

# Intelligent data Data Processing



Maria-Amparo Vila  
vila@decsai.ugr.es

Investigaci'on group Databases and Intelligent  
Systems Informaci'on <https://idbis.ugr.es/>

Department of Sciences of the computation and  
artificial intelligence  
University of Granada

## Introduction to the topic

*Presentaci' structure on*

---

one. Introduction b'asicas ideas about data

two. Type of data

3. Quality problems

Four. Exploraci'on data

4.1 exploration Estad'istica

4.2 Visualizaci'on data

5. Data Transformations

6. Reduccion problems of variables.

6.1 Selecci'on of variables

6.2 Change of coordinates: main components

7. Scaling problems.



## Introduction- data

### Input data

---

M's data structure common to work with DM is the

### Dataset

items variables	$V_{one}$	$V_{two}$	$\dots\dots\dots$	$V_N$
$or_{one}$	$d_{eleven}$	$d_{12}$	$\dots\dots\dots$	$d_{one\ N}$
$\vdots$	$\vdots$	$\vdots$	$\dots\dots\dots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\dots\dots\dots$	$\vdots$
$or_M$	$d_{M\ one}$	$d_{M\ two}$	$\dots\dots\dots$	$d_{M\ N}$



## Introduction- data

### Input data

M's data structure common to work with DM is the

### Dataset

items variables	$V_{one}$	$V_{two}$	$\dots\dots\dots$	$V_N$
$or_{one}$	$d_{eleven}$	$d_{12}$	$\dots\dots\dots$	$d_{one\ N}$
$\vdots$	$\vdots$	$\vdots$	$\dots\dots\dots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\dots\dots\dots$	$\vdots$
$or_M$	$d_{M\ one}$	$d_{M\ two}$	$\dots\dots\dots$	$d_{MN}$

- items represent cases, objects etc.
- Variables can be of many types. Also they called factors
- There may be missing data



## Introduction- data

### *Input data*

---

*The data set may be obtained from previous data, through transformations, resampling, etc. In some cases This is a key point (selection of factors, text mining etc.)*



## Introduction- data

### *Input data*

---

*The data set may be obtained from previous data, through transformations, resampling, etc. In some cases This is a key point (selection of factors, text mining etc.)*

There are problems that the structure of data set is not suitable:

- transactional structures
- Miner's graph (structural patterns are sought)
- Miner's sequence (Biocomputation)



## Introduction- data

### *Input data*

---

*The data set may be obtained from previous data, through transformations, res' umenes etc. In some cases This is a key point (selecci'on of factors, text mining etc.)*

There are problems that the structure of data set is not suitable:

- transactional structures
- Miner'ia graph (structural patterns are sought)
- Miner'ia sequence (Biocomputaci'on)

*In mayor'ia cases they can be transformed into another one representations in order to apply the appropriate t'ecnica*

*From now on, except indicaci'on against, we will focus on the structure of data set*



## Type of data

*num'ericos attributes*

---

Num'ericos has a domain which allows arithmetical operations.

You can classify in:

- **Discretos :**

They are integers or natural numbers. Usually results count. Num'ericos computations allow, seg'

SU domain.

Not to be confused with categorical attributes changed.

**Game 1 level 2 or 3 It is not one attribute num'ericos**





# Type of data

## *num'ericos attributes*

---

- **continued** :

correspond to real numbers.

Advanced computations allowed statisticians to use them. Sometimes they are too detailed and may have rounding problems.

always they have a starting point (zero) and a scale factor. Segments that can be classified in:

- **"Interval"** Zero and scale are arbitrary. (Time in milliseconds and arbitrary starting point, the temperature in Celsius and Fahrenheit etc.)
- **"Ratio"** Zero is not chosen but if the scale factor. (Distances, height, weight, volume etc.) have sense proportion
- **"Absolute"** Both the zero and the scale factor is determined. (Any form of percentage, frequency etc.)

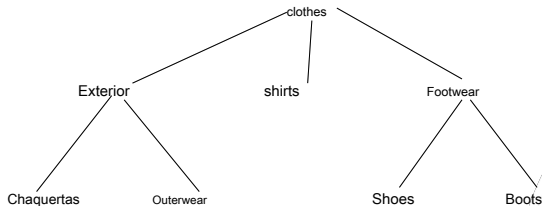
## Type of data

*simb' attributes      olicos, categ'oricos or nominal*

---

A discrete domain have no num'ericos values. Do not allow aritm'eticas operations. In principle equality can only support a jer'arquica structure. with different **granularity levels**

*Example entities simb'olicas*



## Type of data

*simb'olicos*      *oligos, categ'oricos or nominal*

---

Some domains simb'olicos est'an ordered attributes (acad'emicos courses pron'ostico a disease etc.). They are called

**ordinal attributes** , and support comparaci'on operators. Binary attributes (presence / absence) are a form of attribute ordinal

M'as clear example of ordinal attribute is the date. Supporting Granularity

*Example*

***DATE - → MONTH - → TRIMESTER - → TO~***

DO NOT



# Data Quality

*Accuracy (correctness) and precision*

---

*Accuracy: Similarity between the data value and the true value of the attribute.*

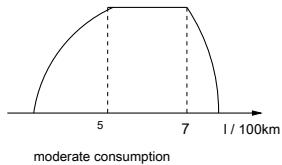
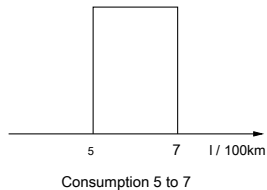
- In the case of numeric attributes:
  - They can exist rounding errors and must unify precision.
  - intervalares or diffuse: Inaccurate valuations may exist. Have to treat them with proper tools: similarities, fuzzy clustering, fuzzy rules association etc.
- In the case of symbolic attributes
  - errors are detected in the data
    - There syntactic detection
    - There semantic detection



# Data Quality

Accuracy (correctness) and precision

## Example inaccurate values



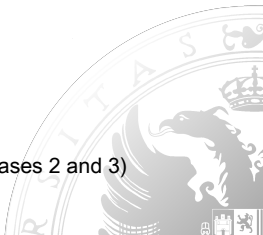
# Data Quality

## *Completeness (Full details)*

---

Completeness Security there is sufficient data and not missing any value of an attribute

- Missing Values in the attributes:
  - . In the case of numericos there t'ecnicas estad'ísticas attributes we will see at the end
  - . In the case of simb'olicos attributes may have prior knowledge that can be used. (Functional dependencies etc.)
- Missing items:
  - one. Lost records or tuples
  - two. Informaci'on biased
  3. scattered data
    - . In this case it is dif fi cult to ensure the quality of the results.
    - . A lot of data does not ensure quality domain t'erminos items (cases 2 and 3)



## Data Quality

*Other problems with data quality*

---

**Anomalies (Outliers)** These are items that really do not belong to the group to be studied and distorting regularity is sought

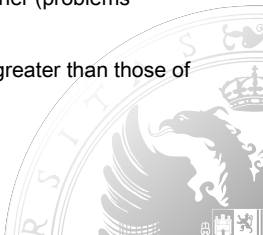
**time lag (timeliness)** Refers to the fact that the data or of them do not have the same "now" than others.

**unbalanced data** It can occur at two levels:

- There are many items of a kind and very few other (problems classifi cation / predicci'on)
- The values of an attribute num'érico are much greater than those of others which has much weight M'as

**Duplicates data** They appear when databases are merged.

Habr'a to clean the data previously



# Exploraci' Data on

*Why and for what*

---

## Motivations to explore data

- Helps you choose the best tools to preprocess and analyze
- Lets make initial hypothesis on patterns extracted as the ability of human being is exploited to recognize patterns

## Exploratory Data analysis (EDA) 1977

- It is because Tuckey
- The EDA est'a focused on visualization. Assumes adequate techniques can be extracted direct knowledge.
- To Tuckey the Clustering and detection of anomalies part of EDA

- More information

<http://www.itl.nist.gov/div898/handbook/eda/eda.htm>





## Exploraci' Data on

*Exploraci' based on descriptive Estadística*

**Frequency Distribuci'on** They can be considered:

- absolute frequencies** :  $F(d_i) n'$       umber of times featured a determined value  $d_i$

Computation of the absolute frequencies

- discrete attributes: Categ'oricos and whole and fi nite num'ericos (with not many domain data). Simple counting on domain values  $D$
- continuous attributes. the discretizaci'on is imposed, the domain becomes discrete intervals through. Usually. Yes  $D = [A, B]$  and we want  $m$  intervals are chosen equal amplitude

$$[to_{one}, to_{two}, \dots, to_{m-one}, to_m], to_1 = A, to_i = to_{i-1} + (B - A) / m \quad \forall i = two, \dots, m$$

The problem of discretization can be complex for visualizaci'on and asociaci'on issues. In some cases it is a matter of preprocessing.

## Exploraci' Data on

*Exploraci' based on descriptive Estad'istica*

Frequency Distribuci'on They can be considered:

- **relative frequencies** :  $f(d_{ij})$  raz'on between absolute frequency and  $n'$   
 UMBER total of items  $f(d_{ij}) = F(d_{ij}) / M$ ,  $M$  It is the  $n'$  UMBER total  
 items. Tambi'en may be given in percentages (  $f(x_{ij}) * 100$ )
- **cumulative frequencies** I s'olo are defined for sorted data. Suppose the set of  
 values  $D = \{d_{one}, \dots, d_m\}$  Y  
 $d_{one} \leq d_{two} \leq \dots \leq d_n$ , it is defined:

$$F(d_{ij}) = \sum_{j=one}^{i=i} f(d_{ij})$$

When percentages are used  $F(d_{ij})$  It indicates the percentage of the poblaci'on  $\leq d_i$  and  
 leads to the concept of:

## Exploraci' Data on

*Exploraci' based on descriptive Estad'istica*

---

### Frequency Distribuci'on

- **percentile** Be  $s$  a value between 0 and 100. We define:

$$p_s = \max \{d \in D / f(d) \leq s\}$$

It is the highest value of the domain having below the  $s$  percent of the poblaci'on.

When  $s = 0, 25, 50, 75, 100$  They are called

quartiles

### Centralizaci'on measures

- For data num'ericos

**Half** :  $d = \frac{\sum_{d \in D} d \cdot M}{M}$

- For all types of data:

**fashion** : The common value  $M$ 'as

**Median** : The 50th percentile.



# Exploraci' Data on

*Exploraci' based on descriptive Estad'istica*

## Measures of dispersion

- For data num'ericos

**variance**  $s^2 = \frac{\sum_{d \in D} (d - \bar{d})^2}{M - 1}$ ,  $s$  is the **Typical deviation**

**Average absolute desviaci'on**  $AAD =$

$$\frac{\sum_{d \in D} |d - \bar{d}|}{M}$$

**Median absolute desviaci'on**

**$MAD = \text{median} \{|d - \bar{d}|; d \in D\}$**

- For all data:

**interquartile range**  $r = p_{75} - p_{25}$



## Exploraci' Data on

*Exploraci' based on descriptive Estad'istica*

Exploraci'on of relaci'on between num'ericos attributes

**Covariance matrix** Be the variables  $V_j, V_k$

$$\text{cov}(V_j, V_k) = \frac{\sum_{i=1}^M (d_{ij} - \bar{d}_j)(d_{ik} - \bar{d}_k)}{M - 1}$$

**Correlation of matrix** Be the variables  $V_j, V_k$

$$\text{corr}(V_j, V_k) = \frac{\text{cov}(V_j, V_k)}{s_j s_k}$$

- When  $\text{corr}(V_j, V_k) \approx 1$  or  $\text{corr}(V_j, V_k) \approx -1$  a linear relaci'on between attributes and one of them can be expressed in funci'on other. This can serve to reduce variables.

# Exploring Data on

Visualization on

---

**Basic ideas** You can be displayed:

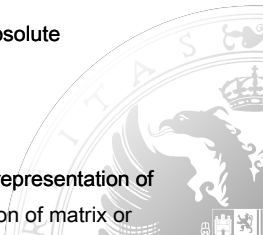
**Objects** an object is represented in a coordinate system, as a value one, two or three attributes. , Point clouds appear different colors etc.

**Attributes** Seg' one are:

- **categorical attributes** : Bar charts, pie charts. Using colors in the representations of objects etc.
- **numeric attributes and ordinal** : Histogram of absolute frequencies, relative cumulative, boxplot (box plot)

**Relations** joint representations of attributes cloud

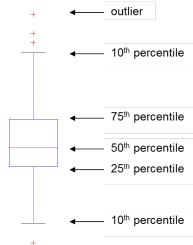
points (scatter plot) dimensional histograms. Joint representation of Cajes diagrams etc. REPRESENTATION Correlation of matrix or covariance function categorical attributes etc.



# Exploraci' Data on

## *The boxplot (Tukey)*

It's another way to view data distribuci'on



For **num'ericos and continuous attributes** tambi'en is usually done using, medium rather than median and  $\pm$  one s, two s, 3 s .. instead of percentiles. Outliers are considered from  $\pm 3$  s

## Exploraci' Data on

*Example*

---

### IRIS dataset



You can be obtained from:

<http://www.ics.uci.edu/ mLearn / MLRepository.html>

- Four numeric attributes: sepal length, sepal width, petal length, petal width
- An attribute categorical: *setosa*, *virginica* and *versicolour*
- 150 items, 50 of each type



## Exploraci' Data on

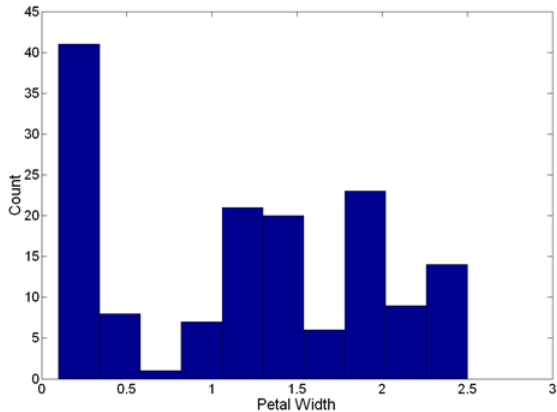
*Example: statistical measures*

Numeric columns		Nominal columns		
Row ID	D sepal l...	D sepal ...	D petal l...	D peta
Minimum	4.3	2	1	0.1
Maximum	7.9	4.4	6.9	2.5
Mean	5.843	3.057	3.758	1.199
Std. deviation	0.828	0.436	1.765	0.762
Variance	0.686	0.19	3.116	0.581
Overall sum	876.5	458.6	563.7	179.9
No. missings	0	0	0	0

## Exploraci' Data on

*Example: histogram with 10 intervals*

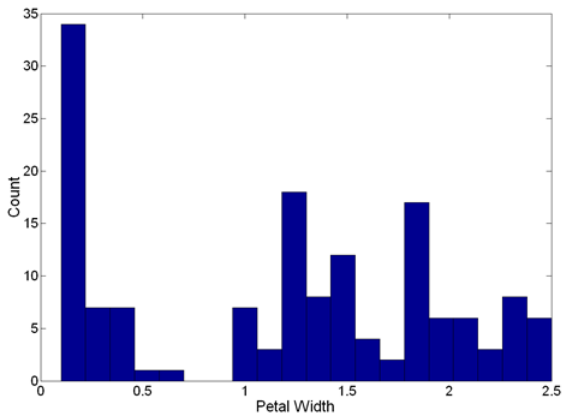
---



## Exploraci' Data on

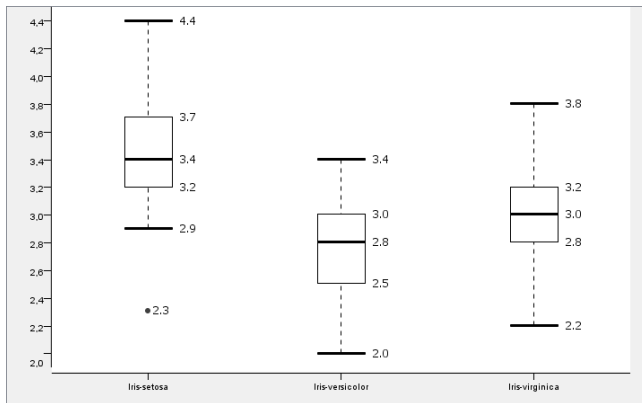
*Example: histograms with 20 intervals*

---



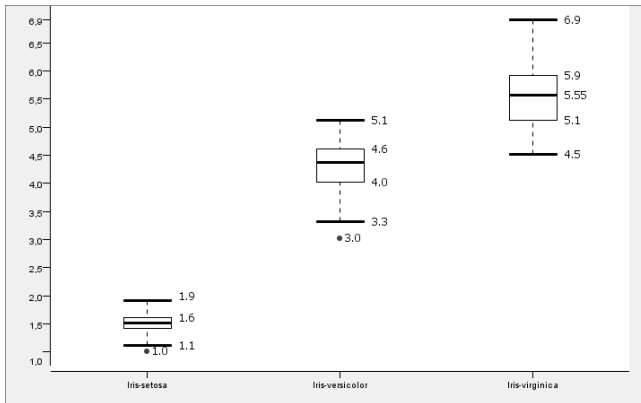
## Exploraci' Data on

*Example: boxplot (petal length)*



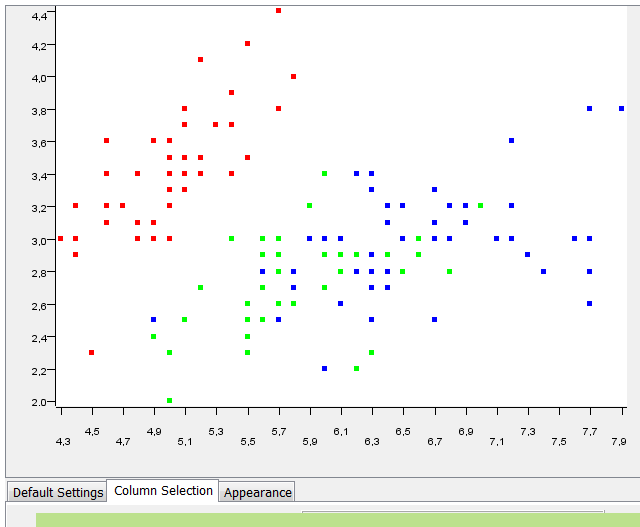
## Exploraci' Data on

*Example: boxplot (petal width)*



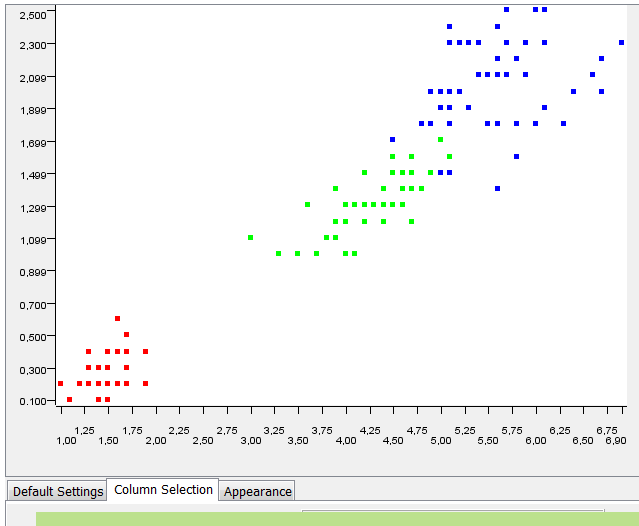
## Exploraci' Data on

*Example: gr'a fi c points (sepal length / width sepal)*



# Exploraci' Data on

*Example: gr'a fi c points (petal length / width petal)*



# Data transformations

*b'asicas ideas*

---

*Before applying t'ecnicas DM in mayor'ia cases it is necessary to preprocess (transform) data*

By qu'e transform data?

- Data needs to be changed because it can be treated directly. (Change of date of birth to age)
- The data are too detailed: **add, sum, discretizar, transforming**
- There are too many variables: **T'ecnicas of reducci'on'o factors selecci'on**
- There is much difference between the ranges. **T'ecnicas scaling**
- There is much lost data **T'ecnicas lost data processing**





# Data transformations

*Agregaci' on, short, discretizaci' on*

---

## Aggregation

Adding data is to combine various objects to get a new one. In general you need to add when the information is too detailed. We have:

- **Horizontal AGGREGATION (abstract)** The data are very detailed level object and must summarize the attributes. Cl'asicas are situations OLAP village-level data are aggregated into zones, etc.
- **Vertical AGGREGATION** An attribute data are very detailed and you need to add to a higher level. The classic example is the time. Other AGGREGATION the semantics of t'erminos seg' ontology

a A

# Data transformations

*Agregaci' on, short, discretizaci' on*

---

## discretization

Discretized data is to substitute a continuous num'érico categ'orico one attribute. It is necessary in extracci'on of asociaci'on rules and classifying certain processes caci'on. Process:

**one.** K intervals are chosen

**two.** Each value is associated with the midpoint of the range where est'a and said midpoint located renames. There are several approaches:

- **Intervals equally distributed** It's the standard
- **Equal frequency intervals** . Try to have equal n' where each interval.

umber of

## Data transformations

*Agregaci' on, short, discretizaci' on*

---

### discretization

- **Intervals obtained by grouping** . one partitional clustering (K-means) considering only the attribute to discretize applies. Centroid give us the values to be replaced.

*There m'etodos discretization LING*

*based u'istica*

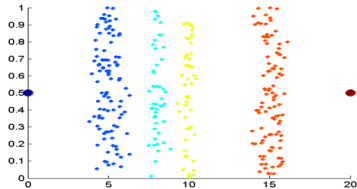
*Fuzzy clustering.*



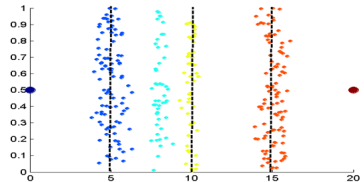
# Data transformations

Example discretizaci' on

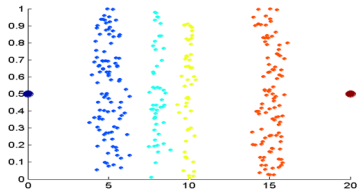
Data



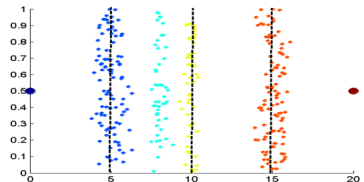
Equal intervals



equal frequency



Grouping



## Reducci' on variables

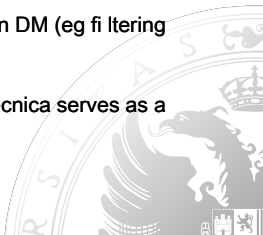
### *Selecci' on variables*

---

It is dimensi'on reduce the problem by taking a smaller set of variables

### T'ecnicas

- **Brute force** Test Test / Error
- **Selecci'on included** The DM will t'ecnica choosing M'as meaningful variables (eg spanning trees of decisi'on)
- **Filtered out** variables are selected prior to the t'ecnicas aplicaci'on DM (eg fi ltering t'erminos in Text Mining)
- **Selecci'on hedging** The goodness of the result of applying DM t'ecnica serves as a criterion for selecci'on.



## Reducci' on variables

### *Principal component analysis*

---

#### Problem to be addressed

##### • *Illustrative example*

Consider a set of students for which the qualifications obtained in five subjects are given, the first two ex'amenenes have been made without notes and the other three with notes. You want to instruct students in funci'on their performance.

#### **Solution :**

Finding a "normalized linear combination" of scores:

$$x_c = \sum_{j=1}^5 x_j \text{ such that } \sum_{j=1}^5 x_j^2 = 1$$

m'axima to collect the variance because "separar'a" to sort items and facil ser'am'as



# Reducci' on variables

*Principal component analysis*

---

Problem to be addressed

*Other examples*

- Ordenaci'on banking customers
- Ordenaci'on of seg' items      A Website
- Overall obtenci'on of res'      umenes attribute



## Reducci' on variables

### *Principal component analysis*

---

#### The model matem'atico

Consider a data set with real num'ericos data, this can be seen as a variable  $N$  dimensional

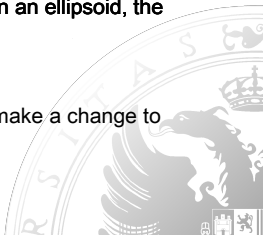
$$X = (x_{one}, \dots, x_N) \quad Y \text{ } m \text{ values}$$

the same  $x_{ij}$ ,  $i \in \{one, \dots, M\}$ ,  $j \in \{one, \dots, N\}$ , we find standard linear transformations (SLC) that "summarize" the best possible data, capturing most variance thereof.

#### *intuitive idea*

If the items are considered as a point cloud  $R_n$ , all can be enclosed in an ellipsoid, the average center, whose matrix is the covariance matrix.

The axes of the ellipsoid are a rectangular coordinate system, if we make a change to this coordinate system, the point spread along axis

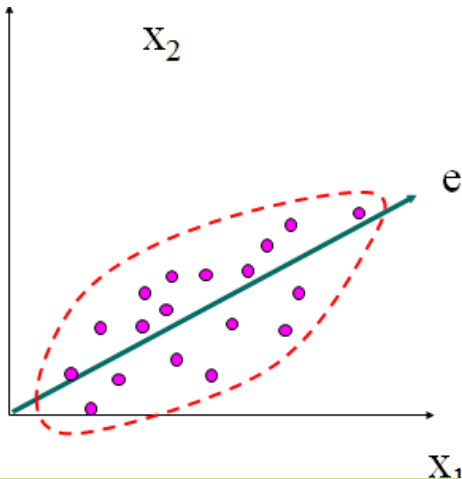




## Reducci' on variables

*Principal component analysis*

*intuitive idea*



## Reducci' on variables

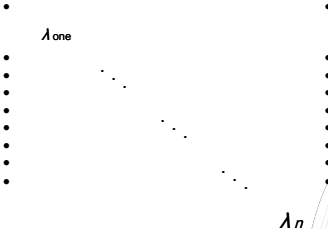
### Principal component analysis

#### The model matem'atico

Be  $\mu$  the average  $\bar{x}$   $Y$   $\Sigma$  its covariance matrix, is  
 find a linear transformaci'on  $y = \Gamma (\bar{x} - \mu)$  such that new  
 coordinate axes are the axes of the ellipsoid. Proves that  $\Gamma$  is a matrix such that:

$$\Gamma \cdot \Sigma \Gamma = \Lambda$$

where,

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) =$$


The diagram shows a diagonal matrix  $\Lambda$  with eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$  on the diagonal. The matrix is represented by a grid of dots, with the diagonal elements labeled  $\lambda_1, \lambda_2, \dots, \lambda_n$ .

## Reducci' on variables

### Principal component analysis

#### The model matem'atico

Vector values  $\bar{x}$   $\lambda$  verified  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$  and are the  
 "Eigenvalue" of the covariance matrix and associated with each  $\lambda_j$  there is a  
 "eigenvector"  $\gamma_j$  that is the  $j$ th column  
 matrix  $\Gamma$  and verify c'andose:

$$\forall j \in \{1, \dots, N\} \text{ and } j = 1 \quad \gamma_j'(\bar{x} - \bar{\mu})$$

$Y_j$  It is called  $j$ th principal component.

$$\forall k, j \in \{1, \dots, N\} \text{ Cov}(and_j, Y_k) = 0 \quad Var(and_j) = \lambda_j$$

$$Var(and_1) \geq \dots \geq Var(and_N)$$

and given  $k \leq N$  there is no SLC which is independent of  $k$  first main component and  
 having a greater variance than the  $k + 1$  main component.

## Reducci' on variables

### *Principal component analysis*

---

#### The model matem'atico

- *Proporci'on explained variance*

?The proporci'on of variance explained by  $k$  factors is

$(\lambda_1 + \dots + \lambda_k) / (\lambda_1 + \dots + \lambda_n)$  and it allows us to reduce the dimensionality of space. That is express fen'omeno less variables.

#### *How many components make ?:*

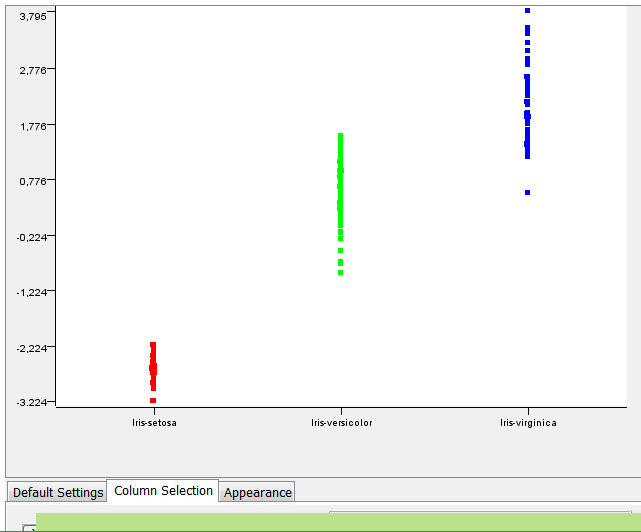
- At least 90% of variance explained
- All eigenvalues that are greater than the average of the same. If the matrix is used instead Correlation of the covariance matrix eigenvalues greater than 1.

*The proporci'on of variaci'on explained variable  $j$  by  $k$  component allows identifying the components and identify one semantics for them.*



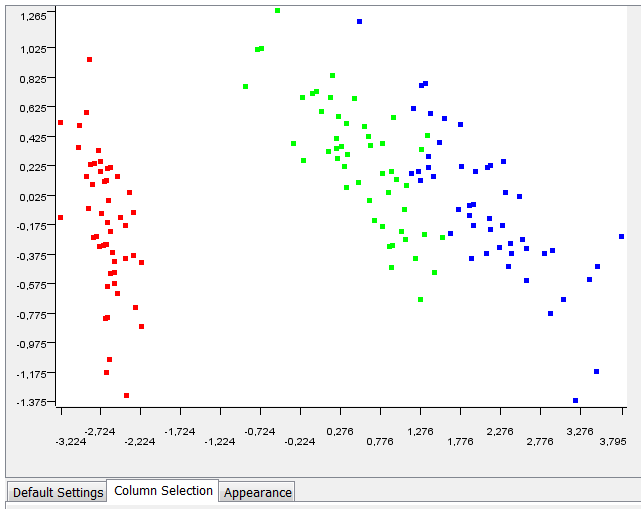
# Main components

*Example: Iris (classes / first factor)*



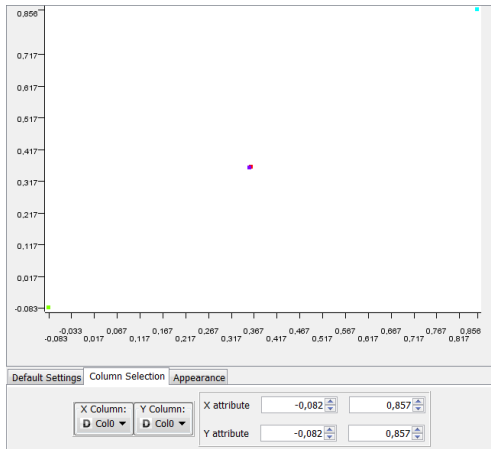
# Main components

*Example: Iris (first factor / second factor)*



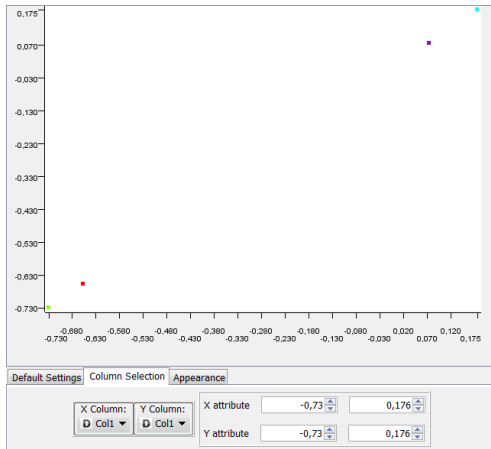
# Main components

*Example: Iris (variables / first factor)*



# Main components

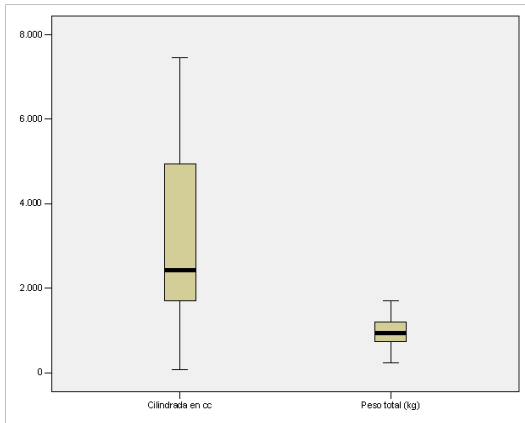
*Example: Iris (variables / second factor)*





## Scaling problems

*Motivación*



## Scaling problems

*b'asicas ideas*

*Num'ericos the values of two attributes are different scales. This makes can not be treated together on issues such as computation of distances etc.*

- *Some expressions to normalize*
- Normalizaci'on in [0,1]

$$V = \frac{V_{to} - \min_{to}}{\max_{to} - \min_{to}}$$

- Tipi fi cac'i'on

$$V = \frac{V_{to} - d}{s}$$

- Tipi Robusta fi cac'i'on

$$V = \frac{V_{to} - \text{median}}{\text{interquartile range}}$$



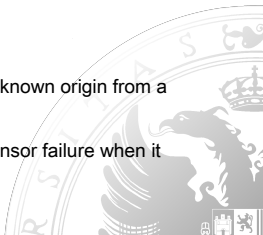
# Problems of missing values

## *b'asicas ideas*

---

*A variable has missing values when its value is not known for a particular data*

- *Origin of missing values*
- Lack of value without random factors. (Someone does not answer something in a survey)
- Property not applicable. (Hair Color frogs). You can not identify with 0 or NO.
- Error random source.
  - random completely: it follows that the data distribuci'on (failure of unknown origin from a sensor) (MCAR)
  - conditional random: follows a conditional distribuci'on, failure of a sensor failure when it rains M'as (MAR)



# Problems of missing values

*To do with missing values*

---

## delete records Eliminating lost data

- If it is a totally random situation.
- If the data volume is not seriously altered

## Replace Replace the lost data with

- a new value is generated within a quantum domain. (NS / NC, NO etc.)
- Missing as follows using the common value segment class registration.
- If quantitative data is replaced by:
  - The average if the error is completely random
  - The conditioned mean to appearance error if we have a type MAR
  - Proximal average values (interpolating) if we know that there is a certain temporal or spatial dependence

how a