

Problem Set #5

Northwestern University EECS 348, Spring 2013

Turn in via Blackboard – it’s okay to upload a PDF scan of a hand-written copy. Due Thursday, May 30th at 11:59PM.

White wire	Red wire	Blue wire	Explosion
0	1	1	1
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	1
1	0	1	1
1	0	0	1

Table 1: Training data

White wire	Red wire	Blue wire	Explosion
0	1	1	1
0	0	1	1
1	1	0	1

Table 2: Validation data

- 1) You’re trying to develop a predictor for whether a particular type of explosive device can be disarmed by cutting specific wires within it. The device has four wires of varying colors, and in experiments on recently seized devices, you’ve observed the training data shown in **Table 1**. (1 indicates the wire has been cut, or an explosion).

- a. (4 points) Draw the decision tree that results when using **Table 1** for training. Always split on the feature (wire) that *minimizes* expected entropy after the split, defined as:

$$E_{After}(Explosion | Feature) = \sum_v P(Feature = v)H(Explosion|Feature = v)$$

where

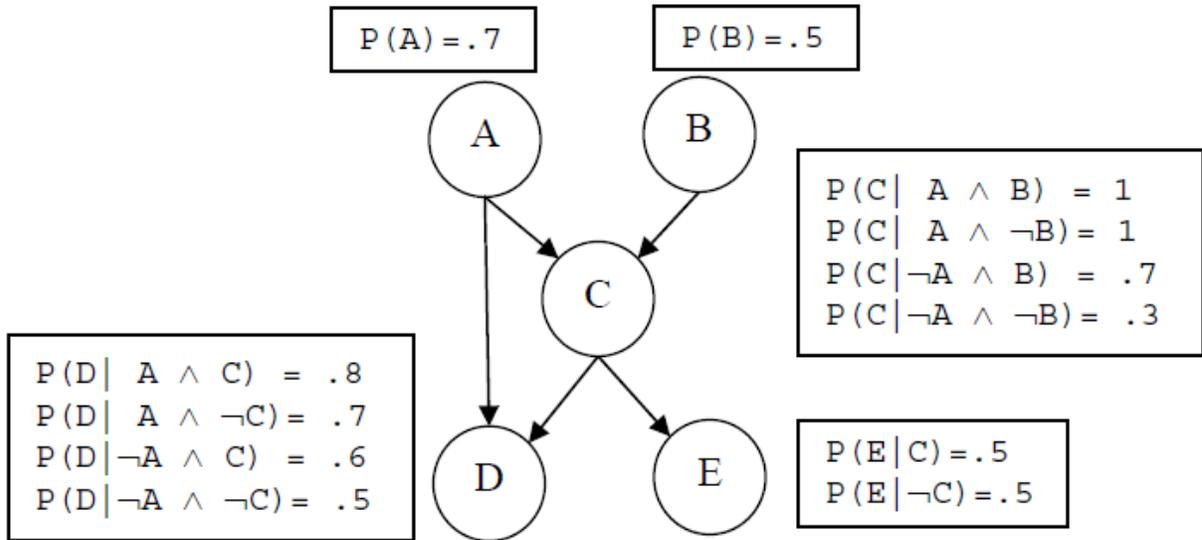
$$H(Explosion|Feature = v) = - \sum_{v'} P(Explosion = v'|Feature = v) \lg P(Explosion = v'|Feature = v)$$

So for example:

$$E_{After}(explosion | red wire) = - (3/7) ((2/3) \lg (2/3) + (1/3) \lg (1/3)) - (4/7) ((3/4) \lg (3/4) + (1/4) \lg (1/4)) = 0.857,$$

Feel free to use the above expressions (e.g. cutting and pasting into a spreadsheet or search engine) to perform your computations quickly. ONLY split if you obtain a decrease in expected entropy by doing so. You should only need to compute about 5 more of these, total.

- b. (1 point) Assume you generate the small set of validation data in Table 2. What is your decision tree's accuracy on the validation data?



- 2) Consider the above Bayesian Network. Here each letter represents a Boolean random variable. The probability distribution for each variable given its parents is shown. In the tables, $P(A)$ indicates the probability that A is 1 (true); to get the probability that A is false, use $1 - P(A)$.
- (1 point) Assume C is known to be true. Will knowing the value of B affect the estimated probability that D is true?
 - (1 point) What is $P(A, \neg B, \neg C, \neg D, E)$? Show your work (in terms of the five factors you multiplied to obtain the final probability).
 - (1/2 point) How many different numbers would you need to specify the full joint distribution of all five variables, if there *weren't* any conditional independencies?
 - (1/2 point) Now let's say you want an even simpler model than the Bayes Net shown above. Instead of using the Bayes Net, you try to predict variable E using a *Naïve Bayes model* of $P(E | A, B, C, D)$. How many different numbers would you need in order to specify the Naïve Bayes model?
- 3) Pick some problem you're interested in modeling that can be described in about five random variables. So for example, you might choose to predict who will win the next Cubs game, based on the ERA of each starting pitcher, whether the game is home or away, and what the winning percentage of the opponent is.
- (1 point) Draw a Bayes Net capturing your domain -- you **don't** need to specify any probability numbers, just the graph.
 - (1 point) Explain one or two conditional independencies exhibited in your graph. If there aren't any, describe the conditional independence that would be most plausible (or "least implausible") in your domain.