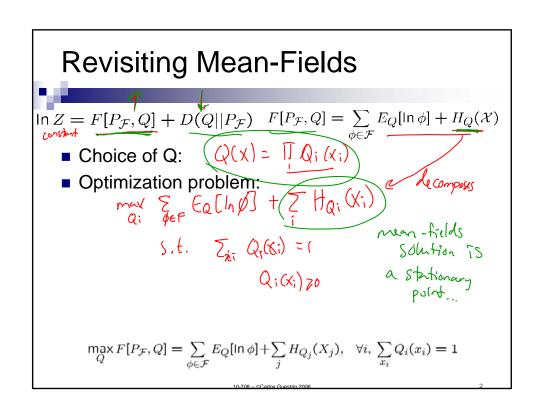
Readings:

K&F: 11.3, 11.5

Yedidia et al. paper from the class website
Chapter 9 – Jordan
K&F: 16.2

Unifying Variational and GBP
Learning Parameters of MNs
EM for BNs

Graphical Models – 10708
Carlos Guestrin
Carnegie Mellon University
November 15th, 2006



Interpretation of energy functional

Energy functional:

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

■ Exact if P=Q: In
$$Z = F[P_{\mathcal{F}}, Q] + D(Q||P_{\mathcal{F}})$$

View problem as an approximation of entropy term:

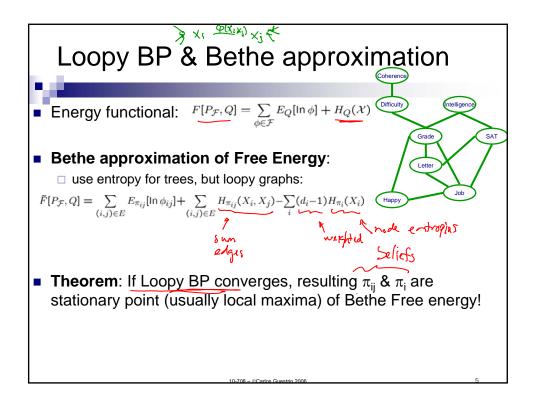
interpretation. Hp(X) & I HQ(X) & Hp(X)

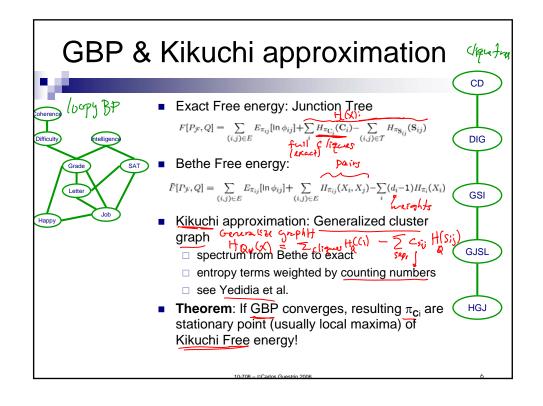
Entropy of a tree distribution

- Entropy term: Hp(X) = -Ep[log P(X)] = -Ep(X) | Difficulty | Difficulty | Ep(X) | Ep
 - Joint distribution: $p(x) = \prod_{i,j} \psi_{i,j}(x_i, x_j)$



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





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

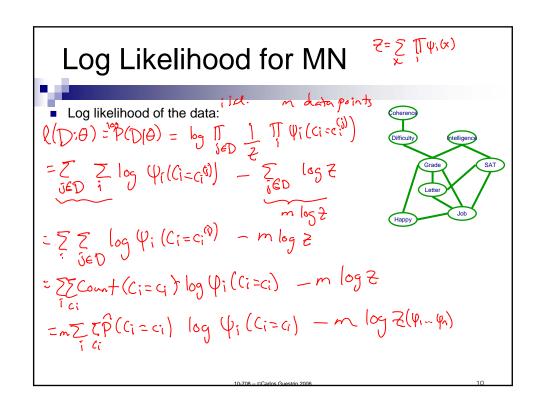
Announcements



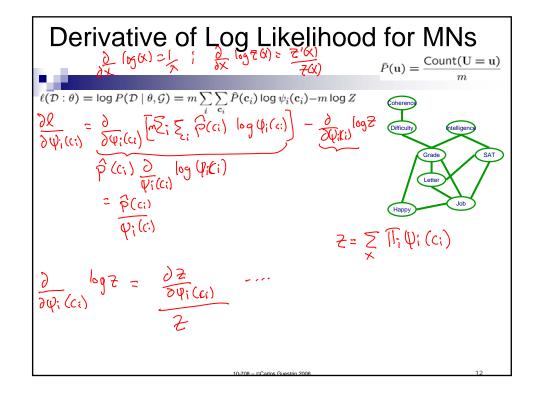
- Tomorrow's recitation
 - ☐ Khalid on learning Markov Networks

10-708 = ©Carlos Guestrin 2006

Learning Parameters of a BN Log likelihood decomposes: $\ell(\mathcal{D}:\theta) = \log P(\mathcal{D} \mid \theta) = m \sum_{i} \sum_{x_i, Pa_{x_i}} \hat{P}(x_i, Pa_{x_i}) \log P(x_i \mid Pa_{x_i}) + \sum_{i} \sum_{x_i, Pa_{x_i}} \hat{P}(x_i, Pa_{x_i}) + \sum_{x_i, Pa_{x_i}} \hat{P}(x_i, Pa_{$

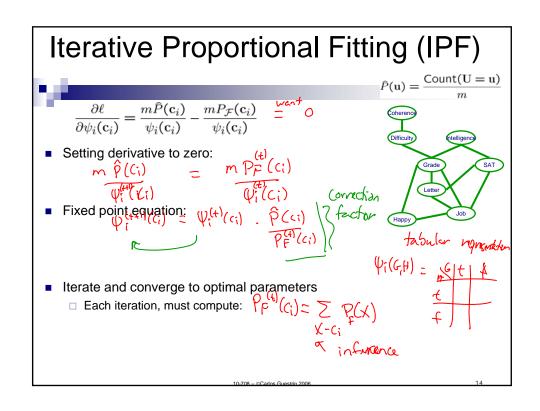


Log Likelihood doesn't decompose for MNs $P(\mathbf{u}) = \frac{\text{Count}(\mathbf{U} = \mathbf{u})}{m}$ • Log likelihood: $\ell(\mathcal{D}:\theta) = \log P(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_{i} \sum_{c_{i}} \hat{P}(c_{i}) \log \psi_{i}(c_{i}) - m \log Z(\Psi_{i} - \Psi_{i}) \log \Psi_{i}(C_{i}) - m \log Z(\Psi_{i}) - m \log Z(\Psi_{i} - \Psi_{i}) \log \Psi_{i}(C_{i}) - m \log Z(\Psi_{i} - \Psi_{i}) - m \log Z(\Psi$



Derivative of Log Likelihood for MNs
$$P(\mathbf{u}) = \frac{\text{Count}(\mathbf{U} = \mathbf{u})}{m}$$

$$P(\mathbf{u})$$



Log-linear Markov network (most common representation)

- **Feature** is some function $\phi[\mathbf{D}]$ for some subset of variables \mathbf{D} e.g., indicator function $\psi(\mathcal{G},\mathcal{I},\mathcal{D}) = \{ \begin{array}{c} \mathcal{G}, \ \mathcal{G}, \$
- Log-linear model over a Markov network H:
 - \square a set of features $\phi_1[\mathbf{D}_1], \ldots, \phi_k[\mathbf{D}_k]$
 - each D_i is a subset of a clique in H
 - two φ's can be over the same variables
 - □ a set of weights w₁,...,w_k

Learning params for log linear models 1 -Generalized Iterative Scaling

P(X_1, \dots, X_n) = $\frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(D_i)\right]$ | $\sum_{x} p(x) \psi_i(x) = \exp chx \text{ value}$ of p_i with p_i with p_i and p_i with p_i wit

- $E_{pur}[\phi_i] = \sum_{di} P_{p}^{(t)}(di) \varphi_i(di)$ $\sum_{i=1}^{n} \varphi_i(x) = 1 \quad \forall x$
 - If conditions violated, equations are not so simple...
 - □ c.f., Improved Iterative Scaling [Berger '97]

Learning params for log linear models 2 – Gradient Ascent

- $P(X_1,...,X_n) = \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{D}_i) \right]$
- Log-likelihood of data: $\begin{cases}
 (D; w) = \log P(D|w) = \log \prod_{j \in D} \frac{1}{2} e^{\sum_{i=1}^{k} w_i \Phi_i(x^{(i)})} \\
 = \log e^{\frac{\pi}{2}} \frac{1}{2} w_i \Phi_i(x^{(i)}) \log 2 \\
 = \frac{\pi}{2} \frac{\pi}{2} v_i \Phi_i(x^{(i)}) m \log 2
 \end{cases}$ Compute derivative & optimize
 - Compute <u>derivative</u> & optimize
 - usually with conjugate gradient ascent
 - ☐ You will do an example in your homework! ☺

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
 - □ applies to both tabular representations and log-linear models

Thus far, fully supervised learning



- We have assumed fully supervised learning:
- Many real problems have missing data:

The general learning problem with missing data



Marginal likelihood –
$$\mathbf{x}$$
 is observed, \mathbf{z} is missing:
$$\ell(\theta:\mathcal{D}) = \log \prod_{j=1}^{m} P(\mathbf{x}_j \mid \theta)$$
$$= \sum_{j=1}^{m} \log P(\mathbf{x}_j \mid \theta)$$
$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} \mid \theta)$$

E-step



- x is observed, z is missing
- Compute probability of missing data given current choice of θ
 □ Q(z|x_i) for each x_i
 - e.g., probability computed during classification step
 - corresponds to "classification step" in K-means

$$Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) = P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})$$

10-708 - ©Carlos Guestrin 2006

21

Jensen's inequality



$$\ell(\theta: \mathcal{D}) = \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{z} \mid \mathbf{x}_{j}) P(\mathbf{x}_{j} \mid \theta)$$

■ Theorem: $\log \Sigma_z P(z) f(z) \ge \Sigma_z P(z) \log f(z)$

10-708 = @Carlos Guestrin 2006

Applying Jensen's inequality



• Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \ge \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\ell(\theta^{(t)}: \mathcal{D}) = \sum_{j=1}^{m} \log \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_{j}) \frac{P(\mathbf{z}, \mathbf{x}_{j} \mid \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_{j})}$$

10-708 - ©Carlos Guestrin 2006

22

The M-step maximizes lower bound on weighted data

Lower bound from Jensen's:

$$\ell(\boldsymbol{\theta}^{(t)}: \mathcal{D}) \geq \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \boldsymbol{\theta}^{(t)}) + H(Q^{(t+1)})$$

- Corresponds to weighted dataset:
 - $\neg < \mathbf{x}_1, \mathbf{z} = 1 > \text{ with weight } \mathbf{Q}^{(t+1)}(\mathbf{z} = 1 | \mathbf{x}_1)$
 - \square < \mathbf{x}_1 , \mathbf{z} =2> with weight Q^(t+1)(\mathbf{z} =2| \mathbf{x}_1)
 - \neg <**x**₁,**z**=3> with weight Q^(t+1)(**z**=3|**x**₁)
 - \neg < \mathbf{x}_2 , $\mathbf{z}=1$ > with weight $Q^{(t+1)}(\mathbf{z}=1|\mathbf{x}_2)$
 - $\neg < \mathbf{x}_2, \mathbf{z} = 2 > \text{ with weight } Q^{(t+1)}(\mathbf{z} = 2|\mathbf{x}_2)$
 - $\neg < \mathbf{x}_2, \mathbf{z} = 3 > \text{ with weight } Q^{(t+1)}(\mathbf{z} = 3 | \mathbf{x}_2)$
 - □ ...

10-708 = ©Carlos Guestrin 200

The M-step



$$\ell(\boldsymbol{\theta}^{(t)}: \mathcal{D}) \; \geq \; \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \boldsymbol{\theta}^{(t)}) + H(Q^{(t+1)})$$

Maximization step:

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_{j}) \log P(\mathbf{z}, \mathbf{x}_{j} \mid \theta)$$

- Use expected counts instead of counts:
 - ☐ If learning requires Count(x,z)
 - \square Use $E_{Q(t+1)}[Count(\mathbf{x},\mathbf{z})]$

10-708 = @Carlos Guestrin 2006

a.e.

Convergence of EM



■ Define potential function $F(\theta,Q)$:

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

- EM corresponds to coordinate ascent on F
 - □ Thus, maximizes lower bound on marginal log likelihood
 - □ As seen in Machine Learning class last semester

10-708 = @Carlos Guestrin 2006

Data likelihood for BNs



Given structure, log likelihood of fully observed data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$

08 = @Carlos Guestrin 2006

27

Marginal likelihood



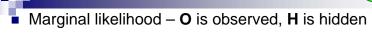
What if S is hidden?

 $\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$



10-708 - ©Carlos Guestrin 2006

Log likelihood for BNs with hidden data



$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^{m} \log P(\mathbf{o}^{(j)} | \theta)$$
$$= \sum_{j=1}^{m} \log \sum_{\mathbf{h}} P(\mathbf{h}, \mathbf{o}^{(j)} | \theta)$$

-708 - ©Carlos Guestrin 2006

29

E-step for BNs





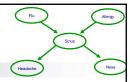
■ E-step computes probability of hidden vars h given o

$$Q^{(t+1)}(\mathbf{x} \mid \mathbf{o}) = P(\mathbf{x} \mid \mathbf{o}, \theta^{(t)})$$

■ Corresponds to inference in BN

0-708 - @Carlos Guestrin 2006

The M-step for BNs





Maximization step:

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{\mathbf{x}} Q^{(t+1)}(\mathbf{h} \mid \mathbf{o}) \log P(\mathbf{h}, \mathbf{o} \mid \theta)$$

- Use expected counts instead of counts:
 - \square If learning requires Count(\mathbf{h} , \mathbf{o})
 - \square Use $E_{Q(t+1)}[Count(\mathbf{h}, \mathbf{o})]$

08 = ©Carlos Guestrio 2006

21

M-step for each CPT





- M-step decomposes per CPT
 - ☐ Standard MLE:

$$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\mathsf{Count}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\mathsf{Count}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

☐ M-step uses expected counts:

$$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\mathsf{ExCount}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\mathsf{ExCount}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

10-708 = @Carlos Guestrin 2006

Computing expected counts

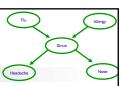


$$P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{z}) = \frac{\mathsf{ExCount}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{z})}{\mathsf{ExCount}(\mathbf{Pa}_{X_i} = \mathbf{z})}$$

- M-step requires expected counts:
 - □ For a set of vars **A**, must compute ExCount(**A**=**a**)
 - \square Some of **A** in example *j* will be observed
 - denote by $\mathbf{A}_{\mathbf{O}} = \mathbf{a}_{\mathbf{O}}^{(j)}$
 - □ Some of **A** will be hidden
 - denote by A_H
- Use inference (E-step computes expected counts):
 - $\square \; \mathsf{ExCount}^{(\mathsf{t+1})}(\mathsf{A}_\mathsf{O} = \mathsf{a}_\mathsf{O}^{(\mathsf{j})}, \, \mathsf{A}_\mathsf{H} = \mathsf{a}_\mathsf{H}) \leftarrow \mathsf{P}(\mathsf{A}_\mathsf{H} = \mathsf{a}_\mathsf{H} \mid \mathsf{A}_\mathsf{O} = \mathsf{a}_\mathsf{O}^{(\mathsf{j})}, \theta^{(\mathsf{t})})$

10-708 - @Carlos Guestrin 2006

Data need not be hidden in the same way



- - When data is fully observed
 - A data point is
 - When data is partially observed
 - □ A data point is
 - But unobserved variables can be different for different data points
 - □ e.g.,
 - Same framework, just change definition of expected counts
 - □ ExCount((t+1)($A_O = a_O(i)$, $A_H = a_H$) ← P($A_H = a_H \mid A_O = a_O(i)$, $\theta(t)$)

10-708 = @Carlos Guestrin 2006

What you need to know



- EM for Bayes Nets
- E-step: inference computes expected counts

 □ Only need expected counts over X_i and Pa_{xi}
- M-step: expected counts used to estimate parameters
- Hidden variables can change per datapoint
- Use labeled and unlabeled data → some data points are complete, some include hidden variables

10-708 = @Carlos Guestrin 2006