

Exercises of Machine Learning and Inductive Inference

See Toledo: “Machine Learning and Inductive Inference: **Lecture**” [H02C1a], folder “Exercise sessions”.

Questions?

Vladimir Dzyuba
E-Mail: Vladimir.Dzyuba@cs.kuleuven.be
Office: 200A 04.034

Davide Nitti
E-Mail: davide.nitti@kuleuven.be
Office: 200A 02.65

Jérôme Renaux
E-Mail: Jerome.Renaux@cs.kuleuven.be
Office: 200A 04.002

1 Exercise Session 1: Concept learning and decision trees

1.1 Generality

Order the following concepts over the boolean attributes A and B according to generality (\geq_g): $A \wedge B$, $A \vee B$, $A \text{ xor } B$, T (true), F (false).

1.2 Version Spaces

Consider the hypothesis space shown in Figure 1 (ordered in a lattice using \geq_g). A Version Space VS is defined by its most general border $G = \{h_1, h_2\}$ and its most specific border $S = \{h_3, h_4\}$.

- Mark all hypotheses that belong to VS (on Figure 1).
- We say that a version space classifies an instance as positive if all the hypotheses in the version space predict it to be positive, negative if all the hypotheses in the version space predict negative, and “don’t know” in all other cases. Let P_i be the instances predicted positive by h_i , and N_i the instances predicted negative by h_i (for $i = 1, 2, 3, 4$). Answer the following questions in terms of P_i and N_i .
 - Which instances are classified positive by VS ?
 - Which instances are classified negative by VS ?
 - Which instances are classified “don’t know” by VS ?

1.3 Decision Trees Representing Logical Concepts

Represent as a decision tree:

- $A \wedge \neg B$
- $\neg A \vee B$
- $A \vee (B \wedge C)$
- $(\neg A \wedge B) \vee (A \wedge \neg B)$
- $(A \vee B) \wedge (C \vee D \vee \neg E)$
- $(A \vee B \vee C) \wedge (D \vee E \vee F)$

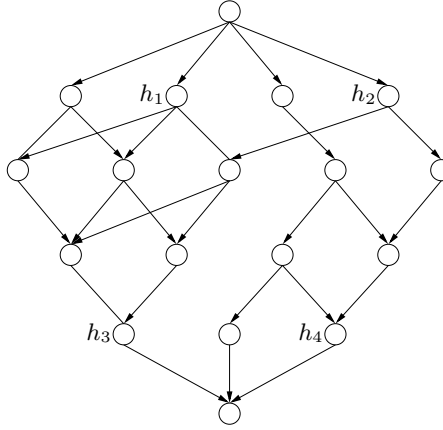


Figure 1: Hypothesis space for exercise 1.2.

1.4 Decision Trees and Generality

Decision tree $D2$ is an *elaboration* of decision tree $D1$ if $D2$ can be constructed out of $D1$ by replacing a leaf of $D1$ with a subtree in $D2$.

1. Which of the trees a and d from Exercise 1.3 is the elaboration of the other?
2. True or false? If a boolean decision tree $D2$ is an elaboration of $D1$ then $D1$ is more general than $D2$.

1.5 Decision Surface

Consider a data set with two *numeric* attributes a_1 and a_2 and one nominal target attribute c with two possible values: \oplus and \ominus . The training examples are shown in Figure 2.

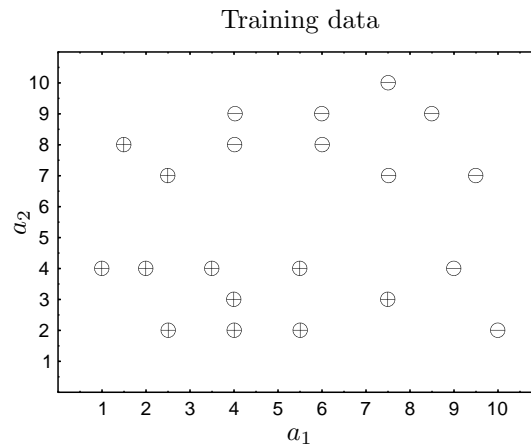


Figure 2: Training data for exercise 1.5

- a. Find a decision tree that classifies all training examples correctly.
- b. Draw the decision surface of this tree on Figure 2.

1.6 Entropy and Information Gain

Consider the following table of training examples:

Instance	Classification	a_1	a_2
1	\oplus	T	T
2	\oplus	T	T
3	\ominus	T	F
4	\oplus	F	F
5	\ominus	F	T
6	\ominus	F	T

- What is the entropy with respect to the “classification” attribute? (Answer without a calculator.)
- What is the information gain of a_2 relative to these training examples? (Answer without a calculator.)
- What is the information gain of a_1 relative to these training examples? (Use a calculator.)

1.7 The ID3 Algorithm

- Show a decision tree that could be learned by ID3 assuming it gets the following examples:

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	sunny	warm	normal	strong	warm	same	yes
2	sunny	warm	high	strong	warm	same	yes
3	rainy	cold	high	strong	warm	change	no
4	sunny	warm	high	strong	cool	change	yes

- Add this example:

5	sunny	warm	normal	weak	warm	same	no
---	-------	------	--------	------	------	------	----

then show how ID3 would induce a decision tree for these 5 examples.

Compute the class entropy for the entire dataset S :

	+	-	Entropy
S	3	2	

Compute the split heuristic for each attribute and select the attribute for the root node:

Attribute	Values	+	-	Entropy	IG
Sky	Sunny	3	1	0.811	0.32
	Rainy	0	1	0.000	
AirTemp	Warm				
	Cold				
Humidity	Normal	1	1	1.000	0.02
	High	2	1	0.918	
Wind	Strong				
	Weak				
Water	Warm				
	Cool				
Forecast	Same				
	Change				

Follow the ID3 algorithm until you obtain a complete decision tree.

1.8 Regression Tree

Consider the following set S of training examples with a numeric target attribute.

Instance	Target	a_1	a_2
1	1.0	T	T
2	1.0	T	T
3	1.0	T	F
4	5.0	F	F
5	6.0	F	T
6	5.5	F	T

We are going to use *weighted average variance of the subsets* as the split heuristic H . For each attribute A , compute:

$$H(A) = \sum_{v \in \text{Values}(A)} \frac{|S_{A=v}|}{|S|} \cdot \text{Var}[S_{A=v}] \quad (\text{heuristic value for attribute } A), \text{ where}$$

$$\text{Var}[S_{A=v}] = \frac{1}{|S_{A=v}|} \cdot \sum_{i \in S_{A=v}} (Target_i - \overline{Target}[S_{A=v}])^2 \quad (\text{variance of } Target \text{ values of instances with } A = v)$$

$$\overline{Target}[S_{A=v}] = \frac{1}{|S_{A=v}|} \cdot \sum_{i \in S_{A=v}} Target_i \quad (\text{average } Target \text{ value of instances with } A = v)$$

Which attribute (a_1 or a_2) will be put in the top node of the regression tree?

1.9 Using Weka: Homework

The explanation of the Weka Tool in this exercise is kept to a minimum. For more information you can look online and download a manual or tutorial from <http://www.cs.waikato.ac.nz/ml/weka/> or you can consult the book: “Data Mining: Practical Machine Learning Tools and Techniques” by Ian H. Witten and Eibe Frank (ISBN 978-0123748560, available in the KU Leuven library).

1. Download the Excel document `weka-zoo.xls` from Toledo (Course Documents/Exercise Sessions/Session 1/Exercise session 1/). Study the animals in the document and without using a data mining tool, draw a decision tree of three to five levels deep that classifies animals into a mammal, bird, reptile, fish, amphibian, insect or invertebrate.
2. Download and install Weka from: <http://www.cs.waikato.ac.nz/ml/weka/>
3. Read about the ARFF-format at <http://www.cs.waikato.ac.nz/ml/weka/arff.html> and construct the header for the animal file.
4. Download and unzip the file `weka-datasets.zip` from Toledo (Course Documents/Exercise Sessions/Session 1/Exercise session 1/). Open the `zoo.arff` file in the Weka Explorer. Find out how many animals this dataset contains.
5. Go to the classifier tab and select the decision tree classifier j48. Click on the line behind the choose button. This shows you the parameters you can set and a button called 'More'. Which algorithm is implemented by j48?
6. Click the start button to run the j48 algorithm. Which percentage of instances is correctly classified by j48? Which families are mistaken for each other? (Hint: Take a look at the confusion matrix.)