

# COMP 307/AIML 420 — Introduction to AI

## Assignment 3: Uncertainty and Probability

10% of Final Mark — Due: 23:59 Wednesday 18 May 2022

### 1 Question Description

In the following, unless explicitly specify, a *capital letter* (e.g.,  $A, B, X, Y$ ) represents a *random variable*, and a *lowercase letter* (e.g.,  $a, b, x, y$ ) represents a value.

#### Part 1: Reasoning Under Uncertainty Basics [20 marks]

This part contains several questions about the basics of reasoning under uncertainty. You need to write your answers to each of these questions in your report, and **Show your working**.

For calculations, you need to show the steps in the form like  $P(A = 0|B = 1) = \frac{P(A=0, B=1)}{P(B=1)}$ , to demonstrate that you *know how to calculate* them.

For proving, you also need to clearly show each step of the proof.

#### Question 1 [10 marks]

The tables below give the prior distribution  $P(X)$ , and two conditional distributions  $P(Y|X)$  and  $P(Z|Y)$ . It is also known that  $Z$  and  $X$  are *conditionally independent given  $Y$* . All the three variables ( $X$ ,  $Y$ , and  $Z$ ) are binary variables.

$X$	$P(X)$	$X$	$Y$	$P(Y X)$	$Y$	$Z$	$P(Z Y)$
0	0.35	0	0	0.10	0	0	0.70
1	0.65	0	1	0.90	0	1	0.30
		1	0	0.60	1	0	0.20
		1	1	0.40	1	1	0.80

1. Compute the table of the joint distribution  $P(X, Y, Z)$ . **Show the rule(s) you used, and the steps of calculating each joint probability.**
2. Create the full joint probability table of  $X$  and  $Y$ , i.e., the table containing the following four joint probabilities  $P(X = 0, Y = 0)$ ,  $P(X = 0, Y = 1)$ ,  $P(X = 1, Y = 0)$ ,  $P(X = 1, Y = 1)$ . **Show the rule(s) used, and the steps of calculating each joint probability.**
3. From the above joint probability table of  $X$ ,  $Y$ , and  $Z$ , calculate the following probabilities. **Show your working.**
  - (a)  $P(Z = 0)$ ,
  - (b)  $P(X = 0, Z = 0)$ ,
  - (c)  $P(X = 1, Y = 0|Z = 1)$ ,
  - (d)  $P(X = 0|Y = 0, Z = 0)$ .

**Question 2 [10 marks]**

Consider three Boolean variables  $A$ ,  $B$ , and  $C$  (can take  $t$  or  $f$ ). We have the following probabilities:

- $P(B = t) = 0.7$
- $P(C = t) = 0.4$
- $P(A = t|B = t) = 0.3$
- $P(A = t|C = t) = 0.5$
- $P(B = t|C = t) = 0.2$

We also know that  $A$  and  $B$  are *conditionally independent given  $C$* . Calculate the following probabilities. **Show your working.**

- (i)  $P(B = t, C = t)$
- (ii)  $P(A = f|B = t)$
- (iii)  $P(A = t, B = t|C = t)$
- (iv)  $P(A = t|B = t, C = t)$
- (v)  $P(A = t, B = t, C = t)$

**Question 3 [for AIML420 ONLY, 10 marks]**

Prove the following statements. **Show your working.**

- (i) If  $P(A|B, C) = P(B|A, C)$ , then  $P(A|C) = P(B|C)$
- (ii) If  $P(A|B, C) = P(A)$ , then  $P(B, C|A) = P(B, C)$
- (iii) If  $P(A, B|C) = P(A|C) * P(B|C)$ , then  $P(A|B, C) = P(A|C)$

## Part 2: Naive Bayes Method [30 marks]

This part is to implement the Naive Bayes algorithm, and evaluate the program on the breast cancer dataset to be described below. The program should build a Naive Bayes classifier from the training dataset and apply it to the test set.

### Dataset Description

The *breast cancer* dataset is obtained from the UCI machine learning library (<https://archive-beta.ics.uci.edu/ml/datasets/breast+cancer>).

The original dataset consists of 286 instances that belong to two classes: *no-recurrence-events* and *recurrence-events*.

Each instance is described by 9 categorical attributes (features). The name and domain of each attribute is described as follows:

1. **age** (9 values): 10-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80-89, 90-99
2. **menopause** (3 values): lt40, ge40, premeno
3. **tumor-size** (12 values): 0-4, 5-9, 10-14, 15-19, 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, 50-54, 55-59
4. **inv-nodes** (13 values): 0-2, 3-5, 6-8, 9-11, 12-14, 15-17, 18-20, 21-23, 24-26, 27-29, 30-32, 33-35, 36-39
5. **node-caps** (2 values): yes, no
6. **deg-malig** (3 values): 1, 2, 3
7. **breast** (2 values): left, right
8. **breast-quad** (5 values): left\_up, left\_low, right\_up, right\_low, central
9. **irradiat** (2 values): yes, no

The dataset has some missing values. After removing the instances with missing values, there are 277 instances remaining. Then, these instances are split into the following training and test datasets as follows.:

- **267 training instances:** 189 *no-recurrence-events* + 78 *recurrence-events*.
- **10 test instances:** 7 *no-recurrence-events* + 3 *recurrence-events*.

The datasets are provided in the `breast-cancer-training.csv` and `breast-cancer-test.csv` files.

### Requirements

Your job is to use the Naive Bayes classifier to classify the test instances in the `breast-cancer-test.csv` file.

The pseudo code of the training is given as follows to obtain the (conditional) probabilities of each feature given the class, and the probabilities of each class.

---

**Algorithm 1:** Training of the Naive Bayes Classifier

---

**Input:** The training set.  
**Output:** A probability table.  
// Initialise the count numbers to 1.

```
1 for each class label  $y$  do
2    $count(y) = 1$ ;
3   for each feature  $X_i$  do
4     for each possible value  $x_i$  of feature  $X_i$  do
5        $count(X_i, x_i, y) = 1$ ;

// Count the numbers of each class and feature value based on the training
// instances.
6 for each training instance  $[X_1 = x_1, \dots, X_n = x_n, Y = y]$  do
7    $count(y) = count(y) + 1$ ;
8   for each feature  $X_i$  do
9      $count(X_i, x_i, y) = count(X_i, x_i, y) + 1$ ;

// Calculate the total/denominators.
10  $class\_total = 0$ ;
11 for each class label  $y$  do
12    $class\_total = class\_total + count(y)$ ;
13   for each feature  $X_i$  do
14      $total(X_i, y) = 0$ ;
15     for each possible value  $x_i$  of feature  $X_i$  do
16        $total(X_i, y) = total(X_i, y) + count(X_i, x_i, y)$ ;

// Calculate the probabilities from the counting numbers.
17 for each class label  $y$  do
18    $prob(y) = count(y)/class\_total$ ;
19   for each feature  $X_i$  do
20     for each possible value  $x_i$  of feature  $X_i$  do
21        $prob(X_i, x_i, y) = count(X_i, x_i, y)/total(X_i, y)$ ;

22 return  $prob$ ;
```

---

For the prediction of each test instance, you need to calculate the score of the test instance for each class, and predict the class with the largest score. The score of a class is calculated as follows.

---

**Algorithm 2:** Calculation of the class score.

---

**Input:** A test instance  $[X_1 = x_1, \dots, X_n = x_n]$ , a class label  $y$ , the probability table  $prob$ .  
**Output:** The score.

```
1  $score = prob(y)$ ;
2 for each feature  $X_i$  do
3    $score = score * prob(X_i, x_i, y)$ ;
4 return  $score$ ;
```

---

**You should implement the Naive Bayes method from scratch (not call it from any machine learning library).** Your program should take two file names as command line arguments, construct a classifier from the data in the first file, and then apply the classifier to the data in the second file.

You may write the program code in **Java**, **Python**, **R**, **C/C++**, or any other programming language.

You should submit the following files electronically and also a report.

- (20 marks) **Program code** for your Naive Bayes Classifier (both the source code and the executable program running on ECS School machines),

- (2 marks) `sampleoutput.txt` containing the output of your program on the test dataset, and
- (8 marks) A **report** in PDF, text or DOC format. The report should include:
  1. The conditional probabilities  $P(X_i = x_i | Y = y)$  for each feature  $X_i$  (e.g., age), its possible value  $x_i$  (e.g., 10-19), and each class label  $Y = y$  ( $y$  can be *no-recurrence-events* or *recurrence-events*).
  2. The class probabilities  $P(Y = y)$  for each class label  $Y = y$ .
  3. For each test instance, given the input vector  $\mathbf{X} = [X_1 = x_1, \dots, X_9 = x_9]$ , give the calculated
    - $\text{score}(Y = \text{no-recurrence-events}, \mathbf{X})$ ,
    - $\text{score}(Y = \text{recurrence-events}, \mathbf{X})$ ,
    - predicted class of the input vector.

### Part 3: Building Bayesian Network [25 marks]

This part is to build a Bayesian Network for the problem described below.

#### Problem Description

Dr. Rachel Nicholson is a Professor, who lives far away from her university. So, she prefer to work at home and she only comes to her office if she has research meetings with her postgraduate students, or teaching lectures for undergraduate students, or she has both meetings and teaching:

- The probability for Rachel to have meetings is 70%, the probability of Rachel has lectures is 60%.
- If Rachel has both meetings and lectures, the probability of Rachel comes to her office is 95%.
- If Rachel only has meetings (without lectures), the probability of Rachel comes to her office is 75% because she can Skype with her students.
- If Rachel only has lectures (without meetings), the probability of Rachel comes to her office is 80%.
- If Rachel has neither meetings nor lectures, there is a only 6% chance that she comes to the office.
- When Rachel is in her office, half the time her light is off (when she is trying to hide from others to get work done quickly).
- When she is not in her office, she leaves her light on only 2% of the time since the cleaners come for cleaning.
- When Rachel is in her office, 80% of the time she logged onto the computer.
- Because she sometimes work from home, 20% of the time she is not in her office, she is still logged onto the computer.

Note regarding the calculation, you should show your *working process* of the calculation to demonstrate that you *know how to calculate* them.

#### Requirements

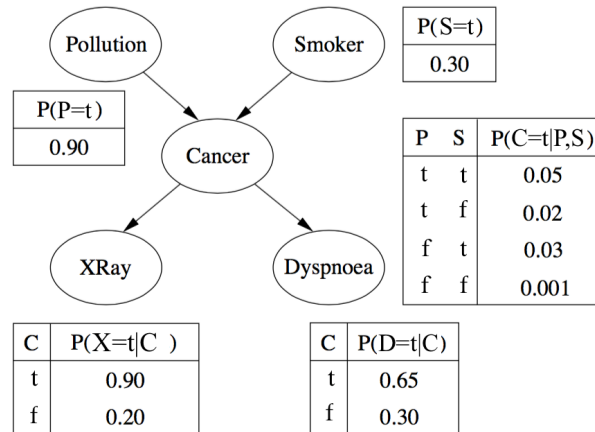
1. Construct a Bayesian network to represent the above scenario. (*Hint: First decide what your domain variables are; these will be your network nodes. Then decide what the causal relationships are between the domain variables and add directed arcs in the network from cause to effect. Finally, you have to add the prior probabilities for nodes without parents, and the conditional probabilities for nodes that have parents.*)

2. Calculate the number of free parameters in your Bayesian network.
3. What is the *joint* probability that Rachel has lectures, has no meetings, she is in her office and logged on her computer but with lights off.
4. Calculate the probability that Rachel is in the office.
5. If we know that Rachel is in the office, what is the *conditional* probability that she is logged on, but her light is off.

## Part 4: Inference in Bayesian Networks [25 marks]

### Problem Description

The following Bayesian Network represents two causes and two effects related to Lung Cancer. Each variable takes the value true (t) or false (f). We will abbreviate the five variable names using their leading letters: P, S, C, X, and D. The probabilities shown are all for the “is true” outcome, e.g. read  $P(P=t) = 0.90$  as the probability that the variable Pollution takes the value true is 0.90. The probability that it is false is not shown, but is easily derived.



### Requirements

Using *inference by enumeration* to calculate the probability  $P(P = t | X = t)$ . You should

1. describe what are the evidence, hidden and query variables in this inference,
2. show each step of the variable elimination in this inference, i.e. to perform the join operation and the elimination operation on the variables in order, and
3. report the final probability.

Note regarding the calculation, you should show your *working process* of the calculation to demonstrate that you *know how to calculate* them.

## Part 5: Bayesian Network: Applications [For AIML420 ONLY, 20 marks]

Identify a real-world application (**different from the examples given in this assignment and the lectures**) that can be described using Bayesian network. There should be at least 5 random variables in this Bayesian network.

**In your report, you should:**

- (i) Clearly define the random variables and their domains.
- (ii) Clearly describe their relationships (using plain language).
- (iii) Draw the Bayesian network that can reflect the described relationship.
- (iv) Write the factorisation of the Bayesian network.

## 2 Submission Guidelines

### 2.1 Submission Requirements

1. Programs (**Executive program file and source files**) for Part 2. Please provide a **readme** file that specifies how to compile and run your program. A script file called **sampleoutput.txt** should also be provided to show how your program run properly. If you programs cannot run properly, you should provide a **buglist** file.
2. A report document that consist of **the answers of all the individual parts**. The document should mark each part clearly. The document can be written in PDF, text or the DOC format.
3. For drawing the diagram such as the Bayesian network, you need to **make the diagram very clear to be marked**. We highly recommend using drawing tools, unless your hand drawing is very clear.

### 2.2 Submission Method

The programs and the PDF report should be submitted through the web submission system (accessible from the COMP307 or AIML420 course web site) **by the due time**. Please ensure you submit to the correct system based on which course you are enrolled in.

### 2.3 Late Penalties

The assignment must be submitted on time unless you have made a prior arrangement with the **course coordinator** or have a valid medical excuse (for minor illnesses it is sufficient to discuss this with the course co-ordinator.) The penalty for assignments that are handed in late without prior arrangement is one grade reduction per day. Assignments that are more than one week late will not be marked.

Remember that you have three late days for this course (which can be used fractionally), but that these apply across the whole course, not per-assignment! Please save some late days in case you need them later on!