

# 15-381 Spring 2007

## Assignment 5: Probability and Uncertainty, Bayes Nets and MDPs

Solutions

Spring 2007

### Problem 1 - Bayes Rule (15 points)

You have just moved to a new town. You were told that there are two types of busses in town: some busses run every 10 minutes, and the others run every 30 minutes (depending on the bus line). Waiting for a bus for the first time in your nearest bus station, you don't yet know which type of bus it serves.

You've been waiting for 20 minutes at the bus station, and the bus hasn't arrived yet.

Assuming that a waiting time interval is modeled by the exponential distribution function:

1. Use Bayes rule to derive your posterior belief about which type of bus this bus stop serves. What is the probability that you are waiting for a bus of the first type (give a number)? of the second type? Show your work, and results.  
Do you find these results intuitive?
2. How would your results be, given that you've been waiting for 5 minutes at the bus stop, and the bus hasn't arrived yet? Give the numerical results.
3. After how many minutes of waiting, would you believe that the bus you're waiting for is more probable to be of the second type?

**About the exponential distribution:** An exponential distribution is often used to model the time between independent events that happen at a constant average rate.

The probability density function of an exponential distribution has the form:

$$f(x = X | \beta) = \frac{1}{\beta} e^{-\frac{x}{\beta}}$$

The cumulative distribution function is given by:

$$F(x \leq X | \beta) = 1 - e^{-\frac{x}{\beta}}$$

where  $\beta$  is the mean interval length (waiting period, here),  $x \geq 0$  and  $\beta > 0$ .

**Hint:** Since the bus hasn't yet arrived, then what you should be estimating here is  $f(\beta | x > X)$ .

**Solution 1**

1. Since we have no prior information about the type of bus this stop serves, the prior is uniform:  
 $Pr(y = 10) = Pr(y = 30) = 0.5$ .

The probability that the waiting time is longer than 20 minutes, for each bus type is:

$$Pr(x = 20 | \beta = 30) = 1 - (1 - e^{-\frac{x}{30}}) = e^{-\frac{x}{30}}$$

Similarly,

$$Pr(x = 20 | \beta = 10) = e^{-\frac{x}{10}}$$

Applying Bayes rule, we get:

$$Pr(\beta = 30 | x = 20) = \alpha Pr(x = 20 | \beta = 30) Pr(y = 30) = 0.5\alpha e^{-\frac{2}{3}}$$

$$\alpha = 0.5e^{-\frac{2}{3}} + 0.5e^{-2} = 0.324$$

We therefore get:

$$Pr(\beta = 30 | x = 20) = 0.791$$

$$Pr(\beta = 10 | x = 20) = 1 - Pr(\beta = 30 | x = 20) = 0.209$$

That is, after 20 minutes of waiting, the posterior probability that the bus stop serves the bus that arrives every 30 minutes on average is much larger than the bus that arrives every 10 minutes on average.

- 2.

$$Pr(\beta = 30 | x = 5) = 0.583, Pr(\beta = 10 | x = 5) = 0.417$$

3. To find the specific waiting time, after which our posterior belief is that the bus line we're waiting for is of the type that arrives every 30 minutes on average, we require:

$$e^{-\frac{x}{30}} > e^{-\frac{x}{10}}$$

The match point is for  $x = 0$ . This means that, as long as we've waited *some* time, the posterior probability that the relevant bus is the one with the longer average wait time is larger. The gap between  $Pr(\beta = 10 | x = X)$  and  $Pr(\beta = 30 | x = X)$  is increased as the waiting time  $x$  increases.

## Problem 2 (15 points)

Suppose we have two sensors -  $x_1, x_2$  - taking samples of the random variable  $y$ . For example, consider two optical sensors, that are intended at measuring an object's location,  $y$ . Each sensor reports a noisy estimate of the true sampled signal. The variability of each sensor is known (it is determined by the specific environment conditions for each sensor). Your job is to determine the best estimate of the object's location.

We will assume the following Gaussian probability distributions:  $f_{(y,a^2)}(x_1), f_{(y,b^2)}(x_2)$ .

The notation  $f_{(m,\sigma^2)}(x)$  denotes a Gaussian density function of variable  $x$ , for which the mean is  $m$  and the variance is  $\sigma^2$ :

$$f_{(m,\sigma^2)}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - m)^2\right)$$

That is, the mean of the measurements given by each sensor is the true sampled location value  $y$  in both cases. But, the variance of  $x_1$  is  $a^2$  and the variance of  $x_2$  is  $b^2$ .

The location  $y$  is also a Gaussian, distributed  $f_{(0,1)}(y)$ .

1. Describe a Bayes net model for this situation.

2. Derive the posterior distribution  $f(y | x_1, x_2)$ .

**Note: following is a recipe for the product of gaussian densities, which you should find helpful:**

The product of two Gaussians is:

$$f_{(m_1, \sigma_1^2)}(x) \times f_{(m_2, \sigma_2^2)}(x) = C_c f_{(m_c, \sigma_c^2)}(x)$$

where

$$\sigma_c^2 = \frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}},$$

$$m_c = \frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}} \left( \frac{m_1}{\sigma_1^2} + \frac{m_2}{\sigma_2^2} \right),$$

and

$$C_c = f_{m_1, \sigma_1^2 + \sigma_2^2}(m_2).$$

**A hint and advice:**

- Due to the symmetry of the Gaussian density function,  $f_{(y, \sigma^2)}(x) = f_{(x, \sigma^2)}(y)$
  - Calculation can be quite neat if you keep  $C_c$  as a “black box” constant (if you have things right, the constants would cancel out :- ) ).
3. What is the maximum likelihood estimate (MLE) for the object’s location given both sensor’s measurements (that is, what is  $E(P(y | x_1, x_2))$  ? (Hint: rather than derive the expression, you may reason about this considering the shape of the posterior distribution function).
4. Interpret the MLE expression, and describe how it weights the information.

## Solution 2

1. The Bayes net structure is:

2.

$$f(y | x_1, x_2) = \frac{f(x_1, x_2 | y)f(y)}{\int f(x_1, x_2 | y)dy} = \frac{f(x_1 | y)f(x_2 | y)f(y)}{\int f(x_1 | y)f(x_2 | y)f(y)dy}$$

$$f_{(0,1)}(y)f_{(y,a^2)}(x_1)f_{(y,b^2)}(x_2) = f_{(0,1)}(y)f_{(x_1,a^2)}(y)f_{(x_1,b^2)}(y) = f_{(0,1)}(y)C_c f_{(x_c, \sigma_c^2)}(y) = C_c C_d f_{(x_d, \sigma_d^2)}(y)$$

where,

$$\frac{1}{\sigma_d^2} = \frac{1}{\sigma_c^2} + 1$$

$$x_d = \sigma_d^2 \frac{x_c}{\sigma_c^2}$$

$$\frac{1}{\sigma_c^2} = \frac{1}{a^2} + \frac{1}{b^2}$$

$$x_c = \sigma_c^2 \left( \frac{x_1}{a^2} + \frac{x_2}{b^2} \right)$$

That is, the posterior distribution is a Gaussian, where:

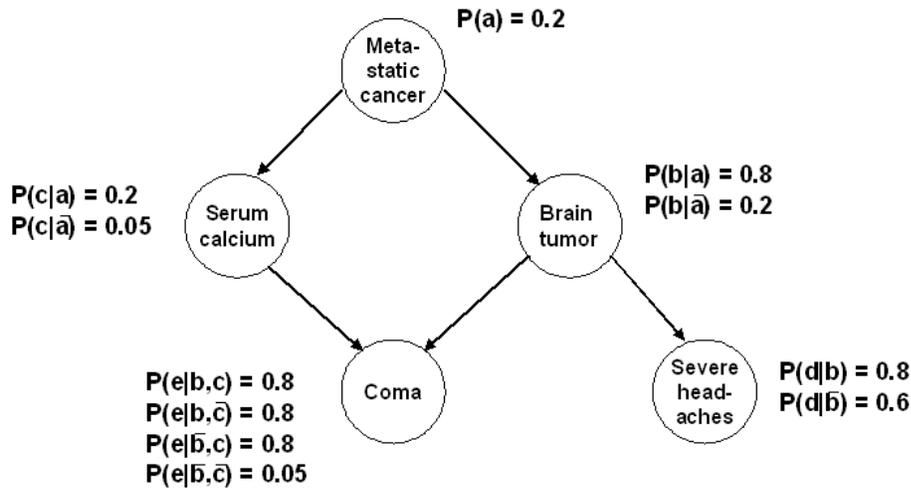
$$\frac{1}{\sigma_d^2} = \frac{1}{a^2} + \frac{1}{b^2} + 1$$

$$x_d = \left( \frac{x_1}{a^2} + \frac{x_2}{b^2} \right) \sigma_d^2 = \left( \frac{x_1}{a^2} + \frac{x_2}{b^2} \right) \frac{1}{\frac{1}{a^2} + \frac{1}{b^2} + 1}$$

3. For a Gaussian, which is symmetric in the space, the MLE is exactly the mean,  $x_d$ .
4. The mean of the posterior distribution weights measurements  $x_1$  and  $x_2$  in inverse proportion to their variance ( $a^2$  and  $b^2$ , respectively).

## Problem 3 - Bayes Nets (35 points)

Given is a simplified version of a network that could be used to diagnose patients arriving at a clinic. Each node in the network corresponds to some condition of the patient. This network demonstrates some causality links. For example, both brain tumor and serum calcium increase the chances of a coma. A brain tumor can cause severe headaches and a comma, and so on.



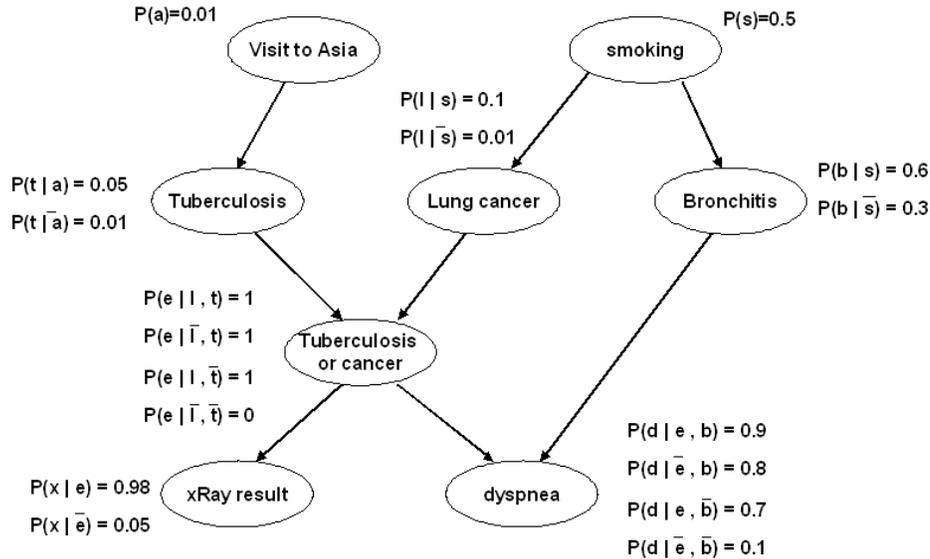
1. (2 points) What is the joint distribution  $P(a, b, c, d, e)$ ? give a factorized expression, according to the network's structure.
2. (3 points) Give an example of 'explaining away' in this Bayes net.
3. (5 points) One of your patients experiences severe headaches, had a comma and serum calcium. What is the probability of him having cancer? show full derivation of this probability, as well as the numerical result.
4. (5 points) What is the probability of a positive serum calcium given severe headaches? Derive this expression. Specifically, start from the joint distribution, factorize it and use variable elimination, so as to lower calculation cost. (Note: a numerical result is *not* required here.)
5. (15 points) Write code for sampling joint and conditional probabilities in Bayes nets. Use the standard name "sample" for your code. The command line arguments should be: FILE-NAME (VAR-NAME=value,VAR-NAME=..) (VAR-NAME=value,VAR-NAME=..). (In case you are to use a non-standard language, be sure to included a README file with runtime instructions.)

The first argument specifies the file describing the Bayes net. The second group of arguments gives the required variables' values, whose probability we'd like to evaluate. Several such values can be specified, using a comma separator within the same round brackets. The last group of arguments gives the variables' values that the requested probability is conditioned on (again, multiple values can be specified, using a comma separator and bounding brackets). If the second group is empty (unfilled brackets), this means we are looking for a joint probability expression, that is conditioned on nothing.

For example: “cancer.text (A=true) (D=true,E=true,C=true)” would give you the probability for question 4.3 .

A text file that includes this net’s information is handed out with this homework, on the class’ website. Using the program, evaluate the probabilities specified in questions 4.3 and 4.4, generating  $n$  samples (where a single sample assigns values to all the nodes in the network). What are the estimated probabilities for  $n = 500$ ?  $n = 1000$ ?  $n = 20,000$ ?

6. (5 points) Consider the following Bayes net, called “Asia”, which can be used to diagnose respiratory diseases (this is a fictitious network).



Use your code to evaluate the probabilities that:

- a patient who smokes, has visited Asia, and got positive xRay results, has lung cancer (give a numerical result).
- a patient who smokes and has dyspnea, has got lung cancer (give a numerical result).

### Solution 3

1.  $p(a, b, c, d, e) = p(e|a, b, c, d)p(a, b, c, d) = p(e|a, b, c, d)p(d|a, b, c)p(c|b, a)p(b|a) = p(e|a, c)p(d|b)p(c|a)p(b|a)p(a)$
2.  $e$  is a common child to  $b$  and  $c$ . In this structure, two causes “compete” to “explain” the observed data. Hence  $b$  and  $c$  become conditionally dependent given that their common child,  $e$ , is observed, even though they are marginally independent. If we know that  $e$  (comma) is true and also  $b$  (brain tumor) is true, this reduces the probability that  $c$  (serum calcium) is true.

3.

$$P(a|d, e, c) = \frac{P(a, c, d, e)}{P(c, d, e)} = \frac{\sum_B P(a, B, c, d, e)}{\sum_{a, B} P(a, B, c, d, e)} \tag{1}$$

$$\propto \sum_B P(a, B, c, d, e) \tag{2}$$

$$\propto \sum_B P(a)P(B|a)P(c|a)P(d|B)P(e|B, c) \tag{3}$$

$$\propto P(a)P(c|a) \sum_B P(B|a)P(d|B)P(e|B, c) \tag{4}$$

$$\propto P(a)P(c|a) (P(b|a)P(d|b)P(e|b, c) + P(\bar{b}|a)P(d|\bar{b})P(e|\bar{b}, c)) \tag{5}$$

$$\propto 0.2 \times 0.2 (0.8 \times 0.8 \times 0.8 + 0.2 \times 0.6 \times 0.8) \quad (6)$$

$$\propto 0.02432 \quad (7)$$

$$P(\bar{a}|d, e, c) = \frac{P(\bar{a}, c, d, e)}{P(c, d, e)} = \frac{\sum_B P(\bar{a}, B, c, d, e)}{\sum_{\bar{a}, B} P(\bar{a}, B, c, d, e)} \quad (8)$$

$$\propto \sum_B P(\bar{a}, B, c, d, e) \quad (9)$$

$$\propto \sum_B P(\bar{a})P(B|\bar{a})P(c|\bar{a})P(d|B)P(e|B, c) \quad (10)$$

$$\propto P(\bar{a})P(c|\bar{a}) \sum_B P(e|B, c)P(B|\bar{a})P(d|B) \quad (11)$$

$$\propto P(\bar{a})P(c|\bar{a}) (P(e|\bar{b}, c)P(b|\bar{a})P(d|b) + P(e|b, c)P(\bar{b}|\bar{a})P(d|\bar{b})) \quad (12)$$

$$\propto 0.8 \times 0.05 (0.8 \times 0.2 \times 0.8 + 0.8 \times 0.8 \times 0.6) \quad (13)$$

$$\propto 0.02048 \quad (14)$$

$$P(a|d, e, c) = \frac{\sum_B P(a, B, c, d, e)}{\sum_B P(a, B, c, d, e) + \sum_B P(\bar{a}, B, c, d, e)} = \frac{0.02432}{0.02432 + 0.02048} = 0.54 \quad (15)$$

4.

$$P(c|d) = \frac{P(c, d)}{P(d)} = \frac{\sum_{A, B, E} P(A, B, c, d, E)}{\sum_{A, B, c, d, E} P(A, B, c, d, E)} \quad (16)$$

$$\propto \sum_{A, B, E} P(A, B, c, d, E) \quad (17)$$

$$\propto \sum_{A, B, E} P(A)P(B|A)P(c|A)P(d|B)P(E|B, c) \quad (18)$$

$$\propto \sum_A P(A)P(c|A) \sum_B P(B|A)P(d|B) \sum_E P(E|B, c) \quad (19)$$

$$\propto \sum_A P(A)P(c|A) \sum_B P(B|A)P(d|B) \quad (20)$$

Likewise,

$$P(\bar{c}|d) \propto \sum_A P(A)P(\bar{c}|A) \sum_B P(B|A)P(d|B) \quad (21)$$

Thus,

$$P(c|d) = \frac{\sum_A P(A)P(c|A) \sum_B P(B|A)P(d|B)}{\sum_A P(A)P(c|A) \sum_B P(B|A)P(d|B) + \sum_A P(A)P(\bar{c}|A) \sum_B P(B|A)P(d|B)} \quad (22)$$

5.  $P(c|d) \sim 0.084$

6.  $P(\text{CANCER} = T | \text{SMOKING} = T, \text{ASIA} = T, \text{XRAY} = T) \sim 0.54$

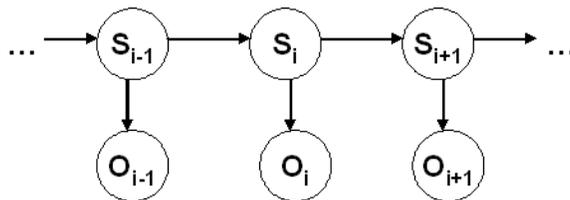
$P(\text{CANCER} = T | \text{DYSPNEA} = T, \text{SMOKING} = T) \sim 0.15$

## Problem 4 - HMM Applications (10 points)

Hidden Markov Models are used for a variety of sequential data processing. As shown in the figure, the model includes a *hidden layer*, which described the data as a sequence of pre-defined states. The output layer maps to the *observed* signals, that are emitted from the unknown states.

For example, in OCR systems, the hidden nodes represent the sequence of the underlying true characters. The observed decoded characters are not necessarily identical to the hand-written ones, due to the high variability of hand writing. For example, the hand-written “i” character is decoded as either “i”, “j”, or “l” with high probabilities. The advantage of using HMM for this problem is that it allows modeling context. For example, if the previous (hidden) state was assumed to be “q”, then the current character is more likely to be a vowel.

Given the observed OCR output sequence, the final output to the user is the sequence of states which maximizes the joint probability of the HMM model.



1. (3 points) In the described OCR scenario, what is the size of the CPT table between two states of the hidden layer? What is the size of the emission probabilities table for every state?
2. (3 points) Suggest a way for obtaining the transition and omission probabilities for the model.
3. (4 points) Suggest an HMM representation for the problem of speech recognition, or automatic Part-of-Speech labelling. (POS labelling is task of assigning every word its grammatical label, e.g., noun, verb, preposition etc.) You may suggest an HMM modeling for other problems as well.

Your suggestion should include the definition of states, the definition of the omitted signals, and a description of the values the hidden and observed nodes can take.

Explain why you think that the addressed problem would benefit using an HMM model.

## Solution 4

1. There are as many states as the size of the alphabet. Thus, the transition table between two states is of size  $\text{Size-of-alphabet} \times \text{Size-of-alphabet}$ .

Given any state, there may be up to  $\text{Size-of-alphabet}$  non-zero entries in the emission table, representing the chance that the OCR recognizes the actual letter as any letter in the alphabet. That is, for every given state, the size of the emission probabilities table is  $\text{Size-of-alphabet} \times 1$ .

2. In many cases, the relevant probabilities are obtained from manually labelled datasets. In the OCR scenario, this means that human annotators first tag a corpus of handwritten text with the correct machine-symbols (the relevant Ascii code, for example). Then, the emission probabilities can be derived by aligning a large annotated corpora with the corresponding OCR output.

As for the transition probabilities between states – these can be obtained more cheaply, by processing machine-readable documents and counting the regularities of the transitions between consecutive letters.

3. In automatic Part-of-Speech (POS) labeling, the states are a pre-defined set of POS tags (e.g., noun, adjective etc.). The emitted signals are the actual words. This representation assumes that there are contextual regularities in the sequences of POS tags (for example, nouns may follow a determiner, but verbs may not.). While most words have only one possible POS tag that can be extracted from a dictionary (e.g., “from” is a Preposition in all cases), there are words with several possible labels (e.g., “walk” can be both a noun and a verb). Using an HMM allows to resolve such ambiguous cases.

HMMs are very dominant in the area of speech processing. There, a separate HMM model is built for every word. A vector of acoustic features is computed every 10 to 30 msec. The states of the HMM are therefore the phonetics of the word, and the observed signals are the acoustic signal representations. As for a human, mapping the measured signals to the correct phonetics is more precise if a sequence is considered rather than every signal in isolation.

## Problem 5 - Markov Decision Processes(25 points)

Jane is a student. Every Sunday, she needs to make a decision - either go out with her friends, in which case she will go to sleep late and start the next week tired. Or, she may choose not go out, that is, go to sleep early. Suppose that Jane understood the material in class the week before. Then, if she chooses to sleep early, she is most likely (80%) to understand next week's lectures. If she starts the week tired, though, the chance that she will understand the material in class next week is only 50%. On the other hand, if Jane didn't understand the material in class last week, then the chance that she understands this week's lectures if she goes out is only 25%. Or, if she chooses to go to sleep early and start the week fresh, her chance to understand the classes next week is 50%.

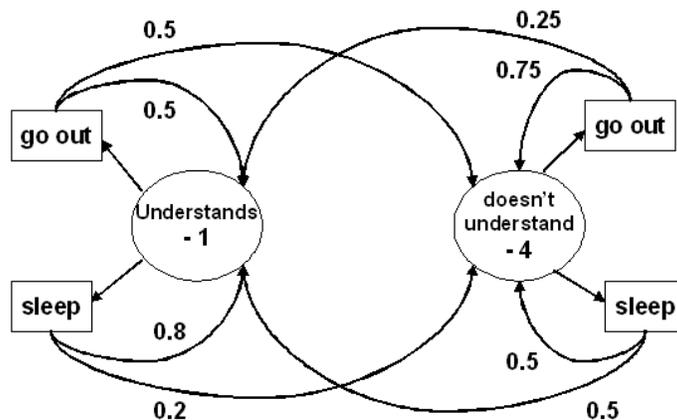
Jane is interested in understanding the material in class, of course. Especially, since in case she doesn't understand the material on a particular week, then she has to put a lot of time into homeworks, which costs her 4 satisfaction points for that week. If she understands the material, it only costs her 1 satisfaction point.

- (4 points) Describe this situation as a Markov Decision Process diagram, where nodes denote the states, and edges are labeled with the transition probabilities. Be sure to include actions in the diagram's edges.
- (4 points) Enumerate all the possible policies that Jane can take. For each policy, give its transition matrix.

In addition to the above, if Jane decides not to go out with her friends on Sunday, then her friends are disappointed. Even worse, they call her a wimp, which costs her 2 satisfaction points. Consider this information in the following questions:

- (5 points) On the first week of school, Jane is equally likely to understand classes, or not. What is her cost value then for every policy, assuming that the term is short, including three weeks? Assume that future satisfaction is discounted by  $\alpha = 0.9$ , per week. Show your work.
- (2 points) What is the optimal policy that Jane should take then ?
- (10 points) What is Jane's optimal policy for the whole school year (40 weeks). Explain how you calculate this. You may use code or a spreadsheet, for computation purposes. Your answer for this question should include a table of the relevant values per week, until convergence. Also, specify what the optimal policy is.

### Solution 5



- Note that it was not incorrect to use 4 states, as if you wanted to set a reward for each state when including the -2 penalty you would need them. It made the problem more difficult though, and was not necessary, as it is valid to accord a reward not just to states but to actions taken from those states.

2. There are four policies (one for each action from each state), that specify the following transition matrices. For convenience, we also included in the tables the modified costs associated with each state given a policy, and the additional information that the action of not going out has a related cost of -2 points. The costs information is used in the next questions.

policy	state	cost	u	d
a: go out either way	u	-1	0.5	0.5
	d	-4	0.25	0.75
b: go to sleep either way	u	-3	0.8	0.2
	d	-6	0.5	0.5
c: go out if understands	u	-1	0.5	0.5
	d	-6	0.5	0.5
d: go out if doesn't understand	u	-3	0.8	0.2
	d	-4	0.25	0.75

3. There were several ways to do this problem, and several ways to interpret the -2 wimp penalties, and we gave credit for any correctly executed combination of approach and assumptions as to wimp penalties.

There were three valid interpretations of the interactions of when wimp penalties were assessed and the discount factor. Assumption 1 was to assume that the wimp penalties were assessed in the first round and subsequently discounted. Assumption 2 was to assume that wimp penalties were not assessed until the second week but were not discounted. Assumption 3 was to assume that wimp penalties were not assessed until the second week and were discounted.

The approach we took was to determine a separate reward for starting from each different state and then to weight them by the initial probabilities of .5 each. Note that you needed to get a single number for the cost of the policy - we took off a bit if you didn't average. Also note that if the likelihood of starting in a particular policy was not 50% you would have needed to weight by likelihoods. It is also ok to start in each state with 50% probability and then evaluate from there. Finally, many students used the  $J^k(s)$  formulation, which is equivalent. Note that in the approach that follows we need to sum across the expected reward each time step:

policy a:

	t	p(u)	p(d)	total cost
Starting from u:	$t_0$	1	0	-1
	$t_1$	0.5	0.5	$0.9(-1 \times 0.5 - 4 \times 0.5)$
	$t_2$	$0.25 + 0.125$	$0.25 + 0.375$	$0.9^2(-1 \times 0.375 - 4 \times 0.625)$
and, from d:	$t_0$	0	1	-4
	$t_1$	0.25	0.75	$0.9(-1 \times 0.25 - 4 \times 0.75)$
	$t_2$	$0.125 + 0.1875$	$0.125 + 0.5625$	$0.9^2(-1 \times 0.3125 - 4 \times 0.685)$

Results for Assumption 1:

$$\text{Total cost for policy } a: \text{Cost}(a, t = 3) = 0.5 \times -5.5788 + 0.5 \times -9.4056 = -7.492$$

Similarly,

$$\text{Cost}(b, t = 3) = -11.316$$

$$\text{Cost}(c, t = 3) = -9.485$$

$$\text{Cost}(d, t = 3) = -9.431$$

Results for Assumption 2:

$$\text{Cost}(a, t = 3) = -7.492$$

$$\text{Cost}(b, t = 3) = -9.664$$

$$\text{Cost}(c, t = 3) = -8.675$$

$$\text{Cost}(d, t = 3) = -8.558$$

Results for Assumption 3:

$$\text{Cost}(a, t = 3) = -7.492$$

$$\text{Cost}(b, t = 3) = -9.316$$

$$\text{Cost}(c, t = 3) = -8.485$$

$$\text{Cost}(d, t = 3) = -8.364$$

4. The optimal policy is then, policy  $a$  - to go out either way (for all assumptions)

(While this is the correct answer, the costs and probabilities may not correspond to CMU student life :-))

5. The results for 40 weeks for Assumption 1 are:

$$\text{Cost}(a, t = 40) = -28.96$$

$$\text{Costs}(b, t = 40) = -38.94$$

$$\text{Costs}(c, t = 40) = -34.53$$

$$\text{Costs}(d, t = 40) = -34.10$$

Other assumptions yielded slightly smaller numbers for policies  $b, c$ , and  $d$ .

Thus policy  $a$  is the optimal one. The relative rankings of the policies is the same as for  $t = 3$  (or earlier). This is the justification for the policy modification method. It was also fine to use policy or value iteration for this problem, as you were only asked to determine the best policy.