

# 343H: Honors AI

Lecture 16: Bayes Nets Inference

3/20/2014

Kristen Grauman

UT Austin

Slides courtesy of Dan Klein, UC Berkeley

# Survey feedback - thank you!

---

- Reading/exercise deadline time
- Web page ease of use
- Programming assignments
  - More project debriefing after deadline
  - Contest rankings beyond top 3
  - Some would like less skeleton, more creativity
  - Python programming standards

# Survey feedback - thank you!

---

- Lecture slides – include answers
- Office hours
- Examples in class lecture
- Textbook

# Announcements

---

- Reading/exercise assignments for next week posted – choose one of the 2 exercises and provide reading response
- PS4 out next week

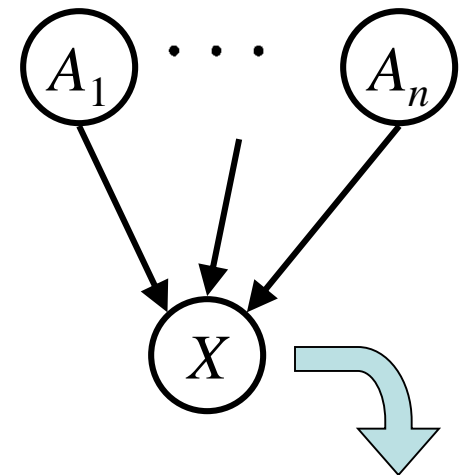
# Bayes' Net Semantics

- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node
  - A collection of distributions over  $X$ , one for each combination of parents' values

$$P(X|a_1 \dots a_n)$$

- Bayes' nets implicitly encode joint distributions
  - As a product of local conditional distributions

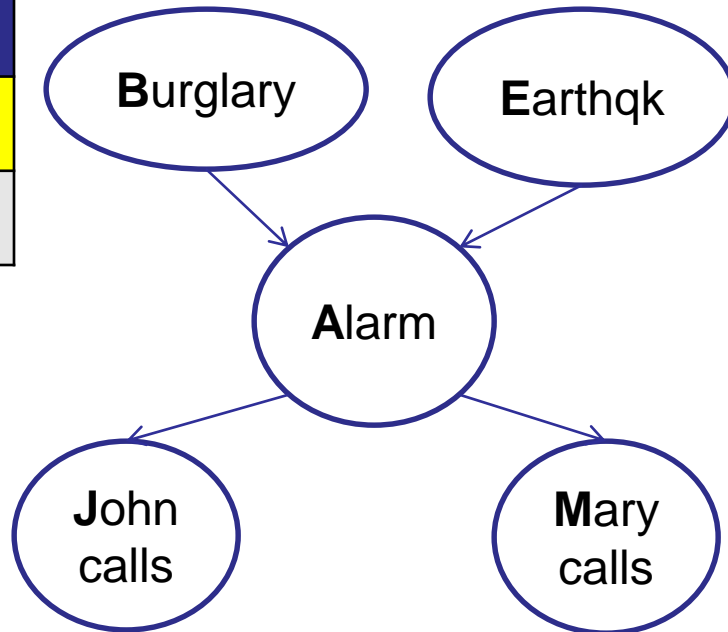
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$



$$P(X|A_1 \dots A_n)$$

# Example: Alarm Network

B	P(B)
+b	0.001
-b	0.999



E	P(E)
+e	0.002
-e	0.998

A	J	P(J A)
+a	+j	0.9

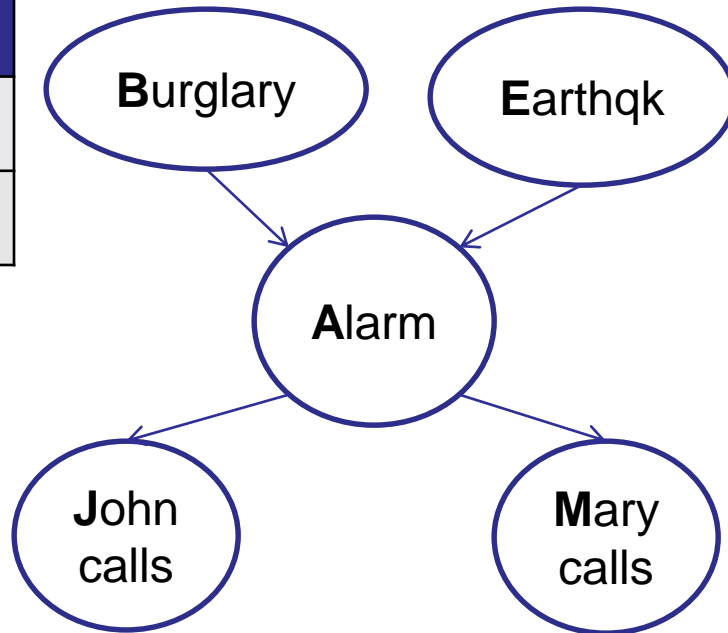
A	M	P(M A)
+a	+m	0.7

B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-e	+a	0.94
+b	-e	-a	0.06
-b	+e	+a	0.29

$$\begin{aligned}
 P(+b, -e, +a, -j, +m) &= \\
 P(+b) P(-e) P(+a \mid +b, -e) P(-j \mid +a) P(+m \mid +a) &= \\
 0.001 \times 0.998 \times 0.94 \times 0.1 \times 0.7 &
 \end{aligned}$$

# Example: Alarm Network

B	P(B)
+b	0.001
-b	0.999



E	P(E)
+e	0.002
-e	0.998

A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-e	+a	0.94
+b	-e	-a	0.06
-b	+e	+a	0.29
-b	+e	-a	0.71
-b	-e	+a	0.001
-b	-e	-a	0.999

# Bayes' Nets

---



Representation



Conditional independences

- Probabilistic inference
  - Enumeration (exact, exponential complexity)
  - Variable elimination (exact, worst-case exponential complexity, often better)
  - Inference is NP-complete
  - Sampling (approximate)
- Learning Bayes' Nets from data



# Inference

---

- Inference: calculating some useful quantity from a joint probability distribution

- Examples:

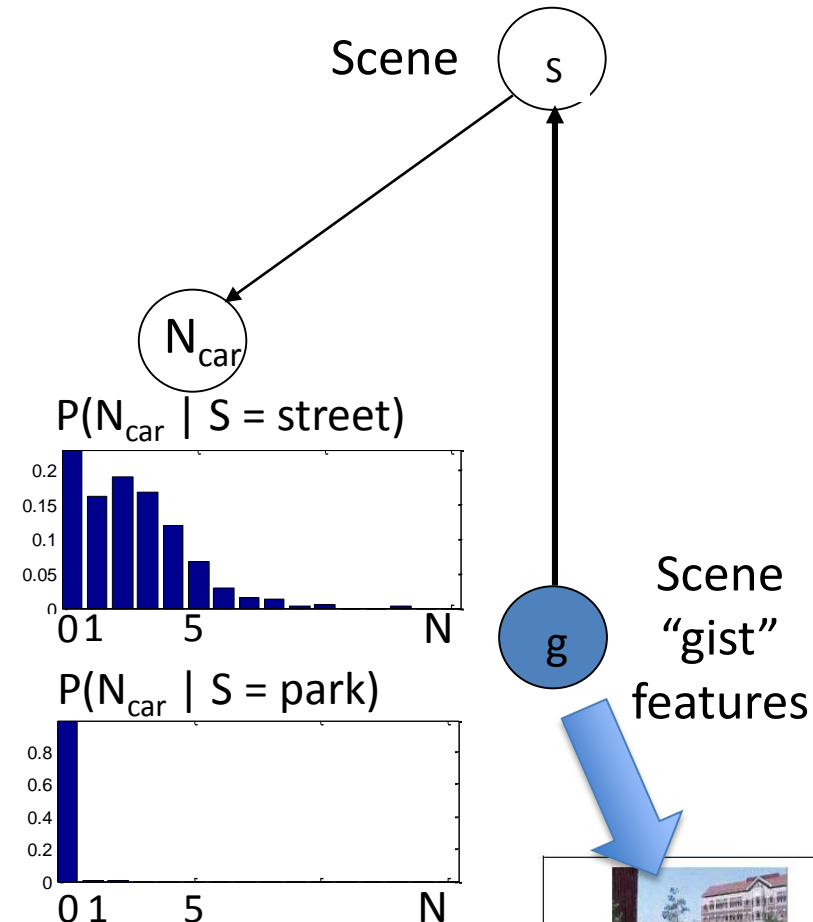
- Posterior probability:

$$P(Q|E_1 = e_1, \dots, E_k = e_k)$$

- Most likely explanation:

$$\operatorname{argmax}_q P(Q = q|E_1 = e_1 \dots)$$

# Recognizing objects in context



# Inference by Enumeration

---

- Given unlimited time, inference in BNs is easy
- Recipe:
  - State the marginal probabilities you need
  - Figure out ALL the atomic probabilities you need
  - Calculate and combine them

# Recall: Inference by Enumeration

- General case:

- Evidence variables:  $E_1 \dots E_k = e_1 \dots e_k$
  - Query\* variable:  $Q$
  - Hidden variables:  $H_1 \dots H_r$
- $\left. \vphantom{\begin{matrix} E_1 \dots E_k = e_1 \dots e_k \\ Q \\ H_1 \dots H_r \end{matrix}} \right\} \begin{matrix} X_1, X_2, \dots, X_n \\ \text{All variables} \end{matrix}$

- We want:  $P(Q|e_1 \dots e_k)$

- Select the entries consistent with the evidence
- Sum out H to get joint of Query and evidence:

$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} \underbrace{P(Q, h_1 \dots h_r, e_1 \dots e_k)}_{X_1, X_2, \dots, X_n}$$

- Normalize

$$Z = \sum_q P(Q, e_1, \dots, e_k)$$

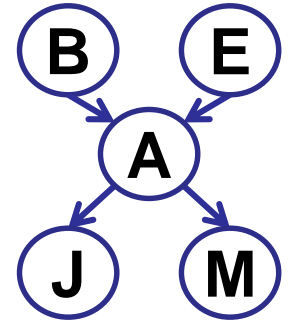
$$P(Q|e_1, \dots, e_k) = \frac{1}{Z} \sum_q P(Q, e_1, \dots, e_k)$$

*\* Works fine with multiple query variables, too*

# Example: Enumeration

---

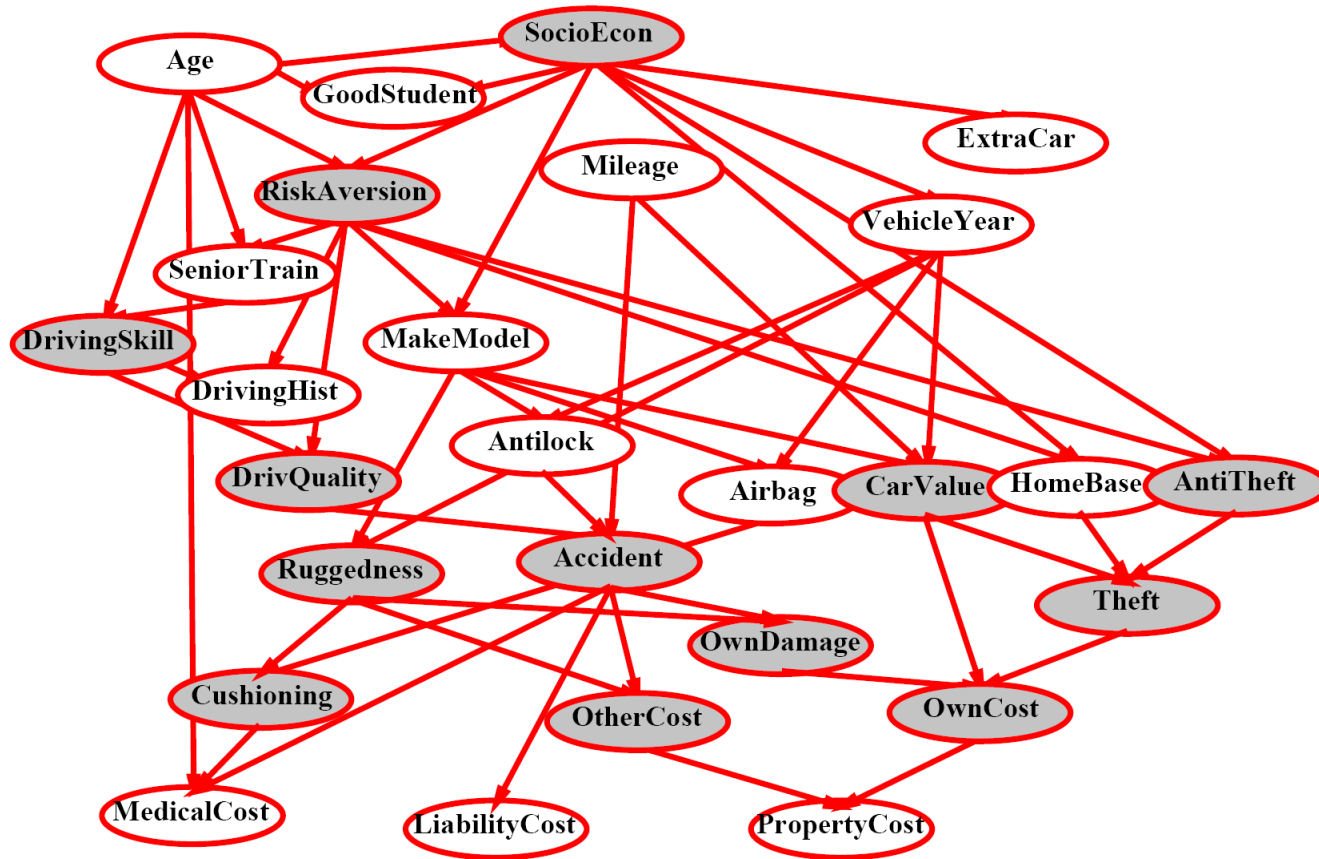
$$P(+b | +j, +m) = \frac{P(+b, +j, +m)}{P(+j, +m)}$$



$$P(+b, +j, +m) =$$

$$\begin{aligned} &P(+b)P(+e)P(+a|+b, +e)P(+j|+a)P(+m|+a) + \\ &P(+b)P(+e)P(-a|+b, +e)P(+j|-a)P(+m|-a) + \\ &P(+b)P(-e)P(+a|+b, -e)P(+j|+a)P(+m|+a) + \\ &P(+b)P(-e)P(-a|+b, -e)P(+j|-a)P(+m|-a) \end{aligned}$$

# Inference by Enumeration?



# Inference by Enumeration vs. Variable Elimination

---

- Why is inference by enumeration so slow?
  - You join up the whole joint distribution before you sum out the hidden variables
- Idea: interleave joining and marginalizing!
  - Called “Variable Elimination”
  - Still NP-hard, but usually much faster than inference by enumeration

# Factor Zoo I

---

- **Joint distribution:  $P(X,Y)$** 
  - Entries  $P(x,y)$  for all  $x, y$
  - Sums to 1
- **Selected joint:  $P(x,Y)$** 
  - A slice of the joint distribution
  - Entries  $P(x,y)$  for fixed  $x$ , all  $y$
  - Sums to  $P(x)$
- **Note: Number of capitals => size of the table**

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

$P(\text{cold}, W)$

T	W	P
cold	sun	0.2
cold	rain	0.3



# Factor Zoo II

$$P(W|T)$$

T	W	P
hot	sun	0.8
hot	rain	0.2
cold	sun	0.4
cold	rain	0.6

$$P(W|hot)$$

$$P(W|cold)$$

- Family of conditionals:  $P(X|Y)$

- Multiple conditionals
- Entries  $P(x|y)$  for all  $x, y$
- Sums to  $|Y|$

- Single conditional:  $P(Y|x)$

- Entries  $P(y|x)$  for fixed  $x$ , all  $y$
- Sums to 1

$$P(W|cold)$$

T	W	P
cold	sun	0.4
cold	rain	0.6

# Factor Zoo III

---

- Specified family:  $P(y | X)$ 
  - Entries  $P(y | x)$  for fixed  $y$ , but for all  $x$
  - Sums to ... who knows!

$$P(\text{rain} | T)$$

T	W	P
hot	rain	0.2
cold	rain	0.6

}  $P(\text{rain} | \text{hot})$   
}  $P(\text{rain} | \text{cold})$

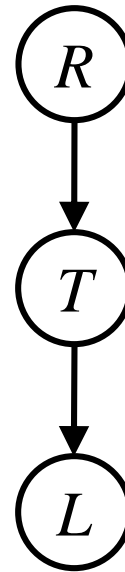
# Factor Zoo Summary

---

- In general, when we write  $P(Y_1 \dots Y_N | X_1 \dots X_M)$ 
  - It is a “factor,” a multi-dimensional array
  - Its values are all  $P(y_1 \dots y_N | x_1 \dots x_M)$
  - Any assigned X or Y is a dimension missing (selected) from the array

# Example: Traffic Domain

- Random Variables
  - R: Raining
  - T: Traffic
  - L: Late for class!



$$P(R)$$

+r	0.1
-r	0.9

$$P(T|R)$$

+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

$$P(L|R)$$

+t	+l	0.3
+t	-l	0.7
-t	+l	0.1
-t	-l	0.9

# Variable Elimination Outline

- Track objects called **factors**
- Initial factors are local CPTs (one per node)

+r	0.1
-r	0.9

+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

+t	+l	0.3
+t	-l	0.7
-t	+l	0.1
-t	-l	0.9

- Any known values are selected
  - E.g. if we know  $L = +\ell$ , the initial factors are

+r	0.1
-r	0.9

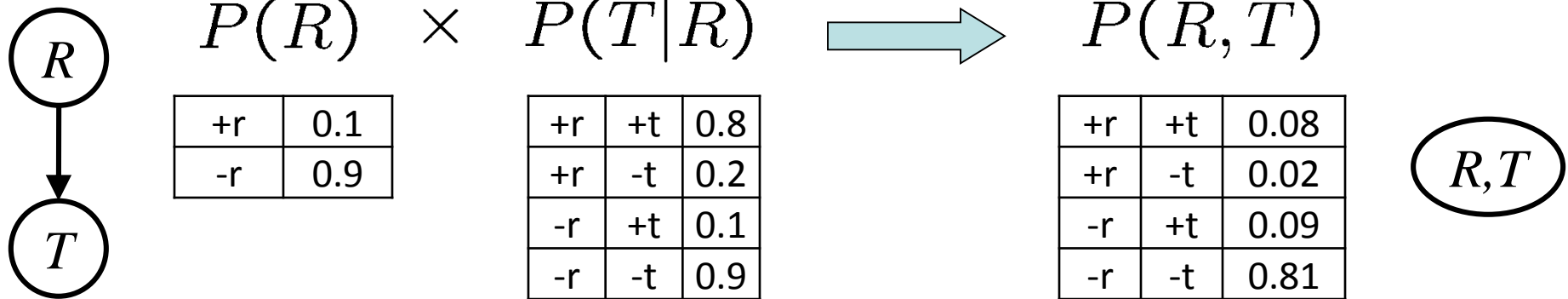
+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

+t	+l	0.3
-t	+l	0.1

- VE: Alternately join factors and eliminate variables

# Operation 1: Join Factors

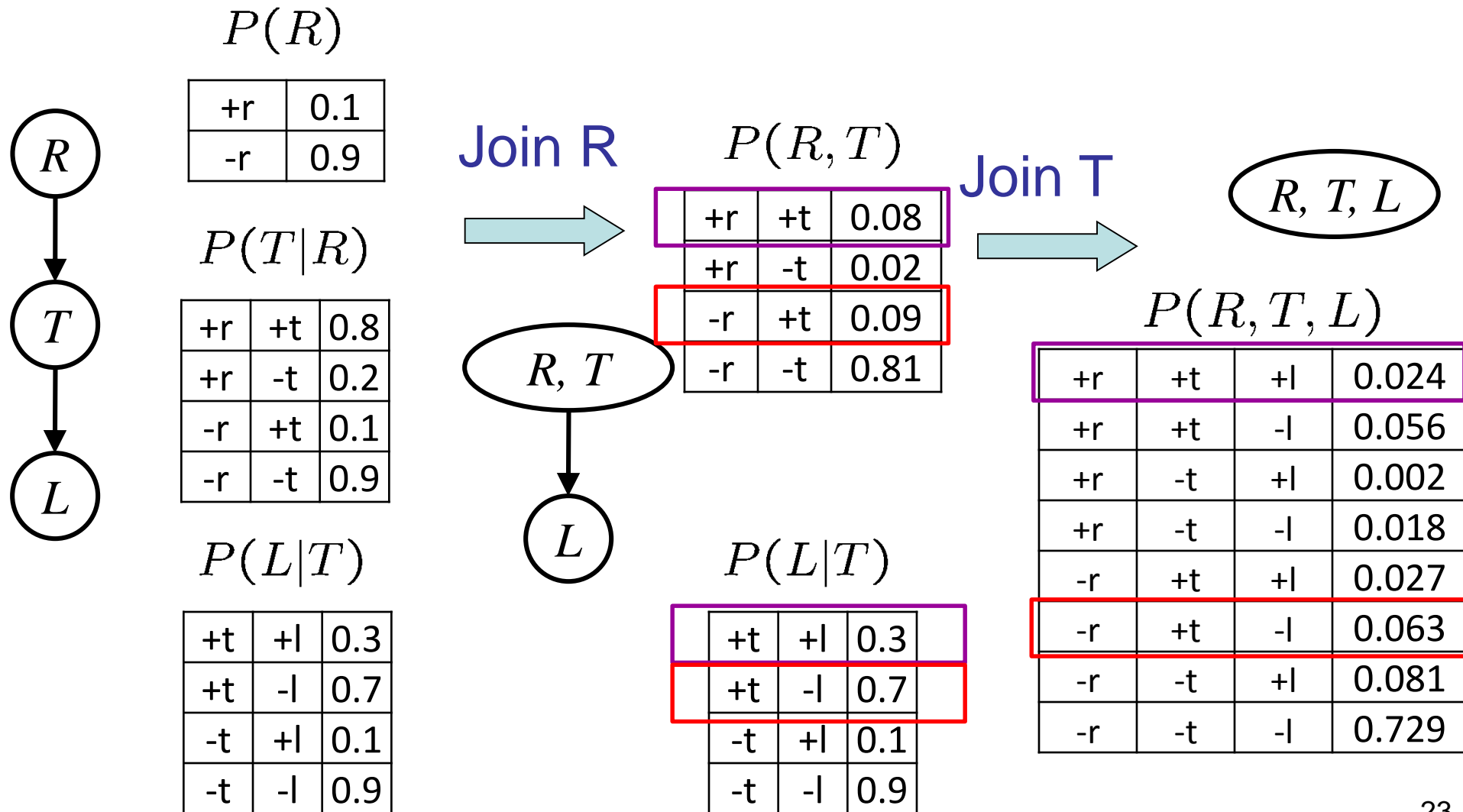
- First basic operation: **joining factors**
- Combining factors:
  - Get all factors over the joining variable
  - Build a new factor over the union of the variables involved
- Example: Join on R



- Computation for each entry: pointwise products

$$\forall r, t : P(r, t) = P(r) \cdot P(t|r)$$


# Example: Multiple Joins



# Operation 2: Eliminate

---

- Second basic operation: **marginalization**
- Take a factor and sum out a variable
  - Shrinks a factor to a smaller one
  - A **projection** operation
- Example:

$P(R, T)$				$P(T)$	
+r	+t	0.08	sum $R$ 	+t	0.17
+r	-t	0.02		-t	0.83
-r	+t	0.09			
-r	-t	0.81			



# Multiple Elimination

$R, T, L$

$P(R, T, L)$

+r	+t	+l	0.024
+r	+t	-l	0.056
+r	-t	+l	0.002
+r	-t	-l	0.018
-r	+t	+l	0.027
-r	+t	-l	0.063
-r	-t	+l	0.081
-r	-t	-l	0.729

Sum  
out R



$T, L$

$P(T, L)$

+t	+l	0.051
+t	-l	0.119
-t	+l	0.083
-t	-l	0.747

Sum  
out T

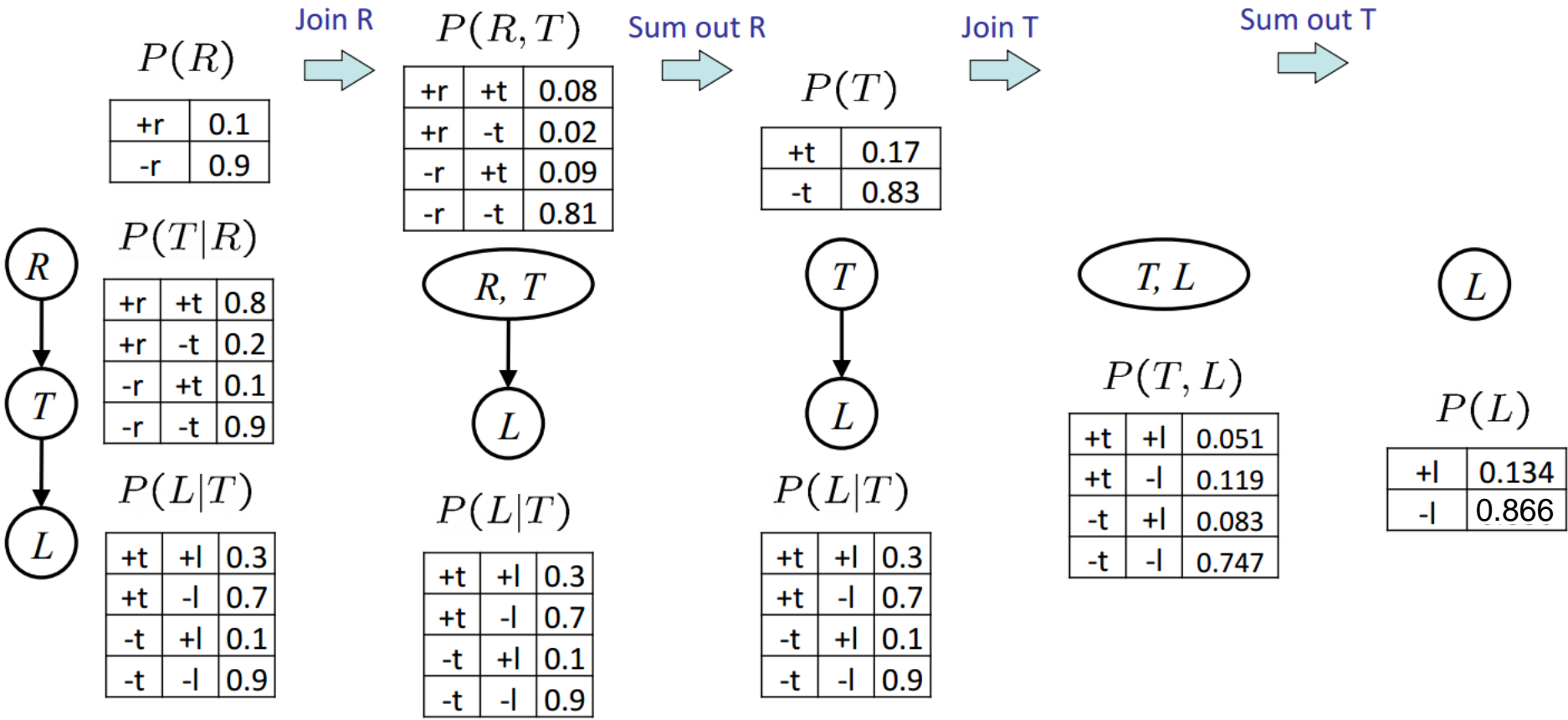


$L$

$P(L)$

+l	0.134
-l	0.886

# Marginalizing early! (aka VE)



# Evidence



- If evidence, start with factors that select that evidence
  - No evidence uses these initial factors:

$$P(R)$$

+r	0.1
-r	0.9

$$P(T|R)$$

+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

$$P(L|T)$$

+t	+l	0.3
+t	-l	0.7
-t	+l	0.1
-t	-l	0.9

- Computing  $P(L|+r)$ , the initial factors become:

$$P(+r)$$

+r	0.1
----	-----

$$P(T|+r)$$

+r	+t	0.8
+r	-t	0.2

$$P(L|T)$$

+t	+l	0.3
+t	-l	0.7
-t	+l	0.1
-t	-l	0.9

- We eliminate all vars other than query + evidence

# Evidence II



- Result will be a selected joint of query and evidence
  - E.g. for  $P(L \mid +r)$ , we'd end up with:

$P(+r, L)$			Normalize	$P(L \mid +r)$	
+r	+l	0.026	→	+l	0.26
+r	-l	0.074		-l	0.74

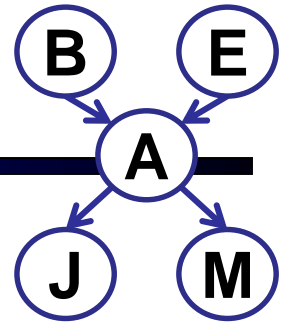
- To get our answer, just normalize this!
- That's it!

# General Variable Elimination

---

- Query:  $P(Q|E_1 = e_1, \dots, E_k = e_k)$
- Start with initial factors:
  - Local CPTs (but instantiated by evidence)
- While there are still hidden variables (not Q or evidence):
  - Pick a hidden variable H
  - Join all factors mentioning H
  - Eliminate (sum out) H
- Join all remaining factors and normalize

# Example



$$P(B|j, m) \propto P(B, j, m)$$

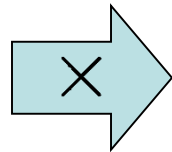
$P(B)$	$P(E)$	$P(A B, E)$	$P(j A)$	$P(m A)$
--------	--------	-------------	----------	----------

Choose A

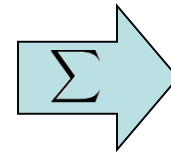
$$P(A|B, E)$$

$$P(j|A)$$

$$P(m|A)$$



$$P(j, m, A|B, E)$$

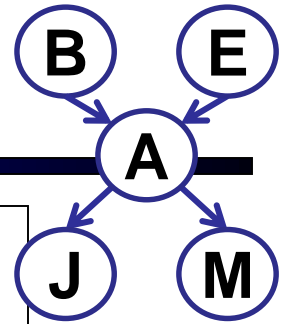


$$P(j, m|B, E)$$

$P(B)$	$P(E)$	$P(j, m B, E)$
--------	--------	----------------

Query:  $P(B|j, m)$

# Example (continued)



$P(B)$	$P(E)$	$P(j, m B, E)$
--------	--------	----------------

Choose E

$$\begin{array}{l} P(E) \\ P(j, m|B, E) \end{array} \xrightarrow{\times} P(j, m, E|B) \xrightarrow{\Sigma} P(j, m|B)$$

$P(B)$	$P(j, m B)$
--------	-------------

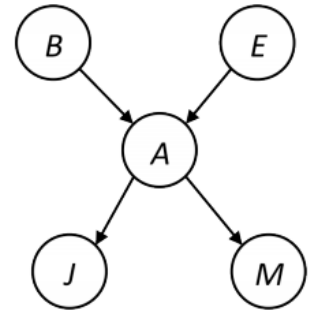
Finish with B

$$\begin{array}{l} P(B) \\ P(j, m|B) \end{array} \xrightarrow{\times} P(j, m, B) \xrightarrow{\text{Normalize}} P(B|j, m)$$

# Same example in equations

$$P(B|j, m) \propto P(B, j, m)$$

$P(B)$	$P(E)$	$P(A B, E)$	$P(j A)$	$P(m A)$
--------	--------	-------------	----------	----------



$$\begin{aligned}
 P(B|j, m) &\propto P(B, j, m) \\
 &= \sum_{e,a} P(B, j, m, e, a) \\
 &= \sum_{e,a} P(B)P(e)P(a|B, e)P(j|a)P(m|a) \\
 &= \sum_e P(B)P(e) \sum_a P(a|B, e)P(j|a)P(m|a) \\
 &= \sum_e P(B)P(e)f_1(B, e, j, m) \\
 &= P(B) \sum_e P(e)f_1(B, e, j, m) \\
 &= P(B)f_2(B, j, m)
 \end{aligned}$$

marginal can be obtained from joint by summing out

use Bayes' net joint distribution expression

use  $x*(y+z) = xy + xz$

joining on a, and then summing out gives  $f_1$

$x*(y+z) = xy + xz$

joining on e, and then summing out gives  $f_2$

We are exploiting:

$$uwy + uwz + uxy + uxz + vwy + vwz + vxy + vxz = (u + v)(w + x)(y+z) \quad 32$$



# Another variable elimination example

Query:  $P(X_3|Y_1 = y_1, Y_2 = y_2, Y_3 = y_3)$

Start by inserting evidence, which gives the following initial factors:

$$p(Z)p(X_1|Z)p(X_2|Z)p(X_3|Z)p(y_1|X_1)p(y_2|X_2)p(y_3|X_3)$$

Eliminate  $X_1$ , this introduces the factor  $f_1(Z, y_1) = \sum_{x_1} p(x_1|Z)p(y_1|x_1)$ , and we are left with:

$$p(Z)f_1(Z, y_1)p(X_2|Z)p(X_3|Z)p(y_2|X_2)p(y_3|X_3)$$

Eliminate  $X_2$ , this introduces the factor  $f_2(Z, y_2) = \sum_{x_2} p(x_2|Z)p(y_2|x_2)$ , and we are left with:

$$p(Z)f_1(Z, y_1)f_2(Z, y_2)p(X_3|Z)p(y_3|X_3)$$

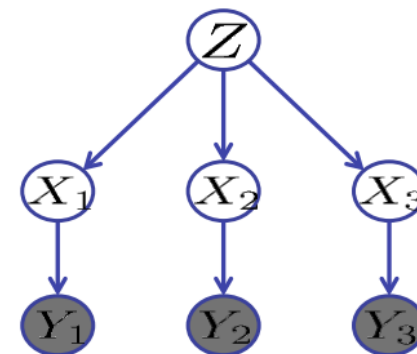
Eliminate  $Z$ , this introduces the factor  $f_3(y_1, y_2, X_3) = \sum_z p(z)f_1(z, y_1)f_2(z, y_2)p(X_3|z)$ , and we are left:

$$p(y_3|X_3), f_3(y_1, y_2, X_3)$$

No hidden variables left. Join the remaining factors to get:

$$f_4(y_1, y_2, y_3, X_3) = P(y_3|X_3)f_3(y_1, y_2, X_3).$$

Normalizing over  $X_3$  gives  $P(X_3|y_1, y_2, y_3)$ .



# Another variable elimination example

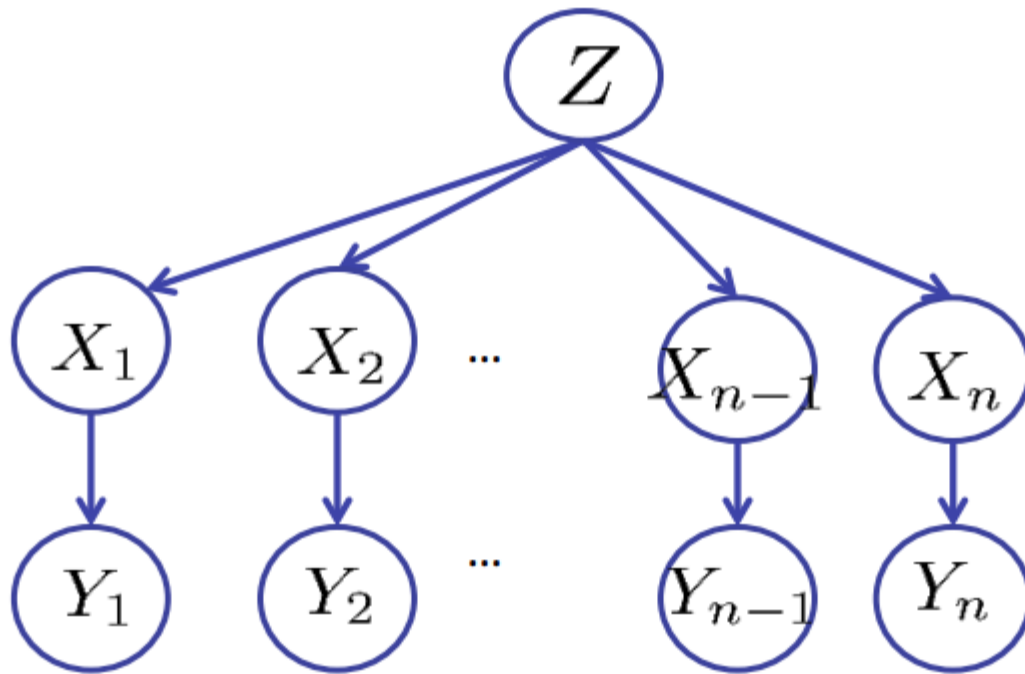
---

- Computational complexity depends on largest factor being generated.
- Size of factor = number of entries in table
- In previous example, assuming all binary variables, all factors are of size 2 – they all have only one variable ( $Z$ ,  $Z$ , and  $X_3$ , respectively)

# Quiz: Variable elimination ordering

---

For the query  $P(X_n \mid y_1, \dots, y_n)$ , what would be a **good** and **bad** ordering for elimination?



# VE: Computational and space complexity

---

- Determined by the largest factor
- Elimination ordering can greatly affect the size of the largest factor
  - e.g., previous example,  $2^n$  vs  $2^2$ .
- Does there always exist an ordering that's good?
  - No.

# Recap: Bayes' Nets

---

- ✓ Representation
- ✓ Conditional independences
  - Probabilistic inference
    - ✓ Enumeration (exact, exponential complexity)
      - Variable elimination (exact, worst-case exponential complexity, often better)
    - ✓ Inference is NP-complete
      - Sampling (approximate)
  - Learning Bayes' Nets from data