

Machine Intelligence

Lecture 6: Inference in Bayesian networks

Thomas Dyhre Nielsen

Aalborg University

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- **Inference in Bayesian networks**
- Machine learning
- Planning
- Multi-agent systems

Exact Inference

Posterior Marginals

Inference Problem:

- Given: a Bayesian network
- Given: an assignment of values to some of the variables in the network: $E_i = e_i$ ($i = 1, \dots, l$)
 - “Instantiation of the nodes \mathbf{E} ”
 - “Evidence $\mathbf{E} = \mathbf{e}$ entered”)
 - “Findings entered”
 - ...
- Want: for variables $A \notin \mathbf{E}$ the *posterior marginal* $P(A \mid \mathbf{E} = \mathbf{e})$.

According to the definition of conditional probability:

$$P(A \mid \mathbf{E} = \mathbf{e}) = \frac{P(A, \mathbf{E} = \mathbf{e})}{P(\mathbf{E} = \mathbf{e})}$$

It is sufficient to compute for each $a \in D_A$ the value

$$P(A = a, \mathbf{E} = \mathbf{e}).$$

Together with

$$P(\mathbf{E} = \mathbf{e}) = \sum_{a \in D_A} P(A = a, \mathbf{E} = \mathbf{e})$$

this gives the desired posterior distribution.

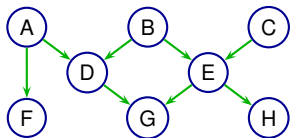
Inference as summation

Let A be the variable of interest, \mathbf{E} the evidence variables, and $\mathbf{Y} = Y_1, \dots, Y_l$ the remaining variables in the network not belonging to $A \cup \mathbf{E}$. Then

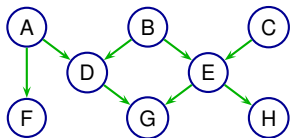
$$P(A = a, \mathbf{E} = \mathbf{e}) = \sum_{y_1 \in D_{Y_1}} \dots \sum_{y_l \in D_{Y_l}} P(A = a, \mathbf{E} = \mathbf{e}, Y_1 = y_1, \dots, Y_l = y_l).$$

Note:

- For each \mathbf{y} the probability $P(A = a, \mathbf{E} = \mathbf{e}, \mathbf{Y} = \mathbf{y})$ can be computed from the network (in time linear in the number of random variables).
- There number of configurations over \mathbf{Y} is exponential in l .



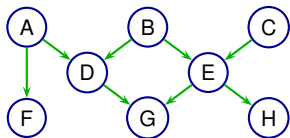
Find $P(B|a, f, g, h) = \frac{P(B, a, f, g, h)}{P(a, f, g, h)}$



Find $P(B|a, f, g, h) = \frac{P(B, a, f, g, h)}{P(a, f, g, h)}$

We can if we have access to $P(A, B, C, D, E, F, G, H)$:

$$P(A, B, C, D, E, F, G, H) = P(A)P(B)P(C)P(D|A, B) \cdot \dots \cdot P(H|E)$$



$$\text{Find } P(B|a, f, g, h) = \frac{P(B, a, f, g, h)}{P(a, f, g, h)}$$

We can if we have access to $P(A, B, C, D, E, F, G, H)$:

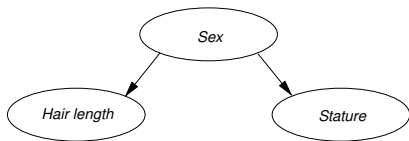
$$P(A, B, C, D, E, F, G, H) = P(A)P(B)P(C)P(D|A, B) \cdot \dots \cdot P(H|E)$$

Inserting evidence we get:

$$P(B, a, f, g, h) = \sum_{C, D, E} P(a, B, C, D, E, f, g, h)$$

and

$$P(a, f, g, h) = \sum_B P(B, a, f, g, h)$$



Conditional probability tables:

<i>Sex</i>	
<i>male</i>	0.49
<i>female</i>	0.51

<i>Hair length</i>	<i>Sex</i>	
	<i>male</i>	<i>female</i>
<i>long</i>	0.05	0.6
<i>short</i>	0.95	0.4

<i>Stature</i>	<i>Sex</i>	
	<i>male</i>	<i>female</i>
≤ 1.68	0.08	0.47
> 1.68	0.92	0.53

Posterior probability inference: Given the value of some observed variables (the evidence) compute the conditional distribution of some other variable:

$$P(\textit{Stature} \mid \textit{Hair length} = \textit{long}) = ?$$

$$P(\textit{Sex} \mid \textit{Hair length} = \textit{short}, \textit{Stature} \leq 1.68) = ?$$

$P(\text{Sex})$

<i>Sex</i>	
<i>male</i>	0.49
<i>female</i>	0.51

$P(\text{Hair length} \mid \text{Sex})$

<i>Hair length</i>	<i>Sex</i>	
	<i>male</i>	<i>female</i>
<i>long</i>	0.05	0.6
<i>short</i>	0.95	0.4

$P(\text{Stature} \mid \text{Sex})$

<i>Stature</i>	<i>Sex</i>	
	<i>male</i>	<i>female</i>
≤ 1.68	0.08	0.47
> 1.68	0.92	0.53

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

$P(\text{Sex})$

Sex	
male	0.49
female	0.51

$P(\text{Hair length} \mid \text{Sex})$

Hair length	Sex	
	male	female
long	0.05	0.6
short	0.95	0.4

$P(\text{Stature} \mid \text{Sex})$

Stature	Sex	
	male	female
≤ 1.68	0.08	0.47
> 1.68	0.92	0.53

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

Step 1: Construct joint distribution

$P(\text{Sex}, \text{Hair length}, \text{Stature})$

	Sex			
	male		female	
	Hair length		Hair length	
Stature	long	short	long	short
≤ 1.68				
> 1.68				

$$P(\text{Sex})$$

Sex	
male	0.49
female	0.51

$$P(\text{Hair length} \mid \text{Sex})$$

Hair length	Sex	
	male	female
long	0.05	0.6
short	0.95	0.4

$$P(\text{Stature} \mid \text{Sex})$$

Stature	Sex	
	male	female
≤ 1.68	0.08	0.47
> 1.68	0.92	0.53

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

Step 1: Construct joint distribution

$$P(\text{Sex}, \text{Hair length}, \text{Stature})$$

	Sex			
	male		female	
	Hair length		Hair length	
Stature	long	short	long	short
≤ 1.68	0.00196			
> 1.68				

$$P(\text{Sex})$$

Sex	
male	0.49
female	0.51

$$P(\text{Hair length} \mid \text{Sex})$$

Hair length	Sex	
	male	female
long	0.05	0.6
short	0.95	0.4

$$P(\text{Stature} \mid \text{Sex})$$

Stature	Sex	
	male	female
≤ 1.68	0.08	0.47
> 1.68	0.92	0.53

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

Step 1: Construct joint distribution

$$P(\text{Sex}, \text{Hair length}, \text{Stature})$$

	Sex			
	male		female	
Stature	Hair length		Hair length	
	long	short	long	short
≤ 1.68	0.00196	0.03724		
> 1.68				

$$P(\text{Sex})$$

Sex	
male	0.49
female	0.51

$$P(\text{Hair length} \mid \text{Sex})$$

Hair length	Sex	
	male	female
long	0.05	0.6
short	0.95	0.4

$$P(\text{Stature} \mid \text{Sex})$$

Stature	Sex	
	male	female
≤ 1.68	0.08	0.47
> 1.68	0.92	0.53

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

Step 1: Construct joint distribution

$$P(\text{Sex}, \text{Hair length}, \text{Stature})$$

	Sex			
	male		female	
	Hair length		Hair length	
Stature	long	short	long	short
≤ 1.68	0.00196	0.03724	0.14382	0.09588
> 1.68	0.02254	0.42826	0.16218	0.10812

Joint distribution $P(\text{Sex}, \text{Hair length}, \text{Stature})$

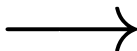
	Sex			
	male		female	
	Hair length		Hair length	
Stature	long	short	long	short
≤ 1.68	0.00196	0.03724	0.14382	0.09588
> 1.68	0.02254	0.42826	0.16218	0.10812

Step 2 “Enter evidence” :

$P(\text{Sex}, \text{Hair length}, \text{Stature})$

	Sex			
	male		female	
	Hair length		Hair length	
Stature	long	short	long	short
≤ 1.68	0.00196	0.03724	0.14382	0.09588
> 1.68	0.02254	0.42826	0.16218	0.10812

Hair length = long

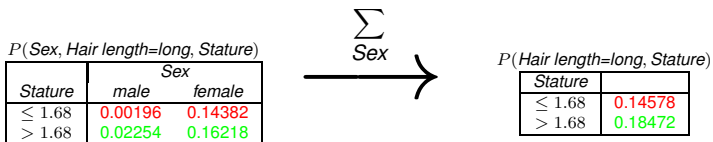


$P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Stature	Sex	
	male	female
≤ 1.68	0.00196	0.14382
> 1.68	0.02254	0.16218

Note: the table on the right shows neither a joint nor a conditional distribution!

Step 3 Marginalize (sum out Sex variable):



Step 3 Marginalize (sum out Sex variable):

Stature	Sex	
	male	female
≤ 1.68	0.00196	0.14382
> 1.68	0.02254	0.16218

$\xrightarrow{\sum_{\text{Sex}}}$

Stature	
≤ 1.68	0.14578
> 1.68	0.18472

Step 4 Normalize

Stature	
≤ 1.68	0.14578
> 1.68	0.18472

$\xrightarrow{\frac{1}{0.14578+0.18472}}$

Stature	
≤ 1.68	0.441
> 1.68	0.559

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) =$
 $P(\text{Sex})P(\text{Hair length} \mid \text{Sex})P(\text{Stature} \mid \text{Sex})$

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) =$
 $P(\text{Sex})P(\text{Hair length} \mid \text{Sex})P(\text{Stature} \mid \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) =$
 $P(\text{Sex})P(\text{Hair length} \mid \text{Sex})P(\text{Stature} \mid \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Marginalize: $P(\text{Hair length}=\text{long}, \text{Stature})=$

$$\sum_{s \in \{\text{male}, \text{female}\}} P(\text{Sex}=s, \text{Hair length}=\text{long}, \text{Stature})$$

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) =$
 $P(\text{Sex})P(\text{Hair length} \mid \text{Sex})P(\text{Stature} \mid \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Marginalize: $P(\text{Hair length}=\text{long}, \text{Stature})=$

$$\sum_{s \in \{\text{male}, \text{female}\}} P(\text{Sex}=s, \text{Hair length}=\text{long}, \text{Stature})$$

Condition: $P(\text{Stature} \mid \text{Hair length}=\text{long}) =$

$$\frac{P(\text{Hair length}=\text{long}, \text{Stature})}{P(\text{Hair length}=\text{long}, \text{Stature} \leq 1.68) + P(\text{Hair length}=\text{long}, \text{Stature} > 1.68)}$$

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} \mid \text{Sex})P(\text{Stature} \mid \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Marginalize: $P(\text{Hair length}=\text{long}, \text{Stature})=$

$$\sum_{s \in \{\text{male}, \text{female}\}} P(\text{Sex}=s, \text{Hair length}=\text{long}, \text{Stature})$$

Condition: $P(\text{Stature} \mid \text{Hair length}=\text{long}) =$

$$\frac{P(\text{Hair length}=\text{long}, \text{Stature})}{P(\text{Hair length}=\text{long}, \text{Stature} \leq 1.68) + P(\text{Hair length}=\text{long}, \text{Stature} > 1.68)}$$

Complexity

Complexity dominated by initial table $P(\text{Sex}, \text{Hair length}, \text{Stature})$ (size 2^3).

For model with n binary variables:

$$O(2^n)$$

Problem

The joint probability distribution will contain exponentially many entries.

Idea

We can use

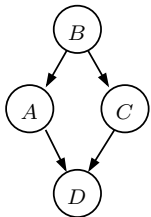
- the form of the joint distribution P
- the law of distributivity

to make the computation of the sum more efficient.

Variable Elimination

Thus, we can adapt our elimination procedure so that:

- we marginalize out variables sequentially
- when marginalizing out a particular variable X , we only need to consider the factors involving X .



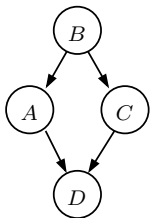
	<i>B</i>	
	<i>t</i>	<i>f</i>
	.5	.5

<i>B</i>	<i>A</i>	
	<i>t</i>	<i>f</i>
<i>t</i>	.7	.3
<i>f</i>	.1	.9

<i>B</i>	<i>C</i>	
	<i>t</i>	<i>f</i>
<i>t</i>	.7	.3
<i>f</i>	.2	.8

<i>A</i>	<i>C</i>	<i>D</i>	
		<i>t</i>	<i>f</i>
<i>t</i>	<i>t</i>	.9	.1
<i>t</i>	<i>f</i>	.7	.3
<i>f</i>	<i>t</i>	.8	.2
<i>f</i>	<i>f</i>	.4	.6

$$P(A, D = f) =$$



	<i>B</i>	
	<i>t</i>	<i>f</i>
	.5	.5

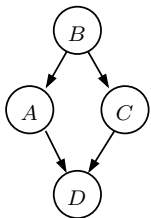
<i>B</i>	<i>A</i>	
	<i>t</i>	<i>f</i>
<i>t</i>	.7	.3
<i>f</i>	.1	.9

<i>B</i>	<i>C</i>	
	<i>t</i>	<i>f</i>
<i>t</i>	.7	.3
<i>f</i>	.2	.8

<i>A</i>	<i>C</i>	<i>D</i>	
		<i>t</i>	<i>f</i>
<i>t</i>	<i>t</i>	.9	.1
<i>t</i>	<i>f</i>	.7	.3
<i>f</i>	<i>t</i>	.8	.2
<i>f</i>	<i>f</i>	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A = c, D = f) =$$



	B	
	t	f
	.5	.5

B	A	
	t	f
t	.7	.3
f	.1	.9

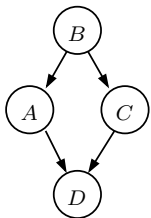
B	C	
	t	f
t	.7	.3
f	.2	.8

A	C	D	
		t	f
t	t	.9	.1
t	f	.7	.3
f	t	.8	.2
f	f	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A = c, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b)P(A = c | B = b)P(D = f | A = c, B = b) =$$



	B	
	t	f
	.5	.5

B	A	
	t	f
t	.7	.3
f	.1	.9

B	C	
	t	f
t	.7	.3
f	.2	.8

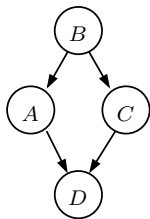
A	C	D	
		t	f
t	t	.9	.1
t	f	.7	.3
f	t	.8	.2
f	f	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b)P(A | B = b)P(C = c | B = b)P(D = f | A, C = c) =$$

$$\sum_{b \in \{t, f\}} P(B = b)P(A | B = b) \sum_{c \in \{t, f\}} P(C = c | B = b)P(D = f | A, C = c)$$



	B	
	t	f
	.5	.5

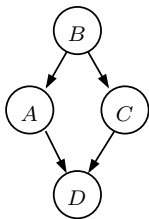
	A	
B	t	f
t	.7	.3
f	.1	.9

	C	
B	t	f
t	.7	.3
f	.2	.8

		D	
	A	C	
	t	t	
	t	f	
	f	t	
	f	f	
		t	f
			.9
			.1
			.7
			.3
			.8
			.2
			.4
			.6

$$\sum_b P(B = b)P(A | B = b) \sum_c P(C = c | B = b)P(D = f | A, C = c) = \\
 \sum_b P(B = b)P(A | B = b)F_1(B = b, A) = F_2(A)$$

Example continued



		B	
	t	f	
	.5	.5	

B	A	
	t	f
t	.7	.3
f	.1	.9

B	C	
	t	f
t	.7	.3
f	.2	.8

		D	
A	C	t	f
t	t	.9	.1
t	f	.7	.3
f	t	.8	.2
f	f	.4	.6

$$\sum_b P(B = b)P(A | B = b) \sum_c P(C = c | B = b)P(D = f | A, C = c) = \sum_b P(B = b)P(A | B = b)F_1(B = b, A) = F_2(A)$$

where

B	C	
	t	f
t	.7	.3
f	.2	.8

A	C	D	
		t	f
t	t	.9	.1
t	f	.7	.3
f	t	.8	.2
f	f	.4	.6

b	a	F ₁ (B, A)	
t	t	.7 · .1	+ .3 · .3 = .16
t	f	.7 · .2	+ .3 · .6 = .32
f	t	.2 · .1	+ .8 · .3 = .26
f	f	.2 · .2	+ .8 · .6 = .52

and

		B	
	t	f	
	.5	.5	

B	A	
	t	f
t	.7	.3
f	.1	.9

b	a	F ₁ (B, A)	
t	t	.16	
⋮	⋮	⋮	

a	F ₂ (A)	
t	...	
f	...	

Calculus of factors

- The procedure operates on **factors**: functions of subsets of variables
- Required operations on factors:
 - *multiplication*
 - *marginalization* (summing out selected variables)
 - *restriction* (setting selected variables to specific values)

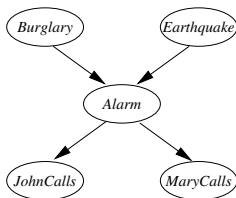
Complexity

Call subsets \mathbf{U} of \mathbf{V} that are the arguments of factors $P(\dots | \dots)$ resp. $F_j(\dots)$ which appear in the elimination process *factor sets*.

The complexity of variable elimination is exponential in the size of the largest factor set.

The size of the largest factor set can depend strongly on the order in which variables are summed out!

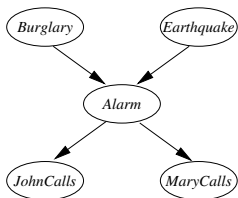
Example



Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

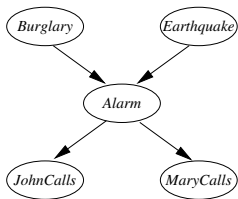
Example



Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

Example

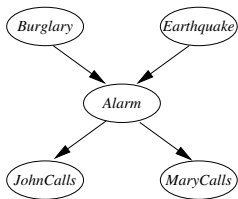


Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)F_1(eq, jc, MC) =$$

Example

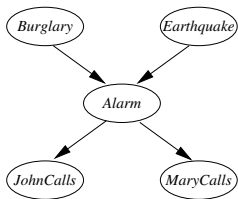


Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)F_1(eq, jc, MC) =$$

Example



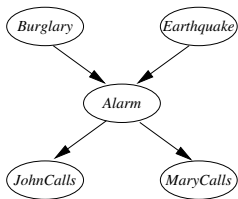
Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)F_1(eq, jc, MC) =$$

$$\sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq)F_2(eq, MC) =$$

Example



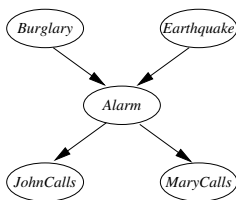
Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)F_1(eq, jc, MC) =$$

$$\sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq)F_2(eq, MC) =$$

Example



Bad ordering for computing $P(MC, B = t)$:

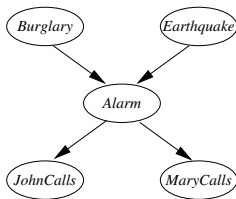
$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)F_1(eq, jc, MC) =$$

$$\sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq)F_2(eq, MC) =$$

$$P(B = t)F_3(MC)$$

Example



Bad ordering for computing $P(MC, B = t)$:

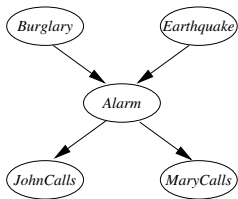
$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} \sum_{a \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

$$\sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)F_1(eq, jc, MC) =$$

$$\sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq)F_2(eq, MC) =$$

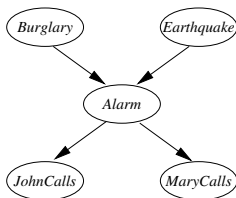
$$P(B = t)F_3(MC)$$

Largest factor (F_1) is function of 3 variables!



Good ordering for computing $P(MC, B = t)$:

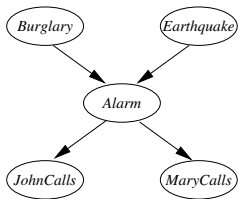
$$\sum_{a \in \{t, f\}} \sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$



Good ordering for computing $P(MC, B = t)$:

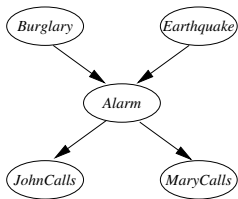
$$\sum_{a \in \{t, f\}} \sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) =$$

$$\sum_{a \in \{t, f\}} \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(MC | A = a)F_1(a) =$$



Good ordering for computing $P(MC, B = t)$:

$$\begin{aligned}
 & \sum_{a \in \{t, f\}} \sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) = \\
 & \sum_{a \in \{t, f\}} \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(MC | A = a)F_1(a) = \\
 & \sum_{a \in \{t, f\}} P(B = t)P(MC | A = a)F_1(a)F_2(a) =
 \end{aligned}$$

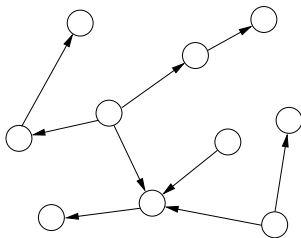


Good ordering for computing $P(MC, B = t)$:

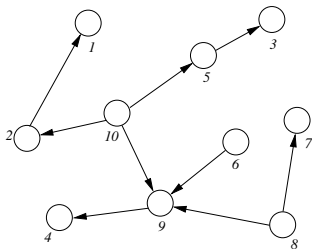
$$\begin{aligned}
 & \sum_{a \in \{t, f\}} \sum_{eq \in \{t, f\}} \sum_{jc \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) = \\
 & \sum_{a \in \{t, f\}} \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq)P(A = a | B = t, EQ = eq)P(MC | A = a)F_1(a) = \\
 & \sum_{a \in \{t, f\}} P(B = t)P(MC | A = a)F_1(a)F_2(a) = \\
 & P(B = t)F_3(MC)
 \end{aligned}$$

Largest factor ($P(A | B = t, EQ)$) is function of 2 variables!

A **singly connected network** is a network in which any two nodes are connected by at most one path of undirected edges:



A **singly connected network** is a network in which any two nodes are connected by at most one path of undirected edges:



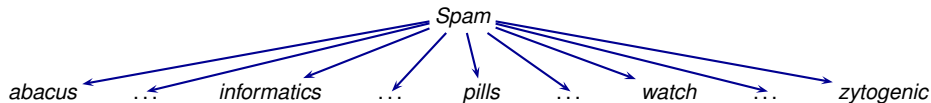
For singly connected network: any elimination order that “peels” variables from outside will only create factors of only one variable.

The complexity of inference is therefore linear in the total size of the network (= combined size of all conditional probability tables).

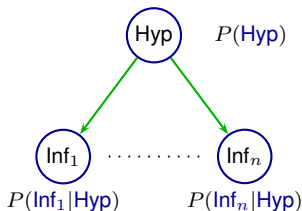
Example: Spam filter

- A single *query variable*: *Spam*
- Many observable features (e.g. words appearing in the body of the message):
abacus, ..., informatics, pills, ..., watch, ..., zytogenic

Network Structure:



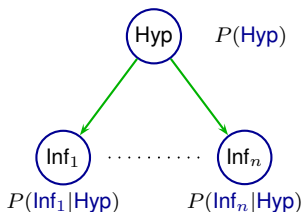
- Inference with large number of variables possible
- Essentially how *Thunderbird* spam filter works



We want the posterior probability of the hypothesis variable **Hyp** given the observations $\{\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n\}$:

$$P(\text{Hyp} | \text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n) = \frac{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n | \text{Hyp}) P(\text{Hyp})}{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n)}$$

Note: The model assumes that the **information variables** are independent given the **hypothesis variable**.



We want the posterior probability of the hypothesis variable **Hyp** given the observations $\{\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n\}$:

$$\begin{aligned}
 P(\text{Hyp} | \text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n) &= \frac{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n | \text{Hyp}) P(\text{Hyp})}{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n)} \\
 &= \mu \cdot P(\text{Inf}_1 = e_1 | \text{Hyp}) \cdot \dots \cdot P(\text{Inf}_n = e_n | \text{Hyp}) P(\text{Hyp})
 \end{aligned}$$

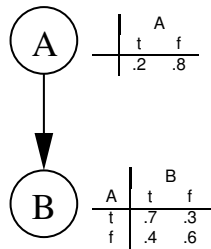
Note: The model assumes that the **information variables** are independent given the **hypothesis variable**.

Approximate Inference

Sample Generator

Observation: can use Bayesian network as random generator that produces states $\mathbf{X} = \mathbf{x}$ according to distribution P defined by the network.

Example:



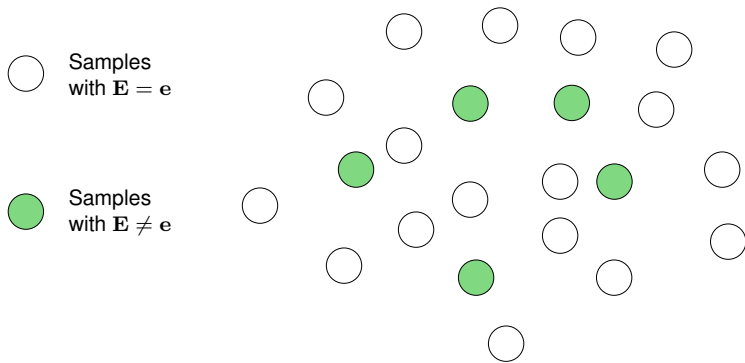
- Generate random numbers r_1, r_2 uniformly from $[0,1]$.
- Set $A = t$ if $r_1 \leq .2$ and $A = f$ else.
- Depending on the value of A and r_2 set B to t or f .

Random generation of one state: linear in size of network.

Approximate Inference from Samples

To compute an approximation of $P(\mathbf{E} = \mathbf{e})$ (\mathbf{E} a subset of the variables in the Bayesian network):

- generate a (large) number of random states
- count the frequency of states in which $\mathbf{E} = \mathbf{e}$.



Hoeffding Bound

- p : true probability $P(\mathbf{E} = \mathbf{e})$
- s : estimate for p from sample of size n
- ϵ : an error bound > 0 .

Then

$$P(|s - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

Hoeffding Bound

- p : true probability $P(\mathbf{E} = \mathbf{e})$
- s : estimate for p from sample of size n
- ϵ : an error bound > 0 .

Then

$$P(|s - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

Required Sample Size

To obtain an estimate that *with probability at most* δ has an *accuracy at least* ϵ , it is sufficient to take

$$n = -\ln(\delta/2)/(2\epsilon^2) \text{ samples.}$$

Hoeffding Bound

- p : true probability $P(\mathbf{E} = e)$
- s : estimate for p from sample of size n
- ϵ : an error bound > 0 .

Then

$$P(|s - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

Required Sample Size

To obtain an estimate that *with probability at most δ* has an *accuracy at least ϵ* , it is sufficient to take

$$n = -\ln(\delta/2)/(2\epsilon^2) \text{ samples.}$$

Example

To get an error ϵ of less than 0.1 in 95% of the cases ($\delta = 0.05$), we need:

$$n > -\ln(0.05/2)/(2 \cdot 0.1^2) \approx 184 \text{ samples}$$

Hoeffding Bound

- p : true probability $P(\mathbf{E} = e)$
- s : estimate for p from sample of size n
- ϵ : an error bound > 0 .

Then

$$P(|s - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

Required Sample Size

To obtain an estimate that *with probability at most δ* has an *accuracy at least ϵ* , it is sufficient to take

$$n = -\ln(\delta/2)/(2\epsilon^2) \text{ samples.}$$

Example

To get an error ϵ of less than 0.1 in 95% of the cases ($\delta = 0.05$), we need:

$$n > -\ln(0.05/2)/(2 \cdot 0.1^2) \approx 184 \text{ samples}$$

How many samples do we need if the error should be less than 0.01?

Hoeffding Bound

- p : true probability $P(\mathbf{E} = e)$
- s : estimate for p from sample of size n
- ϵ : an error bound > 0 .

Then

$$P(|s - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

Required Sample Size

To obtain an estimate that *with probability at most δ* has an *accuracy at least ϵ* , it is sufficient to take

$$n = -\ln(\delta/2)/(2\epsilon^2) \text{ samples.}$$

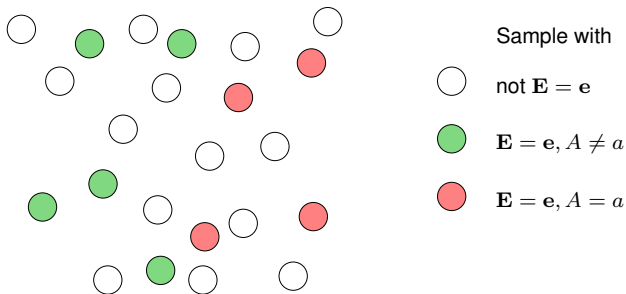
Example

To get an error ϵ of less than 0.1 in 95% of the cases ($\delta = 0.05$), we need:

$$n > -\ln(0.05/2)/(2 \cdot 0.1^2) \approx 184 \text{ samples}$$

How many samples do we need if the error should be less than 0.01? 18444 samples

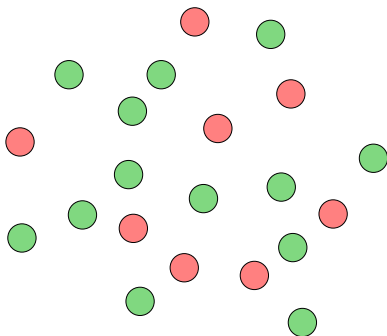
The simplest approach: **Rejection Sampling**



Approximation for $P(A = a \mid E = e)$:
$$\frac{\# \text{ (red circle)}}{\# \text{ (green circle)} \cup \# \text{ (red circle)}}$$

Problem with rejection sampling: samples with $\mathbf{E} \neq \mathbf{e}$ are useless!

Ideally: would draw samples directly from the conditional distribution $P(\mathbf{A} \mid \mathbf{E} = \mathbf{e})$.



First idea (*not to be followed*)

- Fix evidence variables to their observed states.
- Sample from non-evidence variables.

First idea (*not to be followed*)

- Fix evidence variables to their observed states.
- Sample from non-evidence variables.

Problem: This gives a sampling distribution

$$\prod_{X \in \mathbf{X} \setminus \mathbf{E}} P(X \mid \text{pa}(X) \setminus \mathbf{E}, \text{pa}(X) \cap \mathbf{E})$$

somewhere between $P(\mathbf{X})$ and $P(\mathbf{X} \mid \mathbf{e})$.

First idea (*not to be followed*)

- Fix evidence variables to their observed states.
- Sample from non-evidence variables.

Problem: This gives a sampling distribution

$$\prod_{X \in \mathbf{X} \setminus \mathbf{E}} P(X \mid \text{pa}(X) \setminus \mathbf{E}, \text{pa}(X) \cap \mathbf{E})$$

somewhere between $P(\mathbf{X})$ and $P(\mathbf{X} \mid \mathbf{e})$.

Likelihood weighting

We would like to sample from

$$P(\mathbf{X}, \mathbf{e}) = \underbrace{\prod_{X \in \mathbf{X} \setminus \mathbf{E}} P(X \mid \text{pa}(X) \setminus \mathbf{E}, \text{pa}(X) \cap \mathbf{E})}_{\text{Part 1}} \cdot \underbrace{\prod_{E \in \mathbf{E}} P(E = e \mid \text{pa}(E) \setminus \mathbf{E}, \text{pa}(E) \cap \mathbf{E})}_{\text{Part 2}}$$

So instead weigh each generated sample with a weight corresponding to Part 2.

Likelihood weighting

Estimate $P(X = e | \mathbf{e})$ as

$$\hat{P}(X = e | \mathbf{e}) = \frac{\sum_{\text{sample}: X=e} w(\text{sample})}{\sum_{\text{sample}} w(\text{sample})},$$

where

$$w(\text{sample}) = \prod_{E \in \mathbf{E}} P(E = e | \text{pa}(E) = \pi) \quad (\text{Part 2 on the previous slide})$$

and π is the values of $\text{pa}(E)$ under *sample* and \mathbf{e} .

Likelihood weighting

Estimate $P(X = e | \mathbf{e})$ as

$$\hat{P}(X = e | \mathbf{e}) = \frac{\sum_{\text{sample}: X=e} w(\text{sample})}{\sum_{\text{sample}} w(\text{sample})},$$

where

$$w(\text{sample}) = \prod_{E \in \mathbf{E}} P(E = e | \text{pa}(E) = \pi) \quad (\text{Part 2 on the previous slide})$$

and π is the values of $\text{pa}(E)$ under *sample* and \mathbf{e} .

Importance sampling

Likelihood weighting is an instance of importance sampling, where

- samples are weighted and can come from (almost) any proposal distribution.