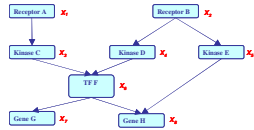


# Approximate Inference: Loopy Belief Propagation and variants

Probabilistic Graphical Models (10-708)

Lecture 14, Nov 7, 2007

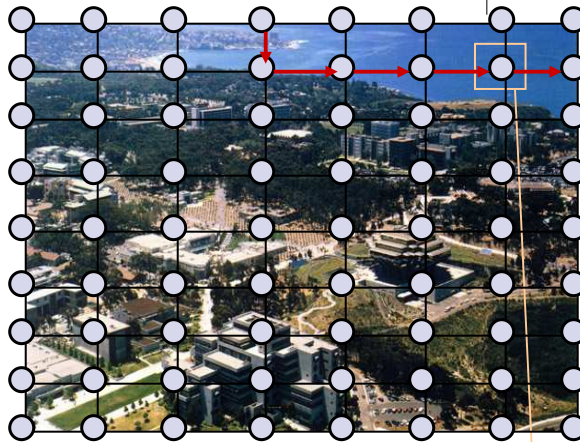


Hetunandan Kamisetty

Reading: KF-Chap. 12, Yedidia et al

## An Ising model on 2-D image

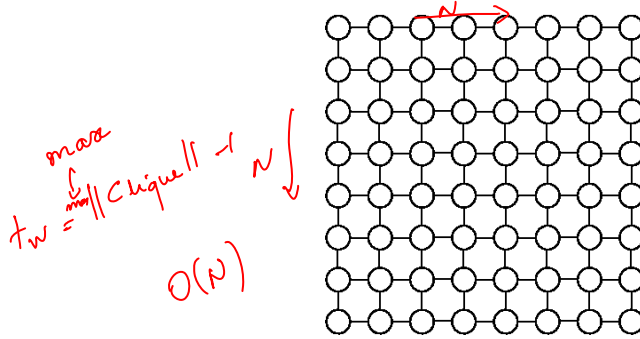
- Nodes encode hidden information (patch-identity).
- They receive local information from the image (brightness, color).
- Information is propagated through the graph over its edges.
- Edges encode 'compatibility' between nodes.



air or water ?

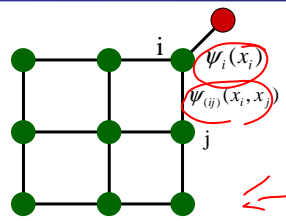


# Why Approximate Inference?



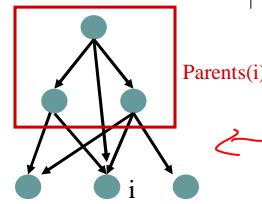
- Why can't we just run junction tree on this graph?
- If  $N \times N$  grid, tree width atleast  $N$ 
  - If  $N \sim O(1000)$ , we have a clique with  $2^{1000}$  entries

# A recap: Factor Graphs



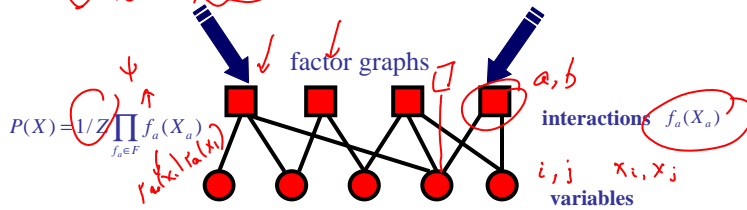
Undirected graph  
(Markov random field)

$$P(x) = \frac{1}{Z} \prod_i \psi_i(x_i) \prod_{(ij)} \psi_{(ij)}(x_i, x_j)$$



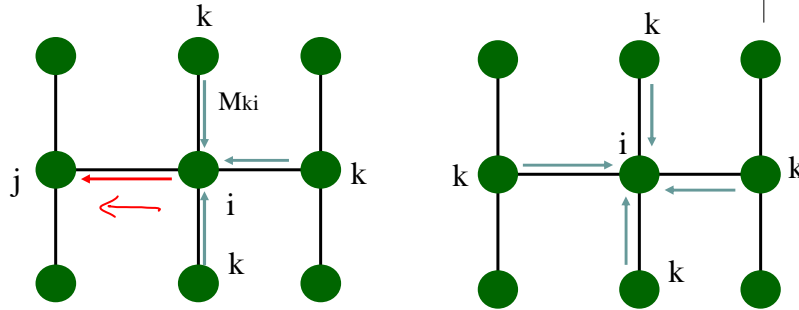
Directed graph  
(Bayesian network)

$$P(x) = \prod_i P(x_i | x_{\text{parents}(i)})$$



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

# Belief Propagation



- BP Message-update Rules

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

external evidence  
Compatibilities (interactions)

$$b_i(x_i)?$$

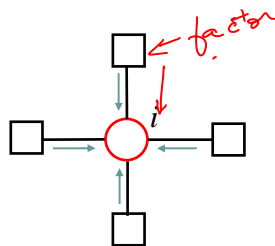
$$\propto \prod_{j \in N(i)} M_{j \rightarrow i}(x_i)$$

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

Eric Xing

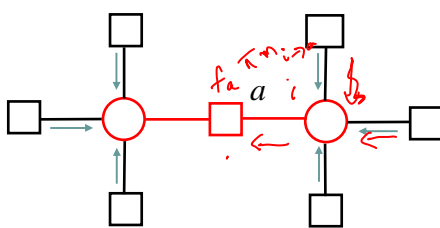
5

# 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"



$$m_{i \rightarrow a}(x_a) = \prod_{c \in N(i), c \neq a} m_{c \rightarrow i}(x_i)$$

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

$$m_{a \rightarrow i}(x_i) = ?$$

Eric Xing

6

# Approximate Inference: What to approximate?



- Let us call the actual distribution  $P$

$$P(X) = 1/Z \prod_{f_a \in F} f_a(X_a) \quad \leftarrow \text{factor graph}$$

- We wish to find a distribution  $Q$  such that  $Q$  is a "good" approximation
- Recall the definition of KL-divergence

$$KL(Q_1 \parallel Q_2) = \int Q_1(x) \log \frac{Q_1(x)}{Q_2(x)} = E_{Q_1} \log \frac{Q_1}{Q_2}$$

- $KL(Q_1 \parallel Q_2) \geq 0$
- $KL(Q_1 \parallel Q_2) = 0$  iff  $Q_1 = Q_2$
- We can therefore use KL as a scoring function to decide a good  $Q$  not sym!
- But,  $KL(Q_1 \parallel Q_2) \neq KL(Q_2 \parallel Q_1)$   $KL(Q_1 \parallel Q_2) \neq KL(Q_2 \parallel Q_1)$

# Which KL?



- Computing  $KL(P \parallel Q)$  requires inference!  $= E_P \log \frac{P}{Q}$
- But  $KL(Q \parallel P)$  can be computed without performing inference on  $P$

$$KL(Q \parallel P) = \sum_x Q(x) \log \left( \frac{Q(x)}{P(x)} \right) = \frac{\sum Q(x) \log Q(x)}{-\sum Q(x) \log P(x)}$$

$$= -H_Q(x) - E_{Q(x)} \log(P(x))$$

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

$$KL(Q \parallel P) = -H_Q(X) - E_Q \log(1/Z \prod_{f_a \in F} f_a(X_a))$$

$$= -H_Q(x) - \log(1/Z) - \sum_a E_Q \log f_a(x_a)$$

# Optimization function



$$KL(Q \parallel P) = -H_Q(X) - \sum_{f_a \in F} E_Q \log f_a(X_a) + \log Z_P$$

*Handwritten notes:*  
 - A red arrow points to the first term:  $-H_Q(X)$   
 - A red arrow points to the second term:  $\sum_{f_a \in F} E_Q \log f_a(X_a)$   
 - A red arrow points to the third term:  $\log Z_P$   
 - A red box is drawn around  $\log Z_P$  with the note " $\leftarrow \text{const}$ ".  
 - A red arrow points to the entire expression with the note " $F(P, Q) \leftarrow$ ".  
 - A red note says " $KL(P, P) = 0$ ".  
 - A red note says " $Q \leftarrow Z_Q$ ".  
 - A red note says " $\text{Textbook } -F(P, Q) \uparrow$ ".

- We will call  $F(P, Q)$  the "Free energy" \*
- $F(P, P) = ? - \log Z$
- $F(P, Q) \geq F(P, P)$

# The Free Energy



- Let us look at the free energy

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

*Handwritten notes:*  
 - A red note says " $\text{approx } \log Z \leftarrow Q(X)$ ".  
 - A red note says " $\text{marginals}$ " with an arrow pointing to the equation.

- $\sum_{f_a \in F} E_Q \log f_a(X_a)$  can be computed if we have marginals over each  $f_a$

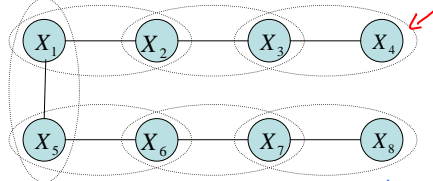
- $H_Q = \sum_X Q(X) \log Q(X)$  is harder! Requires summation over all possible values

- Computing  $F$ , is therefore hard in general.

- Approach 1: Approximate  $F(P, Q)$  with easy to compute  $\hat{F}(P, Q)$

# Easy free energies

- Consider a tree-structured distribution



$\emptyset \rightarrow \textcircled{X_1} \rightarrow \textcircled{X_3}$   
 $b(x) = \prod_a b(x_a) \prod_i b(x_i)^{d_i-1}$   
 $\frac{P(x_1, x_2, x_3)}{P(x_2, x_3)} = \frac{P(x_1, x_2) P(x_2, x_3)}{P(x_2, x_3)}$

- The probability can be written as:  $b(x) = \prod_a b(x_a) \prod_i b(x_i)^{d_i-1}$

$$H_{tree} = -\sum_a \sum_{x_a} b_a(x_a) \log b_a(x_a) + \sum_i (d_i - 1) \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

- $$F_{Tree} = \sum_a \sum_{x_a} b_a(x_a) \log \frac{b_a(x_a)}{f_a(x_a)} + \sum_i (1 - d_i) \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

$$= \uparrow F_{15} + F_{12} + F_{23} + F_{34} \dots - F_3 - F_2 - F_1 \dots$$
- involves summation over edges and vertices and is therefore easy to compute

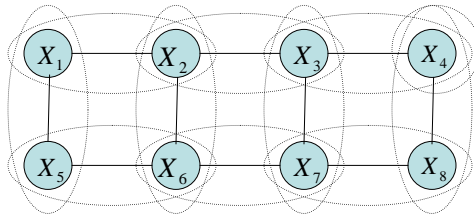
# Bethe Approximation to Gibbs Free Energy

- For a general graph, choose

$$F(P, Q) = \sum_a \sum_{x_a} b_a(x_a) \log \frac{b_a(x_a)}{f_a(x_a)} + \sum_i (1 - d_i) \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

- Called "Bethe approximation" after the physicist Hans Bethe
- Equal to the exact Gibbs free energy when the factor graph is a tree
- Note: This is **not** the same as the entropy of a tree

← true general graph



$f_i = \sum_x q(x) \log q(x)$   
 $\approx q_{12} + q_{23} \dots - 2q_2 \dots$

$$F_{bethe} = F_{15} + F_{12} + F_{23} \dots - F_1 - 2F_2 \dots - (F_4) - F_{24}$$

$$= F_{BETHE}$$

# Bethe Approximation



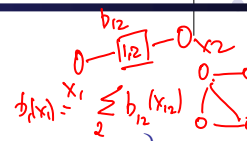
- Pros:
  - Easy to compute, since entropy term involves sum over pairwise and single variables
- Cons:
  - $\hat{F}(P, Q) = F_{\text{bethe}}$  **may or may not** be well connected to  $F(P, Q)$
  - It could, in general, be greater, equal or less than  $F(P, Q)$
- Optimize each  $b(x_a)$ 's.
  - For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

$|F - \hat{F}| < \epsilon$   
 $F > \hat{F}$   
 $F < \hat{F}$

# Minimizing the Bethe Free Energy



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



- Set derivative to zero

$$F_{\text{Bethe}} = \sum_a \sum_{x_a} b_a(x_a) \log \frac{b_a(x_a)}{f_a(x_a)} - \sum_i (1 - d_i) \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

$$\frac{\partial F_b}{\partial b_a(x_a)} = \frac{\partial}{\partial b_a(x_a)} \left[ b_a(x_a) \log \frac{b_a(x_a)}{f_a(x_a)} - \sum_{i \in N(a)} b_i(x_i) \log b_i(x_i) \right]$$

$$= 1 + \log \frac{b_a(x_a)}{f_a(x_a)} - \log f_a - \sum_{i \in N(a)} \frac{\partial b_i(x_i)}{\partial b_a(x_a)} b_i(x_i) \log b_i(x_i)$$





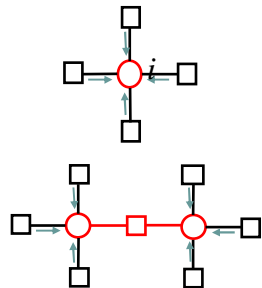
$$b_i(x_i) \propto \exp\left(\frac{1}{d_i-1} \sum_{a \in N(i)} \log m_{i \rightarrow a}\right)$$

## Bethe = BP

- We had

$$b_i(x_i) \propto \exp\left(\frac{1}{d_i-1} \sum_{a \in N(i)} \lambda_{ai}(x_i)\right) \quad b_a(X_a) \propto \exp\left(-\log f_a(X_a) + \sum_{i \in N(a)} \lambda_{ai}(x_i)\right)$$

- Identify  $\lambda_{ai}(x_i) = \log(m_{i \rightarrow a}(x_i)) = \log \prod_{b \in N(i) \neq a} m_{b \rightarrow i}(x_i)$
- to obtain BP equations:



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

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

Eric Xing

17

## Loopy Belief Propagation

- A fixed point iteration procedure that tries to minimize  $F_{\text{bethe}}$
- Start with random initialization of messages and beliefs

- While not converged do

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

$$m_{i \rightarrow a}^{\text{new}}(x_i) = \prod_{c \in N(i) \setminus a} m_{c \rightarrow i}(x_i) \quad m_{a \rightarrow i}^{\text{new}}(x_i) = \sum_{X_a \setminus x_i} f_a(X_a) \prod_{j \in N(a) \setminus i} m_{j \rightarrow a}(x_j)$$

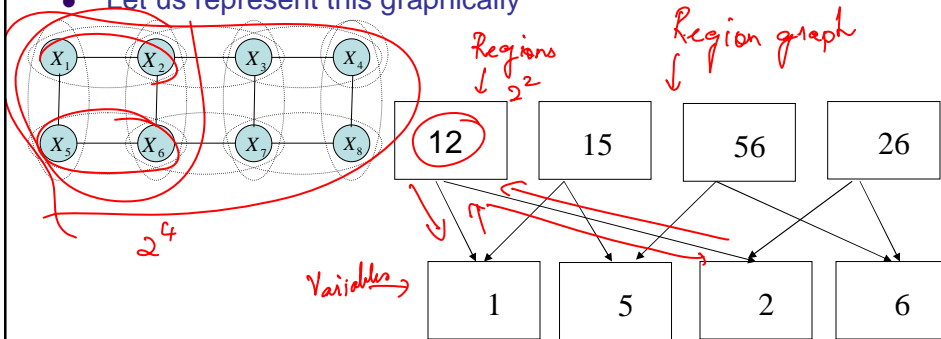
- At convergence, stationary <sup>new</sup> properties are guaranteed <sup>new</sup> Why?
- However, not guaranteed to converge!

Eric Xing

18

# Region graphs

- It will be useful to look explicitly at the messages being passed
  - Messages from variable to factors
  - Messages from factors to variables
- Let us represent this graphically

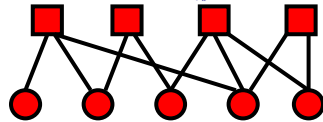


Eric Xing

19

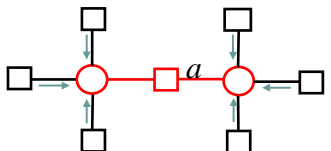
# Summary so far

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



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

$$\hat{F}(P, Q) = \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \log \frac{f_a(\mathbf{x}_a)}{b_a(\mathbf{x}_a)} + \sum_i (1-d_i) \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \log b_i(\mathbf{x}_i)$$



$$b_a(X_a) \propto \exp \left( -\log f_a(X_a) + \sum_{i \in N(a)} \lambda_{ai}(x_i) \right)$$

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

Eric Xing

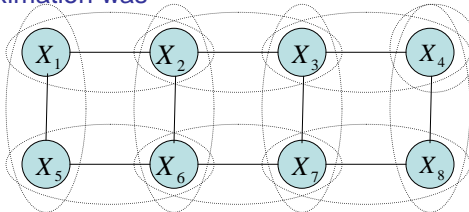
20

# Better approximations?

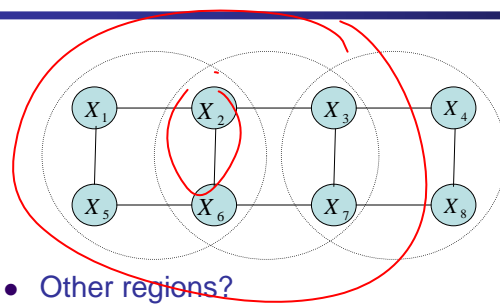


- Recall that Bethe approximation was

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



- We could construct bigger regions
- Rules:
  - If a region includes a factor, it must include the vertices as well
  - Each factor and vertex must appear in atleast one region
  - Associate a weight with each region so that each vertex and factor is counted exactly once



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

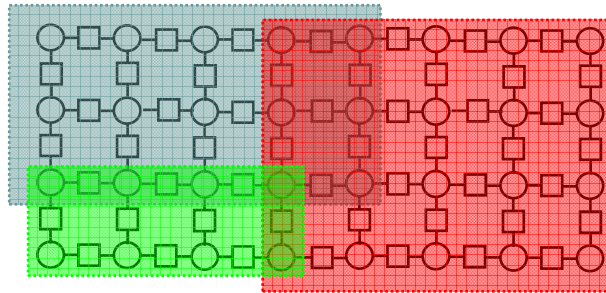
- Other regions?

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



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

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



Eric Xing

23

# Generalized Belief Propagation



- 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.
- Message-update rules obtained by enforcing marginalization constraints.

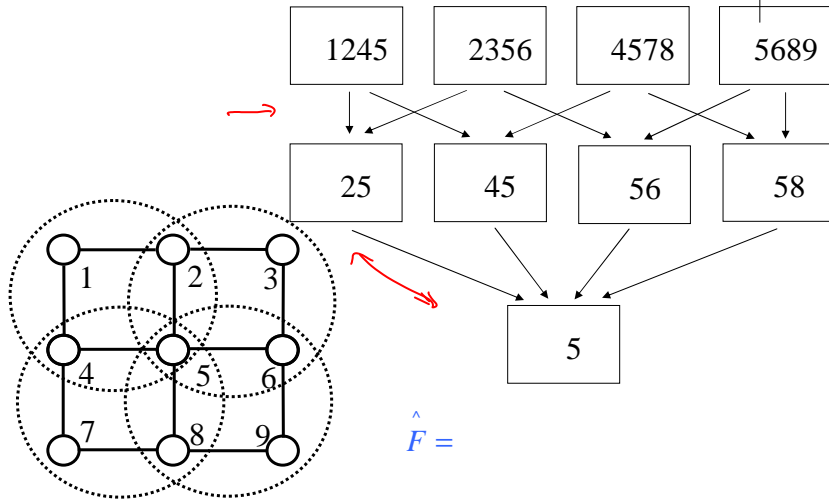
$$F \sim \hat{F} \rightsquigarrow L$$

$$\frac{\partial L}{\partial b_a} \rightarrow \text{terms}$$

Eric Xing

24

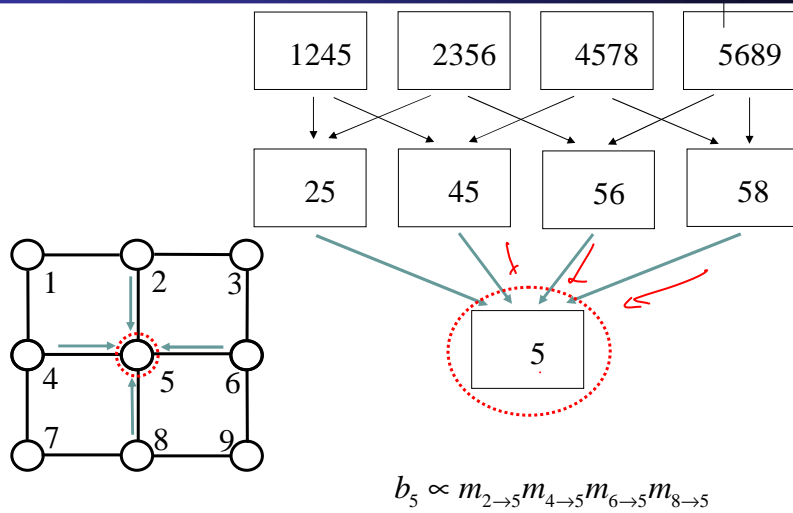
# Generalized Belief Propagation



Eric Xing

25

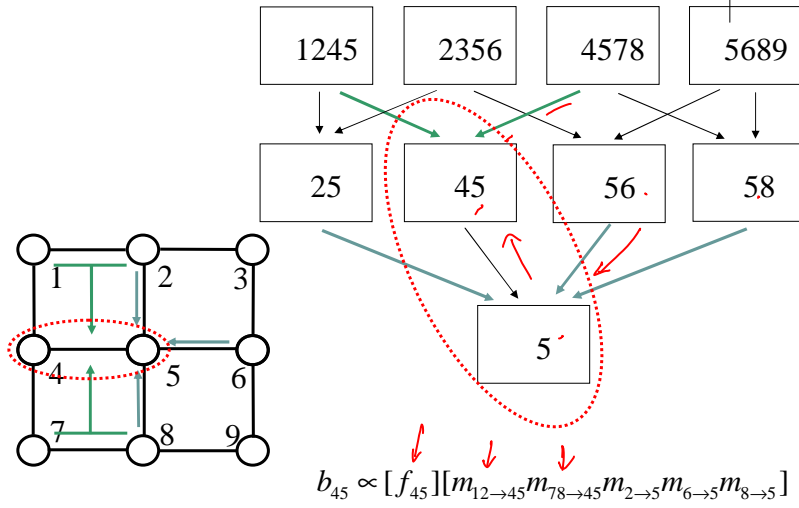
# Generalized Belief Propagation



Eric Xing

26

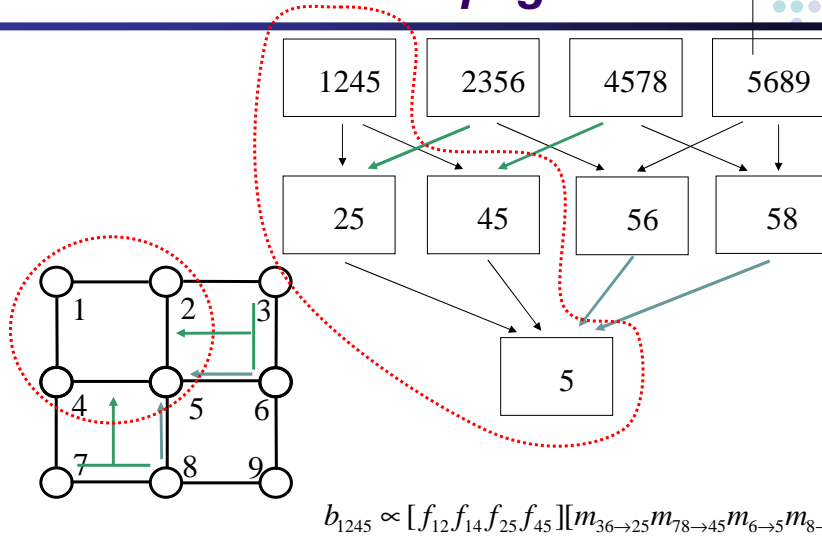
# Generalized Belief Propagation



Eric Xing

27

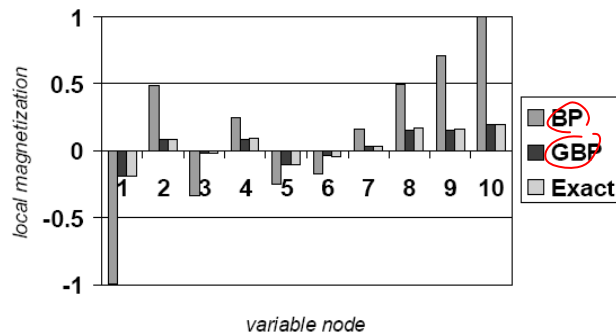
# Generalized Belief Propagation



Eric Xing

28

## Some results



Eric Xing

29

## Summary

- We defined an objective function ( $F$ ) for approximate inference
- However, we found that optimizing this function was hard
- We first approximated objective function  $F$  to simpler  $F_{\text{bethe}}$ 
  - Minima of  $F_{\text{bethe}}$  turned out to be fixed points of BP
- Then we extended this to more complicated approximations
  - The resulting algorithms come under a family called Generalized Belief Propagation
- Next class, we will cover other methods of approximations

Eric Xing

30