

STA 302 / 1001

Note Title

9/29/2011

Model:
$$Y_i = \beta_0 + \beta_1 X_i + e_i, \quad i=1, \dots, n$$

Gauss-Markov conditions

$$E(e_i) = 0$$

$$E(e_i e_j) = 0, \quad i \neq j$$

$$\text{Var}(e_i) = \sigma^2 \quad \text{for all } i$$

+ Distributional assumption e_i 's have a Normal distⁿ

So e_i 's iid $N(0, \sigma^2)$

So Distribution of $Y_i | X_i = x_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

Now Distribution of $\hat{\beta}_1$

$$\hat{\beta}_1 | X_1 = x_1, \dots, X_n = x_n = \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\sum (x_i - \bar{x}) Y_i}{S_{XX}}, \quad S_{XX} = \sum (x_i - \bar{x})^2$$

So $\hat{\beta}_1 | X$ is a linear combination of the Y_i 's
 so $\hat{\beta}_1 | X$ has a normal distribution

$$\hat{\beta}_1 | X_1 = x_1, \dots, X_n = x_n \sim N\left(\beta_1, \frac{\sigma^2}{S_{XX}}\right)$$

Estimate $\text{Var}(\hat{\beta}_1 | X)$ by $\frac{S^2}{S_{XX}}$ where $S^2 = \frac{\sum \hat{e}_i^2}{n-2}$

Standard error - estimate of s.d. of a parameter

$$\text{s.e.}(\hat{\beta}_1 | X) = \frac{S}{\sqrt{S_{XX}}}$$

$$\frac{\hat{\beta}_1 - \beta_1}{\frac{S}{\sqrt{S_{XX}}}} \sim t_{n-2}$$

Confidence interval for $\hat{\beta}_1$

$100(1-\alpha)\%$ CI for β_1

$$b_1 \pm t_{\alpha/2, n-2} \frac{S}{\sqrt{S_{xx}}}, \quad S = \sqrt{\text{MSE}}$$

Similarly, $100(1-\alpha)\%$ CI for β_0

$$b_0 \pm t_{\alpha/2, n-2} S \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

Pre-MP data CFC data:

90% CI for β_1 (9.63, 9.88)

Width \uparrow 95% CI for β_1 (9.61, 9.81)

99% CI for β_1 (9.58, 9.85)

Common test in Simple Linear Regression

Want to test: $H_0: \beta_1 = 0$

vs $H_a: \beta_1 \neq 0$

Practically, is there a relationship between X_i & Y_i ?

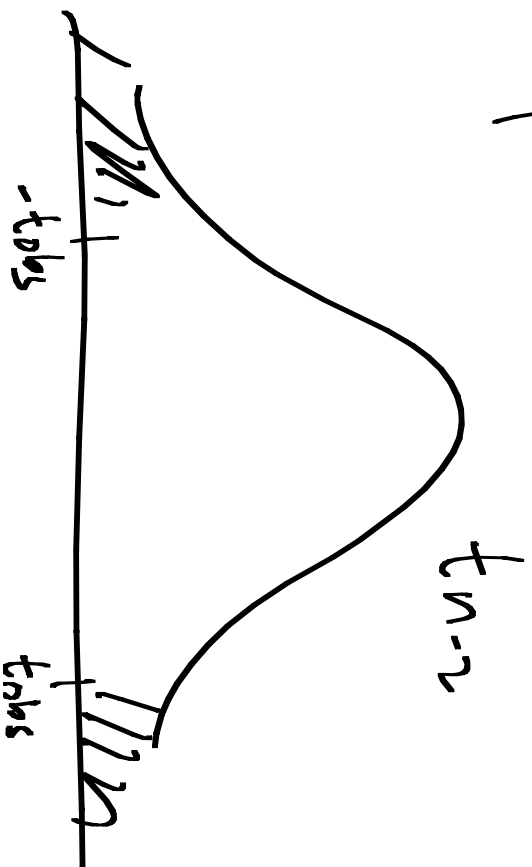
$$\frac{\hat{\beta}_1 - \beta_1}{s / \sqrt{S_{XX}}} \sim t_{n-2}$$

Test statistic for test $H_0: \beta_1 = 0$

$$t_{obs} = \frac{b_1}{s/\sqrt{SXX}}$$

If H_0 is true, t_{obs} is an observation from a t_{n-2} distribution.

Measure strength of evidence against H_0 with p-value



Shaded area = p-value

Pre-MR CFC data $b_1 = 9.71152$

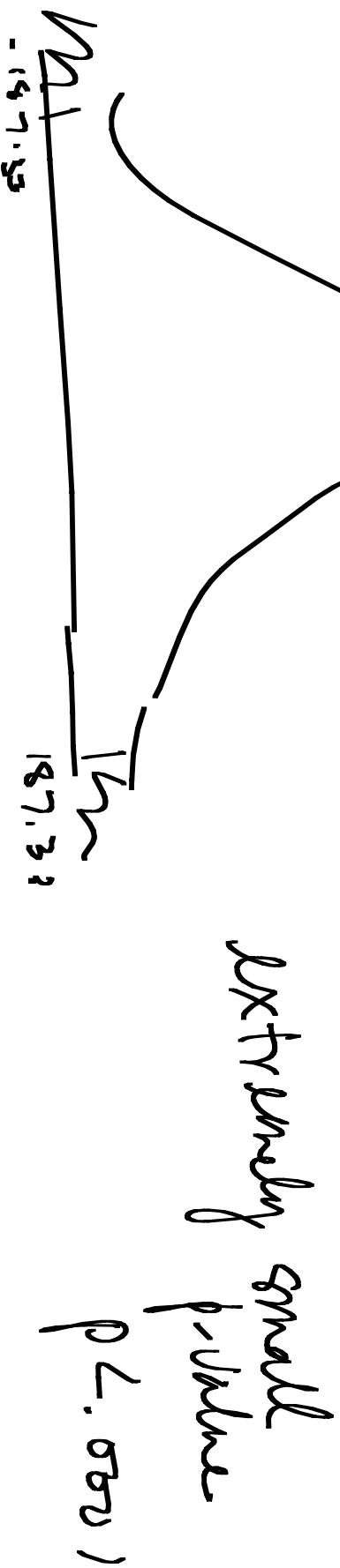
$$\text{s.e.}(\hat{b}_1) = .05184$$

Test statistic for $H_0: \beta_1 = 0$ vs $H_a: \beta_1 \neq 0$

$$t_{\text{obs}} = \frac{9.71152}{.05184} = 187.33$$

$$H_a: \beta_1 > 0$$
$$= 187.33$$

P-value t_{151} ($n=153$)



Very strong evidence that slope is different from 0

95% CI for slope was (9.61, 9.81) \Rightarrow know that
2-sided test for $\beta_1 = 0$ will
have a p-value $< .05$

Practical conclusion: strong evidence that CFC
concentration is changing over time

Suppose want to test $H_0: \beta_1 = 0$ vs $H_a: \beta_1 > 0$
test statistic calculation same



If want to test $H_0: \beta_1 = 0$ vs $H_a: \beta_1 < 0$



But do
 only do
 1-sided
 tests if
 you have
 a ~~really~~
 good
 or priori
 reason.

If want to test $H_0: \beta_1 = \beta_1^*$

$$t_{obs} = \frac{b_1 - \beta_1^*}{s/\sqrt{S_{xx}}}$$

(assume H_0 is true)

Can construct similar tests for β_0

Regression Analysis of Variance

How well does the regression line summarize the data?

Decomposition of Sums of Squares

$$\begin{aligned}y_i &= \hat{y}_i + \hat{e}_i \\ &= b_0 + b_1 x_i + \hat{e}_i \\ &= \bar{y} - b_1 \bar{x} + b_1 x_i + \hat{e}_i \\ y_i - \bar{y} &= b_1 (x_i - \bar{x}) + \hat{e}_i\end{aligned}$$

Square both sides: $(y_i - \bar{y})^2 = b_1^2 (x_i - \bar{x})^2 + 2\hat{e}_i b_1 (x_i - \bar{x}) + \hat{e}_i^2$

Sum:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n \hat{e}_i^2$$

$$\sum_{i=1}^n \hat{e}_i (x_i - \bar{x}) = 0$$

because

$$\sum \hat{e}_i x_i = 0$$

$$\sum \hat{e}_i = 0$$

"Analysis of Variance"
"Decomposition of SS"

SS = Sum of Squares

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

(Corrected SS)

(Uncorrected SS = $\sum y_i^2$)

$\sum_{i=1}^n e_i^2 =$ Residual or Error SS = RSS

- method of least squares minimised this

The term $b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$ is called

model SS or regression SS = SSR_{reg}

If it is the amount of variation in y 's explained by regression line.

Exercise Show $b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 = \sum (y_i - \hat{y})^2$

Merely summarized in ANOVA table:

<u>Source</u>	<u>SS</u>	<u>df</u>	<u>Mean Square</u>
regression	$\sum (y_i - \hat{y})^2$	1	$\frac{\sum (y_i - \hat{y})^2}{1}$
error	$\sum e_i^2$	$n-2$	$\frac{\sum e_i^2}{n-2} = s^2$

<u>Total</u>	<u>$\sum (y_i - \bar{y})^2$</u>	<u>$n-1$</u>	<u>()</u>
--------------	--	-------------------------	------------

leave blank it since doesn't total anymore

Coefficient of Determination

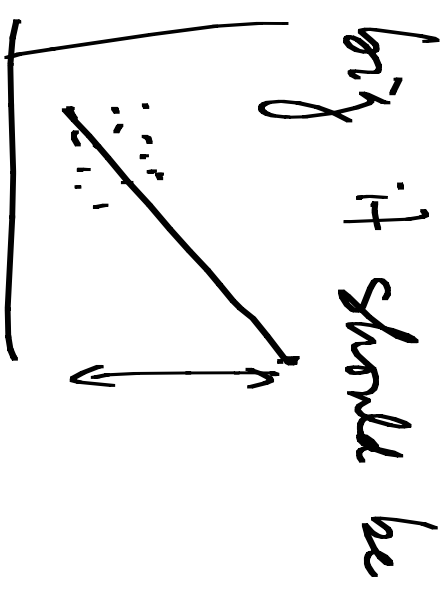
$$R^2 = \frac{SS_{\text{reg}}}{SS_T} = 1 - \frac{RSS}{SS_T}, \quad 0 \leq R^2 \leq 1$$

R^2 gives percent of variation in y 's that is explained by regression line

For MP JCLs data $R^2 = \frac{203119}{203993} = 99.57\%$

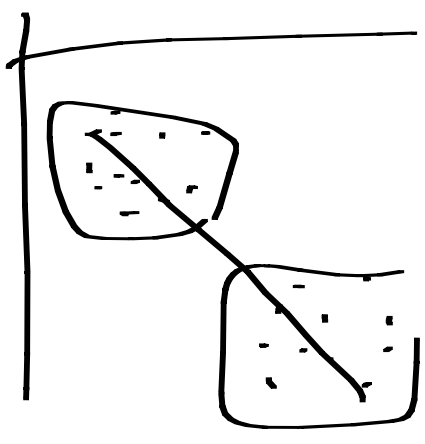
R^2 is useful BUT:

- no absolute rules about how big it should be
- not resistant to outliers
- not meaningful for models with no intercept



- to compare 2 models, R^2 is only useful if:

1. same n (same observations, no transformations)
2. one set of predictor variables is subset of other



- can get very high R^2 by overfitting

(complicated model, may fit well
for data you have but won't
work well on other data)

Some more distribution theory

If $U \sim \text{chi}^2_{\nu_1}$, $V \sim \text{chi}^2_{\nu_2}$,
 U, V independent

$$\text{Then } \frac{U/\nu_1}{V/\nu_2} \sim F_{\nu_1, \nu_2}$$

Treat χ^2 's as non-random

$$\text{Mean Square of Regression} = \text{MS Reg} = \text{SS Reg} / 1 = b_1^2 \sum (x_i - \bar{x})^2$$

Think of MS Reg as estimator $b_1^2 \sum (x_i - \bar{x})^2$

$$E(\text{MS Reg}) = \sigma^2 + b_1^2 S_{xx}$$

MSE "Mean Square Error"

$$= \text{RSS} / n-2 = \sum e_i^2 / n-2$$

$$E(\text{MSE}) = \sigma^2$$

$$\text{If } b_1 = 0, E(\text{MS Reg}) = E(\text{MSE})$$

$$\left. \begin{aligned} E(b_1^2 S_{xx}) &= S_{xx} E(b_1^2) \\ &= S_{xx} [\text{Var}(b_1) + E(b_1)^2] \\ &= S_{xx} \left[\frac{\sigma^2}{S_{xx}} + b_1^2 \right] \end{aligned} \right\}$$

Moreover, If $\beta_1 = 0$, $\frac{MS_{\text{Reg}}}{\sigma^2} \sim \text{chi square } (1)$

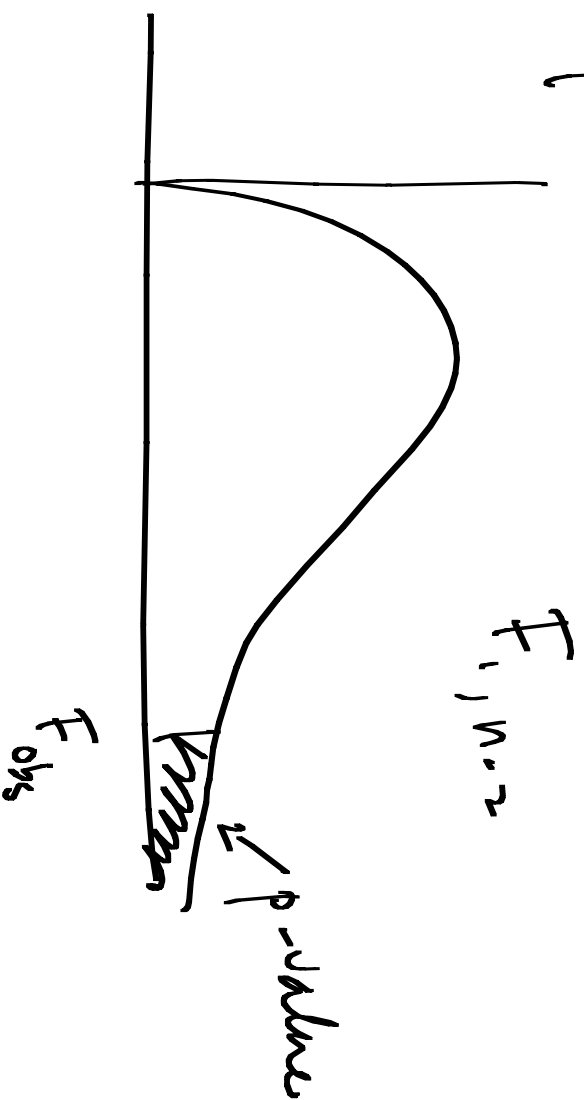
and $\frac{MSE(n-2)}{\sigma^2} \sim \text{chi square } (n-2)$

If $\beta_1 \neq 0$, $\frac{MS_{\text{Reg}}/\sigma^2}{\frac{MSE(n-2)}{\sigma^2/n-2}} \sim F_{1, n-2}$

Another test of $H_0: \beta_1 = 0$ vs $H_a: \beta_1 \neq 0$
Use as test statistic $F_{\text{obs}} = \frac{MS_{\text{Reg}}}{MSE}$

Under H_0 , this is an observation from a F distribution with 1 and $n-2$ degrees of freedom

$b_1 \neq 0$ gives larger values of F_{obs} so deviations from $b_1 = 0$ are in the right tail of F -distribution



Pre-MP CELs data $H_0: \beta_1 = 0$ vs $H_a: \beta_1 \neq 0$

$$F_{obs} = 203119 / 5.788 = 35092.6$$

$$p < .0001$$

Strong evidence that slope is not 0

Note $t_{obs}^2 = (187.33)^2 = 35092.6$

Have 2 tests for $H_0: \beta_1 = 0$ vs $H_a: \beta_1 \neq 0$

Show, in general, $(t_{obs})^2 = F_{obs}$

In general, the square of a r.v. with a t_m distribution results in a r.v. with a $F_{1,m}$ distribution.