

Multivariate

①

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_{p-1} x_{ip-1} + \varepsilon_i$$

→ if 5 basic assumptions hold \Rightarrow LS fit OK

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p-1} \\ 1 & x_{21} & x_{22} & \dots & x_{2p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np-1} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\begin{matrix} y \\ n \times 1 \end{matrix} = \begin{matrix} X \\ n \times p \end{matrix} \begin{matrix} \beta \\ p \times 1 \end{matrix} + \begin{matrix} \varepsilon \\ n \times 1 \end{matrix}$$

$$\text{Minimize LS } Q = \sum (y_i - \hat{y}_i)^2 \text{ where } \hat{y}_i = (1, x_{i1}, \dots, x_{ip-1})'$$

$$= \sum e_i^2$$

$$= e'e = (y - \hat{y}\beta)'(y - \hat{y}\beta)$$

Scalar $|X|$

$$\Rightarrow \frac{\partial Q}{\partial \beta} = 0 \Rightarrow \boxed{\text{Normal Equations } (\hat{X}'\hat{X})\beta = \hat{X}'y}$$

$$(A) \text{ In simple linear regression } \rightarrow \text{solution } \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\text{cov}(x, y)}{\text{var}(x)} \sim \text{var}(x, y)$$

$$(B) \text{ Multiple case } \rightarrow (\hat{X}'\hat{X}) \sim \text{cov}(\hat{X}) \text{ where } \text{var}(x_1), \text{var}(x_2), \dots, \text{var}(x_{p-1}) \text{ on the diagonal, and cross-covariances } \text{cov}(x_1, x_2), \dots, \text{cov}(x_{p-1}, x_p) \text{ off-diagonal}$$

$$\rightarrow (X'y) = \begin{pmatrix} \text{cov}(x_1, y) \\ \text{cov}(x_2, y) \\ \vdots \\ \text{cov}(x_{p-1}, y) \end{pmatrix}$$

(2)

- If $(\bar{X}'\bar{X}) \sim \Lambda$ (diagonal), all x 's are uncorrelated.

$$\rightarrow \text{then } \hat{\beta} = (\bar{X}'\bar{X})^{-1}\bar{X}'y \sim \begin{pmatrix} \text{corr}(x_1, y) \\ \vdots \\ \text{corr}(x_p, y) \end{pmatrix}$$

i.e. each coefficient has a direct interpretation as
the dependency between an individual x -variable and y .

- If $(\bar{X}'\bar{X})$ not diagonal \rightarrow no direct interpretation

\rightarrow can't separate influence of one x on y
from another (correlated) x .

\rightarrow That is, $\hat{\beta}_k$ involves both $\text{corr}(x_k, y)$ and
 $\text{corr}(x_k, y)$, & $\text{corr}(x_k, x_k)$

- Moreover \rightarrow if x 's are highly correlated, $(\bar{X}'\bar{X})^{-1}$
may not exist! ($\det(\bar{X}'\bar{X}) \rightarrow 0$)

or inverse $(\bar{X}'\bar{X})^{-1}$ numerically unstable

\rightarrow When this happens, magnitude & even sign of $\hat{\beta}$'s
may become meaningless. [Many models fit the data equally
well]

Ex : $x_2 = a + b x_1$, (3)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

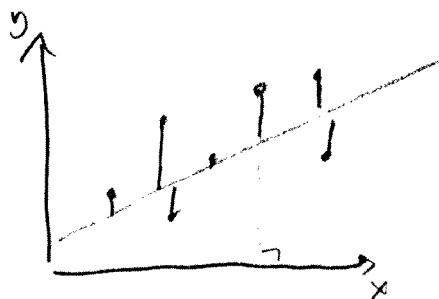
$\left. \right\}$ infinite set of β_0 & β_1 , fit these data equally well

The Hat-matrix

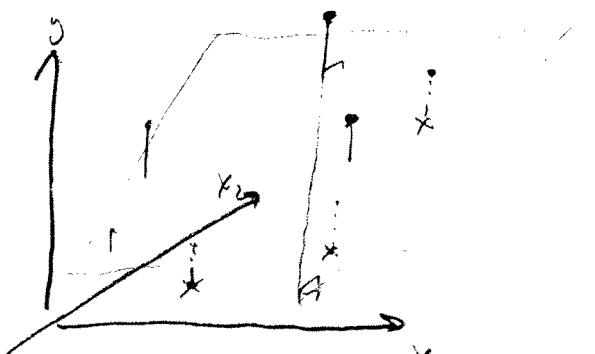
Assume we can solve the normal equations & obtain

$$\hat{\beta} = (\bar{X}' \bar{X})^{-1} \bar{X}' y$$

- Fitted values $\hat{y} = \bar{X} \hat{\beta} = \underbrace{\bar{X} (\bar{X}' \bar{X})^{-1} \bar{X}'}_{\text{Hat-matrix}} y$
- $H = \bar{X} (\bar{X}' \bar{X})^{-1} \bar{X}'$
- H is a projection matrix — projects the data y onto the plane spanned by \bar{X}



Simple case



projection \perp to both x_1 & x_2
in multiple case.

Facts

(4)

- $H^2 = H \quad \left\{ \begin{array}{l} (H \text{ is an idempotent matrix}) \\ (H^T = H) \quad (\text{NO point in redoing regression}) \\ (\text{Symmetric}) \end{array} \right.$
- $e = \hat{\epsilon} = y - \hat{y} = y - Hy = (I - H)y$
- $X'e = X'(I - X(X'X)^{-1}X')y = 0$
 $\quad (\text{orthogonal projection})$
- $y'e = y'H(I - H)y = 0$
 $\quad (\text{fitted values } \perp \text{ to residuals})$

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{pmatrix} h_{11} & h_{12} & \cdots & h_{1n} \\ h_{21} & h_{22} & \cdots & h_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ h_{n1} & h_{n2} & \cdots & h_{nn} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

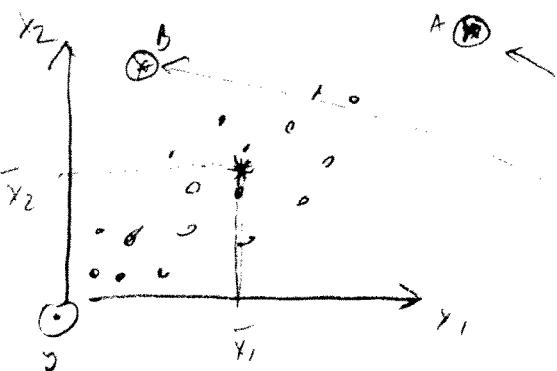
↓
ith row of X

leverage $h_{ii} = x_i'(X'X)^{-1}x_i'$

Note, h_{ii} is large when x_i is extreme compared with \bar{x}

/
 vector of all x-variable values for observation i
 vector, average
 x-values.

Extreme how? wrt mass of data



Extreme in both \$x_1\$ & \$x_2\$

Extreme cmp. w. mass of data
(extreme in \$x_2\$ only.)

Both ① observations have high leverage.

Properties

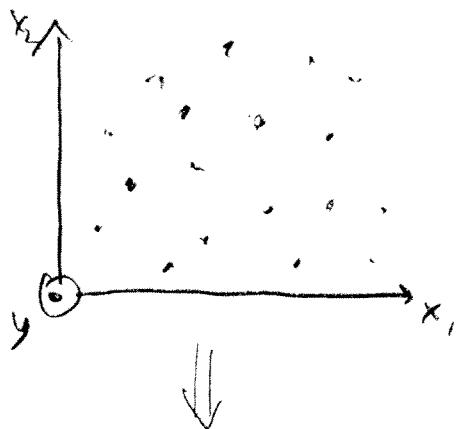
• As before, $E(\hat{\beta}) = \beta$ (unbiased)

$$\cdot V(\hat{\beta}) = V((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'V(\mathbf{y})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$$

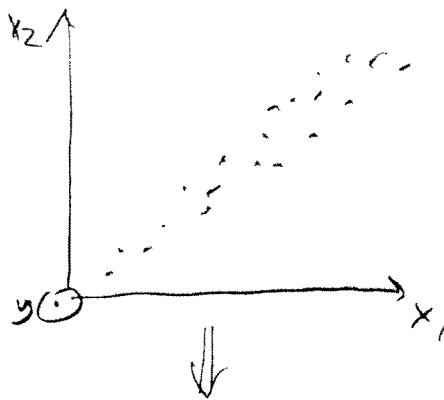
Constant $\sigma^2 I$

Meaning? \rightarrow More spread in any x ($\text{var}(x_i)$ large)
 \rightarrow better estimate of $\hat{\beta}$.

$\rightarrow \sigma^2 \uparrow \rightarrow \text{variance} \uparrow$
 Unless $(\mathbf{X}'\mathbf{X})$ is diagonal, $\hat{\beta}$'s are correlated \Rightarrow Think about impact on inference!
 The more dependent x 's are, the larger the $V(\hat{\beta})$.



$\hat{\beta}_1, \hat{\beta}_2$ nearly uncorrelated



$\hat{\beta}_1, \hat{\beta}_2$ heavily dependent

- large variance associated with estimates
- difficult to make direct inferences.

$$\text{Fitted values } \hat{y} = Hy \Rightarrow E(\hat{y}) = E(y)$$

$$V(\hat{y}) = H V(y) H^T = \sigma^2 H$$

(i.e. large variance @ locations of high leverage)

$$\text{Residuals } e = (I - H)y \Rightarrow E(e) = 0$$

$$V(e) = \sigma^2 (I - H)$$

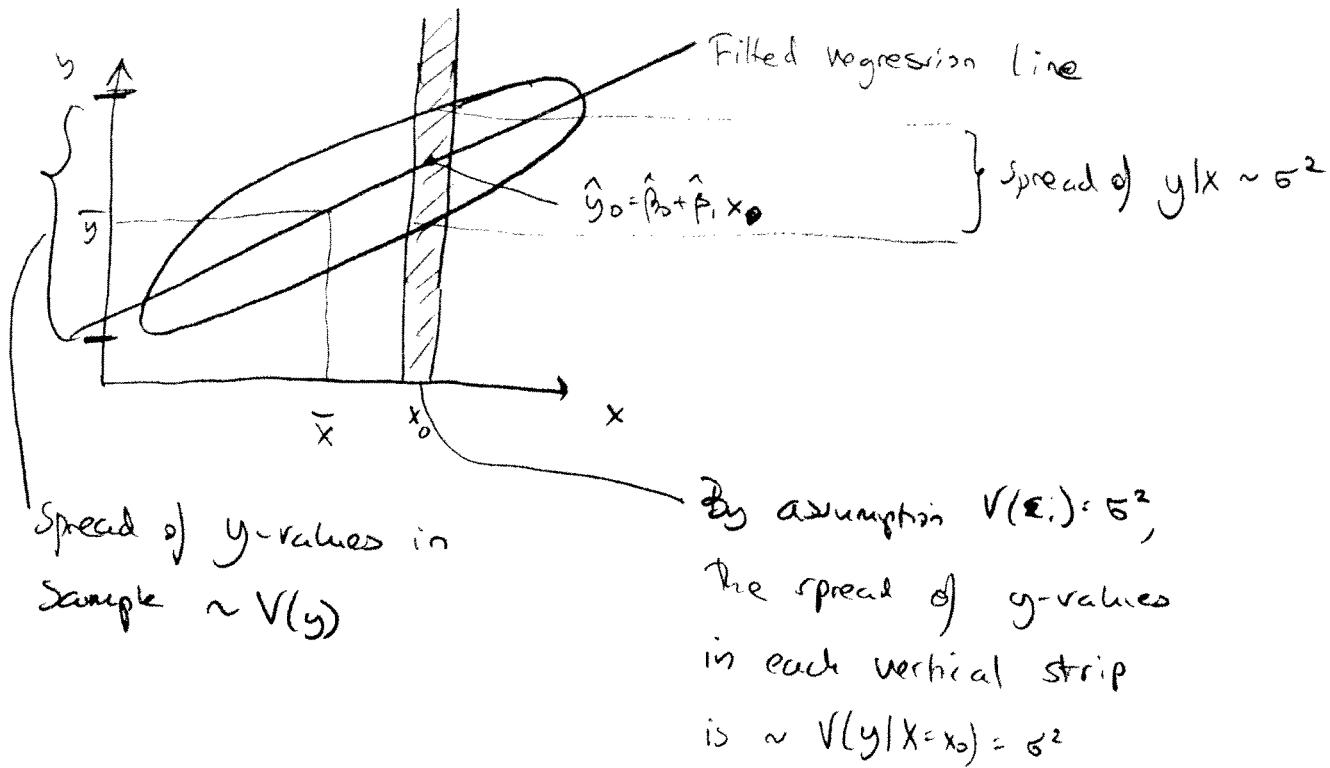
↓ smaller variance @ locations w. high leverage

e's are correlated!

Variance Decomposition

Assume $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$.

①



- Now, if x and y are unrelated, then $V(y) \approx V(y|X=x)$.
- If $y = \beta_0 + \beta_1 x$ exactly, $V(y|X=x) = 0$!

So how much did linear regression reduce the uncertainty in y ?

① $RSS = SSE = \sum (y_i - \hat{y}_i)^2$ (spread of y 's around the regression line)

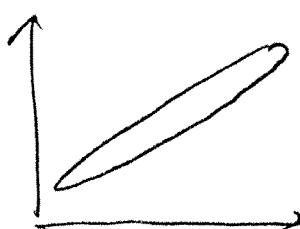
② $SST = \sum (y_i - \bar{y})^2$ (spread around mean of y)

- Since \hat{y} is obtained from the data (the best possible line is fit to the n sample points), RSS is always less than SST

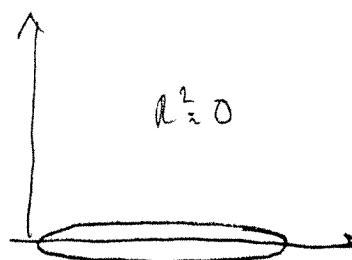
→ Numerical summary; $R^2 = \frac{SST - RSS}{SST} = \frac{\text{reduction in spread}}{\text{original spread}}$

$$= 1 - \frac{RSS}{SST}$$

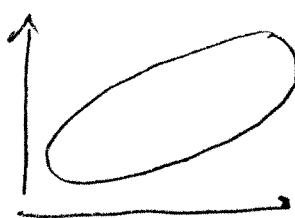
R^2 : the % of variance (spread) in y explained by the regression line.



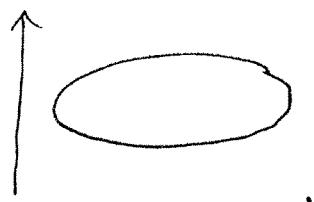
$$R^2 \approx 1$$



$$R^2 \approx 0$$



$$R^2 \approx 0.4$$

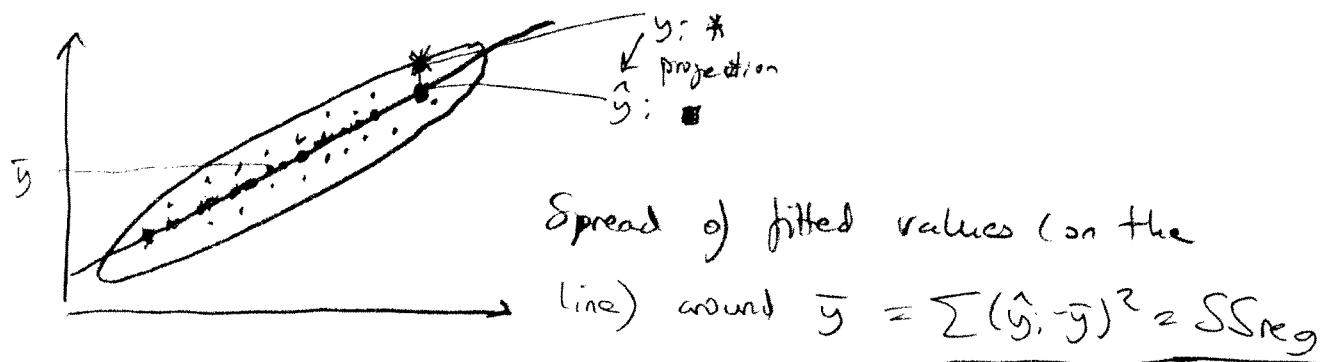


$$R^2 \approx 0$$

Variance decomposition \Rightarrow

③

$$\begin{aligned}
 SS_T - SSE &= \sum (y_i - \bar{y})^2 - \sum (y_i - \hat{y}_i)^2 \\
 &= \underbrace{\sum y_i^2 + \sum \bar{y}^2 - 2\bar{y} \cdot n \bar{y}}_{= n \bar{y}^2} - \sum y_i^2 - \sum \hat{y}_i^2 + 2 \sum y_i \hat{y}_i \\
 &= -n \bar{y}^2 - \sum \hat{y}_i^2 + 2 \sum \hat{y}_i^2 + 2 \sum e_i \hat{y}_i \\
 &= \sum \hat{y}_i^2 - n \bar{y}^2 = \sum (\hat{y}_i - \bar{y})^2 \quad \stackrel{= 0}{\text{by}} \perp \text{projection}
 \end{aligned}$$



So, we have $\boxed{SS_T = SSE + SS_{Reg}}$

$$R^2 = 1 - \frac{SSE}{SS_T} = \frac{SS_{Reg}}{SS_T}$$

How large should R^2 be for regression to "make sense"?



More formally

(9)

⇒ The F-test Goodness-of-fit test (or lack-of-fit)

Assume y and x are not linearly related, i.e. $\beta_1 = 0$!

Then, what would happen to SS_{reg} & SS_E ? (R^2)

- Well, if $\beta_1 = 0$ then true model for y is $y_i = \beta_0 + \varepsilon_i$ ($V(\varepsilon_i) = \sigma^2$)

and our best estimate for β_0 is $\hat{\beta}_0 = \bar{y}$

The expected value of SS_T under the null ($\beta_1 = 0$, \dagger) is

$$E_0 \left(\sum (y_i - \bar{y})^2 \right) = (n-1)\sigma^2 \quad (\text{just our basic estimate of the variance, remember})$$

- What if $\beta_1 \neq 0$ (or maybe it is, but we fit the line to the data anyway), then

$$E \left(\sum (y_i - \hat{y}_i)^2 \right) = E \left(\sum e_i^2 \right) = \sum_i E(y_i^2) + E(\hat{y}_i^2) - 2 E(y_i \hat{y}_i) =$$

$$= \begin{bmatrix} V(\hat{y}_i) = \sigma^2 h_{ii} \\ \sum e_i \hat{y}_i = 0 \\ \sum h_{ii} = \sum_{i=1}^n \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right) = 2 \end{bmatrix} = \sum_i V(y_i) + E(y_i)^2 + V(\hat{y}_i) + E(\hat{y}_i)^2 - 2 \underbrace{(\text{or } (\hat{y}_i + e_i, \hat{y}_i))}_{-2 E(y_i) E(\hat{y}_i)} = -2 V(\hat{y}_i)$$

$$= \sum V(y_i) - \sum V(\hat{y}_i) = \sigma^2 \cdot n - \sigma^2 \cdot 2 = (n-2)\sigma^2$$

[Note, i) we have p parameters in the regression
model $E(SS_E) = \sigma^2(n-p)$]

• What about SS_{reg} ? (5)

Under the null / $\beta_1 = 0$, $E(y_i) = \beta_0$

$$\Rightarrow E_0 \left(\sum (y_i - \bar{y})^2 \right) = E_0 \left(\sum \hat{y}_i^2 + \bar{y}^2 - 2\bar{y}\hat{y}_i \right) = E_0(\sum \hat{y}_i^2) + \underset{(1)}{E_0(\bar{y}^2)} - 2\underset{(2)}{E_0(\bar{y}\hat{y}_i)} + \underset{(3)}{E_0(\bar{y}^2)}$$

$$= \sigma^2 \underbrace{\sum h_{ii}}_{(1)} + n\beta_0^2 + \underbrace{\bar{y}^2}_{(2)} + n\beta_0^2 - 2\bar{y}^2 - 2n\beta_0^2 = \\ (E(y_i) + \sum E(\hat{y}_i)^2) (n(\bar{y}) + E(\bar{y})^2)) (E(\bar{y}\sum \hat{y}_i) = \sum (nw(y_i, \beta_0) + E(\hat{y}))E(\hat{y}_i)),$$

$$= \sigma^2 (\sum h_{ii} - 1) = \sigma^2$$

Note, in general w. p parameters in model, $\sum h_{ii} = p$ and $E_0(SS_{reg}) = \sigma^2(p-1)$

MEANING:

(1) $\beta_1 = 0$ (x and y are not related)

all SS's provide estimates of the same parameter, σ^2 !

Under the null ($\beta_1 = 0$)

$$\frac{SST}{n-1} = \hat{\sigma}^2, \quad \frac{SSE}{n-p} = \hat{\sigma}^2 \quad \text{and} \quad \frac{SS_{reg}}{p-1} = \hat{\sigma}^2$$

Extra assumption : $\epsilon_i \sim N(0, \sigma^2)$ (6)

Now, i) ϵ_i ~ normally distributed and $\beta_1 = 0$

$$\frac{\text{Ssreg}}{\sigma^2} \sim \chi^2_{p-1} \quad \perp \quad \frac{\text{SSE}}{\sigma^2} \sim \chi^2_{n-p}$$

Independent (\perp projection)

sum of sq's
a) normals.

Fact: $\frac{\chi^2_{m/m}}{\text{indep} - \frac{\chi^2_{n/n}}{\chi^2_{n/n}}} \equiv F_{m,n}$

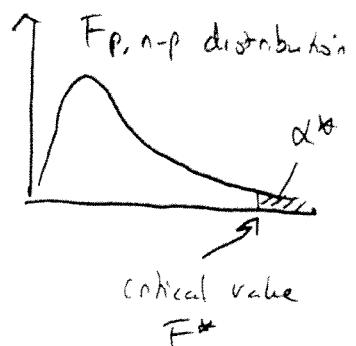
Our test statistic: $F = \frac{\frac{\text{Ssreg}/(p-1)}{\text{SSE}/(n-p)}}{\frac{(SST - SSE)/(n-1 - (n-p))}{\text{SSE}/(n-p)}}$

$$= \frac{\left(\frac{\text{Reduction in SS}}{\text{Reduction in # parameters}} \right)}{\left(\text{MSE of full model} \right)}$$

{ (1) $\beta_1 = 0$, $F \sim F_{p-1, n-p}$

{ (2) $\beta_1 \neq 0$, F is inflated

Hypothesis testing: $\begin{cases} \text{Null: } \beta_1 = 0 \\ \text{Alt: } \beta_1 \neq 0 \end{cases}$ compare F_{observed} to



We reject $\{\beta_1 = 0\}$ @ the α^* -level if

F_{observed} exceeds F^* .

Discussion
Significant is not the same thing as important!
If n is large enough, we reject almost any hypothesis.

Some more distribution theory, assuming $\varepsilon_i \sim N(0, \sigma^2)$

- $\hat{\beta}_1 = \sum k_i y_i$
 - ①) n is large, by CLT $\hat{\beta}_1$ ~ normally distributed
 - ②) $\varepsilon_i \sim N(0, \sigma^2) \Rightarrow y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ $\Rightarrow \hat{\beta}_1 \sim N(\beta_1, \text{Var}(\hat{\beta}_1))$

(also ②) i) $\varepsilon_i \sim N(0, \sigma^2) \Rightarrow y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

$$\Rightarrow \hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}) \quad \text{or} \quad \boxed{\frac{\hat{\beta}_1 - \beta_1}{\frac{\sigma^2}{\sum (x_i - \bar{x})^2}} \sim N(0, 1)}$$

(but) we don't know $\sigma^2 \Rightarrow$ we plug-in estimator $\hat{\sigma}^2$

From before we know $\frac{SSE}{n-2} = \hat{\sigma}^2$ and if ε_i 's are

normal $\frac{SSE}{\hat{\sigma}^2} \sim \chi^2_{n-2}$ or equivalently $(n-2) \frac{\hat{\sigma}^2}{\hat{\sigma}^2} \sim \chi^2_{n-2}$

Fact [Review basic stats]

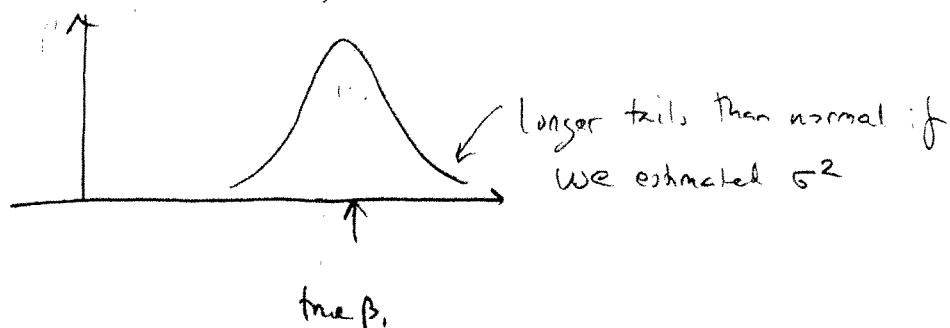
(8)

$$\text{index} - \frac{N(0,1)}{\sqrt{\chi^2_{\nu}/\nu}} \equiv t_{\nu} \Rightarrow \text{(Here)}; \quad \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{G^2 / \sum(x_j - \bar{x})^2}}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{\sum(x_j - \bar{x})^2}}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 / \sum(x_j - \bar{x})^2}} = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} = t\text{-statistic}$$

The standard error

$$\sqrt{V(\hat{\beta}_1)} \leftarrow \hat{\sigma}^2 \text{ estimated.}$$

So, sampling distribution of estimated $\hat{\beta}_1$,



\Rightarrow Back to Basics now!

- Interval estimate $[\hat{\beta}_1 \pm t_{n-p}(1-\alpha/2) \sqrt{\frac{\hat{\sigma}^2}{\sum(x_j - \bar{x})^2}}] = I_{\beta_1}$

- I_{β_1} covers true β_1 w. probability $1-\alpha$.

- Question: Does I_{β_1} cover 0 (the null)?