## Lecture 8 - Reasoning under Uncertainty (Part I)

Jesse Hoey
School of Computer Science
University of Waterloo

June 2, 2022

Readings: Poole & Mackworth (2nd ed.)Chapt. 8 up to 8.4

## Uncertainty

Why is uncertainty important?

- Agents (and humans) don't know everything ,
- but need to make decisions anyways!
- Decisions are made in the absence of information ,
- or in the presence of noisy information (sensor readings)

The best an agent can do:
know how uncertain it is, and act accordingly

## Probability: Frequentist vs. Bayesian



Frequentist view:
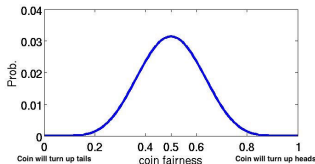probability of heads = # of heads / # of flips
probability of heads this time = probability of heads (history)
Uncertainty is ontological : pertaining to the world

Bayesian view:
probability of heads this time = agent's belief about flip
belief of agent A : based on previous experience of agent A
Uncertainty is epistemological : pertaining to knowledge

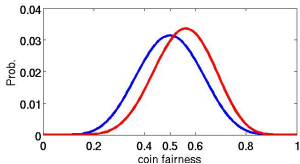## Probability: Bayesian
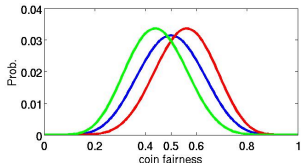
Bayesian probability
all else being equal (Prior)
before any flips

Bayesian probability
all else being equal (Prior)
  after 2 flips heads, heads (Posterior)

Bayesian probability
all else being equal (Prior)
  after 2 flips tails,tails (Posterior)

Should you wear your seatbelt ?
estimate $P(injury)$ given you do/don't wear it

Should you wear your seatbelt ?
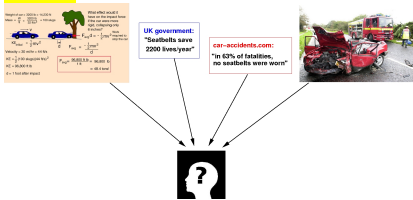
estimate $P(injury)$ given you do/don't wear it

Frequentist:

| test | day | result | P(fatality) |
|---|---|---|---|
| - | Sunday (prior to start) | - | ? |
| 1 | Monday | | 0.0 |
| 2 | Tuesday | | 0.0 |
| 3 | Tuesday | | 0.33333 |
| 4 | Thursday | | 0.25 |
| 5 | Friday | | 0.2 |
| ... | ... | ... | ... |
| N | | | Number of injuries / N |

Should you wear your seatbelt ?

estimate $P(injury)$ given you do/don't wear it

Bayesian:

UK government: "Seatbelts save 2200 lives/year"

car-accidents.com: "in 63% of fatalities, no seatbelts were worn"

if $X$ is a random variable (feature, attribute),
it can take on values $x$, where $x \in Domain(X)$ (or $Dom(X)$)
Assume $x$ is discrete

$\mathbf{P}(\mathbf{x})$ is the probability that $X = x$

joint probability $\mathbf{P}(\mathbf{x}, \mathbf{y})$ is the
probability that $X = x$ and $Y = y$ at the same time

Joint probability distribution:
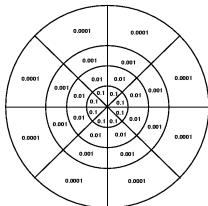


Where is the robot?
features: X,Y

Axioms are things we have to assume about probability:

- $P(X) \geq 0$
- $\sum_x P(X = x) = 1.0$
- $P(a \vee b) = P(a) + P(b)$ if $a$ and $b$ are contradictory - can't both be true at the same time e.g.
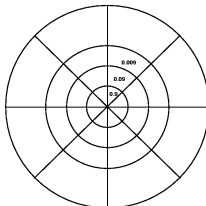  $P(win \vee lose) = P(win) + P(lose) = 1.0$

Some notes:

- probability between 0-1 is purely convention
- $P(a) = 0$ means you think a is definitely false
- $P(a) = 1$ means you think a is definitely true
- $0 < P(a) < 1$ means you have belief about the truth of $a$. It does not mean that $a$ is true to some degree, just that you are ignorant of its truth value.
- Probability = measure of ignorance

- describe a system with $n$ features: $2^n - 1$ probabilities
- Use independence to reduce number of probabilities
- e.g. radially symmetric dartboard, P(hit a sector)
- $P(sector) = P(r, \theta)$ where $r = 1, \ldots, 4$ and $\theta = 1, \ldots, 8$.
- 32 sectors in total - need to give 31 numbers

- describe a system with $n$ features: $2^n - 1$ probabilities
- Use independence to reduce number of probabilities
- e.g. radially symmetric dartboard, P(hit a sector)
- assume radial independence : $P(r, \theta) = P(r)P(\theta)$
- only need 7+3=10 numbers

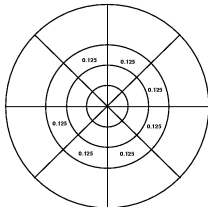- describe a system with $n$ features: $2^n - 1$ probabilities
- Use independence to reduce number of probabilities
- e.g. radially symmetric dartboard, P(hit a sector)
- assume radial independence : $P(r, \theta) = P(r)P(\theta)$
- only need 7+3=10 numbers

if $X$ and $Y$ are random variables, then

$P(x|y)$ is the probability that $X = x$ given that $Y = y$.

e.g.
$P(flies|is\_bird)$ is different than $P(flies)$
$P(flies|is\_a\_penguin, is\_bird)$ is different again

incorporate independence:
$P(flies|is\_bird, has\_feathers) = P(flies|is\_bird)$

Product rule (Chain rule):
$P(flies, is\_bird) = P(flies|is\_bird)P(is\_bird)$
$P(flies, is\_bird) = P(is\_bird|flies)P(flies)$

leads to : Bayes' rule
$P(is\_bird|flies) = \frac{P(flies|is\_bird)P(is\_bird)}{P(flies)}$

## Sum Rule

We know (an Axiom):

$$\sum_x P(X = x) = 1.0 \text{ and therefore that } \sum_x P(X = x|Y) = 1.0$$

This means that (Sum Rule)

$$\sum_x P(X = x, Y) = P(Y)$$

proof:

$$\sum_x P(X = x, Y) = \sum_x P(X = x|Y)P(Y)$$
$$= P(Y) \sum_x P(X = x|Y)$$
$$= P(Y)$$

We call $P(Y)$ the marginal distribution over $Y$

## Conditional Probability

- $X$ and $Y$ are independent iff

  $$P(X) = P(X|Y)$$
  $$P(Y) = P(Y|X)$$
  $$P(X, Y) = P(X)P(Y)$$

  so learning $Y$ doesn't influence beliefs about $X$

- $X$ and $Y$ are conditionally independent given $Z$ iff

  $$P(X|Z) = P(X|Y, Z)$$
  $$P(Y|Z) = P(Y|X, Z)$$
  $$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

  so learning $Y$ doesn't influence beliefs about $X$ if you already know $Z$ ...does not mean $X$ and $Y$ are independent

## Expected Values

expected value of a function on $X$, $V(X)$:

$$\mathbb{E}(V) = \sum_{x \in Dom(X)} P(x)V(x)$$

where $P(x)$ is the probability that $X = x$.

This is useful in decision making, where $V(X)$ is the *utility* of situation $X$.

Bayesian decision making is then

$$\mathbb{E}(V(\text{decision})) = \sum_{outcome} P(outcome|decision)V(outcome)$$

## Value of Independence

- complete independence reduces both representation and inference from $O(2^n)$ to $O(n)$
- Unfortunately, complete mutual independence is rare
- Fortunately, most domains do exhibit a fair amount of conditional independence
- Bayesian Networks or Belief Networks (BNs) encode this information

## Belief Networks

Bayesian network or belief network

- Directed Acyclic graph
- Encodes independencies in a graphical format
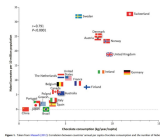- Edges give $P(X_i|parents(X_i))$

Cancer diagnosis example:



- Two tests A and B
- Test A is quick and cheap, but imprecise
- Test A results are read directly
- Test B uses a machine that sometimes malfunctions, but is more precise
- Test B results are not read directly,
- a Report is written (by a human who makes mistakes)
- the Report is entered into a database (by another human who makes mistakes)

## Correlation and Causality

- Directed links in Bayes' net $\approx$ causal
- However, not always the case: chocolate → Nobel or Nobel → chocolate?
- In a Bayes net, it doesn't matter!
- But, some structures will be easier to specify



In this example, its probably
$chocolate \leftarrow$ "$Switzerland - ness$" $\rightarrow Nobel$

## Bayesian networks - example

*If Jesse's alarm doesn't go off (A), Jesse probably won't get coffee (C); if Jesse doesn't get coffee, he's likely grumpy (G). If he is grumpy then it's possible that the lecture won't go smoothly L. If the lecture does not go smoothly then the students will likely be sad S.*



A=Jesse's alarm doesn't go off
C=Jesse doesn't get coffee
G=Jesse is grumpy
L=lecture doesn't go smoothly
S=students are sad

all variables binary (true/false)

## Conditional Independence



- If you learned any of $A$, $C$, $G$, or $L$, would your assessment of $P(S)$ change?
    - If any of these are seen to be true, you would increase $P(s)$ and decrease $P(\bar{s})$.
    - So $S$ is not independent of $A$, $C$, $G$, $L$.
- If you knew the value of $L$, would learning the value of $A$, $C$, or $G$ influence $P(S)$?
    - Influence that these factors have on $S$ is mediated by their influence on $L$.
    - Students aren't sad because Jesse was grumpy, they are sad because of the lecture.
    - Therefore, $S$ is conditionally independent of $A$, $C$, and $G$ (given $L$).

- We say: *S* is **independent** of *A*, *C*, and *G*, given *L*
- (this is **conditional independence**)
- Similarly, we can say
  - *S* is **independent** of *A* and *C*, given *G*
  - *G* is **independent** of *A*, given *C*
  - ...
- This means that:
  - $P(S|L, G, C, A) = P(S|L)$
  - $P(L|G, C, A) = P(L|G)$
  - $P(G|C, A) = P(G|C)$
  - $P(C|A)$ and $P(A)$ don't "simplify"

Chain rule ( **product rule** ):

$P(S, L, G, C, A) =$
$\quad P(S|L, G, C, A)P(L|G, C, A)P(G|C, A)P(C|A)P(A)$

**Independence:**

$P(S, L, G, C, A) = P(S|L)P(L|G)P(G|C)P(C|A)P(A)$

So we can specify the full **joint probability**
using the five local **conditional probabilities**:

$P(S|L), P(L|G), P(G|C), P(C|A), P(A)$

A **Bayesian Network** (Belief Network, Probabilistic Network) or BN over variables $\{X_1, X_2, \ldots, X_N\}$ consists of:
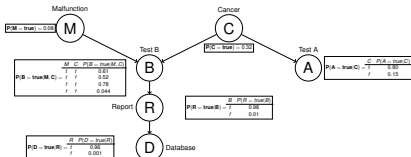
- a **DAG** whose nodes are the variables
- a set of **Conditional Probability tables** (CPTs) giving $P(X_i|Parents(X_i))$ for each $X_i$

**example probability tables** for the Coffee Bayes Net:

Cancer diagnosis:

## Semantics of a Bayes' Net

The structure of the BN means that :

every $X_i$ is conditionally independent of all its nondescendants given its parents:

$$P(X_i|S, Parents(X_i)) = P(X_i|Parents(X_i))$$

for any subset $S \subseteq NonDescendants(X_i)$

The BN defines a factorization of the joint probability distribution. The joint distribution is formed by multiplying the conditional probability tables together.

## Constructing belief networks

To represent a domain in a belief network, you need to consider:

- What are the relevant variables?
  - ▶ What will you observe? - this is the evidence
  - ▶ What would you like to find out? - this is the query
  - ▶ What other features make the model simpler? - these are the other variables
- What values should these variables take?
- What is the relationship between them? This should be expressed in terms of local influence.
- How does the value of each variable depend on its parents? This is expressed in terms of the conditional probabilities.

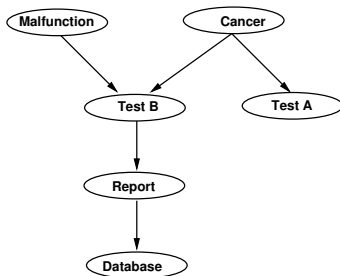## Bayesian Networks - Independence assumptions



- Test B depends on Cancer and Malfunction
- Test A depends only on Cancer
- Report depends only on Test B
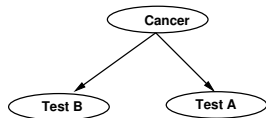- Database depends only on Report

What are the independencies?
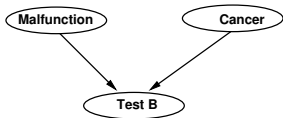
## Three Basic Bayesian Networks

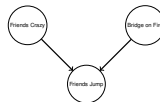**Database and Test B independent if Report is observed**
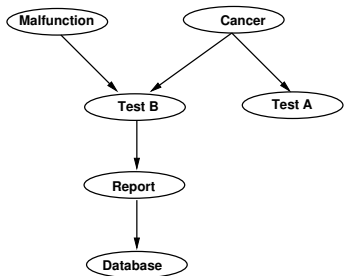
**Test B and Test A are independent if Cancer is observed**

**Malfunction and Cancer are independent if Test B is not observed**

http://imgs.xkcd.com/comics/bridge.png

- Given a BN, how do we determine if two variables X,Y are independent ( given evidence E )?
- D-separation : A set of variables $E$ d-separates X and Y if it blocks every undirected path in the BN between X and Y
- But what does block mean?

**(1)** X any undir. ⟶ Z ⟶ any undir. path Y

If Z in evidence, the path between X and Y blocked

**(2)** X any undir. path ⟵ Z ⟶ any undir. path Y

If Z in evidence, the path between X and Y blocked

**(3)** X any undir. path ⟶ Z ⟵ any undir. path Y

**Descendents(Z)**

If Z is **not** in evidence and **no** descendent of Z is in evidence, then the path between X and Y is blocked

...

The Markov Blanket of a node (variable) $V$ is:

- the parents, children, and the (other) parents of children
- the minimal set of nodes that d-separates $V$ from all other variables

The joint distribution over the Markov Blanket allows for the computation of the distribution $P(V)$.

- TravelSubway and Thermometer (given no evidence)?
- TravelSubway and Thermometer (given Flu or Fever)?
- TravelSubway and Malaria (given Fever)?
- TravelSubway and Exotic Trip (given Jaundice)?
- TravelSubway and Exotic Trip (given Jaundice and Thermometer)?
- TravelSubway and Exotic Trip (given Malaria and Thermometer)?
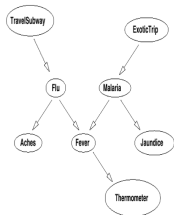
Agent has a prior belief in a hypothesis , $h$, $P(h)$,

Agent observes some evidence $e$
that has a likelihood given the hypothesis: $P(e|h)$.

The agent's posterior belief about $h$ after observing $e$, $P(h|e)$,

is given by Bayes' Rule:

$$P(h|e) = \frac{P(e|h)P(h)}{P(e)} = \frac{P(e|h)P(h)}{\sum_h P(e|h)P(h)}$$

- Often you have causal knowledge :
  $P(symptom \mid disease)$
  $P(light\ is\ off \mid status\ of\ switches\ and\ switch\ positions)$
  $P(alarm \mid fire)$
  $P(image\ looks\ like\ 📌 \mid a\ tree\ is\ in\ front\ of\ a\ car)$
- and want to do evidential reasoning :
  $P(disease \mid symptom)$
  $P(status\ of\ switches \mid light\ is\ off\ and\ switch\ positions)$
  $P(fire \mid alarm)$.
  $P(a\ tree\ is\ in\ front\ of\ a\ car \mid image\ looks\ like\ 📌)$

Before you get any information
- $P(Cancer) = 0.32$
- $P(Malfunction) = 0.08$

## Probabilistic Inference



Suppose the doctor reads a <mark>positive Test B in the Database</mark>
evidence gives Database=true (not directly Test B= true)
we want to know <mark>$P(Cancer = true|Database = true)$</mark>

- $P(Cancer = true|Database = true) = 0.80$
- $P(Malfunction = true|Database = true) = 0.14$

(we will see how to get these numbers later)

## Probabilistic Inference



Suppose <mark>Test A is positive</mark> as well
we want $P(Cancer = true|Database = true \wedge TestA = true)$

- $P(Cancer = true|Database = true \wedge TestA = true) = 0.95$
- $P(M = true|Database = true \wedge TestA = true) = 0.08$

(we will see how to get these numbers later)

## Probabilistic Inference



Suppose <mark>Test A is negative</mark>, though!
we want $P(Cancer = true|Database = true \wedge TestA = false)$

- $P(Cancer = true|Database = true \wedge TestA = false) = 0.48$
- $P(M = true|Database = true \wedge TestA = false) = 0.27$

(we will see how to get these numbers later)

## Simple Forward Inference (Chain)

Computing marginal requires simple forward propagation of
probabilities



- $P(J) = \sum_{M,ET} P(J, M, ET)$
  (marginalisation - sum rule)
- $P(J) = \sum_{M,ET} P(J|M, ET)P(M|ET)P(ET)$
  (chain rule)
- $P(J) = \sum_{M,ET} P(J|M)P(M|ET)P(ET)$
  (conditional indep.)
- $P(J) = \sum_M P(J|M) \sum_{ET} P(M|ET)P(ET)$
  (distribution of sum)

Note: all terms on the last line are CPTs in the BN
Note: only ancestors of J considered. Why?

Same idea when evidence "upstream"



- $P(J|et) = \sum_M P(J, M|et)$
  (marginalisation)
- $P(J|et) = \sum_M P(J|M, et)P(M|et)$
  (chain rule)
- $P(J|et) = \sum_M P(J|M)P(M|et)$
  (conditional indep.)

With multiple parents the evidence is "pooled"



$$P(Fev) = \sum_{Flu,M,TS,ET} P(Fev, Flu, M, TS, ET)$$
$$= \sum_{Flu,M} P(Fev|M, Flu)[\sum_{TS} P(Flu|TS)P(TS)][\sum_{ET} P(M|ET)P(ET)]$$

also works with "upstream" evidence



$$P(Fev|ts, \overline{m}) = \sum_{Flu} P(Fev, Flu|\overline{m}, ts)$$
$$= \sum_{Flu} P(Fev|Flu, ts, \overline{m})P(Flu|ts, \overline{m})$$
$$= \sum_{Flu} P(Fev|Flu, \overline{m})P(Flu|ts)$$

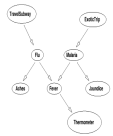When evidence is downstream of query, then we must reason "backwards". This requires Bayes' rule



$$P(ET|j) = P(j|ET)P(ET)/P(J) \propto P(j, ET)$$
$$= P(j|ET)P(ET)$$
$$= \sum_M P(j, M|ET)P(ET)$$
$$= \sum_M P(j|M, ET)P(M|ET)P(ET) \text{ (chain rule)}$$
$$= \sum_M P(j|M)P(M|ET)P(ET)$$

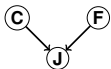normalising constant is $\frac{1}{P(j)}$, but this can be computed as

$$P(j) = \sum_{ET} P(ET, j)$$

http://imgs.xkcd.com/comics/bridge.png

F: Bridge on Fire
C: All friends Crazy
J: All friends Jump
What is $P(F|J = true)$?

| P(C − true) − 0.0001 | | P(F − true) − 0.1 |



| | | F | C | $P(J − true|F, C)$ |
|---|---|---|---|---|
| $P(J − true|F, C)$ | | t | t | 0.95 |
| | | t | f | 0.99 |
| | | f | t | 0.99 |
| | | f | f | 0.01 |

- intuitions above : polytree algorithm
- works for simple networks without loops
- more general algorithm: Variable Elimination
- applies sum-out rule repeatedly
- distributes sums

A factor is a representation of a function from a tuple of random variables into a number.
We will write factor $f$ on variables $X_1, \ldots, X_j$ as $f(X_1, \ldots, X_j)$.
We can assign some or all of the variables of a factor
→(this is restricting a factor):

- $f(X_1 = v_1, X_2, \ldots, X_j)$, where $v_1 \in dom(X_1)$, is a factor on $X_2, \ldots, X_j$.
- $f(X_1 = v_1, X_2 = v_2, \ldots, X_j = v_j)$ is a number that is the value of $f$ when each $X_i$ has value $v_i$.

The former is also written as $f(X_1, X_2, \ldots, X_j)_{X_1 = v_1}$, etc.

$r(X, Y, Z)$:

| X | Y | Z | val |
|---|---|---|---|
| t | t | t | 0.1 |
| t | t | f | 0.9 |
| t | f | t | 0.2 |
| t | f | f | 0.8 |
| f | t | t | 0.4 |
| f | t | f | 0.6 |
| f | f | t | 0.3 |
| f | f | f | 0.7 |

$r(X=t, Y, Z)$:

| Y | Z | val |
|---|---|---|
| t | t | 0.1 |
| t | f | 0.9 |
| f | t | 0.2 |
| f | f | 0.8 |

$r(X=t, Y, Z=f)$:

| Y | val |
|---|---|
| t | 0.9 |
| f | 0.8 |

$r(X=t, Y=f, Z=f) = 0.8$

The $\textcolor{yellow}{\text{product}}$ of factor $f_1(X, Y)$ and $f_2(Y, Z)$, where $Y$ are the variables in common, is the factor $(f_1 \times f_2)(X, Y, Z)$ defined by:

$$(f_1 \times f_2)(X, Y, Z) = f_1(X, Y) f_2(Y, Z).$$

$f_1$:

| A | B | val |
|---|---|-----|
| t | t | 0.1 |
| t | f | 0.9 |
| f | t | 0.2 |
| f | f | 0.8 |

$f_2$:

| B | C | val |
|---|---|-----|
| t | t | 0.3 |
| t | f | 0.7 |
| f | t | 0.6 |
| f | f | 0.4 |

$f_1 \times f_2$:

| A | B | C | val |
|---|---|---|------|
| t | t | t | 0.03 |
| t | t | f | 0.07 |
| t | f | t | 0.54 |
| t | f | f | 0.36 |
| f | t | t | 0.06 |
| f | t | f | 0.14 |
| f | f | t | 0.48 |
| f | f | f | 0.32 |

We can $\textcolor{yellow}{\text{sum out}}$ a variable, say $X_1$ with domain $\{v_1, \ldots, v_k\}$, from factor $f(X_1, \ldots, X_j)$, resulting in a factor on $X_2, \ldots, X_j$ defined by:

$$(\sum_{X_1} f)(X_2, \ldots, X_j)$$
$$= f(X_1 = v_1, \ldots, X_j) + \cdots + f(X_1 = v_k, \ldots, X_j)$$

$f_3$:

| A | B | C | val |
|---|---|---|------|
| t | t | t | 0.03 |
| t | t | f | 0.07 |
| t | f | t | 0.54 |
| t | f | f | 0.36 |
| f | t | t | 0.06 |
| f | t | f | 0.14 |
| f | f | t | 0.48 |
| f | f | f | 0.32 |

$\sum_B f_3$:

| A | C | val |
|---|---|------|
| t | t | 0.57 |
| t | f | 0.43 |
| f | t | 0.54 |
| f | f | 0.46 |

## Evidence

If we want to compute the posterior probability of $Z$ given evidence $Y_1 = v_1 \wedge \ldots \wedge Y_j = v_j$:

$$P(Z|Y_1 = v_1, \ldots, Y_j = v_j)$$
$$= \frac{P(Z, Y_1 = v_1, \ldots, Y_j = v_j)}{P(Y_1 = v_1, \ldots, Y_j = v_j)}$$
$$= \frac{P(Z, Y_1 = v_1, \ldots, Y_j = v_j)}{\sum_Z P(Z, Y_1 = v_1, \ldots, Y_j = v_j)}.$$

The computation reduces to the joint probability of

$P(Z, Y_1 = v_1, \ldots, Y_j = v_j).$

normalize at the end.

## Probability of a conjunction

Suppose the variables of the belief network are $X_1, \ldots, X_n$.
To compute $P(Z, Y_1 = v_1, \ldots, Y_j = v_j)$, we
sum out the variables other than query $Z$ and evidence $Y$,
$Z_1, \ldots, Z_k = \{X_1, \ldots, X_n\} - \{Z\} - \{Y_1, \ldots, Y_j\}$.
We order the $Z_i$ into an elimination ordering $Z_1 \ldots Z_k$.

$$P(Z, Y_1 = v_1, \ldots, Y_j = v_j)$$
$$= \sum_{Z_k} \cdots \sum_{Z_1} P(X_1, \ldots, X_n)_{Y_1 = v_1, \ldots, Y_j = v_j}.$$
$$= \sum_{Z_k} \cdots \sum_{Z_1} \prod_{i=1}^{n} P(X_i | parents(X_i))_{Y_1 = v_1, \ldots, Y_j = v_j}.$$

## Computing sums of products

Computation in belief networks reduces to
computing the sums of products.

- How can we compute $ab + ac$ efficiently?

## Computing sums of products

Computation in belief networks reduces to
computing the sums of products.

- How can we compute $ab + ac$ efficiently?
- Distribute out the $a$ giving $a(b + c)$

Computation in belief networks reduces to
computing the sums of products.
- How can we compute $ab + ac$ efficiently?
- Distribute out the $a$ giving $a(b + c)$
- How can we compute $\sum_{Z_1} \prod_{i=1}^{n} P(X_i | parents(X_i))$ efficiently?

Computation in belief networks reduces to
computing the sums of products.
- How can we compute $ab + ac$ efficiently?
- Distribute out the $a$ giving $a(b + c)$
- How can we compute $\sum_{Z_1} \prod_{i=1}^{n} P(X_i | parents(X_i))$ efficiently?
- Distribute out those factors that don't involve $Z_1$.

To compute $P(Z | Y_1 = v_1 \wedge \ldots \wedge Y_j = v_j)$:
- Construct a factor for each conditional probability .
- Restrict the observed variables to their observed values
- Sum out each of the other variables (the $\{Z_1, \ldots, Z_k\}$ from slide 45) according to some elimination ordering :
  for each $Z_i$ in order starting from $i = 1$:
  - collect all factors that contain $Z_i$
  - multiply together and sum out $Z_i$
  - add resulting new factor back to the pool
- Multiply the remaining factors.
- Normalize by dividing
  the resulting factor $f(Z)$ by $\sum_Z f(Z)$.

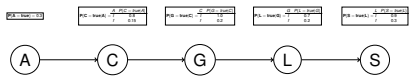To sum out a variable $Z_j$ from a product $f_1, \ldots, f_k$ of factors:
- Partition the factors into
  - those that don't contain $Z_j$, say $f_1, \ldots, f_i$,
  - those that contain $Z_j$, say $f_{i+1}, \ldots, f_k$

We know:

$$\sum_{Z_j} f_1 \times \cdots \times f_k = f_1 \times \cdots \times f_i \times \left( \sum_{Z_j} f_{i+1} \times \cdots \times f_k \right).$$

- Explicitly construct a representation of the rightmost factor $\left( \sum_{Z_j} f_{i+1} \times \cdots \times f_k \right)$.
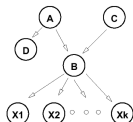- Replace the factors $f_{i+1}, \ldots, f_k$ by the new factor.

Example I

see note `variableelim.pdf`

- Complexity is linear in number of variables, and exponential in the size of the largest factor
- When we create new factors: sometimes this blows up
- Depends on the elimination ordering
- For polytrees : work outside in
- For general BNs this can be hard
- simply finding the optimal elimination ordering is NP-hard for general BNs
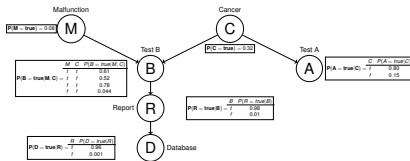- inference in general is NP-hard

- eliminate singly-connected nodes ($D, A, C, X_1, \ldots, X_k$) first
- Then no factor is ever larger than original CPTs
- If you eliminate $B$ first, a large factor is created that includes $A, C, X_1, \ldots, X_k$



- Certain variables have no impact
- In ABC network above, computing $P(A)$ does not require summing over $B$ and $C$

$$P(A) = \sum_{B,C} P(C|B)P(B|A)P(A)$$
$$= P(A) \sum_B P(B|A) \sum_C P(C|B) = P(A) * 1.0 * 1.0$$

- Can restrict attention to relevant variables:
- Given query $Q$ and evidence **E**, complete approximation is:
  - Q is relevant
  - if any node is relevant, its parents are relevant
  - if $E \in$ **E** is a descendent of a relevant variable, then $E$ is relevant
- irrelevant variable: a node that is not an ancestor of a query or evidence variable
- this will only remove irrelevant variables, but may not remove them all



see note variableelim.pdf

---

## Other Representations for Probability distributions      Next:

- Decision Tree or Graph:



- Noisy Or : $P(x | Y_1, \ldots, Y_k)$
- Logistic Regression

$$P(x | Y_1, \ldots, Y_k) = sigmoid(\sum_i w_i Y_i)$$

- Any deep differentiable function – see A. Stassopoulou and M. Petrou

  Obtaining the correspondence between Bayesian and Neural Networks, *International journal of pattern*

  *recognition and artificial intelligence* 12.07 (1998): 901-920.

  https://doi.org/10.1142/S021800149800049X

- Reasoning under Uncertainty Part II (Poole & Mackworth (2nd ed.)Chapter 8.5-8.9)