

## Unit 3: Inductive machine learning

M. A. Gutiérrez Naranjo

F. J. Martín Mateos

A. Riscos Núñez

J. L. Ruiz Reina

Dpto. Ciencias de la Computación e Inteligencia Artificial  
Universidad de Sevilla

- Introduction
- Learning decision trees
- Learning rules
- Instance based learning: kNN
- Clustering: k-means

## Section 1

### Introduction

- Definitions of *learning*:
  - Learning is any process by which a system *improves performance* from experience (*H. Simon*)
  - Modify the representation of perceived world (*R. Michalski*)
  - Making useful changes in our minds (*M. Minsky*)
- Machine learning: build programs that automatically improve with experience
- Tasks examples:
  - Build knowledge bases from experience
  - Classification and diagnosis
  - Data mining, uncover unknown patterns in big databases
  - Resolutions of problems, planning and acting

# Types of problems in machine learning and paradigms

- Types of learning
  - Supervised learning
  - Unsupervised learning
  - Reinforcement learning
- Paradigms
  - Rote learning
  - Classification (*Clustering*)
  - Inductive learning
  - Learning by analogy
  - Discover
  - Genetic algorithms, neural networks

# Example of learning

- Training set
  - Examples: days appropriate (or not) for playing tennis
  - Data represented as a list of pairs attribute-value

EJ.	OUTLOOK	TEMPERATURE	HUMIDITY	WIND	PLAYTENNIS
$D_1$	SUNNY	HIGH	HIGH	WEAK	-
$D_2$	SUNNY	HIGH	HIGH	STRONG	-
$D_3$	OVERCAST	HIGH	HIGH	WEAK	+
$D_4$	RAINY	MILD	HIGH	WEAK	+
...					

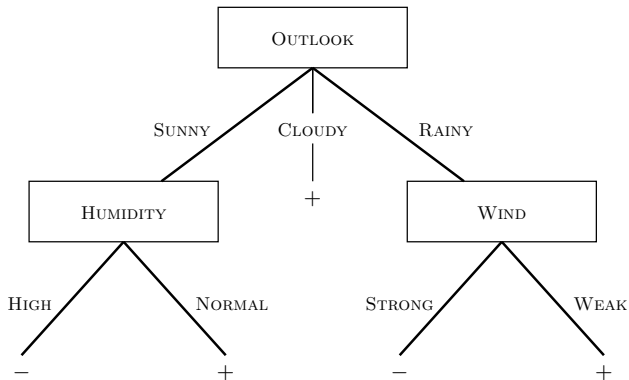
- Goal: given the *training set*, *learn* the concept “Day appropriate for playing tennis”.
  - This is *supervised learning*
- Problem: ¿How to express what has been learned?
  - In this unit, we will see algorithms for learning *decision trees*, *rules*, *neural networks*,...

## Section 2

### Learning decision trees

# Decision trees

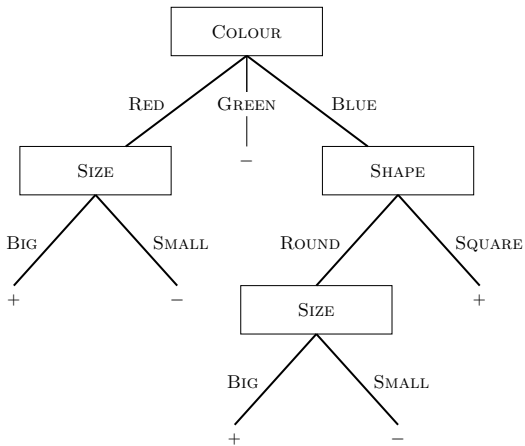
- Decision trees examples





# Decision trees

- Decision trees examples



- Decision trees
  - Internal nodes: attributes
  - Outgoing edges labeled with values of the node
  - Leaves: classification value (usually + or -, although it could be any set of classes, not necessarily binary)
  - It represents a classification function
- Disjunction of propositional rules
$$\begin{aligned} & (\text{OUTLOOK}=\text{SUNNY} \wedge \text{HUMIDITY}=\text{HIGH} \rightarrow \text{PLAYTENNIS}=-) \\ \vee & (\text{OUTLOOK}=\text{SUNNY} \wedge \text{HUMIDITY}=\text{NORMAL} \rightarrow \text{PLAYTENNIS}=+) \\ \vee & (\text{OUTLOOK}=\text{OVERCAST} \rightarrow \text{PLAYTENNIS}=+) \\ \vee & (\text{OUTLOOK}=\text{RAINY} \wedge \text{WIND}=\text{STRONG} \rightarrow \text{PLAYTENNIS}=-) \\ \vee & (\text{OUTLOOK}=\text{RAINY} \wedge \text{WIND}=\text{WEAK} \rightarrow \text{PLAYTENNIS}=+) \end{aligned}$$
- Powerful enough to represent any subset of positive instances

# Learning decision trees

- Goal: Learn a decision tree, consistent with the examples in the training set, and then use it to classify new examples
- Training set example:

EJ.	OUTLOOK	TEMPERATURE	HUMIDITY	WIND	PLAYTENNIS
$D_1$	SUNNY	HIGH	HIGH	WEAK	-
$D_2$	SUNNY	HIGH	HIGH	STRONG	-
$D_3$	OVERCAST	HIGH	HIGH	WEAK	+
$D_4$	RAINY	MILD	HIGH	WEAK	+
...					

## ID3 algorithm

ID3(Examples, Class, Attributes)

1. If all the Examples are positive, return a leaf node labeled with +
2. If all the Examples are negative, return a leaf node labeled with -
3. If Attributes is empty, return a node labeled with the most frequent class in Examples
4. Otherwise:
  - 4.1. Let A be the BEST attribute of Attributes, as for the classifying Examples
  - 4.2. Create Tree, with a node labeled with A
  - 4.3. For every possible value v of A:
    - \* Add an outgoing edge to Tree, labeled with v
    - \* Let Examples(v) the subset of Examples with value v in attribute A
    - \* If Examples(v) is empty:
      - Then put a leaf labeled with the most frequent class in Examples
      - Else, extend the outgoing edge with the subtree ID3(Examples(v), Class, Attributes-{A}).
  - 4.4 Return Tree

# Which is the BEST attribute for an internal node?

- Entropy of a set of examples  $D$  (w.r.t. a classification):

$$Ent(D) = -\frac{|P|}{|D|} \cdot \log_2 \frac{|P|}{|D|} - \frac{|N|}{|D|} \cdot \log_2 \frac{|N|}{|D|}$$

where  $P$  and  $N$  are, respectively, the subsets of positive and negative examples in  $D$

- Notation:  $Ent([p^+, n^-])$ , where  $p = |P|$  and  $n = |N|$
- Intuition:
  - Quantify the “heterogeneity” of a classification
  - Information Theory: average amount of information (in bits) needed to codify the class of an example
- Examples:
  - $Ent([9^+, 5^-]) = -\frac{9}{14} \cdot \log_2 \frac{9}{14} - \frac{5}{14} \cdot \log_2 \frac{5}{14} = 0,94$
  - $Ent([k^+, k^-]) = 1$  (no homogeneity at all)
  - $Ent([p^+, 0^-]) = Ent([0^+, n^-]) = 0$  (total homogeneity)

- Nodes with less entropy are better (small trees)
- Expected entropy after using attribute  $A$  in the tree:

$$\sum_{v \in \text{Valores}(A)} \frac{|D_v|}{|D|} \cdot \text{Ent}(D_v)$$

where  $D_v$  es the subset of  $D$  with value  $v$  in attribute  $A$

- Expected information gain después after using attribute  $A$ :

$$\text{Gain}(D, A) = \text{Ent}(D) - \sum_{v \in \text{Valores}(A)} \frac{|D_v|}{|D|} \cdot \text{Ent}(D_v)$$

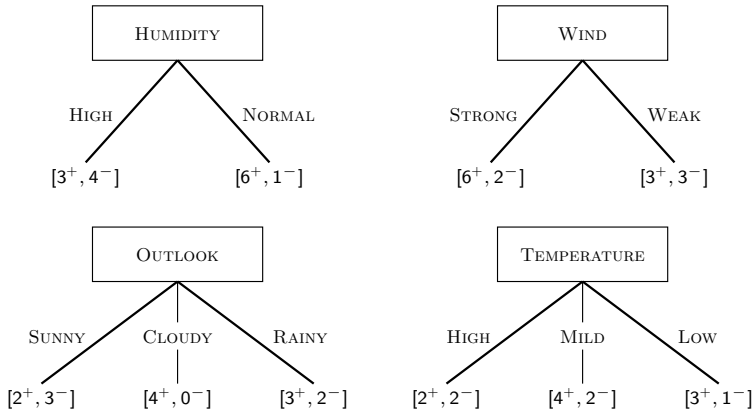
- In ID3, we use *information gain* as the criteria for selecting the best attribute in every internal node (taking into account the examples at that node)

# Algorithm ID3 (example 1)

- Training set:

EJ.	OUTLOOK	TEMPERATURE	HUMIDITY	WIND	PLAYTENNIS
$D_1$	SUNNY	HIGH	HIGH	WEAK	-
$D_2$	SUNNY	HIGH	HIGH	STRONG	-
$D_3$	OVERCAST	HIGH	HIGH	WEAK	+
$D_4$	RAINY	MILD	HIGH	WEAK	+
$D_5$	RAINY	LOW	NORMAL	WEAK	+
$D_6$	RAINY	LOW	NORMAL	STRONG	-
$D_7$	OVERCAST	LOW	NORMAL	STRONG	+
$D_8$	SUNNY	MILD	HIGH	WEAK	-
$D_9$	SUNNY	LOW	NORMAL	WEAK	+
$D_{10}$	RAINY	MILD	NORMAL	WEAK	+
$D_{11}$	SUNNY	MILD	NORMAL	STRONG	+
$D_{12}$	OVERCAST	MILD	HIGH	STRONG	+
$D_{13}$	OVERCAST	HIGH	NORMAL	WEAK	+
$D_{14}$	RAINY	MILD	HIGH	STRONG	-

# Algorithm ID3 (example 1)



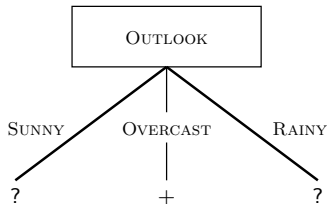


# Algorithm ID3 (example 1)

- Initial entropy:  $Ent([9^+, 5^-]) = 0,94$
- Selecting the attribute for the root node:
  - $Gain(D, \text{HUMIDITY}) =$   
 $0,94 - \frac{7}{14} \cdot Ent([3^+, 4^-]) - \frac{7}{14} \cdot Ent([6^+, 1^-]) = 0,151$
  - $Gain(D, \text{WIND}) =$   
 $0,94 - \frac{8}{14} \cdot Ent([6^+, 2^-]) - \frac{6}{14} \cdot Ent([3^+, 3^-]) = 0,048$
  - $Gain(D, \text{OUTLOOK}) =$   
 $0,94 - \frac{5}{14} \cdot Ent([2^+, 3^-]) - \frac{4}{14} \cdot Ent([4^+, 0^-])$   
 $- \frac{5}{14} \cdot Ent([3^+, 2^-]) = 0,246$  (best attribute)
  - $Gain(D, \text{TEMPERATURE}) =$   
 $0,94 - \frac{4}{14} \cdot Ent([2^+, 2^-]) - \frac{6}{14} \cdot Ent([4^+, 2^-])$   
 $- \frac{4}{14} \cdot Ent([3^+, 1^-]) = 0,02$
- The selected attribute is OUTLOOK

# Algorithm ID3 (example 1)

- Partial tree:



# Algorithm ID3 (example 1)

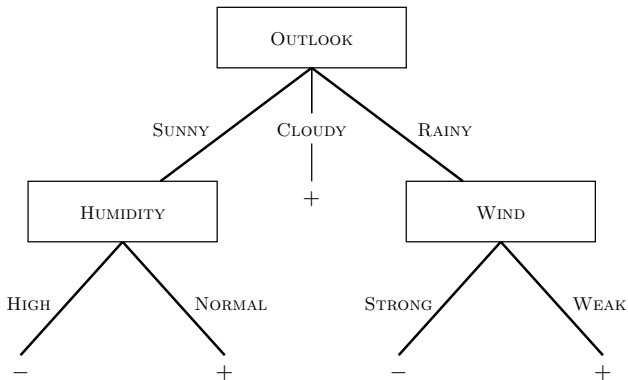
- Selecting the attribute for node  $\text{OUTLOOK}=\text{SUNNY}$
- $D_{\text{SUNNY}} = \{D_1, D_2, D_8, D_9, D_{11}\}$  with entropy  
 $\text{Ent}([2^+, 3^-]) = 0,971$ 
  - $\text{Gain}(D_{\text{SUNNY}}, \text{HUMIDITY}) =$   
 $0,971 - \frac{3}{5} \cdot 0 - \frac{2}{5} \cdot 0 = 0,971$  (best attribute)
  - $\text{Gain}(D_{\text{SUNNY}}, \text{TEMPERATURE}) =$   
 $0,971 - \frac{2}{5} \cdot 0 - \frac{2}{5} \cdot 1 - \frac{1}{5} \cdot 0 = 0,570$
  - $\text{Gain}(D_{\text{SUNNY}}, \text{WIND}) =$   
 $0,971 - \frac{2}{5} \cdot 1 - \frac{3}{5} \cdot 0,918 = 0,019$
- The selected attribute is  $\text{HUMIDITY}$

# Algorithm ID3 (example 1)

- Selecting the attribute for node  $\text{OUTLOOK}=\text{RAINY}$ :
- $D_{\text{RAINY}} = \{D_4, D_5, D_6, D_{10}, D_{14}\}$  with entropy  
 $\text{Ent}([3^+, 2^-]) = 0,971$ 
  - $\text{Gain}(D_{\text{RAINY}}, \text{HUMIDITY}) =$   
 $0,971 - \frac{2}{5} \cdot 1 - \frac{3}{5} \cdot 0,918 = 0,020$
  - $\text{Gain}(D_{\text{RAINY}}, \text{TEMPERATURA}) =$   
 $0,971 - \frac{3}{5} \cdot 0,918 - \frac{2}{5} \cdot 1 = 0,020$
  - $\text{Gain}(D_{\text{RAINY}}, \text{WIND}) =$   
 $0,971 - \frac{3}{5} \cdot 0 - \frac{2}{5} \cdot 0 = 0,971$  (best attribute)
- The selected attribute is  $\text{WIND}$

# Algoritmo ID3 (ejemplo 1)

- The tree finally learned:



# Algoritmo ID3 (ejemplo 2)

- Conjunto de entrenamiento:

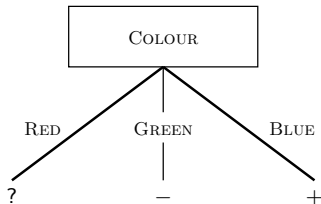
EJ.	COLOUR	SHAPE	SIZE	CLASS
$O_1$	RED	SQUARE	BIG	+
$O_2$	BLUE	SQUARE	BIG	+
$O_3$	RED	ROUND	SMALL	-
$O_4$	GREEN	SQUARE	SMALL	-
$O_5$	RED	ROUND	BIG	+
$O_6$	GREEN	SQUARE	BIG	-

# Algorithm ID3 (example 2)

- Initial entropy,  $Ent([3^+, 3^-]) = 1$
- Selecting the attribute for the root node:
  - $Gain(D, COLOUR) = 1 - \frac{3}{6} \cdot Ent([2^+, 1^-]) - \frac{1}{6} \cdot Ent([1^+, 0^-]) - \frac{2}{6} \cdot Ent([0^+, 2^-]) = 0,543$
  - $Gain(D, SHAPE) = 1 - \frac{4}{6} \cdot Ent([2^+, 2^-]) - \frac{2}{6} \cdot Ent([1^+, 1^-]) = 0$
  - $Gain(D, SIZE) = 1 - \frac{4}{6} \cdot Ent([3^+, 1^-]) - \frac{2}{6} \cdot Ent([0^+, 2^-]) = 0,459$
- The selected attribute is COLOUR

# Algorithm ID3 (example 2)

- Partial tree:



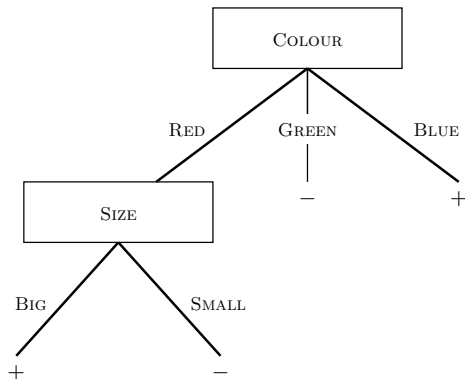


# Algorithm ID3 (example 2)

- Selecting attribute for the node COLOUR=RED:
- $D_{\text{RED}} = \{O_1, O_3, O_5\}$  with entropy  $Ent([2^+, 1^-]) = 0,914$ 
  - $Gain(D_{\text{RED}}, \text{SHAPE}) = 0,914 - \frac{1}{3} \cdot Ent([1^+, 0^-]) - \frac{2}{3} \cdot Ent([1^+, 1^-]) = 0,247$
  - $Gain(D_{\text{RED}}, \text{SIZE}) = 0,914 - \frac{2}{3} \cdot Ent([2^+, 0^-]) - \frac{1}{3} \cdot Ent([0^+, 1^-]) = 0,914$
- The selected attribute is SIZE

# Algorithm ID3 (example 2)

- Final tree learned:



# Searching and inductive bias

- This is a searching problem (local search)
  - In the space of all the decision trees
  - In each step we refine the current tree
  - No backtracking (local optima), hill-climbing search
  - Decisions made by analyzing sets of examples
- Inductive Bias
  - We prefer smaller trees
  - Preference bias, implicit in the search
  - Occam's razor principle

# Some practical questions in machine learning

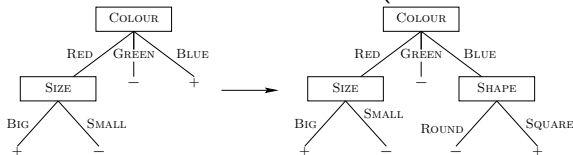
- Validate the learned hypothesis
  - Can we *quantify* how good is what we have learned, compared to the real explanation?
- Overfitting
  - Is it *too close* to the training examples?

# Assesing whta we have learned

- Training and test sets (*test*)
  - Learn with the training set
  - Measure the performance using the test set:
    - *Classification accuracy*: percent of correctly classified examples in the test set
- Repeat this process
  - Learning curve
  - Stratification: each class has to be properly represented both in the training set and in the test set
- Sometimes, it is not possible to “waste” the examples in the test set and we do *cross-validation*:
  - Partition the set of examples in  $k$  equal sized parts. Apply the learning algorithm  $k$  times, each one taking as test set one of the parts and learning with the rest. Return the average of the  $k$  classification accuracies.
  - In practice: cross-validation, with  $k = 10$  and stratification

# Overfitting and noise

- We say that a learned model is *overfitted* if it is so fitted to the training set that it has a poor predictive performance on the rest of unseen examples
- *Noise*: examples in the training set, incorrectly classified. This causes overfitting
- Example: assume that by mistake, we include the following example  
<BLUE,ROUND,SMAL> as a negative example
- The learned model in this case would be (overfitted to the data):



# Overfitting and noise

- Other sources of overfitting:
  - Attributes presenting an apparent regularity in the training set, but actually irrelevant
  - Small training sets
- Ways to avoid overfitting in decision trees:
  - Early stopping
  - *A posteriori* tree pruning
- *A posteriori* pruning, two approaches:
  - Transformation to rules, prune rule conditions
  - Prune the tree
  - Only prune when improving classification accuracy on a test set

## Reduced error pruning

1. Partition the set of examples into Training and Test sets
2. Tree=tree obtained using ID3 on the Trainig set
3. Continue=True
4. While Continue:
  - \* Measure = classification accuracy of the leraned tree on Test
  - \* For every internal node N of Tree:
    - Temporarily prune the subtree of Tree at node N, replacing it by a leaf labeled with the majority clase at that node
    - Compute the classification accuracy of the pruned tree, on the test set
  - \* Let K the node with best accuracy
  - \* If its accuracy is better than Measure, then  
Tree = permanently prune Tree at node K
  - \* Otherwise, Continue=False
5. Return Tree



# More practical questions for tree learning

- ID3 extensions:
  - Real-valued attributes
  - Other measures for selecting the best node
  - Error estimates
  - Missing values in attributes
  - Attributes with cost
- Algorithms C4.5 and C5.0 (Quinlan)

## Section 3

### Learning rules

# Changing the formalism

- Classification rules:
  - R1: If  $\text{OUTLOOK}=\text{SUNNY} \wedge \text{HUMIDITY}=\text{HIGH}$   
Then  $\text{PLAYTENNIS} = -$
  - R2: If  $\text{ASTIGMATISM}=+ \wedge \text{TEAR\_RATE}=\text{NORMAL}$   
Then  $\text{LENSES}=\text{HARD}$
- Advantages of using rules:
  - Clarity
  - Modularity
  - Expressiveness: rules can represent every subset of instances
  - Methods can naturally be extended to first-order rules,
  - Formalism used in knowledge based systems
- Rules and decision trees
  - Easy translation from trees to rules, but not in the other way

# Rule learning

- Goal: learn a set of rules consistent with the examples
  - A rule *covers* an example if the example satisfies its conditions
  - It *correctly covers* the example if the class of the example is equal to the class in the rule conclusion
- Measuring the accuracy of a rule  $R$  with respect to a set of examples  $D$ :
  - *Relative frequency*:  $p/t$  (where  $t$  = examples in  $D$  covered by  $R$ ,  $p$  = examples correctly covered). Notation:  $RF(R, D)$
- Rule learning algorithms:
  - ID3 + translation to rules
  - Covering
  - Genetic algorithms

# Training set

EX.	AGE	PRESCRIPTION	ASTIGMATISM	TEAR_RATE	LENSES
$E_1$	YOUNG	MYOPE	-	REDUCED	NONE
$E_2$	YOUNG	MYOPE	-	NORMAL	SOFT
$E_3$	YOUNG	MYOPE	+	REDUCED	NONE
$E_4$	YOUNG	MYOPE	+	NORMAL	HARD
$E_5$	YOUNG	HYPERMETROPE	-	REDUCED	NONE
$E_6$	YOUNG	HYPERMETROPE	-	NORMAL	SOFT
$E_7$	YOUNG	HYPERMETROPE	+	REDUCED	NONE
$E_8$	YOUNG	HYPERMETROPE	+	NORMAL	HARD
$E_9$	PRE-PRESBYOPIC	MYOPE	-	REDUCED	NONE
$E_{10}$	PRE-PRESBYOPIC	MYOPE	-	NORMAL	SOFT
$E_{11}$	PRE-PRESBYOPIC	MYOPE	+	REDUCED	NONE
$E_{12}$	PRE-PRESBYOPIC	MYOPE	+	NORMAL	HARD
$E_{13}$	PRE-PRESBYOPIC	HYPERMETROPE	-	REDUCED	NONE
$E_{14}$	PRE-PRESBYOPIC	HYPERMETROPE	-	NORMAL	SOFT
$E_{15}$	PRE-PRESBYOPIC	HYPERMETROPE	+	REDUCED	NONE
$E_{16}$	PRE-PRESBYOPIC	HYPERMETROPE	+	NORMAL	NONE

# Un conjunto de entrenamiento

EJ.	AGE	PRESCRIPTION	ASTIGMATISM	TEAR_RATE	LENSES
$E_{17}$	PRESBYOPIC	MYOPE	-	REDUCED	NONE
$E_{18}$	PRESBYOPIC	MYOPE	-	NORMAL	NONE
$E_{19}$	PRESBYOPIC	MYOPE	+	REDUCED	NONE
$E_{20}$	PRESBYOPIC	MYOPE	+	NORMAL	HARD
$E_{21}$	PRESBYOPIC	HYPERMETROPE	-	REDUCED	NONE
$E_{22}$	PRESBYOPIC	HYPERMETROPE	-	NORMAL	SOFT
$E_{23}$	PRESBYOPIC	HYPERMETROPE	+	REDUCED	NONE
$E_{24}$	PRESBYOPIC	HYPERMETROPE		NORMAL	NONE

- R2: If  $ASTIGMATISM = + \wedge TEAR\_RATE = NORMAL$   
Then  $LENSES = HARD$
- R2 covers  $E_4, E_8, E_{12}, E_{16}, E_{20}$  y  $E_{24}$ , but only  $E_4, E_8, E_{12}$  y  $E_{20}$  are correctly covered

# Learning rules to cover examples

- Learn a rule to classify  $\text{LENSES}=\text{HARD}$

- If ?

Then  $\text{LENSES}=\text{HARD}$

- Different choices for ?, and the relative frequency of the resulting rule:

AGE=YOUNG	2/8
AGE=PREPRESBYOPIC	1/8
AGE=PRESBYOPIC	1/8
PRESCRIPTION=MYOPE	3/12
PRESCRIPTION=HYPERMETROPE	1/12
ASTIGMATISM=-	0/12
ASTIGMATISM=+	4/12 *
TEAR_RATE=REDUCED	0/12
TEAR_RATE=NORMAL	4/12 *

- Partially learned rule:

- If  $\text{ASTIGMATISM}=+$

Then  $\text{LENSES}=\text{HARD}$

# Learnig rules to cover examples

- We continue, to exclude incorrectly covered examples:

- If  $\text{ASTIGMATISM} = + \wedge ?$

Then  $\text{LENSES} = \text{HARD}$

- Choices for ?, and relative frequency of the resulting rule:

$\text{AGE} = \text{YOUNG}$	$2/4$
$\text{AGE} = \text{PREPRESBYOPIC}$	$1/4$
$\text{AGE} = \text{PRESBYOPIC}$	$1/4$
$\text{PRESCRIPTION} = \text{MYOPE}$	$3/6$
$\text{PRESCRIPTION} = \text{HYPERMETROPE}$	$1/6$
$\text{TEAR\_RATE} = \text{REDUCED}$	$0/6$
$\text{TEAR\_RATE} = \text{NORMAL}$	$4/6$ *

- Partially learned rule:

- If  $\text{ASTIGMATISM} = + \wedge \text{TEAR\_RATE} = \text{NORMAL}$

Then  $\text{LENSES} = \text{HARD}$



# Learning rules to cover examples

- We continue, to exclude incorrectly covered examples:
  - If  $\text{ASTIGMATISM} = + \wedge \text{TEAR\_RATE} = \text{NORMAL} \wedge ?$   
Then  $\text{LENSES} = \text{HARD}$
  - Choices for ?, and relative frequency of the resulting rule:

AGE=YOUNG	2/2 *
AGE=PREPRESBYOPIC	1/2
AGE=PRESBYOPIC	1/2
PRESCRIPTION=MYOPE	3/3 *
PRESCRIPTION=HYPERMETROPE	1/3

- Rule finally learned (no incorrectly covered examples)
  - If  $\text{ASTIGMATISM} = + \wedge \text{TEAR\_RATE} = \text{NORMAL} \wedge$   
     $\text{PRESCRIPTION} = \text{MYOPE}$   
Then  $\text{LENSES} = \text{HARD}$

# Learning rules to cover examples

- Still an uncovered example with  $\text{LENSES}=\text{HARD}$ 
  - We start again “If ? Then  $\text{LENSES}=\text{HARD}$ ”, but now with  $D' = D \setminus \{E_4, E_{12}, E_{20}\}$
- Rules finally learned for  $\text{LENSES}=\text{HARD}$ :
  - R1: If  $\text{ASTIGMATISM} = + \wedge \text{TEAR\_RATE} = \text{NORMAL} \wedge \text{PRESCRIPTION} = \text{MYOPE}$   
Then  $\text{LENSES} = \text{HARD}$
  - R2: If  $\text{AGE} = \text{YOUNG} \wedge \text{ASTIGMATISM} = + \wedge \text{TEAR\_RATE} = \text{NORMAL}$   
Then  $\text{LENSES} = \text{HARD}$
  - The rules correctly cover the 4 examples of  $\text{LENSES}=\text{HARD}$  (actually, they overlap)
- We now could continue learning rules for:
  - $\text{LENSES}=\text{SOFT}$
  - $\text{LENSES}=\text{NONE}$

# Sequential covering algorithm for learning rules

## Sequential covering algorithm

Sequential-Covering( $D, \text{Class}, c$ )

1. Initialize Learned-Rules to the empty set
2. Initialize  $E$  to  $D$
3. While  $E$  contains examples with Class  $c$ , do:
  - 3.1 Create a rule  $R$  with no conditions and conclusion  $\text{Class}=c$
  - 3.2 While there are examples in  $E$  incorrectly covered by  $R$  and there are attributes to use, do:
    - 3.2.1 Choose the BEST condition  $A=v$  to add to  $R$ , where  $A$  is an attribute not appearing in  $R$  and  $v$  is a value of  $A$
    - 3.2.2 Update  $R$  adding condition  $A=v$  to  $R$
  - 3.3 Add  $R$  to Learned-Rules
  - 3.4 Update  $E$ , removing the examples covered by  $R$
4. Return Learned-Rules

- Algorithm to learn rules from a training set  $D$ 
  - Rules predicting situations where the class is a given value  $c$

# Sequential covering algorithm: control

- Outer loop:
  - Adds rules (the model is *generalized*)
  - Each rule covers some examples, all of them correctly
  - In every step, the covered examples are removed
  - And rules are added while there are examples not covered
- Inner loop:
  - Adds conditions to the rule (the rule is *specialized*)
  - Each new condition excludes examples incorrectly covered
  - And this is done while there are examples incorrectly covered
- Sequential covering vs ID3
  - It learns one rule at a time; in contrast, ID3 learns globally
  - ID3: choose attributes
  - Sequential covering: choose pairs attribute-value

# Sequential covering (properties)

- There are several criteria to choose the “best” condition in the inner loop
  - Choose the condition with *the greatest relative frequency* (as shown in the example)
  - Choose the one with *greatest information gain*:

$$p \cdot (\log_2 \frac{p'}{t'} - \log_2 \frac{p}{t})$$

wher  $p'/t'$  is the relative frequency *after* adding the condition and  $p/t$  is the relative frequency *before* adding the condition

- The set of rules learned by the sequential covering algorithm *perfectly* fit the training set (*danger of overfitting*)
  - *A posteriori* rule pruning: incrementally remove conditions, until no *improvement* is produced
  - *A priori* rule pruning: early stopping rule generation

## Sección

### Instance-based learning: kNN

# Classification by means of nearest neighbour

- An alternative to building probabilistic models is to find a way to calculate the classification for a new instance directly from previous examples (*instance-based learning*)
- Idea: assign the classification of a new example by observing the categories of the “closest” examples.
  - We thus need a notion of “distance” between examples.
  - Most of the time we will handle examples belonging to  $R^n$ , and by default we may assume the Euclidean distance.
  - But of course any other particular distance function might be defined.
- Example of application: documents classification

- $k$ -NN algorithm (stands for “ $k$  nearest neighbours”):
  - Given a training set (numerical vectors with an associated category) and a new example
  - Return the category of the majority within the  $k$  examples of the training set which are the closest to the new example that we want to classify



# Distances for $k$ -NN

- Possible definitions of distance functions:
  - Euclidean:  $d_e(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
  - Manhattan:  $d_m(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n |x_i - y_i|$
  - Hamming: number of components which differ.
- Euclidean is suitable when every dimension refers to similar properties, and Manhattan otherwise; Hamming distance can be used even when handling non-numerical values.
- Normalization: when not all dimensions fall within the same order of magnitude, components should be normalized (by subtracting the average value and dividing by the standard deviation)

# Some observations on $k$ -NN

- Choice of  $k$ :
  - Usually, based on some specific knowledge about the classification problem
  - May be tuned after tests over smaller sets (validation sets)
  - When classification is binary, odd values avoid ties ( $k=5$ , for example)
- Variant of  $k$ -NN: for each class  $c$ , sum the similarity of every example of this class belonging to the  $k$  closests to the new example. Return the class which gets the highest sum.
  - In this way, the closer an example is, the more its classification is taken into account.

## Section 5

### Clustering

- As the last application of statistical learning, let us now overview *clustering* techniques
- The goal is to divide an input data set into several subsets (*clusters*), in such a way that elements within the same subset share some kind of pattern or some characteristics which are previously unknown
- We shall focus on the case when data are expressed as numbers or as numerical vectors, and *the number of clusters is given in advance*
- *Unsupervised* learning: we have no prior information about which is the “correct” cluster that corresponds to each data.
- Application of clustering:
  - Data mining
  - Digital image processing
  - Bioinformatics

# Two examples

- *Color quantization:*
  - A digital image stored with a resolution of 24 bits/pixel (aprox. 16 million colors) is to be displayed on a screen that allows only 8 bits/pixel (256 colors)
  - What is the best correspondence between the colors of the original image and the colors that can be displayed on the screen?
- Merge of distributions:
  - We have some data collection about the weight of people; data are not labelled so for each value we do not know if they refer to men or women. However, we know (or we assume) that there are two different weight distributions for men and women, and both of them follow a normal distribution
  - By just looking at the data, is it possible to *learn* which are the aforementioned two distributions?

# Clustering based on distance

- Idea: given the number  $k$  of *clusters*, find  $k$  points or *centers* which will represent each of the clusters, in such a way that each example will be assigned to the cluster corresponding to the “nearest” center.
- Just like before, distance can be specific for each problem:
  - Expressing the degree of similarity
  - Usually Euclidean

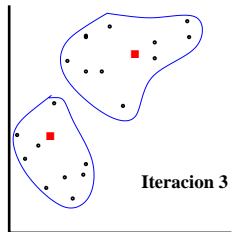
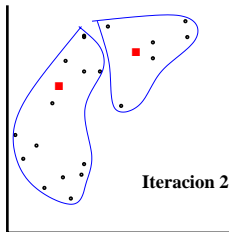
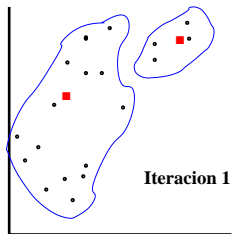
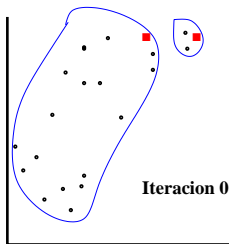
# Classic algorithm: $k$ -means

- Input: a number  $k$  of clusters, a data set  $\{x_i\}_{i=1}^N$  and a distance function
- Output: a set of  $k$  centers  $m_1, \dots, m_k$

## k-means( $k$ ,data,distance)

1. Initialize  $m_i$  ( $i=1, \dots, k$ ) (randomly or with some heuristic criterion)
2. REPEAT (until  $m_i$  do not vary):
  - 2.1 FOR  $j=1, \dots, N$ , DO:  
Calculate the cluster corresponding to  $x_j$ , by selecting, among all  $m_i$ , the center  $m_h$  such that the  $\text{distance}(x_j, m_h)$  is minimum
  - 2.2 FOR  $i=1, \dots, k$ . DO:  
Update  $m_i$  assigning to it the arithmetic mean of the examples assigned to  $i$ -th cluster
3. Return  $m_1, \dots, m_k$

# Graphical intuition for $k$ -means algorithm





# Example of $k$ -means algorithm

- Data about weight of the population: 51, 43, 62, 64, 45, 42, 46, 45, 45, 62, 47, 52, 64, 51, 65, 48, 49, 46, 64, 51, 52, 62, 49, 48, 62, 43, 40, 48, 64, 51, 63, 43, 65, 66, 65, 46, 39, 62, 64, 52, 63, 64, 48, 64, 48, 51, 48, 64, 42, 48, 41
- The algorithm, applied over  $k = 2$  and with Euclidean distance, finds two centers  $m_1 = 63,63$  and  $m_2 = 46,81$  on three iterations
- 19 data belong to the first cluster and 32 to the second one

# Other considerations on the $k$ -means algorithm

- Local search:
  - It can be seen as a local search algorithm, trying to find the centers  $m_i$  which optimize  $\sum_j \sum_i b_{ij} d(x_j, m_i)^2$ , where  $b_{ij}$  is 1 if  $x_j$  has  $m_i$  as its closest center, or 0 otherwise
  - As any other local search algorithm, does not guarantee finding the global optimum
- Initialization: random or with some heuristic technique (for example, randomly split the data into  $k$  clusters, and start the algorithm execution using as input the centers of such clusters)
- In practice, the centers provided as input may make a great difference with respect to the quality of the obtained results

# Another example of $k$ -means algorithm

- The file **iris.arff** of WEKA contains 150 data vectors about the petal and sepal width and length of iris plants, which are classified into three different subtypes (setosa, versicolor and virginica)
  - Example of instance within **iris.arff** file: **5.1,3.5,1.4,0.2,Iris-setosa**
- We can apply  $k$ -means, with  $k = 3$  and Euclidean distance, ignoring the last attribute (pretending it is unknown):
  - After 6 iterations it stabilizes
  - Of the three resulting clusters, the first of them is composed exactly by the 50 instances which were originally classified as “iris setosa”
  - The second cluster includes 47 versicolor and 3 virginicas
  - The third includes 14 versicolor and 36 virginicas
  - It was not able to correctly tell apart versicolor and virginica specimens

- Mitchell, T.M. *Machine Learning* (McGraw-Hill, 1997)
  - Chapters 3,8 and 10
- Russell, S. and Norvig, P. *Artificial Intelligence (A Modern Approach)* (3rd edition) (Prentice Hall, 2010)
  - Sections 18.1, 18.2, 18.3 and 18.8
- Witten, I.H. y Frank, E. *Data mining* (Third edition) (Morgan Kaufmann Publishers, 2011)
  - Chapters 3, 4, 5 and 6.
- Alpaydin, E. *Introduction to Machine Learning* (third edition) (The MIT Press, 2014)
  - Chapter 9, and Section 7.3