

Artificial Intelligence (CS303)

Lecture 12: Representing and Inference with Uncertainty

Hints for this lecture

- An agent can seldom precisely know the state, knowledge should be represented such that wise decisions/actions can still be made.

I. Uncertainty and Rational Decisions

The World is Uncertain

- Things don't always happen with simple true or false.
- We never know what “state” we are in exactly, because the world is only partially observable.
- We (agents) seldom make decisions with full certainty, while more often make **rational** decision based on **utility**.

Language	Ontological Commitment	Epistemological Commitment
Propositional logic	facts	true/false/unknown
First-order logic	facts, objects, relations	true/false/unknown
Probability theory	facts	degree of belief
Fuzzy logic	facts	degree of truth known interval value

Alternative to Logic

- Utility theory: Assign utility to each state/actions
- Probability theory: Summarize the uncertainty associated with each state
- Rational Decisions: Maximize the **expected utility** (Probability + Utility)
- Thus we need to represent states in the language of probability.

II. Basic Probability Theory and Its Use

Basic Probability Theory and Its Use

- Joint probability distribution specifies probability of every atomic event.
- Queries can be answered by summing over atomic events.

	<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

Basic Probability Theory and Its Use

Prior probability

Prior or unconditional probabilities of propositions

e.g., $P(\text{Cavity} = \text{true}) = 0.1$ and $P(\text{Weather} = \text{sunny}) = 0.72$

correspond to belief prior to arrival of any (new) evidence

Probability distribution gives values for all possible assignments:

$\mathbf{P}(\text{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$ (normalized, i.e., sums to 1)

Joint probability distribution for a set of r.v.s gives the probability of every atomic event on those r.v.s (i.e., every sample point)

$\mathbf{P}(\text{Weather}, \text{Cavity}) =$ a 4×2 matrix of values:

<i>Weather =</i>	<i>sunny</i>	<i>rain</i>	<i>cloudy</i>	<i>snow</i>
<i>Cavity = true</i>	0.144	0.02	0.016	0.02
<i>Cavity = false</i>	0.576	0.08	0.064	0.08

Every question about a domain can be answered by the joint distribution because every event is a sum of sample points

Basic Probability Theory and Its Use

Conditional probability

Conditional or posterior probabilities

e.g., $P(\text{cavity}|\text{toothache}) = 0.8$

i.e., **given that toothache is all I know**

NOT “if *toothache* then 80% chance of *cavity*”

(Notation for conditional distributions:

$\mathbf{P}(\text{Cavity}|\text{Toothache}) = 2\text{-element vector of 2-element vectors})$

If we know more, e.g., *cavity* is also given, then we have

$P(\text{cavity}|\text{toothache}, \text{cavity}) = 1$

Note: the less specific belief **remains valid** after more evidence arrives, but is not always **useful**

New evidence may be irrelevant, allowing simplification, e.g.,

$P(\text{cavity}|\text{toothache}, \text{49ersWin}) = P(\text{cavity}|\text{toothache}) = 0.8$

This kind of inference, sanctioned by domain knowledge, is crucial

Basic Probability Theory and Its Use

Conditional probability

Definition of conditional probability:

$$P(a|b) = \frac{P(a \wedge b)}{P(b)} \text{ if } P(b) \neq 0$$

Product rule gives an alternative formulation:

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

A general version holds for whole distributions, e.g.,

$$\mathbf{P}(\textit{Weather}, \textit{Cavity}) = \mathbf{P}(\textit{Weather}|\textit{Cavity})\mathbf{P}(\textit{Cavity})$$

(View as a 4×2 set of equations, **not** matrix mult.)

Chain rule is derived by successive application of product rule:

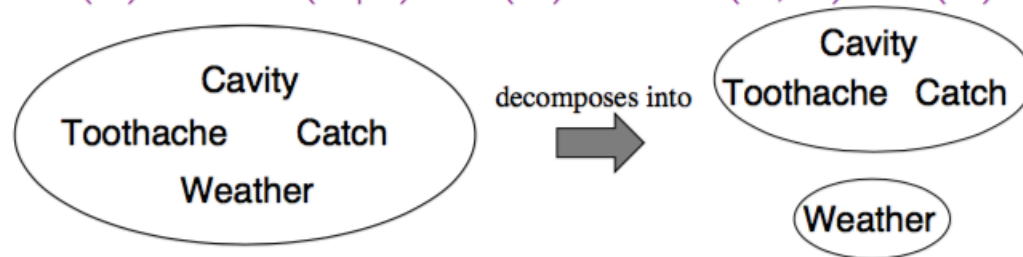
$$\begin{aligned} \mathbf{P}(X_1, \dots, X_n) &= \mathbf{P}(X_1, \dots, X_{n-1}) \mathbf{P}(X_n|X_1, \dots, X_{n-1}) \\ &= \mathbf{P}(X_1, \dots, X_{n-2}) \mathbf{P}(X_{n-1}|X_1, \dots, X_{n-2}) \mathbf{P}(X_n|X_1, \dots, X_{n-1}) \\ &= \dots \\ &= \prod_{i=1}^n \mathbf{P}(X_i|X_1, \dots, X_{i-1}) \end{aligned}$$

Basic Probability Theory and Its Use

Independence

A and B are independent iff

$$P(A|B) = P(A) \quad \text{or} \quad P(B|A) = P(B) \quad \text{or} \quad P(A, B) = P(A)P(B)$$



$$P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}, \textit{Weather}) \\ = P(\textit{Toothache}, \textit{Catch}, \textit{Cavity})P(\textit{Weather})$$

32 entries reduced to 12; for n independent biased coins, $2^n \rightarrow n$

Absolute independence powerful but rare

Dentistry is a large field with hundreds of variables,
none of which are independent. What to do?

Basic Probability Theory and Its Use

Bayes' Rule

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

$$\Rightarrow \text{Bayes' rule } P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

or in distribution form

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \alpha P(X|Y)P(Y)$$

Useful for assessing **diagnostic** probability from **causal** probability:

$$P(Cause|Effect) = \frac{P(Effect|Cause)P(Cause)}{P(Effect)}$$

Basic Probability Theory and Its Use

Bayes' Rule and conditional independence

$$\begin{aligned} & \mathbf{P}(\textit{Cavity}|\textit{toothache} \wedge \textit{catch}) \\ &= \alpha \mathbf{P}(\textit{toothache} \wedge \textit{catch}|\textit{Cavity})\mathbf{P}(\textit{Cavity}) \\ &= \alpha \mathbf{P}(\textit{toothache}|\textit{Cavity})\mathbf{P}(\textit{catch}|\textit{Cavity})\mathbf{P}(\textit{Cavity}) \end{aligned}$$

This is an example of a **naive Bayes** model:

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

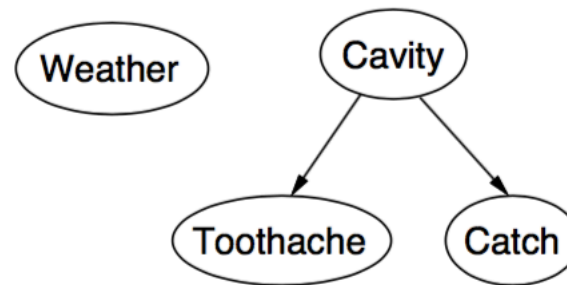


Total number of parameters is **linear** in n

III. Bayesian Networks

What is a BN?

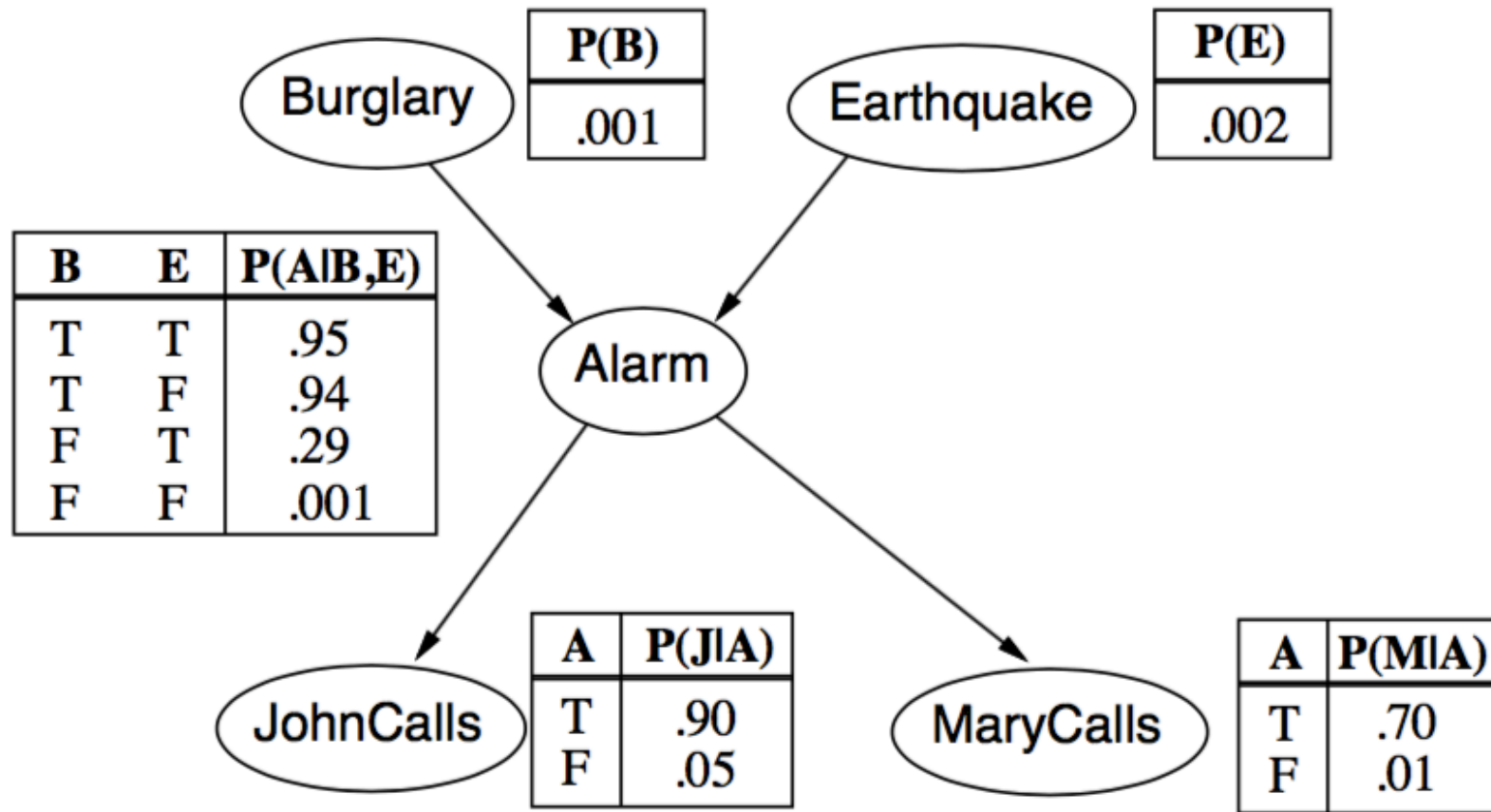
- A Directed Acyclic Graph (DAG).
- Each node is a random variable, associated with conditional distribution.
- Each arc (link) represent **direct influence** of a parent node to a child node.



Weather is independent of the other variables

Toothache and *Catch* are conditionally independent given *Cavity*

A Simple Example of BN



BN in the form of Conditional Probability Table

- More compact representation (compared to propositional logic).

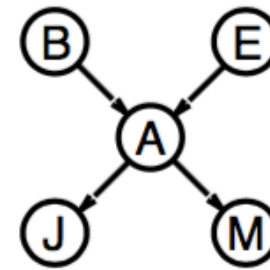
A CPT for Boolean X_i with k Boolean parents has 2^k rows for the combinations of parent values

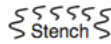


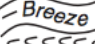
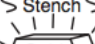

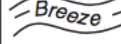
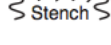
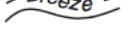



Each row requires one number p for $X_i = \text{true}$ (the number for $X_i = \text{false}$ is just $1 - p$)

If each variable has no more than k parents, the complete network requires $O(n \cdot 2^k)$ numbers

I.e., grows linearly with n , vs. $O(2^n)$ for the full joint distribution

For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers (vs. $2^5 - 1 = 31$)



4	 Stench		 Breeze	PIT
3		 Breeze  Stench  Gold	PIT	 Breeze
2	 Stench		 Breeze	
1	 START	 Breeze	PIT	 Breeze
	1	2	3	4

- Easier to utilize independence and conditional dependence relations to define the joint distribution.

How to construct a CPT for BN?

1. Choose an ordering of variables X_1, \dots, X_n
2. For $i = 1$ to n
 - add X_i to the network
 - select parents from X_1, \dots, X_{i-1} such that
$$\mathbf{P}(X_i | \text{Parents}(X_i)) = \mathbf{P}(X_i | X_1, \dots, X_{i-1})$$

This choice of parents guarantees the global semantics:

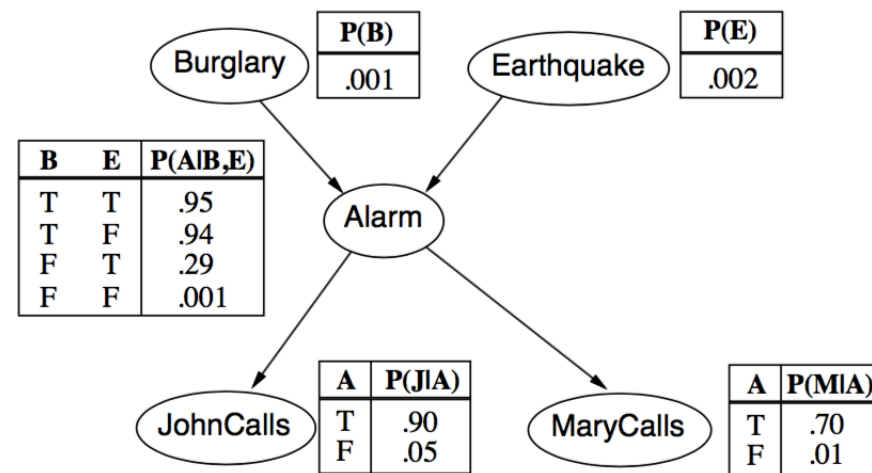
$$\begin{aligned}\mathbf{P}(X_1, \dots, X_n) &= \prod_{i=1}^n \mathbf{P}(X_i | X_1, \dots, X_{i-1}) \quad (\text{chain rule}) \\ &= \prod_{i=1}^n \mathbf{P}(X_i | \text{Parents}(X_i)) \quad (\text{by construction})\end{aligned}$$

IV. Exact Inference with Bayesian Networks

Inference with BN

- Given a Bayesian Network, and an (or some) observed events, which specifies the value for **evidence variables**, we want to know the probability distribution of one (or several) **query variables X**, $P(X \mid \text{events})$.

- E. g. : $P(\text{Burglary} \mid \text{JohnCalls} = \text{true}, \text{MaryCalls} = \text{true})$

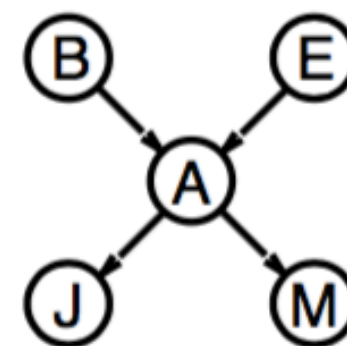


Enumeration

Simple query on the burglary network:

$$\begin{aligned} & \mathbf{P}(B|j, m) \\ &= \mathbf{P}(B, j, m) / P(j, m) \\ &= \alpha \mathbf{P}(B, j, m) \\ &= \alpha \sum_e \sum_a \mathbf{P}(B, e, a, j, m) \end{aligned}$$

Need to consider all values
of “hidden variables”,
e.g., *alarm*=true, *alarm*=false



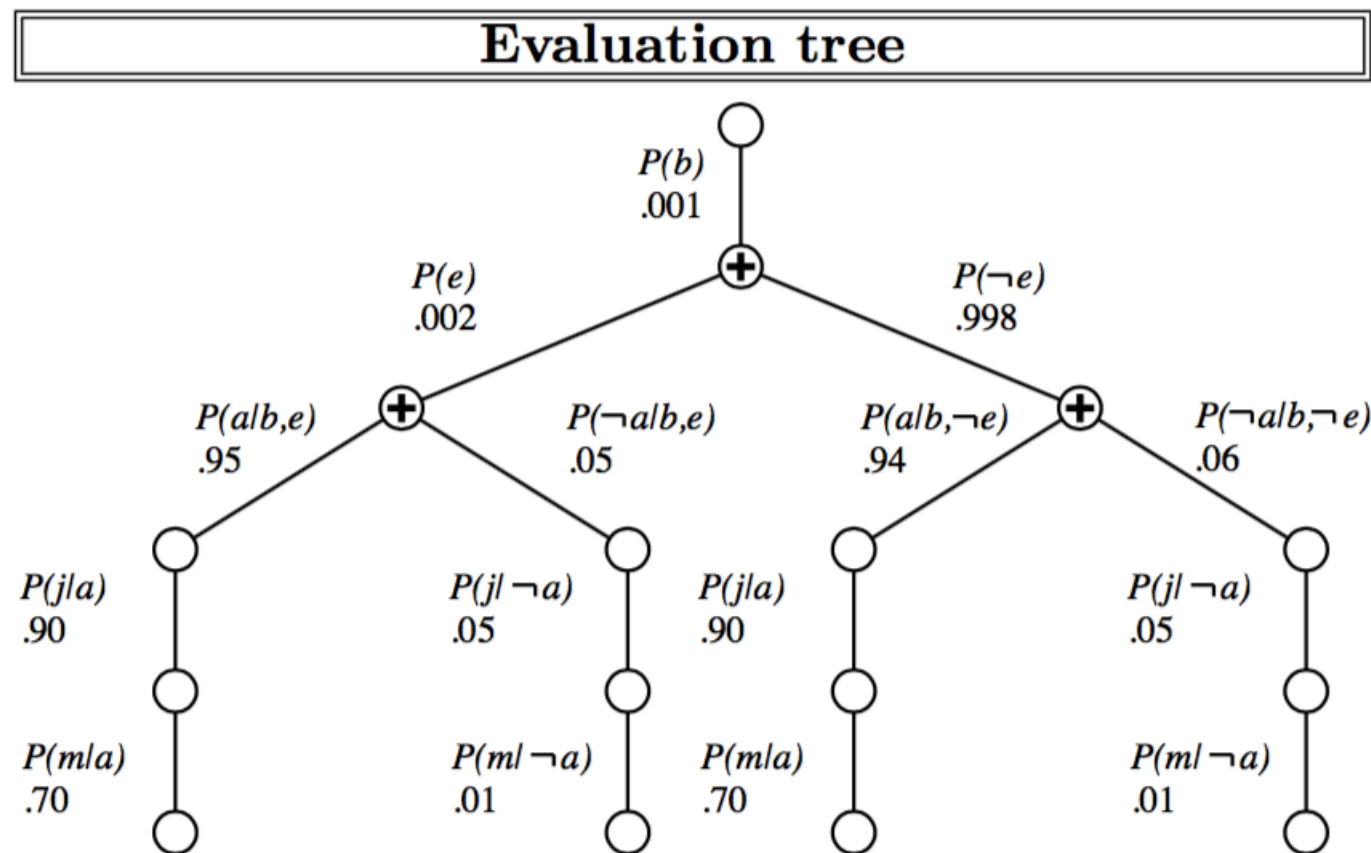
Rewrite full joint entries using product of CPT entries:

$$\begin{aligned} & \mathbf{P}(B|j, m) \\ &= \alpha \sum_e \sum_a \mathbf{P}(B)P(e)\mathbf{P}(a|B, e)P(j|a)P(m|a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e)P(j|a)P(m|a) \end{aligned}$$

Recursive depth-first enumeration: $O(n)$ space, $O(d^n)$ time

$d=2$ for Boolean
variables

Enumeration



Enumeration is inefficient: repeated computation
e.g., computes $P(j|a)P(m|a)$ for each value of e

Enumeration by Variable Elimination

Variable elimination: carry out summations right-to-left, storing intermediate results (**factors**) to avoid recomputation

$$\begin{aligned}\mathbf{P}(B|j, m) &= \alpha \underbrace{\mathbf{P}(B)}_B \sum_e \underbrace{P(e)}_E \sum_a \underbrace{\mathbf{P}(a|B, e)}_A \underbrace{P(j|a)}_J \underbrace{P(m|a)}_M \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e) P(j|a) f_M(a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e) f_J(a) f_M(a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a f_A(a, b, e) f_J(a) f_M(a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) f_{\bar{A}JM}(b, e) \text{ (sum out } A) \\ &= \alpha \mathbf{P}(B) f_{\bar{E}\bar{A}JM}(b) \text{ (sum out } E) \\ &= \alpha f_B(b) \times f_{\bar{E}\bar{A}JM}(b)\end{aligned}$$

V. Approximate Inference with Bayesian Networks

Basic Idea

- Sampling/Monte Carlo/Stochastic Simulation...

Basic idea:

- 1) Draw N samples from a sampling distribution S
- 2) Compute an approximate posterior probability \hat{P}
- 3) Show this converges to the true probability P

Outline:

- Sampling from an empty network
- Rejection sampling: reject samples disagreeing with evidence
- Likelihood weighting: use evidence to weight samples
- Markov chain Monte Carlo (MCMC): sample from a stochastic process whose stationary distribution is the true posterior



Sampling from an empty network

```
function PRIOR-SAMPLE(bn) returns an event sampled from bn  
  inputs: bn, a belief network specifying joint distribution  $P(X_1, \dots, X_n)$   
   $\mathbf{x} \leftarrow$  an event with  $n$  elements  
  for  $i = 1$  to  $n$  do  
     $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$   
    given the values of  $\text{Parents}(X_i)$  in  $\mathbf{x}$   
  return  $\mathbf{x}$ 
```

Assume the joint distribution could be easily sampled

Rejection Sampling

$\hat{P}(X|e)$ estimated from samples agreeing with e

```
function REJECTION-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $N$ , a vector of counts over  $X$ , 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( $N[X]$ )
```

In case distribution of variable e is difficult to sample

E.g., estimate $P(Rain|Sprinkler = true)$ using 100 samples
27 samples have $Sprinkler = true$
Of these, 8 have $Rain = true$ and 19 have $Rain = false$.

$\hat{P}(Rain|Sprinkler = true) = \text{NORMALIZE}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle$

Similar to a basic real-world empirical estimation procedure

$$\begin{aligned}\hat{P}(X|e) &= \alpha N_{PS}(X, e) && \text{(algorithm defn.)} \\ &= N_{PS}(X, e) / N_{PS}(e) && \text{(normalized by } N_{PS}(e)) \\ &\approx P(X, e) / P(e) && \text{(property of PRIOR-SAMPLE)} \\ &= P(X|e) && \text{(defn. of conditional probability)}\end{aligned}$$

Hence rejection sampling returns consistent posterior estimates

Problem: hopelessly expensive if $P(e)$ is small

$P(e)$ drops off exponentially with number of evidence variables!

Likelihood Weighting

Idea: fix evidence variables, sample only nonevidence variables, and weight each sample by the likelihood it accords the evidence

Fix the evidence variables to reduce the sampling space

```
function LIKELIHOOD-WEIGHTING( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $P(X|\mathbf{e})$   
  local variables:  $\mathbf{W}$ , a vector of weighted counts over  $X$ , initially zero  
  for  $j = 1$  to  $N$  do  
     $\mathbf{x}, w \leftarrow \text{WEIGHTED-SAMPLE}(bn)$   
     $\mathbf{W}[x] \leftarrow \mathbf{W}[x] + w$  where  $x$  is the value of  $X$  in  $\mathbf{x}$   
  return NORMALIZE( $\mathbf{W}[X]$ )
```

```
function WEIGHTED-SAMPLE( $bn, \mathbf{e}$ ) returns an event and a weight  
   $\mathbf{x} \leftarrow$  an event with  $n$  elements;  $w \leftarrow 1$   
  for  $i = 1$  to  $n$  do  
    if  $X_i$  has a value  $x_i$  in  $\mathbf{e}$   
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$   
      else  $x_i \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$   
  return  $\mathbf{x}, w$ 
```

Markov Chain Monte Carlo (MCMC)

“State” of network = current assignment to all variables.

Generate next state by sampling one variable given Markov blanket
Sample each variable in turn, keeping evidence fixed

```
function MCMC-Ask( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $P(X|\mathbf{e})$ 
  local variables:  $\mathbf{N}[X]$ , a vector of counts over  $X$ , initially zero
                   $\mathbf{Z}$ , the nonevidence variables in  $bn$ 
                   $\mathbf{x}$ , the current state of the network, initially copied from  $\mathbf{e}$ 

  initialize  $\mathbf{x}$  with random values for the variables in  $\mathbf{Y}$ 
  for  $j = 1$  to  $N$  do
    for each  $Z_i$  in  $\mathbf{Z}$  do
      sample the value of  $Z_i$  in  $\mathbf{x}$  from  $P(Z_i|mb(Z_i))$ 
        given the values of  $MB(Z_i)$  in  $\mathbf{x}$ 
       $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $\mathbf{N}[X]$ )
```

Further reduce the sampling space by only considering variables in Markov blanket

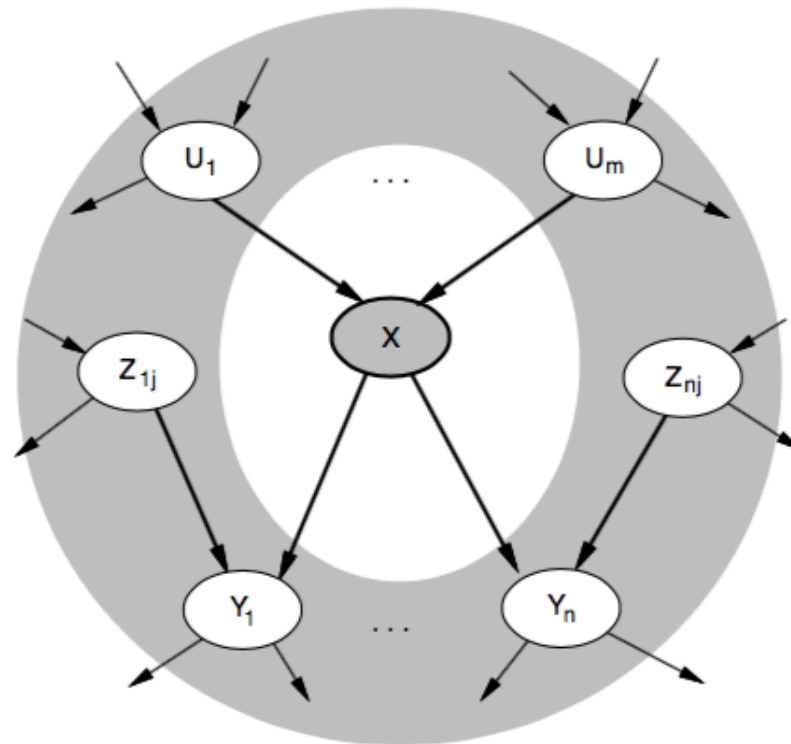
Markov Blanket

Can also choose a variable to sample at random each time

Markov Chain Monte Carlo (MCMC)

Markov blanket

Each node is conditionally independent of all others given its
Markov blanket: parents + children + children's parents



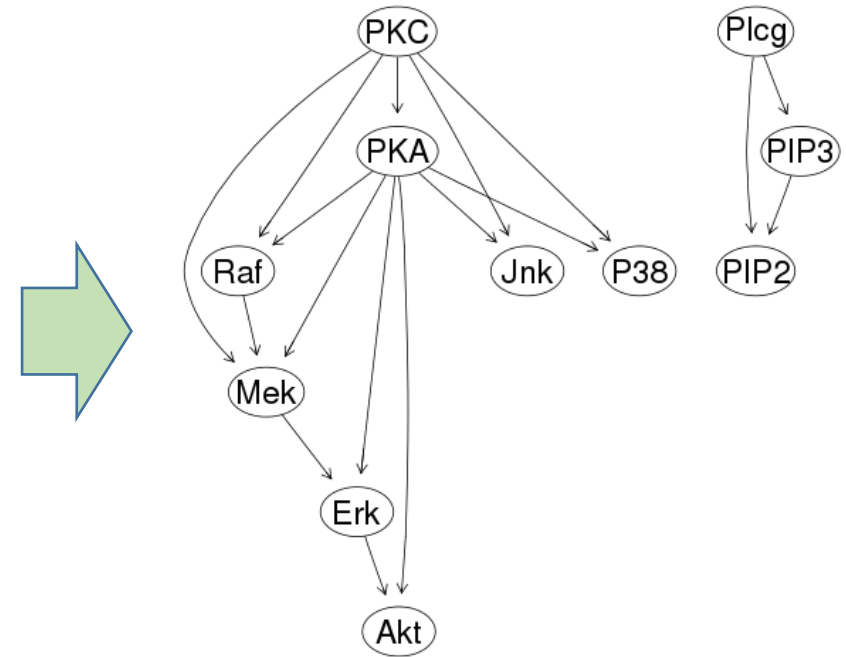
VI. How to construct a BN (or KB in general)

The challenge brought by KB

- Manual implement a KB is tedious (sometimes unaffordable)
- Can we obtain it **automatically** from raw data?

1	Raf	Mek	Plcg	PIP2	PIP3	Erk	Akt	PKA	PKC	P38	Jnk
2	26.4	13.2	8.82	18.3	58.8	6.61	17	414	17	44.9	40
3	35.9	16.5	12.3	16.8	8.13	18.6	32.5	352	3.37	16.5	61.5
4	59.4	44.1	14.6	10.2	13	14.9	32.5	403	11.4	31.9	19.5
5	73	82.8	23.1	13.5	1.29	5.83	11.8	528	13.7	28.6	23.1
6	33.7	19.8	5.19	9.73	24.8	21.1	46.1	305	4.66	25.7	81.3
7	18.8	3.75	17.6	22.1	10.9	11.9	25.7	610	13.7	49.1	57.8
8	44.9	36.5	10.4	132	16.3	8.66	17.9	835	15	35.9	18.1
9	47.4	15	14.6	30.5	17.5	20.2	45.3	466	6.44	24.4	20
10	104	61.5	10.6	21.1	41.8	11.5	23.5	445	29.2	61	25.3
11	21.1	21.5	1.88	205	43.7	13.2	135	213	14.6	26.7	101
12	16.4	16.4	14.5	17	11.2	21.9	34.6	449	20.4	44.9	24.1
13	74.3	22.9	7.5	15.5	26.2	20.9	36.5	389	31.9	71	35.5
14	85.1	39.6	8.9	64.9	11.7	6.67	12.2	528	17.9	44.1	118
15	36.8	29.2	5	9.06	15.5	17.9	17.9	400	14.6	41.4	151
16	29.2	24.1	10.2	16	46.6	8.82	7.23	500	9.73	19.1	7.91
17	50	13.8	11.9	13.2	11.3	18.1	27.9	392	56.2	77	1
18	26.2	26.7	21.3	10.9	14.7	9.06	37.9	89	40	65.5	1.42
19	39.6	38.2	10.5	92.2	22.7	27.6	31.3	223	16.5	30.8	7.64
20	30.5	19.8	7.5	133	15.7	19.1	36.2	319	24.1	37.2	17.2
21	49.1	16.5	13.3	141	25.7	31.1	62.1	710	14.2	27.4	22.1

Raw data

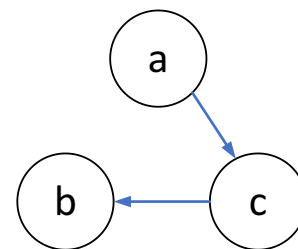
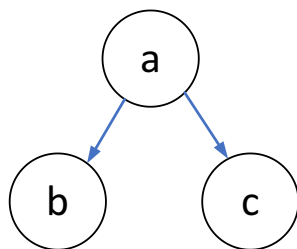
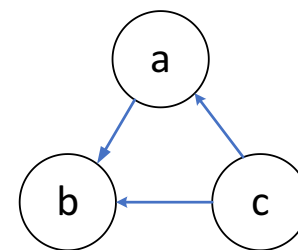
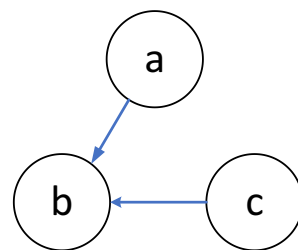
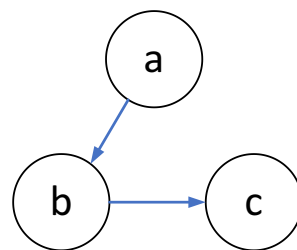
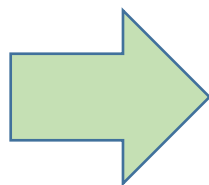


BN

Constructing a BN from data

- Structural Learning
- Parameter Estimation

1	a	b	c
2	0.366329	0.458928	0.937441
3	0.492245	0.829342	0.657999
4	0.521917	0.119465	0.956325
5	0.741545	0.490508	0.063304
6	0.60749	0.785813	0.62433
7	0.935561	0.542337	0.464547
8	0.771979	0.501432	0.25861
9	0.505258	0.398753	0.8417
10	0.147234	0.44209	0.986153
11	0.250739	0.646287	0.167685
12	0.149013	0.907165	0.162593
13	0.127727	0.451099	0.49503
14	0.382873	0.629639	0.016354



.....



Similar to Neural Networks

- Structural Learning: Identify the network structure
- Parameter Estimation: find VALUES for parameters associated with an edge
 - Depending on how you define the relationship between events/nodes
 - values in a CPT
 - parameters of a probability density function
- A machine learning or search problem again.

To be continued