

Lecture 8

Inference of Probabilistic Graphical Models

Instructor: Shibo Li

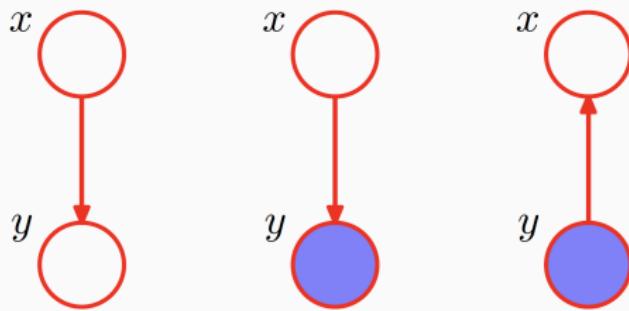
shiboli@cs.fsu.edu



Department of Computer Science
Florida State University

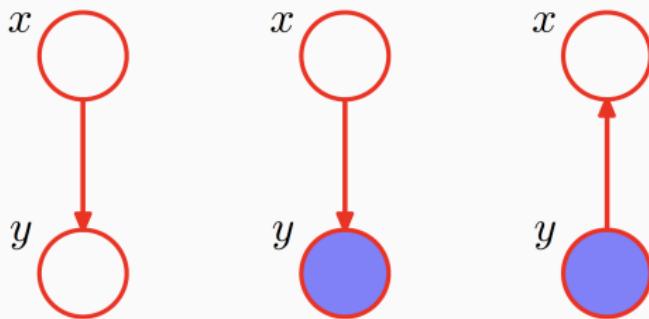
- Bayesian networks
- Markov random fields
- Inference

- Compute the posterior of one or more subsets of nodes given the observed nodes



$$p(y) = \sum_{x'} p(y|x')p(x') \quad p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

- Compute the posterior of one or more subsets of nodes given the observed nodes



$$p(y) = \sum_{x'} p(y|x')p(x')$$

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

Key: compute the marginal of one or a subset of nodes!

Let's Start with a Chain



Let's Start with a Chain



$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$



$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

Each node takes K states and so each potential is a $K \times K$ table

Let's Start with a Chain



$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

Each node takes K states and so each potential is a $K \times K$ table

Let us consider to infer the marginal of a node x_n

Let's Start with a Chain



$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

Each node takes K states and so each potential is a $K \times K$ table

Let us consider to infer the marginal of a node x_n

$$p(x_n) = \sum_{x_1} \cdots \sum_{x_{n-1}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\mathbf{x})$$



$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

Each node takes K states and so each potential is a $K \times K$ table

Let us consider to infer the marginal of a node x_n

$$p(x_n) = \sum_{x_1} \cdots \sum_{x_{n-1}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\mathbf{x}) \quad \text{How much cost?}$$

Let's Start with a Chain



$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

Each node takes K states and so each potential is a $K \times K$ table

Let us consider to infer the marginal of a node x_n

$$p(x_n) = \sum_{x_1} \cdots \sum_{x_{n-1}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\mathbf{x})$$

How much cost?
 $O(K^*K^{N-1})$

$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

$$p(x_n) = \sum_{x_1} \cdots \sum_{x_{n-1}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\mathbf{x})$$

How to reduce the cost?

Key observations: many terms are repeated in the calculation, so we can use the distributive law to save products and sums

$$a_1 b_1 + a_1 b_2 + a_2 b_1 + a_2 b_2 = a_1(b_1 + b_2) + a_2(b_1 + b_2) = (a_1 + a_2)(b_1 + b_2)$$

Let's Start with a Chain

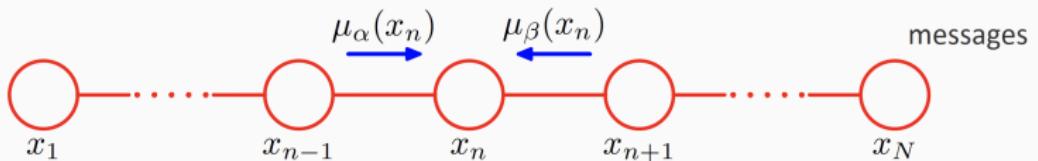
$$p(x_n) = \frac{1}{Z} \underbrace{\left[\sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, \boxed{x_n}) \cdots \left[\sum_{x_2} \psi_{2,3}(x_2, x_3) \left[\sum_{x_1} \psi_{1,2}(x_1, x_2) \right] \right] \cdots \right]}_{\mu_\alpha(x_n)}$$
$$\underbrace{\left[\sum_{x_{n+1}} \psi_{n,n+1}(\boxed{x_n}, x_{n+1}) \cdots \left[\sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) \right] \cdots \right]}_{\mu_\beta(x_n)}$$

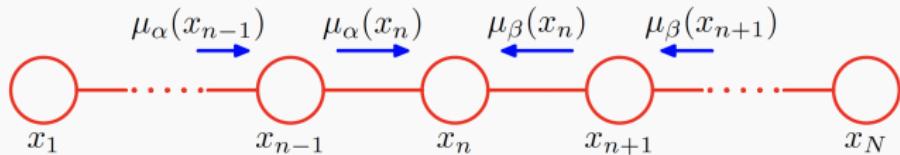
Let's Start with a Chain

$$p(x_n) = \frac{1}{Z}$$

$$\underbrace{\left[\sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, \boxed{x_n}) \cdots \left[\sum_{x_2} \psi_{2,3}(x_2, x_3) \left[\sum_{x_1} \psi_{1,2}(x_1, x_2) \right] \right] \cdots \right]}_{\mu_\alpha(x_n)}$$

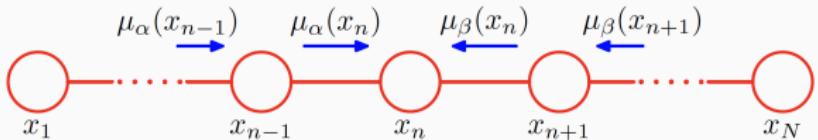
$$\underbrace{\left[\sum_{x_{n+1}} \psi_{n,n+1}(\boxed{x_n}, x_{n+1}) \cdots \left[\sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) \right] \cdots \right]}_{\mu_\beta(x_n)}$$





$$\begin{aligned}\mu_\alpha(x_n) &= \sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_n) \left[\sum_{x_{n-2}} \dots \right] \\ &= \sum \psi_{n-1,n}(x_{n-1}, x_n) \mu_\alpha(x_{n-1}).\end{aligned}$$

$$\begin{aligned}\mu_\beta(x_n) &= \sum_{x_{n+1}} \psi_{n+1,n}(x_{n+1}, x_n) \left[\sum_{x_{n+2}} \dots \right] \\ &= \sum \psi_{n+1,n}(x_{n+1}, x_n) \mu_\beta(x_{n+1}).\end{aligned}$$



$$\mu_\alpha(x_2) = \sum_{x_1} \psi_{1,2}(x_1, x_2) \quad \mu_\beta(x_{N-1}) = \sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N)$$

$$p(x_n) = \frac{1}{Z} \mu_\alpha(x_n) \mu_\beta(x_n)$$

Question: What is Z?

$$Z = \sum_{x_n} \mu_\alpha(x_n) \mu_\beta(x_n)$$

- To compute local marginals:

- Compute and store all forward messages, $\mu_\alpha(x_n)$
 - Compute and store all backward messages, $\mu_\beta(x_n)$
 - Compute Z at any node x_m
 - Compute

$$p(x_n) = \frac{1}{Z} \mu_\alpha(x_n) \mu_\beta(x_n)$$

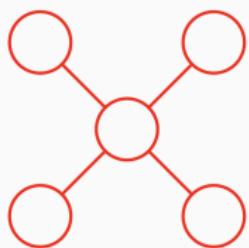
for all variables required

What is the cost? $O(NK^2)$

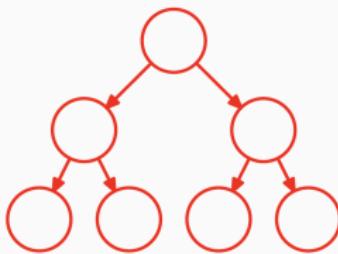
Question: how to infer the marginal of two neighboring variables?

Let us generalize the idea to trees

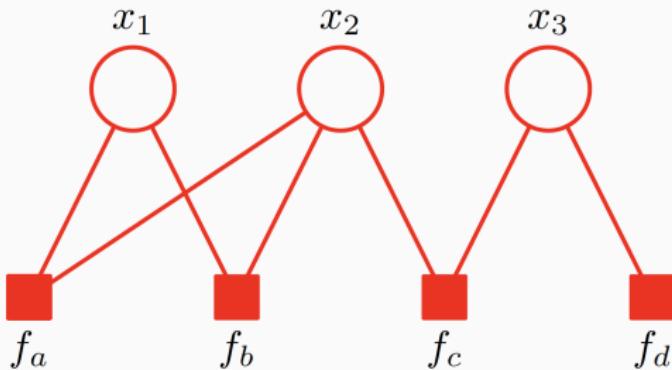
Tree-structured MRF



Tree-structured Bayesian network



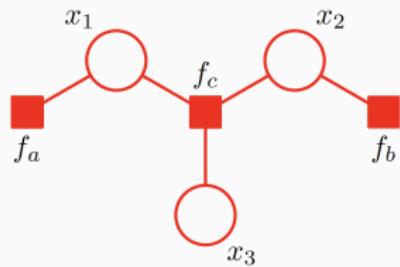
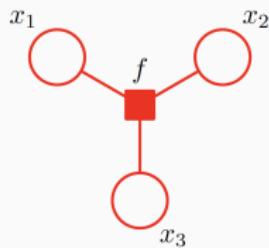
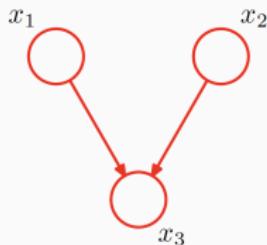
Why trees: tree structures can guarantee exact inference (we will see it later)



$$p(\mathbf{x}) = f_a(x_1, x_2) f_b(x_1, x_2) f_c(x_2, x_3) f_d(x_3)$$

$$p(\mathbf{x}) = \prod_s f_s(\mathbf{x}_s)$$

Factor graphs - multiple choices



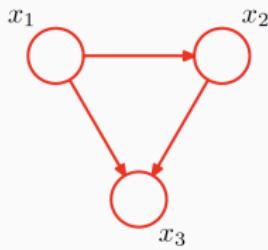
$$p(\mathbf{x}) = p(x_1)p(x_2) \\ p(x_3|x_1, x_2)$$

$$f(x_1, x_2, x_3) = \\ p(x_1)p(x_2)p(x_3|x_1, x_2)$$

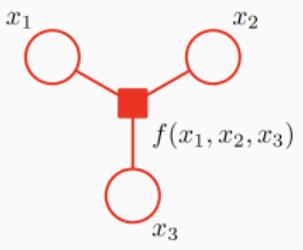
$$f_a(x_1) = p(x_1) \\ f_b(x_2) = p(x_2)$$

$$f_c(x_1, x_2, x_3) = p(x_3|x_1, x_2)$$

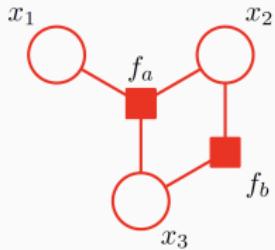
Factor graphs for undirected graphs



$$\psi(x_1, x_2, x_3)$$



$$f(x_1, x_2, x_3) = \psi(x_1, x_2, x_3)$$



$$f_a(x_1, x_2, x_3) f_b(x_2, x_3) = \psi(x_1, x_2, x_3)$$

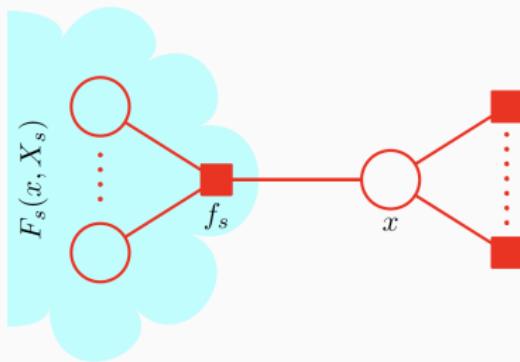
- Objective:
 - efficient, exact inference to find marginals
 - When several marginals are required, allow computations to be shared

Key idea: Distributive Law

$$a_1b_1 + a_1b_2 + a_2b_1 + a_2b_2 = a_1(b_1 + b_2) + a_2(b_1 + b_2) = (a_1 + a_2)(b_1 + b_2)$$

The Sum-Product Algorithm

Given a tree-structured graphical model



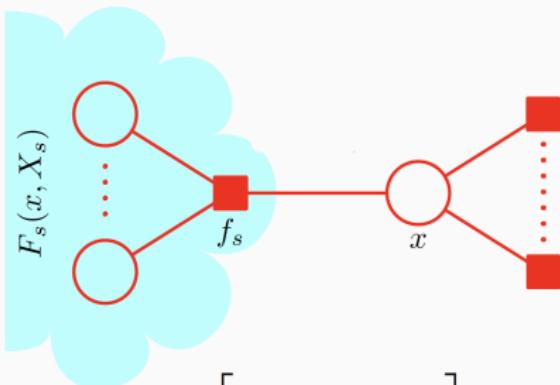
$$p(x) = \sum_{\mathbf{x} \setminus x} p(\mathbf{x})$$

$\text{ne}(x)$: factor nodes that are neighbors of x

$$p(\mathbf{x}) = \prod_{s \in \text{ne}(x)} F_s(x, X_s)$$

X_s : variables in the subtree that connect to x via the factor node f_s

The Sum-Product Algorithm

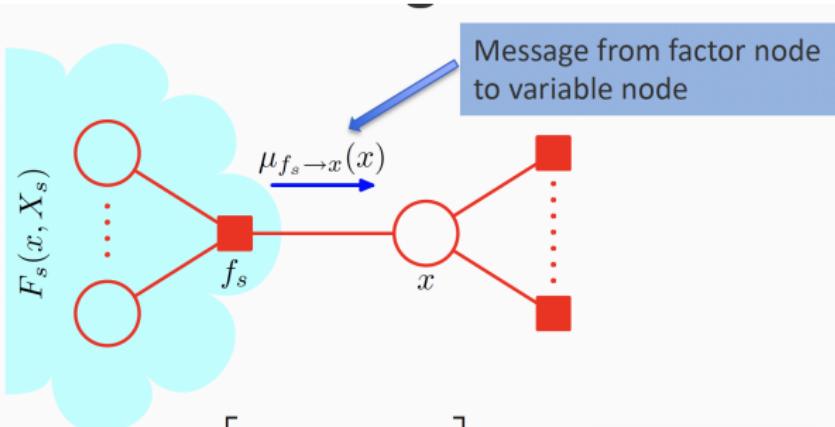


$$\begin{aligned} p(x) &= \prod_{s \in \text{ne}(x)} \left[\sum_{X_s} F_s(x, X_s) \right] \\ &= \prod_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x). \end{aligned}$$

← Why this is true?

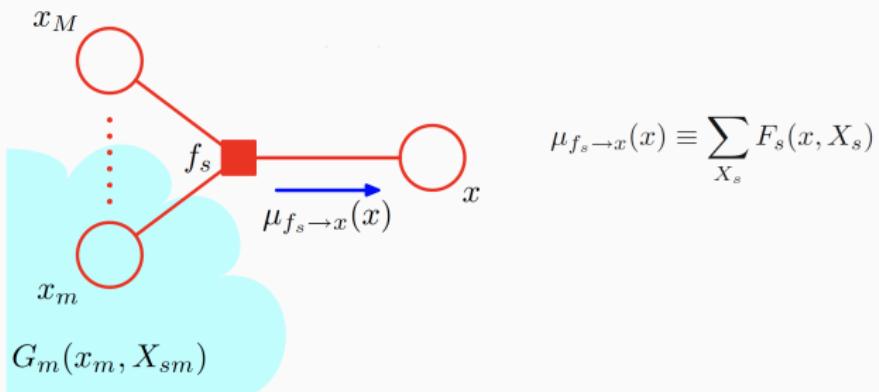
$$\mu_{f_s \rightarrow x}(x) \equiv \sum_{X_s} F_s(x, X_s)$$

The Sum-Product Algorithm



$$\begin{aligned} p(x) &= \prod_{s \in \text{ne}(x)} \left[\sum_{X_s} F_s(x, X_s) \right] && \text{Why this is true?} \\ &= \prod_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x). \end{aligned}$$
$$\mu_{f_s \rightarrow x}(x) \equiv \sum_{X_s} F_s(x, X_s)$$

The Sum-Product Algorithm

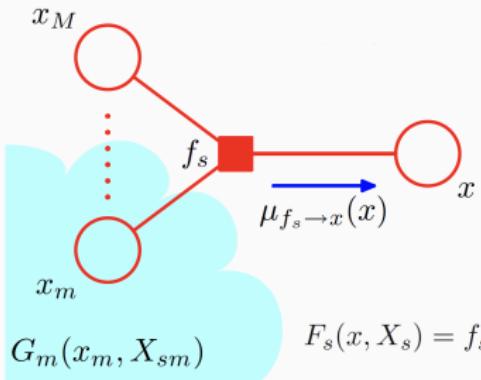


$$F_s(x, X_s) = f_s(x, x_1, \dots, x_M) G_1(x_1, X_{s1}) \dots G_M(x_M, X_{sM})$$

X_{sm} : variables in the subtree that connect to x *first* through x_m and *then* through the factor node f_s

Different X_{sm} do not overlap. Why?

The Sum-Product Algorithm



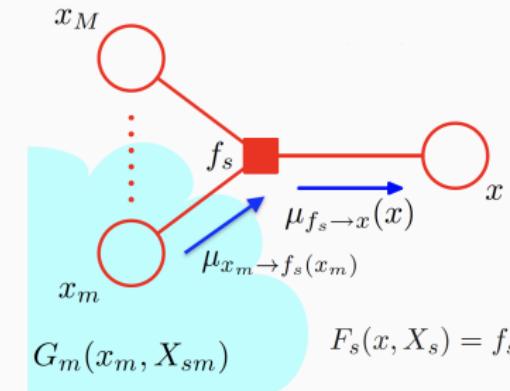
$$\mu_{f_s \rightarrow x}(x) \equiv \sum_{X_s} F_s(x, X_s)$$

$$F_s(x, X_s) = f_s(x, x_1, \dots, x_M) G_1(x_1, X_{s1}) \dots G_M(x_M, X_{sM})$$

distribute the summation

$$\begin{aligned}\mu_{f_s \rightarrow x}(x) &= \sum_{x_1} \dots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m \in \text{ne}(f_s) \setminus x} \left[\sum_{X_{xm}} G_m(x_m, X_{sm}) \right] \\ &= \sum_{x_1} \dots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m \in \text{ne}(f_s) \setminus x} \boxed{\mu_{x_m \rightarrow f_s}(x_m)}\end{aligned}$$

The Sum-Product Algorithm



$$\mu_{f_s \rightarrow x}(x) \equiv \sum_{X_s} F_s(x, X_s)$$

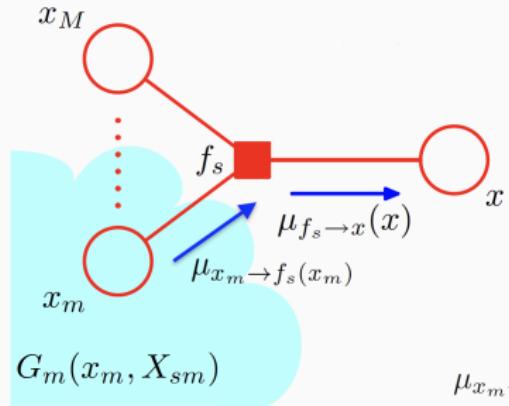
$$F_s(x, X_s) = f_s(x, x_1, \dots, x_M) G_1(x_1, X_{s1}) \dots G_M(x_M, X_{sM})$$

distribute the summation

$$\begin{aligned}\mu_{f_s \rightarrow x}(x) &= \sum_{x_1} \dots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m \in \text{ne}(f_s) \setminus x} \left[\sum_{X_{xm}} G_m(x_m, X_{sm}) \right] \\ &= \sum_{x_1} \dots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m \in \text{ne}(f_s) \setminus x} \boxed{\mu_{x_m \rightarrow f_s}(x_m)}\end{aligned}$$

Message from a
variable node to a
factor node

The Sum-Product Algorithm

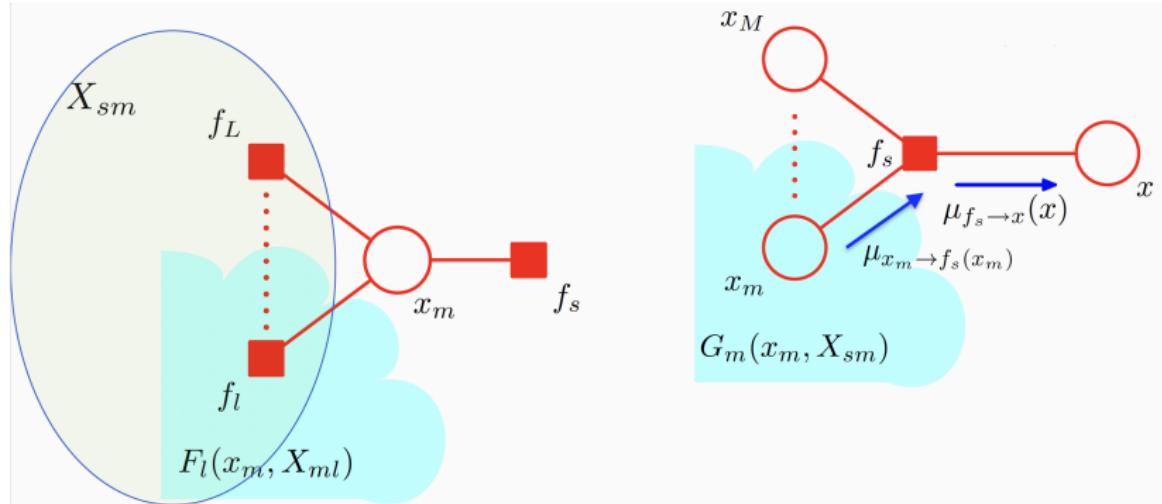


How to compute the message from a variable to a factor?

$$\mu_{x_m \rightarrow f_s}(x_m) \equiv \sum_{X_{sm}} G_m(x_m, X_{sm})$$

Very similar to how we compute $p(x)$,
but with a small difference

The Sum-Product Algorithm



$$\begin{aligned}\mu_{x_m \rightarrow f_s}(x_m) &\equiv \sum_{X_{sm}} G_m(x_m, X_{sm}) &= \sum_{X_{sm}} \prod_{l \in \text{ne}(x_m) \setminus f_s} F_l(x_m, X_{ml}) \\ &= \prod_{l \in \text{ne}(x_m) \setminus f_s} \mu_{f_l \rightarrow x_m}(x_m)\end{aligned}$$

- Now we have two message passing rules

- From a factor node to a variable node

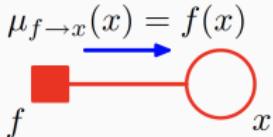
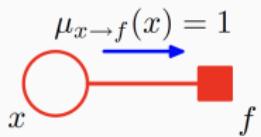
$$\mu_{f_s \rightarrow x}(x) = \sum_{x_1} \dots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f_s}(x_m)$$

- From a variable node to a factor node

$$\mu_{x_m \rightarrow f_s}(x_m) = \prod_{l \in \text{ne}(x_m) \setminus f_s} \mu_{f_l \rightarrow x_m}(x_m)$$

Alternately pass messages!

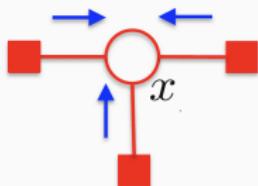
- Initial messages on the leaves



- How to conduct the order of message passing?
 1. Pick an arbitrary node as the root
 2. Compute and propagate messages *from the leaf nodes to the root*, and store received messages at every node
 3. Compute and propagate messages *from the root to the leaf nodes*, storing the messages at every node

- After the message passing done, how to compute the marginals?

$$\begin{aligned} p(x) &= \prod_{s \in \text{ne}(x)} \left[\sum_{X_s} F_s(x, X_s) \right] \\ &= \prod_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x). \end{aligned}$$



Just multiple the received messages of the variable, and normalize as necessary!

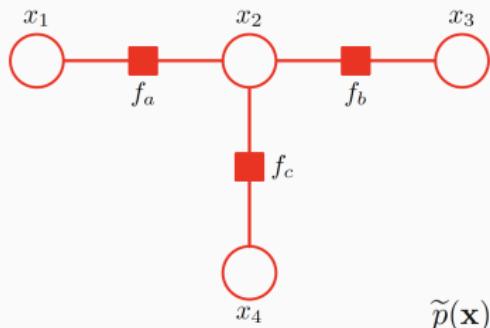
- Why do we need normalization
 - Undirected graphical models (MRF)

$$p(\mathbf{x}) = \frac{1}{Z} \prod_C \psi_C(\mathbf{x}_C) \quad \text{The potentials are not normalization}$$

- Some nodes have been observed

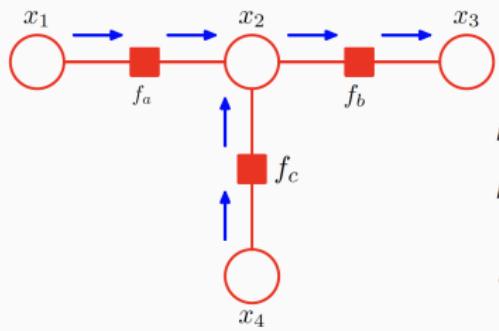
We actually fix the value of the observed nodes in message computation

The Sum-Product Algorithm



$$\tilde{p}(\mathbf{x}) = f_a(x_1, x_2) f_b(x_2, x_3) f_c(x_2, x_4)$$

The Sum-Product Algorithm: example



$$\mu_{x_1 \rightarrow f_a}(x_1) = 1$$

$$\mu_{f_a \rightarrow x_2}(x_2) = \sum_{x_1} f_a(x_1, x_2)$$

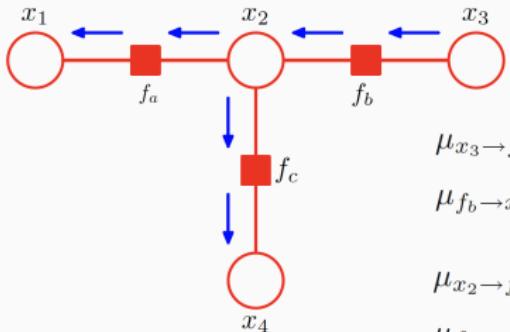
$$\mu_{x_4 \rightarrow f_c}(x_4) = 1$$

$$\mu_{f_c \rightarrow x_2}(x_2) = \sum_{x_4} f_c(x_2, x_4)$$

$$\mu_{x_2 \rightarrow f_b}(x_2) = \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2)$$

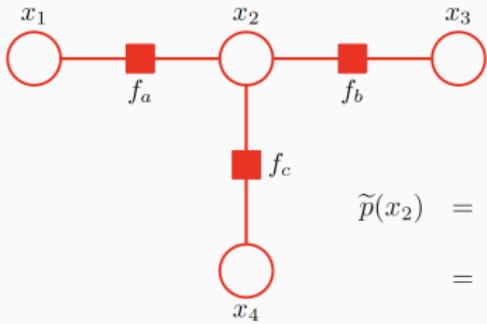
$$\mu_{f_b \rightarrow x_3}(x_3) = \sum_{x_2} f_b(x_2, x_3) \mu_{x_2 \rightarrow f_b}.$$

The Sum-Product Algorithm: example



$$\begin{aligned}\mu_{x_3 \rightarrow f_b}(x_3) &= 1 \\ \mu_{f_b \rightarrow x_2}(x_2) &= \sum_{x_3} f_b(x_2, x_3) \\ \mu_{x_2 \rightarrow f_a}(x_2) &= \mu_{f_b \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2) \\ \mu_{f_a \rightarrow x_1}(x_1) &= \sum_{x_2} f_a(x_1, x_2) \mu_{x_2 \rightarrow f_a}(x_2) \\ \mu_{x_2 \rightarrow f_c}(x_2) &= \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_b \rightarrow x_2}(x_2) \\ \mu_{f_c \rightarrow x_4}(x_4) &= \sum_{x_2} f_c(x_2, x_4) \mu_{x_2 \rightarrow f_c}(x_2).\end{aligned}$$

The Sum-Product Algorithm: example



$$\begin{aligned}\tilde{p}(x_2) &= \mu_{f_a \rightarrow x_2}(x_2) \mu_{f_b \rightarrow x_2}(x_2) \mu_{f_c \rightarrow x_2}(x_2) \\ &= \left[\sum_{x_1} f_a(x_1, x_2) \right] \left[\sum_{x_3} f_b(x_2, x_3) \right] \left[\sum_{x_4} f_c(x_2, x_4) \right] \\ &= \sum_{x_1} \sum_{x_2} \sum_{x_4} f_a(x_1, x_2) f_b(x_2, x_3) f_c(x_2, x_4) \\ &= \sum_{x_1} \sum_{x_3} \sum_{x_4} \tilde{p}(\mathbf{x})\end{aligned}$$

- Step 1. Pick a root node x and arrange the graph into a tree

- Step 2.
 - For each child factor f of x
$$\mu_{f \rightarrow x}(x) = \text{Collect}(f, x)$$

- Step 3.
 - For each child factor f of x
Distribute (x, f)

Collect (x, f)

if x is a leaf, return 1

for each child factor f_j of x (*note: not including f*)

$\mu_{f_j \rightarrow x}(x) = \text{Collect}(f_j, x)$

return $\prod_j \mu_{f_j \rightarrow x}(x)$

Collect (f, x)

if f is a leaf, return $f(x)$

for each child variable x_j of f (*note: not including x*)

$\mu_{x_j \rightarrow f}(x) = \text{Collect}(x_j, f)$

return $\sum_{x_1, \dots, x_M} f(x, x_1, \dots, x_M) \prod_j \mu_{x_j \rightarrow f}(x)$

Distribute (x, f)

- compute and store $\mu_{x \rightarrow f}(x)$ directly
- if f is a leaf, return
- for each child variable x_j of f (*note: not including x*)
 - Distribute (f, x_j)

Distribute (f, x)

- compute and store $\mu_{f \rightarrow x}(x)$ directly
- if x is a leaf, return
- for each child factor f_j of x (*note: not including f*)
 - Distribute (x, f_j)

- In general graphs that contain cycles, sum-product *cannot guarantee exact inference*
- The exact inference on general graphs is called Junction tree algorithm
 - It first merges factors and turns the initial graph into a junction tree and then run a sum-product-like algorithm
 - Intractable on graphs with large factors

- We can still apply sum-product on general graphs as an *approximate* inference algorithm
- First initialize all the messages with 1 (or random)
- Run sum-product (with any message passing order) repeatedly until convergence (not guaranteed!)
- Often works really well, sometimes totally fail
- Striking connections between LBP and decoding (turbo codes) in information theory

- A simple variant of the sum-product algorithm
- Objective: an efficient algorithm to find
 - The value \mathbf{x}_{\max} that maximizes $p(\mathbf{x})$
 - The value of $p(\mathbf{x}_{\max})$
- Very important in many tasks, e.g., structure prediction, decision,

- In general, maximum marginals \neq joint maximum

		$x = 0$	$x = 1$
$y = 0$	$x = 0$	0.3	0.4
	$x = 1$	0.3	0.0

$$\arg \max_x p(x, y) = 1$$

$$\arg \max_x p(x) = 0$$



$$\begin{aligned} p(\mathbf{x}^{\max}) &= \max_{\mathbf{x}} p(\mathbf{x}) = \max_{x_1} \dots \max_{x_M} p(\mathbf{x}) \\ &= \frac{1}{Z} \max_{x_1} \dots \max_{x_N} [\psi_{1,2}(x_1, x_2) \dots \psi_{N-1,N}(x_{N-1}, x_N)] \\ &= \frac{1}{Z} \max_{x_1} \left[\max_{x_2} \left[\psi_{1,2}(x_1, x_2) \left[\dots \max_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) \right] \dots \right] \right] \end{aligned}$$

- We still have the distributive law

$$\max(ab, ac) = a \max(b, c)$$

So we can simply replace sum by max in the sum-product algorithm!

- Generalizes to tree-structured factor graph

$$\max_{\mathbf{x}} p(\mathbf{x}) = \max_{x_n} \prod_{f_s \in \text{ne}(x_n)} \max_{X_s} f_s(x_n, X_s)$$

- To enhance numerical stability, we take log

$$\ln \left(\max_{\mathbf{x}} p(\mathbf{x}) \right) = \max_{\mathbf{x}} \ln p(\mathbf{x}).$$

The distributive law still holds

$$\max(a + b, a + c) = a + \max(b, c)$$

So we only need to replace sum by max, product by sum in the sum-product algorithm

Initialization message (leaf nodes)

$$\mu_{x \rightarrow f}(x) = 0 \quad \mu_{f \rightarrow x}(x) = \ln f(x)$$

Message passing (recursively)

$$\mu_{f \rightarrow x}(x) = \max_{x_1, \dots, x_M} \left[\ln f(x, x_1, \dots, x_M) + \sum_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f}(x_m) \right]$$

$$\mu_{x \rightarrow f}(x) = \sum_{l \in \text{ne}(x) \setminus f} \mu_{f_l \rightarrow x}(x)$$

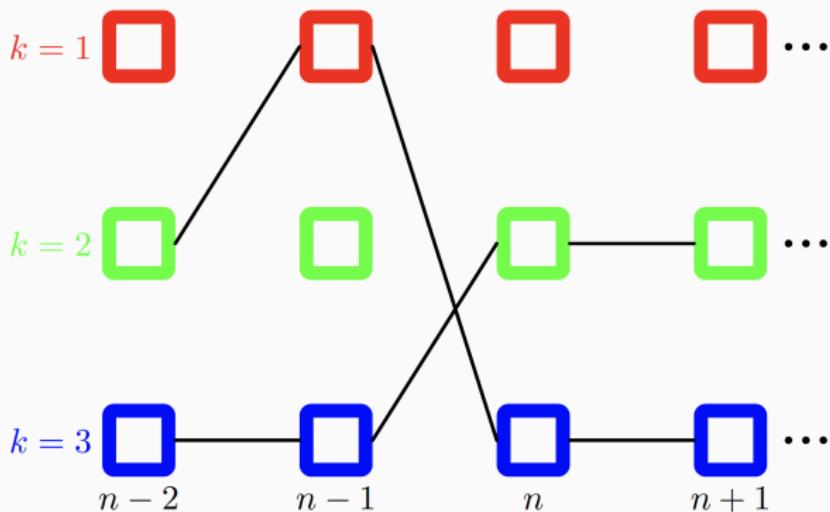
- First pass from leaves to the root and the second pass from the root to leaves
- Termination

$$p^{\max} = \max_x \left[\sum_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x) \right]$$

$$x^{\max} = \arg \max_x \left[\sum_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x) \right]$$

- How to find the global configuration x_{\max} that gives the maximum probability?
- We need to store a quantity to tell us how to trace back to the variable value that maximizes the previous sub-problem (back-tracking)
- So each message can contain two component: (1) the max-sum value (2) the variable value that gives the max-sum (i.e., argmax)

The max-sum algorithm



- This is essentially dynamic programming
- For hidden Markov models, this is known as Viterbi algorithm

- Factor graph definition
- Sum-product algorithm
- Message-passing
- Accurate for tree-structured graphs, not guaranteed to be accurate for graphs with cycles
- Loopy belief propagation
- Max-product algorithm, max-sum
- Be able to implement the algorithms!