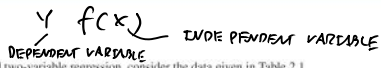




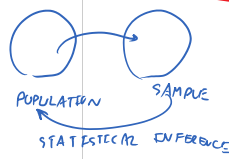
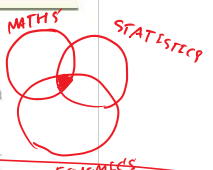
**2. TWO-VARIABLE REGRESSION ANALYSIS**



In order to understand two-variable regression, consider the data given in Table 2.1. The data in the below table refer to a total **Population** of 42 families with their weekly income (X) and weekly consumption expenditure (Y).

Table 2.1: Weekly family Expenditure (Y), Baht and Income (X), Baht

	500	600	700	800	900	1000
Y= Weekly Family Expenditure	360	313	322	310	390	315
	390	400	2800	2870	2820	3710
Conditional means of Y, E(Y X)	350	410	470	530	590	650

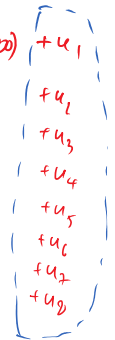


$E(Y) = 9$   
 $E(Y|X=500) = 350$   
 $E(Y|X=900) = 590$   
 $= 498.5714$

① AS INCOME INCREASES,  $E(Y|X)$  INCREASES.

② AT A GIVE LEVEL OF INCOME, SOME FAMILIES SPEND HIGHER THAN ITS OWN (CONDITIONAL MEAN OF Y ( $E(Y|X)$ ), SOME SPEND LESS THAN  $E(Y|X)$

$Y_i = E(Y|X=500) + u_i$   
 $360 = 350 + u_1 = \beta_1 + \beta_2(500) + u_1$   
 $313 = 350 + u_2 =$   
 $322 = 350 + u_3 =$   
 $310 = 350 + u_4 =$   
 $390 = 350 + u_5 =$   
 $315 = 350 + u_6 =$   
 $390 = 350 + u_7 =$   
 $400 = 350 + u_8 =$   
 $= \beta_1 + \beta_2(500)$



$E(u_i | X=500) = 0,$   
 (SEE THE PROOF IN GUJARATI)

Table 2.2: Conditional Probabilities  $p(Y|X)$  for the Weekly Family Income (X) and Expenditure (Y)

	500	600	700	800	900	1000
Y= Weekly Family Expenditure	1/8	1/7	1/6	1/7	1/6	1/8
	1/8	1/7	1/6	1/7	1/6	1/8
	1/8	1/7	1/6	1/7	1/6	1/8
	1/8	1/7	1/6	1/7	1/6	1/8
	1/8	1/7	1/6	1/7	1/6	1/8
	1/8	1/7	1/6	1/7	1/6	1/8
	1/8	1/7	1/6	1/7	1/6	1/8
Conditional means of Y, E(Y X)	350	410	470	530	590	650

Conditional expected value of weekly consumption expenditure given the income level =X,  $E(Y|X)$

$E(Y|X=500) = 350$   
 $E(Y|X=1000) = 650.$

Unconditional expected value,  $E(Y)$

$E(Y) =$

Figure 2.1: Conditional Distribution of Expenditure for Various Levels of Income

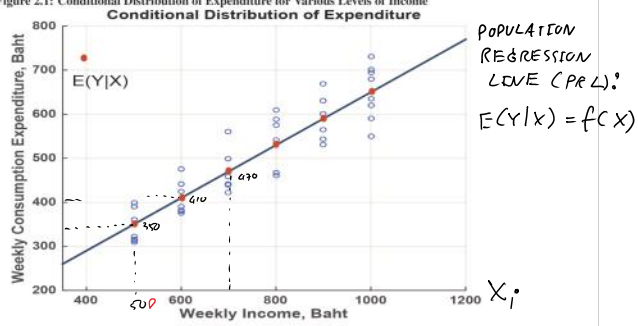
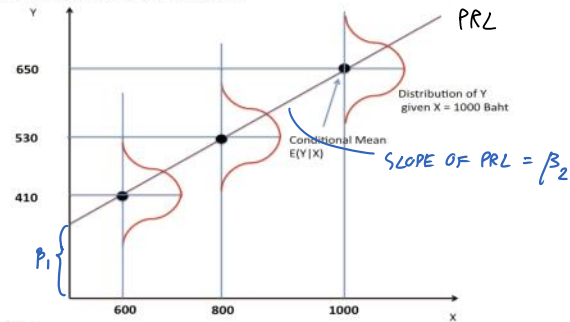


Figure 2.2: Population Regression Line (PRL)



## 2.1 The Concept of Population Regression Function (PRF)

41

## 2.1 The Concept of Population Regression Function (PRF)

The population regression function (PRF) can be written as the function of  $X_i$ :

$$E(Y|X_i) = f(X_i)$$

2.1.1 What form does the function  $f(X)$  assume?If we assume the PRF  $E(Y|X_i)$  is a linear function of  $X_i$ , we get

$$E(Y|X_i) = \beta_1 + \beta_2 X_i$$

 $\beta_1$  = INTERCEPT Y-AXIS OF PRL $\beta_2$  = SLOPE OF THE PRL

$$\frac{\Delta E(Y|X_i)}{\Delta X_i} = \beta_2 = \text{MARGINAL EFFECT OF } X_i \text{ ON } E(Y|X_i)$$

2.1.2 What is the meaning of the term LINEAR?

LINEARITY in the variables

EX:  $E(Y|X_i) = \beta_1 + \beta_2 X_i \rightarrow$  LINEARITY IN VARIABLE

$E(Y|X_i) = \beta_1 + \beta_2 X_i^2 \rightarrow$  NOT LINEARITY IN VARIABLE SINCE  $X_i$  IS RAISED TO THE POWER OF 2

SUMMARY: TO HAVE LINEARITY IN VARIABLES,  $X_i$  MUST BE RAISED TO THE POWER OF 1 ONLY.

LINEARITY in the parameters

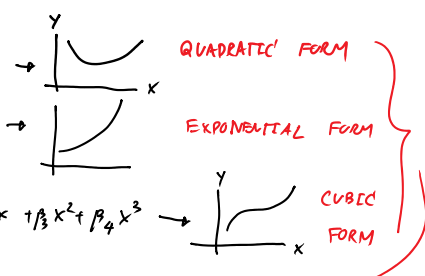
$E(Y|X_i) = \beta_1 + \beta_2 \sqrt{X_i} \rightarrow$  NOT LINEARITY IN VARIABLE

EX:  $E(Y|X_i) = \beta_1 + \beta_2^2 X_i$   
 $E(Y|X_i) = \beta_1 + \sqrt{\beta_2} X_i$   
 $E(Y|X_i) = \beta_1 + \beta_2 \cdot \beta_3 X_i$  } NOT LINEARITY IN PARAMETERS.

SUMMARY: THE TERM "LINEAR REGRESSION MODEL" REFERS TO A MODEL THAT IS "LINEAR IN PARAMETER". IT MAY BE OR MAY NOT BE LINEAR IN VARIABLE.

	LINEAR IN VARIABLES	
LINEAR IN PARAMETERS	Y	N
	Y	N
	LRM	LRM
	NLRM	NLRM

$Y = \beta_1 + \beta_2 X + \beta_3 X^2$   
 $Y = e^{\beta_1 + \beta_2 X}$



LRM = LINEAR REGRESSION MODEL  
 NLRM = NON LINEAR REGRESSION MODEL  
 $Y = \beta_1 + \beta_2 X + \beta_3 X^2 + \beta_4 X^3$

THESE ARE LINEAR-IN-PARAMETER FUNCTIONS!

PAGE 43 STOCHASTIC FORM OF PRF

TAKE A LOOK AT 42 FAMILIES AGAIN... FOR AN INDIVIDUAL FAMILY, WE OBSERVE "A DEVIATION" FROM ITS CONDITIONAL MEAN,  $E(Y|X_i)$ .

LET US CALL "THE DEVIATION" FOR A FAMILY AS  $u_i$ .

EX: AT  $X = 500$ , FOR  $Y_i = 390$ , WE CAN WRITE:

$Y_i = E(Y|X=500) + u_i$   
 $390 = 350 + u_i$   
 $u_i = 390 - 350 = +40$

IN GENERAL,  $u_i = Y_i - E(Y|X_i)$

FOR A GIVEN INCOME LEVEL ( $X_i$ ), AN INDIVIDUAL FAMILY'S WEEKLY EXPENDITURE ( $Y_i$ ) CAN BE DECOMPOSED INTO 2 COMPONENTS

$Y_i = E(Y|X_i) + u_i$   
 (Mean expenditure of all families w/ the same income level) + (RANDOM COMPONENT OR STOCHASTIC COMPONENT)

$Y_i = E(Y|X_i) + u_i$

SYSTEMATIC OR DETERMINISTIC COMPONENT OR STOCHASTIC COMPONENT

2.2.1 The roles of the stochastic disturbance term

1. Vagueness of theory

$Y = f(X) + u$

2.2.1 The roles of the stochastic disturbance term

1. Vagueness of theory

$$Y_i = f(X_i) \rightarrow \text{KEYSIAN CONSUMPTION FUNCTION}$$

THE THEORY ABOVE IS NOT COMPLETE. THERE ARE MANY VARIABLES THAT MAY AFFECT  $Y_i$  BUT WE EXCLUDE THESE VARIABLES FROM THE MODEL.

2. Unavailability of data

IT IS ABOUT DATA AVAILABILITY.

3. Core variables versus peripheral variables

$$Y_i = f(x_1, x_2, x_3, x_4, x_5, \dots, x_n)$$

$$Y_i = f(x_1, x_2, x_3)$$

4. Intrinsic randomness in human behavior

NO MATTER HOW HARD TO EXPLAIN  $Y_i$ , YOU CANNOT "FULLY" EXPLAIN

5. Poor proxy variable

EX: PERMANENT CONSUMPTION =  $f(\text{PERMANENT INCOME})$

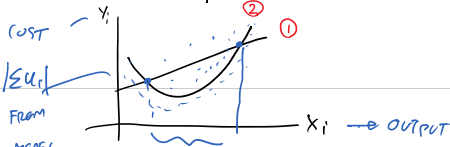
6. Principle of parsimony

"KEEP YOUR REGRESSION MODEL AS SIMPLE AS POSSIBLE"

7. Wrong functional form

$$Y_i = \beta_1 + \beta_2 X_i + u_i \quad \text{--- ① LINEAR}$$

$$Y_i = \beta_1 + \beta_2 X_i + \beta_3 X_i^2 + u_i \quad \text{--- ② NON LINEAR}$$



$| \sum u_i | > | \sum u_i |$   
FROM MODEL ①  
FROM MODEL ②

OR DISTURBANCE COMPONENT  
OR DISTURBANCE COMPONENT

OF THESE VARIABLES

PERIPHERAL VARIABLES

IT DUE TO INTRINSIC RANDOMNESS OF BEHAVIOR.

IF POOR PROXY VARIABLE IS BEING USED, ERROR IN MEASUREMENT MAY ARISE AND IT WILL BE REFLECTED BY THE SIZE OF  $u_i$ .

SUPPOSE THIS IS THE CORRECT FUNCTIONAL FORM

2.3 The Sample Regression Function (SRF)

2.3 The Sample Regression Function (SRF)

As mentioned, in the real situation, we cannot find out all the population of Y values corresponding to the fixed X's. We only have a sample of Y values corresponding to some fixed X's.

Therefore, our goal in this section is to estimate the population regression line (PRF) on the basis of the SAMPLE INFORMATION.

As a result, for the fixed X's as given in table 2.1, we only have a randomly selected sample of Y values. For example, table 2.3 and table 2.4 show a random sample from the population of table 2.1

Table 2.3: A Random Sample From the Population

X	Y
500	390
600	425
700	560
800	575
900	630
1000	679

Table 2.4: Another Random Sample From the Population

X	Y
500	360
600	390
700	440
800	575
900	670
1000	730

POPULATION

$N = 42$   
PRF

$$E(Y|X_i) = f(X_i) = \beta_1 + \beta_2 X_i$$

$\beta_1, \beta_2$ : PARAMETERS  
THEY ARE TRUE PARAMETERS BUT UNKNOWN!

SAMPLE 1

$n = 6$

$$\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$$

SLOPE OF SRF  
INTERCEPT OF SRF

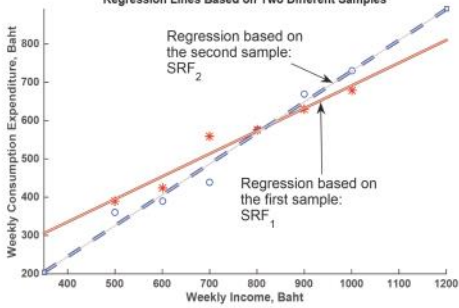
SAMPLE 2

$n = 6$

$$\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$$

STATISTICS

Figure 2.3: Regression lines based on two different samples  
 Regression Lines Based on Two Different Samples



The sample regression function (SRF) can be written as:

$$\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$$

where  $\hat{Y}$  is read as "Y-hat"

$\hat{Y}_i$  = estimator of  $E(Y|X_i)$

$\hat{\beta}_1$  = estimator of  $\beta_1$

$\hat{\beta}_2$  = estimator of  $\beta_2$

We can express the SRF in its stochastic form as follows:

$$Y_i = \beta_1 + \beta_2 X_i + u_i$$

2.3 The Sample Regression Function (SRF)

In sum, our ultimate goal is to estimate the PRF

$N=42$       $E(Y|X_i) = \beta_1 + \beta_2 X_i$

on the basis of the SRF

$n=6$       $Y_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$

STOCHASTIC FORM OF PRF

$$Y_i = E(Y|X_i) + u_i$$

$$= \beta_1 + \beta_2 X_i + u_i$$

SYSTEMATIC     RANDOM

Figure 2.4: Sample and Population Regression Lines

