

Data Mining

Association Rules

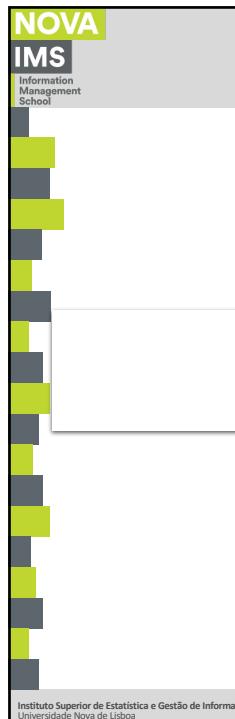
9/12/2021

NOVA-IMS

Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>

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

1



Association Rules

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



2

- **Association Rules:**

- Aims at the extraction of compact patterns that describe subsets of data
 - Events that occur together (market basket analysis);
 - The main purpose is to establish relationships between fields;
 - Are rules of the form if X then Y;
 - Association rules provide information about things that tend to happen together.

- **Association Rules:**

- Table with a set of purchases from a supermarket, with 5 purchases and 5 items.

Cliente	Itens
1	Orange juice, soda
2	Milk, Orange juice, Glass cleaner
3	Orange juice
4	Orange juice, detergent, soda
5	Glass cleaner, soda

- With these data we can create a table of co-occurrences with the number of times that any pair of products was purchased together.

Association Rules

	Orange Juice	Glass Cleaner	Milk	Soda	Detergent
Orange Juice	4	1	1	2	1
Glass Cleaner	1	2	1	1	0
Milk	1	1	1	0	0
Soda	2	1	0	3	1
Detergent	1	0	0	1	2

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

5

Association Rules

- From the analysis, one can conclude that:
 - Orange juice and soda are more likely purchased together than any other two items;
 - Detergent is never purchased with milk or window cleaner;
 - Milk is never purchased with soda or detergent.

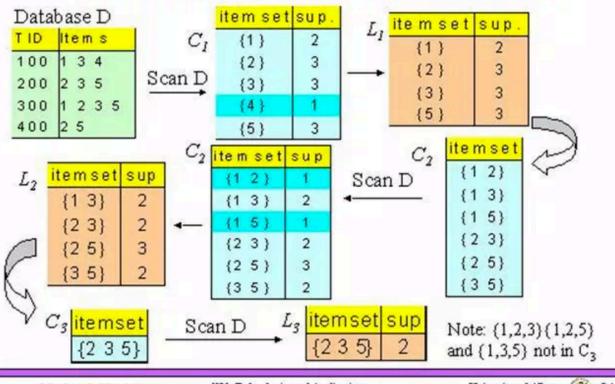
	OJ	GC	Milk	Soda	Det.
OJ	4	1	1	2	1
GC	1	2	1	1	0
Milk	1	1	1	0	0
Soda	2	1	0	3	1
Det.	1	0	0	1	2

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

6

- **Association Rules:**

The Apriori Algorithm -- Example



© Dr. Olave R. Zelmer, 2001

Web Technologies and Applications

University of Alberta

54

- **Association Rules:**

- The rules are expressed in the form of:

"if the item A is part of an event, then the item B will also be part of the event X percent of the time"

- The rules should not be interpreted as a direct causation, but only as an association;
- It is not legitimate to infer rules of causality.

- **How association rules work:**

*If a customer buys **shoes**, then 10% of the time also buys **socks***

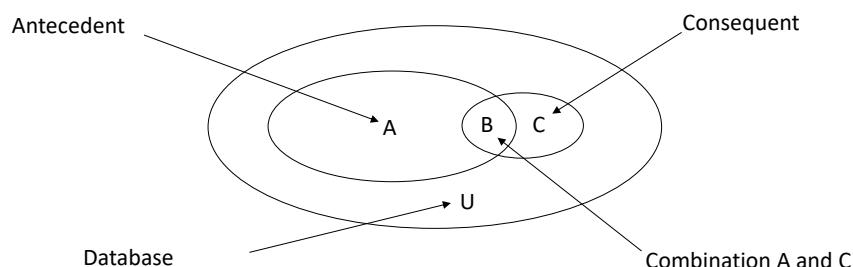
*If a customer buys **plastic paint**, then 85% of the time also buys **brushes***

- All rules have an antecedent and a consequent

*Shoes and plastic paint – **antecedents***

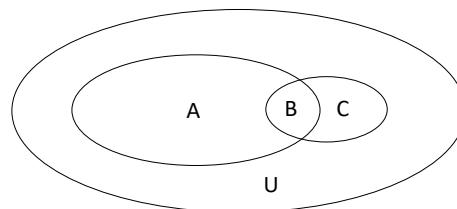
*Socks and brushes - **consequents***

- **Venn Diagram:**



- **How association rules work:**
 - *Evaluating the quality of Association Rules:*
 - Confidence*
 - Support*
 - Expected Confidence*
 - Lift*
 - *The most important measures to assess the rules*

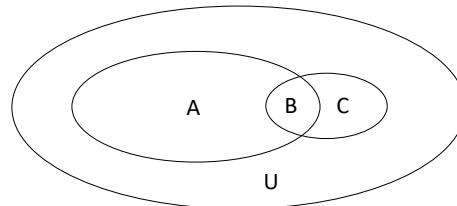
- **How association rules work:**
 - *Evaluating the quality of Association Rules:*
 - **Confidence – the strength of an association, the percentage of a consequent appears given that the antecedent has occurred**



$$\text{Confidence} = B/A$$

- **How association rules work:**

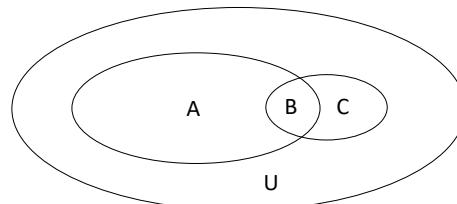
- *Evaluating the quality of Association Rules:*
- **Support** – shows how frequently the combination occurs in the database.



$$\text{Support} = B/U$$

- **How association rules work:**

- *Evaluating the quality of Association Rules:*
- **Expected confidence** – equal to the number of consequent transactions, divided by the total number of transactions



$$\text{ExpectedConfidence} = C/U$$

- **How association rules work:**

- *Evaluating the quality of Association Rules:*
- **Lift** – equal to the confidence factor divided by the expected confidence. Lift is a factor by which the **likelihood** of consequent increases given an antecedent

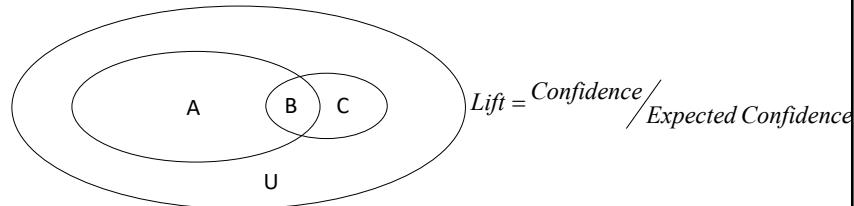


Tabela de Transações	
1,000,000	Total de Transações
200,000	Sapatos
50,000	Meias
20,000	Sapatos e Meias

$$\text{Confidence} = \frac{B}{A}$$

$$\text{Support} = \frac{B}{U}$$



$$\text{Expected Confidence} = \frac{C}{U}$$

$$\text{Lift} = \frac{\text{Confidence}}{\text{Expected Confidence}}$$

Association Rules

Tabela de Transações

1,000,000	Total de Transações
200,000	Sapatos
50,000	Meias
20,000	Sapatos e Meias

Regra

Se um cliente compra sapatos, então 10% das vezes ele comprará meias.

Diagrama de Venn



Critérios de Avaliação

Association Rules

Tabela de Transações

1,000,000	Total de Transações
200,000	Sapatos
50,000	Meias
20,000	Sapatos e Meias

Regra

Se um cliente compra sapatos, então 10% das vezes ele comprará meias.

Diagrama de Venn



Critérios de Avaliação

Confiança: $20,000/200,000 = 10\%$
 Suporte: $20,000/1,000,000 = 2\%$
 Confiança Esperada: $50,000/1,000,000 = 5\%$
 Lift: Confiança/Confiança Esperada = 2

- **How association rules work:**

- *Evaluating the quality of Association Rules:*

Lift – a value of 2 means that, if one bought shoes, one is twice as likely to buy socks than those who have not bought shoes

A **credible rule** must have a good confidence factor, a high level of support and lift higher than 1

Rules with a **high level of confidence** but with a **low support** should be interpreted with caution, as it can result in **idiosyncrasies** that rise due to the small number of cases to support the rule

- **Use:**

The Wine Bible
by Karen MacNeil



List Price: \$19.95
Price: \$13.97 & eligible for FREE Super Saver Shipping on orders over \$25. [See details.](#)

You Save: \$5.98 (30%)

Availability: Usually ships within 24 hours

58 used & new from \$13.00

Edition: Paperback

[Look inside this book](#)

▶ [See more product details](#)

Better Together

Buy this book with [Wine for Dummies](#) by Ed McCarthy, et al today!



Total List Price: \$41.94

Buy Together Today: \$29.36

[Buy both now!](#)

Association Rules

- **Use:**

Customers who bought this book also bought:

- [How to Taste : A Guide to Enjoying Wine](#) by Jancis Robinson ([Rate it](#))
- [The World Atlas of Wine](#) by Hugh Johnson, Jancis Robinson ([Rate it](#))
- [The Wall Street Journal Guide to Wine: New and Improved: How to Buy, Drink, and Enjoy Wine](#) by Dorothy J. Gaiter, John Brecher ([Rate it](#))
- [The University Wine Course: A Comprehensive Text and Tutorial](#) by Marian W. Baldy Ph.D. ([Rate it](#))
- [Parker's Wine Buyer's Guide 6th Edition : The Complete, Easy-to-Use Reference on Recent Vintages, Prices, and Ratings for More Than 8,000 Wines from All the Major Wine Regions](#) by Robert M. Parker ([Rate it](#))

You could win a
\$50 Amazon gift certificate
or the Grand Prize. Start by
adding five items
to your Wish List.
[Learn more.](#)

► [Explore Similar Items: 20 in Books, 20 in DVD, and 16 in Magazine Subscriptions](#)

Association Rules

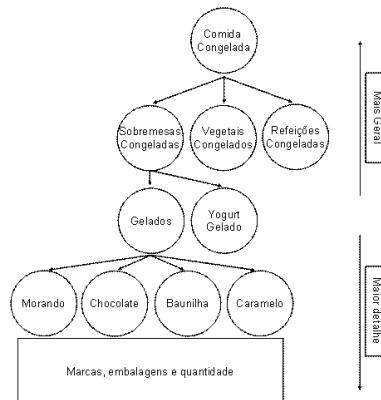
- **Types of rules**

- Trivial rules;
- Inexplicable Rules;
- Actionable Rules;
- Rules that result of promotions made.

- **Additional aspects**

- Available data are essential to success;
- It takes large amounts of data;
- What is an item
 - Frozen pizza, or frozen pizza 4 seasons;
- Setting the appropriate level of detail;
- Classifications.

- **Additional aspects**



- **Additional aspects**

- Use of virtual items;
- Allow you to get information from data beyond the information expressed in the products;
- May include information about the purchase (paid in cash, credit card...), about the customer (new, old...) day of the week, hour of the day.



Data Mining

S3

NOVA-IMS 2021/2022

Fernando Lucas Bação

bacao@isegi.unl.pt

<http://www.isegi.unl.pt/fbacao>

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

1



Agenda

- Data Mining
 - Statistics vs data science
 - The canonical tasks in data mining
 - Exercise
 - The data mining process
 - General aspects of problem definition
 - General aspects of data collection

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

2

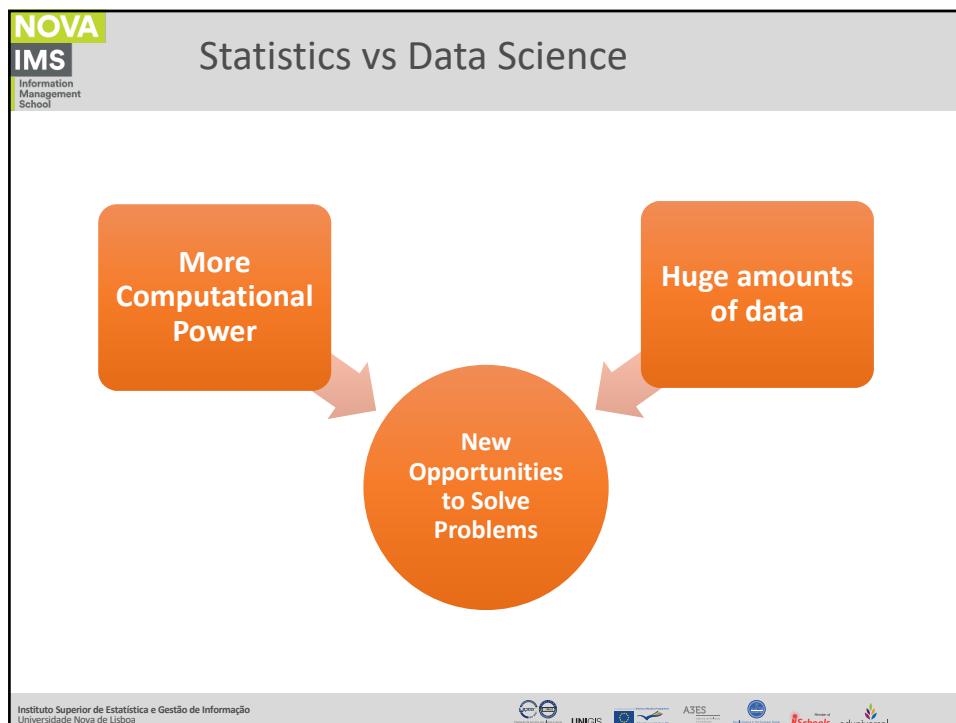


Statistics vs Data Science

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

AACSB UNICIS A3ES iSchools eduniversal

3



4

What new about Data Mining?

(Data Mining: Statistics and More? David J. Hand)

- Statistics might be described as being characterized by **data sets which are small and clean**, which are **static**, which were **sampled in an iid manner**, which were often **collected to answer the particular problem** being addressed, and which are **solely numeric**.
- None of these apply in the data mining context.....

What new about Data Mining?

(Data Mining: Statistics and More? David J. Hand)

- Size of the data sets
 - Data will **not all fit into the main memory** of the computer this means that, if all of the data is to be processed during an analysis, **adaptive or sequential techniques** have to be developed (nonstatistical communities especially to those working in pattern recognition and machine learning).
 - Data sets may be large because the number of records is large or because the number of variables is large (**deep and large**).
 - Data may not be **stored as the single flat** file so but as multiple interrelated flat files.

Incremental (online)

- Examples are presented one at a time and the structure of representation changes.
- In the online learning, the system will handle each instance incrementally, the algorithm itself is updatable, and the knowledge will be updated by every instance in time.

Non incremental (batch)

- Examples are presented all at the same time and are considered together.

Incremental learning algorithms are usually faster than non-incremental algorithms, and for extremely large data sets, non-incremental algorithms may not be applicable at all.

• Size of the data sets

- In the past, in many situations where statisticians have classically worked, the problem has been one of **lack of data** rather than abundance.
- However, when data exists in the **superabundance** the results of tests (significance tests) will lead to **very strong evidence that even tiny effects exist**, effects which are so minute as to be of doubtful practical value.
- In place of statistical significance, we need to consider more carefully substantive significance: **is the effect important or valuable or not?**

- Contaminated data

- In the data mining context, when the analysis is necessarily **secondary**, data is always “dirty”.
- When the data sets are large, it is practically certain that some of the data will be **invalid in some way**.
- This is especially true when the **data describe human interactions** of some kind, such as marketing data, financial transaction data, or human resource data.

- Nonstationarity, Selection Bias, and Dependent Observations

- Very large data sets are **unlikely to arise in an iid manner**;
- **Population drift**, can arise because the underlying population is changing (for example, the population of applicants for bank loans may evolve as the economy heats and cools). Supermarket transactions or Telco phone calls occur every day, not just one day, so that the database is a constantly evolving entity
- **Selection bias**, it arises when developing scoring rules, typically in this situation comprehensive data is available only for those previous applicants who were graded good risk by some previous rule. Those graded bad would have been rejected and hence their true status never discovered.

What new about Data Mining?

(Data Mining: Statistics and More? David J. Hand)

- Spurious Relationships and Automated Data Analysis

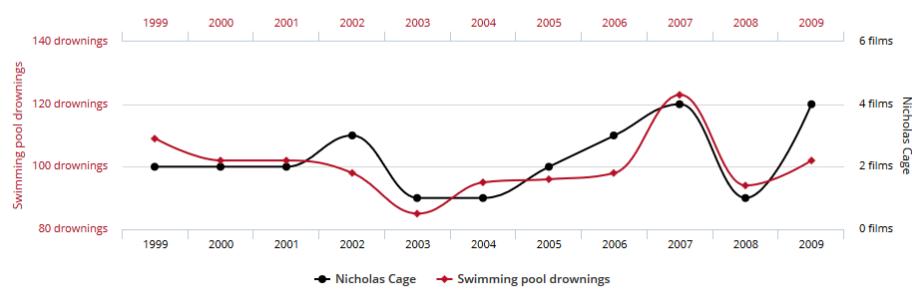
- Because the pattern searches will throw up a **large numbers of candidate patterns**, there will be a **high probability that spurious** (chance) data configurations will be identified as patterns.
- The bottom line is that those patterns and structures identified as potentially interesting will be presented to a **domain expert for consideration to be accepted or rejected** in the context of the substantive domain and objectives, and not merely on the basis of internal statistical structure.

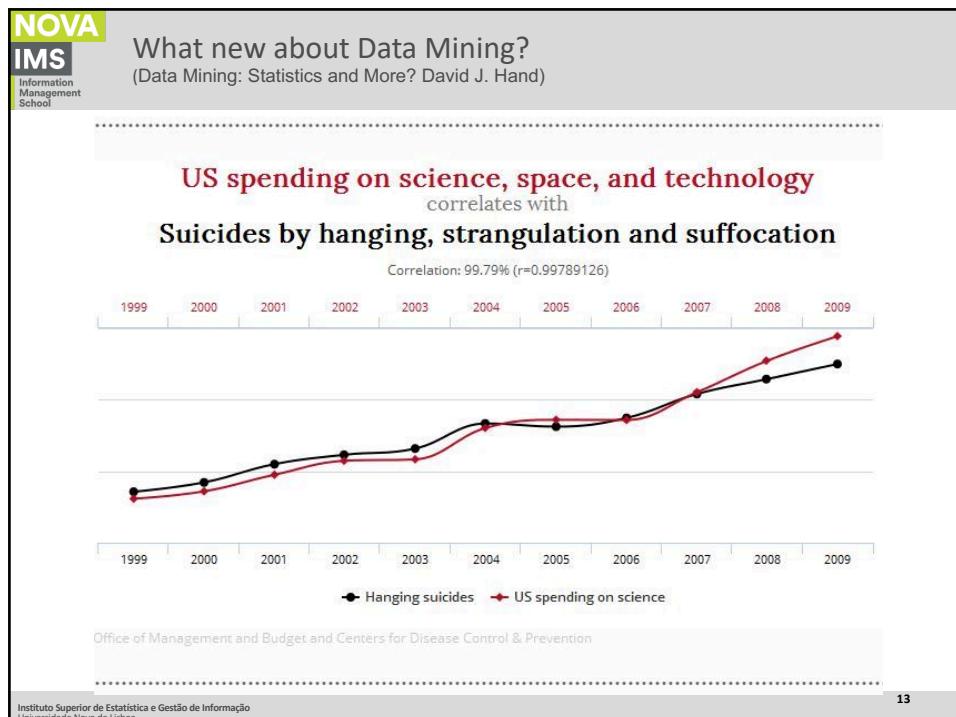
What new about Data Mining?

(Data Mining: Statistics and More? David J. Hand)

Number of people who drowned by falling into a pool correlates with Films Nicolas Cage appeared in

Correlation: 66.6% ($r=0.666004$, $p>0.05$)





13

**NOVA
IMS**
Information Management School

Statistics vs Data Science

	Experimental Primary	Opportunistic Secondary
Purpose	Research	Operational
Value	Scientific	Commercial
Origin	Controlled	Passively observed
Size	Small	Massive
Hygiene	Clean	Dirty
Status	Static	Dynamic

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

14

14

What is Data Mining?

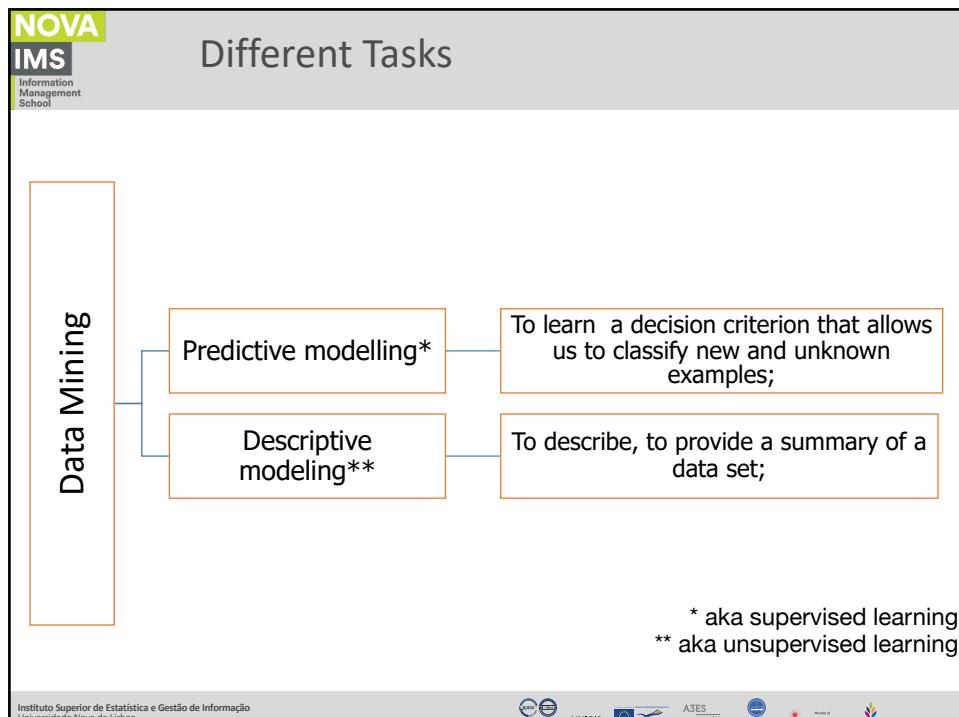
*Data mining is a new **discipline lying at the interface** of statistics, database technology, pattern recognition, machine learning, and other areas.*

*It is concerned with **the secondary analysis of large databases** in order to find previously unsuspected relationships which are of interest or value to the database owners.*

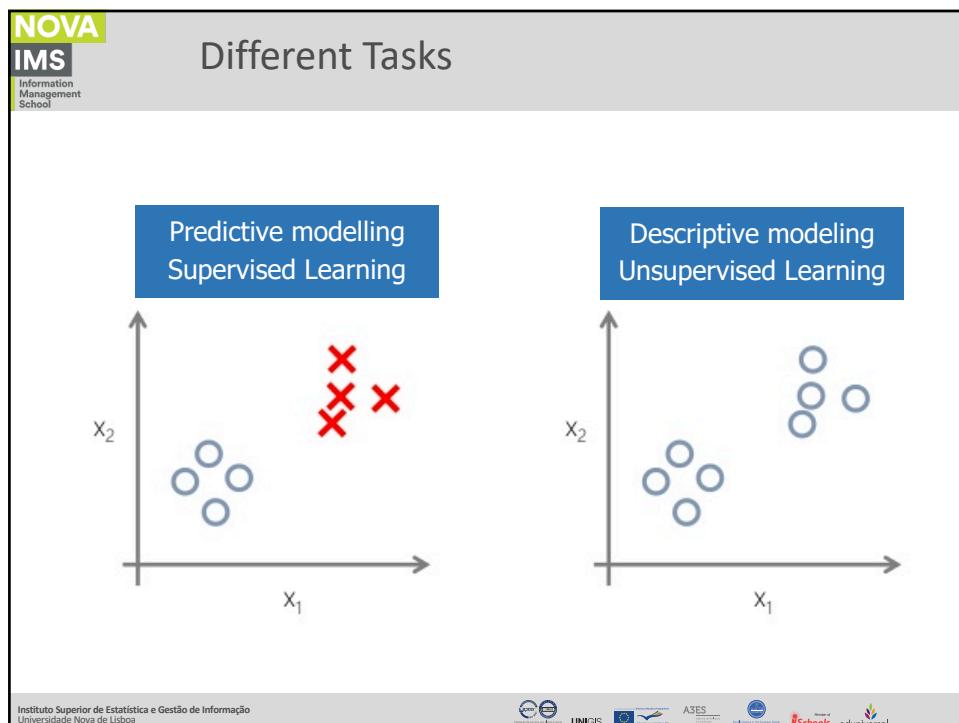
*Data analysis is as much an **art as a science**.*

David J. Hand

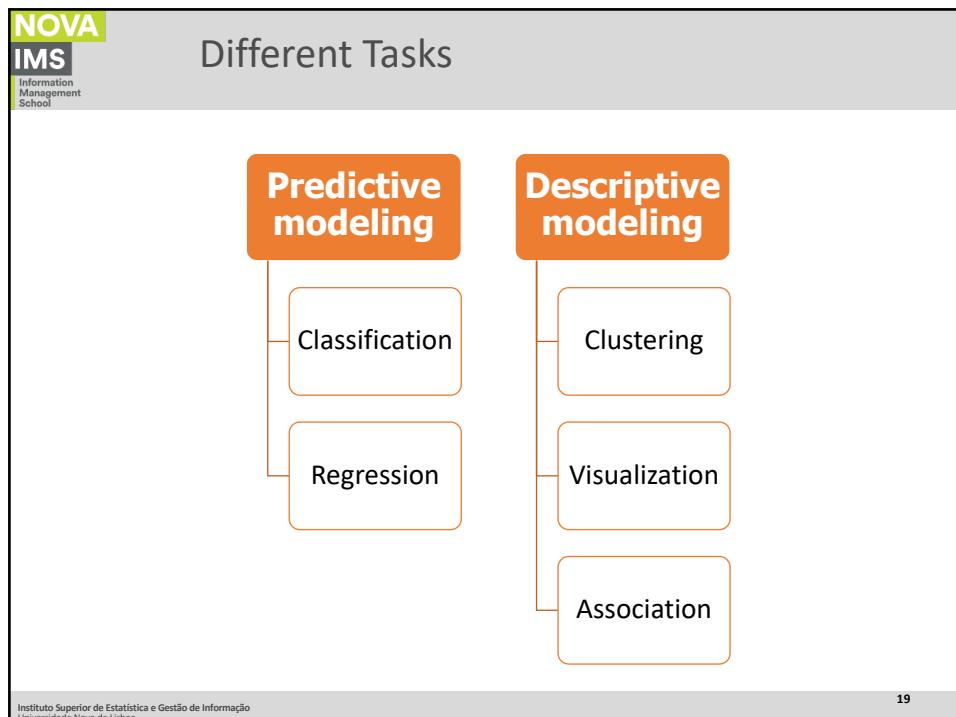
The Canonical Tasks in Data Mining



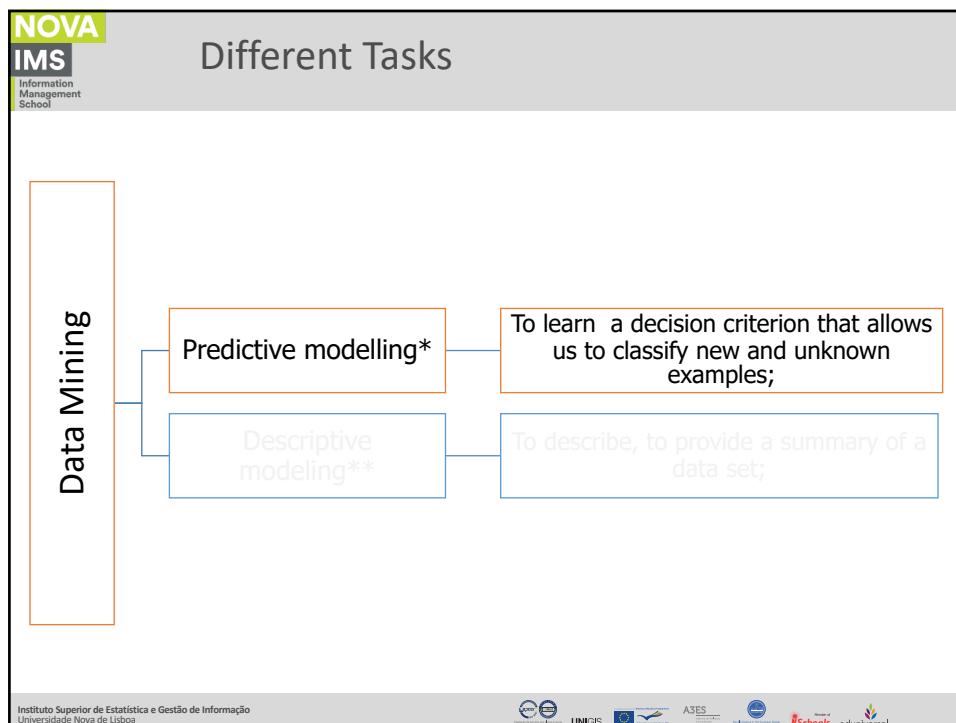
17



18



19



20

NOVA
IMS
Information Management School

Predictive Modeling

The diagram shows a table with 7 columns. The first 6 columns are labeled "Feature": Height, Weight, Sex, Age, Income, and Physical Activity. The 7th column is labeled "Label": Insurance Costs. Blue arrows point from the labels above each column to their respective column headers. The data rows are:

Height	Weight	Sex	Age	Income	Physical Activity	Insurance Costs
1.60	79	M	41	3000	S	N
1.72	82	M	32	4000	S	N
1.66	65	F	28	2500	N	N
1.82	87	M	35	2000	N	S
1.71	66	F	42	3500	N	S

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

Accredited by EQUIS, UNIGIS, A3ES, AACSB, iSchools, eduniversal

21

NOVA
IMS
Information Management School

Predictive Modeling

The diagram illustrates the machine learning process. It is divided into two main sections: "Learning" and "Classification".

Learning: A blue box labeled "Examples (training)" has an arrow pointing to an orange oval labeled "Algorithm". An arrow points from the "Algorithm" to a blue box labeled "Knowledge". A dashed line connects the "Knowledge" box to the "Classification" section.

Classification: A dashed line connects the "Knowledge" box to the "Classification" section. In the "Classification" section, a blue box labeled "Examples (new)" has an arrow pointing to an orange oval labeled "Classifier". An arrow points from the "Classifier" to a blue box labeled "Classification".

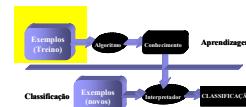
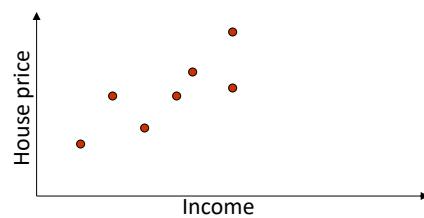
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

Accredited by EQUIS, UNIGIS, A3ES, AACSB, iSchools, eduniversal

22

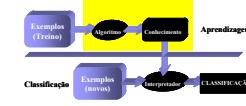
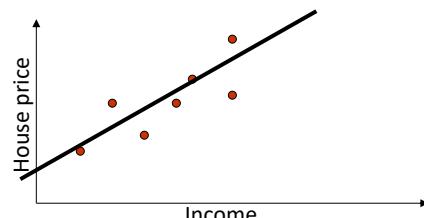
Predictive Modeling

- A real estate agency wants to estimate the price range for each customer based on their income;
- Training examples:
 - Historical data;
 - Income vs sold house prices.



Predictive Modeling

- Algorithm
 - Linear regression
- Knowledge representation
 - Regression line (slope and origin)



NOVA
IMS
Information Management School

Predictive Modeling

- New examples
 - A customer with an income of x
- Interpretation
 - Use the line (prediction method) to obtain an estimate

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

A3ES

25

NOVA
IMS
Information Management School

Predictive Modeling

Classification

Regression

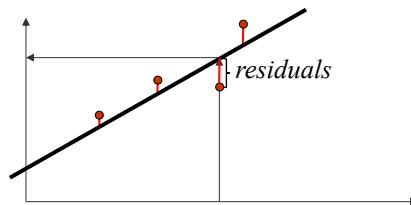
Source: <http://ipython-books.github.io/featured-04/>

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

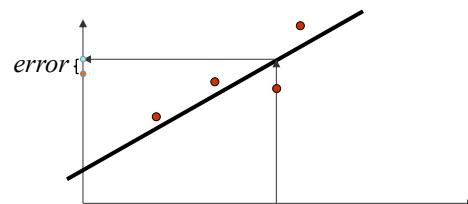
26

26

- If we are facing a regression problem, then we have to consider the average deviations produced by the model.



- If we are facing a regression problem, then we have to consider the average deviations produced by the model.



Predictive Modeling

- If we are facing a classification problem, then we only need to count the number of times that the model was wrong;

Confusion Matrix

	Actually Positive (1)	Actually Negative (0)
Predicted Positive (1)	True Positives (TPs)	False Positives (FPs)
Predicted Negative (0)	False Negatives (FNs)	True Negatives (TNs)

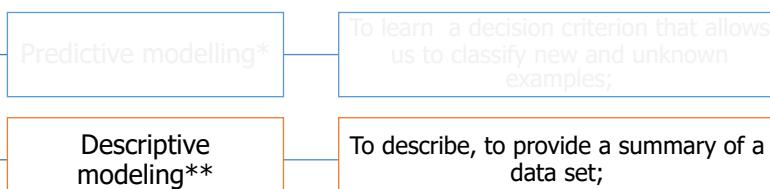
$$\text{Accuracy} = \frac{tp + tn}{tp + tn + fp + fn}$$

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

Different Tasks

Data Mining



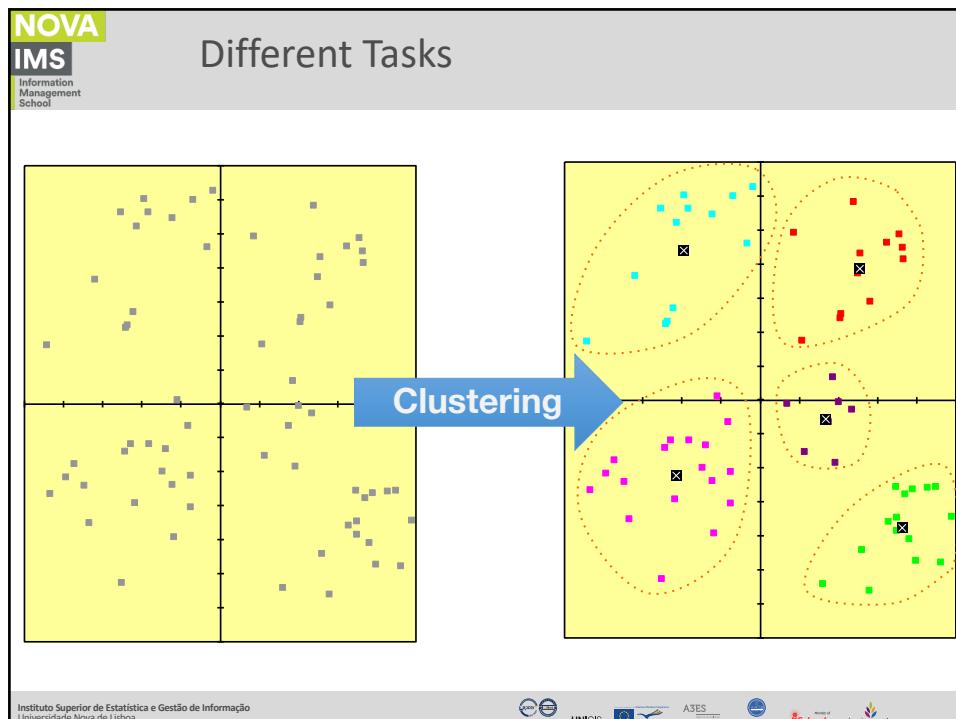
Different Tasks

Feature

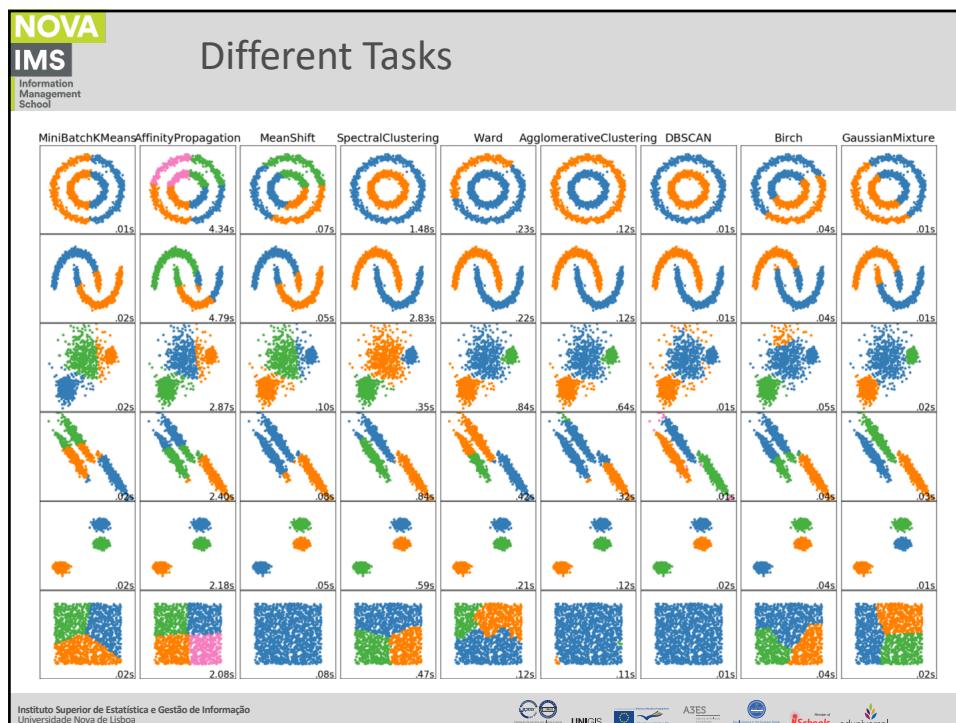


Height	Weight	Sex	Age	Income	Physical Activity
1.60	79	M	41	3000	S
1.72	82	M	32	4000	S
1.66	65	F	28	2500	N
1.82	87	M	35	2000	N
1.71	66	F	42	3500	N

Clustering



33



34

NOVA
IMS
Information Management School

Association Rules

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

AES



35

NOVA
IMS
Information Management School

Different Tasks

Transaction Table

1,000,000 Total Transactions
200,000 Shoes
50,000 Socks
20,000 Shoes and Socks

Rule
If a customer purchases shoes, then 10% of the time he or she will purchase socks.

Evaluation Criteria:
Confidence: $20,000/200,000 = 10\%$
Support $20,000/1,000,000 = 2\%$
Expected Confidence $= 50,000/1,000,000 = 5\%$
Lift = Confidence/Expected Confidence = 2



Note: The confidence factor with socks on the left-hand side and shoes on the right-hand side is 40% ($20,000/50,000$).
The lift value of two implies that you are twice as likely to buy socks if you bought shoes than if you did not buy shoes.

Figure 1. Association Discovery Statistics Example

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

AES

36

NOVA
IMS
Information Management School

Visualization

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

A3ES UNIGIS iSchools eduniversal

37

NOVA
IMS
Information Management School

Different Tasks

Flattening a 3D Chart

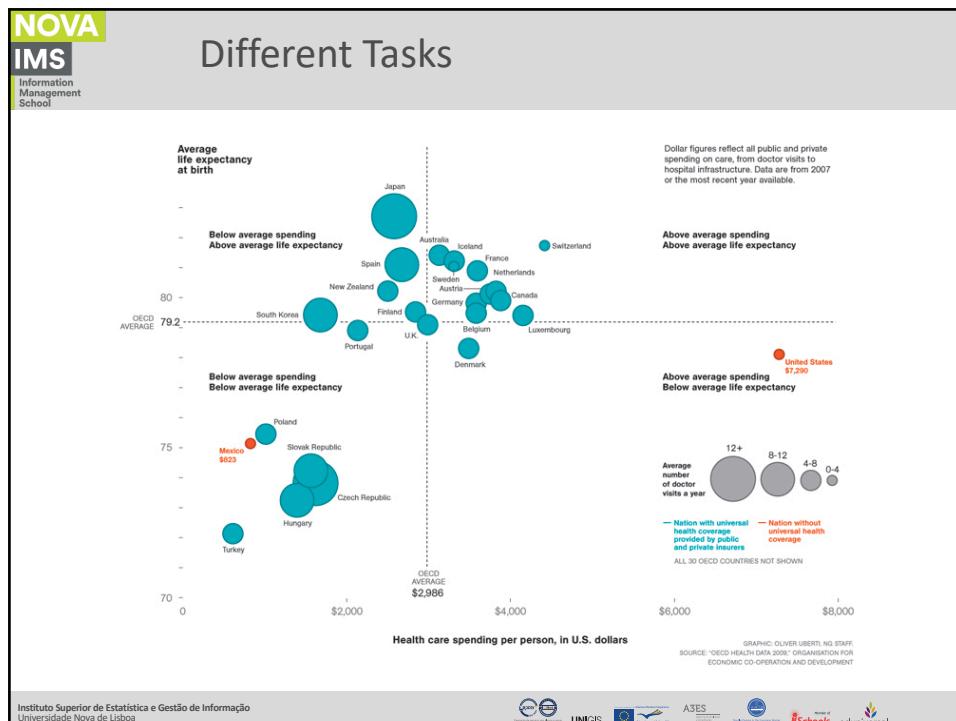
Change in real weekly wages of US-born workers by group, 1990–2006
(Percent)

Education Group	Experience Group	Change in Real Weekly Wages (Percent)
Some High School	Young (experience below 20 years)	0.4
	Old (experience above 20 years)	-5.4
High School Graduate	Young (experience below 20 years)	-1.2
	Old (experience above 20 years)	-1.3
Some College	Young (experience below 20 years)	-1.2
	Old (experience above 20 years)	-3.0
College Graduate	Young (experience below 20 years)	11.3
	Old (experience above 20 years)	6.0

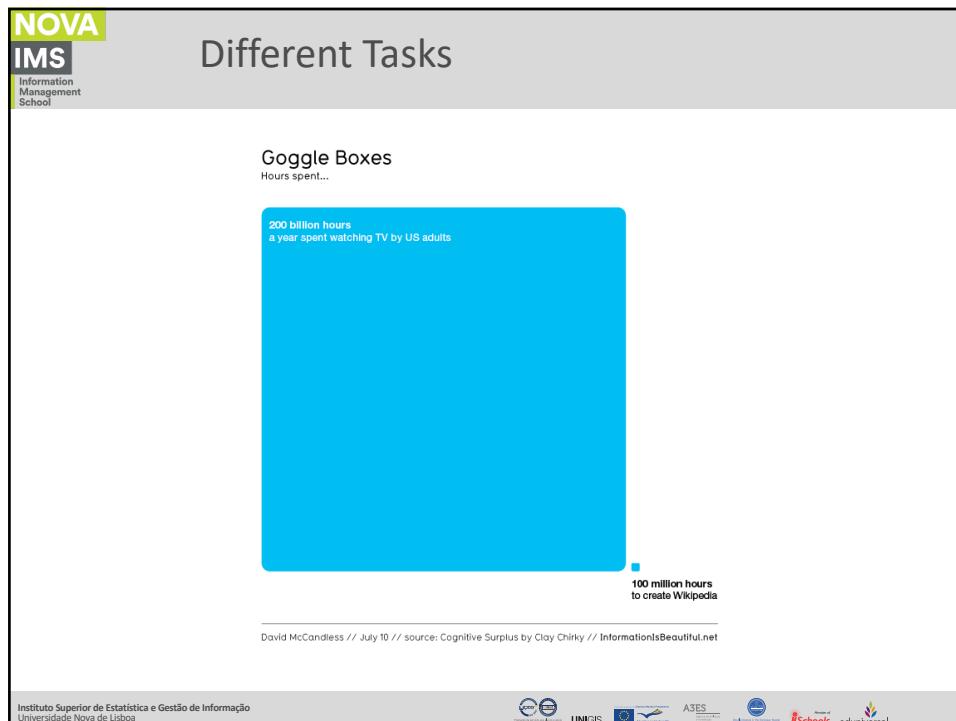
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

A3ES UNIGIS iSchools eduniversal

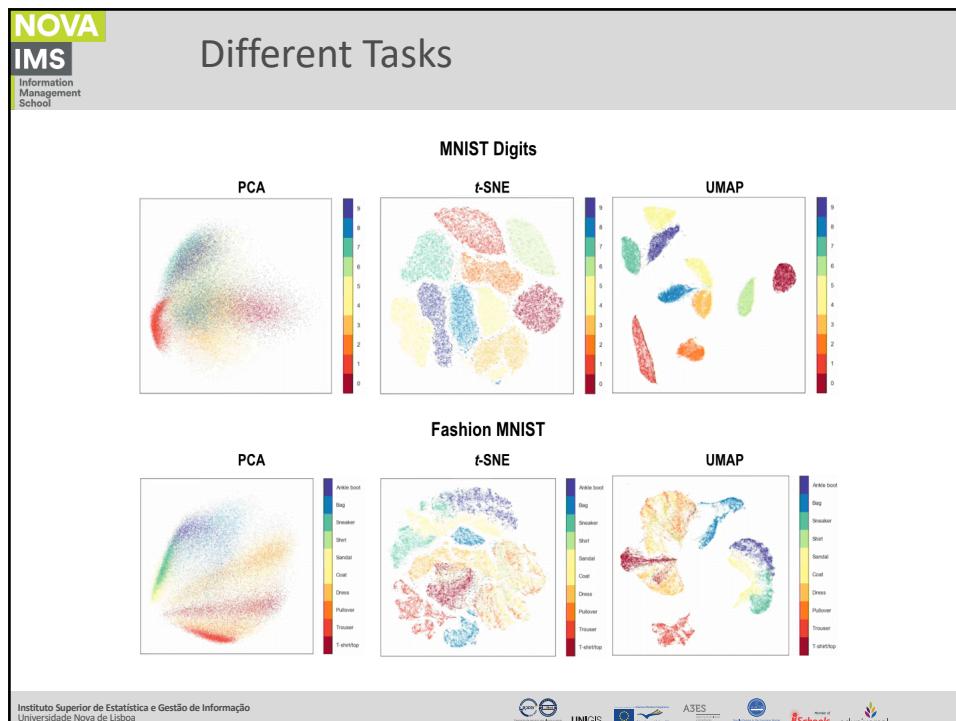
38



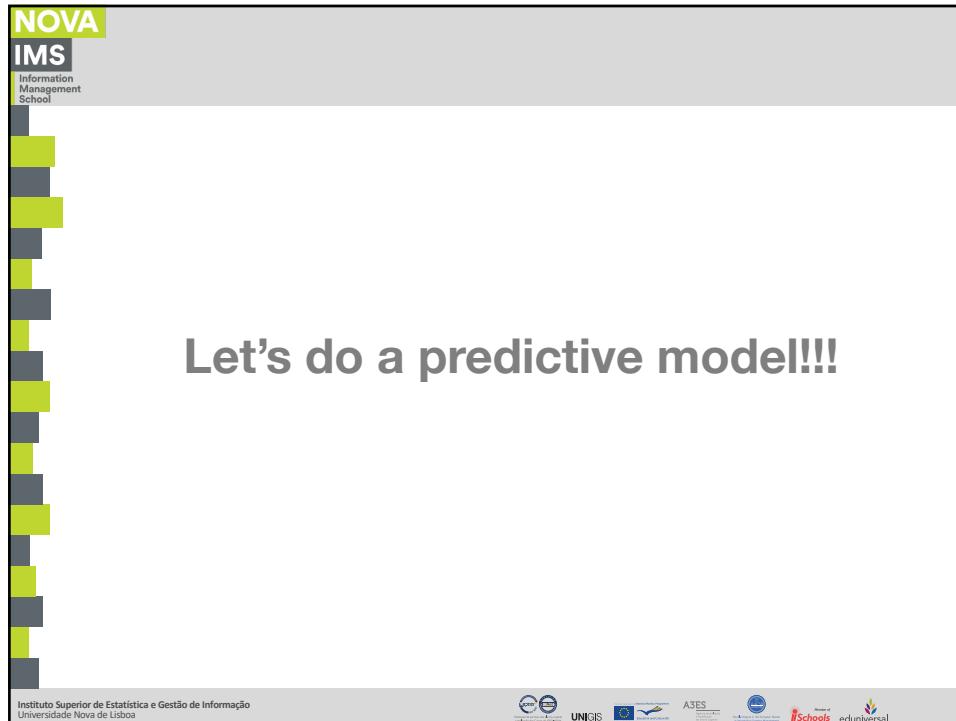
39



40



41

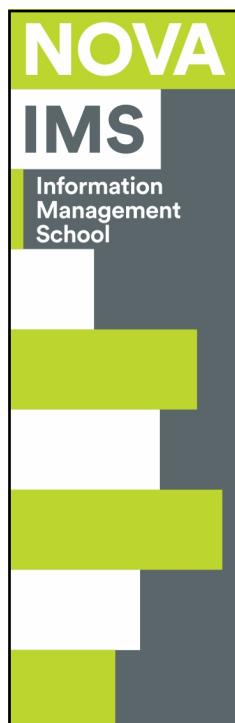


42



Questions?

43



Data Mining

S3

NOVA-IMS 2021/2022

Fernando Lucas Bação

bacao@isegi.unl.pt

<http://www.isegi.unl.pt/bacao>

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

1



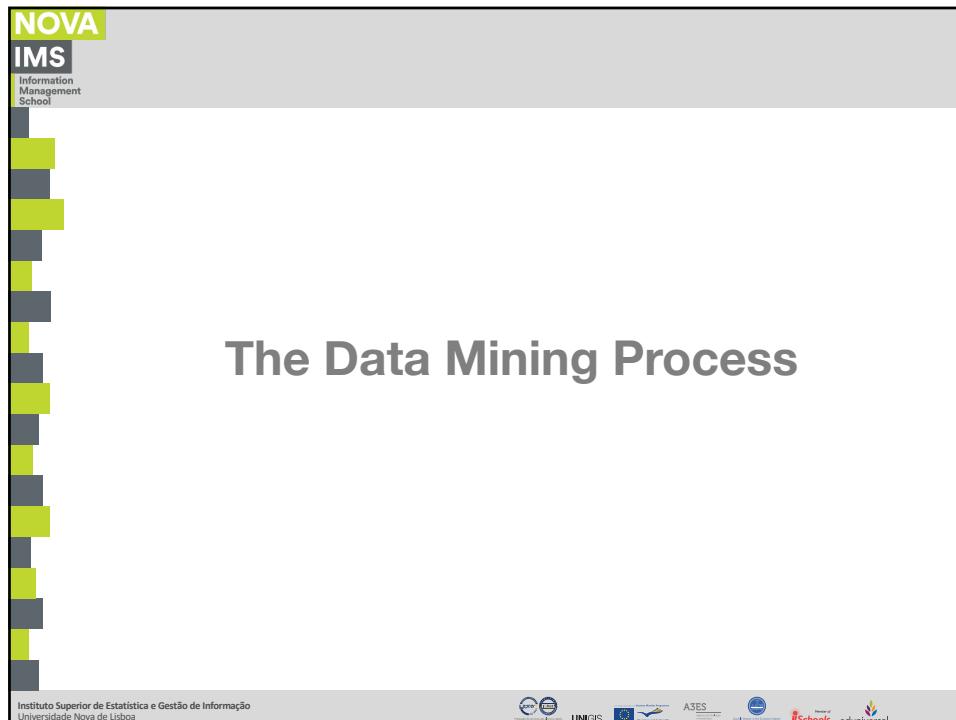
Agenda

- Data Mining
 - The data mining process
 - General aspects of problem definition
 - Input space
 - The curse of dimensionality
 - Input space coverage
 - Binary and multiclass classification
 - Separability and Bayes error
 - Different types of variables
 - Spurious correlations and confounding variables
 - Performance evaluation

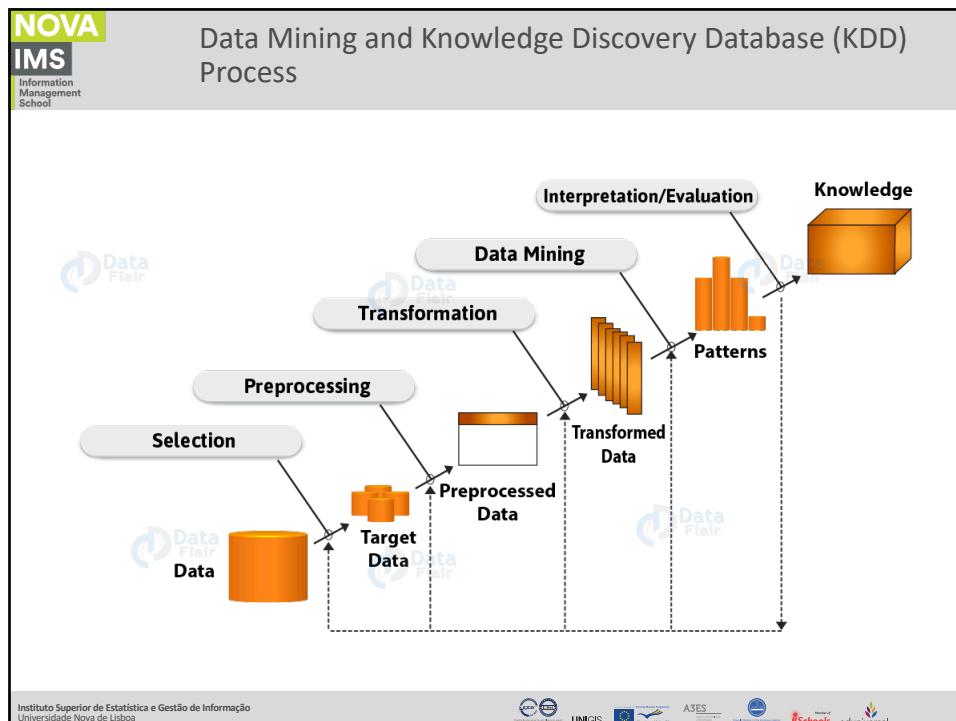
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

UNIGIS A3ES Schools eduniversal

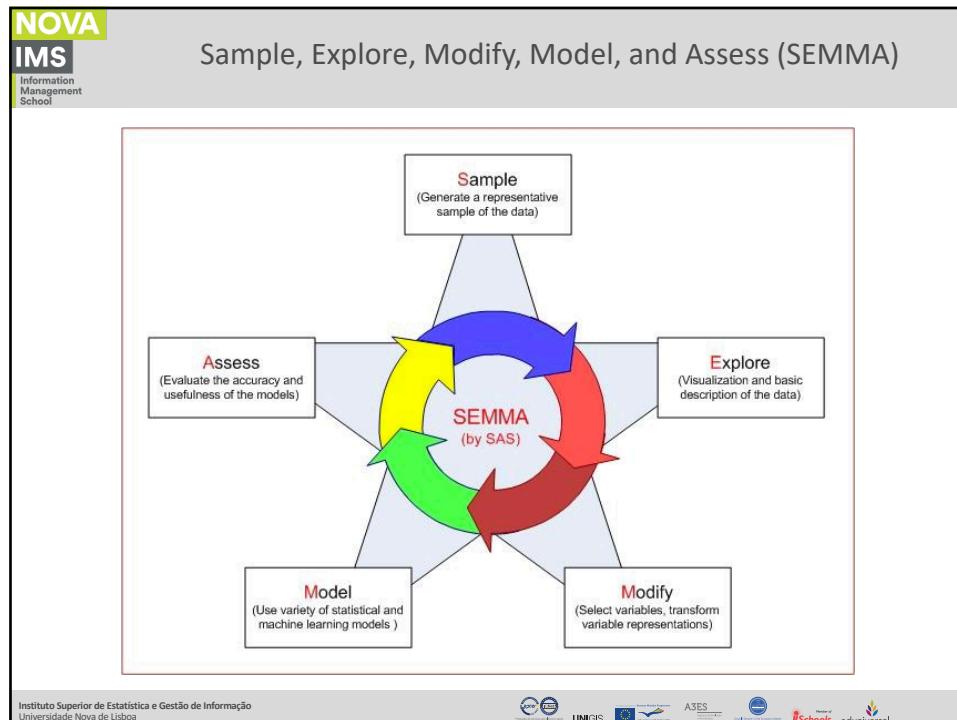
2



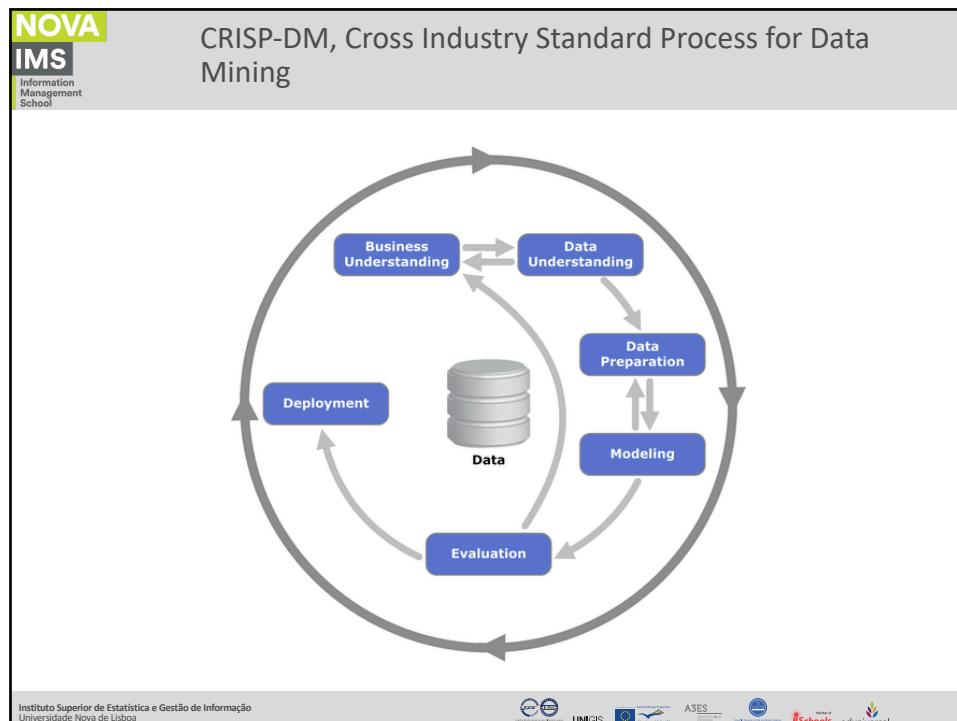
3



4



5



6

NOVA
IMS
Information Management School



General aspects of the problem definition

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



7

NOVA
IMS
Information Management School

Problem definition

- “We’re doomed to complex theories that will never have the **elegance of physics equations**. But if that’s so, we should stop acting as if our goal is to author extremely elegant theories, and instead **embrace complexity** and make use of the best ally we have: **the unreasonable effectiveness of data.**”
- Invariably, **simple models** and a **lot of data** trump more elaborate models based on less data.

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



8

Problem definition

- So, follow the data. Choose a representation that can use **unsupervised learning** on **unlabeled data**, which is so **much more plentiful than labeled data**.
- Represent all the data with a nonparametric model rather than trying to summarize it with a parametric model: "**let the data speak for themselves**"

Problem definition

- Suppose you've constructed the **best set of features** you can, but the **classifiers you're getting are still not accurate enough**.
- What can you do now? There are **two main choices**:
 - design a **better learning algorithm**,
 - or **gather more data** (more examples, and possibly more raw features, subject to the curse of dimensionality).
- Machine learning researchers are mainly concerned with the former, but **pragmatically the quickest path to success is often to just get more data**.
- As a rule of thumb, **a dumb algorithm with lots and lots of data beats a clever one with modest amounts of it**. (After all, machine learning is all about **letting data do the heavy lifting**.)

NOVA
IMS
Information Management School



Input Space

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

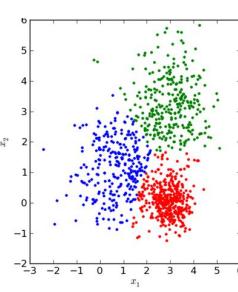


11

NOVA
IMS
Information Management School

Problem definition

- **Input Space**
 - The input space is defined by the input feature vectors.
 - Where the algorithms will try to find a solution to the problem



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



12

NOVA
IMS
Information Management School



The Curse of Dimensionality

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



13

NOVA
IMS
Information Management School

Problem definition

- Number of attributes to be used
 - Few attributes
 - We are unable to distinguish classes.
 - Many attributes
 - Common case in Data Mining;
 - The curse of dimensionality;
 - Difficult visualization and "weird" effects.
 - Important vs. redundant attributes
 - What are the most important attributes for the task?

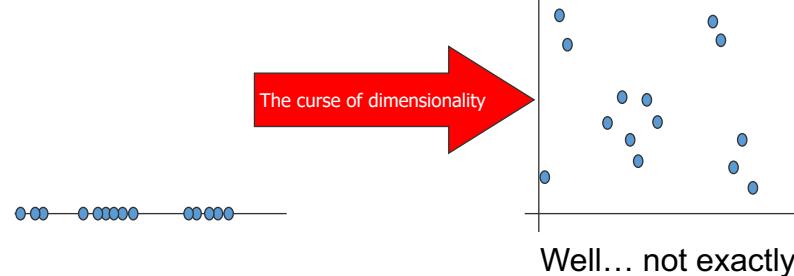
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa



14

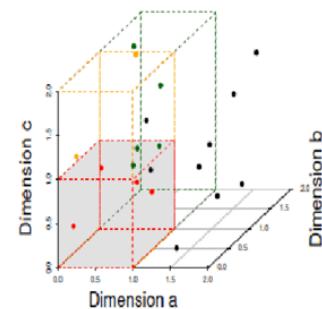
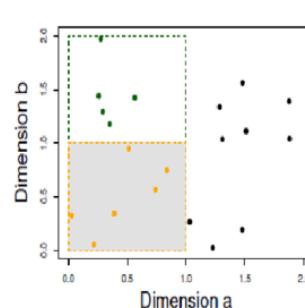
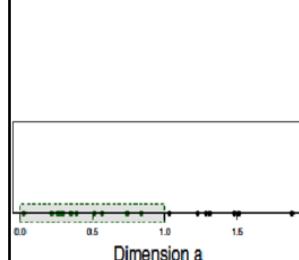
Problem definition

Three groups, right?



When the dimensionality increases, the space becomes more sparse and it becomes more difficult to find groups (you need even more data)

Problem definition

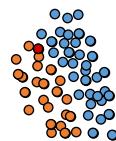


Problem definition

- The curse of dimensionality
 - Generalizing correctly becomes exponentially harder as the dimensionality (number of features) of the examples grows, because a fixed-size training set covers a dwindling fraction of the input space.
 - With a dimension of 100 and a huge training set of a trillion examples, the latter covers only a fraction of about 10^{-18} of the input space. This is what makes machine learning both necessary and hard.

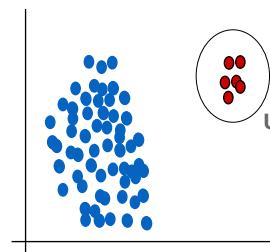
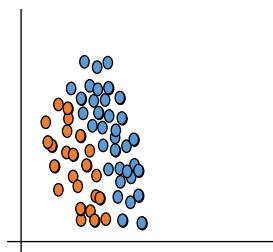
Input Space Coverage

Problem definition



- Good coverage of the problem space increases confidence in the results and in its quality.

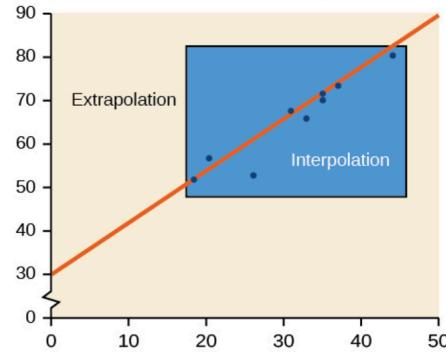
Space coverage



Unknown area where
there are no
training examples

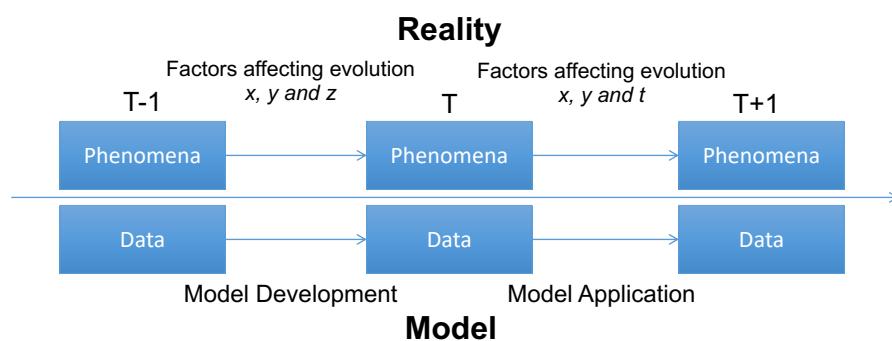
- If a model is developed based on a set of examples, but in fact the examples would be very different, it is natural that the results will be bad (@men women photo).

Extrapolation vs Interpolation



- **Interpolation** involves predicting a value inside the domain and/or range of the data.
- **Extrapolation** involves predicting a value outside the domain and/or range of the data.

General aspects of data collection



NOVA
IMS
Information Management School

Binary and Multiclass Classification

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

A3ES UNIGIS A Schools eduniversal

23

NOVA
IMS
Information Management School

Binary and Multi-class

Binary classification:

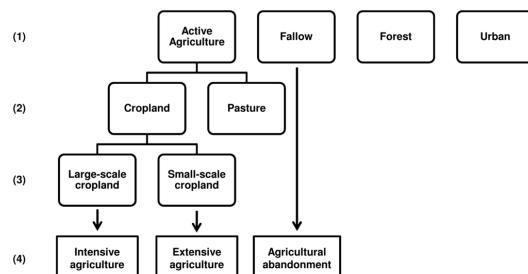
Multi-class classification:

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

A3ES UNIGIS A Schools eduniversal

24

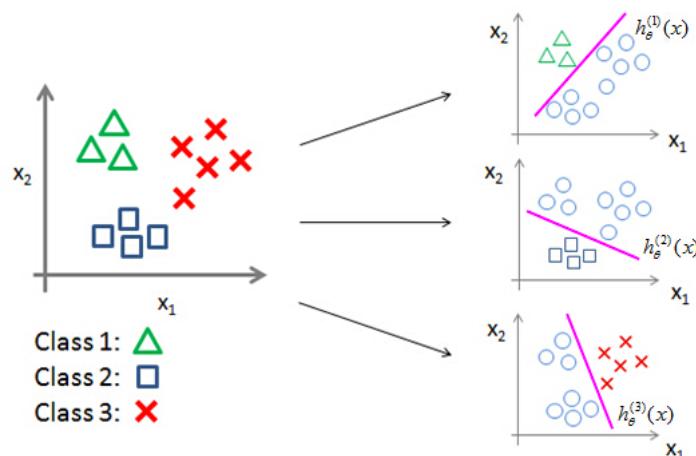
- **Granularity of classification problems:**
 - Initially a small number of output classes;
 - The more output classes, the greater the number of data needed for training.



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

UNIGIS A3ES

25



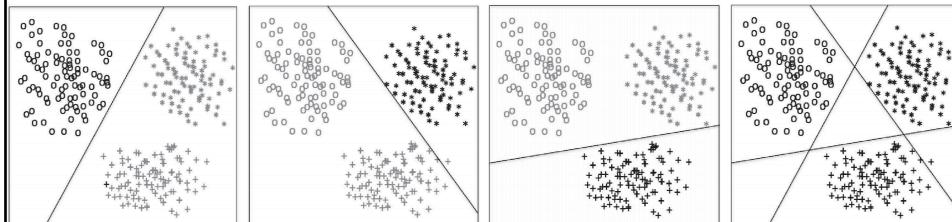
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

UNIGIS A3ES

26

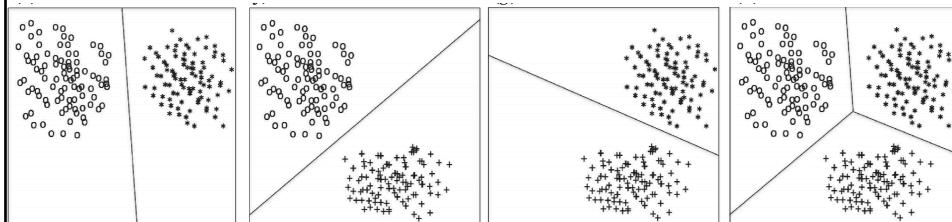
- **One-vs-rest (OVR) strategy:**

- Breaks the multi-class classification down into a series of binary classification
- N-class classification problem is decomposed into N binary classification problems.



- **One-vs-one (OVO) strategy:**

- Enumerating all possible pairs of classes and then to develop a binary classifier for each pair of classes
- Classification is then done by inputting the data point into each particular binary classifier and labelling by majority voting $\frac{1}{2}N(N-1)$



NOVA
IMS
Information Management School

Separability and Bayes Error

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

UNIGIS A3ES iSchools eduniversal

29

NOVA
IMS
Information Management School

Separation and error

not linearly separable

linearly separable

petal width (cm)

petal length (cm)

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

UNIGIS A3ES iSchools eduniversal

30

NOVA
IMS
Information Management School

Problem definition

Separable

- Ø error possible

Not separable

- Always error > Ø
- Bayes error
 - Lowest possible error for a classifier

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

A3ES UNIGIS

31

NOVA
IMS
Information Management School

Different Types of Variables

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

A3ES UNIGIS

32

- **Different types of variables:**
 - **Nominal** are just labels, e.g. ‘red’, ‘green’, ‘blue’, no particular order. Think in classes.
 - **Ordinal** have an order, e.g. ‘satisfied’, ‘very satisfied’, ‘extremely satisfied’. Think in ranks.
 - **Discrete** are just counting data, e.g. 0, 1, 2, ...
 - **Continuous** are just measurement data, e.g. 1.23, 0.001, etc

- **Different types of variables:**
 - **Interval** data are measured and have constant, equal distances between values, but the zero point is arbitrary. The zero isn't meaningful, it doesn't mean a true absence of something.
 - When a **ratio** between two values of a quantitative variable is meaningful, it's a ratio scaled variable. Ratio measurement assumes a zero point where there is no measurement.

Problem definition

Provides:	Nominal	Ordinal	Interval	Ratio
The "order" of values is known		✓	✓	✓
"Counts," aka "Frequency of Distribution"	✓	✓	✓	✓
Mode	✓	✓	✓	✓
Median		✓	✓	✓
Mean			✓	✓
Can quantify the difference between each value			✓	✓
Can add or subtract values			✓	✓
Can multiple and divide values				✓
Has "true zero"				✓

Summary of data types and scale measures

Metadata

- **Metadata:**

- Metadata is "data [information] that provides information about other data".
- Three distinct types of metadata exist:
 - Descriptive metadata describes a resource for purposes such as discovery and identification.
 - Structural metadata is metadata about containers of data and indicates how compound objects are put together, for example, how pages are ordered to form chapters.
 - Administrative metadata provides information to help manage a resource, such as when and how it was created, file type and other technical information, and who can access it.

NOVA
IMS
Information Management School



Spurious Correlations and Confounding Variables

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



37

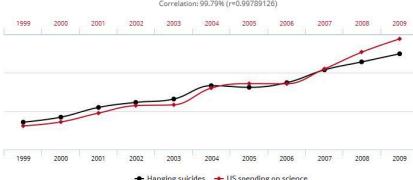
NOVA
IMS
Information Management School

Problem definition

- Input variables should be causally related to the outputs**
 - Spurious correlations
 - Low number of training examples;
 - Large number of input variables.
 - It is important that there is a plausible reason to choose the input variables.

US spending on science, space, and technology correlates with **Suicides by hanging, strangulation and suffocation**

Correlation: 99.79% ($r=0.99789126$)



Office of Management and Budget and Centers for Disease Control & Prevention

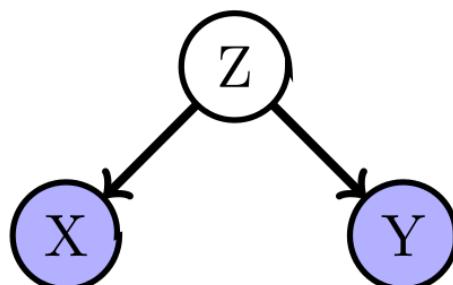
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa



38

- **Input variables should be causally related to the outputs**
 - Confounding variables
 - In statistics, a confounding variable (also confounding factor) is an **extraneous variable** in a statistical model that **correlates** (directly or inversely) with **both the dependent variable and the independent variable**.
 - A **spurious relationship** is a perceived relationship between an independent variable and a dependent variable that has been **estimated incorrectly** because the estimate **fails to account for a confounding factor**.

- **Input variables must be causally related to the outputs**
 - Confounding variables



• Spurious correlations

- Example n.º 1: Ice cream sales and the number of drowning's;
- Example n.º 2: Correlation between the measuring of a patient's temperature on admission to hospital and the probability of his survival;

Performance Evaluation

Problem definition

- Result Evaluation:**

		Actual Value	
		Positives	Negatives
Predicted Value	Positives	TP True Positives	FP False Positives
	Negatives	FN False Negatives	TN True Negatives

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Population}}$$

$$\text{Specificity (true negative rate)} = \frac{\text{True Negatives}}{\text{Total Negatives}}$$

Problem definition

- Result Evaluation:**

		Actual Value	
		Positives	Negatives
Predicted Value	Positives	TP True Positives	FP False Positives
	Negatives	FN False Negatives	TN True Negatives

$$\text{Recall} = \frac{\text{True Positives}}{\text{Total Positives}}$$

$$\text{Precision} = \frac{\text{True Positives}}{\text{Predicted Positives}}$$

Problem definition

- Result Evaluation:**

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$$

$$R^2 = 1 - \frac{MSE(model)}{MSE(baseline)}$$

$$MSE(baseline) = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$$

Problem definition

- Result Evaluation:**

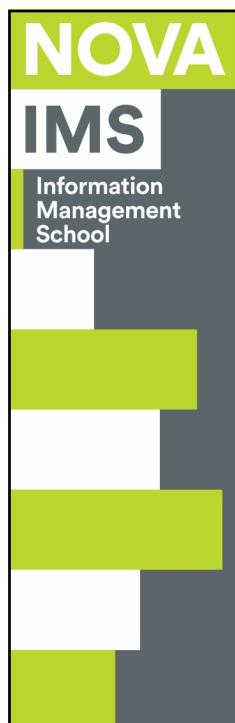
- What is the level of accuracy required to consider the application a success?
- How to compare the quality of an obtained solution?
- What are the existing alternatives that can serve as a standard of comparison?
- What type of data to use to evaluate the various models?



Questions?

47

47



Data Mining

S4

NOVA-IMS 2019/2020

Fernando Lucas Bação

bacao@isegi.unl.pt

<http://www.isegi.unl.pt/bacao>

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

1



**Data Preparation and
Pre-processing**

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

AACSB Accredited

UNIGIS

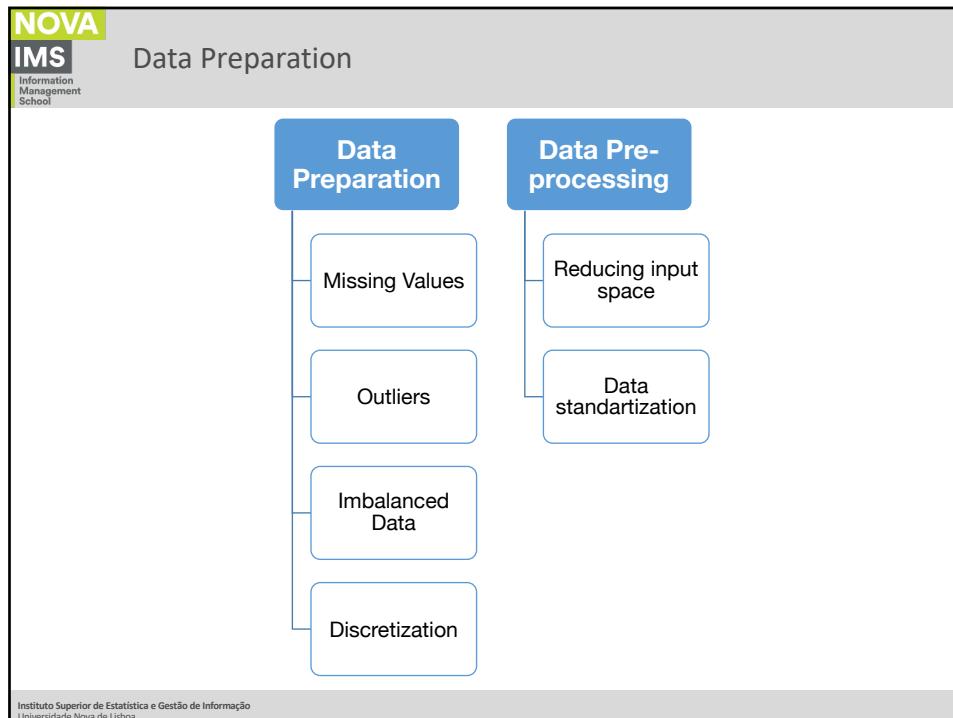
A3ES

EFMD

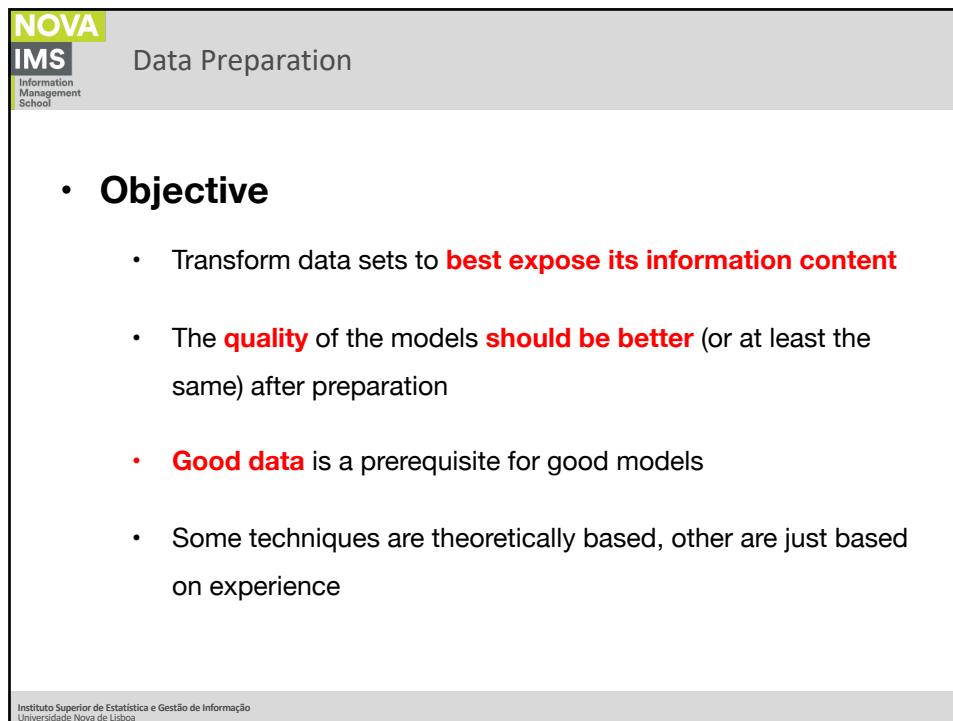
Network of Schools

eduniversal

2



3

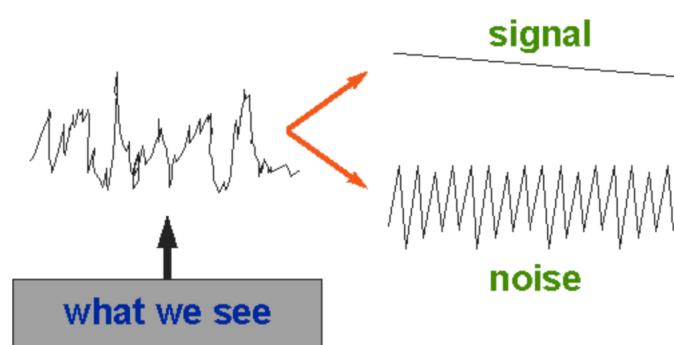


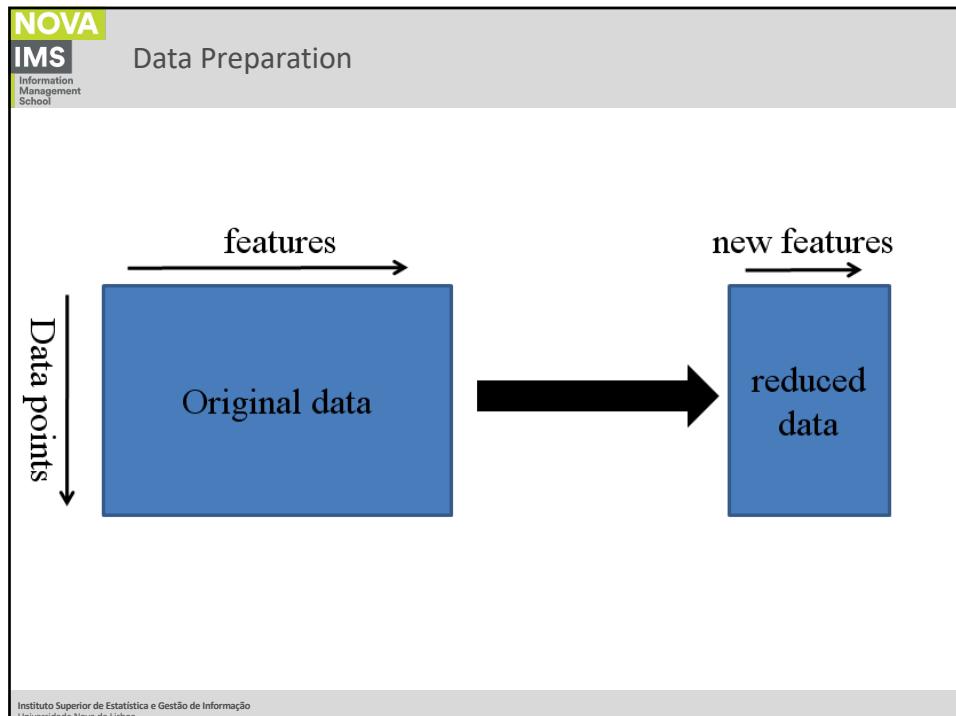
4

- **Signal vs. Noise**

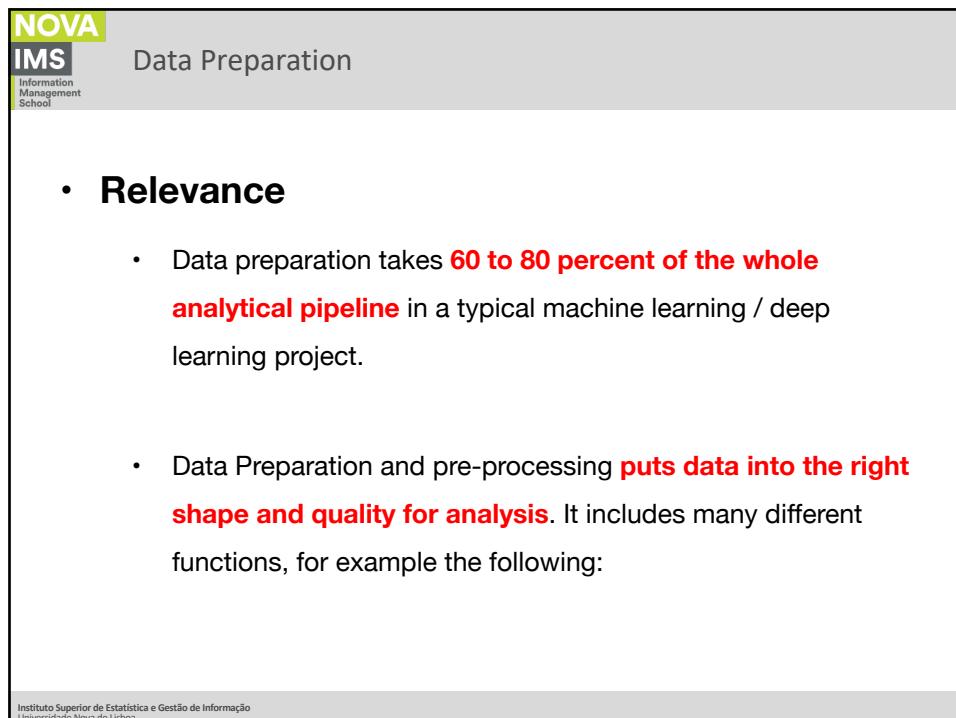
- In science (physics and telecommunications) noise is defined as **fluctuations and external disturbances** in the flow of information (signal) received;
- An **undesired disturbance** in relevant information;
- A disturbance that affects a signal and that may distort the information carried by the signal.

What we observe can be divided into:



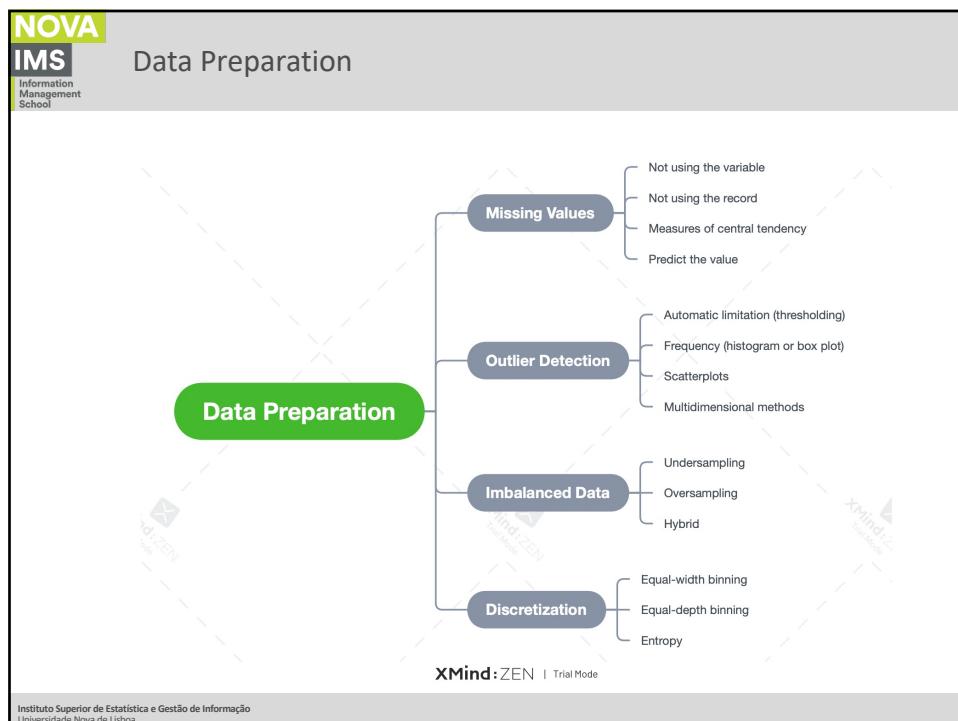


7

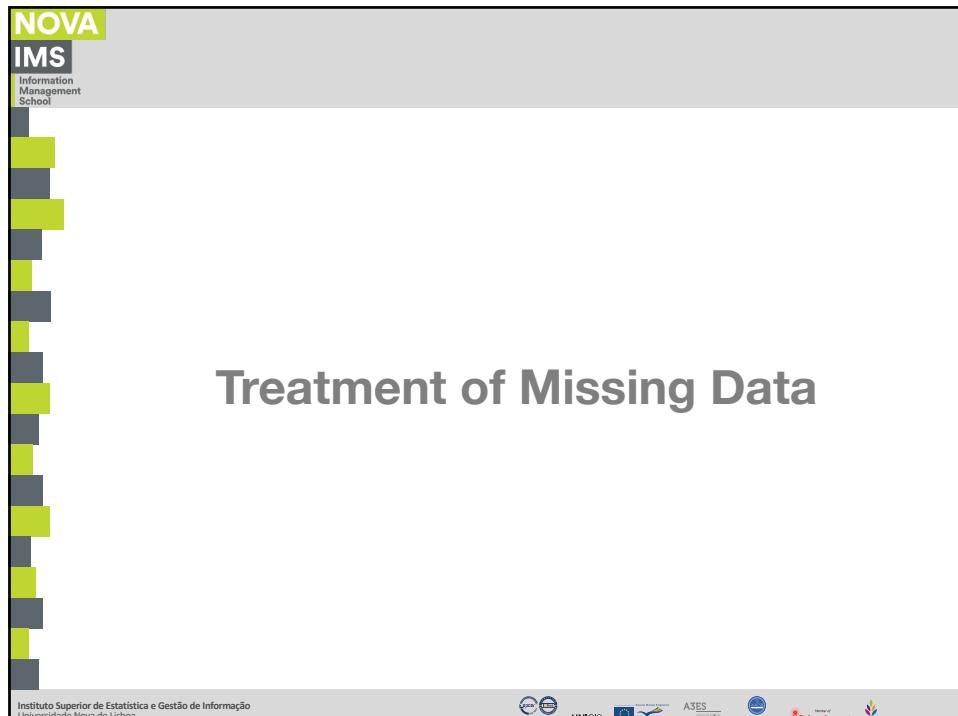


8

- **Real data suffers from several problems**
- Incomplete
 - Missing values, lacking attributes of interest, levels of aggregation
- Noisy
 - Errors and outliers
- Inconsistent
 - E.g. Age=42 Birthday=31/07/1997
 - Changes in scales
 - Duplicate records with different values



11



12

NOVA
IMS
Information Management School

Data Preparation

- Missing Data

Inputs

Inputs				
Records	1	2	3	4
1	?		?	
2		?	?	?
3	?		?	?
4		?		?
5	?		?	

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

13

NOVA
IMS
Information Management School

Data Preparation

- Missing Data

Inputs

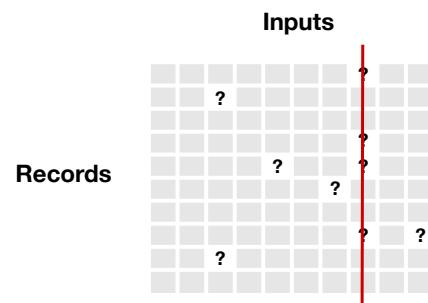
Inputs				
Records	1	2	3	4
1	?		?	
2	?		?	
3	?	?	?	
4	?		?	?
5	?		?	

10 to 3

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

14

- Missing Data



OS	Revenue
Android	1,804
iOS	3,027
iOS	8,788
Android	NA
Android	3,735
Android	1,056
iOS	9,319
Android	6,199
Android	2,235
iOS	NA
Android	1,146

NOVA**IMS**Information
Management
School**Data Preparation**

OS	Revenue
Android	1,804
iOS	3,027
iOS	8,788
Android	NA
Android	3,735
Android	1,056
iOS	9,319
Android	6,199
Android	2,235
iOS	NA
Android	1,146

OS	Global Mean
Android	1,804
iOS	3,027
iOS	8,788
Android	4,145
Android	3,735
Android	1,056
iOS	9,319
Android	6,199
Android	2,235
iOS	4,145
Android	1,146

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

17

NOVA**IMS**Information
Management
School**Data Preparation**

OS	Revenue
Android	1,804
iOS	3,027
iOS	8,788
Android	NA
Android	3,735
Android	1,056
iOS	9,319
Android	6,199
Android	2,235
iOS	NA
Android	1,146

OS	Global Mean	Group Mean
Android	1,804	1,804
iOS	3,027	3,027
iOS	8,788	8,788
Android	4,145	2,696
Android	3,735	3,735
Android	1,056	1,056
iOS	9,319	9,319
Android	6,199	6,199
Android	2,235	2,235
iOS	4,145	7,045
Android	1,146	1,146

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

18

NOVA
IMS
Information Management School

Data Preparation

Missing value record

Other dataset records

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

19

NOVA
IMS
Information Management School

Data Preparation

CUSTID	DAYSWU	AGE	EDU	INCOME
1138	951	29	17	55637
6313	1231	22	16	32795
9185	1186	31	18	40717
7735	1097	27	18	29184
7918	1196	19	13	16479
1585	849	34	19	55403
7897	1120	46	15	63603
7389	792	39	16	57402
6764	1060	missing	15	69803
8416	1217	31	18	56028
6541	1074	54	15	91934
8263	1185	46	18	78262
8052	1007	52	16	91292
2448	840	51	14	62699
10978	553	44	19	62675
1409	1191	missing	14	64449
6063	1222	32	18	55048
9767	961	34	20	53217
4489	1224	69	16	103191
10738	1190	missing	18	64412
10381	939	54	16	73709
8999	1054	45	15	73074
4482	1089	50	19	63983
4257	1074	missing	20	75006
7722	876	50	19	55332
7491	1165	49	14	63150
1037	1127	47	16	69713
10419	1112	48	15	74961
5124	1217	47	17	72890
8140	1229	54	16	74670
3218	1159	54	15	83956
4124	1247	53	15	81937
8990	1025	54	17	79364
7155	1235	54	18	81948
6346	1170	42	17	71182
9547	1125	missing	16	81442
8250	989	12	16	75566
2527	1234	51	15	77635
6663	1185	40	14	61974

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

20

- **Missing data**

- Delete variables (you loose information)
- Delete records (potential bias);
- To examine and manually enter a probable value (tedious + infeasible);
- Automatically fill in with a measure of central tendency (i.e. mean, median, mode);
- Automatically fill in with a measure of central tendency of a subset (e.g. men and women);
- To fill in with values from similar individuals (nearest neighbours);
- Predictive model (linear regression, multiple linear regression);
- Code the missing data explicitly.

- **Missing data**

- The most practical approach to the problem is to initially use the **quickest and simplest option**;
- After achieving some **preliminary results** we can comparatively analyze the performance of the model in the full sample patterns and in those where there was a need to estimate missing values;
- In the event that the **error is significantly higher** than in other data, then we will try to use another method in order to improve results.



Outlier treatment

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



23



Data Preparation

- **Outliers**
 - In statistics, an outlier is an observation point that is distant from other observations.
 - An outlier may be due to variability in the measurement or it may indicate experimental error; the latter are sometimes excluded from the data set.
 - Extreme cases in one or more variables and with great impact on the interpretation of results;
 - Outliers may come from:
 - Unusual but correct situations (the Bill Gates effect),
 - Incorrect measurements,
 - Errors in data collection;
 - Lack of code for missing data.

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

24

**NOVA
IMS**
Information Management School

Data Preparation

- Outliers (leverage effect)**

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

25

**NOVA
IMS**
Information Management School

Data Preparation

- Remove Data Outliers**
 - Automatic limitation (thresholding)
 - The imposition of maximum and minimum values for the variables (age – 0 e 100)

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

26

**NOVA
IMS**
Information Management School

Data Preparation

- Remove Data Outliers

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

27

**NOVA
IMS**
Information Management School

Data Preparation

Q1: Quartile 1, or median of the *left* data subset
after dividing the original data set into 2 subsets via the median
(25% of the data points fall below this threshold)

Q3: Quartile 3, median of the *right* data subset
(75% of the data points fall below this threshold)

IQR: Interquartile-range, $Q_3 - Q_1$

Outliers: Data points are considered to be outliers if
value < $Q_1 - 1.5 \times IQR$ or
value > $Q_3 + 1.5 \times IQR$

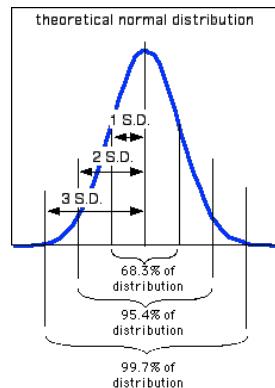
fmi:
http://sebastianraschka.com/Articles/2014_dixon_test.html
http://www.itl.nist.gov/div898/handbook/prc/section1/prc1_6.htm

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

28

- Remove Data Outliers

- Normal data distribution
 - $3\sigma +/\text{-} \text{Average}$

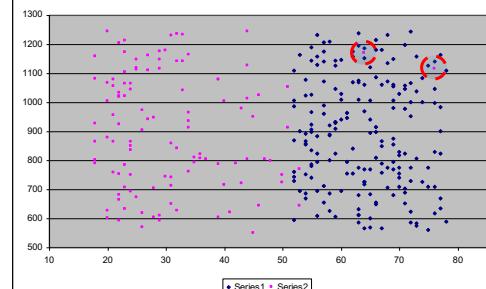
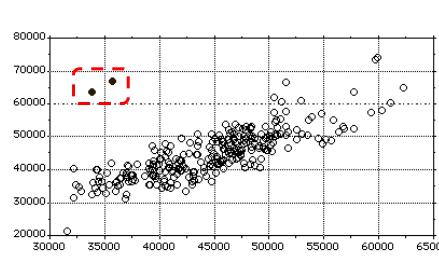


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

29

- Remove Data Outliers

- In two dimensions...



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

30

- Remove Data Outliers
 - Cluster Analysis (K-means)
 - Self-Organizing Maps

Imbalanced Learning

Quick quizz:

	Yes	No
Is it possible to have a useless classifier with 99% accuracy?		
Is it possible to achieve a 99.9% accuracy with a trivial classifier?		

Quick quizz:

	Yes	No
Is it possible to have a useless classifier with 99% accuracy?	x	
Is it possible to achieve a 99.9% accuracy with a trivial classifier?	x	

Quick quizz:

	Yes	No
Is it possible to have a useless classifier with 99% accuracy? If the minority class is 1%	X	
Is it possible to achieve a 99.9% accuracy with a trivial classifier? If the minority class is 0,1%	X	

- Class Imbalance
 - Is this frequent in real-world applications?
 - **Credit Card** frauds - ~2% per year.
 - **HIV prevalence** in the USA - ~0.4%.
 - **Disk drive** failures - ~1% per year.
 - **Factory production** defects - ~ 0.1%.
 - **Business churn** - ~3%

- Imbalanced Learning

- An **imbalanced learning** problem is defined as a classification task for binary or multi-class datasets where a **significant asymmetry** exists between the **number of instances for the various classes**.
- The dominant class is called the **majority class (negative cases)** while the rest of the classes are called the **minority classes (positive cases)**
- The **Imbalance Ratio** (IR), is the ratio between the majority class and the minority class, (depends on the type of application and for binary problems values between 100 and 100.000 have been observed)

- Imbalanced Learning

- **Standard learning methods induce a bias** in favor of the majority class during training.
- This happens because the **minority classes contribute less to the maximization** of the objective function, which is usually accuracy.

$$\text{Accuracy} = \frac{\text{TruePositive} + \text{TrueNegative}}{\text{TotalPopulation}}$$

Imbalanced Learning

- Imbalanced Learning
 - Standard learning methods induce a bias in favor of the majority class during training.
 - This happens because the minority classes contribute less to the maximization of the objective function, which is usually accuracy.

		Actual Value	
		Positives	Negatives
Predicted Value	Positives	TP True Positives	FP False Positives
	Negatives	FN False Negatives	TN True Negatives

Imbalanced Learning

- Imbalanced Learning
 - Standard learning methods induce a bias in favor of the majority class during training.
 - This happens because the minority classes contribute less to the maximization of the objective function, which is usually accuracy.

		Actual Value	
		Positives	Negatives
Predicted Value	Positives	TP True Positives	FP False Positives
	Negatives	FN False Negatives	TN True Negatives

**NOVA
IMS**
Information Management School

Imbalanced Learning

- Imbalanced Learning
 - By optimizing classification accuracy, **most algorithms assume a balanced class distribution**

The figure consists of two side-by-side pie charts. The left pie chart, titled 'Negative' and 'Positive', shows a 50/50 split between the two classes. The right pie chart, also titled 'Negative' and 'Positive', shows a 99% majority for the negative class and a 1% minority for the positive class.

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

A3ES UNIGIS UNI-Schools eduniversal

41

**NOVA
IMS**
Information Management School

Approaches to the Imbalanced Learning Problem

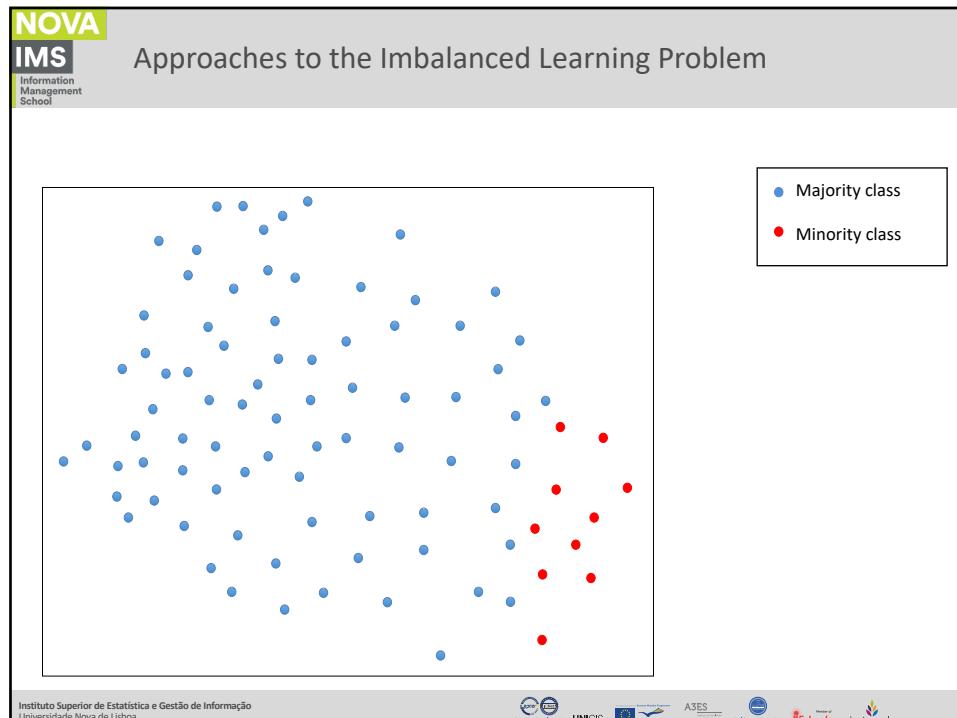
```

graph TD
    A[Solutions to Imbalanced Learning] --> B[Undersampling]
    A --> C[Oversampling]
    A --> D[Hybrid approaches]
  
```

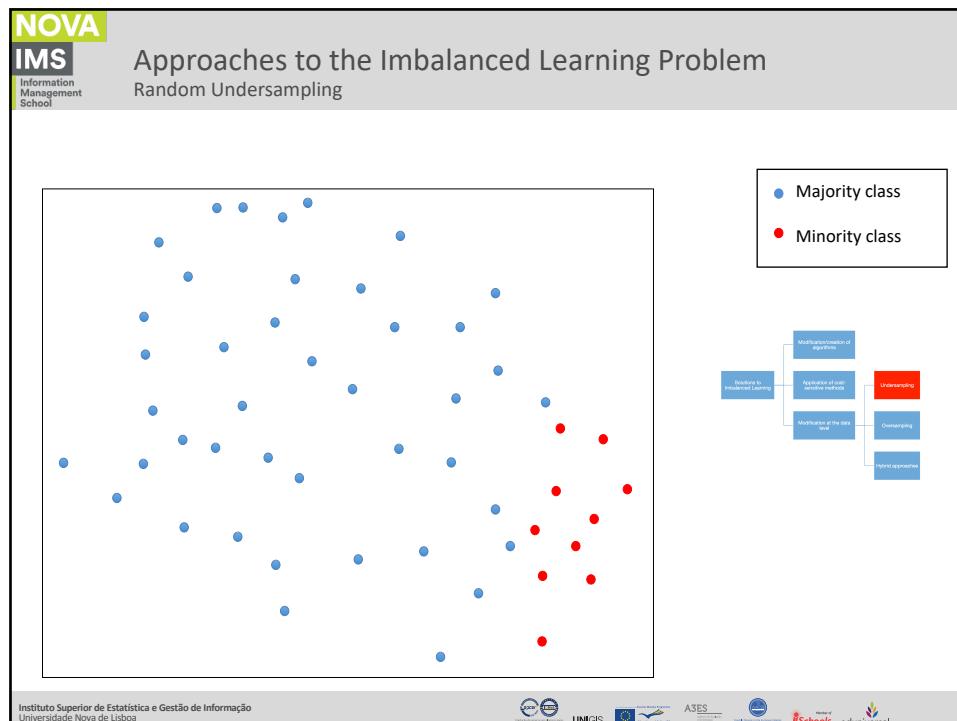
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

A3ES UNIGIS UNI-Schools eduniversal

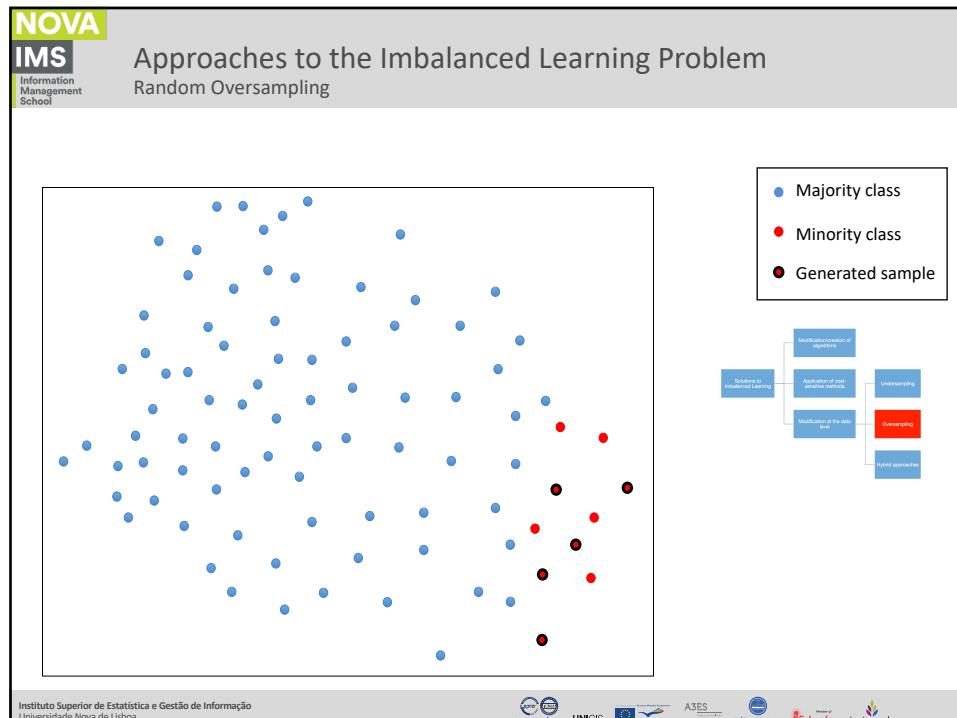
42



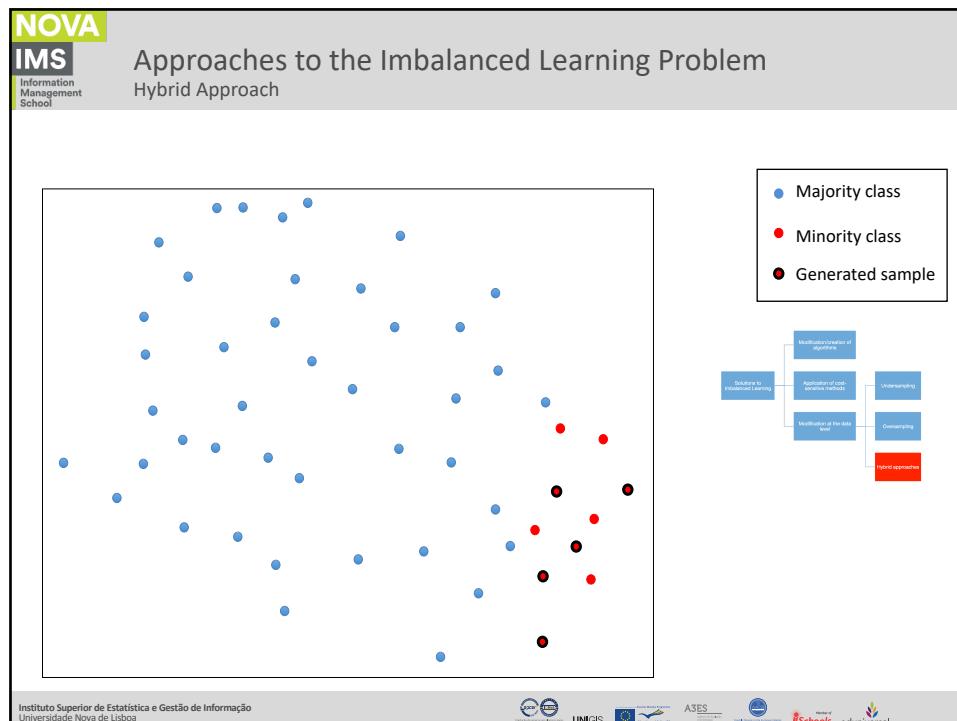
43



44



45



46

NOVA
IMS
Information Management School



SMOTE: Synthetic Minority Over-sampling TEchnique

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

47

Accreditation Logos: EQUIS, UNIGIS, A3ES, AACSB, iSchools, eduniversal

47

NOVA
IMS
Information Management School

Imbalanced Learning

- SMOTE
 - The idea underlying SMOTE is **as simple as it is clever**.
 - The **basic steps** are:
 - randomly selecting a minority class instance x ;
 - then it defines the set of k-nearest neighbors (x_{knn});
 - randomly selects another minority class sample x' from the x_{knn} set.
 - x_{gen} is generated by using a linear interpolation of x and x' , which can be expressed as:

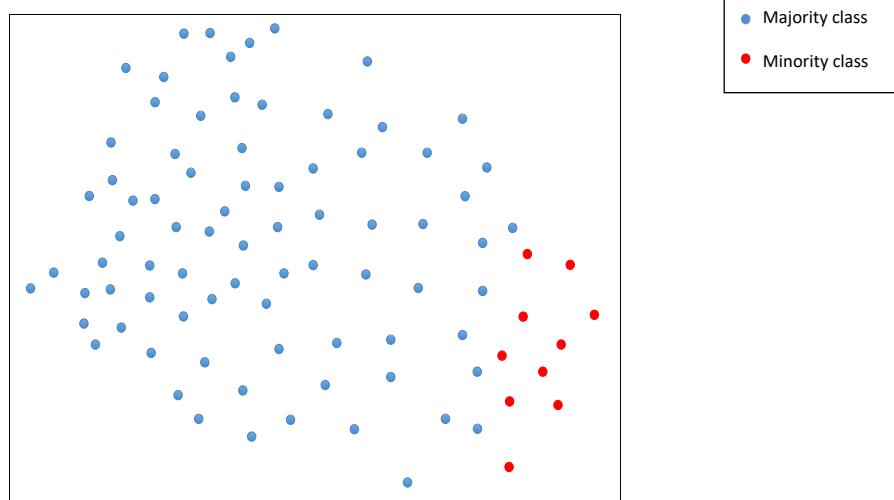
$$x_{gen} = x + a \cdot (x' - x)$$

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

Accreditation Logos: EQUIS, UNIGIS, A3ES, AACSB, iSchools, eduniversal

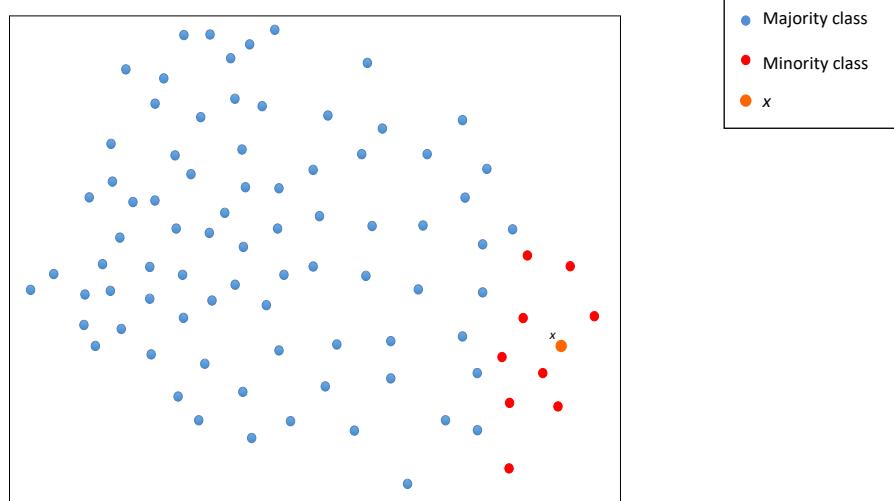
48

General aspects of data collection – acquiring knowledge

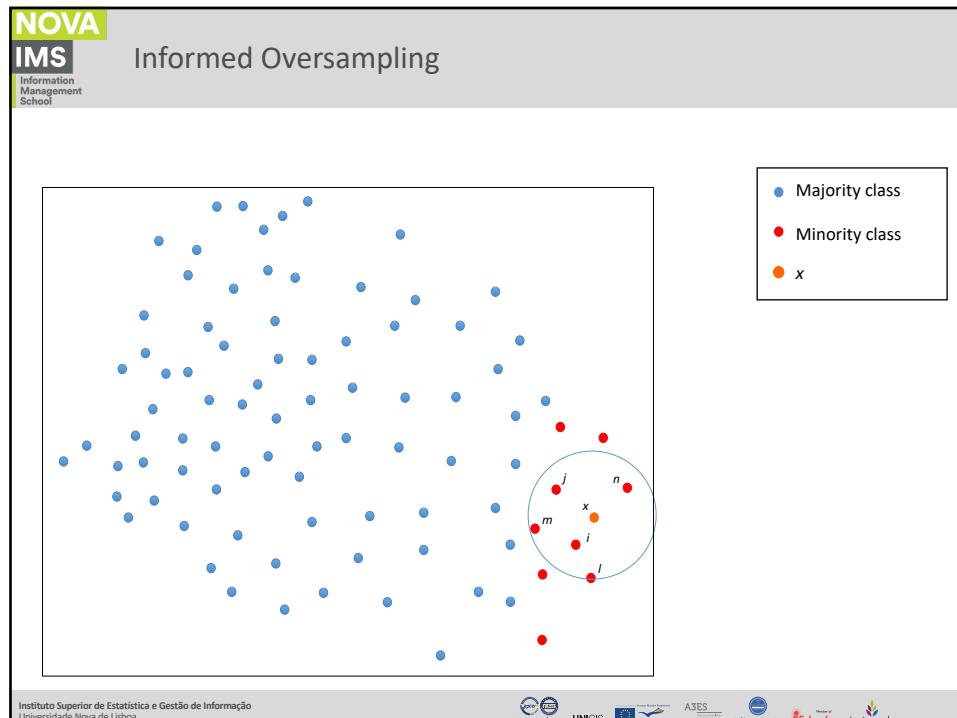
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

49

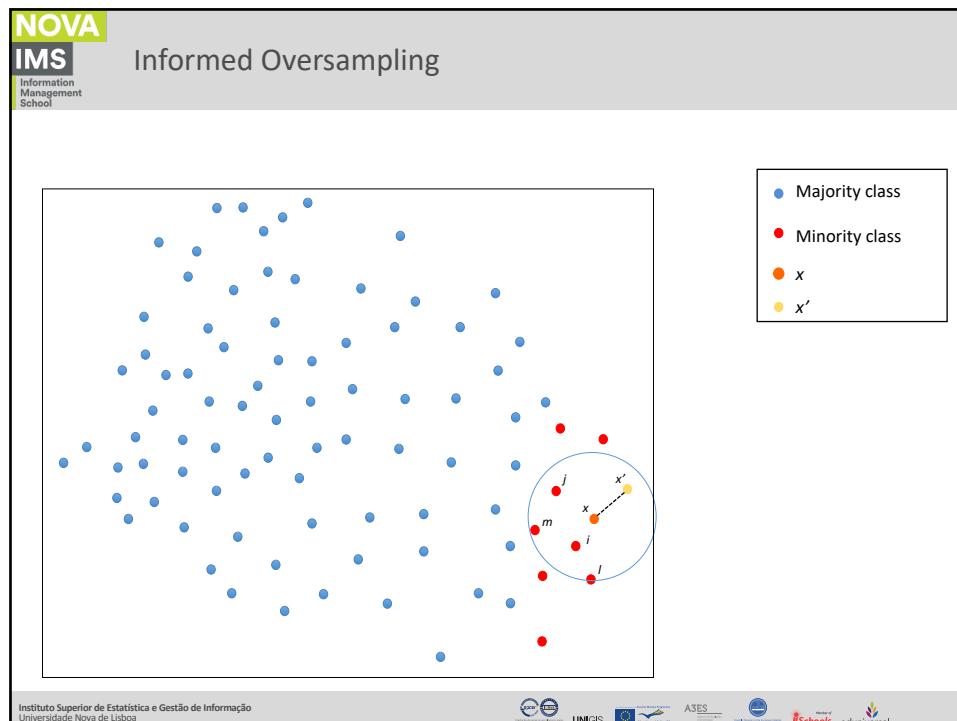
Informed Oversampling

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

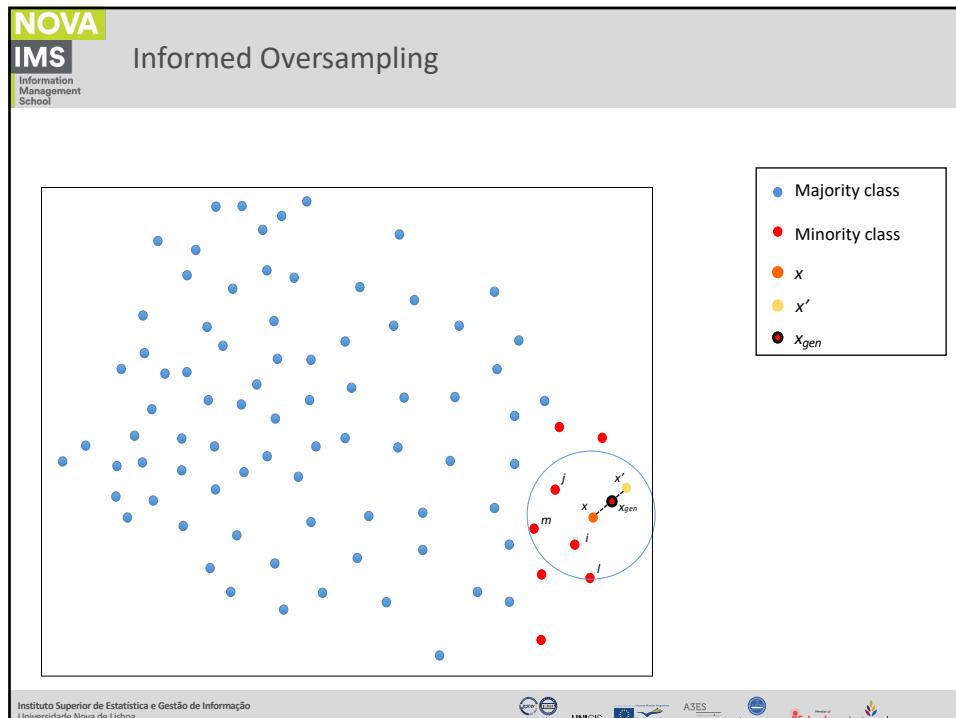
50



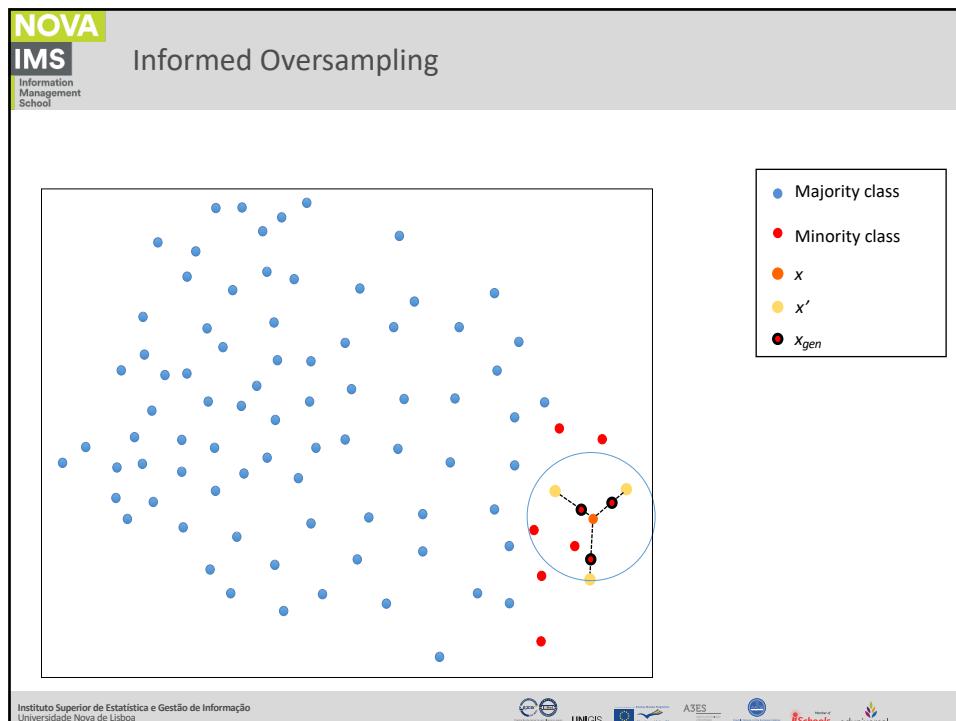
51



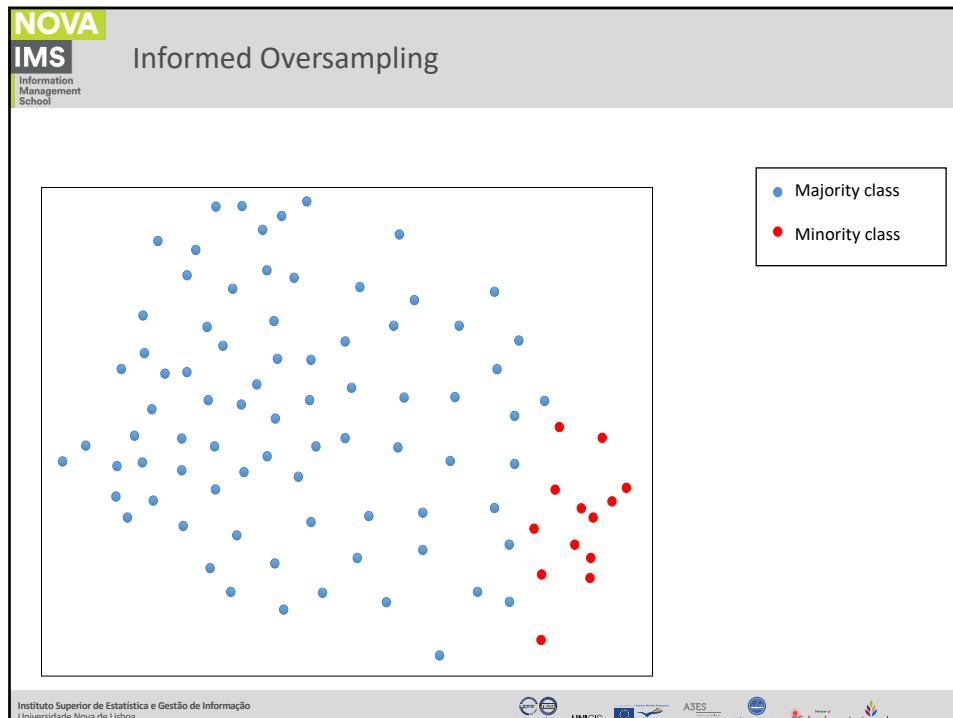
52



53



54



55

NOVA
IMS
Information Management School

General aspects of data collection

- Use of artificial data:
 - It is always preferable to use real data;
 - Create data as realistic as possible;
 - Make artificial data as representative as possible.
 - The quality of the model is constrained by the quality of the data;
 - Creating artificial data translates into the introduction of some noise.

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

56

NOVA
IMS
Information Management School



Discretization

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



57

NOVA
IMS
Information Management School

Data Preprocessing

- **Discretization**
 - Divide the range of a continuous variable into intervals
 - Some classification algorithms only accept discrete attributes
 - Reduce data size
 - Prepare for further analysis
 - Frequently called binning

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

58



Discretization

Unsupervised

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

59



Data Preprocessing

- **Discretization**
 - Unsupervised binning methods transform numerical variables into categorical counterparts but do not use the target (class) information.
 - Equal Width
 - Equal Depth (or frequency)
 - Other methods such as quantiles

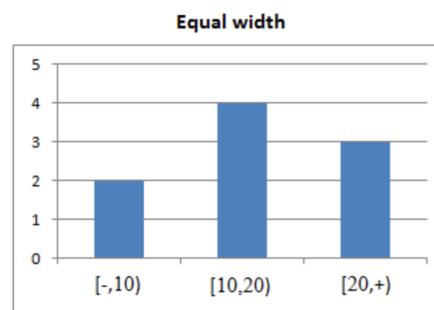
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

60

- **Discretization**

- Equal-width binning

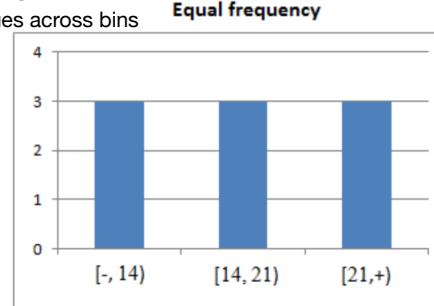
- Divides the range into n intervals of equal size
- If A and B are the minimum and the maximum values of the attribute, the width of the intervals will be: $w=(B-A)/N$
- Most simple method
- Outliers may dominate



- **Discretization**

- Equal-depth binning

- Divides the range into n intervals, each containing approximately the same number of samples
- Generally preferred avoids clumps
- Gives more intuitive breakpoints
- Shouldn't break frequent values across bins



- **Discretization**

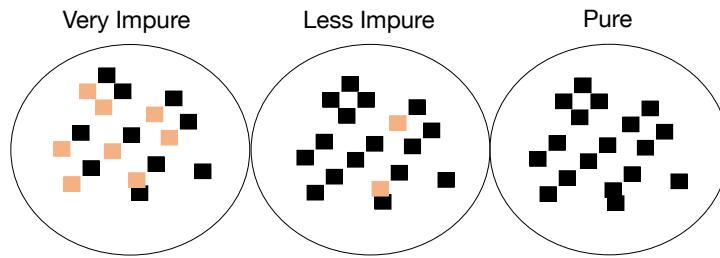
- Class-independent methods (unsupervised)
 - Equal Width is simpler, good for many classes
 - can fail miserably for unequal distributions
 - Equal Height gives better results
- Class-dependent methods can be better for classification
 - Decision tree methods build discretization on the fly
 - Naïve Bayes requires initial discretization
- Many other methods exist ...

Discretization

Supervised

- **Discretization**

- Entropy (also called Expected Information) based discretization

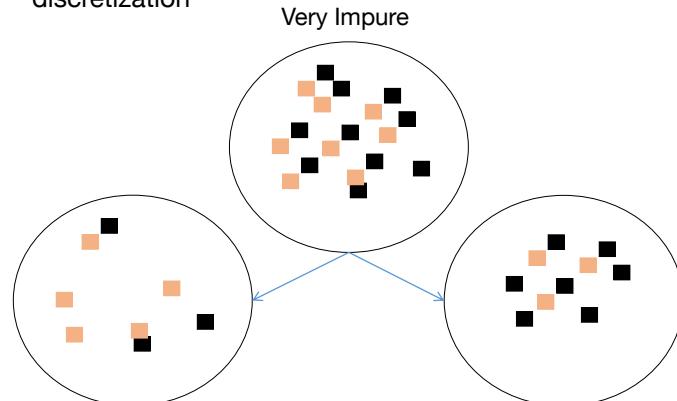


- **Discretization**

- Entropy (also called Expected Information) based discretization
 - Sort examples in increased order
 - Each value forms an interval (m intervals)
 - Calculate the entropy measure of each discretization
 - Find the binary split boundary that minimizes the entropy function over all possible partitions. The split is selected as a binary discretization
 - Apply the process recursively until some stopping criteria is met

- **Discretization**

- Entropy (also called Expected Information) based discretization



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

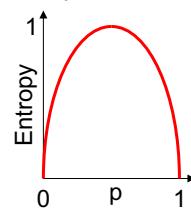
67

- **Discretization**

- Entropy based discretization
 - Entropy
 - Idea: maximize info
 - It measures the purity of a partition:

$$E = -p \log_2(p)$$

- Where p is the probability of the examples belong to a specific class



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

68

- **Discretization**

- Entropy based discretization

	Income <= 50K	Income > 50K
Age < 25	4	6
	Income <= 50K	Income > 50K
Age < 25	9	1

- **Discretization**

- Entropy based discretization

- **Partition entropy:**

$$Ent(S) = - \sum_{i=1}^{\#C} p_i \log_2(p_i)$$

- **Gain in choosing A attribute:**

$$Gain(Ent_{new}) = Ent_{initial} - Ent_{new}$$

$$Gain(S, A) = Ent(S) - \sum_{v \in Valores(A)} \frac{\#S_v}{\#S} Ent(S_v)$$

- **Discretization**

- Entropy based discretization

	Income <= 50K	Income > 50K
	13	7

$$Ent(S) = -(13/20 \log_2(13/20) + 7/20 \log_2(7/20)) = 0.403 + 0.530 = 0.934$$

- **Discretization**

- Entropy based discretization

	Income <= 50K	Income > 50K
Age < 25	4	6

$$Ent(Age < 25) = -(4/10 \log_2(4/10) + 6/10 \log_2(6/10)) = 0.529 + 0.442 = 0.971$$

	Income <= 50K	Income > 50K
Age \geq 25	9	1

$$Ent(Age \geq 25) = -(9/10 \log_2(9/10) + 1/10 \log_2(1/10)) = 0.137 + 0.332 = 0.469$$

- **Discretization**

- Entropy based discretization

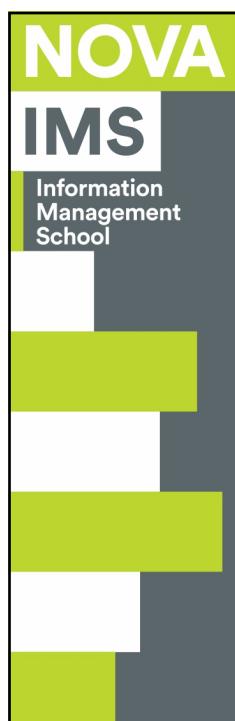
	Income <= 50K	Income > 50K
Age < 25	9	1
Age ≥ 25	4	6

$$\sum_{v \in Valores(A)} \frac{\#S_v}{\#S} Ent(S_v) = \frac{1}{2}(0.469) + \frac{1}{2}(0.971) = 0.72$$

$$Gain(S, A) = Ent(S) - \sum_{v \in Valores(A)} \frac{\#S_v}{\#S} Ent(S_v)$$

$$Ent(S) = -(13/20 \log_2(13/20) + 7/20 \log_2(7/20)) = 0.403 + 0.530 = 0.934$$

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



Data Mining

211020

NOVA-IMS 2019/2020

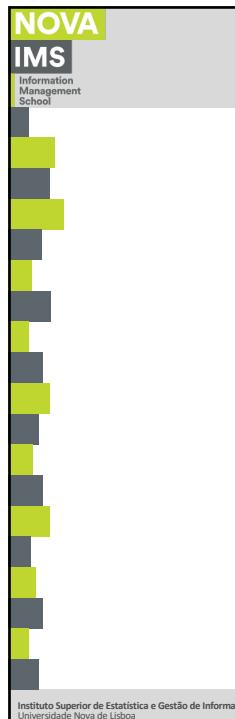
Fernando Lucas Bação

bacao@isegi.unl.pt

<http://www.isegi.unl.pt/bacao>

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

1



Data Pre-processing

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

AACSB Accredited

UNIGIS

A3ES

EFMD

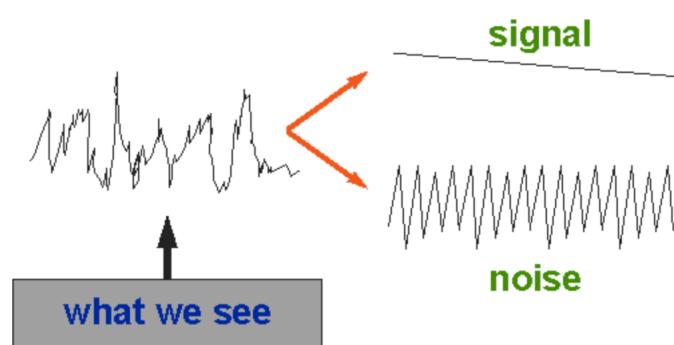
Business Schools

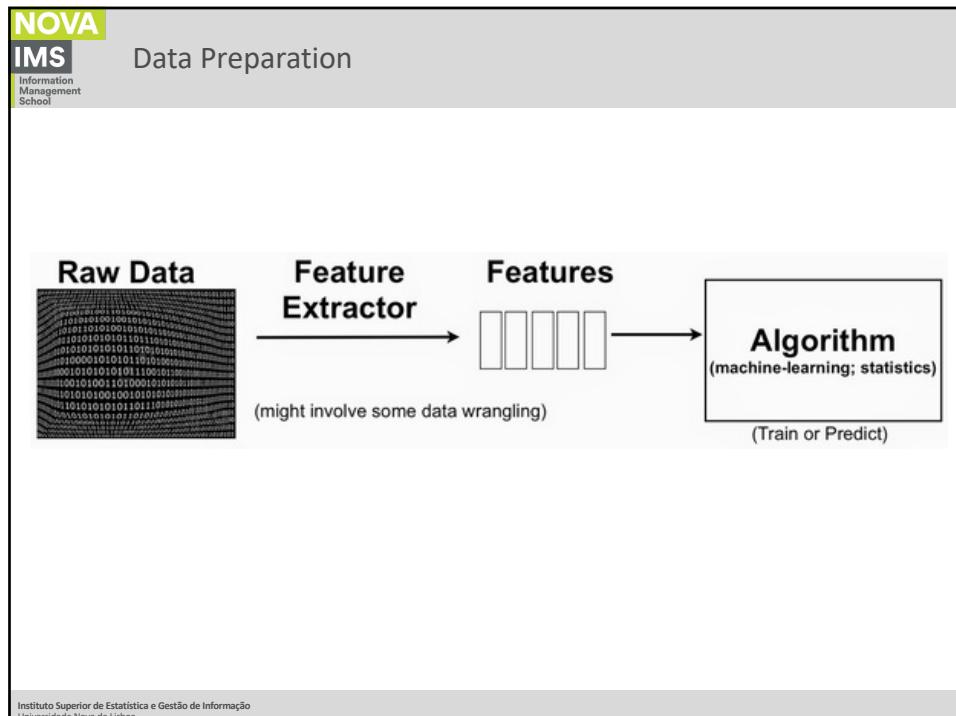
eduniversal

2

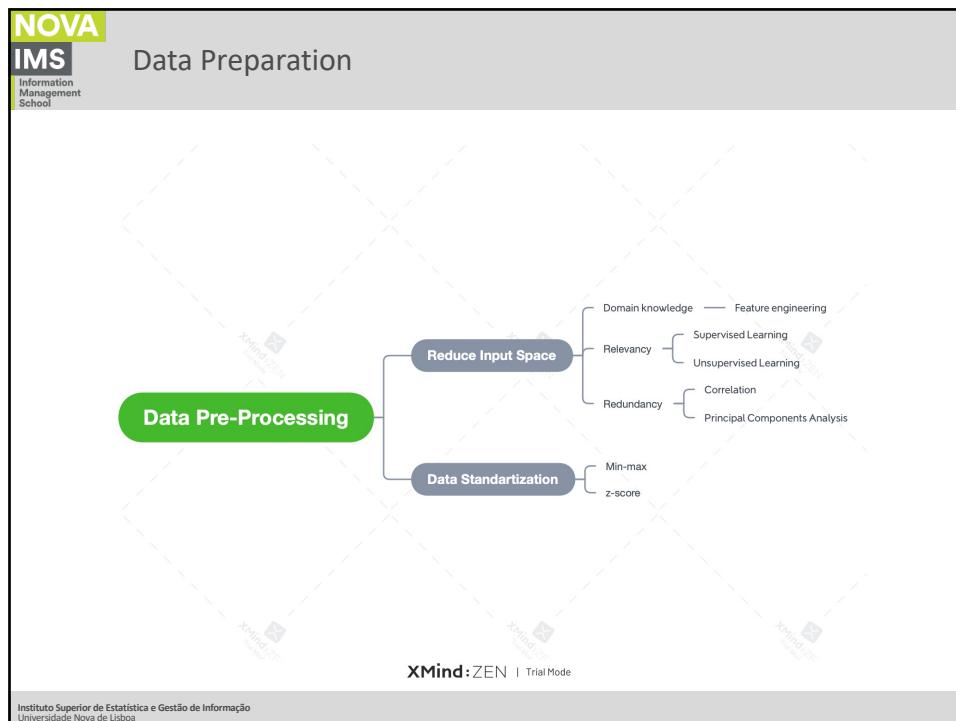
- **Reasons:**
 - Noise Reduction;
 - Signal amplification;
- **Tasks:**
 - Domain-specific knowledge application;
 - Constructing ratios and derived variables
 - Size Reduction of the Input Space;
 - Remove correlated variables
 - Remove irrelevant variables
 - Normalization;

What we observe can be divided into:





5



6



Reducing Input Space

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



7



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.

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

8

NOVA
IMS
Information Management School

Data Preprocessing

Three groups, right?

The curse of dimensionality

Not exactly...

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

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

9

NOVA
IMS
Information Management School

Data Preprocessing

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

E(Target)

$Input_1$ $Input_2$

$Input_1$

$Input_5$

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

10

NOVA
IMS
Information Management School



Reducing Input Space

Feature Engineering

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



11

NOVA
IMS
Information Management School

Data Preprocessing

- **Size Reduction of the Input Space:**
 - To create input combinations
 - Height²/weight (obesity index)
 - Population/area (density)
 - Euros spent/nº of purchases (average buy)
 - Euros spent/time as customer
 - Debt/income
 - Average number of different products purchased per transaction
 - Relative spend on each product

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

12

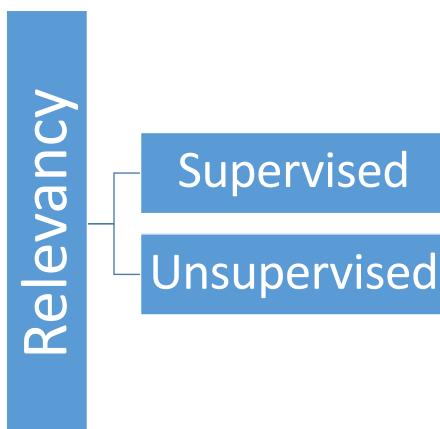
- **Size Reduction of the Input Space:**

- To create input combinations
 1. Average time between transactions (transaction interval)
 2. Variance of transaction interval
 3. Customer stability index (ratio of (2)/(1))

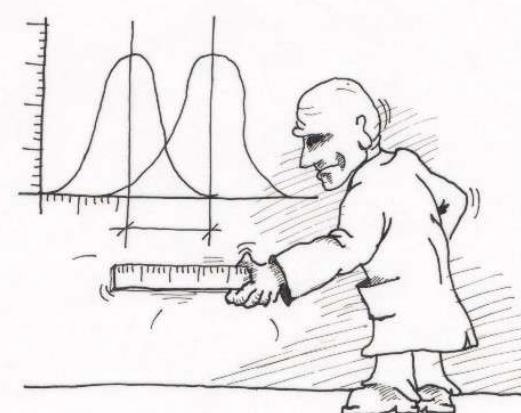


Reducing Input Space

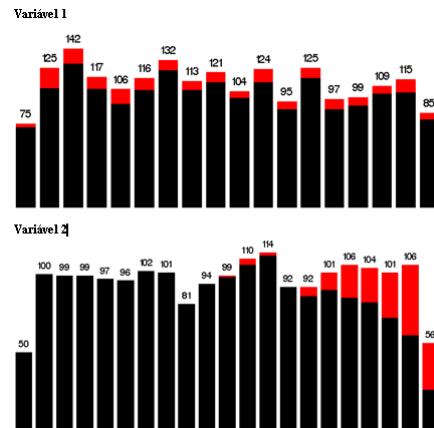
Relevancy



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



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



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

17

- **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



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

18

NOVA
IMS
Information Management School

Reducing Input Space

Redundancy

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

A3ES UNIGIS A Schools eduniversal

19

NOVA
IMS
Information Management School

Data Preprocessing

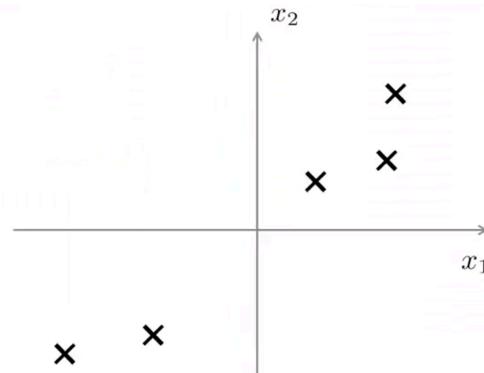
- Size Reduction of the Input Space:**

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

20

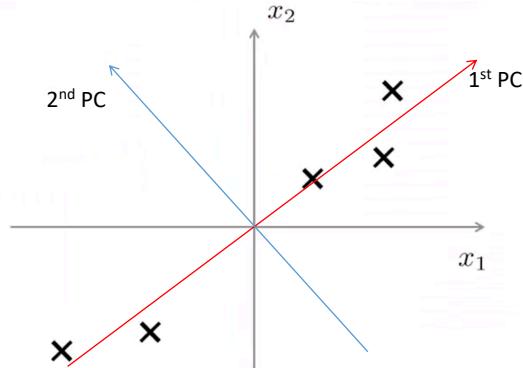
- **Size Reduction of the Input Space:**
 - Principal Component Analysis
 - A procedure that uses an **orthogonal transformation** to convert a set of observations of possibly correlated variables into a set of values of **linearly uncorrelated variables called principal components**.
 - The number of 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.

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



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

- Principal Component Analysis

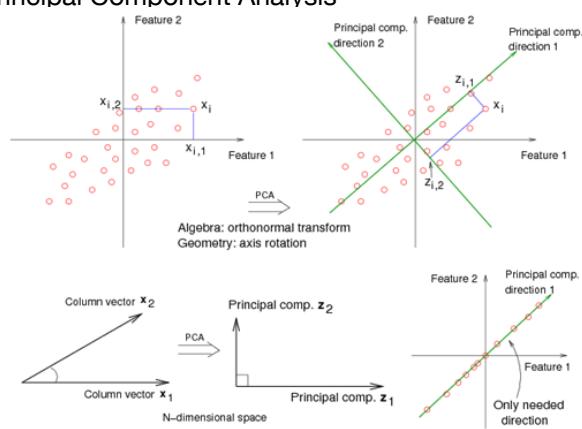


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

23

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

- Principal Component Analysis

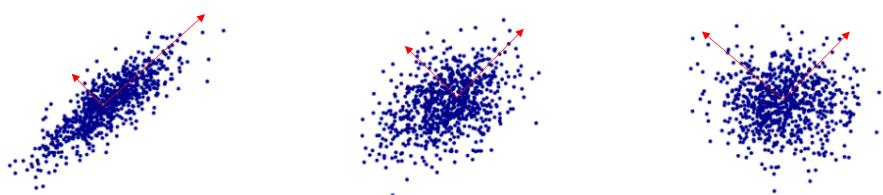


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

24

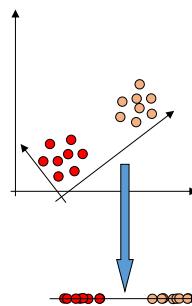
- **Size Reduction of the Input Space:**

- Principal Component Analysis



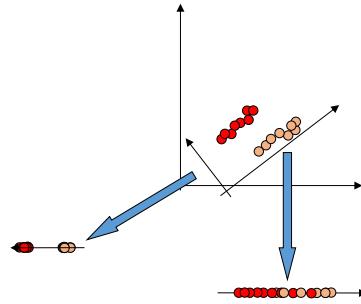
- **Size Reduction of the Input Space:**

- Principal Component Analysis (careful)



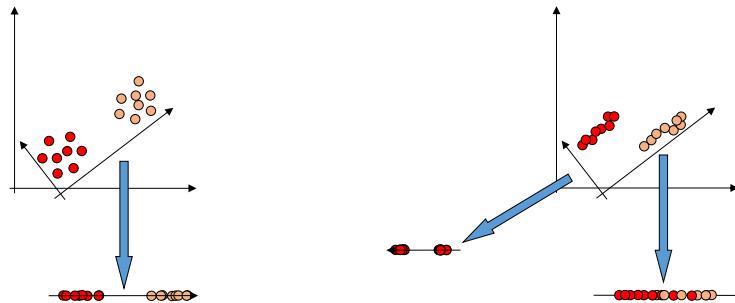
- **Size Reduction of the Input Space:**

- Principal Component Analysis (careful)



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

- Principal Component Analysis (careful)



NOVA
IMS
Information Management School



Data Standardization

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

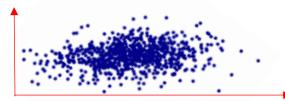


29

NOVA
IMS
Information Management School

Data Preprocessing

- **Normalization:**
 - Models assume that the distances in different directions of the input space have the same importance.
 - Variables come in many **different scales** (percentages, euros, kilos, meters, days...)
 - Normalization: is about adjusting values measured on different scales to a common scale



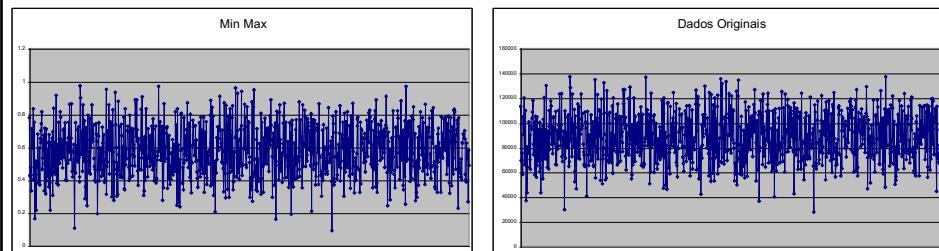
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

30

- **Normalization:**

- Min-Max $y' = \left(\frac{y - \min 1}{\max 1 - \min 1} \right) \underbrace{(\max 2 - \min 2)}_{\text{optional}} + \min 2$
- Zscore $y' = \frac{y - \mu}{std}$

- **Normalization:**

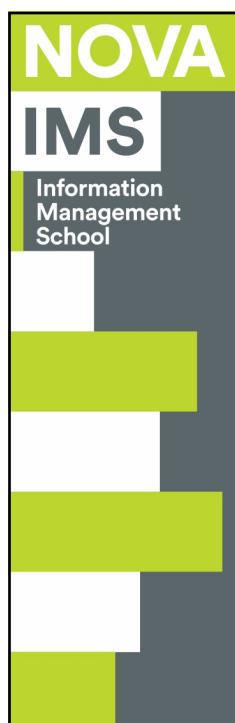




Questions?

33

33



Data Mining

NOVA-IMS
02/11/2021
Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>

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

1



Clustering

Agenda

- Cluster analysis
- Variables to use
- Similarity criterion
- Clustering algorithms
 - A Priori Grouping
 - RFM Analysis

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

2



Cluster Analysis

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



3



Clustering

- **Cluster Analysis**
 - Cluster Analysis is a **basic conceptual activity of human beings**;
 - A **fundamental process**, common to many sciences, essential to the development of scientific theories;
 - The possibility of **reducing the infinite complexity of real** to sets of objects or similar phenomena, is one of the most powerful tools in the service of mankind.

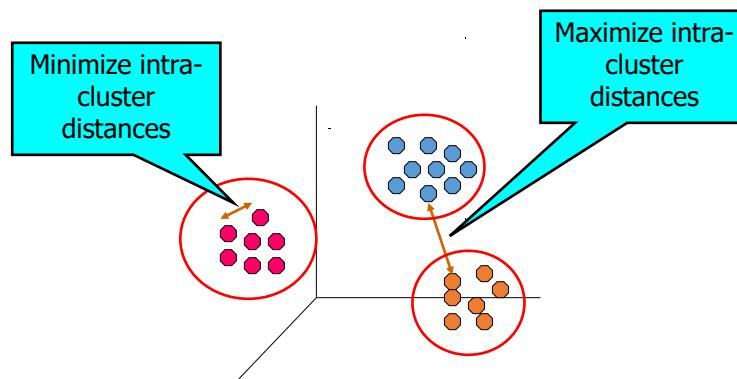
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

4

- **Cluster Analysis**

- Cluster analysis is a generic name for a variety of methods that are used to **group entities**;
- Objective: **To form groups of objects that are similar to each other**;
- From a data collection about a group of entities, seeks to organize them **in homogeneous groups**, assessing a "frame" of similarities/differences between units.

- **Cluster Analysis**



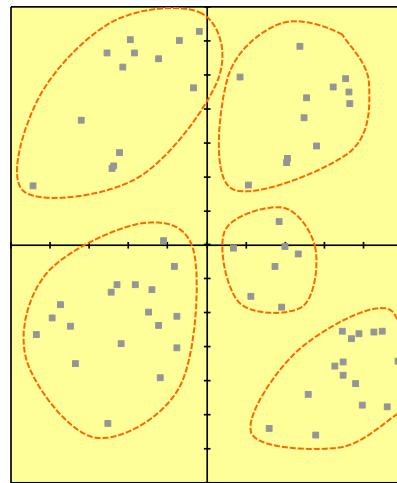
- **Cluster Analysis**

- Classification:
 - Starts out with a **pre-classified training set**, that is, the method has a set of data which contains not only the variables to use in classification but also the class to which each of the records belongs;
 - Attempts to develop a model capable of predicting how a new record will be classified.

- **Cluster Analysis**

- Clustering:
 - There is **no pre-classified data**;
 - We search for groups of records (clusters) that are similar to one another;
 - Underlying is the expectation that similar customers in terms of the variables used will behave in similar ways.

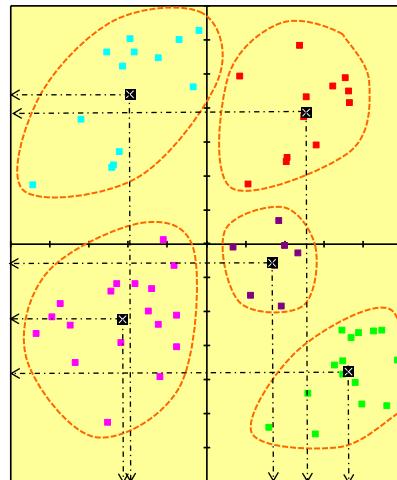
Clustering



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

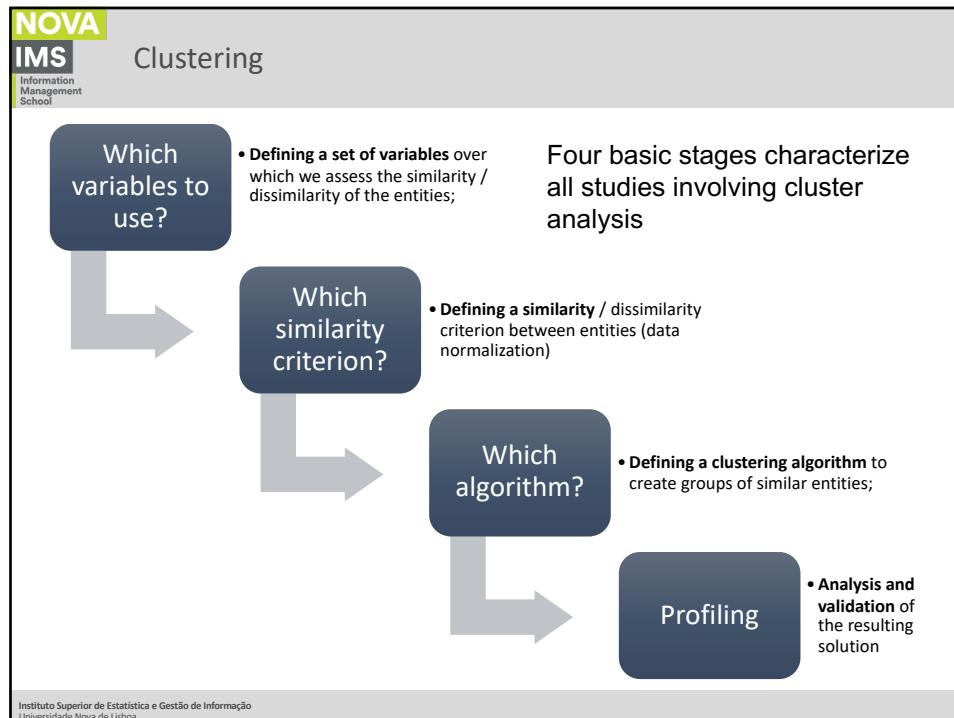
9

Clustering

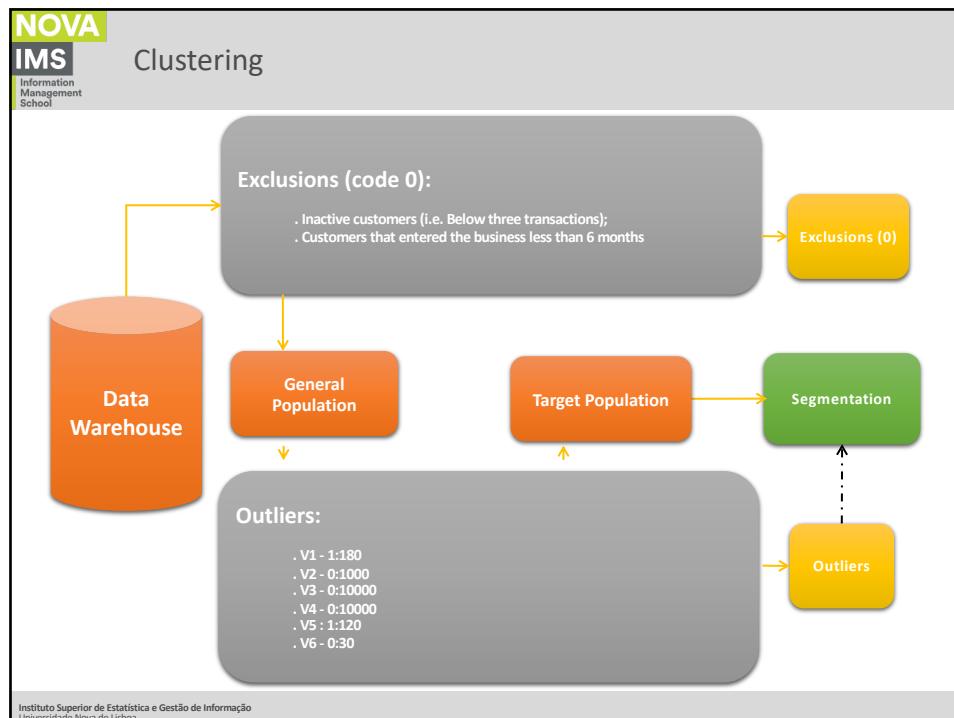


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

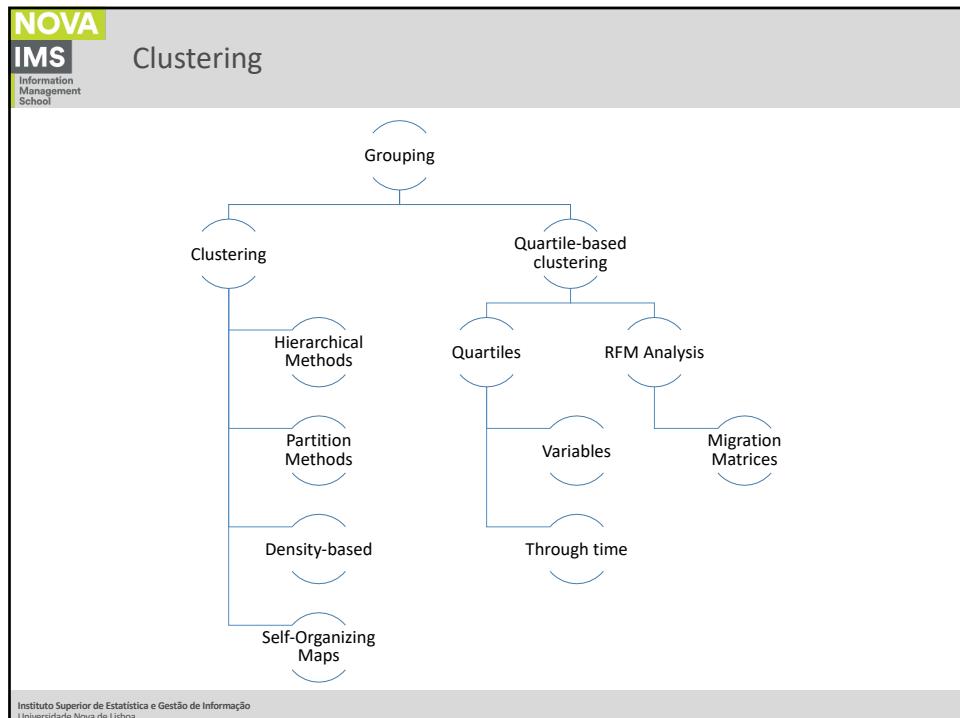
10



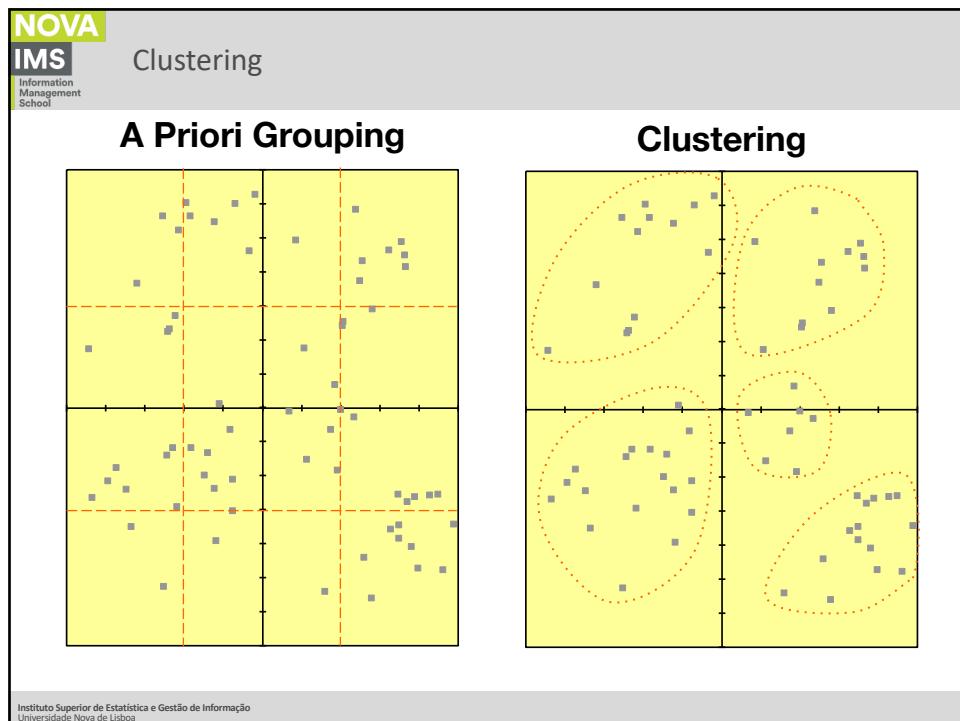
11



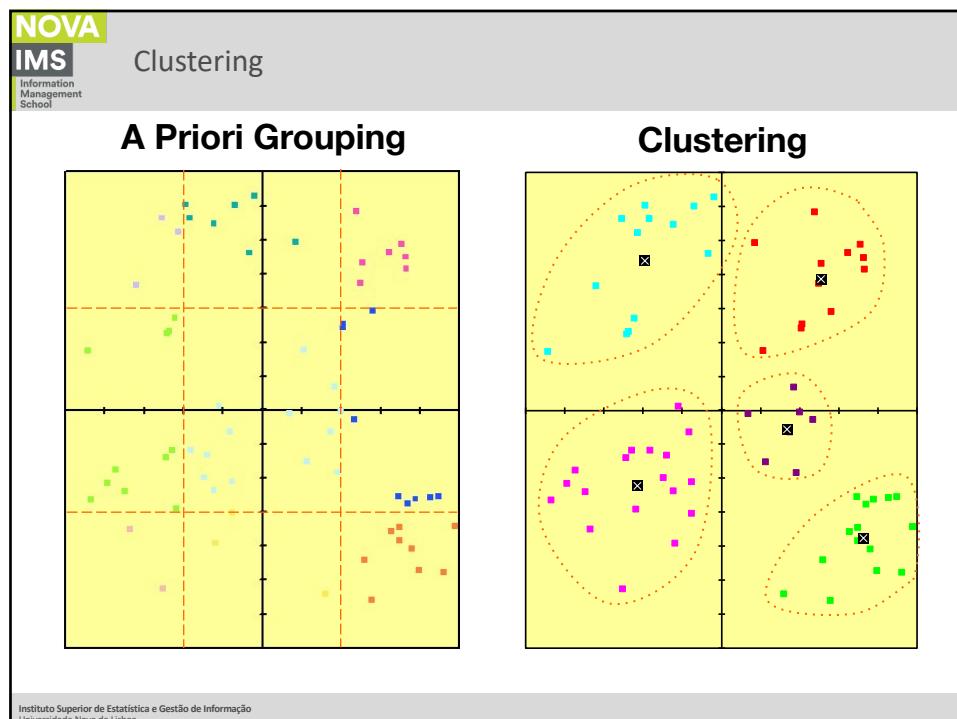
12



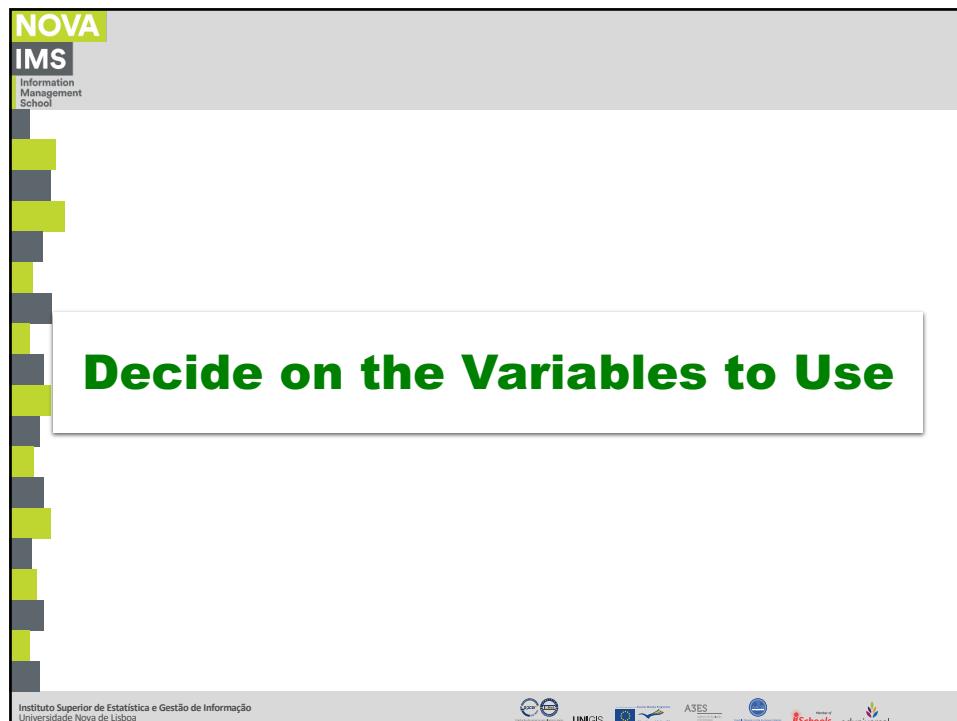
13



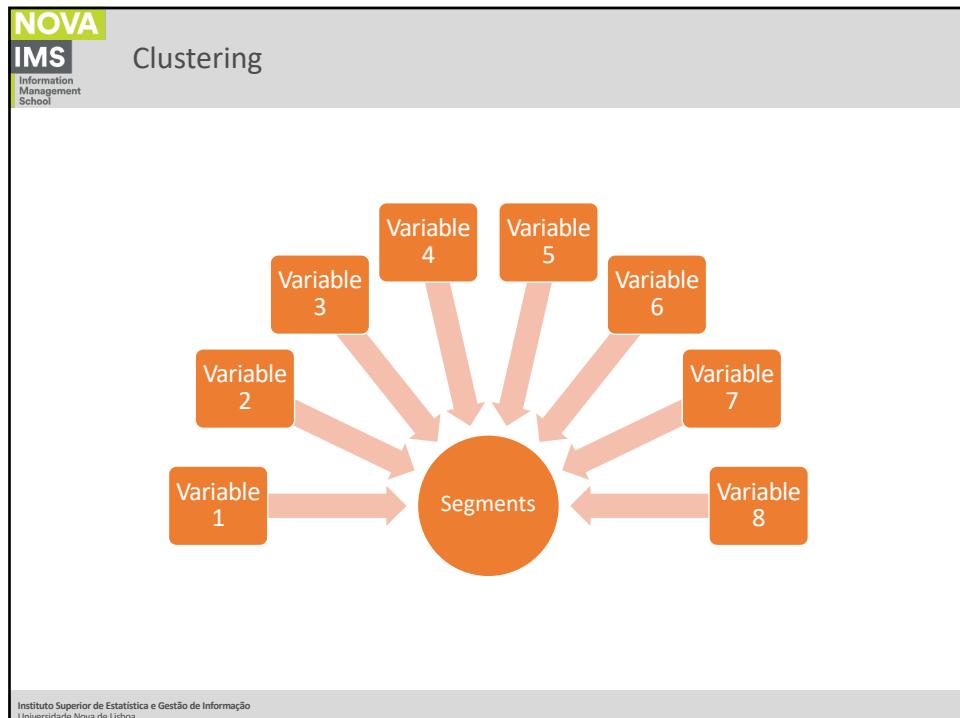
14



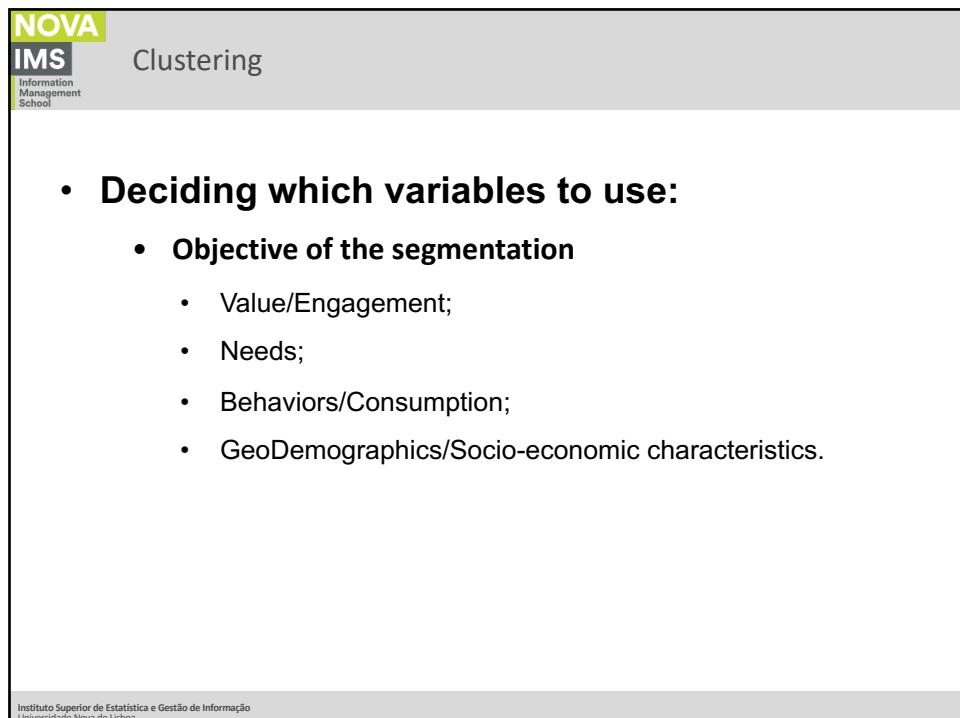
15



16



17



18

- **Deciding which variables to use:**

- The type of problem determines the variables to choose;
- If the purpose is to group objects, the choice of variables with discrimination ability is crucial;
- The quality of any cluster analysis is, first of all, conditioned by the variables used.

- **Deciding which variables to use:**

- The choice of variables should replicate a theoretical context, a reasoning;
- This process is carried out based on a set of variables that we know to be good discriminators for the problem at hand;
- First of all, the quality of the cluster analysis reflects the discrimination ability of the variables we decided to use in our study.



Similarity criterion

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



21



Clustering

- **Similarity criterion:**
 - The analysis of similarity relations has been dominated by metrics based on Euclidean Spaces;
 - Objects as points in a multidimensional space, in a way that the observed dissimilarities between the objects correspond to distances between the respective points;
 - Thus, the use of clustering methods most times means the use of similarity ratios that respect these metrics:

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

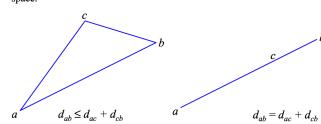
22

- **Similarity criterion:**

- In mathematics, a true measure of distance, called a *metric*, obeys three properties. These metric axioms are as follows, where d_{ab} denotes the distance between objects a and b :

1. $d_{ab} = d_{ba}$ (*measure is symmetric*)
2. $d_{ab} \geq 0$ and $= 0$ if and only if $a = b$ (*distances are always positive except when the objects are identical*)
3. $d_{ab} \leq d_{ac} + d_{cb}$ (*triangle inequality*)

Exhibit 5.1 Illustration of the triangle inequality for distances in Euclidean space.



- **Similarity criterion:**

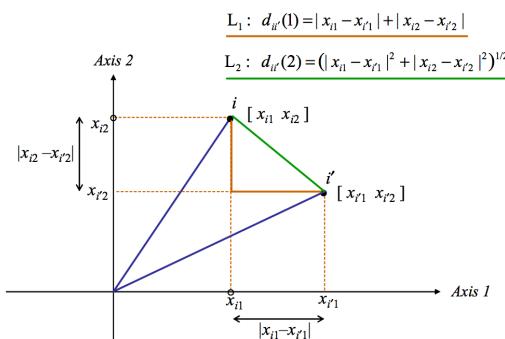
- Euclidian distance: the distance between two elements (i,j) is the square root of the sum of the squares of the differences between i and j values for all variables ($v=1, 2, \dots, p$):

$$d_{ij} = \sqrt{\sum_{v=1}^p (X_{iv} - X_{jv})^2} \text{ euclidean also known as } L_2$$

$$d_{ij} = \sum_{v=1}^p |X_{iv} - X_{jv}| \text{ City Block or } L_1$$

- **Similarity criterion:**

- Euclidian distance: the distance between two elements (i,j) is the square root of the sum of the squares of the differences between i and j values for all variables ($v=1, 2, \dots, p$):



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

25

- **Similarity criterion:**

- Minkowski distance: is defined from the absolute distance, and can be considered as a generalization of both the Euclidean distance and the Manhattan distance. It coincides with Euclidean distance when $r=2$ and with Manhattan distance when $r=1$:

$$d_{ij} = \left(\sum_{v=1}^p |X_{iv} - X_{jv}|^r \right)^{1/r}$$

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

26

- **Similarity criterion:**

- If a weight is assigned to each variable, according to their importance for the analysis, the weighted Euclidean distance takes the following form:

$$d_{ij} = \sqrt{\sum_{v=1}^p w_v (X_{iv} - X_{jv})^2}$$

- **Similarity criterion:**

- Pearson correlation coefficient: its function is to measure the degree of linear correlation between two elements, for a number of variables:

Correlation

$$\rho(x_i, x_{i'}) = \frac{\sum_j (x_{ij} - \bar{x}_i)(x_{i'j} - \bar{x}_{i'})}{\sqrt{\sum_j (x_{ij} - \bar{x}_i)^2 \sum_j (x_{i'j} - \bar{x}_{i'})^2}}$$

NOVA
IMS
Information Management School

Choose the Algorithm

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

A3ES UNIGIS A Schools eduniversal

29

NOVA
IMS
Information Management School

Clustering

```

graph TD
    Grouping --> Clustering
    Grouping --> QuartileBasedClustering
    Clustering --> HierarchicalMethods
    Clustering --> PartitionMethods
    Clustering --> DensityBased
    Clustering --> SOM
    QuartileBasedClustering --> Quartiles
    QuartileBasedClustering --> RFMAnalysis
    Quartiles --> Variables
    Quartiles --> ThroughTime
    RFMAnalysis --> MigrationMatrices
  
```

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

30

**NOVA
IMS**
Information Management School

A Priori Grouping

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

Accreditation Logos: EQUIS, UNIGIS, AACSB, ASES, AMBA, iSchools, eduniversal

31

**NOVA
IMS**
Information Management School

Clustering

```

graph TD
    Grouping --> Clustering
    Clustering --> HierarchicalMethods
    Clustering --> PartitionMethods
    Clustering --> DensityBased
    Clustering --> SOM
    HierarchicalMethods --- QuartileBased
    PartitionMethods --- RFMAnalysis
    DensityBased --- Variables
    DensityBased --- ThroughTime
    RFMAnalysis --- MigrationMatrices
    
```

The diagram illustrates the classification of clustering methods. It starts with a general category 'Grouping' at the top, which branches down to 'Clustering'. 'Clustering' further divides into four main categories: 'Hierarchical Methods', 'Partition Methods', 'Density-based', and 'Self-Organizing Maps'. The 'Hierarchical Methods' category is highlighted with a light blue background. The 'Partition Methods' category is shown in a light grey box. The 'Density-based' and 'Self-Organizing Maps' categories are also shown in light grey boxes. The 'Hierarchical Methods' category leads to 'Quartile-based clustering', which then branches into 'Quartiles' and 'RFM Analysis'. 'Quartiles' further branches into 'Variables' and 'Through time'. 'RFM Analysis' branches into 'Migration Matrices'.

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

32

NOVA
IMS
Information Management School



Quartile-based clusters

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



33

NOVA
IMS
Information Management School

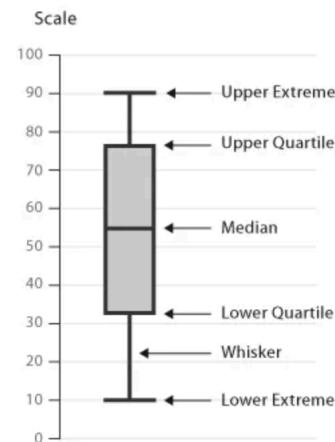
Clustering

- **Quartile-based clusters**
 - A percentile is a measure used in statistics indicating the value below which a given percentage of observations in a group of observations fall.
 - For example, the 20th percentile is the value (or score) below which 20 percent of the observations may be found.
 - The 25th percentile is also known as the first quartile (Q1), the 50th percentile as the median or second quartile (Q2), and the 75th percentile as the third quartile (Q3).
 - In general, percentiles and quartiles are specific types of quantiles.

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

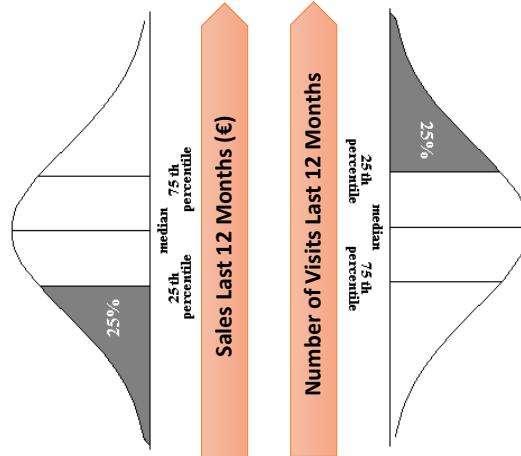
34

- Quartile-based clusters



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

35



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

36

Clustering

Sales last 12M_VL - X - Visits Last 12M_NR

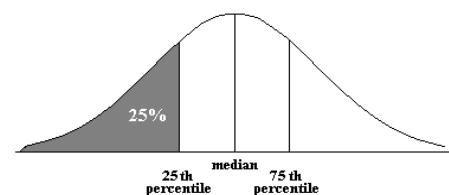
Tabela de Contingência		TICKETS_ULT_12M_NR				Totais
		<Q1	[Q1,Q2[[Q2,Q3[>=Q3	
VENDAS_ULT_12M_VL	<Q1				Buy very small amounts very frequently	
	[Q1,Q2[
	[Q2,Q3[
	>=Q3	A lot of money few times				
Totais						

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

37

Clustering

Migration Matrix



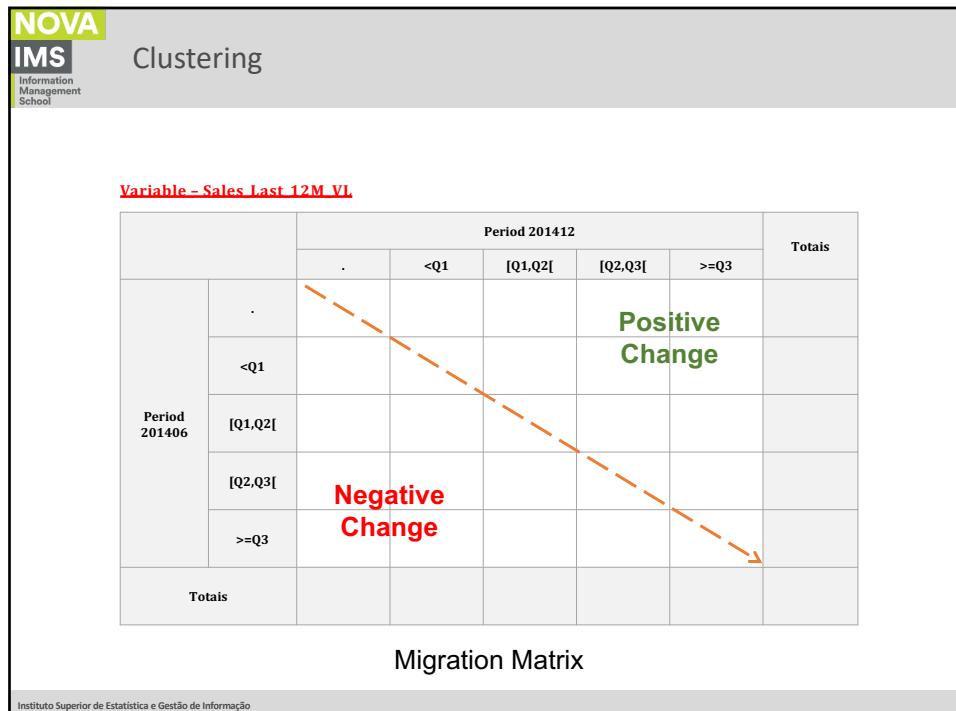
12/2012 – 12/2013

06/2013 – 06/2014

Evolution Sales Last 12 Months (12.2012 – 06.2014)

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

38



39

NOVA
IMS
Information Management School

Clustering

Tabela de Contingência

		Period 06.2013/06.2014					Totais
		. 0	<Q1 a b	[Q1,Q2[f g h	[Q2,Q3[k l m n o	>=Q3 p q r s t	
Periodo 12.2012/ 12.2013	.	30603	17050	8815	6427	62895 13.5%	
	<Q1	28734	55411	13600	2772	178	100695 21.6%
	[Q1,Q2[14834	15506	52421	16838	1097	100696 21.6%
	[Q2,Q3[6540	3450	19346	59756	11608	100700 21.6%
	>=Q3	3765	760	1835	14733	79604	100697 21.6%
	Totais		53873 11.6%	105730 22.7%	104252 22.4%	102914 22.1%	98914 21.2%

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

40

RFM Analysis

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

Accreditation Logos: EQUIS, UNICIS, A3ES, AACSB, iSchools, eduniversal

41

Clustering

- **RFM**
 - Based on the following principles:
 - Customers who have purchased more recently are more likely to purchase again;
 - Customers who have made more purchases are more likely to purchase again;
 - Customers who have made larger purchases are more likely to purchase again.

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

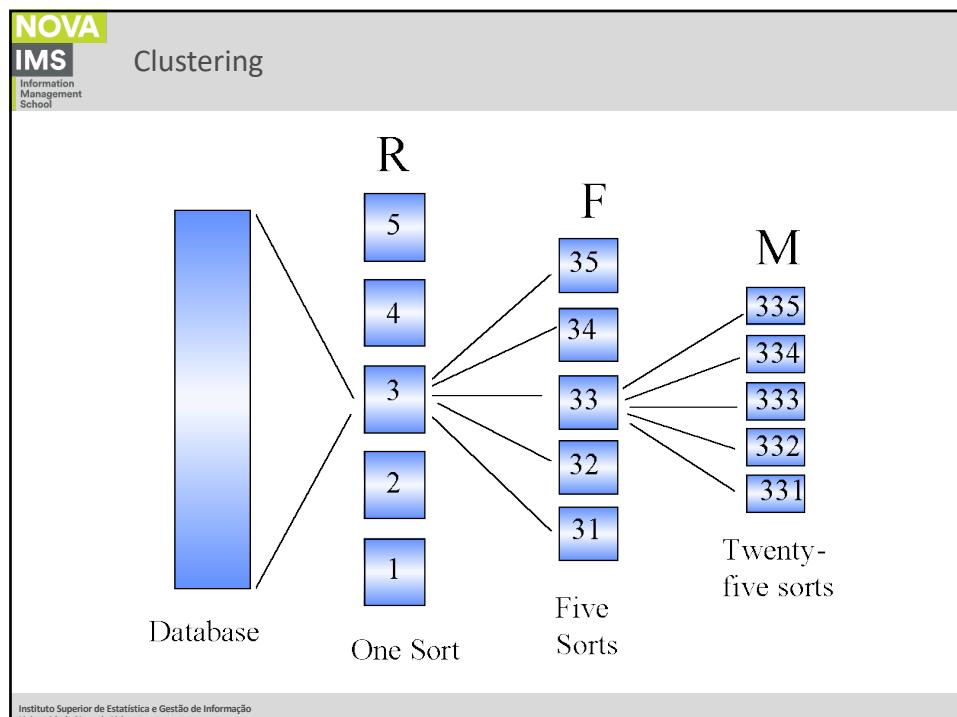
42

- **RFM**

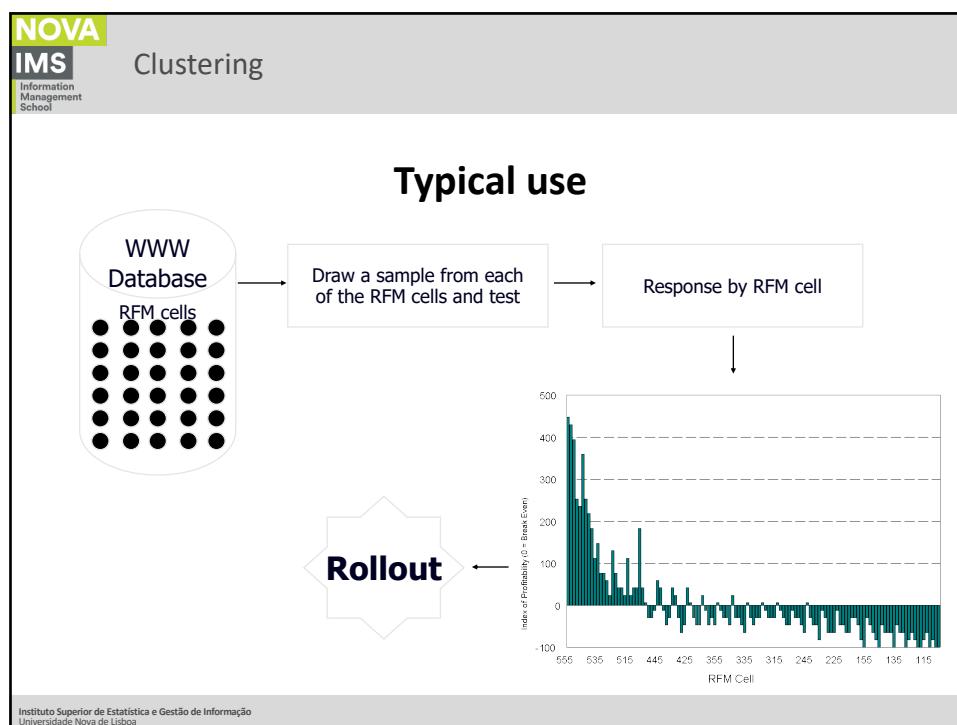
- Has been in active use in Direct Marketing for more than 40 years;
- It can be used only for customer files that contain purchase history;
- There are two methods:
 - Exact Quintiles;
 - Hard coding;

- **RFM**

- How to do it (Exact Quintiles)?
 - We sort the database according to recency and divide into 5 quintiles (5 equal segments);
 - Do the same for the variables frequency and monetary;
 - Result: 125 cells of equal size (5*5*5).



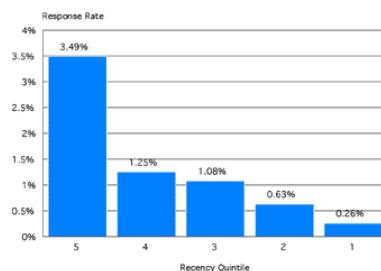
45



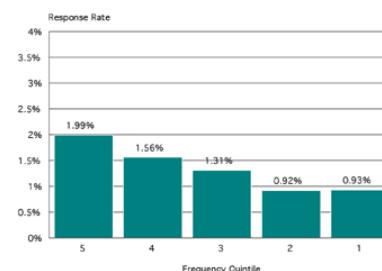
46

Clustering

Response By Recency Quintile



Response By Frequency Quintile



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

47

Clustering

Migration Matrix

Segment in YY/YY/YYYY	Segment 1 XX/XX/XXXX							Total YYYY
	44	45	51	52	53	54	55	
44	.	41914	209	4362	1862	253	22	48622
45	34200	58714	7875	14961	8968	1652	128	126498
51	505	9089	7823	4895	5420	30	.	27762
52	9109	7044	7151	83963	11103	8820	208	127398
53	3572	5758	4211	5578	29736	3691	9	52555
54	382	124	93	6507	2190	36300	4128	49724
55	69	10	22	156	62	4089	14446	18854
Total XXXX	47837	122653	27384	120422	59341	54835	18941	

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

48

- **RFM**

- Hard coding
 - Categories are divided by exact values (0-3 months; 4-6 months; 7-9 months; etc.);
 - More expensive in terms of programming, categories tend to change over time;
 - Very different quantities from cell to cell.

- **RFM**

- Its popularity comes from its simplicity, low cost and capacity to classify customers based on their behavior;
- Opportunity to carry out tests in small, representative groups of each cell;
- A more sophisticated modeling is almost always better, but is it worth it? Not always.



Questions?

51



Data Mining

NOVA-IMS
10/11/2021
Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>

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

1



AGENDA

- Cluster analysis
 - Clustering techniques
 - Hierarchical Methods (agglomerative)
 - Partitioning Methods (kmeans and k-medoids)

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

2

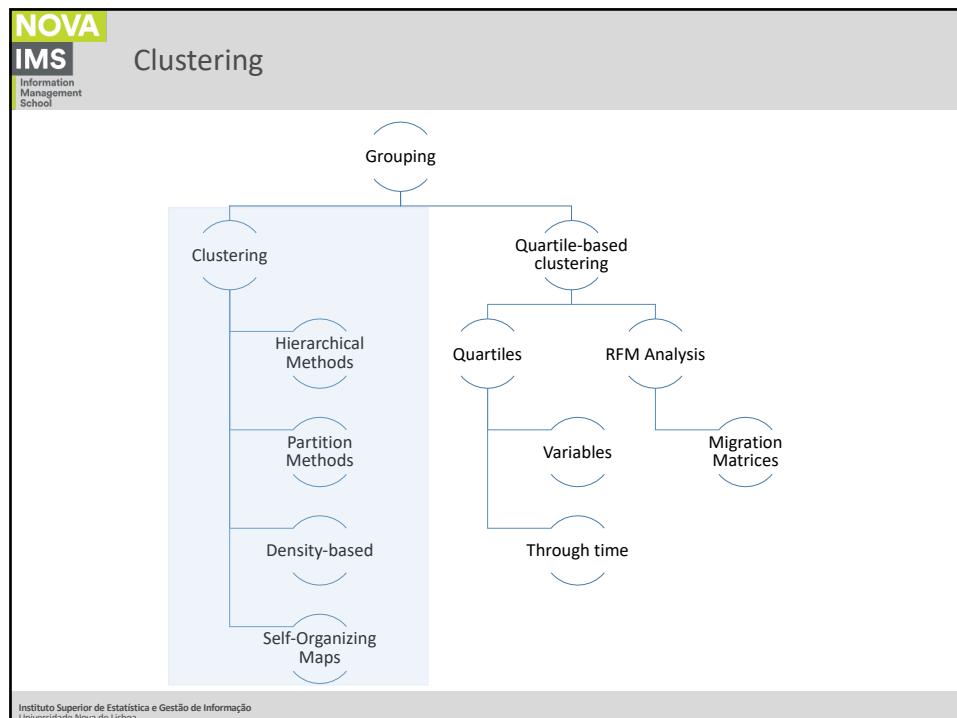
NOVA
IMS
Information Management School

Clustering

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

A3ES
UNIGIS
A3ES
iSchools
eduniversal

3



4

NOVA
IMS
Information Management School



Hierarchical Methods

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



5

NOVA
IMS
Information Management School

Clustering

- **Hierarchical Clustering**

Data Matrix				
	X_1	X_2	...	X_p
I_1				
I_2				
...				
I_n				

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

6

NOVA
IMS
Information Management School

Clustering

- Hierarchical Clustering

	X ₁	X ₂
l ₁		
l ₂		
...		
l _n		

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

7

NOVA
IMS
Information Management School

Clustering

- Hierarchical Clustering

$$d_{ij} = \sqrt{\sum_{v=1}^p (x_{iv} - x_{jv})^2}$$

	X ₁	X ₂
l ₁		
l ₂		
...		
l _n		

	l ₁	l ₂	...	l _n
l ₁	0			
l ₂	d(l ₂ , l ₁)	0		
...	d(l _{..} , l ₁)	d(l _{..} , l ₂)	0	
l _n	d(l _n , l ₁)	d(l _n , l ₂)	d(l _n , l _{..})	0

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

8

Clustering

	BA	FI	MI	NA	RM	TO
BA	0	662	877	255	412	996
FI	662	0	295	468	268	400
MI	877	295	0	754	564	138
NA	255	468	754	0	219	869
RM	412	268	564	219	0	669
TO	996	400	138	869	669	0



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

9

Clustering

- Hierarchical Clustering
 - Linkage or Aggregation Rules

- Single Linkage
 $D(c_i, c_j) = \min D(x_i, x_j)$
 Minimum distance or distance between closest elements in clusters



- Centroid Method
 Combining clusters with minimum distance between the centroids of the two clusters



- Complete Linkage
 $D(c_i, c_j) = \max D(x_i, x_j)$
 Maximum distance between elements in clusters



- Ward's Method
 • Combining clusters where increase in within cluster variance is to the smallest degree.
 • Objective is to minimize the total within cluster variance



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

10

Clustering

	BA	FI	MI/TO	NA	RM
BA	0	662	877	255	412
FI	662	0	295	468	268
MI/TO	877	295	0	754	564
NA	255	468	754	0	219
RM	412	268	564	219	0



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

11

Clustering

	BA	FI	MI/TO	NA/RM
BA	0	662	877	255
FI	662	0	295	268
MI/TO	877	295	0	564
NA/RM	255	268	564	0



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

12

Clustering

	BA/NA/RM	FI	MI/TO
BA/NA/RM	0	268	564
FI	268	0	295
MI/TO	564	295	0



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

13

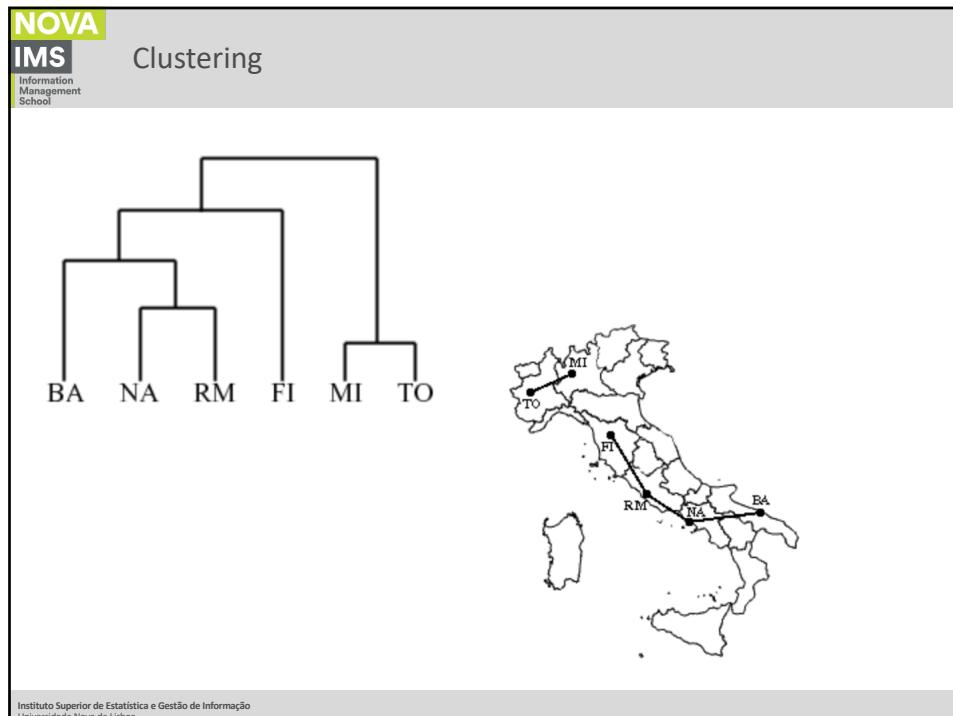
Clustering

	BA/FI/NA/RM	MI/TO
BA/FI/NA/RM	0	295
MI/TO	295	0

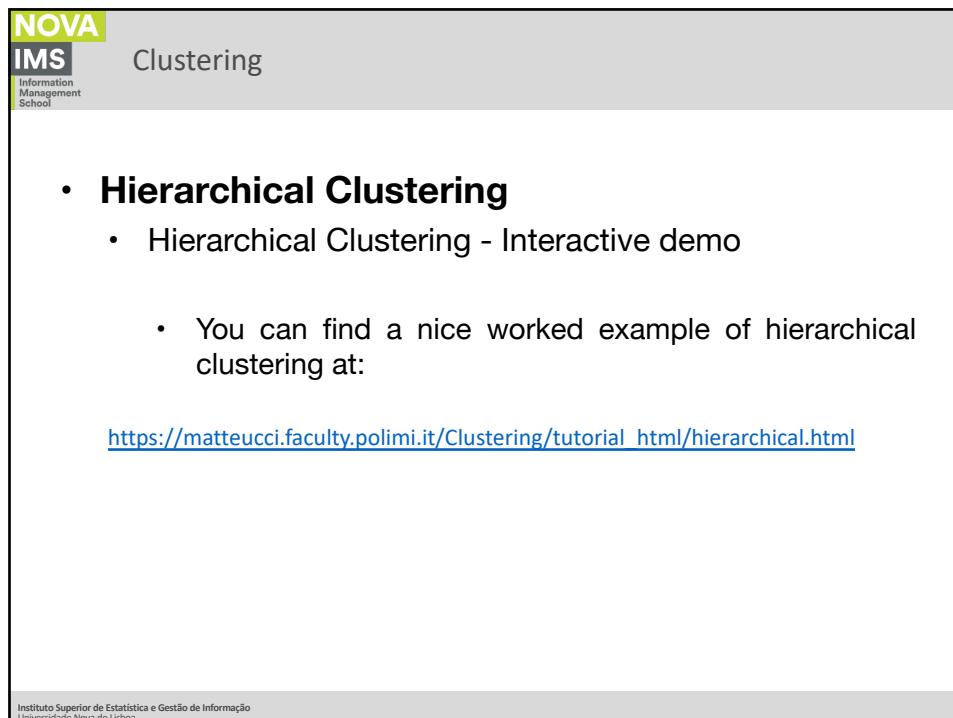


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

14



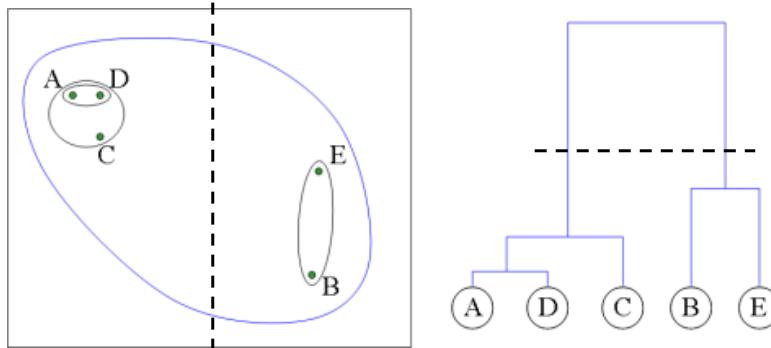
15



16

- **Hierarchical Clustering**

- Dendrogram



- **Hierarchical Clustering (disadvantages)**

- Have an essential problem, once an interaction is performed (merge or separation) we **cannot go back**;
 - This strictness is useful in terms of **computational costs**, because it avoids the cost that originates from the different combinatorial choices;
 - However, this low cost is related to the **impossibility to correct wrong decisions**;
- Time complexity: not suitable for large datasets
 - Due to its calculation needs, several operations with large matrices **do not always work very well with large datasets**.

- **Hierarchical Clustering (other variants)**

- There are two ways of improving the performance of hierarchical methods:
 - To perform a careful analysis of the links produced in each hierarchical partition (CURE and Chameleon methods);
 - To integrate hierarchical clustering and optimization, first using an agglomerative algorithm and then refining the results by using iterative optimization (BIRCH method).



K-Means Algorithm

- **K-means algorithm**

- K-means is a partitional clustering algorithm
- Let the set of data points (or instances) D be

$$\{x_1, x_2, \dots, x_n\},$$

where $x_i = (x_{i1}, x_{i2}, \dots, x_{ir})$ is a vector in a real-valued space $X \in R^r$, and r is the number of attributes (dimensions) in the data.

- The k-means algorithm partitions the given data into k clusters.
 - Each cluster has a cluster center, called centroid.
 - k is specified by the user
 - $k \ll n$.

- **K-means algorithm**

- Classifies the data into K groups, by satisfying the following requirements:
 - each group contains at least one point;
 - each point belongs to exactly one cluster.

- **K-means algorithm**

- Given k , the partition method creates an initial partition (typically randomly);
- Next, uses an iterative relocation technique that tries to improve the partition, moving objects from one group to another;
- Generically, the criterion for a good partitioning is that of objects belonging to the same cluster should be close or related to each other.

- **K-means algorithm**

- Algorithm:
 1. Choose the seeds;
 2. Each individual is associated with the nearest seed;
 3. Calculate the centroids of the formed clusters;
 4. Go back to step 2;
 5. End when the centroids cease to be recentered.

- **K-means algorithm**

- The goal is to minimize intra-group variance (sum of squared error):

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x)$$

- x is a data point in cluster C_i and m_i is the representative point (centroid) for cluster C_i
- One easy way to reduce SSE is to increase K (number of clusters)
- A good clustering with smaller K can have a lower SSE than a poor clustering with higher K



K-Means Algorithm in figures

NOVA
IMS
Information Management School

Clustering

- **K-means**

- Individuals measured based on two variables;
- The goal is to group them in homogeneous sets.

Source: Fiona Cameron, Techniques for Neighbourhood Classification
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

27

NOVA
IMS
Information Management School

Clustering

- **K-means**

- Seeds
- Randomly chosen

Source: Fiona Cameron, Techniques for Neighbourhood Classification
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

28

NOVA
IMS
Information Management School

Clustering

- **K-means**
- Allocate individuals to the nearest seed

Source: Fiona Cameron, Techniques for Neighbourhood Classification
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

29

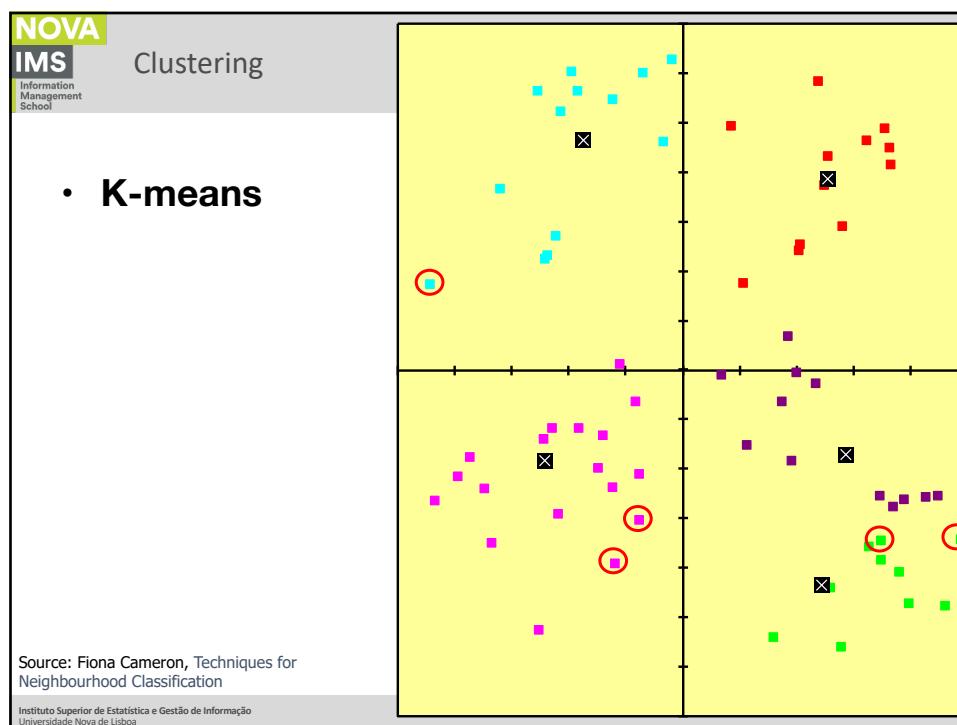
NOVA
IMS
Information Management School

Clustering

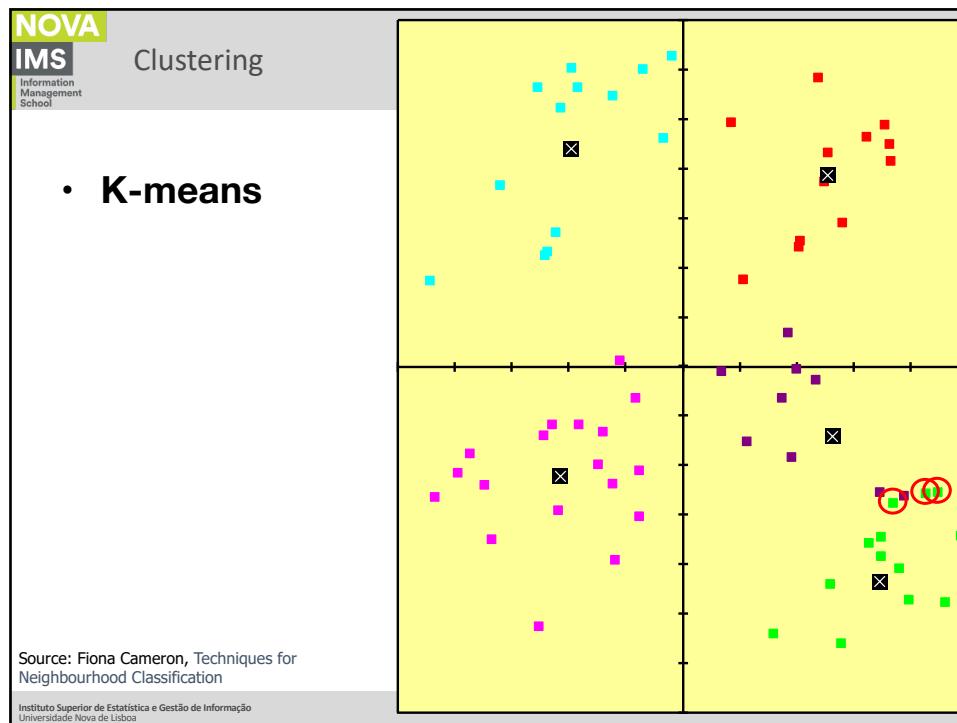
- **K-means**
- Recenter the seed so that it stays in the center of the cloud of points (called centroid)
- Some individuals change cluster

Source: Fiona Cameron, Techniques for Neighbourhood Classification
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

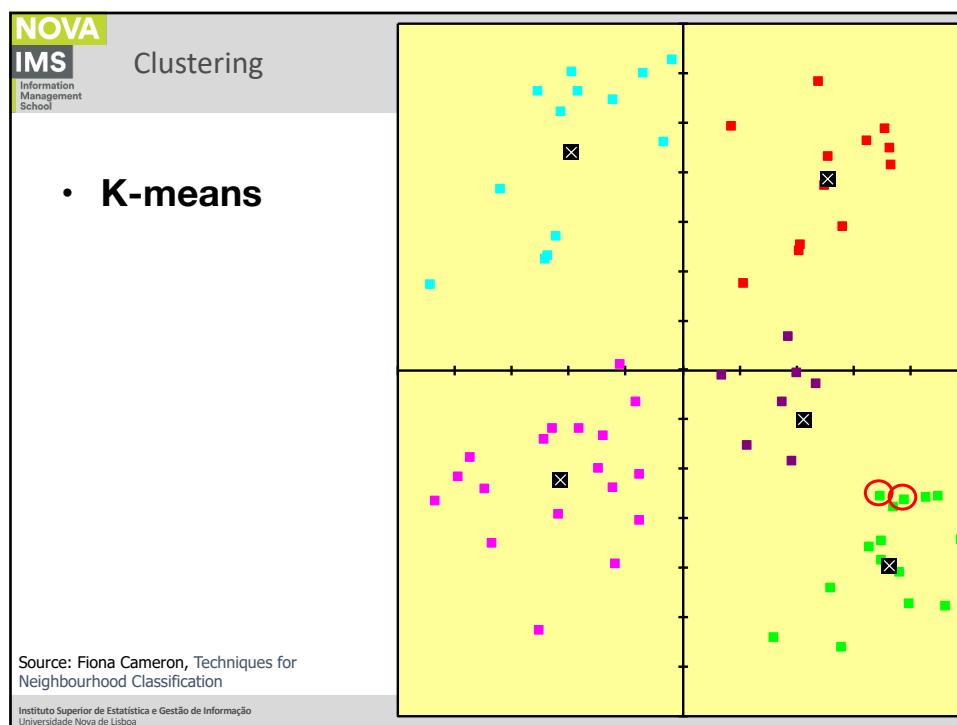
30



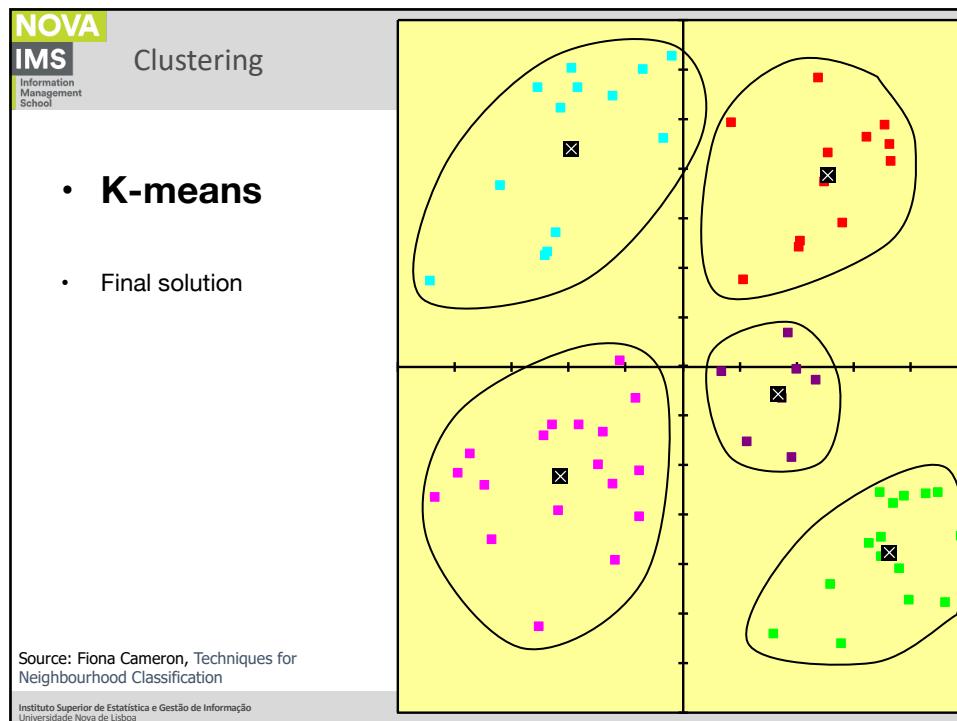
31



32



33



34

NOVA
IMS
Information Management School

Clustering

- **K-means**
- Movement of centroids during optimization process

Source: Fiona Cameron, Techniques for Neighbourhood Classification
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

35

NOVA
IMS
Information Management School

Clustering

- **K-means algorithm (strengths)**
 - Simple: easy to understand and to implement
 - Efficient: Time complexity $O(tkn)$,
where n is the number of data points,
 k is the number of clusters, and
 t is the number of iterations.
 - Since both k and t are small, k-means is considered a linear algorithm.
 - K-means is the most popular clustering algorithm.
 - Note that: it terminates at a local optimum if SSE is used. The global optimum is hard to find due to complexity.

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

36

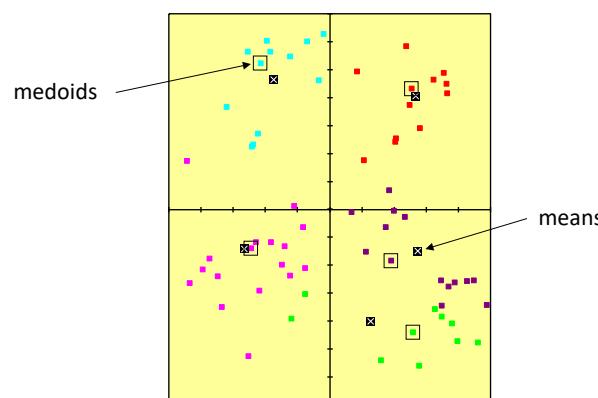
- **K-means algorithm (weaknesses)**
 - Very sensitive to the existence of outliers;
 - Very sensitive to the initial positions of the seeds;
 - Partitioning methods work well with spherical-shaped clusters;
 - Partitioning methods are not the most suitable to find clusters with complex shapes and different densities;
 - The need to set from the start the number of clusters to create.

The Algorithm variant

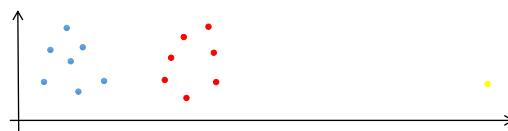
- **K-means and k-medoids algorithms**

- Most algorithms adopt one of two very popular heuristics:
 - k-means algorithm, where each cluster is represented by the average of the values of the points in a cluster;
 - k-medoids algorithm, where each cluster is represented by one of the points located near the center of the cluster.

- **K-means and k-medoids algorithms**

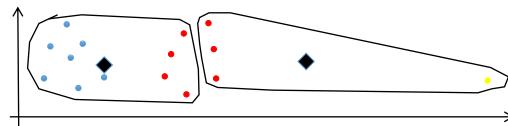


- K-means and k-medoids algorithms



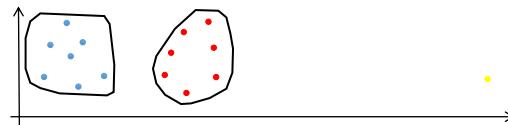
Source: CS583, Bing Liu, UIC, Chapter 4:
Unsupervised Learning

- K-means and k-medoids algorithms



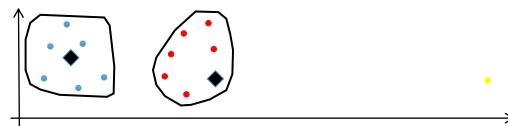
Source: CS583, Bing Liu, UIC, Chapter 4:
Unsupervised Learning

- **K-means and k-medoids algorithms**



Source: CS583, Bing Liu, UIC, Chapter 4:
Unsupervised Learning

- **K-means and k-medoids algorithms**



Source: CS583, Bing Liu, UIC, Chapter 4:
Unsupervised Learning

**NOVA
IMS**
Information Management School

Clustering

- K-means and k-medoids algorithms

Source: CS583, Bing Liu, UIC, Chapter 4:
Unsupervised Learning

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

45

**NOVA
IMS**
Information Management School

The initialization problem

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

46

**NOVA
IMS**
Information Management School

Clustering

- **K-means algorithm (weaknesses)**
 - The algorithm is sensitive to initial seeds.

Source: Tan, Steinbach, Kumar, Introduction to Data Mining 2004

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

47

**NOVA
IMS**
Information Management School

Clustering

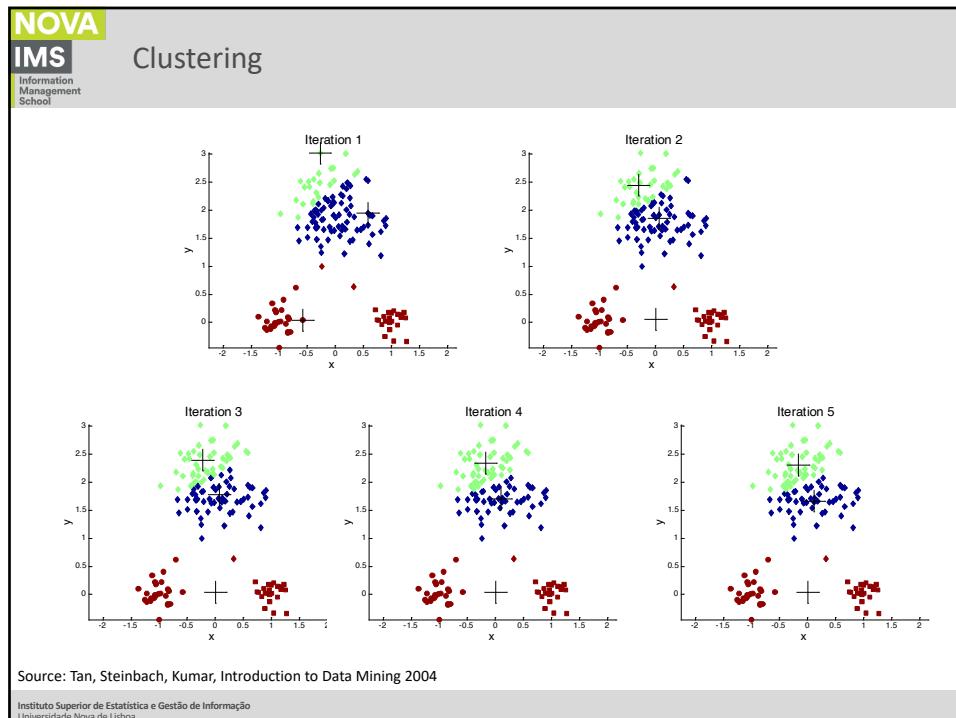
Iteration 1 Iteration 2 Iteration 3

Iteration 4 Iteration 5 Iteration 6

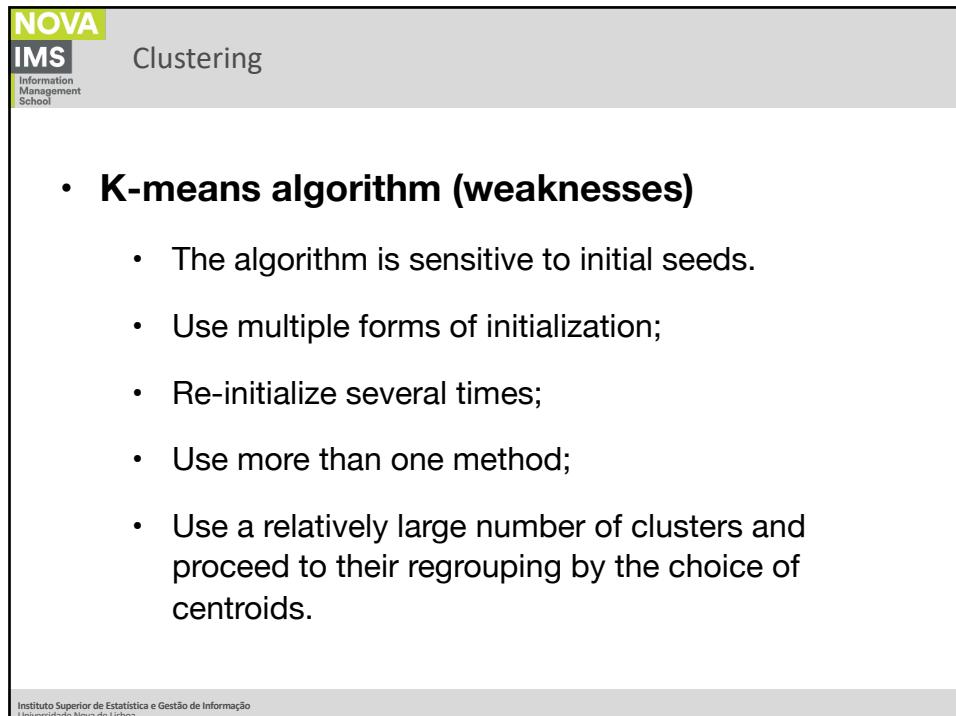
Source: Tan, Steinbach, Kumar, Introduction to Data Mining 2004

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

48



49



50

**NOVA
IMS**
Information Management School

Shape and density

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

UNIGIS A3ES iSchools eduniversal

51

**NOVA
IMS**
Information Management School

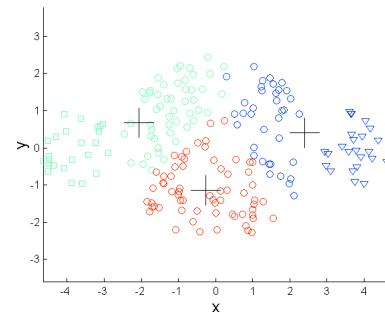
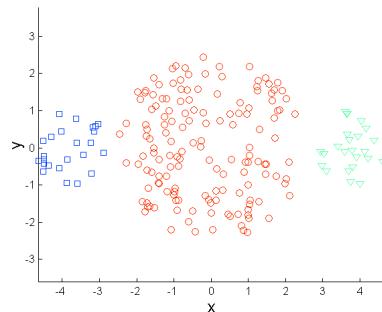
Clustering

- **K-means algorithm (weaknesses)**

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

52

- **K-means algorithm (weaknesses)**
 - Have difficulties in dealing with clusters of different size and density;



- **K-means algorithm (weaknesses)**
 - Each individual either belongs or does not belong to the cluster, having no notion of probability of belonging, in other words, there is no consideration of the quality of the representation of a particular individual in a given cluster.

**NOVA
IMS**
Information Management School

Clustering

- K-means algorithm (weaknesses)

Source: Fiona Cameron, Techniques for Neighbourhood Classification

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

55

**NOVA
IMS**
Information Management School

The number of clusters

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

AES

UNIGIS

A3ES

Schools

eduniversal

56

**NOVA
IMS**
Information Management School

Clustering

- K-means algorithm the number of clusters

The figure consists of four scatter plots arranged in a 2x2 grid. The top row shows two clusters (black dots) and six clusters (black stars, green triangles, cyan circles, orange diamonds, and yellow squares). The bottom row shows two clusters (red squares), four clusters (blue triangles), and six clusters (red stars, blue inverted triangles, yellow diamonds, and yellow squares). Each plot includes a question 'How many clusters?'.

How many clusters?

Six Clusters

Two Clusters

Four Clusters

Source: Tan, Steinbach, Kumar, Introduction to Data Mining 2004

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

57

**NOVA
IMS**
Information Management School

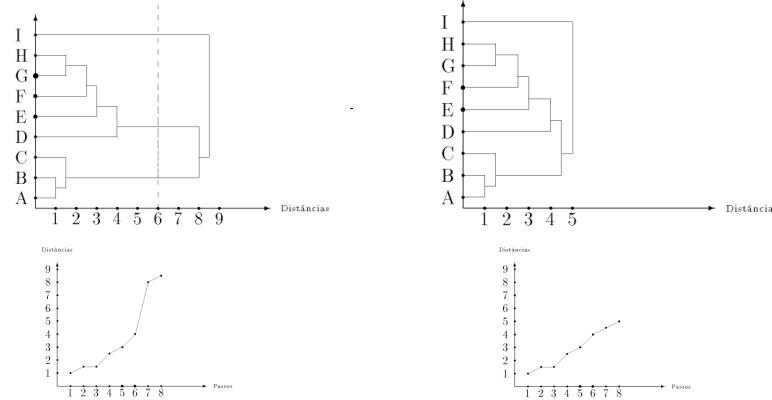
Clustering

- K-means algorithm the number of clusters
 - This is always a difficult problem to solve, and there are no recipes to fix this.
 - One way to minimize the problem is to create various classifications with different K, and choose the best.
 - Use a hierarchical method in order to choose the number of clusters based on the dendrogram.

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

58

- **K-means algorithm the number of clusters**



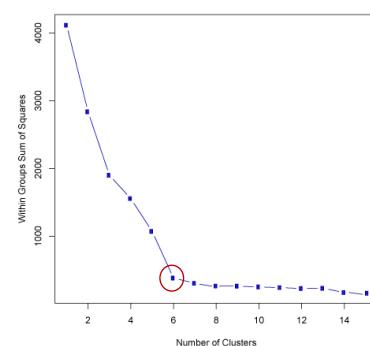
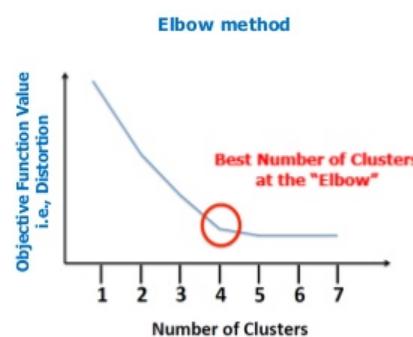
- **K-means algorithm the number of clusters**

- The choice should be guided by three fundamental criteria:
 - intra-cluster variance,
 - evaluation of the profile of the cluster (subjective),
 - operational considerations.

- **K-means algorithm the number of clusters**

- Regarding the first criterion, the analysis is simple and not too subjective, since we know that the lower the intra-cluster variability the greater the cohesion of the cluster, a highly desirable feature in this type of analysis. However, as k increases, variability decreases;
- Regarding the second criterion, the question is not as simple in the sense that it requires much more subjective assessments, which relate to the interpretation of the obtained clusters;
- The third criterion is relatively simple in the sense that these issues are imposed on the analyst.

- **K-means algorithm the number of clusters**



- **K-means algorithm the number of clusters**
 - To test the results by varying k (number of clusters);
 - This procedure allows a series of analyzes that can instruct the choice of the number of clusters;
 - To compare the totals of the distances of the different solutions.

- **K-means algorithm the number of clusters**
 - Operational considerations are related to business environment and usually affect the decisions of the analyst:
 - A number small enough for developing a specific strategy;
 - A number of individuals large enough to be worth it to develop a specific strategy;
 - A good way to accomodate these considerations is the use of a high initial k and then proceed to the grouping of clusters.

**NOVA
IMS**
Information Management School

Clustering

- K-means algorithm the number of clusters**

CLUSTER	Cluster 1	Cluster 2	Cluster 3	Cluster 4
1	0	46.88621549	53.629114781	51.059735073
2	46.88621549	0	35.424488679	46.409611185
3	53.629114781	35.424488679	0	58.950971223
4	51.059735073	46.409611185	58.950971223	0

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

65

**NOVA
IMS**
Information Management School

Clustering

- K-means algorithm the number of clusters**

- This evaluation is carried out by comparing the mean values for each variable in each cluster with the mean values of the population for each variable;
- In this case, it is particularly relevant to take into account the most important differences within the different clusters and the mean population.
- That is why profiling is so important.

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

66



Data Mining

Partitioning Clustering (k-means)

Note on Profiling

17/11/2021

NOVA-IMS

Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/bacao>

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

1



AGENDA

- Cluster analysis
 - Clustering techniques
 - Partitioning Methods (kmeans and k-meadoids)
 - **A Note on Profiling**

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

2

NOVA
IMS
Information Management School



Profiling

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



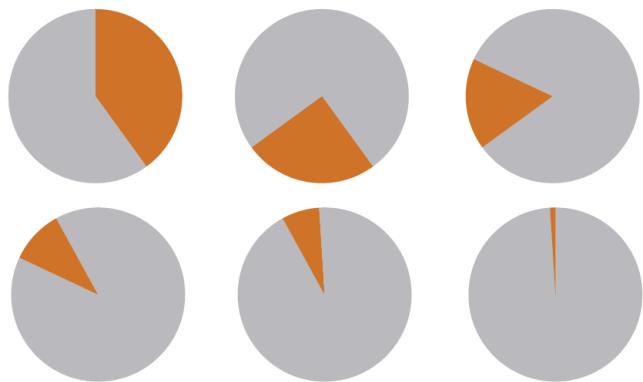
3

NOVA
IMS
Information Management School

Clustering

- Profiling (size of the clusters)**

Part-to-Whole Mini-Pie Charts



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

4

**NOVA
IMS**
Information Management School

Clustering

- Profiling (comparing averages)

A scatter plot illustrating the comparison of averages for four variables (A, B, C, D) across two clusters. The x-axis represents the 'Normalized Mean' from 0 to 1, and the y-axis represents the 'Variables' A, B, C, and D. Two cluster averages are shown: a 'Database average' (represented by a grey vertical bar) and a 'Cluster average' (represented by a dark red square). Data points are represented by blue squares.

Variable	Database Average (Normalized Mean)	Cluster Average (Normalized Mean)
A	~0.85	~0.95
B	~0.15	~0.15
C	~0.35	~0.35
D	~0.25	~0.25

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

5

**NOVA
IMS**
Information Management School

Clustering

- Profiling (comparing profiles)

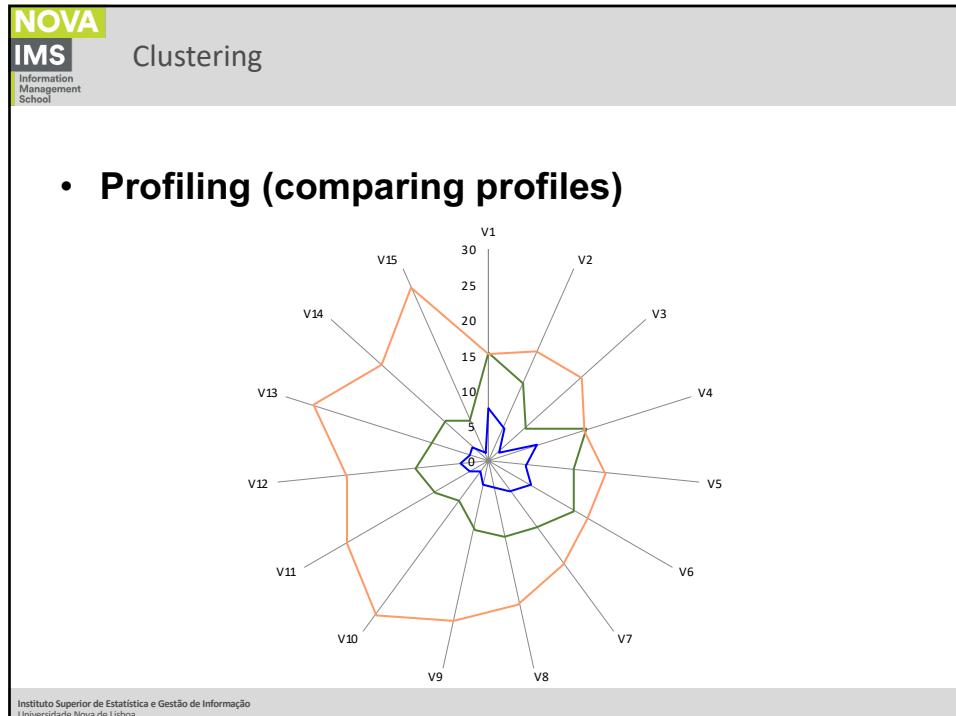
Four line graphs comparing trends in different health categories from 1975 to 2010. Each graph shows a central blue line with data points and surrounding grey lines representing confidence intervals. The graphs are labeled: Circulatory, Mental, Musculoskeletal, and Cancer.

Category	Year	Value
Circulatory	1975	32
	2010	11
Mental	1975	11
	2010	23
Musculoskeletal	1975	17
	2010	26
Cancer	1975	10
	2010	14

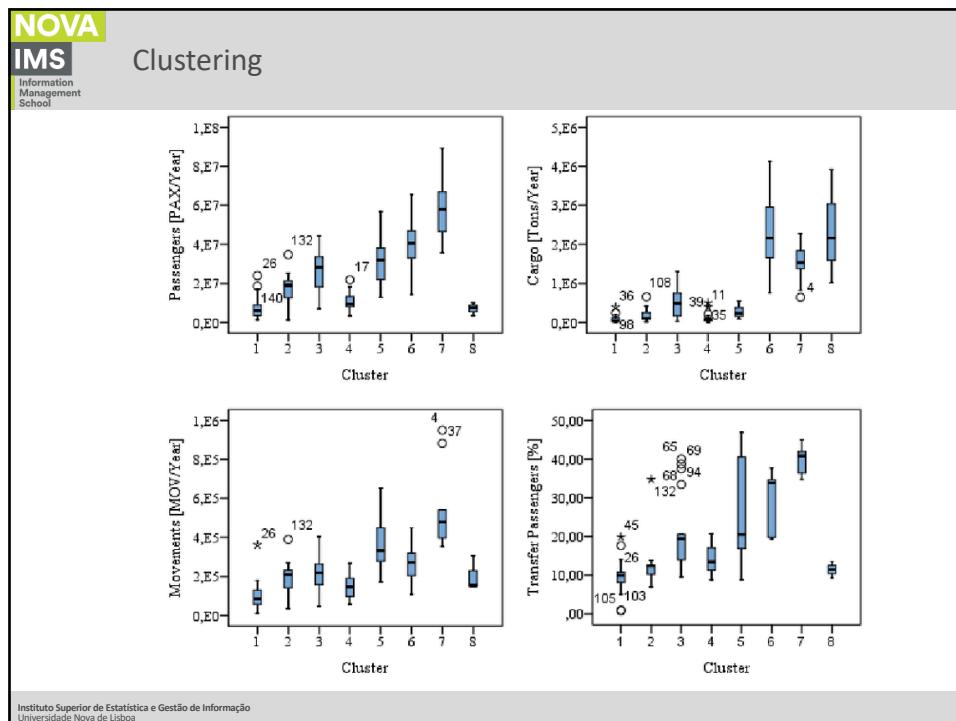
JA Schwabish "An Economist's Guide to Visualizing Data", The Journal of Economic Perspectives, 2014

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

6

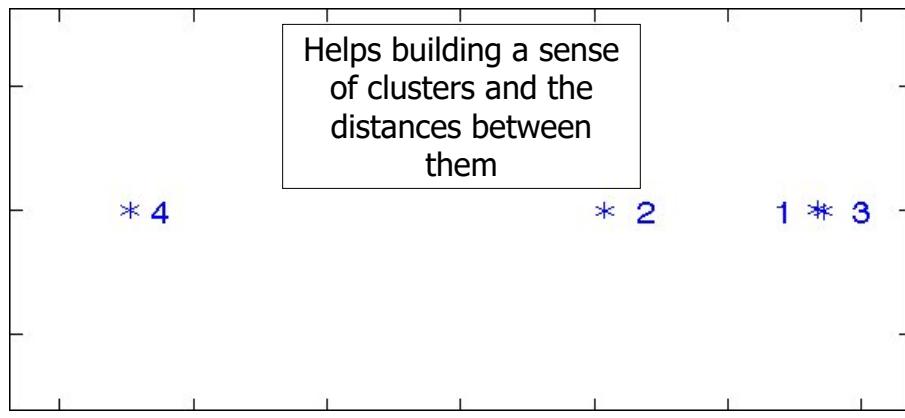


7

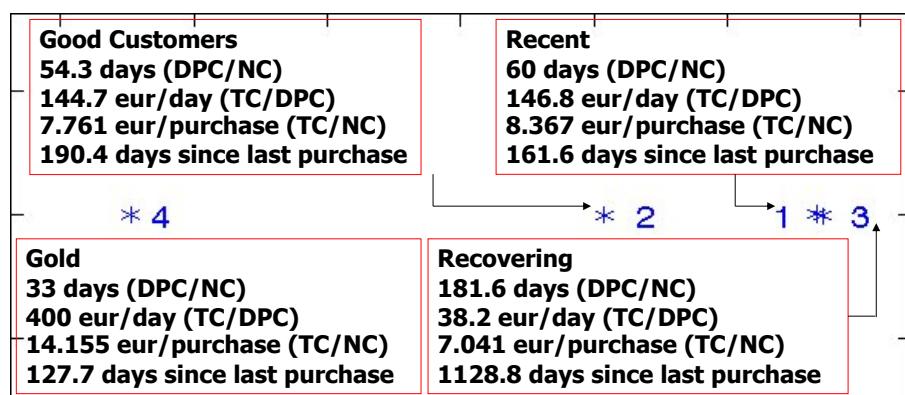


8

- Profiling (multidimensional scaling)



- Profiling (multidimensional scaling)



**NOVA
IMS**
Information Management School

Clustering

- Profiling (comparing variables not used)

Used and unused variable distribution

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

11

**NOVA
IMS**
Information Management School

Clustering

- Profiling (leverage)

Leverage
Ratio
Sales(%) / Individuals(%)

Cluster	Leverage Ratio
C1	8.36
C2	1.96
C3	0.52
C4	0.48

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

12

- Exploring two solutions

		Sofist Digital	Variedade Serviço	Preço	Status_PC	ProAtividade	Tempo Internet e RS	Sustentabilidade	Life Status	Propensao_Pfi delizacao	Switching	
k=5	FREQ	F1_SD	F2_VS	F3_P	F4_SPC	F5_PA	F6_IRS	F7_CC	F8_LS	F9_FI	F10_AB	
1	222	-0,40	-0,54	-0,31	-0,29	0,13	0,65	-1,67	-0,20	-0,09	-0,03	Pouco Envolvidos e pouco sensíveis ao preço
2	463	0,70	0,44	0,38	0,13	-0,32	-0,37	-0,19	0,31	-0,58	0,03	Social Customers
3	469	0,58	-0,16	0,15	0,48	0,23	0,00	0,09	-0,22	0,98	-0,04	Sofisticado com Propensao Fidelizacao
4	428	-1,47	0,27	0,25	-0,13	0,01	-0,25	0,29	0,15	0,15	0,04	Info Excluidos
5	420	0,29	-0,30	-0,68	-0,39	0,02	0,31	0,69	-0,16	-0,56	-0,02	Sustentaveis
k=6	FREQ	F1_SD	F2_VS	F3_P	F4_SPC	F5_PA	F6_IRS	F7_CC	F8_LS	F9_FI	F10_AB	
1	241	0,12	-0,30	-0,09	-0,20	-0,04	1,34	-1,29	-0,03	-0,13	-0,24	Pouco Envolvidos
2	134	-0,16	0,11	0,02	-0,15	-0,10	0,13	0,10	0,00	0,07	2,84	Switchers
3	285	0,43	0,48	-1,28	-0,16	0,52	-0,28	0,05	0,04	0,23	-0,07	Exigentes
4	426	-1,42	0,22	0,30	-0,20	-0,03	-0,27	0,10	0,01	0,21	-0,27	Info Excluidos
5	414	0,39	-0,77	0,20	0,92	0,19	-0,23	0,15	0,21	0,15	-0,16	Status/Conservadores
6	502	0,63	0,30	0,34	-0,37	-0,38	-0,10	0,35	-0,19	-0,39	-0,24	Social Customers

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

13

- Exploring two solutions with different k

K=7 VS K=8	Pouco envolvidos	Switchers	Informados / Social Customers	Fidelizáveis Sofisticados	Info Excluidos	Status	Proativos	Desprendido	Grand Total
Pouco envolvidos	71		23	13	45	1	1	49	203
Switchers	2	106	4	4	2	1	3	122	
Informados / Social Customers	4		219	68	2	20	1	64	378
Fidelizáveis Sofisticados	72		105	191	4	1	2	69	444
Status/Conservadores	7		5	39	13	157	42	263	
Proativos	5		15	20	3		153	19	215
Info Excluidos		12	4	309	1	1	50	50	377
Grand Total	161	106	383	339	378	181	158	296	2002

K=6 VS K=8	Pouco envolvidos	Switchers	Informados / Social Customers	Fidelizáveis Sofisticados	Info Excluidos	Status	Proativos	Desprendido	Grand Total
Pouco Envolvidos	146		16	24	8	17	7	23	241
Switchers	2	105	6	11	1	2	4	3	134
Info Excluidos	1	1	26	56	13	27	50	101	285
Status/Conservadores	3		13	177	5	98	58	60	414
Social Customers	9		322	49	1	20	18	83	502
Grand Total	161	106	383	339	378	181	158	296	2002

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

14

**NOVA
IMS**
Information Management School

Clustering

- **Consolidation of two solutions**

The diagram illustrates the process of consolidating two segments into a larger number of consolidations. It features three main components: 'Value Segment (4)' at the top left, 'Buying Segment (4)' below it, and 'Consolidations (16)' at the bottom right. Each segment is represented by a blue rounded rectangle containing four small colored squares (blue, orange, red, green). Arrows point from both the 'Value Segment' and the 'Buying Segment' towards the 'Consolidations' box, indicating their merging into a total of 16 consolidations.

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

15



Data Mining

Density-based Clustering (dbSCAN)

16/11/2021

NOVA-IMS

Fernando Lucas Bação

bacao@isegi.unl.pt

<http://www.isegi.unl.pt/fbacao>

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

1



AGENDA

- Cluster analysis
 - Clustering techniques
 - Density-based clustering (DBscan)

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

2

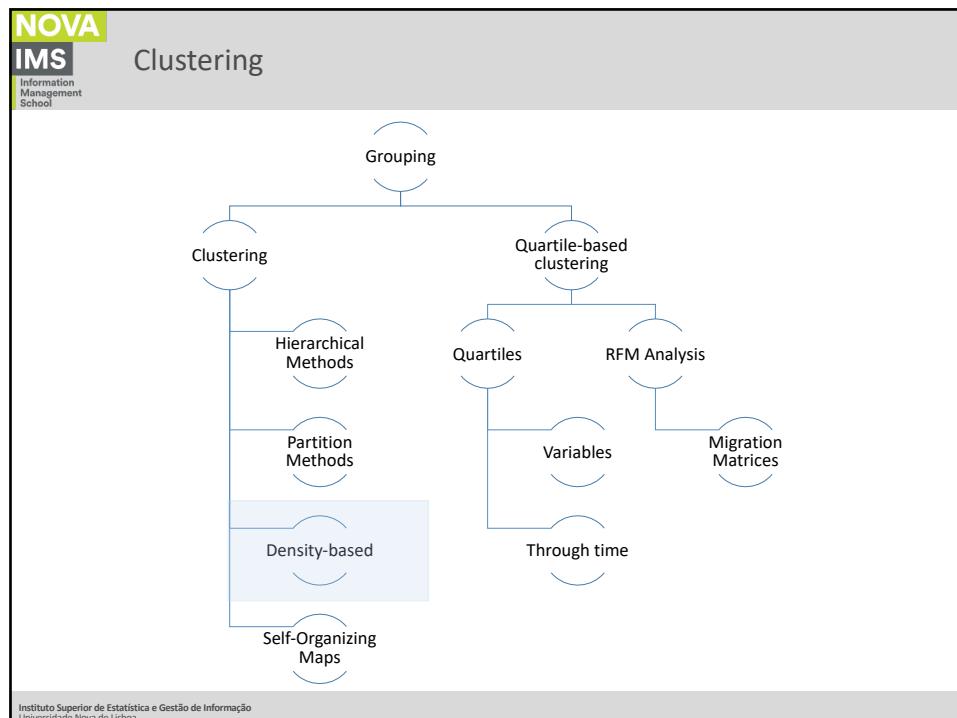
NOVA
IMS
Information Management School

Clustering

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

A3ES
UNIGIS
A3ES
iSchools
eduniversal

3



4

**NOVA
IMS**
Information Management School



Density-Based Clustering (DBSCAN)

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



5

**NOVA
IMS**
Information Management School

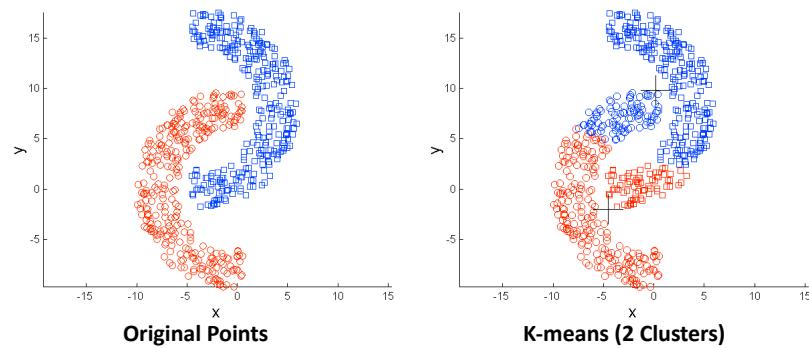
Density-Based Clustering (DBSCAN)

- **The idea**
 - Previous clustering methods have some limitations. They are based on a particular set of assumptions that if not true, the process yields suboptimal results.
 - That is they would likely inaccurately identify non-convex regions, and where noise or outliers are included in the clusters.

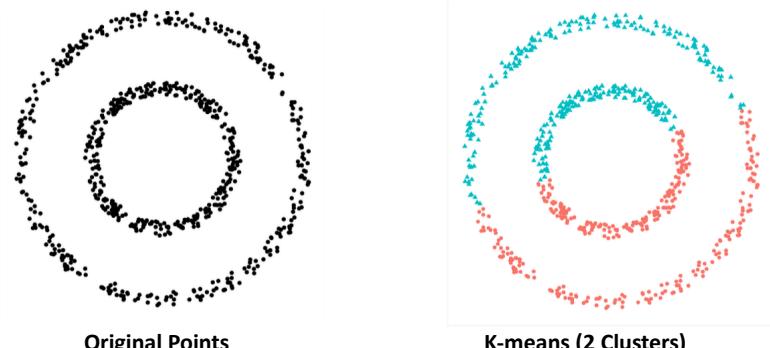
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

6

Density-Based Clustering (DBSCAN)



Density-Based Clustering (DBSCAN)



- **The idea**

- Density-based clustering algorithms try to **find clusters in data without assuming a particular shape**. Thus solving correctly cases like the ring example.
- To find clusters of **arbitrary shape**, these methods model clusters as **dense regions in data space**, separated by sparse regions.

- **The idea**

- We want to be able to cluster data like this:



- **Characteristics**

- What do talk about when we talk about density?
- What is our physical intuition behind density?
 - We say a point p is in a dense region if there are **many points** in the neighborhood of p . In this sense, the density around a point can be measured by the number of points surrounding it.
 - To measure the density around a point p we use the tradition topological definition of neighborhood. The ε -neighborhood of a point p is the space within a radius $\varepsilon > 0$ centered in p .

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

- DBSCAN takes two input parameters:
 - ε - the radius defining the neighborhood
 - $MinPts$ - the minimum of points in ε -neighborhood

Density-Based Clustering (DBSCAN)

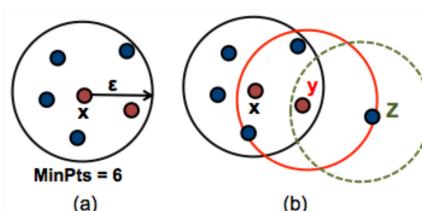
- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density
 - DBSCAN importante concepts:
 - If z is a point that have at least $MinPts$ in its ε -neighborhood is called **core point**.
 - x is **border point**, if the number of its neighbors is less than $MinPts$, but it belongs to the ε -neighborhood of some core point z
 - If a point is neither a **core** nor a **border** point, then it is called a **noise point** or an **outlier**.

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

13

Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density



- Assuming $MinPts = 6$.
 - x is a **core point** because ε -neighborhood (x) = 6,
 - y is a **border point** because ε -neighborhood (y) < $MinPts$, but it belongs to the ε -neighborhood of the core point x .
 - z is a **noise point**.

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

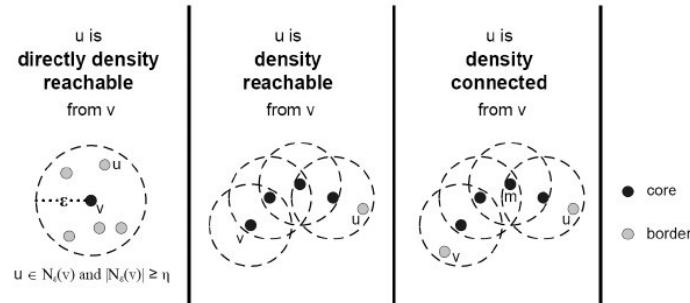
14

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density
 - We define 3 terms, required for understanding the DBSCAN algorithm:
 - **Direct density reachable:** A point “A” is directly density reachable from another point “B” if:
 - i) “A” is in the ε –neighborhood of “B” and
 - ii) “B” is a core point.

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density
 - We define 3 terms, required for understanding the DBSCAN algorithm:
 - **Density reachable:** A point “A” is density reachable from “B” if there are a set of core points leading from “B” to “A”.
 - **Density connected:** Two points “A” and “B” are density connected if there are a core point “C”, such that both “A” and “B” are density reachable from “C”.

Density-Based Clustering (DBSCAN)

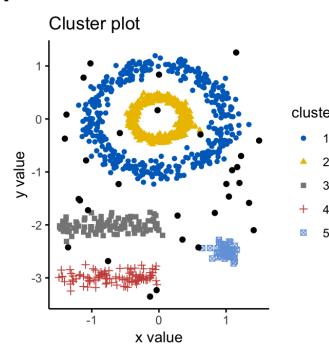
- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density



Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

- A **density-based cluster** is defined as a **group of density connected points**.

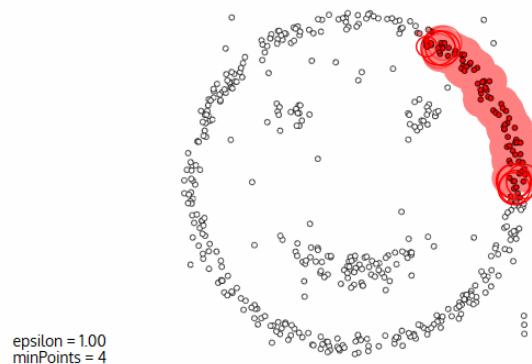


• Algorithm

1. DBSCAN begins with an arbitrary data point that has not been visited. The ε – neighborhood of this point is extracted (all points which are within the ε distance are neighborhood points).
2. If there are a sufficient number of points (according to $MinPts$) within this neighborhood then the clustering process starts and the current data point becomes the first point in the new cluster. Otherwise, the point will be labeled as noise. In both cases that point is marked as “visited”.
3. For this first point in the new cluster, the points within its ε distance neighborhood also become part of the same cluster. This procedure of making all points in the ε –neighborhood belong to the same cluster is then repeated for all of the new points that have been just added to the cluster group.
4. This process of steps 2 and 3 is repeated until all points in the cluster are determined i.e all points within the ε –neighborhood of the cluster have been visited and labeled.

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

19



Restart



Pause

<https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>

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

20

- **Advantages**

- It does not require a pre-set number of clusters at all.
- It identifies outliers as noises (not affected by),
- it can finds arbitrarily sized and arbitrarily shaped clusters quite well.

- **Disadvantages**

- It doesn't perform well when the clusters are of varying density
- Setting of the distance threshold ε and $MinPts$ for identifying the neighborhood points will vary from cluster to cluster when the density varies
- Doesn't work well in high-dimensional spaces



Questions?

Density-Based Clustering (DBSCAN)

- **DBSCAN:** Density-Based Clustering Based on Connected Regions with High Density

1. For each point x_i , compute the distance between x_i and the other points.
Finds all neighbor points within ε –neighborhood of the starting point x_i . Each point, with a neighbor count greater than or equal to $MinPts$, is marked as **core point** or visited.
2. For each **core point**, if it's not already assigned to a cluster, create a new cluster. Find recursively all its density connected points and assign them to the same cluster as the **core point**.
3. Iterate through the remaining unvisited points in the dataset.

Those points that do not belong to any cluster are treated as outliers or noise.



Data Mining

Mean-shift algorithm

16/11/2021

NOVA-IMS

Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>

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

1



AGENDA

- Cluster analysis
 - Clustering techniques
 - Density-based clustering (mean shift clustering)

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

2

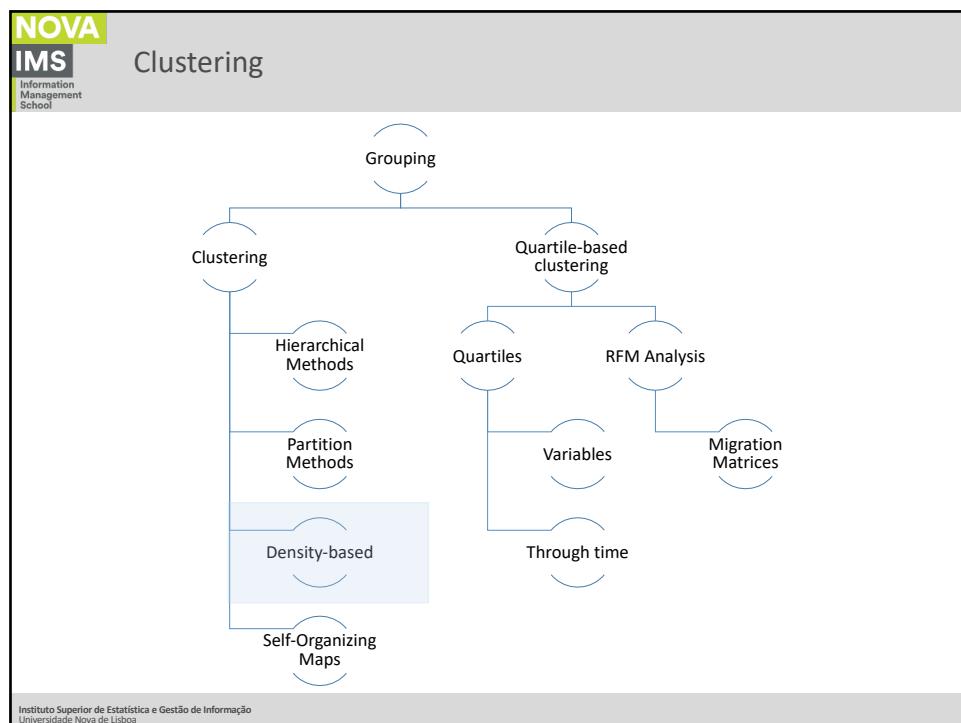
NOVA
IMS
Information Management School

Clustering

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

A3ES
UNIGIS
A3ES
iSchools
eduniversal

3



4



Mean-Shift Clustering

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



5



Mean-Shift Clustering

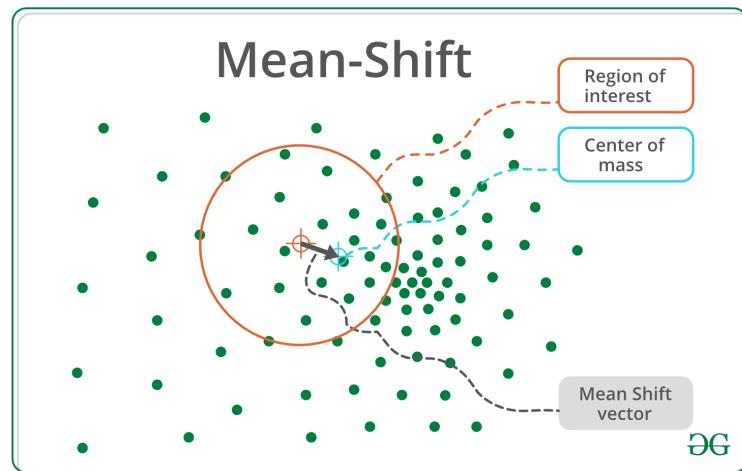
- **Characteristics**
 - Mean shift clustering is a sliding-window-based algorithm that attempts to find dense areas of data points.
 - It is a centroid-based algorithm meaning that the goal is to locate the center points of each group/class,
 - Works by updating candidates for center points to be the mean of the points within the sliding-window.
 - These candidate windows are then filtered in a post-processing stage to eliminate near-duplicates, forming the final set of center points and their corresponding groups.

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

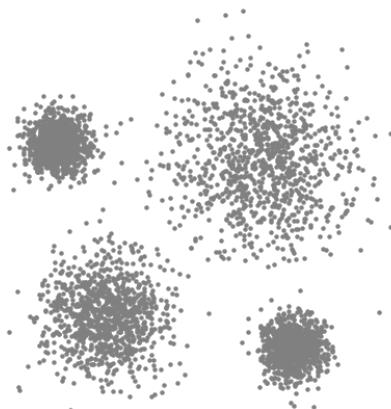
6

- **Algorithm**

1. Begin with a circular sliding window centered at a point C (randomly selected) and having radius r as the kernel.
2. At every iteration, the sliding window is shifted towards regions of higher density by shifting the center point to the mean of the points within the window (will gradually move towards areas of higher point density).
3. We continue shifting the sliding window according to the mean until there is no direction at which a shift can accommodate more points inside the kernel.
4. This process of steps 1 to 3 is done with many sliding windows until all points lie within a window. When multiple sliding windows overlap the window containing the most points is preserved. The data points are then clustered according to the sliding window in which they reside.



Mean-Shift Clustering



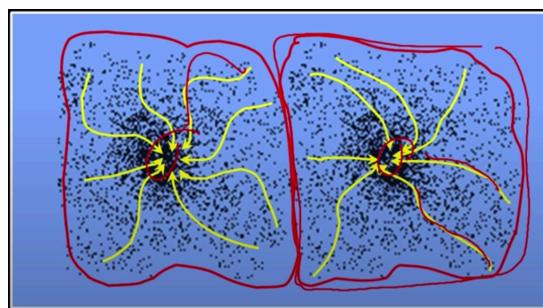
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

9

Mean-Shift Clustering

- **Characteristics**

- **Cluster:** all the data points in the attraction basin of a mode
- **Attraction basin:** the region where all trajectories lead to the same mode



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

10



Questions?

11



Data Mining

Notes on Clustering Strategies

24/11/2021

NOVA-IMS

Fernando Lucas Bação

bacao@isegi.unl.pt

<http://www.isegi.unl.pt/fbacao>

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

1



AGENDA

- Cluster analysis
 - Clustering Strategies
 - Reintroducing the Outliers or New Observations

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

2

NOVA
IMS
Information Management School



Clustering Strategies

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

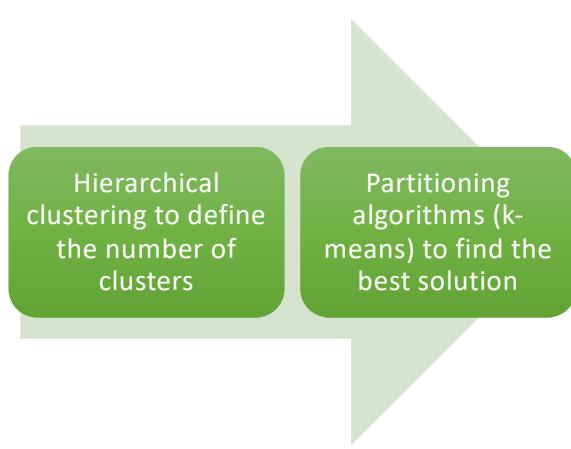


3

NOVA
IMS
Information Management School

Clustering

- Profiling (size of the clusters)



Hierarchical clustering to define the number of clusters

Partitioning algorithms (k-means) to find the best solution

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

4

NOVA
IMS
Information Management School

Clustering

- Profiling (size of the clusters)
 - Partitioning algorithm (k-means) with a large number of clusters
 - Hierarchical algorithm to find the appropriate number of clusters

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

5

NOVA
IMS
Information Management School

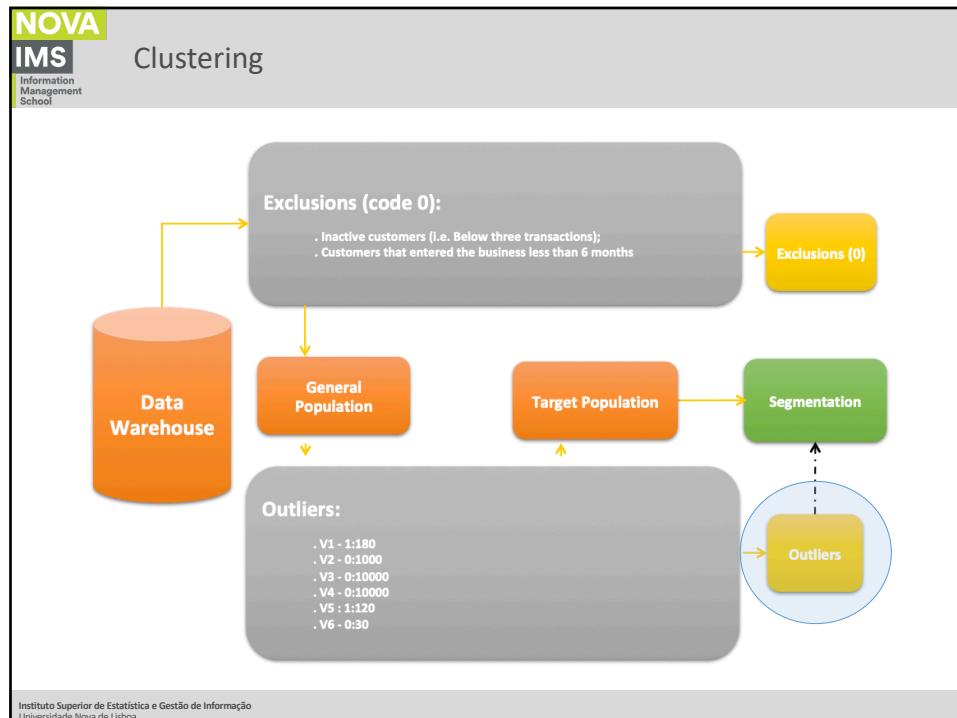


Reintroducing the Outliers or New Observations

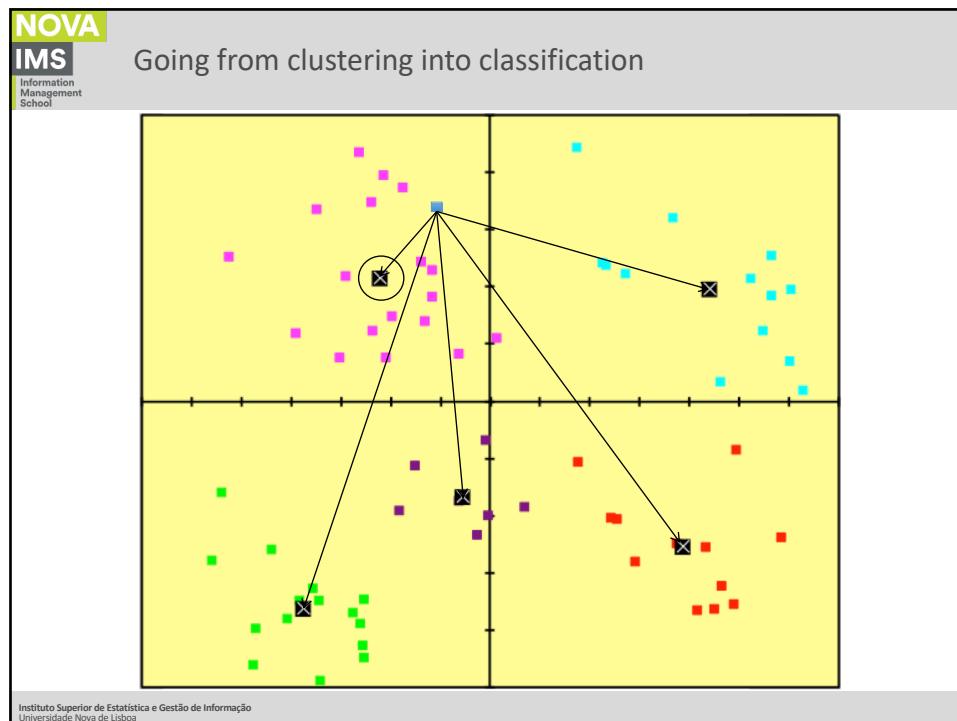
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

A3ES
UNIGIS
EQUIS
AACSB
SAC
Nursery of Schools
eduniversal

6

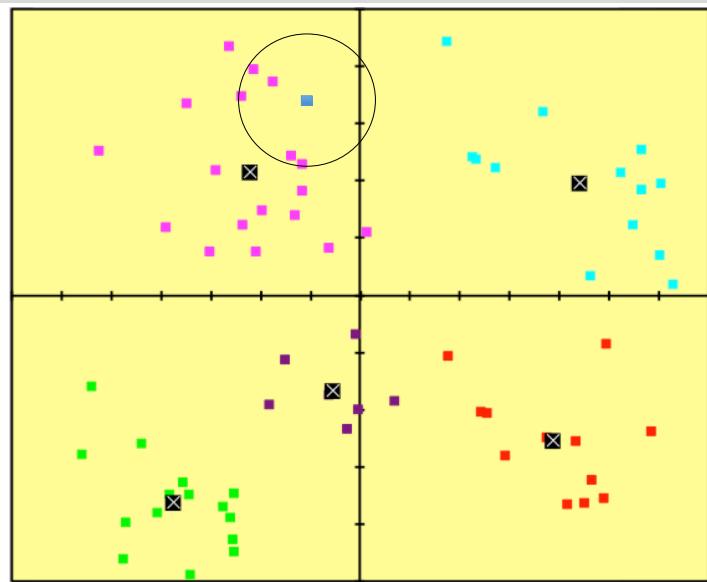


7



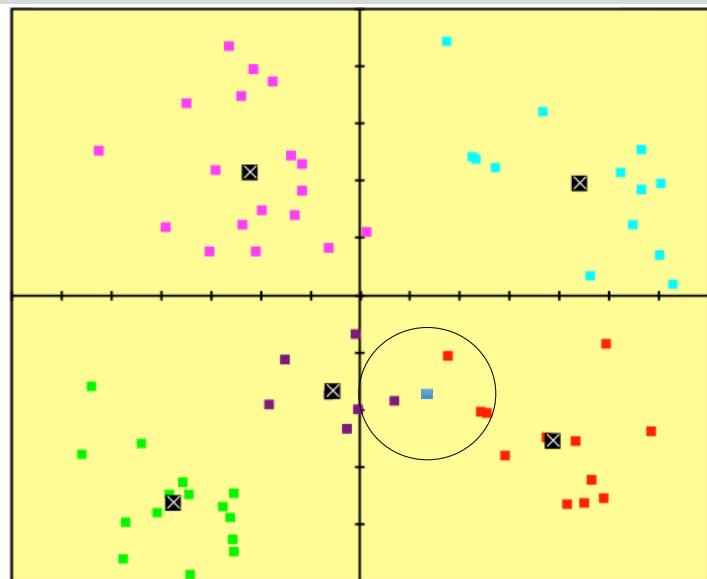
8

Going from clustering into classification

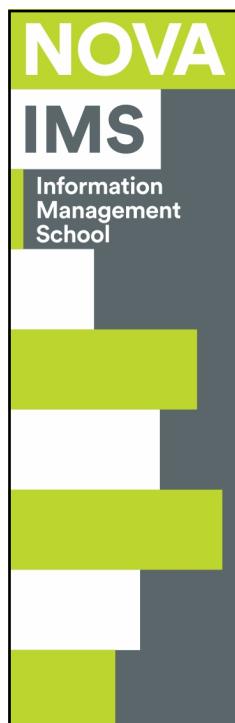
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

9

Going from clustering into classification

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

10



Data Mining

Semi-supervised classification

Nearest Neighbors

24/11/2021

NOVA-IMS

Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>
 Instituto Superior de Estatística e Gestão de Informação
 Universidade Nova de Lisboa

1



AGENDA

- Cluster analysis
 - Semi-supervised classification
 - Nearest Neighbors

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

2



Semi-supervised classification

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

A3ES
UNIGIS
AEGIS
ISchools
eduniversal

3



Going from clustering into classification

Unlabeled Data → Clustering - Labels → Classification – KNN and Trees

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

4

Going from clustering into classification



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

5

***k*-nearest neighbors (*k*-NN)**

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

6

k-nearest neighbors

- Instance based classification:
 - Simplest form of learning;
 - Training instances are searched for instances that most closely resembles new instance;
 - The instances themselves represent the knowledge;
 - Also called instance-based learning
 - Similarity function defines what's "learned"

k-nearest neighbors

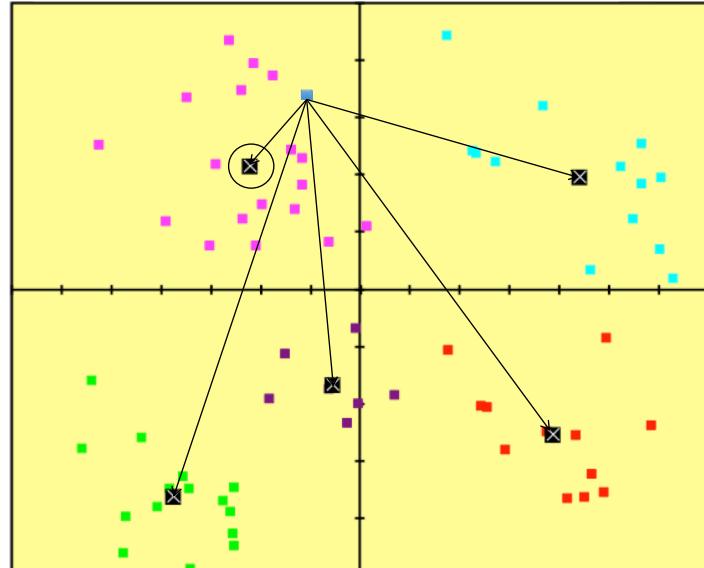
- Requires three things:
 1. The set of stored records (with labels)
 2. A distance metric to compute distance between records (can use Euclidean distance)
 3. The value of k, the number of nearest neighbors to retrieve

k-nearest neighbors

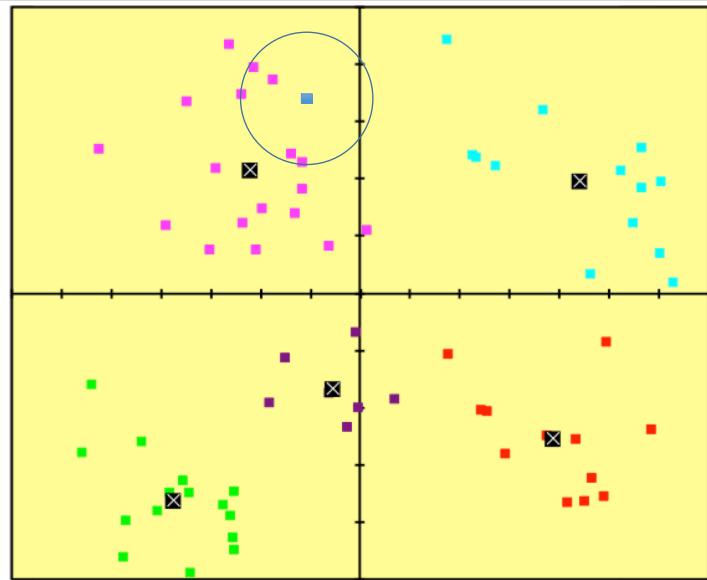
- To classify an unknown record:

1. Compute distance to other training records
2. Identify k nearest neighbors
3. Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

Going from clustering into classification

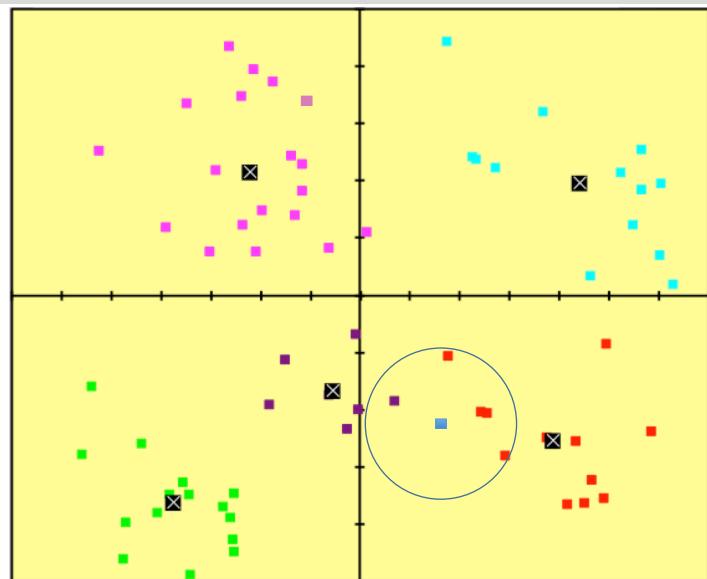


Going from clustering into classification

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

11

Going from clustering into classification

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

12

k-nearest neighbors

- Compute distance between two points:
 - Euclidean distance
- $$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$
- Determine the class from nearest neighbor list
 - Take the majority vote of class labels among the k-nearest neighbors
 - Weigh the vote according to distance
 - weight factor, $w = 1/d^2$

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

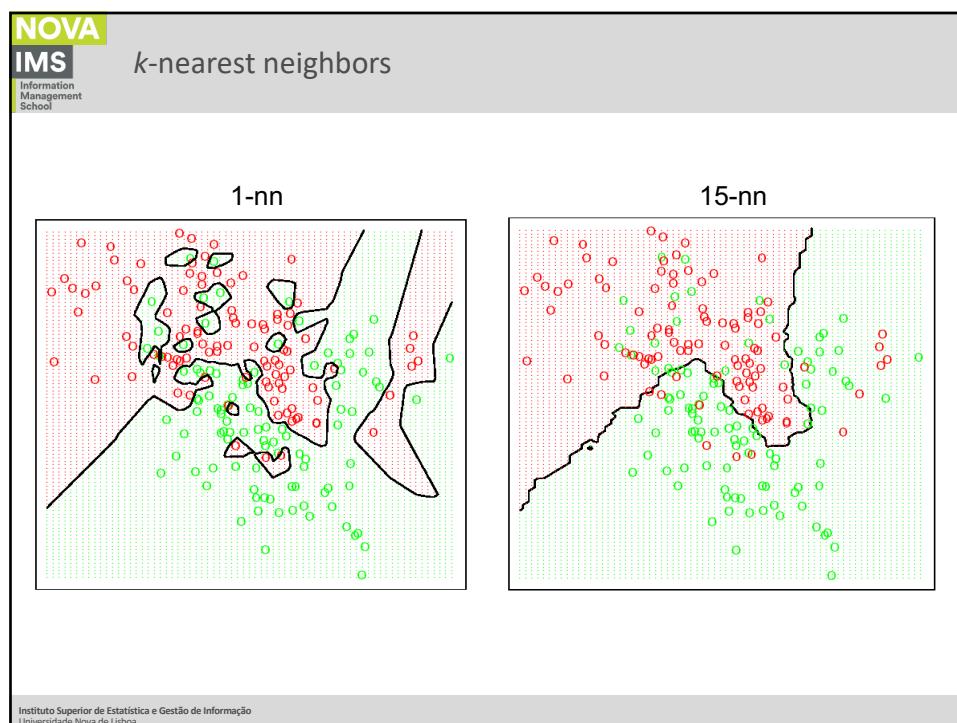
13

k-nearest neighbors

- k-nn frontiers (and the number k):
 - Large k
 - Smooth frontiers
 - Unable to detect small variations
 - Small k
 - Very sensitive to outliers
 - Crisp frontiers

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

14



15



Data Mining

Semi-supervised classification

Classification trees

24/11/2021

NOVA-IMS

Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>
 Instituto Superior de Estatística e Gestão de Informação
 Universidade Nova de Lisboa

1



AGENDA

- Cluster analysis
 - Semi-supervised classification
 - Classification Trees

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

2

NOVA
IMS
Information Management School



Classification Trees

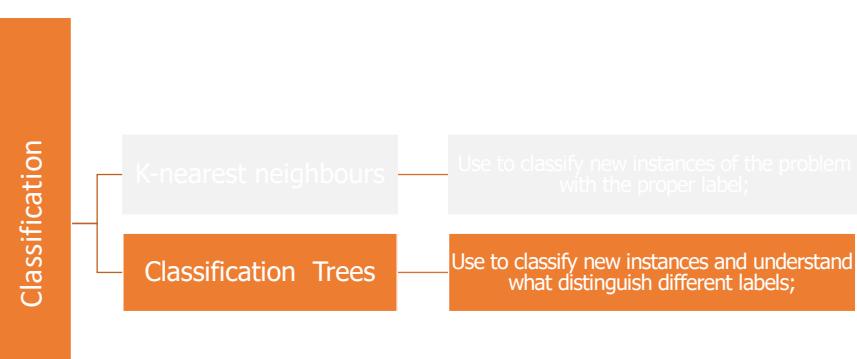
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa



3

NOVA
IMS
Information Management School

Going from clustering into classification



Classification

- K-nearest neighbours
- Classification Trees
 - Use to classify new instances of the problem with the proper label;
 - Use to classify new instances and understand what distinguish different labels;

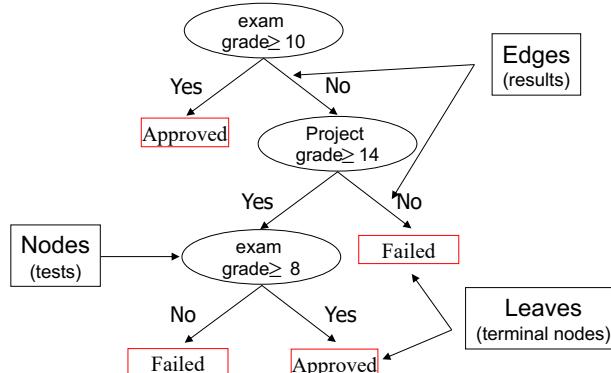
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

4

- **Classification Trees:**

- Classification tree are typically considered to be classification and regression tools
- One of its most important advantages relates with the simplicity of the interpretation of its results
- Thus, the end result of a classification tree can easily be expressed in English or SQL.

- A classification tree is a decisional algorithm
- It can be seen as a way of storing knowledge
- The objective is to discriminate between Class
- Obtain leaves as pure as possible
- If possible each leave should represent only individuals from a specific class



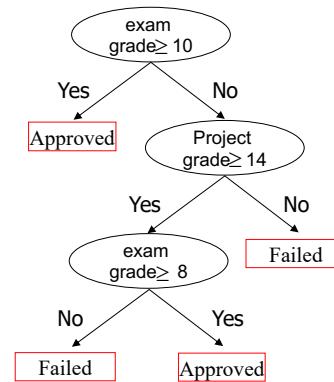
Classification Trees

$$aprovado \Leftrightarrow (exame \geq 10)$$

$$aprovado \Leftrightarrow (exame < 10) \wedge (projeto \geq 14) \wedge (exame \geq 8)$$

$$reprovado \Leftrightarrow (exame < 10) \wedge (projeto < 14)$$

$$reprovado \Leftrightarrow (exame < 10) \wedge (projeto > 14) \wedge (exame < 8)$$



Classification Trees

- **Classification Trees (strengths):**

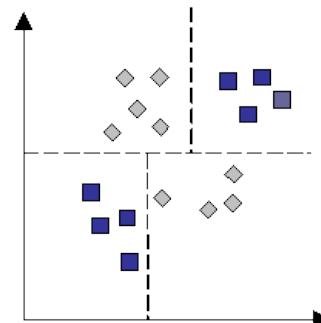
- Interpretation
 - We can easily understand the reasons behind a specific classification decision
- May use different types of data
 - Interval, ordinal, nominal, etc.
- Insensitive to scale factors
 - Variables measured in different scales may be used without any type of normalization

- **Classification Trees (strengths):**

- Automatically defines the most relevant variables
 - These are the variables used at the top of the tree
- Can be adapted to a regression
 - Each leave becomes a linear model

- **Classification Trees (weaknesses):**

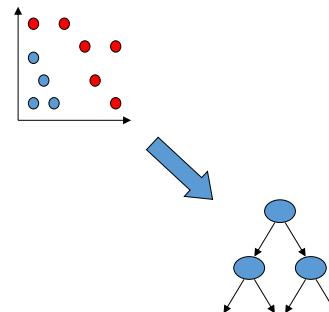
- Boundaries are linear and perpendicular to the variables axys
- Sensitive to small perturbations in the data



From Gahegan and West
http://divcom.otago.ac.nz/SIRC/GeoComp98/61/gc_61.htm

- **Classification Trees:**

- Build (induce) a tree from data
- Problems:
 - What to do?
 - Which variable to use?
 - What partition to use?
 - Which node to split?
 - How many edges per node?
 - When to stop?



- **Classification Trees:**

- ID3, C4.5 e C5 [Quinlan 86,93]
 - Iterative Dichotomizer 3
- CART
 - Classification and regression trees [Breiman 84]
- CHAID [Hartigan 75]
 - Used in SPSS and SAS...
- Muitas (mesmo muitas) outras variantes...
 - In SAS you can choose different parameters to build your tree.

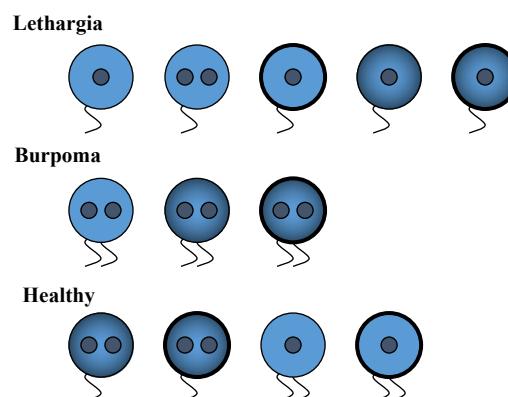
Worked-Example

General Idea

Langley, P: 1996, Elements of Machine Learning, Morgan and Kaufmann Publishers.

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

13



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

14

NOVA
IMS
Information Management School

Classification Trees

Table

# Nucleus	# Tails	Color	Membrane	<i>Class</i>
1	1	Light	Thin	<i>Lethargia</i>
2	1	Light	Thin	<i>Lethargia</i>
1	1	Light	Thick	<i>Lethargia</i>
1	1	Dark	Thin	<i>Lethargia</i>
1	1	Dark	Thick	<i>Lethargia</i>
2	2	Light	Thin	<i>Burpoma</i>
2	2	Dark	Thin	<i>Burpoma</i>
2	2	Dark	Thick	<i>Burpoma</i>
2	1	Dark	Thin	<i>Healthy</i>
2	1	Dark	Thick	<i>Healthy</i>
1	2	Light	Thin	<i>Healthy</i>
1	2	Light	Thick	<i>Healthy</i>

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

15

NOVA
IMS
Information Management School

Classification Trees

- **Classification Trees:**
 - Measure to discriminate the attribute

$$f(A) = \frac{1}{n} \sum_{i=1}^{|A|} C_i$$

- n is the total number of examples and C_i the number of examples correctly classified based on the most frequent class.
- This is a measure of “dominance” or “purity”

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

16

NOVA
IMS
Information Management School

Classification Trees

Table		# Nucleus	# Tails	Color	Membrane	Class
# Nucleus	1	1	Light	Thin	Lethargia	
Lethargia	4	1	Light	Thin	Lethargia	
Burpoma	0	3	Dark	Thin	Lethargia	
Healthy	2	2	Dark	Thin	Lethargia	
		1	1	Light	Thin	Lethargia
		2	2	Light	Thin	Burpoma
		2	2	Dark	Thin	Burpoma
		2	1	Dark	Thin	Healthy
		2	1	Dark	Thick	Healthy
		1	2	Light	Thin	Healthy
		1	2	Dark	Thick	Healthy

Discrimination:
 $(4 + 3) / 12 = 0.58$

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

17

NOVA
IMS
Information Management School

Classification Trees

Table		# Nucleus	# Tails	Color	Membrane	Class
# Tails	1	1	Light	Thin	Lethargia	
Lethargia	5	0	Light	Thin	Lethargia	
Burpoma	0	3	Light	Thin	Lethargia	
Healthy	2	2	Dark	Thin	Lethargia	
		1	1	Light	Thin	Lethargia
		2	2	Light	Thin	Burpoma
		2	2	Dark	Thin	Burpoma
		2	1	Dark	Thin	Healthy
		2	1	Dark	Thick	Healthy
		1	2	Light	Thin	Healthy
		1	2	Light	Thick	Healthy

Discrimination:
 $(5 + 3) / 12 = 0.67$

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

18

NOVA
IMS
Information Management School

Classification Trees

Table

Color	Light	Dark
Lethargia	3	2
Burpoma	1	2
Healthy	2	2

Discrimination:
 $(3 + 2) / 12 = 0.41$

# Nucleus	# Tails	Color	Membrane	Class
1	1	Light	Thin	Lethargia
2	1	Light	Thin	Lethargia
1	1	Light	Thick	Lethargia
1	1	Dark	Thin	Lethargia
1	1	Dark	Thick	Lethargia
2	2	Light	Thin	Burpoma
2	2	Dark	Thin	Burpoma
2	1	Dark	Thin	Healthy
2	1	Dark	Thick	Healthy
1	2	Light	Thin	Healthy
1	2	Light	Thick	Healthy

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

19

NOVA
IMS
Information Management School

Classification Trees

Table

Membrane	Thin	Thick
Lethargia	3	2
Burpoma	2	1
Healthy	3	1

Discrimination:
 $(3 + 2) / 12 = 0.41$

# Nucleus	# Tails	Color	Membrane	Class
1	1	Light	Thin	Lethargia
2	1	Light	Thin	Lethargia
1	1	Light	Thick	Lethargia
1	1	Dark	Thin	Lethargia
1	1	Dark	Thick	Lethargia
2	2	Light	Thin	Burpoma
2	2	Dark	Thin	Burpoma
2	1	Dark	Thin	Healthy
2	1	Dark	Thick	Healthy
1	2	Light	Thin	Healthy
1	2	Light	Thick	Healthy

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

20

NOVA
IMS
Information Management School

Classification Trees

Choice: # Tails

# Nucleus	1	2
Lethargia	4	1
Burpoma	0	3
Healthy	2	2

0.58

Color	Light	Dark
Lethargia	3	2
Burpoma	1	2
Healthy	2	2

0.41

# Tails	1	2
Lethargia	5	0
Burpoma	0	3
Healthy	2	2

0.67

Membrane	Thin	Thick
Lethargia	3	2
Burpoma	2	1
Healthy	3	1

0.41

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

21

NOVA
IMS
Information Management School

Classification Trees

Initial Partition

```

graph TD
    Tails((Tails)) -- one --> TableOne
    Tails -- two --> TableTwo
    
```

# Nucleus	Color	Membrane	Class
1	Light	Thin	Lethargia
2	Light	Thin	Lethargia
1	Light	Thick	Lethargia
1	Dark	Thin	Lethargia
1	Dark	Thick	Lethargia
2	Dark	Thin	Healthy
2	Dark	Thick	Healthy

# Nucleus	Color	Membrane	Class
2	Light	Thin	Burpoma
2	Dark	Thin	Burpoma
2	Dark	Thick	Burpoma
1	Light	Thin	Healthy
1	Light	Thick	Healthy

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

22

NOVA
IMS
Information Management School

Classification Trees

Second Partition

```

graph TD
    Tails((Tails)) -- one --> Nucleus((Nucleus))
    Tails -- two --> Data2[Data]
    Nucleus -- one --> Data1[Data]
    Nucleus -- two --> Data3[Data]
  
```

# Nucleus	Color	Membrane	Class
2	Light	Thin	Burpoma
2	Dark	Thin	Burpoma
2	Dark	Thick	Burpoma
1	Light	Thin	Healthy
1	Light	Thick	Healthy

Color	Membrane	Class
Light	Thin	Lethargia
Light	Thick	Lethargia
Dark	Thin	Lethargia
Dark	Thick	Lethargia

Color	Membrane	Class
Light	Thin	Lethargia
Dark	Thin	Healthy
Dark	Thick	Healthy

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

23

NOVA
IMS
Information Management School

Classification Trees

Tree (cont.2)

```

graph TD
    Tails((Tails)) -- one --> Nucleus((Nucleus))
    Tails -- two --> Data2[Data]
    Nucleus -- one --> Lethargia4[Lethargia (4)]
    Nucleus -- two --> Data3[Data]
  
```

# Nucleus	Color	Membrane	Class
2	Light	Thin	Burpoma
2	Dark	Thin	Burpoma
2	Dark	Thick	Burpoma
1	Light	Thin	Healthy
1	Light	Thick	Healthy

Color	Membrane	Class
Light	Thin	Lethargia
Dark	Thin	Healthy
Dark	Thick	Healthy

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

24

NOVA
IMS
Information Management School

Classification Trees

Tree (cont.3)

# Nucleus	Color	Membrane	Class
2	Light	Thin	Burpoma
2	Dark	Thin	Burpoma
2	Dark	Thick	Burpoma
1	Light	Thin	Healthy
1	Light	Thick	Healthy

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

25

NOVA
IMS
Information Management School

Classification Trees

Final Tree

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

26

The description of the three Classss

$$(tails = 1) \wedge (nucleous = 1) \\ \vee \\ (tails = 1) \wedge (nucleous = 2) \wedge (color = light) \rightarrow Lethargic$$

$$(tails = 2) \wedge (nucleous = 1) \\ \vee \\ (tails = 1) \wedge (nucleous = 2) \wedge (color = dark) \rightarrow Healthy$$

$$(tails = 2) \wedge (nucleous = 2) \rightarrow Burpoma$$

- **Classification Trees:**
 - In each level it divides the set into alternative partitions.
 - Using a measure of quality selects the best partition.
 - The process is repeated for each element of the partition.
 - Stops when a given criteria is reached

- **Classification Trees:**

- It assumes the existence of a target variable “Class” meaning the examples were previously classified.
- Each node specifies a unique attribute which is used as test.
- N – node N
- ASET – Atribute Set
- ISET – Instance Set

If Se the ISET is empty then the terminal node N is an unknown class
if not

If all the examples of ISET are of the same class
then the terminal node N has the name of the class
if not

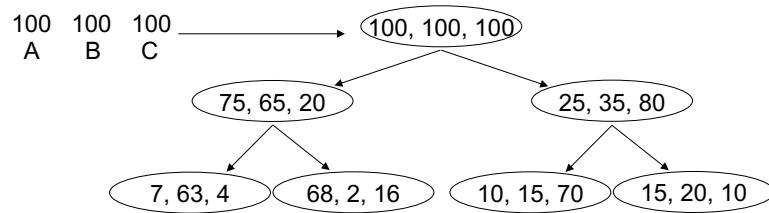
For each attribute A of the set of attribute ASET
Evaluate A according to its capability to discriminate a class
Select the attribute B which has the best discriminate value
For each value V of the best attribute B
Create a new node C from node N
Place the par attribute value (B, V) in C
Let JSET be the set of examples of ISET with value V in B
Let KSET be the set of attributes of ASET with B removed
DDT(C, KSET, JSET)

Worked-Example

Tree Accuracy

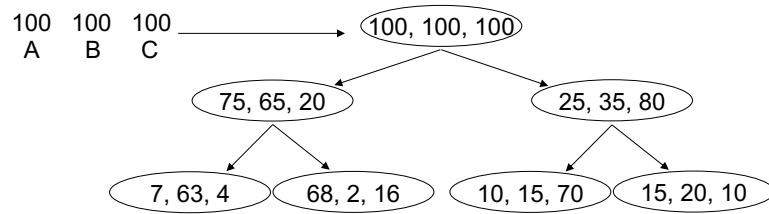
Quality of the results

Error rate



Quality of the results

Error rate



Error Calculation

$$TEA = \frac{(74 * 14,9\% + 86 * 20,9\% + 95 * 26,3\% + 45 * 55,5\%)}{300} = 26,3\%$$

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



Data Mining

Self-Organizing Maps

24/11/2021

NOVA-IMS

Fernando Lucas Bação

bacao@isegi.unl.pt

<http://www.isegi.unl.pt/fbacao>

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

1



AGENDA

- Cluster analysis
 - Clustering techniques
 - Self-organizing maps

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

2

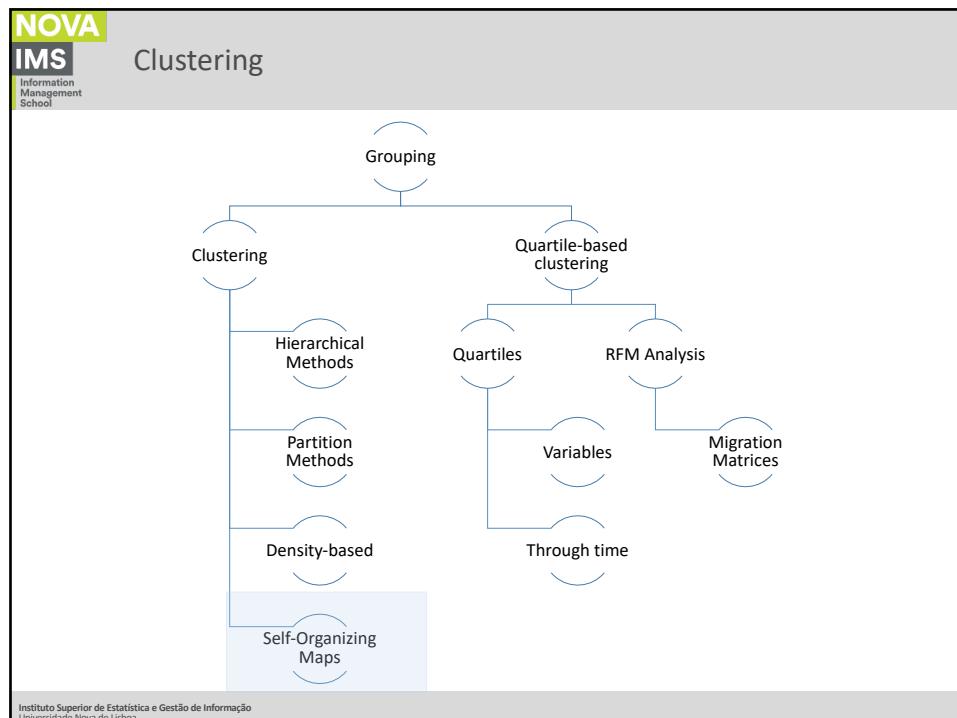
NOVA
IMS
Information Management School

Clustering

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

A3ES
UNIGIS
AEGIS
ISchools
eduniversal

3



4



Self-Organizing Maps

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



5



Self-Organizing Maps

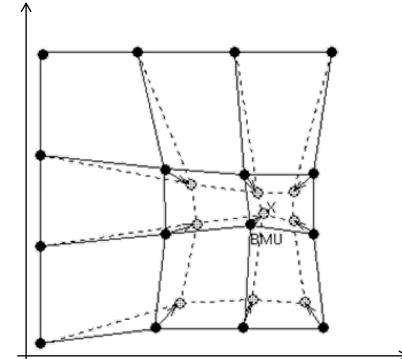
- Unsupervised neural networks;
- Closely related to clustering;
- The inputs are connected to a two-dimensional (it may have several dimensions) matrix of units (neurons);
- Each unit is connected to its neighbors.
- What is its use?
 - Multidimensional data visualization;
 - Cluster detection;
 - Market segmentation;
 - Outlier detection;
 - Solve TSP, robot control, alarm detection, etc., etc., etc.

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

6

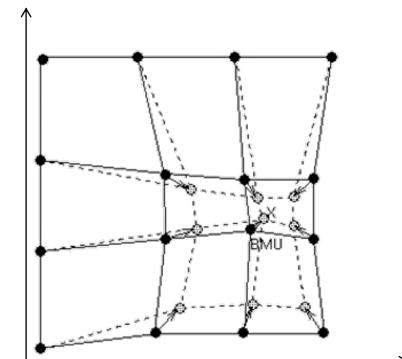
Self-Organizing Maps

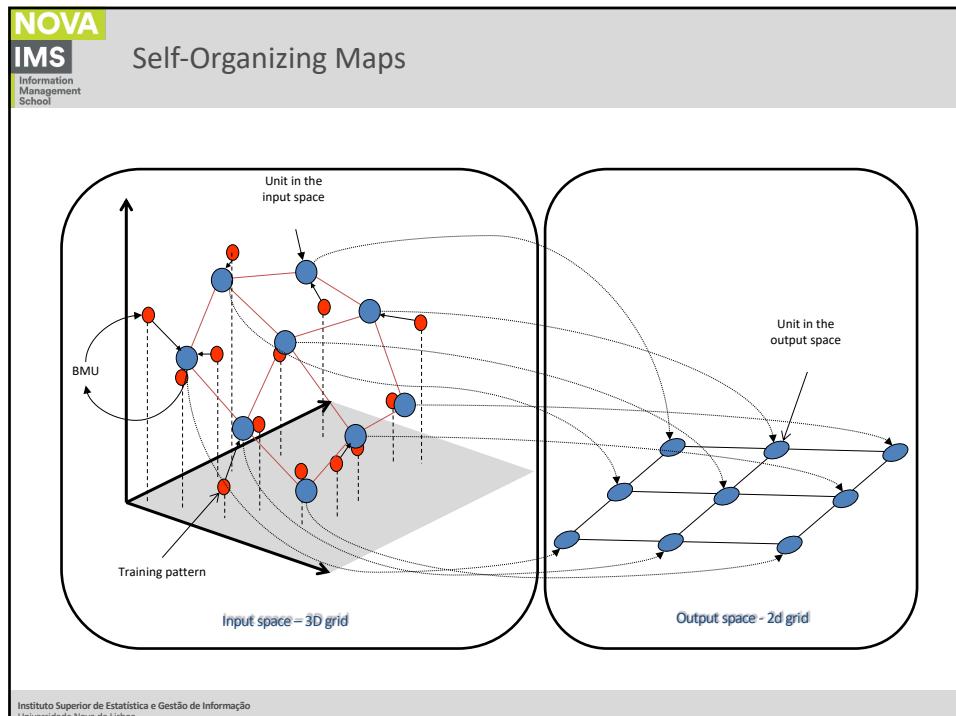
- Each neuron is a **vector in the input space**, just as the data patterns;
- During training, neurons are **pulled** to the positions of the input data, **dragging** with them their neighbors in the output space;
- The map can be seen as a **rubber sheet**, stretched and twisted, so that it passes in (or at least near) the data patterns.



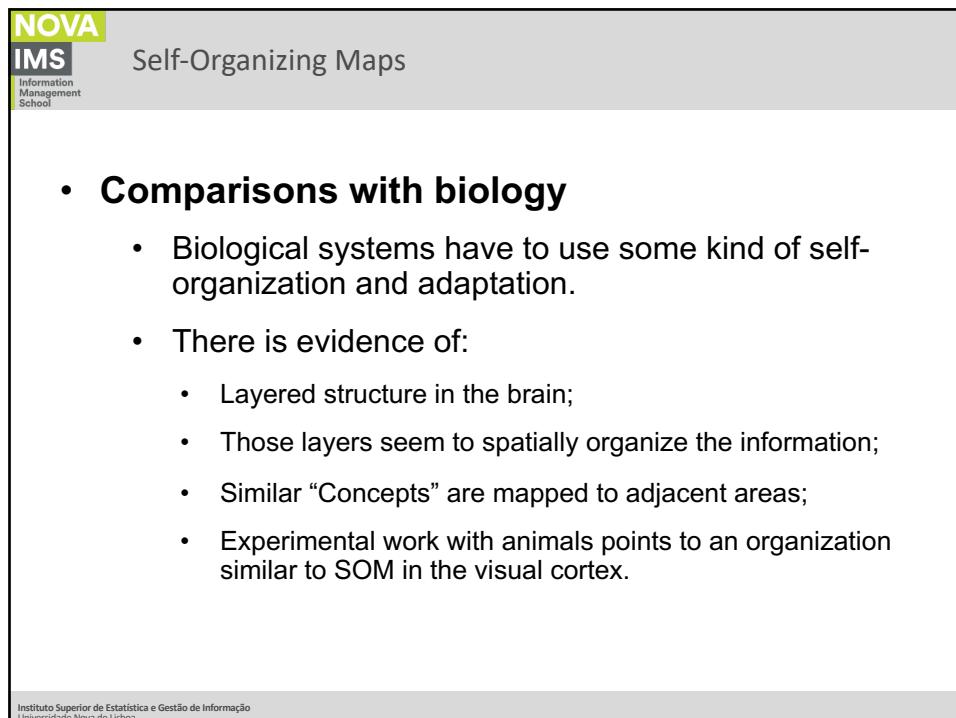
Self-Organizing Maps

- Input patterns are compared with all neurons and the **closest is considered to be the winning neuron**.
- We consider that the input pattern is **mapped to the winning neuron**.
- The **winner is updated** (so that it resembles even more the data pattern that it represents), and its neighbors are also updated a little.
- There is always a slight difference between the data and the neurons that represent them. That difference is the **quantization error**.





9



10

The SOM Algorithm

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

Accredited by EQUIS, UNIGIS, A3ES, AACSB, iSchools, eduniversal

11

Self-Organizing Maps

- **SOM Algorithm:**
 - How does SOM processes data?

Video

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

12

- **SOM Algorithm:**

- What happens in the Input space?
- If we only have two variables, optimization will look like this

SOM 1-dimensional

SOM 2-dimensional

- **SOM Algorithm:**

- What happens in the output space

Color Demo

- Suppose we want to group cells according to their rgb code (red, green and blue)
- Each individual (color) represents a particular combination of rgb intensities
- In this demo, we can see how the output space is transformed as individuals are presented to the network

Self-Organizing Maps

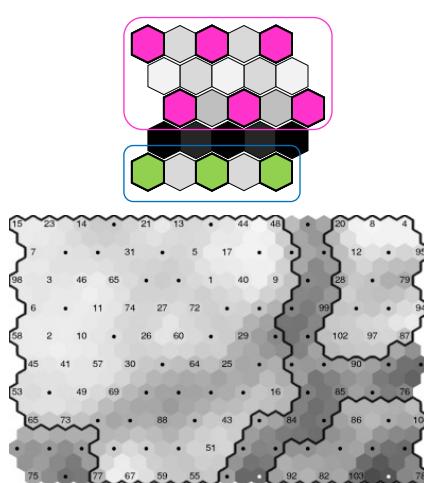
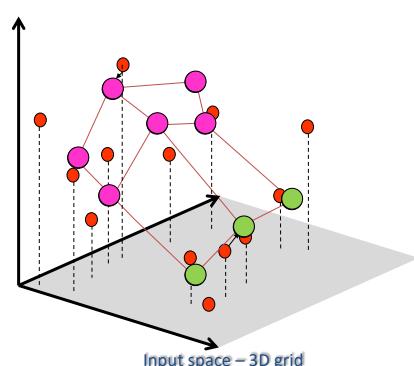
- **SOM Algorithm:**

- Key outputs of the SOM:
 - U-Matrices;
 - Component plans;
 - Hit Plots.

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

15

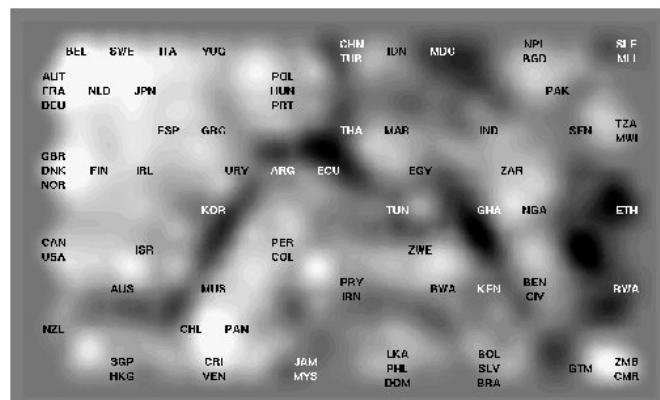
Self-Organizing Maps – U-Matrix



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

16

Self-Organizing Maps – U-Matrix

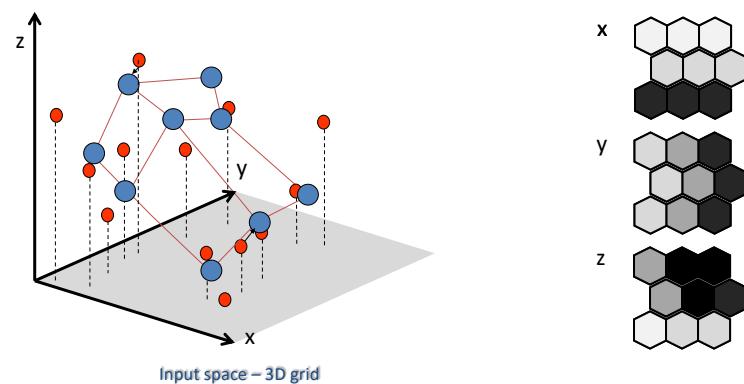


Kaski, S. and Kohonen, T. 1996. Exploratory Data Analysis by the Self-Organizing Map: Structures of Welfare and Poverty in the World. In: Apostolos, P. N., Refenes, Y. A., Moody, J. and Weigend, A. (eds.) Neural Networks in Financial Engineering. Singapore: World Scientific, 498-507.

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

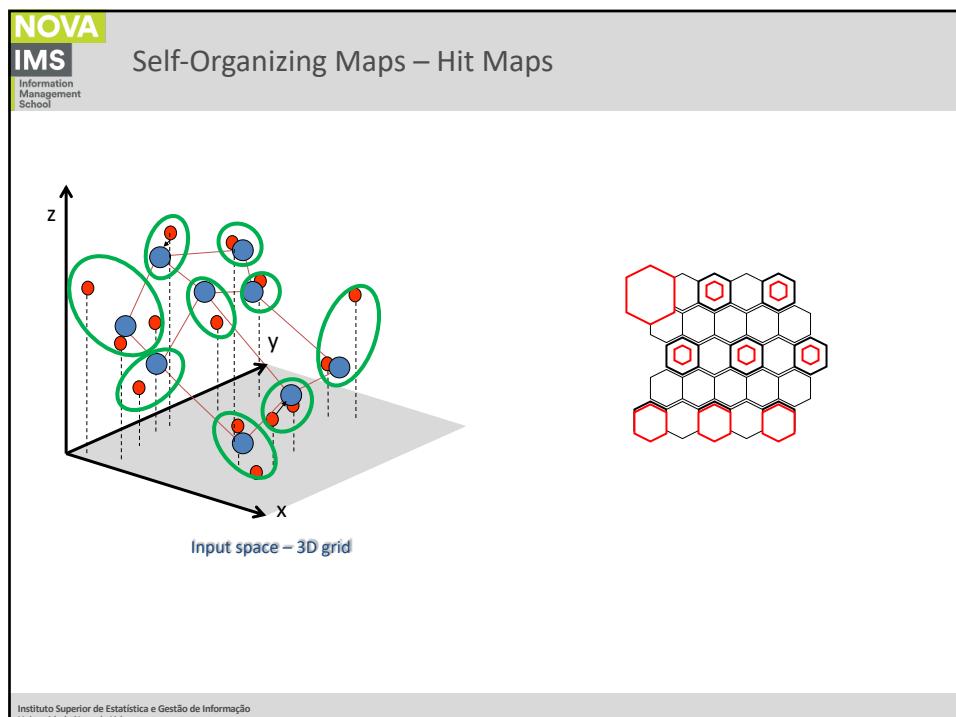
17

Self-Organizing Maps – Component Planes

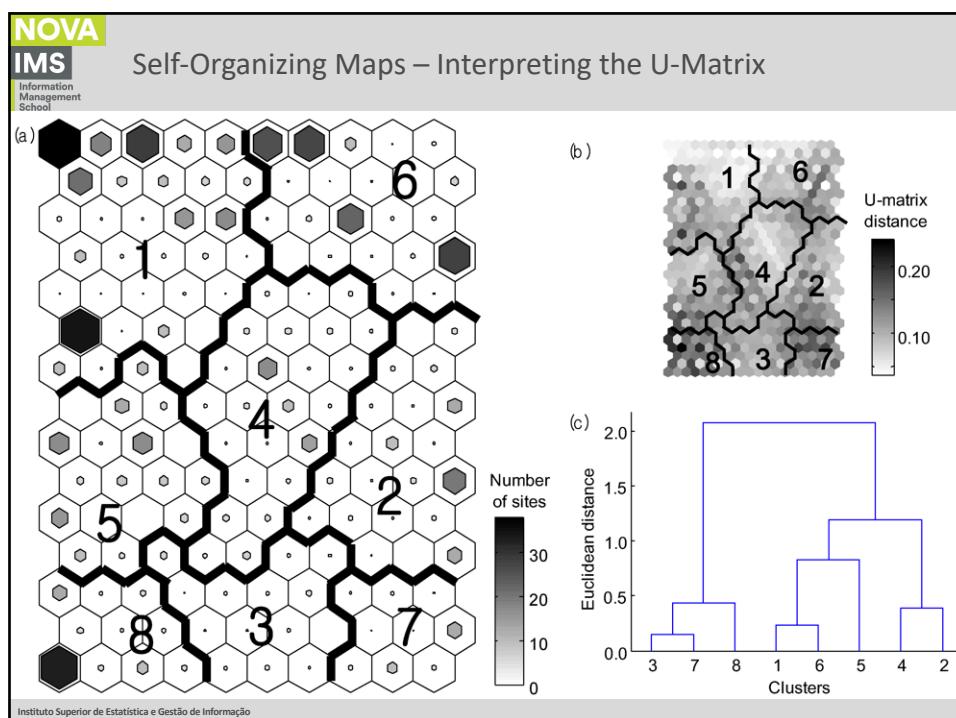


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

18



19



20

- **SOM Algorithm:**

Step 0: Randomly initialize the weights w_{ij}
Set the neighborhood topological parameters
Set the learning rate

Step 1: While stop condition false, do steps 2-7

Step 2: For each input vector x , do steps 3-5
Step 3: For each j , execute:

$$D(j) = \sum (w_{ij} - x_j)^2$$

Step 4: Find the unit that minimizes $D(j)$

Step 5: For every j unit within the predefined and for all the i :

$$w_{ij} (\text{new}) = w_{ij} (\text{old}) + \alpha [x_i - w_{ij} (\text{old})]$$

Step 6: Update the learning rate

Step 7: Update (reduce) the radius of the topological neighborhood

- **SOM Algorithm:**

Vectors to classify:

(1, 1, 0, 0); (0, 0, 0, 1); (1, 0, 0, 0); (0, 0, 1, 1)

Maximum number of clusters to form:

$m = 2$

Learning rate:

$$\alpha(0) = .6, \quad \alpha(t+1) = .5 \alpha(t)$$

- **SOM Algorithm:**

Step 0: Initial matrix of weights

(0.2, 0.6, 0.5, 0.9);

(0.8, 0.4, 0.7, 0.3);

Initial radius: R = 0

Initial learning rate:

$$\alpha(0) = .6$$

- **SOM Algorithm:**

Step 1: Initialize training

Step 2: first vector (1, 1, 0, 0);

Step 3:

$$D(1) = (.2-1)^2 + (.6-1)^2 + (.5-0)^2 + (.9-0)^2$$

$$= 1.86;$$

$$D(2) = (.8-1)^2 + (.4-1)^2 + (.7-0)^2 + (.3-0)^2$$

$$= 0.98$$

Vectors of weights:

(0.2, 0.6, 0.5, 0.9)

(0.8, 0.4, 0.7, 0.3)

- **SOM Algorithm:**

Step 4: The input vector is closest to unit 2, therefore

$$j = 2 = (0.8, 0.4, 0.7, 0.3)$$

Step 5: The weights of the winning unit are adjusted

$$\begin{aligned} w_{i2} (\text{new}) &= w_{i2} (\text{old}) + .6 [x_i - w_{i2} (\text{old})] \\ &= .4 w_{i2} (\text{old}) + .6x_i \end{aligned}$$

- **SOM Algorithm:**

Step 4: The input vector is closest to unit 2, therefore

$$j = 2 = (0.8, 0.4, 0.7, 0.3)$$

Step 5: The weights of the winning unit are adjusted

$$\begin{aligned} w_{i2} (\text{new}) &= w_{i2} (\text{old}) + .6 [x_i - w_{i2} (\text{old})] \\ &= .4 w_{i2} (\text{old}) + .6x_i \end{aligned}$$

- **SOM Algorithm:**

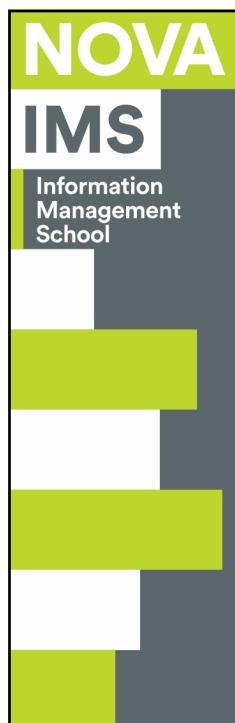
Thus, the weight matrix is adjusted:

.2	.92
.6	.76
.5	.28
.9	.12

- **SOM Algorithm:**

Thus, the weight matrix is adjusted:

.08	.92
.24	.76
.20	.28
.96	.12



Data Mining

Visualization of Multidimensional Data

2/12/2021

NOVA-IMS

Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>

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

1

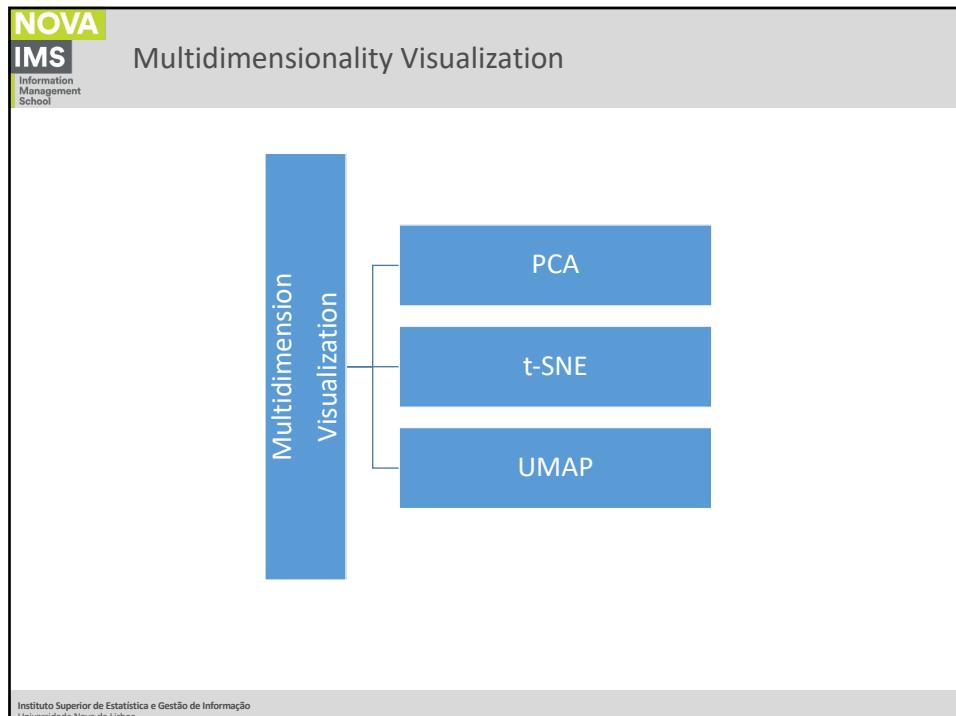


Multidimensionality Visualization

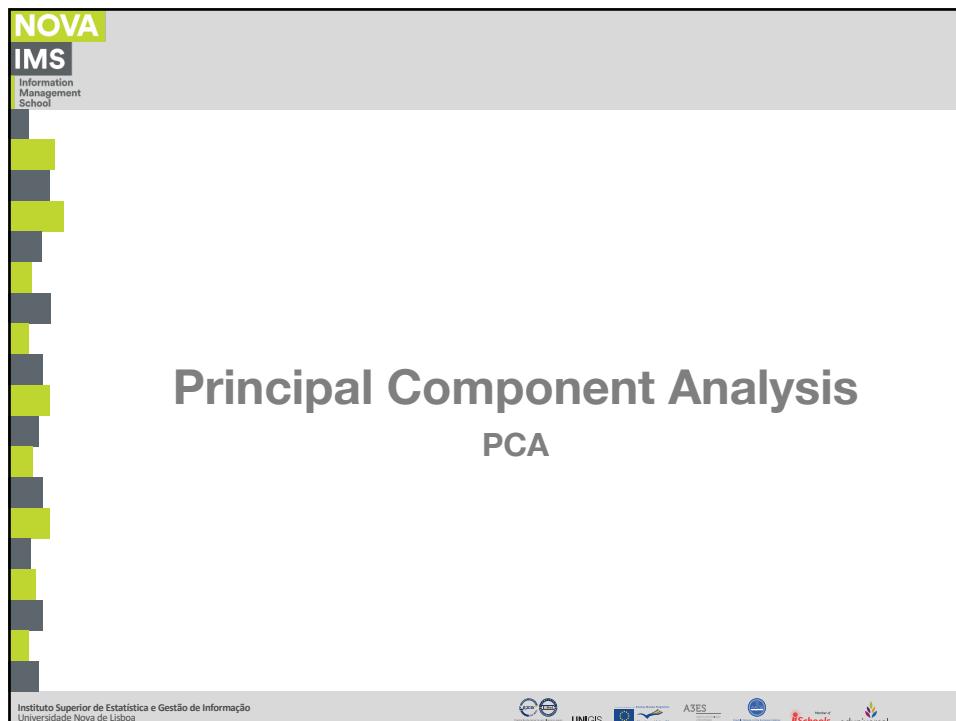
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

AQUIS
UNIGIS
A3ES
e-Schools
eduniversal

2



3



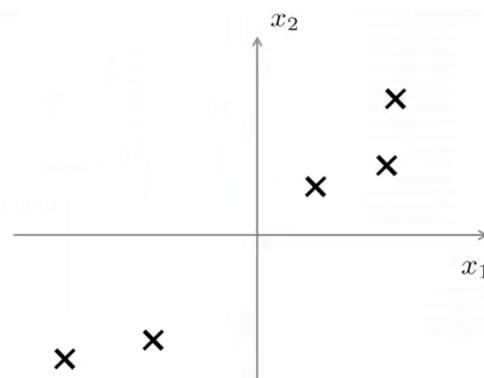
4

- **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.

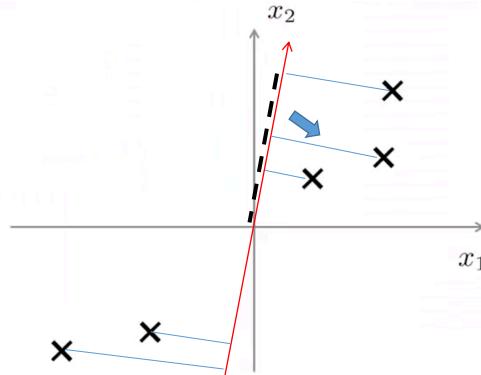
- **Size Reduction of the Input Space:**

- Principal Component Analysis



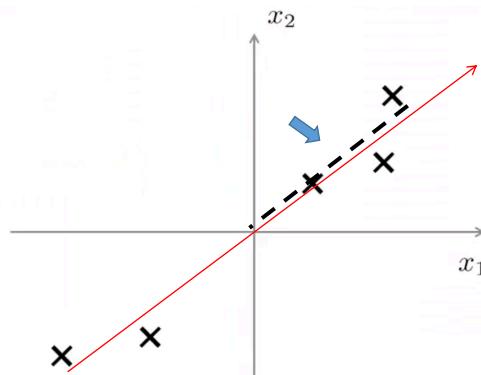
- **Size Reduction of the Input Space:**

- Principal Component Analysis (SVD)



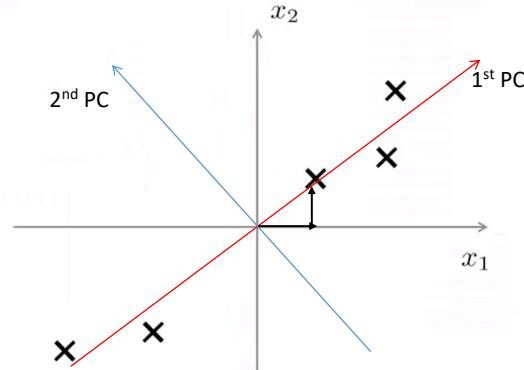
- **Size Reduction of the Input Space:**

- Principal Component Analysis (SVD)



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

- Principal Component Analysis (SVD)

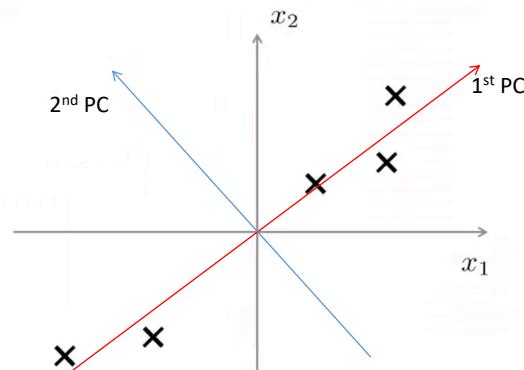


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

9

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

- Principal Component Analysis (SVD)

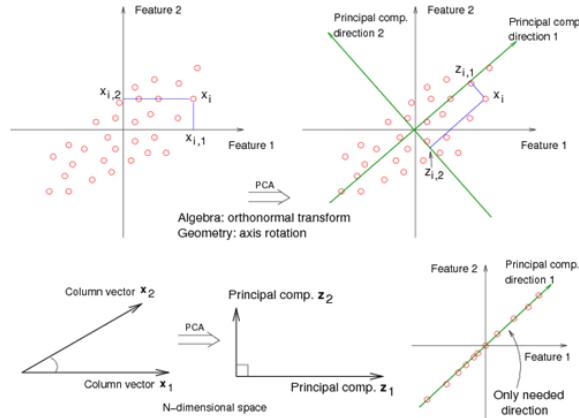


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

10

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

- Principal Component Analysis

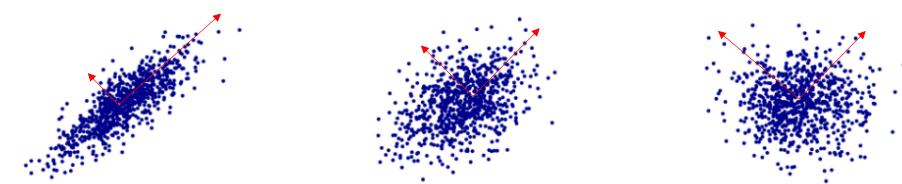


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

11

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

- Principal Component Analysis



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

12

NOVA
IMS
Information Management School

t-distributed Stochastic Neighbor Embedding

t-SNE

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

EQUIS UNIGIS A3ES AACSB iSchools eduniversal

13

NOVA
IMS
Information Management School

t-SNE

- **t-SNE:**
 - t-distributed stochastic neighbor embedding (t-SNE) is a statistical method for visualizing high-dimensional data by **giving each datapoint a location in a two or three-dimensional map.**
 - It is a **nonlinear dimensionality reduction technique** well-suited for embedding high-dimensional data for visualization in a low-dimensional space of two or three dimensions.
 - Specifically, it **models each high-dimensional object** by a two- or three-dimensional point in such a way that **similar objects are modeled by nearby points** and dissimilar objects are modeled by distant points with high probability.

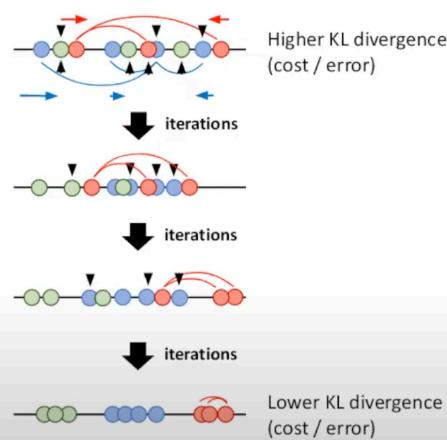
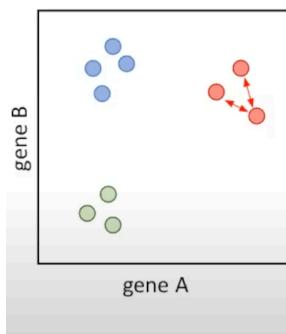
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

14

- **t-SNE:**

- The t-SNE algorithm comprises two stages:
 - First, t-SNE constructs a **probability distribution over pairs of high-dimensional objects** in such a way that similar objects are assigned a higher probability while dissimilar points are assigned a lower probability.
 - Second, t-SNE defines a **similar probability distribution over the points in the low-dimensional map**, and it minimizes the Kullback–Leibler divergence (KL divergence) between the two distributions with respect to the locations of the points in the map.
- While the original algorithm uses the Euclidean distance between objects as the base of its similarity metric, this can be changed as appropriate.

- **t-SNE:**



NOVA
IMS
Information Management School

t-SNE

- **t-SNE:**
 - Optimizing the KL divergence is a measure of how one probability distribution is different from a second, reference probability distribution

Normal distribution vs. t-distribution

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

17

NOVA
IMS
Information Management School

t-SNE

- **t-SNE:**
 - t-SNE has been used for **visualization in a wide range of applications**, including genomics, computer security research, natural language processing, music analysis, cancer research, bioinformatics, etc
 - While t-SNE plots often seem to display clusters, **the visual clusters can be influenced strongly by the chosen parameterization** and therefore a good understanding of the parameters for t-SNE is necessary.
 - Interactive exploration may thus be necessary to choose parameters and validate results.

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

18

**NOVA
IMS**
Information Management School

t-SNE

- **t-SNE:**

 - Perplexity is the main t-SNE parameter, “perplexity,”
 - Basically defines how to balance attention between local and global aspects of your data.
 - Perplexity is a measure for information that is defined as 2 to the power of the Shannon entropy.
 - In t-SNE, the perplexity may be viewed as a knob that sets the number of effective nearest neighbors.
 - It is comparable with the number of nearest neighbors k that is employed in many manifold learners.
 - The original paper says, “The performance of SNE is fairly robust to changes in the perplexity, and typical values are between 5 and 50.”

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

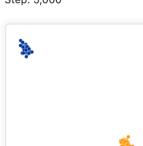
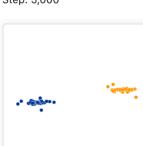
19

**NOVA
IMS**
Information Management School

t-SNE

- **t-SNE:**

 - Perplexity is the main t-SNE parameter, “perplexity,”

					
Original	Perplexity: 2 Step: 5,000	Perplexity: 5 Step: 5,000	Perplexity: 30 Step: 5,000	Perplexity: 50 Step: 5,000	Perplexity: 100 Step: 5,000
					
Original	Perplexity: 30 Step: 10	Perplexity: 30 Step: 20	Perplexity: 30 Step: 60	Perplexity: 30 Step: 120	Perplexity: 30 Step: 1,000

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

20

**NOVA
IMS**
Information Management School

t-SNE

- **t-SNE:**
 - Perplexity is the main t-SNE parameter, “perplexity,”

Perplexity	Step
Original	
2	5,000
5	5,000
30	5,000
50	5,000
100	5,000

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

21

**NOVA
IMS**
Information Management School

Uniform Manifold Approximation and Projection

UMAP

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

Accredited by:

22

- **UMAP:**

- Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualisation similarly to t-SNE.
- The algorithm is founded on three assumptions about the data
 - The data is uniformly distributed on Riemannian manifold (a topological space that locally resembles Euclidean space near each point);
 - The manifold is locally connected.
 - From these assumptions it is possible to model the manifold with a fuzzy topological structure. The embedding is found by searching for a low dimensional projection of the data that has the closest possible equivalent fuzzy topological structure.
- For a more detailed explanation see <https://youtu.be/nq6iPZVUxZU>

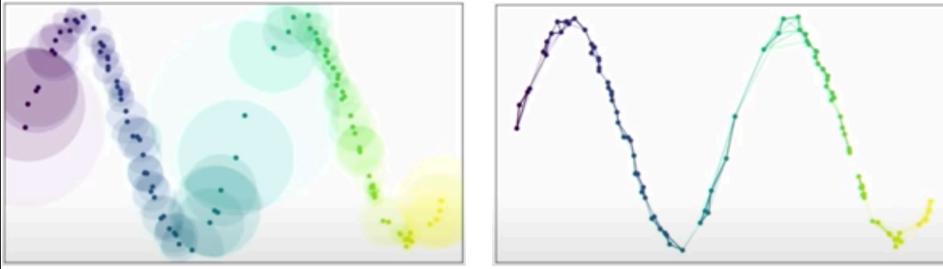
- **UMAP:**

- Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualisation similarly to t-SNE.
- Two phases in the UMAP algorithm
 - The first phase consists of constructing a fuzzy topological representation in the original space;
 - The second phase is simply optimizing (through stochastic gradient descent) the low dimensional representation to have as close a fuzzy topological representation as possible as measured by cross entropy.

NOVA
IMS
Information Management School

UMAP

- **UMAP:**
 - The first phase consists of constructing a fuzzy topological representation;



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

25

NOVA
IMS
Information Management School

UMAP

- **UMAP:**
 - The second phase is simply optimizing (through stochastic gradient descent) the low dimensional representation to have as close a fuzzy topological representation as possible as measured by cross entropy.

Get the clumps right

$$\sum_{a \in A} \mu(a) \log \left(\frac{\mu(a)}{\nu(a)} \right) + (1 - \mu(a)) \log \left(\frac{1 - \mu(a)}{1 - \nu(a)} \right)$$

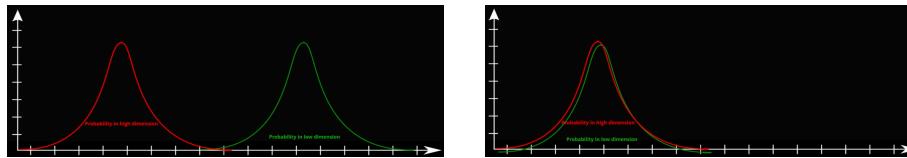
Get the gaps right

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

26

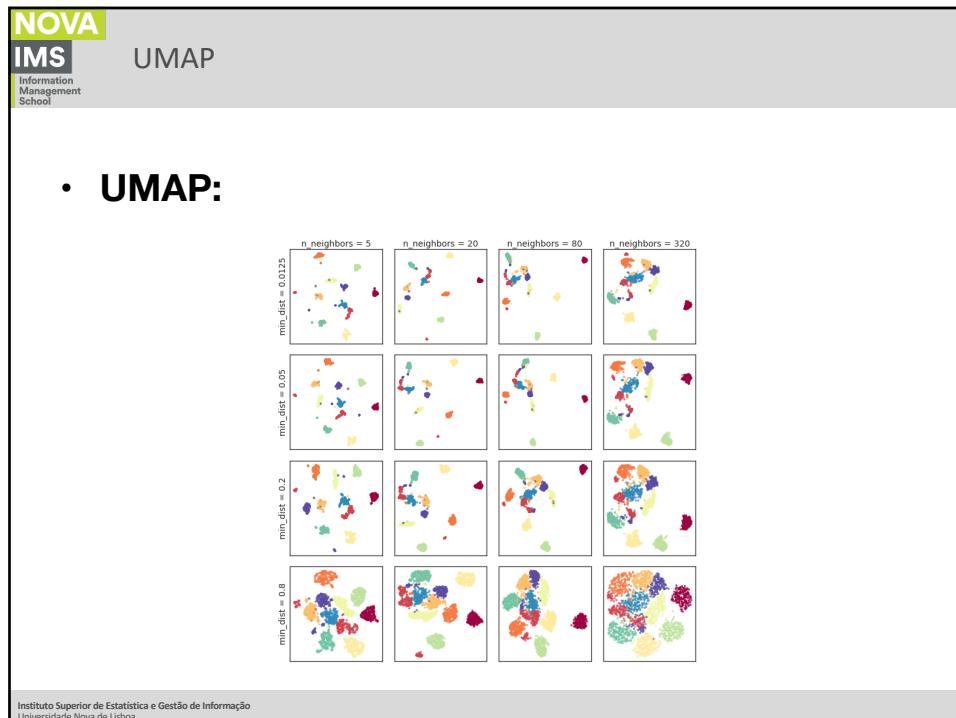
- **UMAP:**

- The second phase is simply optimizing (through stochastic gradient descent) the low dimensional representation to have as close a fuzzy topological representation as possible as measured by cross entropy.

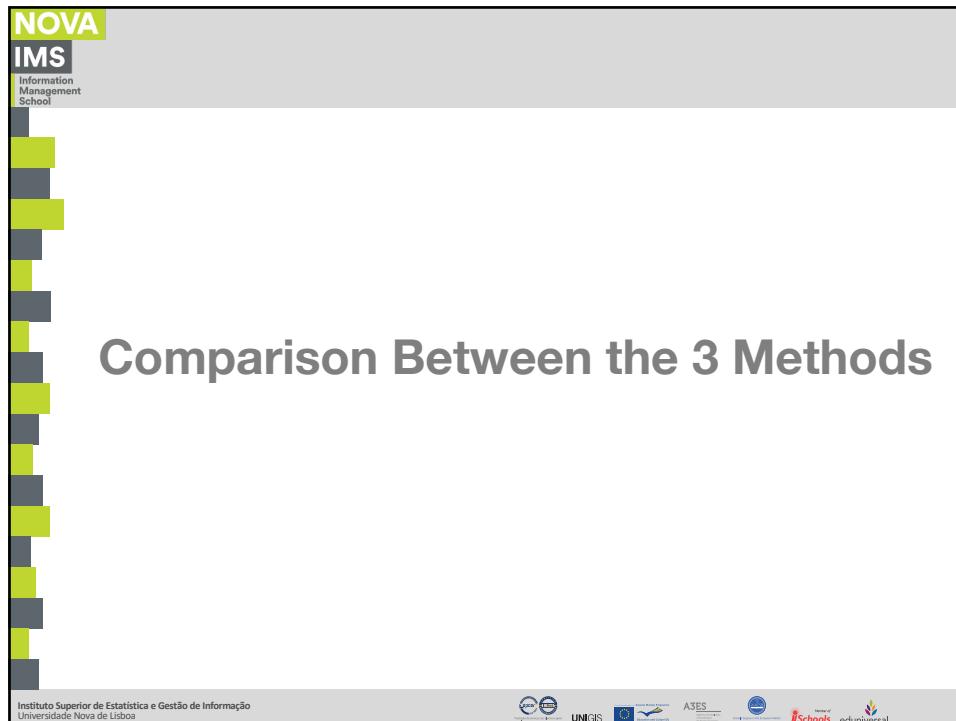


- **UMAP:**

- Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualisation similarly to t-SNE.
- Parameters
 - Number of nearest neighbors - controls how UMAP balances local versus global structure in the data;
 - Minimum distance - controls how tightly UMAP is allowed to pack points together. It provides the minimum distance apart that points are allowed to be in the low dimensional representation



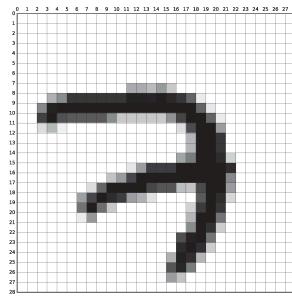
29



30

- **Mnist Digits:**

- 28x28 image (784 dimensions)
- grayscale images of handwritten single digits
- 60,000 examples, and a test set of 10,000 examples



(a) MNIST sample belonging to the digit '7'.

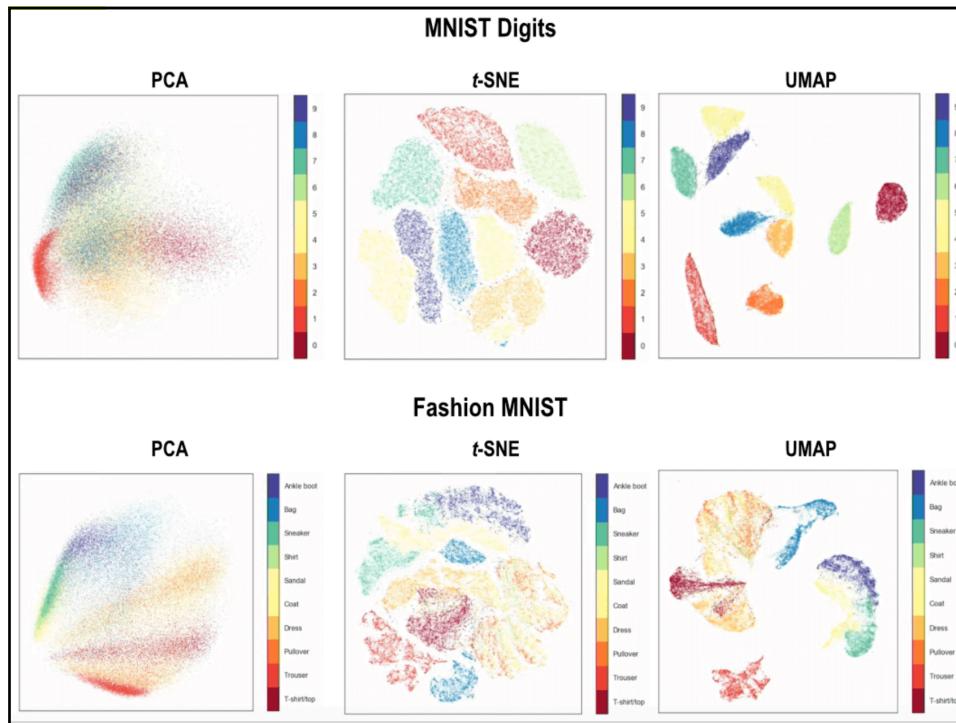


(b) 100 samples from the MNIST training set.

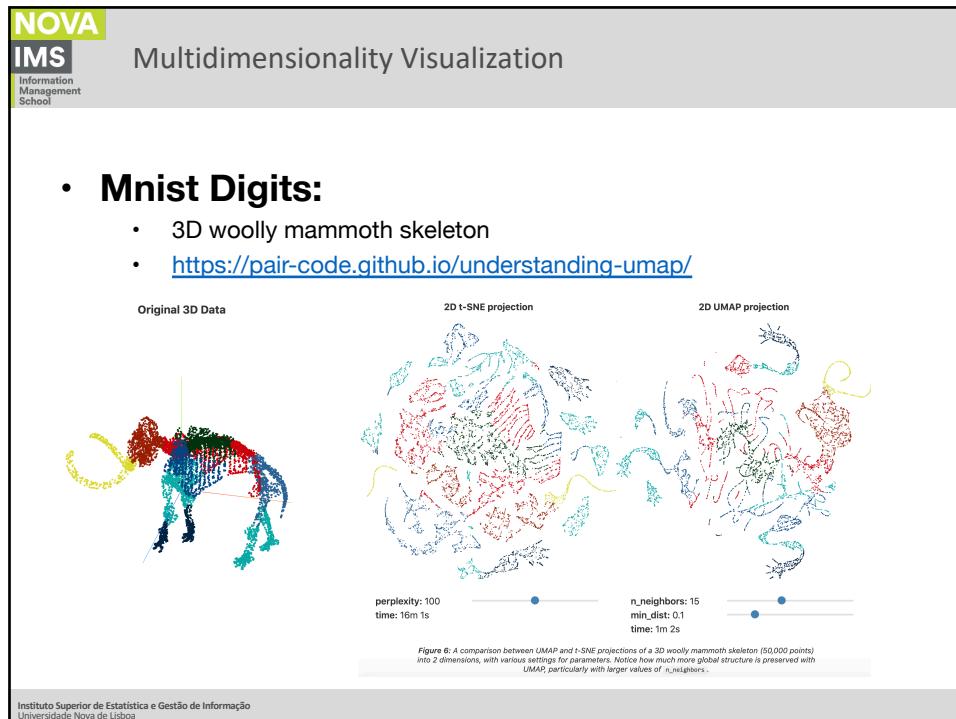
- **Mnist Fashion:**

- 28x28 image (784 dimensions)
- grayscale images of handwritten single digits
- The dataset has 60,000 images





33



34

**NOVA
IMS**
Information Management School

UMAP

- **UMAP:**

	t-SNE	UMAP
COIL20	20 seconds	7 seconds
MNIST	22 minutes	98 seconds
Fashion MNIST	15 minutes	78 seconds
GoogleNews	4.5 hours	14 minutes

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

35



36