

ChE 3240

MATLAB

Statistics: Regression Analysis

Lecture-02

Least Squares of Linear Regression

For points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ the least square regression can be given by-

$$f(x) = b + mx$$

$$\text{Error, } e_i = y_i - f(x_i)$$

$$\text{Sum of squared error, } SSE = \sum (y_i - f(x_i))^2$$

Therefore,

$$y_1 = (b + mx_1) + e_1$$

$$y_2 = (b + mx_2) + e_2$$

.

.

$$y_n = (b + mx_n) + e_n$$

Now let's set up an matrix equation. Let

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} \quad A = \begin{bmatrix} b \\ m \end{bmatrix} \quad E = \begin{bmatrix} e_1 \\ e_2 \\ \dots \\ e_n \end{bmatrix}$$

matrix equation: $Y = XA + E$.

We now just need to solve this for **A**.

The solution to least squares regression equation
 $Y = XA + E$ is:

$$A = (X^T X)^{-1} X^T Y$$

The sum of the squared errors is:

$$SSE = E^T E$$

$$X^T X = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}$$

$$X^T Y = \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix}$$

Quadratic Regression, in General

- In general, the design matrix and the normal equations for a **quadratic model** based on data $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, are:

Design Matrix

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix}$$

Normal Equations

$$X^T X \mathbf{b} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ x_1^2 & x_2^2 & \dots & x_n^2 \end{bmatrix} \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ x_1^2 & x_2^2 & \dots & x_n^2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = X^T \mathbf{y}$$

Normal Equations

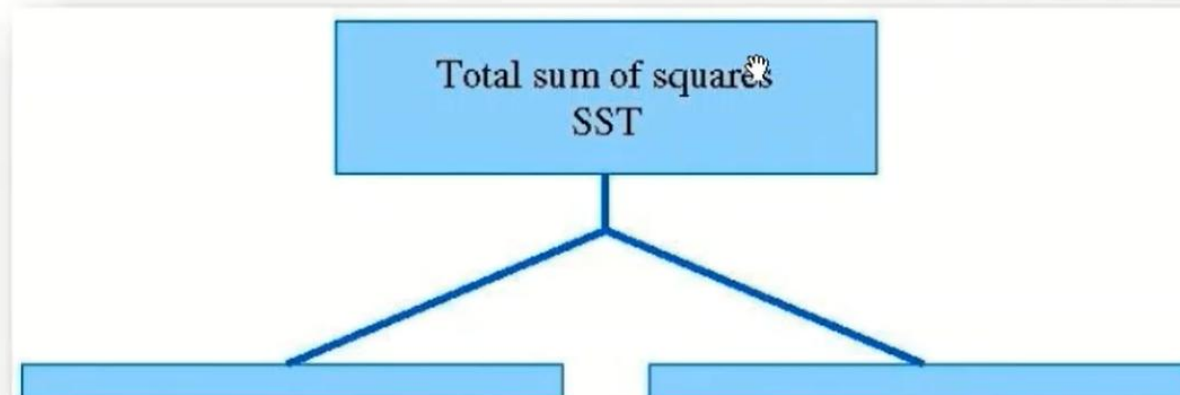
$$X^T X \mathbf{b} = \begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \\ \sum x_i^2 y_i \end{bmatrix} = X^T \mathbf{y}$$

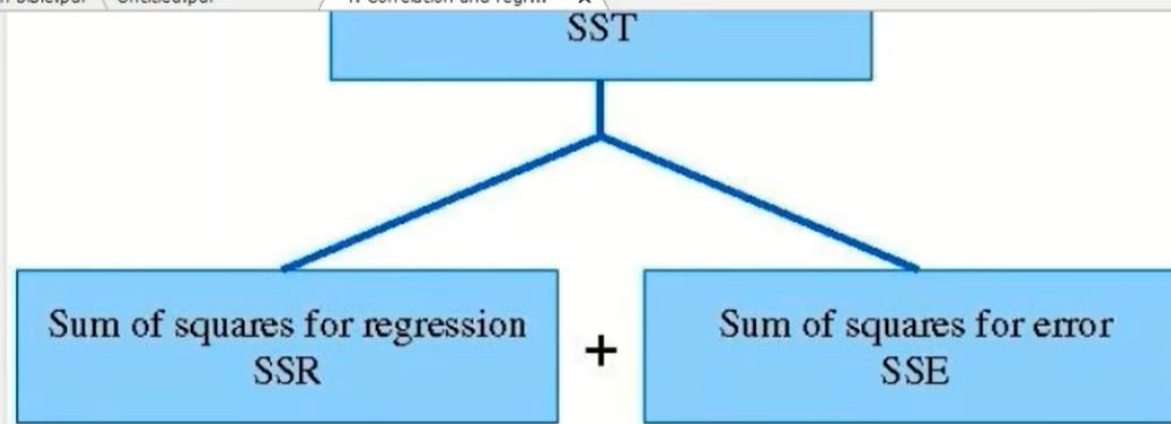
Solution Vector

$$\mathbf{b} = (X^T X)^{-1} (X^T \mathbf{y})$$

Goodness of Fit in Regression and Coefficient of Determination:

$Y_i - \bar{Y}$ is called the total deviation and corresponding $\sum(Y_i - \bar{Y})^2$ is called total sum of squares (SST), $\hat{Y}_i - \bar{Y}$ is called explained deviation and corresponding $\sum(\hat{Y}_i - \bar{Y})^2$ is called sum of squares for regression (SSR) and $Y_i - \hat{Y}_i$ is called unexplained deviation and corresponding $\sum(Y_i - \hat{Y}_i)^2$ is called sum of squares for error (SSE). The relation among SSE, SST, SSR is





$$SST = SSE + SSR$$

Symbolically

$$\sum (Y_i - \bar{Y})^2 = \sum (Y_i - \hat{Y}_i)^2 + \sum (\hat{Y}_i - \bar{Y})^2$$

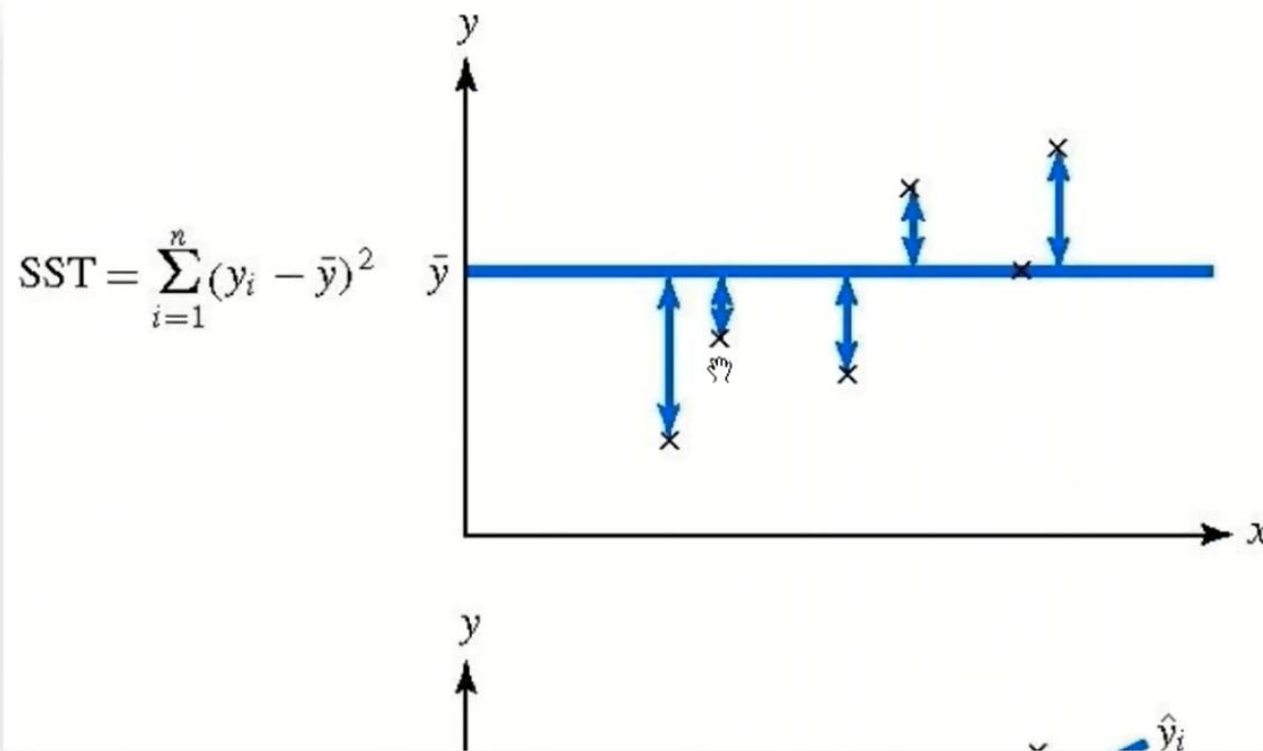
In deterministic relationship $SSE = 0$ i.e. for a perfect fitting estimation line $SST = SSR$ and hence $SSR/SST=1$. For the worst case of data $SSR = 0$ i.e. $SSE = SST$ and hence $SSR/SST=0$

Start

nova-gre-math-bible.pdf

Untitled.pdf

4. Correlation and regr... x



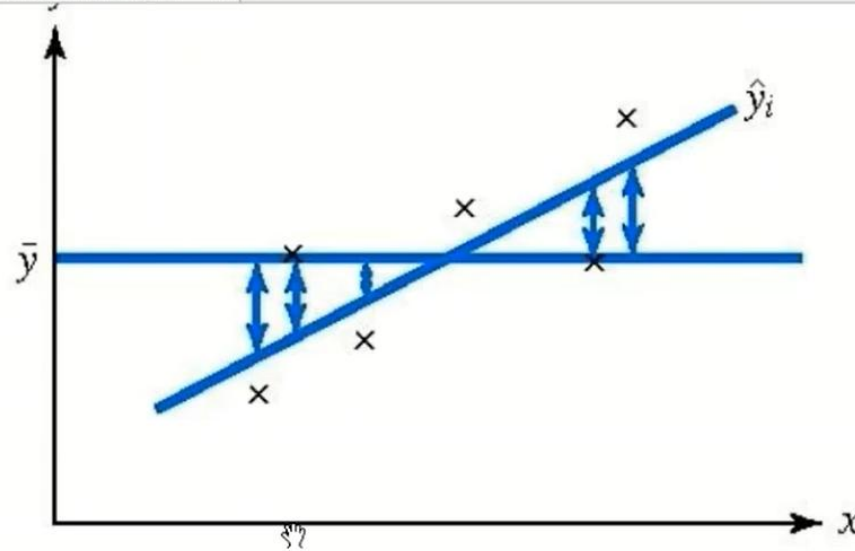
Start

nova-gre-math-bible.pdf

Untitled.pdf

4. Correlation and regr...

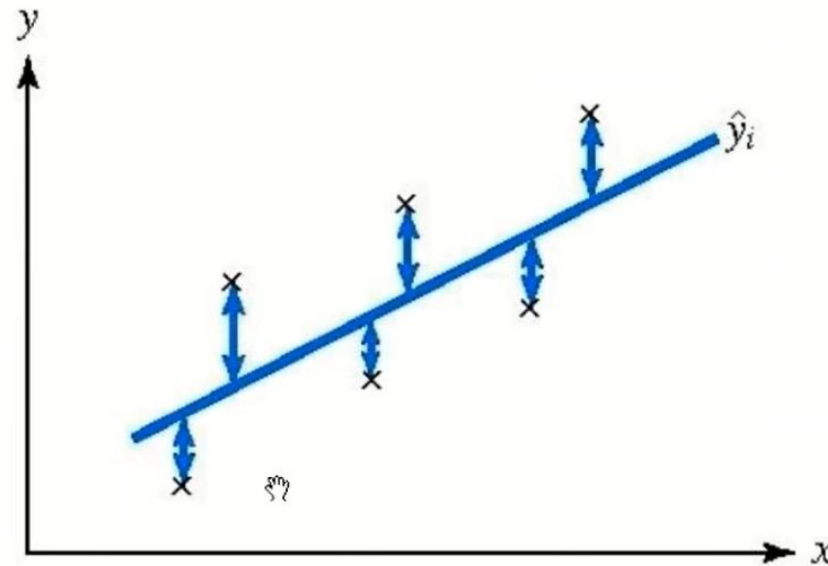
$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$



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



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



So the ratio SSR/SST evaluate how good the estimated regression line is, values of this ratio closer to 1 would imply better fitting estimated line. Thus the ratio SSR/SST is known as the coefficient of determination.

So the ratio SSR/SST evaluate how good the estimated regression line is, values of this ratio closer to 1 would imply better fitting estimated line. Thus the ratio SSR/SST is known as the coefficient of determination.

$$r^2 = \frac{SSR}{SST} = \frac{\sum(\hat{Y}_i - \bar{Y})^2}{\sum(Y_i - \bar{Y})^2} = \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST}$$

Handwritten notes: $= 1 - \frac{SSE/n}{SST/n} = var$ (circled)

r^2 is a non-negative value and it's limit are $0 \leq r^2 \leq 1$. Verbally, r^2 measures the percentage of the total variation in the dependent variable explained by the regression model.