

Linear Regression

(Least Square Error Fit)

TIET, PATIALA

Linear Regression

- In machine learning and statistics, regression attempts to determine the strength and character of the relationship between one dependent variable (usually denoted by Y) and a series of other variables (known as independent variables).
- Mathematically, regression analysis uses an algorithm to learn the mapping function from the input variables to the output variable (Y) i.e. **$Y = f(x)$ where Y is a continuous or real valued variable.**
- Regression is said to be linear regression if the output dependent variable is a linear function of the input variables.

Regression Example

- **House Value Prediction-** The example below shows that the price variable (output dependent continuous variable) depends upon various input (independent) variables such as plot size, number of bedrooms, covered area, granite flooring, distance from city, age, upgraded kitchen, etc.

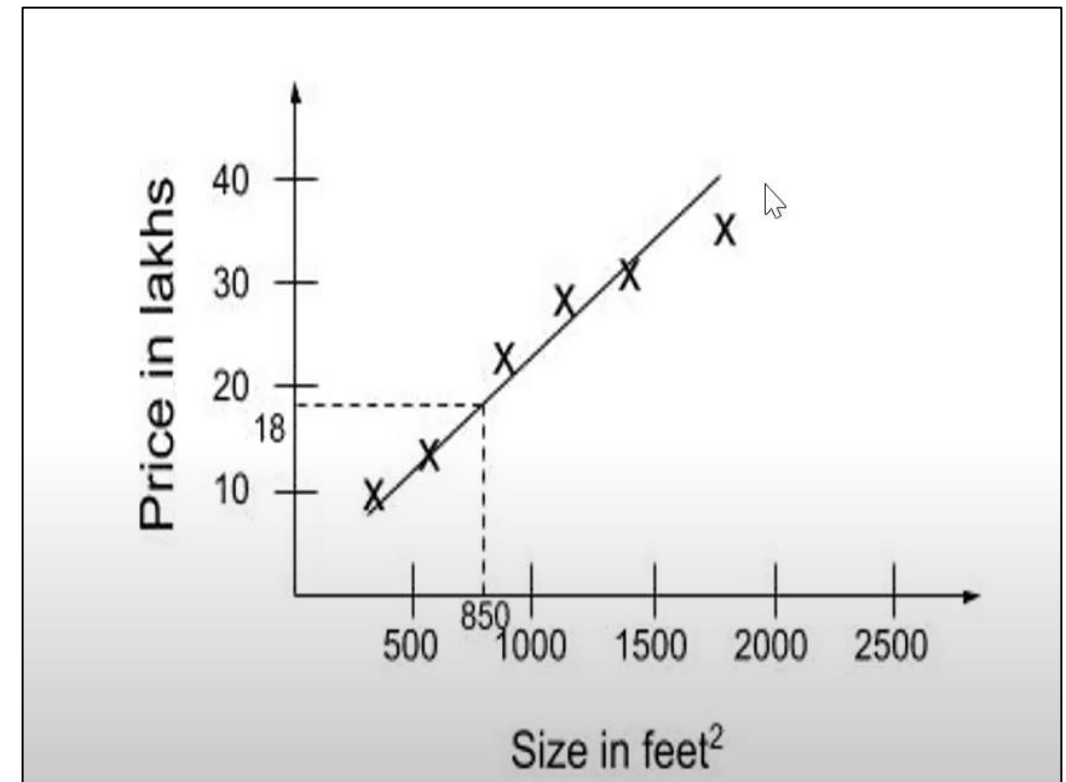
Input Attributes							Output or Class	
Instances	Plot Size	Number of Bedrooms	Covered Area in yards	Granite Flooring	Upgraded Kitchen	Distance from City in Km	Age of flat in years	Price in lakhs
	500	3	150	Y	Y	2	2	70
	1000	2	250	Y	Y	1	1	140
	1800	4	320	N	Y	2	1	200
	300	2	130	Y	Y	3	2	60
	2000	4	500	Y	N	5	3	200
	250	3	160	N	N	1	2	60

Simple Linear Regression (SLR)

- Simple linear regression is a linear regression model with a single explanatory variable.
- It concerns two-dimensional sample points with one independent variable and one dependent variable and finds a linear function (a non-vertical straight line) that, as accurately as possible, predicts the dependent variable values as a function of the independent variable.
- The adjective *simple* refers to the fact that the outcome variable is related to a single predictor.

Simple Linear Regression (SLR) Contd....

- Simple linear regression finds a linear function (a non-vertical straight line) that, as accurately as possible, predicts the dependent variable values as a function of the independent variable.
- For instance, in the house price predicting problem (with only one input variable-plot size), a linear regressor will fit a straight line with x-axis representing plot size and y-axis representing price.



Fitting the Straight Line for SLR

- The linear function that binds the input variable x with the corresponding predicted value of (\hat{y}) can be given by the equation of straight line(slope-intercept form) as:

$$\hat{y} = \beta_0 + \beta_1 x$$

- where β_1 is the slope of line (i.e. it measures change in output variable y with unit change in independent variable x).
- β_0 represents y-intercept i.e. the point at which the line touch x-axis
- \hat{y} is the predicted value of the output for the particular value of input variable x .

Cost/Error function for SLR

- The major goal of SLR model is to fit the straight line that predicts the output variable value quite close to the actual value.
- But, in real world scenario, there is always some error (regression residual) in predicting the values, i.e.

$$\text{actual value}_i = \text{predicted value}_i + \text{error}$$

$$y_i = \hat{y}_i + \epsilon_i$$

$$\text{Residual Error} = \epsilon_i = y_i - \hat{y}_i$$

This error may be positive or negative, as it may predict values greater or lesser than actual values. So we consider **square of each error value**.

Cost/Error function for SLR

- The total error for all the n points in the dataset is given by:

$$\begin{aligned} \text{Total Square Error} &= \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \end{aligned}$$

- The mean of square error is called the cost or error function for simple linear function denoted by $J(\beta_0, \beta_1)$ and given by:

$$J(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

- There exist many methods to optimize (minimize) this cost/error function to **find line of best fit**.

Least Square Method for Line of Best Fit

- The least square method aims to find values $\hat{\beta}_0$ and $\hat{\beta}_1$ for β_0 and β_1 for which the square error between the actual and the predicted values is minimum i.e. least (So, the name is least square error fit).
- The values $\hat{\beta}_0$ and $\hat{\beta}_1$ for β_0 and β_1 for which the square error function ($J(\beta_0, \beta_1)$) is minimum are computed using second derivative test as below:
 1. Compute partial derivatives of $J(\beta_0, \beta_1)$ w.r.t β_0 and β_1 i.e. $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_0}$ and $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1}$
 2. Find values $\hat{\beta}_0$ and $\hat{\beta}_1$ for which $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_0} = 0$ and $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1} = 0$
 3. Find second partial derivative $\frac{\partial^2 J(\beta_0, \beta_1)}{\partial \beta_0^2}$ and $\frac{\partial^2 J(\beta_0, \beta_1)}{\partial \beta_1^2}$; and prove it be minimum for $\hat{\beta}_0$ and $\hat{\beta}_1$.

Least Square Error Fit- Contd.....

$$\text{Total Sqaure Error} = J(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Step 1: Compute partial derivatives of $J(\beta_0, \beta_1)$ w.r.t β_0 and β_1 i.e. $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_0}$ and $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1}$

$$\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (x_i) = -2 \sum_{i=1}^n (x_i y_i - \beta_0 x_i - \beta_1 x_i^2)$$

Step 2: Find values $\hat{\beta}_0$ and $\hat{\beta}_1$ for which $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_0} = 0$ and $\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1} = 0$

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (1)$$

$$\text{and } \sum_{i=1}^n (x_i y_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2) = 0 \quad (2)$$

Least Square Error Fit- Contd.....

From equation 1:

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\beta}_0 - \sum_{i=1}^n \hat{\beta}_1 x_i = 0$$

$$\sum_{i=1}^n y_i - n\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i = 0$$

$$n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \quad (3)$$

From equation 2:

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n \hat{\beta}_0 x_i - \sum_{i=1}^n \hat{\beta}_1 x_i^2 = 0$$

$$\sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0$$

$$\hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \quad (4)$$

Least Square Error Fit- Contd.....

Multiply equation 3 with $\sum_{i=1}^n x_i$ and equation 4 by n

$$n\hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 (\sum_{i=1}^n x_i)^2 = \sum_{i=1}^n x_i \sum_{i=1}^n y_i \quad (5)$$

$$n\hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 n \sum_{i=1}^n x_i^2 = n \sum_{i=1}^n x_i y_i \quad (6)$$

Subtracting Equation 5 from 6, we get,

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

From Equation (3),

$$n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$

$$\hat{\beta}_0 = \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \hat{\beta}_1 \sum_{i=1}^n x_i$$

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

Least Square Error Fit- Contd.....

Step 3: Find second partial derivative $\frac{\partial^2 J(\beta_0, \beta_1)}{\partial \beta_0^2}$ and $\frac{\partial^2 J(\beta_0, \beta_1)}{\partial \beta_1^2}$; and prove it be minimum for $\hat{\beta}_0$ and $\hat{\beta}_1$.

$$\frac{\partial^2 J(\beta_0, \beta_1)}{\partial \beta_0^2} = \frac{\partial(-2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i))}{\partial \beta_0} = -2 \times (-1) = 2$$

$$\frac{\partial^2 J(\beta_0, \beta_1)}{\partial \beta_1^2} = \frac{\partial(-2 \sum_{i=1}^n (x_i y_i - \beta_0 x_i - \beta_1 x_i^2))}{\partial \beta_1} = -2 \times (-x_i^2) = 2x_i^2$$

Therefore, both are positive i.e. the cost function continuously increases with respect to β_0 and β_1 beyond $\hat{\beta}_0$ and $\hat{\beta}_1$. But attains it minimum value at $\hat{\beta}_0$ and $\hat{\beta}_1$

Least Square Error Fit- Summary

- The linear function that binds the input variable x with the corresponding predicted value of (\hat{y}) can be given by the equation of straight line(slope-intercept form) as:

$$\hat{y} = \beta_0 + \beta_1 x$$

- The square error in prediction is minimized when

$$\begin{aligned}\hat{\beta}_1 &= \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \\ &= r_{xy} \frac{\sigma_y}{\sigma_x} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}\end{aligned}$$

$$\text{and } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Least Square Error Fit- Example

The data set (shown in table) gives average masses for women as a function of their height in a sample of American women of age 30–39.

(a) Fit a square line for average mass as function of height using least square error method.

(b) Predict the average mass of women whose height is 1.40 m

Height (m), x_i	Mass (kg), y_i
1.47	52.21
1.50	53.12
1.52	54.48
1.55	55.84
1.57	57.20
1.60	58.57
1.63	59.93
1.65	61.29
1.68	63.11
1.70	64.47
1.73	66.28
1.75	68.10
1.78	69.92
1.80	72.19
1.83	74.46

Least Square Error Fit- Example

i	x_i	y_i	x_i^2	$x_i y_i$
1	1.47	52.21	2.1609	76.7487
2	1.50	53.12	2.25	79.68
3	1.52	54.48	2.3104	82.8096
4	1.55	55.84	2.4025	86.552
5	1.57	57.20	2.4649	89.804
6	1.60	58.57	2.56	93.712
7	1.63	59.93	2.6569	97.6859
8	1.65	61.29	2.7225	101.1285
9	1.68	63.11	2.8224	106.0248
10	1.70	64.47	2.89	109.599
11	1.73	66.28	2.9929	114.6644
12	1.75	68.10	3.0625	119.175
13	1.78	69.92	3.1684	124.4576
14	1.80	72.19	3.24	129.942
15	1.83	74.46	3.3489	136.2618
Total	24.76	931.17	41.0532	1548.2453

Least Square Error Fit- Example

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

$$\hat{\beta}_1 = \frac{15 \times 1548.2453 - 24.76 \times 931.17}{15 \times 41.0532 - 24.76^2} = 61.19$$

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

$$\hat{\beta}_0 = \frac{931.17}{15} - 61.19 \times \frac{24.76}{15} = -38.88$$

Therefore the line of best fit is given by: $\hat{y} = -38.88 + 61.19x$

Predicted value of y when x is 1.4 is

$$\hat{y} = -38.88 + 61.19 \times 1.4 = 46.78$$

Multiple Linear Regression (MLR)

- Multiple regression models describe how a single response variable Y depends linearly on a number of predictor variables.
- Examples:
 - The selling price of a house can depend on the desirability of the location, the number of bedrooms, the number of bathrooms, the year the house was built, the square footage of the lot and a number of other factors.
 - The height of a child can depend on the height of the mother, the height of the father, nutrition, and environmental factors.

Multiple Linear Regression Model

- A multiple linear regression model with k independent predictor variables x_1, x_2, \dots, x_k predicts the output variable as:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots \dots \dots + \beta_k x_k$$

- There is always some error (regression residual) in predicting the values, i.e.

$$\text{actual value}_i = \text{predicted value}_i + \text{error}$$

$$y_i = \hat{y}_i + \epsilon_i$$

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \cdots \dots \dots + \beta_k x_{ik} + \epsilon_i$$

The total error can be computed from all the values in dataset i.e. $i=1, 2, \dots, n$

$$\text{Total Error} = \sum_{i=1}^n \epsilon_i = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^k \beta_j x_{ij}) \quad (7)$$

Multiple Linear Regression Model

- Equation (7) presented in the previous slide, can be represented in matrix form as:

$$\epsilon = y - X\beta$$

- Where $\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$; $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$; $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$

and $X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}$

Least Square Error Fit for MLR

- According to Least Square Error method, we have to find the values of the matrix β for which total square error is minimum.

$$\begin{aligned} \text{Total Square Error} = J(\beta) &= \sum_{i=1}^n \epsilon_i^2 = \epsilon^T \epsilon \\ &= (y - X\beta)^T (y - X\beta) \\ &= (y^T - \beta^T X^T)(y - X\beta) \\ J(\beta) &= y^T y - \beta^T X^T y - y^T X\beta + \beta^T X^T X\beta \\ J(\beta) &= \mathbf{y^T y - 2y^T X\beta + \beta^T X^T X\beta} \end{aligned}$$

[Because $y^T X\beta$ and $\beta^T X^T y$ is always equal with only one entry]

- The square error function is minimized using **second derivative test**.

Least Square Error Fit for MLR

- **Step 1: Compute the partial derivate of $J(\beta)$ w.r.t β**

$$\begin{aligned}\frac{\partial J(\beta)}{\partial \beta} &= \frac{\partial (y^T y - 2y^T X\beta + \beta^T X^T X\beta)}{\partial \beta} \\ &= \frac{\partial y^T y}{\partial \beta} - \frac{\partial 2y^T X\beta}{\partial \beta} + \frac{\partial \beta^T X^T X\beta}{\partial \beta} \\ &= 0 - 2X^T y \frac{\partial \beta}{\partial \beta} + \frac{\partial \beta^T X^T X\beta}{\partial \beta} \\ &\quad [Because \frac{\partial AX}{\partial X} = A^T] \\ &= -2X^T y + 2X^T X\beta \\ &\quad [Because \frac{\partial X^T AX}{\partial X} = 2AX]\end{aligned}$$

Least Square Error Fit for MLR

- **Step 2:** Compute $\hat{\beta}$ for β for which $\frac{\partial J(\beta)}{\partial \beta} = 0$

$$-2X^T y + 2X^T X \hat{\beta} = 0$$

$$X^T X \hat{\beta} = X^T y$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

- **Step 3:** Compute $\frac{\partial^2 J(\beta)}{\partial \beta^2}$ and prove it to be minimum for $\hat{\beta}$

$$\frac{\partial^2 J(\beta)}{\partial \beta^2} = \frac{\partial (-2X^T y + 2X^T X \beta)}{\partial \beta} = 0 + 2 X X^T = +ve$$

Least Square Error Fit for MLR- Example

Example: The Delivery Times Data A soft drink bottler is analyzing the vending machine serving routes in his distribution system. He is interested in predicting the time required by the distribution driver to service the vending machines in an outlet. It has been suggested that the two most important variables influencing delivery time (y in min) are the number of cases of product stocked (x_1) and the distance walked by the driver (x_2 in feet). 3 observations on delivery times, cases stocked and walking times have been recorded.

number of cases of product stocked (x_1)	the distance walked by the driver (x_2)	Delivery time (in min) y
7	560	16.68
3	220	11.50
3	340	12.03

- (a) Fit a multiple regression line using least square error fit.
- (b) Compute the delivery time when 4 cases are stocked and the distance traveled by driver is 80 feet.

Least Square Error Fit for MLR- Example Soln

- The multiple linear regression equation is: $y = \beta_1^{\wedge} + \beta_2^{\wedge}x_1 + \beta_3^{\wedge}x_2$

Where $\beta_1^{\wedge}, \beta_2^{\wedge}, \beta_3^{\wedge}$ or $\beta^{\wedge} = \begin{bmatrix} \beta_1^{\wedge} \\ \beta_2^{\wedge} \\ \beta_3^{\wedge} \end{bmatrix}$ are regression coefficients for line of best fit.

We know, $\beta^{\wedge} = (X^T X)^{-1} X^T y$

$$X = \begin{bmatrix} 1 & 7 & 560 \\ 1 & 3 & 220 \\ 1 & 3 & 340 \end{bmatrix} \text{ and } X^T = \begin{bmatrix} 1 & 1 & 1 \\ 7 & 3 & 3 \\ 560 & 220 & 340 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 3 & 13 & 1120 \\ 13 & 67 & 5600 \\ 1120 & 5600 & 477600 \end{bmatrix}$$

Least Square Error Fit for MLR- Example Soln

$$(X^T X)^{-1} = \begin{bmatrix} 799/288 & 79/288 & -7/720 \\ 79/288 & 223/288 & -7/720 \\ -7/720 & -7/720 & 1/7200 \end{bmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T y = \begin{bmatrix} 7.7696 \\ 0.9196 \\ 0.0044 \end{bmatrix}$$

The line of best fit is, $y = 7.7696 + 0.9196x_1 + 0.0044x_2$

When $x_1 = 4$, $x_2 = 80$

$$y = 7.7696 + 0.9196 \times 4 + 0.0044 \times 80 = 11.80 \text{min}$$