

CS761 Artificial Intelligence

17. Inference with Bayesian Networks

Jiamou Liu
The University of Auckland

Recall: Probability and Uncertainty

- We introduced **probability theory** as a language that extends from propositional logic to capture uncertainty.
- The probability of a proposition expresses the **degree of belief** (belief measure over a sample space) held about the proposition being true.
- When describing a world with uncertainty, we use a set of **atomic propositions** with a joint probability distribution.
 - Causal factor relation
 - Independence
- Probabilistic inference:
 $\text{Posterior probability} \propto \text{Likelihood} \times \text{Prior Probability}$



Diagnosing a Complex Problem

Example. [coughing?]

- When a person gets **influenza (INF)**, he may develop a **sore throat (SOR)**, a **fever (FEV)**, or maybe **bronchitis (BRO)**.
- When a person is a **smoker (SMO)**, he may also get bronchitis.
- When a person gets bronchitis, he may develop **coughing (COU)** and **wheezing (WEE)**.



Prior knowledge:

- *INF* and *SMO* are independent.
- *SOR*, *FEV*, and *BRO* are directly affected by *INF*, and are conditionally independent given *INF*. (*INF* → *SOR*, *INF* → *FEV*, *INF* → *BRO*)
- *BRO* is also affected by *SMO*. (*SMO* → *BRO*)
- *COU* and *WEE* are directly affected by *BRO*, and are conditionally independent given *BRO*. (*BRO* → *COU*, *BRO* → *WEE*)

For each atom, there is a conditional probability table:

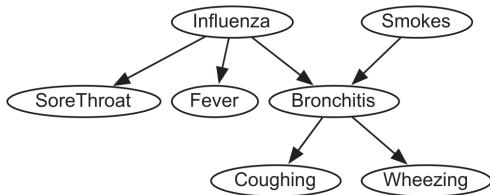
<i>INF</i>	P_{INF}	<i>SMO</i>	P_{SMO}	<i>INF</i>	<i>SOR</i>	P_{SOR}	<i>INF</i>	<i>FEV</i>	P_{FEV}
0	0.95	0	0.8	0	0	0.9999	0	0	0.95
1	0.05	1	0.2	0	1	0.0001	0	1	0.05
				1	0	0.7	1	0	0.1
				1	1	0.3	1	1	0.9

<i>INF</i>	<i>SMO</i>	<i>BRO</i>	P_{BRO}	<i>BRO</i>	<i>COU</i>	P_{COU}	<i>BRO</i>	<i>WEE</i>	P_{WEE}
0	0	0	0.9999	0	0	0.93	0	0	0.999
0	0	1	0.0001	0	1	0.07	0	1	0.05
0	1	0	0.3	1	0	0.2	1	0	0.4
0	1	1	0.7	1	1	0.8	1	1	0.6
1	0	0	0.1						
1	0	1	0.9						
1	1	0	0.01						
1	1	1	0.99						

Question. Suppose the person is coughing. What is the probability that he has influenza?

Probabilistic Model

The problem can be represented by a **directed graph** of the problem domain:



Note.

- The graph is acyclic (i.e. does not contain any directed cycles).
- A node X connects to another node Y by a directed edge if the X is a direct **causal factor** of Y .

Bayesian Networks

Definition [graph terminology]

Let $G = (V, E)$ be a directed graph.

- A **root** of G is any node with no parent;
- $\text{parents}(X) = \{Y \in V \mid (Y, X) \in E\}$ is the set of **parents** of X .
- a **leaf** is any node that is not a parent of any others; the other nodes are called the **intermediate nodes**;
- A node X is a **descendent** of Y if there is a directed path from Y to X .

Bayesian Networks

Definition [graph terminology]

Let $G = (V, E)$ be a directed graph.

- A **root** of G is any node with no parent;
- $\text{parents}(X) = \{Y \in V \mid (Y, X) \in E\}$ is the set of **parents** of X .
- a **leaf** is any node that is not a parent of any others; the other nodes are called the **intermediate nodes**;
- A node X is a **descendent** of Y if there is a directed path from Y to X .

Definition. [Bayesian network]

A **Bayesian network** is a tuple $(V, E, (\mathbf{P}_X)_{X \in V})$ where

- (V, E) forms a directed acyclic graph; each node $X \in V$ represents an **atomic proposition**.
- For every node X , \mathbf{P}_X is a **probability distribution of X conditioning on $\text{parents}(X)$** , i.e., $\mathbf{P}(X \mid \text{parents}(X))$.
- **Local Markov property**: Each node is independent from its non-descendants conditioned on its parents.

Example. [coughing?] A Bayesian network: $(V, E, (\mathbf{P}_X)_{X \in V})$ where

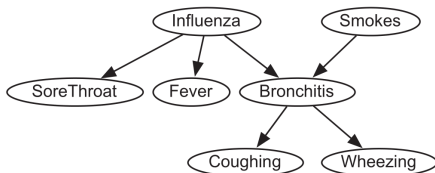
- V contains atoms $INF, SMO, SOR, FEV, BRO, COU, WEE$
- (Syntax) E represents all the dependencies between atoms:

$INF \rightarrow SOR, INF \rightarrow FEV, INF \rightarrow BRO$

$SMO \rightarrow BRO, BRO \rightarrow COU, BRO \rightarrow WEE$

- (Semantics) The CPT given earlier.

Note. Local Markov property ensures all observed conditional independence.



Bayesian Network Inference

For any atom X , we use ℓ_X to denote a literal, i.e. X or $\neg X$.

Bayesian Network Inference Problem

INPUT Given a Bayesian network $(V, E, (\mathbf{P}_X)_{X \in V})$, two nodes $X, Y \in V$, and literal $\ell_Y \in \{Y, \neg Y\}$,

OUTPUT compute the probability $P(X \mid \ell_Y)$ and $P(\neg X \mid \ell_Y)$.

E.g. Given the earlier example Bayesian network, query for $P(INF \mid COU)$.

We now present methods to solve this problem.

Tools:

- ① CPT operations
- ② Chain rule

CPT Operations

CPT Operation 1: Restriction

For a conditional probability table $f(Y_1, \dots, Y_k, X)$, for $\ell_X \in \{X, \neg X\}$, the ℓ_X -restriction operation produces a table $g(Y_1, \dots, Y_k)$ such that for any $\ell_{Y_i} \in \{Y_i, \neg Y_i\}$,

$$g(\ell_{Y_1}, \dots, \ell_{Y_k}) = f(\ell_{Y_1}, \dots, \ell_{Y_k}, \ell_X).$$

E.g.

Y	X	$f(Y, X)$		Y	$g(Y)$
0	0	0.99	\Rightarrow	0	0.01
0	1	0.01		1	0.95
1	0	0.05			
1	1	0.95			

CPT Operation 2: Multiplication

For CPT $f(Y_1, \dots, Y_k, X)$ and $g(X, Z_1, \dots, Z_m)$, the **multiplication** operation produces a new CPT $h(Y_1, \dots, Y_k, X, Z_1, \dots, Z_m)$ such that for any literals $\ell_{Y_i} \in \{Y_i, \neg Y_i\}$, $\ell_X \in \{X, \neg X\}$, $\ell_{Z_j} \in \{Z_j, \neg Z_j\}$,

$$h(\ell_{Y_1}, \dots, \ell_{Y_k}, \ell_X, \ell_{Z_1}, \dots, \ell_{Z_m}) = f(\ell_{Y_1}, \dots, \ell_{Y_k}, \ell_X) \times g(\ell_X, \ell_{Z_1}, \dots, \ell_{Z_m}).$$

E.g.

A	B	f	\times	B	C	g	$=$	A	B	C	h
0	0	0.1		0	0	0.2		0	0	0	0.02
0	1	0.9		0	1	0.8		0	0	1	0.72
1	0	0.4		1	0	0.3		0	1	0	0.27
1	1	0.6		1	1	0.7		0	1	1	0.63
								1	0	0	0.08
								1	0	1	0.32
								1	1	0	0.18
								1	1	1	0.42

CPT Operation 3: Normalisation

For a conditional probability table $f(Y_1, \dots, Y_k, X)$, the **X-normalisation** operation produces a new CPT $g(Y_1, \dots, Y_k, X)$ such that for any $\ell_{Y_i} \in \{Y_i, \neg Y_i\}$, $\ell_X \in \{X, \neg X\}$,

$$g(\ell_{Y_1}, \dots, \ell_{Y_k}, \ell_X) = \frac{f(\ell_{Y_1}, \dots, \ell_{Y_k}, \ell_X)}{f(\ell_{Y_1}, \dots, \ell_{Y_k}, X) + f(\ell_{Y_1}, \dots, \ell_{Y_k}, \neg X)}.$$

E.g.

X	$f(X)$	\Rightarrow	X	$g(X)$
0	0.0096		0	0.2017
1	0.038		1	0.7983

CPT Operation 4: Summation

For a conditional probability table $f(Y_1, \dots, Y_k, X)$, the **X-sum** operation produces a new CPT $g(Y_1, \dots, Y_k)$ such that for any $\ell_{Y_i} \in \{Y_i, \neg Y_i\}$,

$$g(\ell_{Y_1}, \dots, \ell_{Y_k}) = f(\ell_{Y_1}, \dots, \ell_{Y_k}, X) + f(\ell_{Y_1}, \dots, \ell_{Y_k}, \neg X).$$

E.g. C-sum of the following table:

A	B	C	f
0	0	0	0.02
0	0	1	0.72
0	1	0	0.27
0	1	1	0.63
1	0	0	0.08
1	0	1	0.32
1	1	0	0.18
1	1	1	0.42

 \Rightarrow

A	B	g
0	0	0.74
0	1	0.9
1	0	0.4
1	1	0.6

Example. [simple diagnosis] To query $\mathbf{P}(Inf|Test)$:

Inf	\mathbf{P}_{Inf}
0	0.96
1	0.04

and

Inf	$Test$	$\mathbf{P}_{Test}(Inf)$
0	0	0.99
0	1	0.01
1	0	0.05
1	1	0.95

Note. $\mathbf{P}(Inf|Test) \propto \mathbf{P}(Inf) \times \mathbf{P}(Test | Inf)$.

Example. [simple diagnosis] To query $P(\text{Inf}|\text{Test})$:

Inf	P_{Inf}	and	Inf	Test	$P_{\text{Test}}(\text{Inf})$
0	0.96		0	0	0.99
1	0.04		0	1	0.01
			1	0	0.05
			1	1	0.95

Note. $P(\text{Inf}|\text{Test}) \propto P(\text{Inf}) \times P(\text{Test} | \text{Inf})$.

Step 1 Restriction. In P_{Test} keep only those rows where *Test* is true.

Inf	$f(\text{Inf})$
0	0.01
1	0.95

Step 2 Multiplication.

Inf	P_{Inf}	\times	Inf	$f(\text{Inf})$	$=$	Inf	$g(\text{Inf})$
0	0.96		0	0.01		0	$0.01 \times 0.96 = 0.0096$
1	0.04		1	0.95		1	$0.95 \times 0.04 = 0.038$

Step 3. Normalisation.

Inf	$g(\text{Inf})$
0	$0.0096 / (0.0096 + 0.038) = 0.0096 / 0.0476 \approx 0.2017$
1	$0.038 / (0.0096 + 0.038) = 0.038 / 0.0476 \approx 0.7983$

Therefore $P(\text{Inf} | \text{Test}) = 0.7983$ and $P(\neg \text{Inf} | \text{Test}) = 0.2017$.

Chain Rule

- For any two propositions a_1, a_2

$$P(a_1 \wedge a_2) = P(a_2 \mid a_1)P(a_1).$$

- For any three propositions a_1, a_2, a_3

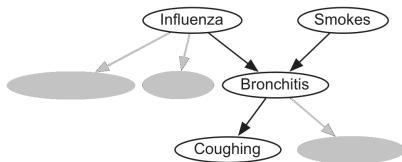
$$P(a_1 \wedge a_2 \wedge a_3) = P(a_1 \wedge a_2 \mid a_3)P(a_3) = P(a_1 \mid a_2 \wedge a_3)P(a_2 \mid a_3)P(a_3).$$

Chain Rule

For any n propositions a_1, a_2, \dots, a_n

$$\begin{aligned} P\left(\bigwedge_{i=1}^n a_i\right) &= P(a_1 \mid a_2 \wedge \dots \wedge a_n)P(a_2 \mid a_3 \wedge \dots \wedge a_n) \dots P(a_{n-1} \mid a_n)P(a_n) \\ &= \prod_{i=1}^n P\left(a_i \mid \bigwedge_{j=i+1}^n a_j\right) \end{aligned}$$

Example. [coughing?]



$$\begin{aligned} & \mathbf{P}(INF, SMO, BRO, COU) \\ &= \mathbf{P}(COU \mid BRO, SMO, INF) \mathbf{P}(BRO \mid SMO, INF) \mathbf{P}(SMO \mid INF) \mathbf{P}(INF) \\ &= \mathbf{P}(COU \mid BRO) \mathbf{P}(BRO \mid SMO, INF) \mathbf{P}(SMO) \mathbf{P}(INF) \quad (\text{by local Markov property}) \end{aligned}$$

Theorem [chain rule on a Bayesian network]

Let $(V, E, (\mathbf{P}_X)_{X \in V})$ be a Bayesian network. Then for any subset $S \subseteq V$

$$\mathbf{P}(S) = \prod_{X \in S} \mathbf{P}(X \mid \text{parents}(X)) = \prod_{X \in S} \mathbf{P}_X$$

Example. [coughing?] What is $P(INF | COU)$?

Idea:

$$\begin{aligned}
 P(INF | COU) &= P(COU | INF)P(INF) = P(COU, INF) \\
 &= \sum_{BRO} \sum_{SMO} P(INF, COU, BRO, SMO) \\
 &= \sum_{BRO} \sum_{SMO} P_{COU} P_{BRO} P_{SMO} P_{INF}
 \end{aligned}$$

Starting from the set of CPT: P_{COU} , P_{BRO} , P_{SMO} , P_{INF}

INF	SMO	BRO	P_{BRO}
0	0	0	0.9999
0	0	1	0.0001
0	1	0	0.3
0	1	1	0.7
1	0	0	0.1
1	0	1	0.9
1	1	0	0.01
1	1	1	0.99

INF	P_{INF}
0	0.95
1	0.05

SMO	P_{SMO}
0	0.8
1	0.2

BRO	COU	P_{COU}
0	0	0.93
0	1	0.07
1	0	0.2
1	1	0.8

Question. How to compute $\sum_{BRO} \sum_{SMO} \mathbf{P}_{COU} \mathbf{P}_{BRO} \mathbf{P}_{SMO} \mathbf{P}_{INF}$?

Naive method.

- ① **Step 1: Restriction.** In \mathbf{P}_{COU} , keep only those rows where COU is true.
Obtain a new CPT $f(BRO)$.
- ② **Step 2: Multiplication.** Compute $f(BRO) \times \mathbf{P}_{BRO} \times \mathbf{P}_{SMO} \times \mathbf{P}_{INF}$
Obtain a new CPT $g(BRO, SMO, INF)$
- ③ **Step 3: Sum (to eliminate variables).** Perform SMO -sum and then BRO -sum, on $g(BRO, SMO, INF)$.
Obtain a new CPT $h(INF)$
- ④ **Step 3: Normalisation.** Extract the conditional probabilities $P(INF \mid COU)$ and $P(\neg INF \mid COU)$.

Complexity. Suppose a CPT has r rows. The algorithm runs in $O(r^n)$.

We now present how to improve the running time of the algorithm.

Variable Elimination

Idea: Instead of multiply all tables and eliminate variables all at once, eliminate variables **one-by-one**.

- ① **Step 1: Restriction.** In \mathbf{P}_{COU} , keep only those rows where COU is true. Obtain a new CPT $f(BRO)$.
- ② **Step 2: Eliminate SMO .**
 - **Step 2.1:** Multiply $\mathbf{P}_{BRO} \times \mathbf{P}_{SMO}$ to get a new table $g(INF, SMO, BRO)$
 - **Step 2.2:** SMO -sum $g(INF, SMO, BRO)$ to get a new table $h(INF, BRO)$.
- ③ **Step 3: Eliminate BRO .**
 - **Step 3.1:** Multiply $f(BRO) \times h(INF, BRO)$ to get a new table $\ell(INF, BRO)$
 - **Step 3.2:** BRO -sum $\ell(INF, BRO)$ to get a new table $r(INF)$.
- ④ **Step 4: Multiplication.** $r(INF) \times \mathbf{P}(INF)$ to a table $s(INF)$.
- ⑤ **Step 5: Normalisation.** Retrieve $\mathbf{P}(INF \mid COU)$.

Illustration.

Step 1: Restriction. In P_{COU} , keep only those rows where COU is true.

BRO	COU	P_{COU}	\Rightarrow	BRO	$f(BRO)$
0	0	0.93		0	0.07
0	1	0.07		1	0.8
1	0	0.2			
1	1	0.8			

Step 2: Eliminate SMO.

2.1	INF	SMO	BRO	P_{BRO}	\times	SMO	P_{SMO}	$=$	INF	SMO	BRO	$g(INF, SMO, BRO)$
	0	0	0	0.9999		0	0.8		0	0	0	0.79992
	0	0	1	0.0001		1	0.2		0	0	1	0.00008
	0	1	0	0.3					0	1	0	0.06
	0	1	1	0.7					0	1	1	0.14
	1	0	0	0.1					1	0	0	0.08
	1	0	1	0.9					1	0	1	0.72
	1	1	0	0.01					1	1	0	0.002
	1	1	1	0.99					1	1	1	0.198

2.2 SMO-Sum	INF	SMO	BRO	$g(INF, SMO, BRO)$	\Rightarrow	INF	BRO	$h(INF, BRO)$
	0	0	0	0.79992		0	0	0.85992
	0	0	1	0.00008		0	1	0.14008
	0	1	0	0.06		1	0	0.082
	0	1	1	0.14		1	1	0.918
	1	0	0	0.08				
	1	0	1	0.72				
	1	1	0	0.002				
	1	1	1	0.198				

Step 3: Eliminate BRO.

$$3.1 \quad \begin{array}{|c|c|} \hline \text{BRO} & \mathbf{P}_{\text{COU}} \\ \hline 0 & 0.07 \\ 1 & 0.8 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline \text{INF} & \text{BRO} & h(\text{INF}, \text{BRO}) \\ \hline 0 & 0 & 0.85992 \\ 0 & 1 & 0.14008 \\ 1 & 0 & 0.082 \\ 1 & 1 & 0.918 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline \text{INF} & \text{BRO} & \ell(\text{INF}, \text{BRO}) \\ \hline 0 & 0 & 0.0601944 \\ 0 & 1 & 0.112064 \\ 1 & 0 & 0.00574 \\ 1 & 1 & 0.7344 \\ \hline \end{array}$$

$$3.2 \text{ BRO-sum} \quad \begin{array}{|c|c|c|} \hline \text{INF} & \text{BRO} & \ell(\text{INF}, \text{BRO}) \\ \hline 0 & 0 & 0.0601944 \\ 0 & 1 & 0.112064 \\ 1 & 0 & 0.00574 \\ 1 & 1 & 0.7344 \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline \text{INF} & r(\text{INF}) \\ \hline 0 & 0.1722584 \\ 1 & 0.74014 \\ \hline \end{array}$$

Step 4: Multiplication.

$$\begin{array}{|c|c|} \hline \text{INF} & r(\text{INF}) \\ \hline 0 & 0.1722584 \\ 1 & 0.74014 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline \text{INF} & \mathbf{P}_{\text{INF}} \\ \hline 0 & 0.95 \\ 1 & 0.05 \\ \hline \end{array} = \begin{array}{|c|c|} \hline & s(\text{INF}) \\ \hline 0 & 0.16364548 \\ 1 & 0.037007 \\ \hline \end{array}$$

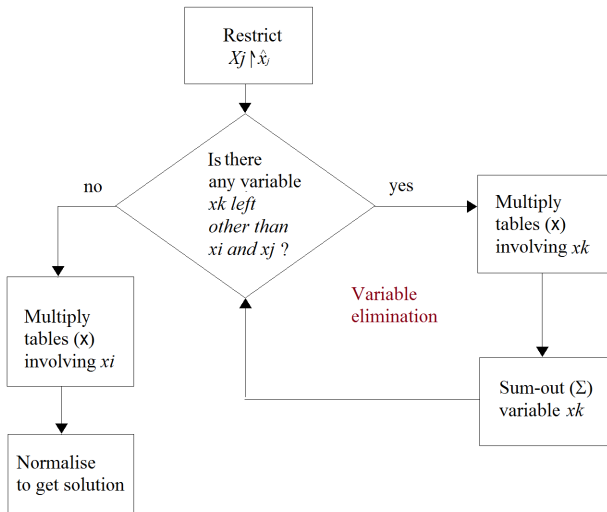
Step 5: Normalisation.

$$\begin{array}{|c|c|} \hline & s(\text{INF}) \\ \hline 0 & 0.16364548 \\ 1 & 0.037007 \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline & \text{result} \\ \hline P(\neg \text{INF} \wedge \text{COU}) & 0.8156 \\ P(\text{INF} \wedge \text{COU}) & 0.1844 \\ \hline \end{array}$$

The probability of a person getting influenza when (only) coughing is observed is **18.44%**!

Variable Elimination Algorithm

Solving the Bayesian network inference problem using the **variable elimination algorithm (VE)**



Variable Elimination (Rough steps)

INPUT:

- CPT P_X for each $X \in V$.
- Query variables X_i, X_j ; Observed value for X_j .

OUTPUT: $P(X_i | X_j)$

- ① Restrict tables that contain X_j to rows where X_j has the observed value.
- ② Eliminate each of the non-query variables X_k (in certain order):
 - Multiply all tables containing X_k .
 - Sum-out X_k from the product table.
- ③ Multiply the remaining tables to get a table over a single variable X_i .
- ④ Normalise and return the result.

Complexity: If every atom appears in at most 2 CPT and one of them has at most 4 rows, then the algorithm runs in $O(n)$.

In the general case, the running time depends on the ordering in which we eliminate the variables, but finding the optimal ordering is NP-complete.

Summary of The Topic

The following are the main knowledge points covered:

- **Bayesian network:** $(V, E, (\mathbf{P}_X)_{X \in V})$ a graph structure denoting causal relations between atoms, CPTs
- **Local markov property**
- **Bayesian network inference:** Given a Bayesian network $(V, E, (\mathbf{P}_X)_{X \in V})$, nodes X, Y and observation $\ell_Y \in \{Y, \neg Y\}$, compute $\mathbf{P}(X \mid \ell_Y)$.
- **CPT Operations:**
 - Restriction
 - Multiplication
 - Normalisation
 - Summation
- **Chain rule:** $\pi(S) = \prod_{X \in S} \pi_X$
- **Naive method** for Bayesian network inference
- **Variable elimination**