

Probabilistic Graphical Models

Undirected Graphical Models

Probabilistic Graphical Models

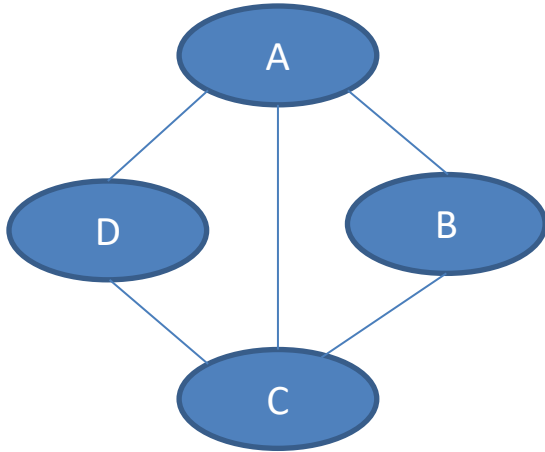
Directed graphical models

- Bayes Nets
- Conditional dependence

Undirected graphical models

- Markov random fields (MRFs)
- Factor graphs

General Markov Networks



$$\Phi = \{\phi_1(\mathbf{D}_1), \dots, \phi_k(\mathbf{D}_k)\}$$

$$\tilde{P}_\Phi(X_1, \dots, X_n) = \prod_i \phi_i(\mathbf{D}_i)$$

$$Z_\Phi = \sum_{X_1, \dots, X_n} \tilde{P}_\Phi(X_1, \dots, X_n)$$

$$P_\Phi(X_1, \dots, X_n) = \frac{1}{Z_\Phi} \prod_i \phi_i(\mathbf{D}_i)$$

a^1	b^1	c^1	0.25
a^1	b^1	c^2	0.35
a^1	b^2	c^1	0.08
a^1	b^2	c^2	0.16
a^2	b^1	c^1	0.05
a^2	b^1	c^2	0.07
a^2	b^2	c^1	0
a^2	b^2	c^2	0
a^3	b^1	c^1	0.15
a^3	b^1	c^2	0.21
a^3	b^2	c^1	0.09
a^3	b^2	c^2	0.18

Log-linear Representation

$$\tilde{P} = \prod_i \phi_i(\mathbf{D}_i)$$

Original parameterization

$$\tilde{P} = \exp \left(- \sum_j w_j f_j(\mathbf{D}_j) \right)$$

Log-linear parameterization

Features (f_j) are functions (like factors) without the non-negativity assumption.

Each feature has a single weight. (coefficient, w_j)

Different features can have the same scope.

Log-linear Representation

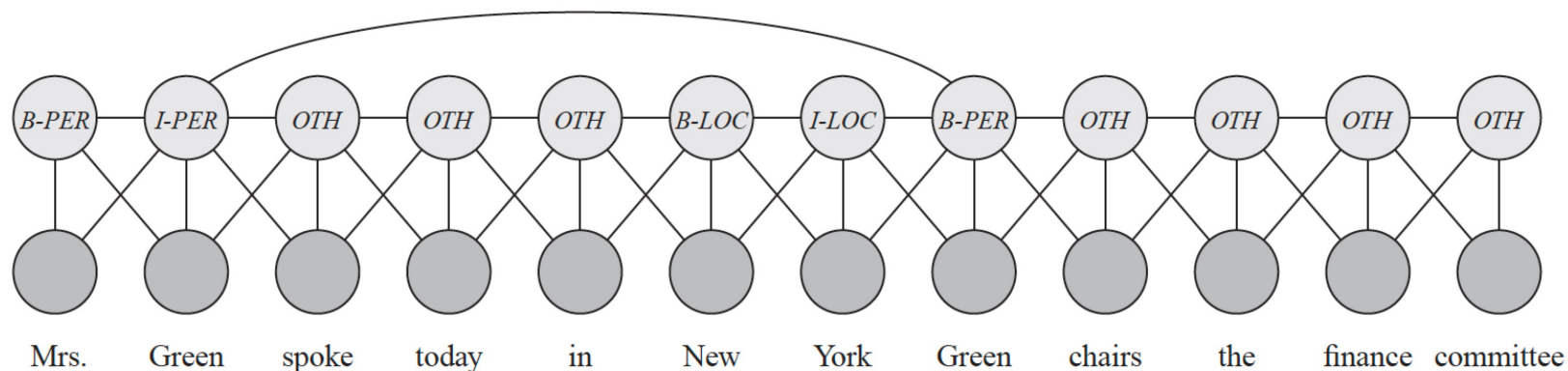
$$\phi(X_1, X_2) = \begin{pmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{pmatrix} \quad f_{12}^{ij} = I(X_1 = i \text{ and } X_2 = j)$$

One feature for each i,j value

$$\phi(X_2, X_3) = \exp\left(-\sum_{kl} w_{kl} f_{12}^{kl}(X_1, X_2)\right)$$

$$w_{kl} = -\log(a_{kl})$$

Feature Example: Text



KEY

<i>B-PER</i>	Begin person name	<i>I-LOC</i>	Within location name
<i>I-PER</i>	Within person name	<i>OTH</i>	Not an entity
<i>B-LOC</i>	Begin location name		

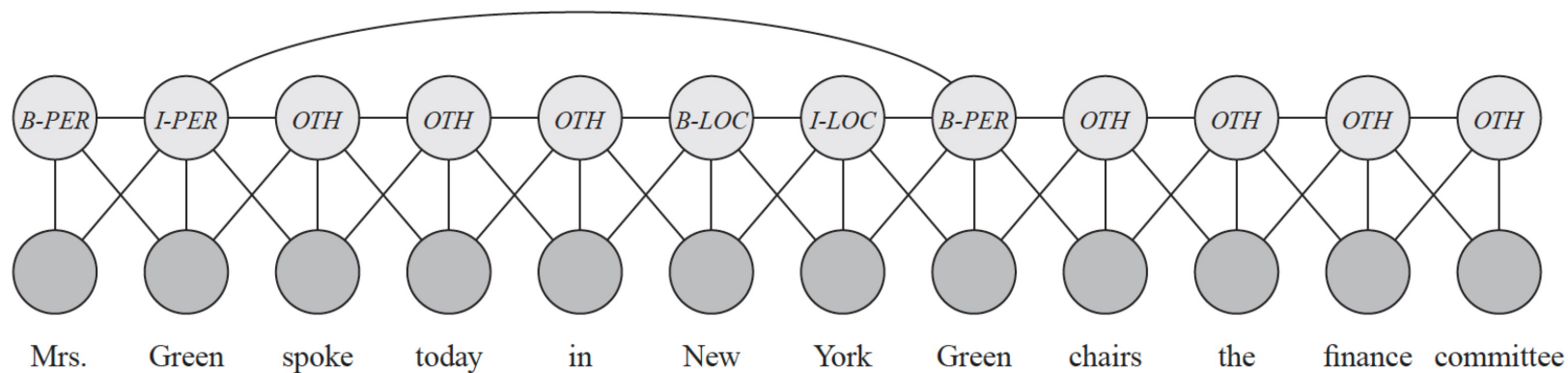
$$\phi(Y_t, Y_{t+1})$$
$$\phi(Y_t | X_{t-1}, X_t, X_{t+1})$$

Problem: Extract entities from a word sequence

For each word: T , a target variable, Y_t , which indicates the entity type of the word.

Possible outcomes of Y_t : B-Person, I-Person, B-Location, I-Location, B-Organization, I-Organization, and Other..

Feature Example: Text



$$\phi(Y_t, Y_{t+1})$$
$$\phi(Y_t | X_{t-1}, X_t, X_{t+1})$$

KEY

<i>B-PER</i>	Begin person name	<i>I-LOC</i>	Within location name
<i>I-PER</i>	Within person name	<i>OTH</i>	Not an entity
<i>B-LOC</i>	Begin location name		

OR

$$f(Y_t, X_t) = 1\{Y_i = Person, X_i \text{ is capitalized}\}$$

Problem: Extract entities from a word sequence

For each word: T , a target variable, Y_t , which indicates the entity type of the word.

Possible outcomes of Y_t : B-Person, I-Person, B-Location, I-Location, B-Organization, I-Organization, and Other..

Example: Ising Models

$$E(x_1, \dots, x_n) = - \sum_{i < j} w_{i,j} x_i x_j - \sum_i u_i x_i$$

$$x_i \in \{-1, 1\}$$

$$f_{i,j}(X_i, X_j) = X_i \cdot X_j$$

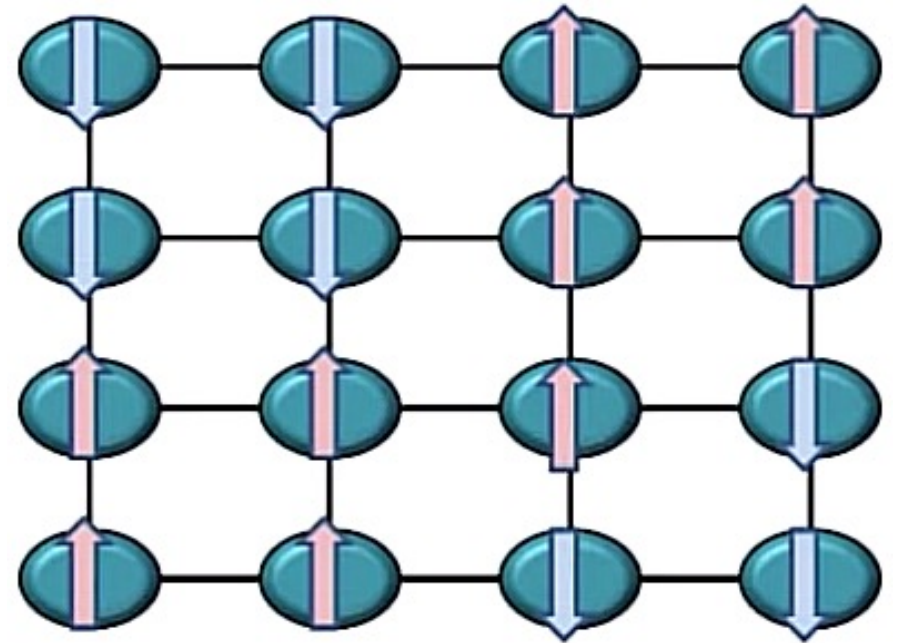
Example: Ising Models

$$E(x_1, \dots, x_n) = - \sum_{i < j} w_{i,j} x_i x_j - \sum_i u_i x_i$$

$$x_i \in \{-1, 1\}$$

$$f_{i,j}(X_i, X_j) = X_i \cdot X_j$$

$$P(\mathbf{X}) \propto e^{-\frac{1}{T}E(\mathbf{X})}$$

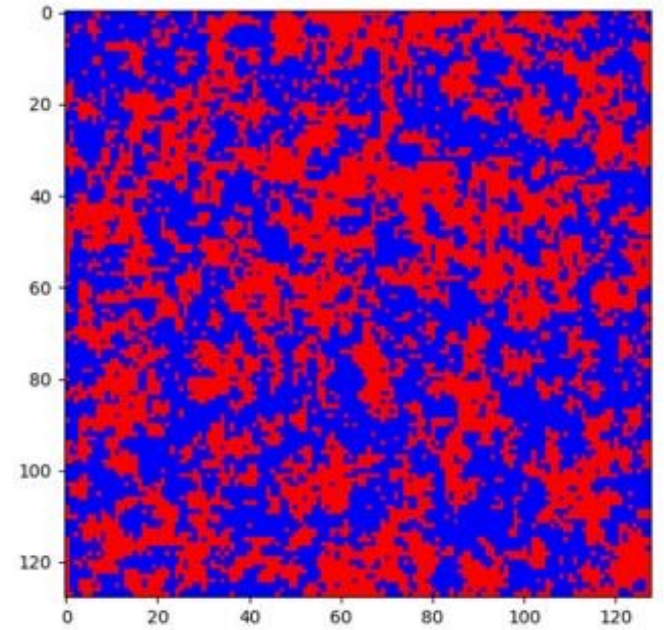
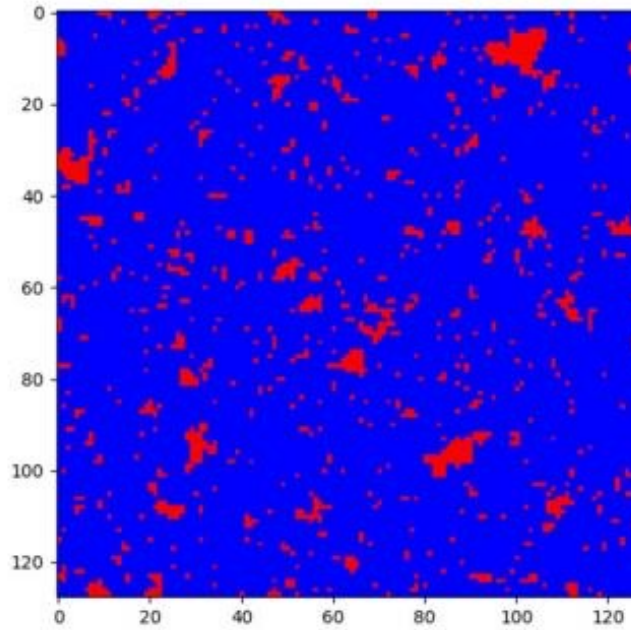


Example: Ising Models

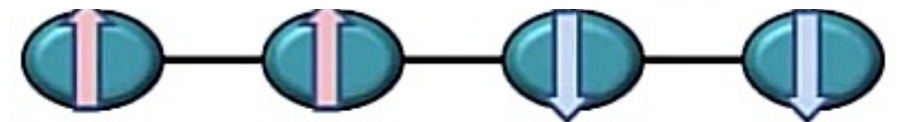
$$E(x_1, \dots, x_n) = - \sum_{i < j} w_{ij} x_i x_j$$

$$x_i \in \{ \pm 1 \}$$

$$f_{i,j}(X_i, X_j) = w_{ij} x_i x_j$$



$$P(\mathbf{X}) \propto e^{-T^{-1} E(\mathbf{X})}$$



As T grows, w_{ij} 's become smaller

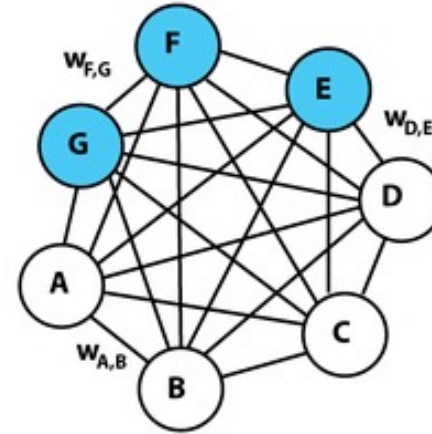
Example: Boltzman machine

$$E = - \sum_{i < j} w_{ij} s_i s_j + - \sum_i \theta_i s_i$$

$$s_i \in \{0, 1\}$$

- w_{ij} is the connection strength between unit j and unit i .
- s_i is the state, $s_i \in \{0,1\}$, of unit i .
- θ_i is the bias of unit i in the global energy function. ($-\theta_i$ is the activation threshold for the unit.)

Model for neural activation



Example: Ising Models

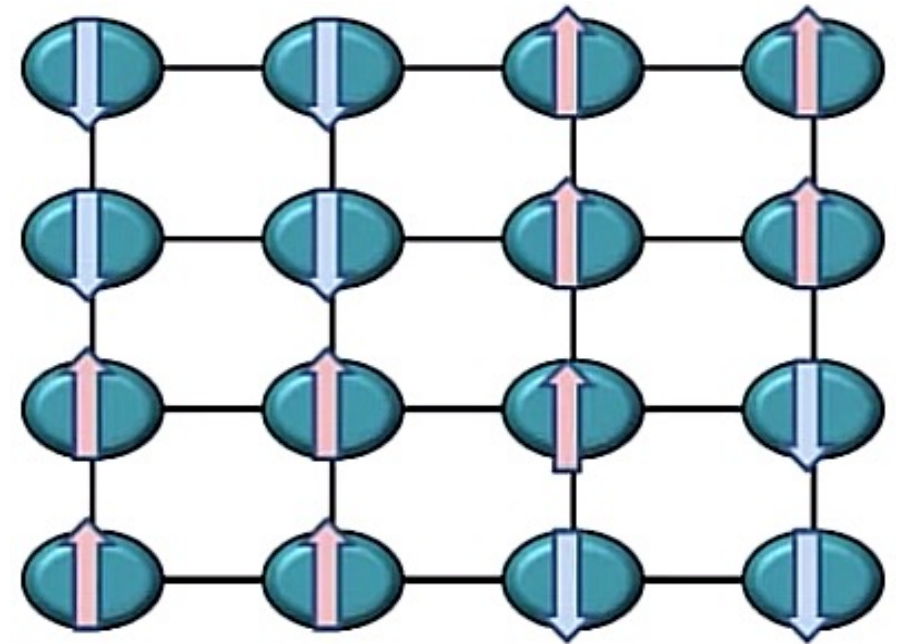
$$E(x_1, \dots, x_n) = - \sum_{i < j} w_{i,j} x_i x_j - \sum_i u_i x_i$$

$$x_i \in \{-1, 1\}$$

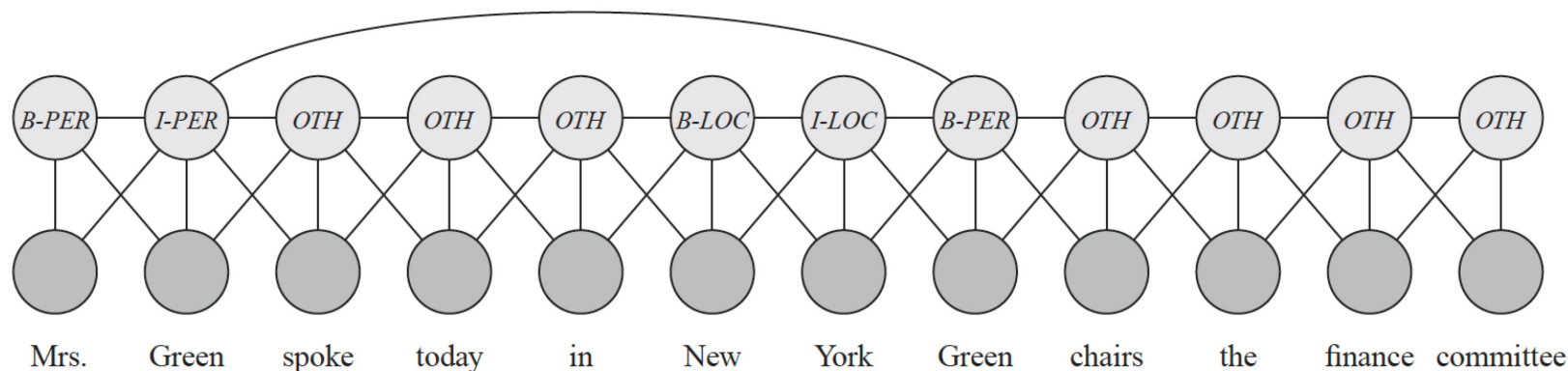
$$f_{i,j}(X_i, X_j) = X_i \cdot X_j$$

$$P(\mathbf{X}) \propto e^{-\frac{1}{T}E(\mathbf{X})}$$

$w_{i,j}$ will in general be the same for every pair i, j



Feature Example: Text



$$\phi(Y_t, Y_{t+1})$$
$$\phi(Y_t | X_{t-1}, X_t, X_{t+1})$$

KEY

<i>B-PER</i>	Begin person name	<i>I-LOC</i>	Within location name
<i>I-PER</i>	Within person name	<i>OTH</i>	Not an entity
<i>B-LOC</i>	Begin location name		

OR

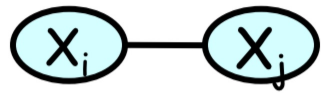
$$f(Y_t, X_t) = 1\{Y_i = \text{Person}, X_i \text{ is capitalized}\}$$

Same energy terms $w_k f_k(X_i, Y_i)$ repeat for all positions i in the sequence

Same energy terms $w_m f_m(Y_i, Y_{i+1})$ also repeat for all positions i

Metric MRFs

- All X_i take values in label space V



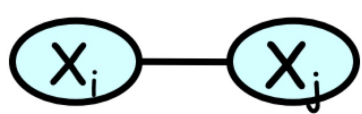
want X_i and X_j to take "similar" values

Distance function $\mu: V \times V \rightarrow \mathbb{R}^+$

- $\mu(v, v) = 0$ for all v
- Symmetry: $\mu(v_1, v_2) = \mu(v_2, v_1)$ for all v_1, v_2
- Triangle inequality: $\mu(v_1, v_2) \leq \mu(v_1, v_3) + \mu(v_3, v_2)$

Metric MRFs

- All X_i take values in label space V



want X_i and X_j to take "similar" values

- Distance function $\mu : V \times V \rightarrow \mathbb{R}$

$$f_{i,j}(X_i, X_j) = \mu(X_i, X_j)$$

lower distance

higher

higher probability

$$\exp(-w_{ij} f_{ij}(X_i, X_j))$$

$$w_{ij} > 0$$

values of X_i and X_j far in μ

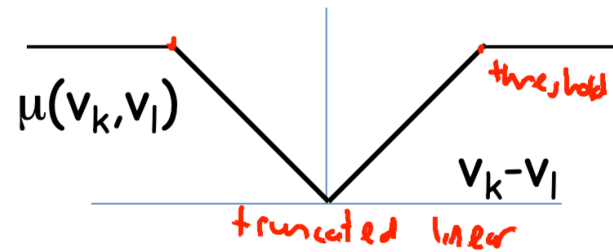
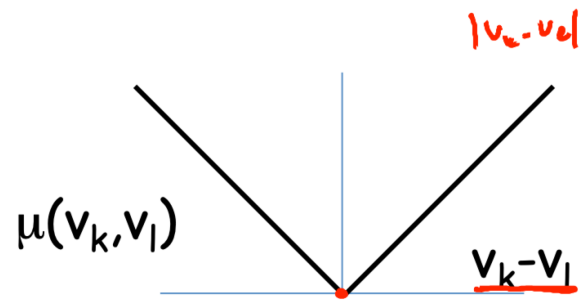


lower probability

Example Metric MRFs

$$\mu(v_k, v_l) = \begin{cases} 0 & v_k = v_l \\ 1 & \text{otherwise} \end{cases}$$

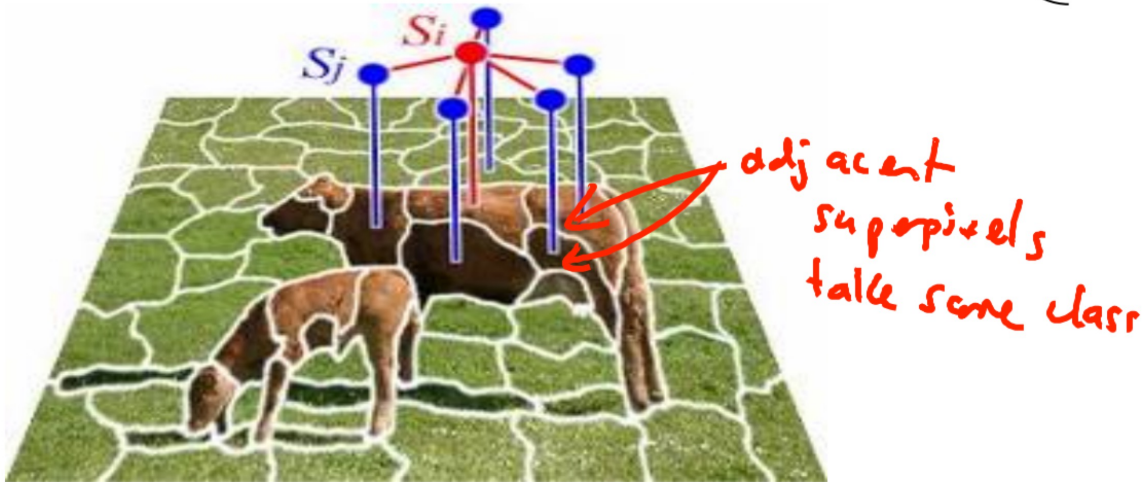
$$\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$



Example: Image Segmentation

$$\mu(v_k, v_l) = \begin{cases} 0 & v_k = v_l \\ 1 & \text{otherwise} \end{cases}$$

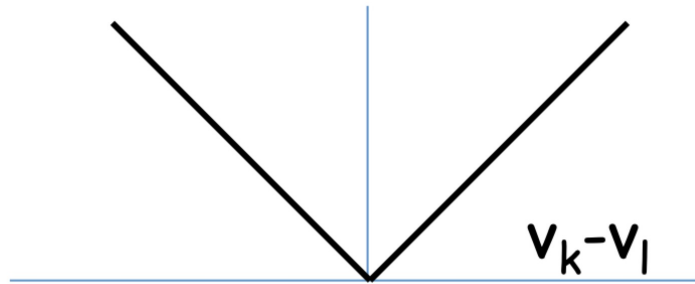
0	1	1	1
1	0	1	1
1	1	0	1
1	1	1	0



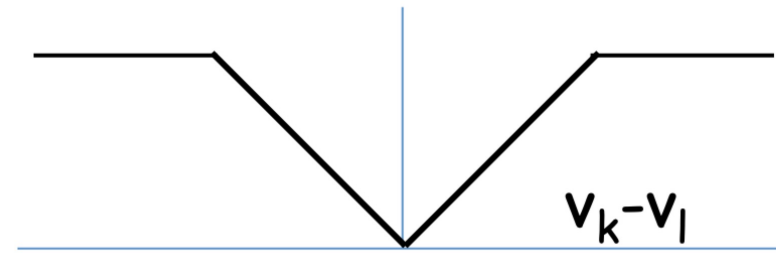
Example: Denoising



$$\mu(v_k, v_l) = |v_k - v_l|$$



$$\mu(v_k, v_l) = \min(|v_k - v_l|, d)$$



Similar idea for stereo reconstruction

Repeated Features

- Need to specify for each feature f_k a set of scopes $Scopes[f_k]$
- For each $\mathbf{D}_k \in Scopes[f_k]$, we have a term $w_k f_k(\mathbf{D}_k)$ in the energy function
- $w_k \sum_{\mathbf{D}_k} f(\mathbf{D}_k)$
- Parameters and structure are reused within an MN and across different MNs

Pt 2: Inference

Queries on PGMs

Conditional Probability Queries

- Evidence: $E = e$
- Query: a subset of variables Y
- Task: compute $P(Y|E = e)$

Applications

- Medical/fault diagnosis
- Pedigree analysis

NP-hardness

Exact Inference is NP hard:

Approximate Inference is also NP hard

Queries on PGMs

Conditional Probability Queries

- Evidence: $E = e$
- Query: a subset of variables Y
- Task: compute $P(Y|E = e)$

Applications

- Medical/fault diagnosis
- Pedigree analysis

NP-hardness

Exact Inference is NP hard:

Approximate Inference is also NP hard

Why is the expression $\sum_{\bar{W}} P(\bar{Y}, \bar{W}, \bar{e})$ hard to compute in general?

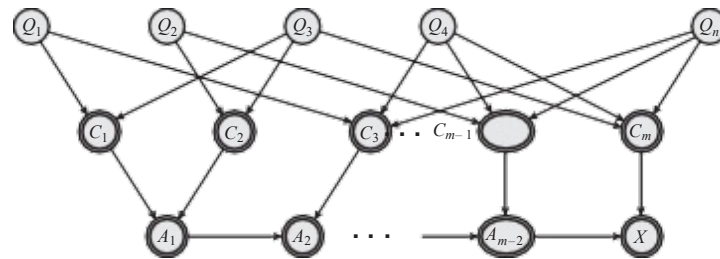
It may be intractable to sum over all the different values that \bar{W} can take.

The summation over all values of \bar{W} is exponential. If \bar{W} has 100 binary variables, then summing will take 2^{100} operations.

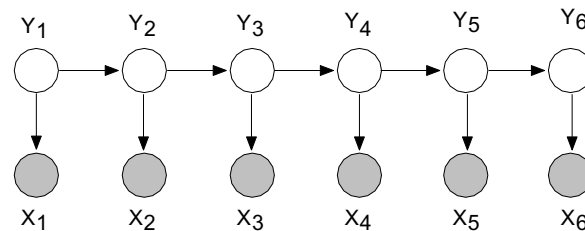
$P(\bar{Y}, \bar{W}, \bar{e})$ is always easy to compute because it is just the product of all CPDs.

Probabilistic inference in practice

- NP-hardness simply says that there **exist** difficult inference problems
- Real-world inference problems are not necessarily as hard as these worst-case instances
- The reduction from SAT created a very complex Bayesian network:



- Some graphs are **easy** to do inference in! For example, inference in hidden Markov models

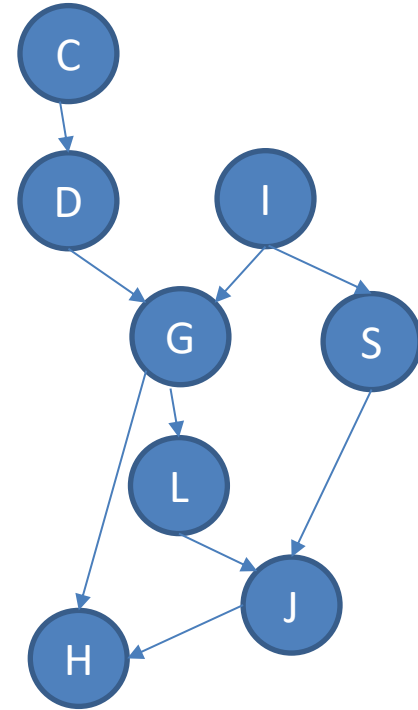


and other tree-structured graphs can be performed in **linear time**

Sum-product Inference

$$\phi_C(C)\phi_D(C, D)\phi_I(I)\phi_G(G, I, D)$$
$$\phi_S(S, I)\phi_L(L, G)\phi_J(J, L, S)\phi_H(H, G, J)$$

Compute $P(J)$

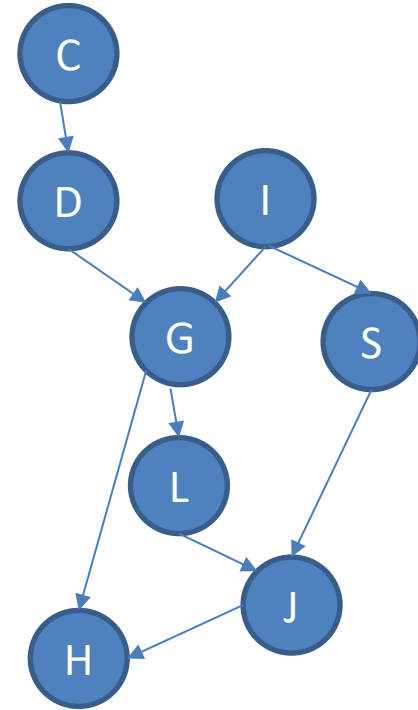


Sum-product Inference

$$\phi_C(C)\phi_D(C, D)\phi_I(I)\phi_G(G, I, D)$$
$$\phi_S(S, I)\phi_L(L, G)\phi_J(J, L, S)\phi_H(H, G, J)$$

Compute $P(J)$

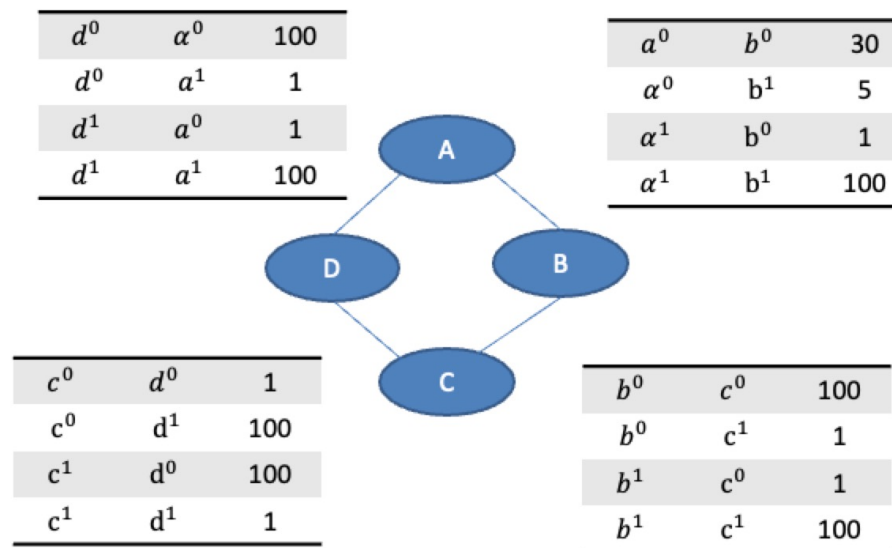
$$\sum_{C, D, I, G, S, L, H} \phi_C(C)\phi_D(C, D)\phi_I(I)\phi_G(G, I, D)\phi_S(S, I)\phi_L(L, G)\phi_J(J, L, S)\phi_H(H, G, J)$$



Sum-product Inference for MNs

$$\tilde{P}(D) = \sum_{A,B,C} \phi_1(A,B)\phi_2(B,C)\phi_3(C,D)\phi_4(D,A)$$

What about the normalization constant?



Introducing Evidence

$$P(Y \mid E = e) = \frac{P(Y, E = e)}{P(E = e)}$$

$$P(Y, E = e) = \sum_W P(Y, W, E = e)$$

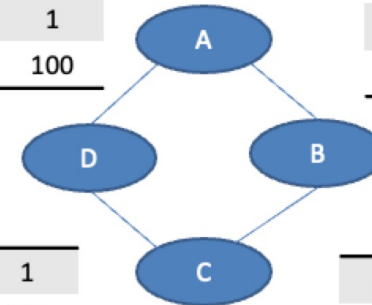
$$W = V \setminus E \cup Y$$

Use the reduced factors:

Example: $A=0$

d^0	a^0	100
d^0	a^1	1
d^1	a^0	1
d^1	a^1	100

a^0	b^0	30
a^0	b^1	5
a^1	b^0	1
a^1	b^1	100



c^0	d^0	1
c^0	d^1	100
c^1	d^0	100
c^1	d^1	1

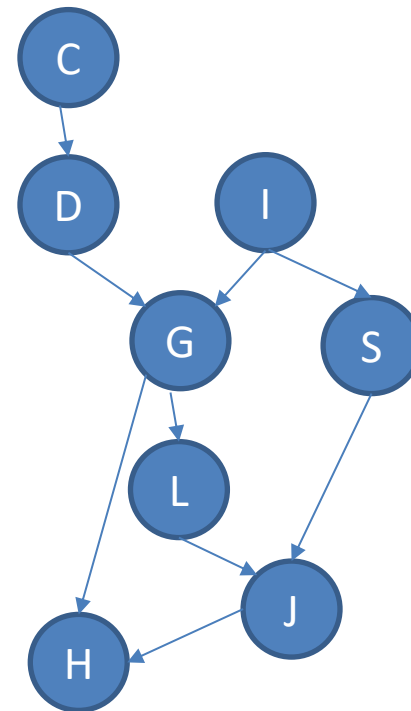
b^0	c^0	100
b^0	c^1	1
b^1	c^0	1
b^1	c^1	100

Sum-product Inference

$$\phi_C(C)\phi_D(C, D)\phi_I(I)\phi_G(G, I, D)$$
$$\phi_S(S, I)\phi_L(L, G)\phi_J(J, L, S)\phi_H(H, G, J)$$

Compute $P(J, i, h)$

$$\sum_{C, D, I, G, S, L, H} \phi_C(C)\phi_D(C, D)\phi_I(I)\phi_G(G, I, D)\phi_S(S, I)\phi_L(L, G)\phi_J(J, L, S)\phi_H(H, G, J)$$



Algorithms: Conditional Probability

Push summations into factor product

- Variable elimination (dynamic programming)

Message passing over a graph

- Belief propagation (exact)
- Variational approximations
- Random sampling instantiations
- Markov chain Monte Carlo (MCMC)
- Importance sampling

Inference in Chains



- We want to compute $p(D)$
- $p(D)$ is a **set** of values, $\{p(D = d), d \in Val(D)\}$.
- Algorithm computes sets of values at a time – an entire distribution

By the chain rule and conditional independence, the joint distribution factors as

$$P(A, B, C, D) = P(A)P(B | A)P(C | B)P(D | C)$$

In order to compute $p(D)$, we have to marginalize over A, B, C :

$$P(A, B, C, D) = \sum_{A, B, C} P(A)P(B | A)P(C | B)P(D | C)$$

Let's be a bit more explicit...

$$\begin{aligned} & P(a^1) P(b^1 | a^1) P(c^1 | b^1) P(d^1 | c^1) \\ + & P(a^2) P(b^1 | a^2) P(c^1 | b^1) P(d^1 | c^1) \\ + & P(a^1) P(b^2 | a^1) P(c^1 | b^2) P(d^1 | c^1) \\ + & P(a^2) P(b^2 | a^2) P(c^1 | b^2) P(d^1 | c^1) \\ + & P(a^1) P(b^1 | a^1) P(c^2 | b^1) P(d^1 | c^2) \\ + & P(a^2) P(b^1 | a^2) P(c^2 | b^1) P(d^1 | c^2) \\ + & P(a^1) P(b^2 | a^1) P(c^2 | b^2) P(d^1 | c^2) \\ + & P(a^2) P(b^2 | a^2) P(c^2 | b^2) P(d^1 | c^2) \end{aligned}$$

$$\begin{aligned} & P(a^1) P(b^1 | a^1) P(c^1 | b^1) P(d^2 | c^1) \\ + & P(a^2) P(b^1 | a^2) P(c^1 | b^1) P(d^2 | c^1) \\ + & P(a^1) P(b^2 | a^1) P(c^1 | b^2) P(d^2 | c^1) \\ + & P(a^2) P(b^2 | a^2) P(c^1 | b^2) P(d^2 | c^1) \\ + & P(a^1) P(b^1 | a^1) P(c^2 | b^1) P(d^2 | c^2) \\ + & P(a^2) P(b^1 | a^2) P(c^2 | b^1) P(d^2 | c^2) \\ + & P(a^1) P(b^2 | a^1) P(c^2 | b^2) P(d^2 | c^2) \\ + & P(a^2) P(b^2 | a^2) P(c^2 | b^2) P(d^2 | c^2) \end{aligned}$$

- There is structure to the summation, e.g., repeated $P(a^1)P(b^1 | a^1) + P(a^2)P(b^1 | a^2)$
- Let's modify the computation to first compute

$$P(a^1)P(b^1 | a^1) + P(a^2)P(b^1 | a^2)$$

Let's be a bit more explicit...

- Let's modify the computation to first compute

$$P(a^1)P(b^1 | a^1) + P(a^2)P(b^1 | a^2)$$

and

$$P(a^1)P(b^2 | a^1) + P(a^2)P(b^2 | a^2)$$

- Then, we get

$$\begin{array}{lll} & (P(a^1)P(b^1 | a^1) + P(a^2)P(b^1 | a^2)) & P(c^1 | b^1) & P(d^1 | c^1) \\ + & (P(a^1)P(b^2 | a^1) + P(a^2)P(b^2 | a^2)) & P(c^1 | b^2) & P(d^1 | c^1) \\ + & (P(a^1)P(b^1 | a^1) + P(a^2)P(b^1 | a^2)) & P(c^2 | b^1) & P(d^1 | c^2) \\ + & (P(a^1)P(b^2 | a^1) + P(a^2)P(b^2 | a^2)) & P(c^2 | b^2) & P(d^1 | c^2) \\ \\ & (P(a^1)P(b^1 | a^1) + P(a^2)P(b^1 | a^2)) & P(c^1 | b^1) & P(d^2 | c^1) \\ + & (P(a^1)P(b^2 | a^1) + P(a^2)P(b^2 | a^2)) & P(c^1 | b^2) & P(d^2 | c^1) \\ + & (P(a^1)P(b^1 | a^1) + P(a^2)P(b^1 | a^2)) & P(c^2 | b^1) & P(d^2 | c^2) \\ + & (P(a^1)P(b^2 | a^1) + P(a^2)P(b^2 | a^2)) & P(c^2 | b^2) & P(d^2 | c^2) \end{array}$$

- We define $\tau_1 : \text{Val}(B) \rightarrow R$, $\tau_1(b^i) = P(a^1)P(b^i | a^1) + P(a^2)P(b^i | a^2)$

Let's be a bit more explicit...

- We now have

$$\begin{aligned} & \tau_1(b^1) P(c^1 | b^1) P(d^1 | c^1) \\ + & \tau_1(b^2) P(c^1 | b^2) P(d^1 | c^1) \\ + & \tau_1(b^1) P(c^2 | b^1) P(d^1 | c^2) \\ + & \tau_1(b^2) P(c^2 | b^2) P(d^1 | c^2) \end{aligned}$$

$$\begin{aligned} & \tau_1(b^1) P(c^1 | b^1) P(d^2 | c^1) \\ + & \tau_1(b^2) P(c^1 | b^2) P(d^2 | c^1) \\ + & \tau_1(b^1) P(c^2 | b^1) P(d^2 | c^2) \\ + & \tau_1(b^2) P(c^2 | b^2) P(d^2 | c^2) \end{aligned}$$

- We can once more reverse the order of the product and the sum and get

$$\begin{aligned} & (\tau_1(b^1)P(c^1 | b^1) + \tau_1(b^2)P(c^1 | b^2)) P(d^1 | c^1) \\ + & (\tau_1(b^1)P(c^2 | b^1) + \tau_1(b^2)P(c^2 | b^2)) P(d^1 | c^2) \end{aligned}$$

$$\begin{aligned} & (\tau_1(b^1)P(c^1 | b^1) + \tau_1(b^2)P(c^1 | b^2)) P(d^2 | c^1) \\ + & (\tau_1(b^1)P(c^2 | b^1) + \tau_1(b^2)P(c^2 | b^2)) P(d^2 | c^2) \end{aligned}$$

- There are still other repeated computations!

Let's be a bit more explicit...

- We define $\tau_2 : \text{Val}(C) \rightarrow R$, with

$$\tau_2(c^1) = \tau_1(b^1)P(c^1 | b^1) + \tau_1(b^2)P(c^1 | b^2)$$

$$\tau_2(c^2) = \tau_1(b^1)P(c^2 | b^1) + \tau_1(b^2)P(c^2 | b^2)$$

- Now we can compute the marginal $p(D)$ as

$$+ \begin{array}{l} \tau_2(c^1) P(d^1 | c^1) \\ \tau_2(c^2) P(d^1 | c^2) \end{array}$$

$$+ \begin{array}{l} \tau_2(c^1) P(d^2 | c^1) \\ \tau_2(c^2) P(d^2 | c^2) \end{array}$$

What Did We Do?

$$P(D) = \sum_C \sum_B \sum_A P(A)P(B | A)P(C | B)P(D | C)$$

Push in the summation of A

$$P(D) = \sum_C \sum_B P(C | B)P(D | C) \sum_A P(A)P(B | A)$$

Push in the summation of B

$$P(D) = \sum_C P(D | C) \sum_B P(C | B) P(B)$$

Push in the summation of C

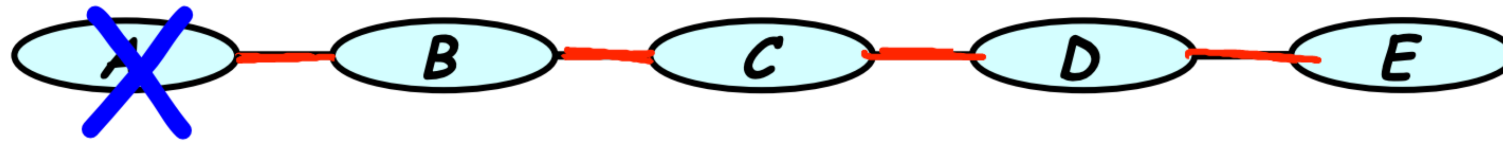
$$P(D) = \sum_C P(D | C)P(C)$$

Rule for Sum-Product VE

If $X \notin \text{Scope}[\phi_1]$, then

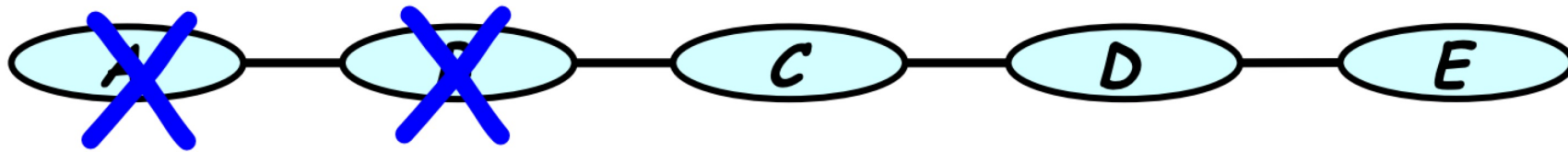
$$\sum_X (\phi_1 \cdot \phi_2) = \phi_1 \cdot \sum_X \phi_2$$

Elimination In Chains: MNs



$$\begin{aligned} P(E) &\propto \sum_D \sum_C \sum_B \sum_A \tilde{P}(A, B, C, D, E) \\ &= \sum_D \sum_C \sum_B \sum_A \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \\ &= \sum_D \sum_C \sum_B \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \sum_A \phi_1(A, B) \tau_1(B) \\ &= \sum_D \sum_C \sum_B \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \tau_1(B) \end{aligned}$$

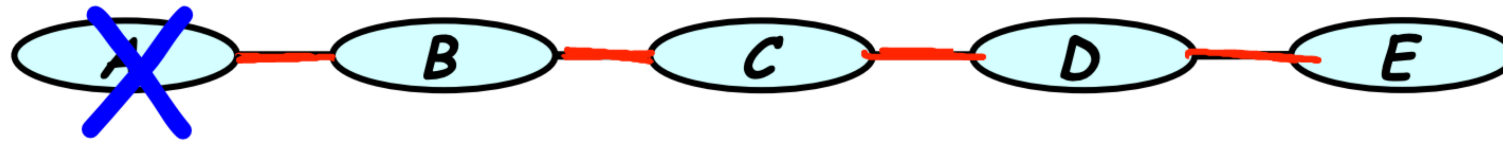
Elimination In Chains: MNs



$$\begin{aligned} P(E) &\propto \sum_D \sum_C \sum_B \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \tau_1(B) \\ &= \sum_D \sum_C \phi_3(C, D) \phi_4(D, E) \left(\sum_B \phi_2(B, C) \tau_1(B) \right) \\ &= \sum_D \sum_C \phi_3(C, D) \phi_4(D, E) \tau_2(C) \end{aligned}$$

Note: Red handwritten annotations in the original image include a bracket under the inner sum in the second line labeled $\tau_2(C)$.

Elimination In Chains: MNs

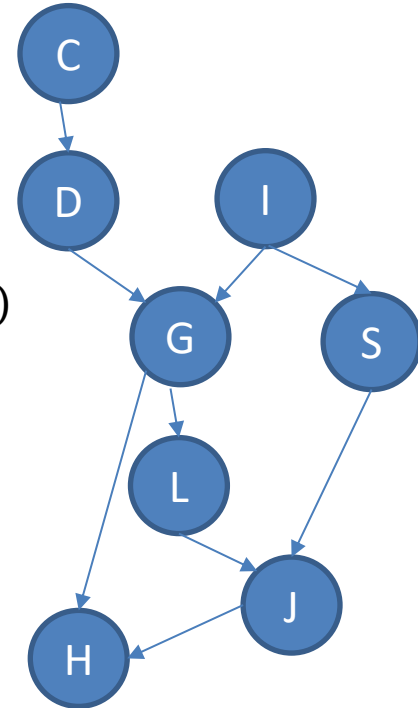


$$\begin{aligned} P(E) &\propto \sum_D \sum_C \sum_B \sum_A \tilde{P}(A, B, C, D, E) \\ &= \sum_D \sum_C \sum_B \sum_A \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \\ &= \sum_D \sum_C \sum_B \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \sum_A \phi_1(A, B) \tau_1(B) \\ &= \sum_D \sum_C \sum_B \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \tau_1(B) \end{aligned}$$

Variable Elimination

- Goal: $P(J)$
- Eliminate: C, D, I, H, G, S, L

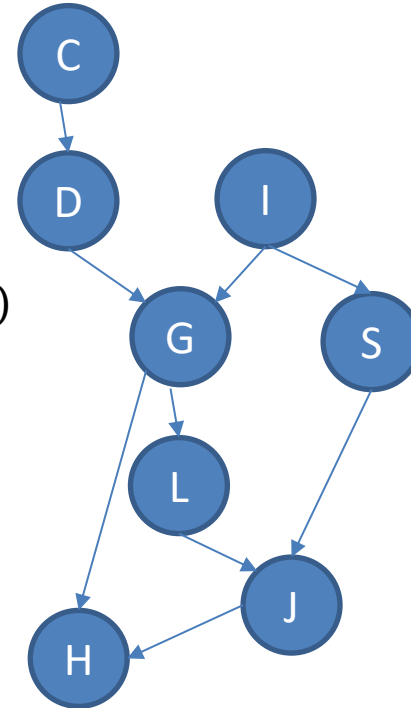
$$\sum_{C, D, I, G, S, L, H} \phi_C(C) \phi_D(C, D) \phi_I(I) \phi_G(G, I, D) \phi_S(S, I) \phi_L(L, G) \phi_J(J, L, S) \phi_H(H, G, J)$$



Variable Elimination with evidence

- Goal: $P(J, I=i, H=h)$
- Eliminate: C, D, G, S, L

$$\sum_{C, D, I, G, S, L, H} \phi_C(C) \phi_D(C, D) \phi_I(I) \phi_G(G, I, D) \phi_S(S, I) \phi_L(L, G) \phi_J(J, L, S) \phi_H(H, G, J)$$



Variable Elimination in MNs

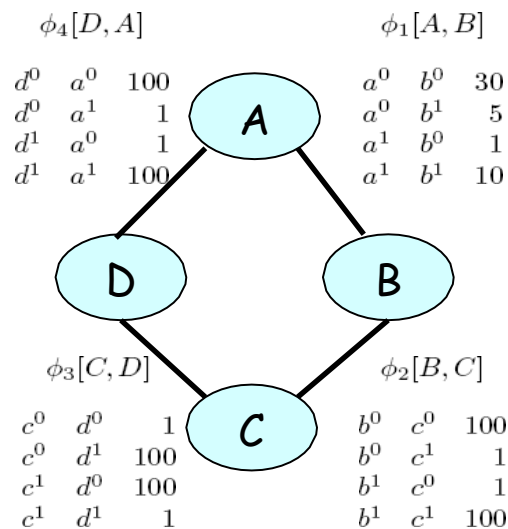
- Goal: $P(D)$
- Eliminate: A, B, C

$$\sum_{A, B, C} \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(A, D)$$

$$\sum_{B, C} \phi_2(B, C) \phi_3(C, D) \sum_A \phi_1(A, B) \phi_4(A, D)$$

$$\sum_{B, C} \phi_2(B, C) \phi_3(C, D) \tau_1(B, D)$$

At the end of elimination get $\tau_3(D) \propto P(D)$
 re-normalize $= \tilde{P}(D)$



- Reduce all factors by evidence
 - Get a set of factors Φ
- For each non-query variable Z
 - Eliminate-Var Z from Φ :

$$\Phi' = \{\phi_i \in \Phi : Z \in \text{Scope}[\phi_i]\}$$

$$\psi = \prod_{\phi_i \in \Phi'} \phi_i$$

$$\tau = \sum_Z \psi$$

$$\Phi := \Phi - \Phi' \cup \{\tau\}$$

- Multiply all remaining factors
- Renormalize to get distribution

Summary

- Simple algorithm
- Works for both BNs and MNs
- Factor product and summation steps can be done in any order, subject to:
 - when Z is eliminated, all factors involving Z have been multiplied in