# Correlation (Pearson & Spearman) & Linear Regression

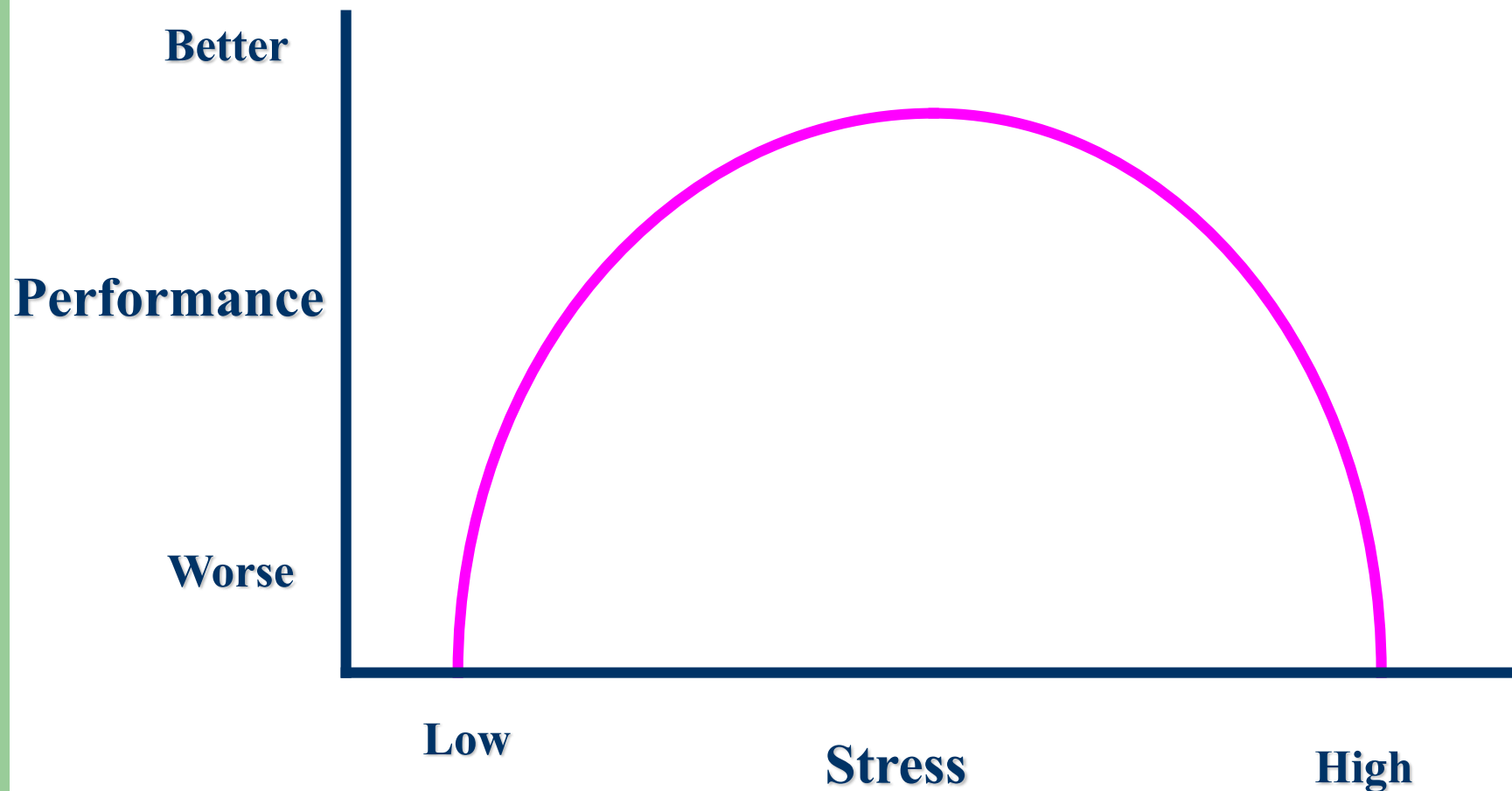**Azmi Mohd Tamil**

# Key Concepts

- Correlation as a statistic
- Positive and Negative Bivariate Correlation
- Range Effects
- Outliers
- Regression & Prediction
- Directionality Problem
- Third Variable Problem (& partial correlation)

# Assumptions

- Related pairs
- Scale of measurement. For Pearson, data should be interval or ratio in nature.
- Normality
- Linearity
- Homocedasticity

# Example of Non-Linear Relationship Yerkes-Dodson Law – not for correlation

# Correlation

X ⟶ Y

**Stress**          **Illness**

# Correlation – parametric & non-para

- **2 Continuous Variables - Pearson**
  - **linear relationship**
  - **e.g., association between height and weight**

- **1 Continuous, 1 Categorical Variable (Ordinal) Spearman/Kendall**
  - **e.g., association between Likert Scale on work satisfaction and work output**
  - **pain intensity (no, mild, moderate, severe) and dosage of pethidine**

# Pearson Correlation

- **2 Continuous Variables**
    - **linear relationship**
    - **e.g., association between height and weight, +**

- measures the degree of linear association between two interval scaled variables

- analysis of the relationship between two quantitative outcomes, e.g., height and weight,
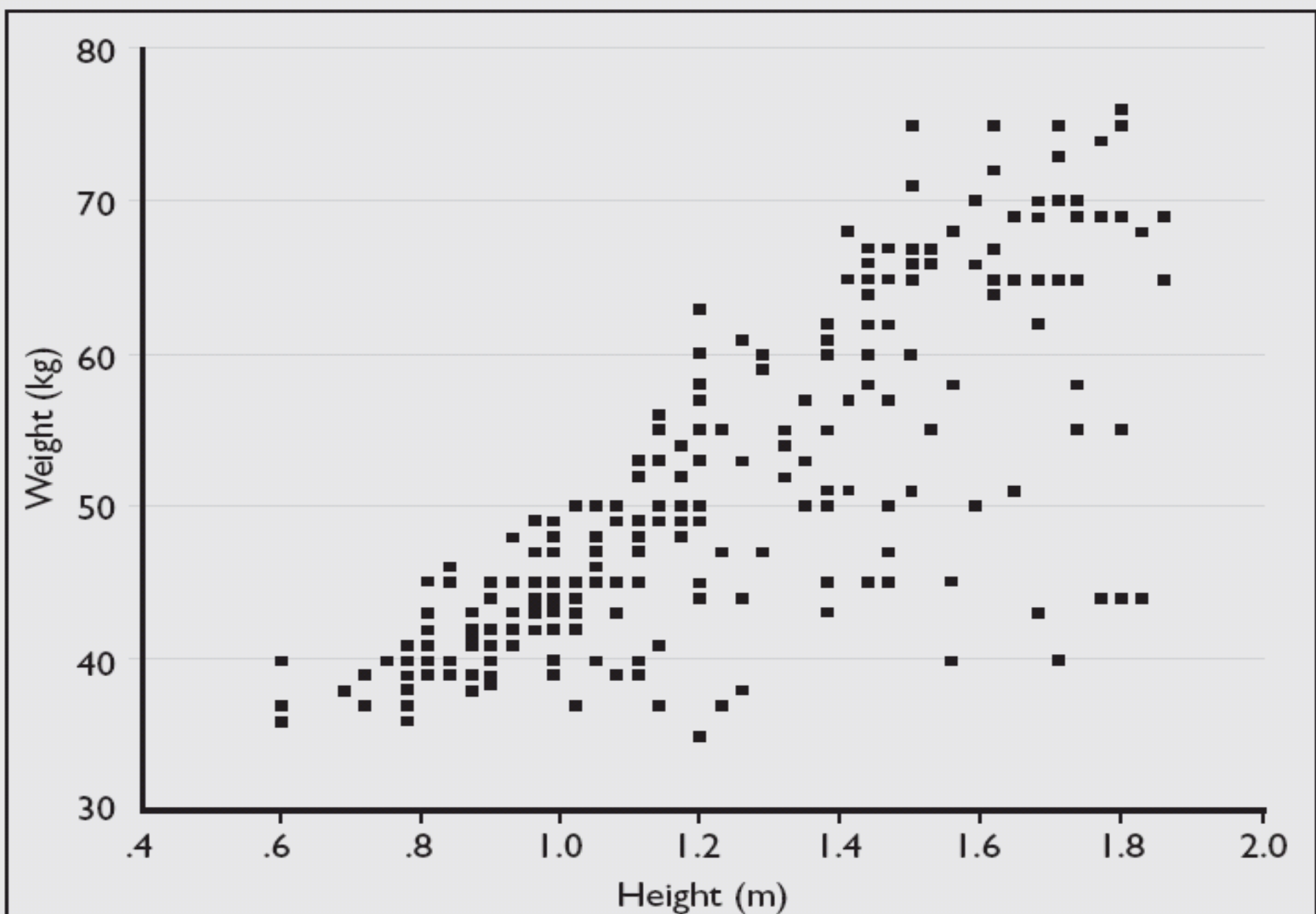
# History of Pearsons' Correlation

- Sir Francis Galton was studying the relationship between the height of the fathers and the height of their sons and discovered a way to mathematically measure this relationship. He called it the "co-efficient of correlation." He gave a specific formula for computing this number from the data he collected. Galton died in 1911. It was his disciple, Karl Pearson, who first formulated the idea in its most complete form in 1895.

- In 1915, Pearson introduced R.A. Fisher to the difficult problem of determining the statistical distribution of **Galton's correlation co-efficient**. Fisher thought about the problem, cast it into a geometric formulation, and within a week had a complete answer. He submitted it for publication in Biometrika; but Pearson & William Sealy Gosset had difficulty understanding the paper. Pearson got his workers to check the calculations. In every case, they agreed with Fisher's more general solution.

# History of Pearsons' Correlation

- Please note that Pearson stated it as **Galton's correlation co-efficient** not **Pearson's correlation co-efficient** to R.A. Fisher. However it is now known as **Pearson's correlation co-efficient** .

- This is an example of what Stephen Stigler, a contemporary historian of science, calls the law of misonomy, that nothing in mathematics is ever named after the person who discovered it. Sir Francis Galton was the one who came out with the co-efficient of correlation theory but Karl Pearson's was the one credited for it.

**Fig. 1** Relationship between height and weight.

# How to calculate r?

$$r = \cfrac{\Sigma XY - \cfrac{\Sigma X \Sigma Y}{N}}{\sqrt{(\Sigma X^2 - \cfrac{(\Sigma X)^2}{N}) \; (\Sigma Y^2 - \cfrac{(\Sigma Y)^2}{N})}}$$

$$t = r\sqrt{\cfrac{n-2}{1-r^2}}$$

$$df = n_p - 2$$

# How to calculate r?

# Example

$$r = \dfrac{\Sigma XY - \dfrac{\Sigma X \Sigma Y}{N}}{\sqrt{\left(\Sigma X^2 - \dfrac{(\Sigma X)^2}{N}\right)\left(\Sigma Y^2 - \dfrac{(\Sigma Y)^2}{N}\right)}}$$

- $\Sigma$ x = 4631　　$\Sigma$ x$^2$ = 688837
- $\Sigma$ y = 2863　　$\Sigma$ y$^2$ = 264527
- $\Sigma$ xy = 424780　　n = 32

- a=424780-(4631*2863/32)=10,450.22
- b=688837-4631$^2$/32=18,644.47
- c=264527-2863$^2$/32=8,377.969
- r=a/(b*c)$^{0.5}$
  =10,450.22/(18,644.47*83,77.969)$^{0.5}$
  =0.836144

$$t = r\sqrt{\dfrac{n-2}{1-r^2}}$$

- t= 0.836144*((32-2)/(1-0.836144$^2$))$^{0.5}$
t = 8.349436 & d.f. = n - 2 = 30,
p < 0.001

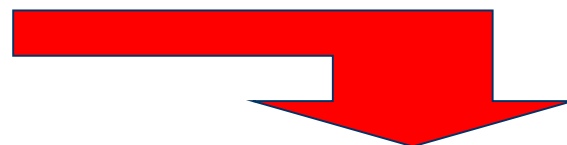| nores | bps1 | bpd1 | x2 | y2 | xy |
|---|---|---|---|---|---|
| 234 | 118 | 67 | 13924 | 4489 | 7906 |
| 235 | 126 | 76 | 15876 | 5776 | 9576 |
| 238 | 105 | 68 | 11025 | 4624 | 7140 |
| 240 | 112 | 71 | 12544 | 5041 | 7952 |
| 243 | 99 | 55 | 9801 | 3025 | 5445 |
| 244 | 99 | 66 | 9801 | 4356 | 6534 |
| 245 | 110 | 75 | 12100 | 5625 | 8250 |
| 274 | 133 | 85 | 17689 | 7225 | 11305 |
| 248 | 134 | 88 | 17956 | 7744 | 11792 |
| 253 | 129 | 83 | 16641 | 6889 | 10707 |
| 255 | 140 | 80 | 19600 | 6400 | 11200 |
| 256 | 117 | 72 | 13689 | 5184 | 8424 |
| 259 | 137 | 86 | 18769 | 7396 | 11782 |
| 231 | 164 | 95 | 26896 | 9025 | 15580 |
| 232 | 164 | 94 | 26896 | 8836 | 15416 |
| 233 | 164 | 89 | 26896 | 7921 | 14596 |
| 236 | 156 | 87 | 24336 | 7569 | 13572 |
| 237 | 147 | 103 | 21609 | 10609 | 15141 |
| 239 | 186 | 108 | 34596 | 11664 | 20088 |
| 241 | 170 | 102 | 28900 | 10404 | 17340 |
| 242 | 170 | 99 | 28900 | 9801 | 16830 |
| 246 | 176 | 121 | 30976 | 14641 | 21296 |
| 247 | 186 | 116 | 34596 | 13456 | 21576 |
| 249 | 157 | 107 | 24649 | 11449 | 16799 |
| 250 | 142 | 91 | 20164 | 8281 | 12922 |
| 251 | 159 | 85 | 25281 | 7225 | 13515 |
| 252 | 144 | 97 | 20736 | 9409 | 13968 |
| 254 | 155 | 113 | 24025 | 12769 | 17515 |
| 257 | 162 | 72 | 26244 | 5184 | 11664 |
| 258 | 151 | 98 | 22801 | 9604 | 14798 |
| 260 | 164 | 109 | 26896 | 11881 | 17876 |
| 261 | 155 | 105 | 24025 | 11025 | 16275 |
|  | 4631 | 2863 | 688837 | 264527 | 424780 |

Table A3 Percentage points of the $t$ distribution.

Adapted from Table 7 of White et al. (1979) with permission of authors and publishers.

| | | | | One-sided $P$ value | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 0.25 | 0.1 | 0.05 | 0.025 | 0.01 | 0.005 | 0.0025 | 0.001 | 0.0005 |
| | | | | Two-sided $P$ value | | | | | |
| d.f. | 0.5 | 0.2 | 0.1 | 0.05 | 0.02 | 0.01 | 0.005 | 0.002 | 0.001 |
| 1 | 1.00 | 3.08 | 6.31 | 12.71 | 31.82 | 63.66 | 127.32 | 318.31 | 636.62 |
| 2 | 0.82 | 1.89 | 2.92 | 4.30 | 6.96 | 9.92 | 14.09 | 22.33 | 31.60 |
| 3 | 0.76 | 1.64 | 2.35 | 3.18 | 4.54 | 5.84 | 7.45 | 10.21 | 12.92 |
| 4 | 0.74 | 1.53 | 2.13 | 2.78 | 3.75 | 4.60 | 5.60 | 7.17 | 8.61 |
| 5 | 0.73 | 1.48 | 2.02 | 2.57 | 3.36 | 4.03 | 4.77 | 5.89 | 6.87 |
| 6 | 0.72 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 | 4.32 | 5.21 | 5.96 |
| 7 | 0.71 | 1.42 | 1.90 | 2.36 | 3.00 | 3.50 | 4.03 | 4.78 | 5.41 |
| 8 | 0.71 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 | 3.83 | 4.50 | 5.04 |
| 9 | 0.70 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 | 3.69 | 4.30 | 4.78 |
| 10 | 0.70 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 | 3.58 | 4.14 | 4.59 |
| 11 | 0.70 | 1.36 | 1.80 | 2.20 | 2.72 | 3.11 | 3.50 | 4.02 | 4.44 |
| 12 | 0.70 | 1.36 | 1.78 | 2.18 | 2.68 | 3.06 | 3.43 | 3.93 | 4.32 |
| 13 | 0.69 | 1.35 | 1.77 | 2.16 | 2.65 | 3.01 | 3.37 | 3.85 | 4.22 |
| 14 | 0.69 | 1.34 | 1.76 | 2.14 | 2.62 | 2.98 | 3.33 | 3.79 | 4.14 |
| 15 | 0.69 | 1.34 | 1.75 | 2.13 | 2.60 | 2.95 | 3.29 | 3.73 | 4.07 |
| 16 | 0.69 | 1.34 | 1.75 | 2.12 | 2.58 | 2.92 | 3.25 | 3.69 | 4.02 |
| 17 | 0.69 | 1.33 | 1.74 | 2.11 | 2.57 | 2.90 | 3.22 | 3.65 | 3.96 |
| 18 | 0.69 | 1.33 | 1.73 | 2.10 | 2.55 | 2.88 | 3.20 | 3.61 | 3.92 |
| 19 | 0.69 | 1.33 | 1.73 | 2.09 | 2.54 | 2.86 | 3.17 | 3.58 | 3.88 |
| 20 | 0.69 | 1.32 | 1.72 | 2.09 | 2.53 | 2.84 | 3.15 | 3.55 | 3.85 |
| 21 | 0.69 | 1.32 | 1.72 | 2.08 | 2.52 | 2.83 | 3.14 | 3.53 | 3.82 |
| 22 | 0.69 | 1.32 | 1.72 | 2.07 | 2.51 | 2.82 | 3.12 | 3.50 | 3.79 |
| 23 | 0.68 | 1.32 | 1.71 | 2.07 | 2.50 | 2.81 | 3.10 | 3.48 | 3.77 |
| 24 | 0.68 | 1.32 | 1.71 | 2.06 | 2.49 | 2.80 | 3.09 | 3.47 | 3.74 |
| 25 | 0.68 | 1.32 | 1.71 | 2.06 | 2.48 | 2.79 | 3.08 | 3.45 | 3.72 |
| 26 | 0.68 | 1.32 | 1.71 | 2.06 | 2.48 | 2.78 | 3.07 | 3.44 | 3.71 |
| 27 | 0.68 | 1.31 | 1.70 | 2.05 | 2.47 | 2.77 | 3.06 | 3.42 | 3.69 |
| 28 | 0.68 | 1.31 | 1.70 | 2.05 | 2.47 | 2.76 | 3.05 | 3.41 | 3.67 |
| 29 | 0.68 | 1.31 | 1.70 | 2.04 | 2.46 | 2.76 | 3.04 | 3.40 | 3.66 |
| 30 | 0.68 | 1.31 | 1.70 | 2.04 | 2.46 | 2.75 | 3.03 | 3.38 | 3.65 |
| 40 | 0.68 | 1.30 | 1.68 | 2.02 | 2.42 | 2.70 | 2.97 | 3.31 | 3.55 |
| 60 | 0.68 | 1.30 | 1.67 | 2.00 | 2.39 | 2.66 | 2.92 | 3.23 | 3.46 |
| 120 | 0.68 | 1.29 | 1.66 | 1.98 | 2.36 | 2.62 | 2.86 | 3.16 | 3.37 |
| ∞ | 0.67 | 1.28 | 1.65 | 1.96 | 2.33 | 2.58 | 2.81 | 3.09 | 3.29 |

**We refer to Table A3.
so we use df=30 .
t = 8.349436 > 3.65 (p=0.001)**

**Therefore if t=8.349436, p<0.001.**

## Two-sided $P$ value

| d.f. | 0.5 | 0.2 | 0.1 | 0.05 | 0.02 | 0.01 | 0.005 | 0.002 | 0.001 |
|---|---|---|---|---|---|---|---|---|---|
| 30 | 0.68 | 1.31 | 1.70 | 2.04 | 2.46 | 2.75 | 3.03 | 3.38 | 3.65 |
| 40 | 0.68 | 1.30 | 1.68 | 2.02 | 2.42 | 2.70 | 2.97 | 3.31 | 3.55 |
| 60 | 0.68 | 1.30 | 1.67 | 2.00 | 2.39 | 2.66 | 2.92 | 3.23 | 3.46 |
| 120 | 0.68 | 1.29 | 1.66 | 1.98 | 2.36 | 2.62 | 2.86 | 3.16 | 3.37 |
| ∞ | 0.67 | 1.28 | 1.65 | 1.96 | 2.33 | 2.58 | 2.81 | 3.09 | 3.29 |

# Correlation

Two pieces of information:

- The strength of the relationship
- The direction of the relationship

# Strength of relationship

- r lies between -1 and 1. Values near 0 means no (linear) correlation and values near ± 1 means very strong correlation.

|  |  |  |
|---|---|---|
| -1.0 | 0.0 | +1.0 |
| Strong Negative | No Rel. | Strong Positive |

# How to interpret the value of r?

## Table II. Strength of linear relationship.

| Correlation Coefficient value | Strength of linear relationship |
|---|---|
| At least 0.8 | Very strong |
| 0.6 up to 0.8 | Moderately strong |
| 0.3 to 0.5 | Fair |
| Less than 0.3 | Poor |

# Correlation ( + direction)

- Positive correlation: high values of one variable associated with high values of the other

- Example: Higher course entrance exam scores are associated with better course grades during the final exam.

Positive and Linear

# Correlation ( - direction)

- Negative correlation: The negative sign means that the two variables are inversely related, that is, as one variable increases the other variable decreases.

- Example:  Increase in body mass index is associated with reduced effort tolerance.



Negative and Linear

# Pearson's *r*

- A 0.9 is a very strong positive association (as one variable rises, so does the other)

- A -0.9 is a very strong negative association
  (as one variable rises, the other falls)

*r*=0.9 has nothing to do with 90%

*r=correlation coefficient*

# Coefficient of Determination Defined

- Pearson's $r$ can be squared , $r^2$, to derive a coefficient of determination.

- Coefficient of determination – the portion of variability in one of the variables that can be accounted for by variability in the second variable
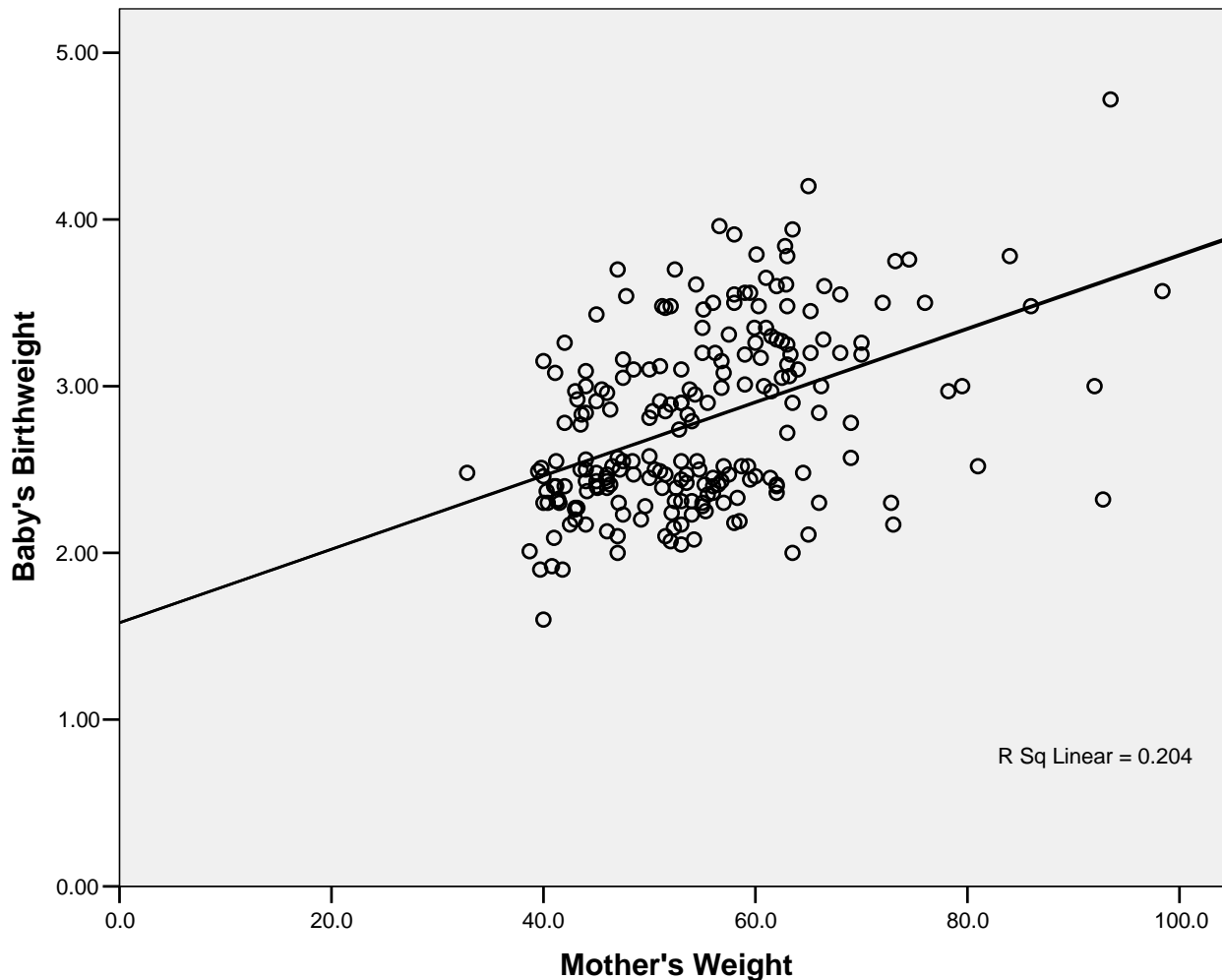
# Coefficient of Determination

- Pearson's $r$ can be squared , $r^2$, to derive a coefficient of determination.

- Example of depression and CGPA
  - Pearson's r shows negative correlation, $r$=-0.5
  - $r^2$=0.25

  - In this example we can say that 1/4 or 0.25 of the variability in CGPA scores can be accounted for by depression (remaining 75% of variability is other factors, habits, ability, motivation, courses studied, etc)

# Coefficient of Determination and Pearson's *r*

- Pearson's *r* can be squared , $r^2$

- If *r*=0.5, then $r^2$=0.25
- If *r*=0.7 then $r^2$=0.49

- Thus while *r*=0.5 versus 0.7 might not look so different in terms of strength, $r^2$ tells us that *r*=0.7 accounts for about twice the variability relative to *r*=0.5

# A study was done to find the association between the mothers' weight and their babies' birth weight. The following is the scatter diagram showing the relationship between the two variables.



The coefficient of correlation (r) is 0.452

The coefficient of determination ($r^2$) is 0.204

Twenty percent of the variability of the babies' birth weight is determined by the variability of the mothers' weight.

# Causal Silence:
# Correlation Does Not Imply Causality

**Causality – must demonstrate that variance in one variable can only be due to influence of the other variable**

- **Directionality of Effect Problem**
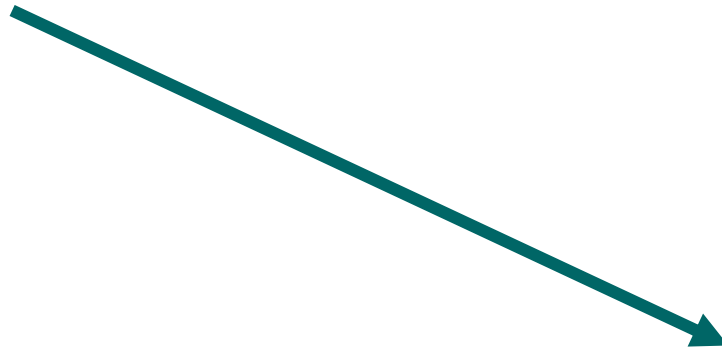
- **Third Variable Problem**

# CORRELATION DOES NOT MEAN CAUSATION

- A high correlation **does not** give us the evidence to make a cause-and-effect statement.
- A common example given is the high correlation between the cost of damage in a fire and the number of firemen helping to put out the fire.
- Does it mean that to cut down the cost of damage, the fire department should dispatch less firemen for a fire rescue!
- The intensity of the fire that is highly correlated with the cost of damage and the number of firemen dispatched.
- The high correlation between smoking and lung cancer. However, one may argue that both could be caused by stress; and smoking does not cause lung cancer.
- In this case, a correlation between lung cancer and smoking may be a result of a cause-and-effect relationship (by clinical experience + common sense?). To establish this cause-and-effect relationship, controlled experiments should be performed.

Big Fire

More
Firemen
Sent

More
Damage

# Directionality of Effect Problem

X ⟶ Y

X ⟵ Y

X ⟷ Y

# Directionality of Effect Problem

**X** → **Y**

**Class Attendance**          **Higher Grades**

**X** ← **Y**

**Class Attendance**          **Higher Grades**

# Directionality of Effect Problem

**X** → **Y**

**Aggressive Behavior**  **Viewing Violent TV**

**X** ← **Y**

**Aggressive Behavior**  **Viewing Violent TV**

Aggressive children may prefer violent programs or
Violent programs may promote aggressive behavior

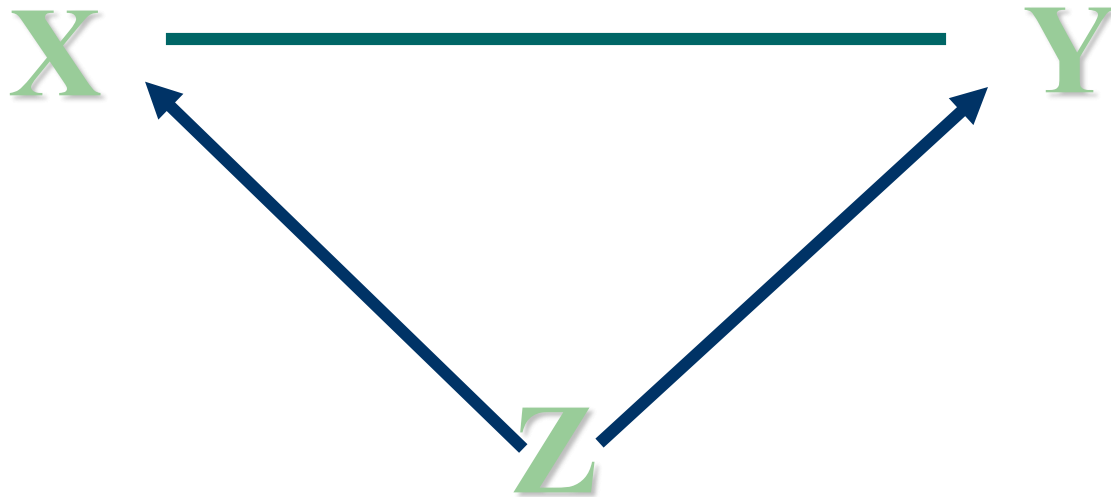# Methods for Dealing with Directionality

- Cross-Lagged Panel design
    - A type of longitudinal design
    - Investigate correlations at several points in time
    - STILL NOT CAUSAL
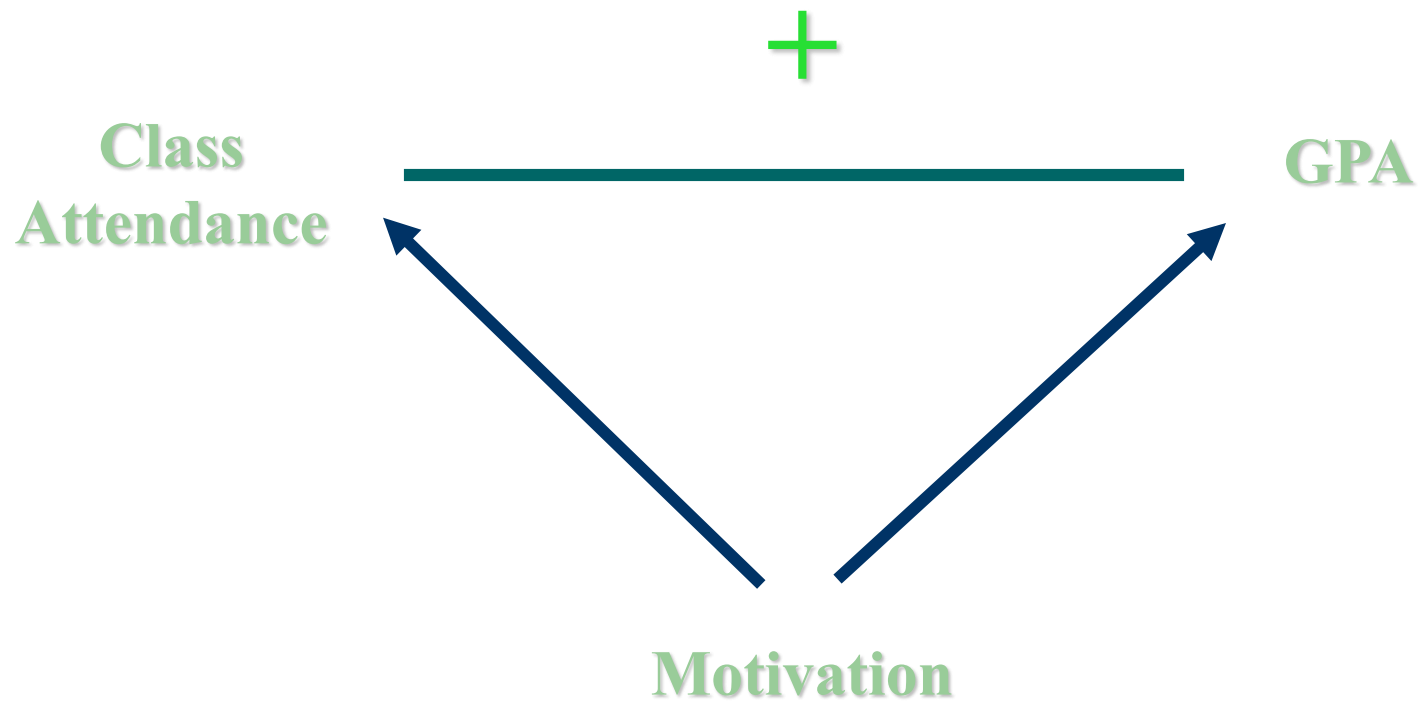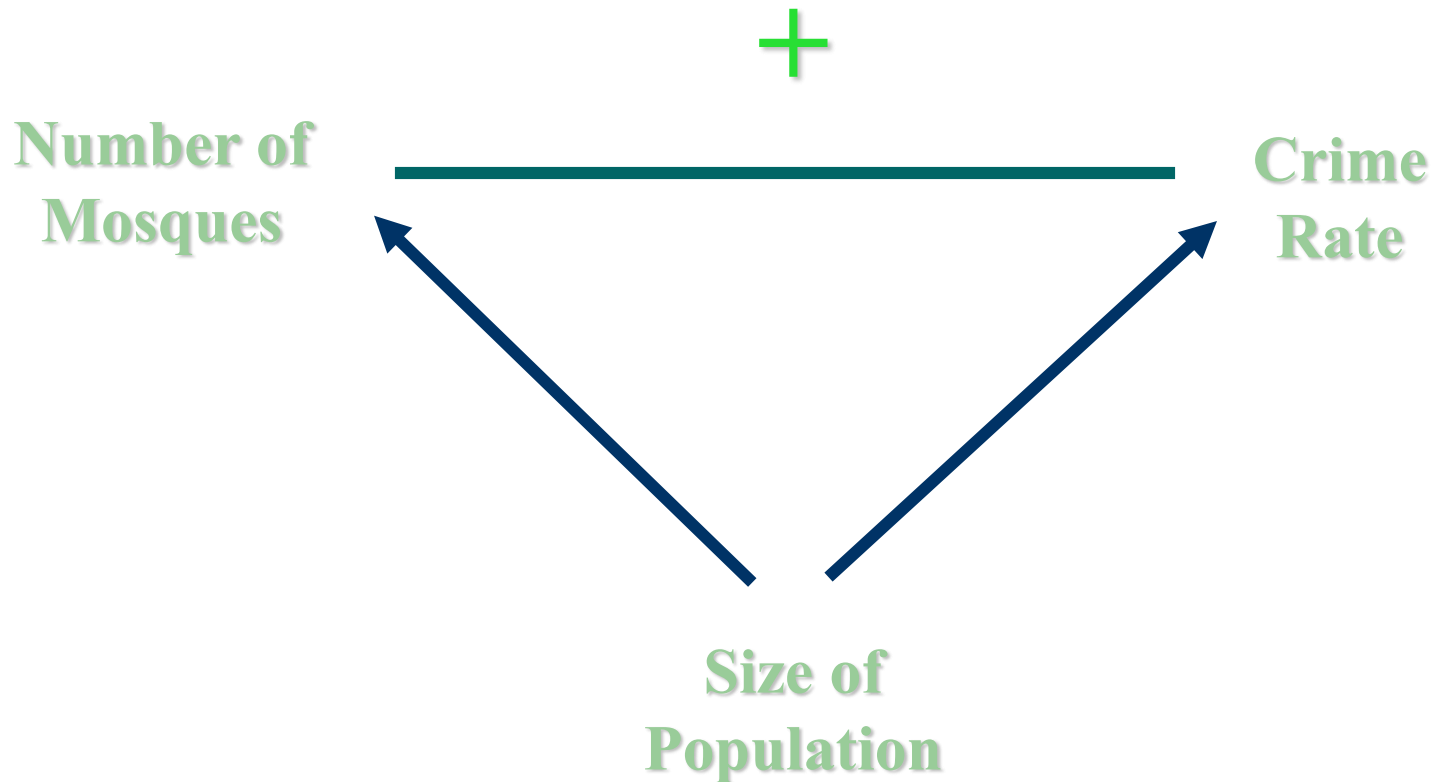
    Example next page

# Cross-Lagged Panel

Pref for violent TV       .05       Pref for violent TV

3rd grade                      13th grade

.21         .31         .01         -.05

TV to later
aggression      Aggression to later TV

Aggression               Aggression

3rd grade       .38       13th grade

# Third Variable Problem

X ——————— Y

Z

# Class Exercise

Identify the

third variable

that influences both X and Y

# Third Variable Problem

# Third Variable Problem

+

Number of
Mosques

Crime
Rate

Size of
Population

# Third Variable Problem

+

**Ice Cream Consumed** ———————————— **Number of Drownings**

**Temperature**

# Third Variable Problem

+

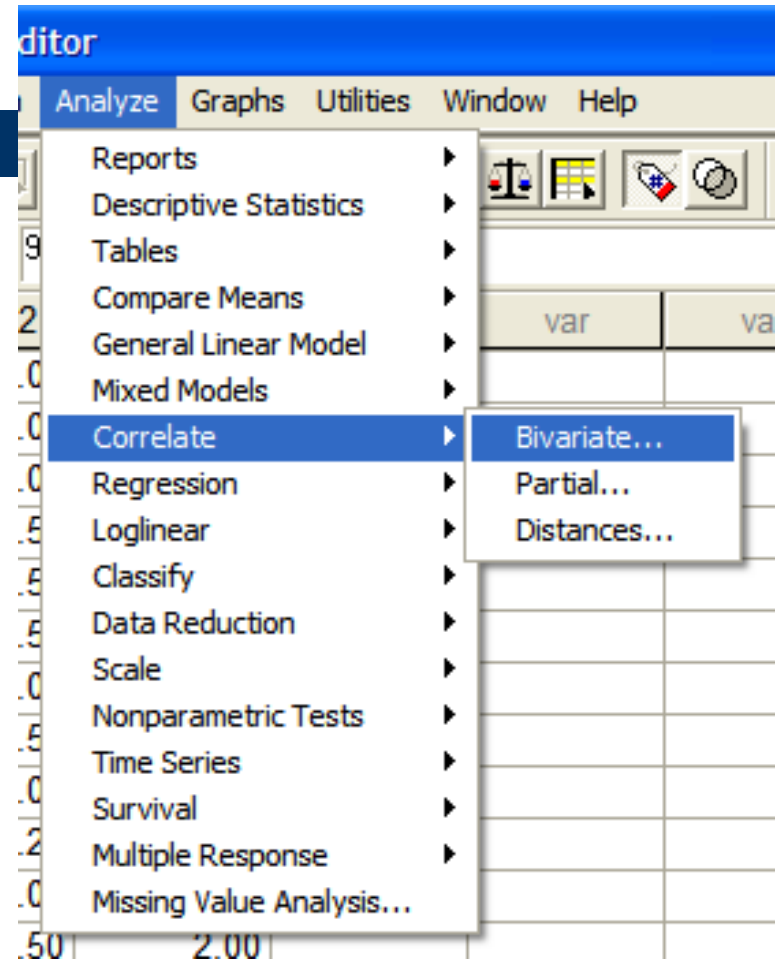**Reading Score** ———————— **Reading Comprehension**

**IQ**

# Data Preparation - Correlation

- Screen data for outliers and ensure that there is evidence of linear relationship, since correlation is a measure of linear relationship.

- Assumption is that each pair is bivariate normal.

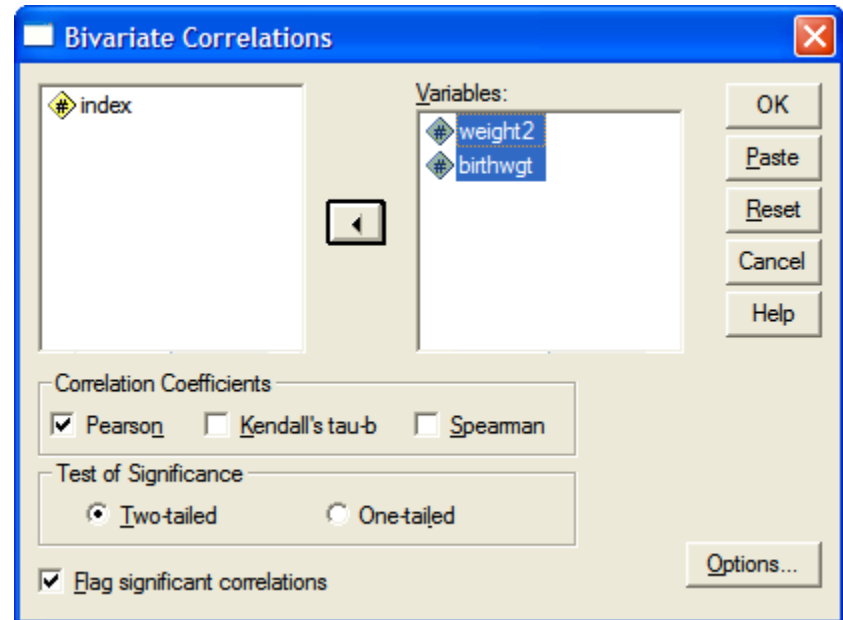- If not normal, then use Spearman.

# Correlation In SPSS

- For this exercise, we will be using the data from the CD, under Chapter 8, korelasi.sav

- This data is a subset of a case-control study on factors affecting SGA in Kelantan.

- Open the data & select -
>Analyse
    >Correlate
        >Bivariate…

# Correlation in SPSS

- We want to see whether there is any association between the mothers' weight and the babies'weight. So select the variables (weight2 & birthwgt) into 'Variables'.

- Select 'Pearson' Correlation Coefficients.
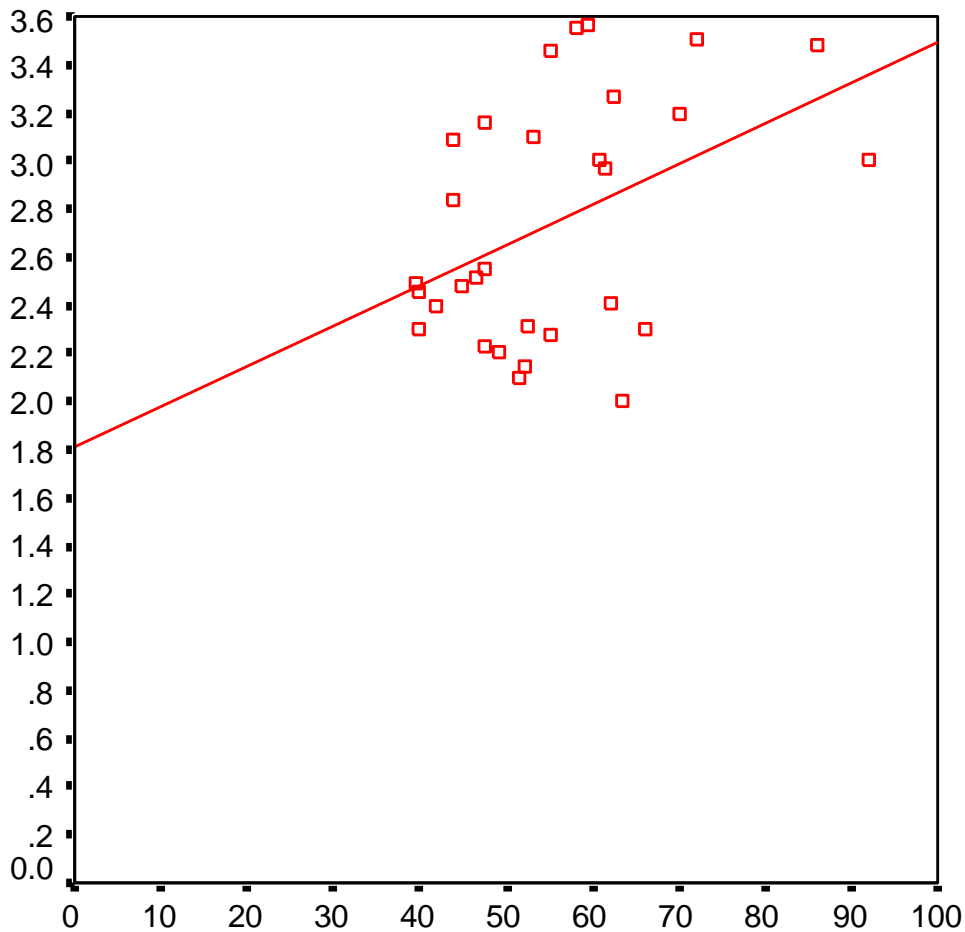
- Click the 'OK' button.

# Correlation Results

**Correlations**

|  |  | WEIGHT2 | BIRTHWGT |
|---|---|---|---|
| WEIGHT2 | Pearson Correlation | 1 | .431* |
|  | Sig. (2-tailed) | . | .017 |
|  | N | 30 | 30 |
| BIRTHWGT | Pearson Correlation | .431* | 1 |
|  | Sig. (2-tailed) | .017 | . |
|  | N | 30 | 30 |

*. Correlation is significant at the 0.05 level (2-tailed).

- The r = 0.431 and the p value is significant at 0.017.

- The r value indicates a fair and positive linear relationship.

# Scatter Diagram



Rsq = 0.1861

MOTHERS' WEIGHT

- If the correlation is significant, it is best to include the scatter diagram.

- The r square indicated mothers' weight contribute 19% of the variability of the babies' weight.
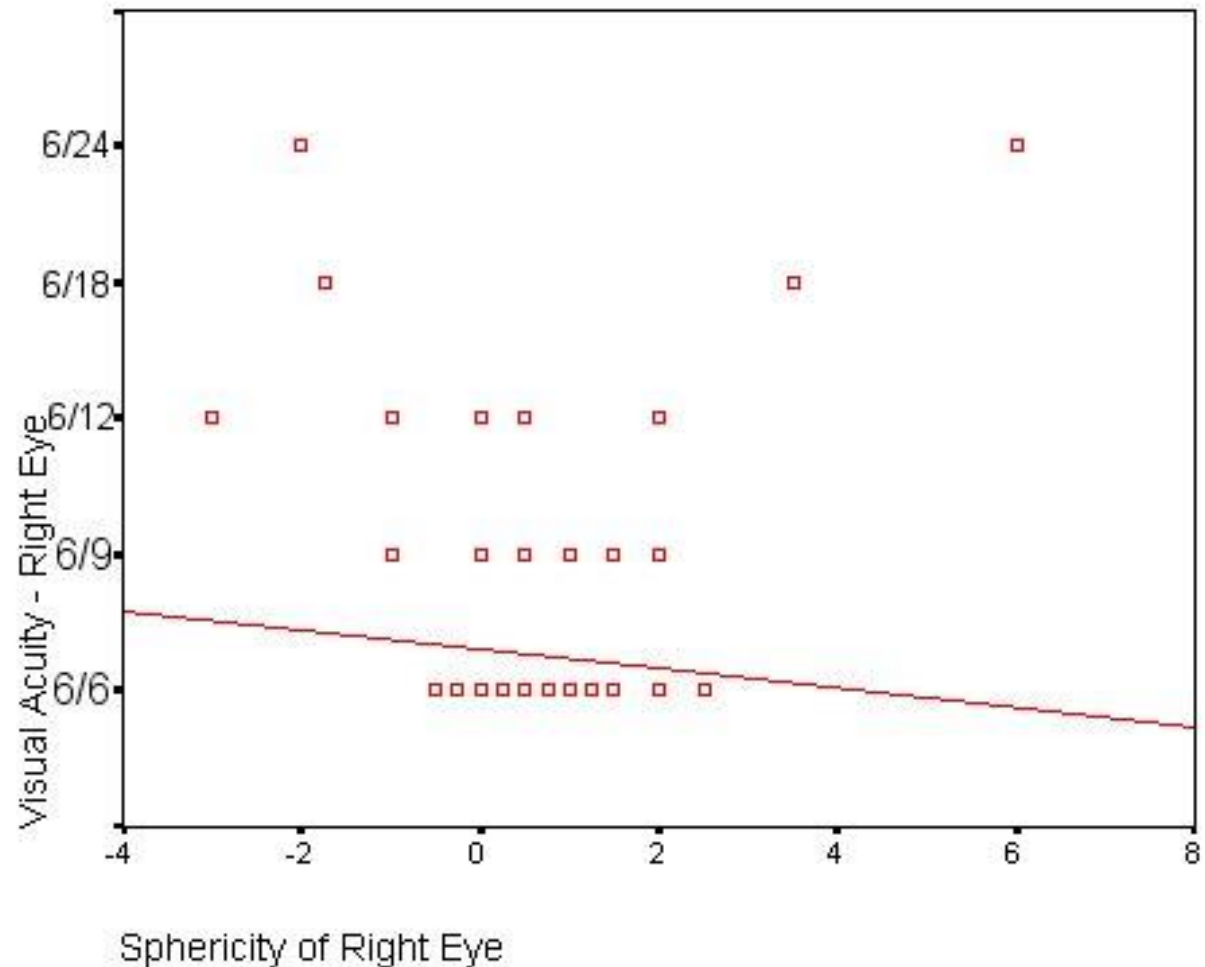
# Spearman/Kendall Correlation

- To find correlation between a related pair of continuous data (not normally distributed); or
- **Between 1 Continuous, 1 Categorical Variable (Ordinal)**
  - **e.g., association between Likert Scale on work satisfaction and work output.**

# Spearman's rank correlation coefficient

- In statistics, **Spearman's rank correlation coefficient**, named for Charles Spearman and often denoted by the Greek letter ρ (rho), is a non-parametric measure of correlation – that is, it assesses how well an arbitrary monotonic function could describe the relationship between two variables, without making any assumptions about the frequency distribution of the variables. Unlike the Pearson product-moment correlation coefficient, it does not require the assumption that the relationship between the variables is linear, nor does it require the variables to be measured on interval scales; it can be used for variables measured at the ordinal level.

# Example

• Correlation between sphericity and visual acuity.

• Sphericity of the eyeball is continuous data while visual acuity is ordinal data (6/6, 6/9, 6/12, 6/18, 6/24), therefore Spearman correlation is the most suitable.

• The Spearman rho correlation coefficient is -0.108 and p is 0.117. P is larger than 0.05, therefore there is no significant association between sphericity and visual acuity.



Sphericity of Right Eye

**Correlations**

|  |  |  | Visual Acuity - Right Eye | Sphericity of Right Eye |
|---|---|---|---|---|
| Spearman's rho | Visual Acuity - Right Eye | Correlation Coefficient | 1.000 | -.108 |
|  |  | Sig. (2-tailed) | . | .117 |
|  |  | N | 215 | 211 |
|  | Sphericity of Right Eye | Correlation Coefficient | -.108 | 1.000 |
|  |  | Sig. (2-tailed) | .117 | . |
|  |  | N | 211 | 211 |

# Example 2

•- Correlation between glucose level and systolic blood pressure.

•Based on the data given, prepare the following table;

•For every variable, sort the data by rank. For ties, take the average.

•Calculate the difference of rank, d for every pair and square it. Take the total.

•Include the value into the following formula;

$$r_s = 1 - \frac{6 \Sigma d^2}{n(n^2 - 1)}$$

•$\Sigma$ d$^2$ = 4921.5     n = 32

•Therefore r$_s$ = 1-((6*4921.5)/(32*(32$^2$-1)))
                = 0.097966.

This is the value of Spearman correlation coefficient (or $\Upsilon$).

•Compare the value against the Spearman table;

•p is larger than 0.05.

•Therefore there is no association between systolic BP and blood glucose level.

| nores | glu | rank x | bps1 | rank y | d | d2 |
|---|---|---|---|---|---|---|
| 231 | 123 | 23 | 164 | 25.5 | -2.5 | 6.25 |
| 232 | 97 | 9 | 164 | 25.5 | -16.5 | 272.25 |
| 233 | 325 | 32 | 164 | 25.5 | 6.5 | 42.25 |
| 234 | 124 | 24 | 118 | 7 | 17 | 289 |
| 235 | 107 | 12.5 | 126 | 8 | 4.5 | 20.25 |
| 236 | 95.7 | 8 | 156 | 20 | -12 | 144 |
| 237 | 122 | 22 | 147 | 16 | 6 | 36 |
| 238 | 112 | 17 | 105 | 3 | 14 | 196 |
| 239 | 119 | 20 | 186 | 31.5 | -11.5 | 132.25 |
| 240 | 132 | 25 | 112 | 5 | 20 | 400 |
| 241 | 105 | 11 | 170 | 28.5 | -17.5 | 306.25 |
| 242 | 219 | 30 | 170 | 28.5 | 1.5 | 2.25 |
| 243 | 141 | 26 | 99 | 1.5 | 24.5 | 600.25 |
| 244 | 93.6 | 4 | 99 | 1.5 | 2.5 | 6.25 |
| 245 | 206 | 29 | 110 | 4 | 25 | 625 |
| 246 | 113 | 18.5 | 176 | 30 | -11.5 | 132.25 |
| 247 | 167 | 28 | 186 | 31.5 | -3.5 | 12.25 |
| 248 | 95.6 | 7 | 134 | 11 | -4 | 16 |
| 249 | 108 | 14.5 | 157 | 21 | -6.5 | 42.25 |
| 250 | 297 | 31 | 142 | 14 | 17 | 289 |
| 251 | 109 | 16 | 159 | 22 | -6 | 36 |
| 252 | 100 | 10 | 144 | 15 | -5 | 25 |
| 253 | 83.3 | 2 | 129 | 9 | -7 | 49 |
| 254 | 145 | 27 | 155 | 18.5 | 8.5 | 72.25 |
| 255 | 90.2 | 3 | 140 | 13 | -10 | 100 |
| 256 | 113 | 18.5 | 117 | 6 | 12.5 | 156.25 |
| 257 | 108 | 14.5 | 162 | 23 | -8.5 | 72.25 |
| 258 | 121 | 21 | 151 | 17 | 4 | 16 |
| 259 | 94.5 | 6 | 137 | 12 | -6 | 36 |
| 260 | 69.4 | 1 | 164 | 25.5 | -24.5 | 600.25 |
| 261 | 94.2 | 5 | 155 | 18.5 | -13.5 | 182.25 |
| 274 | 107 | 12.5 | 133 | 10 | 2.5 | 6.25 |
|  |  |  |  |  |  | 4921.5 |

# Spearman's table

- 0.097966 is the value of Spearman correlation coefficient (or $\rho$).
- Compare the value against the Spearman table;
- $0.098 < 0.364$ (p=0.05)
- p is larger than 0.05.
- Therefore there is no association between systolic BP and blood glucose level.

| N (the number of pairs of scores): | 0.05 | 0.02 | 0.01 |
|---|---|---|---|
| 5 | 1 | 1 | |
| 6 | 0.886 | 0.943 | 1 |
| 7 | 0.786 | 0.893 | 0.929 |
| 8 | 0.738 | 0.833 | 0.881 |
| 9 | 0.683 | 0.783 | 0.833 |
| 10 | 0.648 | 0.746 | 0.794 |
| 12 | 0.591 | 0.712 | 0.777 |
| 14 | 0.544 | 0.645 | 0.715 |
| 16 | 0.506 | 0.601 | 0.665 |
| 18 | 0.475 | 0.564 | 0.625 |
| 20 | 0.45 | 0.534 | 0.591 |
| 22 | 0.428 | 0.508 | 0.562 |
| 24 | 0.409 | 0.485 | 0.537 |
| 26 | 0.392 | 0.465 | 0.515 |
| 28 | 0.377 | 0.448 | 0.496 |
| 30 | 0.364 | 0.432 | 0.478 |

# SPSS Output

**Correlations**

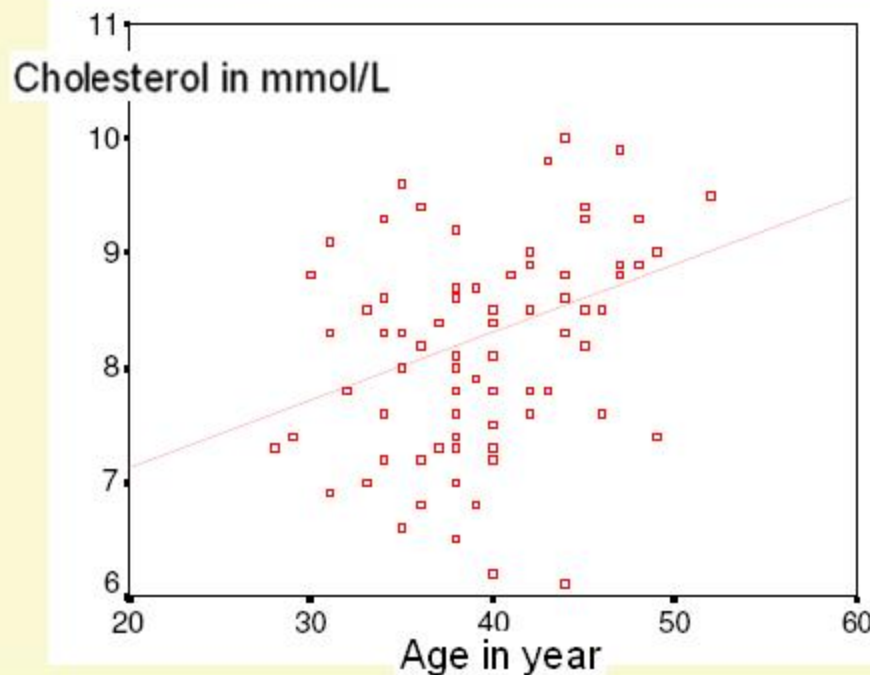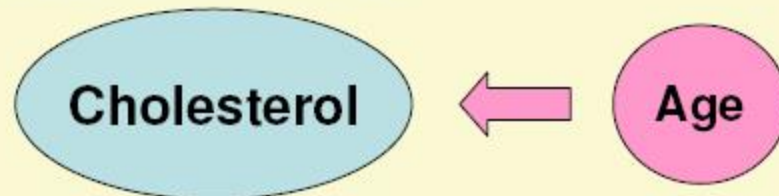| | | | GLU | BPS1 |
|---|---|---|---|---|
| Spearman's rho | GLU | Correlation Coefficient | 1.000 | .097 |
| | | Sig. (2-tailed) | . | .599 |
| | | N | 32 | 32 |
| | BPS1 | Correlation Coefficient | .097 | 1.000 |
| | | Sig. (2-tailed) | .599 | . |
| | | N | 32 | 32 |

# Linear Regression

# Simple Linear Regression

- To determine the relationship between age and blood cholesterol level



▶ Here, we may use either 'correlation analysis' or 'regression analysis', as both cholesterol and age are numerical variables.

▶ *Correlation* can give the strength of relationship, but *regression* can describe the relationship in more detail.

▶ In above example, if we decide to do *regression*, cholesterol will be our outcome (dependent) variable, because age may determine cholesterol but cholesterol cannot determine age.

# Simple Linear Regression

- **To determine the relationship between age and blood cholesterol level**

# Simple Linear Regression

# Simple Linear Regression

**Data Editor**

| rm | Analyze | Graphs | Utilities | Window | Help |
|----|---------|--------|-----------|--------|------|

Reports ▶
Descriptive Statistics ▶
Custom Tables ▶
Compare Means ▶
General Linear Model ▶
Mixed Models ▶
Correlate ▶
Regression ▶    Linear… ①
Loglinear ▶    Curve Estimat
Classify ▶

se_stat

**Linear Regression**

- # age
- # diet
- # exercise
- # se_stat

② → Dependent: # chol

Previous    Block 1 of 1

③ → Independent(s): # age

Method: Enter

Selection Variabl

Case Labels:

WLS >>      ④ Statistics…   Plots…

**Linear Regression: Statistics**

Regression Coefficients
- ✔ Estimates
- ☑ Confidence intervals ⑤
- ☐ variance matrix

$$Y = a + bX$$

$$Chol = 5.9 + (0.058*age)$$

**Coefficients**[a]

| Model | | Unstandardized Coefficients | | Standardized Coefficients | | *P* value | 95% Confidence Interval for B | |
|-------|--|------------------------------|--|----------------------------|--|-----------|-------------------------------|--|
| | | B | Std. Error | Beta | t | Sig. | Lower Bound | Upper Bound |
| 1 | (Constant) | 5.895 | .735 | | 8.026 | .000 | 4.434 | 7.357 |
| | AGE  age in year | 5.776E-02 | .018 | .331 | 3.134 | .002 | .021 | .094 |

$H_o: \beta = 0$

a. Dependent Variable: CHOL  cholesterol in mmol/L
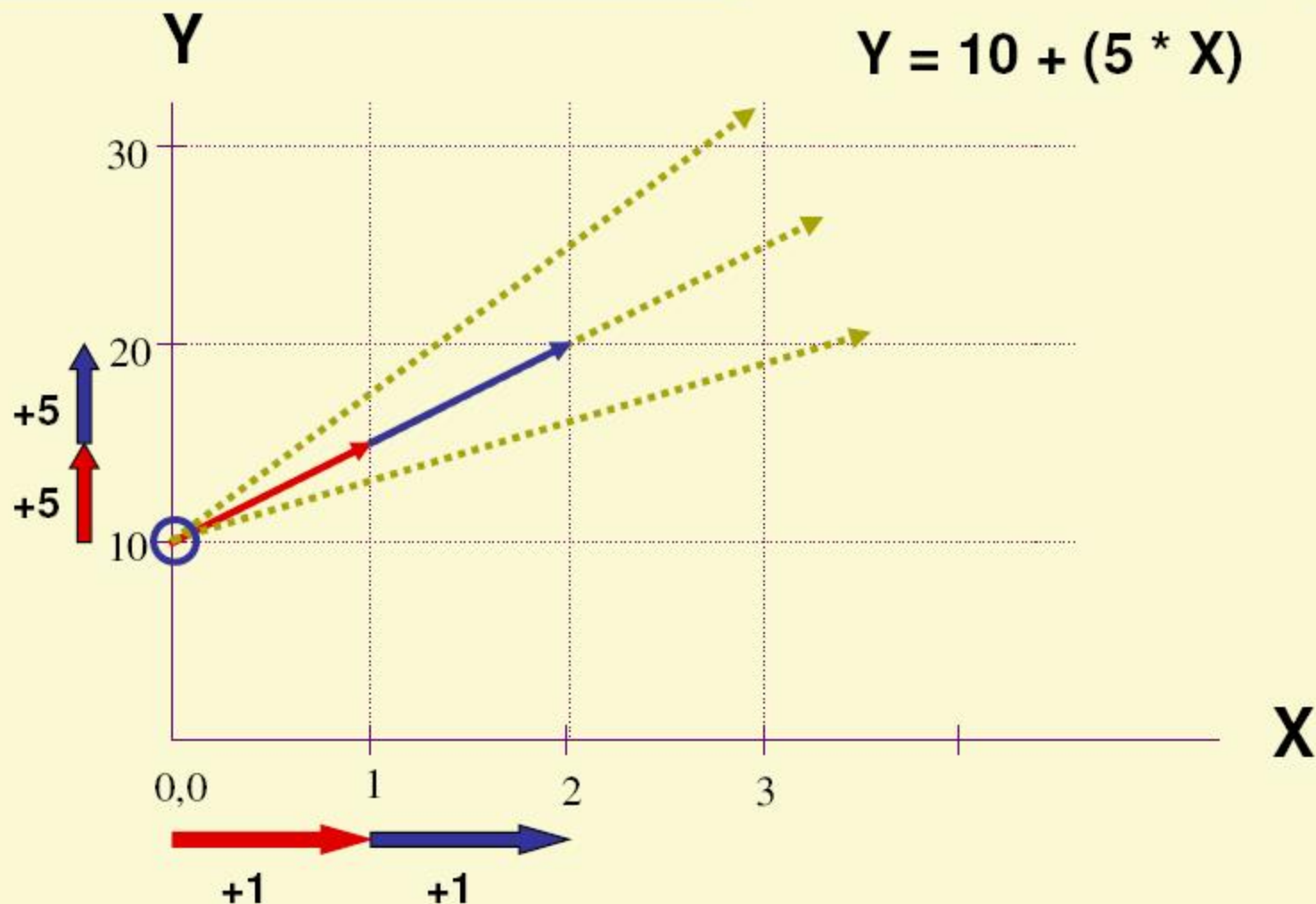
Slope (b) = 0.058 (95% CI: 0.021, 0.094)

**The Linear line is described by the "Linear Equation".**

$Y = a + (b * X)$

$Y = \text{Constant} + (\text{slope} * X)$

$Y = 10 + (5 * X)$

# The Least Squares (Regression) Line

A good line is one that minimizes
the sum of squared differences between the
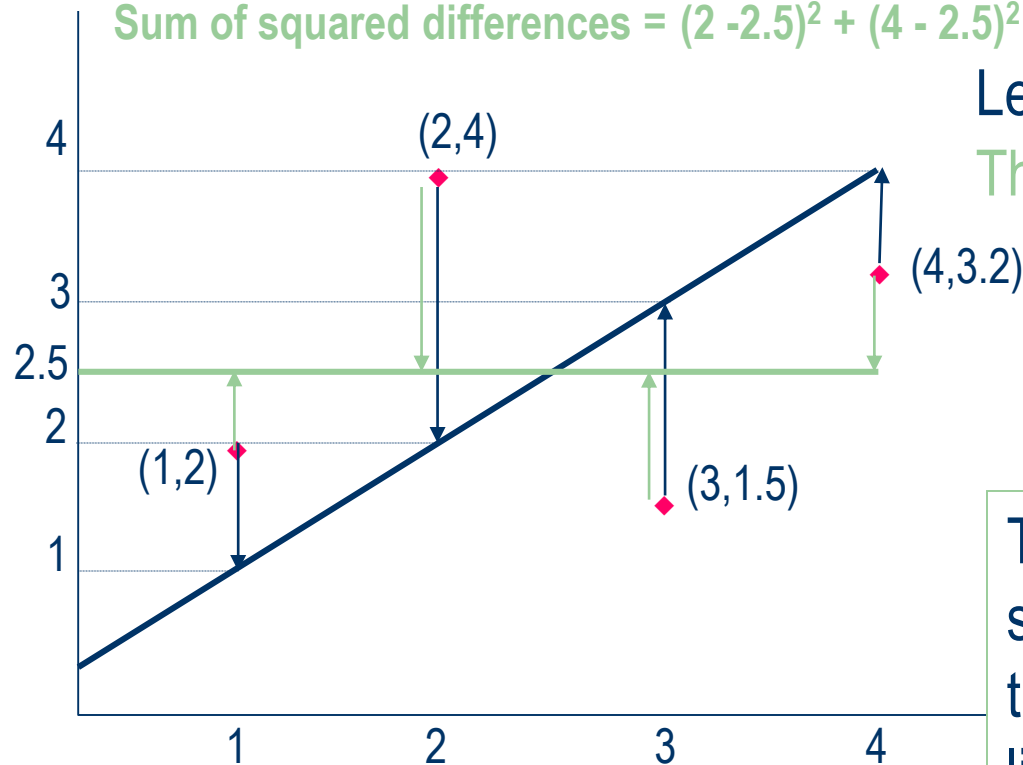points and the line.

# The Least Squares (Regression) Line

**Sum of squared differences = $(2 - 1)^2 + (4 - 2)^2 + (1.5 - 3)^2 + (3.2 - 4)^2 = 6.89$**

**Sum of squared differences = $(2 - 2.5)^2 + (4 - 2.5)^2 + (1.5 - 2.5)^2 + (3.2 - 2.5)^2 = 3.99$**
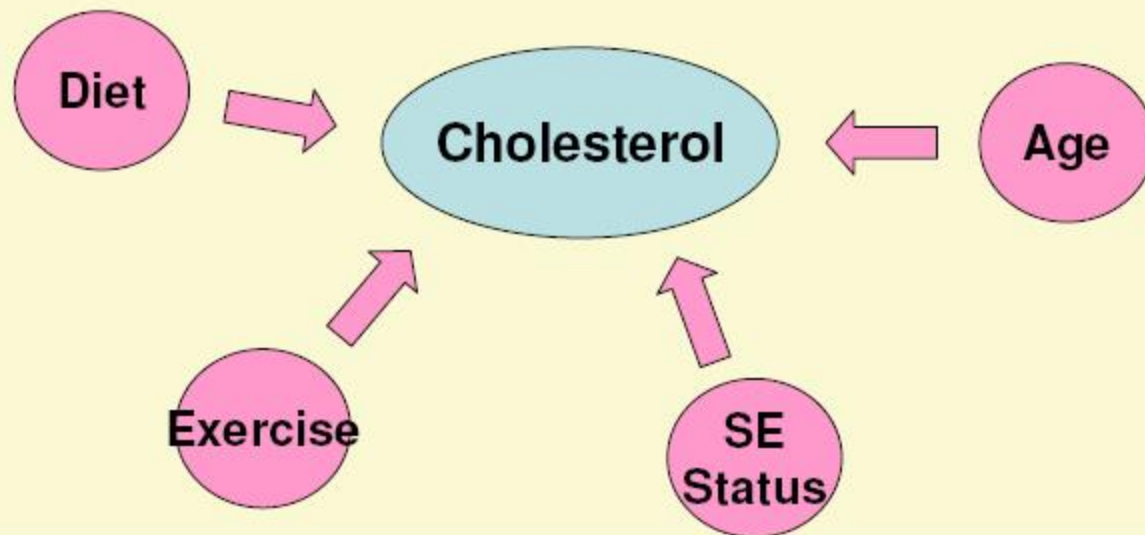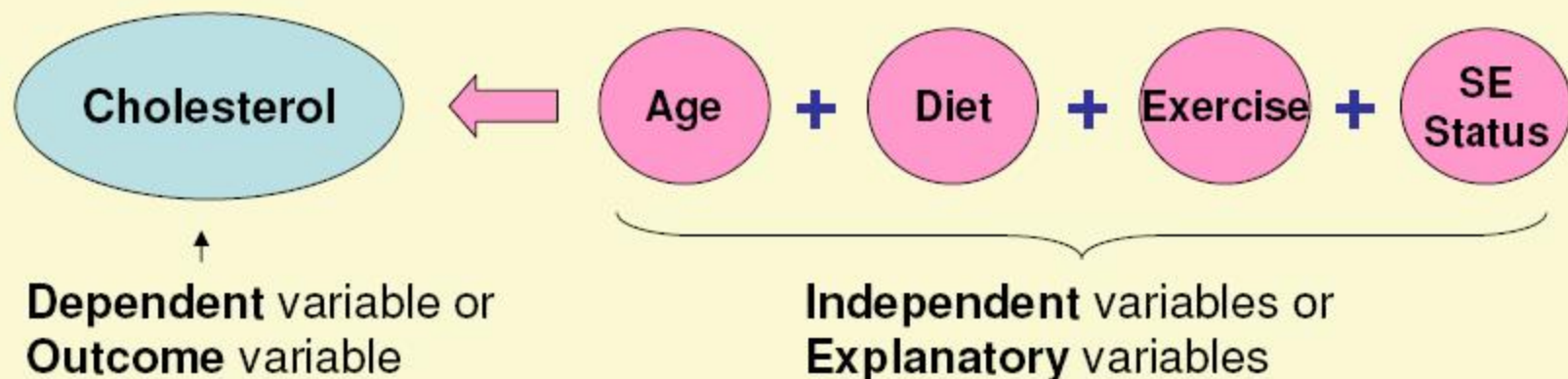
Let us compare two lines

The second line is horizontal



The smaller the sum of squared differences the better the fit of the line to the data.

# Basic Theory of MLR

- **Most of the outcomes (events) are determined (influenced) by more than one factors (e.g. blood pressure, cholesterol level, etc.)**
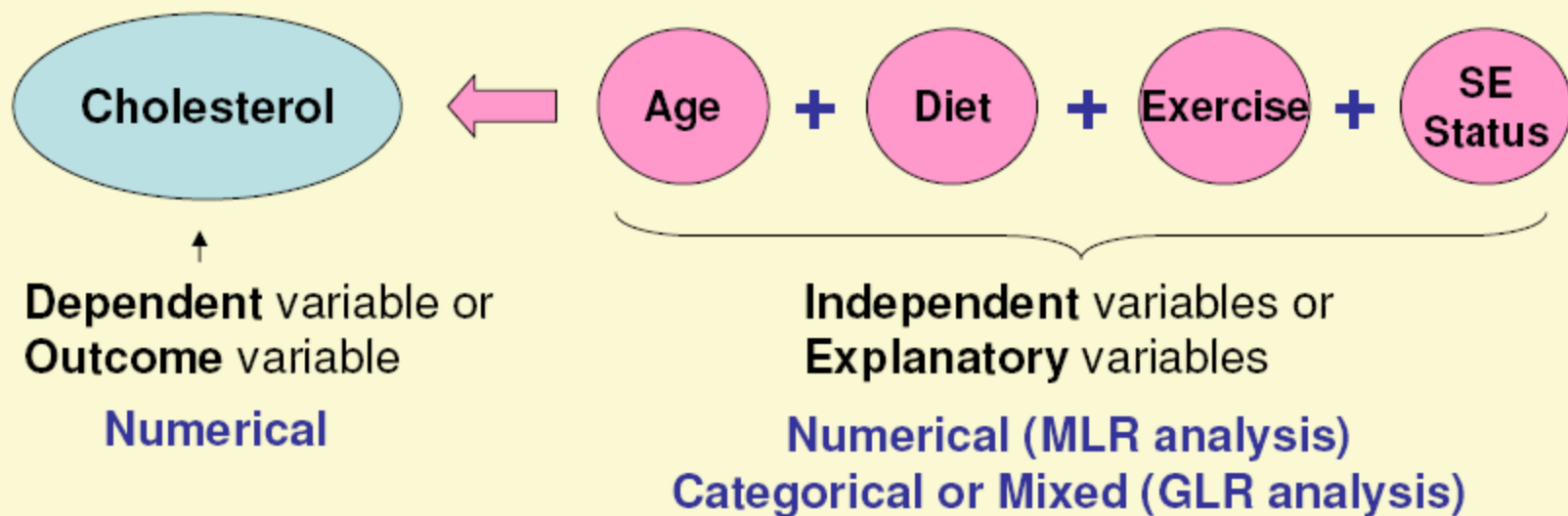
# Basic Theory

Cholesterol ⟸ Age + Diet + Exercise + SE Status

**Dependent** variable or
**Outcome** variable

**Independent** variables or
**Explanatory** variables

- **This analysis is used for ….**
  - **Exploring associated / influencing / risk factors to outcome (exploratory study)**
  - **Developing prediction model (exploratory study)**
  - **Confirming a specific relationship (confirmatory study)**

# Basic Theory

**Cholesterol** ⟸ **Age** **+** **Diet** **+** **Exercise** **+** **SE Status**

**Dependent** variable or
**Outcome** variable

**Numerical**

**Independent** variables or
**Explanatory** variables

**Numerical (MLR analysis)**
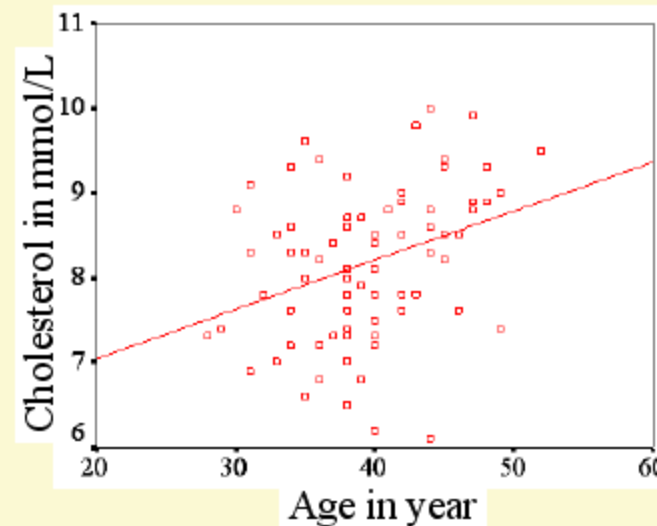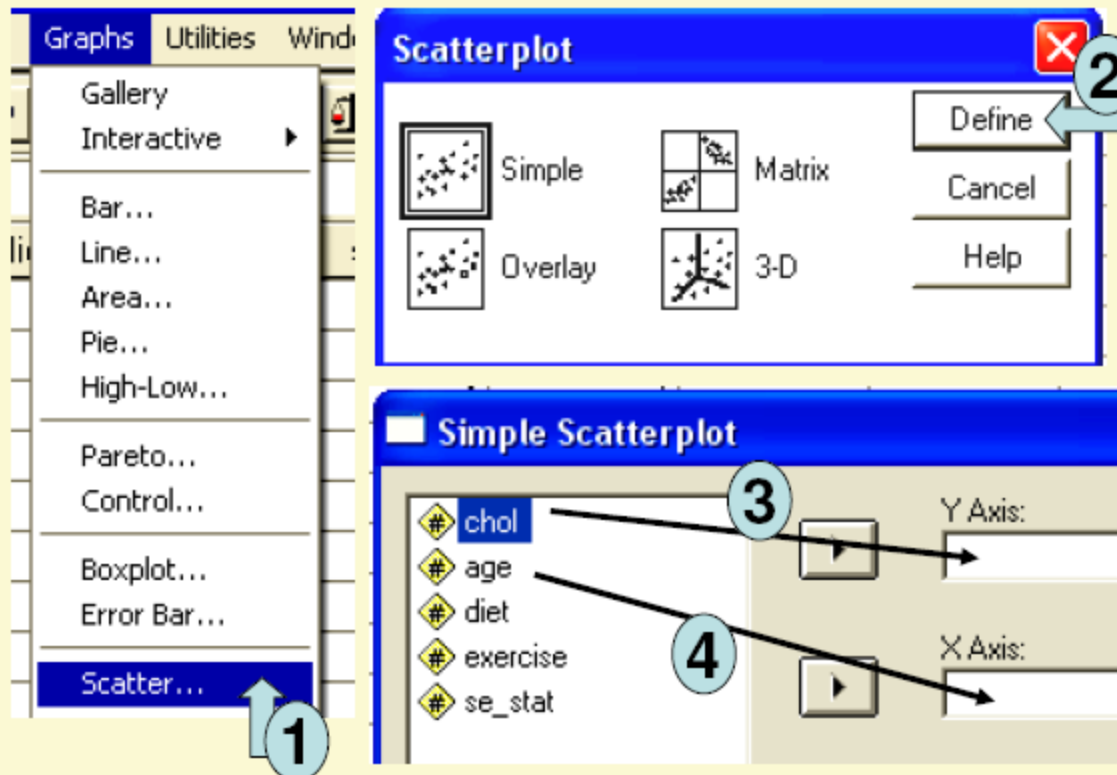**Categorical or Mixed (GLR analysis)**

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \ldots\ldots + \beta_n X_n$$

- If the dependent variable is numerical and independent variables are numerical, it will be called <u>Multiple Linear Regression</u> (MLR) analysis.
- MLR can be with categorical independent variables, but special name is given as <u>General Linear Regression</u> analysis.

# Step 2: Simple Linear Regression

**Two main reasons:**

1) To check the 'gross' relationship between dependent and each independent variable
2) Later this result will be compared with multiple linear regression result. This comparison indicates the confounding effects if it is present.

# Step 2: Simple Linear Regression



**Coefficients**[a]

| Model | | Unstandardized Coefficients | | Standardized Coefficients | | P value | 95% Confidence Interval for B | |
|---|---|---|---|---|---|---|---|---|
| | | B | Std. Error | Beta | t | Sig. | Lower Bound | Upper Bound |
| 1 | (Constant) | 5.895 | .735 | | 8.026 | .000 | 4.434 | 7.357 |
| | AGE  age in year | 5.776E-02 | .018 | .331 | 3.134 | .002 | .021 | .094 |

a. Dependent Variable: CHOL  cholesterol in mmo **Slope (b) = 0.058 (95% CI: .021 .094)**

Table 3: Factors associated with blood cholesterol level (mmol/L) among the study population ($n$=82) using simple linear regression

| Independent Variable | SLR[a] | | |
|---|---|---|---|
| | $b$ ( 95%CI ) | | $P$ value |
| Age (year) | 0.06 ( 0.02, 0.09) | | 0.002 |
| Duration of exercise (hrs/wk) | - 0.62 (- 0.79,- 0.46) | | <0.001 |
| Diet inventory score | 0.45 ( 0.30, 0.61) | | <0.001 |
| Socio-economic index | 0.21 ( 0.17, 0.25) | | <0.001 |

[a] Simple linear regression (Outcome as Cholesterol mmol/L)
$b$ = crude regression coefficient

# Regression Line

- In a scatterplot showing the association between 2 variables, the regression line is the "best-fit" line and has the formula

y=a + bx

a=place where line crosses Y axis

b=slope of line (rise/run)

Thus, given a value of X, we can predict a value of Y

# Linear Regression

- Come up with a **Linear Regression Model** to predict a continous outcome with a continuous risk factor, i.e. predict BP with age. Usually LR is the next step after correlation is found to be strongly significant.

- y = a + bx; a = y - bx

  – e.g. BP = **constant (a) + regression coefficient (b) * age**

- b=

$$b = \frac{\sum xy - \frac{(\sum x)(\sum y)}{n}}{\sum x^2 - \frac{(\sum x)^2}{n}}$$

# Example

$$b = \cfrac{\sum xy - \cfrac{(\sum x)(\sum y)}{n}}{\sum x^2 - \cfrac{(\sum x)^2}{n}}$$

$\sum x = 6426$        $\sum x^2 = 1338088$

$\sum y = 4631$        $\sum xy = 929701$

n = 32

b = (929701-(6426*4631/32))/
(1338088-(6426²/32)) = -0.00549

Mean x =   6426/32=200.8125
mean y =   4631/32=144.71875

y = a + bx

a = y – bx  (replace the x, y & b value)

a = 144.71875+(0.00549*200.8125)
= 145.8212106

Systolic BP = 145.82121 - 0.00549.chol

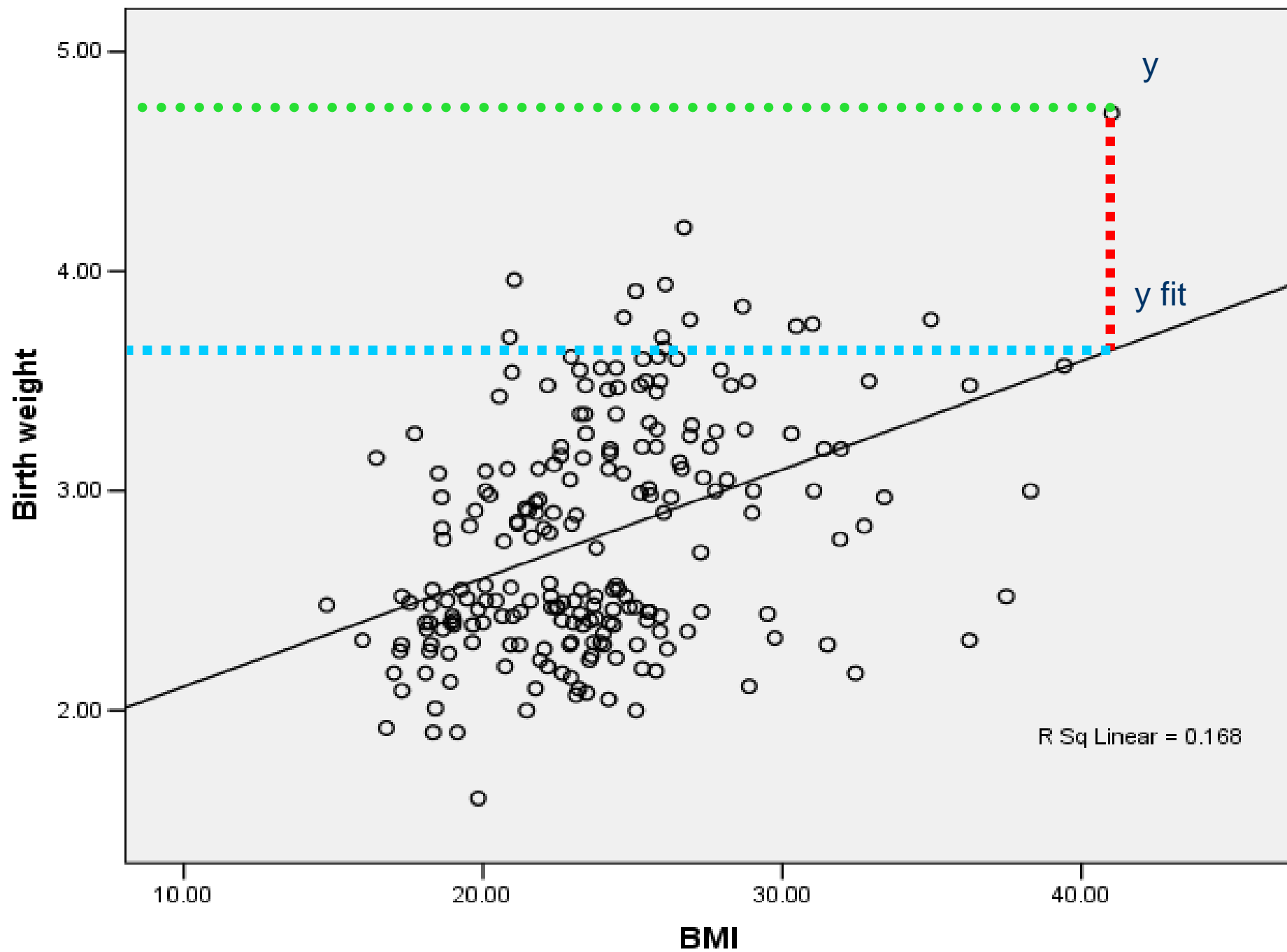| nores | chol x | bps1 y | x2 | y2 | xy |
|---|---|---|---|---|---|
| 234 | 162 | 118 | 26244 | 13924 | 19116 |
| 235 | 210 | 126 | 44100 | 15876 | 26460 |
| 238 | 239 | 105 | 57121 | 11025 | 25095 |
| 240 | 187 | 112 | 34969 | 12544 | 20944 |
| 243 | 181 | 99 | 32761 | 9801 | 17919 |
| 244 | 180 | 99 | 32400 | 9801 | 17820 |
| 245 | 156 | 110 | 24336 | 12100 | 17160 |
| 274 | 191 | 133 | 36481 | 17689 | 25403 |
| 248 | 203 | 134 | 41209 | 17956 | 27202 |
| 253 | 169 | 129 | 28561 | 16641 | 21801 |
| 255 | 221 | 140 | 48841 | 19600 | 30940 |
| 256 | 223 | 117 | 49729 | 13689 | 26091 |
| 259 | 269 | 137 | 72361 | 18769 | 36853 |
| 231 | 151 | 164 | 22801 | 26896 | 24764 |
| 232 | 151 | 164 | 22801 | 26896 | 24764 |
| 233 | 249 | 164 | 62001 | 26896 | 40836 |
| 236 | 206 | 156 | 42436 | 24336 | 32136 |
| 237 | 252 | 147 | 63504 | 21609 | 37044 |
| 239 | 219 | 186 | 47961 | 34596 | 40734 |
| 241 | 129 | 170 | 16641 | 28900 | 21930 |
| 242 | 150 | 170 | 22500 | 28900 | 25500 |
| 246 | 194 | 176 | 37636 | 30976 | 34144 |
| 247 | 164 | 186 | 26896 | 34596 | 30504 |
| 249 | 223 | 157 | 49729 | 24649 | 35011 |
| 250 | 264 | 142 | 69696 | 20164 | 37488 |
| 251 | 232 | 159 | 53824 | 25281 | 36888 |
| 252 | 165 | 144 | 27225 | 20736 | 23760 |
| 254 | 232 | 155 | 53824 | 24025 | 35960 |
| 257 | 286 | 162 | 81796 | 26244 | 46332 |
| 258 | 180 | 151 | 32400 | 22801 | 27180 |
| 260 | 198 | 164 | 39204 | 26896 | 32472 |
| 261 | 190 | 155 | 36100 | 24025 | 29450 |
|  | 6426 | 4631 | 1338088 | 688837 | 929701 |

# Testing for significance

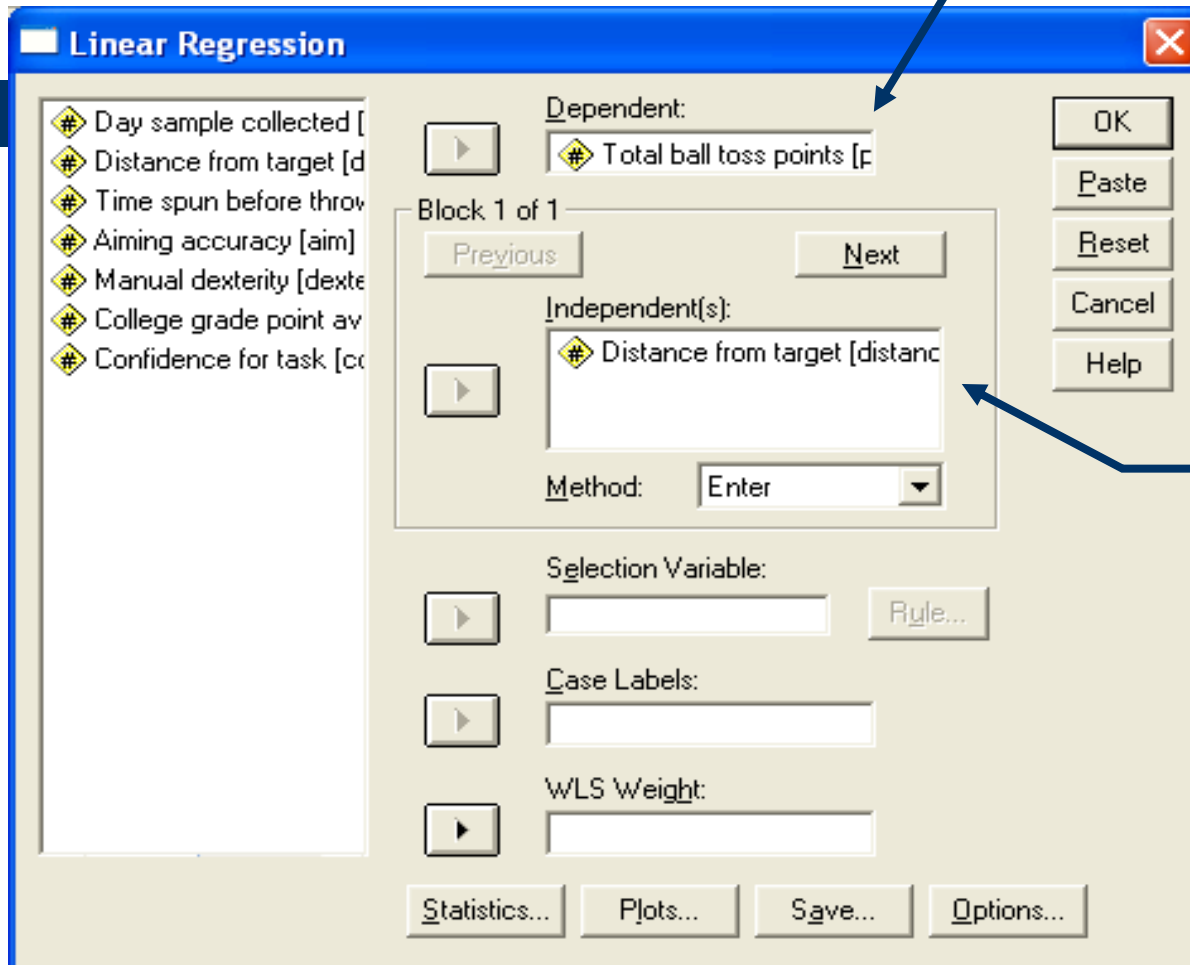test whether the slope is significantly different from zero by:

$$t = b/SE(b)$$

$$SE_{(b)} = \frac{S_{res}}{\sqrt{\Sigma(x - \bar{x})^2}}$$

$$S_{res} = \sqrt{\frac{\Sigma(y - y_{fit})^2}{n - 2}}$$

| | index | BMI | birth wgt | yfit | ytola kyfit | var |
|---|---|---|---|---|---|---|
| 1 | 1 | 32.44 | 2.17 | 3.20 | 1.07 | |
| 2 | 2 | 20.74 | 2.20 | 2.63 | .19 | |
| 3 | 3 | 22.04 | 2.28 | 2.70 | .17 | |
| 4 | 4 | 14.77 | 2.48 | 2.34 | .02 | |
| 5 | 5 | 18.33 | 1.90 | 2.51 | .38 | |
| 6 | 6 | 19.03 | 2.41 | 2.55 | .02 | |
| 7 | 7 | 27.29 | 2.45 | 2.95 | .25 | |
| 8 | 8 | 21.00 | 2.43 | 2.64 | .05 | |
| 9 | 9 | 18.92 | 2.40 | 2.54 | .02 | |

## SPSS Regression Set-up

**"Criterion,"**
- y-axis variable,
- what you're trying to predict

**"Predictor,"**
- x-axis variable,
- what you're basing the prediction on



**Linear Regression**

- Day sample collected [
- Distance from target [d
- Time spun before throw
- Aiming accuracy [aim]
- Manual dexterity [dexte
- College grade point av
- Confidence for task [co

Dependent:
- Total ball toss points [p

Block 1 of 1

Previous | Next

Independent(s):
- Distance from target [distanc

Method: Enter

Selection Variable:
Rule...

Case Labels:

WLS Weight:

OK | Paste | Reset | Cancel | Help

Statistics... | Plots... | Save... | Options...

# Getting Regression Info from SPSS

**Model Summary**

| Model | R | R Square | Adjusted R Square | Std. Error of the Estimate |
|---|---|---|---|---|
| 1 | .777[a] | .603 | .581 | 18.476 |

a. Predictors: (Constant), Distance from target

$y' = a + b (x)$

$y' = 125.401 - 4.263(20)$

a

**Coefficients[a]**

| Model | | Unstandardized Coefficients | | Standardized Coefficients | t | Sig. |
|---|---|---|---|---|---|---|
| | | B | Std. Error | Beta | | |
| 1 | (Constant) | 125.401 | 14.265 | | 8.791 | .000 |
| | Distance from target | -4.263 | .815 | -.777 | -5.230 | .000 |

a. Dependent Variable: Total ball toss points

b

# Birthweight=1.615+0.049mBMI

**Coefficients**[a]

| Model | | Unstandardized Coefficients | | Standardized Coefficients | t | Sig. |
|---|---|---|---|---|---|---|
| | | B | Std. Error | Beta | | |
| 1 | (Constant) | 1.615 | .181 | | 8.909 | .000 |
| | BMI | .049 | .007 | .410 | 6.605 | .000 |

a. Dependent Variable: Birth weight