

Applied Statistics I

(IMT224 β /AMT224 β)

Department of Mathematics
University of Ruhuna

A.W.L. Pubudu Thilan

Linear Regression

Introduction

- ▶ The goal of a correlation analysis was to quantify the strength of the linear relationship between the variables, whereas regression expresses the relationship in the form of an equation.
- ▶ Both correlation and linear regression assume that the relationship between the variables is linear.

Example

- ▶ In students taking a Maths and English test, we could use correlation to determine whether students who are good at Maths tend to be good at English.
- ▶ Regression can be used to determine whether the marks in English can be predicted for given marks in Maths.

Regression

- ▶ The purpose of running the regression is to find a formula that fits the relationship between the variables.
- ▶ Then you can use that formula to predict values for the dependent variable when only the independent variables are known.
- ▶ **Eg:** A doctor could prescribe the proper dose based on a person's body weight.

Independent and dependent variables

- ▶ Independent and dependent variables are related to one another.
- ▶ The independent part is what you, the experimenter, changes or enacts in order to do your experiment.
- ▶ The dependent variable is what changes when the independent variable changes.
- ▶ The dependent variable depends on the outcome of the independent variable.

Examples

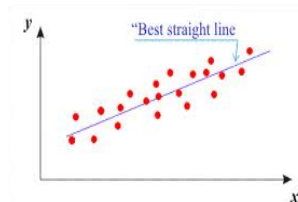
- ▶ The profit made by the manufacturing unit, which is dependent on the sales volumes of the company. Here sales volume is the independent variable and profit is the dependent variable.
- ▶ You are interested in how stress affects heart rate in humans. Your independent variable would be the stress and the dependent variable would be the heart rate.

Extraneous and Confounding Variables

- ▶ The independent and dependent variables are not the only variables present in many experiments.
- ▶ Any variables in your experiment that are not part of your manipulation is called extraneous variables.
- ▶ They are factors you haven't controlled.
- ▶ Extraneous variables affect your results, but usually they affect all your conditions equally and so they do not create any biases in your results.

Simple linear regression

- ▶ When there is only one explanatory (predictor) variable, the prediction method is called **simple regression**.
- ▶ In **simple linear regression**, the predictions of Y when plotted as a function of X form a straight line.



Simple linear regression model

- ▶ The relationship between the dependent variable and the explanatory variable is represented by the following equation:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

β_0 – constant term

β_1 – coefficients of explanatory variable

ϵ_i – error term

The model assumptions

The model

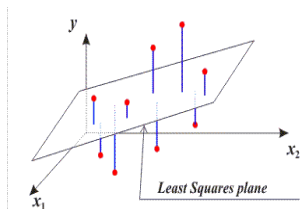
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, 2, 3, \dots, n$$

is used with the following assumptions,

1. $E(\epsilon_i) = 0$.
2. The distribution of the errors in prediction of the value of y is constant regardless of the value of x .
 $\Rightarrow \text{var}(\epsilon_i) = \sigma^2$
3. Errors in prediction of the value of y are all independent of one another.
 $\Rightarrow \text{cov}(\epsilon_i, \epsilon_j) = 0$ for $i \neq j$

Multiple linear regression

- ▶ In multiple linear regression, a linear combination of two or more explanatory variables is used to explain the variation in a response.
- ▶ When there are more than one explanatory variable, the method is quite similar, but instead of a scatterplot in two dimensions, we have to imagine a space with as many dimensions as there are variables.



Multiple linear regression model

The relationship between the dependent variable and the p explanatory variables is represented by the following equation:

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_p x_{pi} + \epsilon_i$$

β_0 – constant term

$\beta_i (i \neq 0)$ – coefficients relating the p explanatory variables

ϵ_i – error term

Remark

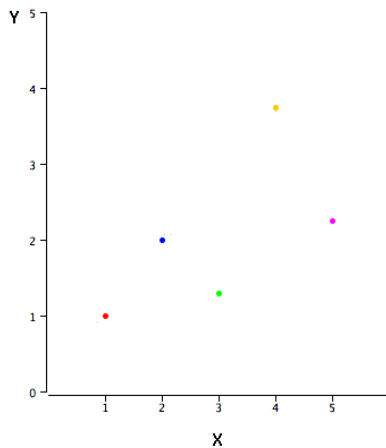
- ▶ Multiple linear regression can be thought of an extension of simple linear regression.
- ▶ Simple linear regression can be thought of as a special case of multiple linear regression, where $p = 1$.
- ▶ The term **linear** is used because in multiple linear regression, we assume that y is directly related to a linear combination of the explanatory variables.

Example

Find the regression line for following data.

X	Y
1.00	1.00
2.00	2.00
3.00	1.50
4.00	3.75
5.00	2.25

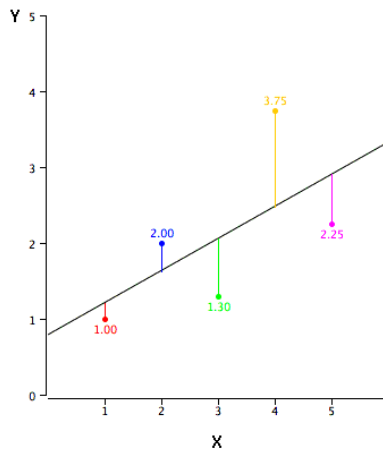
Cont...



Cont...

- ▶ Simple linear regression consists of finding the best-fitting straight line through the points.
- ▶ The best-fitting line is called a **regression line**.
- ▶ The black diagonal line in Figure is the regression line and consists of the predicted score on Y for each possible value of X .

Cont...



Cont...

- ▶ The vertical lines from the points to the regression line represent the errors of prediction.
- ▶ The red point is very near the regression line; its error of prediction is small.
- ▶ The yellow point is much higher than the regression line and therefore its error of prediction is large.

Cont...

- ▶ The error of prediction for a point is the value of the point minus the predicted value (the value on the line).
- ▶ For example, the first point has a y of 1.00 and a predicted \hat{y} of 1.21. Therefore its error of prediction is -0.21.
- ▶ The below Table shows the predicted values (\hat{y}_i) and the errors of prediction ($y_i - \hat{y}_i$).

Cont...

x_i	y_i	\hat{y}_i	$(y_i - \hat{y}_i)$	$(y_i - \hat{y}_i)^2$
1.00	1.00	1.210	-0.210	0.044
2.00	2.00	1.635	0.365	0.133
3.00	1.30	2.060	-0.760	0.578
4.00	3.75	2.485	1.265	1.600
5.00	2.25	2.910	-0.660	0.436

Cont...

- ▶ The most commonly used criterion for the **best fitting line** is the line that minimizes the sum of the squared errors of prediction.
- ▶ That is the criterion that was used to find the line in Figure.
- ▶ The last column in Table shows the squared errors of prediction.
- ▶ The sum of the squared errors of prediction shown in Table is lower than it would be for any other regression line.

Least square method

- ▶ In the least square method, we minimize the sum of square of differences of observed y_i and \hat{y}_i .
- ▶ Let us consider n pairs (x_i, y_i) for $i = 1, 2, 3, \dots, n$.
- ▶ The liner regression model is $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ and we wish to compute a line of the form $\hat{y} = \beta_0 + \beta_1 x$ and is termed as regression line of y on x .

Parameter estimation

- ▶ In the regression line of y on x , that is $\hat{y} = \beta_0 + \beta_1 x$, the unknown parameters β_0 and β_1 should be estimated.
- ▶ As an estimation method we use **least square method**.

Cont...

- ▶ Let

y_i – observed value

\hat{y}_i – estimated value

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

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

- ▶ Our intension is to estimate values of β_0 and β_1 by minimizing S .

Cont...

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

$$S = \sum_{i=1}^n [y_i - \beta_0 - \beta_1 x_i]^2$$

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

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

Cont...

- Equating these partial derivatives to zero we get normal equations.

$$\begin{aligned}\frac{\partial S}{\partial \beta_0} &= 0 \\ \Rightarrow 2 \sum_{i=1}^n [y_i - \beta_0 - \beta_1 x_i] (-1) &= 0 \\ \sum_{i=1}^n y_i - n\beta_0 - \beta_1 \sum_{i=1}^n x_i &= 0\end{aligned}\tag{1}$$

Cont...

$$\begin{aligned}\frac{\partial S}{\partial \beta_1} &= 0 \\ \Rightarrow \sum_{i=1}^n x_i y_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 &= 0\end{aligned}\quad (2)$$

Cont...

$$\begin{aligned}
 & (1) \times \sum_{i=1}^n x_i \Rightarrow \\
 & \sum_{i=1}^n x_i \sum_{i=1}^n y_i - n\beta_0 \sum_{i=1}^n x_i - \beta_1 \left(\sum_{i=1}^n x_i \right)^2 = 0 \quad (3)
 \end{aligned}$$

$$\begin{aligned}
 & (2) \times n \Rightarrow \\
 & n \sum_{i=1}^n x_i y_i - n\beta_0 \sum_{i=1}^n x_i - \beta_1 n \sum_{i=1}^n x_i^2 = 0 \quad (4)
 \end{aligned}$$

Cont...

(4)-(3) \Rightarrow

$$\begin{aligned}
 \hat{\beta}_1 &= \frac{\sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n}}{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}} \\
 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
 \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}
 \end{aligned}$$

Note: $\hat{\beta}_0$ and $\hat{\beta}_1$ are termed as least square estimates . The regression line is $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.

Cont...

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

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

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

Cont...

$$\hat{\beta}_1 = \frac{rS_y}{S_x} \text{ in the case of sample}$$

$$\hat{\beta}_1 = \frac{r\sigma_y}{\sigma_x} \text{ in the case of population.}$$

$$\text{where } \sigma_x = \sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \quad \sigma_y = \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}$$

$$S_x = \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \quad S_y = \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

Example

- (a) Fit a least square line to the following data.
- (b) Estimate y if $x = 7$.

x_i	y_i
1	2
2	5
3	3
4	8
5	7

Solution

x_i	y_i	$x_i y_i$	x_i^2
1	2	2	1
2	5	10	4
3	3	9	9
4	8	32	16
5	7	35	25
15	25	88	55

Cont...

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n}}{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}} \\&= \frac{88 - \frac{15 \times 25}{5}}{55 - \frac{(15)^2}{5}} \\&= 1.3\end{aligned}$$

Cont...

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ &= \frac{25}{5} - 1.3 \frac{15}{5} \\ &= 1.1\end{aligned}$$

The equation of least square line becomes $y = 1.1 + 1.3x$.

The value of y when $x = 7$ is $= 1.1 + 1.3 \times 7 = 10.2$.

Coefficient of determination

- ▶ It is used to check the suitability of the model.
- ▶ The coefficient of determination is defined as,

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

- ▶ If $R^2 \rightarrow 1$ the regression line is appropriate for the given set of data.
- ▶ Otherwise it is necessary to find some other models.

Properties

1. $0 \leq R^2 \leq 1$.
2. $R^2 = r^2$.
3. $R^2 = \frac{\hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$

Proof of (3)

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \text{ by definition.}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

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

$$\hat{y}_i - \bar{y} = \hat{\beta}_1 (x_i - \bar{x})$$

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

Example

The data below give details of x and y such that

x — water content of skew on 1st may in a river

y — water yield in the river after 3-months from that day

The data have been recorded for 10 years.

- (a) Obtain the regression line of y on x .
- (b) Find the coefficient of determination.
- (c) Discuss the suitability of the model.

Cont...

x_i	y_i
23.1	10.5
32.8	16.7
31.8	18.2
32.0	17.0
30.4	16.3
24.0	10.5
39.5	23.1
24.2	12.4
52.5	24.9
37.9	22.8

Solution

x_i	y_i	$x_i y_i$	x_i^2
23.1	10.5	242.55	533.61
32.8	16.7		
31.8	18.2		
32.0	17.0		
30.4	16.3		
24.0	10.5		
39.5	23.1		
24.2	12.4		
52.5	24.9		
37.9	22.8		
328.2	172.4	6044.49	11483.4

Cont...

$$\begin{aligned}
 \hat{\beta}_1 &= \frac{\sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n}}{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}} \\
 &= \frac{6044.49 - \frac{328.2 \times 172.4}{10}}{11483.4 - \frac{(328.2)^2}{10}} \\
 &= \frac{386.322}{711.876} \\
 &= 0.543
 \end{aligned}$$

Cont...

$$\bar{x} = \frac{\sum_{i=1}^{10} x_i}{10} = \frac{328.2}{10} = 32.82$$

$$\bar{y} = \frac{\sum_{i=1}^{10} y_i}{10} = \frac{172.4}{10} = 17.24$$

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ &= 17.24 - (0.543) \times 32.82 \\ &= -0.58\end{aligned}$$

The regression line of y on x is $y = -0.58 + 0.543x$.

Con...

x_i	y_i	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$
23.1	10.5	94.4784	45.4276
32.8	16.7		
31.8	18.2		
32.0	17.0		
30.4	16.3		
24.0	10.5		
39.5	23.1		
24.2	12.4		
52.5	24.9		
37.9	22.8		
328.2	172.4	711.876	240.364

Con...

$$\begin{aligned}\hat{\beta}_1 &= 0.543 \\ R^2 &= \frac{\hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \\ &= \frac{(0.543)^2 (711.876)}{240.367} \\ &= 0.87\end{aligned}$$

The value is closed to 1. Therefore the obtained line is suitable.

Thanks