



**Data Mining**

**S3**

**NOVA-IMS 2021/2022**

Fernando Lucas Bação

[bacao@isegi.unl.pt](mailto:bacao@isegi.unl.pt)

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

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

1



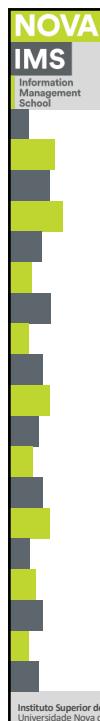
## Agenda

- 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

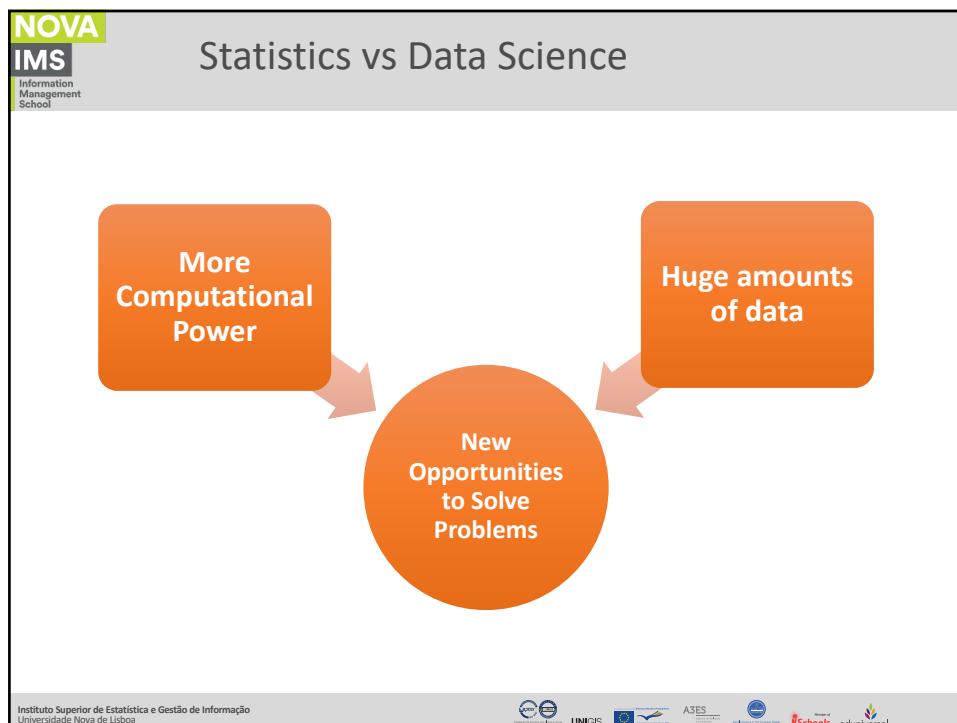


# Statistics vs Data Science

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

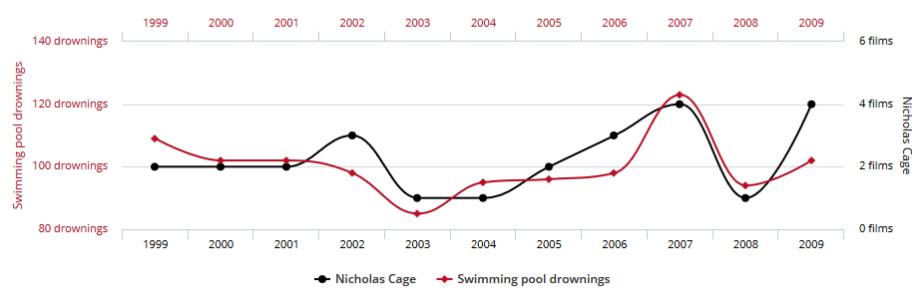
- Because the pattern searches will throw up a **large numbers of candidate patterns**, there will be a **high probability that spurious** (chance) data configurations will be identified as patterns.
- The bottom line is that those patterns and structures identified as potentially interesting will be 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.

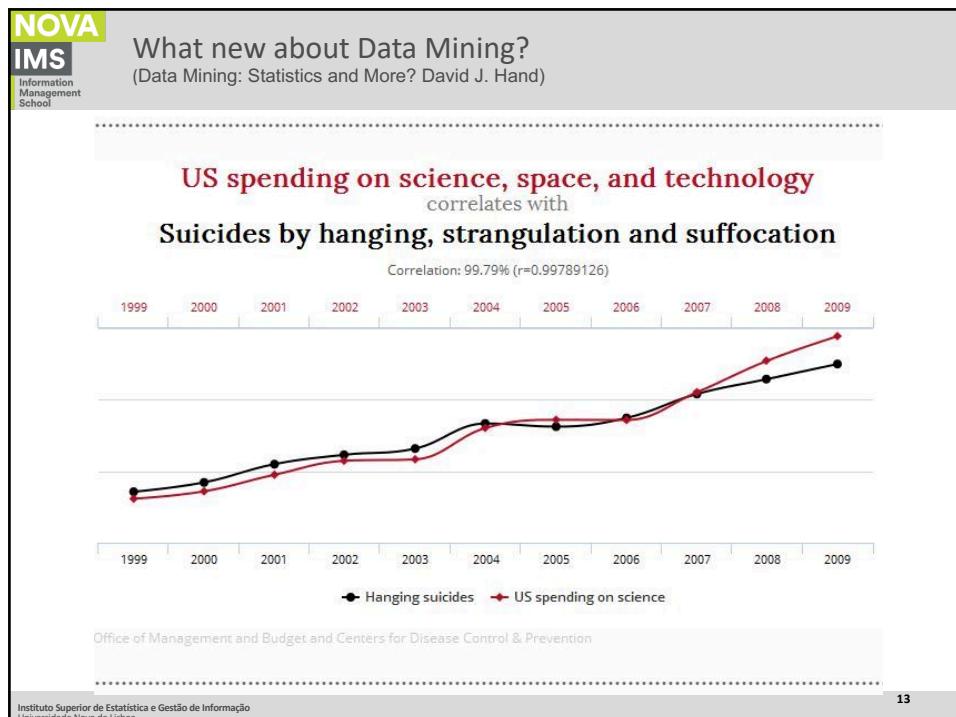
## What new about Data Mining?

(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?

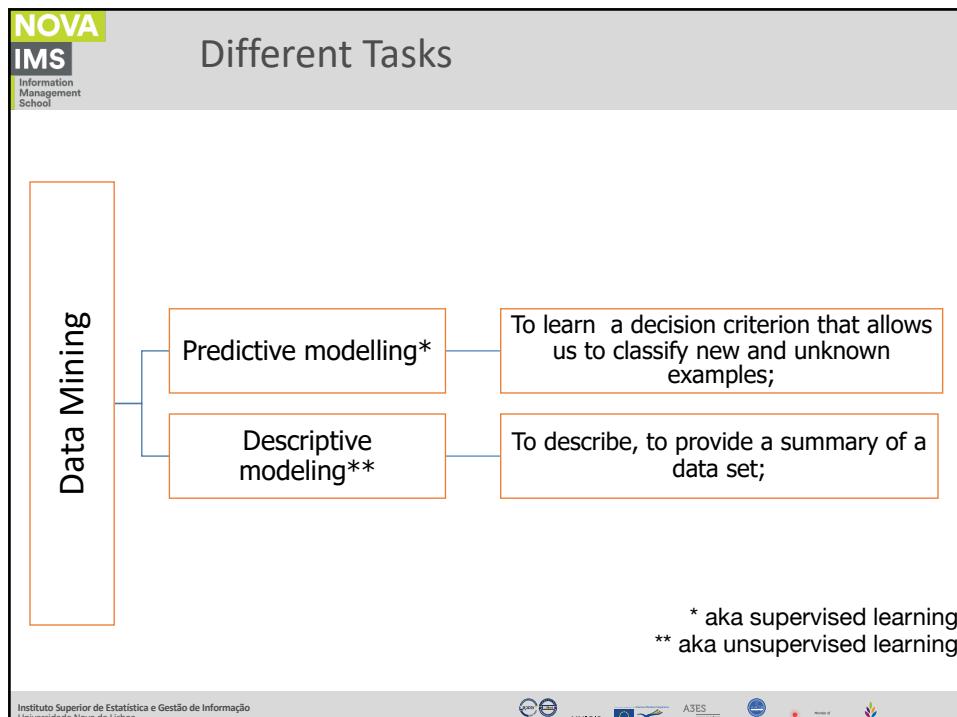
*Data mining is a new **discipline lying at the interface** of statistics, database technology, pattern recognition, machine learning, and other areas.*

*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.*

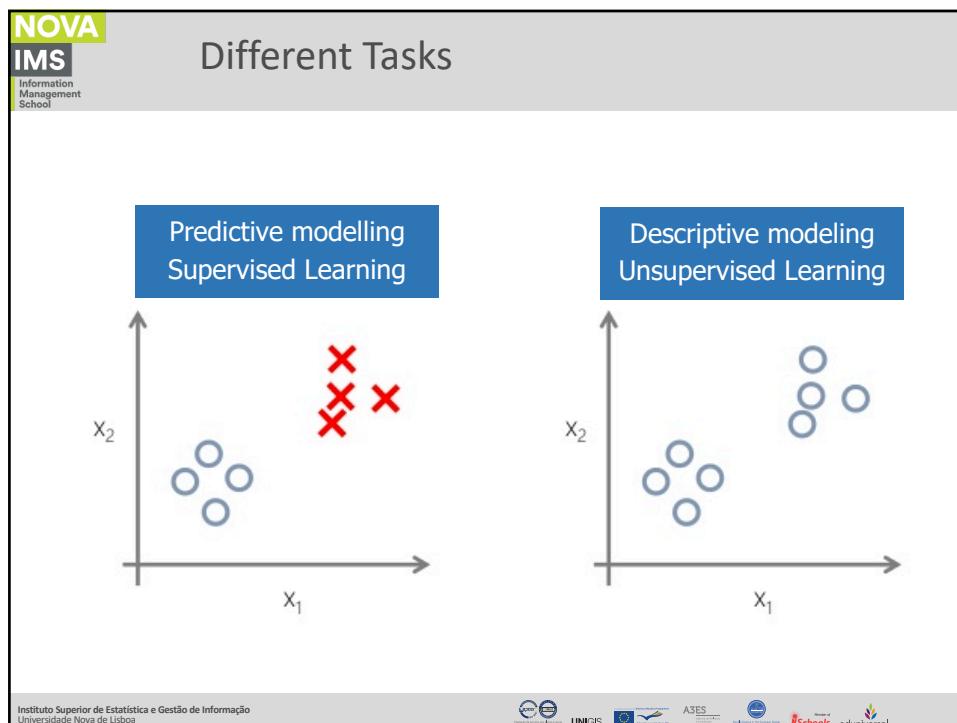
*Data analysis is as much an **art as a science**.*

David J. Hand

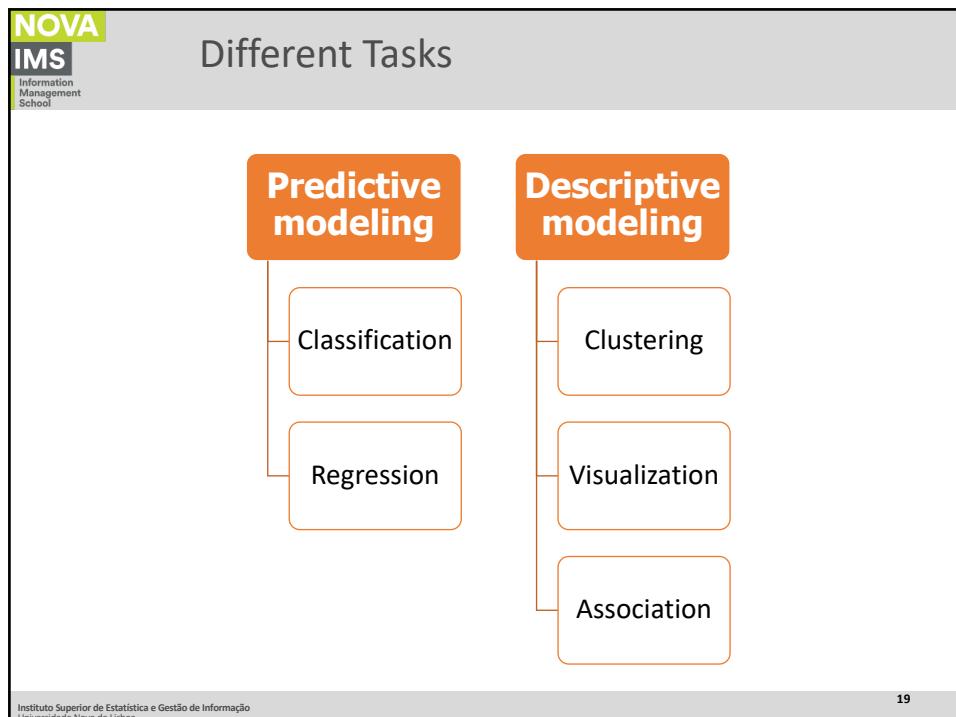
## The Canonical Tasks in Data Mining



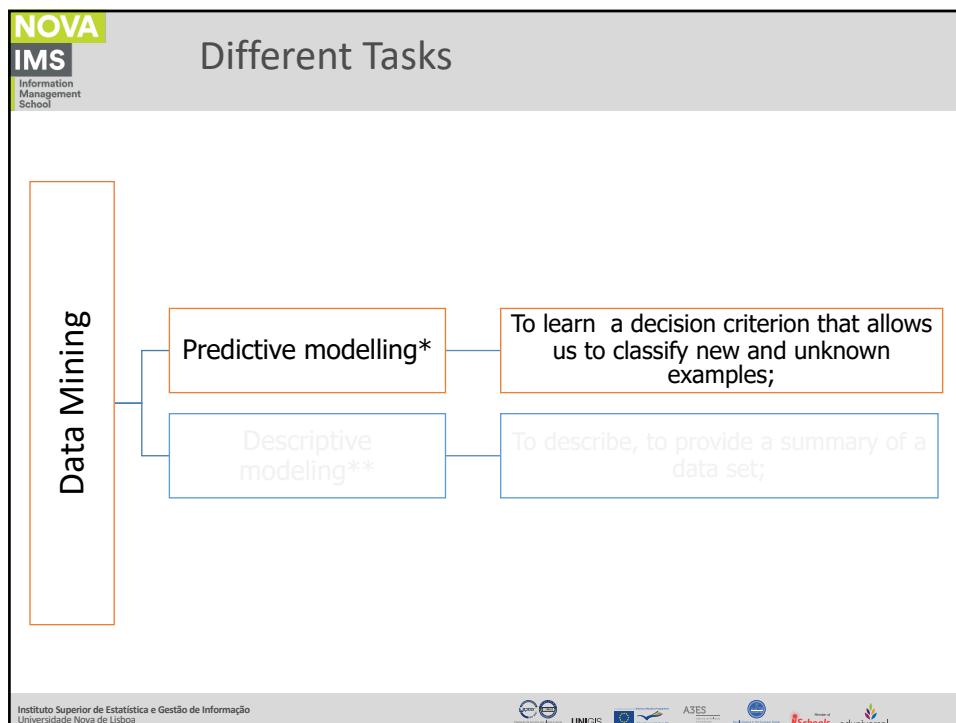
17



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".

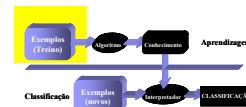
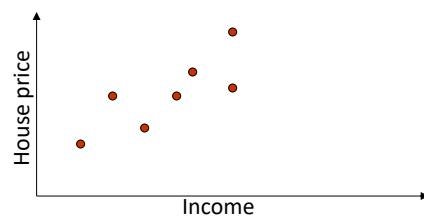
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

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

22

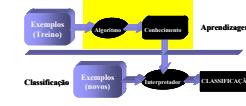
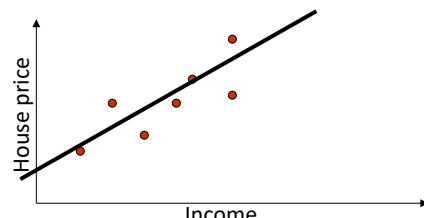
## Predictive Modeling

- 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**

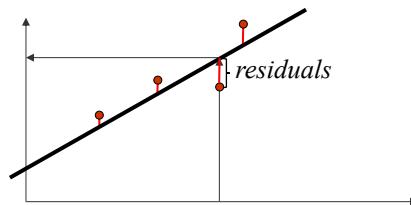
Source: <http://ipython-books.github.io/featured-04/>

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

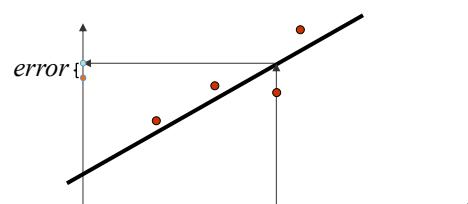
26

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)

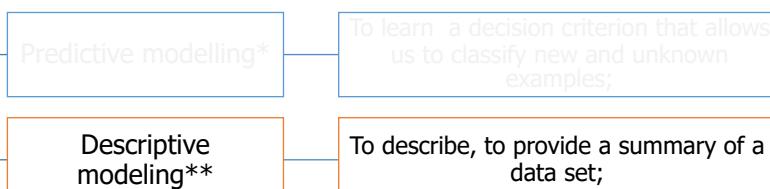
$$\text{Accuracy} = \frac{tp + tn}{tp + tn + fp + fn}$$

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

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

## Different Tasks

### Data Mining



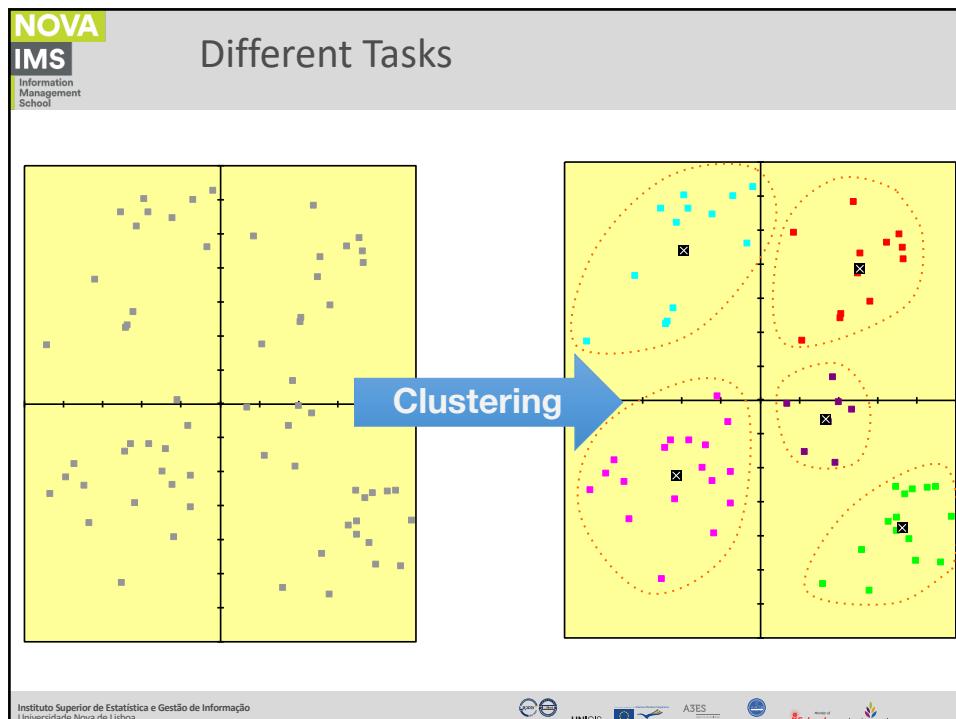
## Different Tasks

Feature

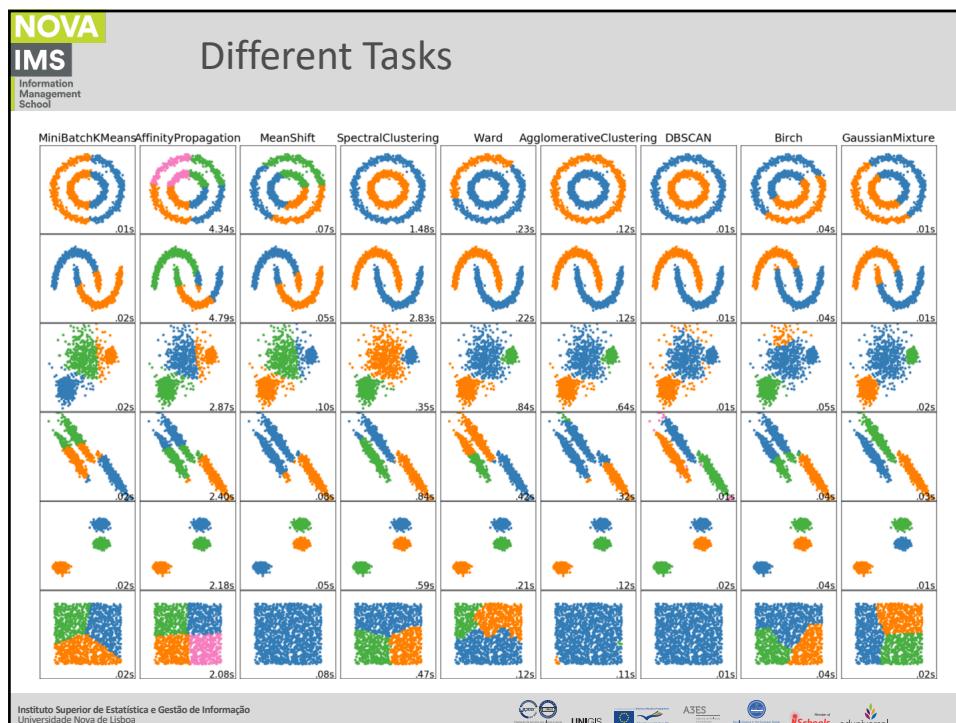


Height	Weight	Sex	Age	Income	Physical Activity
1.60	79	M	41	3000	S
1.72	82	M	32	4000	S
1.66	65	F	28	2500	N
1.82	87	M	35	2000	N
1.71	66	F	42	3500	N

## 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

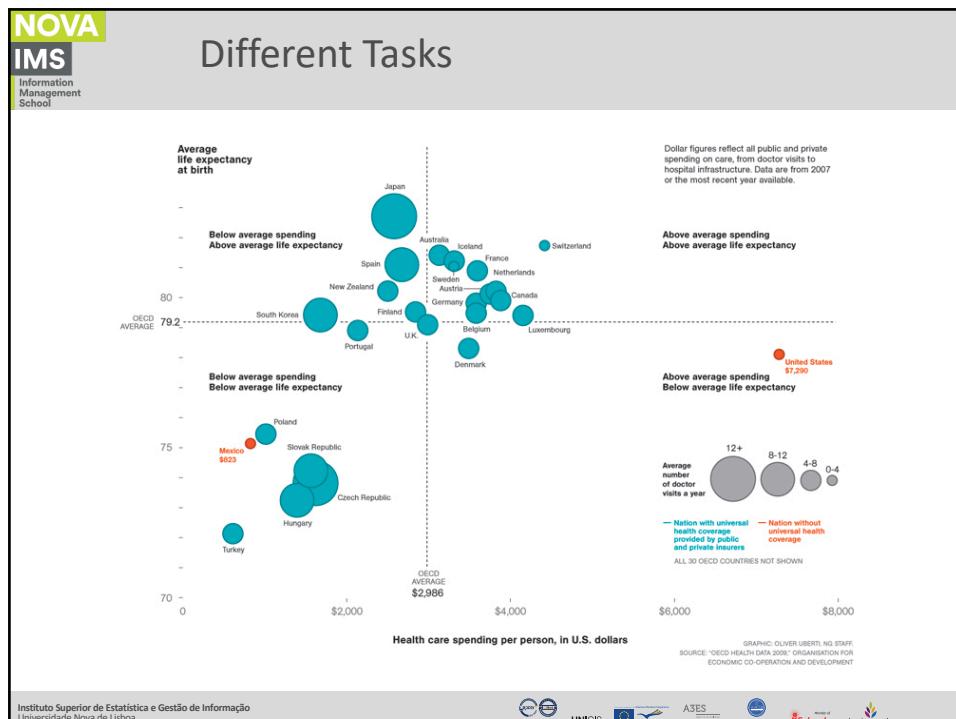
Change in real weekly wages of US-born workers by group, 1990–2006  
(Percent)

Education Level	Young (experience below 20 years)	Old (experience above 20 years)
Some High School	0.4	-5.4
High School Graduate	-1.2	-1.3
Some College	-1.2	-3.0
College Graduate	11.3	6.0

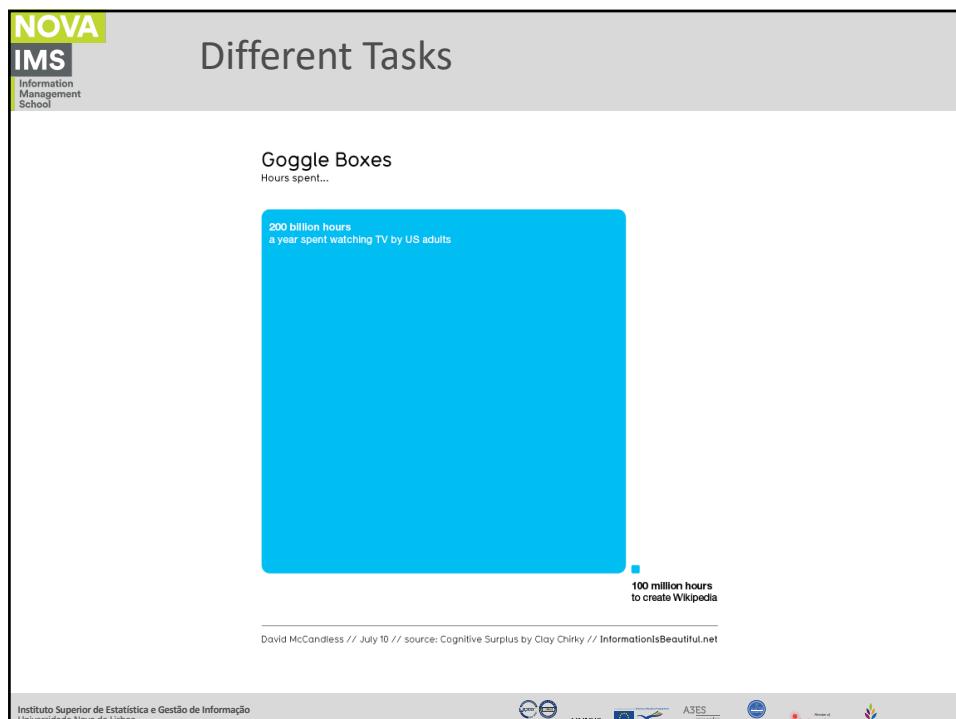
Instituto Superior de Estatística e Gestão de Informação  
Universidade Nova de Lisboa

A3ES UNIGIS iSchools eduniversal

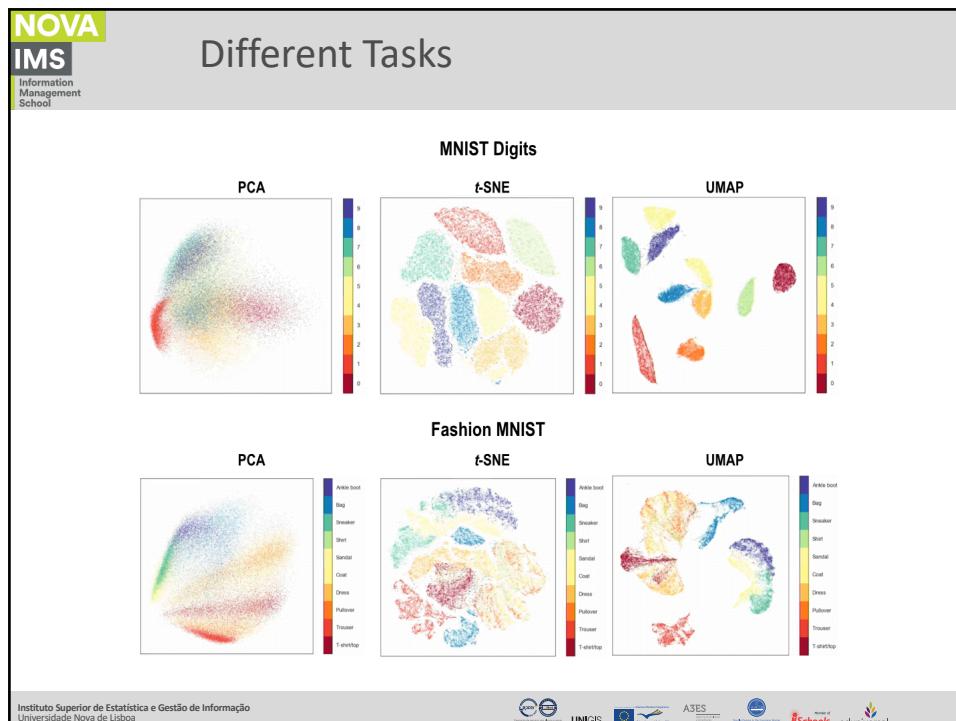
38



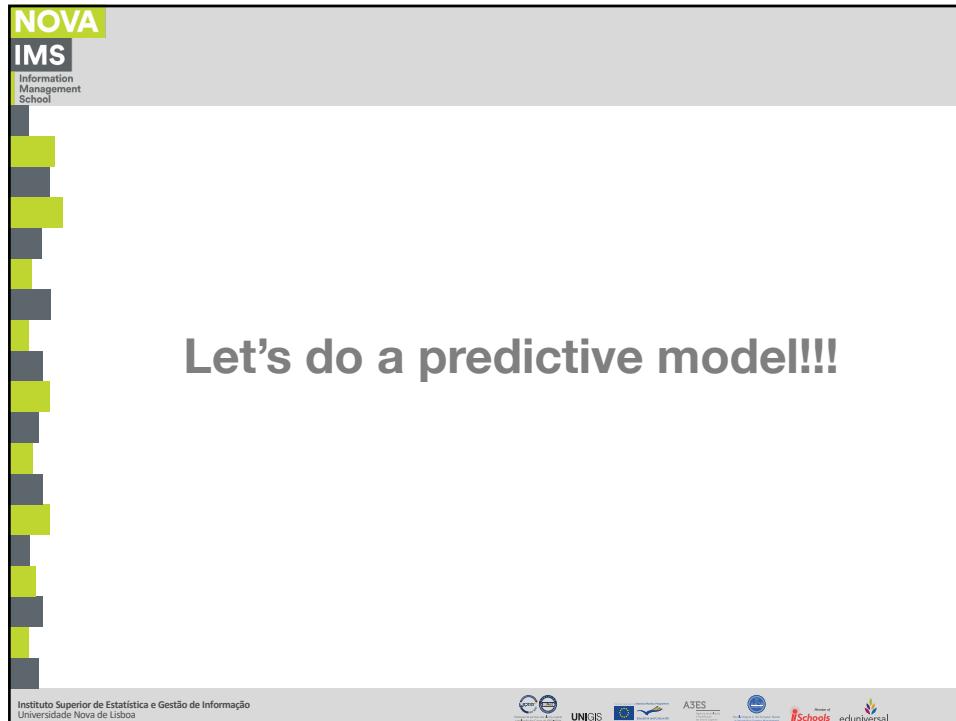
39



40



41



42



Questions?

43