

## MACHINE LEARNING (SS2012)

Prof. Dr. M. Riedmiller, Manuel Blum

### Exercise Sheet 1

#### Exercise 1.1: Introduction to Machine Learning

- (a) Visit the following website:  
<http://www.pacman-vs-ghosts.net>  
Think of different aspects of the Pacman Task, that could be solved using Machine Learning algorithms.
- (b) Given an appropriate dataset, the problems given below can be solved by Machine Learning algorithms. Which of the problems would you apply supervised learning to?
- (c) There are two types of supervised learning tasks, classification and regression. Decide for the supervised learning problems given below, whether it is a classification problem or not.

Task	Supervised Learning	Classification
Predict tomorrow's price of a particular stock.	<input type="checkbox"/>	<input type="checkbox"/>
Discover whether there are different types of spam mail and what categories there are.	<input type="checkbox"/>	<input type="checkbox"/>
Predict your life expectancy.	<input type="checkbox"/>	<input type="checkbox"/>
Predict if it is going to rain tomorrow.	<input type="checkbox"/>	<input type="checkbox"/>
Learn to grasp an object by trial and error.	<input type="checkbox"/>	<input type="checkbox"/>

Table 1: Problems that can be solved using Machine Learning techniques.

#### Exercise 1.2: Version Spaces and Conjunctive Hypotheses

- (a) What are the elements of the version space? How are they ordered? What can be said about the meaning and sizes of  $S$  and  $G$ ?
- (b) In the following, it is desired to describe whether a person is *ill*. We use a representation based on conjunctive constraints (three per subject) to describe individual person. These constraints are “running nose”, “coughing”, and “reddened skin”, each of which can take the value true (+) or false (-). We say that somebody is

ill, if he is coughing and has a running nose. Each single symptom individually does not mean that the person is ill.

Specify the space of hypotheses that is being managed by the version space approach. To do so, arrange all hypotheses in a graph structure using the more-specific-than relation.

- (c) Apply the candidate elimination (CE) algorithm to the sequence of training examples specified in Table 2 and name the contents of the sets  $S$  and  $G$  after each step.

Training	running nose	coughing	reddened skin	Classification
$d_1$	+	+	+	positive (ill)
$d_2$	+	+	-	positive (ill)
$d_3$	+	-	+	negative (healthy)
$d_4$	-	+	+	negative (healthy)
$d_5$	-	-	+	negative (healthy)
$d_6$	-	-	-	negative (healthy)

Table 2: List of training instances for the medical diagnosis task.

- (d) Does the order of presentation of the training examples (according to Table 2) to the learner affect the finally learned hypothesis?
- (e) Assume a domain with two attributes, i.e. any instance is described by two constraints. How many positive and negative training examples are minimally required by the candidate elimination algorithm in order to learn an arbitrary concept?
- (f) We are now extending the number of constraints used for describing training instances by one additional constraint named “fever”. We say that somebody is ill, if he has a running nose and is coughing (as we did before), or if he has fever.

How does the version space approach using the CE algorithm perform now, given the training examples specified in Table 3? What happens, if the order of presentation of the training examples is altered?

Training	running nose	coughing	reddened skin	fever	Classification
$d_1$	+	+	+	-	positive (ill)
$d_2$	+	+	-	-	positive (ill)
$d_3$	-	-	+	+	positive (ill)
$d_4$	+	-	-	-	negative (healthy)
$d_5$	-	-	-	-	negative (healthy)
$d_6$	-	+	+	-	negative (healthy)

Table 3: List of training instances using the extended representation.

### Exercise 1.3: Decision Tree Learning with ID3

- (a) Apply the ID3 algorithm to the training data provided in Table 4.
- (b) Does the resulting decision tree provide a disjoint definition of the classes?

Training	fever	vomiting	diarrhea	shivering	Classification
$d_1$	no	no	no	no	healthy (H)
$d_2$	average	no	no	no	influenza (I)
$d_3$	high	no	no	yes	influenza (I)
$d_4$	high	yes	yes	no	salmonella poisoning (S)
$d_5$	average	no	yes	no	salmonella poisoning (S)
$d_6$	no	yes	yes	no	bowel inflammation (B)
$d_7$	average	yes	yes	no	bowel inflammation (B)

Table 4: Multi-class training examples.

- (c) Consider the use of real-valued attributes, when learning decision trees, as described in the lecture. Table 5 shows the relationship between the body height and the gender of a group of persons (the records have been sorted with respect to the value of *height* in cm). Calculate the information gain for potential splitting thresholds and determine the best one.

<i>Height</i>	161	164	169	175	176	179	180	184	185
<i>Gender</i>	F	F	M	M	F	F	M	M	F

Table 5: Data on the correlation between body height and gender.