

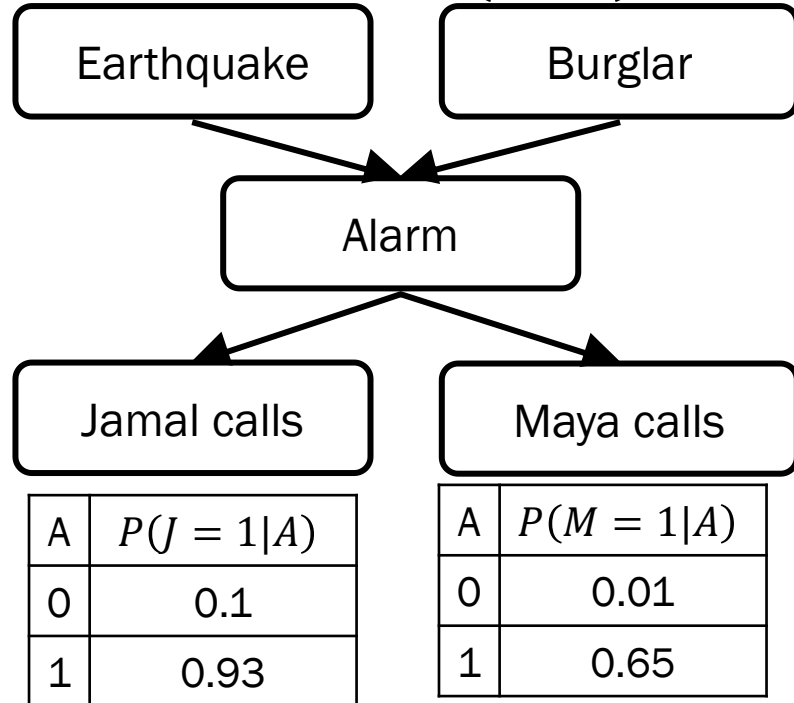


FUDAN SUMMER SCHOOL INTRODUCTION TO AI

Probabilistic Reasoning and Decision Making
Day 7 – variable elimination exercise, MLE learning

Inference in Bayes Nets: Variable Elimination

$$P(E = 1) = 0.001 \quad P(B = 1) = 0.005$$



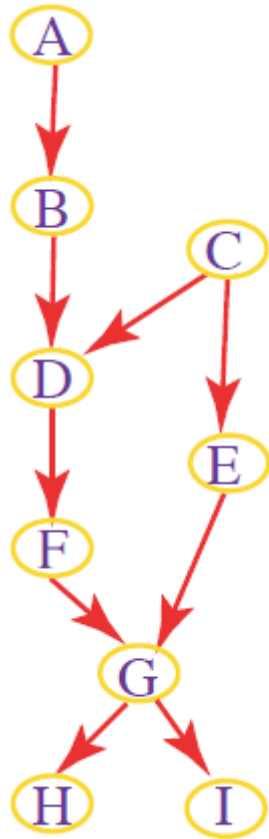
B	E	$P(A = 1 B, E)$
0	0	0.002
0	1	0.35
1	0	0.96
1	1	0.98

Compute $P(B|J = 1, M = 1)$

$$= \frac{\sum_e \sum_a P(B, J = 1, M = 1, A = a, E = e)}{\sum_b P(B = b, J = 1, M = 1)}$$

1. Create a factor for each CPT
2. Choose an elimination ordering for non-query variables (we'll use M, J, A, E)
3. Eliminate each variable in order by applying factor operations.
 1. Product
 2. Summing out
 3. Condition

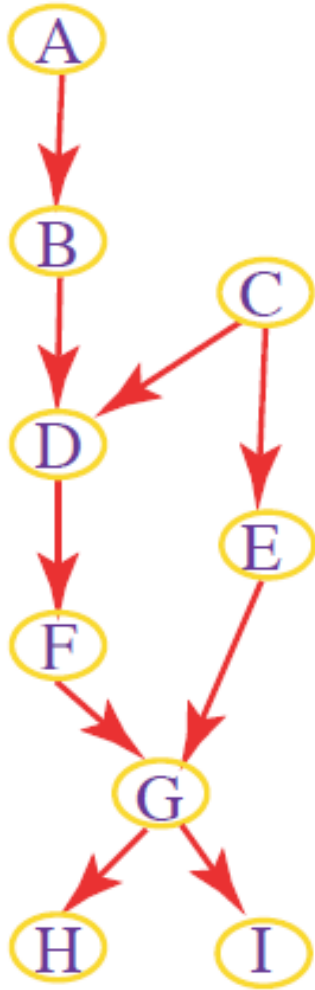
Calculate $P(G \mid H = 1)$



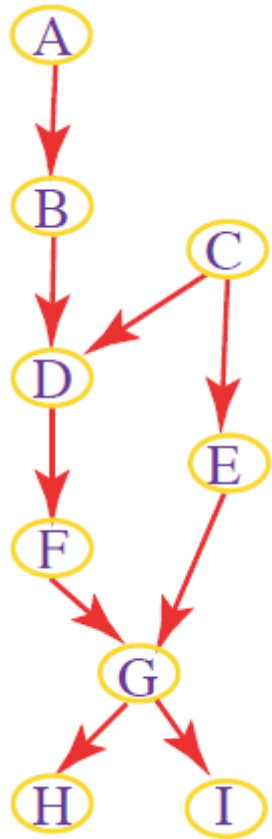
If I simply use joint probability to calculate the result, what is the number of ops that we need to do? Assume each variable is binary and we have all the CPTs

- A. 2^{14} multiplications and 2^8 additions
- B. 2^{15} multiplications and 2^7 additions
- C. 2^{11} multiplications and 2^9 additions
- D. 2^{12} multiplications and 2^{10} additions
- E. None of the above

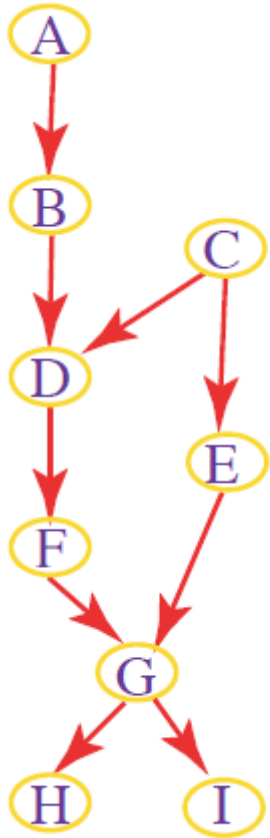
Calculate $P(G \mid H = 1)$



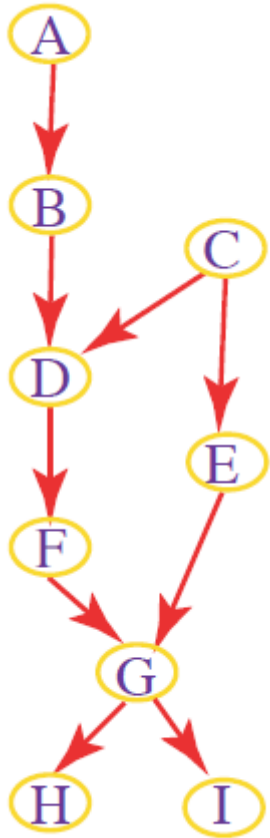
Calculate $P(G \mid H = 1)$



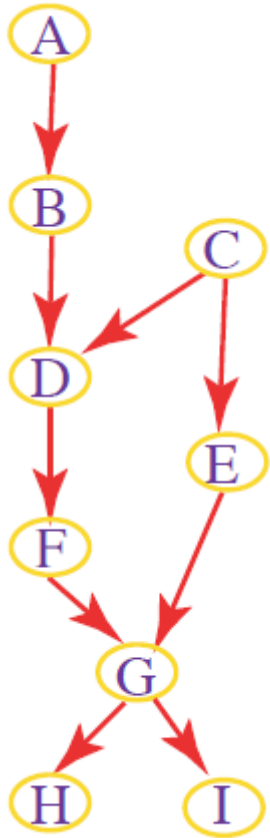
Calculate $P(G \mid H = 1)$



Calculate $P(G \mid H = 1)$

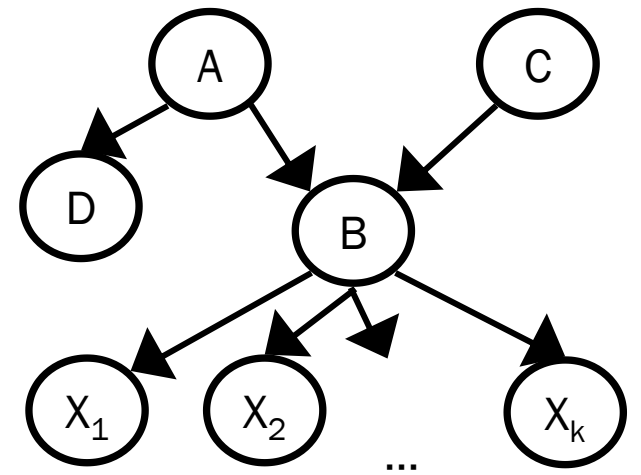


Calculate $P(G \mid H = 1)$



Does elimination order matter?

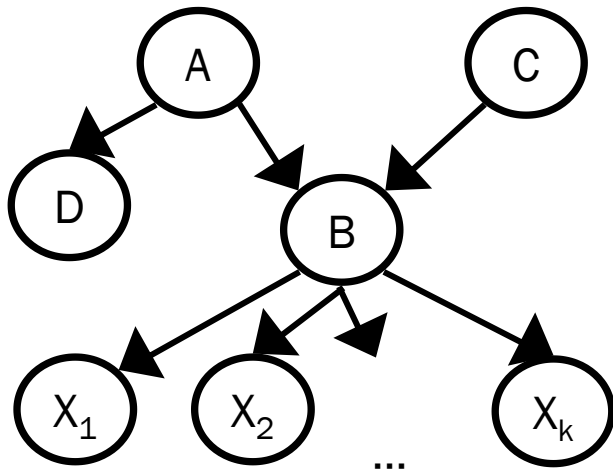
- In general, yes (but not in the trivial graphs we've been considering)
- Time and space of VE is dominated by the largest factor created
- Heuristic: Eliminate the variable that will lead to the smallest next factor being created
 - *In a polytree this leads to linear time inference (in size of largest CPT)*



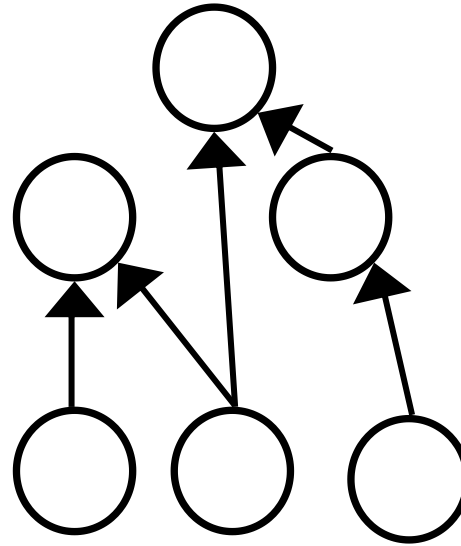
What is a polytree? A graph with no *undirected* loops

■ Which are polytrees?

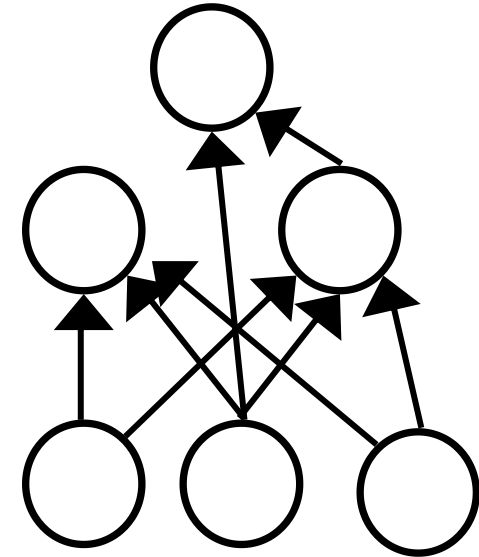
A. None of these B. I only C. I and II D. I, II and III



I.

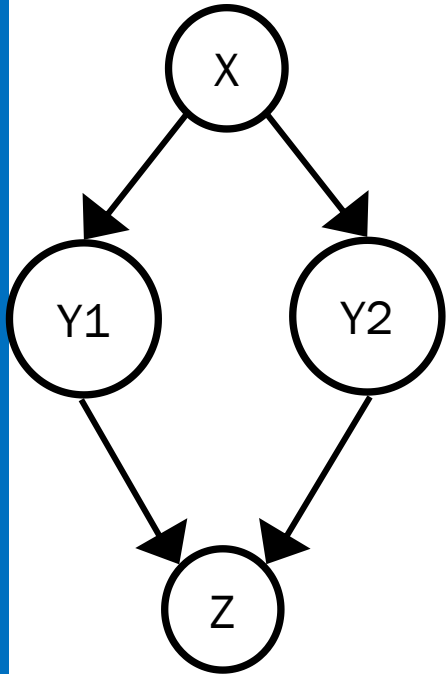


II.



III.

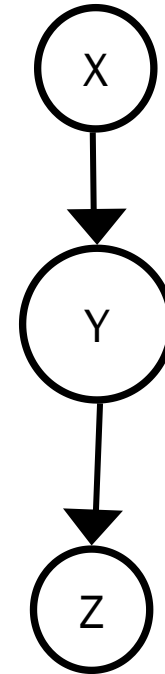
Making a polytree by collapsing nodes



X	$P(Y1=1 X)$
0	0.1
1	0.5

X	$P(Y2=1 X)$
0	0.9
1	0.7

Y1	Y2	$P(Z=1 Y1, Y2)$
0	0	0.2
0	1	0.3
1	0	0.6
1	1	0.8



Y1	Y2	Y	X	$P(Y X)$	$P(Z=1 Y)$
0	0	0	0		
0	1	1	0		
1	0	2	0		
1	1	3	0		
0	0	0	1		
0	1	1	1		
1	0	2	1		
1	1	3	1		
1	1	3	1		

What we've learned so far

- Pause and spend 5 minutes writing down the main ideas we've learned so far, and how they connect.

What we've learned so far

- Basics of probability and how to use the product rule, Bayes rule and marginalization to do inference.
- How to represent relationships in the world with Bayes nets:
 - *Graph to represent variables and their direct dependencies*
 - *CPTs to represent strength of dependencies*
- Noisy-OR model for representing CPTs
- Reasoning about conditional independence between variables in a Bayes' net using d-separation
- General algorithms for inference in Bayes nets:
 - *Enumeration (exponential in number of undefined variables)*
 - *Variable Elimination (linear in # of undefined variables and size of largest CPT for polytrees)*
- How to turn a graph into a polytree (at the cost of the size of the CPT)

Up next: Learning in Bayes nets (focused on learning CPTs)

Learning Bayes Nets

- Aspects of the Bayes Net might not be known or easy to elicit from experts. This could include:
 - *the structure of the DAG*
 - *the CPTs*
- In this course we will assume a given DAG (structure) and focus on learning CPTs from data:
 - *With complete data (all variables observed)*
 - *With incomplete data (some variables not observed) – LATER*

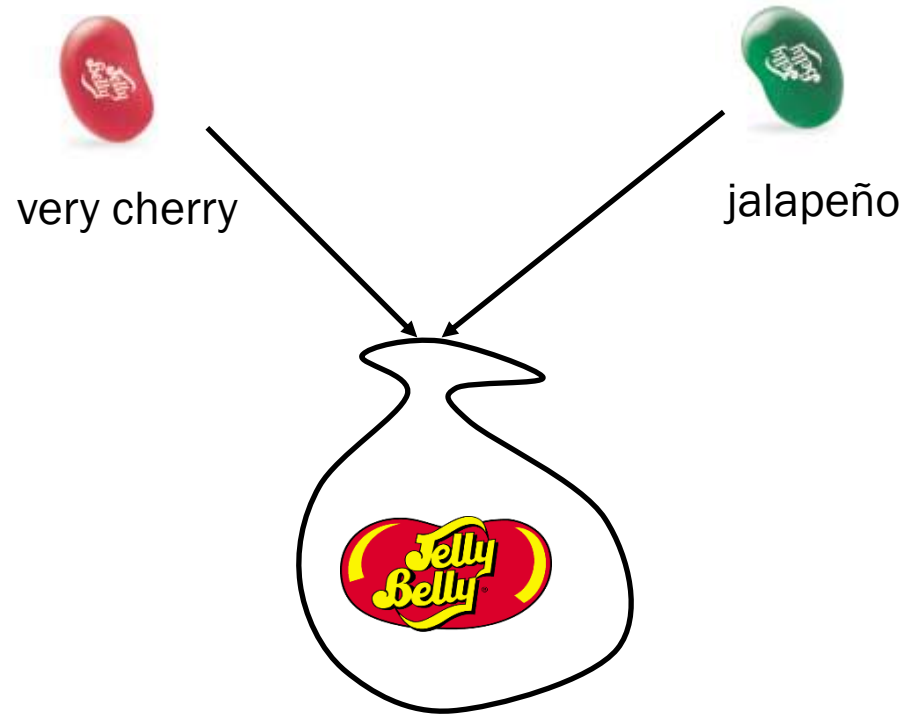
Learning from data via Maximum Likelihood (ML) Learning

- ML is the simplest form of learning in BNs
- Idea: Choose (learn) the model (CPTs) that maximizes the probability of the data

$$P_{model}(observed\ data)$$

Simple Example: Estimating the proportion of candies in a bag

Goal: Estimate proportion of cherries/jalapenos by drawing samples



Simple Example: Estimating the proportion of candies in a bag

Goal: Estimate proportion of cherries/jalapenos by drawing samples



Simple Example: Estimating the proportion of candies in a bag

Goal: Estimate $p = P(X = \text{cherry})$ by drawing samples



Data: T samples $\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}$ where each $x^{(t)}$ is either cherry or jalapeño

The probability of selecting these samples given the parameter $p = P(X = \text{cherry})$ is:

$$P(\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\} | p)$$

The **likelihood** of p for this data

Assumption: The samples are independently drawn and identically distributed (i.i.d.)

In other words, each sample is drawn using the same p and don't depend on each other:

$$P(\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\} | p) = P(X = x^{(1)} | p) P(X = x^{(2)} | p) \dots P(X = x^{(T)} | p) = \prod_{t=1}^T P(X = x^{(t)} | p)$$

Simple Example: Estimating the proportion of candies in a bag

Goal: Estimate $p = P(X = \text{cherry})$ by drawing samples



Data: T i.i.d. samples $\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}$ where each $x^{(t)}$ is either cherry or jalapeño

$$\text{likelihood}(p) = P(\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}) = \prod_{t=1}^T P(X = x^{(t)})$$

For a given p , how many possible values can $P(X = x^{(t)})$ take?

- A. 1
- B. 2
- C. 4
- D. There is no way to know

Simple Example: Estimating the proportion of candies in a bag

Goal: Estimate $p = P(X = \text{cherry})$ by drawing samples



Data: T i.i.d. samples $\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}$ where each $x^{(t)}$ is either cherry or jalapeño

$$\text{likelihood}(p) = P(\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}) = \prod_{t=1}^T P(X = x^{(t)})$$

How many terms in $\prod_{t=1}^T P(X = x^{(t)})$ will equal p ?

- A. The number of cherry candies in the sample
- B. The number of jalapeno candies in the sample
- C. The total number of candies in the sample
- D. You cannot tell from the data given

Simple Example: Estimating the proportion of candies in a bag

Goal: Estimate $p = P(X = \text{cherry})$ by drawing samples



Data: T i.i.d. samples $\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}$ where each $x^{(t)}$ is either cherry or jalapeño

$$\text{likelihood}(p) = P(\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}) = \prod_{t=1}^T P(X = x^{(t)}) = p^{N_c}(1 - p)^{N_j}$$

We want to choose p that maximizes this function

Log-Likelihood: An easier function

Goal: Estimate $p = P(X = \text{cherry})$ by drawing samples



Data: T i.i.d. samples $\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}$ where each $x^{(t)}$ is either cherry or jalapeño

log-likelihood: $\mathcal{L}(p) = \log P(\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}) = \log \prod_{t=1}^T P(X = x^{(t)}) =$

For this problem $\mathcal{L}(p) =$

Maximize the log-likelihood by taking the derivative

Goal: Estimate $p = P(X = \text{cherry})$ by drawing samples



Data: T i.i.d. samples $\{x^{(1)}, x^{(2)}, \dots, x^{(T)}\}$ where each $x^{(t)}$ is either cherry or jalapeño

$$\text{log-likelihood: } \mathcal{L}(p) = N_c \log p + N_j \log(1 - p)$$

Estimating Parameters of a Bayes Net

Goal: Estimate $p = P(X = \text{cherry})$ by drawing samples



$$P(X = \text{cherry}) = p = \frac{N_c}{N_c + N_j}$$
$$P(X = \text{jalapeno}) = 1 - p$$

But what about a more complex Bayes Net?

Estimating Parameters of a Bayes Net

Given: A fixed DAG with n discrete nodes $\{X_1, X_2, \dots, X_n\}$

Goal: Estimate the values in the CPTs to maximize probability of observed data

Estimating Parameters of a Bayes Net

Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

Data set: $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)})\}_{t=1}^T$

$$\mathcal{L} = \log P(data) = \log \prod_{t=1}^T P(X_1 = x_1^{(t)}, X_2 = x_2^{(t)}, \dots, X_n = x_n^{(t)})$$

Under what assumption(s) is the above log-likelihood equation true?

A. Data is i.i.d.

B. $P(X_i | X_1, X_2, \dots, X_{i-1}) = P(X_i | Parents(X_i))$

C. $T > 1$ (you have more than one data point)

D. More than one of the above

E. None of the above (it is generally true)