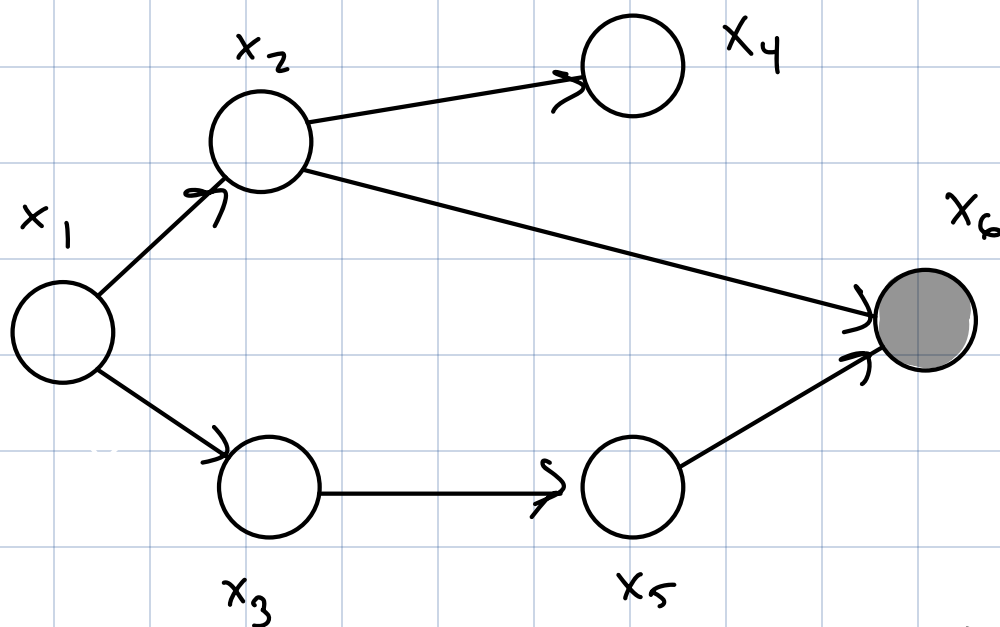# Outline

# Variable elimination

- model $P(x_1 \ldots x_n) \equiv P(x)$, $x_i \in \{1 \ldots k\}$

- 3 subsets of variables
  - $X_E$ "evidence" nodes
  - $X_F$ "query" nodes
  - $X_R = (X_E \cup X_F)^c$ "everything else"

- Goal: $\underline{P(x_F \mid X_E)}$

$$\frac{P(X_E, X_F)}{P(X_E)} = \frac{\sum_{X_R} P(x)}{\sum_{X_F} \sum_{X_R} P(x)}$$

- Naively: $O(k^{|R|})$ to compute $P(x_E, x_F)$

---

- Example:



- compute $P(x_1 \mid \bar{x}_6) \equiv P(X_1 = x_1 \mid X_6 = \bar{x}_6)$

- Define the "evidence potential"

$$\delta_{\bar{x}_6}(x_6) = \begin{cases} 1 & \text{if } x_6 = \bar{x}_6 \\ 0 & \text{otherwise} \end{cases}$$

- $P(x_1, \bar{x_6})$

$$= \sum_{x_2} \cdots \sum_{x_6} P(x_1) P(x_2|x_1) \cdots P(x_6|x_2, x_5) \, \delta_{\bar{x_6}}(x_6)$$

$$= P(x_1) \sum_{x_2} P(x_2|x_1) \sum_{x_3} \cdots \underbrace{\sum_{x_6} P(x_6|x_2, x_5) \, \delta_{\bar{x_6}}(x_6)}_{\triangleq \; m_6(x_2, x_5)}$$

- "intermediate factor"
- function of $x_2, x_5$
- variable $x_6$ is <u>eliminated</u>
- $m_6(\cdots)$ is a "message" from $x_6$

$$= P(x_1) \sum_{x_2} \cdots \underbrace{\sum_{x_5} P(x_5|x_3) \, m_6(x_2, x_5)}_{\triangleq \; m_5(x_2, x_3)}$$

This factor does not depend on $x_4$, so ...

$$= P(x_1) \cdots \underbrace{\sum_{x_3} P(x_3|x_1) \, m_5(x_2, x_3)}_{\triangleq \; m_3(x_1, x_2)} \underbrace{\sum_{x_4} P(x_4|x_2)}_{\triangleq \; m_4(x_2) \; = 1}$$

$$= P(x_1) \sum_{x_2} P(x_2|x_1) \, m_3(x_1, x_2) \, m_4(x_2)$$

$$= P(x_1) \, m_2(x_1)$$

This computes $P(x_1, \bar{x_6})$. The <u>evidence</u> is then:
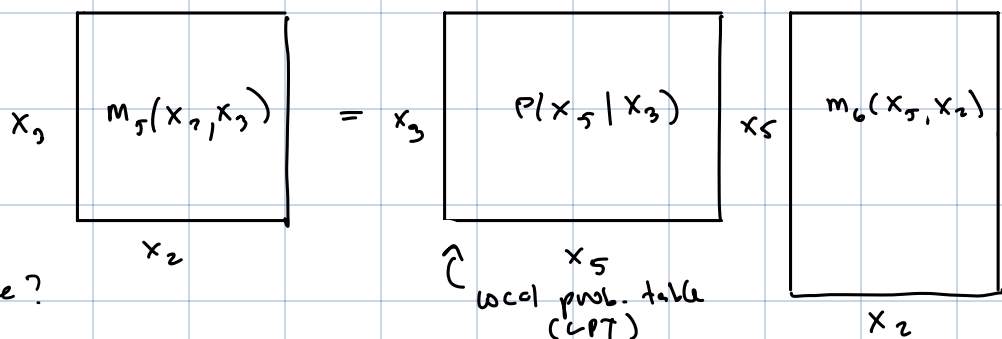
$$P(\bar{x_6}) = \sum_{x_1} P(x_1) \, m_2(x_1)$$

And the posterior is then:

$$P(x_1 | \bar{x_6}) = \frac{P(x_1, \bar{x_6})}{P(\bar{x_6})}$$

# Discussion:

- iterative product-sum procedure

- eliminated variables leave behind <u>messages</u>

- the message $m_i(\ldots)$ from $x_i$ is a <u>function</u> which takes as arguments the values of <u>non-eliminated</u> variables

- recall that $m_5(x_2, x_3) = \sum_{x_5} P(x_5 | x_3) m_6(x_5, x_2)$

- think of these as tables:

$$x_3 \boxed{m_5(x_2, x_3)} = x_3 \boxed{P(x_5|x_3)} x_5 \boxed{m_6(x_5,x_2)}$$

$x_2$     $\underset{\substack{\text{local prob. table} \\ \text{(CPT)}}}{x_5}$     $x_2$

- what <u>is</u> a message?

- $m_5(x_2, x_3) = \sum_{x_5} P(x_5 | x_3) \underbrace{\sum_{x_6} P(x_6 | x_2, x_5) \, \delta_{\bar{x}_6}(x_6)}_{= P(\bar{x}_6 | x_2, x_5)}$

$$\underbrace{\phantom{= \sum_{x_5} P(x_5|x_3) P(\bar{x}_6|x_2,x_5)}}_{= P(x_5, \bar{x}_6 | x_2, x_3)}$$

$$= P(\bar{x}_6 | x_2, x_3) \quad (\text{a conditional evidence table})$$

- look at the graph and imagine eliminating (i.e., marginalizing out) $x_3$

- <u>Question</u>: what if we had eliminated $x_2$ first?

$$m_2(\ldots) = \sum_{x_2} P(x_2 | x_1) P(x_4 | x_2) P(x_6 | x_2, x_5)$$

$$\equiv m_2(x_1, x_4, x_5, x_6) \quad (\text{4d table})$$

- <u>Intuition</u>: complexity depends on the <u>ordering</u>, and a good ordering exploits the cond. indep. structure

# The variable elimination (VE) algorithm

Input: graph $G$, target $P(X_F \mid \bar{X}_E)$

① ORDER $X_{1...n}$ with $\underline{X_F \text{ last}}$

② initialize "active list" of functions

$$A = \left[\, P(x_i \mid x_{\pi_i}) \text{ for } i = 1...n \,\right]$$
$$\underset{(CPTs)}{}$$

$$+ \left[\, \delta_{\bar{x}_i}(x_i) \text{ for } i \in E \,\right]$$
$$\underset{(\text{evidence potentials})}{}$$

③ for $i$ in ORDER:

• $A_i = \left[\, \text{all } f(x_i, ...) \text{ in } A \,\right]$

• $S_i = \left[\, \text{all } j \text{ s.t. } f(x_i, x_j, ...) \text{ in } A_i \,\right]$

• $\phi_i(x_i, s_i) = \prod_{f \in A_i} f(x_i, ...)$

• $M_i(S_i) = \sum_{x_i} \prod_{f \in A_i} f(x_i, ...)$ ← some subset of $S_i$

"eliminate" ⟶ • $A = A - A_i$

• $A = A + m_i(S_i)$

④ Return $\dfrac{\phi_F(x_F, S_F)}{m_F(S_F)} \equiv \dfrac{P(x_F, \bar{x}_E)}{\sum_{x_F} P(x_F, \bar{x}_E)}$

- Notice that the complexity depends on the sizes of the intermediate $S_i$ sets that are created to eliminate $x_i$

- $S_i$ is called the __elimination clique__

- for a given ORDER, we can see the sequence of $S_i$'s graphically, by converting to an undirected graph:
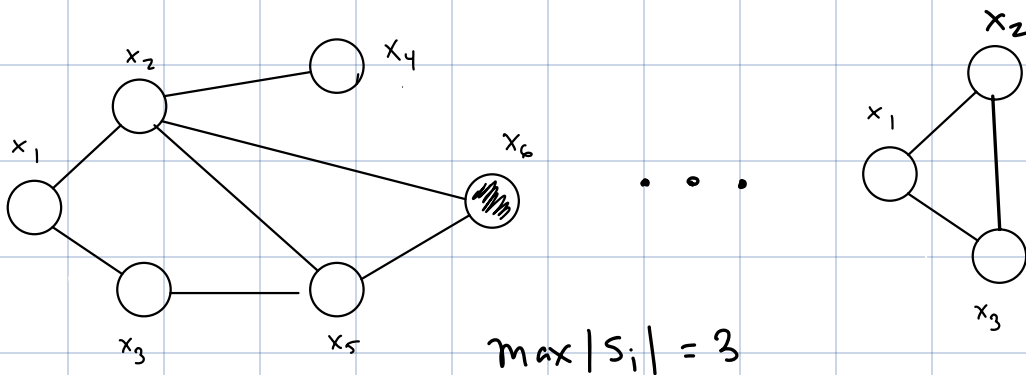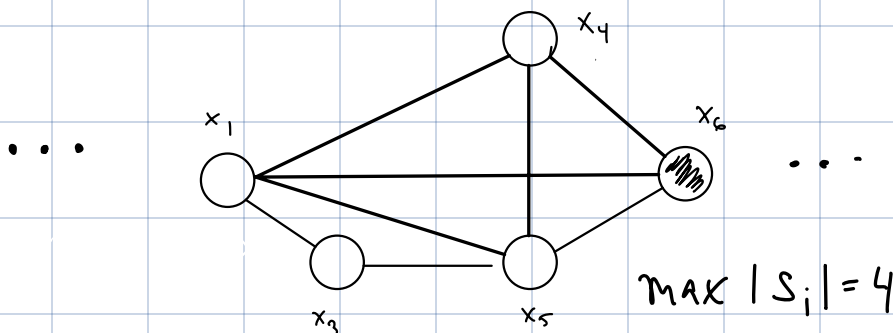
  ① moralize

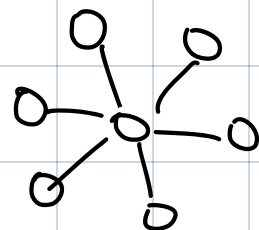  ② for i in ORDER:
     remove i
     connect all neighbors(i)

e.g. ORDER = $\{6, 4, 5, 2, 3\}$



$\max |S_i| = 3$

e.g. ORDER = $\{2, 5, 3, 4, 6\}$

consider also:



$\max |S_i| = 4$

# VE for undirected graphs

- Recall the joint in an UGM:

$$P(x_1 \ldots x_n) = \frac{1}{Z} \prod_{c \in C} \psi(x_c)$$

- Note that

$$P(x_F \mid \bar{x}_E) = \frac{\sum_{x_R} P(x_F, x_R, \bar{x}_E)}{\sum_{x_F} \quad ''}$$

$$= \frac{\frac{1}{Z} \sum_{x_R} \prod_{c \in C} \psi(x_c)}{\frac{1}{Z} \sum_{x_F} \quad ''}$$

$$\propto \sum_{x_F} \sum_{x_R} \prod_{c \in C} \psi(x_c)$$

- So computing arbitrary marginals boils down to iterative sum-product of potentials $\psi(x_c) > 0$

- **VE algo:** same as above except initialize $A$ with $[\psi_c(x_c)$ for $c \in C]$ instead of CPTs

- for convenience also singleton potentials $\psi(x_i)$

$$P(x_1 \ldots x_n) \propto \prod_i \psi(x_i) \prod_c \psi(x_c)$$
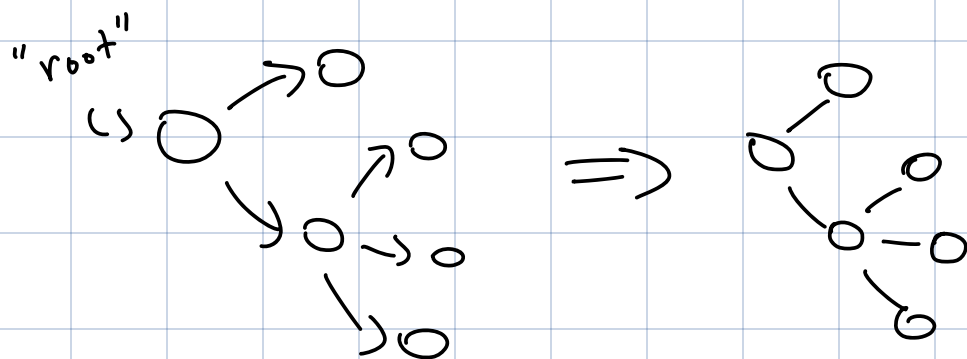
- This makes conditioning simple:

$$\psi(x_i) := \psi(x_i) \, \delta_{\bar{x}_i}(x_i) \quad \text{if } i \in E$$

# Trees

Def'n (directed): $|par(x_i)| = 1 \quad \forall i \neq r$

Def'n (undirected): $|paths(x_i, x_j)| = 1 \quad \forall i, j$

Fact: "Moralized" directed tree = undirected tree

"root"

$\circlearrowleft$



Parameterization directed $\Rightarrow$ undirected:

$$P(x_1 \ldots x_n) = \underbrace{P(x_r)}_{\text{"root"}} \prod_{i \to j} P(x_j | x_i)$$

$\Downarrow$

$$\psi(x_i, x_j) \overset{\Delta}{=} P(x_j | x_i) \quad \forall \; i \to j$$
$$\psi(x_r) \overset{\Delta}{=} P(x_r)$$
$$\psi(x_i) \overset{\Delta}{=} 1 \quad \forall \; i \neq r$$

(where again we will redefine the singletons to
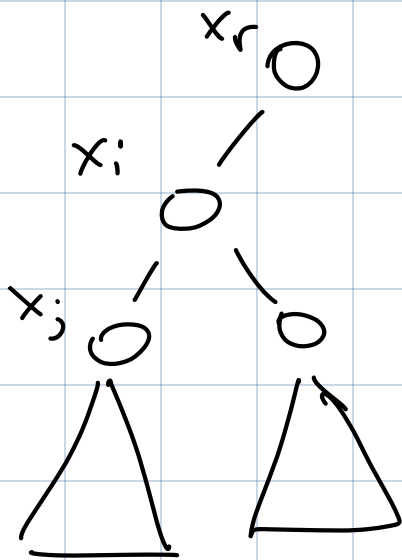condition on evidence $\psi(x_i) = \psi(x_i) \delta_{\bar{x}_i}(x_i)$ )

# VE for trees

- Query node : $X_F$

- To determine the ORDER :
  - ① Set $X_F$ to the root $X_r$
  - ② Direct edges away from $X_F$
  - ③ Order by depth first

$X_F$

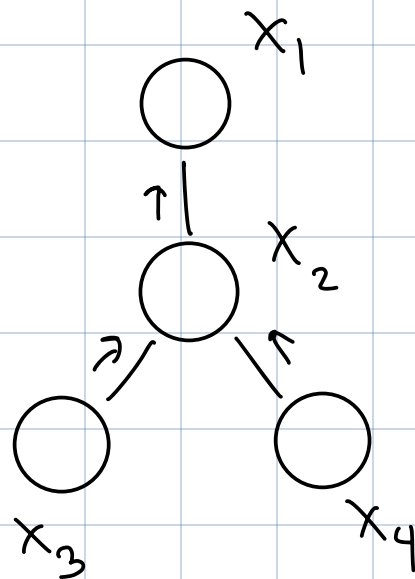$\Rightarrow$

$\max(|S_i|) = 2$

- Now eliminate bottom-up

$X_r$

$X_i$

$X_j$

$$m_j(\cdots) = \sum \prod_{\substack{j \quad f \in A_j}} f(m_j \cdots)$$

$$= \sum_j \psi(x_j) \psi(x_i, x_j) \prod_{k \in \text{child}(j)} m_k(\cdots)$$

$$= m_j(x_i) \equiv m_{j \rightarrow i}(x_i)$$

"message"

Example:

Goal: $P(X_1)$

$$m_{3 \to 2}(x_2) = \sum_{x_3} \psi(x_3)\psi(x_2, x_3)$$

$$m_{4 \to 2}(x_2) = \cdots$$

$$m_{2 \to 1}(x_1) = \sum_{x_2} \psi(x_2)\psi(x_1, x_2)\, m_{3 \to 2}(x_2)\, m_{4 \to 2}(x_2)$$

$$P(X_1) \propto \psi(x_1)\, m_{2 \to 1}(x_1)$$
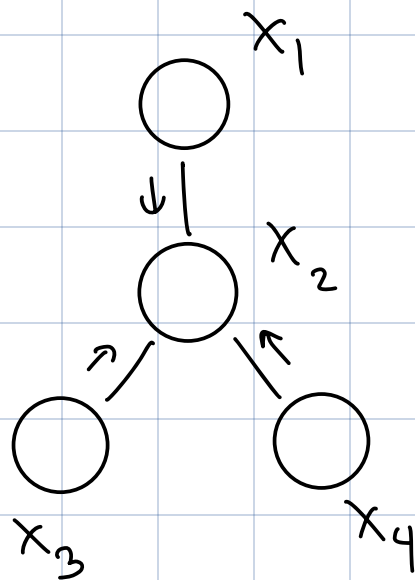
Goal: $P(X_2)$

Notice that:

$$m_{3 \to 2} = \cdots$$

$$m_{4 \to 2} = \cdots$$

are the same as above.

Goal: $P(X_4)$

$$m_{1 \to 2} = \cdots \quad (\text{same})$$

$\longrightarrow$ wasting computation for multiple queries

# The sum-product algorithm
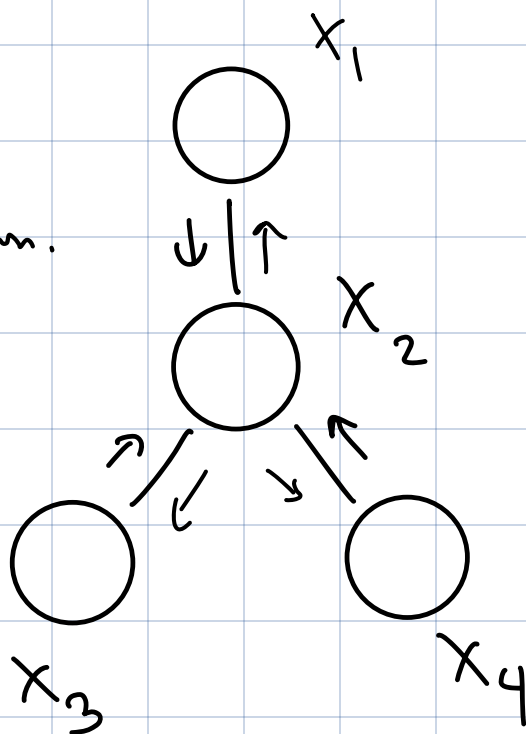
- computes <u>all</u> marginals in a tree

- only 2 messages per node

- <u>schedule</u> : choose <u>any</u> node as root.
  Messages up the tree, and messages down.

- <u>Claim</u> : all marginals can be computed
  from $\{ m_{i \to j}, m_{j \to i} \}$ using this
  schedule.

- $P(x_i \mid \bar{x}_E) \propto \varphi(x_i) \prod\limits_{j \in Neigh(i)} m_{j \to i}(x_i)$

- <u>Proof</u> : Def'n of the undirected tree...

- Another perspective : asynchronous <u>message-passing</u>
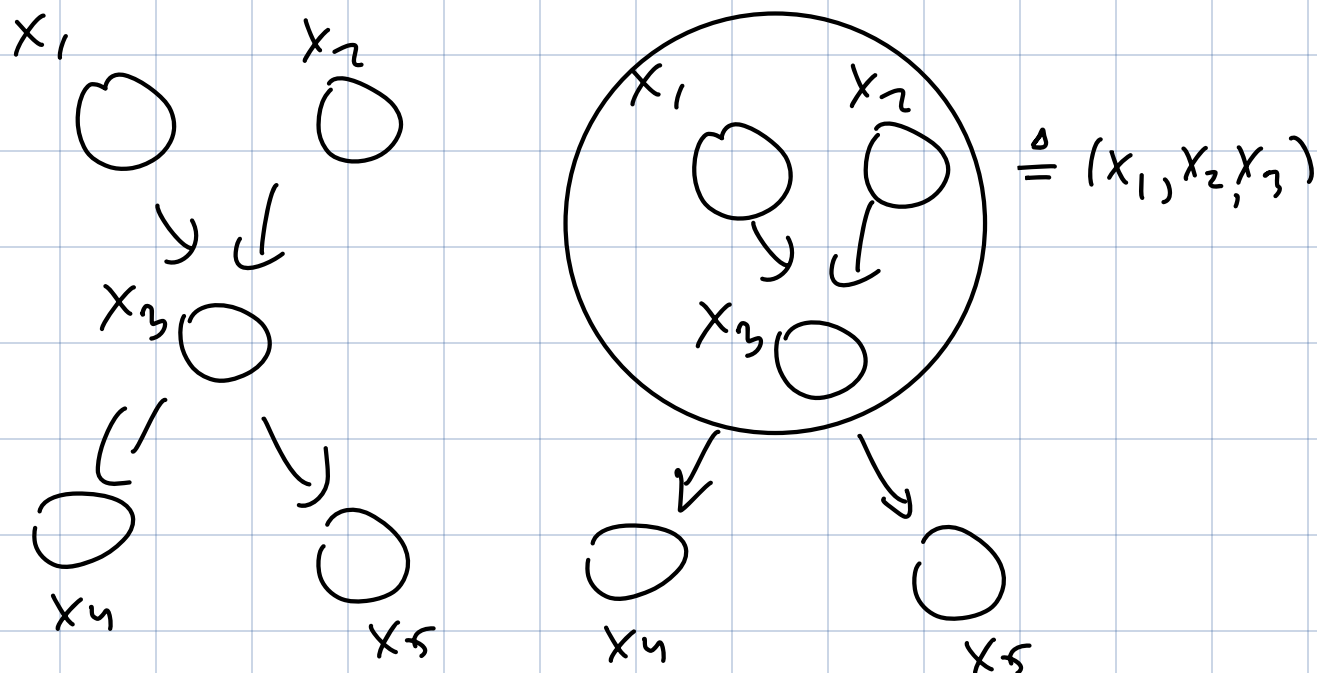
- <u>Message-passing protocol</u> :
  A node $i$ sends message to $j$ $m_{i \to j}$ only after
  it has received messages from other neighbors $k \in Neigh(i) \setminus j$.

- For <u>trees</u>, we can confirm that a fixed schedule
  produces messages that satisfy the protocol.

# Junction tree algorithm

- Convert any graph into a tree of super nodes, then run SP.

$$\stackrel{\Delta}{=} (X_1, X_2, X_3)$$

- Exponential in size of super nodes

# Belief propagation (BP)

- SP = BP on trees
- for more graphs BP is "loopy" and must be iterated until convergence

- Convergence not guaranteed

- When LBP converges it often produces good solutions

- One can show that LBP is minimizing the
  **Bethe free energy**

$$KL(Q \| P) = \sum_x Q(x) \log P(x) - \sum_x Q(x) \log P(x)$$

$$= -H_Q(x) - \mathbb{E}_Q[\log P(x)]$$

$$= -H_Q(x) - \mathbb{E}_Q[\log f(x)] + \log Z$$

- So we can minimize $KL(Q\|P)$ without
  doing inference in $P(x)$

$$\quad \hookrightarrow \min_Q KL(Q\|P) = \max_Q \underbrace{\mathbb{E}_Q[\log f(x)] + H_Q(x)}$$

$$\equiv \text{"ELBO"}$$

- $H_Q(x) = \sum_{x_1} \cdots \sum_{x_n} Q(x) \log \frac{1}{Q(x)}$ $\qquad \equiv -\text{"Free Energy"}$

  often intractable

- for trees, this is tractable, for the same reasons BP is.

- $H^{\text{Bethe}}(x)$ computes $H_Q(x)$ as if the graph were a tree:

$$Q^{\text{bethe}}(x) = \prod_{j \to i} \frac{Q(x_i, x_j)}{Q(x_j) Q(x_i)} \prod_i Q(x_i)$$

  (products of singleton and pairwise marginals)

- (This is exact for a tree.)
- $F^{Bethe} = -\mathbb{E}_{Q^B}[\log f(x)] + H^{Bethe}(x)$
- $Q^{Bethe}(x)$ is only "legal" if it is in the <u>marginal polytope</u>

$$M(G) = \{(Q(x_i), Q(x_i, x_j)) : \exists \, Q \text{ with those marginals}\}$$

- Searching over this space is hard. Instead, LBP searches in $L(G) \supset M(G)$
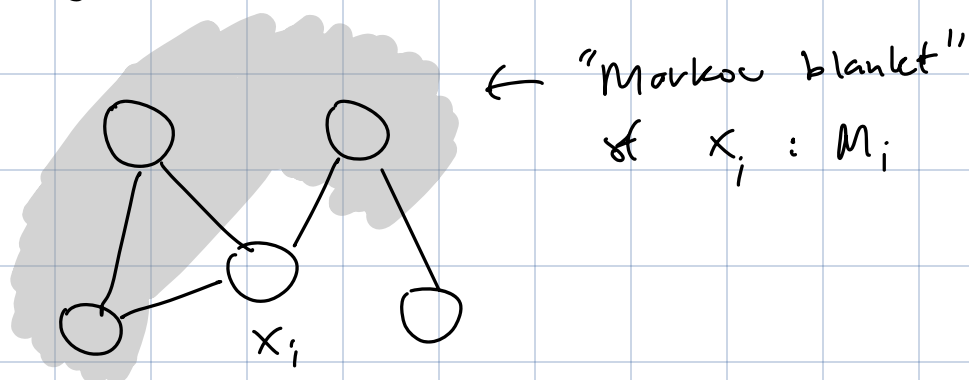
- A fixed point of LBP is

$$\text{argmin } F^{Bethe} (P, Q)$$
$$Q \in L(G)$$

- often works well, but Bethe free energy is not guaranteed to be near the free energy and LBP is not guaranteed to converge

---

## Connections

- LBP is a form of <u>variational inference</u>

- More well-known VI is based on

$$\text{argmax } \mathbb{E}_Q[\log f(x)] + H(x)$$
$$Q \in Q$$

for a simple class of tractable $Q$

- EM is a special case of VI

- Many other inference algos are a form of iterative message-passing

- e.g., Gibbs sampling



$\leftarrow$ "Markov blanket" of $x_i$ : $M_i$

$x_i$

for iter in $1 \ldots S$:
    for $i$ in nodes:
$$x_i \sim p(x_i \mid x_{M_i})$$
        "complete conditional"

- For exact inference, discrete or Gaussian

- $$P(x) = \frac{1}{Z} \prod_{(i,j)} \varphi(x_i, x_j) \prod_i \psi_i(x_i)$$

$$\varphi(x_i, x_j) \overset{\circ}{=} \exp(x_i^T V_{ij} x_j)$$

$$\psi(x_i) \overset{\circ}{=} \exp\left(\frac{1}{\delta_i}(x_i - \mu_i)^2\right)$$

- Same BP messages as above...

$$P(x_i) \propto \psi(x_i) \prod_{j \in neigh(j)} m_{j \to j}(x_i)$$
$$\propto \mathcal{N}(x_i, \ldots)$$

- Generalizes the "Kalman filter"

# Max product algorithm

Goal: $x^* = \underset{x}{\text{argmax}} \ P(X)$  (MAP)

$$P(x^*) = \underset{x}{\max} \ P(X)$$

$$= \underset{x_1}{\max} \ P(X_1) \ \underset{x_2}{\max} \ P(x_2 | x_1) \ldots$$

(max is also distributive)

## Messages (for trees):

### up the tree:

$$m_{j \to i}(x_i) = \underset{x_j}{\max} \ \varphi(x_j) \ \varphi(x_i, x_j) \ \overline{\prod_{\kappa \in neigh(j) \setminus i}} m_{\kappa \to j}(x_j)$$

$$\delta_{j \to i}(x_i) = \underset{x_j}{\text{argmax}} \ \text{"} \quad \text{"}$$

### at the root:

$$\max_x p(x) = \underset{x_r}{\max} \ \varphi(x_v) \ \overline{\prod_{\kappa \in neigh(r)}} m_{\kappa \to r}(x_r)$$

$$x_r^* \leftarrow \underset{x_r}{\text{argmax}} \ \text{"} \quad \text{"}$$

### down the tree:

$$x_j^* \leftarrow \delta_{j \to i}(x_i^*)$$