

Readings:
K&F: 11.1, 11.5

Mean Field and Variational Methods

First approximate inference

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 1st, 2006

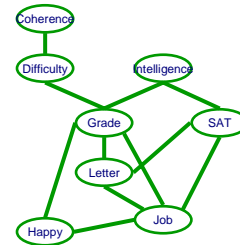
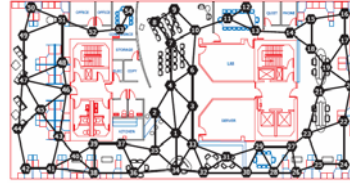
1

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 *a week from Mon.*

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?**



10-708 - ©Carlos Guestrin 2006

3

KL divergence: Distance between distributions

- Given two distributions p and q KL divergence:

$$D(p||q) = \sum_x p(x) \cdot \log \frac{p(x)}{q(x)}$$

- $D(p||q) = 0$ iff $p=q$
- Not symmetric – p determines where difference is important

- $p(x)=0$ and $q(x) \neq 0$

$$p(x) \cdot \log \frac{p(x)}{q(x)} = 0 \cdot \log \frac{0}{q(x)} = 0$$

doesn't appear in distance
 - $p(x) \neq 0$ and $q(x)=0$

$$\overset{\approx \epsilon}{p(x)} \log \frac{\overset{\approx \epsilon}{p(x)}}{\overset{\approx 0}{q(x)}} = \epsilon \log \frac{\epsilon}{0} = +\infty$$

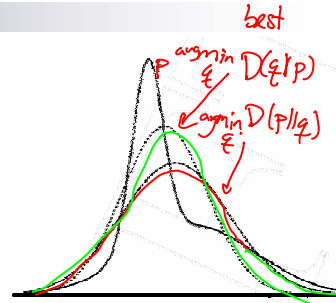
if p thinks $p(x) > 0$, then q better think $q(x) > 0$

10-708 - ©Carlos Guestrin 2006

4

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 q defines support
 - tends to give overconfident results
 - will be tractable



p complex - Gaussian
w. 2 bumps

q simple - Gaussian
w. 1 bump

10-708 - ©Carlos Guestrin 2006

5

Back to graphical models

- Inference in a graphical model:

- $P(\mathbf{x}) = \frac{1}{Z} \prod_i \phi_i(c_i)$
- want to compute $P(X_i | \mathbf{e}) \propto P(X_i, \mathbf{e})$
- our p : $P(\mathbf{x}) = \frac{1}{Z} \prod_i \phi_i(c_i)$, want $P(x_i)$

- What is the simplest q ?

- every variable is independent: $q(\mathbf{x}) = \prod_i q_i(x_i)$
- mean field approximation
- can compute any prob. very efficiently

$$P(\mathbf{x}, \mathbf{e}) = \frac{1}{Z} \prod_i \phi_i(c_i)$$

notation drop \mathbf{e}

$q_i(x_i)$

$q_2(x_2)$ $q_3(x_3)$

10-708 - ©Carlos Guestrin 2006

6

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

- p: $P(x) = \frac{1}{Z} \prod_i \phi_i(c_i)$
- q: $Q(x) = \prod_i Q_i(x_i)$
- $D(p||q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} = \sum_x P(x) \log P(x) - \sum_x P(x) \log Q(x)$

$= -H_p(x) - \sum_x P(x) \log \prod_i Q_i(x_i) = -H_p(x) - \sum_i \sum_x P(x) \log Q_i(x_i)$

$\hookrightarrow = \sum_{x_i} P(x_i) \underbrace{\sum_{x_{\setminus i}} P(x_{\setminus i} | x_i)}_1 \log Q_i(x_i) = \sum_{x_i} P(x_i) \log Q_i(x_i)$

*doesn't depend on z
no role in approx.*
-H_p(x)
 *$\prod_i Q_i$
↓
 Q_i*
*need P(x_i) which is
what was intractable
in the first place*

10-708 – ©Carlos Guestrin 2006

7

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

- p: $P(x) = \frac{1}{Z} \prod_i \phi_i(c_i)$
- q: $Q(x) = \prod_i Q_i(x_i)$
- $D(q||p) = \sum_x Q(x) \log \frac{Q(x)}{P(x)} = \sum_x Q(x) \log Q(x) - \sum_x Q(x) \log P(x)$

$H_q(x) = -\sum_x Q(x) \log \prod_i Q_i(x_i) = -\sum_i \sum_{x_i} Q_i(x_i) \log Q_i(x_i) = -\sum_i H_{q_i}(x_i)$

$= -\sum_i \sum_x Q(x) \log Q_i(x_i) = -\sum_i \sum_{x_i} \underbrace{Q(x)}_{Q_i(x_i)} \log Q_i(x_i) = \sum_i \underbrace{E_{Q_i}[\log Q_i(x_i)]}_{H_{q_i}(x_i)}$

$\sum_x Q(x) \log P(x) = \sum_x Q(x) \log \frac{1}{Z} \prod_i \phi_i(c_i) = \sum_i \sum_x Q(x) \log \phi_i(c_i) + \sum_x Q(x) \log Z$

$\hookrightarrow \sum_x Q(x) \log \phi_i(c_i) = \sum_{c_i} \sum_y Q(y | c_i) Q(c_i) \log \phi_i(c_i) = \sum_{c_i} Q(c_i) \log \phi_i(c_i) \sum_y Q(y | c_i) = E_{Q_i}[\log \phi_i(c_i)] = +\log Z$

$Q_i(x_i) = \sum_{x_{\setminus i}} Q(x)$
 $H_q(x) = \sum_i H_{q_i}(x_i)$
tractable model
chain rule

10-708 – ©Carlos Guestrin 2006

8

What you need to know so far

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

10-708 – ©Carlos Guestrin 2006

9

Announcements

- Tomorrow's recitation
 - Khalid on Variational Methods
- Monday's special recitation
 - Khalid on Dirichlet Processes
 - An exciting way to deal with model selection using graphical models, e.g., selecting (or averaging over) number of clusters in unsupervised learning
 - you even get to see a BN with infinitely many variables

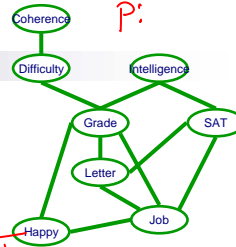
10-708 – ©Carlos Guestrin 2006

10

Reverse KL & The Partition Function

Back to the general case

$$E_Q(f) = \sum_x q(x) f(x)$$



- Consider again the defn. of $D(q||p)$:

$$p \text{ is Markov net } P_F = \frac{1}{Z} \prod_i \phi_i(c_i)$$

$$KL(q||p) = \sum_x q(x) \log \frac{q(x)}{p(x)} = -H_q(x) - \sum_x q(x) \log p(x)$$

$$= -H_q(x) - \sum_i \sum_x q(x) \log \phi_i(c_i) + \log Z \quad \left| \quad \sum_x q(x) \log \phi_i(c_i) = E_q[\log \phi_i] \right.$$

Theorem: $\ln Z = F[P_F, Q] + D(Q||P_F)$

$$D(q||p) = -H_q(x) - \sum_i E_q[\log \phi_i] + \log Z$$

- where energy functional:

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

$$q: \begin{matrix} C \\ | \\ D-I \\ | \\ H-G \\ | \\ L-J \end{matrix}$$

10-708 - ©Carlos Guestrin 2006

11

Understanding Reverse KL, Energy Function & The Partition Function

$$\ln Z = F[P_F, Q] + D(Q||P_F)$$

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

- Maximizing Energy Functional \Leftrightarrow Minimizing Reverse KL

$$\uparrow F[P_F, Q] \Leftrightarrow \downarrow D(Q||P_F) \quad , \forall Q \ D(Q||P_F) \geq 0$$

- Theorem:** Energy Function is lower bound on partition function

$$\ln Z \geq F[P_F, Q] \quad \uparrow \text{maximize}$$

$$\ln Z = F[P_F, Q] \quad \text{iff } Q = P_F$$

- Maximizing energy functional corresponds to search for tight lower bound on partition function

10-708 - ©Carlos Guestrin 2006

12

Structured Variational Approximate Inference

$$\ln Z = F[P_{\mathcal{F}}, Q] + D(Q || P_{\mathcal{F}})$$

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

- Pick a family of distributions Q that allow for exact inference

e.g., graphical model for Q

- e.g., fully factorized (mean field) $Q(x) = \prod_i Q_i(x_i)$

- Find $Q \in \mathcal{Q}$ that maximizes $F[P_{\mathcal{F}}, Q]$

\mathcal{F} is graphical model for P

- For mean field: $F[P_{\mathcal{F}}, Q] = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + \sum_i H_{Q_i}(x_i)$

$$\forall x_i: \sum_{x_i} Q_i(x_i) = 1 \quad Q_i(x_i) \geq 0 \quad \forall x_i$$

10-708 - ©Carlos Guestrin 2006

13

Optimization for mean field

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

- Constrained optimization, solved via Lagrangian multiplier

- $\exists \lambda$, such that optimization equivalent to:

$$L(Q, \lambda) = \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi] + \sum_j H_{Q_j}(x_j) + \lambda \left(\sum_{x_j} Q_j(x_j) - 1 \right)$$

- Take derivative, set to zero

- **Theorem:** Q is a stationary point of mean field approximation iff for each i :

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

10-708 - ©Carlos Guestrin 2006

14

conditional expectation

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

$$E_Q[\ln \phi \mid x_i] = \sum_x Q(x_i) \log \phi(c)$$

$$= \sum_{i \in I} Q(I=i \mid G=g) \log \phi(G=g, I=i)$$

$$= \sum_{i \in I} Q(I=i) \log \phi(G=g^t, I=i)$$

$$= 0.8 \log 1 + 0.2 \log 2$$

$$\phi(G, I) = \begin{array}{c|cc} & t & f \\ \hline t & 1 & 2 \\ \hline f & 3 & 4 \end{array}$$

meanfield

$$Q(I=i \mid G=g) = Q(I=i)$$

$$Q(I) = \begin{array}{c|cc} & t & f \\ \hline & 0.8 & 0.2 \end{array}$$

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

```
graph TD; Coherence((Coherence)) --> Difficulty((Difficulty)); Difficulty --> Grade((Grade)); Intelligence((Intelligence)) --> SAT((SAT)); Grade --> Letter((Letter)); Grade --> Job((Job)); SAT --> Job; Letter --> Happy((Happy)); Letter --> Job; Happy --> Job;
```

10-708 – ©Carlos Guestrin 2006

16

Q_i only needs to consider factors that intersect X_i

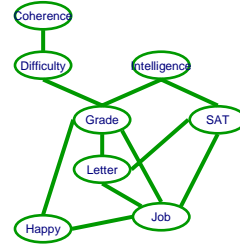
- **Theorem:** The fixed point:

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

is equivalent to:

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

□ where the $\text{Scope}[\phi_j] = \mathbf{U}_j \cup \{X_i\}$



10-708 – ©Carlos Guestrin 2006

17

There are many stationary points!

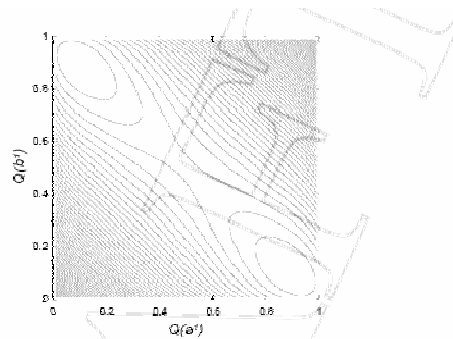


Figure 11.18 An example of a multi-modal mean field energy functional landscape. In this network, $P(a, b) = 0.25 - \epsilon$ if $a \neq b$ and ϵ if $a = b$. The axes correspond to the mean field marginal for A and B and the contours show equi-values of the energy functional.

10-708 – ©Carlos Guestrin 2006

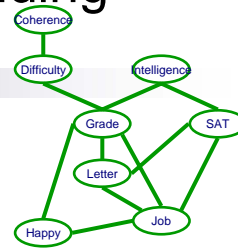
18

Very simple approach for finding one stationary point

- Initialize Q (e.g., randomly or smartly)
- Set all vars to unprocessed
- Pick unprocessed var X_i
 - update Q_i :

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

- set var i as processed
 - if Q_i changed
 - set neighbors of X_i to unprocessed
- Guaranteed to converge



10-708 - ©Carlos Guestrin 2006

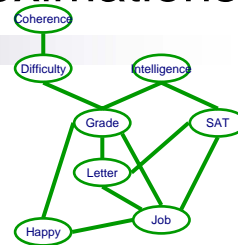
19

More general structured approximations

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

- **Theorem:** stationary point of energy functional:

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



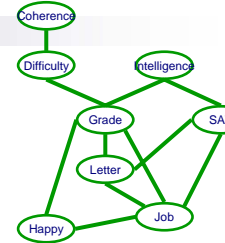
10-708 - ©Carlos Guestrin 2006

20

Computing update rule for general case

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

- Consider one ϕ :



10-708 - ©Carlos Guestrin 2006

21

Structured Variational update requires inference

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

- 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:
 - sequential: compute marginals, update ψ_j , recompute marginals
 - parallel: compute marginals, update all ψ 's, recompute marginals

10-708 - ©Carlos Guestrin 2006

22

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