# Assignment2 – Bayesian Networks

**Given: May 4    Due: May 14**

**Problem 2.1 (Is your TA in the office?)**                                  0 pt

You want to discuss something with your TA. You know that

1. the probability of your TA being in the office, assuming it is morning, is $\frac{1}{5}$,

2. if your TA is in the office, there is a $\frac{1}{3}$ probability it is morning,

3. the probabilities that it is morning or afternoon are both $\frac{1}{2}$

Your tasks:

1. Write down the probabilities mentioned above as formulas

2. Compute the full joint probability distribution

3. What's the probability you'll meet your TA, if you come to the office in the afternoon?

---

*Solution:*   Let $m$ denote that it is morning and $o$ denote that the TA is in the office.

1.   (a)  $P(o|m) = \frac{1}{5}$

 (b)  $P(m|o) = \frac{1}{3}$

 (c)  $P(m) = P(\neg m) = \frac{1}{2}$

2.   • $P(o, m) = P(o|m) \cdot P(m) = \frac{1}{10}$ *(product rule)*

 • $P(\neg o, m) = \frac{4}{10}$, because $P(m) = P(o, m) + P(\neg o, m)$ *(marginalization)*

 Now, from $P(m|o) \cdot P(o) = P(m, o)$ it follows that $P(o) = \frac{\frac{1}{10}}{\frac{1}{3}} = \frac{3}{10}$. So we get

 • $P(o, \neg m) = \frac{2}{10}$, because $P(o) = P(o, m) + P(o, \neg m)$ *(marginalization)*

 • $P(\neg o, \neg m) = \frac{3}{10}$, because $1 - P(o) = P(\neg o) = P(\neg o, m) + P(\neg o, \neg m)$ *(marginalization)*

3.  $P(o|\neg m) = \frac{P(o, \neg m)}{P(\neg m)} = \frac{4}{10}$

---

**Problem 2.2 (AFT Tests)**                                  30 pt

Trisomy 21 (*Down syndrome*) is a genetic anomaly that can be diagnosed during pregnancy using an amniotic fluid test.

The probability of a foetus having Down syndrome is strongly correlated with the age of the pregnant parent. We will only consider the following two age groups.

1. For 25 year olds the probability is one in 1250,

2. for 43 year old parents it increases to one in fifty.

However, diagnostic tests are never perfect. We distinguish two kinds of errors:

3. Type I Error (False Positive): The test result is positive even though the child is healthy.

4. Type II Error (False Negative): The test result is negative even though the child has trisomy 21.

The probabilities of Type I and Type II Errors are both merely 1% for amniotic fluid tests for Down syndrome.

1. Express the four items above in the form of conditional probabilities. Use the random variable $F$ with domain $\{Age_{25}, Age_{43}\}$ for the age of the pregnant person and the Boolean random variables $Pos$ and $Down$ for the propositions "*The amniotic fluid test is positive*" and "*The child has Down syndrome*" respectively.

2. Assume that we have a 25 year old pregnant person. Using Bayes' theorem, express and compute the probability that their child has Down syndrome, given that the amniotic fluid test is positive. What can we conclude from the result?

---

*Solution:*

1. $P(Down \mid F = Age_{25}) = 0.0008$, $P(Down \mid F = Age_{43}) = 0.02$, $P(Pos \mid \neg Down) = 0.01$, $P(\neg Pos \mid Down) = 0.01$.

2. We normalize to $F = Age_{25}$, making $P(Down) = 0.0008$ and compute:

$$
\begin{aligned}
P(Down \mid Pos) &= \frac{P(Pos \mid Down) \cdot P(Down)}{P(Pos)} = \frac{P(Pos \mid Down) \cdot P(Down)}{P(Pos \wedge Down) + P(Pos \wedge \neg Down)} \\
&= \frac{P(Pos \mid Down) \cdot P(Down)}{P(Pos \mid Down) \cdot P(Down) + P(Pos \mid \neg Down) \cdot P(\neg Down)} \\
&= \frac{(1 - P(\neg Pos \mid Down)) \cdot P(Down)}{(1 - P(\neg Pos \mid Down) \cdot P(Down)) + P(Pos \mid \neg Down) \cdot (1 - P(Down))} \\
&= \frac{0.99 \cdot 0.0008}{0.99 \cdot 0.0008 + 0.01 \cdot 0.9992} \approx 0.07
\end{aligned}
$$

So, even with a positive test result, the probability of the child actually having Down syndrome is still only 7%, simply due to Down syndrome being relatively rare in young parents. Consequently, there is little point in applying this particular test without exceptional cause for concern.
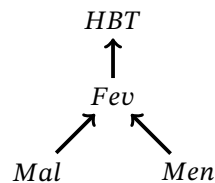
**Problem 2.3 (Medical Bayesian Network)** 30 pt

  Both Malaria and Meningitis can cause a fever, which can be measured by check-
ing for a high body temperature. Of course you may also have a high body temper-
ature for other reasons. We consider the following random variables for a given
patient:

- *Mal*: The patient has malaria.

- *Men*: The patient has meningitis.

- *HBT*: The patient has a high body temperature.

- *Fev*: The patient has a fever.

Consider the following Bayesian network for this situation:

$$HBT$$
$$\uparrow$$
$$Fev$$
$$\nearrow \quad \nwarrow$$
$$Mal \qquad Men$$

1. Explain the purpose of the edges in the network regarding the conditional
   probability table.

2. What would have happened if we had constructed the network using the vari-
   able order $Mal, Men, HBT, Fev$? Would that have led to a better network?

3. How do we compute the probability distribution for the patient having malaria,
   given that he has high body temperature? State the query variables, hidden
   variables and evidence and write down the equation for the probability we
   are interested in.

---

*Solution:*

1. The parents (i.e, nodes from which there are incoming edges) of $X$ are the
   variables that $X$ may depend on. The conditional probability table for $X$ must
   take all of those as additional inputs.

2. We would have obtained additional edges from $Mal$ and $Men$ to $Fev$ because
   they affect the probability of fever. That would be a worse network because
   more edges increase the complexity.

3. Query variable: $Mal$. Evidence: $HBT$. Hidden variables: $Men, Fev$. We get:

   - start
     $$P(Mal|HBT = true)$$

3

- normalization to turn the conditional distribution into an unconditional one
$$= \alpha P(Mal, HBT = true)$$
where $\alpha = 1/P(HBT = true)$ is the constant factor that normalizes the vector $\langle P(Mal = true, HBT = true), P(Mal = false, HBT = true)\rangle$
- marginalization to bring in the hidden variables
$$= \alpha \sum_{m,f} P(Mal, HBT = true, Men = m, Fev = f)$$
where $m$ and $f$ range over the possible values of $Men$ and $Fev$
- chain rule to turn the joint distribution into a product of conditional ones, ordering the variables according to the structure of the network
$$= \alpha \cdot \sum_{m,f} P(Mal) \cdot P(Men = m | Mal) \cdot P(Fev = f | Mal, Men = m) \cdot$$
$$P(HBT = true | Fev = f, Mal, Men = m)$$
Note that each factor is a vector with two entries, one for $Mal = true$ and one for $Mal = false$. These vectors are multiplied component-wise.
- use the structure of the network to drop redundant conditions
$$= \alpha \cdot \sum_{m,f} P(Mal) \cdot P(Men = m) \cdot P(Fev = f | Mal, Men = m) \cdot P(HBT = true | Fev = f)$$
Now all factors are entries of the conditional probability table of the network, which can be plugged in to compute the result.

Because $Mal$ is boolean, we could have started with $P(Mal = true | HBT = true)$ right away. Then we could have skipped the normalization step and would not have to multiply vectors. But for variables with many values, the above is practical because it derives the entire distribution in one go.

---

**Problem 2.4 (Bayesian Networks in Python)** 40 pt

The goal of this exercise is to implement inference by enumeration in Bayesian networks in Python. You can find the necessary files at `https://kwarc.info/teaching/AI/resources/AI2/bayes/`.

Your task is to implement the `query` function in `bayes.py`. Use `test.py` for testing your implementation.

**Important:** We will test your code automatically. So please make sure that:

- The tests in `test.py` work on your code (without any modifications to `test.py`)

- You use a recent Python version ($\geq 3.5$)

- You don't use any libraries

- You only upload a single file `bayes.py` with your implementation of `query`

Otherwise you risk getting no points.

*Hint:* First implement a function for the full joint probability distribution.