Introduction to Belief Networks¹

Salvador Ruiz Correa

Centro de Investigación en Matemáticas (CIMAT)

э.

¹These slides are *adapted* from those that accompany the book *Bayesian Reasoning and Machine Learning*. The book and demos can be downloaded from www.cs.ucl.ac.uk/staff/D.Barber/brml. We acknowledge David Barber for providing the original slides.

Bayesian network structure

A Bayesian network structure \mathcal{G} is a directed acyclic graph whose nodes represent variables X_1, \ldots, X_n . The structure encodes the following set of conditional independence assumptions, called local independences, and denoted by $\mathcal{I}_l(\mathcal{G})$:

For each variable X_i : $X_i \perp \mathsf{ND}(X_i) | \mathsf{Pa}(X_i)$,

where $ND(X_i)$ and $Pa(X_i)$ denote the non-descendants and the parents of X_i , respectively.

Independencies in P

Let P be a distribution over X_1, \ldots, X_n . $\mathcal{I}(P)$ denotes the set of independence assertions of the form $X \perp \!\!\!\perp Y \mid Z$ that hold in P.

I-Map

Let \mathcal{G} be a graph associated with a set of induced pendences $\mathcal{I}(\mathcal{G})$. \mathcal{G} is an I-map for a set of dependencies \mathcal{I} if $I(\mathcal{G}) \subseteq \mathcal{I}$. \mathcal{G} is win I-map for P if $I(\mathcal{G})$ is an I-map for P. If $I(\mathcal{G}) \subseteq \mathcal{I}(P)$, any independence that \mathcal{G} asserts must hold in P. P may have additional independences that are not reflected in \mathcal{G} .

CPD parametrization



 $p(A, B, C) = p(B \mid A, B)p(A)p(B)$

C	A	B	p(C A,B)
0	0	0	$ heta_1$
0	0	1	θ_2
0	1	0	$ heta_3$
0	1	1	$ heta_4$
1	0	0	$1- heta_1$
1	0	1	$1- heta_2$
1	1	0	$1- heta_3$
1	1	1	$1- heta_4$

A	p(A)
0	$ heta_5$
1	$1-\theta_5$

B	p(B)
0	θ_6
1	$1-\theta_6$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Factorization (chain rule)

Let \mathcal{G} be a Bayesian network structure over the variables X_1, \ldots, X_n . A distribution P over the same variables factorizes according to \mathcal{G} if P can be expressed as a product:

 $P(X_1,\ldots,X_n) = \prod_{i=1}^n P(X_i \mid \mathsf{Pa}(X_i)).$

This equation is called chain rule for Bayesian networks. The terms $P(X_i | Pa(X_i))$ are called conditional probability distributions (CPDs), or local probability models.

Bayesian Network

A Bayesian network is a pair (\mathcal{G}, P) where P factorizes over \mathcal{G} and where P is specified as a set of CPDs associated with \mathcal{G} 's nodes.

$\text{I-Map} \Rightarrow \text{factorization}$

Let \mathcal{G} be a Bayesian network structure over X_1, \ldots, X_n and let P be a joint distribution over the same variables. If \mathcal{G} is an I-map for P, then P factorizes over \mathcal{G} .

$\mathsf{Factorization} \Rightarrow \mathsf{I}\text{-}\mathsf{Map}$

Let \mathcal{G} be a Bayesian network structure over X_1, \ldots, X_n and let P be a joint distribution over the same variables. If P factorizes over \mathcal{G} , then \mathcal{G} is an I-map for P.

Completness

For all distributions P that factorize over \mathcal{G} ,

$$\mathcal{I}(P) \stackrel{a.e.}{=} \mathcal{I}(\mathcal{G});$$

that is, for all distributions except for a set of measure zero in the space of CPD parametrizations, we have that $\mathcal{I}(P) = \mathcal{I}(\mathcal{G})$:

Belief Networks (Bayesian Networks)

A belief network is a directed acyclic graph in which each node has associated the conditional probability of the node given its parents.

The joint distribution is obtained by taking the product of the conditional probabilities:

p(A,B,C,D,E) = p(A)p(B)p(C|A,B)p(D|C)p(E|B,C)



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Example – Part I

Sally's burglar **A**larm is sounding. Has she been **B**urgled, or was the alarm triggered by an **E**arthquake? She turns the car **R**adio on for news of earthquakes.

Choosing an ordering Without loss of generality, we can write

$$\begin{split} p(A,R,E,B) &= p(A|R,E,B)p(R,E,B) \\ &= p(A|R,E,B)p(R|E,B)p(E,B) \\ &= p(A|R,E,B)p(R|E,B)p(E|B)p(B) \end{split}$$

Assumptions:

- The alarm is not directly influenced by any report on the radio, p(A|R,E,B) = p(A|E,B)
- The radio broadcast is not directly influenced by the burglar variable, p(R|E,B) = p(R|E)

 $\bullet\,$ Burglaries don't directly 'cause' earthquakes, p(E|B)=p(E)

Therefore

p(A,R,E,B) = p(A|E,B)p(R|E)p(E)p(B)

Example – Part II: Specifying the Tables



m(D|F)

p(A|B, E)

Alarm $= 1$	Burglar	Earthquake		$p(\mathbf{n})$	
0.9999	1	1		Radio = 1	Earthquake
0.99	1	0		1	1
0.99	0	1	1	0	0
0.0001	0	0			

The remaining tables are p(B = 1) = 0.01 and p(E = 1) = 0.000001. The tables and graphical structure fully specify the distribution.

Example Part III: Inference

Initial Evidence: The alarm is sounding

$$p(B = 1|A = 1) = \frac{\sum_{E,R} p(B = 1, E, A = 1, R)}{\sum_{B,E,R} p(B, E, A = 1, R)}$$
$$= \frac{\sum_{E,R} p(A = 1|B = 1, E)p(B = 1)p(E)p(R|E)}{\sum_{B,E,R} p(A = 1|B, E)p(B)p(E)p(R|E)} \approx 0.99$$

Additional Evidence: The radio broadcasts an earthquake warning:

A similar calculation gives $p(B = 1|A = 1, R = 1) \approx 0.01$.

Initially, because the alarm sounds, Sally thinks that she's been burgled. However, this probability drops dramatically when she hears that there has been an earthquake.

The earthquake 'explains away' to an extent the fact that the alarm is ringing.

Uncertain Evidence

In soft/uncertain evidence the variable is in more than one state, with the strength of our belief about each state being given by probabilities. For example, if y has the states dom $(y) = \{$ red, blue, green $\}$ the vector (0.6, 0.1, 0.3) could represent the probabilities of the respective states.

hard evidence

We are certain that a variable is in a particular state. In this case, all the probability mass is in one of the vector components, (0, 0, 1).

inference

Inference with soft-evidence can be achieved using Bayes' rule. Writing the soft evidence as \tilde{y} , we have

$$p(x|\tilde{y}) = \sum_{y} p(x|y) p(y|\tilde{y})$$

where $p(y = \mathbf{i}|\tilde{y})$ represents the probability that y is in state i under the soft-evidence.

Jeffrey's rule

For variables x, y, and $p_1(x, y)$, how do we form a joint distribution given soft-evidence \tilde{y} ?

Form the conditional We first define

$$p_1(x|y) = \frac{p_1(x,y)}{\sum_x p_1(x,y)}$$

Define the joint

The soft evidence $p(y|\tilde{y})$ then defines a new joint distribution

 $p_2(x, y|\tilde{y}) = p_1(x|y)p(y|\tilde{y})$

One can therefore view soft evidence as defining a new joint distribution. We use a dashed circle to represent a variable in an uncertain state.

Uncertain evidence example



Revisiting the earthquake scenario, we think we hear the burglar alarm sounding, but are not sure, specifically p(A = tr) = 0.7. For this binary variable case we represent this soft-evidence for the states (tr, fa) as $\tilde{A} = (0.7, 0.3)$. What is the probability of a burglary under this soft-evidence?

$$\begin{split} p(B = \mathsf{tr} | \tilde{A}) &= \sum_{A} p(B = \mathsf{tr} | A) p(A | \tilde{A}) \\ &= p(B = \mathsf{tr} | A = \mathsf{tr}) \times 0.7 + p(B = \mathsf{tr} | A = \mathsf{fa}) \times 0.3 \approx 0.6930 \end{split}$$

This value is lower than 0.99, the probability of being burgled when we are sure we heard the alarm. The probabilities p(B = tr|A = tr) and p(B = tr|A = fa) are calculated using Bayes' rule from the original distribution, as before.

Unreliable evidence (likelihood evidence)

Under potentially confusing reports, you decide to replace the influence of the radio variable with your own model. You decide that you want the radio evidence to influence the inference 80% towards being an earthquake and 20% to not being an earthquake.

$$p(R|E) \rightarrow p(\mathsf{R}|E) = \begin{cases} 0.8 & E = \mathsf{tr} \\ 0.2 & E = \mathsf{fa} \end{cases}$$

This then gives a distribution, with R in an arbitrary fixed state,

 $p(B, E, A, \mathsf{R}) = p(A|B, E)p(B)p(E)p(\mathsf{R}|E)$

This can then be used to form inference.

Examples of Belief Networks in Machine Learning

```
Prediction (discriminative)
p(class|input)
```

Prediction (generative) $p(class|input) \propto p(input|class)p(class)$

Time-series Markov chains, Hidden Markov Models.

Unsupervised learning $p(data) = \sum_{latent} p(data|latent)p(latent).$

And many more

Personally I find the framework very useful for understanding and rationalising the many different approaches in machine learning and related areas.

Independence ⊥⊥ in Belief Networks – Part I

All belief networks with three nodes and two links:



• In (a), (b) and (c), A, B are conditionally independent given C.

(a) $p(A, B|C) = \frac{p(A,B,C)}{p(C)} = \frac{p(A|C)p(B|C)p(C)}{p(C)} = p(A|C)p(B|C)$ (b) $p(A, B|C) = \frac{p(A)p(C|A)p(B|C)}{p(C)} = \frac{p(A,C)p(B|C)}{p(C)} = p(A|C)p(B|C)$ (c) $p(A, B|C) = \frac{p(A|C)p(C|B)p(B)}{p(C)} = \frac{p(A|C)p(B,C)}{p(C)} = p(A|C)p(B|C)$

• In (d) the variables A, B are conditionally dependent given C, $p(A, B|C) \propto p(C|A, B)p(A)p(B)$.

Independence ⊥⊥ in Belief Networks – Part II



- In (a), (b) and (c), the variables A, B are marginally dependent.
- In (d) the variables A, B are marginally independent.

 $p(A,B) = \sum_C p(A,B,C) = \sum_C p(A)p(B)p(C|A,B) = p(A)p(B)$

Collider

A collider contains two or more incoming arrows along a chosen path. Summary of two previous slides:



If C has more than one incoming link, then $A \perp\!\!\!\perp B$ and $A \not\!\!\perp B \mid C$. In this case C is called collider.



If C has at most one incoming link, then $A \perp B \mid C$ and $A \not\perp B$. In this case C is called non-collider.

The 'connection'-graph

All paths in the connection graph need to be blocked to obtain $\bot\!\!\!\bot$:



General Rule for Independence in Belief Networks

Given three sets of nodes $\mathcal{X}, \mathcal{Y}, \mathcal{C}$, if all paths from any element of \mathcal{X} to any element of \mathcal{Y} are blocked by \mathcal{C} , then \mathcal{X} and \mathcal{Y} are conditionally independent given \mathcal{C} .

A path \mathcal{P} is blocked by \mathcal{C} if at least one of the following conditions is satisfied:

- 1. there is a collider \mathcal{P} in the path \mathcal{P} such that neither the collider nor any of its descendants is in the conditioning set \mathcal{C} .
- 2. there is a non-collider in the path \mathcal{P} that is in the conditioning set \mathcal{C} .

d-connected/separated

We use the phrase 'd-connected' if there is a path from \mathcal{X} to \mathcal{Y} in the 'connection' graph – otherwise the variable sets are 'd-separated'. Note that d-separation implies that $\mathcal{X} \perp \!\!\!\perp \! \mathcal{Y} | \mathcal{Z}$, but d-connection does not necessarily imply conditional dependence.

Markov Equivalence

skeleton

Formed from a graph by removing the arrows

immorality

An immorality in a DAG is a configuration of three nodes, A,B,C such that C is a child of both A and B, with A and B not directly connected.

Markov equivalence

Two graphs represent the same set of independence assumptions if and only if they have the same skeleton and the same set of immoralities.



Limitations of expressibility



$$\begin{split} p(t_1,t_2,y_1,y_2,h) &= p(t_1)p(t_2)p(y_1|t_1,h)p(y_2|t_2,h) \\ t_1 \amalg t_2,y_2, \qquad t_2 \amalg t_1,y_1 \end{split}$$

$$p(t_1, t_2, y_1, y_2) = p(t_1)p(t_2)\sum_h p(y_1|t_1, h)p(y_2|t_2, h)$$



Still holds that:

 $t_1 \perp\!\!\!\perp t_2, y_2, \qquad t_2 \perp\!\!\!\perp t_1, y_1$

No Belief network on t_1, t_2, y_1, y_2 can represent all the conditional independence statements contained in $p(t_1, t_2, y_1, y_2)$. Sometimes we can extend the representation by adding for example a bidirectional link, but this is no longer a Belief Network.

Causality

Males	Recovered	Not Recovered	Rec. Rate
Given Drug	18	12	60%
Not Given Drug	7	3	70%

Females	Recovered	Not Recovered	Rec. Rate
Given Drug	2	8	20%
Not Given Drug	9	21	30%

Combined	Recovered	Not Recovered	Rec. Rate
Given Drug	20	20	50%
Not Given Drug	16	24	40%

Simpson's paradox

For the males, it's best not to give the drug. For the females, it's also best not to give the drug. However, for the combined data, it's best to give the drug!

Resolving the paradox

We can write the distribution as



observational calculation

p(G, D, R) = p(R|G, D)p(D|G)p(G)

Our observational calculation computed p(R|G,D) and p(R|D) using the above distribution.

Sampling from the distribution

The above formula suggests that we would first chose a gender (the term p(G)) then decide whether or not to give the drug (the term p(D|G)).

Resolving the paradox

interventional calculation

We must use a distribution that is consistent with an interventional experiment. In this case, the term p(D|G) should play no role. That is, we need to consider a modified distribution (conditioned on the drug)

R

$$\tilde{p}(G, R|D) = p(R|G, D)p(G)$$

$$p(R||D) = \sum_{G} \tilde{p}(G, R|D) = \sum_{G} p(R|G, D)p(G)$$

This gives the non-paradoxical result:

 $p(\text{recovery}|\text{drug}) = 0.6 \times 0.5 + 0.2 \times 0.5 = 0.4$ $p(\text{recovery}|\text{no drug}) = 0.7 \times 0.5 + 0.3 \times 0.5 = 0.5$

The moral of the story is that you have to make the distribution match the experimental conditions, otherwise apparent paradoxes may arise.