

Lecture 3: Efficient probabilistic reasoning

Frederik Mallmann-Trenn

Department of Informatics
King's College London

(Version 1.8)



Today

- Introduction
- Probabilistic Reasoning I
- Probabilistic Reasoning II
- Sequential Decision Making
- Temporal Probabilistic Reasoning
- Game Theory
- Argumentation I
- Argumentation II
- (A peek at) Machine Learning
- AI & Ethics



Recap

- In Lecture 2 we talked about using probability theory to represent uncertainty in an agent's knowledge of the world.
- **Probability distribution** gives values for all possible values of a random variable:

$$\mathbf{P}(\textit{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$$

- **Product rule** relates joint probabilities to priors:

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

- **Bayes rule** relates conditional probabilities:

$$P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$



- **Chain rule** comes from applying the product rule multiple times:

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

- In terms of a more concrete example:

$$P(a, b, c) = P(a)P(b|a)P(c|b, a)$$

Recap

- Important step in probabilistic inference is establishing joint probabilities over all variables:

$$P(a, b, c)$$

$$P(\neg a, b, c)$$

$$\vdots$$

- With a full joint probability distribution over all the state variables
 - which we can either measure directly

$$P(\text{toothache}, \text{cavity}, \neg \text{catch})$$

or we can compute from conditionals

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

we can compute any specific values we want.



- Typically, we want the joint distribution of the **query variables** **Y** given specific values **e** for the **evidence variables** **E**

$$\begin{aligned} \mathbf{P}(\mathbf{Y}|\mathbf{E} = \mathbf{e}) &= \alpha \mathbf{P}(\mathbf{Y}, \mathbf{E} = \mathbf{e}) \\ &= \alpha \sum_{\mathbf{h}} \mathbf{P}(\mathbf{Y}, \mathbf{E} = \mathbf{e}, \mathbf{H} = \mathbf{h}) \end{aligned}$$

- The **H** are **hidden** variables, the ones we don't care about for this query.

- For

$$\begin{aligned} \mathbf{P}(\mathbf{Y}|\mathbf{E}=\mathbf{e}) &= \alpha \mathbf{P}(\mathbf{Y}, \mathbf{E}=\mathbf{e}) \\ &= \alpha \sum_{\mathbf{h}} \mathbf{P}(\mathbf{Y}, \mathbf{E}=\mathbf{e}, \mathbf{H}=\mathbf{h}) \end{aligned}$$

- The **H** are **hidden** variables, the ones we don't care about for this query.

- For example, we have:

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

- We want to know the probability of having a cavity, given we have a toothache.
- *Toothache* = *toothache* is the **evidence** variable.
- $\mathbf{P}(\textit{Cavity})$ is the **query**
- *Catch* is the **hidden** variable.

- For example, we have:

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

- We want to know the probability of having a cavity, given we have a toothache.
- *Toothache* = *toothache* is the **evidence** variable.
- $\mathbf{P}(\textit{Cavity})$ is the **query**
- *Catch* is the **hidden** variable.

$$\mathbf{P}(\textit{Cavity}|\textit{toothache}) = \mathbf{P}(\textit{Cavity}, \textit{catch}|\textit{toothache}) \\ + \mathbf{P}(\textit{Cavity}, \neg \textit{catch}|\textit{toothache})$$

- Computationally this is awkward.
 - n binary variables requires 2^n probabilities.
- **Conditional independence** is the mechanism we exploit to do reduce this number.
- *Catch* is **conditionally independent** of *Toothache* given *Cavity*

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

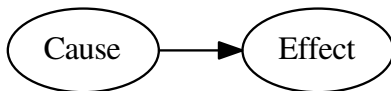
Catch does not depend on *Toothache* if we know that *Cavity* is true.

Recap

- Often easier to assess causal probabilities.

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

- Can visualise this as:



- **Bayes rule** allows us to use these easier probabilities.

Bayes' Rule & conditional independence

- So, in our running dentist example

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

- This is an example of a **naive Bayes** model, in which one assumes that

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



Naive Bayes

- Models n conditionally independent effects

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

- Total number of parameters is **linear** in n
- It is called 'naive', because it is oversimplifying: in many cases the 'effect' variables aren't actually conditionally independent given the cause variable. Example:
 - Cause*: it rained last night
 - Effect*₁: the street is wet
 - Effect*₂: I'm late for my class
 - If the streets were still wet, then an accident was more likely to happen and the caused traffic jam could be the reason for being late

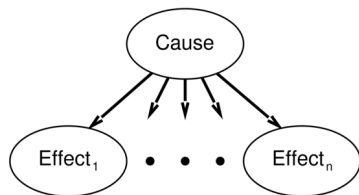
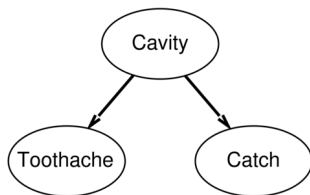


Naive Bayes

- Models two conditionally independent effects

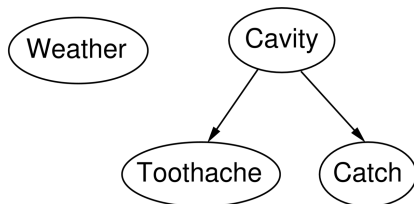
$$\mathbf{P}(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = \mathbf{P}(\text{Cause}) \prod_i \mathbf{P}(\text{Effect}_i | \text{Cause})$$

- Visualise as:



Bayesian networks

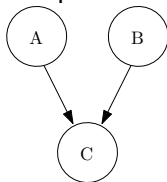
- A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions
- Topology of network encodes conditional independence assertions:



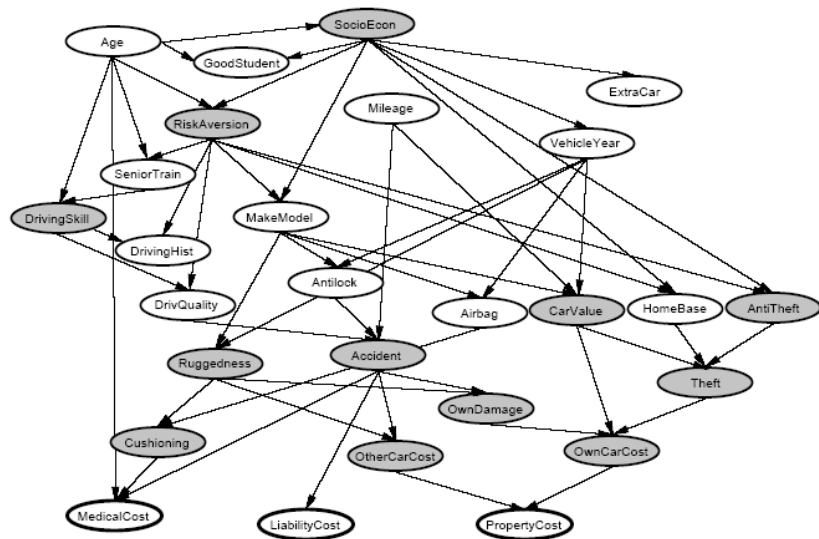
- *Weather* is independent of the other variables
- *Toothache* and *Catch* are conditionally independent given *Cavity*

Bayesian Networks

- **Bayesian networks** are a way to represent these dependencies:
 - Each node corresponds to a random variable (which may be discrete or continuous)
 - A directed edge (also called link or arrow) from node u to node v means that u is the **parent** of v .
 - The graph has no directed cycles (and hence is a directed acyclic graph, or DAG).
 - Each node v has a conditional probability $Pr(u \mid Parents(u))$ that quantifies the effect of the parent nodes
- Example: C depends on A and B , and A and B are independent.



Bayesian networks



(<http://www.igi.tugraz.at>)

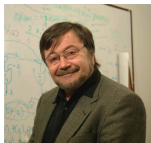
- How can represent the knowledge about the probabilities?
- Conditional distribution represented as a **conditional probability table** (CPT) giving the distribution over u for each combination of parent values

A	B	$P(C \mid A, B)$
T	T	0.2
T	F	0.123
F	T	0.9
F	F	0.51

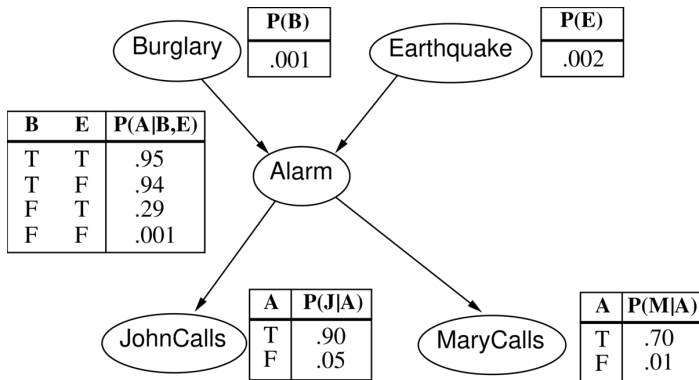
- Bayesian networks \neq Naive Bayes
- These are somewhat orthogonal. Naive Bayes might be used in Bayesian networks.
- Also don't confuse them with Bayes' rule

Bayesian networks

- An example (from California):
I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?
- Variables: *Burglar, Earthquake, Alarm, JohnCalls, MaryCalls*
- Network topology reflects “causal” knowledge:
 - A burglar can set the alarm off
 - An earthquake can set the alarm off
 - The alarm can cause Mary to call
 - The alarm can cause John to call



Bayesian networks



A note on CPTs

- The CPTs in the previous slide appear to be missing some values:

A	$P(J A)$
T	0.90
F	0.05

has two values rather than the four which would completely specify the relation between J and A .

- The table tells us that:

$$P(J = T|A = T) = 0.9$$

which means:

$$P(J = F|A = T) = 0.1$$

because $P(J = T|A = T) + P(J = F|A = T) = 1$



- Or, writing the values of J and A the other way:

$$P(j|a) = 0.9$$

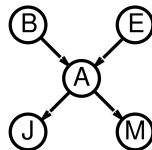
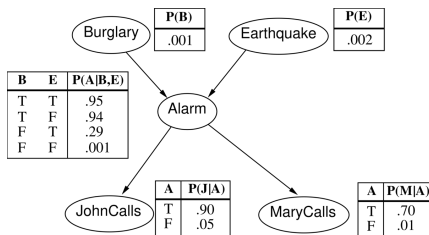
$$P(\neg j|a) = 0.1$$

because $P(j|a) + P(\neg j|a) = 1$

Global semantics

- **Global** semantics defines the full joint distribution as the product of the local conditional distributions:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$



- Compute: $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$
- <https://pollev.com/frederikm106>

Global semantics

- **Global** semantics defines the full joint distribution as the product of the local conditional distributions:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

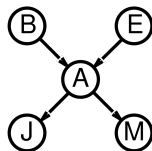
- $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e)$$

$$= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998$$

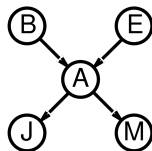
$$\approx 0.00063$$

- Application of the chain rule.



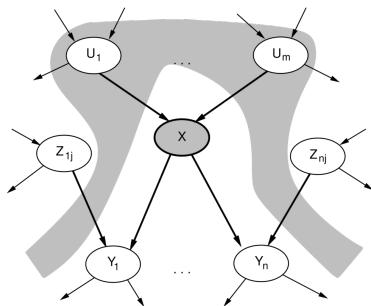
Compactness

- 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
- Grows linearly with n , vs. $O(2^n)$ for the full joint distribution
- For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers (vs. $2^5 - 1 = 31$)



Local semantics

- **Local** semantics: each node is conditionally independent of its nondescendants given its parents **define descendants**

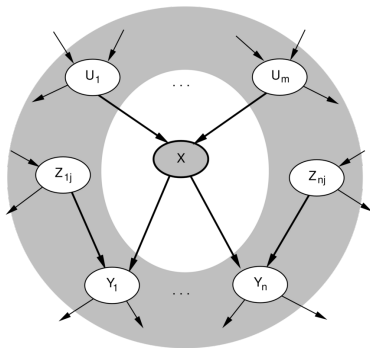


- Theorem:

Local semantics \Leftrightarrow global semantics

Markov blanket

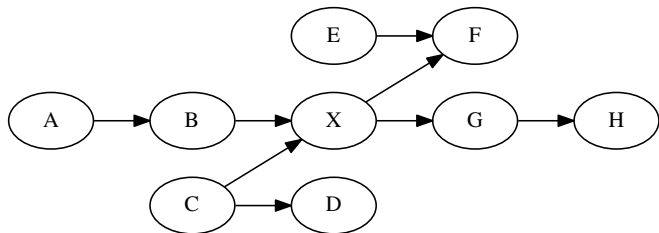
- Each node is conditionally independent of all others given its **Markov blanket**: parents + children + children's parents



Andrey Markov

Markov blanket

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



- Markov blanket of X?

<https://pollev.com/frederikm106>

Last year's success rate: 0.72

Constructing Bayesian networks

- Build Bayesian networks like anyother form of **knowledge representation**.
- First figure out the variables that describe the world.
- Then decide how they are connected.
Conditional independence.
- Then work out the values in the CPTs.



Kathy Laskey

Compact conditional distributions

- CPT grows exponentially with number of parents
 - Use **canonical** distributions that are defined compactly
- **Deterministic** nodes are the simplest case.
- $X = f(\text{Parents}(X))$ for some function f
 - Boolean functions:

$$\text{NorthAmerican} \Leftrightarrow \text{Canadian} \vee \text{US} \vee \text{Mexican}$$

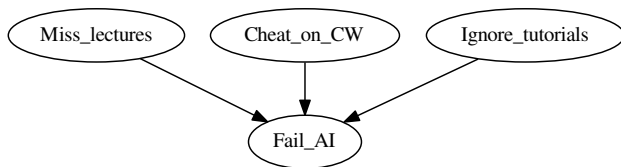
- Numerical relationships among continuous variables

$$\frac{\partial \text{Level}}{\partial t} = \text{inflow} + \text{precipitation} - \text{outflow} - \text{evaporation}$$



Noisy-OR

- Say we have this network



- The noisy-OR model allows for uncertainty about the ability of each parent to cause the child to be true—the causal relationship between parent and child may be inhibited: a student could miss a lecture, but not fail.
- This is different from conventional logic
- Compare Noisy-OR to Naive Bayes
- We would like to compute things such as:

$$P(\text{fail_AI} \mid \text{miss_lectures}, \text{ignore_tutorials}, \neg \text{cheat_on_CW})$$

Compact conditional distributions

- **Noisy-OR** distributions model multiple noninteracting causes

- 1 Parents $U_1 \dots U_k$ include all causes
(Can add **leak node** that covers “miscellaneous causes.”)
- 2 It assumes that inhibition of each parent is independent of inhibition of any other parents. Formally,

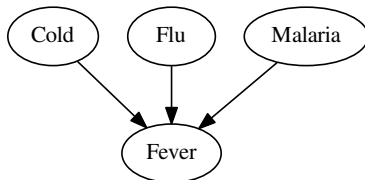
$$\Rightarrow P(X|U_1 \dots U_j, \neg U_{j+1} \dots \neg U_k) = 1 - \prod_{i=1}^j q_i,$$

where $q_i = P(X | U_i)$

(only U_i is true)



Example



$$q_{cold} := P(\neg fever \mid cold) = 0.6$$

- $q_{flu} := P(\neg fever \mid flu) = 0.2$

$$q_{malaria} := P(\neg fever \mid malaria) = 0.1$$

- In this model, we then get

$$P(\neg fever \mid cold, flu, \neg malaria) = 1 - q_{cold} \cdot q_{flu} = 1 - 0.12$$

- Given that you have the flu and a cold, the only way you cannot have fever is if neither of them gave you fever. (Think of all possible outcomes)

More examples

$$q_{cold} := P(\neg fever \mid cold) = 0.6$$

$$q_{flu} := P(\neg fever \mid flu) = 0.2$$

$$q_{malaria} := P(\neg fever \mid malaria) = 0.1$$

<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	$P(\text{Fever})$	$P(\neg \text{Fever})$
F	F	F	0.0	1.0
F	F	T	0.9	0.1
F	T	F	0.8	0.2
F	T	T	0.98	$0.02 = 0.2 \times 0.1$
T	F	F	0.4	0.6
T	F	T	0.94	$0.06 = 0.6 \times 0.1$
T	T	F	0.88	$0.12 = 0.6 \times 0.2$
T	T	T	0.988	$0.012 = 0.6 \times 0.2 \times 0.1$

- Number of parameters **linear** in number of parents

- The required space is small!

- Okay, so what can we do with Bayesian Networks?

- Okay, so what can we do with Bayesian Networks?
- They are useful for inference (a conclusion reached on the basis of evidence and reasoning)

- **Simple queries**: compute posterior marginal $\mathbf{P}(X_i|\mathbf{E} = \mathbf{e})$
 $\mathbf{P}(\text{Brexit} | PM = \text{Boris}, Economy = \text{not_great})$

- **Conjunctive queries**

$$\mathbf{P}(X_i, X_j | \mathbf{E} = \mathbf{e}) = \mathbf{P}(X_i | \mathbf{E} = \mathbf{e}) \mathbf{P}(X_j | X_i, \mathbf{E} = \mathbf{e})$$

- **Optimal decisions**: decision networks include utility information; probabilistic inference required for $P(\text{outcome} | \text{action}, \text{evidence})$
- **Value of information**: which evidence to seek next?
- **Sensitivity analysis**: which probability values are most critical?
- **Explanation**: why do I need a new starter motor?

- We will focus on simple queries:

- Compute posterior marginal

$$\mathbf{P}(X_i | \mathbf{E} = \mathbf{e})$$

$$\mathbf{P}(\textit{Brexit} | \textit{PM} = \textit{Boris}, \textit{Economy} = \textit{not_great})$$

- We will look a several ways of doing this.

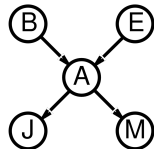
- 1 Enumeration
- 2 Rejection sampling (using prior sampling)
- 3 Likelihood weighting
- 4 Gibbs sampling

- Simplest approach to evaluating the network is to do just as we did for the dentist example.
- Difference is that we use the structure of the network to tell us which sets of joint probabilities to use.
 - Thanks Professor Markov
- Gives us a slightly intelligent way to sum out variables from the joint without actually constructing its explicit representation

Inference by enumeration

- Simple query on the burglary network.

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



- Rewrite full joint entries taking network into account:

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

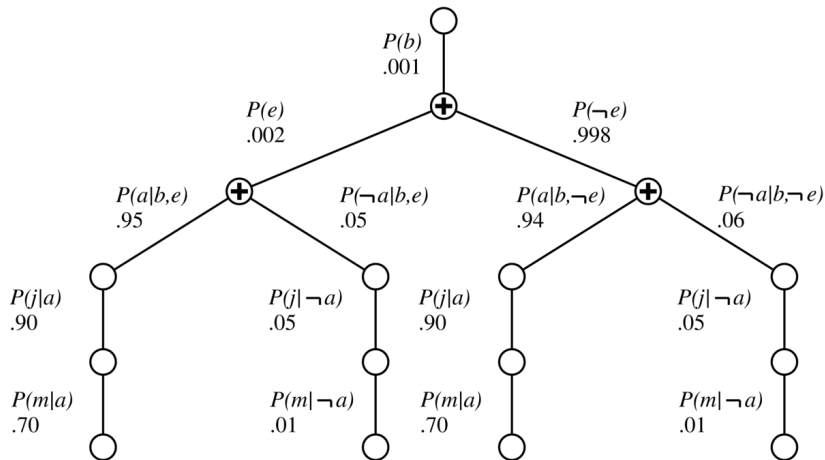
- We evaluate this expression

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

by going through the variables in order, multiplying CPT entries along the way.

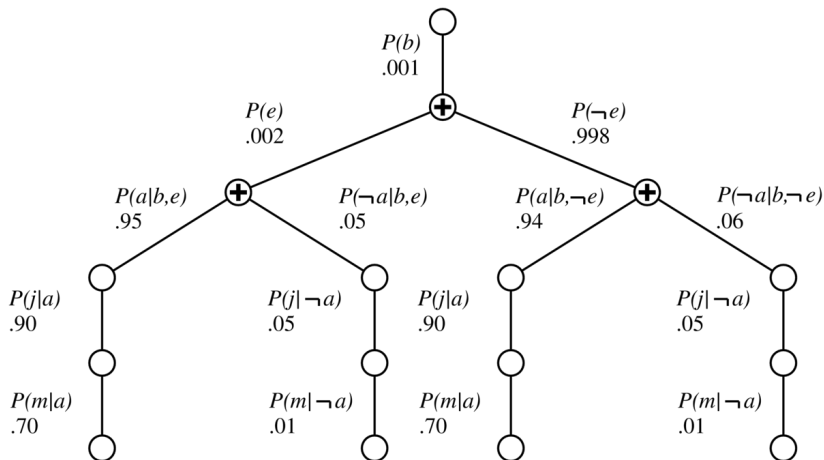
- At each point, we need to loop through the possible values of the variable.
- Involves a lot of repeated calculations.

Evaluation tree



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

Evaluation tree



Inefficient: computes $P(j|a)P(m|a)$ for each value of e

Enumeration algorithm

function **ENUMERATION-Ask**(X, \mathbf{e}, bn) **returns** a distribution over X

inputs: X , the query variable

\mathbf{e} , observed values for variables \mathbf{E}

bn , a Bayesian network $\{X\} \cup \mathbf{E} \cup \mathbf{Y}$

$\mathbf{Q}(X) \leftarrow$ a distribution over X , initially empty

for each value x_i of X **do**

 extend \mathbf{e} with value x_i for X

$\mathbf{Q}(x_i) \leftarrow \text{ENUMERATE-ALL}(\text{VARS}[bn], \mathbf{e})$

return $\text{NORMALIZE}(\mathbf{Q}(X))$



Enumeration algorithm

```
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  $P(y \mid Pa(Y))$ 
         $\times$  ENUMERATE-ALL(REST(vars), e)
else return  $\sum_y P(y \mid Pa(Y))$ 
         $\times$  ENUMERATE-ALL(REST(vars), ey)
    where ey is e extended with Y = y
```


Other exact approaches

- We can improve on enumeration.
- **Variable elimination** evaluates the enumeration tree bottom up, remembering intermediate values.
 - Simple and efficient for single queries
- **Clustering algorithms** can be more efficient for multiple queries
 - Group variables together strategically.
- However, *all* exact inference can be computationally intractable.



Complexity of exact inference

- Singly connected networks (or polytrees)
- Any two nodes are connected by at most one (undirected) path
- Time and space cost of variable elimination are:

$$O(d^k n)$$

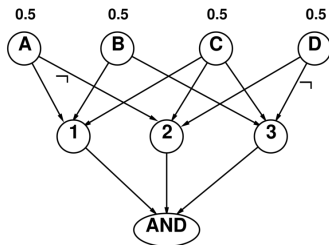
for k parents, d values.



Complexity of exact inference

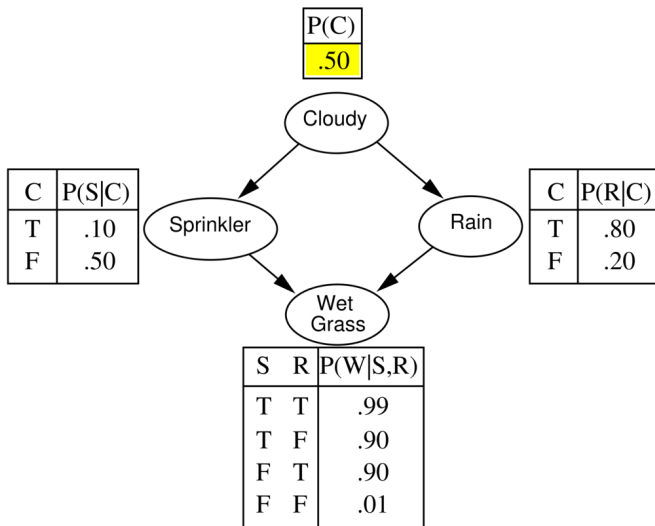
- **Multiply connected** networks.
- Exponential time and space complexity, even when number of parents of a node is bounded.
- Inference is NP-hard.
- In fact, #P-complete.

1. $A \vee B \vee C$
2. $C \vee D \vee \neg A$
3. $B \vee C \vee \neg D$



- Let's use stochastic sampling instead!

Prior sampling



- How would you estimate $P(c, \neg s, r, w)$?

Take a step back



- How would you estimate $P(\text{die shows } 7)$?

- How would you estimate $P(\text{die shows } 7)$?
- Simple: you take n random samples from the network
- Let X_i be the binary r.v. that is 1 if the event sampled in the i th run is 7
- Simply output

$$\hat{P}(7) = \frac{\sum_{i=1}^n X_i}{n}$$

- Law of large numbers says that $\lim_{n \rightarrow \infty} \hat{P} = P$.

- Back to our Bayesian network with cloudy, sprinkler, rain, and wet grass.
- How would you estimate $P(c, \neg s, r, w)$?

Stochastic simulation

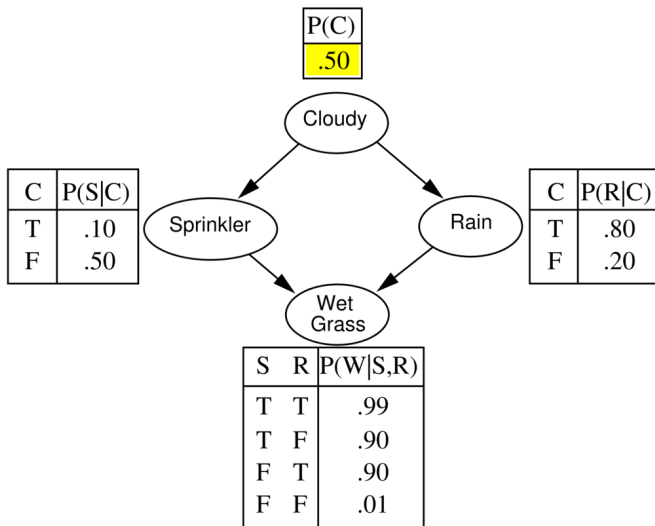
- Back to our Bayesian network with cloudy, sprinkler, rain, and wet grass.
- How would you estimate $P(c, \neg s, r, w)$?
- Simple you just take say n random samples from the network
- Let X_i be the binary r.v. that is 1 if the event sampled in the i th run is $c, \neg s, r, w$
- Simply output

$$\hat{P}(c, \neg s, r, w) = \frac{\sum_{i=1}^n X_i}{n}$$

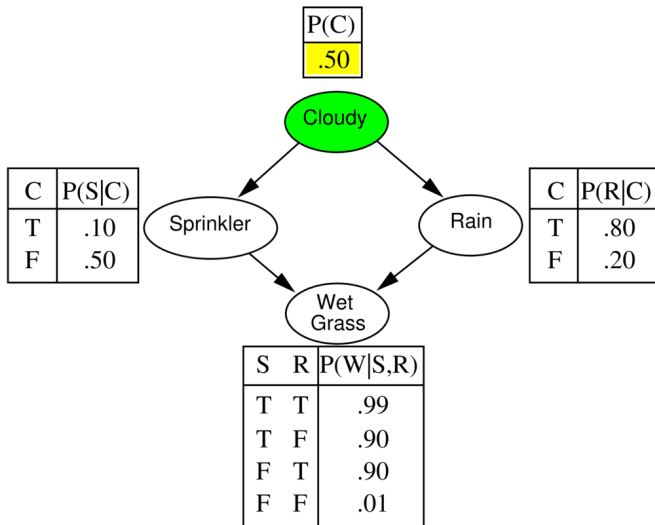
- Law of large numbers says that $\lim_{n \rightarrow \infty} \hat{P} = P$.



Prior sampling

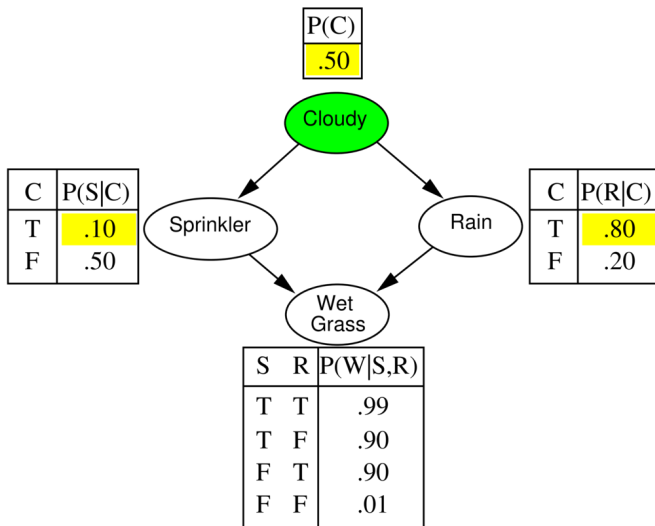


Prior sampling

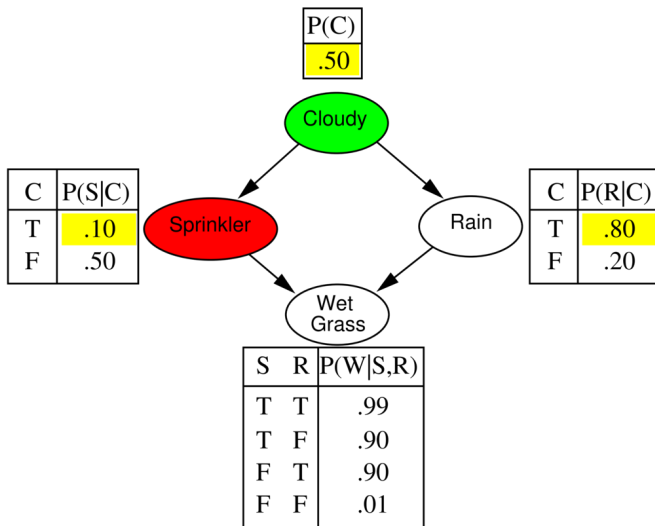


- *Cloudy = true*

Prior sampling

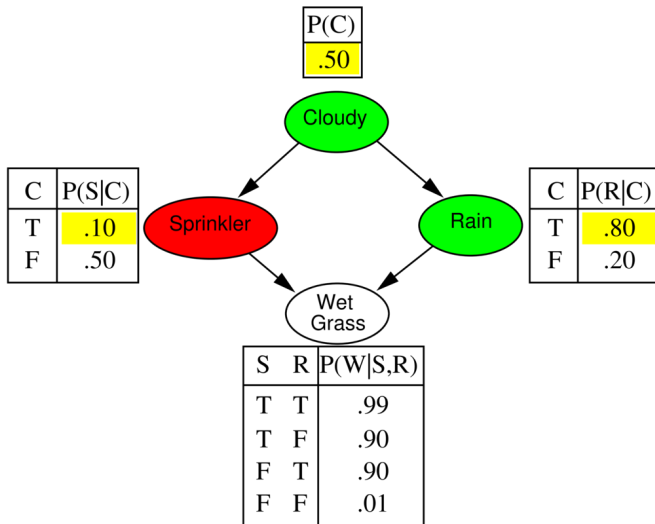


Prior sampling



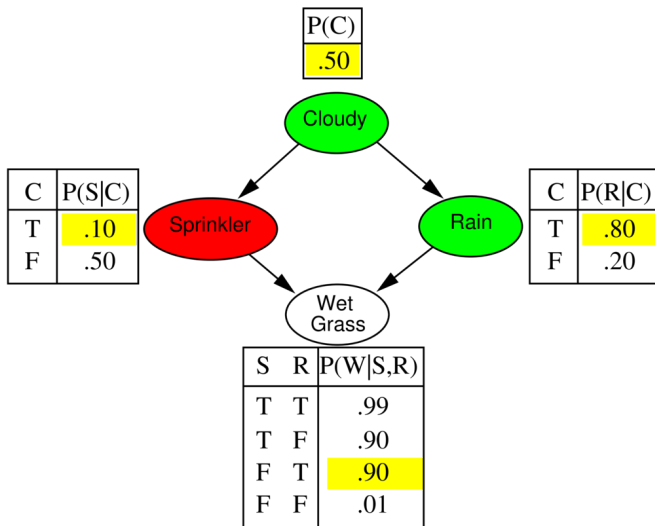
- *Sprinkler = false*

Prior sampling

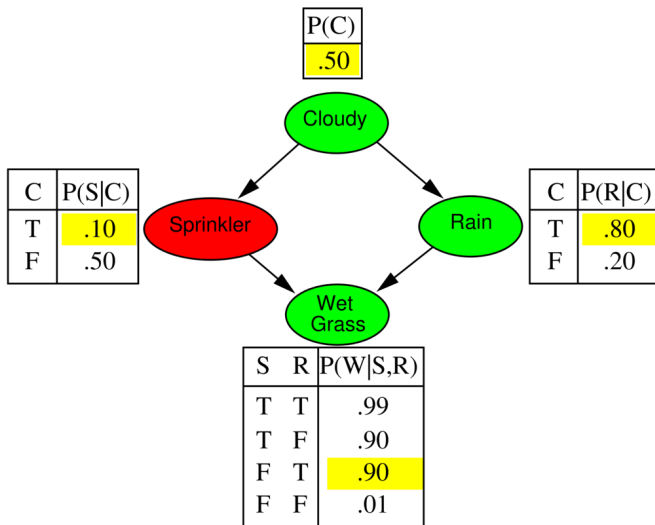


- *Rain = true*

Prior sampling



Prior sampling



- *WetGrass = true*

- So, this time we get the event

$[Cloudy = true, Sprinkler = false, Rain = true, WetGrass = true]$

- Will write this:

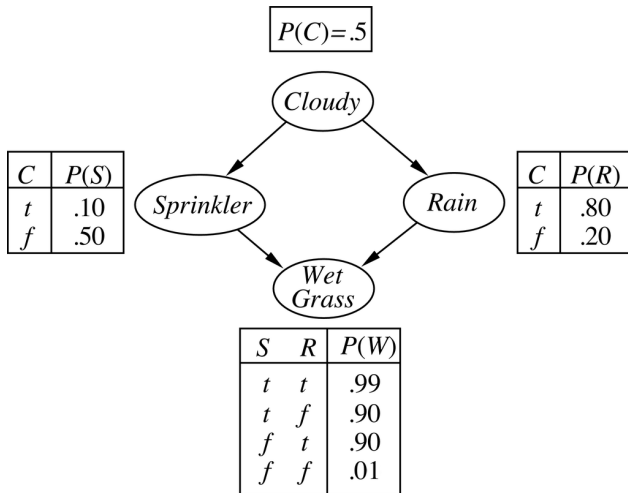
$[true, false, true, true]$

- If we repeat the process many times, we can count the number of times $[true, false, true, true]$ is the result.
- The proportion of this to the total number of runs is:

$$P(c, \neg s, r, w)$$

- The more runs, the more accurate the probability.
- Similarly for other joint probabilities.

Again?



What is $P(c, \neg s, r, w)$?

Prior sampling

function **PRIOR-SAMPLE**(*bn*) **returns** an event sampled from *bn*

inputs: *bn*, a belief network specifying joint distribution $\mathbf{P}(X_1, \dots, X_n)$

x \leftarrow an event with *n* elements

for *i* = 1 **to** *n* **do**

x_i \leftarrow a random sample from $\mathbf{P}(X_i \mid \text{parents}(X_i))$
 given the values of *Parents*(*X_i*) in **x**

return **x**



- How would you get the following marginal distribution?

$$\mathbf{P}(X|\mathbf{e})$$

Rejection sampling

- Use **Rejection sampling**.
- Generate samples as before (using prior sampling)
- If the sample is such that \mathbf{e} holds use it to build an estimate
- Otherwise, ignore it

Rejection sampling

function REJECTION-SAMPLING(X, \mathbf{e}, bn, N) **returns** an estimate of $P(X|\mathbf{e})$

local variables: \mathbf{N} , a vector of counts over X , initially zero

for $j = 1$ to N **do**

$\mathbf{x} \leftarrow \text{PRIOR-SAMPLE}(bn)$

if \mathbf{x} is consistent with \mathbf{e} **then**

$\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$ where x is the value of X in \mathbf{x}

return NORMALIZE($\mathbf{N}[X]$)

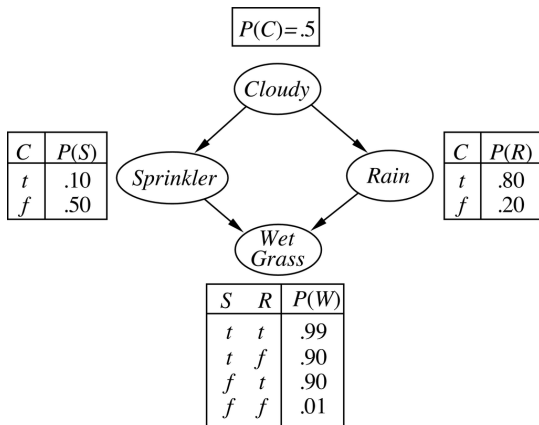
Rejection sampling

- More efficient than prior sampling, but
- For unlikely events, may have to wait a long time to get enough matching samples.
- Still inefficient.
- So, use **likelihood weighting**

- Version of **importance sampling**.
- Fix evidence variable to *true*, so just sample relevant events.
- Have to weight them with the likelihood that they fit the evidence.
- Use the probabilities we know to weight the samples.

Likelihood weighting

- Consider we have the following network:



- Say we want to establish
 $\mathbf{P}(\text{Rain} | \text{Cloudy} = \text{true}, \text{WetGrass} = \text{true})$

- We want $\mathbf{P}(Rain|Cloudy = true, WetGrass = true)$
- We pick a variable ordering, say:
Cloudy, Sprinkler, Rain, WetGrass.
as before.
- Set the weight $w = 1$ and we start.
- Deal with each variable in order.

- Remember, we want $\mathbf{P}(\text{Rain} | \text{Cloudy} = \text{true}, \text{WetGrass} = \text{true})$
- Cloudy* is true, so:

$$w \leftarrow w \times P(\text{Cloudy} = \text{true})$$

$$w \leftarrow 0.5$$

- Cloudy* = true, *Sprinkler* = ?, *Rain* = ?, *WetGrass* = ?.
- $w = 0.5$

Likelihood weighting

- *Sprinkler* is not an evidence variable, so we don't know whether it is true or false.
- Sample a value just as we did for prior sampling:

$$\mathbf{P}(\textit{Sprinkler} | \textit{Cloudy} = \textit{true}) = \langle 0.1, 0.9 \rangle$$

- Let's assume this returns *false*.
- w remains the same.
- $\textit{Cloudy} = \textit{true}$, $\textit{Sprinkler} = \textit{false}$, $\textit{Rain} = ?$, $\textit{WetGrass} = ?$.
- $w = 0.5$



- *Rain* is not an evidence variable, so we don't know whether it is true or false.
- Sample a value just as we did for prior sampling:

$$\mathbf{P}(Rain|Cloudy = true) = \langle 0.8, 0.2 \rangle$$

- Let's assume this returns *true*.
- *w* remains the same.
- *Cloudy* = *true*, *Sprinkler* = *false*, *Rain* = *true*, *WetGrass* = ?.
- *w* = 0.5

- *WetGrass* is an evidence variable with value *true*, so we set:

$$w \leftarrow w \times P(WetGrass = true | Sprinkler = false, Rain = true)$$

$$w \leftarrow 0.45$$

- *Cloudy* = *true*, *Sprinkler* = *false*, *Rain* = *true*,
WetGrass = *true*.
- $w = 0.45$

- So we end with the event $[true, false, true, true]$ and weight 0.45.
- To find a probability we tally up all the relevant events, weighted with their weights.
- The one we just calculated would tally under

$$Rain = true$$

- As before, more samples means more accuracy.

Likelihood weighting

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($\mathbf{W}[X]$)



Likelihood weighting

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

- A rather different approach to sampling.
- Part of the Markov Chain Monte Carlo (MCMC) family of algorithms.
- Don't generate each sample from scratch.
- Generate samples by making a random change to the previous sample.

Gibbs sampling

- Gibbs sampling for Bayesian networks starts with an arbitrary state.
- So pick state with evidence variables fixed at observed values.
(If we know *Cloudy* = *true*, we pick that.)
- Generate next state by randomly sampling from a non-evidence variable.
- Do this sampling conditional on the current values of the Markov blanket.
- “The algorithm therefore wanders randomly around the state space . . . flipping one variable at a time, but keeping the evidence variables fixed”.



- Consider the query:

$$\mathbf{P}(\textit{Cloudy} | \textit{Sprinkler} = \textit{true}, \textit{WetGrass} = \textit{true})$$

- The evidence variables are fixed to their observed values.
- The non-evidence variables are initialised randomly.

Cloudy = true

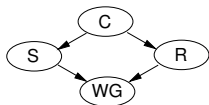
Rain = false

- State is thus:

[*Cloudy* = true, *Sprinkler* = true,
Rain = false, *WetGrass* = true].

Gibbs sampling

- First we sample *Cloudy* given the current state of its Markov blanket.
- Markov blanket is *Sprinkler* and *Rain*.
- So, sample from:



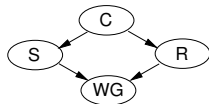
$\mathbf{P}(\textit{Cloudy} | \textit{Sprinkler} = \textit{true}, \textit{Rain} = \textit{false})$

- Suppose we get *Cloudy* = *false*, then new state is:

[*Cloudy* = *false*, *Sprinkler* = *true*,
Rain = *false*, *WetGrass* = *true*].

Gibbs sampling

- Next we sample *Rain* given the current state of its Markov blanket.
- Markov blanket is *Cloudy*, *Sprinkler* and *WetGrass*.
- So, sample from:



$\mathbf{P}(\text{Rain} | \text{Cloudy} = \text{false}, \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{false})$

- Suppose we get *Rain* = *true*, then new state is:

[*Cloudy* = *false*, *Sprinkler* = *true*,
Rain = *true*, *WetGrass* = *true*].

Gibbs sampling

- Each state visited during this process contributes to our estimate for:

$$\mathbf{P}(\textit{Cloudy} | \textit{Sprinkler} = \textit{true}, \textit{WetGrass} = \textit{true})$$

- Say the process visits 80 states.
- In 20, *Cloudy* = *true*
- In 60, *Cloudy* = *false*
- Then

$$\begin{aligned}\mathbf{P}(\textit{Cloudy} | \textit{Sprinkler} = \textit{true}, \textit{WetGrass} = \textit{true}) \\ &= \textit{Normalize}(\langle 20, 60 \rangle) \\ &= \langle 0.25, 0.75 \rangle\end{aligned}$$



Gibbs sampling

function **GIBBS-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{Z}

for $j = 1$ to N **do**

for each Z_i in \mathbf{Z} **do**

sample the value of Z_i in \mathbf{x} from $\mathbf{P}(Z_i | mb(Z_i))$
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 $\text{NORMALIZE}(\mathbf{N}[X])$



- All of this begs the question:
“How do we sample a variable given the state of its Markov blanket?”
- For a value x of a variable X :

$$\mathbf{P}(X|mb(X)) = \alpha \mathbf{P}(X|parents(X)) \prod_{Y \in Children(X)} \mathbf{P}(Y|parents(Y))$$

where $mb(X)$ is the Markov blanket of X .

- Given $\mathbf{P}(X|mb(X))$, we can sample from it just as we have before.

To summarize

- Bayesian networks exploit conditional independence to create a more compact set of information.
- Reasonably efficient computation for some problems.
- Five approaches to inference in Bayesian networks.
 - Exact: Inference by enumeration.
 - Approximate: Prior sampling
 - Approximate: Rejection sampling
 - Approximate: Importance sampling/likelihood weighting
 - Approximate: Gibbs sampling
- Can answer a simple query for any BN.

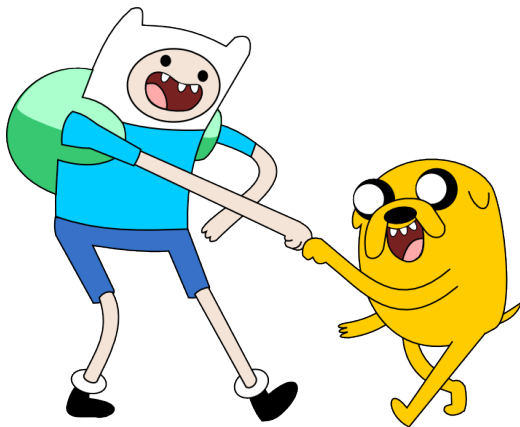
From probability to decision making

- What we have covered allows us to compute probabilities of interesting events.
- But **beliefs** alone are not so interesting to us.
- In Pacman don't care so much if there is a ghost in (2, 2), so much as we care whether we should go WEST or SOUTH.
- This is complicated in an uncertain world.
 - Don't know the outcome of actions.
 - Non-deterministic as well as partially observable



- Next time!

Mathematical!



(Pendleton Ward/Cartoon Network)

Summary

- This lecture started with exact probabilistic inference in Bayesian networks.
 - Inference by enumeration.
- Efficiency prompted us to look at approximate techniques based on sampling:
 - Prior sampling
 - Rejection sampling
 - Likelihood weighting (Importance sampling)
 - Gibbs sampling