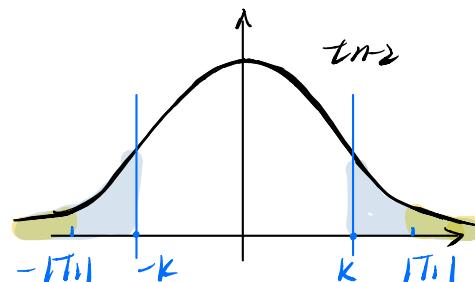
$$\Leftrightarrow \frac{\hat{\beta}_1 - c}{\text{ese}(\hat{\beta}_1)} \leq -k \text{ or } \frac{\hat{\beta}_1 + c}{\text{ese}(\hat{\beta}_1)} \geq k$$

$$\Leftrightarrow |T_1| \geq k \equiv \frac{\alpha}{2}, n-2 \quad T_1 = \frac{\hat{\beta}_1 - c}{\text{ese}(\hat{\beta}_1)} \quad t \text{ test statistic}$$

Wald Test or t-test

$$\Leftrightarrow \Pr(|T_1| > |T_{\alpha/2}|) < \alpha$$

p-value



For  $\beta_0$

$\begin{cases} H_0: \beta_0 = c \\ H_1: \beta_0 \neq c \end{cases}$

We reject  $H_0$  if  $|T_0| \geq k \equiv \frac{\alpha}{2}, n-2 \quad T_0 = \frac{\hat{\beta}_0 - c}{\text{ese}(\hat{\beta}_0)}$

### 3.1 Testing Significance of Regression

Failing to reject  $H_0: \beta_1 = 0$  implies that there is no linear relationship between  $X_1$  and  $Y$ .

- either that  $X_1$  is of little value in explaining the variation in  $Y$ .
- or that the true relationship between  $X$  and  $Y$  is not linear.

We reject  $H_0$  with significance level  $\alpha$  if the  $100(1-\alpha)\%$  confidence interval  $(\hat{\beta}_1)$  doesn't cover 0.

The corresponding test statistic is

$$T_1 = \frac{\hat{\beta}_1 - 0}{\text{ese}(\hat{\beta}_1)} = \frac{\hat{\beta}_1}{\text{ese}(\hat{\beta}_1)}$$

We reject  $H_0: \beta_1 = 0$  if  $|T_1| > k \equiv z_{\frac{\alpha}{2}, n-2}$

### 3.2 Statistical Significance

- What  $\alpha$  should we use? Conventionally set  $\alpha=0.05$ . If you're very conservative and cannot afford false rejection, then you should use smaller  $\alpha$ .
- If we test the hypothesis that  $H_0: \beta_1 = 0$  and reject, we say that the difference between  $\beta_1$  and 0 is **statistically significant**.
- statistical significance  $\neq$  importance or magnitude.  
When  $H_0$  is attained, it should be interpreted as  
"the effect of  $\beta_1$  is statistically undetectable." or  
" $\beta_1$  is statistically undistinguishable from 0."

No strong evidence to show.....

e.g. A common fallacy is

" I tested whether  $H_0: \beta_1 = 0$  or not, and I retain  $H_0$ .

Therefore,  $\beta_1$  is insignificant, and it is unimportant and you can ignore it."

Because we retain  $H_0: \beta_1 = 0$  when  $| \frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} | \leq k$ , there are two cases that might result in a retain of  $H_0$ :

- $\hat{\beta}_1$  is very close to zero, which implies  $\beta_1$  may be as well as zero, thus unimportant
- $\text{se}(\hat{\beta}_1)$  is so large that we can't tell anything about  $\beta_1$  with any confidence  $\Rightarrow \beta_1$  may be either large or small.

## ④ Analysis of Variance (ANOVA)

The distance from any point  $y_i$  in a collection of data to  $\bar{y}$ , is the deviation written as  $y_i - \bar{y}$ .

The ANOVA is based on a partition of total variability in  $y$ .

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$$

where  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1}$

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

$SST$                      $SS_{\text{Res}}$                      $SS_R$

Total sum of squares.

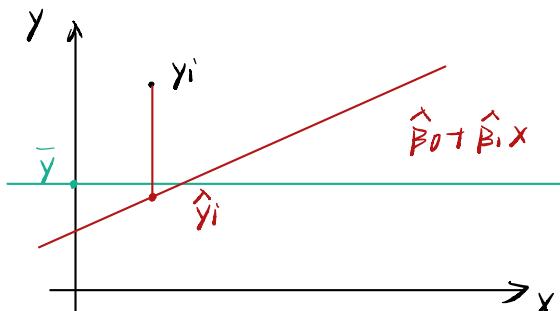
It measures the "total variance" in  $y$  that can be decomposed to two components.

Residual sum of squares

It measures amount of "residual variance" in  $y$ .

It is an aggregate measure of mis-fitted regression line.

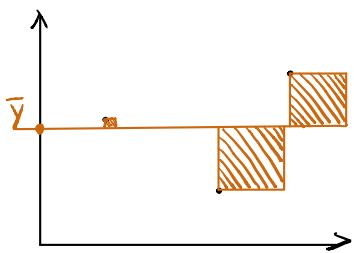
Regression sum of squares  
It measures amount of "systematic variation" in  $y$  due to the  $y \sim x$  linear relationship.



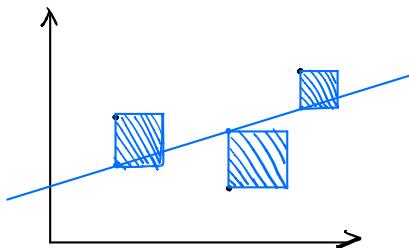
$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$$

cannot be explained by regression

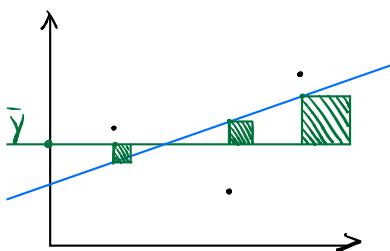
variation explainable by regression



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



$$SSRES = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



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

$$\begin{aligned}
 \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\
 &= \sum_{i=1}^n \left( \frac{(y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})}{e_i} \right)^2 \\
 &= \sum_{i=1}^n ((\hat{y}_i - \bar{y})^2 + 2e_i(\hat{y}_i - \bar{y}) + e_i^2) \\
 &= \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{SSR} + \frac{\sum_{i=1}^n e_i^2}{SSRES} + 2 \sum_{i=1}^n e_i (\hat{y}_0 + \hat{y}_1 x_{ii} - \bar{y}) \\
 &= SSR + SSRES + 2(\hat{y}_0 - \bar{y}) \sum_{i=1}^n e_i + 2\hat{y}_1 \sum_{i=1}^n e_i x_{ii} \quad 0
 \end{aligned}$$

## (5) F-test

(Intercept vs. whole model)

The idea is to compare two models

$$H_0: Y = \beta_0 + \varepsilon$$

$$H_1: Y = \beta_0 + \beta_1 X_1 + \varepsilon$$

or equivalent to

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0.$$

If we fit the first model, the least squares estimator is  $\hat{\beta}_0 = \bar{y}$

The idea is to create a statistic that measures how much better the second model is than the first model.

Under AN-SLR, and under  $H_0: \beta_1 = 0$ , we show that

$$\frac{SS_T}{\sigma^2} \sim \chi^2_{n-1} \quad df = n-1$$

$$\frac{SS_{RES}}{\sigma^2} \sim \chi^2_{n-p} \quad df_{RES} = n-p = n-2$$

$$\frac{SS_R}{\sigma^2} \sim \chi^2_{p-1} \quad df_R = p-1 = 2-1 = 1$$

for AN-SLR,  $p=2$

The degree of freedom has an additive property

$$df_T = df_R + df_{RES}$$

Consider now the F-test statistic.

$$F_0 = \frac{\frac{SSR}{\sigma^2} / df_R}{\frac{SS_{RES}}{\sigma^2} / df_{RES}} = \frac{\frac{SSR}{\sigma^2} / (p-1)}{\frac{SS_{RES}}{\sigma^2} / (n-p)} = \frac{MSR}{MS_{RES}}$$

under  $H_0: \beta_1 = 0$

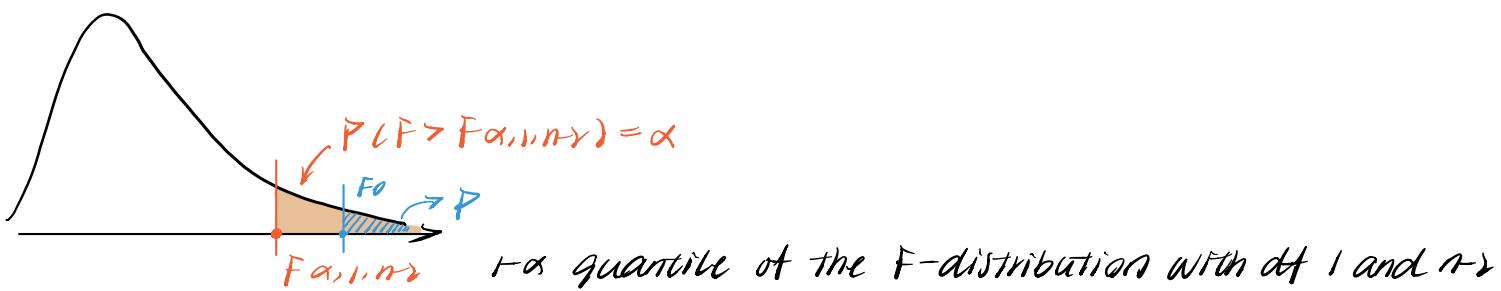
The statistic  $F_0$  follows F distribution, with degrees of freedom  $p-1$  and  $n-p$

$$F_0 \sim F_{p-1, n-p} = F_{1, n-2}$$

We reject  $H_0$  if  $F_0$  is "large enough"

We set significance level  $\alpha$  0.01 or 0.05 and reject  $H_0$  if  $F_0 > F_{\alpha, p-1, n-p}$

$\Leftrightarrow$  We reject  $H_0$  if p-value  $< \alpha$        $\Pr(F > F_0) < \alpha = \Pr(F > F_{\alpha, 1, n-2})$



Source	SS	df	MS	F
Regression	$SSR$	$p-1$	$MSR = \frac{SSR}{p-1}$	$\frac{SSR/(p-1)}{SSRes/(n-p)}$
Residual	$SSRes$	$n-p$	$MSRes = \frac{SSRes}{n-p}$	
Total	$SST$	$n-1$		

## 5.1 Equivalence of F-test and t-test for SLR.

$$(T_1)^2 = \left( \frac{\hat{\beta}_1}{\text{ese}(\hat{\beta}_1)} \right)^2 = F_0$$

The p-value for F-test is exactly the same as the p-value for t-test

$$F > F_0 \Leftrightarrow T^2 > (T_1)^2$$

$$\text{Thus, } P(F > F_0) = P(T > T_1)$$

Therefore, the test result based on F and t are the same.

Now we show that  $(T_1)^2 = F_0$ , since by

$$\begin{aligned} SSR &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n [(\hat{\beta}_0 + \hat{\beta}_1 x_{i1}) - (\hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1)]^2 \\ &= \hat{\beta}_1^2 \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2 \\ &= \hat{\beta}_1^2 S_{xx} \end{aligned}$$

$$\text{and } \text{ese}(\hat{\beta}_1) = \sqrt{\frac{\hat{\beta}_1^2}{S_{xx}}}$$

$$F = \frac{SSR/(p-1)}{SSRes/(n-p)} = \frac{SSR/1}{SSRes/(n-p)} = \frac{SSR}{\hat{\beta}_1^2} = \frac{\hat{\beta}_1^2 S_{xx}}{\text{ese}^2(\hat{\beta}_1) S_{xx}} = \left[ \frac{\hat{\beta}_1}{\text{ese}(\hat{\beta}_1)} \right]^2 = (T_1)^2$$

## 5.2 What the F-test really tests?

1. If retain  $H_0: \beta_1 = 0$  - do not find any significant share of variance associated with the regression.
  - (a) The intercept-only model is correct.
  - (b)  $\beta_1 \neq 0$ , but the test doesn't have enough power to detect it if  $\frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)}$  small
  - (c) The relationship is non-linear, but the best linear approximation has zero slope.
2. If reject  $H_0$  - doesn't mean the SLR model is right, only that the non-zero slope linear model predict better than the intercept-only model.

## ⑥ Prediction and Predictive Inference

### 6.1 Prediction.

The conditional mean of  $y$  is

$$\mu = m(x_i) = E(Y|X_i=x_i) = \beta_0 + \beta_1 x_i$$

while our estimate of the conditional mean is

$$\hat{\mu} = \hat{m}(x_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

Note:  $m(x)$  is deterministic, but unknown.

$\hat{m}(x)$  is a function of data, random, known. It inherits the randomness from  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , which in turn inherits from  $y$ .

We can show that

$$E(\hat{m}(x_i)) = \hat{\beta}_0 + \hat{\beta}_1 x_i = m(x_i)$$

$$\text{var}(\hat{m}(x_i)) = \sigma^2 \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}} \right)$$

Under the AN-SLR assumption,  $\hat{m}(x_1)$  is Gaussian.

$$\hat{m}(x_1) \sim N(m(x_1), \sigma^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{s_{xx}} \right))$$

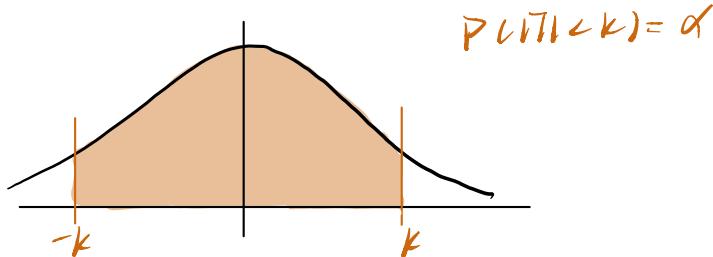
A  $100(1-\alpha)$  percent CI on the mean response, at the point  $x_1 = x_{01}$

$$CI(m(x_{01})) = \left[ \hat{m}(x_{01}) \pm k \sqrt{\text{MSRes} \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{s_{xx}} \right)} \right]$$

$$= \left[ \hat{m}(x_{01}) \pm k \cdot \frac{\text{ese}(\hat{m}(x_{01}))}{\hat{s}_s} \right]$$

$\sqrt{\text{var}(\hat{m}(x_{01}))}$  with  $\sigma^2$  replaced by  $\hat{\sigma}^2$ .

where  $k = T \frac{\alpha}{2}, n-2$ .



A  $100(1-\alpha)$  percent CI on the future observation is

$$CI(y) = \left[ \hat{m}(x_{01}) \pm k \cdot \sqrt{\text{MSRes} \left( 1 + \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{s_{xx}} \right)} \right]$$

We should distinguish mean response  $CI(m(x_{01}))$  and new observation  $CI(y)$

- A  $100(1-\alpha)\%$  CI on the mean response is an interval  $[l, u]$  where  $P(l \leq m(x_1) \leq u | x_1 = x_{01}) = 1-\alpha$

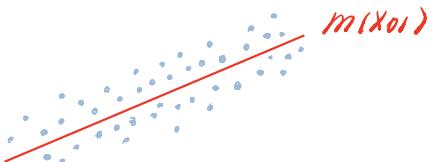
recall:  $Y = \beta_0 + \beta_1 x_1 + \varepsilon$

$$E(Y | X_1 = x_1) = m(x_1)$$

- A  $100(1-\alpha)\%$  CI on the new observation  $Y$  is an interval  $[l', u']$  where

$$P(l' \leq Y_{x_{01}} \leq u' | x_1 = x_{01}) = 1-\alpha$$

$$Y_{x_{01}} = \beta_0 + \beta_1 x_1 + \varepsilon$$



# Proof

## Mean Response

We predict on  $X_1 = x_{01}$

$$\begin{aligned}\hat{m}(x_{01}) &= \beta_0 + \beta_1 x_{01} \\ &= (\bar{y} - \hat{\beta}_1 x_{01}) + \hat{\beta}_1 x_{01} \\ &= \bar{y} + (x_{01} - \bar{x}_1) \hat{\beta}_1\end{aligned}$$

using  $\hat{\beta}_1 = \sum_{i=1}^n c_i y_i$ ,  $c_i = \frac{x_{i1} - \bar{x}_1}{S_{xx}}$   $\sum_{i=1}^n c_i = 0$

Now plug-in  $\hat{\beta}_1$

$$\begin{aligned}\hat{m}(x_{01}) &= \bar{y} \sum_{i=1}^n c_i y_i + (x_{01} - \bar{x}_1) \sum_{i=1}^n c_i y_i \\ E(\hat{m}(x_{01})) &= E \left[ \bar{y} \sum_{i=1}^n c_i y_i + (x_{01} - \bar{x}_1) \sum_{i=1}^n c_i y_i \right] \\ \textcircled{1} \quad E \left[ \bar{y} \sum_{i=1}^n c_i y_i \right] &= E \left[ \bar{y} \sum_{i=1}^n (\beta_0 + \beta_1 x_{i1} + \varepsilon_i) c_i \right] = \beta_0 + \beta_1 \bar{x}_1 \\ \textcircled{2} \quad E[(x_{01} - \bar{x}_1) \sum_{i=1}^n c_i y_i] &= (x_{01} - \bar{x}_1) \beta_1 \quad E(\hat{\beta}_1) = \beta_1\end{aligned}$$

$$\Rightarrow E(\hat{m}(x_{01})) = \beta_0 + \beta_1 \bar{x}_1 + (x_{01} - \bar{x}_1) \beta_1 = \beta_0 + \beta_1 x_{01}$$

$$\begin{aligned}\text{var}(\hat{m}(x_{01})) &= \text{var} \left( \bar{y} \sum_{i=1}^n c_i y_i + (x_{01} - \bar{x}_1) \sum_{i=1}^n c_i y_i \right) \\ &= \text{var} \left( \bar{y} \sum_{i=1}^n (1 + n(x_{01} - \bar{x}_1) c_i) y_i \right) \quad y_i, y_j | x_i = x_{01} \text{ independent} \\ &= \bar{y}^2 \sum_{i=1}^n (1 + n(x_{01} - \bar{x}_1) c_i)^2 \text{var}(y_i) \\ &= \frac{\sigma^2}{n^2} \sum_{i=1}^n (1 + 2n(x_{01} - \bar{x}_1) c_i + n^2(x_{01} - \bar{x}_1)^2 c_i^2) \\ &= \frac{\sigma^2}{n^2} (n + 0 + n^2(x_{01} - \bar{x}_1)^2 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}{S_{xx}}) \\ &= \frac{\sigma^2}{n^2} (n + 0 + \frac{n^2(x_{01} - \bar{x}_1)^2}{S_{xx}}) \\ &= \sigma^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)\end{aligned}$$

cannot eliminate.

NOTE: The variance of  $m(x_{01})$  grows as  $\sigma^2$  goes.  
The larger the  $n$ , the smaller the variance

The further our prediction point  $x_{01}$  is from the center of data, the bigger the variance.

Thus, under AN-SLR,  $\hat{m}(x_{01})$  is gaussian

$$\begin{aligned}\hat{m}(x_{01}) &\sim N(E(\hat{m}(x_{01})), \text{Var}(\hat{m}(x_{01}))) \\ &= N(m(x_{01}), \sigma^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right))\end{aligned}$$

$$\frac{\hat{m}(x_{01}) - m(x_{01})}{\sqrt{\sigma^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)}} \sim N(0, 1)$$

$$\frac{\hat{m}(x_{01}) - m(x_{01})}{\text{ese}(\hat{m}(x_{01}))} = \frac{\hat{m}(x_{01}) - m(x_{01})}{\sqrt{\sigma^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)}} \sim t_{n-2}$$

Thus, a  $100(1-\alpha)\%$  CI on the mean response at point  $x_1 = x_{01}$  is

$$CI(m(x_{01})) = \left[ \hat{m}(x_{01}) \pm k \sqrt{\sigma^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)} \right]$$

### Future Observation

Using AN-SLR model, we know that at  $x_1 = x_{01}$ ,  
the future observation of  $y_0$  is

$$y_0 = m(x_{01}) + \varepsilon = \beta_0 + \beta_1 x_{01} + \varepsilon$$

Now we could construct a confidence interval for  $y_0$   
or prediction interval

We call it the prediction interval for the future observation  $y_0$ .  
(which is predicted at  $x_1 = x_{01}$ )

$$\begin{aligned}PI(y_0) &= [\hat{y}_0 \pm k \cdot \text{ese}(\hat{y}_0)] \quad \xrightarrow{\text{wider than } CI(m(x_{01}))} \\ &= \left[ \hat{m}(x_{01}) \pm k \cdot \sqrt{MS_{\text{Res}} \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)} \right]\end{aligned}$$

Since  $Y_0 | X_1 = x_{01} = m(x_{01}) + \varepsilon \sim N(m(x_{01}), \sigma^2)$

The variance of  $Y_0$  is  $\text{Var}(Y_0) = \text{Var}(m(x_{01}) + \varepsilon)$   
=  $\text{Var}(m(x_{01})) + \text{Var}(\varepsilon)$   
=  $\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_{01}) + \sigma^2$   
=  $\sigma^2$ .

In the formula,  $\sigma^2, \hat{\beta}_0, \hat{\beta}_1$  are unknown, replace them with  $\hat{\sigma}^2, \hat{\beta}_0, \hat{\beta}_1$

$$\begin{aligned}\hat{\text{Var}}(Y_0) &= \hat{\text{Var}}(\hat{m}(x_{01}) + \varepsilon) \\ &= \hat{\text{Var}}(\hat{\beta}_0 + \hat{\beta}_1 x_{01}) + \hat{\text{Var}}(\varepsilon) \\ &= \hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right) + \hat{\sigma}^2 \\ &= \hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)\end{aligned}$$

Therefore, the estimated standard error of  $Y_0$  is

$$\text{esr}(Y_0) = \sqrt{\hat{\text{Var}}(Y_0)} = \sqrt{1 + \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}}}$$

## ⑦ $R^2$ (crubbish)

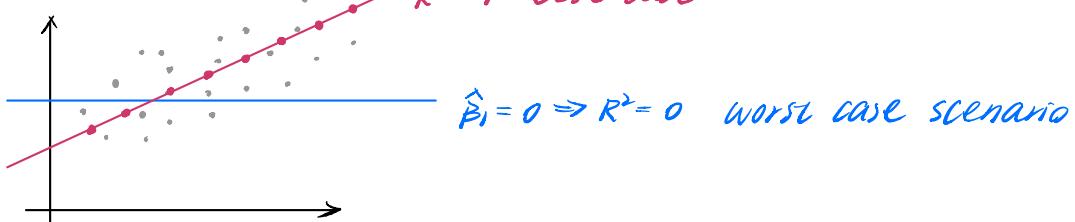
If we want to get a global measure of how well  $X_1$  predict  $Y$ , you may consider the proportion of variation explained the regression model.

$$R^2 = \frac{SSR}{SST} = \frac{SST - SS_{\text{Res}}}{SST} = 1 - \frac{SS_{\text{Res}}}{SST}$$

Because  $SSR = \hat{\beta}_1 S_{xy} = \hat{\beta}_1^2 S_{xx}$

- $R^2$  will be 0 if  $\hat{\beta}_1 = 0$
- $R^2$  will be 1 if all the residuals are zero  $SS_{\text{Res}} = 0$ .

$R^2 = 1$  best case scenario



## Interpretation of $R^2$ .

Distinguish  $R^2$  from t-test (P-values)

### ① Low $R^2$ and low P-value ( $P\text{-value} < \alpha$ )

The model doesn't explain much of variation of the data but the model is significant.

(Including  $X_1$  now the model is necessary but not sufficient).

### ② Low $R^2$ and high P-value

The model doesn't explain much of variation of the data and  $X_1$  is not useful < worst cases

### ③ High $R^2$ and low P-value

Model sufficient,  $X_1$  is useful < best cases

### ④ High $R^2$ and high P-value

## Adjusted $R^2$

We have  $E(SS_{\text{Res}}) = \sigma^2(n-p)$  # variables

If we increase the model complexity  $P \rightarrow n$

$$SS_{\text{Res}} \rightarrow 0 \quad R^2 = 1 - \frac{SS_{\text{Res}}}{SS_T} \rightarrow 1$$

$R^2$  doesn't account for model "complexity".

## adjusted $R^2$

$$R^{\text{adj}} = 1 - \frac{SS_{\text{Res}} / (n-p)}{SS_T / (n-1)}$$

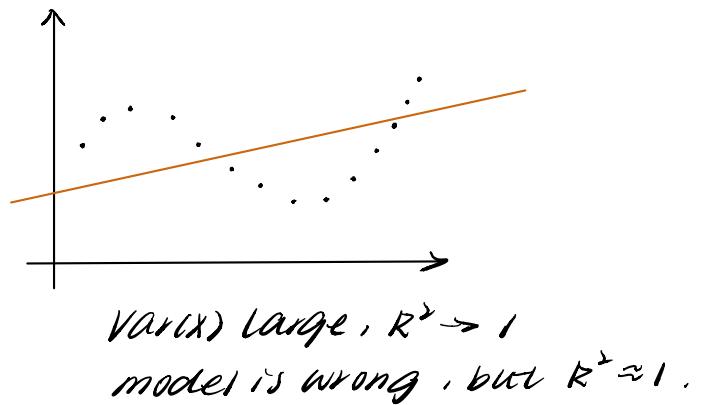
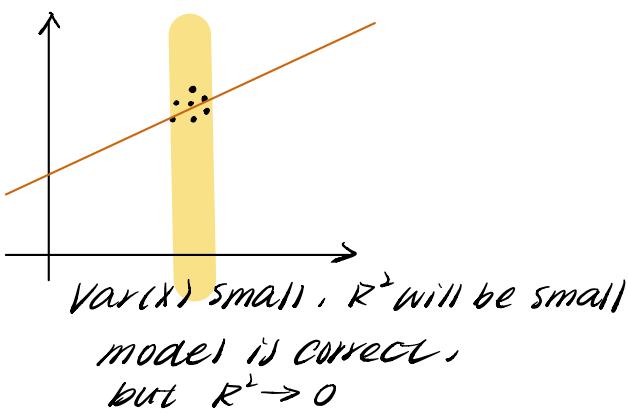
Note :

- When  $p=1$ ,  $R^{\text{adj}} = R^2$
- Adjusted  $R^2$  is better than  $R^2$  in term of model selection.

## Limitation of $R^2$ (or $R_{adj}^2$ )

$$R^2 = \frac{SSR}{SST} = \frac{SSR}{SSR + SSR_{\text{res}}} = \frac{\hat{\beta}_1^2 S_{xx}}{\hat{\beta}_1^2 S_{xx} + SSR_{\text{res}}} = \frac{\hat{\beta}_1^2 \frac{S_{xx}}{n}}{\hat{\beta}_1^2 \frac{S_{xx}}{n} + \frac{SSR_{\text{res}}}{n}} \xrightarrow{n \rightarrow \infty} \frac{\hat{\beta}_1^2 \text{Var}(X)}{\hat{\beta}_1^2 \text{Var}(X) + 0}$$

$$\hat{\alpha} = \frac{SSR_{\text{res}}}{n} \xrightarrow{n \rightarrow \infty} 0$$



Note:

- (a) We can manipulate the data that  $R^2 = 0$  but the model is correct.  
 $R^2 = 1$  but the model is wrong
- (b)  $R^2$  can be compared when different models are fit to the same dataset with the same, untransformed response variable.  
 $R^2$  cannot be compared across different datasets.