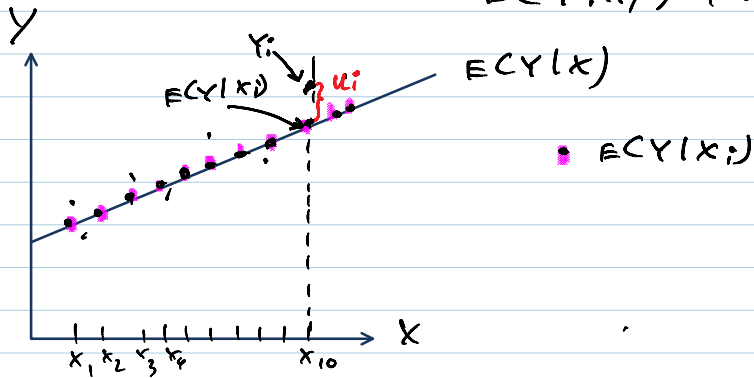


WEEK 5 (7, 9 FEB 2012)

THE SIGNIFICANCE OF THE STOCHASTIC DISTURBANCE TERM

LAST CLASS :  $u_i = Y_i - E(Y|x_i)$   
 OR  $Y_i = \underbrace{E(Y|x_i)}_{\text{SYSTEMATIC}} + \underbrace{u_i}_{\text{NON-SYSTEMATIC}}$



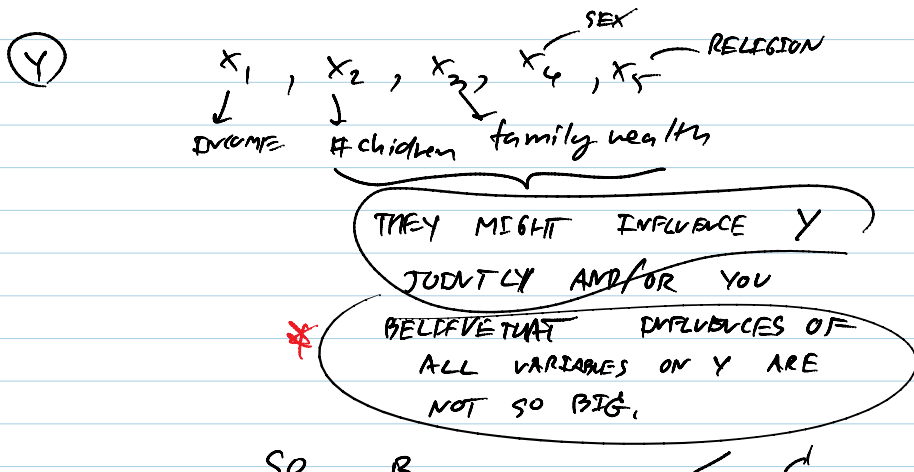
Q: WHAT ARE THE ROLES OF THE ERROR TERM?

A: REASONS ARE MANY:

① VALUES OF THEORY : EVEN THOUGH INCOME IS A CRUCIAL VARIABLE INFLUENCING CONSUMPTION BEHAVIOR, BUT THERE ARE SOME OTHER VARIABLES THAT MIGHT BE IMPORTANT AND WE IGNORE OR OMIT THEM. SO WE USE  $u_i$  AS "A SURROGATE" FOR ALL THOSE VARIABLES (SUBSTITUTE) THAT ARE OMITTED FROM THE MODEL BUT THAT COLLECTIVELY AFFECT  $Y$ .

② UNAVAILABILITY OF DATA : YOU CANNOT INCLUDE SOME VARIABLES THAT ARE IMPORTANT, SIMPLY BECAUSE THE CONSTRAINT OVER DATA AVAILABILITY, EX: FAMILY WEALTH IS GENERALLY UNAVAILABLE OR COSTLY TO GET A PRECISE DATA SET.

③ CORE VARIABLES VS. PERIPHERAL VARIABLES



NOT SO BIG.

SO B INCLUDING THEM < C EXCLUDING THEM

SO, WE INCORPORATE ALL VARIABLES INTO  $u_i$

④ "INTRINSIC RANDOMNESS" IN HUMAN BEHAVIOR

⇒  $Y$  HAS "INTRINSIC RANDOMNESS" THAT CANNOT BE EXPLAINED NO MATTER HOW HARD WE TRY. THE  $u_i$  TERM REFLECTS THIS INTRINSIC RANDOMNESS.

⑤ POOR PROXY VARIABLES.

TAKE MELTON FRIEDMAN'S MODEL AS AN EXAMPLE:

PERMANENT CONSUMPTION ( $Y^p$ ) IS A FUNCTION OF PERMANENT INCOME ( $X^p$ )

PROXY  $Y$  →  $Y^p$  → USES CURRENT CONSUMPTION  
 PROXY  $X$  →  $X^p$  → USES CURRENT INCOME

⇓  
 ERRORS OF MEASUREMENT

AS A RESULT,  $u_i$  REPRESENTS THIS MEASUREMENT ERRORS.

⑥ WRONG FUNCTIONAL FORM?

$$Y_i = a + bX_i + u_i$$

$$Y_i = a + bX_i + cX_i^2 + u_i$$

EVEN IF WE HAVE THEORETICALLY CORRECT VARIABLES TO BE USED IN EXPLAINING  $Y$ , BUT QUITE OFTEN WE DON'T EXACTLY KNOW THE FUNCTIONAL RELATIONSHIP BET.  $Y$  AND  $X(s)$ .

⑦ PRINCIPLE OF PARSIMONY

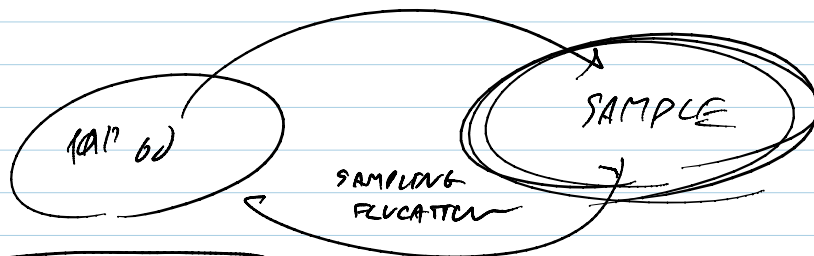
OCCAM'S RAZOR : KEEP THINGS AS SIMPLE AS POSSIBLE.

IF BEHAVIOR OF  $Y$  CAN BE "SUBSTANTIALLY"

EXPLAINED BY TWO OR THREE VARIABLES +  
OUR THEORY IS NOT STRONG ENOUGH TO  
SUGGEST WHAT ELSE SHOULD BE PUT IN,  
WHY PUT MORE?

ANOTHER INTERPRETATION: IF YOU HAVE TWO  
COMPETING THEORIES OR MODELS THAT MAKE  
EXACTLY THE SAME PREDICTION, THE SIMPLER  
ONE IS THE BEST.  
(SIMPLICITY IS A VIRTUE.)

NOW, THE SAMPLE REGRESSION FUNCTION (SRF)



UNKNOWN PARAMETERS

FROM PRF:  $Y_i = a + bX_i + u_i$

$E(Y_i | X_i) = a + bX_i$

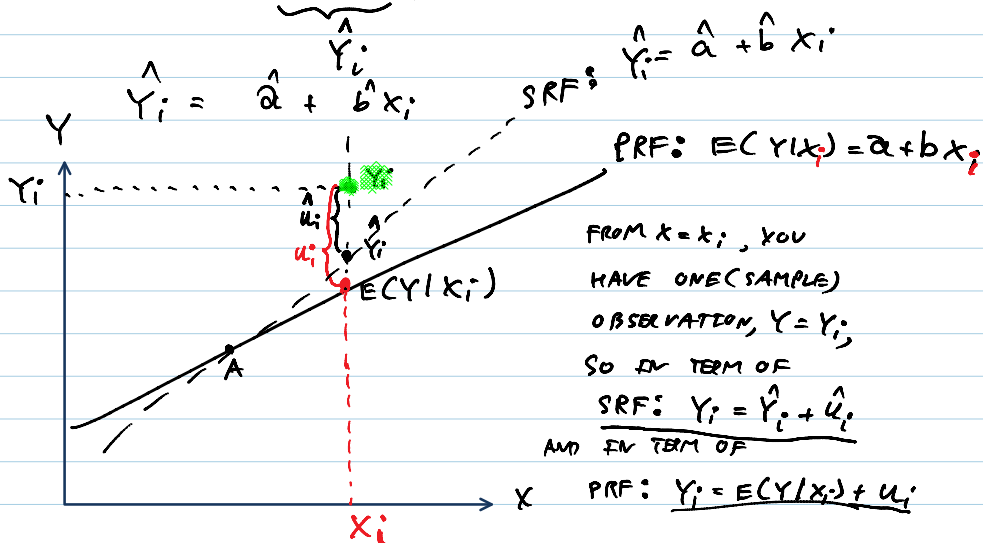
SRF:  $Y_i = \hat{a} + \hat{b}X_i + u_i$

(TRUE)  
UNKNOWN PARAMETERS:

$a, b,$

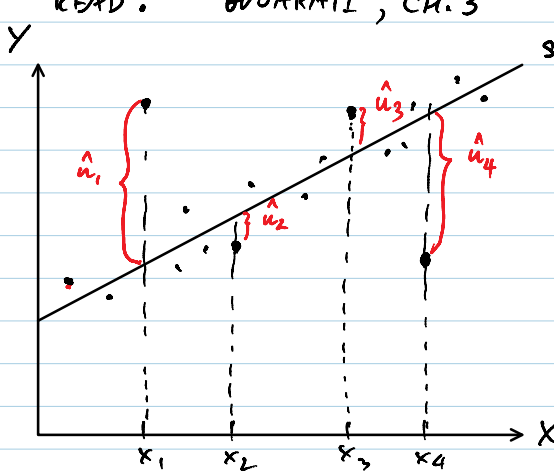
ESTIMATORS:

$\hat{a}, \hat{b},$



(9 FEB 2012)

PRINCIPLES OF ORDINARY LEAST SQUARES



SRF:  $\hat{Y}_i = \hat{a} + \hat{b}X_i$

TWO-VARIABLE REGRESSION MODEL

$Y_i = a + bX_i + u_i$  (PRF)

- $a$  and  $b$  are parameters to be estimated
- $Y$  = dependent variable
- $X$  = independent variable

TASK: WE WANT TO ESTIMATE PRF ON THE BASIS OF SRF AS ACCURATELY AS POSSIBLE.

HOW TO ACCOMPLISH THIS?

- WE WANT TO ESTIMATE  $a$  AND  $b$  FROM SRF:

$Y_i = \hat{a} + \hat{b}X_i + \hat{u}_i$

$\hat{u}_i = Y_i - \hat{Y}_i$

$= Y_i - \hat{a} - \hat{b}X_i$

- WE WANT THE LINE TO BE AS CLOSE TO  $Y_i$  AS POSSIBLE.

- WE WANT TO FIND ESTIMATORS'  $\hat{a}$  AND  $\hat{b}$  SUCH THAT

THE SUM OF SQUARED ERRORS IS THE LEAST

LET  $\hat{u}_i$  BE ESTIMATED ERRORS.

∴ 
$$\left. \begin{aligned} \hat{u}_1 &= Y_1 - \hat{Y}_1 \\ \hat{u}_2 &= Y_2 - \hat{Y}_2 \\ &\vdots \\ \hat{u}_n &= Y_n - \hat{Y}_n \end{aligned} \right\} \begin{array}{l} \text{some are positives} \\ \text{some are negatives} \end{array}$$

$$\sum_{i=1}^n (\hat{u}_i) \quad \text{OR} \quad \sum_{i=1}^n (Y_i - \hat{Y}_i) = 0$$

SUPPOSE  $10, -2, +2, -10 \rightarrow \sum \hat{u}_i = 0$

$\uparrow$                      $\uparrow$                      $\uparrow$                      $\uparrow$   
 $u_1$                      $u_2$                      $u_3$                      $u_4$

$$\sum (\hat{u}_i)^2 \quad \text{OR} \quad \sum (Y_i - \hat{Y}_i)^2$$

$$\sum (u_i) \quad \text{OR} \quad \sum (y_i - \hat{y}_i)$$

- FIND  $\hat{a}$ ,  $\hat{b}$  (ESTIMATORS) SUCH THAT THE SUM OF SQUARED OF THE DIFFERENCE BETWEEN THE ACTUAL  $y$  AND THE ESTIMATED  $\hat{y}$  ( $\hat{y}_i = \hat{a} + \hat{b}x_i$ ) AT LEAST.

$$\therefore \hat{y}_i = \hat{a} + \hat{b}x_i \Rightarrow \hat{u}_i = y_i - \hat{y}_i = y_i - (\hat{a} + \hat{b}x_i)$$

$$\hat{u}_i^2 = (y_i - \hat{a} - \hat{b}x_i)^2$$

$$\sum_{i=1}^n \hat{u}_i^2 = \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i)^2$$

NOTICE THAT  $\sum_{i=1}^n \hat{u}_i^2 = f(\hat{a}, \hat{b})$   
 is a function of

MINIMIZE  $\sum_{i=1}^n \hat{u}_i^2$  WITH RESPECT TO  $\hat{a}$ ,  $\hat{b}$  or depends on

LET 'S DENOTE  $\hat{u}_i = e_i$  FOR CONVENIENCE.

$$\frac{\partial \sum e_i^2}{\partial \hat{a}} = \frac{\partial \sum (y_i - \hat{a} - \hat{b}x_i)^2}{\partial \hat{a}} = 2 \sum (y_i - \hat{a} - \hat{b}x_i)(-1) = 0 = -2 \sum (y_i - \hat{a} - \hat{b}x_i) = 0$$

IT IMPLIES THAT  $\sum (y_i - \hat{a} - \hat{b}x_i) = 0$

OR

$$\sum (y - \hat{y}_i) = 0$$

OR

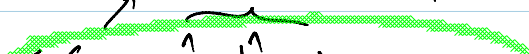
$$\sum e_i = 0$$

ALWAYS.

$$\frac{\partial \sum e_i^2}{\partial \hat{b}} = -2 \sum (y_i - \hat{a} - \hat{b}x_i)x_i = 0$$



ACTUALLY ESTIMATED



o b

IT IMPLIES THAT  $\sum (y_i - \hat{a} - \hat{b}x_i) x_i = 0$

$\nwarrow$  ↑ ↑  
 ACCURATELY ESTIMATED

$\sum (y_i - \hat{y}_i) x_i = 0$

$$\sum e_i x_i = 0$$

THE EXPLANATORY VARIABLE (X) AND ERRORS ( $e_i$ ) MUST BE INDEPENDENT.

FROM  $\sum (y_i - \hat{a} - \hat{b}x_i) = 0 \Leftrightarrow \sum y_i - \sum \hat{a} - \hat{b} \sum x_i = 0$

$$\sum y_i = \sum \hat{a} + \hat{b} \sum x_i$$

$$\sum y_i = n\hat{a} + \hat{b} \sum x_i$$

$$n\hat{a} = \sum y_i - \hat{b} \sum x_i$$

$$\hat{a} = \frac{\sum y_i - \hat{b} \sum x_i}{n}$$

$$= \frac{\sum y_i}{n} - \frac{\hat{b} \sum x_i}{n}$$

$$= \bar{y} - \hat{b} \bar{x}$$

∴

$$\hat{a} = \bar{y} - \hat{b} \bar{x}$$

FROM  $\sum (y_i - \hat{a} - \hat{b}x_i) x_i = 0$

$$\sum x_i y_i - \sum \hat{a} x_i - \sum \hat{b} x_i^2 = 0$$

$$\sum x_i y_i = \hat{a} \sum x_i + \hat{b} \sum x_i^2$$

$$\hat{b} \sum x_i^2 = \sum x_i y_i - \hat{a} \sum x_i$$

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

$$= \sum x_i y_i - \left( \frac{\sum y_i}{n} - \hat{b} \frac{\sum x_i}{n} \right) \sum x_i$$

$$= \sum x_i y_i - \sum x_i \bar{y} + \hat{b} \left( \sum x_i \right)^2 / n$$

$$\hat{b} \sum x_i^2 - \frac{\hat{b} (\sum x_i)^2}{n} = \sum x_i y_i - \frac{\sum x_i y_i}{n}$$

$$\frac{\hat{b} n \sum x_i^2 - \hat{b} (\sum x_i)^2}{n} = \frac{n \sum x_i y_i - \sum x_i y_i}{n}$$

$$\hat{b} n \sum x_i^2 - \hat{b} (\sum x_i)^2 = n \sum x_i y_i - \sum x_i y_i$$

$$\hat{b} = \frac{n \sum x_i y_i - \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

### HOMEWORK

SHOW THAT  $\hat{b} = \frac{\sum_{i=1}^n x_i y_i}{\sum x_i^2}$

LOWERCASE LETTER  
 LOWERCASE LETTER

WHERE

$$x_i = X_i - \bar{X}$$

(DEVIATION OF  $X_i$   
FROM ITS MEAN)

AND  $y_i = Y_i - \bar{Y}$

(DEVIATION OF  $Y_i$   
FROM ITS MEAN)