# Data Management & Analysis: Intermediate PASW Topics II Workshop

1

## PASW STATISTICS V 17.0 (SPSS FOR WINDOWS)

Beginning, Intermediate & Advanced Applied Statistics

Zayed University Office of Research SPSS for Windows®
Workshop Series

Presented by

Dr. Maher Khelifa
Associate Professor
Department of Humanities and Social Sciences
College of Arts and Sciences

PASW Statistics 17 (2).lnk

# Data Analysis Using SPSS

2

## MEASURING RELATIONSHIPS BETWEEN VARIABLES (CORRELATIONS)

# Topics to be Covered

- Correlational vs. Regression model.
- Correlational model
  - Bivariate Correlations
    - Pearson Product Moment Correlation Coefficient
    - Spearman Rank Order Correlations
    - Kendall's tau
  - Partial correlation
    - Partial correlations & Zero-order correlations

# Correlational Analysis

- Correlational analysis: Any statistical procedures used for assessing the extent of a relationship between 2 variables.

- A correlational index is typically used to identify the extent to which two variables correlate or co-vary.

- The correlational method can be applied when data have been collected on two different variables X & Y.

- Each case has 2 scores, one on X and one on Y.

# Correlational Analysis

- The correlational model is different from the regression model.
  - Correlational model
    - Seeks to establish relationships.
    - Does not distinguish between variables X & Y.
  - Regression model
    - Seeks prediction and causality
      - Causal relationships are investigated using experimental research.
    - Distinguishes between X (predictor) and Y (criterion) variables.
      - Changes in X are used to predict changes in Y.

# Correlational Model

6

## BIVARIATE CORRELATIONS

# Bivariate Correlations

- Bivariate correlations explores pairwise relationships for a set of scale or ordinal variables.

- In Bivariate Correlations the test measures if variables or rank orders are related. Results are displayed in a correlation matrix.

- In SPSS, depending on the level of measurement, three different tests of significance can be computed under Bivariate Correlations procedure:
  - Pearson Product Moment Correlation Coefficient,
  - Spearman Rank Order Correlations, and
  - Kendall's tau.

- Bivariate tests of significance are very useful in determining the strength and direction of associations between variables (for scale and ordinal variables only).

# Bivariate Correlation

8

## THE PEARSON PRODUCT MOMENT CORRELATION COEFFICIENT

# The Pearson Product-Moment Correlation ($r$)

- The Pearson Product-Moment Correlation Coefficient ($r$) assesses the degree that two quantitative variables are linearly related.

- The $r$ provides information on the direction and magnitude of an observed correlation between two variables (X and Y).

- $rxy = \dfrac{\sum zxi\ zyi}{n\text{-}1}$

- This formula requires transforming X and Y to Z scores, then summing the cross-products over all individuals in the sample. The sum is divided by n-1.

- $z = \dfrac{X - \overline{X}}{s}$

# The Pearson Product-Moment Correlation (*r*)

- The correlation coefficient, *r,* ranges from -1 to +1.

- **Value of r Interpretation**
    - r= 0          The two variables do not vary together at all.
    - 0 > r > 1          The two variables tend to increase or decrease together.
    - r = 1.0          Perfect correlation.
    - -1 > r > 0          One variable increases as the other decreases.
    - r = -1.0          Perfect negative or inverse correlation.

# The Pearson Product-Moment Correlation (*r*)

- Guidelines for interpreting *r*

| Range of Positive Coefficient | Type of Relationship | Range of Negative Coefficient |
|---|---|---|
| .00 to .20 | low | .00 to -.20 |
| .21 to .40 | Low to moderate | -.21 to -.40 |
| .41 to .60 | Moderate | -.41 to -.60 |
| .61 to .80 | Moderate to high | -.61 to -.80 |
| .81 to 1.00 | High | -.81 to -1.00 |

•What is large or small relationship depends on the researcher's discipline.

•In Social Sciences correlation coefficients of .10, .30, and .50 (irrespective of sign) are interpreted as small, medium, and large relationships, respectively.

# How to Obtain a Person *r*

- Go to Analyze, Correlate, Bivariate.

- A correlation dialogue box will appear.

# How to Obtain a Person *r*

- Move the variables you want to correlate to the variable box. Then select Pearson.
- Select options and select means and SD and press continue.
- Press OK.

# How to Obtain a Person *r*

A correlation matrix will appear in the Output window.

Check the Pearson correlation coefficient and the level of significance.

**Correlations**

|  |  | quads | gluts |
|---|---|---|---|
| quads | Pearson Correlation | 1 | .484** |
|  | Sig. (2-tailed) |  | .000 |
|  | Sum of Squares and Cross-products | 9211.640 | 2671.520 |
|  | Covariance | 93.047 | 26.985 |
|  | N | 100 | 100 |
| gluts | Pearson Correlation | .484** | 1 |
|  | Sig. (2-tailed) | .000 |  |
|  | Sum of Squares and Cross-products | 2671.520 | 3311.360 |
|  | Covariance | 26.985 | 33.448 |
|  | N | 100 | 100 |

**. Correlation is significant at the 0.01 level (2-tailed).

# Effect Size

- To properly interpret the value of the Pearson $r$, *the $r$* is squared to calculate $r^2$ (*r* squared) also called the *coefficient of determination.*

- *$r^2$ ranges from 0 to 1.*

- *$r^2$ is the* variance shared between X and Y.
  - For example, if $r^2 = 0.59$, then 59% of the variance in X can be explained by variation in Y, and 59% of the variance in Y can be explained by variation in X.

  - So *$r^2$ is the amount of variance accounted for by the linear relationship of the 2 variables.*

# Effect Size

- The Pearson *r* scale values are not of equal intervals.
- The correlation scale is <span style="color:red">ordinal</span> not interval or ratio.
- For example:
  - A value of r = .80 does not represent twice the relationship as r = .40.
  - The difference in a relationship of .40 to .50 does not represent the same magnitude of relationship difference as .70 to .80.
- If we use r² descriptive <span style="color:red">comparative statements</span> can be made.
- The r² provides an added dimension for interpreting and understanding the magnitude of a linear relationship.

# Assumptions of the Pearson *r*

- There are 2 assumptions underlying the Pearson *r* test of significance.

- 1. Assumption of Independence:
  - Cases represent a random sample from the population and the scores on variables for one case are independent of scores of other cases.
    - The test is not robust to violations of this assumption and therefore the Pearson *r* test of significance should not be computed.

# Assumptions of the Pearson *r*

- 2. <span style="color:red">Normality Assumptions</span>:

- The variables are bivariately normally distributed.

  ○ Pearson correlation calculations are based on the assumption that both X and Y values are sampled from populations that follow a normal distribution, at least approximately.

  ○ the Pearson correlation coefficient works best when the variables are approximately normally distributed and have no outliers.

  ○ With large samples, the test is robust to violations of this assumption.

  ○ If the normality assumption is met, the only type of statistical relationship that can exist between 2 variables is a linear association.

# Assumptions of the Pearson *r*

- 2. Normality Assumptions: (Continued)

- The Pearson's *r* is a measure of linear association only. However, two variables can be perfectly related but not in a linear fashion. The Pearson's correlation coefficient is not an appropriate statistic for measuring non-linear associations.

- Before calculating *r*, data should be screened for outliers and evidence of a linear relationship. A scatterplot can reveal possible problems.

- Outliers may cause misleading results.

- If the normality assumption is not made or met, alternate distribution-free tests should be used instead.

# How to Obtain A Scatter Plot

- Go to Graphs,
- Chart builder.

# How to Obtain A Scatter Plot

- From Gallery, chose Scatter/dot, then drag and drop the type of graph you want into the Chart preview
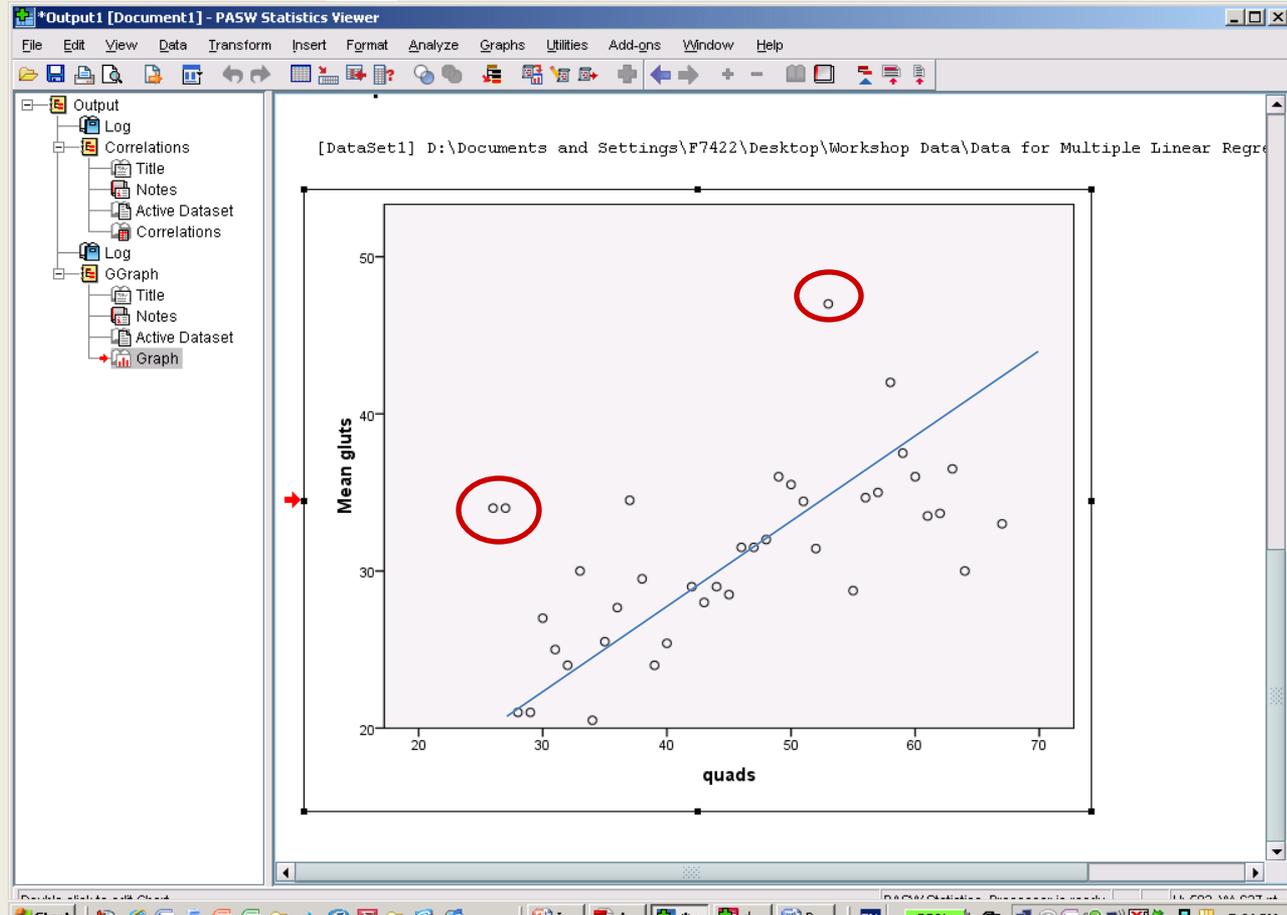
# How to Obtain A Scatter Plot

- Drag and drop the selected variables to correlate into the X and Y axes. Then Press OK.
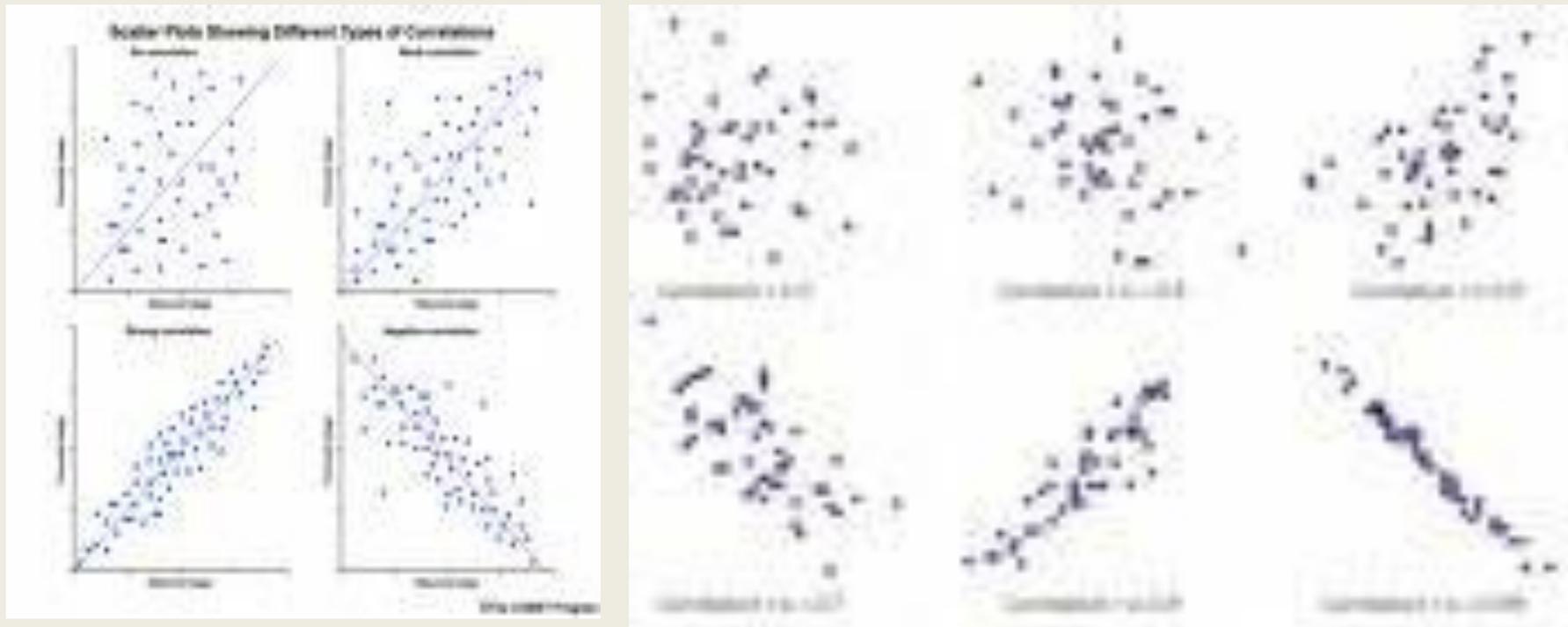
# How to Obtain A Scatter Plot

- A scatter plot is obtained in the Output window.
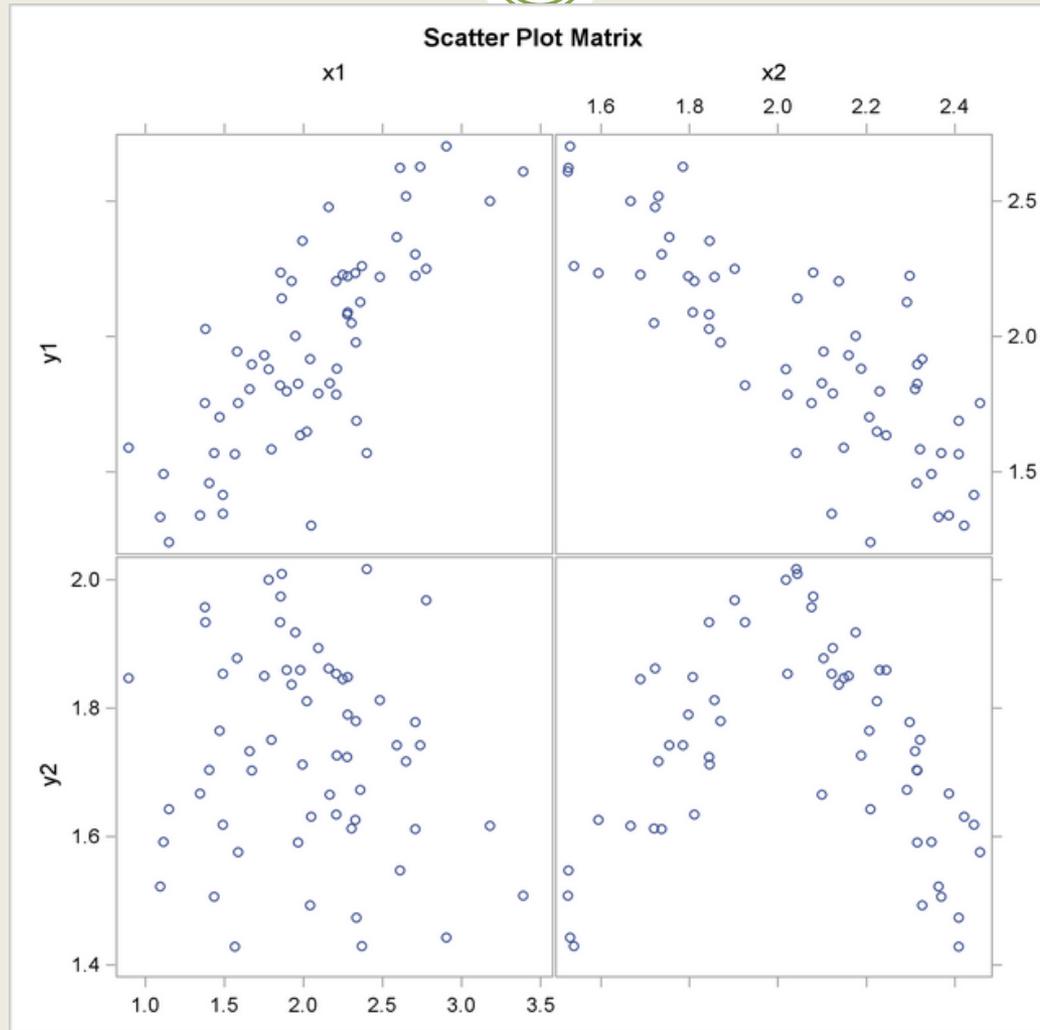- Look for outliers and serious deviations from a linear relationship pattern.
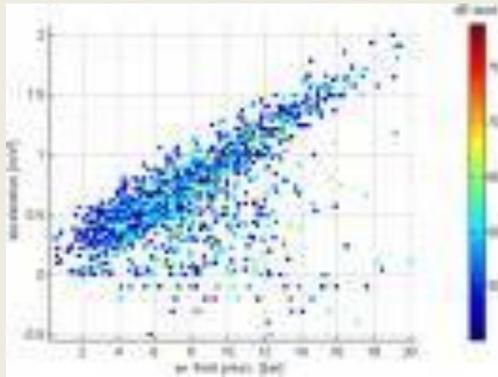
# Example of Relationship Scatterplots

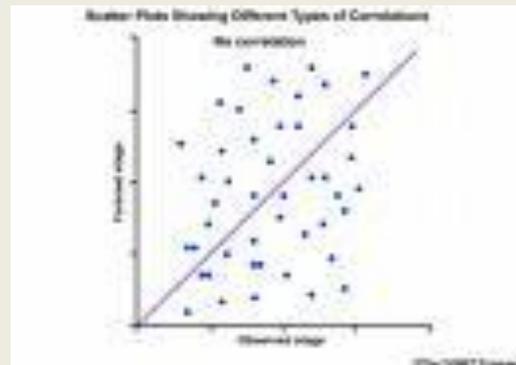# Example of Relationship Scatterplots
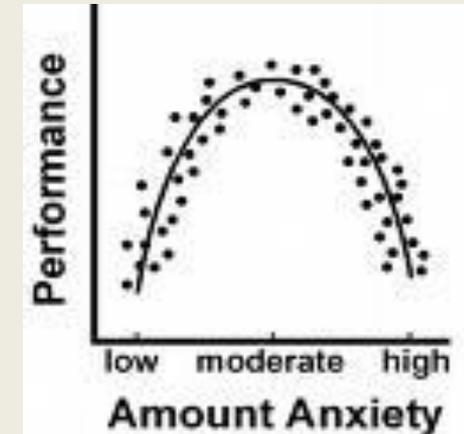
# Example of Relationship Scatterplots

•If a relationship is less than perfect, the scatterplot of data points will deviate from a perfect straight line.

# Non Linear Relationships

- Two variables can be perfectly related, but in a non-linear fashion.

- When this is the case, the Pearson's correlation coefficient is not an appropriate statistic for measuring their association.

- Scan the plots for relationships other than linear.

- The attached plots indicate non-linear relationships (curvilinear).

# Bivariate Correlation

28

## THE SPEARMAN RANK CORRELATION COEFFICIENT
## &
## THE KENDALL'S RANK CORRELATION COEFFICIENT

# The Spearman's rho & Kendall's tau-*b*

- If the data are not normally distributed, or have ordered categories, the Spearman's rho and Kendall's tau-*b* statistics may offer good alternatives to the Pearson *r*.

- The Spearman's rho and Kendall's tau-b measure the association between rank-orders (ranked variables).

- The Spearman's rho and Kendall's tau-b work regardless of the distributions of the variables. They are considered distribution-free or non-parametric tests.

- Spearman's rho and Kendall's tau-b can be computed on quantitative variables or variables with ordered categories, i.e.:
  - on 2 scale variables,
  - on 2 ordinal variables, on
  - on one scale and one ordinal variable.

# The Spearman's rho ($r_s$)

- The Spearman correlation coefficient is a correlation coefficient between 2 ranked variables.

- Raw scores in the 2 variables are converted into ranks and differences between the ranks for each observation on the 2 variables are calculated using the following formula.

$$( R ) = 1 - \frac{6 \sum d^2}{n^3 - n}$$

# Spearman rho: Example

$$(R) = 1 - \frac{6\sum d^2}{n^3 - n}$$

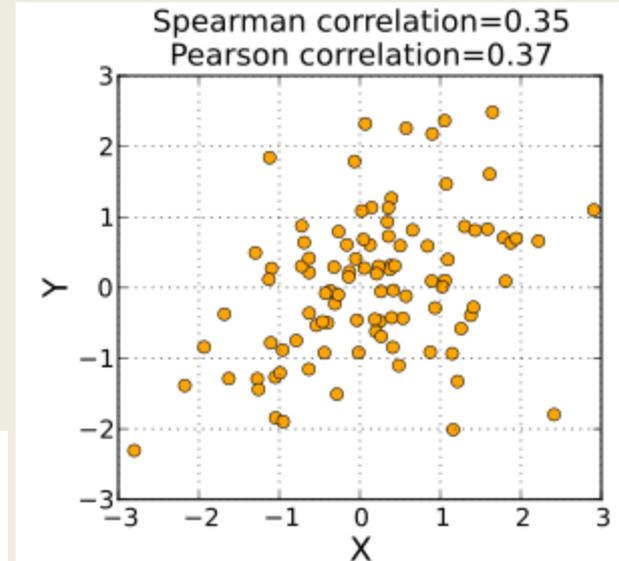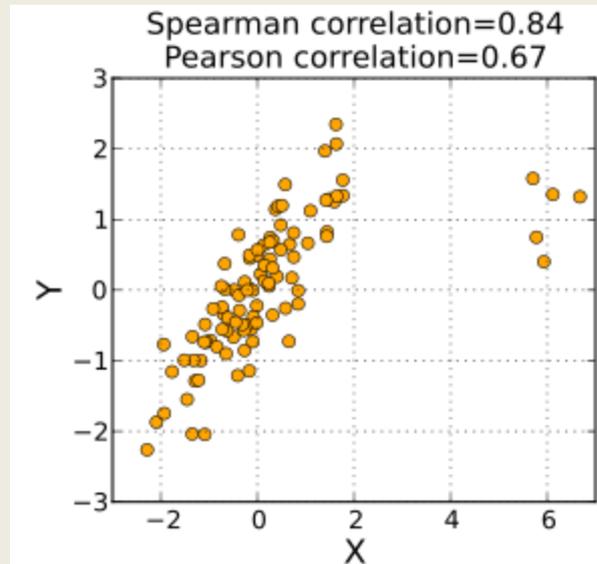| Convenience Store | Distance from CAM (m) | Rank variable distance | Price of 50cl bottle (€) | Rank variable price | Difference between ranks (d) | d² |
|---|---|---|---|---|---|---|
| 1 | 50 | 10 | 1.80 | 2 | 8 | 64 |
| 2 | 175 | 9 | 1.20 | 3.5 | 5.5 | 30.25 |
| 3 | 270 | 8 | 2.00 | 1 | 7 | 49 |
| 4 | 375 | 7 | 1.00 | 6 | 1 | 1 |
| 5 | 425 | 6 | 1.00 | 6 | 0 | 0 |
| 6 | 580 | 5 | 1.20 | 3.5 | 1.5 | 2.25 |
| 7 | 710 | 4 | 0.80 | 9 | -5 | 25 |
| 8 | 790 | 3 | 0.60 | 10 | -7 | 49 |
| 9 | 890 | 2 | 1.00 | 6 | -4 | 16 |
| 10 | 980 | 1 | 0.85 | 8 | -7 | 49 |
| | | | | | | **d² = 285.5** |

© Dr. Maher Khelifa

# Interpreting the Spearman rho

- The Spearman rho varies between -1 and +1 and is interpreted in the same fashion as the Pearson *r*.

# Spearman rho & Outliers

- Outliers have less of an effect on Spearman's rho.



Spearman correlation=0.35
Pearson correlation=0.37



Spearman correlation=0.84
Pearson correlation=0.67

# The Kendall's tau-*b*

- Like the Spearman rho, the Kendall's tau ($\tau$) coefficient is a distribution-free statistic used to measure the statistical association between two variables.

- The Kendall's tau ($\tau$) is a measure of rank correlation.

- The test examines the similarity of the orderings of the ranked data.

  - If the two variable rankings are identical, the coefficient has value 1.

  - If one ranking is the perfect reverse of the other, the coefficient has value −1.

  - Increasing coefficient values indicate more concordance between the rankings.

  - If the coefficient has value of 0, this implies that the rankings are completely independent.

# Kendall's tau-*b*

- $\tau = \dfrac{n_c - n_d}{n(n-1)/2}$

- *Where $n_c$* is the number of <span style="color:red">concordant pairs</span>, and $n_d$ is the number of <span style="color:red">discordant pairs</span> in the data set.

- The range of the Kendall's tau is -1 to +1, and the tau coefficient is interpreted the same way as the Spearman rho and Pearson r.

# How to obtain the Spearman's rho and Kendall's tau-*b*

- In PASW, Spearman's rho and Kendall's tau-b are obtained the same way as the Pearson *r*. Just select the test you are interested in (Pearson, rho, and tau-*b*).

- PASW allows you to select all tests at the same time. The three tests generally yield similar alphas (not identical values though).

# Comapring the Bivariate Tests

- The Pearson $r$ is more powerful if the data satisfies the test assumptions (normality and independence).

- More predictive models are available for linear relationships, and the linear models are generally easier to implement and interpret.

- The Pearson $r$ is also used extensively in measurement, in test construction, and in research to determine instrument reliability:
  - test-retest reliability,
  - equivalence reliability, and
  - internal consistency reliability.

- The distribution-free measures of associations are handy for discovering whether there is any kind of association between two variables with ranked data.

- They distribution-free measures of associations are especially useful if the data is not normally distributed or there are outliers.

# Alternate tests to Bivariate Correlations

- The Bivariate Correlation procedures discussed above are useful for studying the pairwise associations for a set of scale or ordinal variables.

- If the variables are nominal, use the Crosstabs procedure to obtain measures of association including Chi square, contingency coefficient, and Phi and Cramer's V.

- If you want to predict the value of a scale variable based on its linear relationship to other variable(s), use Linear Regression procedure.

- If you want to explore the variation in your data to look for underlying patterns, use Factor Analysis.

# Partial Correlations

# Partial Correlations ($r_p$)

- Any time the relationship between 2 variables are examined, one must be concerned with the effects of other variables on the relationship of interest.

- Example: Education and income (controlling age and work experience).

- Partial correlation coefficients is a procedure that describes the linear relationship between two variables while controlling for the effects of one or more additional variables.

- $r_p$ is an effect size index which indicates the degree that two variables are linearly related in a sample, partialling out the effects of one or more variables.

# Partial Correlations ($r_p$)

- The partial correlation attempts to estimate the correlation between 2 variables holding constant the values of the control variable(s).

- The significance test evaluates whether in the population the partial correlation is equal to zero.

- To interpret properly the partial correlations results, zero-order correlations must be considered.

- The data file for a partial correlation should contain scores on at least 3 variables. All the variables should be at the scale level of measurement.

- Results of the partial correlation, like the Pearson r, range in value between -1 and +1 and are interpreted in a similar manner.

# Usefulness of Partial Correlation

- Partial correlation is very useful for:
  - uncovering spurious relationships, and
  - Detecting hidden relationships

- Spurious correlation: When two or more variables are statistically related but are not in fact causally linked—usually because the statistical relation is caused by a third variable (usually the true predictor). When the effects of the third variable are removed, they are said to have been partialed out.

- Spurious correlations mean that 2 things look like they cause each other, but in reality they don't (a correlation between A & B is actually caused by a third variable C, usually referred to as lurking variable, confounding and moderator variable).

- Hidden correlation: Theory, intuition and common sense may suggest that there should be a relationship between two variables even though the data suggests no correlation. It may be the case that one or more variables are suppressing the expected relationship.

# How to obtain a Partial Correlation

- Select Analyze, Correlate, Partial

# How to obtain a Partial Correlation

- Move the 2 variables to correlate to the variable box.

- Select the control variable.

- Select Options, and select zero-order correlations

# How to obtain a Partial Correlation

- A Partial correlation matrix appears in the Output Window.

- Look for zero-order correlations and compare them to the values of the partial correlation.

# Alternate tests

- The Partial Correlations procedure is only appropriate for scale variables.

- If you have categorical (nominal or ordinal) data, use the Crosstabs procedure. Layer variables in Crosstabs are similar to control variables in Partial Correlations.

- If you want to predict the value of a scale variable based on its linear relationship to other variables, try Linear Regression.

# Thank you for your Attention