

CS 540 (Shavlik) HW 3 – Probabilistic Reasoning

Assigned: 10/14/15

Due: 11/5/15 at 11:55pm (not accepted after 11:55pm on 11/12/15)

Points: 150

Suggestion: you should consider doing this HW in a word processor since cut-and-pasting is likely to be useful.

You should do Problem 1 before the midterm since it involves material from Lecture 14, the last lecture whose topics might be on the midterm.

Your solutions to Problems 1-5 should be placed in HW3_p1_to_p5.pdf.

Problem 1: Full Joint Probability Distributions (15 points)

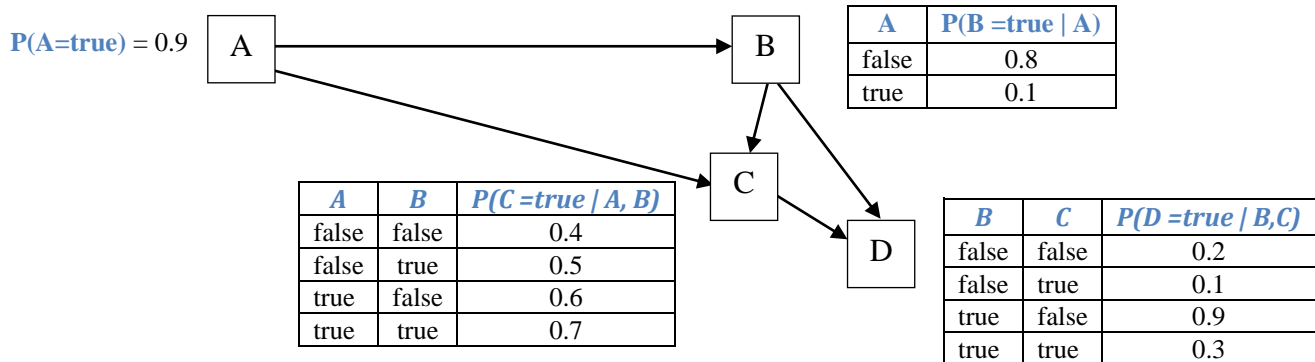
Consider this *full joint probability distribution* involving four Boolean-valued random variables (A-D). For this problem you can simply walk through the cells in the table below one-by-one and add up those probabilities where the expression for a given question below is true. I.e., you do not need to create “complete world states.”

<u>A</u>	<u>B</u>	<u>C</u>	<u>D</u>	<u>Prob</u>
F	F	F	F	0.012
F	F	F	T	0.013
F	F	T	F	0.010
F	F	T	T	0.025
F	T	F	F	0.015
F	T	F	T	0.025
F	T	T	F	0.025
F	T	T	T	0.025
T	F	F	F	?
T	F	F	T	0.020
T	F	T	F	0.030
T	F	T	T	0.040
T	T	F	F	0.050
T	T	F	T	0.060
T	T	T	F	0.070
T	T	T	T	0.080

- Compute $P(A = \text{true and } B = \text{false and } C = \text{false and } D = \text{false})$.
- Compute $P(A = \text{true and } C = \text{false and } D = \text{false})$.
- Compute $P(B = \text{true or } D = \text{false})$.
- Compute $P(D = \text{false} \mid A = \text{false and } B = \text{true and } C = \text{true})$.
- Compute $P(A = \text{false and } B = \text{true and } C = \text{true} \mid D = \text{false})$.

Problem 2: Bayesian Networks (20 points)

Consider the following Bayesian Network, where variables **A-D** are all Boolean-valued:



Show your work for the following calculations.

- i. Compute $P(A = \text{true} \text{ and } B = \text{true} \text{ and } C = \text{true} \text{ and } D = \text{true})$.
- ii. Compute $P(A = \text{false} \text{ and } C = \text{false} \text{ and } D = \text{false})$.
- iii. Compute $P(C = \text{true} \mid A=\text{true} \text{ and } B=\text{true} \text{ and } D=\text{true})$.
- iv. Compute $P(D=\text{false} \mid B = \text{true} \text{ and } C = \text{true})$. // Use the *Markov Blanket* property here.
- v. Compute $P(D=\text{false} \mid A=\text{true} \text{ and } B = \text{true} \text{ and } C = \text{true})$.
 // Do this one without employing the *Markov Blanket* property (i.e., to algebraically or
 // numerically confirm that knowing $A=\text{true}$ doesn't change the result from iv).
- vi. Compute $P((A = \text{false} \text{ and } C = \text{false}) \text{ or } (B = \text{true} \text{ and } D = \text{true}))$.

Problem 3: BNs and Full Joint Probability Tables (5 points)

Use the Bayes Net from Problem 2 to fill in the top quarter of the full joint probability table below (feel free to fill it all in, but to reduce the tedium we are only requiring you to fill in the top quarter). Show your work.

<u>A</u>	<u>B</u>	<u>C</u>	<u>D</u>	<u>Prob</u>
F	F	F	F	?
F	F	F	T	?
F	F	T	F	?
F	F	T	T	?
F	T	F	F	?
F	T	F	T	?
F	T	T	F	?
F	T	T	T	?
T	F	F	F	?
T	F	F	T	?
T	F	T	F	?
T	F	T	T	?
T	T	F	F	?
T	T	F	T	?
T	T	T	F	?
T	T	T	T	?

Problem 4: Bayes' Rule (20 points)

Define the following two variables about people:

shot = got flu shot last year

flu = caught flu this year

Assume we know from past experience that:

$$P(\textit{shot}) = 0.6$$

$$P(\textit{flu}) = 0.4$$

$$P(\textit{flu} \mid \textit{shot}) = 0.2$$

- i. Given someone did *not* get a *shot*, what is the probability he or she does gets the *flu*?
- ii. Given you find out someone got the *flu*, what's the probability he or she got a *shot*?

Be sure to show and explain your calculations for both parts (i) and (ii). Start by writing out the above questions as conditional probabilities.

In the general population, 10 in a 100,000 people have the dreaded *DislikeHomework* disease. Fortunately, there is a test (*test4it*) for this disease that is 97.5% accurate. That is, if one has the disease, 975 times out of 1000 *test4it* will turn out positive; if one does *not* have the disease, 25 times out of 1000 the test will turn out positive.

- iii. You take *test4it* and the results come back true. Use Bayesian reasoning to calculate the probability that you actually have *DislikeHomework*. That is, compute:

$$P(\textit{DislikeHomework} = \textit{true} \mid \textit{test4it} = \textit{true})$$

Show your work; you may use *DH* for *DislikeHomework* and *T4* for *test4it* if you wish.

Problem 5: Naïve Bayes and ‘Bag-of-Words’ Text Processing (15 points)

Sally has divided her books into two groups, those she likes and those she doesn't. For simplicity, assume no book contains a given word more than once.

The 8 books that Sally likes contain (only) the following words:

animal (7 times), *vegetable* (1 time), *see* (6 times), *eat* (2 times)

The 12 books that Sally dislikes contain (only) the following words:

animal (11 times), *mineral* (3 times), *vegetable* (8 times), *see* (9 times), *eat* (5 times)

Using Bayes Rule and the Naive Bayes assumption, determine whether it is more probable that Sally likes the following book (call it book23) than that she dislikes it.

see animal eat vegetable // These four words are the entire contents of book23.

That is, compute the ratio:

$$\frac{\text{Prob}(\text{Sally likes book23} \mid \text{'see' in book} \wedge \text{'animal' in book} \wedge \text{'eat' in book} \wedge \text{'vegetable' in book})}{\text{Prob}(\text{Sally dislikes book23} \mid \text{'see' in book} \wedge \text{'animal' in book} \wedge \text{'eat' in book} \wedge \text{'vegetable' in book})}$$

Be sure to show and explain your work. Ensure that none of your probabilities are zero by starting all your counters at 1 instead of 0 (the counts above result from starting at 0, i.e., imagine that there is one more book (book21) that Sally likes that contains each of the five unique words above exactly once and also one more book (book22) she dislikes that also contains each of the above five words exactly once).

Technical note: in the *evidence* (ie, the 'given') part of the above conditional probabilities we are ignoring the words not present in the book (i.e., *mineral*). This is different from our usual 'fixed-length feature vector' view of data, where instances of *randomVar=false* impact the calculation. Since the average vocabulary is around 50,000 words, when dealing with text as our data we want to focus on the words present rather than having the absence of tens of thousands of words impact the calculation.

You can visualize the random process as follows. There is one bag of all the books that Sally would like (this bag contains many more than the books mentioned above; these 'liked' books are just a sample drawn from this bag). There is a second bag, which contains all the books that Sally would dislike. If we draw a book from the 'likes' bag, what is the probability this book contains the word *see*? What would this probability be if we draw a book from the 'dislikes' bag? Since we are assuming conditional independence of the words given the type of book (liked vs. disliked), we repeat this process for *animal*, then repeat for *eat*, and finally repeat for *vegetable*. We also need to remember to include the answer to: if we combined these two bags of books and then (uniformly) randomly drew a book from the combination, what is the probability we drew a book Sally liked?

Problem 6: Creating Probabilistic Reasoners that Play Nannon (75 points)

This problem involves writing Java code that implements three probabilistic reasoners to play the two-person board game called Nannon (<http://nannon.com>), which is a simplified version of the well-known game Backgammon (<http://en.wikipedia.org/wiki/Backgammon>). Instructions for Nannon are available at <http://nannon.com/rules.html>.

Here is how we will formulate the task. At each turn, whenever there is more than one legal move, your *chooseMove* method is given access to:

- 1) The *current board configuration*.
- 2) A *list of legal moves*; each move's effect is also provided, and you are able to determine the next board configuration from each move. (Explanations of which effects are computed for you appear in the *chooseMove* method of the provided *RandomNannonPlayer* and in the *ManageMoveEffects.java* file).

Your *chooseMove* method needs to return one of the legal moves. It should do this by using Bayes' Rule to estimate the odds each move will lead to a winning game, returning the one with the highest odds. That is, it should compute for each possible move:

$$Prob(\text{will } \underline{\text{win}} \text{ game} \mid \text{current board, move, next board, and move's effect})$$

$$Prob(\text{will } \underline{\text{lose}} \text{ game} \mid \text{current board, move, next board, and move's effect})$$

Your solution need not use ALL of these given's to estimate these probabilities, and you can choose to define whichever random variables you wish from the provided information. The specific design is up to you and we expect each student's solution to be unique.

You need to create three solutions. In one, you will create a full joint probability table. In the other two you will create two (Bayesian Networks (BNs); neither can be a BN equivalent to your full joint probability table. One BN should be Naive Bayes (NB) and the other needs to somehow go beyond the NB conditional-independence assumption (see Lecture 17). It is up to you to decide the specific random variables used and, for the non-naive Bayesian Network, which conditional independence assumptions you wish to make. The random variables in your three solutions need not be the same.

You need to place your three solutions in these files, where *YourMoodleLoginName* is your actual Moodle (i.e., your UWisc) login:

FullJointProbTablePlayer_YourMoodleLoginName.java

NaiveBayesNetPlayer_YourMoodleLoginName.java

BayesNetPlayer_YourMoodleLoginName.java

Copy all the Java files in <http://pages.cs.wisc.edu/~shavlik/cs540/HWs/HW3/> to your working space. The provided *PlayNannon.java*, *NannonPlayer.java*, and *RandomNannonPlayer.java* files contain substantial details on what you need to do. You should start by reading the comments in them; I suggest you read the files in the order they appear in the previous sentence.

So how do you get the necessary information to compute these probabilities? After each game, your players' *updateStatistics* method is given information about the sequence of board

configurations encountered and the moves chosen by your player in that game, as well as whether or not your player won that game. You should not try to figure out which moves were good or bad in any one specific game; instead, if a game was won, all moves in it should be considered good (i.e., led to a win) and if a game is lost all moves should be considered bad (led to a loss). Obviously some moves in losing games were good and vice versa, but because we are using statistics, we are robust to this ‘noise.’

Your three players need to implement the *reportLearnedModel* method, which reports the value of the random variable (or values of the combination of random variables) where the following ratio is *largest* (i.e., most indicative of a win) and the *smallest* (i.e., most indicative of a loss):

$$\text{prob}(\text{randomVariable(s)} \mid \text{win}) / \text{prob}(\text{randomVariable(s)} \mid \text{loss})$$

For your full-joint-prob table, *randomVariables* should be a setting of all the random variables other than the ‘win’ variable (i.e., $\text{loss} = \neg \text{win}$). For your NB player, *randomVariable(s)* should be *one* variable other than *win*. For your Bayes Net approach, *randomVariable(s)* should be one of the non-NB entries in the product of probabilities you compute. (Recall that if we want to allow some dependencies among our random variables, the product in a non-naive Bayes Net calculation will include something like $p(A \mid B \wedge \text{win}) \times p(B \mid \text{win})$, which is equivalent to $p(A \wedge B \mid \text{win})$, as explained in class.)

The *reportLearnedModel* method is automatically called (by code we have written) after a run of *k* games completes when `Nannon.reportLearnedModels` is set to true.

It is fine to print more about what was learned, but the above requested information should be easy to find in your printout.

What needs to be turned in, in addition to the three Java files listed above (be sure to comment your code)? A report (in HW3_p6.pdf) containing:

1. A description of your design for your `FullJointProbTablePlayer`; should be 1-2 pages. Include an illustrative picture (the picture can be hand drawn). Don’t draw the complete full-joint table, but provide a visual of its contents.
2. A description of your design for your `NaiveBayesNetPlayer`; should be 1-2 pages. An illustrative picture of your Bayes Net is required (the picture can be hand drawn).
3. A description of your design for your `BayesNetPlayer`; should be 1-2 pages. An illustrative picture of your Bayes Net is required (the picture can be hand drawn). Focus your discussion on the differences from your NB player.
4. A 1-2 page report presenting and discussing how your three players performed against each other, as well as against the provided ‘hand coded’ and ‘random moves’ players.

Create a table where the columns are (change “Your” to “My” in your report):

YourFullJoint YourNB YourBN HandCoded Random

And the rows are

YourFullJoint YourNB YourBN

Feel the ‘upper triangle’ portion of this table with the results (the percentage of the time the column’s label beat the row’s label) on the 6 cells x 3 pieces per player board using 1M games after 1K ‘burn-in’ games.

Please type your report.

Here some statistics from my solution on the 6 cells x 3 pieces per player board. Note that results can vary depending on the ‘seed’ for the random-number generator. Also note that several students’ solutions beat my solutions in previous cs540 classes.

33% random-moves vs 67% my Full-Joint Player
47% hand-coded vs 53% my Full-Joint Player

33% random-moves vs 67% my Bayes Net
46% hand-coded vs 54% my Bayes Net

39% random-moves vs 61% my Naive Bayes
52% hand-coded vs 48% my Naive Bayes
(so don't worry if your Bayes Net cannot beat the hand-coded player)

Note that we will test your `NaiveBayesNetPlayer` and your `BayesNetPlayer` under various conditions. The `PlayNannon.java` file provides details. Also note that you should not modify any of the provided Java files since when we test your solution we will be using the original versions of these files.

A class-wide tournament will be conducted soon after this homework is due, using your non-naive BN player. It will be for fun and a learning experience (e.g., what did the winner(s) do that worked so well?); performance in this tournament will not be part of your grade on this homework. If your code is too slow, uses too much RAM, or crashes, it won’t be able to participate in the class-wide tournament. Hopefully the tournament-running code can handle solutions that use no more than 8GB RAM and play one million games in less than a minute (these ‘specs’ subject to change; my goal is to include all students’ solution, but the first time I did this in cs540, dealing with 3-4 ‘outlier’ solutions took up about half my time).

Feel free to create additional players based on genetic algorithms and/or a nearest-neighbor approach (or even decision trees/forests), but that is not required nor will any extra credit be given. I will be quite happy to talk to students about such approaches, though.

A GUI-Based Player

There is a GUI-based way to play against a trained computer opponent. If you haven’t already, download the `GUI_player.java`, `MiscForGUI.java`, `NannonGUI.java`, and `PlayingField.java`. (I crudely hacked up some old GUI code I had written around 1998 called the `AgentWorld`, so don't look at the code - it has a lot of unused junk in it.)

To invoke this player (as X, use arg2 to be O) do this in PlayNannon:

```
String arg1 = "gui";
```

It will first silently play the specified number of burn-in games against your chosen computer-based player; the GUI-based player will make random moves. It will then pop up a Java window and allow you to use the mouse to play against the computerized player.

Legal moves are visually highlighted. If you press down on the mouse but don't release, you will see where a moveable piece will go. If you release quickly enough (ie, "click on") a moveable player, the move will happen. If you press down on a moveable piece but do NOT want to make that move, move the mouse off that player and then release it (or hold it down long enough so the press+release is not viewed as click).

Note you only get to select a move for those board states where you have a CHOICE of moves, so games might seem a bit haphazard.

I have only tested this in Windows, but it is simple Java and should run anywhere Java does.

In a related note, at the end of a match you can manually watch two computerized players play. The file PlayNannon.java's main has a line like this:

```
Nannon.useGUIToWatch = true;
```

Change true to false to have the GUI turned off.

One known bug: you can't kill the GUI using the upper-right X in Windows, but there is a QUIT button (or kill the PlayNannon java process).

Some Additional Tips

Here are some miscellaneous notes and tips about Nannon that came up in student discussions when I used Nannon in previous semesters.

1) In Java when you spec arrays, you need not use constants but can instead do something like this:

```
myRandomVar_win
= new int[NannonGameBoard.piecesPerPlayer][NannonGameBoard.cellsOnBoard];
// I am simply illustrating variable-sized arrays here and don't use
// something exactly like this in my code.
// (The actual code uses 'accessor' methods for reading and writing these static vars.)
```

You should use NannonGameBoard.piecesPerPlayer and NannonGameBoard.cellsOnBoard in your BayesNet players so that it can handle any game size, but your FullJoint player only needs to handle NannonGameBoard.piecesPerPlayer = 3 and NannonGameBoard.cellsOnBoard = 6 (in Java the SIZE of an individual dimension can vary, but the NUMBER of dimensions cannot - at least as far as I know; possible some advanced feature of Java allows this - and this limitation makes it hard for the FullJoint player to handle various game sizes, plus for larger games, the FullJoint table might cause you to run out of memory since the table grows exponentially with game size).

2) You can vary `Nannon.numberOfGamesInBurnInPhase` to run experiments via an 'accessor' function. I made this and other variables private and provided accessor functions, so that in the class tourney I can alter the accessor functions to only allow the TAs or me to change this variable. In the class tourney, your player has to 'live with' the settings of this variable (it is ok to make random moves after the burn-in phase, but random moves are likely to increase the odds your player will lose).

Ditto for `Nannon.gamesToPlay` and `Nannon.printProgress` (be sure your turned-in player 'runs silently' unless a TA or I set `Nannon.printProgress = true` because excess printing will slow down things and clutter our output).

Ditto for `NannonGameBoard.piecesPerPlayer` and `NannonGameBoard.cellsOnBoard`. It would be rather chaotic if players could change the game in the middle of playing!

Good programming practice is to ALWAYS use accessor functions (ie, setters and getters) rather than making public variables, for reasons like these.

3) If a full joint prob table has K cells and we initially place a '1' in each cell to avoid having $\text{prob} = 0$, then technically the *numberWins* and *numberLosses* should initially be set to K because world states do not overlap (ie, we imagined K wins and K losses). But in the full joint table approach, K will be very large and the real wins and losses will likely be small compared to K (and, hence "washed out"). So I suggest you initially set *numberWins* and *numberLosses* to 1 (and hence the initial `numberGames` = 2 if you keep that as an explicit counter). In the Bayes Net approaches, the various local probability tables can overlap, so one could initially set *numberWins* and *numberLosses* based on the LARGEST such table; this is still a heuristic rather than necessarily correct, but a reasonable approach. Initially setting *numberWins* = 1 and *numberLosses* = 1 is also fine here, given we are playing millions of games.

4) Remember that we COUNT things in our cells, but our formula uses PROB's. So don't forget that $\text{prob}(\text{win}) = \text{numberWins} / (\text{numberWins} + \text{numberLosses})$ and if you have some conditional prob that, say, depends on three random variables ('features'), then you need to do $\text{prob}(F1=a \text{ and } F2=b \text{ and } F3=c \mid \text{win}) = \text{countWins}(F1=a \text{ and } F2=b \text{ and } F3=c) / \text{numberWins}$ and $\text{prob}(F1=a \text{ and } F2=b \text{ and } F3=c \mid \text{not win}) = \text{countLosses}(F1=a \text{ and } F2=b \text{ and } F3=c) / \text{numberLosses}$.

5) I recommend you first create and debug your NB player, then cut-paste-and-extend it to create your non-NB player. Whether you do the FullJoint player first or last is up to you.

Be sure to monitor the HW3 Forum in Moodle for suggestions, extra information, and (hopefully rarely, if at all) bug fixes.