

**NOVA**  
**IMS**  
Information Management School

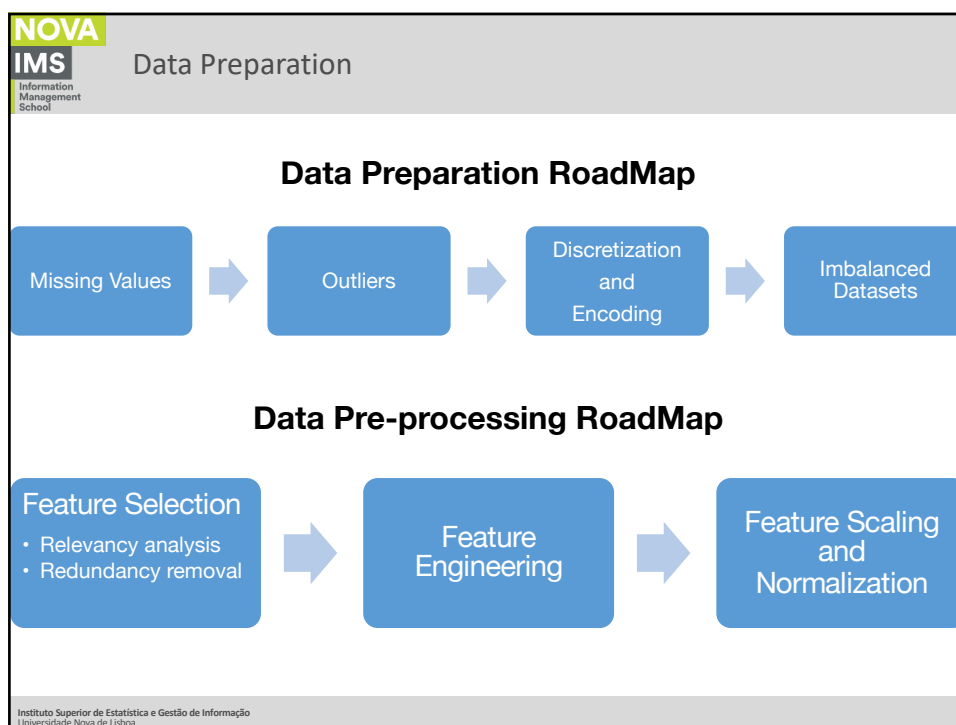
# Data Mining

## S6

Fernando Lucas Bação  
[bação@isegi.unl.pt](mailto:bação@isegi.unl.pt)  
<http://www.isegi.unl.pt/fbação>

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

1



2

**NOVA**  
**IMS**  
Information Management School

# Data Preprocessing

Preparing data for mining models

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

UNICIS

ASES

Schools

eduniversal

3

**NOVA**  
**IMS**  
Information Management School

## Data Preparation

Through preprocessing, we transform original data into a more compact and informative representation.

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

4

**NOVA**  
**IMS**  
Information  
Management  
School

## Data Preprocessing

- **Reasons:**
  - Noise reduction and signal amplification;
  - Size Reduction of the Input Space;
    - Domain-specific knowledge application;
    - Feature selection
      - Remove correlated variables
      - Remove irrelevant variables
    - Feature engineering
      - Constructing ratios and derived variables
  - Scaling (normalization);

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa









5

**NOVA**  
**IMS**  
Information  
Management  
School

## Reducing Input Space

**Why fewer variables often means better models**

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

6

**NOVA**  
**IMS**  
Information  
Management  
School

## Data Preprocessing

- **Additional considerations about data:**
  - Curse of dimensionality – the input space grows exponentially with the number of input variables;
  - The larger the input space, the more data and computing power we need.
  - In high dimensions, distance metrics lose meaning.

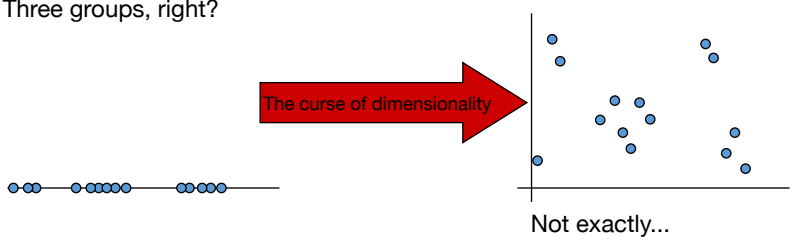
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

7

**NOVA**  
**IMS**  
Information  
Management  
School

## Data Preprocessing

Three groups, right?



Not exactly...

**When the dimensionality increases, the space becomes more sparse and it becomes more difficult to find groups**

**As dimensionality increases, apparent clusters disappear.**

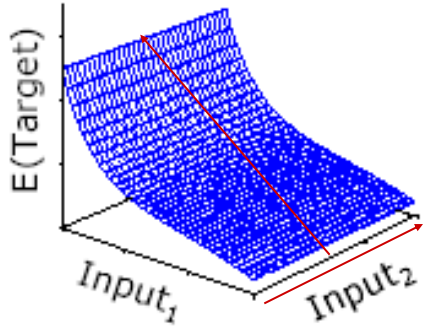
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

8

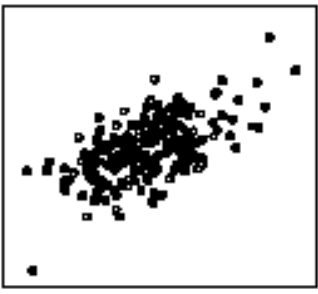
**NOVA**  
**IMS**  
Information Management School

## Data Preprocessing

- **Size Reduction of the Input Space** (or feature selection):  
Two major principles:  
**Irrelevance and Redundancy**



**Irrelevant variable:** doesn't affect target.



**Redundant variable:** highly correlated with another

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa


9

**NOVA**  
**IMS**  
Information Management School

## Reducing Input Space

Relevancy

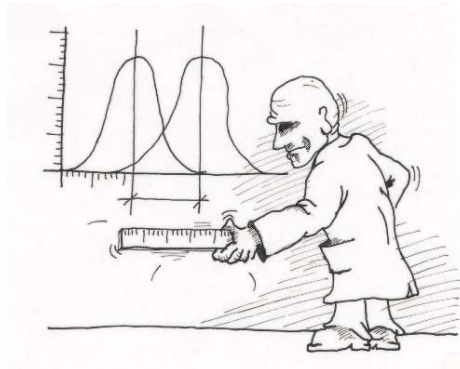
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa



10

## • Size Reduction of the Input Space:

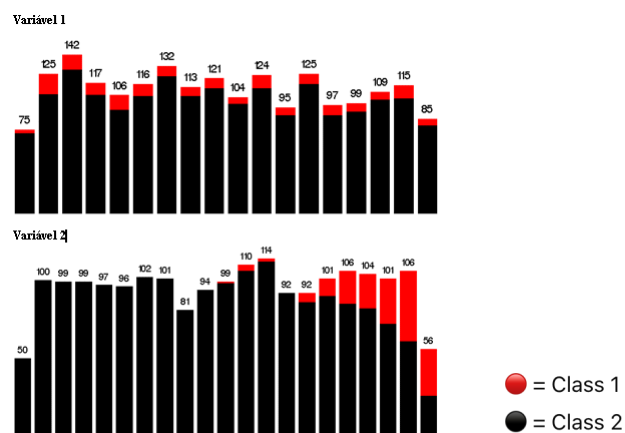
- **Feature selection: measuring which variable best separates the classes**



Feature selection helps us find the variables that best separate the target classes. We are, in essence, measuring which 'ruler' (feature) allows us to distinguish one group from another.

## • Size Reduction of the Input Space:

- **Feature selection: measuring which variable best separates the classes**



**NOVA**  
**IMS**  
Information  
Management  
School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - To create input combinations
    - Height2/weight (Body Mass Index, BMI)
    - Population/area
    - Euros spent/nº of purchases
    - Euros spent/time as customer
    - Debt/income

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

13

**NOVA**  
**IMS**  
Information  
Management  
School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - To create input combinations
 

1. Customer ID (could be anonymous)	8. Average number of different products purchased per transaction
2. Total revenue for the customer	9. Relative spend on each product
3. Number of transactions per customer (frequency)	10. NRS on each product (and where a product taxonomy exists):
4. Average time between transactions (transaction interval)	11. Relative spend in each product subgroup
5. Variance of transaction interval	12. NRS in each product subgroup
6. Customer stability index (ratio of (5)/(4))	13. NRS in each product group
7. Days since last visit (recency)	

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

14

**NOVA**  
**IMS**  
Information Management School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - Heuristic feature selection methods:
    - Best single features
      - Choose by information gain measures (e.g. entropy)
      - A feature is interesting if it reduces uncertainty

No improvement

Perfect Split

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

15

**NOVA**  
**IMS**  
Information Management School

# Reducing Input Space

## Redundancy

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

UNICIS

16



**NOVA**  
**IMS**  
Information Management School

## Data Preprocessing

- **Size Reduction of the Input Space:**

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

17

**NOVA**  
**IMS**  
Information Management School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis
    - A procedure that uses an **orthogonal transformation** to convert a set of observations of possibly correlated variables into a set of values of **linearly uncorrelated variables called principal components**.
    - The number of principal components is **equal to the number of original variables**.
    - This transformation is defined in such a way that the first principal component has the **largest possible variance** (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance.

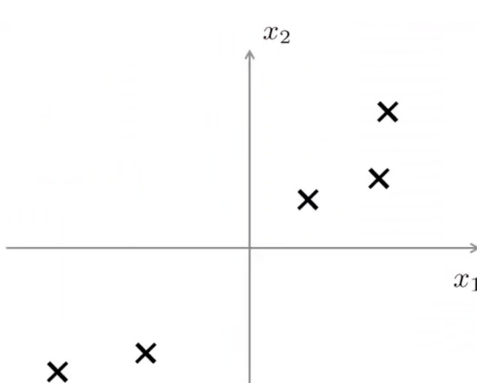
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

18

**NOVA**  
**IMS**  
Information  
Management  
School

Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis



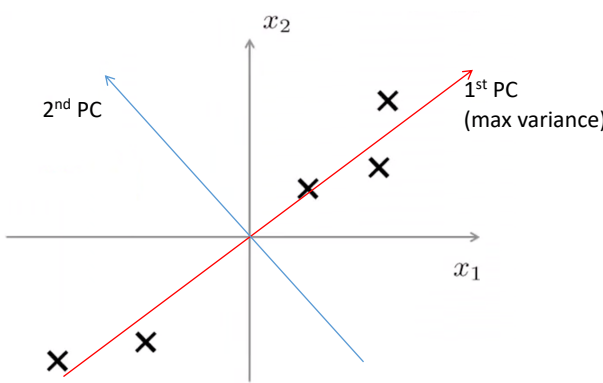
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

19

**NOVA**  
**IMS**  
Information  
Management  
School

Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis



Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

20

**NOVA**  
**IMS**  
Information Management School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis

Algebra: orthonormal transform  
Geometry: axis rotation

Column vector  $x_2$   
Column vector  $x_1$   
N-dimensional space  
Principal comp.  $z_2$   
Principal comp.  $z_1$   
Only needed direction

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

21

**NOVA**  
**IMS**  
Information Management School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - Each principal component (PC) captures a portion of the total variance in the dataset — that is, how much of the information (or variability) in the original variables it represents.
    - PC1 captures the largest possible variance.
    - PC2 captures the next largest, uncorrelated with PC1.
    - Subsequent PCs each explain less and less variance.
  - By summing the variance explained by each component, we get the cumulative variance explained, which tells us how much of the total information in the data is retained when using only the first k components.
  - Example:
    - PC1 explains 60% of the variance,
    - PC2 explains 25%,
    - PC3 explains 10%,
  - then the first three components together explain 95% of the total variance.

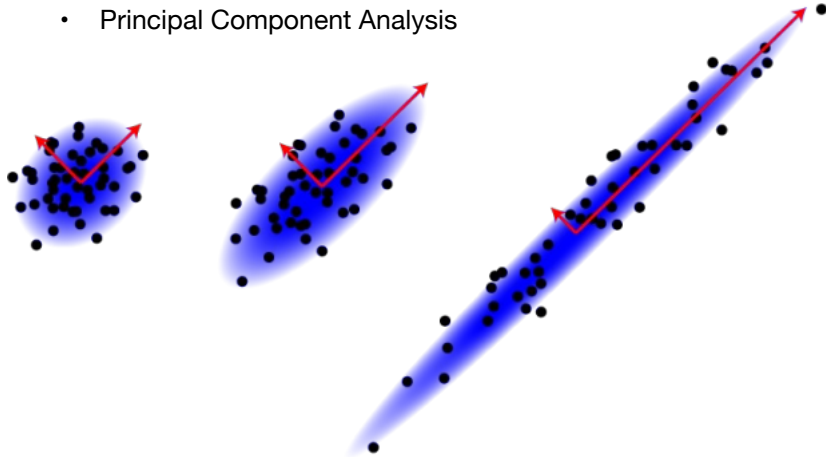
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

22

**NOVA**  
**IMS**  
Information Management School

Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis



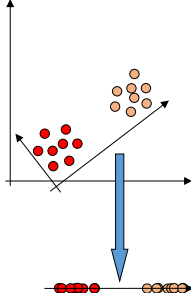
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

23

**NOVA**  
**IMS**  
Information Management School

Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis (careful)



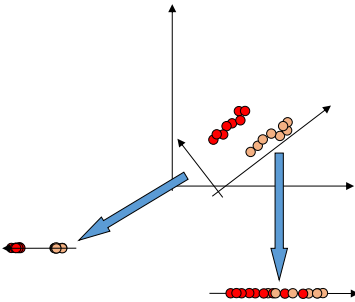
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

24

**NOVA**  
**IMS**  
Information  
Management  
School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis (careful)



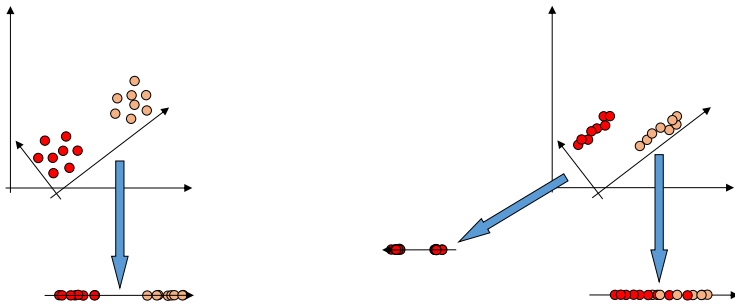
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

25

**NOVA**  
**IMS**  
Information  
Management  
School

## Data Preprocessing

- **Size Reduction of the Input Space:**
  - Principal Component Analysis (careful)



PCA preserves variance, not class separability  
Dimensionality reduction can hurt supervised learning

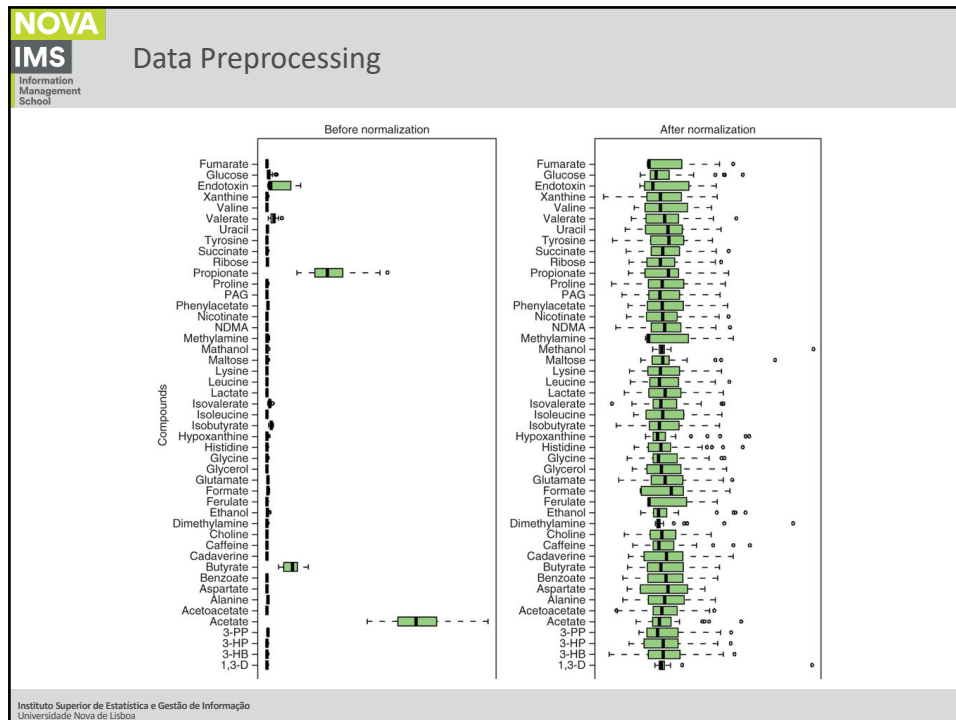
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

26

27

- **Standardization or Normalization:**
  - Many algorithms (e.g., k-means, SVMs) assume features are on comparable scales.
    - Variables come in many different scales (percentages, euros, kilos, meters, days...)
    - Normalization: is about adjusting values measured on different scales to a common scale

28



29

**NOVA**  
**IMS**  
Information Management School

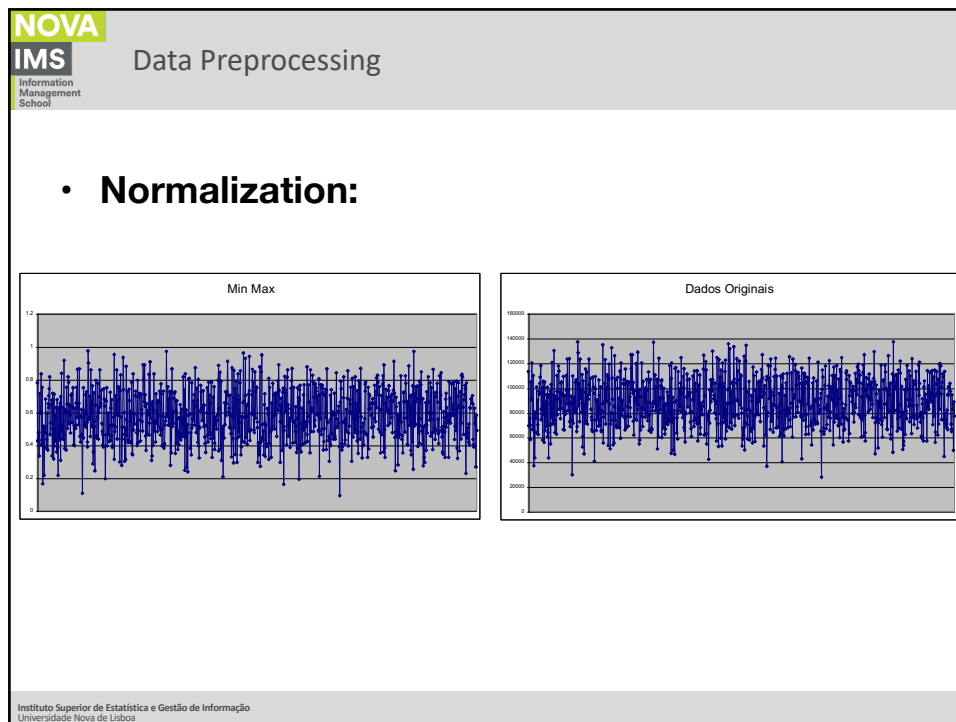
## Data Preprocessing

- Standardization or Normalization:**
  - Min-Max:  $x' = \frac{x - \min(x)}{\max(x) - \min(x)} (\max2 - \min2) + \min2$   

optional
  - Zscore:  $z = \frac{x - \mu}{\sigma}$

Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

30



31



32