

Probabilistic Graphical Models

Lecture 2 – Bayesian Networks Representation

CS/CNS/EE 155
Andreas Krause

Announcements

- Will meet in Steele 102 for now
- Still looking for another 1-2 TAs..
- Homework 1 will be out soon. Start early!! 😊

Multivariate distributions

- Instead of random variable, have random vector

$$\mathbf{X}(\omega) = [X_1(\omega), \dots, X_n(\omega)]$$

- Specify $P(X_1=x_1, \dots, X_n=x_n)$
- Suppose all X_i are Bernoulli variables.
- How many parameters do we need to specify?

Marginal distributions

- Suppose we have joint distribution $P(X_1, \dots, X_n)$
- Then

$$P(\underline{X_i = x_i}) = \sum_{\underline{x_1, \dots, x_{i-1}}, \underline{x_{i+1}, \dots, x_n}} P(\underline{x_1, \dots, x_n})$$

- If all X_i binary: How many terms?

Rules for random variables

- Chain rule

$$P(x_1 \dots x_n) = P(x_1) P(x_2 | x_1) \dots P(x_n | x_1 \dots x_{n-1})$$

- Bayes' rule

$$P(x|y) = \frac{P(y|x) P(x)}{P(y)}$$

↑
How do we get $P(y)$?

Key concept: Conditional independence

- Events α, β conditionally independent given γ if

$$P(\alpha \wedge \beta \mid \gamma) = P(\alpha \mid \gamma) P(\beta \mid \gamma)$$

- Random variables X and Y cond. indep. given Z if
for all $x \in \text{Val}(X)$, $y \in \text{Val}(Y)$, $z \in \text{Val}(Z)$

$$P(\underline{X = x}, \underline{Y = y} \mid \underline{Z = z}) = \underline{P(X = x \mid Z = z)} \underline{P(Y = y \mid Z = z)}$$

- If $P(Y=y \mid Z=z) > 0$, that's equivalent to

$$P(X = x \mid Z = z, Y = y) = P(X = x \mid Z = z)$$

Similarly for sets of random variables **X, Y, Z**

We write: $P \models \underline{\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z}}$

Why is conditional independence useful?

- $P(X_1, \dots, X_n) = P(X_1) P(X_2 | X_1) \dots P(X_n | X_1, \dots, X_{n-1})$

How many parameters?

$$2^0 + 2^1 + 2^2 + \dots + 2^{n-1} = 2^n - 1$$

- Now suppose $X_1 \dots X_{i-1} \perp X_{i+1} \dots X_n | X_i$ for all i

Then

$$P(X_1, \dots, X_n) = \underbrace{P(X_1)}_1 \cdot \underbrace{P(X_2 | X_1)}_2 \cdot \underbrace{P(X_3 | X_2)}_2 \cdot \dots \cdot \underbrace{P(X_n | X_{n-1})}_2$$

How many parameters?

$$2^{n-1} \ll 2^n$$

Exponential reduction in # params

- Can we compute $P(X_n)$ more efficiently? *Yes (often)*

Properties of Conditional Independence

- **Symmetry**

- $X \perp Y \mid Z \Rightarrow Y \perp X \mid Z$

- **Decomposition**

- $X \perp Y, W \mid Z \Rightarrow X \perp Y \mid Z$

- **Contraction**

"Inverse Decomposition"

- $(X \perp Y \mid Z) \wedge (X \perp W \mid Y, Z) \Rightarrow X \perp Y, W \mid Z$

- **Weak union**

- $X \perp Y, W \mid Z \Rightarrow X \perp Y \mid Z, W$

- **Intersection**

- $(X \perp Y \mid Z, W) \wedge (X \perp W \mid Y, Z) \Rightarrow X \perp Y, W \mid Z$

- Holds only if distribution is positive, i.e., $P > 0$

Key questions

- How do we specify distributions that satisfy particular independence properties?

→ Representation

- How can we exploit independence properties for efficient computation?

→ Inference

- How can we identify independence properties present in data?

→ Learning

Will now see example: Bayesian Networks

Key idea

- Conditional parameterization
(instead of joint parameterization)
- For each RV, specify $P(X_i | \mathbf{X}_A)$ for set \mathbf{X}_A of RVs
- Then use chain rule to get joint parametrization

$$P(x_1 \dots x_n) = \prod P(x_i | x_{A_i})$$

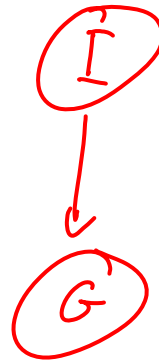
- Have to be careful to guarantee legal distribution...

$$P(X|Y), P(Y|X)$$

Does there exist $P(X,Y)$ with above ~~dist~~ cond. distributions

Not in general

Example: 2 variables

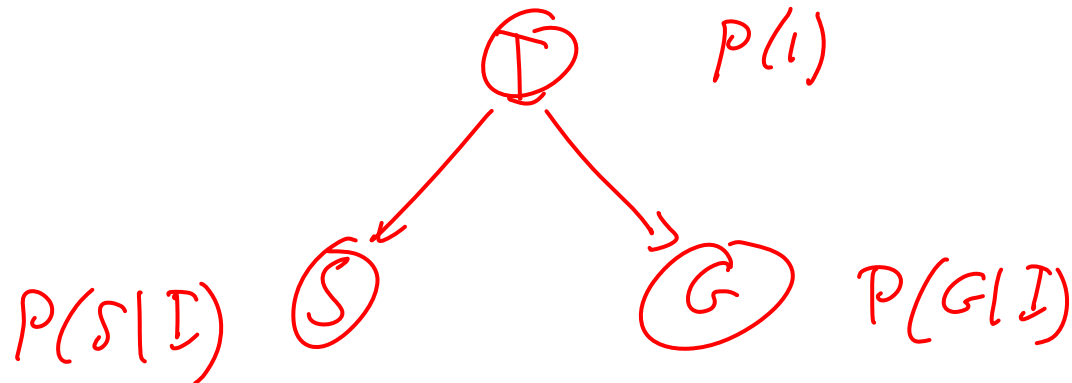


$$P(I = VH) = 0.8$$

$P(G|I)$

I \ G	A	B	
VH	0.8	0.2	$\Sigma_i = 1$
H	0.6	0.4	$\Sigma_i = 1$

Example: 3 variables



$$P(I, S, G) = P(I) P(G|I) P(S|I)$$

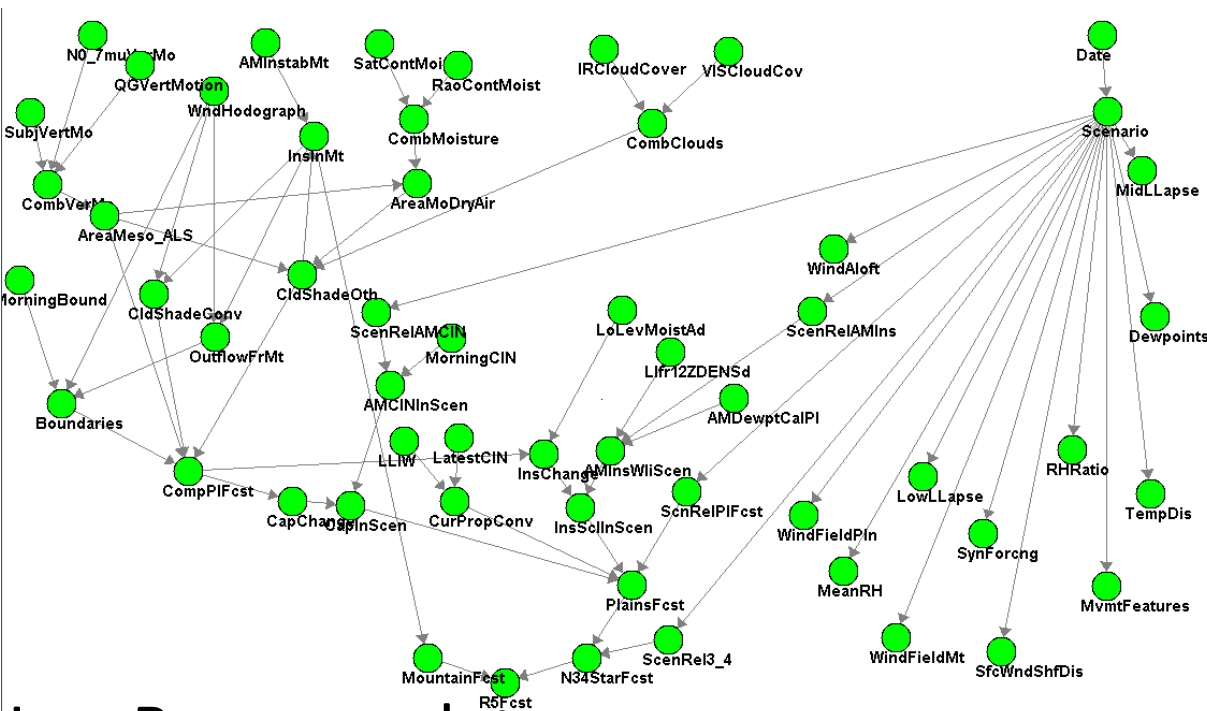
Example: Naïve Bayes models

- Class variable Y
- Evidence variables X_1, \dots, X_n
- Assume that $X_A \perp X_B \mid Y$
for all subsets X_A, X_B of $\{X_1, \dots, X_n\}$
- Conditional parametrization:
 - Specify $P(Y)$
 - Specify $P(X_i \mid Y)$
- Joint distribution

$$\underline{P(x_1, \dots, x_n, y) = P(y) \prod_i P(x_i \mid y)}$$

Today: Bayesian networks

- Compact representation of distributions over large number of variables
- (Often) allows efficient exact inference (computing marginals, etc.)



HailFinder

56 vars

~ 3 states each

→ $\sim 10^{26}$ terms

> 10.000 years

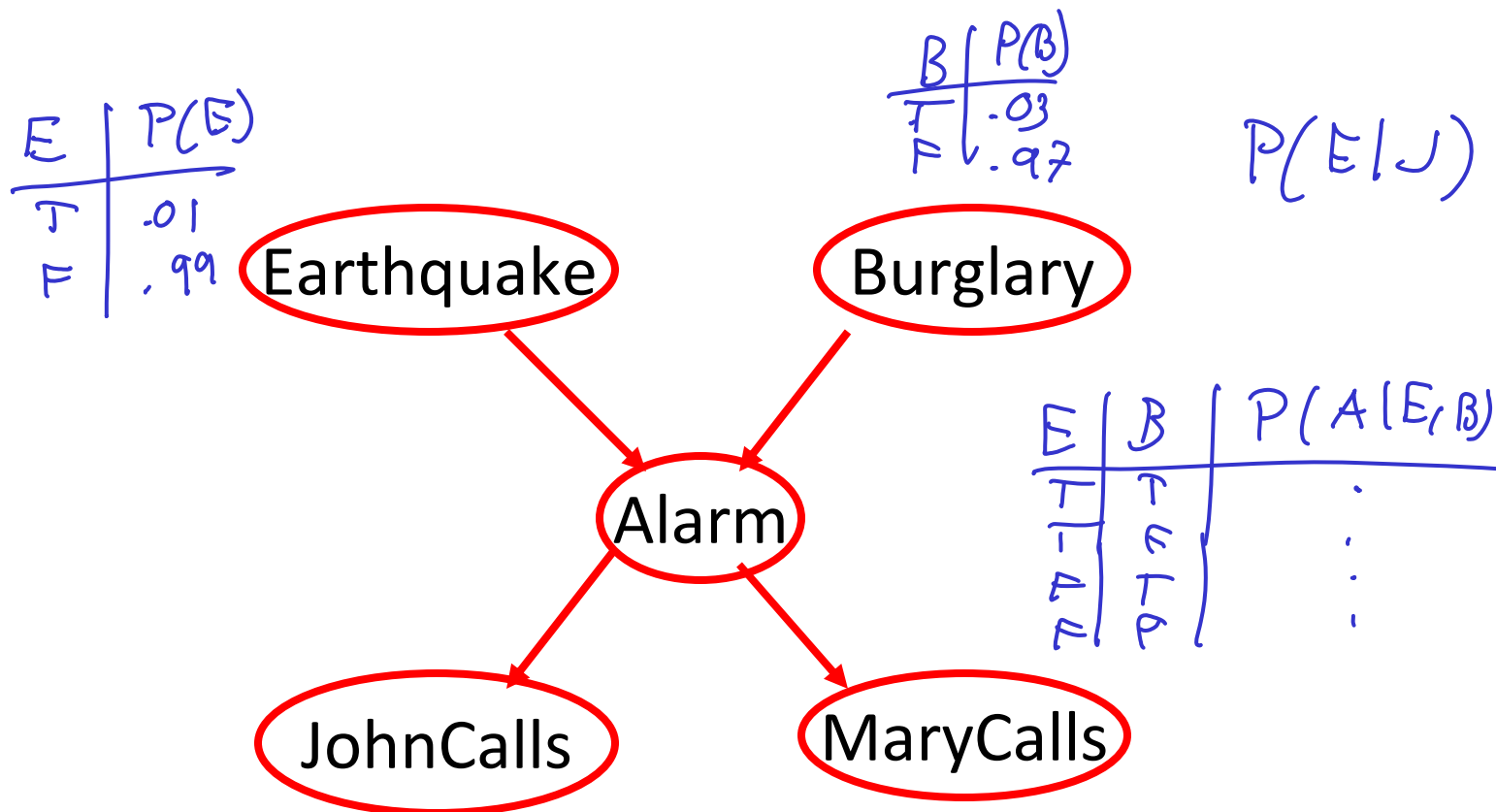
on Top

supercomputers

JavaBayes applet

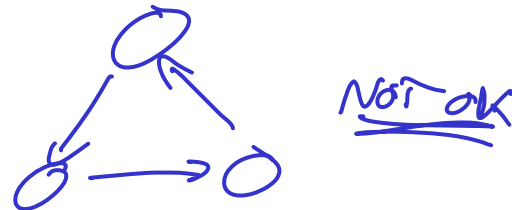
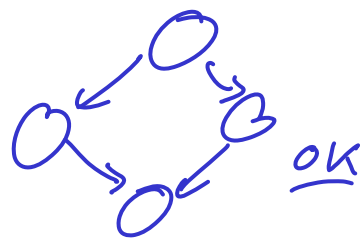
Causal parametrization

- Graph with directed edges from (immediate) causes to (immediate) effects



Bayesian networks

- A **Bayesian network structure** is a directed, acyclic graph G , where each vertex s of G is interpreted as a random variable X_s (with unspecified distribution)



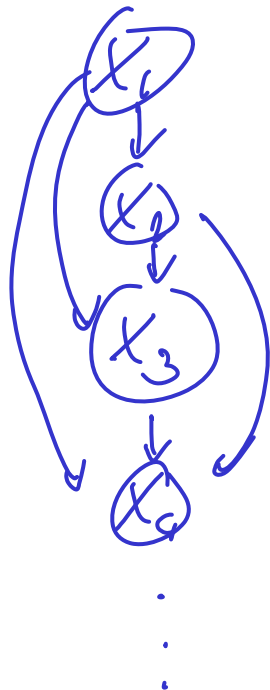
- A **Bayesian network** (G, P) consists of
 - A BN structure G and ..
 - ..a set of conditional probability distributions (CPTs) $P(X_s \mid \mathbf{Pa}_{X_s})$, where \mathbf{Pa}_{X_s} are the parents of node X_s such that
 - (G, P) defines joint distribution

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \mathbf{Pa}_{X_i})$$

Bayesian networks

- Can every probability distribution be described by a BN?

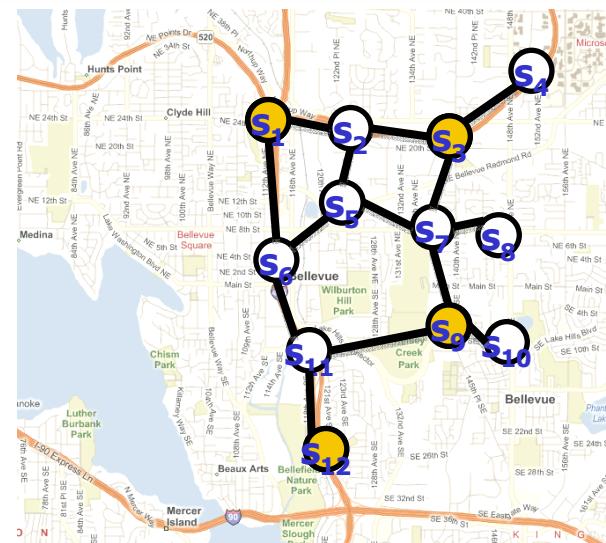
$$P(x_1, \dots, x_n) = P(x_1) P(x_2 | x_1) \dots P(x_n | x_1, \dots, x_{n-1})$$



Representing the world using BNs



represent

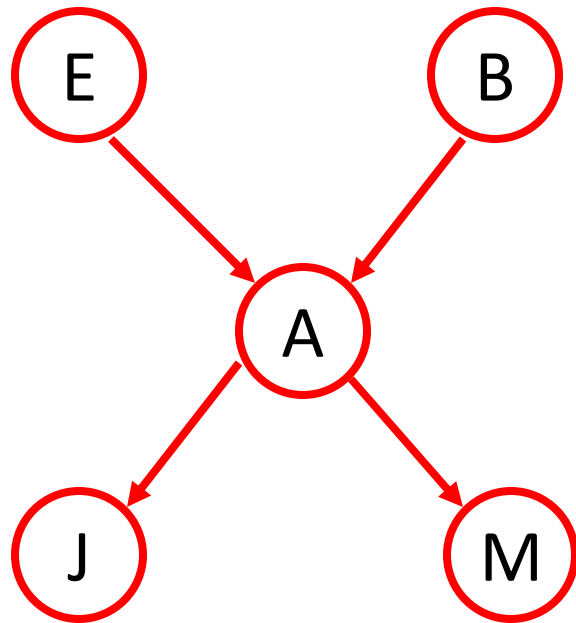


True distribution P'
with cond. ind. $I(P')$

Bayes net (G, P)
with $I(P)$

- Want to make sure that $I(P) \subseteq I(P')$
- Need to understand CI properties of BN (G, P)

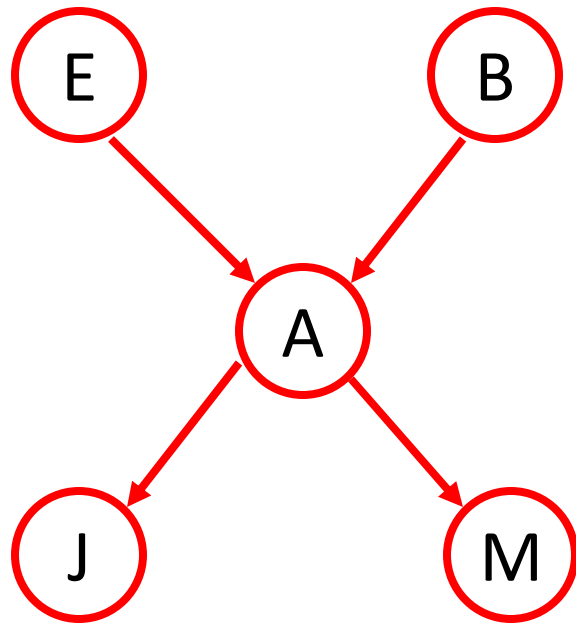
Which kind of CI does a BN imply?



$$E \perp B$$

$$\begin{aligned} P(E, B) &= \sum_{A, J, M} P(E, B, A, J, M) \\ &= \sum_{A, J, M} \underbrace{P(E) P(B) P(A|E, B)}_{P(E) P(B)} \underbrace{P(J|A) P(M|A)}_{P(J, M|A)} \\ &= P(E) P(B) \sum_{A, J, M} P(A|E, B) P(J, M|A) \\ &= P(E) P(B) \end{aligned}$$

Which kind of CI does a BN imply?



$$J \perp M \mid A$$

$$P(J|AM) = \frac{P(J, A, M)}{P(A, M)}$$

$$P(J, A, M) = \sum_{E, B} P(J, A, M, E, B)$$

$$= \sum_{E, B} P(E) P(B) P(A|E, B) P(J|A) P(M|A)$$

$$= P(J|A) P(M|A) \underbrace{\sum_{E, B} P(E) P(B) P(A|E, B)}_{= P(A)} = P(J, M, A)$$

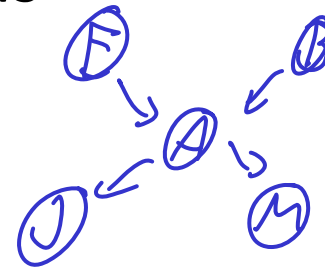
$$\Rightarrow P(J|AM) = \frac{P(J, A, M)}{P(A, M)} = \frac{P(J|A) \cancel{P(A, M)}}{\cancel{P(A, M)}} = P(J|A) \quad \square$$

Local Markov Assumption

- Each BN Structure G is associated with the following conditional independence assumptions

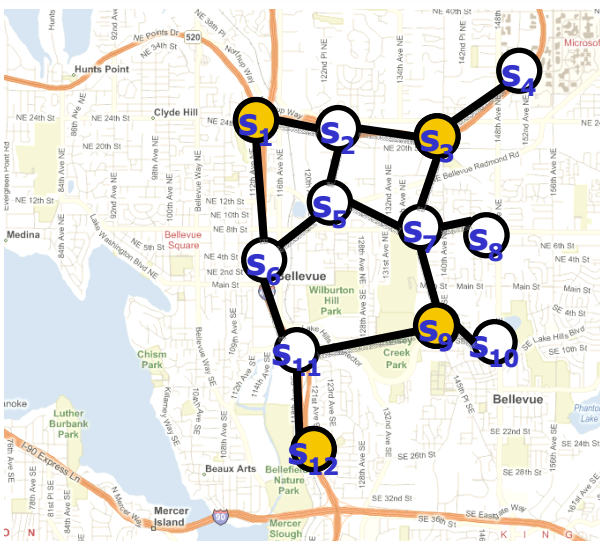
$$J \perp B \mid A$$

$$X \perp \text{NonDescendants}_X \mid \text{Pa}_X$$



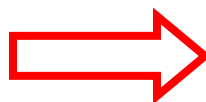
- We write $I_{loc}(G)$ for these conditional independences
- Suppose (G, P) is a Bayesian network representing P
Does it hold that $I_{loc}(G) \subseteq I(P)$?
If this holds, we say G is an I-map for P.

Factorization Theorem



True distribution P
can be represented exactly as

$$I_{\text{loc}}(G) \subseteq I(P)$$

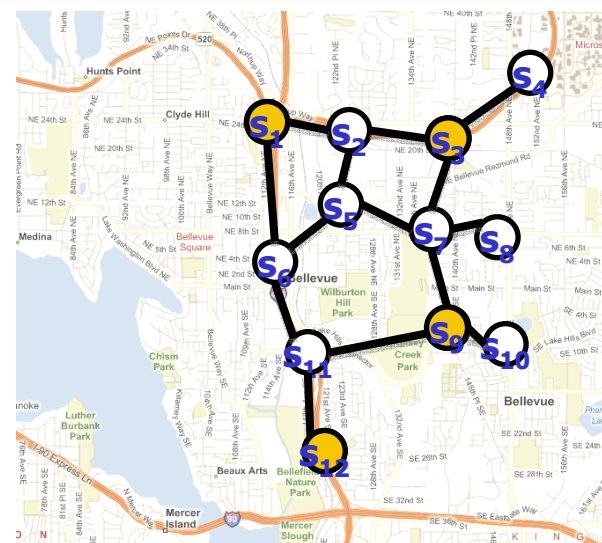


$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \mathbf{Pa}_{X_i})$$

G is an **I-map** of P
(independence map)

i.e., P can be represented as
a Bayes net (G, P)

Factorization Theorem



True distribution P
can be represented exactly as
a Bayes net (G, P)

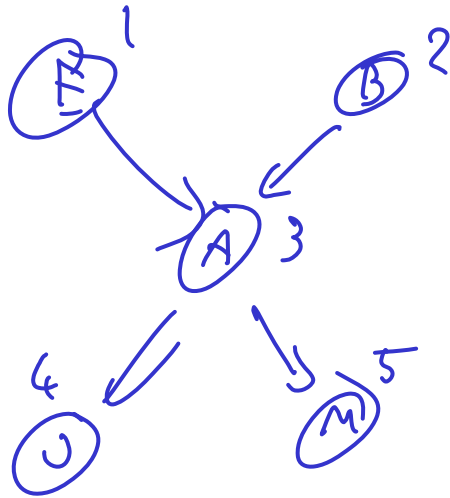
$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \mathbf{Pa}_{X_i})$$



$$I_{\text{loc}}(G) \subseteq I(P)$$

G is an **I-map** of P
(independence map)

Proof: I-Map to factorization



Ordering $\pi : \{1 \dots n\} \rightarrow \{1 \dots n\}$ topological

If: $\forall X_j$ descendent of X_i

$$\pi(j) > \pi(i)$$

Can find topological ordering in linear time

$$P(X_1 \dots X_n) = \prod \underbrace{P(X_{\pi(i)} \mid X_{\pi(i_1)} \dots X_{\pi(i_{i-1})})}_{P(X_{\pi(i)} \mid \text{Pa}_{X_{\pi(i)}})} \quad \text{Chain rule}$$

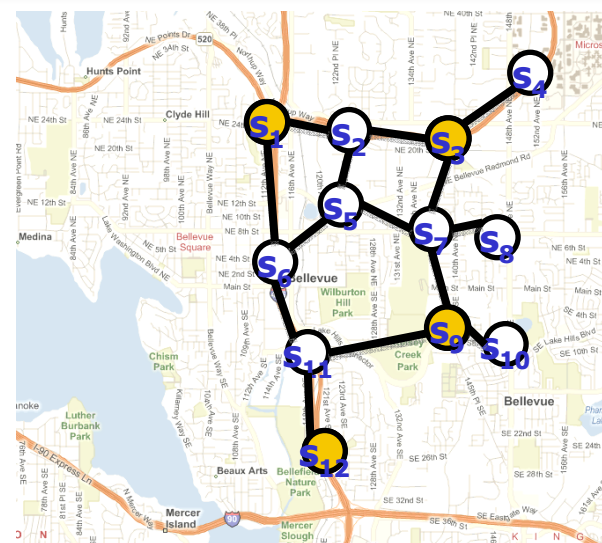
□

Factorization Theorem



True distribution P
can be represented exactly as
a Bayes net (G, P)

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \mathbf{Pa}_{X_i})$$



$$I_{\text{loc}}(G) \subseteq I(P)$$

G is an **I-map** of P
(independence map)

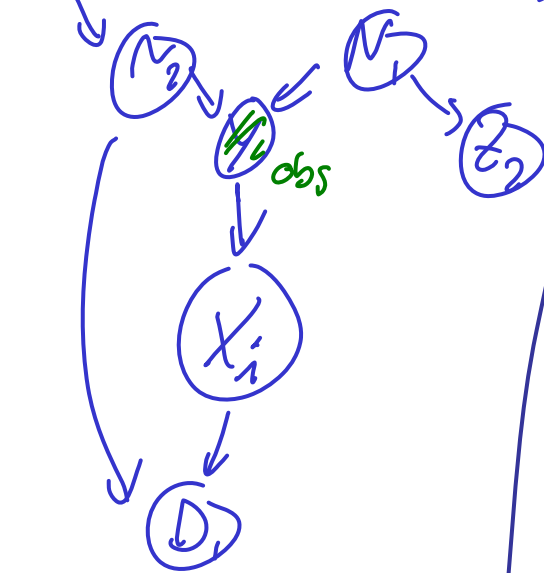
The general case

$$P(x_1 \dots x_n) = \prod P(x_i | Pa_{x_i}) \stackrel{!}{\Rightarrow} \forall x_i: \forall N \subseteq \underline{\text{Nondesc}_{x_i}} \\ x_i \perp N \mid Pa_{x_i}$$

$$Y = Pa_{x_i}$$

$$Z = \text{Nondesc}_{x_i} \setminus (Y \cup N)$$

$$\textcircled{2} D = \text{Desc}(x_i)$$



$$P(x_i | Y, N) = \frac{P(x_i, Y, N)}{\underline{P(Y, N)}}$$

$$P(x_i, Y, N) = \sum_{Z, D} P(x_i, Y, N, Z, D)$$

$$= \sum_{Z, D} P(x_i | Y) \prod_{x \in D} P(x | Pa_x) \prod_{x' \in (N \cup Z \cup Y)} P(x' | Pa_{x'})$$

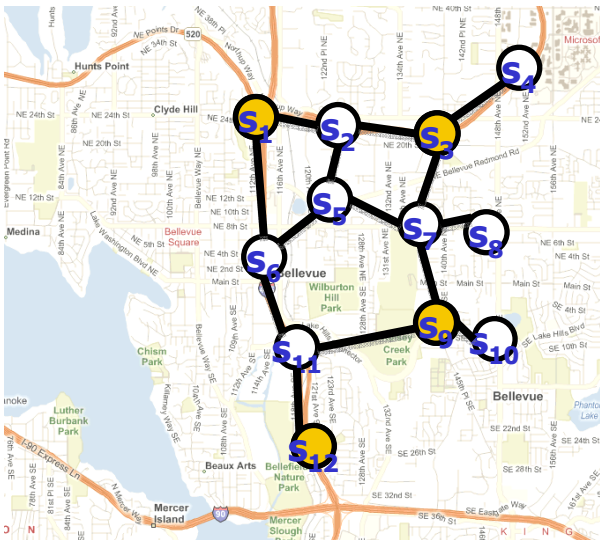
$$= P(x_i | Y) \underbrace{\sum_Z \prod_{x' \in (N \cup Z \cup Y)} P(x' | Pa_{x'})}_{=1} \underbrace{\sum_D \prod_{x \in D} P(x | Pa_x)}_{=1}$$

$$P(Y, N) = \sum_{x_i} P(x_i, Y, N) = \sum_Z \prod_{x' \in (N \cup Z \cup Y)} P(x' | Pa_{x'}) \underbrace{\sum_{x_i} P(x_i | Y)}_{=1}$$

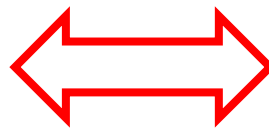
$$\Rightarrow P(x_i | Y, N) = P(x_i | Y)$$

□

Factorization Theorem



$$I_{\text{loc}}(G) \subseteq I(P)$$



G is an I-map of P
(independence map)

True distribution P
can be represented exactly as
Bayesian network (G,P)

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \text{Pa}_{X_i})$$

Defining a Bayes Net

- Given random variables and known conditional independences
- Pick ordering X_1, \dots, X_n of the variables
- For each X_i
 - Find minimal subset $A \subseteq \{X_1, \dots, X_{i-1}\}$ such that $X_i \perp X_{\neg A} \mid A$,
where $\neg A = \{X_1, \dots, X_n\} \setminus A$
Ensure local Markov property holds!
 - Specify / learn $\text{CPD}(X_i \mid A)$ \nwarrow *A parents of X_i*
- Ordering matters a lot for compactness of representation! More later this course.

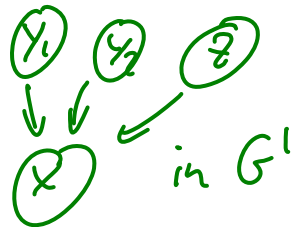
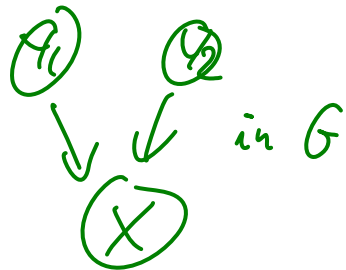
Adding edges doesn't hurt

Theorem:

Let G be an I-Map for P , and G' be derived from G by adding an edge. Then G' is an I-Map of P
(G' is strictly more expressive than G)

• **Proof** : wtp: $I_{loc}(G') \subseteq I_{loc}(G)$

Then $I_{loc}(G') \subseteq I(P)$ since $I_{loc}(G) \subseteq I(P)$



wtp: $X \perp \text{Nondesc}(X; G) \mid \text{Pa}_G(X; G)$
 $\Rightarrow X \perp \text{Nondesc}(X; G') \mid \text{Pa}_{G'}(X; G')$

$X \perp N_i Z \mid Y \Rightarrow X \perp N \mid Y, Z$

holds b/c: Weak Union property of C.I. \square

Additional conditional independencies

- BN specifies joint distribution through conditional parameterization that satisfies Local Markov Property
- But we also talked about additional properties of CI
 - Weak Union, Intersection, Contraction, ...
- Which additional CI does a particular BN specify?
 - All CI that can be derived through algebraic operations

Local Markov prop. $I_{loc}(G) \subseteq I(G)$

What you need to know

- Bayesian networks
- Local Markov property
- I-Maps
- Factorization Theorem

Tasks

- Subscribe to Mailing list
<https://utils.its.caltech.edu/mailman/listinfo/cs155>
- Read Koller & Friedman Chapter 3.1-3.3
- Form groups and think about class projects. If you have difficulty finding a group, email Pete Trautman