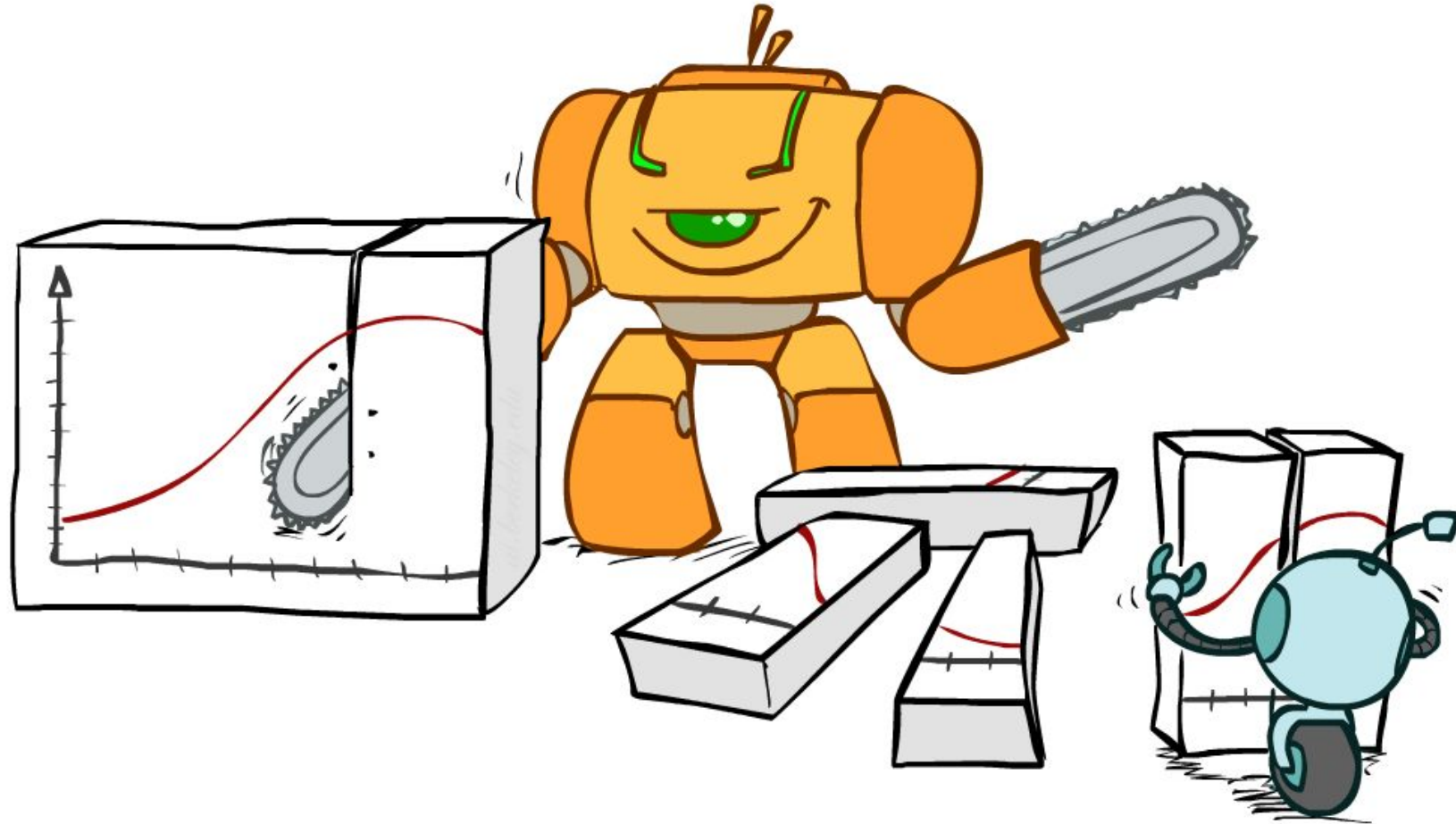


Reminder: elementary probability

- Basic laws: $0 \leq P(\omega) \leq 1$ $\sum_{\omega \in \Omega} P(\omega) = 1$
- Events: subsets of Ω : $P(A) = \sum_{\omega \in A} P(\omega)$
- Random variable $X(\omega)$ has a value in each ω
 - Distribution $P(X)$ gives probability for each possible value x
 - Joint distribution $P(X,Y)$ gives total probability for each combination x,y
- Summing out/marginalization: $P(X=x) = \sum_y P(X=x, Y=y)$
- Conditional probability: $P(X|Y) = P(X,Y)/P(Y)$
- Product rule: $P(X|Y)P(Y) = P(X,Y) = P(Y|X)P(X)$
 - Generalize to chain rule: $P(X_1, \dots, X_n) = \prod_i P(X_i | X_1, \dots, X_{i-1})$

Bayes Rule



Bayes' Rule

- Write the product rule both ways:

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

- Dividing left and right expressions, we get:

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

- Why is this at all helpful?
 - Lets us build one conditional from its reverse
 - Often one conditional is tricky but the other one is simple
 - Describes an “update” step from prior $P(a)$ to posterior $P(a | b)$
 - Foundation of many systems we'll see later (e.g. ASR, MT)
- In the running for most important AI equation!

That's my rule!



Inference with Bayes' Rule

- Example: Diagnostic probability from causal probability:

$$P(\text{cause} \mid \text{effect}) = \frac{P(\text{effect} \mid \text{cause}) P(\text{cause})}{P(\text{effect})}$$

- Example:

- M: meningitis, S: stiff neck

$$\left. \begin{array}{l} P(s \mid m) = 0.8 \\ P(m) = 0.0001 \\ P(s) = 0.01 \end{array} \right\} \text{Example gives}$$

$$P(m \mid s) = \frac{P(s \mid m) P(m)}{P(s)} = \frac{0.8 \times 0.0001}{0.01}$$

- Note: posterior probability of meningitis still very small: 0.008 (80x bigger – why?)
- Note: you should still get stiff necks checked out! Why?

Probabilistic Inference

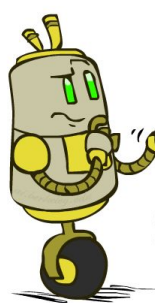
- Probabilistic inference: compute a desired probability from a probability model
 - Typically for a *query variable* given *evidence*
 - E.g., $P(\text{airport on time} \mid \text{no accidents}) = 0.90$
 - These represent the agent's *beliefs* given the evidence
- Probabilities change with new evidence:
 - $P(\text{airport on time} \mid \text{no accidents, 5 a.m.}) = 0.95$
 - $P(\text{airport on time} \mid \text{no accidents, 5 a.m., raining}) = 0.80$
 - Observing new evidence causes *beliefs to be updated*



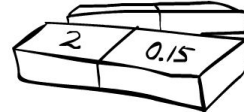
Inference by Enumeration

- Probability model $P(X_1, \dots, X_n)$ is given
 - Partition the variables X_1, \dots, X_n into sets as follows:
 - Evidence** variables: $E = e$
 - Query** variables: Q
 - Hidden** variables: H
- We want: $P(Q | e)$

- Step 1: Select the entries consistent with the evidence



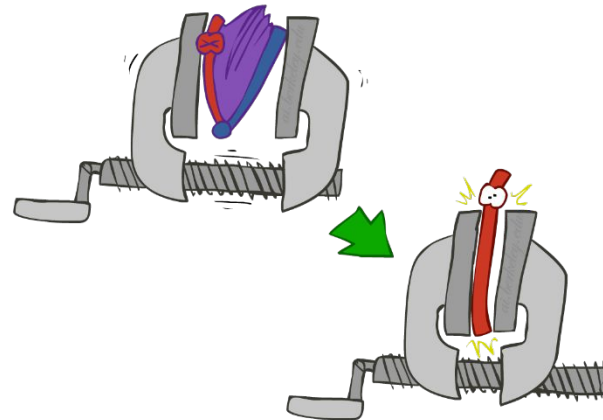
x	P(x)
-3	0.05
-1	0.25
0	0.07
1	0.2
5	0.01



2	0.15
---	------

- Step 2: Sum out H from model to get joint of query and evidence

$$P(Q, e) = \sum_h \underbrace{P(Q, h, e)}_{X_1, \dots, X_n}$$



- Step 3: Normalize

$$P(Q | e) = \alpha P(Q, e)$$

Inference by Enumeration

- $P(W \mid \text{winter})?$

Season	Temp	Weather	P
summer	hot	sun	0.35
summer	hot	rain	0.01
summer	hot	fog	0.01
summer	hot	meteor	0.00
summer	cold	sun	0.10
summer	cold	rain	0.05
summer	cold	fog	0.09
summer	cold	meteor	0.00
winter	hot	sun	0.10
winter	hot	rain	0.01
winter	hot	fog	0.02
winter	hot	meteor	0.00
winter	cold	sun	0.15
winter	cold	rain	0.20
winter	cold	fog	0.18
winter	cold	meteor	0.00

Inference by Enumeration

- Obvious problems:
 - Worst-case time complexity $O(d^n)$
 - Space complexity $O(d^n)$ to store the joint distribution
 - $O(d^n)$ data points to estimate the entries in the joint distribution

Independence

- Two variables X and Y are (absolutely) **independent** if

$$\forall x, y \quad P(x, y) = P(x) P(y)$$

- I.e., the joint distribution **factors** into a product of two simpler distributions

- Equivalently, via the product rule $P(x, y) = P(x | y) P(y)$,

$$P(x | y) = P(x) \quad \text{or} \quad P(y | x) = P(y)$$

- Example: two dice rolls $Roll_1$ and $Roll_2$

- $P(Roll_1=5, Roll_2=3) = P(Roll_1=5) P(Roll_2=3) = 1/6 \times 1/6 = 1/36$

- $P(Roll_2=3 | Roll_1=5) = P(Roll_2=3)$



Example: Independence

- n fair, independent coin flips:

$P(X_1)$

H	0.5
T	0.5

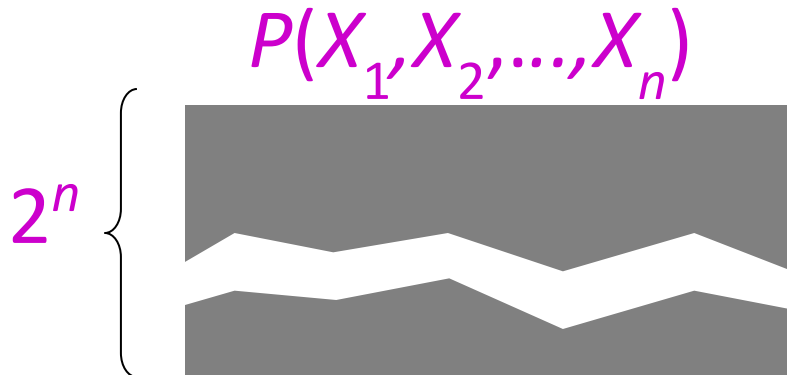
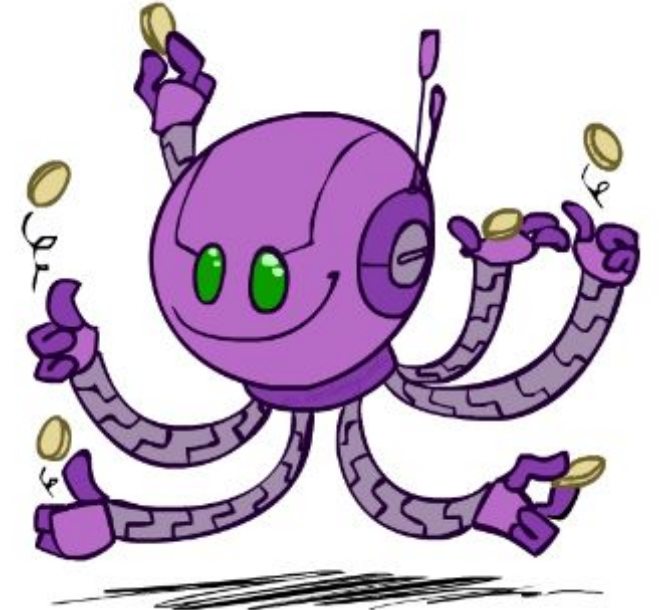
$P(X_2)$

H	0.5
T	0.5

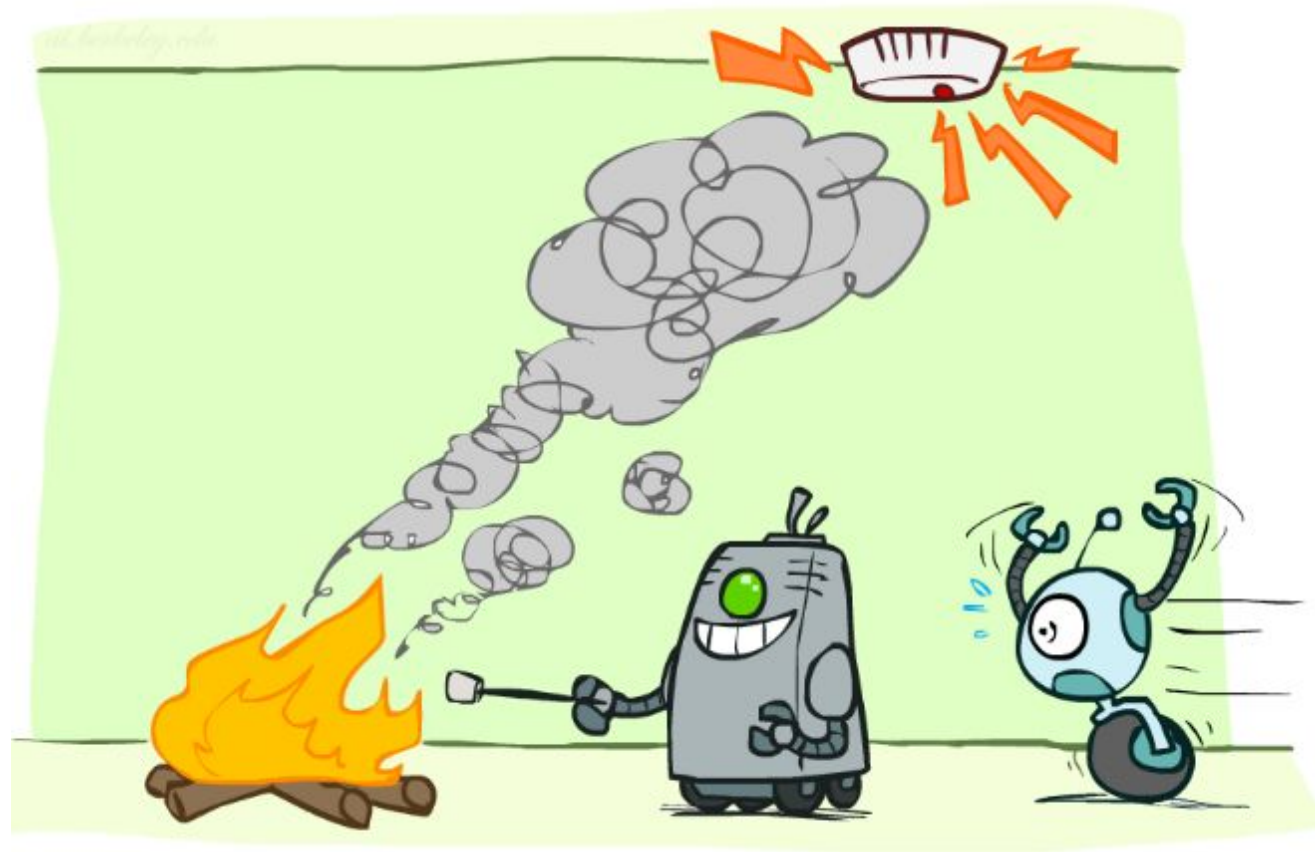
...

$P(X_n)$

H	0.5
T	0.5

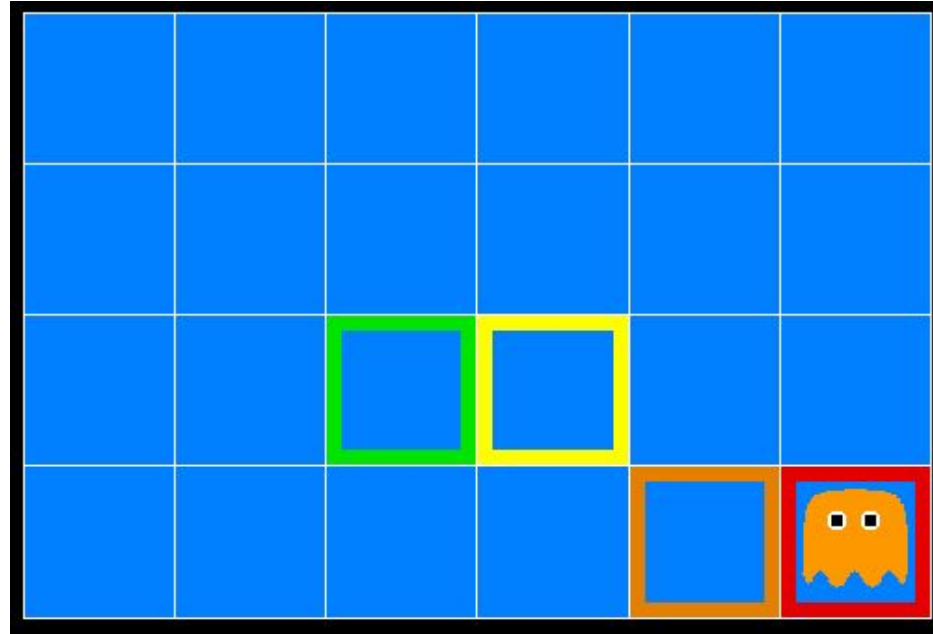


Conditional Independence



Ghostbusters

- A ghost is in the grid somewhere
- Sensor readings tell how close a square is to the ghost
 - On the ghost: usually red
 - 1 or 2 away: mostly orange
 - 3 or 4 away: typically yellow
 - 5+ away: often green
- Click on squares until confident of location, then “*bust*”



Video of Demo Ghostbusters with Probability



Ghostbusters model

- Variables and ranges:

- G (ghost location) in $\{(1,1), \dots, (3,3)\}$
- $C_{x,y}$ (color measured at square x,y) in $\{\text{red, orange, yellow, green}\}$

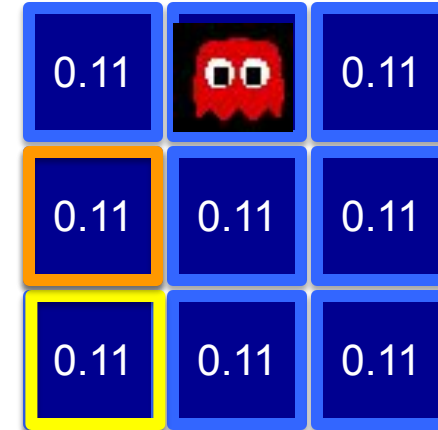
0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11

- Ghostbuster physics:

- **Uniform prior distribution** over ghost location: $P(G)$
- **Sensor model**: $P(C_{x,y} \mid G)$ (depends only on distance to G)
 - E.g. $P(C_{1,1} = \text{yellow} \mid G = (1,1)) = 0.1$

Ghostbusters model, contd.

- $P(G, C_{1,1}, \dots, C_{3,3})$ has $9 \times 4^9 = 2,359,296$ entries!!!
- Ghostbuster independence:
 - Are $C_{1,1}$ and $C_{1,2}$ independent?
 - E.g., does $P(C_{1,1} = \text{yellow}) = P(C_{1,1} = \text{yellow} \mid C_{1,2} = \text{orange})$?
- Ghostbuster physics again:
 - $P(C_{x,y} \mid G)$ **depends only on distance to G**
 - So $P(C_{1,1} = \text{yellow} \mid \underline{G = (2,3)}) = P(C_{1,1} = \text{yellow} \mid \underline{G = (2,3)}, C_{1,2} = \text{orange})$
 - I.e., $C_{1,1}$ is **conditionally independent** of $C_{1,2}$ **given G**



Ghostbusters model, contd.

- Apply the chain rule to decompose the joint probability model:
- $P(G, C_{1,1}, \dots, C_{3,3}) = P(G) P(C_{1,1} | G) P(C_{1,2} | G, C_{1,1}) P(C_{1,3} | G, C_{1,1}, C_{1,2}) \dots P(C_{3,3} | G, C_{1,1}, \dots, C_{3,2})$
- Now simplify using conditional independence:
- $P(G, C_{1,1}, \dots, C_{3,3}) = P(G) P(C_{1,1} | G) P(C_{1,2} | G) P(C_{1,3} | G) \dots P(C_{3,3} | G)$
- I.e., conditional independence properties of ghostbuster physics simplify the probability model from **exponential** to **quadratic** in the number of squares
- This is called a **Naïve Bayes** model:
 - One discrete query variable (often called the **class** or **category** variable)
 - All other variables are (potentially) evidence variables
 - Evidence variables are all conditionally independent given the query variable

Conditional Independence

- **Conditional independence** is our most basic and robust form of knowledge about uncertain environments.

- X is conditionally independent of Y given Z if and only if:

$$\forall x, y, z \quad P(x \mid y, z) = P(x \mid z)$$

or, equivalently, if and only if

$$\forall x, y, z \quad P(x, y \mid z) = P(x \mid z) P(y \mid z)$$

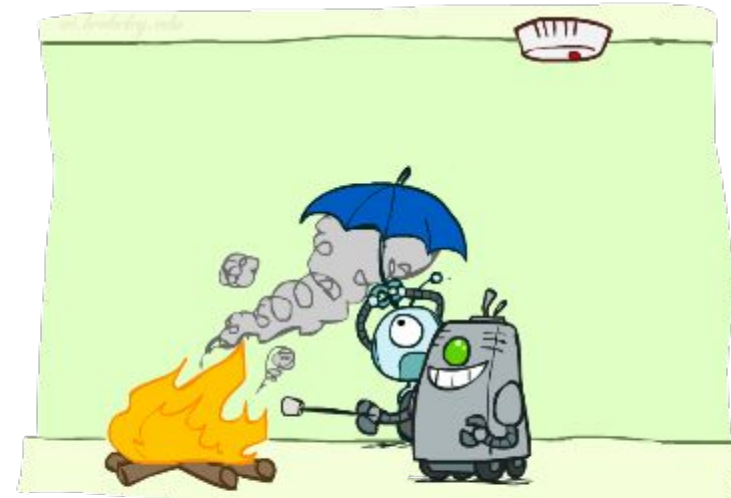
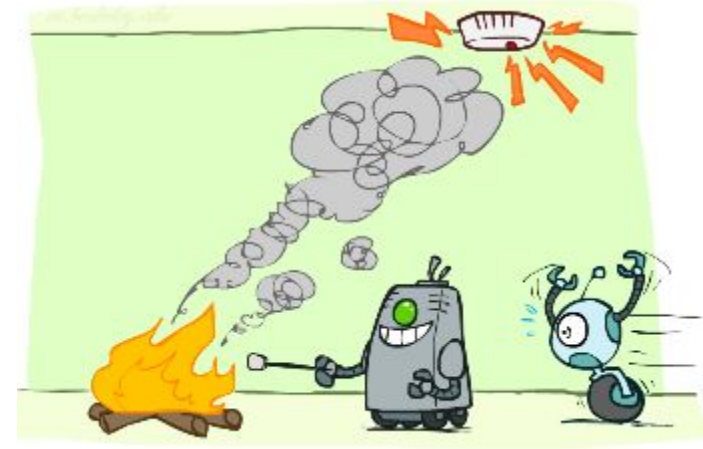
Conditional Independence

- What about this domain:
 - Traffic
 - Umbrella
 - Raining

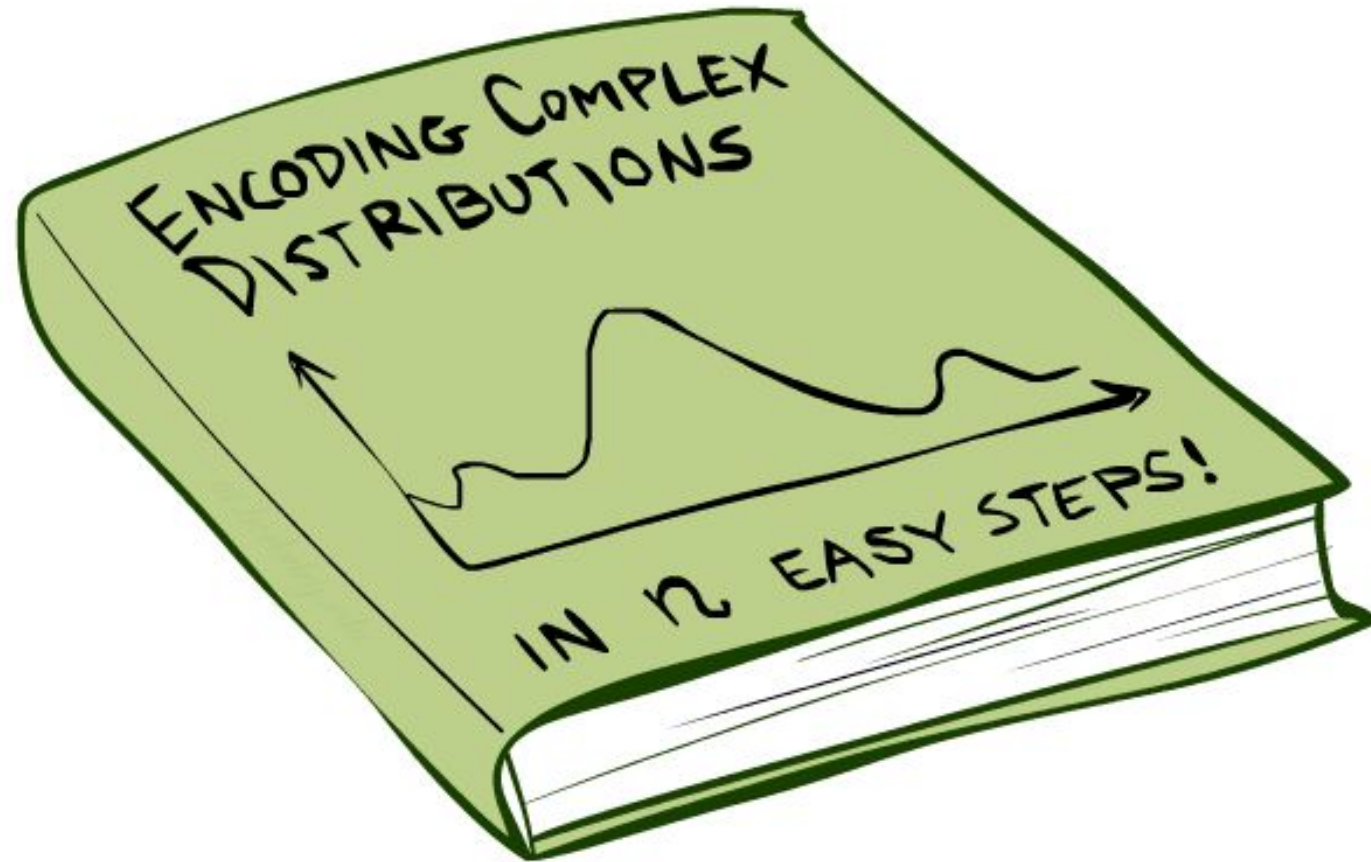


Conditional Independence

- What about this domain:
 - Fire
 - Smoke
 - Alarm

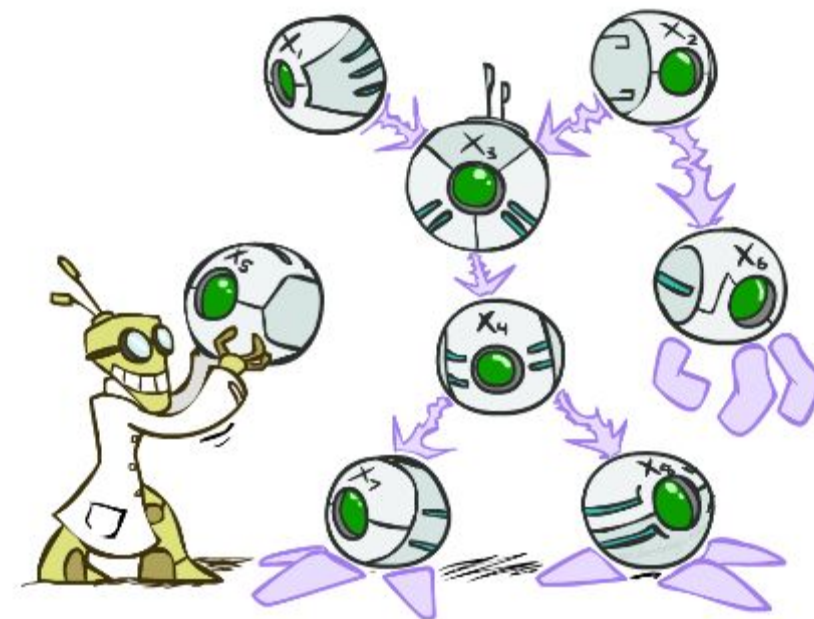


Bayes Nets: Big Picture



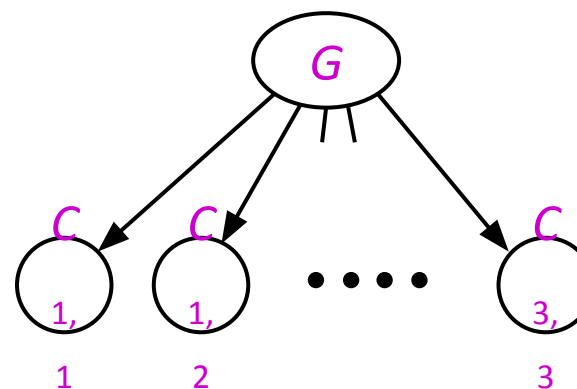
Bayes Nets: Big Picture


- **Bayes nets:** a technique for describing complex joint distributions (models) using simple, conditional distributions
 - A subset of the general class of **graphical models**
- Use local causality/conditional independence:
 - the world is composed of many variables,
 - each interacting locally with a few others



Graphical Model Notation

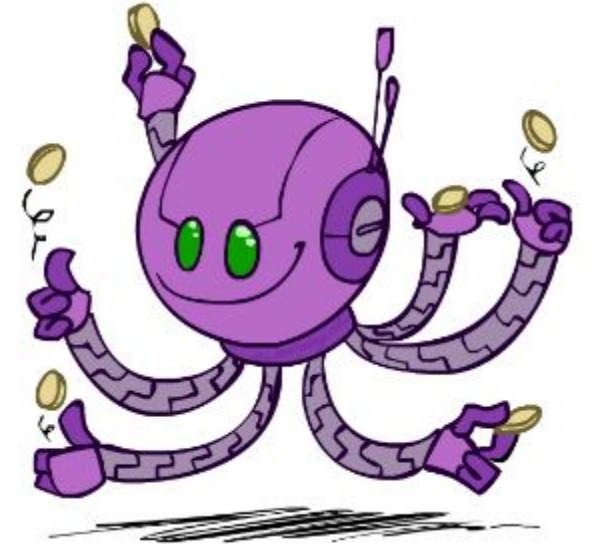
- **Nodes: variables (with domains)**
 - Can be assigned (observed) or unassigned (unobserved)
- **Arcs: interactions**
 - Indicate “direct influence” between variables
 - Formally: encode conditional independence (more later)



0.11		0.11
0.11	0.11	0.11
0.11	0.11	0.11

Example: Coin Flips

- N independent coin flips

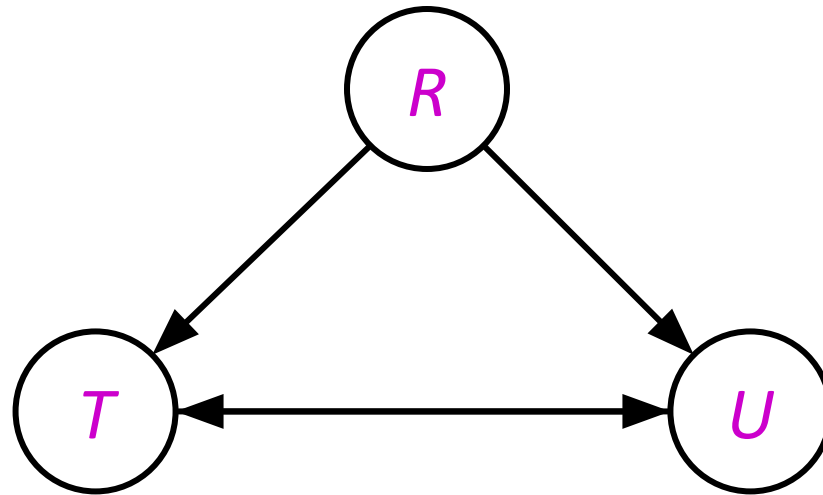
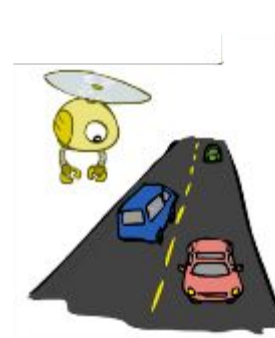


- No interactions between variables: absolute independence

Example: Traffic

- Variables:

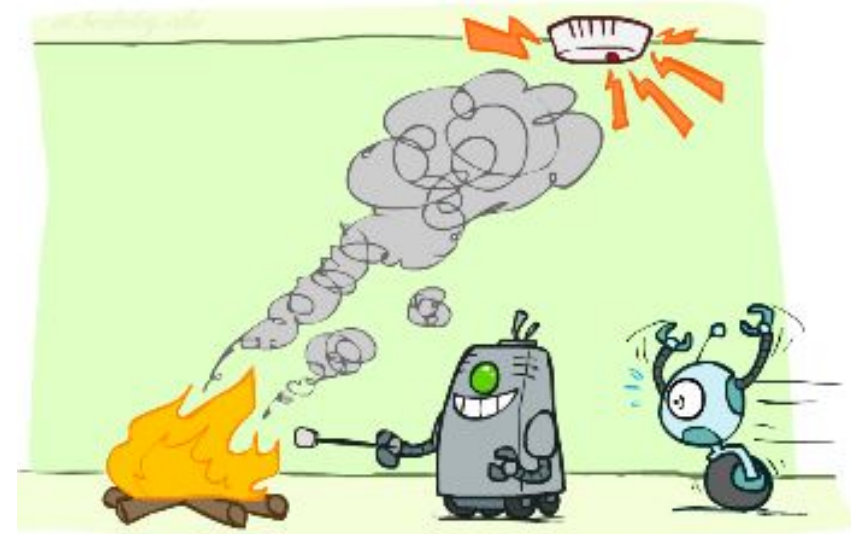
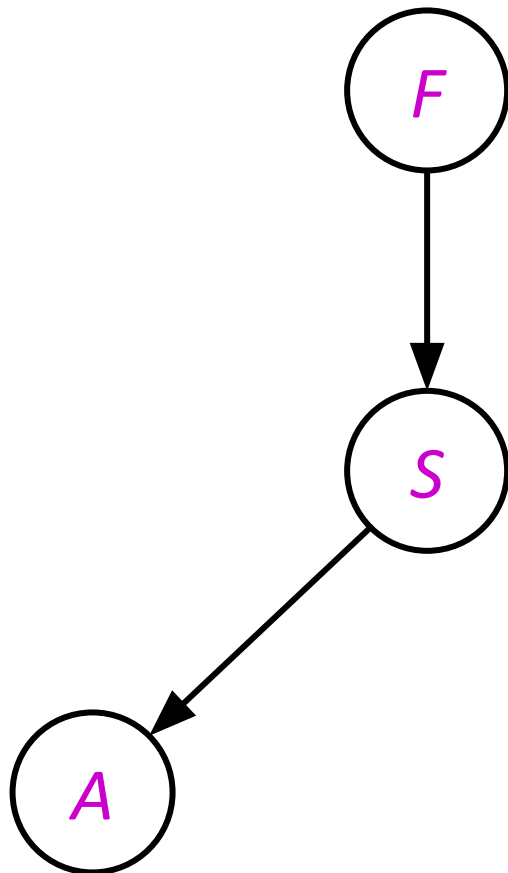
- T: There is traffic
- U: I'm holding my umbrella
- R: It rains



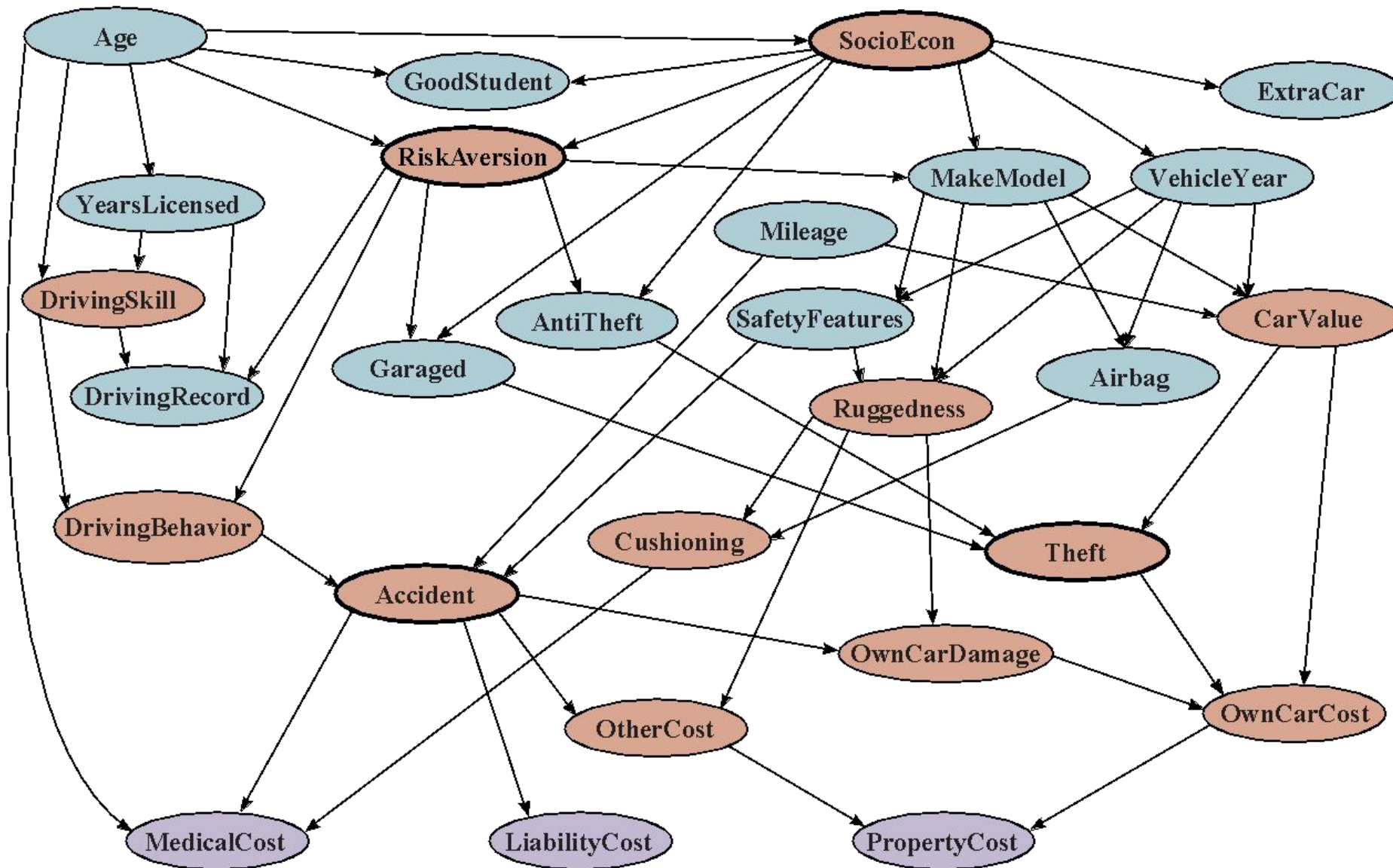
Example: Smoke alarm

- Variables:

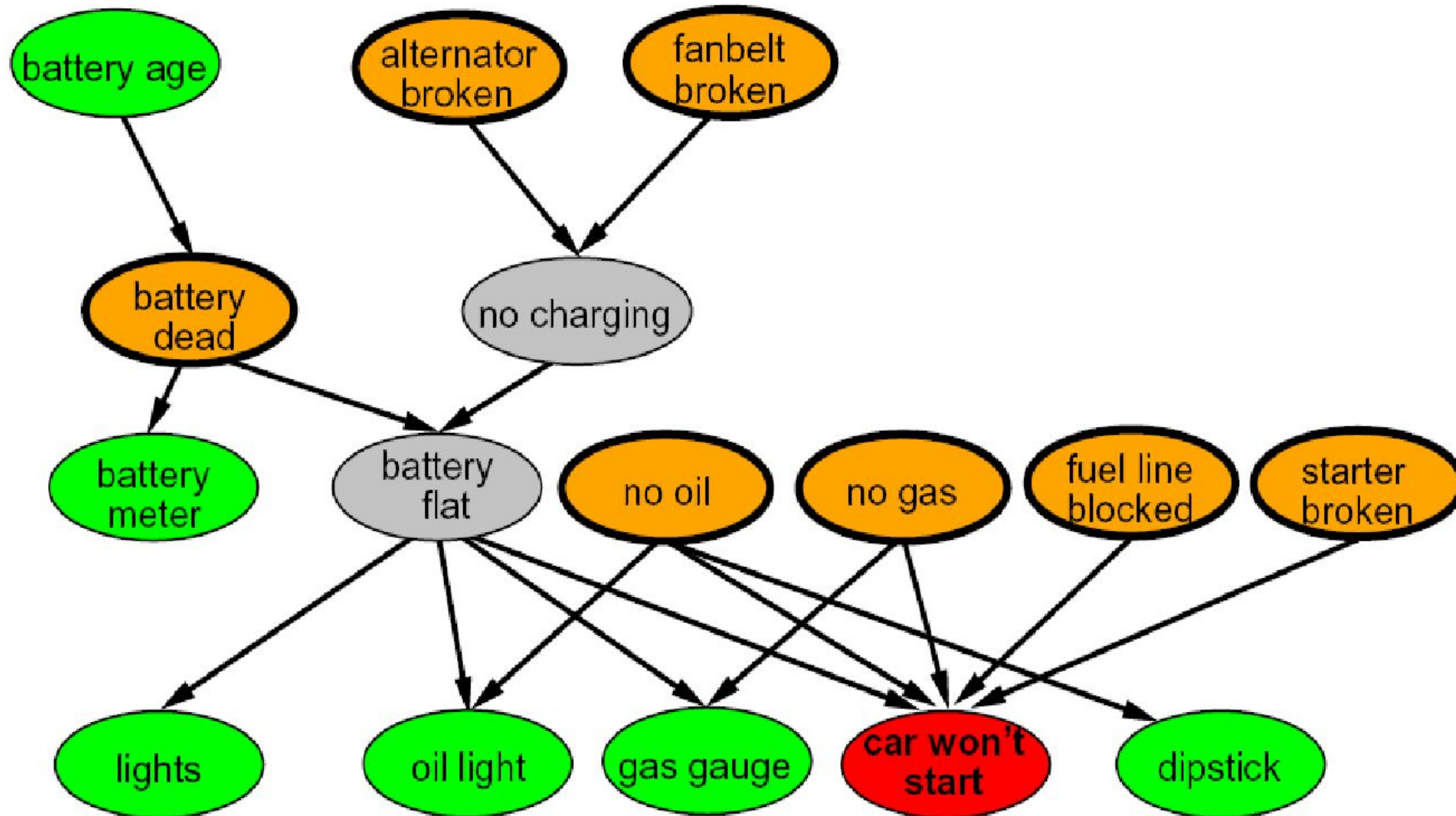
- F: There is fire
- S: There is smoke
- A: Alarm sounds



Example Bayes' Net: Car Insurance



Example Bayes' Net: Car Won't Start



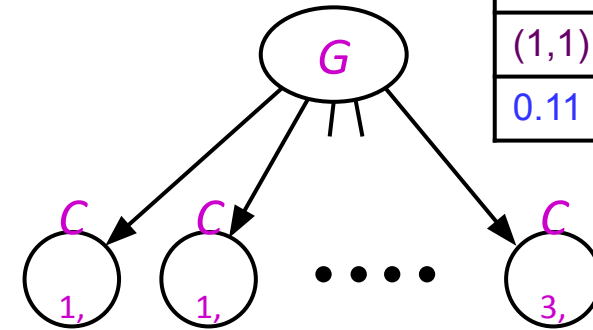
Bayes Net Syntax and Semantics



Bayes Net Syntax



- A set of nodes, one per variable X_i
- A directed, acyclic graph
- A conditional distribution for each node given its **parent variables** in the graph
 - **CPT** (conditional probability table); each row is a distribution for child given values of its parents



P(G)			
(1,1)	(1,2)	(1,3)	...
0.11	0.11	0.11	...

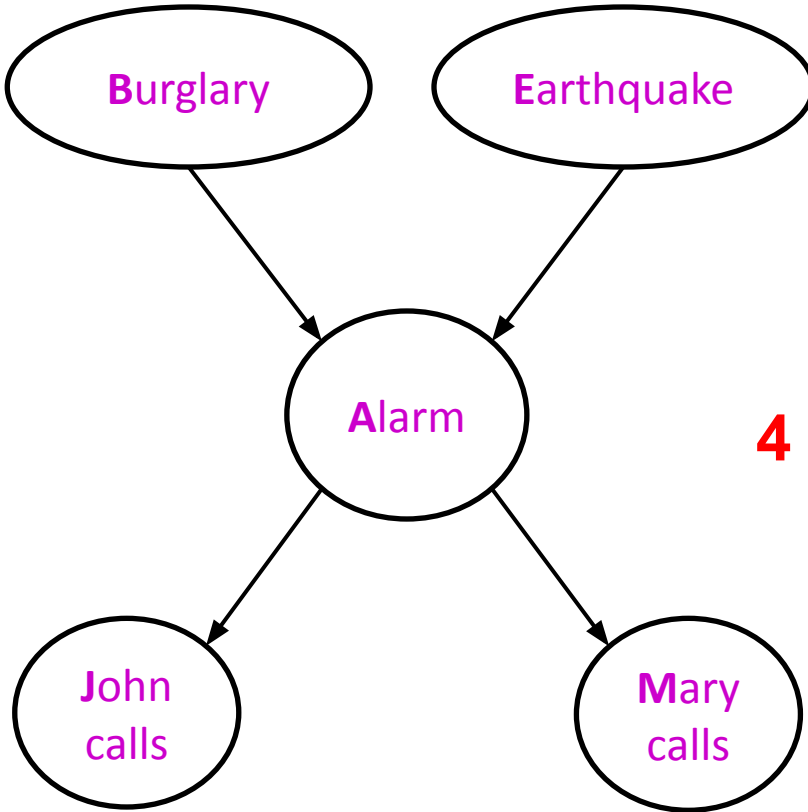
G	P(C _{1,1} G)			
	g	y	o	r
(1,1)	0.01	0.1	0.3	0.59
(1,2)	0.1	0.3	0.5	0.1
(1,3)	0.3	0.5	0.19	0.01
...				

Bayes net = Topology (graph) + Local Conditional Probabilities

Example: Alarm Network

P(B)	
true	false
0.001	0.999

1

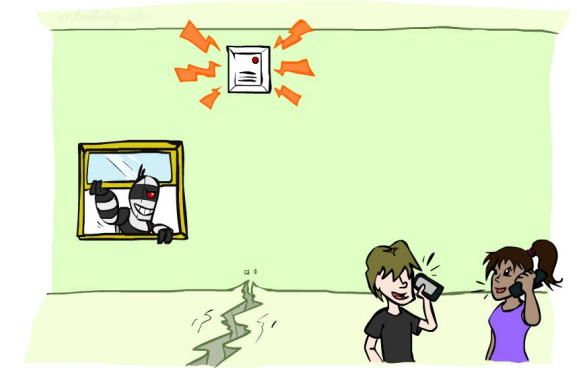


P(E)	
true	false
0.002	0.998

1

B	E	P(A B,E)	
		true	false
true	true	0.95	0.05
true	false	0.94	0.06
false	true	0.29	0.71
false	false	0.001	0.999

4



A	P(J A)	
	true	false
true	0.9	0.1
false	0.05	0.95

2

A	P(M A)	
	true	false
true	0.7	0.3
false	0.01	0.99

2

Number of *free parameters* in each CPT:

Parent range sizes d_1, \dots, d_k

Child range size d

Each table row must sum to 1

$$(d-1) \prod_i d_i$$

General formula for sparse BNs

- Suppose
 - n variables
 - Maximum range size is d
 - Maximum number of parents is k
- Full joint distribution has size $O(d^n)$
- Bayes net has size $O(n \cdot d^k)$
 - Linear scaling with n as long as causal structure is local

Bayes net global semantics

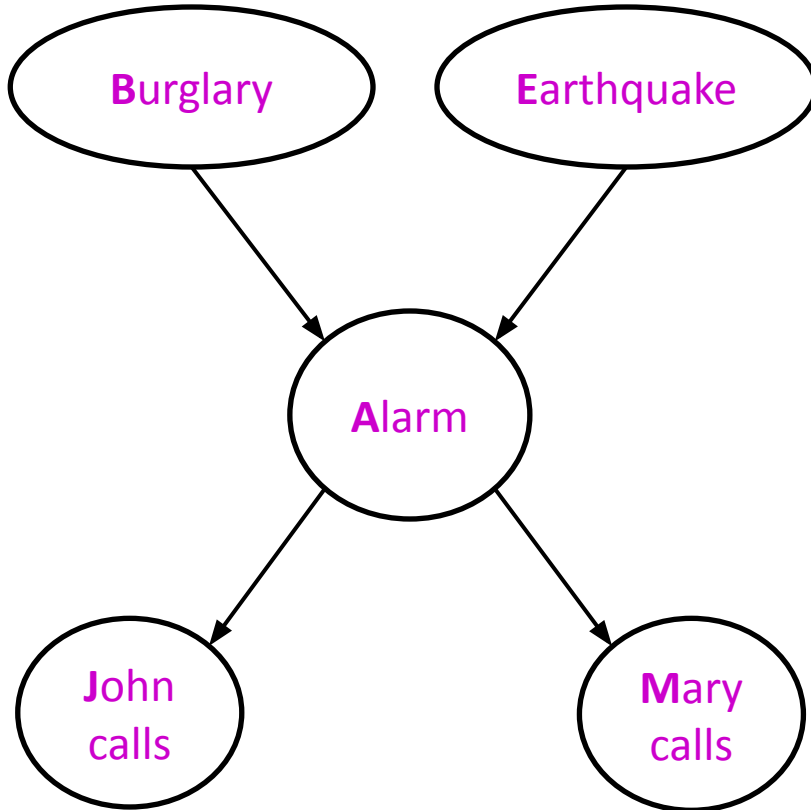


- Bayes nets encode joint distributions as product of conditional distributions on each variable:

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \text{Parents}(X_i))$$

Example

P(B)	
true	false
0.001	0.999



P(E)	
true	false
0.002	0.998

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

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

$$= .001 \times .998 \times .94 \times .1 \times .3 = .000028$$

B	E	P(A B,E)	
		true	false
true	true	0.95	0.05
true	false	0.94	0.06
false	true	0.29	0.71
false	false	0.001	0.999

A	P(J A)	
	true	false
true	0.9	0.1
false	0.05	0.95

A	P(M A)	
	true	false
true	0.7	0.3
false	0.01	0.99

Conditional independence in BNs



- Compare the Bayes net global semantics

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \text{Parents}(X_i))$$

with the chain rule identity

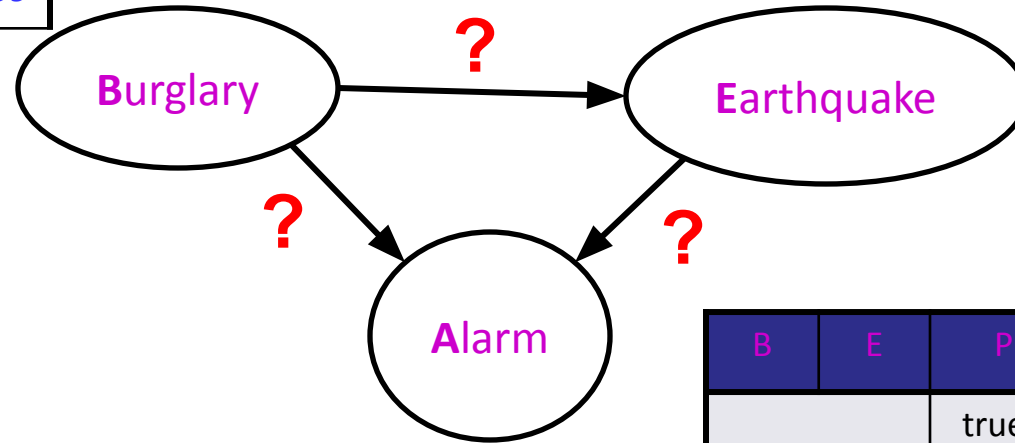
$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid X_1, \dots, X_{i-1})$$

- Assume (without loss of generality) that X_1, \dots, X_n sorted in topological order according to the graph (i.e., parents before children), so $\text{Parents}(X_i) \subseteq X_1, \dots, X_{i-1}$
- So the Bayes net asserts conditional independences $P(X_i \mid X_1, \dots, X_{i-1}) = P(X_i \mid \text{Parents}(X_i))$
 - To ensure these are valid, choose parents for node X_i that “shield” it from other predecessors

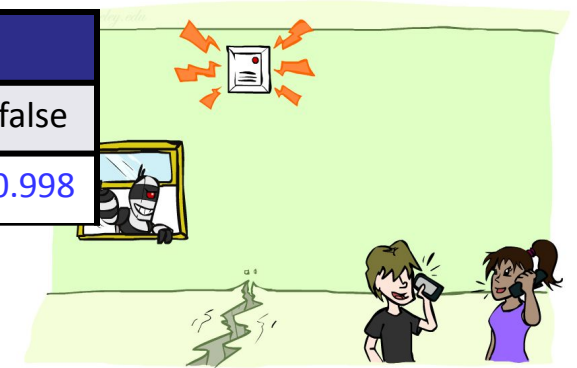
Example: Burglary

- Burglary
- Earthquake
- Alarm

P(B)	
true	false
0.001	0.999



P(E)	
true	false
0.002	0.998

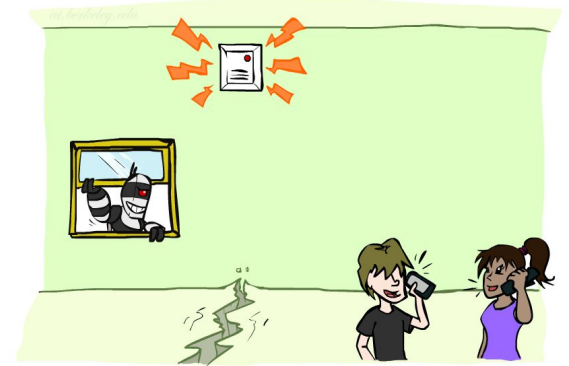
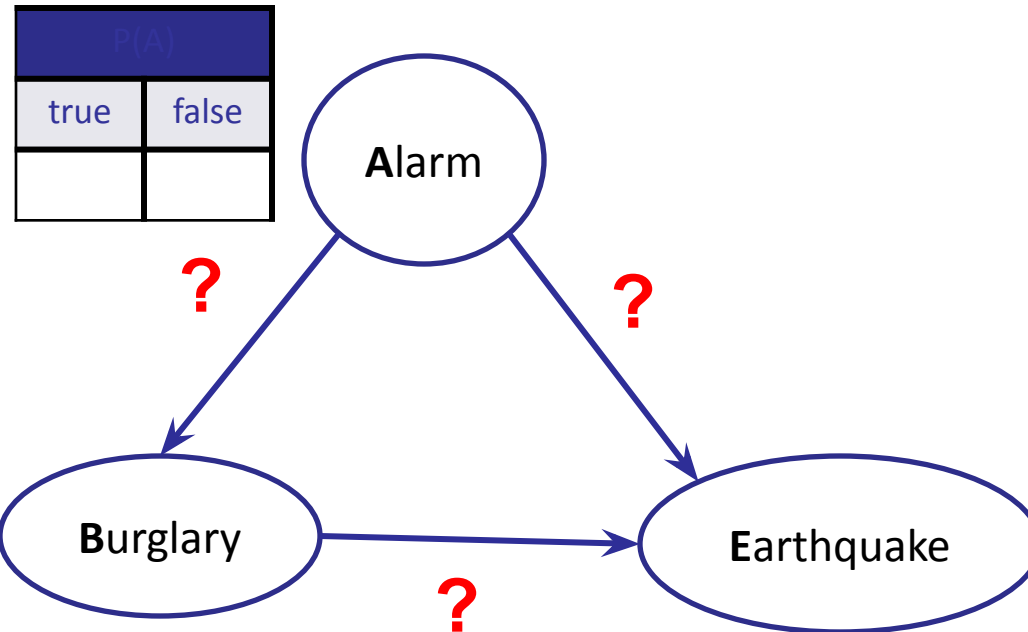


B	E	P(A B,E)	
		true	false
true	true	0.95	0.05
true	false	0.94	0.06
false	true	0.29	0.71
false	false	0.001	0.999

Example: Burglary

- Alarm
- Burglary
- Earthquake

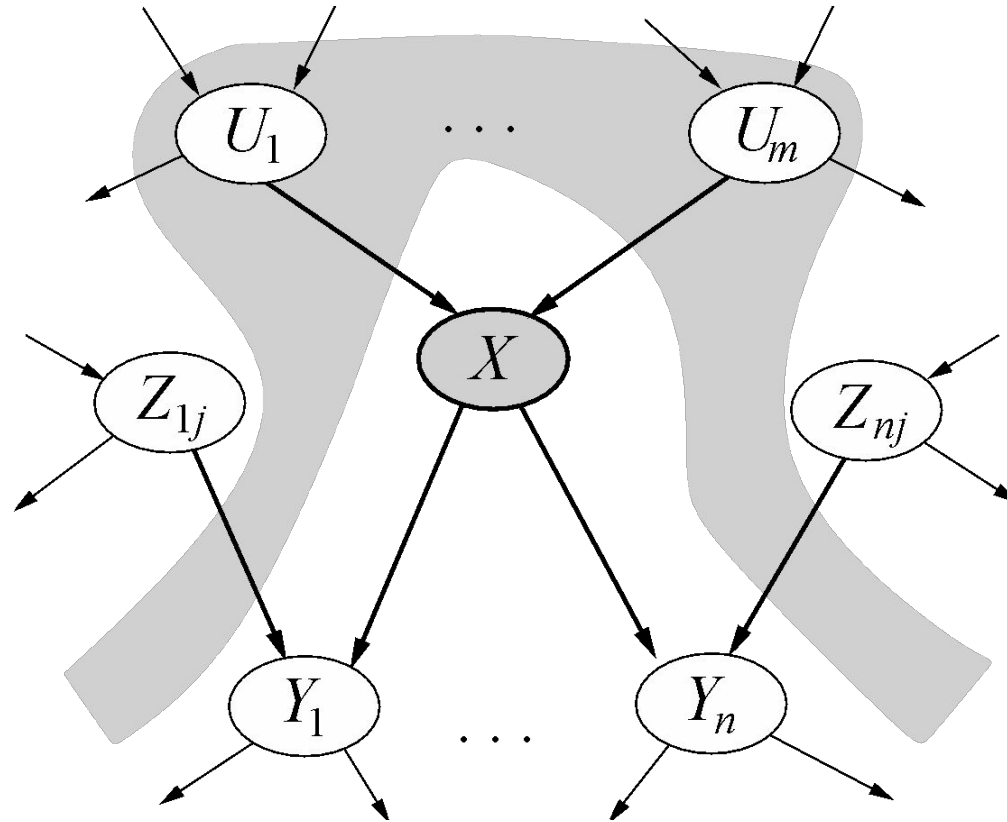
A	P(E A)	
	true	false
true	?	
false		



A	B	P(E A,B)	
		true	false
true	true	?	
true	false		
false	true		
false	false		

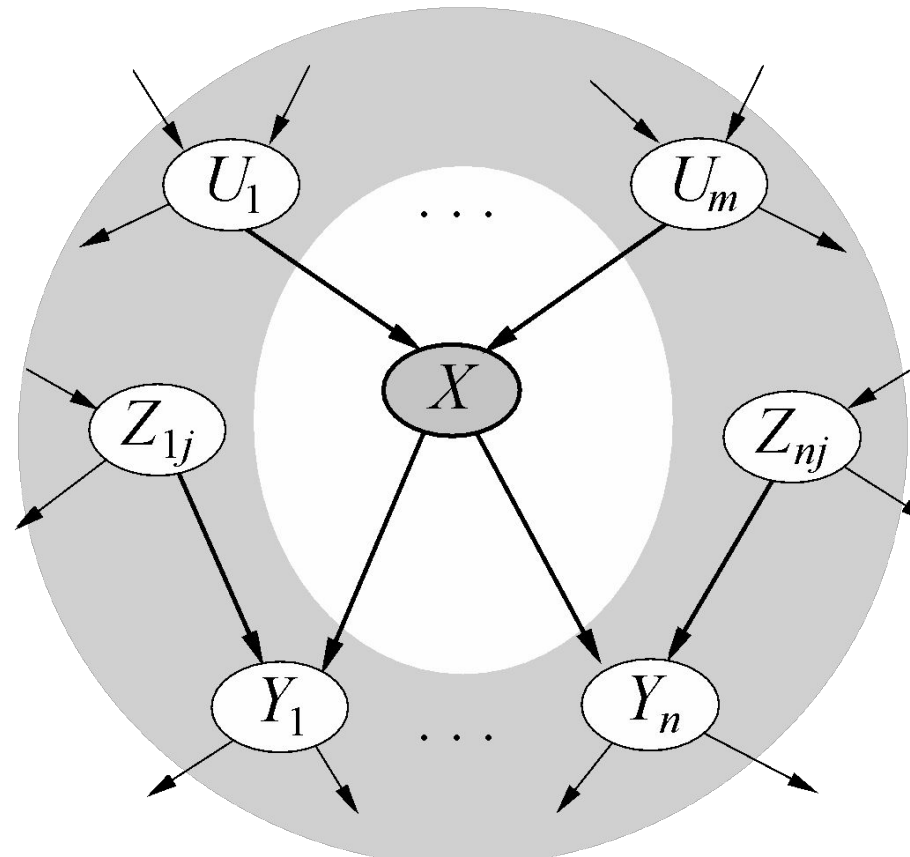
Conditional independence semantics

- *Every variable is conditionally independent of its non-descendants given its parents*
- Conditional independence semantics \Leftrightarrow global semantics



Markov blanket

- A variable's Markov blanket consists of parents, children, children's other parents
- ***Every variable is conditionally independent of all other variables given its Markov blanket***



Summary

- Independence and conditional independence are important forms of probabilistic knowledge
- Bayes net encode joint distributions efficiently by taking advantage of conditional independence
 - Global joint probability = product of local conditionals
- Next: how to answer queries, i.e., compute conditional probabilities of queries given evidence

