

Lecture 14: March 9, 2015
cs 573: Probabilistic Reasoning
Professor Nevatia
Spring 2015

Review

- Assignment # 4 due today
- Assignment #5 to be posted Mar 11, due Mar 25
- Exam 1 graded by March 11
 - How to provide access?
- Exam 2, April 29, class period; NOT cumulative
- Previous lecture:
 - Gaussian networks
 - Representation of Gaussian distributions
 - Gaussian BNs and MRFs
 - Inference algorithms
- Today's objective
 - MAP Inference

Canonical Form

- Co-variance Form

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

- Canonical form:

$$\begin{aligned} K &= \Sigma^{-1} \\ h &= \Sigma^{-1} \boldsymbol{\mu} \\ g &= -\frac{1}{2} \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu} - \log \left((2\pi)^{n/2} |\Sigma|^{1/2} \right) \end{aligned}$$

$$\mathcal{C}(\mathbf{X}; K, \mathbf{h}, g) = \exp \left(-\frac{1}{2} \mathbf{X}^T K \mathbf{X} + \mathbf{h}^T \mathbf{X} + g \right) \quad \text{Eq 14.1}$$

Marginalize/Conditionalize a Gaussian Distribution

$$p(\mathbf{X}, \mathbf{Y}) = \mathcal{N} \left(\begin{pmatrix} \mu_{\mathbf{X}} \\ \mu_{\mathbf{Y}} \end{pmatrix}; \begin{bmatrix} \Sigma_{\mathbf{X}\mathbf{X}} & \Sigma_{\mathbf{X}\mathbf{Y}} \\ \Sigma_{\mathbf{Y}\mathbf{X}} & \Sigma_{\mathbf{Y}\mathbf{Y}} \end{bmatrix} \right)$$

- Marginal over \mathbf{Y} (sum out \mathbf{X}) is given by $\mathcal{N}(\mu_{\mathbf{Y}}, \Sigma_{\mathbf{Y}\mathbf{Y}})$
- Conditioning is easier in the canonical form:

$$\begin{aligned} K' &= K_{\mathbf{X}\mathbf{X}} \\ h' &= h_{\mathbf{X}} - K_{\mathbf{X}\mathbf{Y}}\mathbf{y} \\ g' &= g + h_{\mathbf{Y}}^T\mathbf{y} - \frac{1}{2}\mathbf{y}^T K_{\mathbf{Y}\mathbf{Y}}\mathbf{y} \end{aligned}$$

Gaussian Bayesian Networks (GBN)

- In a GBN, all variables are continuous; all CPDs are linear Gaussian

Let Y be a continuous variable with continuous parents X_1, \dots, X_k . We say that Y has a linear Gaussian model if there are parameters β_0, \dots, β_k and σ^2 such that

$$p(Y \mid x_1, \dots, x_k) = \mathcal{N}(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k; \sigma^2).$$

In vector notation,

$$p(Y \mid \mathbf{x}) = \mathcal{N}(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}; \sigma^2).$$

– Thm 7.3

Given $p(Y \mid \mathbf{x}) = \mathcal{N}(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}; \sigma^2)$; X_1, \dots, X_k distributed as $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

We can show that

$$\begin{aligned}\mu_Y &= \beta_0 + \boldsymbol{\beta}^T \boldsymbol{\mu} \\ \sigma_Y^2 &= \sigma^2 + \boldsymbol{\beta}^T \boldsymbol{\Sigma} \boldsymbol{\beta}\end{aligned}$$

$$\text{Cov}[X_i; Y] = \sum_{j=1}^k \beta_j \Sigma_{i,j}.$$

Follows that in a GBN, joint distribution is Gaussian

Distribution to GBN

- Previous slide shows how given a GBN (network parameters), we can get joint distribution parameters.
- Reverse: given the joint distribution, recover the linear model (thm 7.4)

Given:

$$p(\mathbf{X}, \mathbf{Y}) = \mathcal{N} \left(\begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}; \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix} \right)$$

Derive:

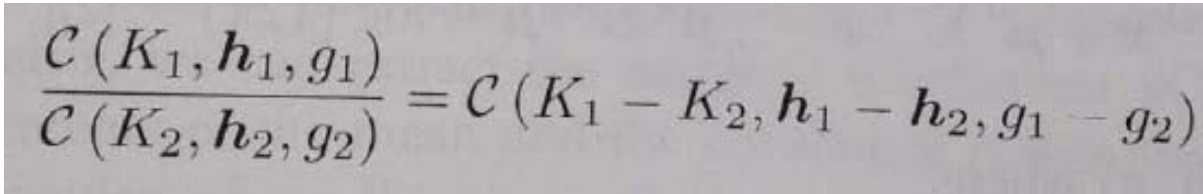
$$p(Y | \mathbf{X}) = \mathcal{N} \left(\beta_0 + \boldsymbol{\beta}^T \mathbf{X}; \sigma^2 \right)$$

$$\begin{aligned} \beta_0 &= \mu_Y - \Sigma_{YX} \Sigma_{XX}^{-1} \mu_X \\ \boldsymbol{\beta} &= \Sigma_{XX}^{-1} \Sigma_{YX} \\ \sigma^2 &= \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY} \end{aligned}$$

Distribution to GBN

- Given a distribution over n variables, we can construct a GBN (BN with linear Gaussian Models) that is an I-map of the distribution
- A key difference with discrete case: number of parameters in GBN is not necessarily smaller than in the joint distribution itself as the joint distribution is compact by itself.
- We can also go from distributions to Markov networks
 - J_{ij} help define pairwise (log) potentials
 - However, additional complexity in case of MN; not every set of potentials induces a valid Gaussian distributions
 - Sufficient but not necessary condition is that each edge potential be normalizable (corresponding information matrix is positive definite)
 - We skip other details of Gaussian Markov Random Fields (sec 7.3)

Inference in Networks with Continuous Variables

- Essentially, all the algorithms for discrete case apply
- Sum-Product algorithm
 - Steps consist of multiplying factors (product) and marginalizing over some variables (sum) and passing messages to other nodes
 - Iterate until convergence (two passes suffice for trees)
 - We already know how to marginalize Gaussians
 - Now consider product and division
- Product of two canonical forms over the same set of variables:
$$C(\mathbf{X}, K_1, \mathbf{h}_1, g_1) \cdot C(\mathbf{X}, K_2, \mathbf{h}_2, g_2) = C(\mathbf{X}, K_1 + K_2, \mathbf{h}_1 + \mathbf{h}_2, g_1 + g_2)$$
- Division:

$$\frac{C(K_1, \mathbf{h}_1, g_1)}{C(K_2, \mathbf{h}_2, g_2)} = C(K_1 - K_2, \mathbf{h}_1 - \mathbf{h}_2, g_1 - g_2)$$
- If scope is different, expand by adding variables and zero entries

SP and BP Algorithms

- Sum-Product algorithm
 - As before; need to show that resulting factors maintain K matrices to be positive definite so integration is not infinite
 - Shown in Proposition 14.1
- Gaussian Belief Propagation
 - Similar steps as in discrete case but the observation that the potentials are quadratic simplifies the equations (eqs 14.7 to 14.9)
 - **Interesting property:** If BP converges, resulting beliefs encode the correct means but estimated variances are underestimates (overconfident estimates)
 - Also convergence guaranteed if pairwise normalizability condition holds.
- We skip other details

Hybrid Networks

- Mix of discrete and continuous variables
- In general, even representation is difficult
 - Restrict to continuous linear Gaussian (CLG) networks
 - No discrete variables with continuous parents
- Even for CLG networks, exact inference is difficult
 - Even though the initial factors are all Gaussians, intermediate factors are not necessarily Gaussians (but mixtures of Gaussians)
 - Number of Gaussian to mix grows exponentially in the size of the network
- Approximation techniques need to be used; we omit details

MAP Queries

- Often, we are interested in joint assignment of variables that have the highest probability, rather than the individual distributions
 - Decoding a signal
 - Word sequence from a speech signal...
 - Classification/diagnosis tasks: debugging a piece of equipment, object recognition.
- MAP (maximum *a posteriori*) Query
 - Maximum a posteriori probability assignment or most probable explanation (MPE)
 - $\text{MAP}(\mathbf{W} | \mathbf{e}) = \arg \max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e}), \mathbf{W} = \mathbf{X} - \mathbf{E}$
 - (Notation: $\arg \max_x f(x)$ = value of x for which $f(x)$ is maximal)
 - Note that maximal joint assignment is not same as maxima of individual assignments, example on next slide.

MAP Example

Consider a two node chain $A \rightarrow B$ where A and B are both binary-valued. Assume that:

$$\begin{array}{cc} a^0 & a^1 \\ \hline 0.4 & 0.6 \end{array} \quad \begin{array}{c|cc} A & b^0 & b^1 \\ \hline a^0 & 0.1 & 0.9 \\ a^1 & 0.5 & 0.5 \end{array} \quad \text{CPD tables} \quad (2.13)$$

$$\text{MAP}(A) = a^1$$

However, $\text{MAP}(A, B) = (a^0, b^1)$:

$$\arg \max_{a,b} P(a, b) \neq (\arg \max_a P(a), \arg \max_b P(b))$$

Overview

- Task is to compute:

$$\xi^{map} = \arg \max_{\xi} P_{\Phi}(\xi) = \arg \max_{\xi} \frac{1}{Z} \tilde{P}_{\Phi}(\xi) = \arg \max_{\xi} \tilde{P}_{\Phi}(\xi).$$

- Note that we don't need to compute partition function if we just want variable assignments but do need it for computing the actual max probability
- Can adapt the *sum-product* algorithm to compute MAP by changing the sum operation to a max operation- *max product* algorithm
- If we work with log of probabilities, we get a *max sum* algorithm
- If we work with energy factors, we get a *min sum* algorithm

Marginal MAP Query

- We may only care about the assignments of a subset of the variables
 - Disease diagnosis: full MAP query would compute joint distribution of diseases, symptoms, and test outcomes; we may only be interested in disease probabilities
 - $\text{MAP}(\mathbf{Y} | \mathbf{e}) = \arg \max_{\mathbf{y}} P(\mathbf{y} | \mathbf{e}) = \arg \max_{\mathbf{y}} P(\mathbf{y}, \mathbf{e})$
 - Let $\mathbf{Z} = \mathbf{X} - \mathbf{Y} - \mathbf{E}$
 - $\text{MAP}(\mathbf{Y} | \mathbf{e}) = \arg \max_{\mathbf{y}} \{ \sum_{\mathbf{z}} P(\mathbf{Y}, \mathbf{Z} | \mathbf{e}) \}$
 - Note: marginal MAP query can not be computed directly from a MAP query (can not reverse summation and maximization operations)
- Complexity of computing MAP query is NP-complete
- Complexity of marginal MAP is even higher (NP^{PP})
- MAP Problem is NP-hard even for a polytree

VE for MAP

- Example 13.1. Simple network $A \rightarrow B$, no evidence variables
- $\max_{a,b} P(a,b) = \max_{a,b} P(a) P(b|a)$
 $= \max_a \max_b P(a) P(b|a)$
 $= \max_a P(a) \max_b P(b|a)$
- Let $\phi(a)$ denote $\max_b P(b|a)$

a^0	a^1
0.4	0.6

A	b^0	b^1
a^0	0.1	0.9
a^1	0.55	0.45

$$\phi(a^1) = \max_b P(b | a^1) = 0.55 \text{ and } \phi(a^0) = \max_b P(b | a^0) = 0.9.$$

$$\max_a P(a) \phi(a) = \max [0.4 \cdot 0.9, 0.6 \cdot 0.55] = 0.36.$$

- How to find the MAP assignments?
 - When we computed $\phi(A)$, we could not know what value of B would lead to MAP solution
 - At the last step, we know that $A = a^1$ is selected, so $B = b^0$

Factor Maximization

- Definition 13.2

Let X be a set of variables, and $Y \notin X$ a variable. Let $\phi(X, Y)$ be a factor. We define the factor maximization of Y in ϕ to be factor ψ over X such that:

$$\psi(X) = \max_Y \phi(X, Y).$$



Max-marginalization over B

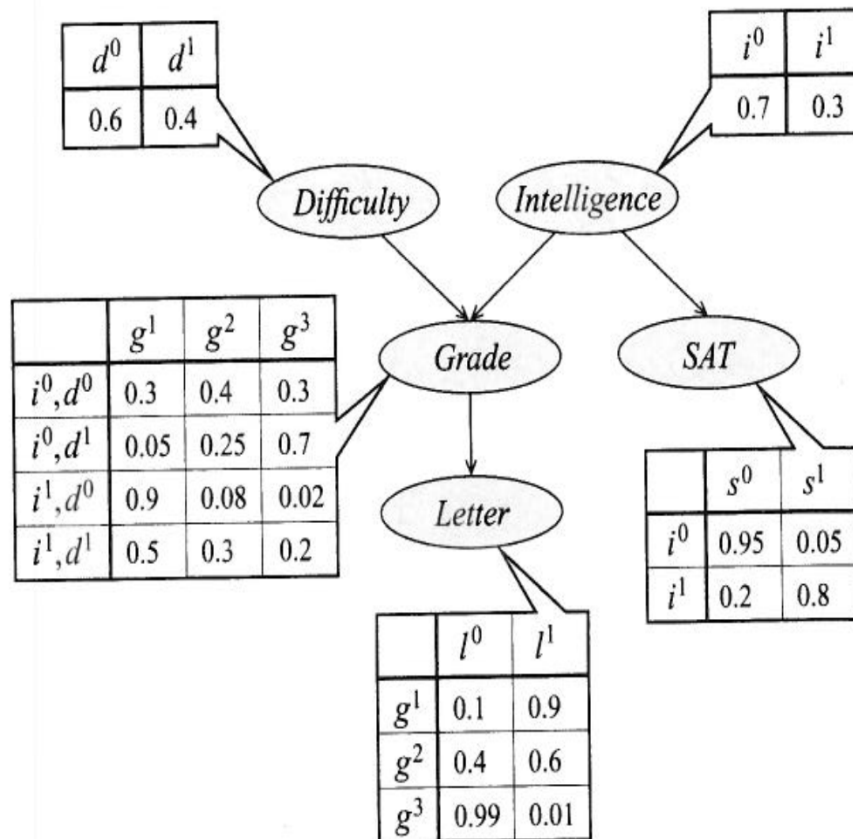
a^1	b^1	c^1	0.25
a^1	b^1	c^2	0.35
a^1	b^2	c^1	0.08
a^1	b^2	c^2	0.16
a^2	b^1	c^1	0.05
a^2	b^1	c^2	0.07
a^2	b^2	c^1	0
a^2	b^2	c^2	0
a^3	b^1	c^1	0.15
a^3	b^1	c^2	0.21
a^3	b^2	c^1	0.09
a^3	b^2	c^2	0.18

a^1	c^1	0.25
a^1	c^2	0.35
a^2	c^1	0.05
a^2	c^2	0.07
a^3	c^1	0.15
a^3	c^2	0.21

Rules for Factor Maximization

- $\sum_x (\phi_1 \cdot \phi_2) = \phi_1 \cdot \sum_x \phi_2$, if X is not in domain of ϕ_1
- $\max_x (\phi_1 \cdot \phi_2) = \phi_1 \cdot \max_x \phi_2$, X is not in domain of ϕ_1
- $\max_x (\phi_1 + \phi_2) = \phi_1 + \max_x \phi_2$, X is not in domain of ϕ_1

Max Product VE: Example



Eliminate variables in
order S,I,D,L,G

Max instead of sum at
each step

$$\tau_1(I) = \max_s \phi_s(I, s)$$

$$\tau_1(i^0) = .95 \quad \tau_1(i^1) = .8$$

Note different values of S
for the two values above

See table on next slide for other
variables.

Example

Step	Variable eliminated	Factors used	Intermediate factor	New factor
1	S	$\phi_S(I, S)$	$\psi_1(I, S)$	$\tau_1(I)$
2	I	$\phi_I(I), \phi_G(G, I, D), \tau_1(I)$	$\psi_2(G, I, D)$	$\tau_2(G, D)$
3	D	$\phi_D(D), \tau_2(G, D)$	$\psi_3(G, D)$	$\tau_3(G)$
4	L	$\phi_L(L, G)$	$\psi_4(L, G)$	$\tau_4(G)$
5	G	$\tau_4(G), \tau_3(G)$	$\psi_5(G)$	$\tau_5(\emptyset)$

Note that this procedure does not provide the values of eliminated variables that give rise to the maximum

Consider elimination of S , we don't consider what value of S gave us $\tau(I)$

TraceBack

- At the last elimination, we do have access to the argument that maximizes the corresponding factor (ψ_5), this gives a value for g^*
- Given g^* we can compute $l^* = \operatorname{argmax}_l \psi_4(g^*, l)$
- Continue, one variable at a time
- - $d^* = \operatorname{argmax}_d \psi_3(g^*, d)$
 - $i^* = \operatorname{argmax}_i \psi_2(g^*, i, d^*)$
 - $s^* = \operatorname{argmax}_s \psi_1(i^*, s)$
- Pseudo-code on following slides
 - Identical to sum-product except that “sum” is replaced by “max” and *traceback* function is added

Procedure Max-Product-VE (

Φ , // Set of factors over X

\prec // Ordering on X

)

1 Let X_1, \dots, X_k be an ordering of X such that

2 $X_i \prec X_j$ iff $i < j$

3 **for** $i = 1, \dots, k$

4 $(\Phi, \phi_{X_i}) \leftarrow \text{Max-Product-Eliminate-Var}(\Phi, X_i)$

5 $x^* \leftarrow \text{Traceback-MAP}(\{\phi_{X_i} : i = 1, \dots, k\})$

6 **return** x^*, Φ // Φ contains the probability of the MAP

Procedure Max-Product-Eliminate-Var (

Φ , // Set of factors

Z // Variable to be eliminated

)

1 $\Phi' \leftarrow \{\phi \in \Phi : Z \in \text{Scope}[\phi]\}$

2 $\Phi'' \leftarrow \Phi - \Phi'$

3 $\psi \leftarrow \prod_{\phi \in \Phi'} \phi$

4 $\tau \leftarrow \max_Z \psi$

5 **return** $(\Phi'' \cup \{\tau\}, \psi)$

```

Procedure Traceback-MAP (
     $\{\phi_{X_i} : i = 1, \dots, k\}$ 
)
1   for  $i = k, \dots, 1$ 
2        $u_i \leftarrow (x_{i+1}^*, \dots, x_k^*) \langle \text{Scope}[\phi_{X_i}] - \{X_i\} \rangle$ 
3       // The maximizing assignment to the variables eliminated after
         $X_i$ 
4        $x_i^* \leftarrow \arg \max_{x_i} \phi_{X_i}(x_i, u_i)$ 
5       //  $x_i^*$  is chosen so as to maximize the corresponding entry in
        the factor, relative to the previous choices  $u_i$ 
6   return  $x^*$ 

```

VE with a Marginal MAP

- We want to find max over some set \mathbf{Y} but don't care about the rest (set \mathbf{W})
- $\mathbf{y}^{\text{m-map}} = \arg \max_{\mathbf{y}} P_{\phi}(\mathbf{y}) = \arg \max_{\mathbf{y}} \sum_{\mathbf{w}} P_{\phi}(\mathbf{y}, \mathbf{w})$
- What we want to compute is a max sum product

$$\max_{\mathbf{Y}} \sum_{\mathbf{W}} \prod_{\phi \in \Phi} \phi.$$

- More difficult: as we must compute both sum and max functions. Can be shown that the sum computations must be done first and max later: the two are not commutative
 - Less freedom in choosing the elimination order
- Example on next slide

$$\max_{S,L} \sum_{G,I,D} P(I, D, G, S, L).$$

$$\psi_1(I, G, D) = \phi_D(D) \cdot \phi_G(G, I, D)$$

$$\tau_1(I, G) = \sum_D \psi_1(I, G, D)$$

$$\psi_2(S, G, I) = \phi_I(I) \cdot \phi_S(S, I) \cdot \tau_1(I, G)$$

$$\tau_2(S, G) = \sum_I \psi_2(S, G, I)$$

$$\psi_3(S, G, L) = \tau_2(S, G) \cdot \phi_L(L, G)$$

$$\tau_3(S, L) = \sum_G \psi_3(S, G, L)$$

$$\psi_4(S, L) = \tau_3(S, L)$$

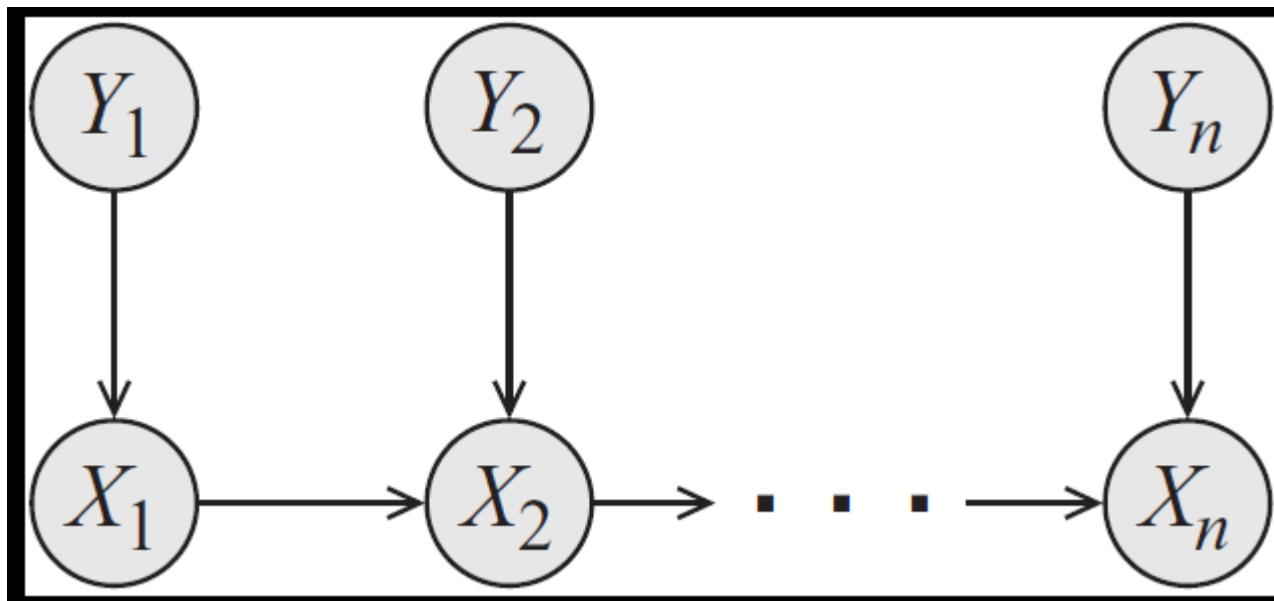
$$\tau_4(L) = \max_S \psi_4(S, L)$$

$$\psi_5(L) = \tau_4(L)$$

$$\tau_5(\emptyset) = \max_L \psi_5(L).$$

Computationally Complex Example

- Goal is to compute $\mathbf{y}^{\text{m-map}} = \arg \max_{\mathbf{Y}} \sum_{\mathbf{X}} P(\mathbf{Y}, \mathbf{X})$



- First eliminate X_i variables in order
- Resulting factor will contain all of the Y_i variables, so exponential in n
- Find max over Y_i variables

Max-Product in Clique Trees

- Compute max-marginals over cliques instead of sum-marginals
- Advantage of this algorithm over VE is less clear as we are not seeking distributions over all variables (just over one assignment)
- May show other assignments that have nearly the same probability as the max assignment
- Method generalizes to loopy BP calculations

Algorithm 13.2 Max-product message computation for MAP

```
Procedure Max-Message (  
    i,    // sending clique  
    j    // receiving clique  
)  
1    $\psi(C_i) \leftarrow \psi_i \cdot \prod_{k \in (\text{Nb}_i - \{j\})} \delta_{k \rightarrow i}$   
2    $\tau(S_{i,j}) \leftarrow \max_{C_i - S_{i,j}} \psi(C_i)$   
3   return  $\tau(S_{i,j})$ 
```

Definition: Max-marginal of a Function

- Max-marginal of a function f , relative to a set of variables \mathbf{Y} is said to be:

$$\text{MaxMarg}_f(\mathbf{y}) = \max_{\xi(\mathbf{Y})=\mathbf{y}} f(\xi), \quad \text{for any assignment } \mathbf{y} \in \text{Val}(\mathbf{Y})$$

- Consider $\text{MaxMarg}_{P \sim \Phi}(\mathbf{Y})$: not a single number but a factor
 - For each assignment of values \mathbf{y} to \mathbf{Y} , the MaxMarg function corresponds to max of $P \sim \Phi(\mathbf{y})$; note that P may contain other variables whose values are being chosen to maximize the expression

Max-marginal

- After convergence, in a clique tree, beliefs in each clique will be max-marginals

$$\beta_i(\mathbf{c}_i) = \text{MaxMarg}_{\tilde{P}_\Phi}(\mathbf{c}_i).$$

- For each assignment \mathbf{c}_i to variables in the clique C_i , belief is the unnormalized max of all assignments elsewhere in the graph, consistent with the clique assignments, \mathbf{c}_i .
- Max-marginals agree over the variables in the sepset, as is the case for the product sum algorithm, i.e. clique tree is *max-calibrated*
- Surprisingly, the resulting beliefs still are a reparameterization of the probability function (multiply all and divide by product of beliefs on sepset variables gives the original, unnormalized distribution)

Decoding the Max-Marginals

- More complex than in the VE algorithm
- Find assignments, ξ^* that are locally optimal (for one clique)
- Can be shown that set of such assignments over all the cliques defines the global MAP
 - Theorem 13.6, Lemma 13.1; rather difficult to follow
- In case of multiple assignments having the same value, choose one and *extend* it
 - Find a locally optimal assignment for one clique, for the neighboring clique, use this assignment for the common variables and iterate

Max Product BP in Loopy Graphs

- Similar to the Sum Product BP algorithm but Sum is replaced by Max
- Convergence is not guaranteed
- Max-marginals at convergence are not necessarily the correct max-marginals of the distribution: call them *pseudo max-marginals*
- Recovering assignments is more complex as the local optimality property may not be satisfied (i.e. globally optimal assignment may not be locally optimal)
 - A globally consistent solution may not exist
 - Can not just extend the locally optimal assignments
 - Search to find an assignment that satisfies all local optimal conditions
 - Equivalent to a constraint satisfaction problem
 - Fallback: take assignments to variables that are unambiguous, do exact inference on rest if needed

Next Class

- Read section 11.5.1 of the KF book