

## Coursework Part 1: The Naive Bayesian Network

In the first coursework we will look at calculating joint and conditional probability tables (link matrices) and making inferences with a naive Bayesian network.

Take a look at the first four functions in `DAPICourseworkLibrary.py` to make sure you understand what is going on. You may find the python syntax self explanatory, but if not you can find several good online tutorials through google. The first function reads a data file in the above format, extracting the data so that it can be used. It is extensively commented to help you understand what is going on. The second writes a data array to a file, appending it to what is there already. The data is written out in a fixed floating point format which is suitable for printing probabilities. The third appends a list (one dimensional array or vector) to a text file and the fourth appends a string to a text file. These allow you to save your results in a text file and add comments to them. Thus you can put together your results in a suitable form for submitting your solution to CATE.

Now take a look at the file `DAPICourseworkSkeleton.py`. The first five functions are to be completed for coursework 1 as detailed below. They are all quite short, the longest part in my solution is 11 lines.

### Task 1.1: 4 marks

Complete the definition of function `Prior`. It should calculate the prior distribution over the states of the variable passed as the parameter 'root' in the data array. The first line simply sets up an array for the result with zero entries.

### Task 1.2: 5 marks

Complete the definition of `CPT` which calculates a conditional probability table (or link matrix) between two variables indicated by parameters `varP` (the parent) and `varC` (the child). In the skeleton the conditional probability table is simply initialised to zeros.

### Task 1.3: 4 marks

Complete the definition of function `JPT` which calculates the joint probability table of any two variables. The joint probability is simply the frequency of a state pair occurring in the data. So for state  $a_0$  of variable  $A$  and state  $b_3$  of variable  $B$ , if  $N(a_0 \& b_3)$  is the number of data points containing both  $a_0$  and  $b_3$  then  $P(a_0 \& b_3) = N(a_0 \& b_3) / noDataPoints$ . The table should be arranged so that for  $P(A \& B)$  the states of  $A$  are rows and the states of  $B$  are columns.

### Task 1.4: 4 marks

Complete the function `JPT2CPT` which calculates a conditional probability table from a joint probability table. This is purely a matter of normalising each column to sum to 1. In effect we use the equation  $P(A|B) = P(A \& B) / P(B)$ .

### Task 1.5: 6 marks (harder)

Finally complete the function `Query` which calculates the probability distribution over the root node of a naive Bayesian network. To represent a naive network in Python we will use a list containing an entry for each node (in numeric order) giving the associated probability table: [`prior`, `cpt1`, `cpt2`, `cpt3`, `cpt4`, `cpt5`]. You can calculate the prior of the root and the conditional probability tables between each child variable and the root using your solutions to Tasks 1 and 2. A query is a list of the instantiated states of the child nodes, for example [1,0,3,2,0]. The returned value is a list (or vector) giving the posterior probability distribution over the states of the root node, for example [0.1,0.3,0.4,0.2].

## Results File

You should finish by writing a main program part at the end of the skeleton file. Some example code is there to help you do this. You should use the `Neurones.txt` data set, and create a results file containing: