

Q1 Score: 7.0

Total Score: 75.0

5B35B33E-ED27-41CE-8E70-267BC6F472CB

561-S16-midterm3

#323 2 of 13



### 1. [10%] General AI Knowledge

For each of the statements below, fill in the box T if the statement is always and unconditionally true, or fill in the box F if it is always false, sometimes false, or just does not make sense.

- a)  X
- b)  T
- c)  F
- d)  F
- e)  X
- f)  T
- g)  F
- h)  X
- i)  T
- j)  F

- a) If A is one of B's k-nearest-neighbors for a given value of k, then B must be one of A's k-nearest-neighbors.
- b) SVM can only classify data that is linearly separable.
- c) Assuming Boolean attributes, the depth of a decision tree, built using common algorithms such as ID3 (Iterative Dichotomiser 3), can never be larger than the number of training examples.
- d) Every Boolean function can be represented by some Bayesian network.
- e) Naive Bayes is a linear classifier.
- f) A Markov process is a random process in which the future is independent of the present, given the past.
- g) A single perceptron cannot compute the XOR function.
- h) For reinforcement learning, we need to know the transition probabilities between states before we start.
- i) In supervised learning, the examples given to the learner are not labeled.
- j) A perceptron is guaranteed to learn a given linearly separable function within a finite number of training steps.

**Q2A Score: 4.0**

357B0FCA-65BD-4C25-9BDE-E4B10316A4B0

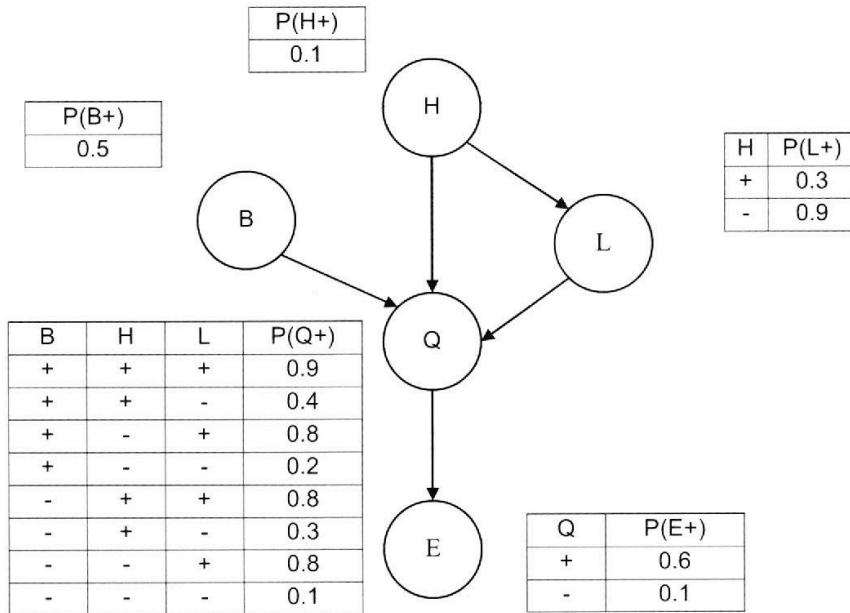
561-S16-midterm3

#323 3 of 13



## 2. [20%] Bayesian Networks

In the network below, the Boolean variables have the semantics: B: Brilliant, H: Honest, L: LotsOfFriends, Q: Qualified, E: Elected.



**2A. [6%]** Which of these, if any, are asserted by the structure of the network (leaving aside the conditional probability tables (CPTs))?

- |   |                                |  |
|---|--------------------------------|--|
| 1. <input checked="" type="checkbox"/> <input type="checkbox"/> F | $P(B, L) = P(B) P(L)$          |  |
| 2. <input type="checkbox"/> <input checked="" type="checkbox"/> T | $P(E   Q, L) = P(E   Q, L, H)$ |  |
| 3. <input type="checkbox"/> <input checked="" type="checkbox"/> T | $P(Q   B, H) = P(Q   B, H, L)$ |  |

**Q2B Score: 7.0**

99BF7DF5-E1CD-45A7-8ADD-43ACE1BA7383

561-S16-midterm3

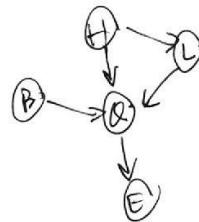
#323 4 of 13

**2B. [7%]** Calculate the value of  $P(B+, H+, L-, Q+, E-)$ . Show your work.

$$\begin{aligned} & P(B+, H+, L-, Q+, E-) \\ &= P(E-|Q+) P(Q+|B+, H+, L-) P(B+) P(H+) P(L-|H+) \\ &= 0.4 \times 0.4 \times 0.5 \times 0.1 \times 0.7 \\ &= 0.0056 \end{aligned}$$



Q2C Score: 5.0



8D064568-4A6A-4034-96AC-58B43FF2E529

561-S16-midterm3

#323 5 of 13

 $B+$ 

**2C. [7%]** Calculate the probability that a candidate is brilliant or not given that she is honest, does not have lots of friends, and gets elected. That is, calculate  $P(B | H+, L-, E+)$ . Show your work. (You need to give both  $P(B+ | H+, L-, E+)$  and  $P(B- | H+, L-, E+)$ )

$$P(B | H+, L-, E+) = P(B+ | H+, L-, E+) + P(B- | H+, L-, E+)$$

$$\textcircled{1} \quad P(B+ | H+, L-, E+)$$

$$= \frac{P(B+, H+, L-, E+)}{P(H+, L-, E+)} = \frac{P(H+) P(B+) P(L- | H+) [P(Q+ | B+, H+, L-) P(E+ | Q+)]}{P(H+) P(L- | H+) [P(B) [P(Q+ | B+, H+, L-) P(E+ | Q+)] + P(Q- | B+, H+, L-) P(E+ | Q-)]} \\ + P(B) [P(Q+ | B-, H+, L-) P(E+ | Q+)] + P(Q- | B-, H+, L-) P(E+ | Q-)]$$

$$= \frac{0.1 \times 0.5 \times 0.7 \times [0.4 \times 0.6 + 0.6 \times 0.1]}{0.1 \times 0.7 \times [0.5 \times (0.4 \times 0.6 + 0.6 \times 0.1) + 0.5 \times (0.3 \times 0.6 + 0.7 \times 0.1)]}$$

$$\textcircled{2} \quad P(B- | H+, L-, E+)$$

$$= \frac{P(H+) P(B-) P(L- | H+) [P(Q+ | B-, H+, L-) P(E+ | Q+)] + P(Q- | B-, H+, L-) P(E+ | Q-)]}{P(H+) P(L- | H+) [P(B+) [P(Q+ | B+, H+, L-) P(E+ | Q+)] + P(Q- | B+, H+, L-) P(E+ | Q-)] + P(B) [P(Q+ | B-, H+, L-) P(E+ | Q+)] + P(Q- | B-, H+, L-) P(E+ | Q-)]}$$

$$\Rightarrow P(B | H+, L-, E+) = \textcircled{1} + \textcircled{2} = \frac{0.000504 + 0.000441}{0.000945}$$

$$= 1$$

-2 for incorrect final answer

Q3A Score: 8.0

DFA0FC4C-7BF4-4A50-AC9F-D3D8B4D5EAA7

561-S16-midterm3

#323 6 of 13

**3. [23%] Decision Tree Learning**

You are given the task of learning to classify first names by gender. You are given a list of names labeled as female (F) or male (M) and you want to learn a classifier based on decision tree learning.

For a given name, let us define **L** as its length, **V** as its number of vowels and **C** as its number of consonants. We will consider that A-E-I-O-U-Y are vowels. The other letters are consonants.

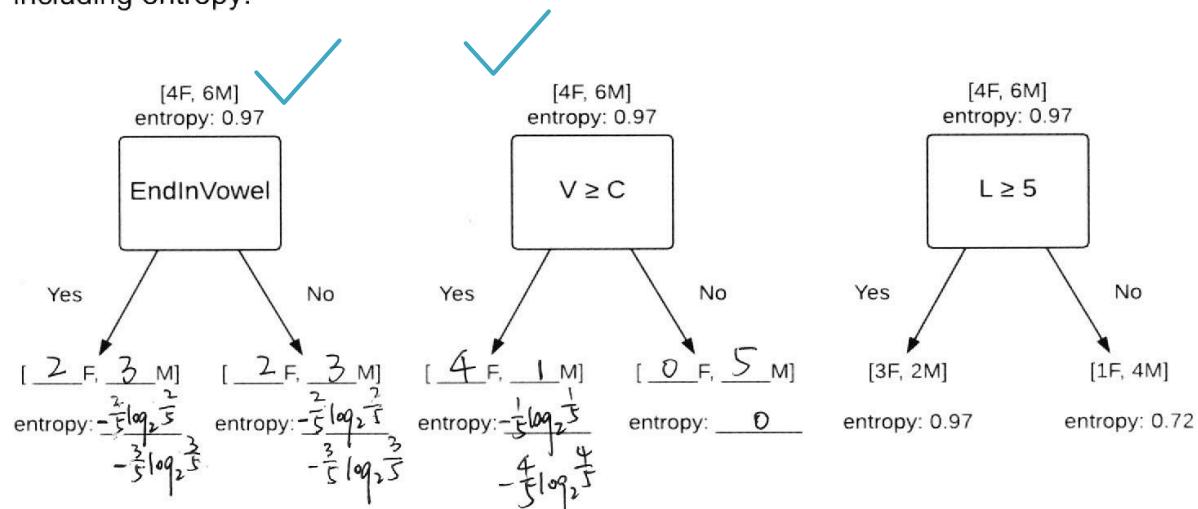
You decide to use the following features to predict the classes:

- *EndInVowel*: The name ends in a vowel.
- $V \geq C$ : The name has more vowels than consonants.
- $L \geq 5$ : The name contains 5 letters or more.

Name			Feature	Gender
	EndInVowel		$V \geq C$	
Annie	✓ Yes F		Yes F	Yes F
Brad	No M		No M	No M
Carl	No M		No M	No M
Daisy	✓ Yes F		Yes F	Yes F
Eleanor	✓ No P		Yes F	Yes F
Fernando	Yes M		No M	Yes M
Gary	✓ Yes M		Yes M	No M
Hans	No M		No M	No M
Isis	✓ No P		Yes F	No P
Jerry	Yes M		No M	Yes M

With 4 Female names and 6 Male names, the entropy of the decision in bits is 0.97.

**3A. [8%]** Consider the following decision trees, splitting on (*EndInVowel*), ( $V \geq C$ ), ( $L \geq 5$ ). The ( $L \geq 5$ ) tree has been filled out. Complete the values for the other features, including entropy.



## Q3BC Score: 5.0

737FAACC-9D14-4B06-A0BE-ADE9E010527F

561-S16-midterm3

#323 7 of 13



**3B. [6%]** Calculate the information gain for splitting on each of the 3 features. Show formulas and steps clearly.

$$\begin{aligned}
 \textcircled{1} \text{ IG(EndInVowel)} &= 0.97 - \left( \frac{1}{2} \cdot \text{Entropy(Yes)} + \frac{1}{2} \cdot \text{Entropy(No)} \right) \\
 &= 0.97 - \frac{1}{2} \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) - \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \\
 &= 0.97 + \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \quad \checkmark
 \end{aligned}$$
  

$$\begin{aligned}
 \textcircled{2} \text{ IG(V>c)} &= 0.97 - \left( \frac{1}{2} \cdot \text{Entropy(Yes)} + \frac{1}{2} \cdot \text{Entropy(No)} \right) \\
 &= 0.97 - \left[ \frac{1}{2} \left( \frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} \right) + 0 \right] = 0.97 + \frac{1}{10} \log_2 \frac{1}{5} + \frac{2}{5} \log_2 \frac{4}{5} \quad \checkmark
 \end{aligned}$$
  

$$\begin{aligned}
 \textcircled{3} \text{ IG(L>5)} &= 0.97 - \left( \frac{1}{2} \cdot \text{Entropy(Yes)} + \frac{1}{2} \cdot \text{Entropy(No)} \right) \\
 &= 0.97 - \frac{1}{2} (0.97 + 0.72) = 0.125 \quad \checkmark
 \end{aligned}$$

-3, result?

**3C. [2%]** Which attribute should you split on first? Justify your answer.

- ① Should use (V>c) to split on first. ✓
- ② After comparing the IG of these 3 attribute, the Information Gain of (V>c) is the highest, so we should use it first.

**Q3D Score: 2.0**

553EC66B-FA1B-4FD7-BABE-1259702CF3EE

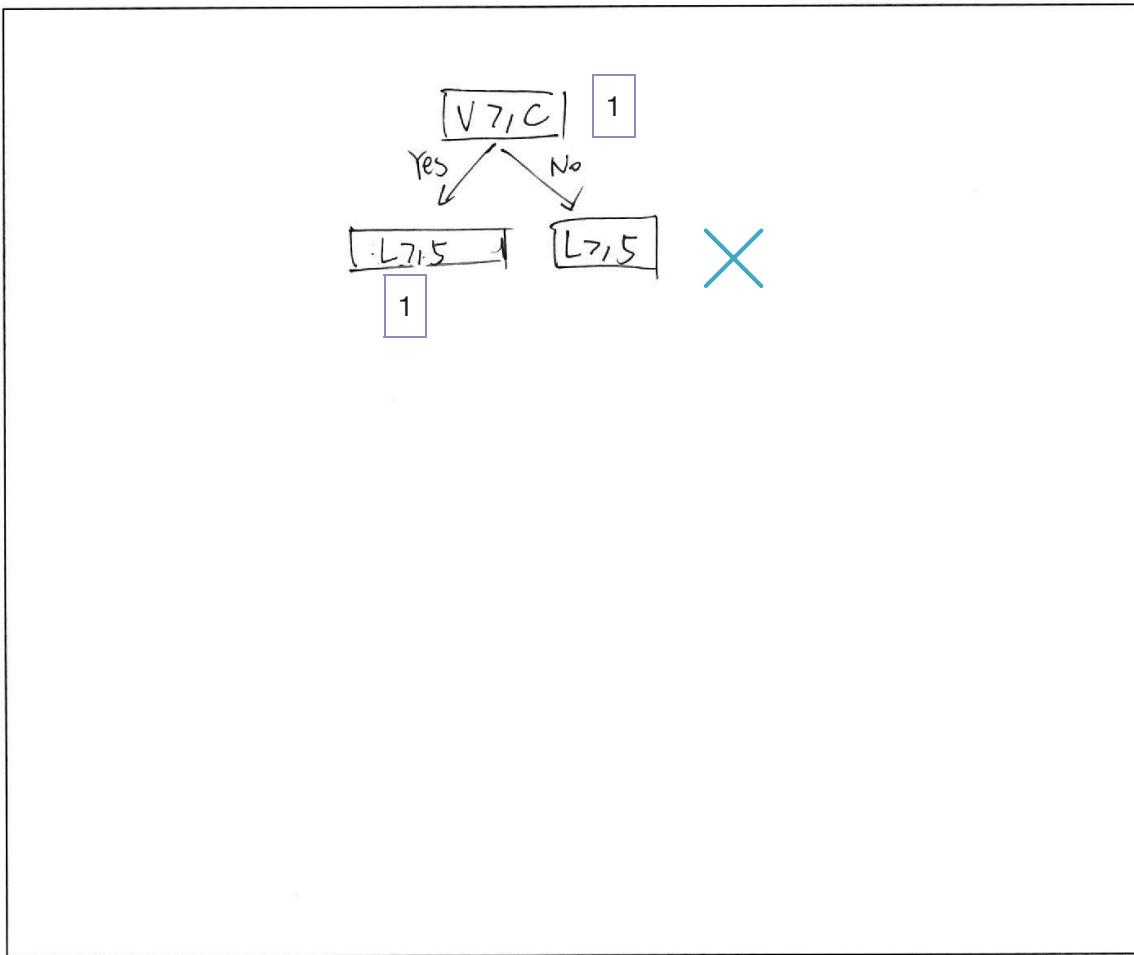
561-S16-midterm3

#323 8 of 13

**3D. [7%]** For the second level of the tree, you decide to use the following rule:

- split on attribute ( $V \geq C$ ) if it was not split on first
- split on attribute ( $L \geq 5$ ) otherwise

Draw the entire decision tree.



Q4A Score: 7.0

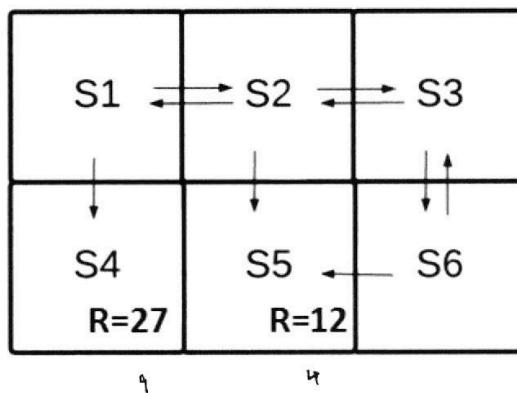
121E623E-CA16-4BCA-B7A5-B2DFF892F85B

561-S16-midterm3

#323 9 of 13

**4. [17%] Markov Decision Process**

Consider the 6-state Markov Decision process below. The goals with rewards are in state S4 and S5. At each state, the possible transitions are deterministic and indicated by the arrows. You get a reward of  $R_4=27$  if you get to the goal S4 and a reward of  $R_5=12$  if you get to the goal S5.



**4A. [7%]** Consider a discount factor of  $\gamma = 2/3$ . On the figure below, show the optimal value  $V^*$  for each state and the arrows corresponding to the set of optimal actions.

$\downarrow : 18$	$\leftarrow : 12$	$\leftarrow : 8$
S1	S2	S3
$\times : 27$	$\times 12$	$\leftarrow : 8$
S4	S5	S6

$$(1-\alpha)\delta + \alpha[R(\gamma) + r \text{Max} \delta]$$

$$R(\gamma) + \frac{2}{3} \times \delta'$$

+6 for correct  $V^*$  value

+1 for ALL correct arrows

**Q4BC Score: 10.0**

DS216F94-2111-47BA-BB5F-90087413BD78

561-S16-midterm3

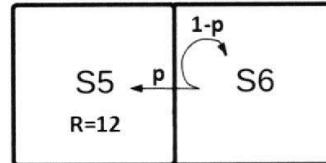
#323 10 of 13



**4B. [5%]** What values of  $\gamma$  would result in a different optimal action in S2? Indicate which policy action changes.

If  $S_2$  choose  $\leftarrow$ : the  ~~$\alpha = r^2 \cdot 27$~~   
else if  $S_2$  choose  $\downarrow$ : the  $\alpha = r \cdot 12$ .  
So, if we want to change policy action  $r^2 \cdot 27 < r \cdot 12$   
 $\therefore 0 < r < \frac{4}{9}$   
So, if  $0 < r < \frac{4}{9}$ , action of  $S_2$  will be  $\downarrow$ , value is  $r \cdot 12$

**4C. [5%]** In this question, you consider only states  $S_5$  and  $S_6$ . The transition is no longer deterministic. When going to  $S_5$  from  $S_6$ , you have a probability  $p$  of succeeding and a probability  $1-p$  of tripping, and staying in state  $S_6$ . What is the optimal value  $V^*$  at state  $S_6$  if the discount factor  $\gamma = 2/3$  and  $p = 1/4$ ?



$$\alpha = r \cdot \max(p \cdot 12, (1-p)r \cdot \alpha)$$

$$= \frac{2}{3} \cdot \max(3, \frac{3}{4} \times \alpha)$$

$$= \frac{2}{3} \times \cancel{6}$$

$$= 4$$

**Q5AB Score: 6.0**

B07356BF-FC88-4A7C-8E6C-4F1C6FCA62E6

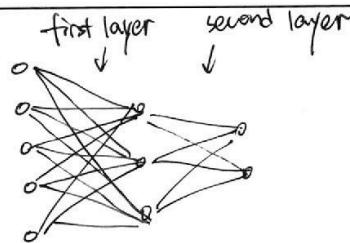
561-S16-midterm3

#323 11 of 13

**5. [20%] Neural Networks**

**5A. [4%]** How many weights does a 2-layer feed-forward neural network with 5 input units, 3 hidden units and 2 output units contain, including the biases (dummy input weights)? Show your work.

$$\begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix}$$



$$\text{first layer: } 5 \times 3 = 15$$

$$\text{second layer: } 3 \times 2 = 6$$

$$\Rightarrow 15 + 6 = 21 \text{ weights.}$$

Partial Credit +2

**5B. [4%] True or False.** T      

1. The back-propagation algorithm, when run until a minimum error is achieved, always converges to the same set of weights no matter what the initial set of weights is.

 T      

2. When choosing between two different neural network structures, we should always prefer the one with the lower error on the training set.

**Q5C Score: 12.0**

2CE41B69-D347-42C4-AB60-DSEC16E8419C

561-S16-midterm3

#323 12 of 13

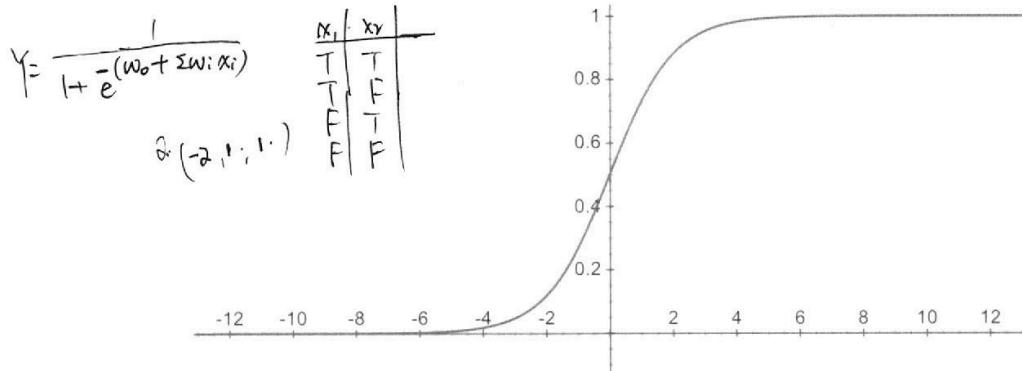


**5C. [12%]** Consider the neural network built out of units with real-valued inputs  $X_1 \dots X_n$ , where the unit output  $Y$  is given by

$$Y = \frac{1}{1 + \exp(-(w_0 + \sum_i w_i X_i))}$$

Here we will explore the expressiveness of neural nets, by examining their ability to represent Boolean functions. Here the inputs  $X_i$  will be 0 or 1. The output  $Y$  will be real-valued, ranging anywhere between 0 and 1. We will interpret  $Y$  as a Boolean value by interpreting it to be a Boolean 1 if  $Y > 0.5$ , and interpreting it to be 0 otherwise.

The figure for  $\frac{1}{1+e^{-x}}$  is:



Give 3 weights for a single unit with two inputs  $X_1$  and  $X_2$ , that implements the logical OR function  $Y = X_1 \vee X_2$  and the logical AND function  $Y = X_1 \wedge X_2$ , respectively.

Functions	$w_0$	$w_1$	$w_2$
Logical OR function $Y = X_1 \vee X_2$	-1	1	1
Logical AND function $Y = X_1 \wedge X_2$	-2	1	1

**Q6 Score: 2.0**

048D6F9C-BE7A-4613-AA1B-2CB07B42DCB4

561-S16-midterm3

#323 13 of 13

**6. [10%] AI Applications.** B

1. [2%] Which statement is true about cognitive architectures?
- a. A cognitive architecture is a hypothesis about the fixed structures that provide a mind.
  - b. A cognitive architecture tries to yield intelligent behavior in complex environments.
  - c. A generically cognitive architecture spans both the creation of artificial intelligence and the modeling of natural intelligence, at a suitable level of abstraction.
  - d. All of the above
  - e. None of the above

 D

2. [2%] In the task of randomly assigning air marshals to flights using game theory, which argument allows us to use an incremental strategy for scaling-up?

- a. The support set size is small: most variables are 0.
- b. The full rewards matrix is sparse.
- c. The computation can be parallelized.
- d. All of the above
- e. None of the above

 A

3. [2%] Which method can be used to solve a problem in which the utility function is not known?

- a. Reinforcement learning
- b. Markov Decision Process
- c. Perceptron learning
- d. All of the above
- e. None of the above

 A

4. [2%] In Natural Language Processing, which of these algorithms takes advantage of grammars to represent sentences as trees?

- a. Conditional Random Field (CRF)
- b. Cocke-Younger-Kasami (CYK)
- c. Hidden Markov Models (HMM)
- d. All of the above
- e. None of the above

 C

5. [2%] In the minimax algorithm, which of the following is the most unrealistic in practice?

- a. The knowledge of the utility values for the terminal states
- b. The generation of the whole game tree
- c. The assumption that the players are rational
- d. All of the above
- e. None of the above

## **1. [10%] True/False Questions**

[1% each no partial credit]

a) F

b) F

c) T (Each node in the decision tree divides the examples into two nonempty sets. Since there are n training examples, the number of nodes from the leaf to the root (i.e., the height) is bounded by n.)

d) T (In the worst case, we can have one Boolean output node with n Boolean parents and an n-dimensional,  $2^n$ -big CPT.)

e) F (The decision rule for Naive Bayes cannot be written as:  $\sum_i w_i x_i > k$ .)

f) F

g) T

h) F

i) F

j) T

## **2. [20%] Bayes Net**

2A [6%]

1. T [2%]
2. T [2%]
3. F [2%]

2B. [7%]

$$P(B+, H+, L-, Q+, E-)$$

$$= P(B+) P(H+) P(L- | H+) P(Q+ | B+, H+, L-) P(E- | Q+) \quad [5\%, \text{ Deduct } 1\% \text{ for each missing term}]$$

$$= 0.5 * 0.1 * (1 - 0.3) * 0.4 * (1 - 0.6) \quad [2\%, -0.5\% \text{ for each wrongly replaced value}]$$

$$= 0.0056.$$

## 2C. [7%]

$$P(B | H+, L-, E+)$$

$$= \alpha P(B, H+, L-, E+) [2\%]$$

$$= \alpha(P(B, H+, L-, Q+, E+) + P(B, H+, L-, Q-, E+)) [1\%]$$

$$= \alpha(<P(B+) P(H+) P(L- | H+) P(Q+ | B+, H+, L-) P(E+ | Q+),$$

$$P(B-) P(H+) P(L- | H+) P(Q+ | B-, H+, L-) P(E+ | Q+)>$$

$$+ <P(B+) P(H+) P(L- | H+) P(Q- | B+, H+, L-) P(E+ | Q-),$$

$$P(B-) P(H+) P(L- | H+) P(Q- | B-, H+, L-) P(E+ | Q-)>$$

$$) [1\%]$$

$$= \alpha P(H+) P(L- | H+) * (<P(B+) P(Q+ | B+, H+, L-) P(E+ | Q+),$$

$$P(B-) P(Q+ | B-, H+, L-) P(E+ | Q+)>$$

$$+ <P(B+) P(Q- | B+, H+, L-) P(E+ | Q-),$$

$$P(B-) P(Q- | B-, H+, L-) P(E+ | Q-)>$$

$$= \beta * (<0.5 * 0.4 * 0.6, 0.5 * 0.3 * 0.6> + <0.5 * 0.6 * 0.1, 0.5 * 0.7 * 0.1>)$$

[where  $\beta := \alpha P(H+) P(L- | H+)$ ] [1%]

$$= \beta * 0.5 * (<0.4 * 0.6, 0.3 * 0.6> + <0.6 * 0.1, 0.7 * 0.1>)$$

$$= \gamma * (<0.24, 0.18> + <0.06, 0.07>)$$

[where  $\gamma := \beta * 0.5$ ]

[Note: =  $\alpha(<0.084, 0.063> + <0.021, 0.0245>)$  ]

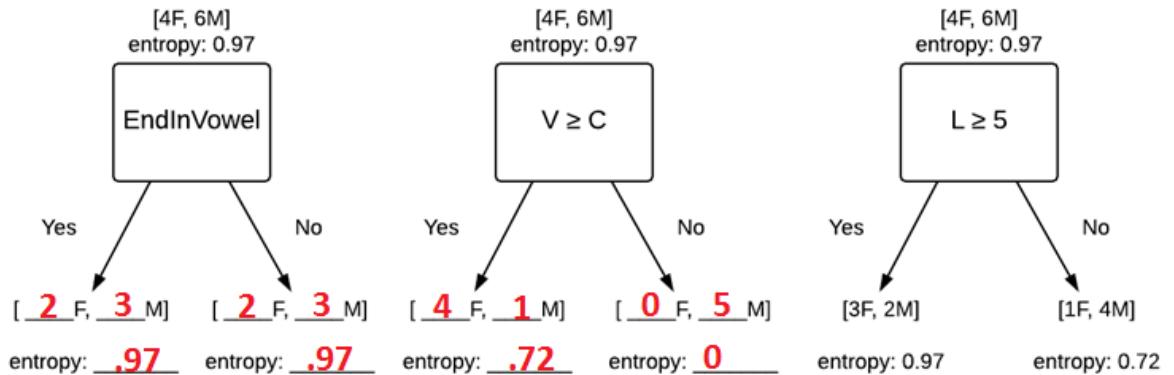
$$= \gamma * <0.3, 0.25> [1\%, any value pair with the correct ratio (6:5) is ok]$$

$$\approx <0.545, 0.455>. [1\%, normalization]$$

[If one of  $P(B+ | H+, L-, E+)$  and  $P(B- | H+, L-, E+)$  is computed correctly while the other is not, 5% partial credit is given]

### 3. [23%] Decision Tree Learning

3A. [8%] 2% per correct entropy with number of F/M



3B. [6%] 2% per feature

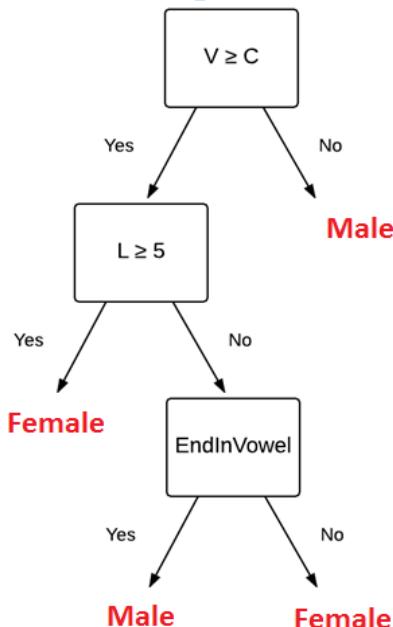
$$\text{IG}(\text{EndInVowel}) = 0.97 - (5/10 * .97 + 5/10 * .97) \\ = 0$$

$$\text{IG}(V \geq C) = 0.97 - (5/10 * .72 + 5/10 * 0) \\ = 0.61$$

$$\text{IG}(L \geq 5) = 0.97 - (5/10 * 0.97 + 5/10 * 0.72) \\ = 0.125$$

3C. [2%] ( $V \geq C$ ) since it has highest information gain.

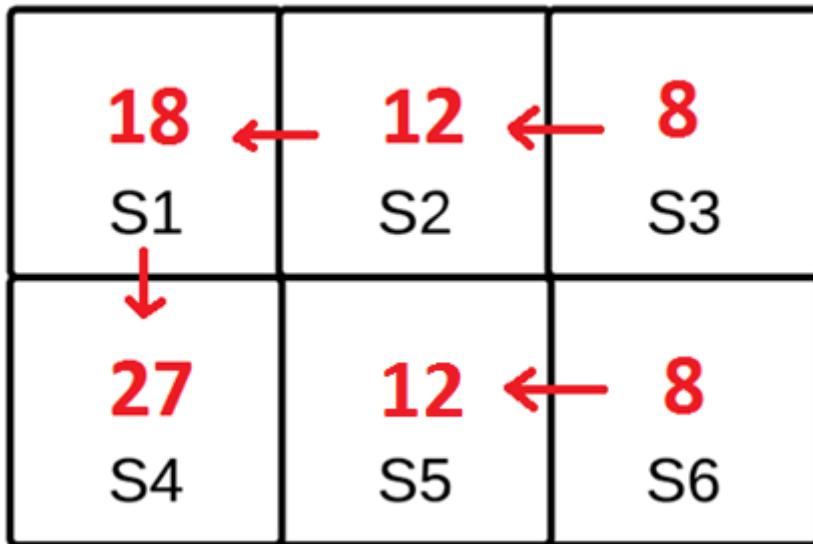
3D. [7%] 1% per correct node



[-0.5%] per leaf that shows  $[_M, _F]/\text{entropy}=0$  instead of taking a decision

#### 4. [17%] Markov Decision Process

4A. [7%] 1% per correct  $V^*$  value (total 6%)  
1% for all correct arrows in S1/S2/S3/S6.



4B. [5%] (1% for changed policy – 4% for  $\gamma$  value)  
 $\gamma^2 R_4 < \gamma R_5$  leads to:  $S2 \Rightarrow S5$  instead of  $S2 \Rightarrow S1$   
 $\gamma < R_5 / R_4 = 4/9$

4C. [5%]

$$V_6^* = p \gamma R_5 + (1-p) \gamma V_6^*$$

therefore  $V_6^* = p \gamma R / (1 - (1-p) \gamma)$  [3% for either form]  
 $= \frac{1}{4} * \frac{2}{3} * 12 / (1 - \frac{3}{4} * \frac{2}{3})$

$$V_6^* = 4 \quad [2%]$$

#### 5. [20%] Neural Networks

5A. [4%]

$$((5 + 1) * 3) + ((3 + 1) * 2) = 18 + 8 = 26$$

[2% partial credit if the students compute the right number not including the biases:  $5 * 3 + 3 * 2 = 21$ ]

5B. [4%]

1. [2%] F.
2. [2%] F. (Overfitting)

### 5C. [12%]

[6%] Note that  $y > 0.5$  if  $x > 0$ , and  $y \leq 0.5$  if  $x \leq 0$ . Given this, we need to choose  $w_i$  so that  $w_0 + w_1 * x_1 + w_2 * x_2$  will be greater than 0 when  $x_1 \vee x_2$  is equal to 1.

$X_1$	$X_2$	Y	Constraints on $w_i$	Credit
0	0	0	$w_0 \leq 0$	[1.5%]
0	1	1	$w_0 + w_2 > 0$	[1.5%]
1	0	1	$w_0 + w_1 > 0$	[1.5%]
1	1	1	$w_0 + w_1 + w_2 > 0$	[1.5%]

[6%] Similar to previous part, we need to choose  $w_i$  so that  $w_0 + w_1 * x_1 + w_2 * x_2$  will be greater than 0 when  $x_1 \wedge x_2$  is equal to 1.

$X_1$	$X_2$	Y	Constraints on $w_i$	Credit
0	0	0	$w_0 \leq 0$	[1.5%]
0	1	0	$w_0 + w_2 \leq 0$	[1.5%]
1	0	0	$w_0 + w_1 \leq 0$	[1.5%]
1	1	1	$w_0 + w_1 + w_2 > 0$	[1.5%]

## 6. [10%] Multiple-Choice Questions

[2% per answer]

1.D    2.A    3.A    4.B    5.B