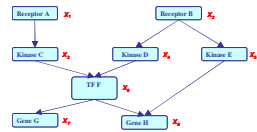


Approximate Inference: Loopy Belief Propagation and variants

Probabilistic Graphical Models (10-708)

Lecture 14, Nov 7, 2007

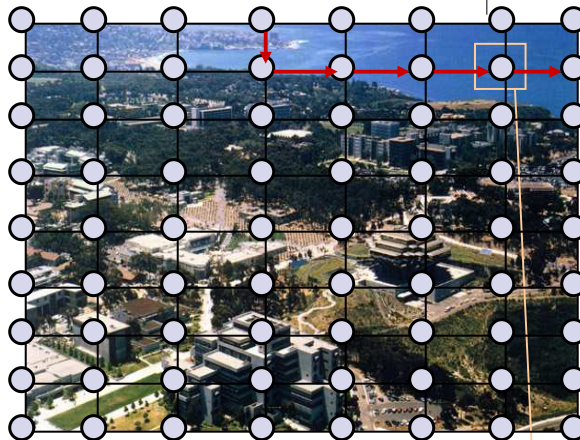


Hetunandan Kamisetty

Reading: J-Chap. 5,6, KF-Chap. 8

An Ising model on 2-D image

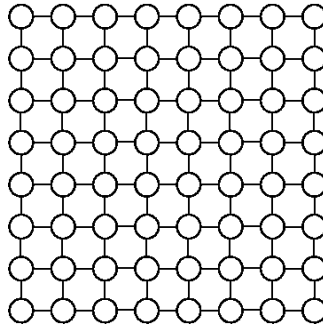
- Nodes encode hidden information (patch-identity).
- They receive local information from the image (brightness, color).
- Information is propagated through the graph over its edges.
- Edges encode 'compatibility' between nodes.



air or water ?



Why Approximate Inference?

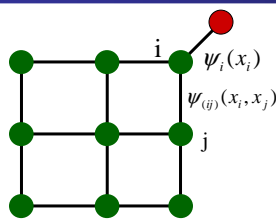


- Why can't we just run junction tree on this graph?
- If NxN grid, tree width atleast N
 - If N~O(1000), we have a clique with 2¹⁰⁰⁰ entries

Eric Xing

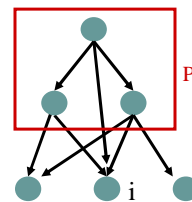
3

A recap: Factor Graphs



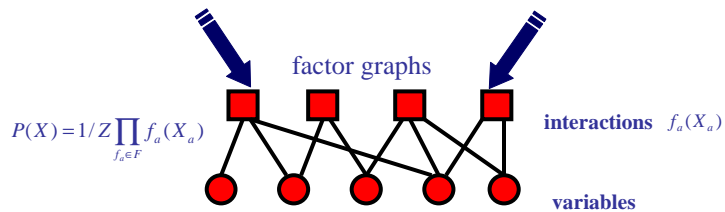
Undirected graph
(Markov random field)

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



Directed graph
(Bayesian network)

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

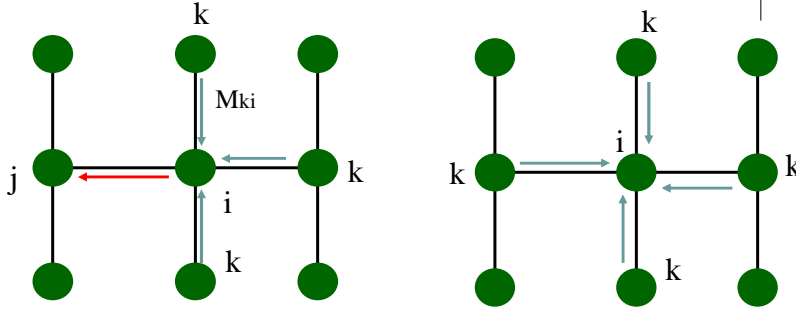


$$P(X) = 1/Z \prod_{f_a \in F} f_a(X_a)$$

Eric Xing

4

Belief Propagation



- BP Message-update Rules

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

↑ Compatibilities (interactions)
 ↑ external evidence

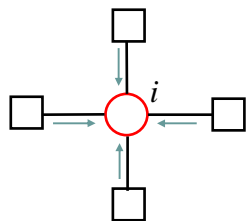
$$b_i(x_i) \propto \psi_i(x_i) \prod_k M_k(x_k)$$

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

Eric Xing

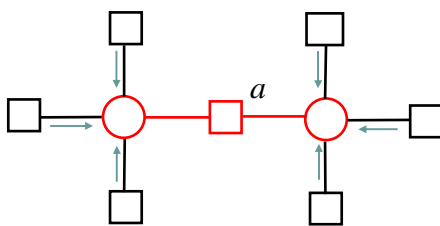
5

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"



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

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

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

Eric Xing

6

Approximate Inference: What to approximate?



- Let us call the actual distribution P

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

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

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

- $KL(Q_1 \parallel Q_2) \geq 0$
- $KL(Q_1 \parallel Q_2) = 0$ iff $Q_1 = Q_2$
- We can therefore use KL as a scoring function to decide a good Q
- But, $KL(Q_1 \parallel Q_2) \neq KL(Q_2 \parallel Q_1)$

Eric Xing

7

Which KL?



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

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

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

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

Eric Xing

8

Optimization function



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

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

The Free Energy



- Let us look at the free energy

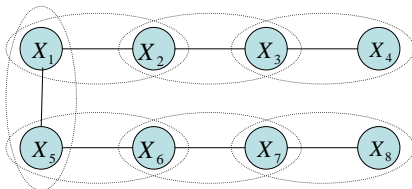
$$F(P, Q) = -H_Q(X) - \sum_{f_a \in F} E_Q \log f_a(X_a)$$

- $\sum_{f_a \in F} E_Q \log f_a(X_a)$ 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



- The probability can be written as: $b(\mathbf{x}) = \prod_a b_a(\mathbf{x}_a) \prod_i b_i(\mathbf{x}_i)^{1-d_i}$
- $H_{tree} = -\sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \log b_a(\mathbf{x}_a) + \sum_i (d_i - 1) \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \log b_i(\mathbf{x}_i)$
- $F_{Tree} = \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \log \frac{b_a(\mathbf{x}_a)}{f_a(\mathbf{x}_a)} + \sum_i (1 - d_i) \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \log b_i(\mathbf{x}_i)$
 $= F_{12} + F_{23} + \dots + F_{67} + F_{78} - F_1 - F_5 - F_2 - F_6 - F_3 - F_7$
 - involves summation over edges and vertices and is therefore easy to compute

Eric Xing

11

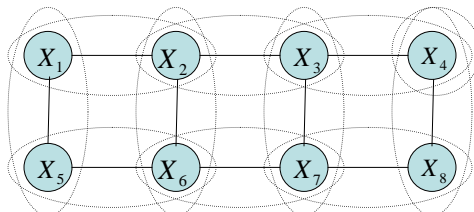
Bethe Approximation to Gibbs Free Energy



- For a general graph, choose

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

- 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



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

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 $b(x_a)$'s.
 - For discrete belief, constrained opt. with *Lagrangian* multiplier
 - For continuous belief, not yet a general formula
 - Not always converge

Eric Xing

13

Minimizing the Bethe Free Energy



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

Eric Xing

14

Proof



-

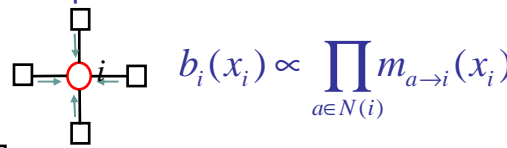


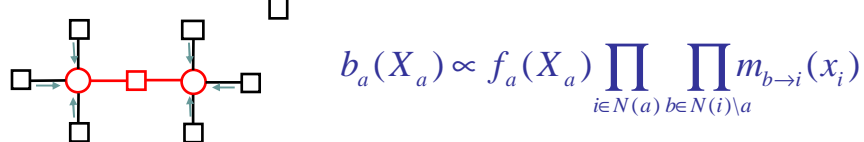
Bethe = BP

- 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:



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


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

Eric Xing

17

Loopy Belief Propagation

- A fixed point iteration procedure that tries to minimize F_{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, stationarity properties are guaranteed
- However, not guaranteed to converge!

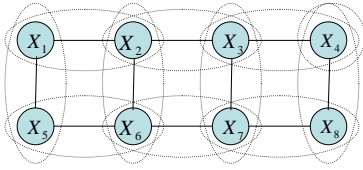
Eric Xing

18

Region graphs



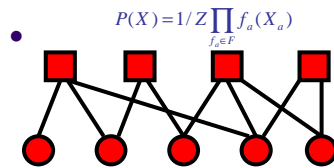
- It will be useful to look explicitly at the messages being passed
 - Messages from variable to factors
 - Messages from factors to variables
- Let us represent this graphically



Eric Xing

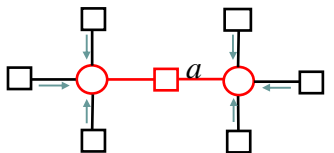
19

Summary so far



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

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



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

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

Eric Xing

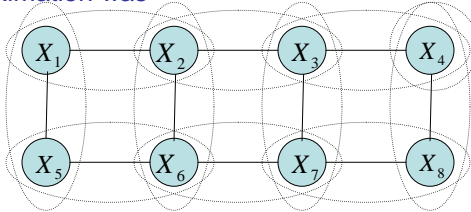
20

Better approximations?

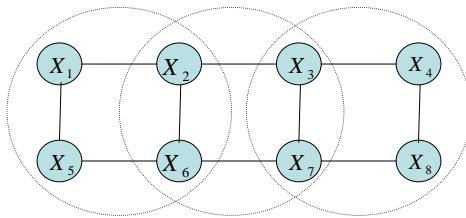


- Recall that Bethe approximation was

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



- 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



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

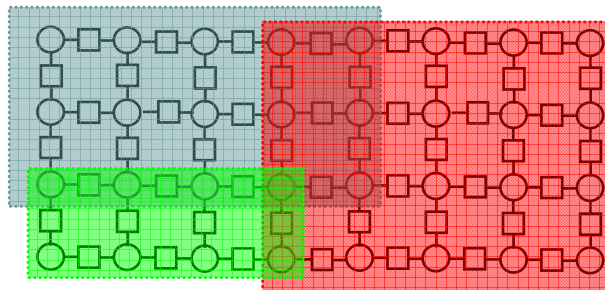
- Other regions?

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



Exact: $\mathcal{G}[q(X)]$ (intractable)

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



Eric Xing

23

Generalized Belief Propagation

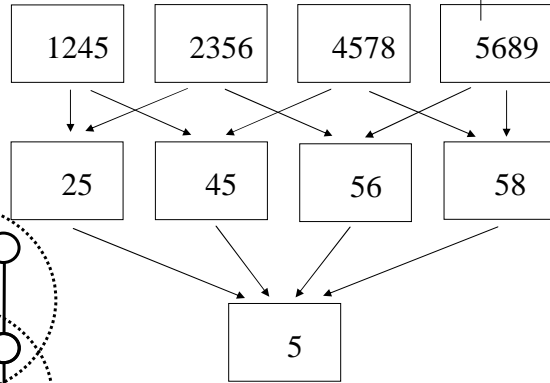
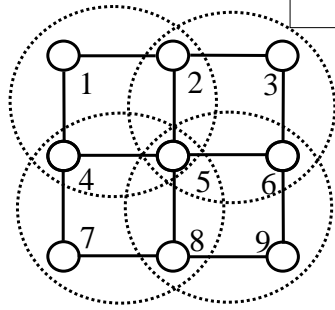


- Belief in a region is the product of:
 - Local information (factors in region)
 - Messages from parent regions
 - Messages into descendant regions from parents who are not descendants.
- Message-update rules obtained by enforcing marginalization constraints.

Eric Xing

24

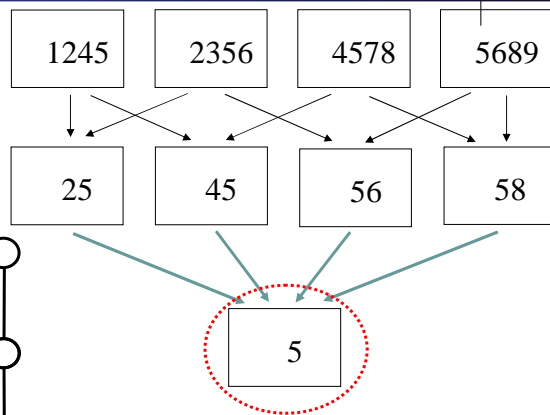
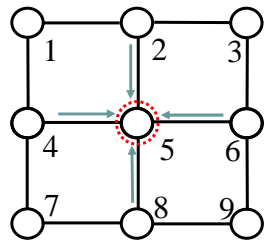
Generalized Belief Propagation



Eric Xing

25

Generalized Belief Propagation

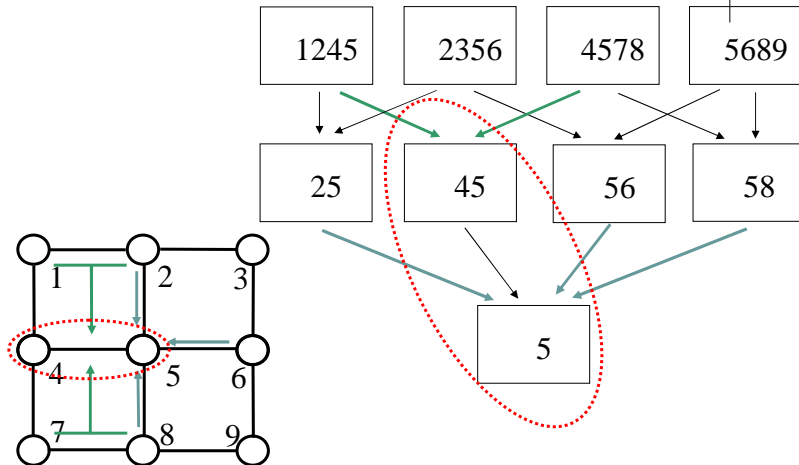


$$b_5 \propto m_{2 \rightarrow 5} m_{4 \rightarrow 5} m_{6 \rightarrow 5} m_{8 \rightarrow 5}$$

Eric Xing

26

Generalized Belief Propagation

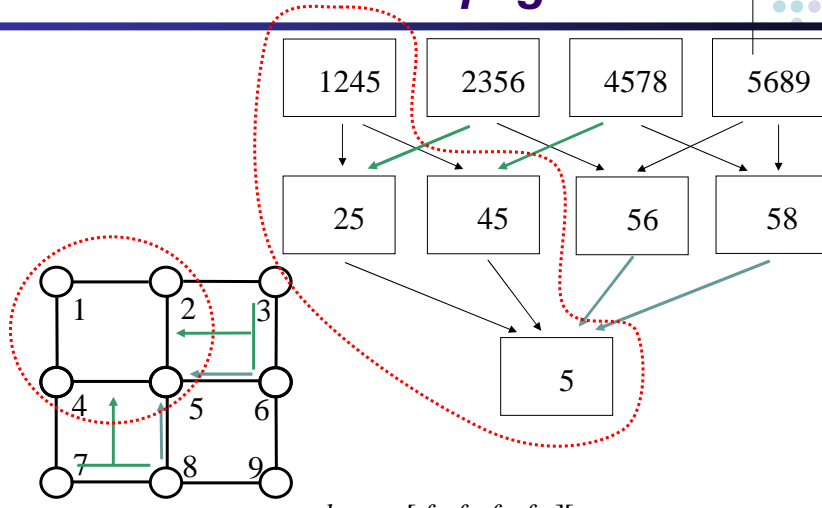


$$b_{45} \propto [f_{45}][m_{12 \rightarrow 45} m_{78 \rightarrow 45} m_{2 \rightarrow 5} m_{6 \rightarrow 5} m_{8 \rightarrow 5}]$$

Eric Xing

27

Generalized Belief Propagation



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

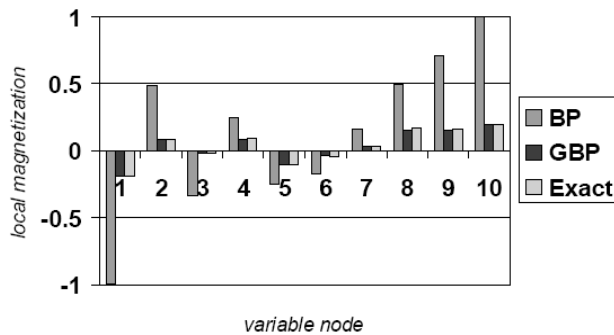
Eric Xing

28

Some results



-



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_{bethe}
 - Minima of F_{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