

Problem 1. (Simpson's Paradox)

A dangerous new virus is sweeping the world. Currently, there are two potential drug treatments (A and B) for patients. Dr. Homer Simpson wants to compare the un-cured rate of patients after receiving either treatment A or B , in order to determine the better drug.

The data indicates that there are 240 patients that are not cured among the 1500 patients who received treatment A . There are 105 patients that are not cured among the 550 who received treatment B . *Note: this is a fictitious scenario and we made up these numbers.*

Problem 1.a. Can you help Homer construct a probabilistic graphical model for the above scenario.

Solution:

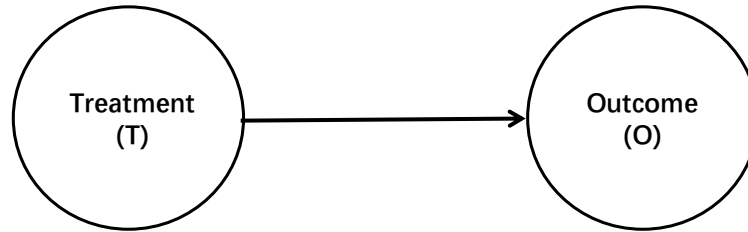


Figure 1: PGM constructed by ourselves.

Define two binary random variables $\text{Treatment}(T)$, which takes the value of A and B , and $\text{Outcome}(O)$, which takes the value of Cured(0) and Not cured(1). We can define a Bayes Net as Fig.1. The joint distribution $p(T, O) = p(O|T)p(T)$. From the data, we can compute the MLE estimates (assuming Categorical distributions):

$$p(O = 1|T = A) = \frac{240}{1500} = 0.16 \quad (1)$$

$$p(O = 0|T = A) = 1 - \frac{240}{1500} = 0.84 \quad (2)$$

Similarly, for the treatment B

$$p(O = 1|T = B) = \frac{105}{550} = 0.19 \quad (3)$$

$$p(O = 0|T = B) = 1 - \frac{105}{550} = 0.81 \quad (4)$$

For prior distribution $p(T)$, we can compute it as

$$p(T = A) = \frac{1500}{1500 + 550} = 0.73 \quad (5)$$

$$p(T = B) = \frac{550}{1500 + 550} = 0.27 \quad (6)$$

Problem 1.b. The data seems to indicate that treatment A is more effective. Can Homer confirm (just from the data) that one of the treatments results in more cures? *Hint:* Consider what happens when there are unobserved variables that could affect the treatment and the outcome.

Solution:

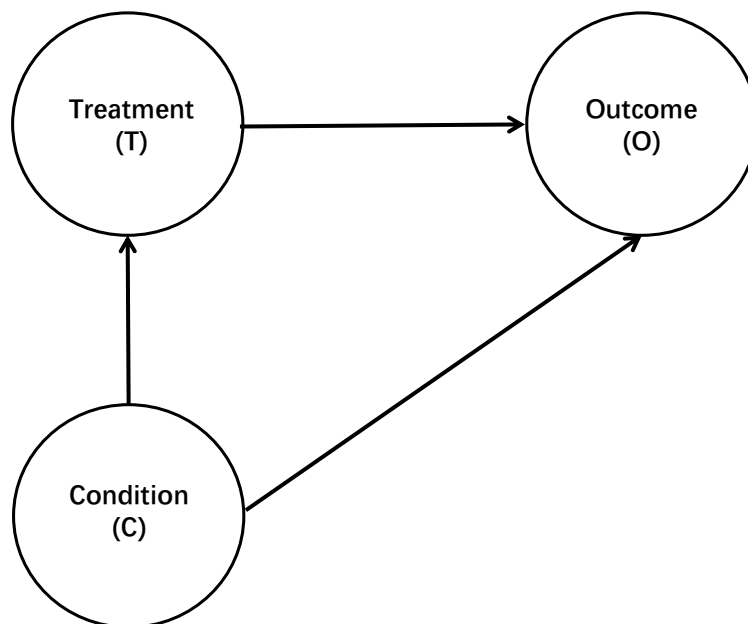


Figure 2: True PGM of the underlying causal process.

We cannot confirm that treatment A is more effective, since we may have some unobserved variables. Suppose the true causal process is represented by Fig.2, where there is another random variable Condition. Condition is a binary random variable that takes value of *Mild*(0) or *Severe*(1). It represents the severity of the patient's sickness. If a patient is in the severe condition, doctors tend to give him treatment B (e.g. B has better a treatment outcome, but is more expensive). If a patient is under the mild condition, doctors tend to give him treatment A (cheaper but less effective).

Suppose among the 1500 patients who received treatment A , there are 1400 patients under the *Mild* condition with 210 patients that are not cured. There are 100 patients under *Severe* condition with 30 patients that are not cured. Among the 550 patients who received treatment B , there are 50 patients under the *Mild* condition with 5 patients that are not cured and 500 patients under the *Severe* condition with 100 patients that are not cured.

Given this information, we can compute the probability of un-cured patients who received treatment A

or B under different conditions.

$$p(O = 1|T = A, C = 0) = \frac{210}{1400} = 0.15 \quad (7)$$

$$p(O = 1|T = A, C = 1) = \frac{30}{100} = 0.3 \quad (8)$$

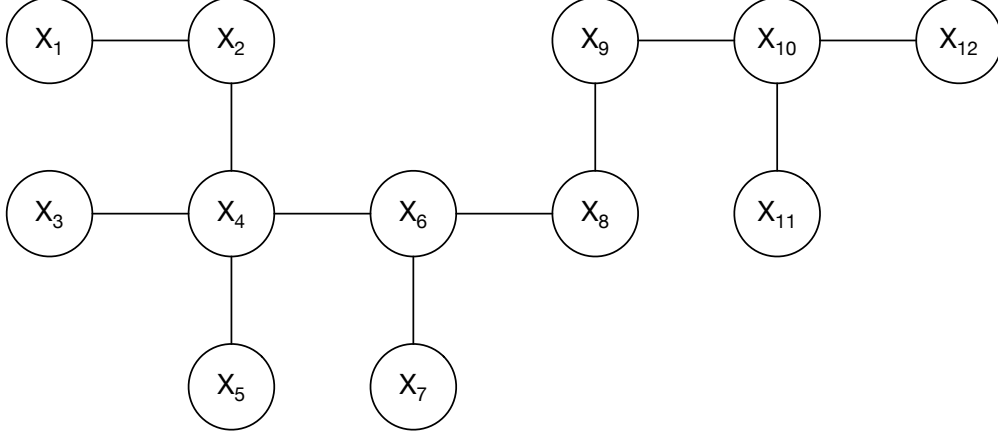
$$p(O = 1|T = B, C = 0) = \frac{5}{50} = 0.1 \quad (9)$$

$$p(O = 1|T = B, C = 1) = \frac{100}{500} = 0.2 \quad (10)$$

We can see that under each condition, the un-cured rate of patients who received treatment B is smaller than that of patients who received treatment A . The reason is that the majority of patients who received treatment B are under *Severe* condition. Therefore, even we know that B is better, the overall un-cured rate is larger than patients who received A because of the allocation of patients to the different treatment groups.

Problem 2. (MRT Inference)

You are given the following *pairwise* undirected graphical model which models the activity (low or high) at 12 MRT stations.



Each node represents a random variable indicating whether the activity at a particular station is low (0) or high (1). Assume the following factorization:

$$p(x_1, x_2, \dots, x_{12}) = \frac{1}{Z} \prod_{i \in V} \psi(x_i) \prod_{(i,j) \in E} \psi(x_i, x_j) \quad (11)$$

where V is the set of nodes, E is the set of edges, and that the unary and pairwise factors are given by:

x_i	$\psi(x_i)$
0	10
1	2

Figure 3: Unary Factors

x_i	x_j	$\psi(x_i, x_j)$
0	0	20
0	1	5
1	0	5
1	1	20

Figure 4: Pairwise Factors

Note that the factors are the same across the nodes. Your task is to compute the following conditional probabilities.

Problem 2.a. Compute $p(x_{12} = 1 | x_1 = 0, x_7 = 0, x_9 = 1, x_{10} = 0)$.

Solution: According to the conditional independence assertions in the MRF.

$$p(x_{12} | x_1, x_7, x_9, x_{10}) = p(x_{12} | x_{10}) = \frac{p(x_{10}, x_{12})}{\sum_{x_{12}} p(x_{10}, x_{12})}$$

Denote

$$m(x_{10}) = \sum_{i,j \in M_{E/x_{10}, x_{12}}} \psi(x_i) \psi(x_i, x_j)$$

Then

$$\begin{aligned} p(x_{12} | x_1, x_7, x_9, x_{10}) &= \frac{p(x_{10}, x_{12})}{\sum_{x_{12}} p(x_{10}, x_{12})} = \frac{\sum_{i,j \in M_{E/x_{10}, x_{12}}} p(x_1 \dots x_{12})}{\sum_{x_{12}} \sum_{i,j \in M_{E/x_{10}, x_{12}}} p(x_1 \dots x_{12})} \\ &= \frac{m(x_{10}) \psi(x_{10}) \psi(x_{10}, x_{12}) \psi(x_{12})}{\sum_{x_{12}} m(x_{10}) \psi(x_{10}) \psi(x_{10}, x_{12}) \psi(x_{12})} \end{aligned}$$

$$\begin{aligned}
&= \frac{\psi(x_{10}, x_{12})\psi(x_{12})}{\sum_{x_{12}} \psi(x_{10}, x_{12})\psi(x_{12})} = \frac{\psi(x_{10} = 0, x_{12} = 1)\psi(x_{12} = 1)}{\sum_{x_{12}} \psi(x_{10} = 0, x_{12})\psi(x_{12})} \\
&= \frac{5 \times 2}{20 \times 10 + 5 \times 2} = \frac{1}{21} = 0.0476
\end{aligned}$$

Problem 2.b. Compute $p(x_1 = 1 | x_3 = 0, x_4 = 1, x_6 = 0)$.

Solution: According to the conditional independence assertions in MRF.

$$p(x_1 | x_3, x_4, x_6) = p(x_1 | x_4) = \frac{p(x_1, x_4)}{\sum_{x_1} p(x_1, x_4)}$$

Denote

$$m(x_4) = \sum_{i,j \in M_{E/x_1, x_2, x_4}} \psi(x_i)\psi(x_i, x_j)$$

Then

$$\begin{aligned}
p(x_1 | x_3, x_4, x_6) &= \frac{\sum_{x_2} \sum_{i,j \in M_{E/x_1, x_2, x_4}} p(x_1 \dots x_{12})}{\sum_{x_2} \sum_{x_1} \sum_{i,j \in M_{E/x_1, x_2, x_4}} p(x_1 \dots x_{12})} \\
&= \frac{\sum_{x_2} \psi(x_1, x_2)\psi(x_2)\psi(x_2, x_4)\psi(x_1)\psi(x_4)m(x_4)}{\sum_{x_1} \sum_{x_2} \psi(x_1, x_2)\psi(x_2)\psi(x_2, x_4)\psi(x_1)\psi(x_4)m(x_4)} \\
&= \frac{\sum_{x_2} \psi(x_1, x_2)\psi(x_2)\psi(x_2, x_4)\psi(x_1)}{\sum_{x_1} \sum_{x_2} \psi(x_1, x_2)\psi(x_2)\psi(x_2, x_4)\psi(x_1)} = \frac{\sum_{x_2} \psi(x_1 = 1, x_2)\psi(x_2)\psi(x_2, x_4 = 1)\psi(x_1 = 1)}{\sum_{x_1} \sum_{x_2} \psi(x_1, x_2)\psi(x_2)\psi(x_2, x_4 = 1)\psi(x_1)} \\
&= \frac{5 \times 10 \times 5 \times 2 + 20 \times 2 \times 20 \times 2}{20 \times 10 \times 5 \times 10 + 5 \times 2 \times 20 \times 10 + 5 \times 10 \times 5 \times 2 + 20 \times 2 \times 20 \times 2} = \frac{7}{47} = 0.1489
\end{aligned}$$

Problem 2.c. Compute $p(x_{10} = 1 | x_9 = 1, x_{12} = 1, x_2 = 0)$.

Solution: According to the conditional independence assertions in MRF.

$$p(x_{10} | x_9, x_{12}, x_2) = p(x_{10} | x_9, x_{12}) = \frac{p(x_{10}, x_9, x_{12})}{p(x_9, x_{12})}$$

Denote

$$m(x_9) = \sum_{i,j \in M_{E/x_9, x_{10}, x_{11}, x_{12}}} \psi(x_i)\psi(x_i, x_j)$$

Then

$$\begin{aligned}
p(x_{10} | x_9, x_{12}, x_2) &= \frac{p(x_{10}, x_9, x_{12})}{p(x_9, x_{12})} = \frac{\sum_{x_{11}} \sum_{i,j \in M_{E/x_9, x_{10}, x_{11}, x_{12}}} p(x_1 \dots x_{12})}{\sum_{x_{10}} \sum_{x_{11}} \sum_{i,j \in M_{E/x_9, x_{10}, x_{11}, x_{12}}} p(x_1 \dots x_{12})} \\
&= \frac{\sum_{x_{11}} \psi(x_9)\psi(x_9, x_{10})\psi(x_{10})\psi(x_{10}, x_{11})\psi(x_{11})\psi(x_{10}, x_{12})\psi(x_{12})m(x_9)}{\sum_{x_{10}} \sum_{x_{11}} \psi(x_9)\psi(x_9, x_{10})\psi(x_{10})\psi(x_{10}, x_{11})\psi(x_{11})\psi(x_{10}, x_{12})\psi(x_{12})m(x_9)} \\
&= \frac{\sum_{x_{11}} \psi(x_9)\psi(x_9, x_{10})\psi(x_{10})\psi(x_{10}, x_{11})\psi(x_{11})\psi(x_{10}, x_{12})\psi(x_{12})m(x_9)}{\sum_{x_{10}} \sum_{x_{11}} \psi(x_9)\psi(x_9, x_{10})\psi(x_{10})\psi(x_{10}, x_{11})\psi(x_{11})\psi(x_{10}, x_{12})\psi(x_{12})m(x_9)}
\end{aligned}$$

$$\begin{aligned}
&= \frac{\sum_{x_{11}} \psi(x_9 = 1, x_{10} = 1) \psi(x_{10} = 1) \psi(x_{10} = 1, x_{11}) \psi(x_{11}) \psi(x_{10} = 1, x_{12} = 1)}{\sum_{x_{10}} \sum_{x_{11}} \psi(x_9 = 1, x_{10}) \psi(x_{10}) \psi(x_{10}, x_{11}) \psi(x_{11}) \psi(x_{10}, x_{12} = 1)} \\
&= \frac{20 \times 2 \times 5 \times 10 \times 20 + 20 \times 2 \times 20 \times 2 \times 20}{5 \times 10 \times 20 \times 10 \times 5 + 20 \times 2 \times 5 \times 10 \times 20 + 5 \times 10 \times 5 \times 2 \times 5 + 20 \times 2 \times 20 \times 2 \times 20} \\
&= \frac{48}{83} = 0.5783
\end{aligned}$$

Problem 2.d. Compute $p(x_6 = 0 | x_4 = 1, x_8 = 1, x_{10} = 0)$.

Solution: According to the conditional independence assertions in MRF.

$$p(x_6 | x_4, x_8, x_{10}) = p(x_6 | x_4, x_8) = \frac{\sum_{x_7} p(x_4, x_6, x_8, x_7)}{\sum_{x_7} \sum_{x_7} p(x_4, x_6, x_8, x_7)}$$

Then

$$\begin{aligned}
p(x_6 = 0 | x_4 = 1, x_8 = 1) &= \frac{\sum_{x_7} \psi(x_4 = 1, x_6 = 0) \psi(x_8 = 1, x_6 = 0) \psi(x_7) \psi(x_7, x_6 = 0) \psi(x_6 = 0)}{\sum_{x_6} \sum_{x_7} \psi(x_4 = 1, x_6) \psi(x_8 = 1, x_6) \psi(x_7) \psi(x_7, x_6) \psi(x_6)} \\
&= \frac{10 \times 5 \times 5 \times 2 \times 5 + 5 \times 5 \times 10 \times 20 \times 10}{10 \times 5 \times 5 \times 2 \times 5 + 5 \times 5 \times 10 \times 20 \times 10 + 20 \times 20 \times 10 \times 5 \times 2 + 20 \times 20 \times 2 \times 20 \times 2} = \frac{35}{83} = 0.4217
\end{aligned}$$

Problem 2.e. Compute $p(x_8 = 1 | x_1 = 0, x_6 = 0, x_9 = 1, x_{12} = 1)$.

Solution: According to the conditional independence assertions in MRF.

$$p(x_8 | x_1, x_6, x_9, x_{12}) = p(x_8 | x_6, x_9) = \frac{p(x_8, x_6, x_9)}{\sum_{x_8} p(x_8, x_6, x_9)}$$

Then

$$\begin{aligned}
p(x_8 = 1 | x_6 = 0, x_9 = 1) &= \frac{\psi(x_8 = 1, x_6 = 0) \psi(x_9 = 1, x_8 = 1) \psi(x_8 = 1)}{\sum_{x_8} \psi(x_8, x_6 = 0) \psi(x_9 = 1, x_8) \psi(x_8)} \\
&= \frac{5 \times 20 \times 2}{5 \times 20 \times 2 + 20 \times 5 \times 10} = \frac{1}{6} = 0.1667
\end{aligned}$$

Problem 2.f. Compute $p(x_2 = 0 | x_1 = 0, x_3 = 1, x_4 = 1, x_7 = 1, x_{11} = 0)$.

Solution: According to the conditional independence assertions in MRF.

$$p(x_2 | x_1, x_3, x_4, x_7, x_{11}) = p(x_2 | x_1, x_4) = \frac{p(x_1, x_2, x_4)}{\sum_{x_2} p(x_1, x_2, x_4)}$$

Then

$$\begin{aligned} p(x_2|x_1, x_4) &= \frac{\psi(x_2 = 0, x_1 = 0)\psi(x_4 = 1, x_2 = 0)\psi(x_2 = 0)}{\sum_{x_2} \psi(x_2, x_1 = 0)\psi(x_4 = 1, x_2)\psi(x_2)} \\ &= \frac{20 \times 5 \times 10}{5 \times 20 \times 2 + 20 \times 5 \times 10} = \frac{5}{6} = 0.8333 \end{aligned}$$

Problem 3. (Image Denoising)

For this problem, you will be working on Image Denoising, taking a noisy image and making it a clean one. Please refer to the provided `Image-Denoising-Pre.ipynb` notebook. You can download the notebook and relevant images in a zipfile from NUS Canvas (in the MRF module under Home).

To use the notebook, you have to install jupyter (<https://jupyter.org/install>) and a python distro; we use anaconda (<https://www.anaconda.com/products/individual>) but you can use whichever distribution you like.