



NOVA
IMS
Information
Management
School

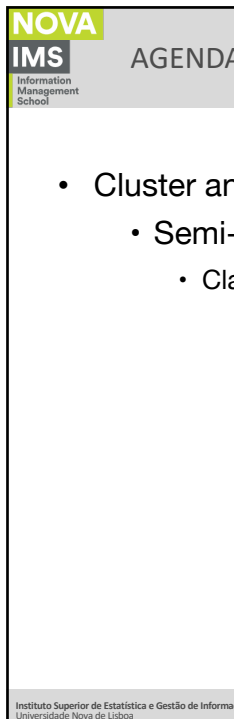
Data Mining

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



NOVA
IMS
Information
Management
School

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

UNIGES, A3ES, Schools, eduniversal

3

NOVA
IMS
Information Management School

Going from clustering into classification

```

graph LR
    Classification[Classification] --- KNN[K-nearest neighbours]
    Classification --- CT[Classification Trees]
    KNN --- KNN_Usage[Use to classify new instances of the problem with the proper label;]
    CT --- CT_Usage[Use to classify new instances and understand what distinguish different labels;]
  
```

The diagram illustrates the transition from clustering to classification. It features a vertical orange bar on the left labeled "Classification". To its right, two boxes are connected by a horizontal line: a light gray box for "K-nearest neighbours" and an orange box for "Classification Trees". Each box has a corresponding description to its right: "Use to classify new instances of the problem with the proper label;" for K-nearest neighbours, and "Use to classify new instances and understand what distinguish different labels;" for Classification Trees.

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

4

NOVA
IMS
Information Management School

Classification Trees

- **Classification Trees:**
 - Classification trees 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.

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

5

NOVA
IMS
Information Management School

Classification Trees

- 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 leaf should represent only individuals from a specific class

```

graph TD
    Node1([exam grade ≥ 10]) -- Yes --> Leaf1[Approved]
    Node1 -- No --> Node2([Project grade ≥ 14])
    Node2 -- Yes --> Node3([exam grade ≥ 8])
    Node2 -- No --> Leaf2[Failed]
    Node3 -- No --> Leaf3[Failed]
    Node3 -- Yes --> Leaf4[Approved]
  
```

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

6

NOVA
IMS
Information Management School

Classification Trees

$aprovado \Leftrightarrow (exame \geq 10)$
 $aprovado \Leftrightarrow (exame < 10) \wedge (projecto \geq 14) \wedge (exame \geq 8)$
 $reprovado \Leftrightarrow (exame < 10) \wedge (projecto < 14)$
 $reprovado \Leftrightarrow (exame < 10) \wedge (projecto > 14) \wedge (exame < 8)$

```

graph TD
    A([exam grade ≥ 10]) -- Yes --> B[Approved]
    A -- No --> C([Project grade ≥ 14])
    C -- Yes --> D([exam grade ≥ 8])
    C -- No --> E[Failed]
    D -- No --> F[Failed]
    D -- Yes --> G[Approved]
  
```

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

7

NOVA
IMS
Information Management School

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

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

8

NOVA
IMS
Information Management School

Classification Trees

- **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

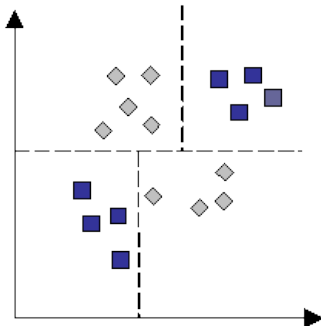
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

9

NOVA
IMS
Information Management School

Classification Trees

- **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/GeoComp/GeoComp98/61/gc_61.htm

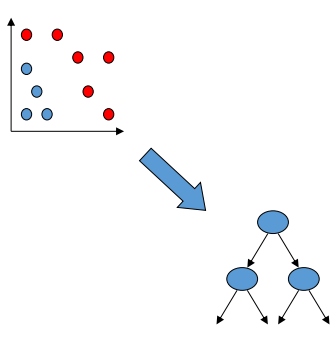
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

10

NOVA
IMS
Information Management School

Classification Trees

- **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?



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

11

NOVA
IMS
Information Management School

Classification Trees

- **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.

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

12

NOVA
IMS
Information
Management
School

Classification Trees

Worked-Example General Idea

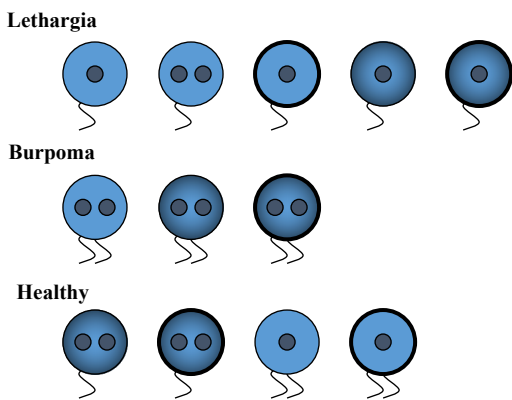
Langley, P: 1996, Elements of Machine Learning, Morgan and Kaufmann Publishers.

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

13

NOVA
IMS
Information
Management
School

Classification Trees



Lethargia

Burpoma

Healthy

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	Class
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	1	2
Lethargia	4	1
Burpoma	0	3
Healthy	2	2

# 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	2	Dark	Thick	Burpoma
2	1	Dark	Thin	Healthy
2	1	Dark	Thick	Healthy
1	2	Light	Thin	Healthy
1	2	Light	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

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

# 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	2	Dark	Thick	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

# 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	2	Dark	Thick	Burpoma
2	1	Dark	Thin	Healthy
2	1	Dark	Thick	Healthy
1	2	Light	Thin	Healthy
1	2	Light	Thick	Healthy

Discrimination:

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

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

# 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	2	Dark	Thick	Burpoma
2	1	Dark	Thin	Healthy
2	1	Dark	Thick	Healthy
1	2	Light	Thin	Healthy
1	2	Light	Thick	Healthy

Discrimination:

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

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

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

0.67

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

0.41

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

Tails

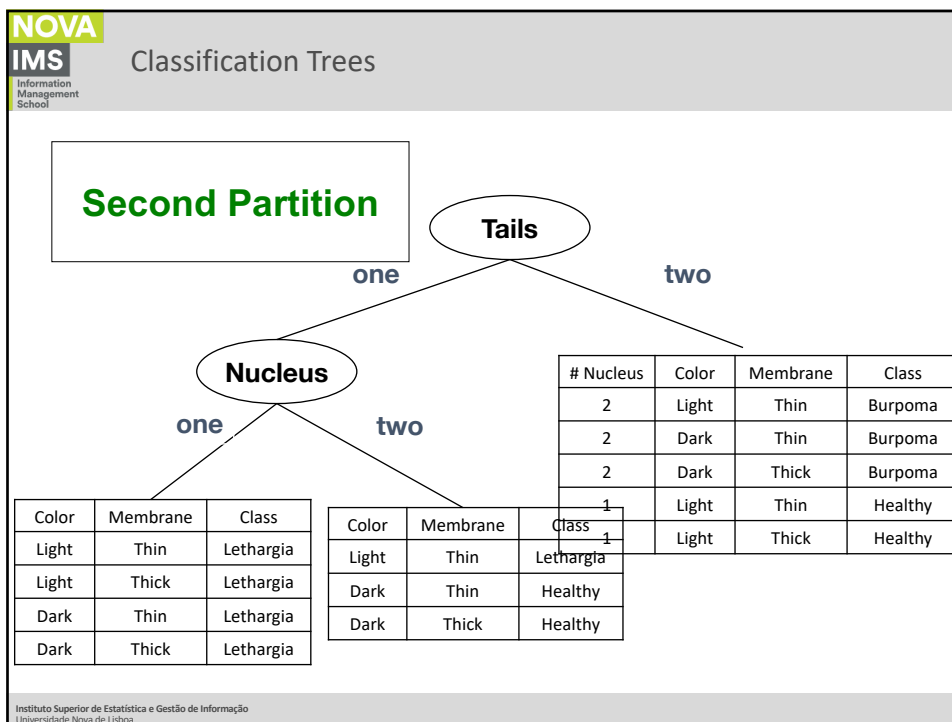
one **two**

# 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

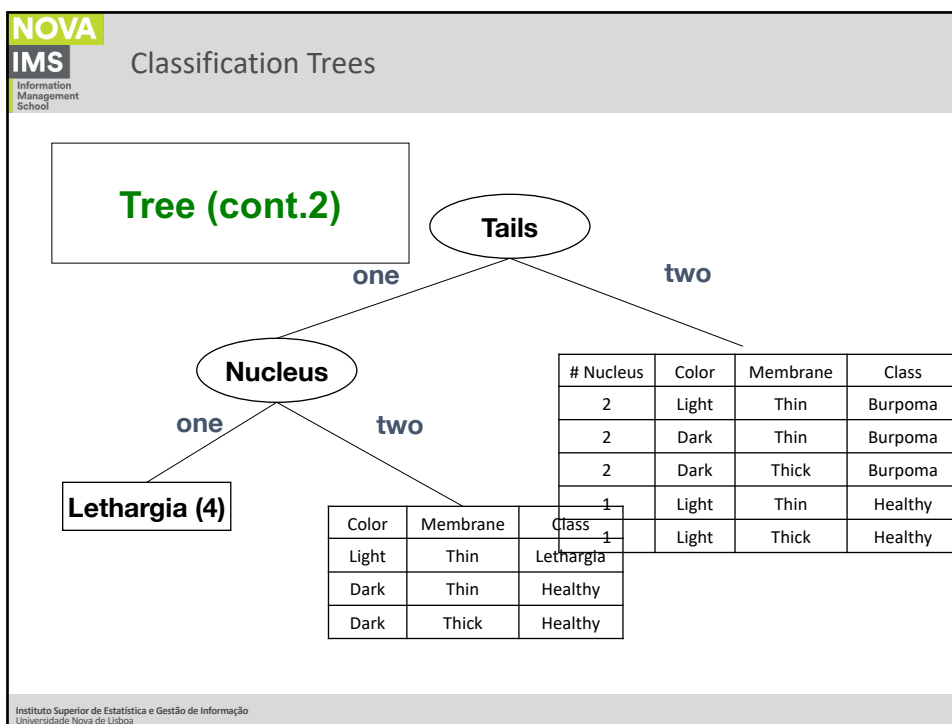
# Nucleus	Color	Membrane	Class
2	Light	Thin	Burpoma
2	Dark	Thin	Burpoma
2	Dark	Thick	Burpoma
1	Light	Thin	Healthy
1	Light	Thick	Healthy

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

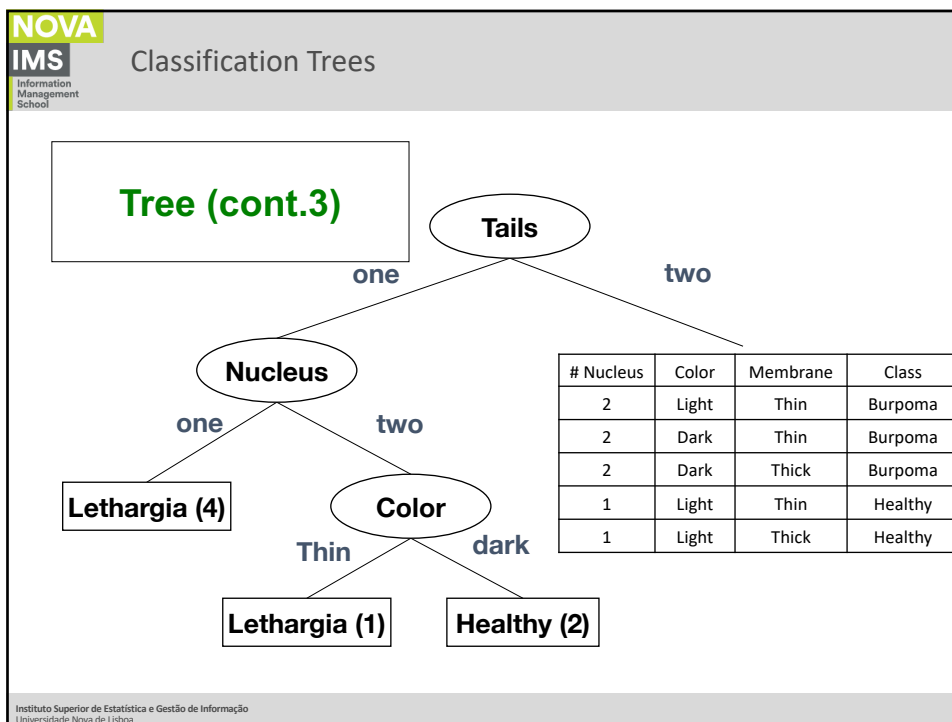
22



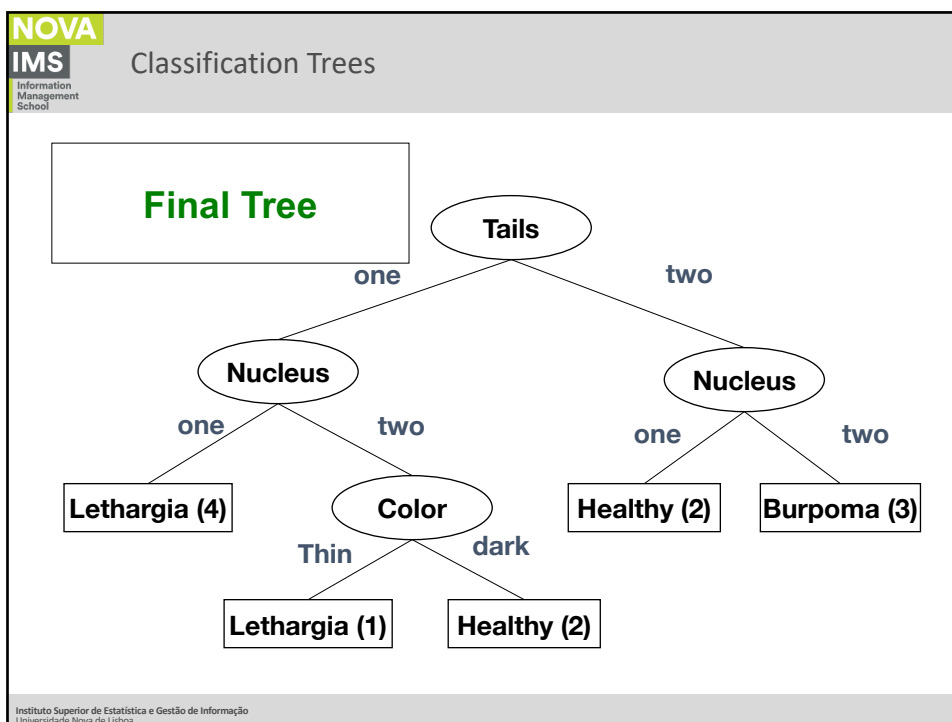
23



24



25



26

NOVA
IMS
Information Management School

Classification Trees

The description of the three Classs

$$\begin{array}{c} (tails = 1) \wedge (nucleous = 1) \\ \vee \\ (tails = 1) \wedge (nucleous = 2) \wedge (color = light) \rightarrow Lethargic \end{array}$$

$$\begin{array}{c} (tails = 2) \wedge (nucleous = 1) \\ \vee \\ (tails = 1) \wedge (nucleous = 2) \wedge (color = dark) \rightarrow Healthy \end{array}$$

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

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

27

NOVA
IMS
Information Management School

Classification Trees

- **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

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

28

NOVA

IMS

Information
Management
School

Classification Trees

- **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 – Attribute Set
 - ISET – Instance Set

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

29

NOVA

IMS

Information
Management
School

Classification Trees

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 pair 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)

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

30

NOVA
IMS
Information Management School

Classification Trees

Worked-Example Tree Accuracy

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

31

NOVA
IMS
Information Management School

Classification Trees

Quality of the results Error rate

100 A 100 B 100 C → 100, 100, 100

75, 65, 20 25, 35, 80

7, 63, 4 68, 2, 16 10, 15, 70 15, 20, 10

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

32

