

STA 302 / 1001

Note Title

9/22/2011

What to Read? Ch 1 - Motivation
Ch 2

Assignment 1 will be posted soon

Andriy will have QUEST hours in RW 107 / 109

Tuesday Sept. 27 11:00 - 13:00

Tuesday Oct. 4 11:00 - 13:00

(so he won't be here on those 2 Tuesdays.)

Try SAS & the assignment before you ask him for help.

Simple Linear Regression Model

$$Y = \beta_0 + \beta_1 X + e$$

r.v.

random var. $E(e) = 0$

sometimes random

parameters of model
Used to estimate from data

Estimators: $\hat{\beta}_0, \hat{\beta}_1$

Estimates: b_0, b_1

How to estimate β_0, β_1 from observed data:

Data: n pairs $(x_i, y_i) \quad i = 1, \dots, n$

Fitted value for each x_i :

$$\hat{y}_i = b_0 + b_1 x_i \quad (\text{estimate } e \text{ by } \delta)$$

" y_i hat"

Deviations from line: $y_i - \hat{y}_i = e_i$ "residual"

Why interested in vertical errors:

- regression treats x and y differently
- trying to predict y from x

Method of Least Squares

- makes no statistical assumptions

Least Squares: find b_0, b_1 that minimize sum of squares of residuals

Why Least Squares? - mean square error is most common way to measure error in statistics

- Gauss-Markov Theorem
- least squares est. are "best" (min. variance)

Define $RSS =$ "residual sum of squares"

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

$$= \sum (y_i - b_0 - b_1 x_i)^2$$

Find b_0, b_1 to minimize this

$$\frac{\partial \text{RSS}}{\partial b_0} = -2 \sum (y_i - b_0 - b_1 x_i)$$

$$\frac{\partial \text{RSS}}{\partial b_1} = -2 \sum x_i (y_i - b_0 - b_1 x_i)$$

Set equal to 0:

b_0, b_1 must satisfy

"Normal equations"

$$\sum y_i = nb_0 + b_1 \sum x_i \quad \textcircled{1}$$

$$\sum x_i y_i = b_0 \sum x_i + b_1 \sum x_i^2 \quad \textcircled{2}$$

Notation: $\bar{x} = \frac{1}{n} \sum x_i, \quad \bar{y} = \frac{1}{n} \sum y_i$

$$\textcircled{1} \Rightarrow b_0 = \bar{y} - b_1 \bar{x}$$

$$\textcircled{2} \Rightarrow \sum x_i y_i = n \bar{x} \bar{y} + b_1 (\sum x_i^2 - n \bar{x}^2)$$

$$b_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

EXERCISE:
 Show
 that 2
 formulas
 for b_1
 are
 equivalent

$$b_1 = \frac{S_{xy}}{S_{xx}} \quad S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$S_{xx} = \sum (x_i - \bar{x})^2$$

For CFCs Example:

Before MP, $b_0 = -19064$
 $b_1 = 9.71152$

After MP, $b_0 = 3929.6775$
 $b_1 = -1.83289$

Fitted model pre-MP:

$$\hat{y}_i = -19064 + 9.71152 x_i$$

x_i : time in years

Fitted model post-MP:

$$\hat{y}_i = 3929.68 + -1.834 x_i$$

y_i : CFE core in ppt

Estimate time of first manufacture of CFCs:

Pre-MP model, want x when $y=0$, gives ≈ 1963

Actually first produced in (193)

Why a discrepancy? Assumes same growth rate from beginning

- we only know what happened 1977 \rightarrow

• When will atmospheric concentrations of CFCs reach 0?

Post-WP model what x when $y=0$

Not a reasonable question
(would have to assume linear decline continues at

same rate)

• What does model say about 2011? A: nothing

• What is the practical interpretation of intercept?

A: nothing, $x=0$ is not in range of data.

Our models were based on data 1977-1990
1995-2004

We have no idea what happens outside that range

EXTRAPOLATION IS DANGEROUS

Properties of Fitted Regression Line (consequences of least squares) (No statistical assumption)

① Residuals

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

$$= y_i - (b_0 + b_1 x_i)$$

$$= y_i - (b_0 + b_1 \bar{x}) - b_1 x_i$$

$$= (y_i - \bar{y}) - b_1 (x_i - \bar{x})$$

$$\sum e_i = 0$$

So $\bar{e} = 0$ (average of residuals is 0)

$$\sum (x_i - \bar{x}) = 0$$

$$\textcircled{2} \sum \hat{\epsilon}_i^2 \neq 0 \quad (\text{unless data fit perfectly on line})$$

$$= \text{RSS}$$

$$\textcircled{3} \sum \hat{\epsilon}_i x_i = 0 \quad \left\{ \begin{array}{l} \text{EXCURSUS} \end{array} \right.$$

$$\textcircled{4} \sum \hat{\epsilon}_i \hat{y}_i = 0 \quad \left(\sum \hat{\epsilon}_i y_i \neq 0 \right.$$

hard to show)

$$\textcircled{5} \sum \hat{y}_i = \sum (b_0 + b_1 x_i)$$

$$= \sum (\bar{y} - b_1 \bar{x} + b_1 x_i)$$

$$= n\bar{y} - b_1 n\bar{x} + b_1 n\bar{x}$$

$$= n\bar{y} = \sum y_i$$

So far: We assumed that linear model is appropriate

Now some statistical assumptions

$$Y_i = \beta_0 + \beta_1 X_i + e_i$$

Gauss-Markov Conditions

$$\textcircled{1} E(e_i) = 0$$

$$\Rightarrow E(Y_i | X_i = x_i) = \beta_0 + \beta_1 x_i$$

$$\textcircled{2} \text{Var}(e_i) \text{ constant}$$

(same for all
obs. values)

$$\textcircled{3} e_i \text{'s uncorrelated}$$

Gauss-Markov Theorem (without proof)

x_i 's known, fixed

linear model, with Gauss-Markov conditions.

Can show: - least squares estimators are unbiased (will show)
- linear combinations of y_i 's

(Exercise: show $b_1 = \sum d_i y_i$

where d_i function of x_i 's)

Theorem says:

Least squares estimators are BLUE

("Best linear unbiased Estimators")

Best = minimum variance

Model: $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ for x_i not random
estimate of $(\beta_0, \beta_1) = (b_0, b_1)$ functions of
observed data (numbers)

estimators of $(\beta_0, \beta_1) = (\hat{\beta}_0, \hat{\beta}_1)$ functions
of random variables,

estimators of r.v.'s,
can talk about their
statistical properties

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$$

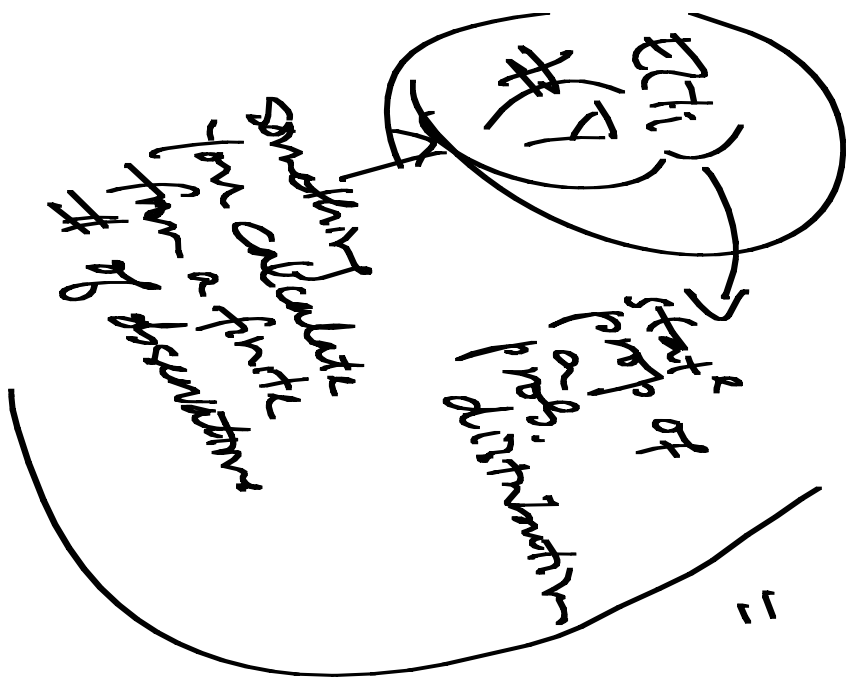
$$\hat{\beta}_1 = \frac{\sum x_i Y_i - n \bar{x} \bar{Y}}{\sum x_i^2 - n \bar{x}^2}$$

Mean and Variance of $\hat{\beta}_1$ Model: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

Going to treat x_i as not random

If x_i 's i.i.d.'s, do everything the same, but
condition on $x_1 = x_1, \dots, x_n = x_n$

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}\right), & S_{XX} &= \sum x_i^2 - n \bar{x}^2 \\ &= \frac{1}{S_{XX}} E\left(\sum x_i y_i - n \bar{x} \bar{y}\right) & &= \sum (x_i - \bar{x})^2 \end{aligned}$$



$$= \frac{1}{S_{XX}} \left[\sum x_i E(Y_i) - n \bar{x} E(\bar{Y}) \right]$$

$$= \frac{1}{S_{XX}} \left[\sum x_i (\beta_0 + \beta_1 x_i) - n \bar{x} (\beta_0 + \beta_1 \bar{x}) \right]$$

$$\begin{aligned} \bar{Y} &= \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) \\ &= \beta_0 + \beta_1 \bar{x} + \bar{\epsilon} \\ E(\bar{\epsilon}) &= 0 \end{aligned}$$

since $E(\epsilon_i) = 0$

$$\begin{aligned}
 \text{So } E(\hat{\beta}_1) &= \frac{1}{S_{XX}} \left[\beta_0 n \bar{x} + \beta_1 \sum x_i^2 - n \bar{x} \beta_0 \right. \\
 &\quad \left. - \beta_1 n \bar{x}^2 \right] \\
 &= \frac{1}{S_{XX}} \beta_1 \left(\sum x_i^2 - n \bar{x}^2 \right) \\
 &= \beta_1
 \end{aligned}$$

So $\hat{\beta}_1$ is an UNBIASED estimator of β_1

Exercise: Show $\hat{\beta}_0$ is an unbiased estimator of β_0

Want: $\text{Var}(\hat{\beta}_1)$

Recall u, v r.v.'s
a, b constants

$$\text{Var}(aU + bV)$$

$$= a^2 \text{Var} U + b^2 \text{Var} V + 2ab \text{Cov}(U, V)$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\sum (x_i - \bar{x})Y_i - \sum (x_i - \bar{x})\bar{Y}}{S_{XX}}$$

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

$$\text{Var}(\hat{\beta}_1) = \text{Var}\left(\frac{\sum (x_i - \bar{x}) y_i}{S_{xx}}\right)$$

(treating x_i as const; otherwise condition on $x_1 = x_1, \dots, x_n = x_n$)

$$= \frac{1}{S_{xx}^2} \sum [(x_i - \bar{x})^2 \text{Var}(y_i)]$$

$$\text{Var}(\sum y_i) = \sum \text{Var}(y_i) ?$$

Yes, because assumed e_i 's uncorr.

$$\underline{y}_i = \beta_0 + \beta_1 x_i + \underline{e}_i$$

$$\text{Var}(y_i) = \text{Var}(e_i)$$

$= \sigma^2$
(σ^2 is the var of the errors)

$$= \frac{1}{S_{XX}} \sum (x_i - \bar{x})^2 \sigma^2$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{XX}} = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

So the more spread out the x_i 's,
the smaller the variance of
the estimator of the slope

Exercise ^{show} $\text{Var}(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}} \right]$

$$\text{Show } \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = - \frac{\sigma^2 \bar{x}}{S_{XX}}$$
