

Introduction to Artificial Intelligence

Roman Barták

Department of Theoretical Computer Science and Mathematical Logic



Agents may need to handle **uncertainty** due to:

- **partial observability** (agent may not know for certain the state, where it is)
- **nondeterminism** (agent may not know where it will end up after performing its action)

Logical agent can:

- work with **belief states** (belief state = a set of possible world states, where it might be in)
- generate **contingency plans** (contingency plan handles every possible eventuality).

But, belief-state representations and contingency plans can be impossible **large** and **complex** as they need to cover every possible explanation of observation and every eventuality.

A logical agent believes each sentence to be true or false or has no opinion.

A **probabilistic agent** may have a numerical degree of belief between 0 (certainly false) and 1 (certainly true).

We have a maze with pits that are detected in neighboring squares via breeze (Wumpus and gold will not be assumed now).

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 B OK	2,2	3,2	4,2
1,1 OK	2,1 B OK	3,1	4,1

Where does the agent should go, if there is breeze at (1,2) and (2,1)?

Each cell – (1,3), (2,2), (3,1) – might contain a pit. Pure logical inference can conclude nothing about which square is most likely to be safe!

=> a logical agent might have to choose randomly

Can we do better?



Like logical assertions, probabilistic assertions are about possible worlds – **sample space** Ω .

- the possible worlds are **mutually exclusive** and **exhaustive**

Each **possible world** ω is associated with a numerical probability $P(\omega)$ such that:

$$0 \leq P(\omega) \leq 1, \sum_{\omega \in \Omega} P(\omega) = 1$$

Example: If we are about to roll two (distinguishable) dice, there are 36 possible worlds to consider: (1,1), (1,2), ..., (6,6), $P(\omega) = 1/36$

The sets of possible worlds are called **events**.

The probability of event is the sum of probabilities of possible worlds in the event.

$$P(\phi) = \sum_{\omega \in \phi} P(\omega)$$

These probabilities are called **unconditional** or **prior probabilities** („priors“ for short).

Example: $P(\text{doubles}) = 1/36 + 1/36 + 1/36 + 1/36 + 1/36 + 1/36 = 1/6$

Frequently, we have some information, called **evidence** (**b**), and we are interested in probability of some event (**a**).

$$P(a \mid b) = P(a \wedge b) / P(b), \text{ whenever } P(b) > 0$$

This is called **conditional** or **posterior probability**

Example: What is the probability of double if we already know that first die rolled to 5?

$$P(\text{doubles} \mid \text{Die}_1 = 5) = 1/36 / (6 \cdot 1/36) = 1/6$$

This can be also written in a different form called the **product rule**:

$$P(a \wedge b) = P(a \mid b) \cdot P(b)$$



We can refer to a possible world using a **factored representation** – a possible world is represented by a set of variable/value pairs.

Variables in probability theory are called **random variables**. Every random variable has a **domain** – the set of possible values it can take on (similarly to a CSP).

Die₁ – represents a value on the first die 1 (1,...,6)

Cavity – describes whether the patient has or has not cavity (true, false)

A possible world is fully identified by values of all random variables.

P(Die₁ = 5, Die₂ = 5)

Probability of all possible worlds can be described using a table called a **full joint probability distribution**

– elements are indexed by values of random variables.

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

To compute posterior probability of a query proposition given observed evidence, we add up probabilities of possible worlds in which the proposition is true (**marginalization** or **summing out**).

$$P(\phi) = \sum_{\omega:\omega|\models\phi} P(\omega), \quad P(Y) = \sum_{z \in Z} P(Y,z)$$

A variant of this rule involves conditional probabilities (**conditioning**):

$$P(Y) = \sum_{z \in Z} P(Y|z)P(z)$$

In most cases, we are interested in **computing conditional probabilities** of some variables, given **evidence** about others:

$$P(Y | E=e) = P(Y, E=e) / P(E=e)$$

We also know $\sum_{y \in Y} P(Y=y | E=e) = 1$, so we do not need to calculate $P(E=e)$ at all and we can do **normalization** instead:

$$P(Y | E=e) = \alpha P(Y, E=e)$$

where $\alpha = 1 / P(E=e) = 1 / \sum_{y \in Y} P(Y=y, E=e)$ (normalization constant).

Probabilistic query answering:

In a typical case, we know values \mathbf{e} of random variables \mathbf{E} from the **observation** and we are looking for probability distribution of random variables \mathbf{Y} from the **query**. The other random variables are **hidden** $\mathbf{H} = \mathbf{X} - \mathbf{Y} - \mathbf{E}$.

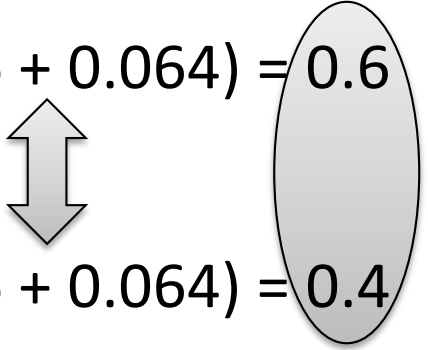
$$P(Y | E=e) = \alpha P(Y, E=e) = \alpha \sum_h P(Y, E=e, H=h)$$

$$P(\text{toothache}) (= P(\text{Toothache}=\text{true})) \\ = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$$

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

$$P(\text{cavity} \vee \text{toothache}) \\ = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$$

$$P(\text{cavity} | \text{toothache}) \\ = P(\text{cavity} \wedge \text{toothache}) / P(\text{toothache}) \\ = (0.108 + 0.012) / (0.108 + 0.012 + 0.016 + 0.064) = 0.6$$

$$P(\neg \text{cavity} | \text{toothache}) \\ = P(\neg \text{cavity} \wedge \text{toothache}) / P(\text{toothache}) \\ = (0.016 + 0.064) / (0.108 + 0.012 + 0.016 + 0.064) = 0.4$$


$$\mathbf{P}(\text{Cavity} | \text{toothache}) = \alpha \mathbf{P}(\text{Cavity}, \text{toothache}) \\ = \alpha [\mathbf{P}(\text{Cavity}, \text{toothache}, \text{catch}) + \mathbf{P}(\text{Cavity}, \text{toothache}, \neg \text{catch})] \\ = \alpha [\langle 0.108, 0.016 \rangle + \langle 0.012, 0.064 \rangle] \\ = \alpha [\langle 0.12, 0.08 \rangle] = [\langle 0.6, 0.4 \rangle]$$

Can we represent full joint probability distribution more compactly?

We can exploit **(absolute) independence** of random variables:

$$\mathbf{P(X|Y) = P(X) \text{ or } P(Y|X) = P(Y) \text{ or } P(X,Y) = P(X).P(Y)}$$

Example: two dice $P(\text{Die}_1 = 5, \text{Die}_2 = 3) = P(\text{Die}_1 = 5).P(\text{Die}_2 = 3)$

More frequently, two variables X and Y are **conditionally independent** given a third variable Z:

$$\mathbf{P(X|Y,Z) = P(X|Z) \text{ or } P(Y|X,Z) = P(Y|Z) \text{ or } P(X,Y|Z) = P(X|Z) P(Y|Z)}$$

Example:

$P(\text{Catch}, \text{Toothache} | \text{Cavity}) = P(\text{Catch} | \text{Cavity}). P(\text{Toothache} | \text{Cavity})$

Toothache and Catch are independent given information about cavity.

One big table can be represented using several smaller tables.

Random Boolean variables:

$P_{i,j}$ – pit at square (i,j)

$B_{i,j}$ – breeze at square (i,j)

(only for the observed squares $B_{1,1}$, $B_{1,2}$ a $B_{2,1}$).

Query to be answered:

$P(P_{1,3} \mid \text{known}, b)$.

where we have **evidence**:

$$b = \neg b_{1,1} \wedge b_{1,2} \wedge b_{2,1}$$

$$\text{known} = \neg p_{1,1} \wedge \neg p_{1,2} \wedge \neg p_{2,1}$$

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 B OK	2,2	3,2	4,2
1,1 OK	2,1 B OK	3,1	4,1

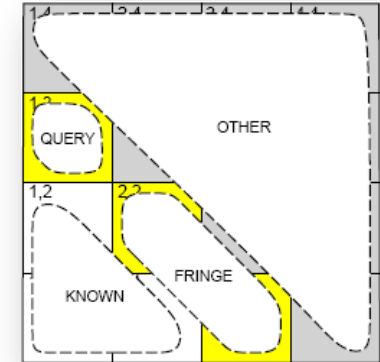
Answer can be computed by enumeration of the full joint probability distribution:

$$P(P_{1,3} \mid \text{known}, b) = \alpha \sum_{\text{unknown}} P(P_{1,3}, \text{unknown}, \text{known}, b)$$

where unknown be the variables $P_{i,j}$ except $P_{1,3}$ and known.

However there are $2^{12} = 4096$ terms! Can we do it better (faster)?

unknown = fringe \cup other
 from conditional independence we get:
 $P(b \mid P_{1,3}, \text{known}, \text{unknown}) = P(b \mid P_{1,3}, \text{known}, \text{fringe})$



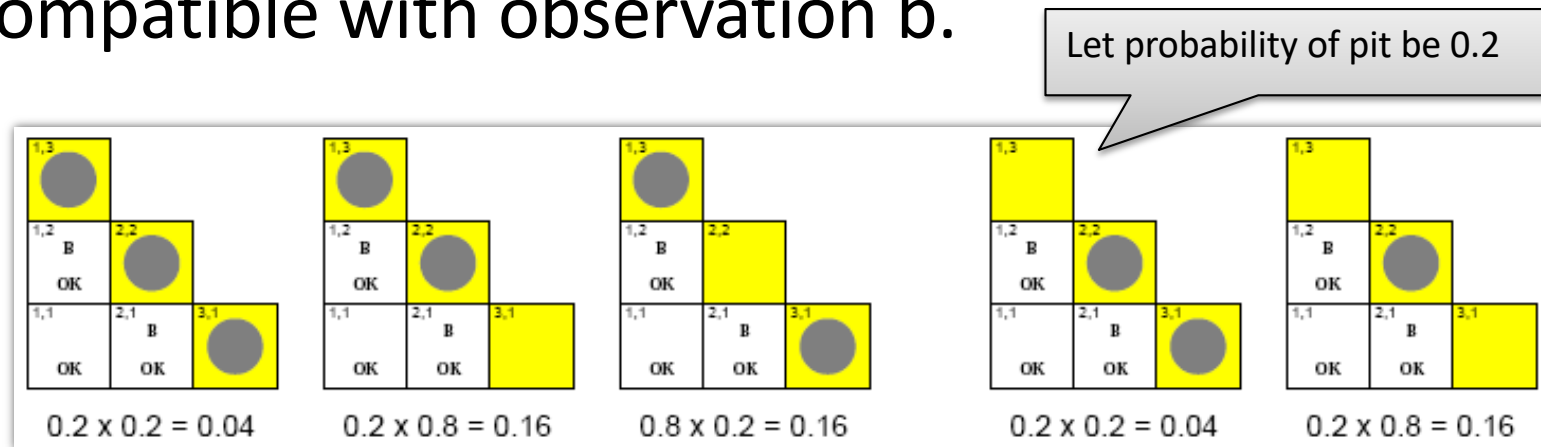
product rule $P(X,Y) = P(X \mid Y) P(Y)$

$$\begin{aligned}
 & \mathbf{P}(P_{1,3} \mid \text{known}, b) \\
 &= \alpha \sum_{\text{unknown}} \mathbf{P}(P_{1,3}, \text{known}, \text{unknown}, b) \\
 &= \alpha \sum_{\text{unknown}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{unknown}) * \mathbf{P}(P_{1,3}, \text{known}, \text{unknown}) \\
 &= \alpha \sum_{\text{fringe}} \sum_{\text{other}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{fringe}, \text{other}) * \mathbf{P}(P_{1,3}, \text{known}, \text{fringe}, \text{other}) \\
 &= \alpha \sum_{\text{fringe}} \sum_{\text{other}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{fringe}) * \mathbf{P}(P_{1,3}, \text{known}, \text{fringe}, \text{other}) \\
 &= \alpha \sum_{\text{fringe}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{fringe}) * \sum_{\text{other}} \mathbf{P}(P_{1,3}, \text{known}, \text{fringe}, \text{other}) \\
 &= \alpha \sum_{\text{fringe}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{fringe}) * \sum_{\text{other}} \mathbf{P}(P_{1,3}) P(\text{known}) P(\text{fringe}) P(\text{other}) \\
 &= \alpha P(\text{known}) \mathbf{P}(P_{1,3}) \sum_{\text{fringe}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{fringe}) P(\text{fringe}) \sum_{\text{other}} P(\text{other}) \\
 &= \alpha' \mathbf{P}(P_{1,3}) \sum_{\text{fringe}} \mathbf{P}(b \mid P_{1,3}, \text{known}, \text{fringe}) P(\text{fringe})
 \end{aligned}$$

$\alpha' = \alpha \cdot P(\text{known})$
 $\sum_{\text{other}} P(\text{other}) = 1$

$$P(P_{1,3} \mid \text{known}, b) = \alpha' P(P_{1,3}) \sum_{\text{fringe}} P(b \mid P_{1,3}, \text{known}, \text{fringe}) P(\text{fringe})$$

Let us explore possible models (values) of **fringe** that are compatible with observation b.



$$P(P_{1,3} \mid \text{known}, b)$$

$$= \alpha' \langle 0.2 (0.04 + 0.16 + 0.16), 0.8 (0.04 + 0.16) \rangle$$

$$= \langle 0.31, 0.69 \rangle$$

$$P(P_{2,2} \mid \text{known}, b) = \langle 0.86, 0.14 \rangle$$

Definitely avoid square (2,2)!





Recall the product rule

$$P(a \wedge b) = P(a | b) P(b) = P(b | a) P(a)$$

We can deduce a so called **Bayes' rule** (law or theorem):

$$P(a | b) = P(b | a) P(a) / P(b)$$

in general:

$$P(Y | X) = P(X | Y) P(Y) / P(X) = \alpha P(X | Y) P(Y)$$

It looks like two steps backward as now we need to know $P(X | Y)$, $P(Y)$, $P(X)$.

But these are the values that we frequently have.

$$P(\text{cause} | \text{effect}) = P(\text{effect} | \text{cause}) P(\text{cause}) / P(\text{effect})$$

- $P(\text{effect} | \text{cause})$ describes the **causal direction**
- $P(\text{cause} | \text{effect})$ describes the **diagnostic direction**

If all the effects are conditionally independent given the cause variable, we get:

$$P(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = P(\text{Cause}) \prod_i P(\text{Effect}_i | \text{Cause})$$

Such a probability distribution is called a **naive Bayes model**

(it is often used even in cases where the “effect” variables are not actually conditionally independent given the value of the cause variable).

Medical diagnosis

- from past cases we know $P(\text{symptoms} | \text{disease})$, $P(\text{disease})$, $P(\text{symptoms})$
- for a new patient we know symptoms and looking for diagnosis $P(\text{disease} | \text{symptoms})$

Example:

- meningitis causes a stiff neck 70% of the time
- the prior probability of meningitis is 1/50 000
- the prior probability of stiff neck is 1%

What is the probability that a patient having a stiff neck has meningitis?

$$P(m|s) = P(s|m).P(m) / P(s) = 0.7 * 1/50000 / 0.01 = 0.0014$$

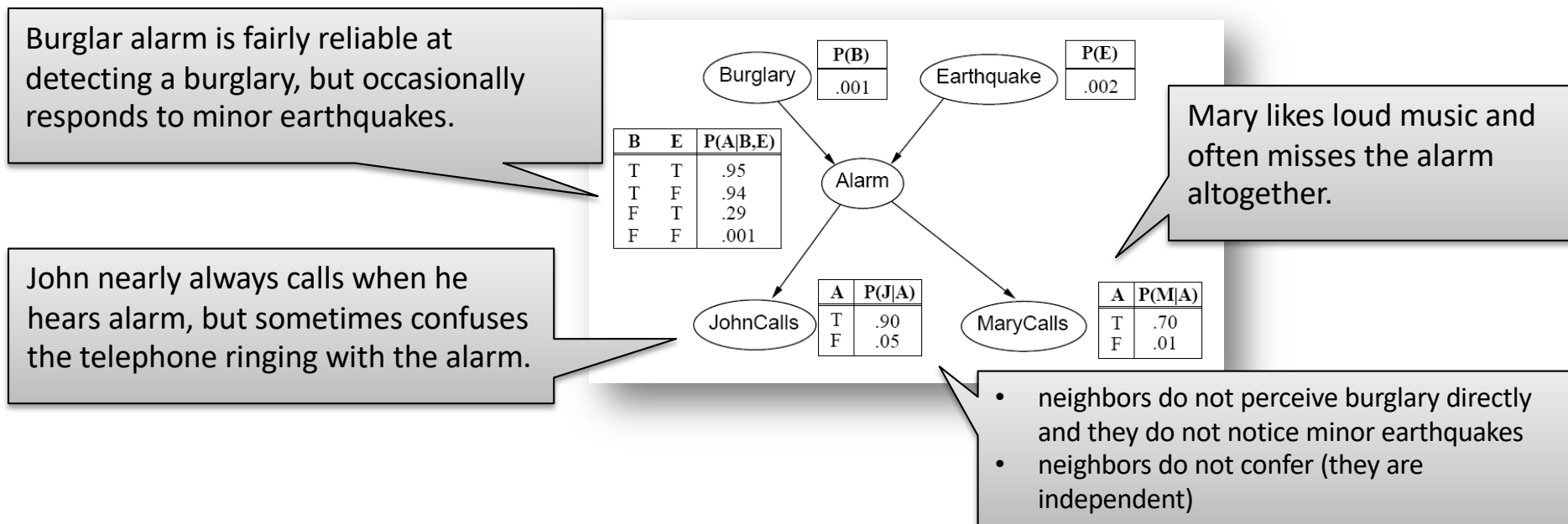
Why the conditional probability for the diagnostic direction is not stored directly?

- diagnostic knowledge is often more fragile than causal knowledge
- for example, if there is a sudden epidemic of meningitis, the unconditional probability of meningitis $P(m)$ will go up so $P(m|s)$ should also go up while the causal relation $P(s|m)$ is unaffected by the epidemic, as it reflects how meningitis works

How to represent efficiently any full joint probability distribution by exploiting conditional independence?

Bayesian network specifies **conditional independence relationships** among random variables using a directed acyclic graph (DAG)

- nodes correspond to random variables
- predecessors of nodes are called parents
- each node X has a **conditional probability distribution** $P(X \mid \text{Parents}(X))$



The Bayesian network represents the full joint probability distribution.

$$P(x_1, \dots, x_n) = \prod_i P(x_i \mid \text{parents}(X_i))$$

Nodes:

determine the set of random variables that are required to model the domain and order them

- any order will work, but the resulting networks will be different
- a recommended order is such that causes precede effects (leads to smaller networks and easier-to-fill CPTs)

Arcs:

choose variables X_i in a given order from 1 to n

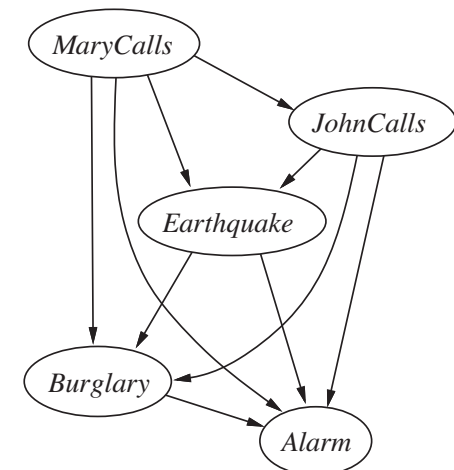
- in the set $\{X_1, \dots, X_{i-1}\}$ choose a minimal set of parents for X_i , such that $\mathbf{P}(X_i \mid \text{Parents}(X_i)) = \mathbf{P}(X_i \mid X_{i-1}, \dots, X_1)$ holds
- for each parent insert a link from the parent to X_i
- write down the conditional probability table $\mathbf{P}(X_i \mid \text{Parents}(X_i))$

Why does it work?

$$P(x_1, \dots, x_n) = \prod_i P(x_i \mid x_{i-1}, \dots, x_1) \quad \text{(chain rule)}$$

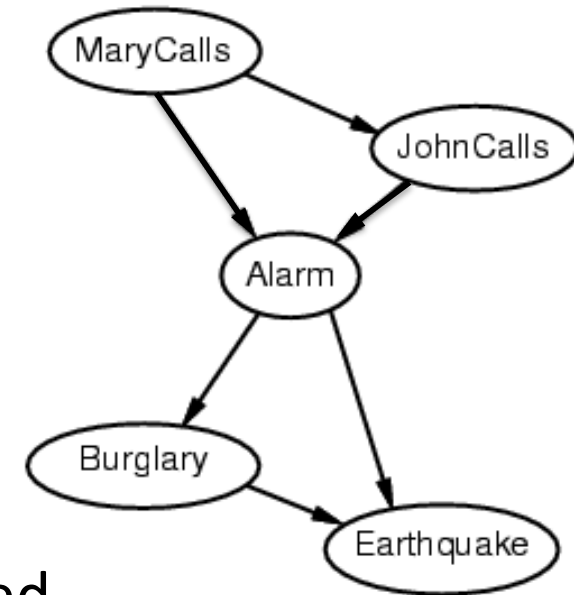
$$\mathbf{P}(X_i \mid X_{i-1}, \dots, X_1) = \mathbf{P}(X_i \mid \text{Parents}(X_i))$$

$$\text{where } \text{Parents}(X_i) \subseteq \{X_{i-1}, \dots, X_1\}$$



Let us use the following order of random variables:
MarryCalls, JohnCalls, Alarm, Burglary, Earthquake

- MarryCalls has no parents
- if Marry calls then the alarm is probably active which would make it more likely that John calls
- alarm is probably active if Marry or John calls
- if we know the alarm state then the calls from Marry and John do not influence whether the burglary happened



$$P(\text{Burglary} \mid \text{Alarm, JohnCalls, MarryCalls}) = P(\text{Burglary} \mid \text{Alarm})$$

- the alarm is an earthquake detector of sorts, but if there was a burglary then it explains the alarm and the probability of an earthquake is only slightly above normal

We introduced Bayesian networks to **do inference** – to deduce posterior probability of some variable(s) **X** from the query given the values **e** of observed variables (evidence), while having the other variables **Y** hidden.

$$P(X | e) = \alpha P(X, e) = \alpha \sum_y P(X, e, y)$$

where $P(X, e, y)$ is computed as follows

$$P(x_1, \dots, x_n) = \prod_i P(x_i | \text{parents}(X_i))$$

Example:

Assume a query about the probability of burglary when both Marry and John calls

$$P(b | j, m)$$

$$= \alpha \sum_e \sum_a P(b) P(e) P(a | b, e) P(j | a) P(m | a)$$

$$= \alpha P(b) \sum_e P(e) \sum_a P(a | b, e) P(j | a) P(m | a)$$

The structure of computation can be described using a tree structure.

- it is very similar to solving CSPs and SAT

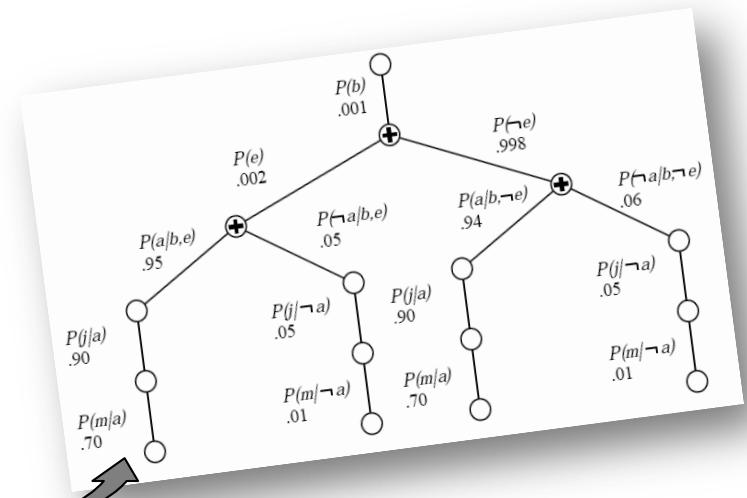
Notice that some parts are repeated!

```

function ENUMERATION-ASK(X, e, bn) returns a distribution over X
inputs: X, the query variable
       e, observed values for variables E
       bn, a Bayes net with variables {X} ∪ E ∪ Y /* Y = hidden variables */

Q(X) ← a distribution over X, initially empty
for each value x_i of X do
  Q(x_i) ← ENUMERATE-ALL(bn.VARS, e_{x_i})
           where e_{x_i} is e extended with X = x_i
return NORMALIZE(Q(X))

function ENUMERATE-ALL(vars, e) returns a real number
if EMPTY?(vars) then return 1.0
Y ← FIRST(vars)
if Y has value y in e
  then return P(y | parents(Y)) × ENUMERATE-ALL(REST(vars), e)
else return ∑_y P(y | parents(Y)) × ENUMERATE-ALL(REST(vars), e_y)
           where e_y is e extended with Y = y
    
```



Enumeration repeats the same parts of the computation.

We can remember the result and reuse it later (**dynamic programming**).

$$\begin{aligned}
 P(B \mid j,m) &= \alpha P(B) \sum_e P(e) \sum_a P(a \mid B,e) P(j \mid a) P(m \mid a) \\
 &= \alpha f_1(B) \sum_e f_2(E) \sum_a f_3(A,B,E) f_4(A) f_5(A)
 \end{aligned}$$

Factors f_i are matrices (tables) corresponding to CPTs.

Evaluation will be done **from right to left**.

- **the product of factors** corresponds to the pointwise product (it is not a multiplication of matrices)
- **summing out a variable** is done by adding up the sub-matrices formed by fixing the variable to each of its values in turn

Notes:

- The algorithm works for any ordering of variables.
- The complexity is given by the size of the largest factor constructed during the operation of the algorithm.
- Eliminate whichever variable minimizes the size of the next factor to be constructed (heuristic).

```

function ELIMINATION-ASK( $X, e, bn$ ) returns a distribution over  $X$ 
inputs:  $X$ , the query variable
        $e$ , observed values for variables  $E$ 
        $bn$ , a Bayesian network specifying joint distribution  $P(X_1, \dots, X_n)$ 

factors  $\leftarrow []$ 
for each var in ORDER( $bn.VARS$ ) do
  factors  $\leftarrow$  [MAKE-FACTOR( $var, e$ ) | factors]
  if var is a hidden variable then factors  $\leftarrow$  SUM-OUT( $var, factors$ )
return NORMALIZE(POINTWISE-PRODUCT(factors))
    
```

The pointwise **product** of two factors yields a new factor whose variables are the union of the variables from the original factors.

$$f(X_1, \dots, X_j, Y_1, \dots, Y_k, Z_1, \dots, Z_l) = f(X_1, \dots, X_j, Y_1, \dots, Y_k) \cdot f(Y_1, \dots, Y_k, Z_1, \dots, Z_l)$$

A	B	f₁(A,B)		B	C	f₂(B,C)		=	A	B	C	f₃(A,B,C)
T	T	0.3		T	T	0.2			T	T	T	0.06 = 0.3*0.2
T	F	0.7	*	T	F	0.8			T	T	F	0.24 = 0.3*0.8
F	T	0.9		F	T	0.6			T	F	T	0.42 = 0.7*0.6
F	F	0.1		F	F	0.4			T	F	F	0.28 = 0.7*0.4
									F	T	T	0.18 = 0.9*0.2
									F	T	F	0.72 = 0.9*0.8
									F	F	T	0.06 = 0.1*0.6
									F	F	F	0.04 = 0.1*0.4

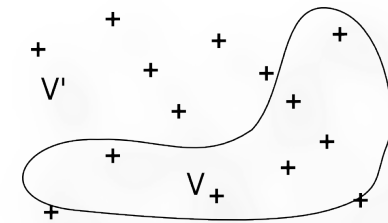
Then we **sum out** a variable to eliminate it: $\sum_a f(A,B,C) = f(B,C)$

A	B	C	f₃(A=T,B,C)		A	B	C	f₃(A=F,B,C)		=	B	C	f₄(B,C)
T	T	T	0.06		F	T	T	0.18			T	T	0.24
T	T	F	0.24	+	F	T	F	0.72			T	F	0.96
T	F	T	0.42		F	F	T	0.06			F	T	0.48
T	F	F	0.28		F	F	F	0.04			F	F	0.31

Exact inference is intractable for large, multiply connected networks so we may need to consider approximate inference methods based on **Monte Carlo** algorithms.

Monte Carlo algorithms are used to estimate quantities that are difficult to calculate exactly.

- generate many samples
- use statistics to estimate the quantity
- more samples = more accuracy



A **sample** corresponds to an instantiation of random variables. Each sample should be generated from a known probability distribution (given by CPTs in the Bayesian network).

- nodes (variables) are taken in topological order
- the probability distribution is conditioned on the values already assigned to parents
- generate a sample value based on this distribution

```

function PRIOR-SAMPLE(bn) returns an event sampled from the prior specified by bn
inputs: bn, a Bayesian network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 

x ← an event with n elements
foreach variable  $X_i$  in  $X_1, \dots, X_n$  do
     $x[i]$  ← a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$ 
return x
  
```

Let N be the number of samples and $N(x_1, \dots, x_n)$ be the number of occurrences of event x_1, \dots, x_n , then

$$P(x_1, \dots, x_n) = \lim_{N \rightarrow \infty} (N(x_1, \dots, x_n) / N)$$

However, we are looking for $P(X | e)$!

Rejection sampling

From all the generated samples, we will select only those consistent with the evidence e (other samples are rejected).

$$P(X | e) \approx N(X, e) / N(e)$$

Major weakness: **rejecting too many samples**

Likelihood weighting

Instead of rejecting inconsistent samples, it seems more efficient to generate only samples consistent with evidence e .

- Fix the values for the evidence variables E and sample only the non-evidence variables.
- The probability of obtaining a sample is $P(z, e) = \prod_i P(z_i | \text{parents}(z_i))$
- But this is not what we want! We miss $w(z, e) = \prod_j P(e_j | \text{parents}(e_j))$.
- Hence each sample is **weighted** as follows:
 $P(X | e) \approx \alpha N(X, e) w(X, e)$

Major weakness: **too small weights**

```
function REJECTION-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
inputs:  $X$ , the query variable
        $e$ , observed values for variables  $E$ 
        $bn$ , a Bayesian network
        $N$ , the total number of samples to be generated
local variables:  $N$ , a vector of counts for each value of  $X$ , initially zero

for  $j = 1$  to  $N$  do
   $x \leftarrow$  PRIOR-SAMPLE( $bn$ )
  if  $x$  is consistent with  $e$  then
     $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
return NORMALIZE( $N$ )
```



```
function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
inputs:  $X$ , the query variable
        $e$ , observed values for variables  $E$ 
        $bn$ , a Bayesian network specifying joint distribution  $P(X_1, \dots, X_n)$ 
        $N$ , the total number of samples to be generated
local variables:  $W$ , a vector of weighted counts for each value of  $X$ , initially zero

for  $j = 1$  to  $N$  do
   $x, w \leftarrow$  WEIGHTED-SAMPLE( $bn, e$ )
   $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
return NORMALIZE( $W$ )
```

```
function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
 $w \leftarrow 1$ ;  $x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
foreach variable  $X_i$  in  $X_1, \dots, X_n$  do
  if  $X_i$  is an evidence variable with value  $x_i$  in  $e$ 
    then  $w \leftarrow w \times P(X_i = x_i | \text{parents}(X_i))$ 
    else  $x[i] \leftarrow$  a random sample from  $P(X_i | \text{parents}(X_i))$ 
return  $x, w$ 
```



Probability theory is a formal tool to handle **uncertainty**.

The **full joint probability distribution** specifies the probability of each complete assignment of values to random variables (possible worlds).

It is usually too large but **independence** relations between subsets of random variables allows us to factor it into smaller joint or conditional distributions.

Bayesian network is such a compact representation.

We can use it to answer **queries $P(\mathbf{X} | \mathbf{E})$** about probability distributions of variables \mathbf{X} under evidence \mathbf{E} .

- exact methods (enumeration, variable elimination)
- approximate methods (rejection sampling, likelihood weighting)



© 2020 Roman Barták

Department of Theoretical Computer Science and Mathematical Logic

bartak@ktiml.mff.cuni.cz