

University of Wisconsin-Madison
Computer Sciences Department

CS 760 — Machine Learning

Spring 1995

Midterm Exam

(two pages of notes allowed)

100 points, 90 minutes

May 3, 1995

Write your answers on these pages and show your work. If you feel that a question is not fully specified, state any assumptions you need to make in order to solve the problem. You may use the backs of these sheets for scratch work. Notice that all questions do not have the same point-value. Divide your time appropriately.

Before starting, write your name on this and all other pages of this exam. Also, make sure your exam contains six (6) problems on seven (7) pages.

Problem	Score	Max Score
1	_____	20
2	_____	15
3	_____	20
4	_____	15
5	_____	10
6	_____	20
Total	_____	100

1. Decision Trees (20 pts)

Part A. Assume you are given the following three nominal features with the possible values shown.

$$\begin{aligned} F1 &\in \{v1, v2\} \\ F2 &\in \{v3, v4, v5\} \\ F3 &\in \{v6, v7, v8, v9\} \end{aligned}$$

Using ID3 and its max-gain formula, produce a decision tree that accounts for the following examples. Break ties by choosing the lowest numbered feature. **Show all your work.**

F1 = v1	F2 = v3	F3 = v7	category = +
F1 = v1	F2 = v5	F3 = v8	category = +
F1 = v1	F2 = v4	F3 = v9	category = -
F1 = v2	F2 = v3	F3 = v9	category = -
F1 = v2	F2 = v4	F3 = v8	category = -

Part B. Imagine you have trained a neural network to recognize a given Boolean concept, and are now interested in a “comprehensible” description of the network’s representation of the concept. People claim that decision trees are comprehensible, at least if they are not too big. Briefly describe how one could use the ID3 algorithm to produce a description of the concept the network learned. State any assumptions you have to make, and discuss one strength and one weakness of your proposed approach.

2. Neural Networks (15 pts)

Part A. Based on the Bayesian analysis of what a neural network should optimize, under what conditions is it appropriate to use the squared-difference (between the network's and the teacher's output vectors) as the error function? In this case, what semantics should one give to the network's output? Finally, briefly describe how the Bayesian interpretation addresses *overfitting avoidance*.

Part B. Consider a variant of the traditional learning-from-examples paradigm where the teacher provides, for each input vector, both the desired output and the *slope* of the desired function with respect to one or more of the inputs. Sketch how this information about the slope of the function being learned could be used in a neural network that is to be trained by backpropagation. Hypothesize about what impact using this extra information might have.

3. Reinforcement Learning (20 pts)

Part A. Imagine an environment containing two light bulbs, where each light bulb has a switch that toggles between being on or off; initially both light bulbs are off. The actor in this environment is able to see both light bulbs and at each step has to flip the switch of one of them. Assume the actor always receives a reward of +1 when turning *on* the first light bulb, +2 when turning *on* the second light bulb, and -1 when turning *off* the first light bulb.

Apply the one-step, Q-learning algorithm to this problem, using a *table* to represent your Q-function (all entries in the table should initially be zero); let $\gamma=0.9$. In the space below, show the state of the Q-table after the first three (3) steps of the learner. For simplicity, always follow the current policy during learning (i.e., *no* exploration) and break ties by choosing the first light bulb. Briefly explain why each of the three steps was made.

Part B. Under Part A's simplifying assumptions, will the optimal policy be learned in the limit? In general, what is needed to ensure that a Q-table learner produces the optimal policy in the limit?

4. Explanation-Based Learning (15 pts)

Part A. Consider the following EBL domain theory. Terms beginning with ?'s are implicitly universally-quantified variables.

$$\begin{array}{lll} A(?x, ?y) \wedge B(?y, ?x, ?z) & \rightarrow & C(?x, ?y, ?z) \\ D(?x, ?x) \wedge E(?y, ?x) & \rightarrow & B(?y, 1, ?x) \\ F(?x, ?x) \wedge F(?y, ?y) & \rightarrow & B(?x, ?y, 1) \end{array}$$

Assume the following problem-specific facts are asserted:

$$\begin{array}{llllll} A(1, 1) & A(3, 1) & D(1, 1) & E(2, 3) & F(2, 2) & F(3, 2) \\ A(2, 2) & A(3, 2) & D(3, 3) & E(1, 2) & F(2, 3) & F(3, 3) \end{array}$$

Explain, with a proof tree, that $C(3, 2, 1)$ is true. Draw to the *right* of your proof tree the corresponding *explanation structure*. Clearly indicate the necessary unifications by marking them with three parallel lines.

Assuming that predicates A, D, E, and F are operational, what rule would the EGGS algorithm produce? Explain your answer.

Part B. Briefly explain the difference between learning macro-operators and search-control rules. Describe a strength of each approach.

5. Computational Learning Theory (10 pts)

Consider an inductive learning task that involves two real-valued input features, x and y . Assume that the chosen hypothesis space is the set of all rectangles whose vertices are integer-valued, i.e., candidate concepts are of the form $(a \leq x \leq b) \wedge (c \leq y \leq d)$ where a , b , c , and d are *integers* between 0 and 9 (inclusive), and $a < b$ and $c < d$.

Part A. Provide an upper bound on the number of training examples needed so that, with probability 0.99, a concept that is consistent with the training set has an error rate of no more than 10%.

Part B. Consider some algorithm, A , that produces a concept consistent with the training examples. In addition to providing A with enough training examples, what else must hold in order for this “rectangle-learning” task to be PAC-learnable by A ?

6. Short Essays (20 pts)

Part A. Briefly explain the importance in machine learning of the following:

minimal description length (MDL) principle

t-tests

crossover

the utility problem

exploration vs. exploitation

Part B. Briefly describe two important roles for *tuning sets* in inductive learning. Can you think of any role for them in speedup learning? Explain your answer.