# Basic Descriptive Analysis

# Numerical Variables

### K. Gibert

*Department of Statistics and Operation Research*

*Knowledge Engineering and Machine Learning group*

*Universitat Politècnica de Catalunya, Barcelona*

*karina.gibert@upc.edu*

*www.eio.upc.edu/homepages/karina*

# Guió

# 0. Descriptive analysis

*Compact and Informative view of the variable structure*

**DATA= FIT+ ERROR**

**General Pattern**          **Deviations**          Characterizacion

Structural Component          Random Component

# Tools

1. Graphical

   Visualitze variable's distribution

   *the best inductor*

2. Numerical

   Quantify what is observed in he graphs

   *Objectivate*

# 1. Graphical tools

1. Performing the graph

       Mechanical                      (software)

2. Reading the graph

       Technical       (statistitian or data miner)

3. Interpretation

       Conceptual            (domain expert)

       Contextualization

# Graphical tools

# for numerical variables

1. Histogram

2. Boxplot

3. Others (dotplot, stem and leaf plot….)

# Histogram

Visualitzation of frequencies distribution table

| Intervalo | Número de Observaciones | Observaciones Acumuladas | Frec Relativas | Frec Acumuladas |
|---|---|---|---|---|
| 45–65 | 1 | 1 | $1/17 = .06$ | 0.06 |
| 65–75 | 5 | 6 | 0.29 | 0.35 |
| 75–85 | 5 | 11 | 0.29 | 0.64 |
| 85–95 | 1 | 12 | 0.06 | 0.70 |
| 95–105 | 3 | 15 | 0.17 | 0.87 |
| 105–115 | 0 | 15 | 0 | 0.87 |
| 115–125 | 1 | 16 | 0.06 | 0.93 |
| 125–135 | 0 | 16 | 0 | |
| 135–145 | 1 | 17 | 0.06 | |

*Frequency Classes?*

*Bars' AREAS PROPORTIONAL to frequencies*

Heuristics: $\begin{cases} 6 \, log_{10}(n) & , si \; n < 100 \\ 1{,}2\sqrt{n} & , si \; n \geq 100 \end{cases}$ $\qquad 3{,}49 \; s \; n^{-\frac{1}{3}} \qquad 2 \, d_i \, n^{-\frac{1}{3}}$

# Histogram

| Intervalo | Número de Observaciones |
|-----------|-------------------------|
| 45–65     | 1                       |
| 65-75     | 5                       |
| 75-85     | 5                       |
| 85-95     | 1                       |
| 95-105    | 3                       |
| 105–115   | 0                       |
| 115–125   | 1                       |
| 125–135   | 0                       |
| 135-145   | 1                       |

Bars' AREAS PROPORTIONAL to frequencies

*© K. Gibert*

# READING HISTOGRAMS

1. Range of variable (max-min)

2. Central trend

3. Dispersion

4. Simmetry

5. Anomalies
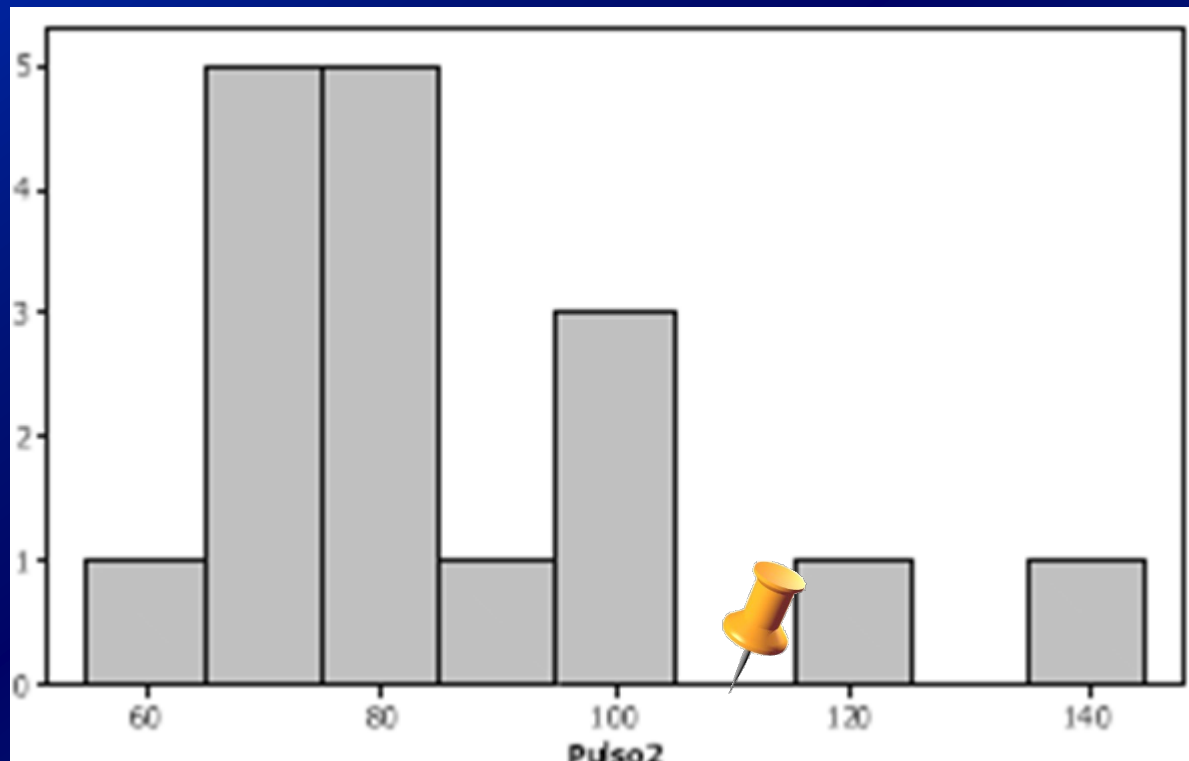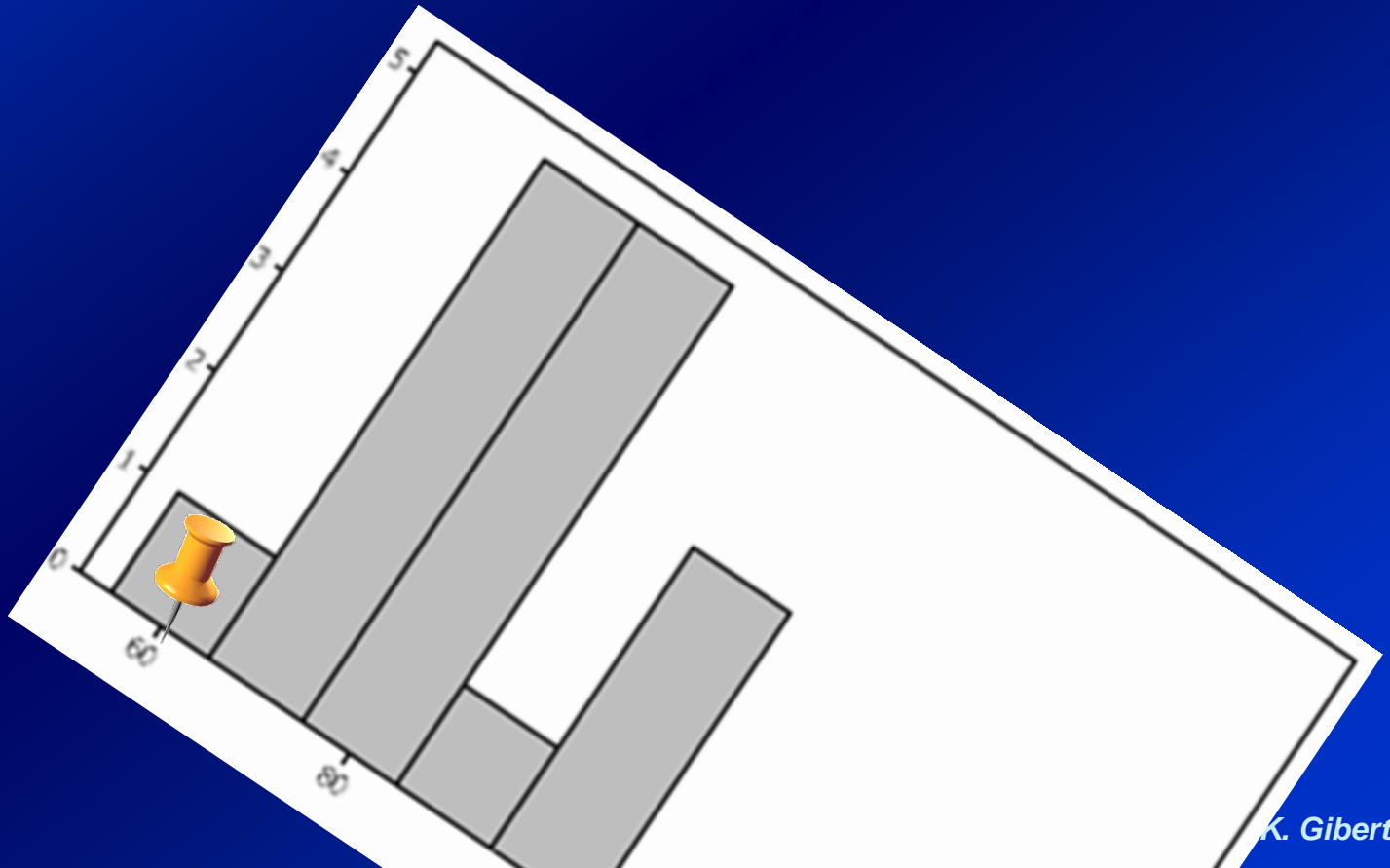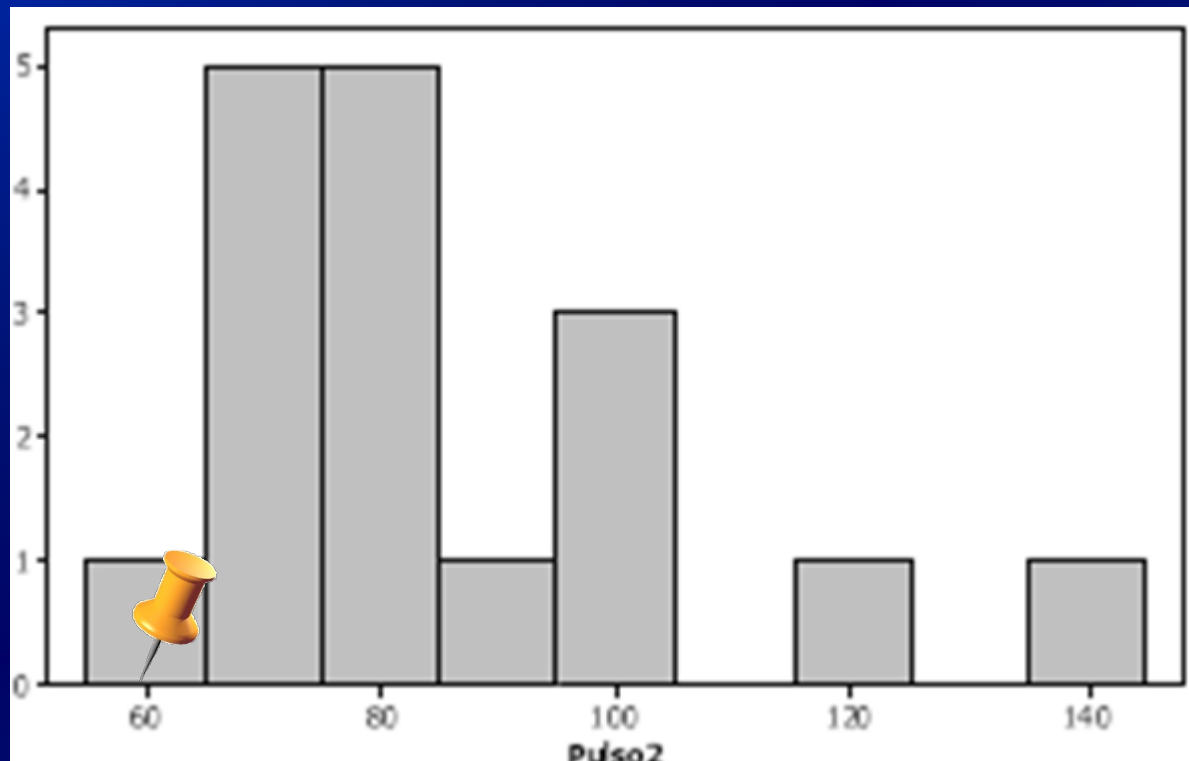
# READING HISTOGRAMS

## Central trend

*Value arround which observations distribute*

# READING HISTOGRAMS

## Central trend

*Value arround which observations distribute*

# READING HISTOGRAMS

## Central trend

*Value arround which observations distribute*

# READING HISTOGRAMS

## Central trend

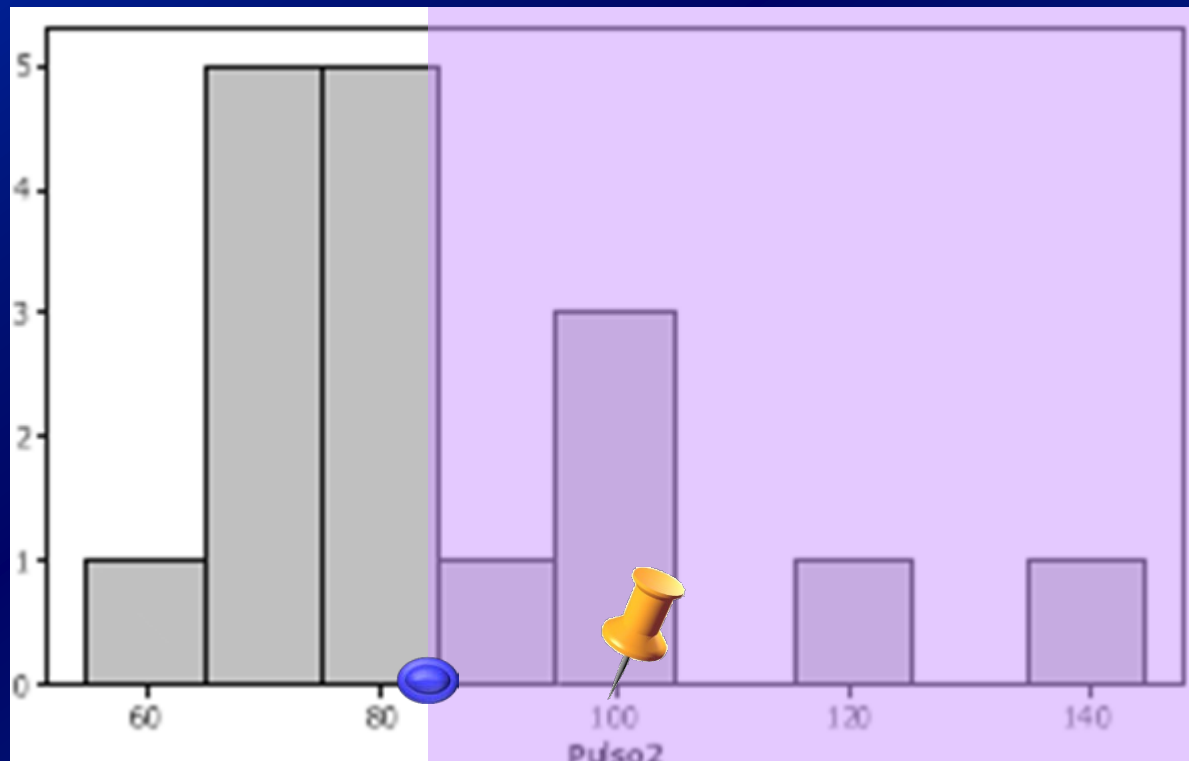*Value arround which observations distribute*



K. Gibert

# READING HISTOGRAMS

## Central trend

*Value arround which observations distribute*

# READING HISTOGRAMS

## Dispersion/Variability
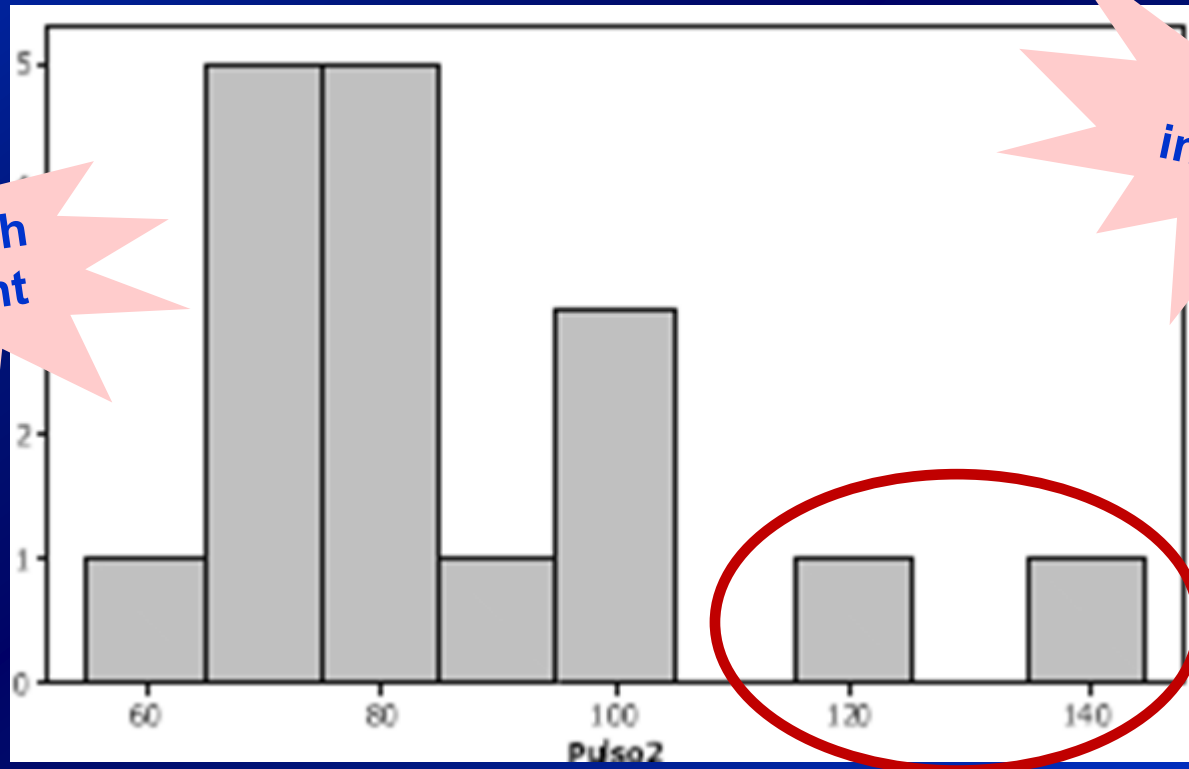
*How observations concentrate arround central trend?*

*Mean distance to central trend*

# READING HISTOGRAMS

## Anomalies
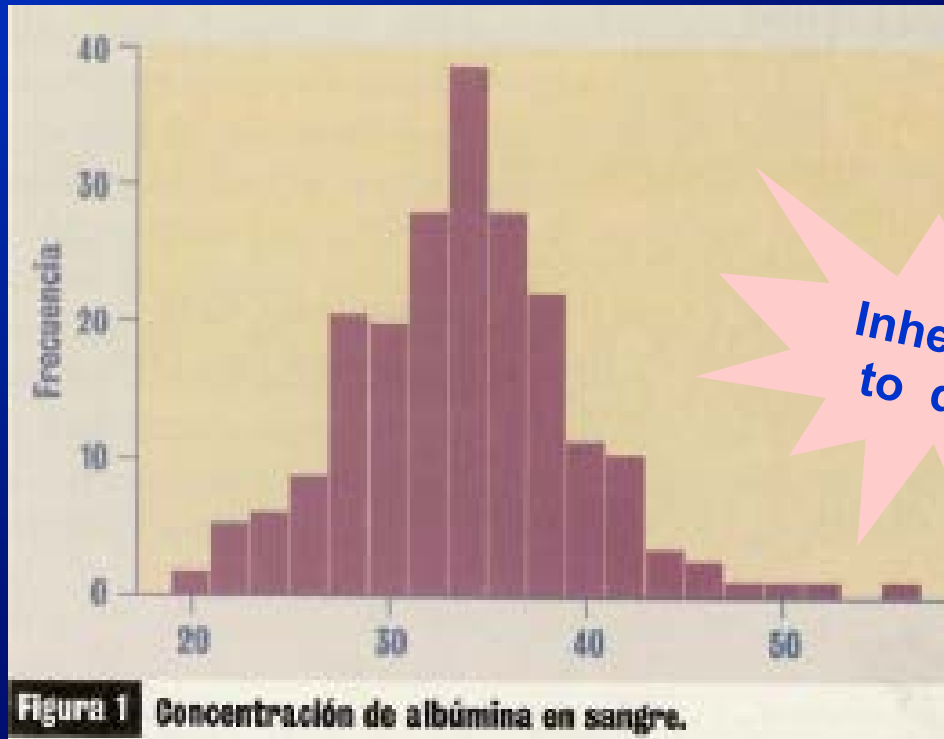
*Outliers: Observations anormaly far from rest*

# READING HISTOGRAMS

## Main Patterns    *[Gibert, JANO1996]*

### *Albumine*



Figura 1 Concentración de albúmina en sangre.

### *Bilirrubine*



F2 Nivel de bilirubina

**Inherent to data**

## *Symmetric*
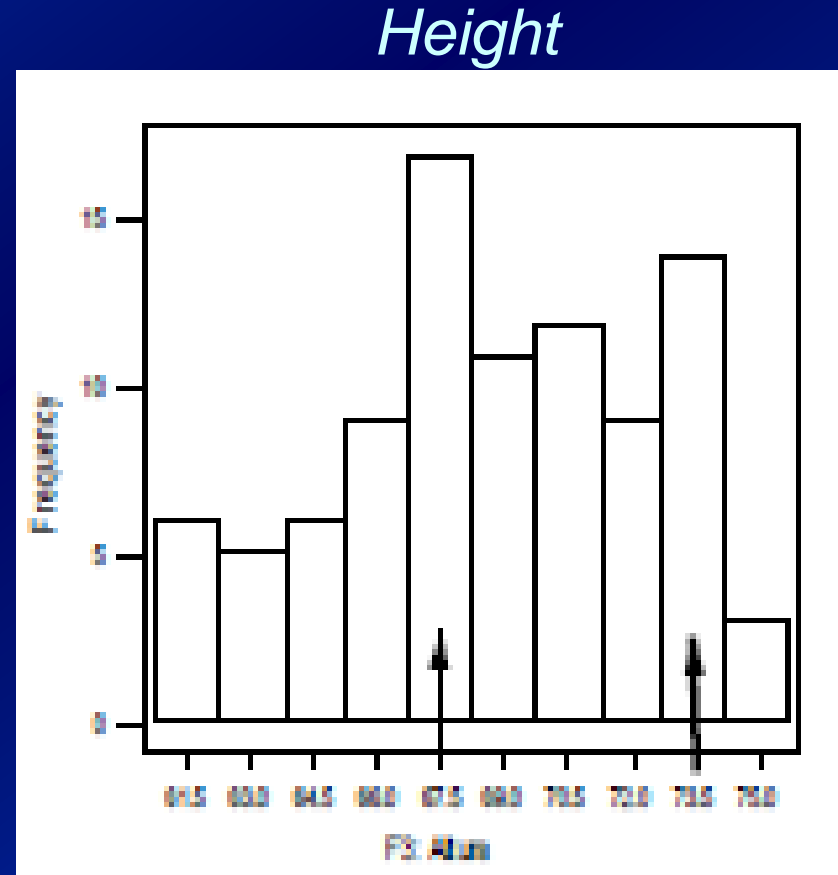
## *Asymmetric*

# READING HISTOGRAMS

## Main Patterns

*Height*

*Multimodality*

*Several central trends!!*

**Find discriminant factor**

# READING HISTOGRAMS
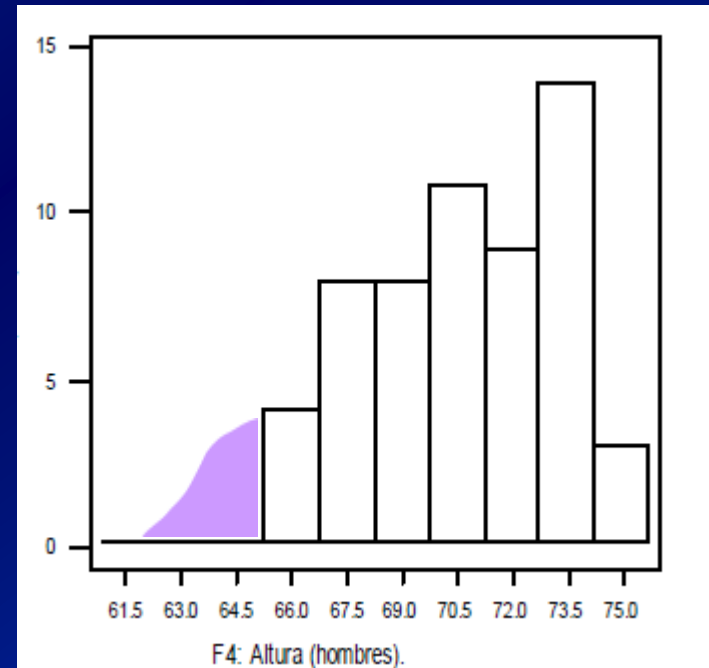
## Main Patterns

*Find discriminant factor*

*Height*

*Women*

*Men*

© K. Gibert

# READING HISTOGRAMS

## Main Patterns

*Height of Men*

*Scarped*

*Part of distribution trunked!*

*(only adult men)*



F4: Altura (hombres).

# READING HISTOGRAMS

## Main Patterns

*Pulse per minute*

*Dentat*

*Measurement approximations!*

*Count one minute*

*Count 10 sec x 6*

*Count 25 sec x 4*

*…*

# Tools

1. Graphical

   Visualitze variable's distribution

   *the best inductor*

2. Numerical

   Quantify what is observed in he graphs

   *Objectivate*

# 2. Numerical tools

*Quantify and synthetize characteristics of a distribution*

1. According to the information provided

    1. Central trend statistics

    2. Variability statistics

2. According to the stability

    1. Classic

    2. Robust

# Numerical tools

# for numerical variables

|  | Robusto | Clásico |
|---|---|---|
| **Posición** | Mediana<br>Cuartiles<br>Percentiles<br>Moda | Media |
| **Dispersión** | Distancia entre cuartiles | $S$<br>Desviación estándar<br>$S^2$<br>Varianza<br>Coef. variación<br>Amplitud |

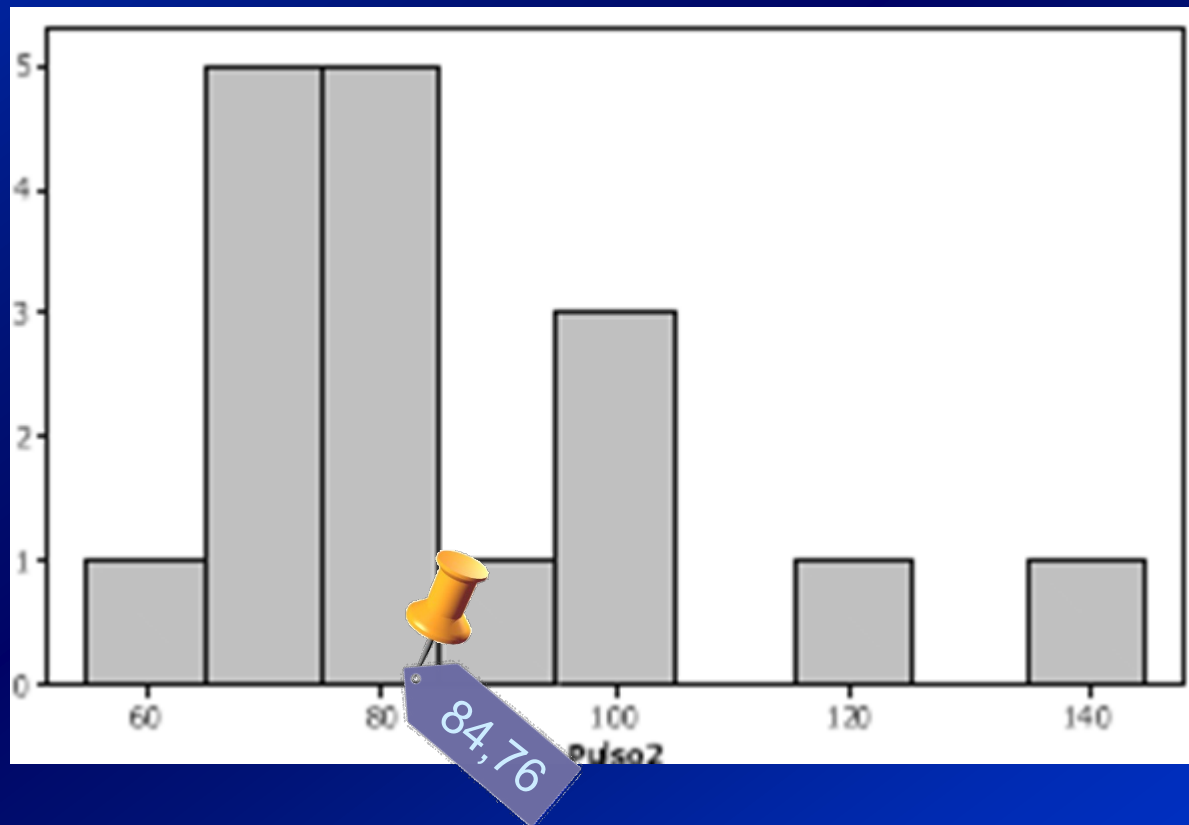# Numerical tools

## Mean

$x_1, x_2, \ldots x_n$ are $n$ observations of a variable X

$$\bar{x} = \frac{x_1 + x_2 + \ldots + x_n}{n} = \frac{\sum_{i=1}^{n} x_i}{n}$$

UPC

# READING HISTOGRAMS

## Central trend

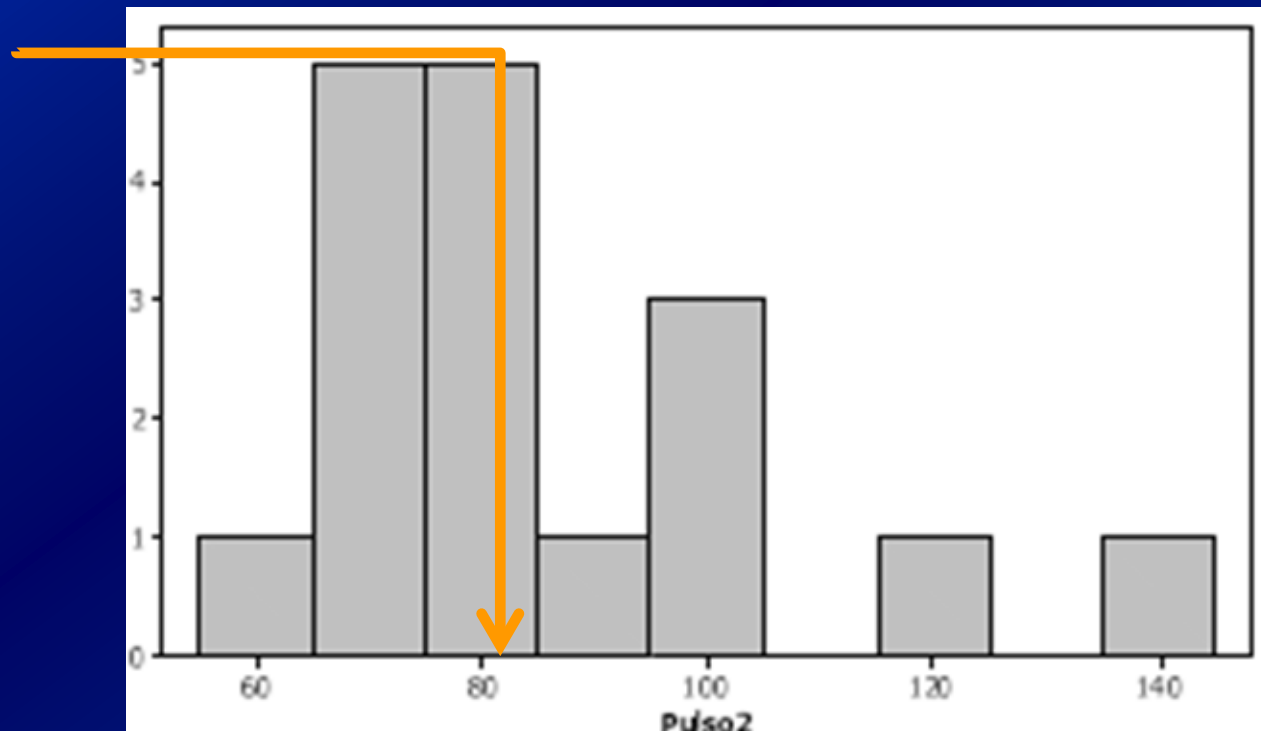*Value arround which observations distribute*

# Numerical tools

## Mean

$$\overline{x} = \frac{\displaystyle\sum_{i=1}^{n} x_i}{n}$$

© K. Gibert

# Numerical tools

# Mode

*The most frequent observation*

# Numerical tools

## Dispersion Measures

$$V(X) = \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n}$$

$$S_X = \sqrt{V(X)} = \sqrt{\frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n}}$$

$$qV(X) = \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}$$

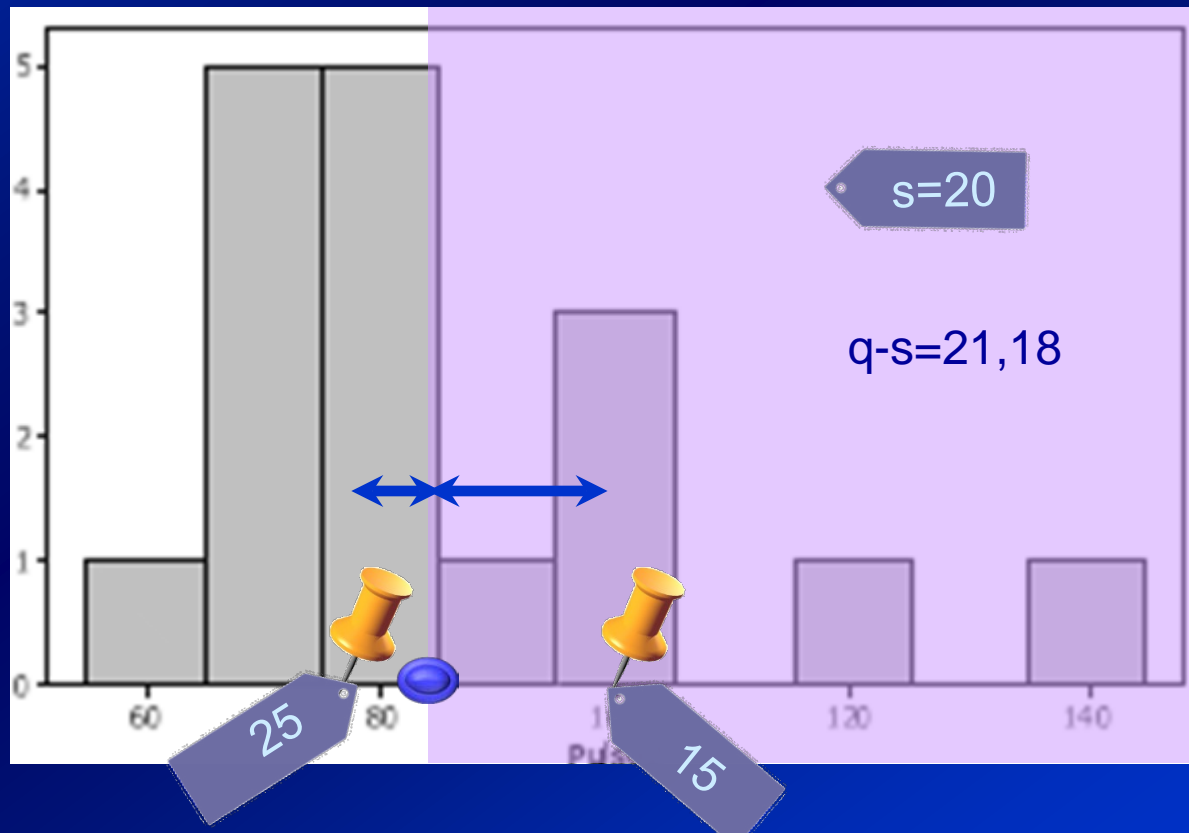$$q\text{-}S_X = \sqrt{qV(X)} = \sqrt{\frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}}$$

**Mean distance to central trend**

UPC

# READING HISTOGRAMS

## Dispersion/Variability

*How observations concentrate arround central trend?*
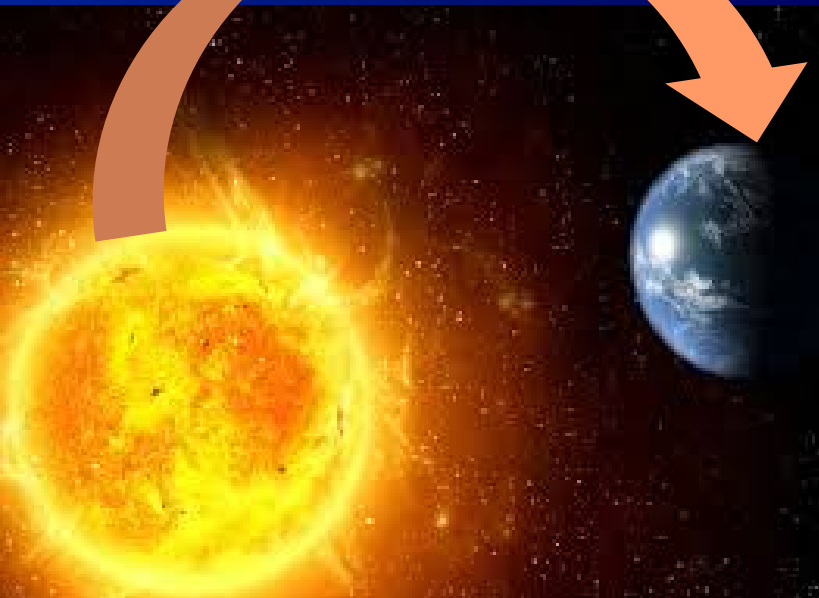
*Mean distance to central trend*

# Numerical tools

## Dispersion Measures

Variation Coeficient: $\dfrac{s}{\bar{x}} \times 100$

$d_i/range$

d(Sun, Earth): $150\times10^6$ Km, $s=10^3$ Km

d(BCN, Moscow): 3000 Km, $s=10^3$ Km

*© K. Gibert*

# Numerical tools

## 5-Number Summary

Robust

*<min, Q1, Me, Q3, max>*

Requires
sorting

# Boxplot   *[Tukey 1956]*

*Symbolic representation of 5-Number Summary*

**Boxplot of N(0,1)**

Outliers

extreme      smooth

(3di)              (1.5di)

**Outliers**

**BOX**

**Whiskers**

**(1.5di)**

**Whiskers**

**(1.5di)**

-4      -3      -2      -1      0      1      2      3

**N(0,1)**

**Min**           **Q1**      **Me**      **Q3**                **Max**   *© K. Gibert*
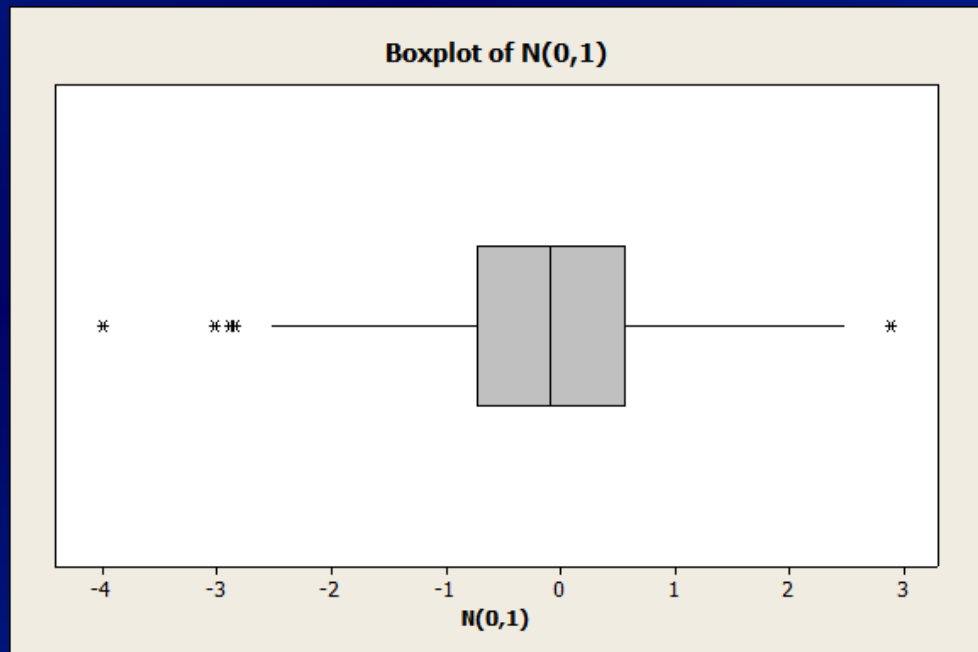
UPC

# READING BOXPLOTS

1. Range of variable (max-min)

2. Central trend

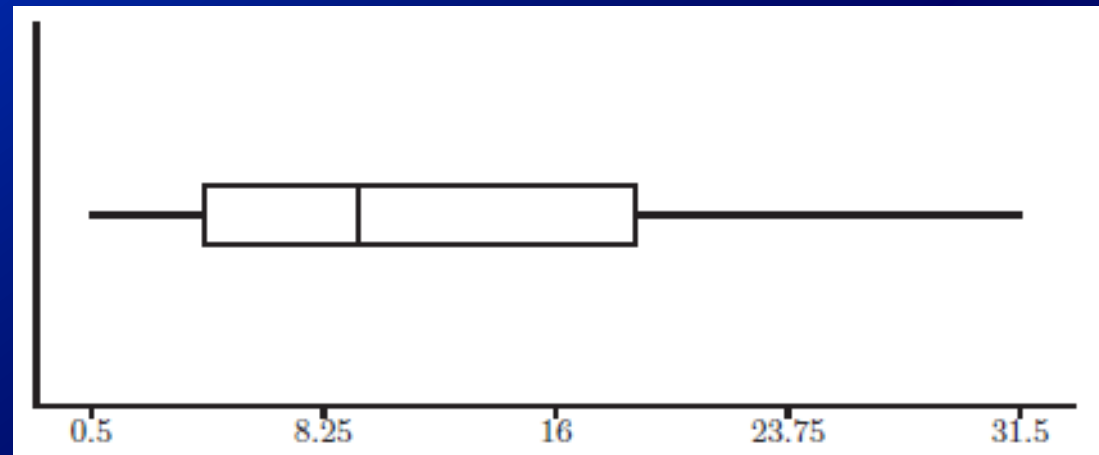3. Dispersion

4. Simmetry

5. Anomalies

# READING BOXPLOTS

1. Range of variable (max-min)
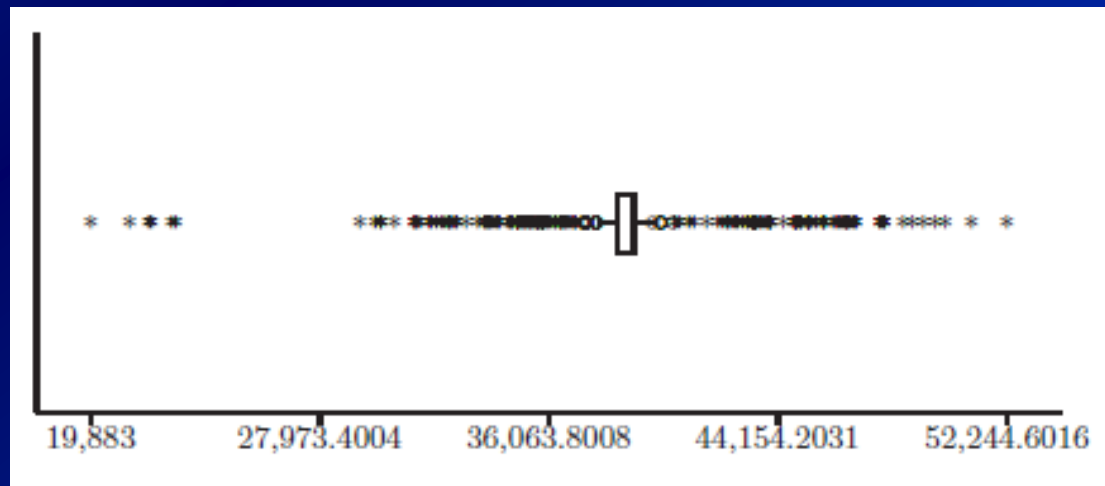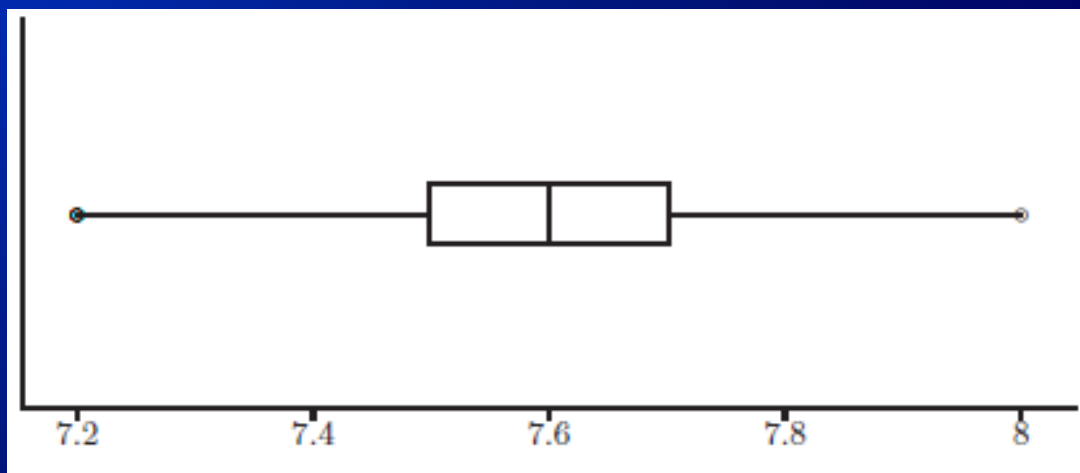
2. Central trend

3. Dispersion

4. Simmetry

5. Anomalies



Boxplot of N(0,1)
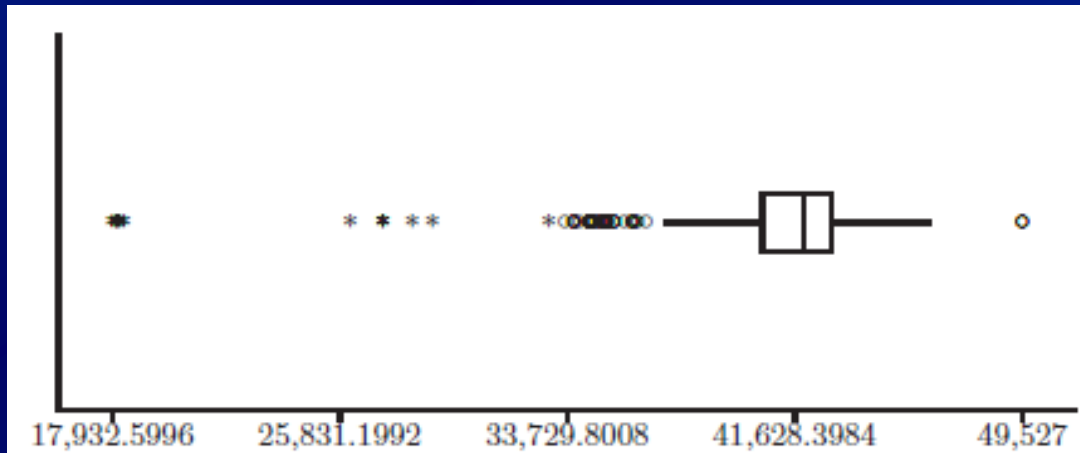
# READING BOXPLOTS
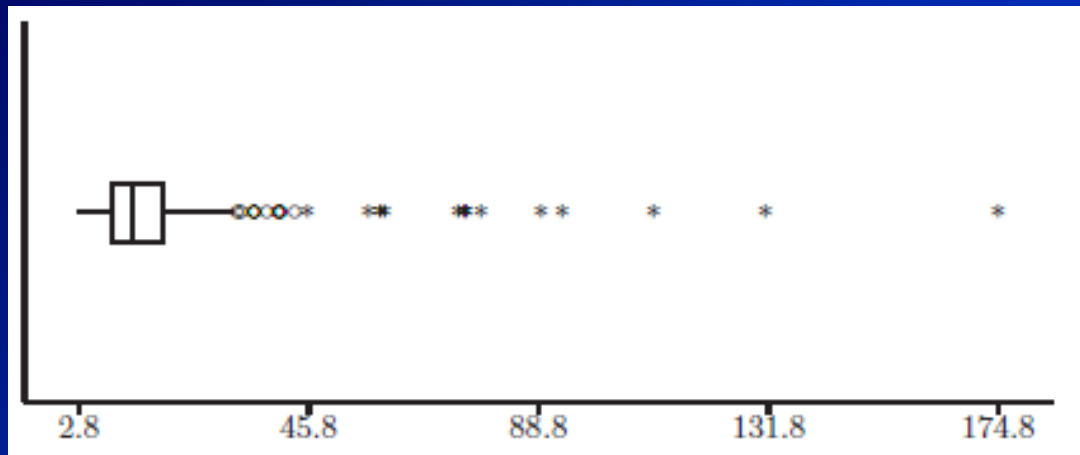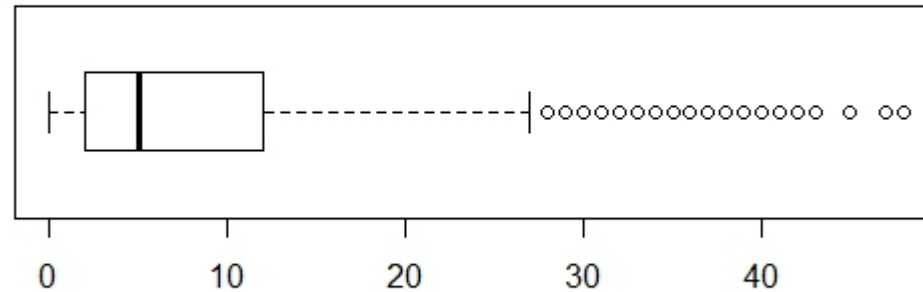
## Dispersion



*Ammonium*
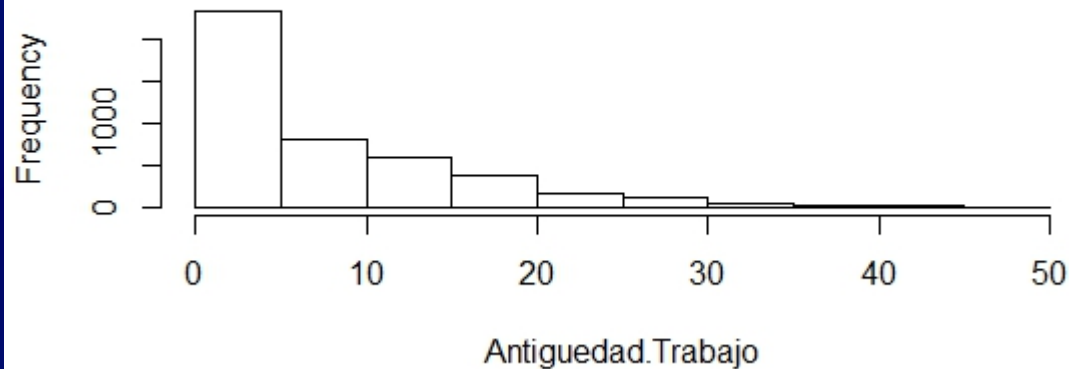
*QB-B*

**READING BOXPLOTS**

Simmetry

# READING BOXPLOTS



Boxplot of Antiguedad.Trabajo

Histogram of Antiguedad.Trabajo

# Symmetry

if Mean <> Median then

      if Me-Q1 < >Q3-Me then  assymmetry

                    else  outliers

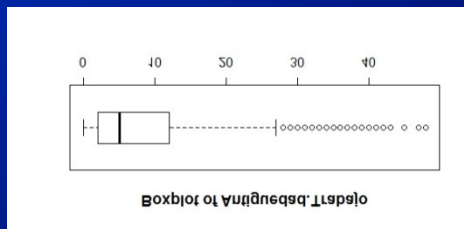else  symmetry without outliers

# Synthesis

1. Descriptive analysis of numerical variable

    1. Central trend and variability (classical/robust)

    2. Graphical and Numerical tools



Boxplot of Antiguedad.Trabajo



*5-Number-Summary:*

*<min, Q1, Me, Q3, max>*

+

*mean, q-stdev, variation coefficient*
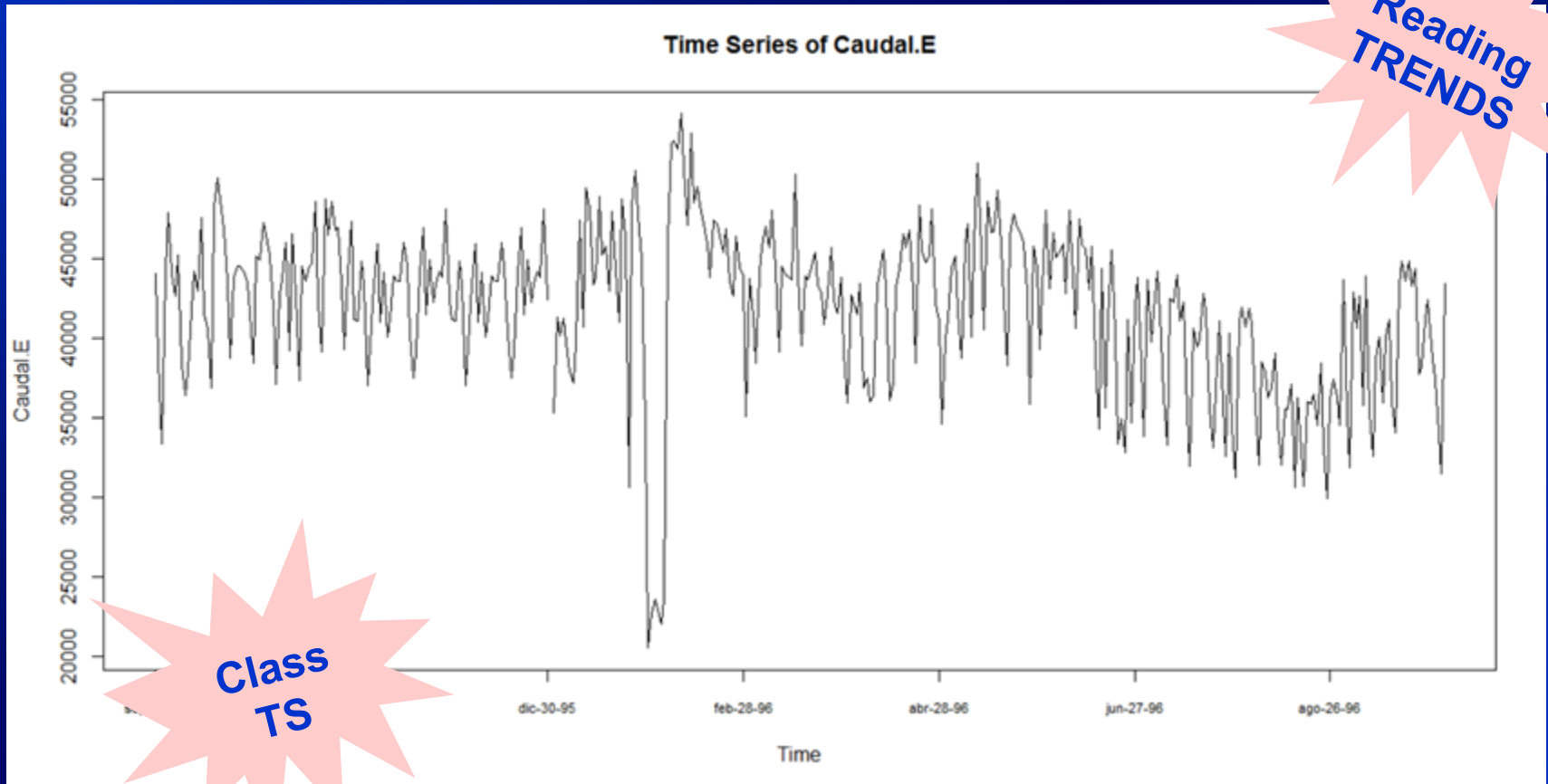
3. Characterize the variable

    *Central trend, variability, symmetry, n-modality….*

© K. Gibert

# 3. Cronological data

*Observations are sequentially sorted in the dataset*

## Time series plot

Reading TRENDS

Class TS



**Time Series of Caudal.E**

© K. Gibert

# Cronological data

*Observations are sequentially sorted in the dataset*

*Evenctually the Date is Available*

Class Date

# Cronological data

## *Date Objects in R*

| Symbo | Meaning | Example |
|-------|---------|---------|
| %d | day as a number (0-31) | 01-31 |
| %D | Date format | |
| %a | abbreviated weekday | Mon |
| %A | unabbreviated weekday | Monday |
| %m | month (00-12) | 00-12 |
| %b | abbreviated month | Jan |
| %B | unabbreviated month | January |
| %y | 2-digit year | 07 |
| %Y | 4-digit year | 2007 |

*31/12/2014  : %d/%m/%Y*

*31-Dic-07: %d-%b-%y*

UPC

# Cronological data

## *Date Objects in R*

| Symbo | Meaning | Example |
|-------|---------|---------|
| %c | Date and time | |
| %C | Century | |
| %H | Hours (00-23) | 15 |
| %I | Hours (1-12 ) | 3 |
| %j | Day of the year (0-365) | 250 |
| %M | minute (00-59) | January |
| %S | Second as integer(0-61) | 07 |

*23:12:59  = %H:%M:%S*

*11 12 59  = %I %M  %S*

# Cronological data

*Observations are sequentially sorted in the dataset*

*Evenctually the Date is Available*

**Class Date**

*To consider time*

**Class POSIXCT**

UPC

# 4. Assessing Normality

1. 68-97-99.5 Rule

   [x +-s], [x+- 2s] [x+-3s]

2. Normality plot (qq-plot, Henri line)

3. Normality assessment test: Shapiro Wilk

$$W = \frac{\left(\sum_{i=1}^{N} a_i y_i\right)^2}{\sum_{i=1}^{N} (y_i - m_1)^2}$$

yi= ith order statistic  m1=sample mean

ai= computed as linear regression to the expected value of standard normal order statistics

# 5. Assessing Exponentiality

## The rule of 70

Time of doubling : 70/R, R growing factor *[Moore, McCabe 93]*

*X has exponential growth with constant factor R*

*if needs 70/R time to pass from X to 2X*

# Basic Descriptive Analysis
# Numerical Variables

*Karina Gibert*

*Dpt. Statistics and Operation Research*

*Knowledge Engineering and Machine Learning Research group*

*Universitat Politècnica de Catalunya-BarcelonaTech (Spain)*
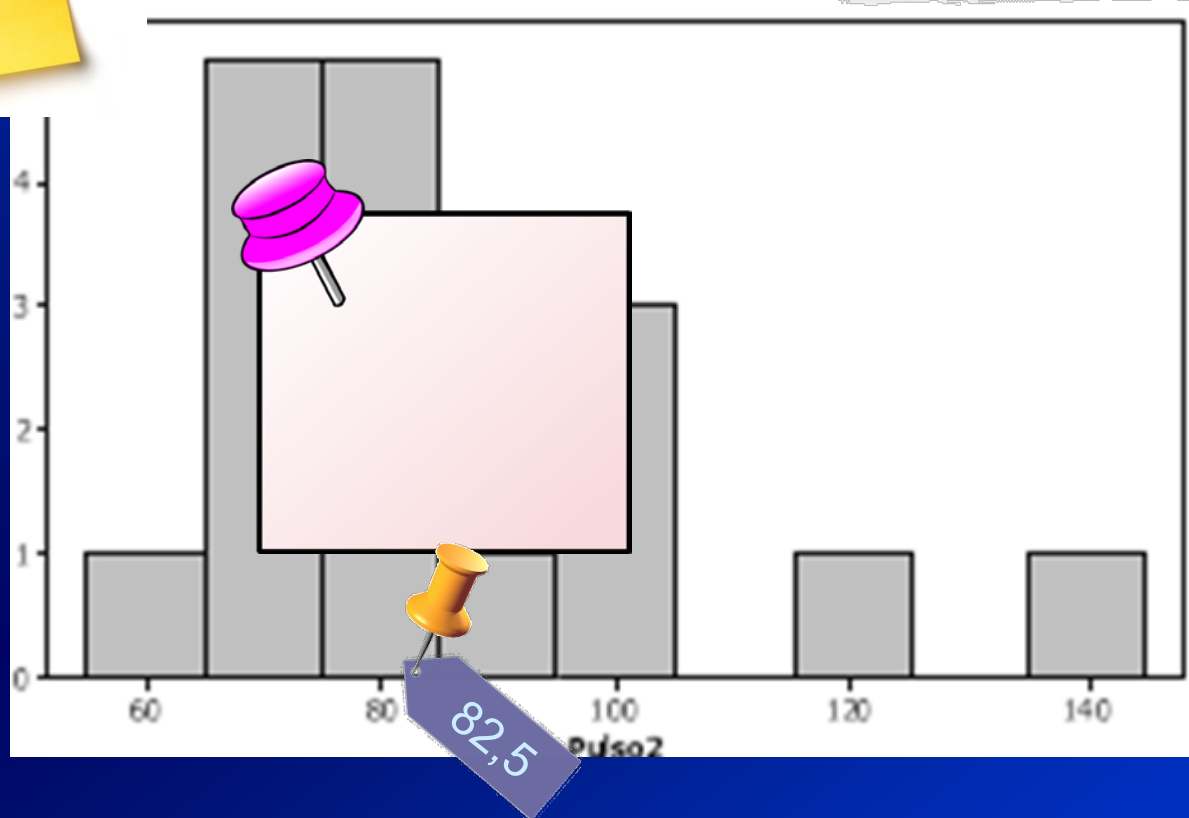
*karina.gibert@upc.edu*

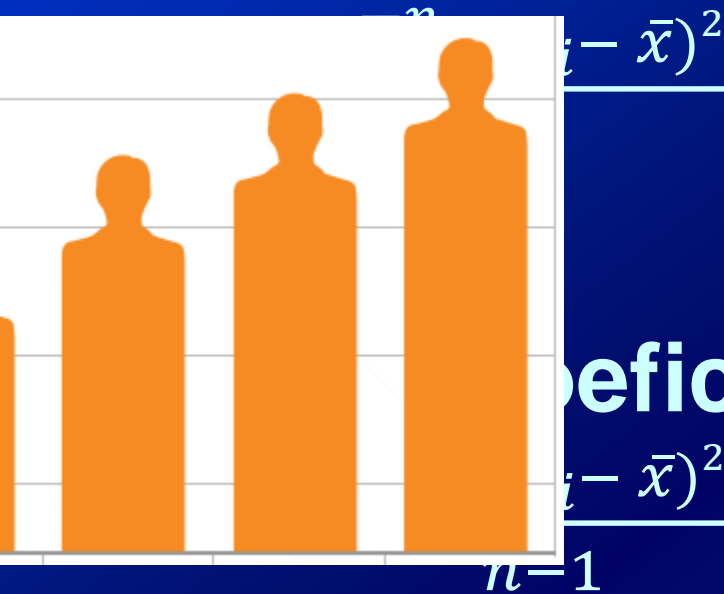*www.eio.upc.edu/homepages/karina*

*Are there any questions?...*

# Numerical tools

## Dispersion

$$\frac{\sum_{i}^{n}(x_i - \bar{x})^2}{}$$

**Coefficient**

$$\frac{(x_i - \bar{x})^2}{n-1}$$

Mean dist central

© *K. Gibert*