
Machine Learning, 2024 Fall

Homework 2

Notice

Due 23:59 (CST), Nov 19, 2024

Plagiarizer will get 0 points.

\LaTeX is highly recommended. Otherwise you should write as legibly as possible.

A Bayesian Network

(24pt) The figure below shows a Bayesian network, illustrating the conditional dependencies among the following variables: Season (S), Flu (F), Dehydration (D), Chills (C), Headache (H), Nausea (N), and Dizziness (Z).

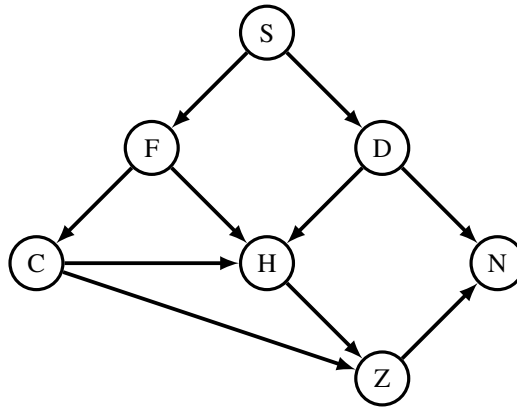


Figure 1: A Bayesian network that represents the conditional dependencies among the variables: Season, Flu, Dehydration, Chills, Headache, Nausea, and Dizziness.

A.1 Independence in Bayesian Network

(12pt) Justify the following independence statements are true or false and give a brief explanation.

- $S \perp N | D$
- $S \perp Z | H$
- $F \perp N | Z, D$
- $F \perp Z | C, H, N$
- $C \perp N | Z$
- $C \perp D | S$

- False: influence can flow along the path $S \rightarrow F \rightarrow H \rightarrow Z \rightarrow N$
- False: influence can flow along the path $S \rightarrow F \rightarrow C \rightarrow Z$
- True: influence cannot flow along any path
- False: influence can flow along the path $F \leftarrow S \rightarrow D \rightarrow N \leftarrow Z$
- False: influence can flow along the path $C \leftarrow F \leftarrow S \rightarrow D \rightarrow N$
- True: influence cannot flow along any path

A.2 Evaluating Probability Queries

(12pt) Given the conditional probability tables for the Bayesian network in the following table, calculate each of the queried probabilities. **Note:** You have to write down the calculation steps if necessary.

(1) Use conditional independency properties, write down the factorized form of the conditional independence over all of the variables, $P(S, F, D, C, H, N, Z)$

$$P(S, F, D, C, H, N, Z) = P(S)P(F|S)P(D|S)P(C|F)P(H|F, D, C)P(N|D, Z)P(Z|C, H)$$

(2) What is the probability that you have the flu, given that it is summer?

The problem can be translated to: $P(\text{Flu}=\text{true}|\text{Season}=\text{summer})$

$$\begin{aligned} & \sum_s P(F = \text{true}, S = s) \\ &= 0.2 \end{aligned}$$

(3) What is the probability that you are feeling chill, when no prior information is known?

The problem can be translated to: $P(\text{Chills}=\text{true})$

$$\begin{aligned} & \sum_{f,s} P(C = \text{true}, F = f, S = s) \\ &= \sum_{f,s} P(C = \text{true}|F = f) P(F = f|S = s) P(S = s) \\ &= 0.5 * 0.2 * 0.7 + 0.5 * 0.8 * 0.8 + 0.5 * 0.4 * 0.7 + 0.5 * 0.6 * 0.8 \\ &= 0.77 \end{aligned}$$

(4) What is the probability that you have the flu, given that it is summer and that you have a headache, and you know that you are dehydrated?

The problem can be translated to: $P(\text{Flu}=\text{true}|\text{Season}=\text{summer}, \text{Headache}=\text{true})$

$$\begin{aligned} & P(F = \text{true}|S = \text{summer}, H = \text{true}) \\ &= \frac{P(F = \text{true}, S = \text{summer}, H = \text{true})}{P(S = \text{summer}, H = \text{true})} \\ &= \frac{\sum_c P(F = \text{true}, S = \text{summer}, H = \text{true}, D = \text{true}, C = c)}{\sum_{f,c} P(F = f, S = \text{summer}, H = \text{true}, D = \text{true}, C = c)} \\ &= \frac{\sum_c P(H = \text{true}|F = \text{true}, D = \text{true}, C = c) P(C = c|F = \text{true}) P(F = \text{true}|S = \text{summer}) P(D = \text{true}|S = \text{summer}) P(S = \text{summer})}{\sum_{f,c} P(H = \text{true}|F = f, D = \text{true}, C = c) P(C = c|F = f) P(F = f|S = \text{summer}) P(D = \text{true}|S = \text{summer}) P(S = \text{summer})} \\ &= \frac{0.5 * 0.1 * 0.2 * 0.7 * 0.7 + 0.5 * 0.1 * 0.2 * 0.3 * 0.3}{0.5 * 0.1 * 0.2 * 0.7 * 0.7 + 0.5 * 0.1 * 0.2 * 0.3 * 0.3 + 0.5 * 0.1 * 0.8 * 0.8 * 0.5 + 0.5 * 0.1 * 0.8 * 0.2 * 0.1} \\ &= \frac{0.0058}{0.0226} \\ &= 0.26 \end{aligned}$$

| $P(S = \text{winter})$ | $P(S = \text{summer})$ | | $P(F = \text{true} S)$ | $P(F = \text{false} S)$ |
|------------------------|------------------------|---------------------|------------------------|-------------------------|
| 0.5 | 0.5 | $S = \text{winter}$ | 0.4 | 0.6 |
| | | $S = \text{summer}$ | 0.2 | 0.8 |

| | $P(D = \text{true} S)$ | $P(D = \text{false} S)$ | | $P(C = \text{true} F)$ | $P(C = \text{false} F)$ |
|---------------------|------------------------|-------------------------|--------------------|------------------------|-------------------------|
| $S = \text{winter}$ | 0.3 | 0.7 | $F = \text{true}$ | 0.7 | 0.3 |
| $S = \text{summer}$ | 0.1 | 0.9 | $F = \text{false}$ | 0.8 | 0.2 |

| | $P(H = \text{true} C, F, D)$ | $P(H = \text{false} C, F, D)$ |
|--|------------------------------|-------------------------------|
| $C = \text{true}, F = \text{true}, D = \text{true}$ | 0.7 | 0.3 |
| $C = \text{true}, F = \text{true}, D = \text{false}$ | 0.6 | 0.4 |
| $C = \text{true}, F = \text{false}, D = \text{true}$ | 0.5 | 0.5 |
| $C = \text{true}, F = \text{false}, D = \text{false}$ | 0.4 | 0.6 |
| $C = \text{false}, F = \text{true}, D = \text{true}$ | 0.3 | 0.7 |
| $C = \text{false}, F = \text{true}, D = \text{false}$ | 0.2 | 0.8 |
| $C = \text{false}, F = \text{false}, D = \text{true}$ | 0.1 | 0.9 |
| $C = \text{false}, F = \text{false}, D = \text{false}$ | 0.8 | 0.2 |

| | $P(N = \text{true} D, Z)$ | $P(N = \text{false} D, Z)$ |
|--------------------------------------|---------------------------|----------------------------|
| $D = \text{true}, Z = \text{true}$ | 0.7 | 0.3 |
| $D = \text{true}, Z = \text{false}$ | 0.8 | 0.2 |
| $D = \text{false}, Z = \text{true}$ | 0.2 | 0.8 |
| $D = \text{false}, Z = \text{false}$ | 0.5 | 0.5 |

| | $P(Z = \text{true} C, H)$ | $P(Z = \text{false} C, H)$ |
|--------------------------------------|---------------------------|----------------------------|
| $C = \text{true}, H = \text{true}$ | 0.1 | 0.9 |
| $C = \text{true}, H = \text{false}$ | 0.3 | 0.7 |
| $C = \text{false}, H = \text{true}$ | 0.4 | 0.6 |
| $C = \text{false}, H = \text{false}$ | 0.8 | 0.2 |

B Variable Elimination

(26pt) Given a Bayesian network in the following figure, which consists of binary variables. We will use variable elimination. The chosen variable elimination ordering is A, C, E, G to compute the query $P(B, D, H|f = 1)$.

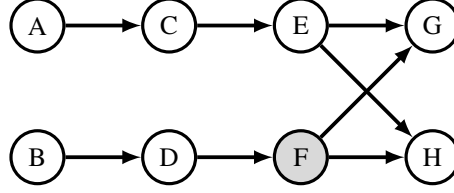


Figure 2: A Bayesian network.

(1) (3pt) What is the corresponding moral graph?

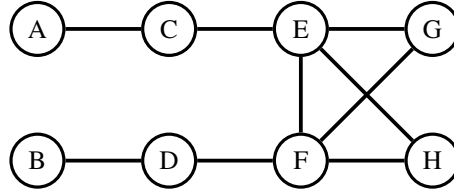


Figure 3: Moral graph.

(2) (3pt) Write down all initial factors after inserting evidence $f = 1$.

$$P(a), P(b), P(c|a), P(d|b), P(e|c), f_{P(g|e,f)}(g, e), f_{P(h|e,f)}(h, e)$$

(3) (15pt) Run variable elimination

1. When eliminating A , please write all factors including the new generated factor f_1 .
2. When eliminating C , please write all factors including the new generated factor f_2 .
3. When eliminating E , please write all factors including the new generated factor f_3 .
4. When eliminating G , please write all factors including the new generated factor f_4 .
5. Compute $P(B, D, H|f = 1)$ from the factors left in 4

(1) $P(b), P(d|b), P(e|c), f_{P(g|e,f)}(g, e), f_{P(h|e,f)}(h, e), f_1(c)(f_a(c) \text{ is ok})$

(2) $P(b), P(d|b), f_{P(g|e,f)}(g, e), f_{P(h|e,f)}(h, e), f_2(e)(f_c(e) \text{ is ok})$

(3) $P(b), P(d|b), f_3(g, h)(f_e(g, h) \text{ is ok})$

(4) $P(b), P(d|b), f_4(h)(f_g(h) \text{ is ok})$

(5)
$$P(B, D, H|f = 1) = \frac{P(b)P(d|b)f_4(h)}{\sum_{B, D, H} P(b)P(d|b)f_4(h)}$$

Note: You have to write all factors in each question, or you will loss some points.

(4) (5pt) Among the factors, f_1, f_2, f_3, f_4 , which is the most largest factor generated, and what is the size of this factor? In this context, we assume that all variables have binary domains, and we determine the factor size by counting the number of rows in the table that represents the factor.

$f_3(g, h)$ is the largest factor generated. It has 2 variables, hence the factor size is $2^2 = 4$.

C Learning in Naive Bayesian Networks

Assume a training dataset $\mathcal{D} = \{\mathbf{x}, y\}$, where $\mathbf{x} \in R^d$ represents different features and y represents categories, and let \mathcal{G} be a network over these variables.

(1) If \mathcal{G} is a Naive Bayesian network, consider a apple quality dataset that describes the features of good and bad apples. The training data is as follows:

| Color | Skin Smoothness | Size | Aroma | Flesh Firmness | Category |
|--------|-----------------|--------|--------|----------------|------------|
| Red | Smooth | Medium | Strong | Hard | Good Apple |
| Green | Rough | Large | Faint | Medium | Bad Apple |
| Red | Smooth | Medium | Faint | Soft | Bad Apple |
| Red | Rough | Large | Strong | Hard | Good Apple |
| Green | Smooth | Small | Strong | Hard | Good Apple |
| Red | Smooth | Medium | Strong | Medium | Good Apple |
| Yellow | Rough | Medium | Faint | Soft | Bad Apple |
| Red | Smooth | Large | Strong | Hard | Good Apple |
| Green | Rough | Small | Faint | Soft | Bad Apple |
| Red | Smooth | Medium | Strong | Hard | Good Apple |
| Yellow | Smooth | Medium | Faint | Medium | Bad Apple |
| Red | Rough | Small | Strong | Soft | Good Apple |

- (5pts) Please draw this Naive Bayesian network
- (5pts) Please write down the CPTs with add-1 Laplace smoothing.
- (5pts) Given an apple with the features *Red, Rough, Small, Faint, and Soft*, Is this apple a good apple?

(2)(10pts) For each local node in \mathcal{G} , the parameters of CPT can be written as θ_i . If we independently maximize the local likelihood and combine each local optimal solution θ_i^* into a result denoted $\theta^P = [\theta_1^*, \theta_2^*, \dots, \theta_n^*]$. Is the result θ^P equal to the optimal maximum likelihood estimate for the global parameter $\theta^* = \arg \max L(\theta)$? If so, please prove it; if not, please explain why.

(1.a)

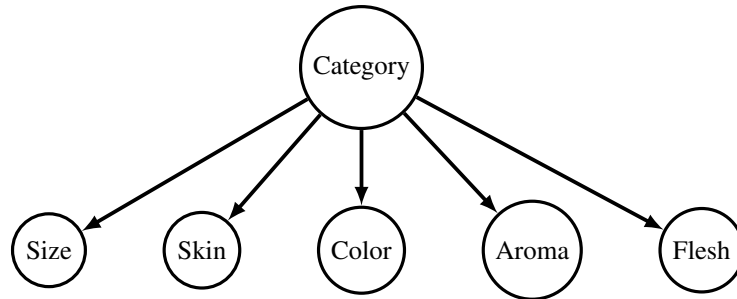


Figure 4: The Naive Bayesian network.

(1.b)

| | | | | | | | |
|---------------------------------|--|---|--|---|--|---|--|
| | | $P(\text{Category} = \text{Good})$ | | $P(\text{Category} = \text{Bad})$ | | | |
| | | 8/14(0.5714) | | 6/14(0.4286) | | | |
| | | $P(\text{Color} = \text{Red} \text{Category})$ | | $P(\text{Color} = \text{Green} \text{Category})$ | | $P(\text{Color} = \text{Yellow} \text{Category})$ | |
| $\text{Category} = \text{Good}$ | | 7/10 | | 2/10 | | 1/10 | |
| $\text{Category} = \text{Bad}$ | | 2/8 | | 3/8 | | 3/8 | |
| | | $P(\text{Skin} = \text{Smooth} \text{Category})$ | | $P(\text{Skin} = \text{Rough} \text{Category})$ | | | |
| $\text{Category} = \text{Good}$ | | 6/9 | | 3/9 | | | |
| $\text{Category} = \text{Bad}$ | | 3/7 | | 4/7 | | | |
| | | $P(\text{Size} = \text{Small} \text{Category})$ | | $P(\text{Size} = \text{Medium} \text{Category})$ | | $P(\text{Size} = \text{Large} \text{Category})$ | |
| $\text{Category} = \text{Good}$ | | 3/10 | | 4/10 | | 3/10 | |
| $\text{Category} = \text{Bad}$ | | 2/8 | | 4/8 | | 2/8 | |
| | | $P(\text{Aroma} = \text{Strong} \text{Category})$ | | $P(\text{Aroma} = \text{Faint} \text{Category})$ | | | |
| $\text{Category} = \text{Good}$ | | 8/9 | | 1/9 | | | |
| $\text{Category} = \text{Bad}$ | | 1/7 | | 6/7 | | | |
| | | $P(\text{Flesh} = \text{Hard} \text{Category})$ | | $P(\text{Flesh} = \text{Medium} \text{Category})$ | | $P(\text{Flesh} = \text{Soft} \text{Category})$ | |
| $\text{Category} = \text{Good}$ | | 6/10 | | 2/10 | | 2/10 | |
| $\text{Category} = \text{Bad}$ | | 1/8 | | 3/8 | | 4/8 | |

(1.c) Bad Apple

(2) Assume that there are n samples in dataset \mathcal{D}

$$\begin{aligned}
 L(\theta) &= \prod_n P_G(\mathbf{x}[n]; \theta) \\
 &= \prod_n \prod_i P(x_i[n] | pa_{x_i}[n]; \theta) \\
 &= \prod_i \left(\prod_n P(x_i[n] | pa_{x_i}[n]; \theta) \right) \\
 &= \prod_i L_i(\theta_i)
 \end{aligned}$$

D Latent Variable Analysis

A latent variable model for T data vectors $\mathbf{x}_0, \dots, \mathbf{x}_T$ is

$$P(\mathbf{x}_0, \dots, \mathbf{x}_T) = P(\mathbf{x}_0) \prod_{t=1}^T P(\mathbf{x}_t | \mathbf{x}_0)$$

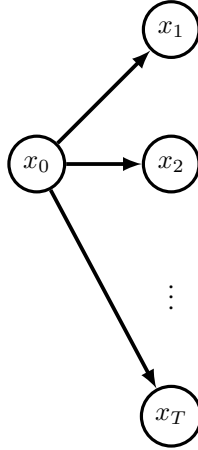
where $P(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_t} \mathbf{x}_0, (1 - \alpha_t) \mathbf{I})$, $\mathcal{N}(\cdot)$ is a Normal distribution, $\alpha_t > 0$ is a scalar parameter, and \mathbf{I} is a identity matrix.

(1)(5pts) Drawing a graphical model to depict the generative process. (Note that each variable \mathbf{x}_i is influenced by its own parameter α_i)

(2)(10pts) If $P(\mathbf{x}_{t-1} | \mathbf{x}_t, x_0)$ follows a Gaussian distribution with variance σ_t^2 , show this distribution.

(3)(10pts) If the condition in (2) is relaxed to a, b , where $0 < a < b < T$, show the distribution of $P(\mathbf{x}_a | \mathbf{x}_b, \mathbf{x}_0)$

(1)



(2) Assume $P(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(m_t \mathbf{x}_t + n_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$, thus

$$\mathbf{x}_{t-1} = m_t \mathbf{x}_t + n_t \mathbf{x}_0 + \sigma_t \mathbf{e}$$

where $\mathbf{e} \sim \mathcal{N}(0, \mathbf{I})$, and there are

$$\begin{aligned}
 \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \mathbf{e}_2 \\
 \mathbf{x}_{t-1} &= \sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1}} \mathbf{e}_3
 \end{aligned}$$

And

$$\begin{aligned}\mathbf{x}_{t-1} &= m_t(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\mathbf{e}_2) + n_t\mathbf{x}_0 + \sigma_t\mathbf{e} \\ &= (m_t\sqrt{\alpha_t} + n_t)\mathbf{x}_0 + m\sqrt{1-\alpha_t}\mathbf{e}_2 + \sigma_t\mathbf{e}\end{aligned}$$

So we can get the equations

$$\begin{cases} m_t\sqrt{\alpha_t} + n_t = \sqrt{\alpha_{t-1}} \\ m_t^2(1-\alpha_t) + \sigma_t^2 = 1 - \alpha_{t-1} \end{cases}$$

So that $m_t = \sqrt{\frac{1-\alpha_{t-1}-\sigma_t^2}{1-\alpha_t}}$ and $n_t = \sqrt{\alpha_{t-1}} - \sqrt{\frac{\alpha_t}{1-\alpha_t}(1-\alpha_{t-1}-\sigma_t^2)}$, and

$$P(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\frac{1-\alpha_{t-1}-\sigma_t^2}{1-\alpha_t}}\mathbf{x}_t + (\sqrt{\alpha_{t-1}} - \sqrt{\frac{\alpha_t}{1-\alpha_t}(1-\alpha_{t-1}-\sigma_t^2)})\mathbf{x}_0, \sigma_t^2\mathbf{I}\right)$$

(3) Similar to (2), we can get

$$P(\mathbf{x}_a|\mathbf{x}_b, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\frac{1-\alpha_a-\sigma_b^2}{1-\alpha_b}}\mathbf{x}_b + (\sqrt{\alpha_a} - \sqrt{\frac{\alpha_b}{1-\alpha_b}(1-\alpha_a-\sigma_b^2)})\mathbf{x}_0, \sigma_b^2\mathbf{I}\right)$$