# Final: CS 7301
# Spring 2016

---

**The exam is closed book (2 cheat sheets allowed). If you run out of room for an answer, use an additional sheet (available from the instructor) and staple it to your exam. Questions marked with ∗∗ are the hardest, followed by questions marked with a ∗. Unmarked questions are relatively easy and a good strategy might be to attempt them before the questions marked with ∗ and ∗∗.**

---

- **NAME** _____

- **UTD-ID if known** _____

- **Time:** 2 hours 40 minutes.

| Question | Points | Score |
|---|---|---|
| Decision Trees | 10 | |
| Linear versus Non-linear functions | 10 | |
| Neural Networks | 10 | |
| Maximum Likelihood Estimation | 8 | |
| Computational Learning Theory | 10 | |
| AdaBoost | 12 | |
| Bayesian Networks: Inference | 10 | |
| Bayesian networks: Learning | 10 | |
| Support Vector Machines | 10 | |
| K-means | 10 | |
| Total: | 100 | |

## Question 1: Decision Trees  (10 points)

(a) (6 points)  **(True/False) You are given a propositional formula in conjunctive normal form (CNF). Recall that a formula is in CNF if it is a conjunction of clauses where a clause is a disjunction of literals and a literal is a propositional variable or its negation. For example, the following formula over four variables $\{X_1, X_2, X_3, X_4\}$ is in CNF:

$$(X_1 \vee \neg X_3) \wedge (X_3 \vee X_4) \wedge (\neg X_4 \vee X_2 \vee X_1)$$

Let $f$ be a CNF formula having $m$ clauses such that the length of each clause is at most $k$ (also called $k$-CNF). Is the following statement true or false?  $f$ can be represented by a decision tree having height $k$ and $m$ leaf nodes.  Explain your answer.  If your answer is True, describe an algorithm for constructing the decision tree given $f$.  If your answer is False, give a counter-example or provide a reasonable argument.

(b) (4 points) Which of the following statements are true or false? Explain your answer in one or two sentences.

- Reduced-Error pruning reduces the bias but increases the variance.

- The Entropy impurity is always preferable (in the sense that it will give better generalization accuracy) to the Gini impurity because the former is based on sound information-theoretic principles.
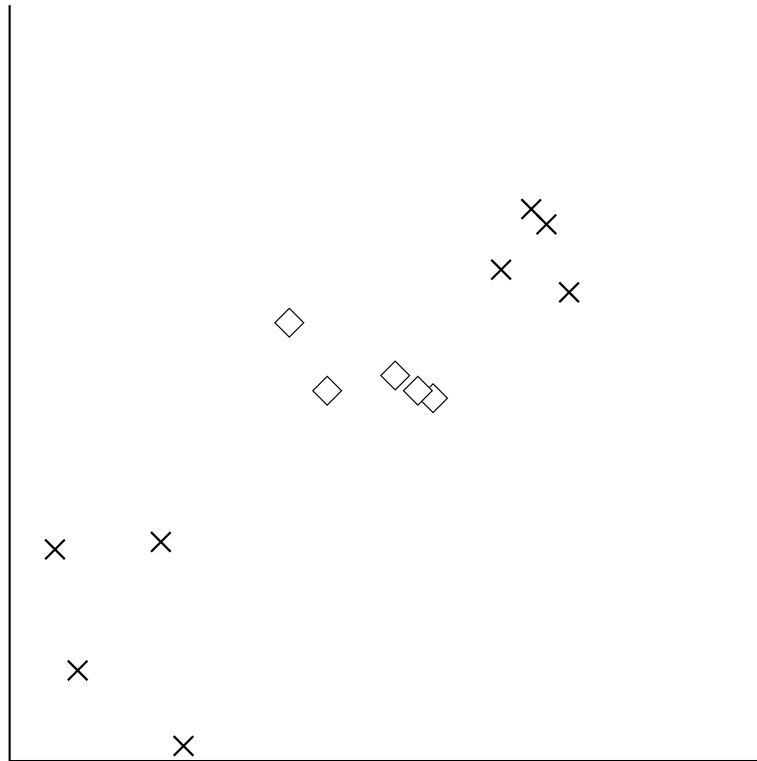
## Question 2: Linear versus Non-linear functions (10 points)

Recall that a linear threshold function or a linear classifier is given by: If $(w_0 + \sum_i w_i x_i) > 0$ then the class is positive, otherwise it is negative. Here, $x_1, \ldots, x_n$ are the attribute values and $w_0, \ldots, w_n$ are the weights ($w_0$ is the bias term). Assume that $1$ is true and $0$ is false.

(a) (5 points) * Consider the Boolean function given below. $x_1$, $x_2$ and $x_3$ are the attributes and $y$ is the class variable.

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | -1 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | -1 |
| 1 | 1 | 0 | -1 |
| 1 | 1 | 1 | -1 |

Can you represent the function using a linear threshold function. If your answer is **YES**, then give a precise numerical setting of the weights. Otherwise, clearly explain, why this function cannot be represented using a linear threshold function. No credit will be given if the explanation is incorrect.
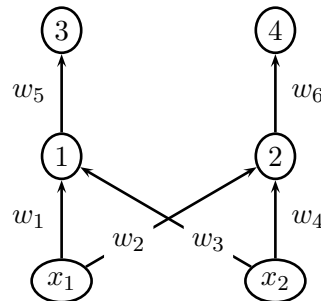
(b) (5 points) Which of the following will classify the data perfectly and why? Draw the decision boundary.

- Logistic Regression
- SVMs using a quadratic kernel.

## Question 3: Neural Networks  (10 points)

Consider the Neural network given below.



**Assume that all internal nodes and output nodes compute the** $tanh$ **function.** In this question, we will derive an explicit expression that shows how back propagation (applied to minimize the least squares error function) changes the values of $w_1$, $w_2$, $w_3$, $w_4$, $w_5$ and $w_6$ when the algorithm is given the example $(x_1, x_2, y_1, y_2)$ with $y_1$ and $y_2$ being outputs at 3 and 4 respectively (there are no bias terms). Assume that the learning rate is $\eta$. Let $o_1$ and $o_2$ be the output of the hidden units 1 and 2 respectively. Let $o_3$ and $o_4$ be the output of the output units 3 and 4 respectively.

Hint: Derivative of $tanh(x) = 1 - tanh^2(x)$.

 (a) (3 points)  Forward propagation. Write equations for $o_1$, $o_2$, $o_3$ and $o_4$.

 (b) (3 points)  Backward propagation. Write equations for $\delta_1$, $\delta_2$, $\delta_3$ and $\delta_4$ where $\delta_1$, $\delta_2$ , $\delta_3$ and $\delta_4$ are the values propagated backwards by the units denoted by 1, 2, 3 and 4 respectively in the neural network.

(c) (4 points) Give an explicit expression for the new (updated) weights $w_1$, $w_2$, $w_3$, $w_4$, $w_5$ and $w_6$ after backward propagation.

## Question 4: Maximum Likelihood Estimation  (8 points)

(a) (8 points)  *You are given a collection of $n$ documents, where the word count of the $i$-th document is $x_i$. Assume that the word count is given by an Exponential Distribution with parameter $\lambda$. In other words, for a non-negative integer $x$,

$$P(wordcount = x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$

Compute $\lambda$ such that the likelihood of observing $\{x_1, x_2, \ldots, x_n\}$ is maximized.

## Question 5: Computational Learning Theory (10 points)

Useful formulas for this section:

$$m \geq \frac{1}{\epsilon} \left( \ln(1/\delta) + \ln(|H|) \right)$$

$$m \geq \frac{1}{2\epsilon^2} \left( \ln(1/\delta) + \ln(|H|) \right)$$

$$m \geq \frac{1}{\epsilon} \left( 4 \log_2(2/\delta) + 8 VC(H) \log_2(13/\epsilon) \right)$$

(a) (5 points) Consider the class C of concepts of the form $(a \leq x_1 \leq b) \wedge (c \leq x_2 \leq d) \wedge (e \leq x_3 \leq f)$. Let $a, b$ be integers in the range $[0, 199]$ and $c, d, e, f$ be integers in the range $[0, 99]$. Give an upper bound on the number of training examples sufficient to assure that for any target concept $c \in C$, any **consistent learner** using $H = C$ will, with probability 0.99, output a hypothesis with error at most 0.05.

(b) (5 points) Consider the class C of concepts of the form $(a \leq x_1 \leq b) \wedge (c \leq x_2 \leq d)$. Suppose that $a, b, c, d$ take on real values instead of integers. Give an upper bound on the number of training examples sufficient to assure that for any target concept $c \in C$, a learner will, with probability 0.99, output a hypothesis with error at most 0.05.
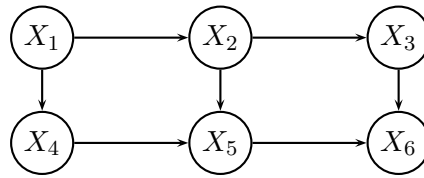
**Question 6: AdaBoost  (12 points)**

Consult the AdaBoost algorithm given on the last page for this question. Suppose you have two weak learners, $h_0$ and $h_1$, and a set of 17 points.

(a) (2 points)  You find that $h_1$ makes one mistake and $h_2$ makes four mistakes on the dataset. Which learner will AdaBoost choose in the first iteration (namely $m = 1$)? Justify your answer.

(b) (2 points)  What is $\alpha_1$?

(c) (2 points)  Calculate the data weighting co-efficients $w_2$ for the following two cases: (1) the points on which the chosen learner made a mistake and (2) the points on which the chosen learner did not make a mistake.

(d) (6 points) Consider a simple modification to the AdaBoost algorithm in which we normalize the data weighting co-efficients. Namely, we replace $w_n^{(m+1)}$ by $w_n^{(m+1)}/Z^{(m+1)}$ where $Z^{(m+1)} = \sum_{n=1}^{N} w_n^{(m+1)}$. Prove that $Z^{(m+1)} = 2(1 - \epsilon_m)$.

Hint: Notice that if the weights are normalized, then $\epsilon_m = \sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)$.
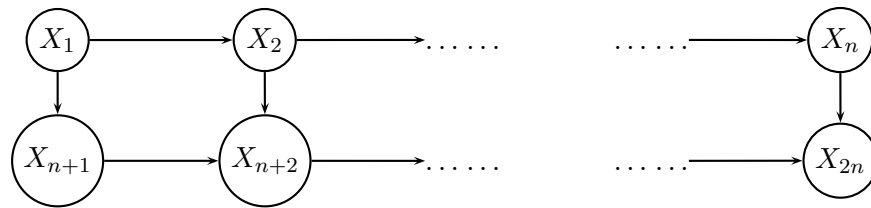
## Question 7: Bayesian Networks: Inference  (10 points)

Consider the Bayesian network given below:



(a) (5 points)  Let $X_6$ be the evidence variable. Trace the operation of **Variable elimination** for computing $\Pr(X_6 = x_6)$ along the order $(X_1, X_2, X_3, X_4, X_5)$ where $x_6$ is a value in the domain of $X_6$. Precisely show the various factors (functions) that will be generated. What is the time and space complexity of the algorithm along the ordering.

Consider a generalization of the Bayesian network given on the previous page:



(b) (5 points) Let $X_{2n}$ be the evidence variable. Assume that all variables have exactly $d$ values in their domain. What is the best (optimal) time and space complexity of computing the probability of evidence using the **Variable elimination algorithm** on the network given above. Be sure to mention the **best elimination order**, briefly explaining why it is optimal.

## Question 8: Bayesian networks: Learning  (10 points)

Consider a Bayesian network with edges $A \to B$ and $A \to C$, and parameters which are given below:

- $P(A = 1) = 0.9$
- $P(B = 1|A = 1) = 0.1$, $P(B = 1|A = 0) = 0.6$
- $P(C = 1|A = 1) = 0.7$, $P(C = 1|A = 0) = 0.3$

Consider the dataset given below:

| A | B | C |
|---|---|---|
| 0 | 1 | ? |
| 0 | 1 | 1 |
| ? | 0 | 1 |
| 1 | 1 | ? |
| 1 | 0 | ? |
| 0 | 0 | 0 |
| 1 | 1 | 1 |

Assume that the CPTs are the CPTs at some iteration of EM. What you are going to do is derive the new set of parameters after running one iteration of EM.

(a) (5 points) Show the calculations involved in the E-step. Recall that in the E-step, you make the dataset bigger (by considering all possible completions) and weigh each new data point appropriately.

(b) (5 points) Show the calculations involved in the M-step. What are the new parameters?

## Question 9: Support Vector Machines  (10 points)

Consider the following 2-D dataset ($x_1$ and $x_2$ are the attributes and $y$ is the class variable).

Dataset:

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 0 | $+1$ |
| 0 | 1 | $-1$ |
| 1 | 0 | $-1$ |
| 1 | 1 | $+1$ |

(a) (5 points)  Precisely write the expression for the dual problem. Let $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ be the lagrangian multipliers associated with the four data points.

(b) (5 points)  **Solve the dual problem. What is the value of $\alpha_1$, $\alpha_2$, $\alpha_3$ and $\alpha_4$.

## Question 10: K-means (10 points)

(a) (10 points) Give a formal proof that K-means always converges to local minima. Recall that K-means minimizes the following objective function:

$$\phi(\{x_i\}, \{a_i\}, \{c_j\}) = \sum_{i=1}^{n} dist(x_i, c_{a_i})$$

where $\{x_i\}$ is the set of $n$ points, $a_i \in \{1, \ldots, k\}$ gives the cluster to which the $i$-th point is assigned to, $c_j$ is the mean of the $j$-th cluster and $dist$ denotes the Euclidean distance.

Page left Blank

Page left Blank

Page left Blank

## AdaBoost

1. Initialize the data weighting coefficients $\{w_n\}$ by setting $w_n^{(1)} = 1/N$ for $n = 1, \ldots, N$.

2. For $m = 1, \ldots, M$:

   (a) Fit a classifier $y_m(\mathbf{x})$ to the training data by minimizing the weighted error function

   $$J_m = \sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \qquad (14.15)$$

   where $I(y_m(\mathbf{x}_n) \neq t_n)$ is the indicator function and equals 1 when $y_m(\mathbf{x}_n) \neq t_n$ and 0 otherwise.

   (b) Evaluate the quantities

   $$\epsilon_m = \frac{\displaystyle\sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\displaystyle\sum_{n=1}^{N} w_n^{(m)}} \qquad (14.16)$$

   and then use these to evaluate

   $$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \qquad (14.17)$$

   (c) Update the data weighting coefficients

   $$w_n^{(m+1)} = w_n^{(m)} \exp \left\{ \alpha_m I(y_m(\mathbf{x}_n) \neq t_n) \right\} \qquad (14.18)$$

3. Make predictions using the final model, which is given by

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^{M} \alpha_m y_m(\mathbf{x}) \right). \qquad (14.19)$$