

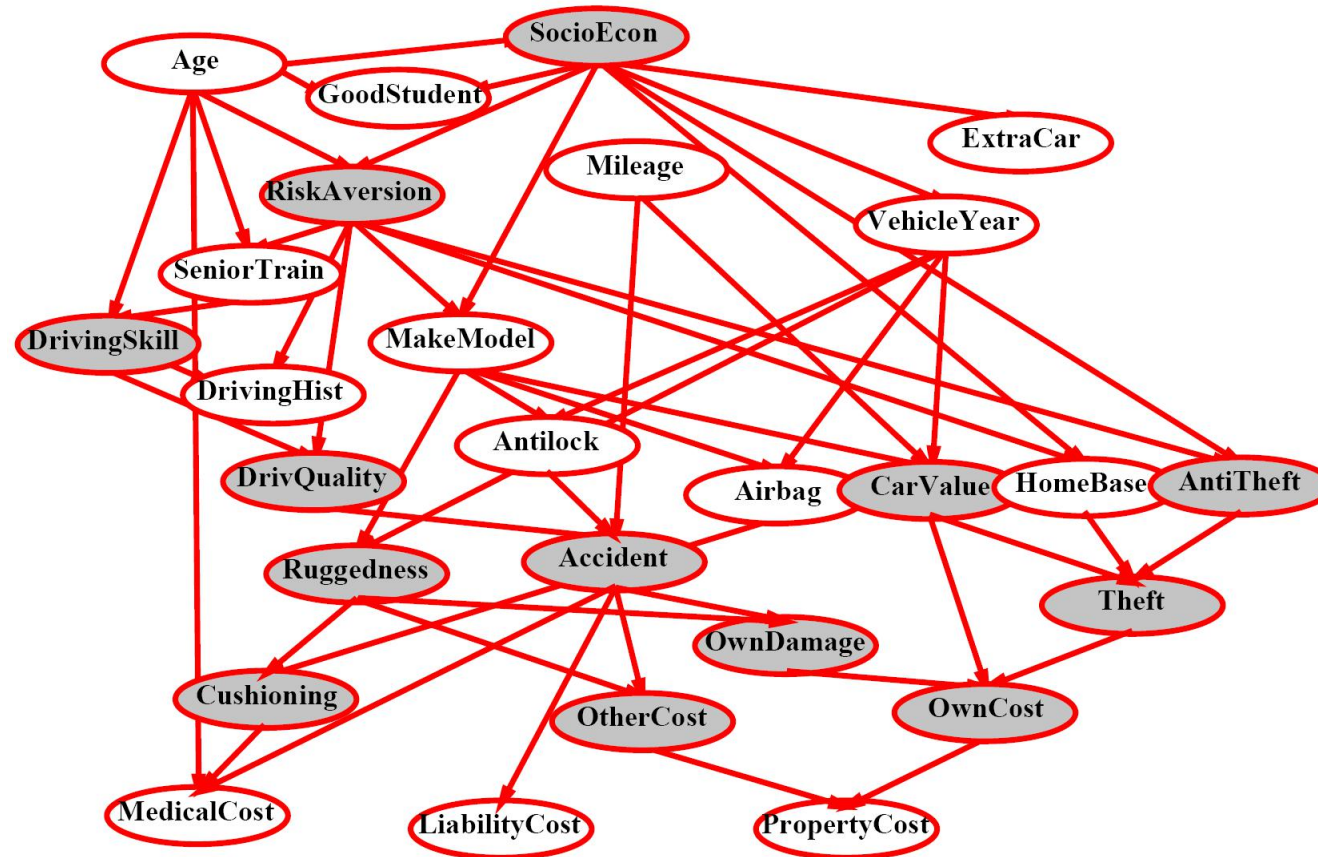
Artificial Intelligence

CS4365 --- Fall 2022

Bayes' Net: Variable Elimination

Instructor: Yunhui Guo

Inference by Enumeration



$$P(\textit{Antilock} | \textit{observed variables}) = ?$$

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 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 $P(x)$
- Number of capitals = dimensionality 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

$P(cold, W)$

T	W	P
cold	sun	0.2
cold	rain	0.3

Factor II

- Single conditional: $P(Y \mid x)$
 - Entries $P(y \mid x)$ for fixed x , all y
 - Sums to 1

$$P(W \mid cold)$$

T	W	P
cold	sun	0.4
cold	rain	0.6

- Family of conditionals: $P(X \mid Y)$
 - Multiple conditionals
 - Entries $P(x \mid y)$ for all x, y
 - Sums to $|Y|$

$$P(W \mid T)$$

T	W	P	
hot	sun	0.8	} $P(W \mid hot)$
hot	rain	0.2	
cold	sun	0.4	} $P(W \mid cold)$
cold	rain	0.6	

Factor III

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

$$P(rain|T)$$

T	W	P	
hot	rain	0.2	} $P(rain hot)$ } $P(rain cold)$
cold	rain	0.6	

Summary

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

Example: Traffic Domain

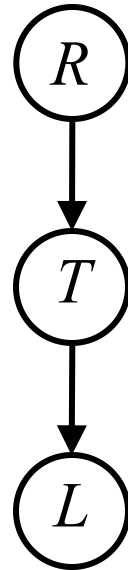
- Random Variables

- R: Raining
- T: Traffic
- L: Late for class

$$P(L) = ?$$

$$= \sum_{r,t} P(r, \mathbf{t}, L)$$

$$= \sum_{r,t} P(r)P(\mathbf{t}|r)P(L|\mathbf{t})$$



$$P(R)$$

+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

Inference by Enumeration: Procedural Outline

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

$$P(R)$$

+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

Any known values are selected

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

$$P(R)$$

+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

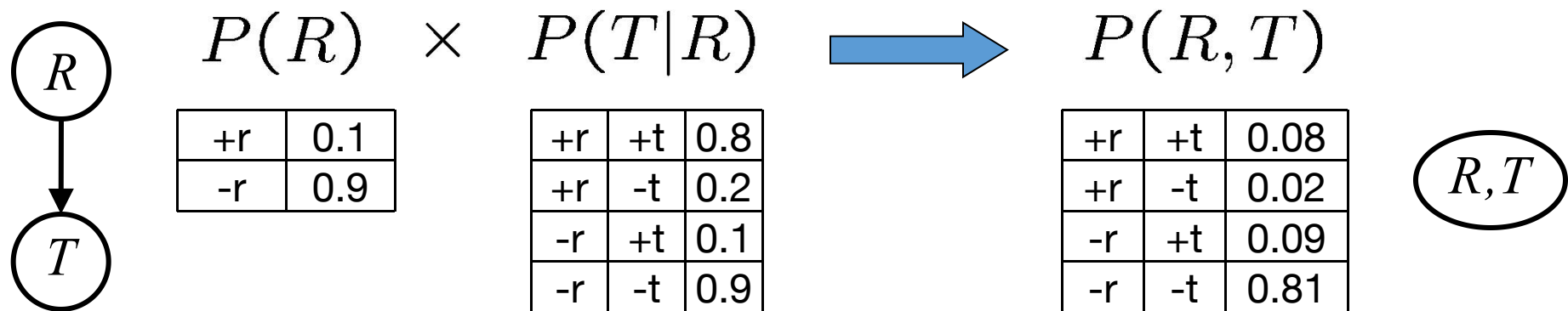
$$P(+\ell|T)$$

+t	+l	0.3
-t	+l	0.1

- Procedure: Join all factors, then eliminate all hidden variables

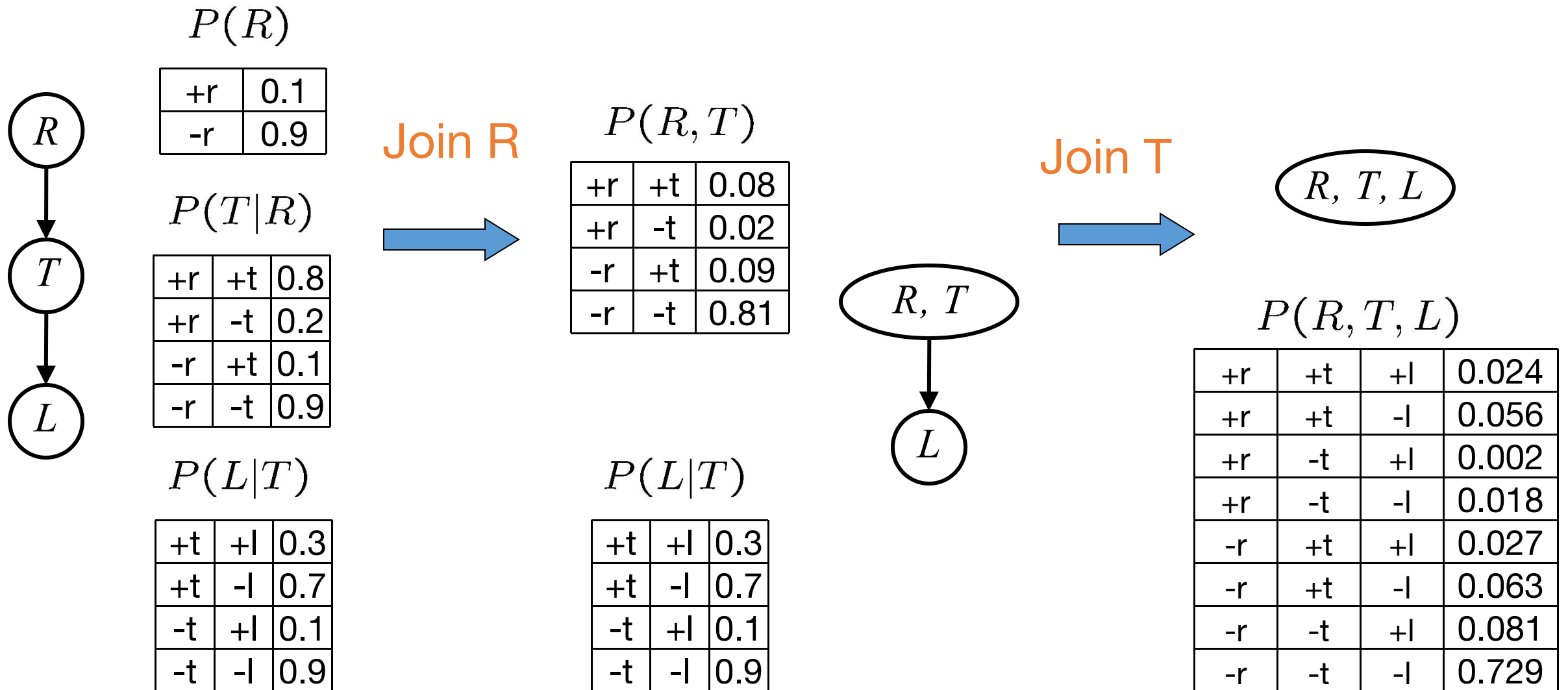
Operation 1: Join Factors

- First basic operation: **joining 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

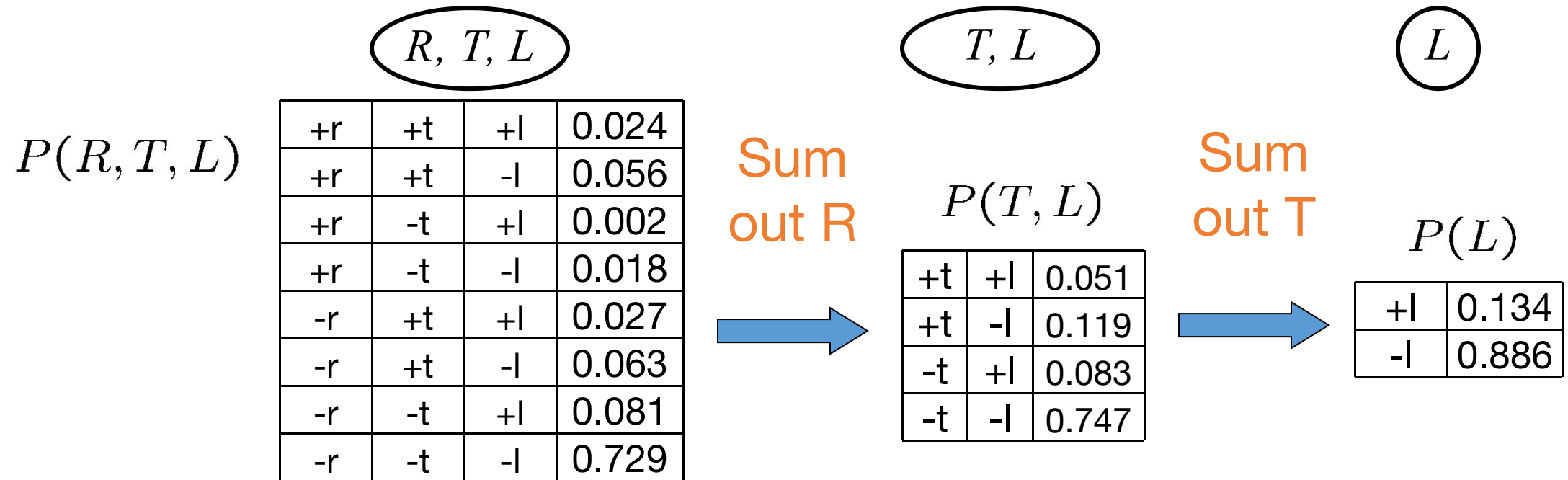


Operation 2: Eliminate

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

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

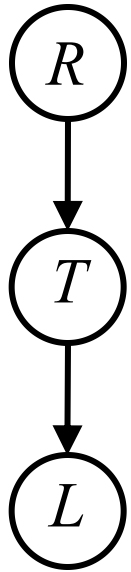
Multiple Elimination



Summary

- Multiple Join, Multiple Eliminate (= Inference by Enumeration)
- Marginalizing Early (= Variable Elimination)

Traffic Domain



$$P(L) = ?$$

• Inference by Enumeration

$$= \sum_t \sum_r \underbrace{P(L|t)P(r)P(t|r)}_{\text{Join on } r}$$

Join on t

Eliminate r

Eliminate t

Variable Elimination

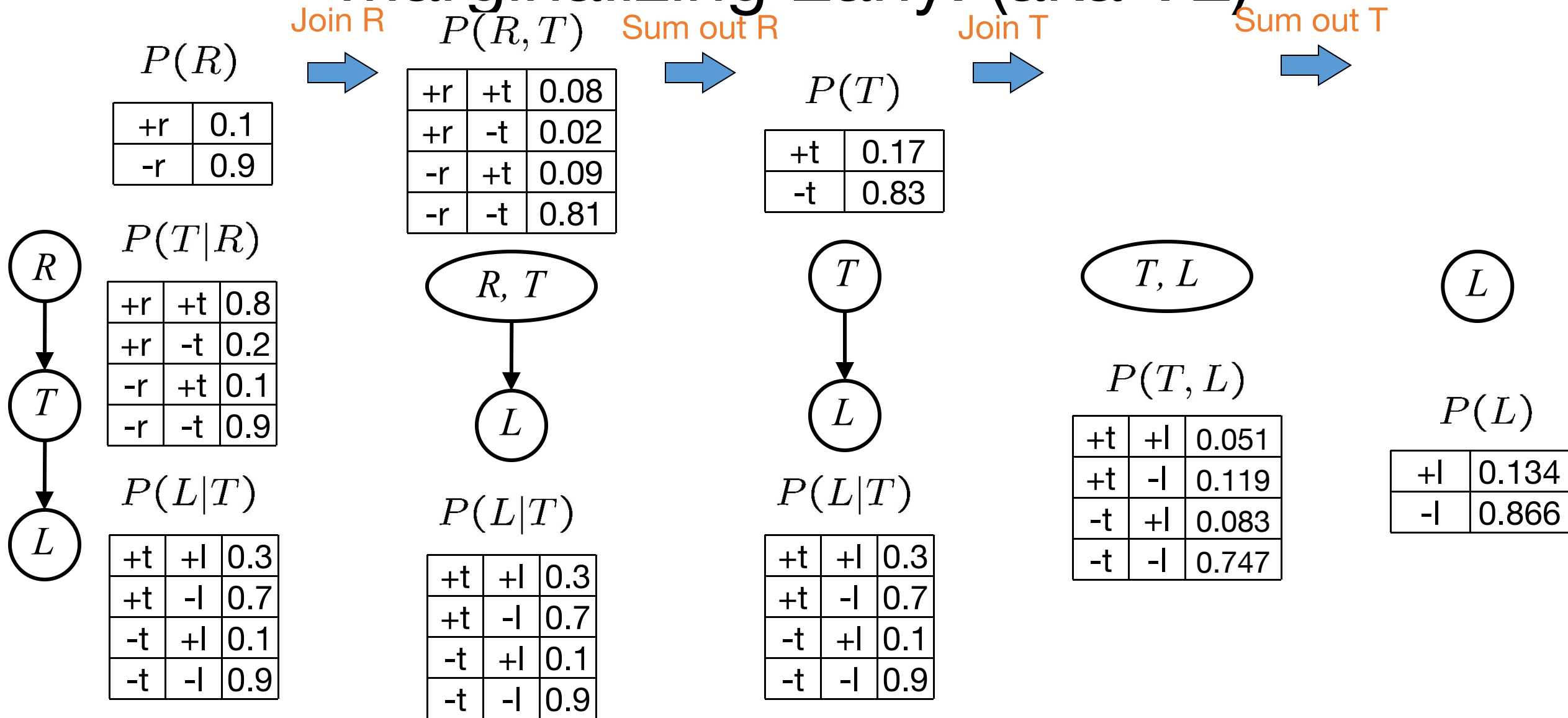
$$= \sum_t P(L|t) \underbrace{\sum_r P(r)P(t|r)}_{\text{Join on } r}$$

Eliminate r

Join on t

Eliminate t

Marginalizing Early! (aka VE)



Evidence

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

$P(R)$		$P(T R)$			$P(L T)$		
+r	0.1	+r	+t	0.8	+t	+l	0.3
-r	0.9	+r	-t	0.2	+t	-l	0.7
		-r	+t	0.1	-t	+l	0.1
		-r	-t	0.9	-t	-l	0.9

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

$P(+r)$		$P(T +r)$			$P(L T)$		
+r	0.1	+r	+t	0.8	+t	+l	0.3
		+r	-t	0.2	+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 would end up with:

$$P(+r, L)$$

+r	+l	0.026
+r	-l	0.074

Normalize



$$P(L \mid +r)$$

+l	0.26
-l	0.74

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

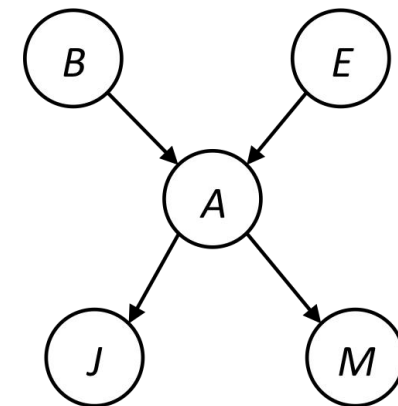
General Variable Elimination

- 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

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

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

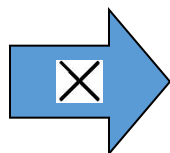


Choose A

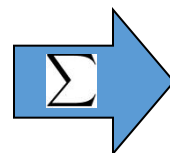
$$P(A|B, E)$$

$$P(j|A)$$

$$P(m|A)$$



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



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

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

Example

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

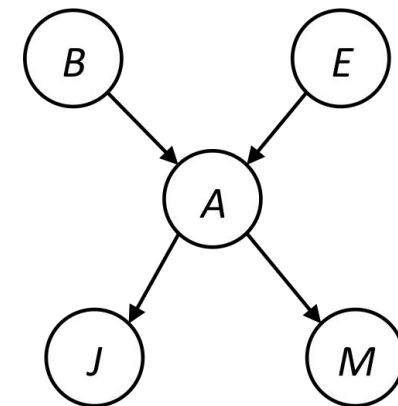
Choose E

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

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

Finish with B

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

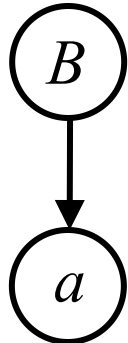


Example 2: $P(B|a)$

Start / Select

$P(B)$

B	P
+b	0.1
¬b	0.9

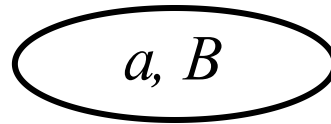


A directed graph with two nodes: a circle labeled B at the top and a circle labeled a at the bottom. A vertical arrow points from B down to a .

$P(A|B) \rightarrow P(a|B)$

B	A	P
+b	+a	0.8
b	¬a	0.2
¬b	+a	0.1
¬b	¬a	0.9

Join on B



$P(a, B)$

A	B	P
+a	+b	0.08
+a	¬b	0.09

Normalize

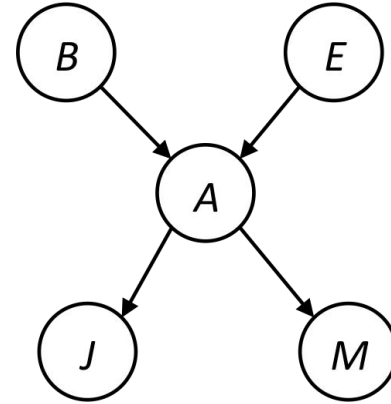
$P(B|a)$

A	B	P
+a	+b	8/17
+a	¬b	9/17

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

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

joining on e, and then summing out gives f_2

All we are doing is exploiting $uw y + uw z + ux y + ux z + vw y + vw z + vx y + vx z = (u+v)(w+x)(y+z)$ to improve computational efficiency!

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:

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

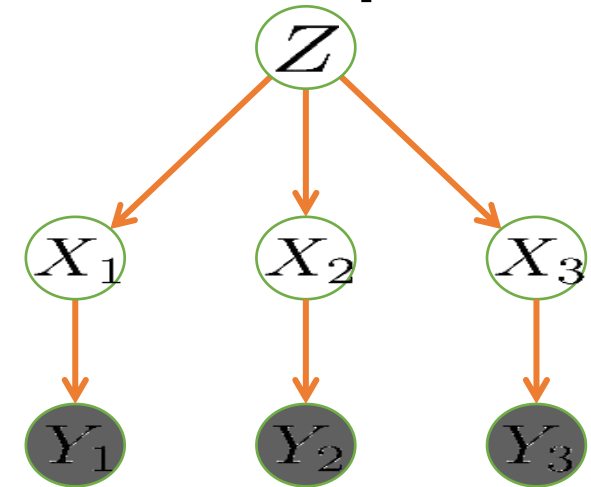
Eliminate Z , this introduces the factor $f_3(y_1, y_2, X_3) = \sum_z p(z)f_1(z, y_1)f_2(z, y_2)p(X_3|z)$, and we are left:

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

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

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

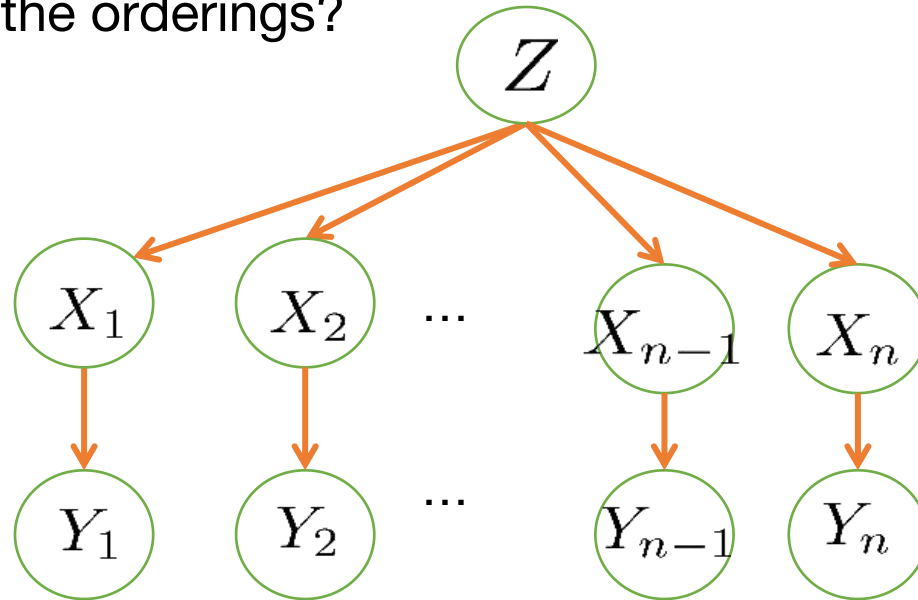
Normalizing over X_3 gives $P(X_3|y_1, y_2, y_3)$.



- Computational complexity critically depends on the **largest factor** being generated in this process.
- **Size of factor = number of entries in table.**
- In example above (assuming binary) all factors generated are of size 2 --- as they all only have one variable (Z , Z , and X_3 respectively).

Variable Elimination Ordering

- For the query $P(X_n | y_1, \dots, y_n)$ work through the following two different orderings as done in previous slide: Z, X_1, \dots, X_{n-1} and X_1, \dots, X_{n-1}, Z . What is the size of the maximum factor generated for each of the orderings?



- Answer: 2^{n+1} versus 2^2 (assuming binary)
- In general: the ordering can greatly affect efficiency.

VE: Computational and Space Complexity

- The **computational** and **space complexity** of **variable elimination** is determined by the **largest factor**
- The elimination ordering can greatly affect the size of the largest factor.
 - E.g., previous slide's example 2^{n+1} vs. 2^2
- Does there always exist an ordering that only results in small factors?
 - **No!**

Worst Case Complexity?

- 3-SAT:

$$(x_1 \vee x_2 \vee \neg x_3) \wedge (\neg x_1 \vee x_3 \vee \neg x_4) \wedge (x_2 \vee \neg x_2 \vee x_4) \wedge (\neg x_3 \vee \neg x_4 \vee \neg x_5) \wedge (x_2 \vee x_5 \vee x_7) \wedge (x_4 \vee x_5 \vee x_6) \wedge (\neg x_5 \vee x_6 \vee \neg x_7) \wedge (\neg x_5 \vee \neg x_6 \vee x_7)$$

$$P(X_i = 0) = P(X_i = 1) = 0.5$$

$$Y_1 = X_1 \vee X_2 \vee \neg X_3$$

$$\dots$$
$$Y_8 = \neg X_5 \vee X_6 \vee X_7$$

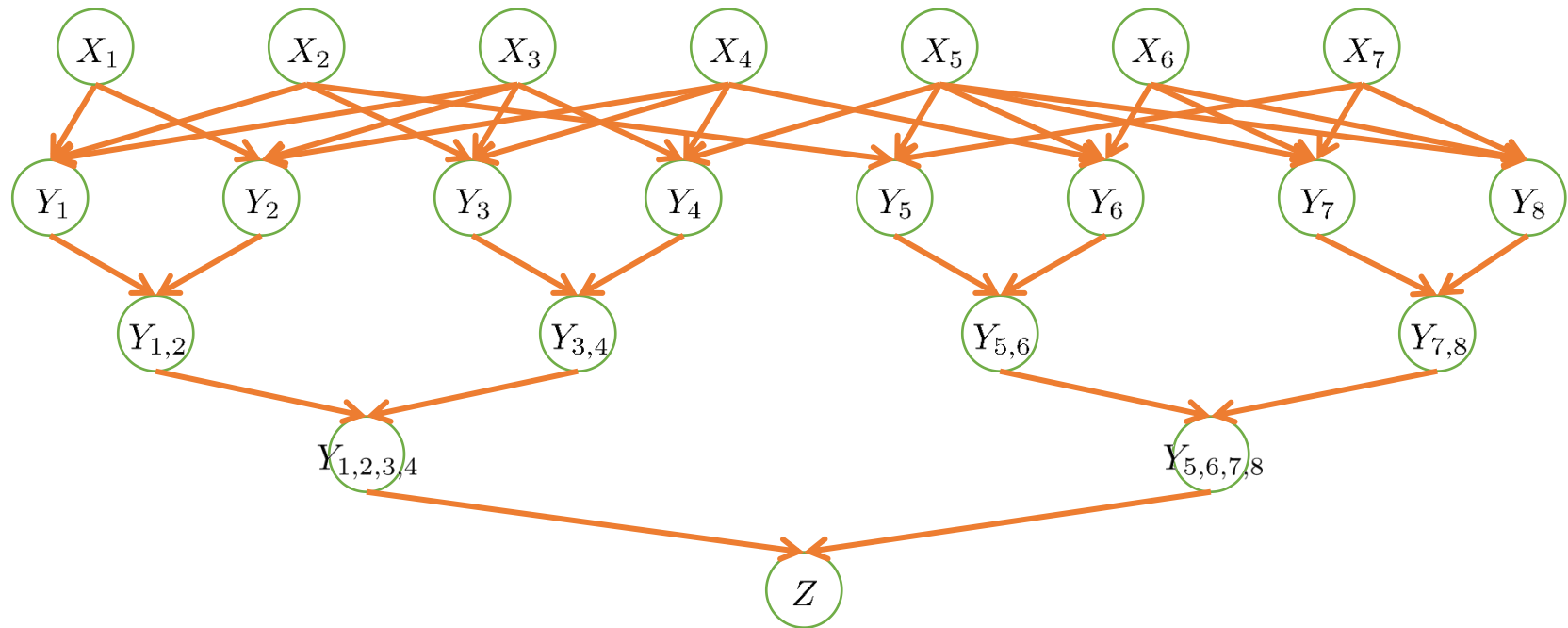
$$Y_{1,2} = Y_1 \wedge Y_2$$

$$\dots$$
$$Y_{7,8} = Y_7 \wedge Y_8$$

$$Y_{1,2,3,4} = Y_{1,2} \wedge Y_{3,4}$$

$$Y_{5,6,7,8} = Y_{5,6} \wedge Y_{7,8}$$

$$Z = Y_{1,2,3,4} \wedge Y_{5,6,7,8}$$

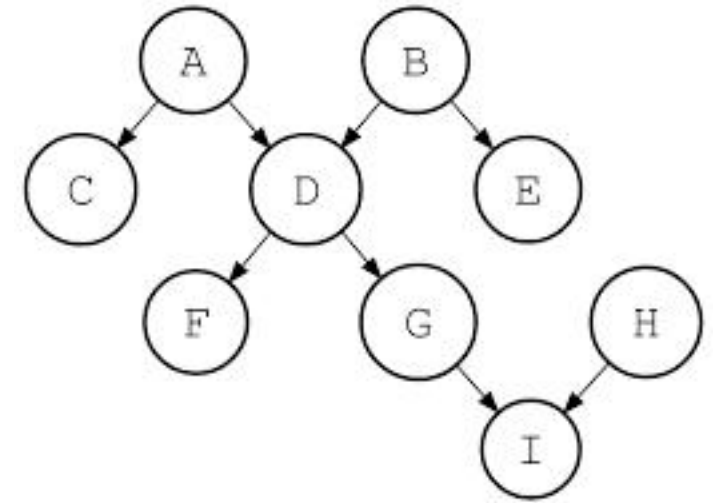


If we can answer $P(X_1, X_2, \dots, X_n | z=1)$, we answered whether the 3-SAT problem has a solution.

Hence **inference in Bayes' nets is NP-hard**. No known efficient probabilistic inference in general

Polytrees

- A **polytree** is a directed graph with **no undirected cycles**
- For poly-trees you can always find an ordering that is efficient
 - Try it!!
- **Cut-set conditioning** for Bayes' net inference
 - Choose set of variables such that if removed only a polytree remains
 - Exercise: Think about how the specifics would work out!



Summary: VE

- Time is **exponential in size** of largest factor
- **Bad elimination order** can generate huge factors
- **NP Hard** to find the best elimination order
- There are reasonable heuristics for picking an elimination order (such as choosing the variable that results in the smallest next factor)
- Inference in polytrees (nets with no cycles) is linear in size of the network (the largest CPT)

Probabilistic Reasoning

- Let p be a formula. We try to prove it from **available information** (evidence) that is known to be true.
- Let q_1, \dots, q_n be this **evidence**. Theorem proving can be used to determine if:

$$(q_1 \wedge q_2 \wedge \dots \wedge q_n) \rightarrow p$$

When the evidence does not imply that p is true it may imply that p is false:

$$(q_1 \wedge q_2 \wedge \dots \wedge q_n) \rightarrow \neg p$$

Probabilistic Reasoning

- But in many practical situations the evidence cannot be used to prove that p is either **true** or **false**. Still, it can always be used for **probabilistic reasoning**.
- The **optimal Bayes decision rule** is:

Take p as true if:

$$P(p \mid q_1 \wedge q_2 \wedge \dots \wedge q_n) \geq P(\neg p \mid q_1 \wedge q_2 \wedge \dots \wedge q_n)$$

This is the same condition as:

Take p as true if:

$$P(p \mid q_1 \wedge q_2 \wedge \dots \wedge q_n) \geq 0.5$$

Probabilistic Reasoning: Generalization of Logical Reasoning

- Observe that if $(q_1 \wedge q_2 \wedge \dots \wedge q_n) \rightarrow p$, then
$$P(p \mid q_1 \wedge q_2 \wedge \dots \wedge q_n) = 1$$

So, **probabilistic reasoning** generalizes **logic reasoning**.

Probabilistic Reasoning

- The **optimal Bayes decision rule** is:

Take p as true if:

$$P(p \mid q_1 \wedge q_2 \wedge \dots \wedge q_n) \geq P(\neg p \mid q_1 \wedge q_2 \wedge \dots \wedge q_n)$$

- **Simplification using Bayes Rule**

- Take p as true if:

$$P(q_1 \wedge q_2 \wedge \dots \wedge q_n \mid p) P(p) \geq P(q_1 \wedge q_2 \wedge \dots \wedge q_n \mid \neg p) P(\neg p)$$

Naïve Bayes assumption

- **Simplification using Naïve Bayes assumption**

- The presence of a particular evidence is unrelated to the presence of any other evidences.
- Take p as true if:

$$P(q_1|p) P(q_2|p) \dots P(q_n|p)P(p) \geq P(q_1|\neg p) P(q_2|\neg p) \dots P(q_3|\neg p) P(\neg p)$$

Naïve Bayes assumption

- Classifying documents by their content, for example into **spam** and **non-spam** e-mails
 - The **document D** is drawn from a number of classes (topics)
 - Each class C is modeled as sets of words, the probability that the i-th word occurs is $P(w_i | C)$
 - The probability of the document D given the classes C is

$$P(D|C) = \prod_i p(w_i|C)$$

Question: what is the probability that a given document D belongs to a given class C ? $P(C|D)$

Naïve Bayes assumption

- $P(D|C) = P(D \wedge C) / P(C)$ conditional distribution
- $P(C|D) = P(D \wedge C) / P(D)$ conditional distribution
 $= P(C)P(D|C) / P(D)$ product rule

Assume that there two classes: **Spam** (**S**) and **not Spam** ($\neg S$)

$$P(D|\mathbf{S}) = \prod_i p(w_i|\mathbf{S}) \quad \text{Assumption}$$

$$P(D|\neg \mathbf{S}) = \prod_i p(w_i|\neg \mathbf{S}) \quad \text{Assumption}$$

$$P(\mathbf{S}|D) = P(\mathbf{S}) \prod_i p(w_i|\mathbf{S}) / P(D) \quad \text{Bayes rule}$$

$$P(\neg \mathbf{S}|D) = P(\neg \mathbf{S}) \prod_i p(w_i|\neg \mathbf{S}) / P(D) \quad \text{Bayes rule}$$

Assume $P(\mathbf{S}) = P(\neg \mathbf{S}) = 0.5$

$$P(\mathbf{S}|D) / P(\neg \mathbf{S}|D) = \prod_i p(w_i|\mathbf{S}) / p(w_i|\neg \mathbf{S})$$