

Approximate inference overview

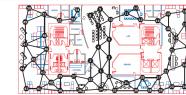
- So far: VE & junction trees
 - □ exact inference
 - □ exponential in tree-width
- There are many many many many approximate inference algorithms for PGMs
- We will focus on three representative ones:
 - sampling
 - □ variational inference
 - □ loopy belief propagation and generalized belief propagation
- There will be a special recitation by Pradeep Ravikumar on more advanced methods

10-708 - @Carlos Guestrin 2006

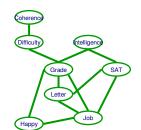
Approximating the posterior v. approximating the prior



- Prior model represents entire world
 - □ world is complicated
 - □ thus prior model can be very complicated



- Posterior: after making observations
 - sometimes can become much more sure about the way things are
 - □ sometimes can be approximated by a simple model
- First approach to approximate inference: find simple model that is "close" to posterior
- Fundamental problems:
 - □ what is close?
 - posterior is intractable result of inference, how can we approximate what we don't have?



0-708 - ©Carlos Guestrin 2006

3

KL divergence:

Distance between distributions



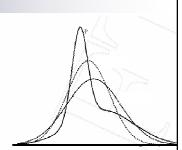
- Given two distributions p and q KL divergence:
- D(p||q) = 0 iff p=q
- Not symmetric p determines where difference is important
 - \Box p(x)=0 and q(x) \neq 0
 - \Box p(x) \neq 0 and q(x)=0

10-708 – ©Carlos Guestrin 2006

Find simple approximate distribution



- Suppose p is intractable posterior
- Want to find simple q that approximates p
- KL divergence not symmetric
- D(p||q)
 - □ true distribution p defines support of diff.
 - □ the "correct" direction
 - □ will be intractable to compute
- D(q||p)
 - □ approximate distribution defines support
 - tends to give overconfident results
 - □ will be tractable



0-708 – ©Carlos Guestrin 200

5

Back to graphical models



- Inference in a graphical model:
 - \square P(x) =
 - \square want to compute $P(X_i|\mathbf{e})$
 - □ our *p*:
- What is the simplest q?
 - □ every variable is independent:
 - □ mean field approximation
 - □ can compute any prob. very efficiently

0.708 = @Carlos Guestrin 200

D(p||q) for mean field – KL the right way

- p:
- **q**:
- D(p||q)=

10-708 = ©Carlos Guestrin 2006

D(q||p) for mean field – KL the reverse direction

- p:
 - **q**:
 - D(p||q)=

10-708 = @Carlos Guestrin 2006

What you need to know so far



- Goal:
 - ☐ Find an efficient distribution that is close to posterior
- Distance:
 - □ measure distance in terms of KL divergence
- Asymmetry of KL:
 - \square D(p||q) \neq D(q||p)
- Computing right KL is intractable, so we use the reverse KL

10-708 - ©Carlos Guestrin 2006

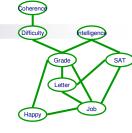
۵

Reverse KL & The Partition Function

Back to the general case



- Consider again the defn. of D(q||p):
 - $\ \square$ p is Markov net P_F



- Theorem: In $Z = F[P_{\mathcal{F}}, Q] + D(Q||P_{\mathcal{F}})$
- where energy functional:

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

0-708 - ©Carlos Guestrin 2006

Understanding Reverse KL, Energy Function & The Partition Function

$$\ln Z = F[P_{\mathcal{F}},Q] + D(Q||P_{\mathcal{F}}) \qquad \qquad F[P_{\mathcal{F}},Q] = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + H_Q(\mathcal{X})$$

- Maximizing Energy Functional ⇔ Minimizing Reverse KL
- **Theorem**: Energy Function is lower bound on partition function
 - Maximizing energy functional corresponds to search for tight lower bound on partition function

10-708 - @Carlos Guestrin 2006

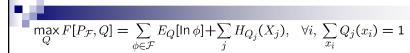
11

Structured Variational Approximate In $Z = F[P_{\mathcal{F}}, Q] + D(Q||P_{\mathcal{F}})$ $F[P_{\mathcal{F}}, Q] = \sum\limits_{\phi \in \mathcal{F}} E_Q[\ln \phi] + H_Q(\mathcal{X})$

- Pick a family of distributions Q that allow for exact inference
 - □ e.g., fully factorized (mean field)
- Find $Q \in Q$ that maximizes $F[P_F, Q]$
- For mean field

10-708 - @Carlos Guestrin 2006

Optimization for mean field



- Constrained optimization, solved via Lagrangian multiplier
 - \Box \exists λ , such that optimization equivalent to:
 - □ Take derivative, set to zero
- **Theorem**: Q is a stationary point of mean field approximation iff for each *i*:

$$Q_i(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi \mid x_i] \right\}$$

10-708 - @Carlos Guestrin 2006

13

Understanding fixed point equation

١,

$$Q_i(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi \mid x_i] \right\}$$

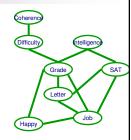
Oherency
Difficulty Itelligency
Grade SAT
Letter
Happy Job

10-708 - ©Carlos Guestrin 2006

Simplifying fixed point equation



$$Q_i(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi \mid x_i] \right\}$$



Q_i only needs to consider factors that intersect X_i

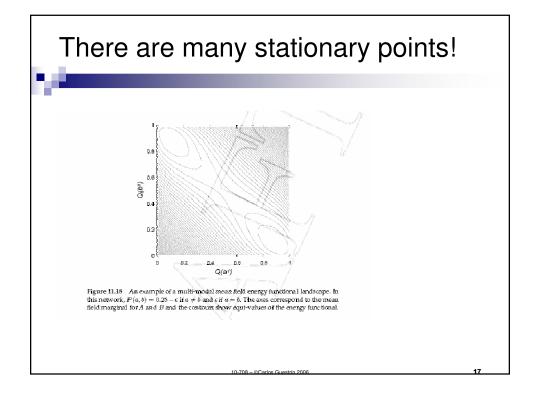


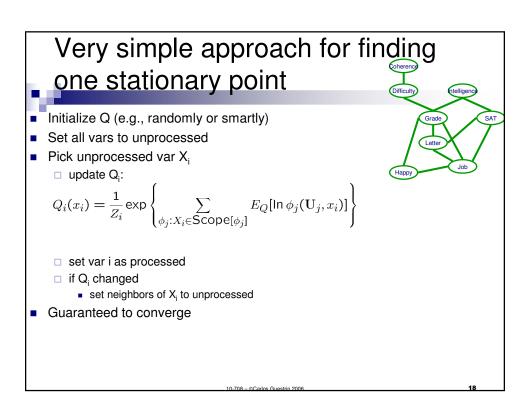
■ **Theorem**: The fixed point:

$$Q_i(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi \mid x_i] \right\}$$

is equivalent to:
$$Q_i(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi_j: X_i \in \mathsf{Scope}[\phi_j]} E_Q[\ln \phi_j(\mathbf{U}_j, x_i)] \right\}$$

 $\quad \ \ \, \square \ \, \text{where the Scope}[\varphi_i] = \boldsymbol{U}_i \cup \{X_i\}$

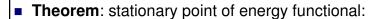




More general structured approximations



- Mean field very naïve approximation
- Consider more general form for Q
 - □ assumption: exact inference doable over Q



$$\psi_j(\mathbf{c_j}) \propto \exp \left\{ \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi \mid \mathbf{c_j}] - \sum_{\psi \in \mathcal{Q} \backslash \{\psi_j\}} E_Q[\ln \psi \mid \mathbf{c_j}] \right\}$$

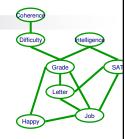
0-708 - ©Carlos Guestrin 2006

40

Computing update rule for general case

$$\psi_j(\mathbf{c_j}) \propto \exp\left\{\sum_{\phi \in \mathcal{F}} E_Q[\ln \phi \mid \mathbf{c_j}] - \sum_{\psi \in \mathcal{Q} \backslash \{\psi_j\}} E_Q[\ln \psi \mid \mathbf{c_j}]\right\} \text{ otherwise}$$

Consider one φ:



I0-708 – ©Carlos Guestrin 2006

Structured Variational update requires inferece

$$\psi_{j}(\mathbf{c_{j}}) \propto \exp \left\{ \sum_{\phi \in \mathcal{F}} E_{Q}[\ln \phi \mid \mathbf{c_{j}}] - \sum_{\psi \in \mathcal{Q} \setminus \{\psi_{j}\}} E_{Q}[\ln \psi \mid \mathbf{c_{j}}] \right\}$$

- Compute marginals wrt Q of cliques in original graph and cliques in new graph, for all cliques
- What is a good way of computing all these marginals?
- Potential updates:
 - $\ \square$ sequential: compute marginals, update ψ_i , recompute marginals

10-708 - ©Carlos Guestrin 2006

21

What you need to know about variational methods



- Structured Variational method:
 - □ select a form for approximate distribution
 - □ minimize reverse KL
- Equivalent to maximizing energy functional
 - searching for a tight lower bound on the partition function
- Many possible models for Q:
 - □ independent (mean field)
 - □ structured as a Markov net
 - cluster variational
- Several subtleties outlined in the book

10-708 - ©Carlos Guestrin 2006