

Readings:
K&F: 10.2, 10.3

Generalized Belief Propagation

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 12th, 2008

10-708 – Carlos Guestrin 2006-2008

1

More details on Loopy BP

Numerical problem:

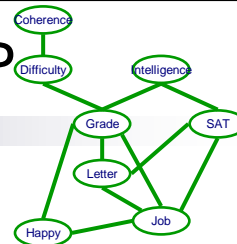
- messages < 1 get multiplied together as we go around the loops
- numbers can go to zero
- normalize messages to one:

$$\delta_{i \rightarrow j}(X_j) = \frac{1}{Z_{i \rightarrow j}} \sum_{x_i} \phi_i(x_i) \phi_{ij}(x_i, X_j) \prod_{k \in \mathcal{N}(i) - j} \delta_{k \rightarrow i}(x_i)$$

- $Z_{i \rightarrow j}$ doesn't depend on X_j , so doesn't change the answer

Computing node "beliefs" (estimates of probs.):

$$b_i(x_i) = \hat{P}(X_i) = \frac{1}{Z_i} \phi_i(X_i) \prod_{k \in \mathcal{N}(i)} \delta_{k \rightarrow i}(X_i)$$



sometimes important to compute in log space

$$\log \delta_{i \rightarrow j}(x_i) = \log$$

10-708 – Carlos Guestrin 2006-2008

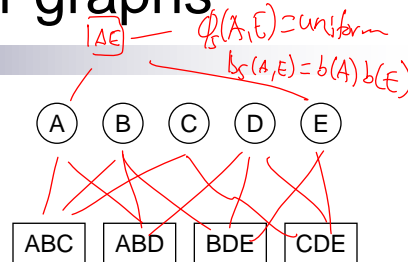
2

Loopy BP in Factor graphs

From node i to factor j :

- $F(i)$ factors whose scope includes X_i

$$\delta_{i \rightarrow j}(X_i) \propto \prod_{k \in F(i) - j} \delta_{k \rightarrow i}(X_i)$$



From factor j to node i :

- $\text{Scope}[\phi_j] = Y \cup \{X_i\}$

$$\delta_{j \rightarrow i}(X_i) \propto \sum_{\underline{y}} \phi_j(X_i, \underline{y}) \prod_{X_k \in \text{Scope}[\phi_j] - X_i} \delta_{k \rightarrow j}(x_k)$$

Belief:

- Node:

$$P(x_i) \approx b_i(x_i) \propto \prod_{\phi_j, x_i \in \text{Scope}[\phi_j]} \delta_{j \rightarrow x_i}(x_i)$$

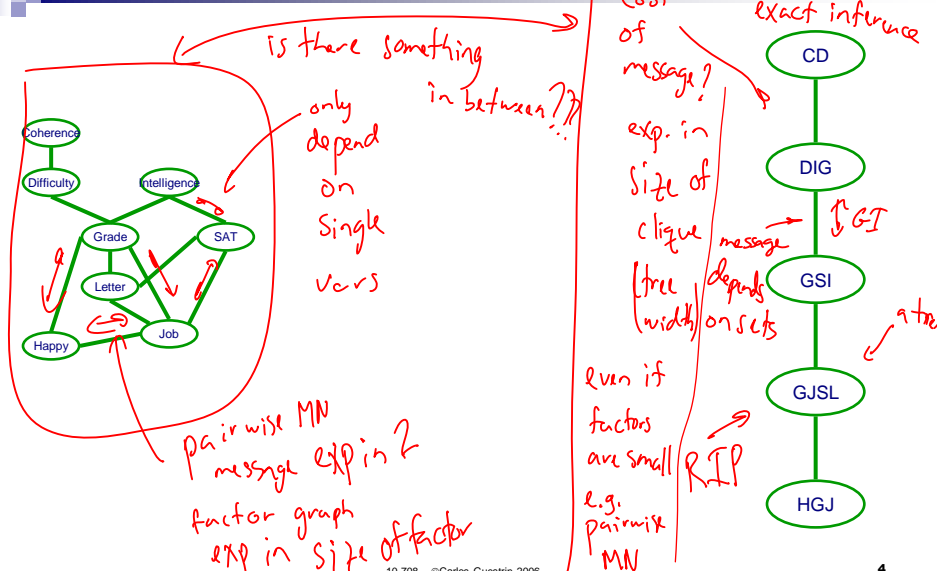
- Factor:

$$P(Y) \approx b_j(Y) \propto \phi_j(Y) \prod_{x_i \in Y} \delta_{x_i \rightarrow \phi_j}(x_i)$$

10-708 © Carlos Guestrin 2006-2008

3

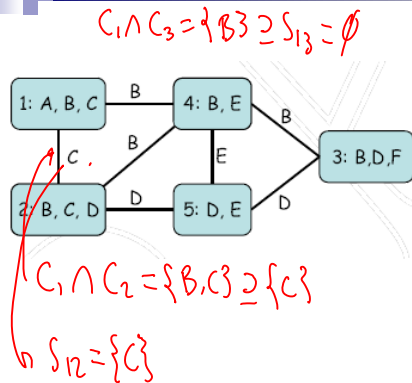
Loopy BP v. Clique trees: Two ends of a spectrum



10-708 © Carlos Guestrin 2006

4

Generalize cluster graph



Generalized cluster graph:

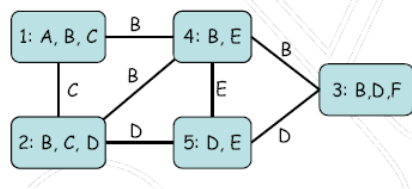
For set of factors F

- Undirected graph *not a tree*
- Each node i associated with a cluster C_i
- *Family preserving*: for each factor $f_j \in F$, \exists node i such that $\text{scope}[f_j] \subseteq C_i$
- Each edge $i - j$ is associated with a set of variables
 $S_{ij} \subseteq C_i \cap C_j$

10/708 - ©Carlos Guestrin 2006

5

Running intersection property



(Generalized) Running intersection property (RIP)

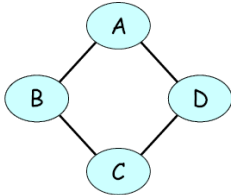
- Cluster graph satisfies RIP if whenever $X \in C_i$ and $X \in C_j$ then *one and only one path* from C_i to C_j where $X \in S_{uv}$ for every edge (u,v) in the path

10/708 - ©Carlos Guestrin 2006

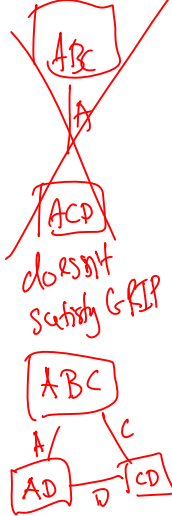
6

Examples of cluster graphs

pairwise MN :



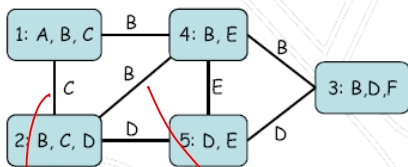
some cluster graph that satisfy RIP



10/708 - ©Carlos Guestrin 2006

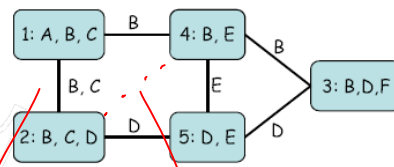
7

Two cluster graph satisfying RIP with different edge sets



$S_{1,2} = \{C\}$

$S_{2,4} = \{B\}$



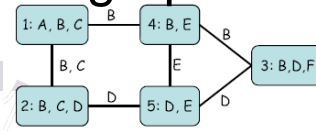
$S_{1,2} = \{BC\}$

$S_{2,4} = \emptyset$

10/708 - ©Carlos Guestrin 2006

8

Generalized BP on cluster graphs satisfying RIP



Initialization:

- Assign each factor ϕ to a clique $\alpha(\phi)$, $\text{Scope}[\phi] \subseteq \mathbf{C}_{\alpha(\phi)}$
- Initialize cliques: $\psi_i^0(\mathbf{C}_i) \propto \prod_{\phi: \alpha(\phi)=i} \phi$
- Initialize messages: $\delta_{j \rightarrow i} = 1$

While not converged, send messages:

$$\delta_{i \rightarrow j}(\mathbf{S}_{ij}) \propto \sum_{\mathbf{C}_i - \mathbf{S}_{ij}} \psi_i^0(\mathbf{C}_i) \prod_{k \in \mathcal{N}(i) - j} \delta_{k \rightarrow i}(\mathbf{S}_{ik})$$

normalize messages
marginalize

Belief:

or $P(\mathbf{C}_i | \mathbf{e})$

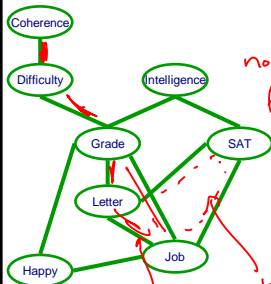
$$P(\mathbf{C}_i) \approx b_i(\mathbf{C}_i) \propto \psi_i^0(\mathbf{C}_i) \prod_{k \in \mathcal{N}(i)} \delta_{k \rightarrow i}(\mathbf{S}_{ik})$$

clique potential
all incoming messages from other nodes

10-708 - ©Carlos Guestrin 2006

9

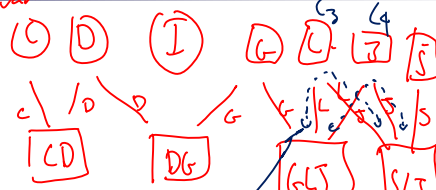
Cluster graph for Loopy BP



node per var

node per factor

triple factor



message and about L or J independently

share L, J

no joint message about L & J

special case of cluster graph w. RIP

because each var forms a "tree"

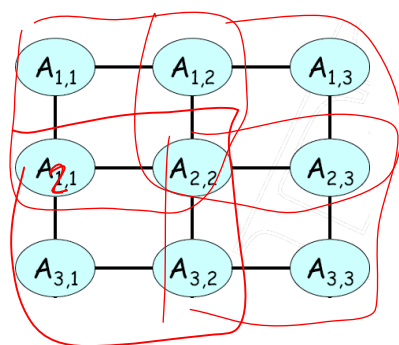
think

XOR

10-708 - ©Carlos Guestrin 2006

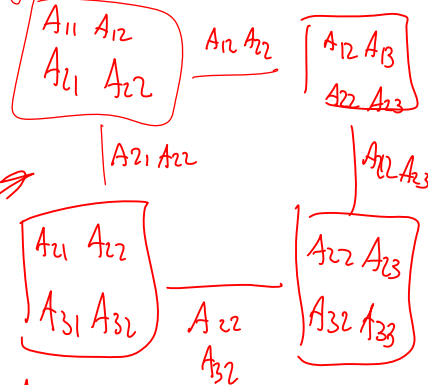
10

What if the cluster graph doesn't satisfy RIP



doesn't satisfy
RIP because of $A_{2,2}$

it is possible to deal with such a cluster graph, but messier see book
intuitive cluster graph



10.708 - ©Carlos Guestrin 2006

11

Region graphs to the rescue

- Can address generalized cluster graphs that don't satisfy RIP using region graphs:

□ Book: 10.3

- Example in your homework! ☺

10.708 - ©Carlos Guestrin 2006

12

Revisiting Mean-Fields

$\ln Z = \underbrace{F[P_{\mathcal{F}}, Q]}_{\text{lower bound}} + D(Q || P_{\mathcal{F}}) \quad F[P_{\mathcal{F}}, Q] = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + H_Q(\mathcal{X})$

- Choice of Q: $Q(x) = \prod_i Q_i(x_i)$ = equal
- Optimization problem:

$$\max_Q \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + \sum_i H_{Q_i}(x_i)$$

$Q_i(x_i) \geq 0$

$\sum_{x_i} Q_i(x_i) = 1$

$Q \approx p$
intuitively
as approx to
 $H_p(x)$

$$\max_Q F[P_{\mathcal{F}}, Q] = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + \sum_j H_{Q_j}(X_j), \quad \forall i, \sum_{x_i} Q_i(x_i) = 1$$

10-708 - ©Carlos Guestrin 2006

13

Announcements

- Recitation tomorrow
- HW5 out soon
- Will not cover relational models this semester
 - Instead, recommend Pedro Domingos' tutorial on Markov Logic
 - Markov logic is one example of a relational probabilistic model
 - November 14th from 1:00 pm to 3:30 pm in Wean 4623

you should go!!

10-708 - ©Carlos Guestrin 2006-2008

14

Interpretation of energy functional

■ Energy functional: $F[P_{\mathcal{F}}, Q] = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + H_Q(\mathcal{X})$

■ ~~Exact if $P=Q$.~~ $\ln Z = \overset{\text{Constant}}{F[P_{\mathcal{F}}, Q]} + D(Q||P_{\mathcal{F}})$

■ View problem as an approximation of entropy term:

$$H_Q(\mathcal{X}) \approx H_P(\mathcal{X})$$

$$F(P_{\mathcal{F}}, Q) = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + H_Q(\mathcal{X}) \approx \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + H_P(\mathcal{X})$$

10-708 - ©Carlos Guestrin 2006

15

Entropy of a tree distribution

■ Entropy term: $H_P(\mathcal{X})$

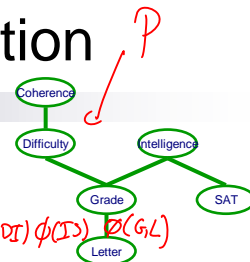
■ Joint distribution: $P(\mathcal{X}) = \frac{1}{Z} \phi(CD) \phi(DG) \phi(DI) \phi(IG) \phi(LG) \phi(LI)$

■ Decomposing entropy term:

$$H(\mathcal{X}) = H(CD) + H(DG) + H(DI) + H(IG) + H(LG) + H(LI) - H(D) - 2H(G) - H(I)$$

for any tree MN

■ More generally: $H_P(\mathcal{X}) = \sum_{(i,j) \in E} H(X_i, X_j) - \sum_i (d_i - 1) H(X_i)$
 □ d_i number neighbors of X_i



10-708 - ©Carlos Guestrin 2006

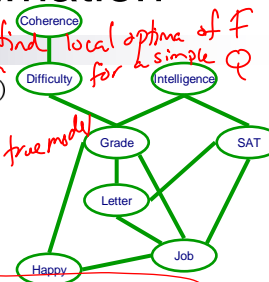
16

Loopy BP & Bethe approximation

- **Energy functional:** $F[P_{\mathcal{F}}, Q] = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + H_Q(\mathcal{X})$

Variational methods try to find local optima of F for a simple Q
 - **Bethe approximation of Free Energy:**
 - use entropy for trees, but loopy graphs: $\approx H_P(X)$

equation for true model
- call edges in loopy graph for tree entropy eqn for a loopy graph*
- $$\tilde{F}[P_{\mathcal{F}}, Q] = \sum_{(i,j) \in E} E_{b_{ij}}[\ln \phi_{ij}] + \sum_{(i,j) \in E} H_{b_{ij}}(X_i, X_j) - \sum_i (d_i - 1) H_{b_i}(X_i)$$
- **Theorem:** If Loopy BP converges, resulting b_{ij} & b_i are stationary point (usually local maxima) of Bethe Free energy!



10.708 - ©Carlos Guestrin 2006

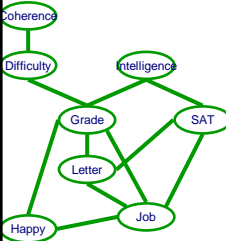
17

GBP & Kikuchi approximation

- **Exact Free energy: Junction Tree**

$$F[P_{\mathcal{F}}, Q] = \sum_{(i,j) \in E} E_{\pi_{ij}}[\ln \phi_{ij}] + \sum_i H_{\pi_{C_i}}(C_i) - \sum_{(i,j) \in T} H_{\pi_{S_{ij}}}(S_{ij})$$
- **Bethe Free energy:**

$$\tilde{F}[P_{\mathcal{F}}, Q] = \sum_{(i,j) \in E} E_{b_{ij}}[\ln \phi_{ij}] + \sum_{(i,j) \in E} H_{b_{ij}}(X_i, X_j) - \sum_i (d_i - 1) H_{b_i}(X_i)$$
- **Kikuchi approximation: Generalized cluster graph**
 - spectrum from Bethe to exact
- **Theorem:** If GBP converges, resulting b_{C_i} are stationary point (usually local maxima) of Kikuchi Free energy!



10.708 - ©Carlos Guestrin 2006

18

What you need to know about GBP

- Spectrum between Loopy BP & Junction Trees:
 - More computation, but typically better answers
- If satisfies RIP, equations are very simple
- General setting, slightly trickier equations, but not hard
- Relates to variational methods: Corresponds to local optima of approximate version of energy functional

10-708 – ©Carlos Guestrin 2006

19

Readings:

K&F: ~~10.2, 10.3~~

Parameter learning in Markov nets

Graphical Models – 10708
Carlos Guestrin
Carnegie Mellon University
November 12th, 2008

10-708 – ©Carlos Guestrin 2006-2008

20

Learning Parameters of a BN

- Log likelihood decomposes:

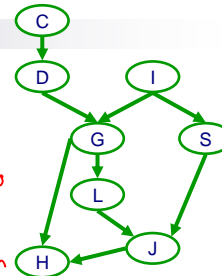
$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta) = m \sum_i \sum_{x_i, \text{Pa}_{x_i}} \hat{P}(x_i, \text{Pa}_{x_i}) \log P(x_i | \text{Pa}_{x_i})$$

parameters
I want
to learn

- Learn each CPT independently

- Use counts

$$P(x_i | \text{Pa}_{x_i} = u) \stackrel{\text{MLE}}{=} \frac{\text{count}(x_i = x_i, \text{Pa}_{x_i} = u)}{\text{count}(\text{Pa}_{x_i} = u)}$$



$$\hat{P}(u) = \frac{\text{Count}(U = u)}{m}$$

10-708 - ©Carlos Guestrin 2006

21

Log Likelihood for MN $\log Z_\theta = \log \sum_x \prod_{i,j} \phi_{ij}(x_i, x_j | \theta_{ij})$

- Log likelihood of the data:

$$\ell(\mathcal{D}; \theta) = \log P(\mathcal{D} | \theta) \stackrel{\text{iid}}{=} \sum_k \log P(x^{(k)} | \theta)$$

$$= \sum_k \log \frac{1}{Z} \prod_{i,j} \phi_{ij}(x_i^{(k)}, x_j^{(k)} | \theta)$$

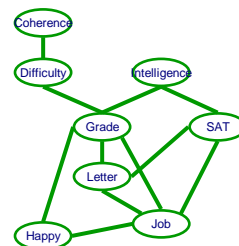
$$= \sum_k \sum_{i,j} \log \phi_{ij}(x_i^{(k)}, x_j^{(k)}) - \sum_{k=1}^m \log Z$$

$$= \sum_{i,j} \sum_{x_i, x_j} \text{count}(x_i = x_i, x_j = x_j) \log \phi_{ij}(x_i = x_i, x_j = x_j) - m \log Z$$

$$= m \sum_{i,j} \sum_{x_i, x_j} \hat{P}(x_i = x_i, x_j = x_j) \log \phi_{ij}(x_i = x_i, x_j = x_j | \theta_{ij}) - m \log Z_\theta$$

decomposes nicely into factors just like BN

doesn't decompose



10-708 - ©Carlos Guestrin 2006

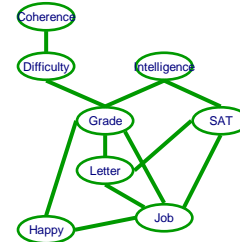
22

Log Likelihood doesn't decompose for MNs

$$\hat{P}(\mathbf{u}) = \frac{\text{Count}(\mathbf{U} = \mathbf{u})}{m}$$

Log likelihood:

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \log \psi_i(\mathbf{c}_i) - m \log Z$$



A convex problem

- Can find global optimum!!

Term log Z doesn't decompose!!

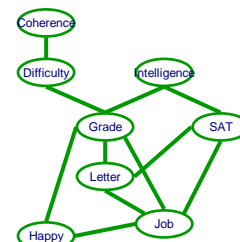
10/708 - ©Carlos Guestrin 2006

23

Derivative of Log Likelihood for MNs

$$\hat{P}(\mathbf{u}) = \frac{\text{Count}(\mathbf{U} = \mathbf{u})}{m}$$

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \log \psi_i(\mathbf{c}_i) - m \log Z$$



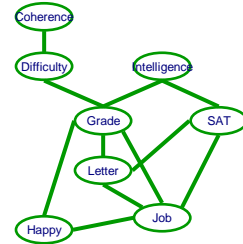
10/708 - ©Carlos Guestrin 2006

24

Derivative of Log Likelihood for MNs 2

$$\hat{P}(\mathbf{u}) = \frac{\text{Count}(\mathbf{U} = \mathbf{u})}{m}$$

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \log \psi_i(\mathbf{c}_i) - m \log Z$$



10-708 - ©Carlos Guestrin 2006

25

Derivative of Log Likelihood for MNs

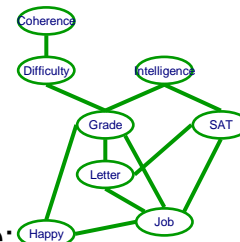
$$\hat{P}(\mathbf{u}) = \frac{\text{Count}(\mathbf{U} = \mathbf{u})}{m}$$

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \log \psi_i(\mathbf{c}_i) - m \log Z$$

- Derivative:

$$\frac{\partial \ell}{\partial \psi_i(\mathbf{c}_i)} = \frac{m \hat{P}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} - \frac{m P_{\mathcal{F}}^{\psi}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)}$$

- Computing derivative requires inference:



- Can optimize using gradient ascent

- Common approach
- Conjugate gradient, Newton's method,...

- Let's also look at a simpler solution

10-708 - ©Carlos Guestrin 2006

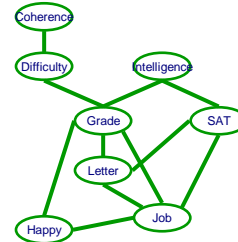
26

Iterative Proportional Fitting (IPF)

$$\hat{P}(u) = \frac{\text{Count}(U = u)}{m}$$

$$\frac{\partial \ell}{\partial \psi_i(\mathbf{c}_i)} = \frac{m \hat{P}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} - \frac{m P_{\mathcal{F}}^{\psi}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)}$$

- Setting derivative to zero:
- Fixed point equation:
- Iterate and converge to optimal parameters
 - Each iteration, must compute:



10-708 - ©Carlos Guestrin 2006

27

What you need to know about learning MN parameters?

- BN parameter learning easy
- MN parameter learning doesn't decompose!
- Learning requires inference!
- Apply gradient ascent or IPF iterations to obtain optimal parameters

10-708 - ©Carlos Guestrin 2006

28