

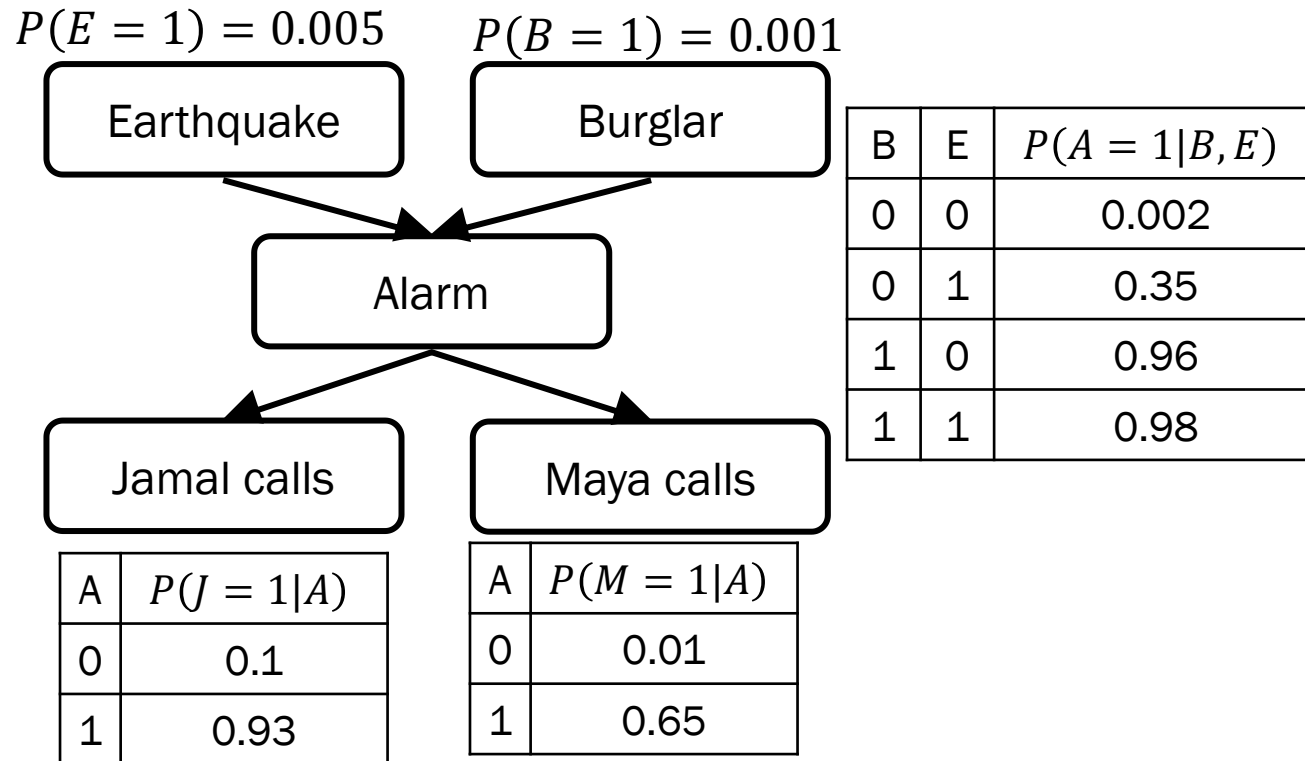


FUDAN SUMMER SCHOOL INTRODUCTION TO AI

Probabilistic Reasoning and Decision Making

Day 6 – variable elimination, MLE learning

Inference in Bayes Nets: Enumeration



$$\begin{aligned}
 &P(B = 1, J = 1, M = 1) \\
 &= P(B = 1)P(E = 1)P(A = 1|E = 1, B = 1)P(J = 1|A = 1)P(M = 1|A = 1) \\
 &+ P(B = 1)P(E = 0)P(A = 1|E = 0, B = 1)P(J = 1|A = 1)P(M = 1|A = 1) \\
 &+ P(B = 1)P(E = 1)P(A = 0|E = 1, B = 1)P(J = 1|A = 0)P(M = 1|A = 0) \\
 &+ P(B = 1)P(E = 0)P(A = 0|E = 0, B = 1)P(J = 1|A = 0)P(M = 1|A = 0)
 \end{aligned}$$

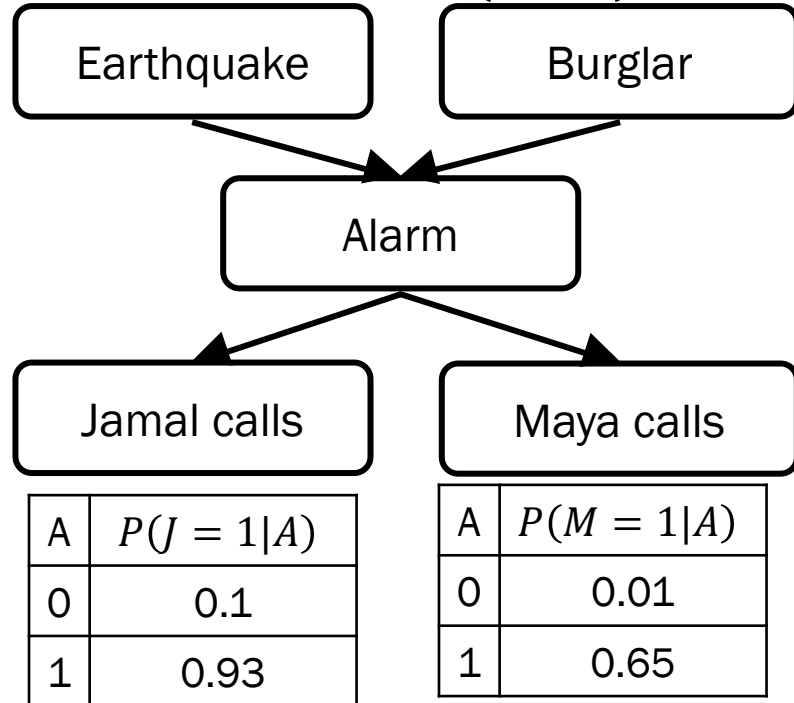
Intuition

- Solve $ab + ac + ad + aeh + afh + agh$
- Can we simplify better?
- $ab + ac + ad + aeh + afh + agh$

Inference in Bayes Nets: Variable Elimination

$$P(E = 1) = 0.005$$

$$P(B = 1) = 0.001$$



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

Eliminate redundant calculations by storing intermediate results in "factors"

A factor is:

Operations on Factors: Multiplication

 $f_0(E)$

E	val
0	0.995
1	0.005

 $f_1(B)$

B	val
0	0.999
1	0.001

 $f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

 $f_3 * f_4$ $f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

 $f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

Operations on Factors: Multiplication

$f_0(E)$

E	val
0	0.995
1	0.005

$f_1(B)$

B	val
0	0.999
1	0.001

$f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

How many values (rows) are in the table for the factor $f_2 * f_3$

- A. 4
- B. 8
- C. 16
- D. 32

$f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

$f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

Operations on Factors: Multiplication

$f_0(E)$

E	val
0	0.995
1	0.005

$f_1(B)$

B	val
0	0.999
1	0.001

$f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

How many total multiplications are needed to produce the table $f_2 * f_3$

- A. 0
- B. 16
- C. 32
- D. 64

$f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

$f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

Operations on Factors: Summing out

 $f_0(E)$

E	val
0	0.995
1	0.005

 $f_1(B)$

B	val
0	0.999
1	0.001

 $f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

$$\sum_e f_2$$

 $f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

 $f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

Operations on Factors: Conditioning

 $f_0(E)$

E	val
0	0.995
1	0.005

 $f_1(B)$

B	val
0	0.999
1	0.001

 $f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

 $f_2(A, B = 1, E)$ $f_3(J, A)$

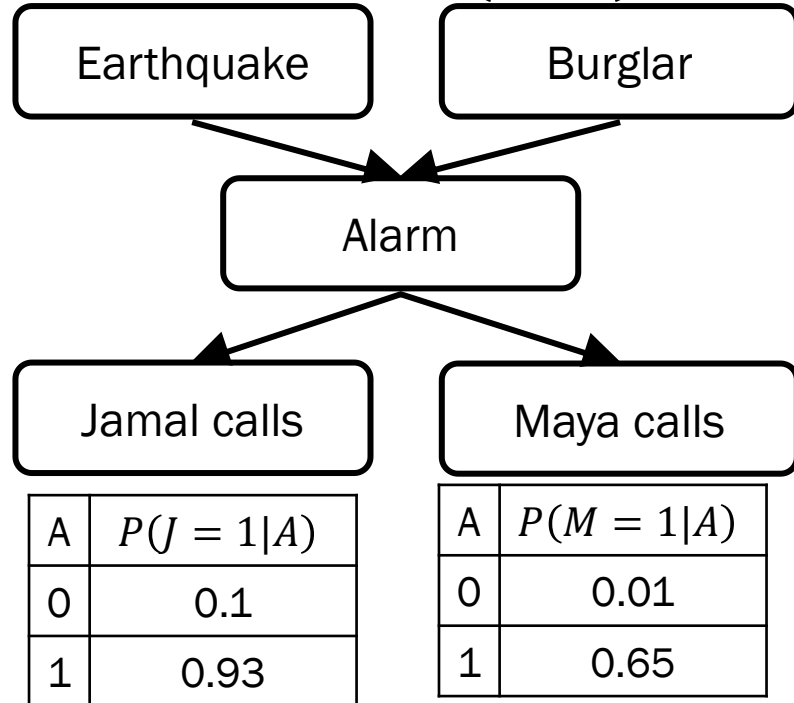
J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

 $f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

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.

Variable Elimination: $\sum_e \sum_a P(B, J = 1, M = 1, A = a, E = e)$

Order: M, J, A, E

$f_0(E)$

E	val
0	0.995
1	0.005

$f_1(B)$

B	val
0	0.999
1	0.001

$f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

$f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

$f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

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.

Which factor(s) do we need to consider to eliminate M?

- A. f_4
- B. f_3 and f_4
- C. f_2 , f_3 and f_4
- D. All of the factors

Variable Elimination: $\sum_e \sum_a P(B, A, E, M = 1, J = 1)$

Order: M, J, A, E

$f_0(E)$

E	val
0	0.995
1	0.005

$f_1(B)$

B	val
0	0.999
1	0.001

$f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

$f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

$f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

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.

$$\sum_e \sum_a P(B, A, E, M = 1, J = 1)$$

$$= \sum_e \sum_a f_0(E) f_1(B) f_2(A, B, E) f'_3(A) f'_4(A)$$

$f_0(E)$

E	val
0	0.995
1	0.005

 $f_1(B)$

B	val
0	0.999
1	0.001

 $f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

 $f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

 $f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

 $f'_3(A)$

A	val
0	0.1
1	0.93

 $f'_4(A)$

A	val
0	0.01
1	0.65

$$\sum_e \sum_a P(B, A, E, M = 1, J = 1)$$

$$= \sum_e \sum_a f_0(E) f_1(B) f_2(A, B, E) f'_3(A) f'_4(A)$$

$f_0(E)$

E	val
0	0.995
1	0.005

 $f_1(B)$

B	val
0	0.999
1	0.001

 $f_2(A, B, E)$

A	B	E	val
0	0	0	0.998
0	0	1	0.65
0	1	0	0.04
0	1	1	0.02
1	0	0	0.002
1	0	1	0.35
1	1	0	0.96
1	1	1	0.98

 $f_3(J, A)$

J	A	val
0	0	0.9
0	1	0.07
1	0	0.1
1	1	0.93

 $f_4(M, A)$

M	A	val
0	0	0.99
0	1	0.35
1	0	0.01
1	1	0.65

 $f'_3(A)$

A	val
0	0.1
1	0.93

 $f'_4(A)$

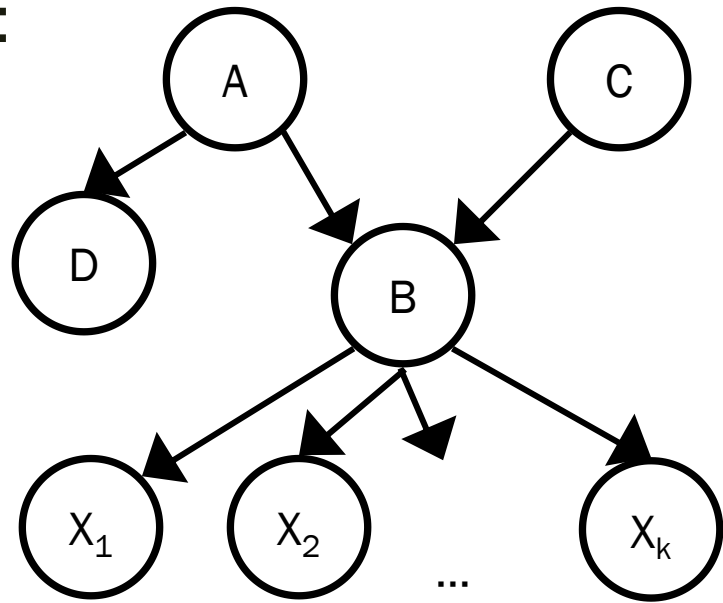
A	val
0	0.01
1	0.65

$$\sum_e \sum_a P(B, A, E, M = 1, J = 1)$$

$$= \sum_e \sum_a f_0(E) f_1(B) f_2(A, B, E) f'_3(A) f'_4(A)$$

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
- Consider:



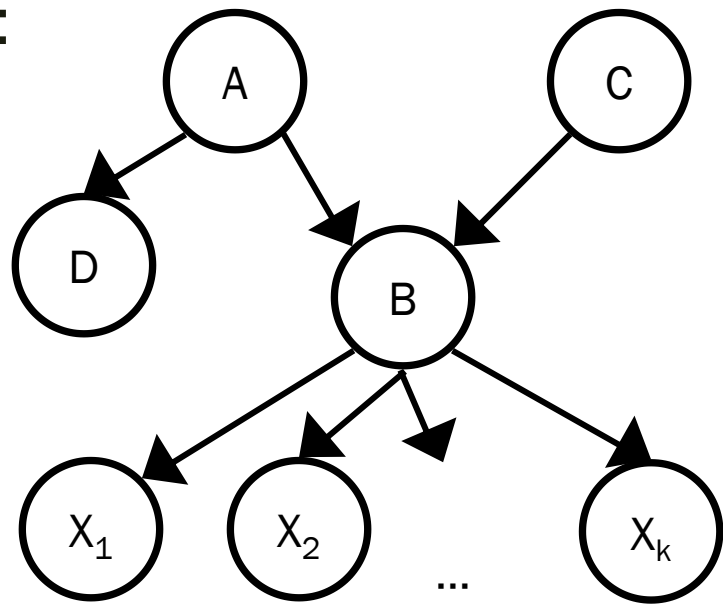
To calculate $P(C, X_1 = 1)$

How large (how many rows) is the factor that is created if you eliminate B first (after accounting for X_1)?

- A. 2
- B. 4
- C. 8
- D. 2^k
- E. 2^{k+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
- Consider:



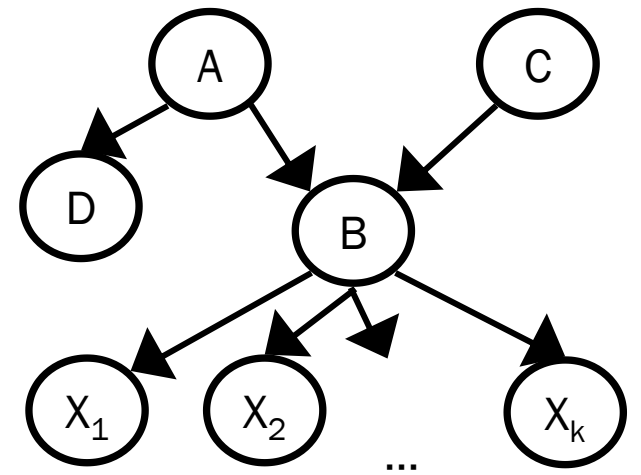
To calculate $P(C, X_1 = 1)$

Which one should I eliminate first?

- A. one of X's
- B. A
- C. C
- D. D
- E. More than one is correct

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)*



Consider the belief network with 5 variables: Earthquake, Burglary, Alarm, JamalCalls, MayaCalls.

Suppose we want to compute the probabilities $P(B = 1|J = 1, M = 0)$ and $P(B = 0|J = 1, M = 0)$.

a) Express the conditional probabilities $P(B = 1|J = 1, M = 0)$ and $P(B = 0|J = 1, M = 0)$ in terms of the probabilities $P(B = 1, J = 1, M = 0)$ and $P(B = 0, J = 1, M = 0)$.

Consider the belief network with 5 variables: Earthquake, Burglary, Alarm, JamalCalls, MayaCalls. Suppose we want to compute the probabilities $P(B = 1|J = 1, M = 0)$ and $P(B = 0|J = 1, M = 0)$.

b) Letting b represent either 0 or 1, express $P(B = b, J = 1, M = 0)$ in terms of the CPTs of the belief network. (Do not get rid of redundant operations. This corresponds to the brute-force enumeration method.)

Consider the belief network with 5 variables: Earthquake, Burglary, Alarm, JamalCalls, MayaCalls. Suppose we want to compute the probabilities $P(B = 1|J = 1, M = 0)$ and $P(B = 0|J = 1, M = 0)$.

c) Count the number of additions and multiplications needed to compute the answer to part (b). (Include operations for both $b = 0$ and $b = 1$.)

Consider the belief network with 5 variables: Earthquake, Burglary, Alarm, JamalCalls, MayaCalls. Suppose we want to compute the probabilities $P(B = 1|J = 1, M = 0)$ and $P(B = 0|J = 1, M = 0)$.

d) Rewrite the answer in (b) by pulling as many constant factors as possible out of sums. (Hint: to get the most efficient answer, you should have a nested sum: there will be a sum over values of E inside a sum over values of A.)

Consider the belief network with 5 variables: Earthquake, Burglary, Alarm, JamalCalls, MayaCalls. Suppose we want to compute the probabilities $P(B = 1|J = 1, M = 0)$ and $P(B = 0|J = 1, M = 0)$.

e) Count the number of additions and multiplications needed to compute the answer to part (d). (Include operations for both $b = 0$ and $b = 1$.)

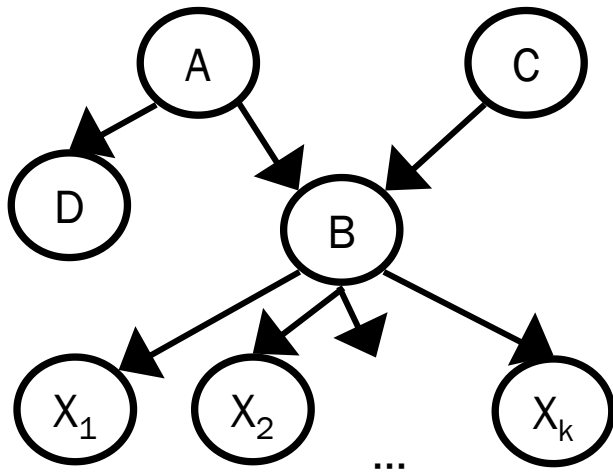
Consider the belief network with 5 variables: Earthquake, Burglary, Alarm, JamalCalls, MayaCalls. Suppose we want to compute the probabilities $P(B = 1|J = 1, M = 0)$ and $P(B = 0|J = 1, M = 0)$.

f) Using the elimination ordering J,M,E,A, write out the sequence of factors that would be introduced when applying the variable elimination algorithm to this computation. This answer should correspond closely to your answer in part (d).

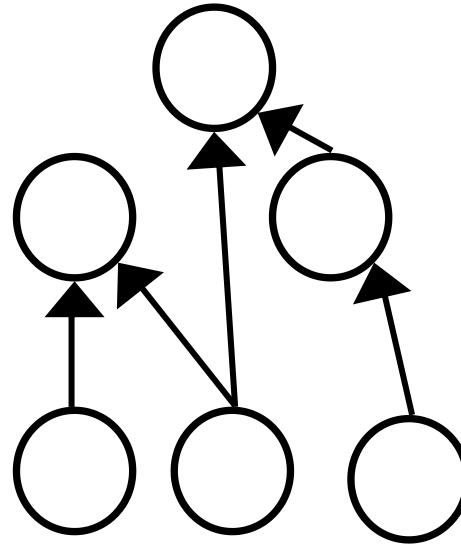
What is a polytree? A graph with no *undirected* loops

■ Which are polytrees?

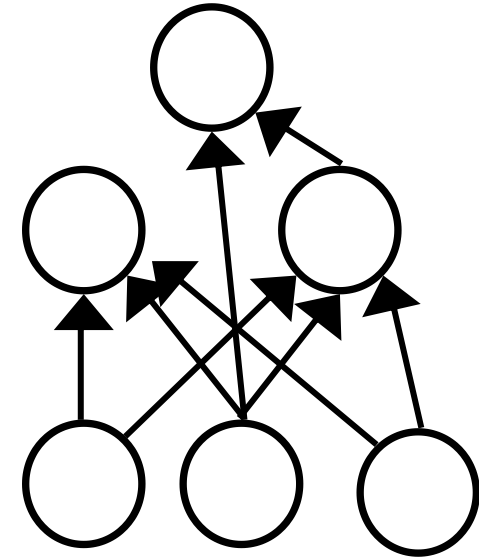
A. None of these B. I only C. I and II D. I, II and III



I.

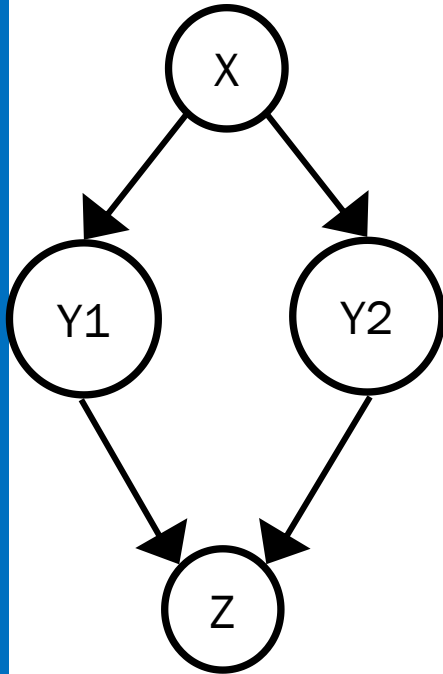


II.



III.

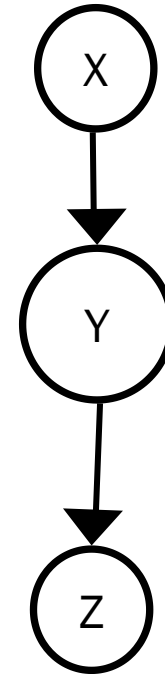
Making a polytree by collapsing nodes



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		

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*

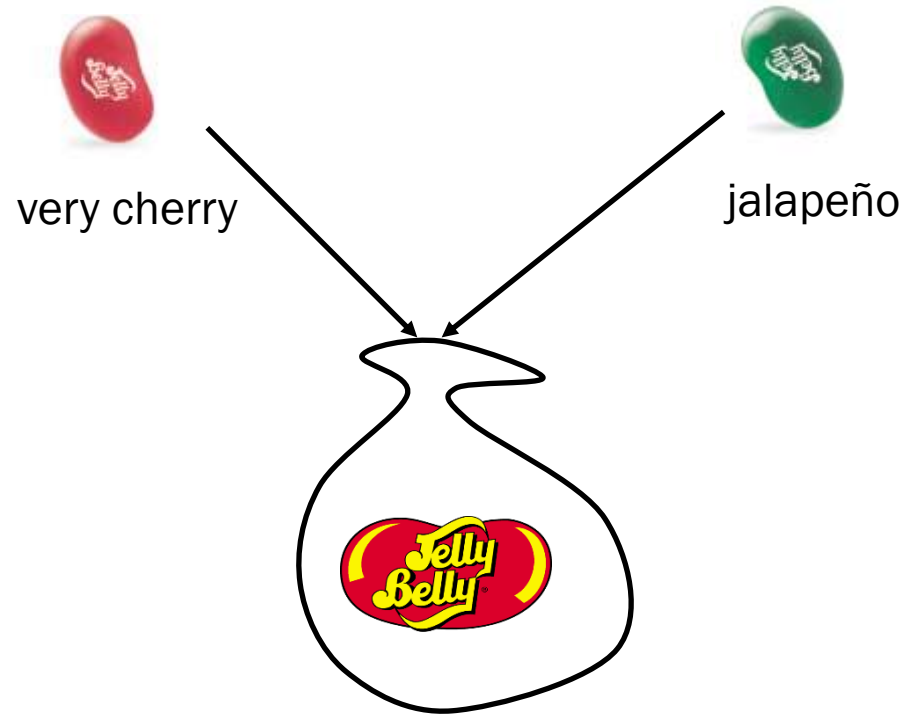
Learning from data via Maximum Likelihood (ML) Learning

- 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)