

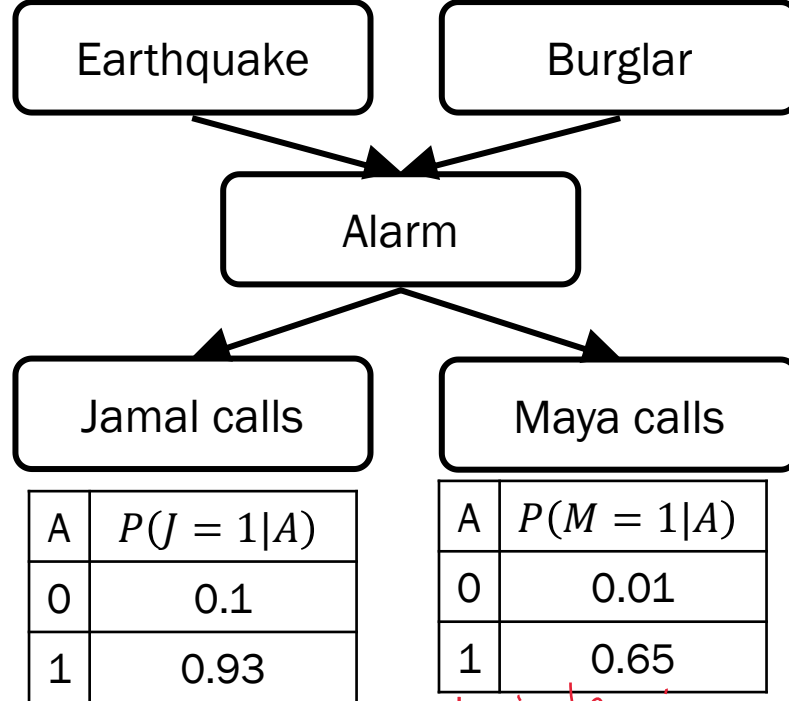


# FUDAN SUMMER SCHOOL INTRODUCTION TO AI

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

Factor: { multiple summing out condition

Handwritten notes and equations:

- $E, A$ : non query, non-evidence
- Marginal  $P(B, J=1, M=1)$
- $\sum_e \sum_a P(B, J=1, M=1, A=a, E=e)$

2<sup>7</sup> terms, each term has 2<sup>3</sup> \*

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

$$P(G=0 \mid H=1) \leftarrow P(G=0, H=1)$$

$$P(G=1 \mid H=1) \leftarrow P(G=1, H=1)$$

$$P(G=1 \mid H=1) = \frac{P(G=1, H=1)}{P(G=0, H=1) + P(G=1, H=1)}$$

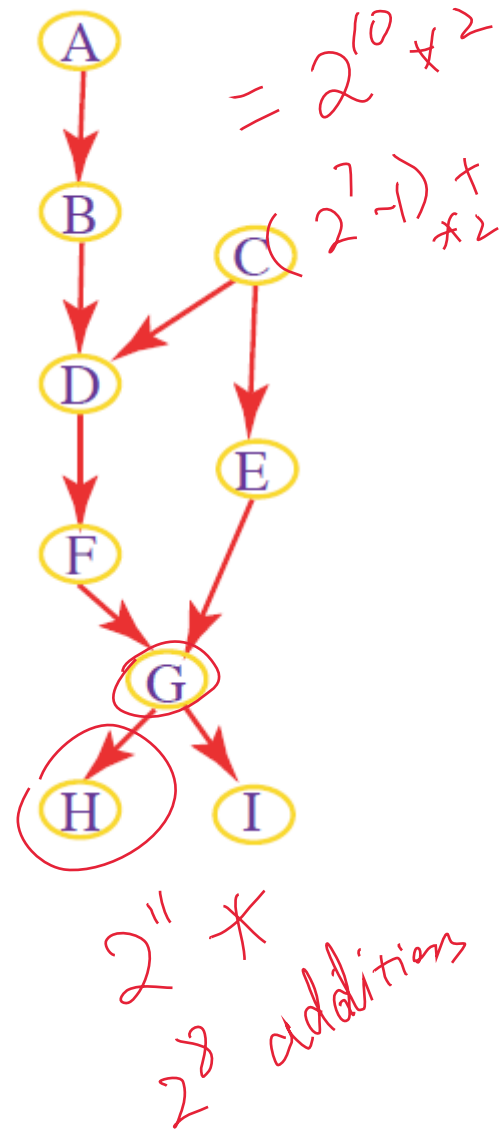
$$P(G=0, H=1) = \sum_a \sum_b \sum_c \sum_d \sum_e \sum_f \sum_i P(G=0, H=1, A=a, B=b, C=c, D=d, E=e, F=f, I=i)$$

$$= \sum_a \sum_b \dots \sum_i P(A=a) P(B=b \mid A=a) P(C=c) P(D=d \mid B=b, C=c) P(E=e \mid C=c)$$

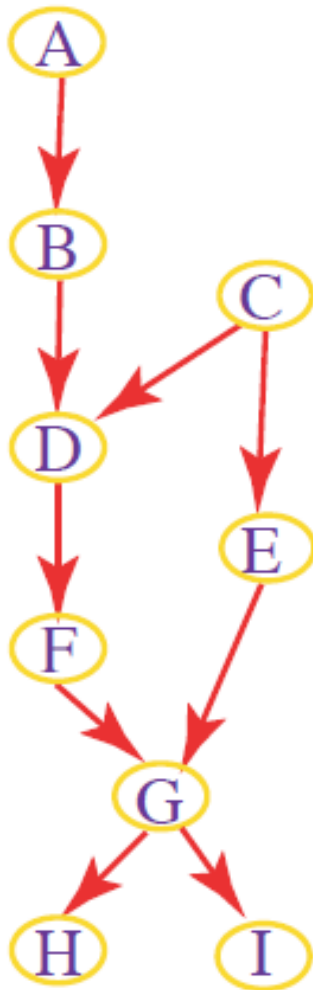
If I simply use joint probability to calculate the result, what is the number of ops that we need to do? Assume each variable is binary and we have all the CPTs

- A. 2<sup>14</sup> multiplications and 2<sup>8</sup> additions
- B. 2<sup>15</sup> multiplications and 2<sup>7</sup> additions
- C. 2<sup>11</sup> multiplications and 2<sup>9</sup> additions
- D. 2<sup>12</sup> multiplications and 2<sup>10</sup> additions
- E. None of the above

$$P(I=i \mid D=d) \cdot P(G=0 \mid I=i, \bar{C}=e) \cdot P(H=1 \mid G=0) \cdot P(I=i \mid G=0)$$



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



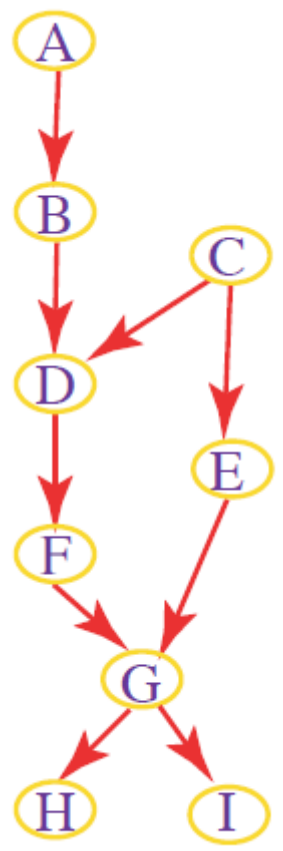
$$P(G, H=1)$$

$$= \sum_{a,b,c,d,e,f,i} P(A=a, B=b, C=c, D=d, E=e, G, H=1, I=i)$$

$$= \sum_{a,b,c,d,e,f,i} P(A=a) P(B=b|A=a) P(C=c) P(D=d|B=b, C=c) \\ P(E=e|C=c) P(F=f|D=d) P(G|F=f, E=e) \\ P(H=1|G) P(I=i|G)$$

elimination order:

$I, H, A, B, C, D, E, F$



$$= \sum_{a,b,c,d,e,f,i} P(A=a) P(B=b|A=a) P(C=c) P(D=d|B=b, C=c) \\ P(E=e|C=c) P(F=f|D=d) P(G|F=f, E=e) \\ P(H=h|G) P(I=i|G)$$

elimination order:

I, H, A, B, C, D, E, F

$$F(G) \begin{array}{c|c} G & F(G) \\ \hline 0 & \sim \\ 1 & \sim \end{array}$$

2+

$$\sum_b \left( P(D=d|B=b, C=c) \sum_a \left( P(A=a) P(B=b|A=a) \right) \right) \cdot \underbrace{P(H=h|G)}_{\text{condition}} \cdot \underbrace{\sum_i P(I=i|G)}_{F(G) = 2*}$$

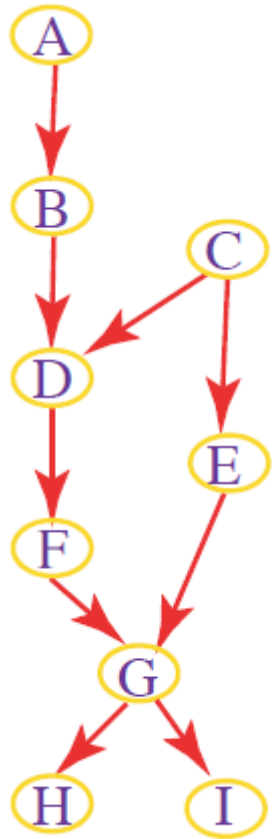
$$F(I, G)$$

I	G	F(I, G)
0	0	
0	1	
1	0	
1	1	

$$F(B, CD) \xrightarrow{4+} F(A, B) \xrightarrow{2+} F(B) \xrightarrow{8*} F(B, CD)$$

$$F(H=1, G)$$

H	G	F(H=1, G)
1	0	~
1	1	~



$$= \sum_{a,b,c,d,e,f,i} \cancel{P(A=a)} \cancel{P(B=b|A=a)} \cancel{P(C=c)} \cancel{P(D=d|B=b, C=c)} \\ P(C=c|C=c) P(F=f|D=d) P(G|F=f, E=e)$$

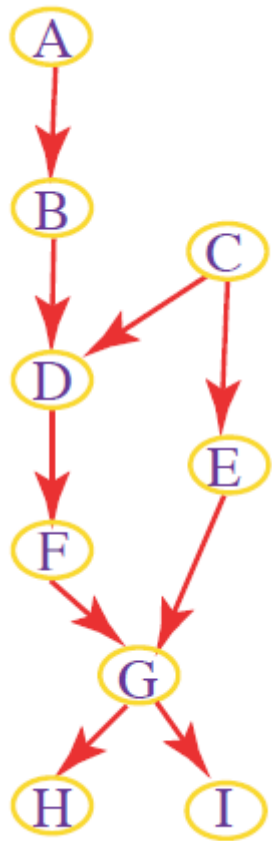
$$\cancel{P(H=h|G)} \cancel{P(I=i|G)}$$

elimination order:

I, H, A, B, C, D, E, F

$$\sum_{c,d,e,f} F(C,D) \cdot F(G) P(C=c|C=c) \\ P(F=f|D=d) \\ P(G|F=f, E=e)$$

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



$$\sum_{cdef} F(C,D) \cdot F(G) P(\bar{E}=e \mid C=c) P(\bar{F}=f \mid D=d) P(G \mid \bar{F}=f, \bar{E}=e)$$

$$\sum_{cdef} F(C,D) F(G) F(\bar{C},E) F(D,F) F(E,F,G)$$

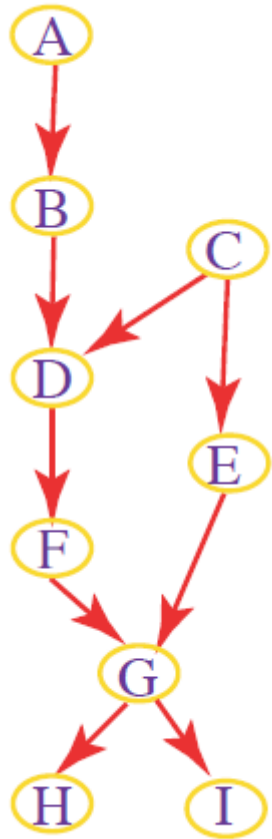
$$= \sum_g \sum_e \sum_d \sum_c \sum_f F(G) F(E,F,G) F(\bar{C},E) F(C,D) F(\bar{C},E)$$

$\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$   
 $F(G)$   $F(E,F,G)$   $F(\bar{C},E)$   $F(C,D)$   $F(\bar{C},E)$   
 $2^*$   $2^+$   $4^+$   $8^*$   $4^+$   $8^*$   $4^+$   $8^*$   
 $F(G)$   $F(E,F,G)$   $F(E,F)$   $F(D,E,F)$   $F(C,D,E)$   $F(C,D)$   $F(\bar{C},E)$

$$40^* \leftarrow 2048^*$$

$$20^+ \leftarrow 256^+$$

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

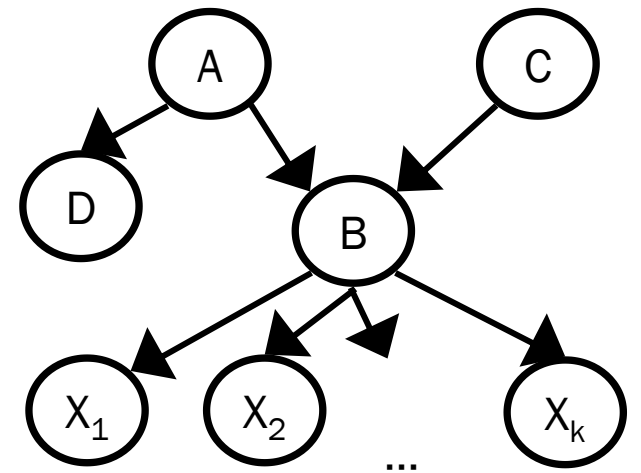




# Does elimination order matter?

- 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)*

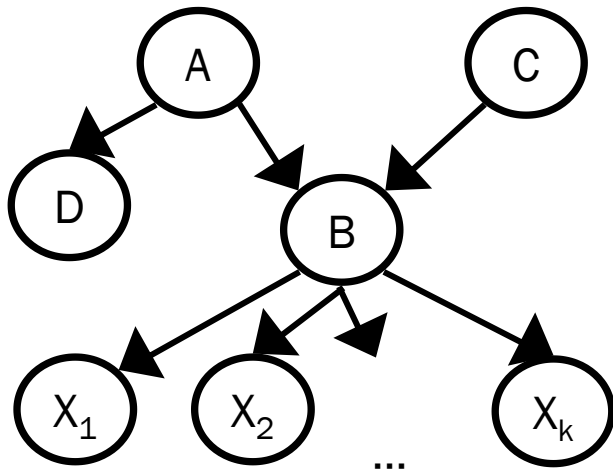
*PAG w/o undirected loops  
no two undirected path between  
any two nodes*



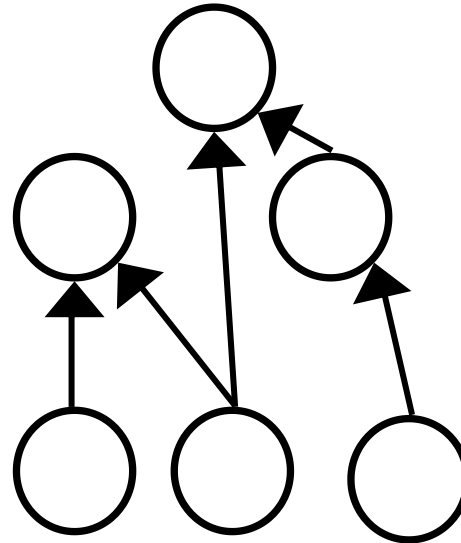
# What is a polytree? A graph with no *undirected* loops

■ Which are polytrees?

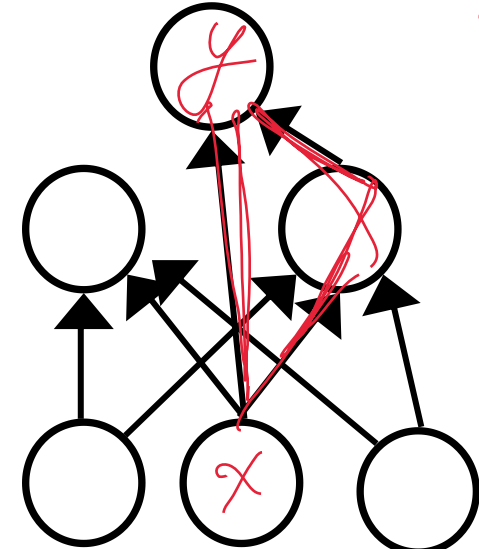
A. None of these B. I only C. I and II D. I, II and III



I.



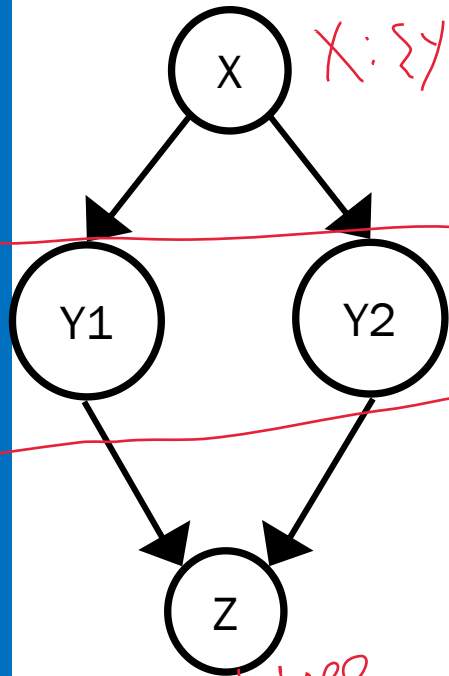
II.



III.

no a polytree

# Making a polytree by collapsing nodes



$$X: \{Y_1\}, E: \{X\}, Y: \{Y_2\}$$

$$Y_1 \leftarrow X \rightarrow Y_2 \text{ condition 2}$$

$$d\text{-separation}$$

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

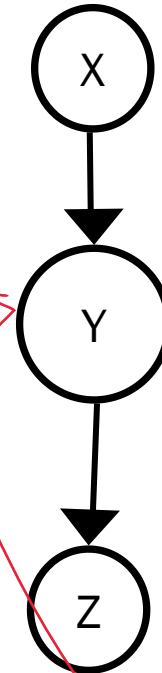
X	P(Y2=1 X)
0	0.9
1	0.7

non-polytree

condition 3

$$Y_1 \rightarrow Z \leftarrow Y_2$$

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



$$P(Y_1=0, Y_2=0|X) = P(Y_1=0|X)$$

Y1	Y2	Y	X	P(Y X)	P(Z=1 Y)
0	0	0	0	0.1	0.2
0	1	1	0	0.9	0.3
1	0	2	0	0.1	0.6
1	1	3	0	0.9	0.8
0	0	0	1	0.5	
0	1	1	1	0.5	
1	0	2	1	0.5	
1	1	3	1	0.5	

$$P(Z=1|Y_1=0, Y_2=0)$$

# 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)

learning  
① Learn parameters (CPTs)  
② learn the structure  
(NP-hard)

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

# Learning Bayes Nets : data

- 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)* : *MLE*
  - *With incomplete data (some variables not observed)* – *LATER*
    - *EM*

# 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

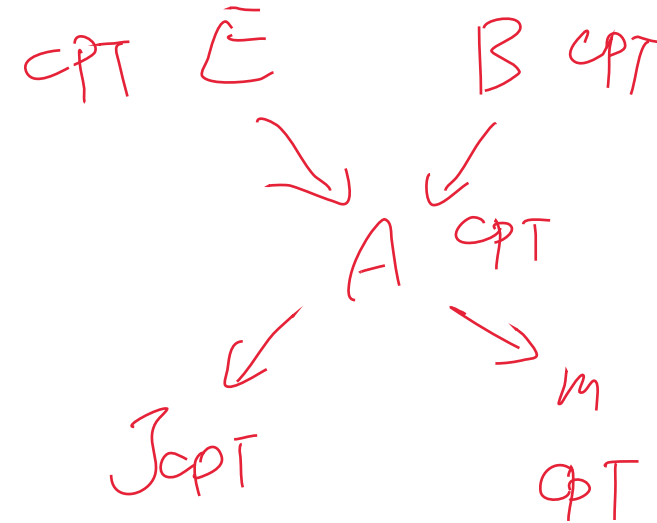
data: E B A J M

sample 1	→	0	0	1	1	1
sample 2	→	1	0	1	1	1
sample 3	→	0	0	1	1	1

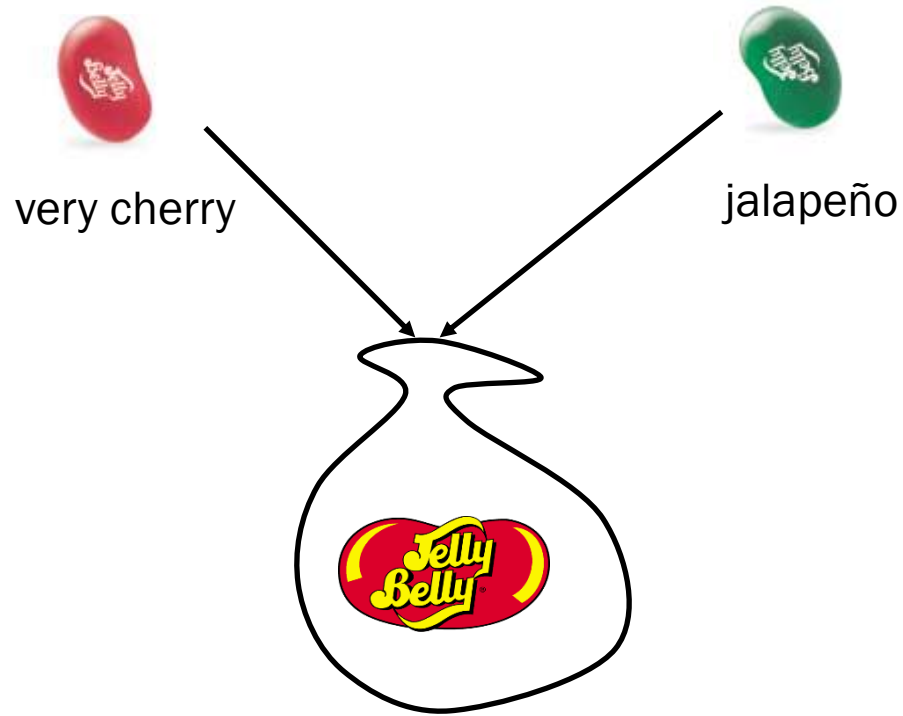
Given the observed data and the BN, ML is to find the values for CPTs that is the most likely to have the data.

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

$\arg \max_{\text{CPTs}} P(\text{data} | \text{CPTs})$

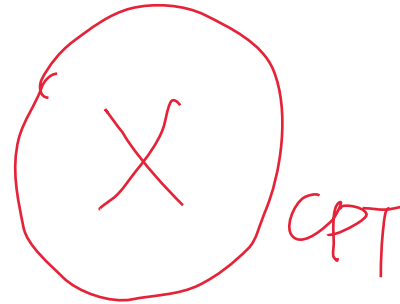


# Simple Example: Estimating the proportion of candies in a bag



Goal: Estimate proportion of cherries/jalapenos by drawing samples

data : drawing	type
1	cherry
2	jalapeño
3	cherry
4	cherry
⋮	⋮
⋮	⋮
⋮	⋮



$$\underline{P(X = \text{cherry}) ??}$$



# Simple Example: Estimating the proportion of candies in a bag

Goal: Estimate proportion of cherries/jalapenos by drawing samples



(X)

$X: \{\text{cherry}, \text{jalapeno}\}$

CPT:  $P(X = \text{cherry}) = p$

$P(X = \text{jalapeno}) = 1 - p$

data:  $T$  samples  $\{X^{(1)}, X^{(2)}, X^{(3)}, \dots, X^{(T)}\}$

$$\begin{aligned} & \underset{p}{\operatorname{argmax}} P(\text{data} | p) \\ &= \underset{p}{\operatorname{argmax}} P(X^{(1)}, X^{(2)}, \dots, X^{(T)} | p) \\ &= \underset{p}{\operatorname{argmax}} \prod_{i=1}^T P(X^{(i)} | p) \end{aligned}$$

$T$  samples  
Where  $X^{(i)}$  is either  
cherry or jalapeno

# 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

cherry  $p$   
jalapeño  $1-p$

# 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

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

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

- ☒ A. The number of cherry candies in the sample (*known to us*)
- ☐ 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



$$N_c + N_j = T$$

$N_c$ : # of cherries in sample

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

$N_j$ : # of Jalapenos in sample

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

We want to choose  $p$  that maximizes this function

$$\begin{aligned} \arg \max_p p^{N_c} \cdot (1-p)^{N_j} &= \arg \max_p \log(p^{N_c} \cdot (1-p)^{N_j}) \\ &= \arg \max_p (N_c \log p + N_j \log(1-p)) \end{aligned}$$

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

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

$$\frac{\partial \mathcal{L}(p)}{\partial p} = N_c \cdot \frac{1}{p} + N_j \cdot \frac{1}{1-p} \cdot (-1)$$

$$\stackrel{\Delta}{=} 0$$

$$\frac{N_c}{p} = \frac{N_j}{1-p}$$

$$(\underbrace{N_c + N_j}_T) p = N_c$$

$$p = \frac{N_c}{T}$$

# 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

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

BIN

CPT	
flavor	$P(\text{flavor})$
cherry	$\frac{N_c}{T}$

Marginal Probability

$$= \frac{\text{Count of } X = x_i}{\text{total \# of samples}}$$

# 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?

: all data are observed

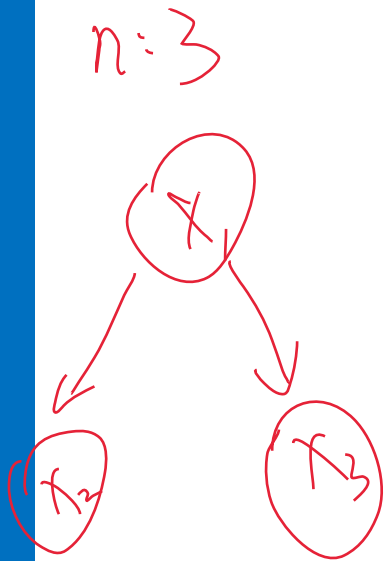


# Estimating Parameters of a Bayes Net

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

5 probabilities



data	$x_1$	$x_2$	$x_3$
①	0	0	0
②	0	0	0
③	0	0	1
④	0	0	0
⋮			
T	1	0	1

CPT

$x_1$	$P(x_1)$
0	??
1	??

$x_1$	$P(x_2=1 x_1)$
0	??
1	??

$x_1$	$P(x_3=1 x_1)$
0	??
1	??

# Estimating Parameters of a Bayes Net

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

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

where each sample is  $\{x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)}\}$   
 $t^{th}$  sample

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

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)