

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

# **Correlation and linear regression**

# Pearson's Correlation Coefficient (r)

$$n \sum XY - (\sum X) (\sum Y)$$

- $r = \frac{\quad}{\quad}$

$$\sqrt{[n \sum X^2 - (\sum X)^2] [n \sum Y^2 - (\sum Y)^2]}$$

**Correlation** is defined  
as the **degree or strength of relationship between two characteristics in a population**

□ **The aim is**

- ❖ to **investigate the linear association between two**
- ❖ **continuous quantitative variables.**

**Correlation therefore measures the closeness of the association.**

**For the correlation to be obtained we need the followings;**

a. One population

b. two characteristics

c. both should be continuous type (quantitative data)

d. both should be changing (variables) (not constant)

e. There must be some sort of relationship between two

✓ in order to obtain the strength of this relationship

□ After that we need to determine

❖ which of the two variables is X and

❖ which one is Y according to the following;



After that we need to **determine which** of the two **variables is X** and **which one is Y** according to the following;

**X Independent:**

**Y Dependent:**

The change in X is independent on the change in Y

The change in Y is dependent on the change in X

Less changing in a short period of time (more constant)

More changing in a short period of time (more changing)

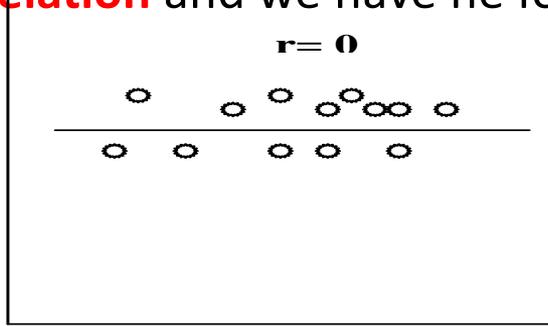
As the cause

As the effect

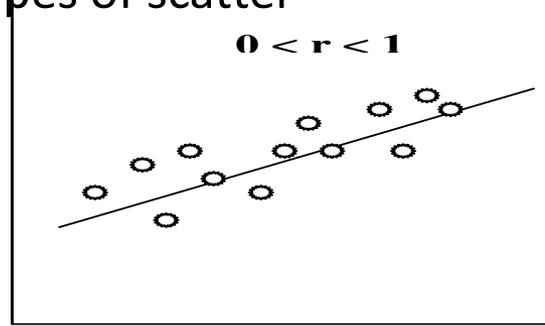
□ After that we need to **draw a scatter diagram** in order to **ascertain** the **presence of correlation** and we have the **following types of scatter**



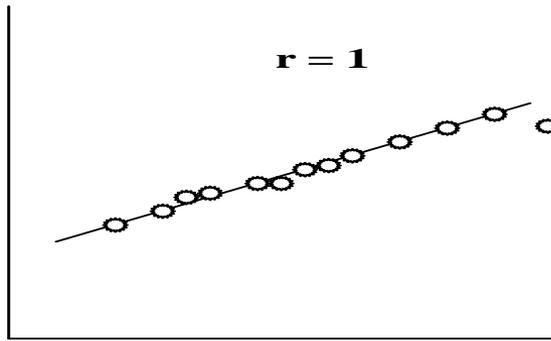
After that we need to **draw a scatter diagram** in order to **ascertain** the **presence of correlation** and we have the following types of scatter



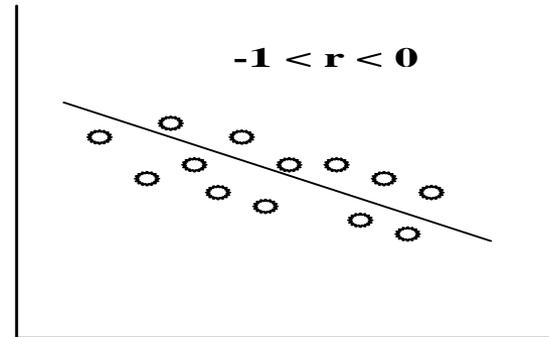
a) No correlation



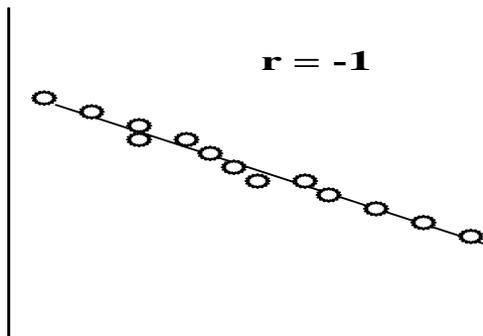
b) Direct positive correlation



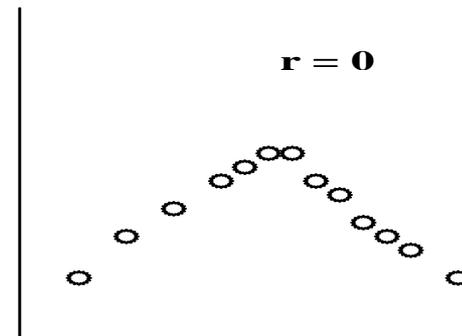
c) Perfect direct positive correlation



d) Inverse negative correlation



e) Perfect inverse negative correlation



f) Strong but non-linear relationship exists

# Pearson's Correlation Coefficient (r)

$$n \sum XY - (\sum X) (\sum Y)$$

- $r = \frac{\quad}{\quad}$

$$\sqrt{[n \sum X^2 - (\sum X)^2] [n \sum Y^2 - (\sum Y)^2]}$$

**r** only measures the linear relationship so we have to draw a scatter diagram first to identify non-linear relationship

## Linear regression;

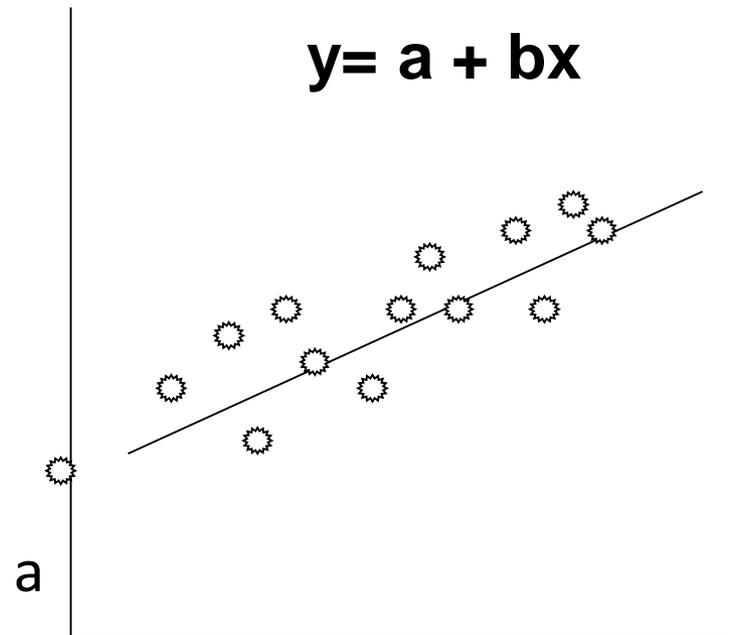
Gives the equation of straight line that best describes it and enables the prediction of one variable from the other.

**a**= constant, called **y-intercept**, it is the place where the regression line intercept with y axis

**b**=regression coefficient

**x**= any value of X variable

**y**= any value of Y variable



## Interpretation of r:

- **r** is always a **number** between **-1 and +1**
- **r** is **positive** if **x and y** tend to **be high or low together**, and the larger its value, the closer the association
- **r** is **negative** if **high value of y** tend to go with **low values of x** and vice versa
- **r** only **measures the linear relationship** so we have to **draw a scatter diagram first** to identify non-linear relationship.

### **Pearson's Correlation Coefficient (r)**

The Pearson correlation coefficient (r) is the most common way of measuring a linear correlation

- **The  $r^2$**  (The coefficient of determination), i.e. when value of  **$r=0.58$ , then  $r^2=0.34$** , this means **that 34% of the variation** in the values of **y may** be accounted for by **knowing values of x** or vice versa

## Pearson's Correlation Coefficient (r)

- **The value of (r) indicates the strength of the relationship**
- **<0.2 : very weak**
- **0.2- <0.4 : weak**
- **0.4- <0.7 : moderate**
- **0.7- <0.9 : strong**
- **$\geq 0.9$  : very strong**

The Pearson correlation coefficient (r) is the most common way of measuring a linear correlation

e.g; The body weight (Kg) and plasma volume (Liter) of 8 healthy men are presented in this table;

<b>N</b>	<b>Body weight (Kg)</b>		<b>Plasma volume (Liter)</b>		
<b>0</b>					
	<b>X</b>	<b>X<sup>2</sup></b>	<b>Y</b>	<b>Y<sup>2</sup></b>	<b>X x Y</b>
1	<b>58</b>	<b>3364</b>	<b>2.75</b>	<b>7.56</b>	<b>159.50</b>
2	<b>70</b>	<b>4900</b>	<b>2.86</b>	<b>8.18</b>	<b>200.20</b>
3	<b>74</b>	<b>5476</b>	<b>3.37</b>	<b>11.36</b>	<b>249.38</b>
4	<b>63.5</b>	<b>4032.25</b>	<b>2.76</b>	<b>7.62</b>	<b>175.26</b>
5	<b>62</b>	<b>3844</b>	<b>2.62</b>	<b>6.86</b>	<b>162.44</b>
6	<b>70.5</b>	<b>4970.25</b>	<b>3.49</b>	<b>12.18</b>	<b>246.05</b>
7	<b>71</b>	<b>5041</b>	<b>3.05</b>	<b>9.30</b>	<b>216.55</b>
8	<b>66</b>	<b>4356</b>	<b>3.12</b>	<b>9.73</b>	<b>205.92</b>
	<b>Σx=535</b>	<b>Σx<sup>2</sup>=35983.5</b>	<b>Σy=24.02</b>	<b>Σy<sup>2</sup>=72.798</b>	<b>Σx.y=1615.292</b>

# Pearson's Correlation Coefficient (r)

$$n \sum XY - (\sum X) (\sum Y)$$

- $r = \frac{\quad}{\quad}$

$$\sqrt{[n \sum X^2 - (\sum X)^2] [n \sum Y^2 - (\sum Y)^2]}$$

$$\sum (X - \bar{X}) (Y - \bar{Y})$$

$$r = \frac{\sum (X - \bar{X}) (Y - \bar{Y})}{\sqrt{[\sum (X - \bar{X})^2 \cdot \sum (Y - \bar{Y})^2]}}$$

$$\sqrt{[\sum (X - \bar{X})^2 \cdot \sum (Y - \bar{Y})^2]}$$

**SP<sub>xy</sub>**

**SP<sub>xy</sub>**

$$r = \frac{\text{SP}_{xy}}{\sqrt{(\text{SS}_x \cdot \text{SS}_y)}} = \frac{\text{SP}_{xy}}{\sqrt{(\text{SQ}_x \cdot \text{SQ}_y)}}$$

**SP = Sum of products of X and Y**

**SS=SQ= Sum of squares of X or of Y**

$$\text{SP}_{xy} = \sum x \cdot y - (\sum x) \cdot (\sum y) / n \Rightarrow 1615.292 - (535 \times 24.02) / 8 \Rightarrow 8.9545$$

$$\text{SQ}_x = \sum x^2 - (\sum x)^2 / n \Rightarrow 35983.5 - (535)^2 / 8 \Rightarrow 0.675$$

$$\text{SQ}_y = \sum y^2 - (\sum y)^2 / n \Rightarrow 72.798 - (24.02)^2 / 8 \Rightarrow 205.375$$

**SP = Sum of products of X and Y**

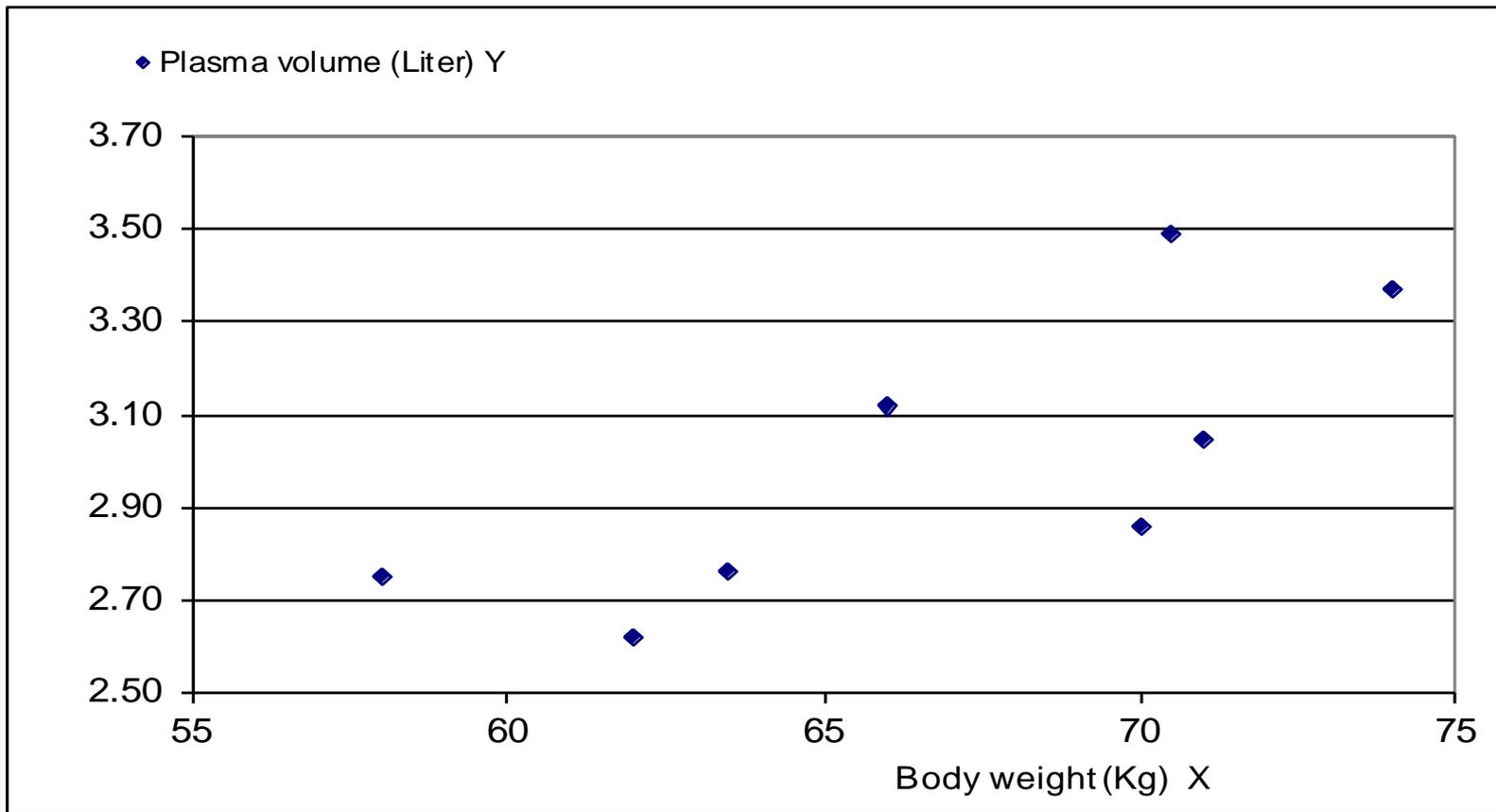
**SS=SQ= Sum of squares of X or of Y**

$$\sum x.y - (\sum x).(\sum y) / n$$

$$r = \frac{\sum x.y - (\sum x).(\sum y) / n}{\sqrt{[\sum x^2 - (\sum x)^2 / n][\sum y^2 - (\sum y)^2 / n]}}$$

$$r = 8.9545 / \sqrt{(0.675 \times 205.375)} \implies + 0.759$$

There is strong direct relationship between weight (Kg) and plasma volume (Liter)



Scatter diagram of plasma volume and body weight showing linear regression line

# Regression:

## Regression:

Regression: it is the best fit for the relationship between two characteristics in a population.

Regression Line: The best line that fits the relationship between two characteristics in a population. Usually determined by the equation of straight line of the first degree which is  $\rightarrow y=a+bx$

$Y \rightarrow$  any value on y axis, the dependent variable

$a \rightarrow$  constant, the intercept, the value of y when x equals to zero, it is the distance between the X axis and the point at which the regression line or its extension cuts the Y axis.

$.b = \text{regression coefficient} = \frac{SP_{xy}}{SSX} \rightarrow \frac{SP_{xy}}{SQx}$

$X \rightarrow$  any value on X axis, the independant variable

We can get the equation according to the following

$$Y = a + bX$$

$$\bar{Y} = \sum y/n = 24.02/8 \rightarrow 3.0025$$

$$\bar{X} = \sum x/n = 535/8 \rightarrow 66.875$$

$$b = SP_{xy}/SQ_x = 8.9545/205.375 \rightarrow 0.0436$$

$$Y = a + bX$$

$$3.0025 = a + 0.0436 \times 66.875$$

$$3.0025 = a + 2.916$$

$$a = 3.0025 - 2.916 \rightarrow 0.0865$$

$$Y = 0.0865 + 0.0436 X$$

Regression coefficient “b” means that; for each one unit change in x axis there is about (b) unit change in y axis

**For each one Kg increase in weight, there is about 0.0867 Liter increase in plasma volume**

$$Y = a + bX$$

$$3.0025 = a + 0.0436 \times 66.875$$

$$3.0025 = a + 2.916$$

$$a = 3.0025 - 2.916 = 0.0865$$

$$Y = 0.0865 + 0.0436 X$$

# Thank you for attention

year 1<sup>ST</sup>  
medical students

