# Data Mining



Lecture 2: data and data preprocessing

# Data mining analysis steps CRISP-dm

1. Determining the purpose of the analysis,
2. Preliminary selection of variables, data cleaning (preprocessing),
3. Analytical form selection, on the basis of descriptive data analysis, correlation matrix, and knowledge,
4. Final variable and model selection,
5. Model parameters estimation,
6. Model verification (statistical, substantive),
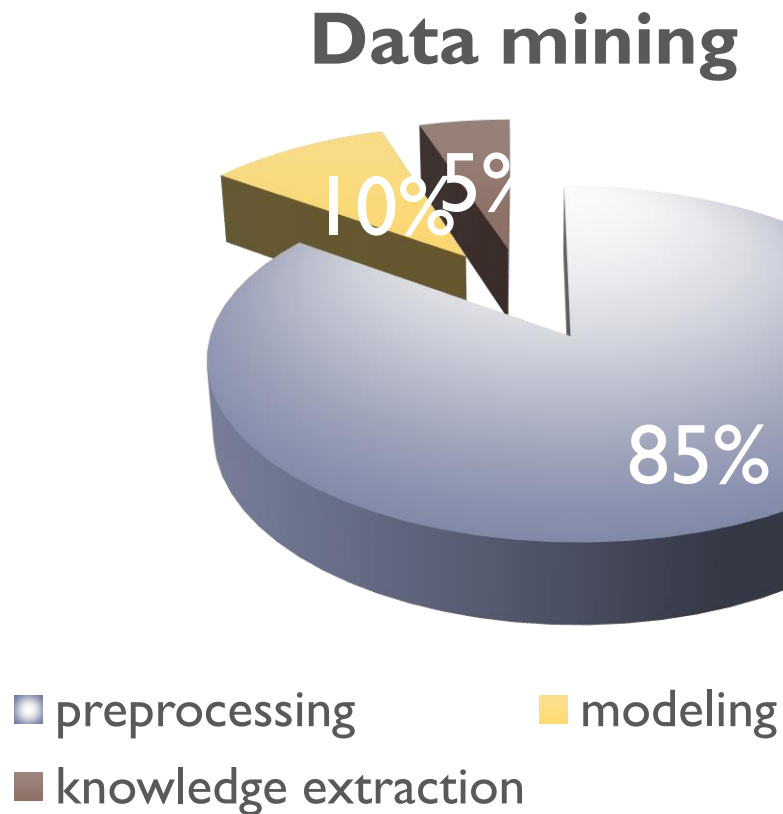7. Practical model uses (dependence analysis, prediction).

# GIGO



Garbage in- Garbage out

# Time-consuming data mining phases

**Data mining**



- preprocessing
- modeling
- knowledge extraction

10% 5%

85%

# Data preprocessing

- Preparation of data variables from the preliminary, raw data and the final data set that will be used in all subsequent phases.

- Cases selection and variables that will be analysed and which are suitable for analysis.

- Variable transformations, if necessary.

- Clear raw data so that it is ready to be used by modelling tools.

# Motivation of data preprocessing

- Preparation of data variables from the preliminary, raw data and the final data set that will be used in all subsequent phases.

- Cases selection and variables that will be analysed and which are suitable for analysis.

- Variable transformations, if necessary.

- Clear the raw data so that it is ready to be used by the modelling tools.
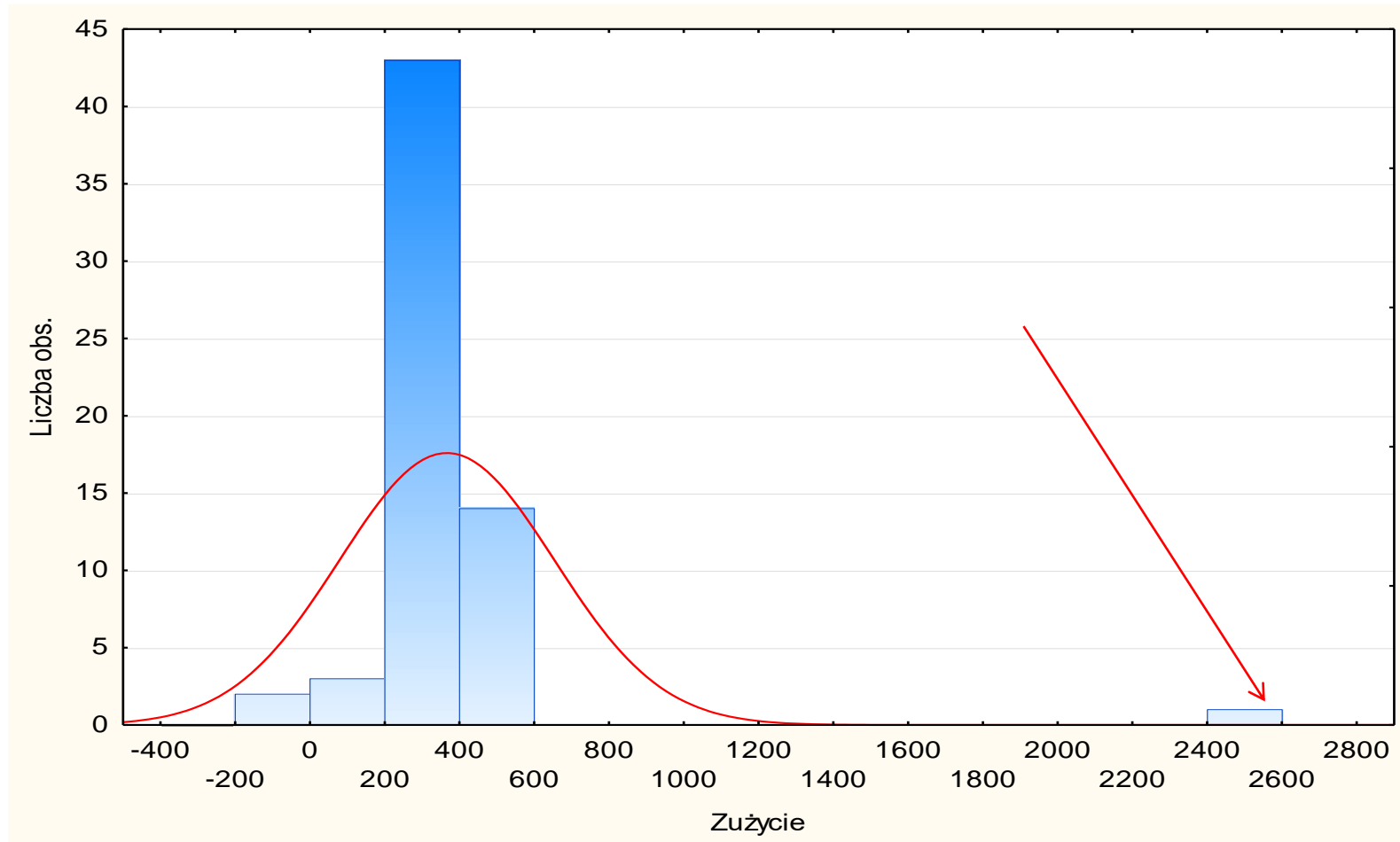
# Data cleaning

- Analytical and graphical detection of erroneous and unusual observations,

- Handling and replacing missing data,

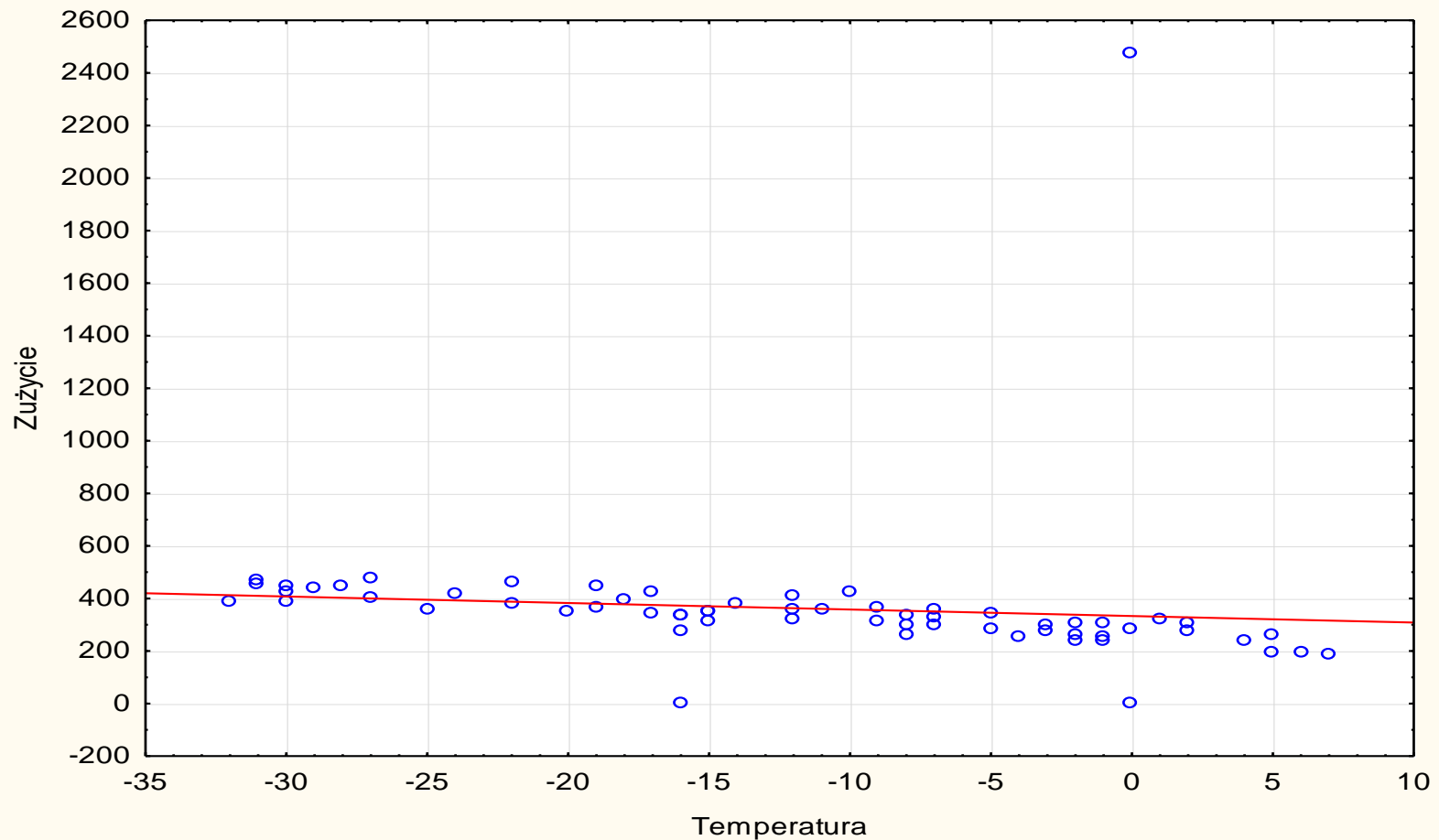- Identification and removal of duplicate records.

# Analytical and graphical detection of erroneous and unusual observations

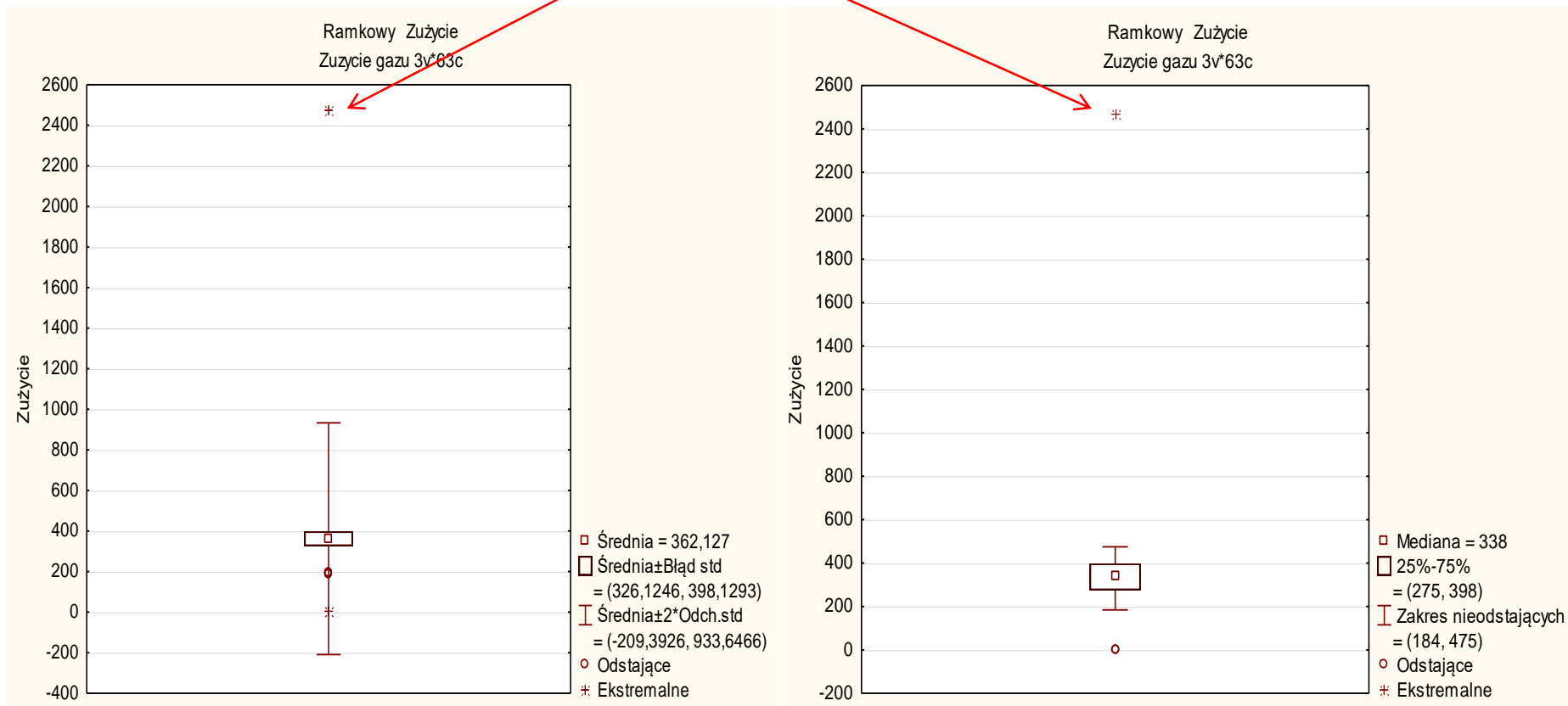# Analytical and graphical detection of erroneous and unusual observations

# Analytical and graphical detection of erroneous and unusual observations

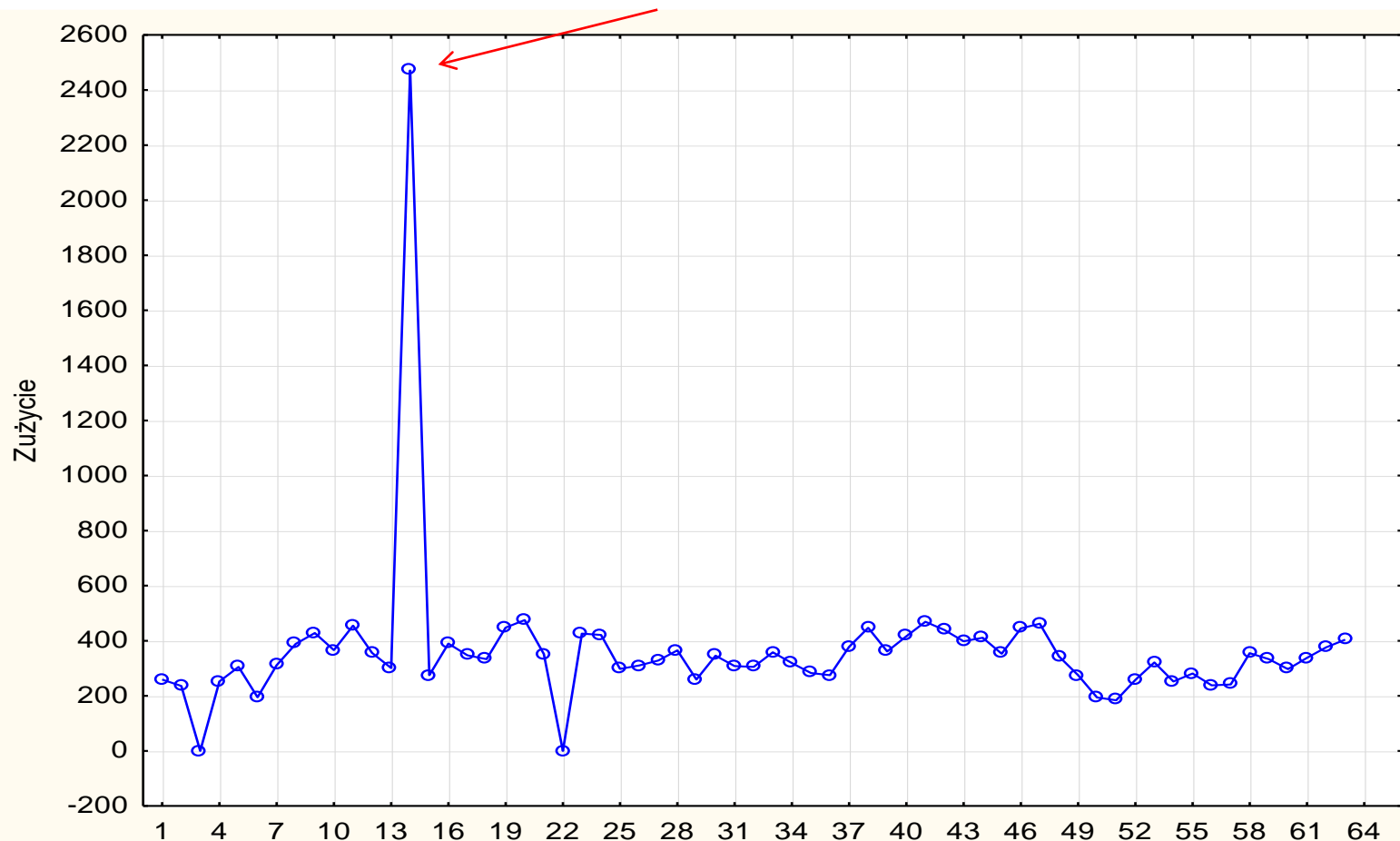| Descriptive statistics | | | | | | |
|---|---|---|---|---|---|---|
| | N | Mean | Median | Min. | Max. | Std. Dev. |
| Dataset | 63 | 362,1270 | 338,0000 | 0,00 | 2471,000 | 285,7598 |

# Analytical and graphical detection of erroneous and unusual observations

**Extreme values**

# Analytical and graphical detection of erroneous and unusual observations
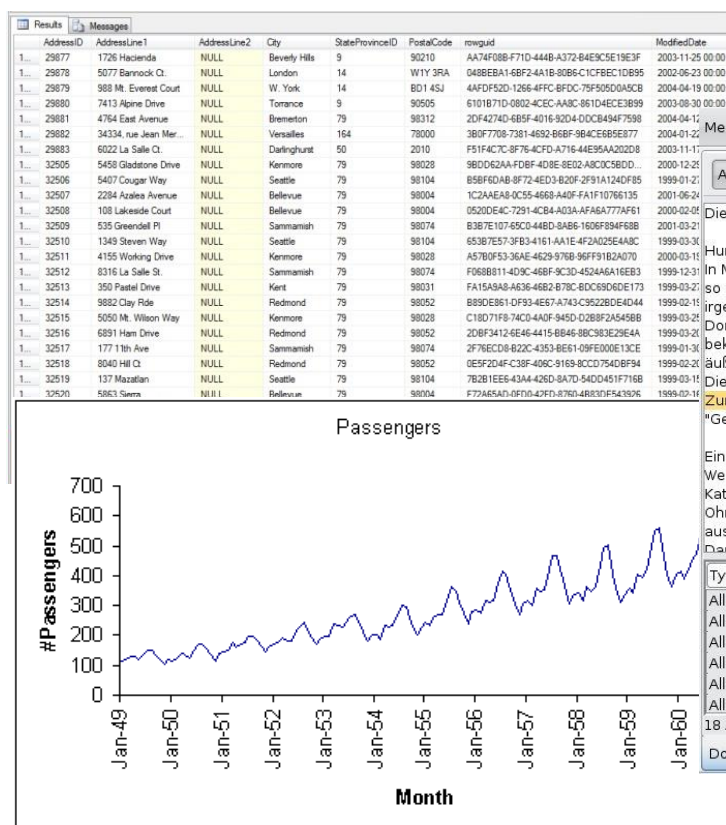
Extreme values

# Missing data

Methods for missing data:

▸ Omit rows with empty values,

▸ Replacement of the missing value determined by the analyst,

▸ Replace with mean or median - for numbers,

▸ Modal value substitution for qualitative variables,

▸ Generation of a random value from the observed distribution of the variable,

▸ Link the blank data to the rest of the object data and give the most probable value.

# Missing data

Omit rows with empty values- only for „table data",
Omit columns with empty values- user criterion,
UDL, LDL (ADL, UDL), right censored, left censored.

# Missing data

Replacement of the missing value determined by the analyst

# Data transformations

Calculation of derived variables:

- Differentiation of time series,

- Logarithms,

- Roots,

- Calculation of indicators.

# Data transformations

Change assignment category (transcoding)

# Data transformations

Rank assignment

# Data transformations

Latent variables

# Data transformations

Text variables operations:

- Counting,
- Conversion to numbers (recoding),
- Complete the information,
- Analysis of descriptions.

# Data analysis

It is a knowledge discovery about:

- pure data,  data type,
- data collection,
- data preparation,
- factors that affect the development of the phenomenon,
- seasonality, trend, events and incidents.

# Data in Matrix Form

This data type having one (dependent) variable described by one or many (independent /predictors) variables.

Example: Runners finish time on 100 meters (dependent variable) described by: height, weight, number of workouts, leg length (independent variables).

Regression problems can also be written as binary variables (they replace all variables that are not quantitative).

# Time-Series Data

Time-series database

- Consists of sequences of values or events changing with time,

- Data is recorded at regular intervals,

- Characteristic time-series components,

    **Trend, cycle, seasonal, irregular.**

Applications

- Financial: stock price, inflation,

- Industry: power consumption,

- Scientific: experiment results,

- Meteorological: precipitation.

# Categories of Time-Series Movements

Categories of Time-Series Movements

- Long-term or trend movements (trend curve): general direction in which a time series is moving over a long interval of time,

- Cyclic movements or cycle variations: long term oscillations about a trend line or curve,

  e.g., business cycles, may or may not be periodic.

- Seasonal movements or seasonal variations,

  almost identical patterns that a time series appears to follow during corresponding months of successive years.

- Irregular or random movements,

Time series analysis: decomposition of a time series into these four basic movements:

- Additive Modal: TS = T + C + S + I

- Multiplicative Modal: TS = T $\times$ C $\times$ S $\times$ I

# Categories of Time-Series Movements

- Long-term or trend movements,

- Average level,

Populacja USA w latach 1790-1980

# Categories of Time-Series Movements

- Cyclic movements or cycle variations

M1: dane sezonowe bez trendu

# Categories of Time-Series Movements

▸ Irregular or random movements

# Aim of the Time series analysis

- Modelling of a certain phenomenon / process based on observed changes in some measurable quantities describing this process,

- Isolation and measurement of time series components (decomposition of time series),

- Prediction of future values using the obtained model.

# How to get information about data?

- Descriptive analysis,

- Data visualizations: graphs, plots, histograms.

Price and market penetration of consumer electronics over the past 50 years

# Data analysis, what to do?

1. Plots: linear, spot, column, box and whiskers,

2. Descriptive statistics (mean, median, max., min., std. dev.,…), tables, crosstabulation tables,

3. Histograms, normality graph,

4. ACF and PACF graphs, Fourier analysis – for time series,

5. Correlation matrices,

6. Cluster analysis.

# Descriptive statistics

| Descriptive statistics (chronologie jesiony) | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | N | Mean | Median | Mode | Min. | Max. | Variance | Std. dev. | Variation coefficient | Skewness |
| ry1fr_n | 116 | 0,977 | 0,961 | Wielokr. | 0,499 | 1,650 | 0,033 | 0,181 | 18,486 | 0,249 |
| ry1fr_r | 116 | 1,002 | 1,003 | Wielokr. | 0,547 | 1,632 | 0,029 | 0,172 | 17,156 | 0,365 |
| ry1fr_c | 116 | 0,963 | 0,954 | Wielokr. | 0,236 | 1,754 | 0,106 | 0,326 | 33,862 | -0,193 |

# Histograms -shapes

# Data analysis



Histogram   QTY
Napój gazowany 21v*3144c
QTY = 3114*50*normal(x; 72,8651; 37,8077)

# Data analysis

## Histogram - numbers of bars

# Histograms- numer of bars

Juran's Quality Control Handbook provides these guidelines for the number of bars and states that they are not "rigid" and should be adjusted when necessary.

| Number of Data Points | Number of Bars | | Number of Data Points | Number of Bars |
|---|---|---|---|---|
| 20-50 | 6 | | 201-500 | 9 |
| 51-100 | 7 | | 501-1000 | 10 |
| 101-200 | 8 | | 1000+ | 11-20 |

$$k < 5\log n \qquad k \approx \sqrt{n} \qquad k \approx 1+3.3\log n \qquad k \leq \frac{n}{10}$$

Szerokość klasy: $\approx (x_{max}-x_{min})/k$

# Data analysis



Liniowy wiele zmiennych
chronologie jesiony 13v*122c

# Data analysis

# Data analysis



Air Passengers from 1949 to 1961

# Data analysis



Wykres rozrzutu   QTY względem Price
Napój gazowany 21v*3144c
QTY = 42,2709+10,0118*x

# Data analysis



Wykres rozrzutu ry1fr_rcs względem ry1fr_ne
chronologie jesiony 13v*122c
ry1fr_rcs = 0,1795+0,8418*x

# Data analysis



Ramkowy  wiele zmiennych
chronologie jesiony 13v*122c
Mediana; Współczynnik: 25%-75%; Wąs: Zakres nieodstających

# Data analysis

# Data analysis



Podsumowanie:QTY: ilość sprzedanyh w tyś.

K-S d=,10848, p<,01 ; Lilliefors p<,01
— Oczekiwana normalna

Wykres normalności: QTY

Statystyki:QTY
N ważnych=3114,000000
Średnia= 72,865125
Minimum=  4,000000
Maksimum=552,000000
Odch.std= 37,807717

Średnia = 72,8651
Średnia±Odch.std
= (35,0574, 110,6728)
Średnia±1,96*Odch.std
= (-1,238, 146,9683)

# Autocorrelation

**Autocorrelation correlogram.** Seasonal patterns of time series can be examined via correlograms. The correlogram (autocorrelogram) displays graphically and numerically the autocorrelation function (*ACF*), that is, serial correlation coefficients (and their standard errors) for consecutive lags in a specified range of lags (e.g., 1 through 30). Ranges of two standard errors for each lag are usually marked in correlograms but typically the size of auto correlation is of more interest than its reliability because we are usually interested only in very strong (and thus highly significant) autocorrelations.

# Partial autocorrelations

**Partial autocorrelations.** Another useful method to examine serial dependencies is to examine the partial autocorrelation function (*PACF*) - an extension of autocorrelation, where the dependence on the intermediate elements (those *within* the lag) is removed. In other words, the partial autocorrelation is similar to autocorrelation, except that when calculating it, the (auto) correlations with all the elements within the lag are partially out.

If a lag of 1 is specified (i.e., there are no intermediate elements within the lag), then partial autocorrelation is equivalent to autocorrelation. In a sense, the partial autocorrelation provides a "cleaner" picture of serial dependencies for individual lags (not confounded by other serial dependencies).

# Autocorelation



ACF

PACF

TREND

# Autocorelation



ACF

PACF

seasonality

# Autocorelation



ACF                                    PACF

TREND + SEASONALITY

# Autocorrelations/ Partial autocorrelations

Brak trendu i sezonowości

# Autocorrelations/ Partial autocorrelations

# Partial autocorrelations



Funkcja autokorelacji cząstkowej

QTY     : ilość sprzedanyh w tyś.

(Błędy std. przy założeniu AR rzędu k-1)

# Partial autocorrelations



Funkcja autokorelacji cząstkowej

PRICE

(Błędy std. przy założeniu AR rzędu k-1)

| Opóźn | Kor. | S.E |
|---|---|---|
| 1 | +,972 | ,0178 |
| 2 | +,059 | ,0178 |
| 3 | +,025 | ,0178 |
| 4 | +,036 | ,0178 |
| 5 | +,004 | ,0178 |
| 6 | -,021 | ,0178 |
| 7 | -,013 | ,0178 |
| 8 | +,126 | ,0178 |
| 9 | +,041 | ,0178 |
| 10 | -,008 | ,0178 |
| 11 | +,053 | ,0178 |
| 12 | +,046 | ,0178 |
| 13 | +,016 | ,0178 |
| 14 | +,015 | ,0178 |
| 15 | +,074 | ,0178 |
| 16 | +,049 | ,0178 |
| 17 | -,010 | ,0178 |
| 18 | +,073 | ,0178 |
| 19 | -,029 | ,0178 |
| 20 | +,014 | ,0178 |
| 21 | +,018 | ,0178 |
| 22 | +,022 | ,0178 |
| 23 | +,029 | ,0178 |
| 24 | +,002 | ,0178 |
| 25 | +,038 | ,0178 |
| 26 | -,023 | ,0178 |
| 27 | +,008 | ,0178 |
| 28 | +,036 | ,0178 |
| 29 | +,053 | ,0178 |
| 30 | +,013 | ,0178 |

-1,0   -0,5   0,0   0,5   1,0   — P. ufności

# Autocorrelations



Funkcja autokorelacji

PRICE

(Błędy standardowe to oceny białego szumu)

| Opóźn | Kor. | S.E | Q | p |
|---|---|---|---|---|
| 1 | +,972 | ,0178 | 2970, | 0,000 |
| 2 | +,948 | ,0178 | 5795, | 0,000 |
| 3 | +,925 | ,0178 | 8489, | 0,000 |
| 4 | +,905 | ,0178 | 111E2 | 0,000 |
| 5 | +,886 | ,0178 | 135E2 | 0,000 |
| 6 | +,865 | ,0178 | 159E2 | 0,000 |
| 7 | +,845 | ,0178 | 181E2 | 0,000 |
| 8 | +,832 | ,0178 | 203E2 | 0,000 |
| 9 | +,820 | ,0178 | 224E2 | 0,000 |
| 10 | +,808 | ,0178 | 245E2 | 0,000 |
| 11 | +,798 | ,0178 | 265E2 | 0,000 |
| 12 | +,790 | ,0178 | 285E2 | 0,000 |
| 13 | +,783 | ,0178 | 304E2 | 0,000 |
| 14 | +,777 | ,0178 | 323E2 | 0,000 |
| 15 | +,773 | ,0178 | 342E2 | 0,000 |
| 16 | +,770 | ,0178 | 361E2 | 0,000 |
| 17 | +,767 | ,0178 | 379E2 | 0,000 |
| 18 | +,766 | ,0178 | 398E2 | 0,000 |
| 19 | +,762 | ,0178 | 416E2 | 0,000 |
| 20 | +,758 | ,0178 | 435E2 | 0,000 |
| 21 | +,755 | ,0178 | 453E2 | 0,000 |
| 22 | +,752 | ,0178 | 470E2 | 0,000 |
| 23 | +,749 | ,0178 | 488E2 | 0,000 |
| 24 | +,746 | ,0178 | 506E2 | 0,000 |
| 25 | +,743 | ,0178 | 523E2 | 0,000 |
| 26 | +,739 | ,0178 | 541E2 | 0,000 |
| 27 | +,735 | ,0178 | 558E2 | 0,000 |
| 28 | +,732 | ,0178 | 575E2 | 0,000 |
| 29 | +,731 | ,0178 | 592E2 | 0,000 |
| 30 | +,730 | ,0177 | 609E2 | 0,000 |

0

-1,0    -0,5    0,0    0,5    1,0

—— P. ufności

# Partial autocorrelations



Funkcja autokorelacji cząstkowej

QTY     : ilość sprzedanyh w tyś.

(Błędy std. przy założeniu AR rzędu k-1)

| Opóźn | Kor. | S.E |
|---|---|---|
| 1 | +,606 | ,0179 |
| 2 | +,120 | ,0179 |
| 3 | -,029 | ,0179 |
| 4 | +,139 | ,0179 |
| 5 | +,137 | ,0179 |
| 6 | +,212 | ,0179 |
| 7 | +,206 | ,0179 |
| 8 | -,050 | ,0179 |
| 9 | -,126 | ,0179 |
| 10 | -,021 | ,0179 |
| 11 | +,014 | ,0179 |
| 12 | -,002 | ,0179 |
| 13 | +,179 | ,0179 |
| 14 | +,113 | ,0179 |
| 15 | -,062 | ,0179 |
| 16 | -,049 | ,0179 |
| 17 | -,007 | ,0179 |
| 18 | +,004 | ,0179 |
| 19 | +,072 | ,0179 |
| 20 | +,102 | ,0179 |
| 21 | +,065 | ,0179 |
| 22 | -,067 | ,0179 |
| 23 | -,051 | ,0179 |
| 24 | -,005 | ,0179 |
| 25 | +,019 | ,0179 |
| 26 | +,032 | ,0179 |
| 27 | +,081 | ,0179 |
| 28 | +,067 | ,0179 |
| 29 | -,051 | ,0179 |
| 30 | -,068 | ,0179 |

-1,0    -0,5    0,0    0,5    1,0    —— P. ufności

# Correlation

Correlation is a measure of the relation between two or more variables. The measurement scales used should be at least interval scales, but other correlation coefficients are available to handle other types of data. Correlation coefficients can range from -1.00 to +1.00. The value of -1.00 represents a perfect *negative* correlation while a value of +1.00 represents a perfect *positive* correlation. A value of 0.00 represents a lack of correlation.

Both datasets should have the same numer of observations.

# Correlation

# Pearson linear correlation coefficient

$$\rho_{X,Y} = corr(X,Y) = \frac{cov(X,Y)}{\sigma_X \sigma_Y} = \frac{E[(X-\mu_X)(Y-\mu_Y)]}{\sigma_X \sigma_Y}$$

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \overline{x})(y_i - \overline{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \overline{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \overline{y})^2}},$$

# Pearson linear correlation coefficient

# Nonparametric correlations

- ▸ Could be used to find nonlinear dependences in two datasets for many data types.

- ▸ Correlation coefficients can range from -1.00 to +1.00.

- ▸ The following are three types of commonly used nonparametric correlation coefficients (Spearman R, Kendall Tau, and Gamma coefficients).

# Spearman R.

$$r_S = 1 - \frac{6 \sum\limits_{i=1}^{n} d_i^2}{n(n^2 - 1)},$$

$$d_i = \mathrm{R}\,x_i - \mathrm{R}\,y_i$$

*di* is a difference between ranks variables *x* and *y* for *i* observation

# Spearman R.

| IQ, $X_i$ | Hours of TV per week, $Y_i$ | Rank $x_i$ | Rank $y_i$ | $d_i$ | $d_i^2$ |
|---|---|---|---|---|---|
| 86 | 0 | 1 | 1 | 0 | 0 |
| 97 | 20 | 2 | 6 | -4 | 16 |
| 99 | 28 | 3 | 8 | -5 | 25 |
| 100 | 27 | 4 | 7 | -3 | 9 |
| 101 | 50 | 5 | 10 | -5 | 25 |
| 103 | 29 | 6 | 9 | -3 | 9 |
| 106 | 7 | 7 | 3 | 4 | 16 |
| 110 | 17 | 8 | 5 | 3 | 9 |
| 112 | 6 | 9 | 2 | 7 | 49 |
| 113 | 12 | 10 | 4 | 6 | 36 |

# Spearman R.

$$\rho = 1 - \frac{6 \times 194}{10(10^2 - 1)}$$

$\rho$ = -29/165 = −0.175757575... With p = 0.627188 (using t distribution)

# tau Kendal correlation

To calculate Kendal correlations, data should be compiled into all possible pairs and then divide these pairs into three possible categories:

- **Concordant pairs-** ordered in the same way (P),

- **Discordant pairs** – ordered differently (Q),

- **Bonded pairs** – the same values in pair for both pairs (T).

# tau Kendal correlation

$$\tau = \frac{P - Q}{P + Q + T}$$

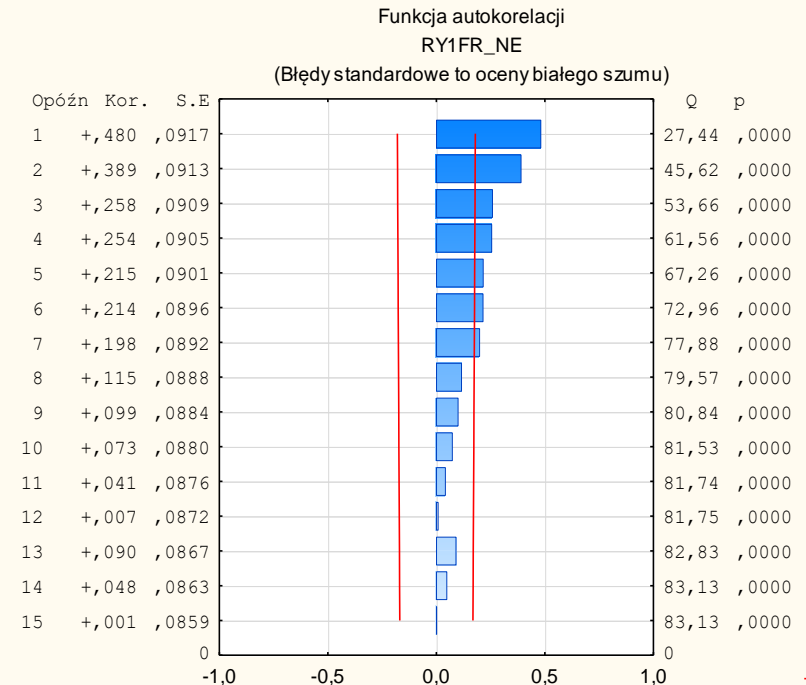$$P + Q + T = \binom{N}{2} = \frac{N(N-1)}{2}$$

$$\tau = 2\frac{P - Q}{N(N-1)}$$

# Correlations

| Korelacje (chronologie jesiony) Oznaczone wsp. korelacji są istotne z p < ,05000 N=116 (Braki danych usuwano przypadkami) | | | | |
|---|---|---|---|---|
| Zmienna | Średnia | Odch.std | ry1fr_ne | ry1fr_rcs |
| ry1fr_ne | 0,977198 | 0,180644 | 1,000000 | 0,884473 |
| ry1fr_rcs | 1,002135 | 0,171929 | 0,884473 | 1,000000 |

| Korelacja porządku rang Spearmana (chronologie jesiony) BD usuwane parami Oznaczone wsp. korelacji są istotne z p <,05000 | | |
|---|---|---|
| Zmienna | ry1fr_ne | ry1fr_rcs |
| ry1fr_ne | 1,000000 | 0,866205 |
| ry1fr_rcs | 0,866205 | 1,000000 |

| Korelacja tau Kendalla (chronologie jesiony) BD usuwane parami Oznaczone wsp. korelacji są istotne z p <,05000 | | |
|---|---|---|
| Zmienna | ry1fr_ne | ry1fr_rcs |
| ry1fr_ne | 1,000000 | 0,731634 |
| ry1fr_rcs | 0,731634 | 1,000000 |

Funkcja autokorelacji
RY1FR_NE
(Błędy standardowe to oceny białego szumu)

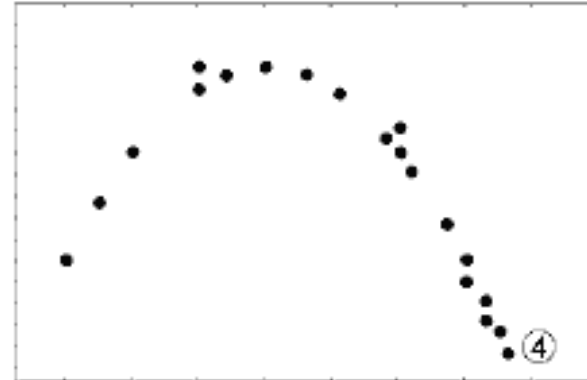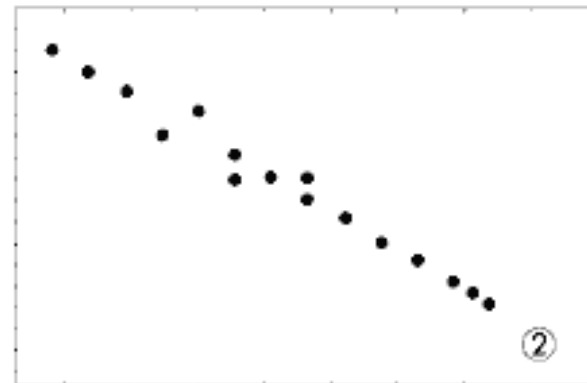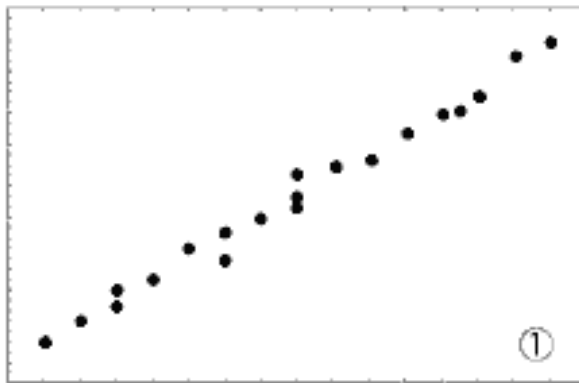| Opóźn | Kor. | S.E | | Q | p |
|---|---|---|---|---|---|
| 1 | +,480 | ,0917 | | 27,44 | ,0000 |
| 2 | +,389 | ,0913 | | 45,62 | ,0000 |
| 3 | +,258 | ,0909 | | 53,66 | ,0000 |
| 4 | +,254 | ,0905 | | 61,56 | ,0000 |
| 5 | +,215 | ,0901 | | 67,26 | ,0000 |
| 6 | +,214 | ,0896 | | 72,96 | ,0000 |
| 7 | +,198 | ,0892 | | 77,88 | ,0000 |
| 8 | +,115 | ,0888 | | 79,57 | ,0000 |
| 9 | +,099 | ,0884 | | 80,84 | ,0000 |
| 10 | +,073 | ,0880 | | 81,53 | ,0000 |
| 11 | +,041 | ,0876 | | 81,74 | ,0000 |
| 12 | +,007 | ,0872 | | 81,75 | ,0000 |
| 13 | +,090 | ,0867 | | 82,83 | ,0000 |
| 14 | +,048 | ,0863 | | 83,13 | ,0000 |
| 15 | +,001 | ,0859 | | 83,13 | ,0000 |

Measures similar to correlation:
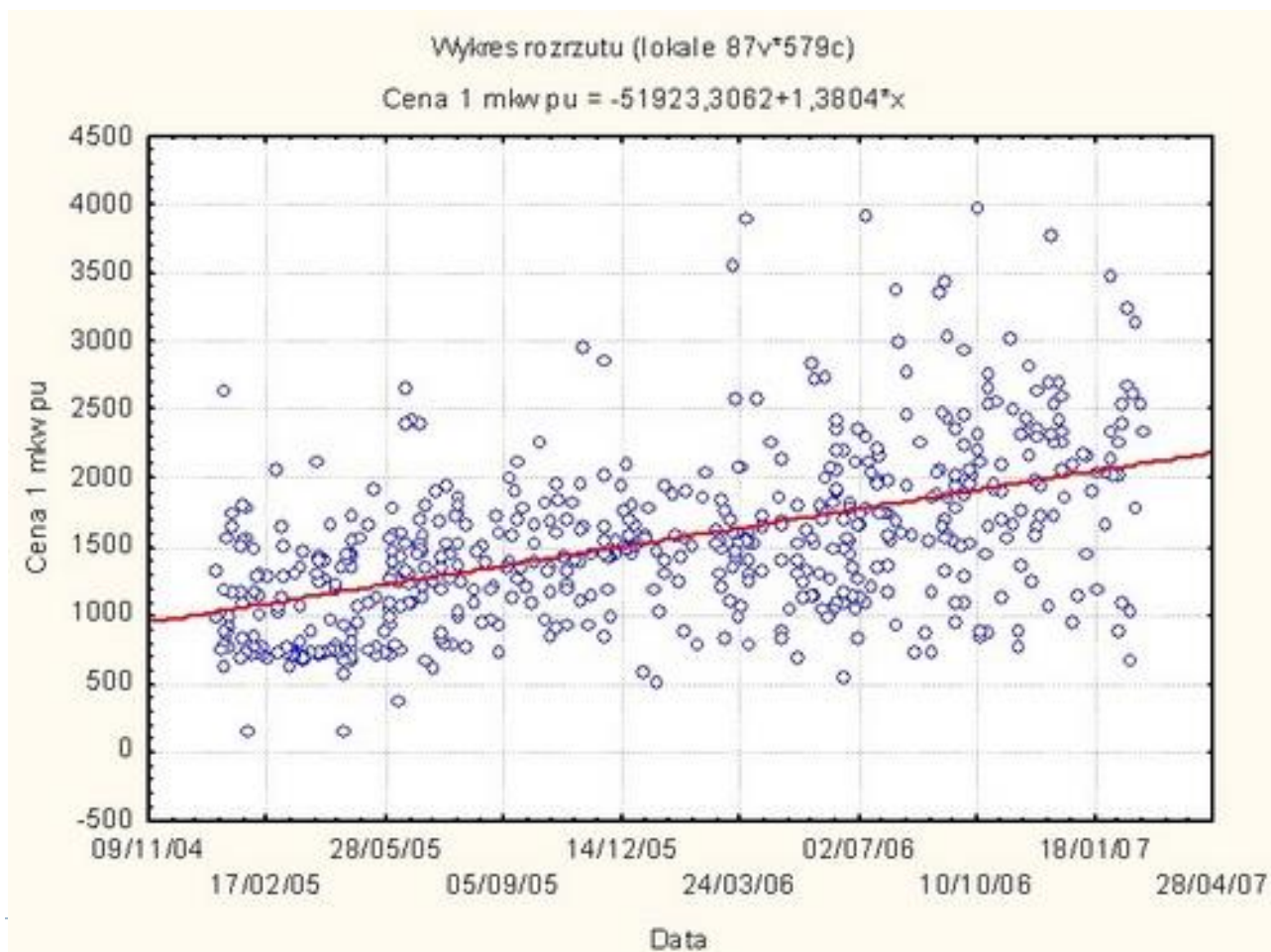- Fi coefficient,
- V Cramer's coefficient.

# Scatter plot- correlation

# Scatter plot- correlation

# More correlation coefficients

The Point-Biserial Correlation Coefficient is a correlation measure of the strength of association between a continuous-level variable (ratio or interval data) and a binary variable,
The phi coefficient (or mean square contingency coefficient) is a measure of association for two binary variables. It is known as the Matthews correlation coefficient (MCC),
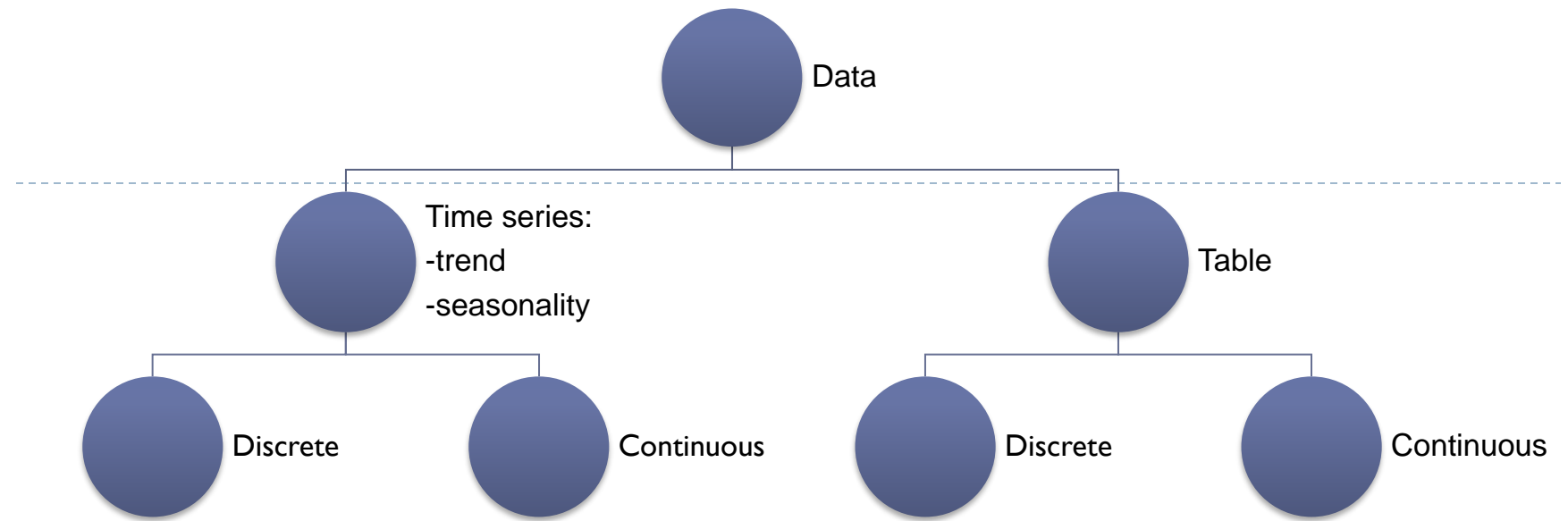Tetrachoric Correlation: Used to calculate the correlation between binary categorical variables.
Polychoric Correlation: Used to calculate the correlation between ordinal categorical variables.
 Cramer's V: Used to calculate the correlation between nominal categorical variables.

# Pytania?

**Do zastanowienia się:**
-Jaki jest wpływ wykonania preprocessingu na dalsze modelowanie
-Jaką różnicę w wynikach możemy uzyskać pracując z próbą zamiast populacji i czym jest ona spowodowana

**Na ćwiczenia wiadomości:**
-Preprocessing w zależności od typu danych i ich formy
-W jaki sposób analizujemy dane (na co zwracamy uwagę, jakie analizy wykonujemy)

# Lecture 3:
# Regression model