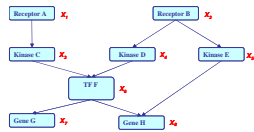


Approximate Inference: Loopy Belief Propagation and Variants

Probabilistic Graphical Models (10-708)

Lecture 16, Nov 7, 2007

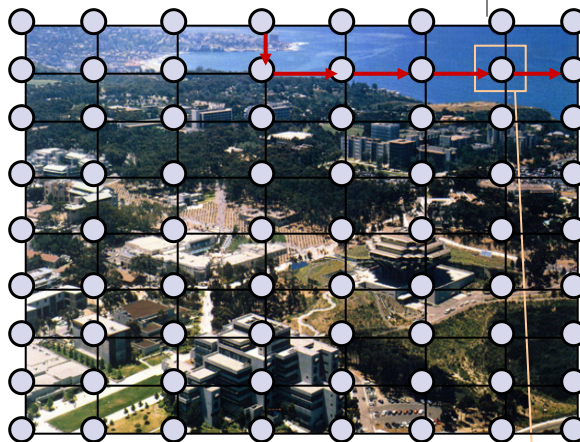


Eric Xing

Reading: KF-Chap. 12

An Ising model on 2-D image

- Nodes encode hidden information (patch-identity).
- They receive local information from the image (brightness, color).
- Information is propagated through the graph over its edges.
- Edges encode 'compatibility' between nodes.



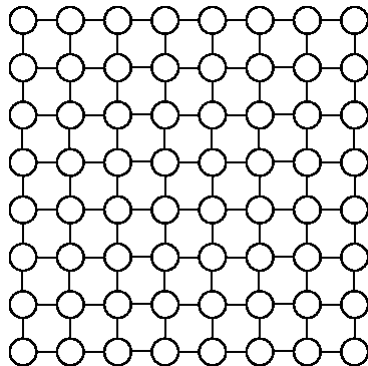
air or water ?



Why Approximate Inference?



- Tree-width of $N \times N$ graph is $O(N)$
- N can be a huge number (~1000s of pixels)
- Exact inference will be too expensive



$$p(X) = \frac{1}{Z} \exp \left\{ \sum_{i < j} \theta_{ij} X_i X_j + \sum_i \theta_{i0} X_i \right\}$$

Eric Xing

3

Variational Methods



- For a distribution $p(X/\theta)$ associated with a complex graph, computing the marginal (or conditional) probability of arbitrary random variable(s) is intractable
- Variational methods
 - formulating probabilistic inference as an optimization problem:

$$\text{e.g. } f^* = \arg \max_{f \in \mathcal{S}} \{ F(f) \}$$

f : a (tractable) probability distribution
or, solutions to certain probabilistic queries

Eric Xing

4



Bethe Energy Minimization



The Objective

- Let us call the actual distribution P

$$P(X) = 1/Z \prod_{f_a \in F} f_a(X_a)$$

- We wish to find a distribution Q such that Q is a “good” approximation to P
- Recall the definition of KL-divergence

$$KL(Q_1 \parallel Q_2) = \sum_x Q_1(X) \log\left(\frac{Q_1(X)}{Q_2(X)}\right)$$

- $KL(Q_1 \parallel Q_2) \geq 0$
- $KL(Q_1 \parallel Q_2) = 0$ iff $Q_1 = Q_2$
- But, $KL(Q_1 \parallel Q_2) \neq KL(Q_2 \parallel Q_1)$

Which KL?



- Computing $KL(P||Q)$ requires inference!
- But $KL(P||Q)$ can be computed without performing inference on P

$$\begin{aligned} KL(Q || P) &= \sum_X Q(X) \log\left(\frac{Q(X)}{P(X)}\right) \\ &= \sum_X Q(X) \log Q(X) - \sum_X Q(X) \log P(X) \\ &= -H_Q(X) - E_Q \log P(X) \end{aligned}$$

- Using $P(X) = 1/Z \prod_{f_a \in F} f_a(X_a)$

$$\begin{aligned} KL(Q || P) &= -H_Q(X) - E_Q \log\left(1/Z \prod_{f_a \in F} f_a(X_a)\right) \\ &= -H_Q(X) - \log 1/Z - \sum_{f_a \in F} E_Q \log f_a(X_a) \end{aligned}$$

Eric Xing

7

The Objective



-

$$KL(Q || P) = \underbrace{-H_Q(X) - \sum_{f_a \in F} E_Q \log f_a(X_a)}_{F(P, Q)} + \log Z$$

- We will call $F(P, Q)$ the “Energy Functional” *
- $F(P, P) = ?$
- $F(P, Q) \geq F(P, P)$

Eric Xing

*also called Gibbs Free Energy

The Energy Functional



- Let us look at the functional

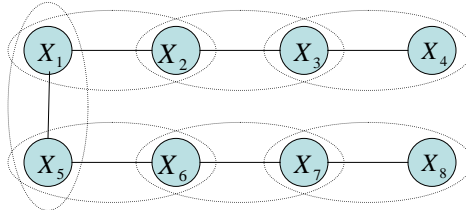
$$F(P, Q) = -H_Q(X) - \sum_{f_a \in F} E_Q \log f_a(X_a)$$

- $\sum_{f_a \in F} E_Q \log f_a(X_a)$ can be computed if we have marginals over each f_a
- $H_Q = -\sum_X Q(X) \log Q(X)$ is harder! Requires summation over all possible values
- Computing F , is therefore hard in general.
- Approach 1: Approximate $F(P, Q)$ with easy to compute $\hat{F}(P, Q)$

Tree Energy Functionals



- Consider a tree-structured distribution



- The probability can be written as: $b(\mathbf{x}) = \prod_{i,j \in E} b_{ij}(x_i, x_j) \prod_i b_i(x_i)^{-1}$
- $H_{tree} = -\sum_{i,j \in E} \sum_{x_i, x_j} b_{ij}(x_i, x_j) \ln b_{ij}(x_i, x_j) + \sum_i 1 \times \sum_{x_i} b_i(x_i) \ln b_i(x_i)$
- $$F_{tree} = -\left(-\sum_{i,j \in E} \sum_{x_i, x_j} b_{ij}(x_i, x_j) \ln b_{ij}(x_i, x_j) + \sum_i 1 \times \sum_{x_i} b_i(x_i) \ln b_i(x_i) \right) - \sum_{i,j \in E} \sum_{x_i, x_j} b_{ij}(x_i, x_j) \ln f_{i,j}(x_i, x_j) - \sum_i \sum_{x_i} b_i(x_i) \ln f_i(x_i)$$

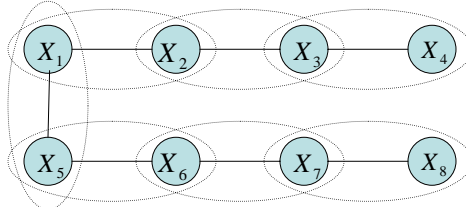
$$= \sum_{i,j \in E} \sum_{x_i, x_j} b_{ij}(x_i, x_j) \ln \frac{b_{ij}(x_i, x_j)}{f_{i,j}(x_i, x_j)} + \sum_i \sum_{x_i} b_i(x_i) \ln \frac{b_i(x_i)}{f_i(x_i)} - \sum_i 2 \times \sum_{x_i} b_i(x_i) \ln b_i(x_i)$$

$$= F_{12} + F_{23} + \dots + F_{67} + F_{78} - F_1 - F_5 - F_2 - F_6 - F_3 - F_7$$
- involves summation over edges and vertices and is therefore easy to compute

Tree Energy Functionals



- Consider a tree-structured distribution



- The probability can be written as: $b(\mathbf{x}) = \prod_a b_a(\mathbf{x}_a) \prod_i b_i(x_i)^{1-d_i}$
- $H_{tree} = -\sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln b_a(\mathbf{x}_a) + \sum_i (d_i - 1) \sum_{x_i} b_i(x_i) \ln b_i(x_i)$
- $F_{Tree} = \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln \frac{b_a(\mathbf{x}_a)}{f_a(\mathbf{x}_a)} + \sum_i (1 - d_i) \sum_{x_i} b_i(x_i) \ln b_i(x_i)$
 $= F_{12} + F_{23} + \dots + F_{67} + F_{78} - F_1 - F_5 - F_2 - F_6 - F_3 - F_7$

- involves summation over edges and vertices and is therefore easy to compute

Eric Xing

11

Bethe Approximation to Gibbs Free Energy

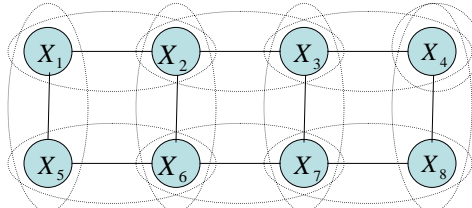


- For a general graph, choose $\hat{F}(P, Q) = F_{Bethe}$

$$H_{Bethe} = -\sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln b_a(\mathbf{x}_a) + \sum_i (d_i - 1) \sum_{x_i} b_i(x_i) \ln b_i(x_i)$$

$$F_{Bethe} = \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln \frac{b_a(\mathbf{x}_a)}{f_a(\mathbf{x}_a)} + \sum_i (1 - d_i) \sum_{x_i} b_i(x_i) \ln b_i(x_i) = -\langle f_a(\mathbf{x}_a) \rangle - H_{Bethe}$$

- Called "Bethe approximation" after the physicist Hans Bethe



$$F_{Bethe} = F_{12} + F_{23} + \dots + F_{67} + F_{78} - F_1 - F_5 - 2F_2 - 2F_6 \dots - F_8$$

- Equal to the exact Gibbs free energy when the factor graph is a tree
- In general, H_{Bethe} is **not** the same as the H of a tree

Eric Xing

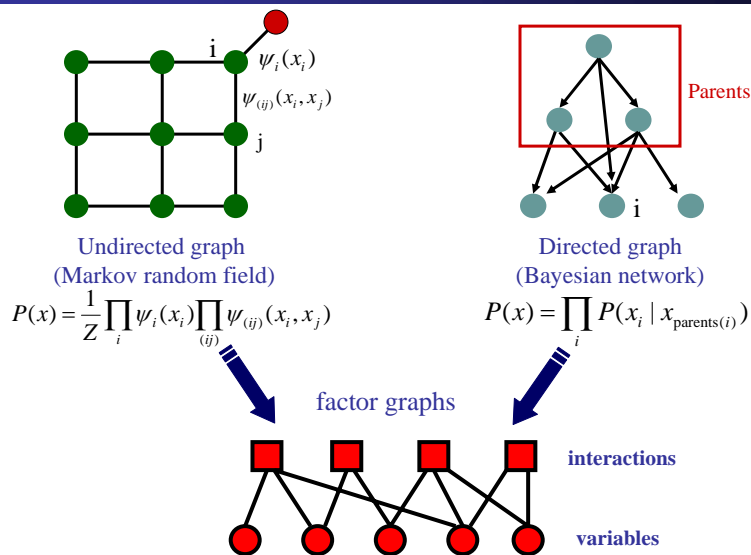
12

Bethe Approximation

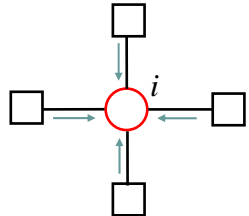


- Pros:
 - Easy to compute, since entropy term involves sum over pairwise and single variables
- Cons:
 - $\hat{F}(P, Q) = F_{\text{bethe}}$ **may or may not** be well connected to $F(P, Q)$
 - It could, in general, be greater, equal or less than $F(P, Q)$
- Optimize each $b(x_{\partial i})$'s.
 - For discrete belief, constrained opt. with *Lagrangian* multiplier
 - For continuous belief, not yet a general formula
 - Not always converge

From GM to factored graphs

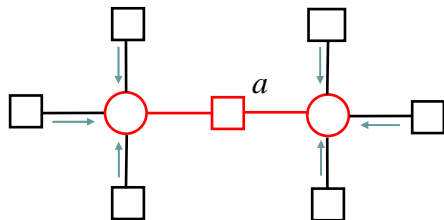


Recall Beliefs and messages in FG



$$b_i(x_i) \propto f_i(x_i) \prod_{a \in N(i)} m_{a \rightarrow i}(x_i)$$

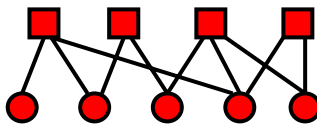
↑ "beliefs" ↑ "messages"



$$b_a(X_a) \propto f_a(X_a) \prod_{i \in N(a)} \prod_{c \in N(i) \setminus a} m_{c \rightarrow i}(x_i)$$

The "belief" is the BP approximation of the marginal probability.

Bethe Free Energy for FG



$$F_{Bethe} = \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln \frac{b_a(\mathbf{x}_a)}{f_a(\mathbf{x}_a)} + \sum_i (1 - d_i) \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \ln b_i(\mathbf{x}_i)$$

$$H_{Bethe} = - \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln b_a(\mathbf{x}_a) + \sum_i (d_i - 1) \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \ln b_i(\mathbf{x}_i)$$

$$F_{Bethe} = - \langle f_a(\mathbf{x}_a) \rangle - H_{Bethe}$$

Constrained Minimization of the Bethe Free Energy



$$L = F_{\text{Bethe}} + \sum_i \gamma_i \{ \sum_{x_i} b_i(x_i) - 1 \} + \sum_a \sum_{i \in N(a)} \sum_{x_i} \lambda_{ai}(x_i) \left\{ \sum_{X_a \setminus x_i} b_a(X_a) - b_i(x_i) \right\}$$

$$\frac{\partial L}{\partial b_i(x_i)} = 0 \implies b_i(x_i) \propto \exp\left(\frac{1}{d_i - 1} \sum_{a \in N(i)} \lambda_{ai}(x_i)\right)$$

$$\frac{\partial L}{\partial b_a(X_a)} = 0 \implies b_a(X_a) \propto \exp\left(-E_a(X_a) + \sum_{i \in N(a)} \lambda_{ai}(x_i)\right)$$

Eric Xing

17

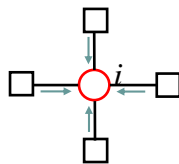
Bethe = BP on FG



- Identify

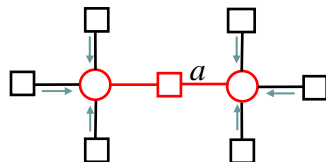
$$\lambda_{ai}(x_i) = \ln \prod_{b \in N(i) \neq a} m_{b \rightarrow i}(x_i)$$

- to obtain BP equations:



$$b_i(x_i) \propto f_i(x_i) \prod_{a \in N(i)} m_{a \rightarrow i}(x_i)$$

↑ "beliefs" ↑ "messages"



$$b_a(X_a) \propto f_a(X_a) \prod_{i \in N(a)} \prod_{c \in N(i) \setminus a} m_{c \rightarrow i}(x_i)$$

The "belief" is the BP approximation of the marginal probability.

Eric Xing

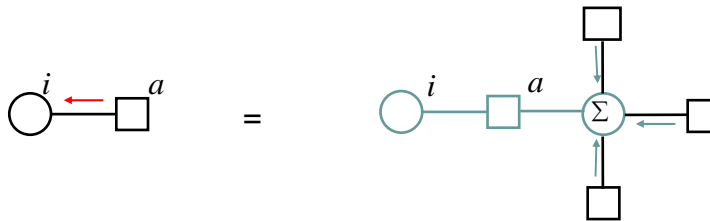
18

BP Message-update Rules

Using $b_{a \rightarrow i}(x_i) = \sum_{X_a \setminus x_i} b_a(X_a)$, we get

$$m_{a \rightarrow i}(x_i) = \sum_{X_a \setminus x_i} f_a(X_a) \prod_{j \in N(a) \setminus i} \prod_{b \in N(j) \setminus a} m_{b \rightarrow j}(x_j)$$

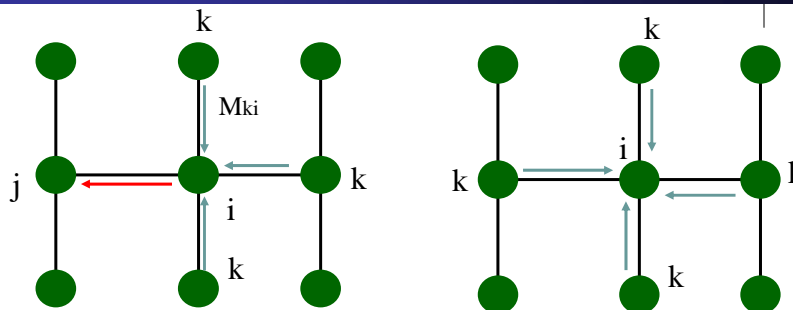
(A sum product algorithm)



Eric Xing

19

Belief Propagation on trees



- BP Message-update Rules

$$M_{i \rightarrow j}(x_j) \propto \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_k M_{k \rightarrow i}(x_i)$$

↑ Compatibilities (interactions)
 ↑ external evidence

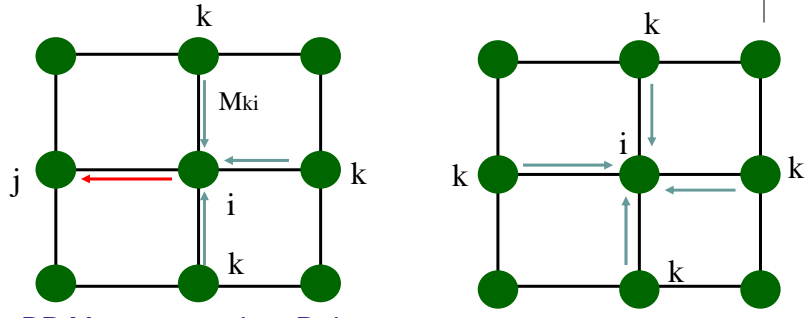
$$b_i(x_i) \propto \psi_i(x_i) \prod_k M_k(x_k)$$

- BP on trees always converges to exact marginals (cf. Junction tree algorithm)

Eric Xing

20

Belief Propagation on loopy graphs



- BP Message-update Rules

$$M_{i \rightarrow j}(x_j) \propto \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_k M_{k \rightarrow i}(x_i)$$

↑ Compatibilities (interactions)
 ↑ external evidence

$$b_i(x_i) \propto \psi_i(x_i) \prod_k M_k(x_k)$$

- May not converge or converge to a wrong solution

Eric Xing

21

Loopy Belief Propagation



- If BP is used on graphs with loops, messages may circulate indefinitely
- Empirically, a good approximation is still achievable
 - Stop after fixed # of iterations
 - Stop when no significant change in beliefs
 - If solution is not oscillatory but converges, it usually is a good approximation

Eric Xing

22

The Theory Behind LBP



- For a distribution $p(X/\theta)$ associated with a complex graph, computing the marginal (or conditional) probability of arbitrary random variable(s) is intractable
- Variational methods
 - formulating probabilistic inference as an optimization problem:

$$q^* = \arg \min_{q \in \mathcal{S}} \{ F_{\text{Betha}}(p, q) \}$$

$$F_{\text{Betha}} = \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln \frac{b_a(\mathbf{x}_a)}{f_a(\mathbf{x}_a)} + \sum_i (1 - d_i) \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \ln b_i(\mathbf{x}_i) = -\langle f_a(\mathbf{x}_a) \rangle - H_{\text{betha}}$$

q : a (tractable) probability distribution

Eric Xing

23

The Theory Behind LBP



- But we do not optimize $q(\mathbf{X})$ explicitly, focus on the set of beliefs
 - e.g., $b = \{b_{i,j} = \tau(x_i, x_j), b_i = \tau(x_i)\}$
- Relax the optimization problem
 - approximate objective: $H_{\text{Betha}} = H(b_{i,j}, b_i)$
 - relaxed feasible set: $\mathcal{M}_o = \{ \tau \geq 0 \mid \sum_{x_i} \tau(x_i) = 1, \sum_{x_i} \tau(x_i, x_j) = \tau(x_j) \}$

$$b^* = \arg \min_{b \in \mathcal{M}_o} \{ \langle E \rangle_b + F(b) \}$$

- The loopy BP algorithm:
 - a fixed point iteration procedure that tries to solve b^*

Eric Xing

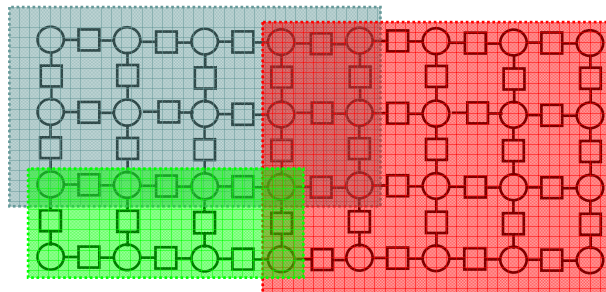
24

Region-based Approximations to the Gibbs Free Energy (Kikuchi, 1951)



Exact: $G[q(X)]$ (intractable)

Regions: $G[\{b_r(X_r)\}]$



Eric Xing

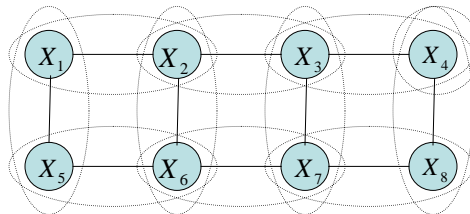
25

Better approximations?



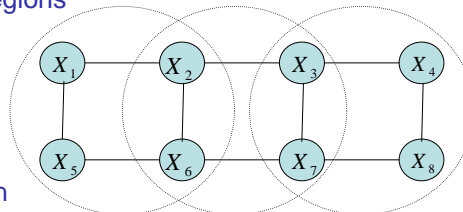
- Recall that Bethe approximation was

$$F_{\text{bethe}} = F_{12} + F_{23} + \dots + F_{67} + F_{78} - F_1 - F_5 - 2F_2 - 2F_6 \dots - F_8$$



- We could construct bigger regions

$$\hat{F} = F_{1256} + F_{2367} + F_{3478} - F_{26} - F_{36} + F_2 + F_6 + F_3 + F_7$$



- Called Kikuchi approximation

Eric Xing

26

Recipe for generalizing BP

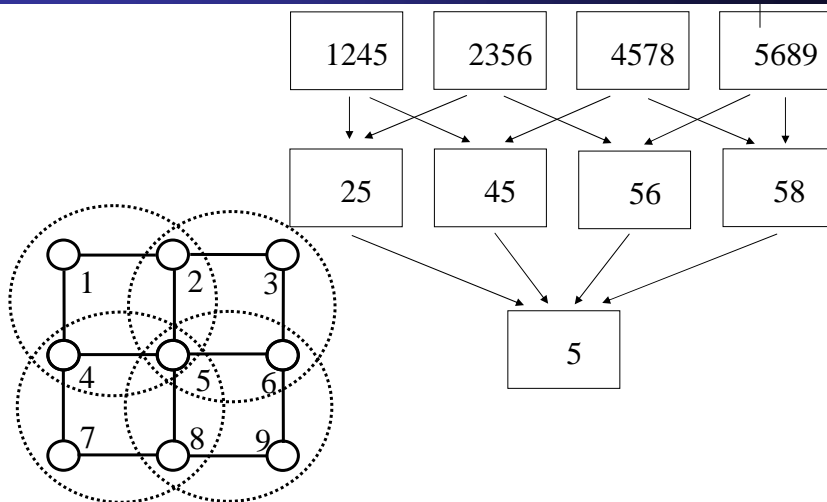


- Define an approximation of the free energy
 - Some restrictions on choice of regions
- Obtain fixed point equations by solving constrained optimization problem using this approximation
- Convert to a message passing algorithm ala BP
 - Message-update rules obtained by enforcing marginalization constraints.
- Belief in a region is the product of:
 - Local information (factors in region)
 - Messages from parent regions
 - Messages into descendant regions from parents who are not descendants.
- Called Generalized BP or GBP

Eric Xing

27

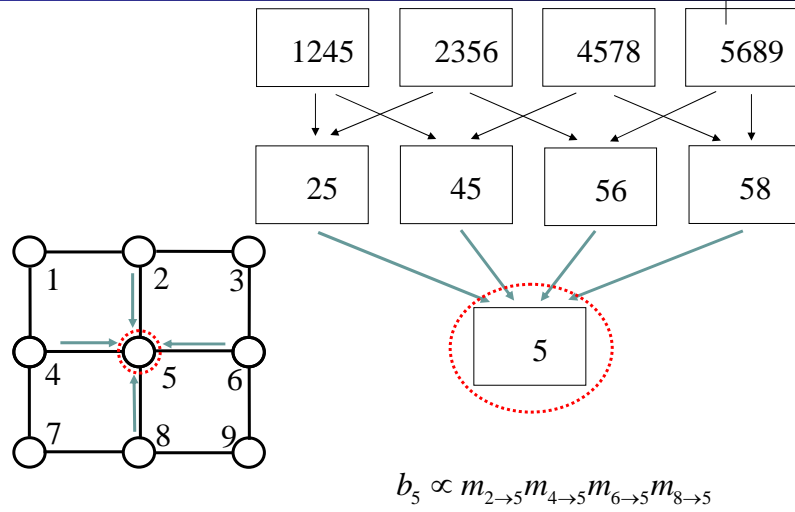
Generalized Belief Propagation



Eric Xing

28

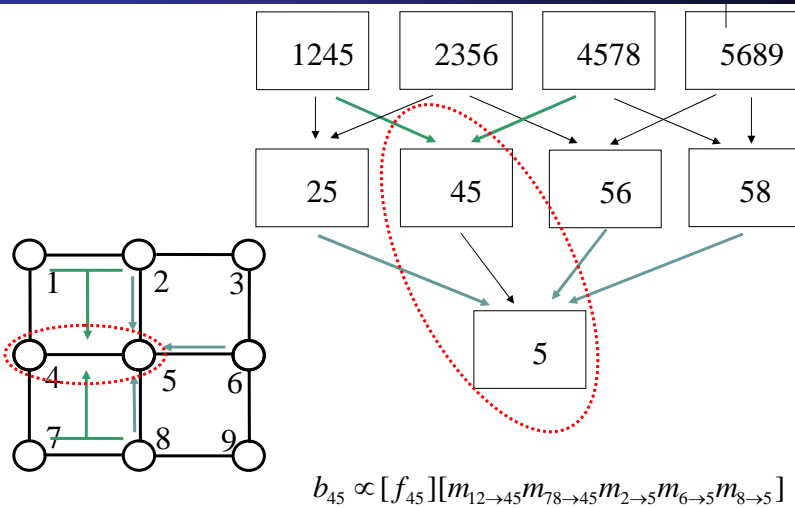
Generalized Belief Propagation



Eric Xing

29

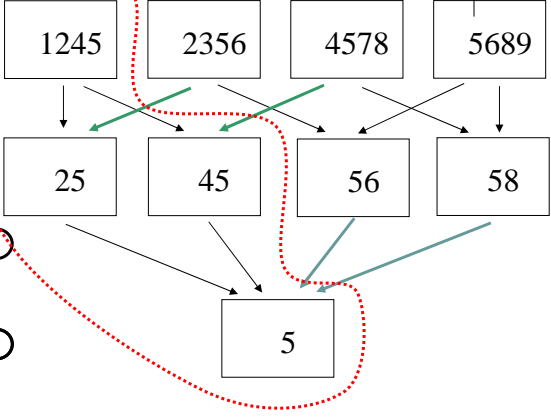
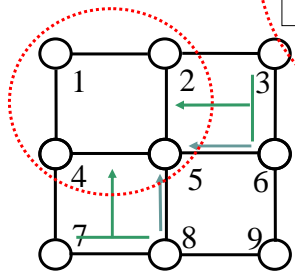
Generalized Belief Propagation



Eric Xing

30

Generalized Belief Propagation



$$b_{1245} \propto [f_{12} f_{14} f_{25} f_{45}] [m_{36 \rightarrow 25} m_{78 \rightarrow 45} m_{6 \rightarrow 5} m_{8 \rightarrow 5}]$$