

LEAST SQUARES REGRESSION

Note Title

7/2/2010

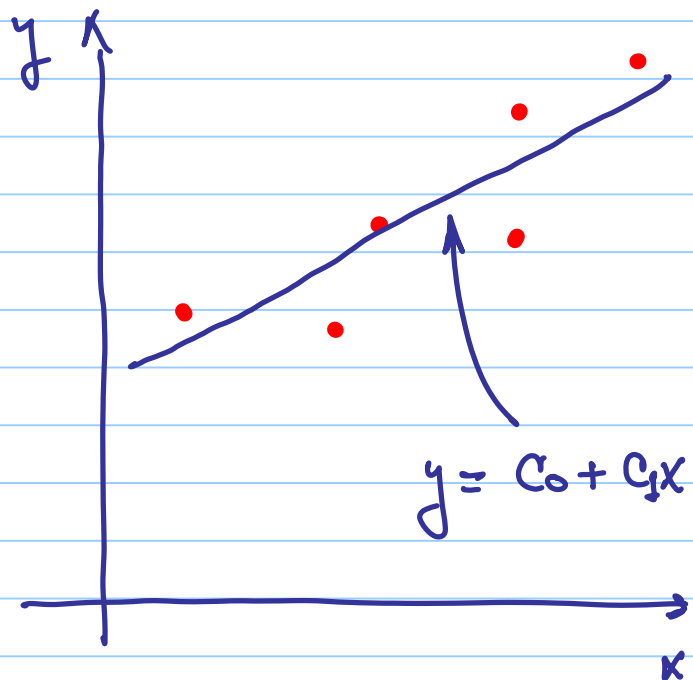
REGRESSION IS THE METHOD OF OBTAINING THE BEST FIT TO A GIVEN SET OF DATA.

THE PROCEDURE OF FITTING THE BEST STRAIGHT LINE (LINEAR EQUATION) TO THE DATA IS KNOWN AS

"LINEAR REGRESSION"

DATA POINTS :

X	y
x_1	y_1
x_2	y_2
x_3	y_3
x_4	y_4
x_5	y_5
\vdots	\vdots
\vdots	\vdots



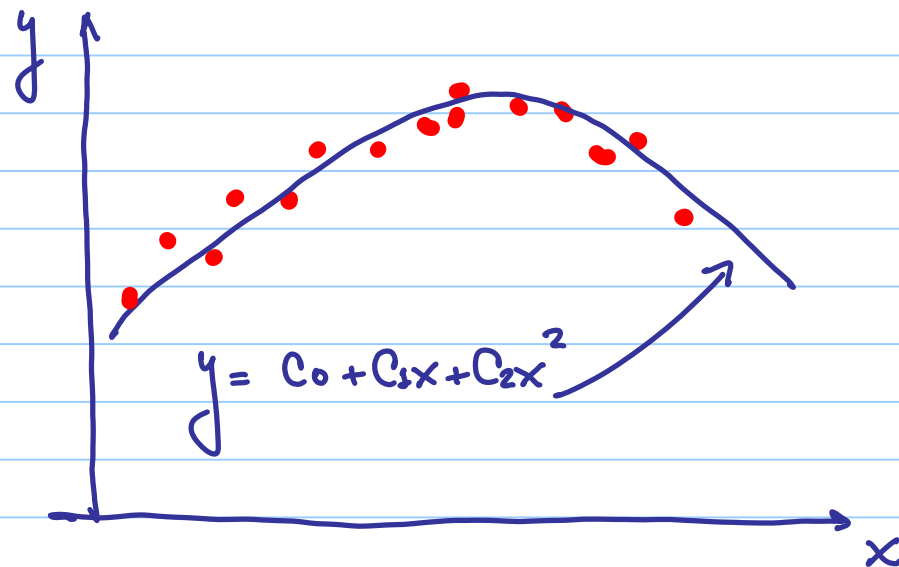
$$y = C_0 + C_1x$$

WE NEED TO FIND C_0 AND C_1 .

IN MANY ENGINEERING APPLICATIONS THE VARIABLES MAY BE RELATED NON LINEARLY.

FOR DATA POINTS EXHIBITING NON-LINEAR BEHAVIOR WE CAN USE

'Polynomial Regression'



POLYNOMIAL PROPOSED :

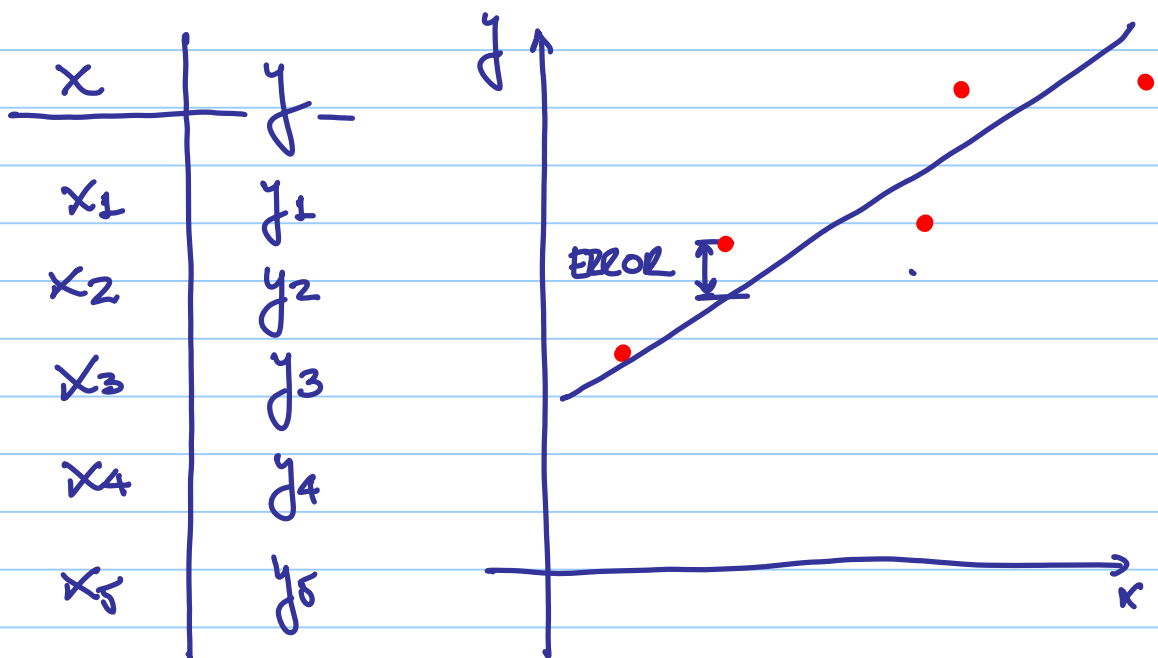
$$\boxed{y = C_0 + C_1x + C_2x^2}$$

WE NEED TO FIND COEFFICIENTS

C_0 , C_1 AND $C_2 \dots$ HOW ??

Well, THE BEST COEFFICIENTS will be THOSE THAT MINIMIZE THE ERROR AT EVERY DATA POINT.

ASSUME WE HAVE THE FOLLOWING DATA :



WE PROPOSE THE LINE

$$y = C_0 + C_1 x$$

SO THAT MEANS THAT THE FOLLOWING EQUATIONS MUST HOLD:

$$y_1 = C_0 + C_1 x_1$$

$$y_2 = C_0 + C_1 x_2$$

$$y_3 = C_0 + C_1 x_3$$

$$y_4 = C_0 + C_1 x_4$$

$$y_5 = C_0 + C_1 x_5$$

WE COULD USE THIS SYSTEM OF LINEAR EQUATIONS TO FIND C_0 AND C_1 ,

HOWEVER, THIS SYSTEM HAS NOT AN EXACT SOLUTION BECAUSE IT IS OVERDETERMINED

IT HAS 5 EQUATIONS AND ONLY TWO UNKNOWN || .
⌒

SINCE NO SOLUTION SATISFIES ALL THE EQUATIONS EXACTLY, WE CAN DEFINE A RESIDUAL ERROR VECTOR AND TRY TO FIND THE COEFFICIENTS C_0 AND C_1 THAT MINIMIZE THIS VECTOR.

USE MATRIX NOTATION TO WRITE THE SYSTEM OF EQUATIONS :

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \\ 1 & x_5 \end{bmatrix} \begin{bmatrix} C_0 \\ C_1 \end{bmatrix}$$

$$\vec{y} = [A] \vec{C}$$

RESIDUAL ERROR VECTOR :

$$\vec{e} = [A] \vec{C} - \vec{y}$$

A commonly used procedure to minimize the error vector involves minimizing the sum of squares of the components of the error vector, known as 'LEAST SQUARE REGRESSION'

THE SUM OF SQUARES OF THE COMPONENTS OF \vec{e} CAN BE EXPRESSED AS

$$S(c, c) = \vec{e}^T \cdot \vec{e} = e_1^2 + e_2^2 + e_3^2 \dots$$

$$S(c, c) = (\vec{c}^T [A]^T - \vec{y}^T) \cdot ([A] \vec{c} - \vec{y})$$

FOR MINIMIZING S , WE SET THE PARTIAL DERIVATIVES OF S WITH RESPECT TO EACH OF THE COMPONENTS OF \vec{c} EQUAL TO ZERO:

$$\frac{\partial S}{\partial c_0} = 0$$

$$\frac{\partial S}{\partial c_1} = 0$$

THIS LEADS TO THE EQUATION

$$[A]^T [A] \vec{c} = [A]^T \vec{y}$$

$$[K] \vec{c} = \vec{b}$$

$$(n \times n) \quad n \times 1 \quad n \times 1$$

→ THIS IS A SYSTEM OF n EQNS. AND n UNKNOWNIS ... HENCE,

it can be solved easily! ☺

Accuracy of the linear regression:

- 1.- Compute first the sum of squares of the deviation of the data around the mean

$$S_0 = \sum_{i=1}^n (y_i - \bar{y})^2$$

\bar{y} denotes the mean value of y_i :

$$\bar{y} = \frac{1}{n} \sum_{j=1}^n y_j$$

- 2.- Compute the sum of squares of the deviation of the data about the straight line (errors)

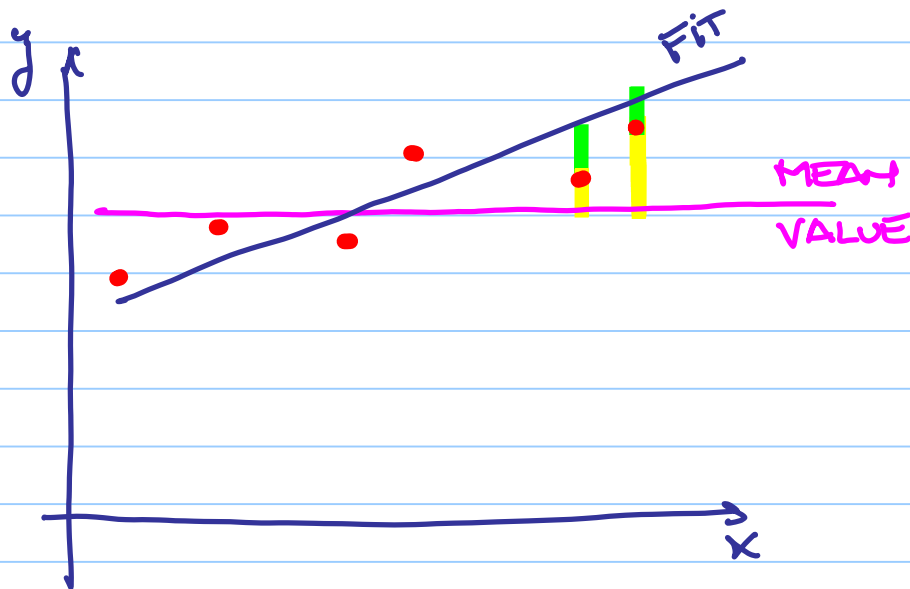
$$S = \sum_{i=1}^n (y_i - (C_0 + C_1 x_i))^2$$

The difference between S_0 and S provides a measure of accuracy of the regression

$$R^2 = \frac{S_0 - S}{S_0}$$

FOR A PERFECT FIT S WILL BE ZERO AND R^2 WILL BE 1.

FOR A POOR FIT $S = S_0$ AND $R^2 = 0$.



IN PRACTICE, A GOOD LEAST SQUARES FIT, IS INDICATED BY A VALUE OF R^2 CLOSE TO 1.

THE SAME APPROACH EXPLAINED HERE CAN BE USED FOR POLYNOMIAL REGRESSION, EVEN FOR OTHER TYPE OF NON-LINEAR REGRESSIONS.

EXAMPLE 1 :

THE NUMBER OF GRADUATE STUDENTS ENROLLED IN THE CIVIL ENGINEERING DEPARTMENT OF FIU OVER A PERIOD OF 10 YEARS IS TABULATED AS FOLLOWS :

YEAR (x_i)	Number (y_i)
1	204
2	216
3	208
4	220
5	240
6	257
7	271
8	270
9	285
10	291

FIT A LINEAR EQUATION TO THE DATA

$$y = C_0 + C_1 x$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ \vdots \\ y_{10} \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \\ \vdots & \vdots \\ 1 & x_{10} \end{bmatrix} \begin{bmatrix} C_0 \\ C_1 \end{bmatrix}$$

\Downarrow \Downarrow \Downarrow
 \vec{y} $[A]$ \vec{C}

EXAMPLE 2 :

THE TEMPERATURE MEASURED AT VARIOUS POINTS ALONG THE THICKNESS OF A WALL SUBJECTED TO A SUDDEN CHANGE IN THE INNER TEMPERATURE YIELDED THE FOLLOWING RESULTS :

x →	x_i (% WALL THICKNESS)	0	0.25	0.5	0.75	1
y →	T_i (°C)	100	70	45	25	15

DEVELOP A SECOND ORDER POLYNOMIAL TO SHOW THE VARIATION OF TEMPERATURE ALONG THE THICKNESS OF THE WALL.

$$y = C_0 + C_1 x + C_2 x^2$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \\ 1 & x_4 & x_4^2 \\ 1 & x_5 & x_5^2 \end{bmatrix} \begin{bmatrix} C_0 \\ C_1 \\ C_2 \end{bmatrix}$$



\vec{C}



\vec{y}



$[A]$