

Readings:

K&F: 4.1, 4.2, 4.3, 4.4, 4.5

Undirected Graphical Models (finishing off)

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 3rd, 2008

10-708 – © Carlos Guestrin 2006-2008

1

What you learned about so far

- Bayes nets
- Junction trees
- (General) Markov networks
- Pairwise Markov networks
- Factor graphs

- How do we transform between them?
- More formally:
 - I give you an graph in one representation, find an **I-map** in the other

10-708 – © Carlos Guestrin 2006-2008

2

BNs ~~→~~ MNs: Moralization

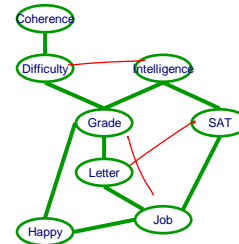
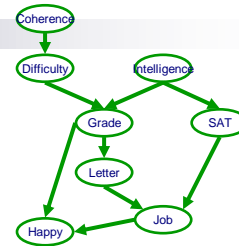
■ **Theorem:** Given a BN G the Markov net H formed by moralizing G is the *minimal I-map* for $I(G)$

■ **Intuition:**

- in a Markov net, each factor must correspond to a subset of a clique
- the factors in BNs are the CPTs
- CPTs are factors over a node and its parents
- thus node and its parents must form a clique

■ **Effect:**

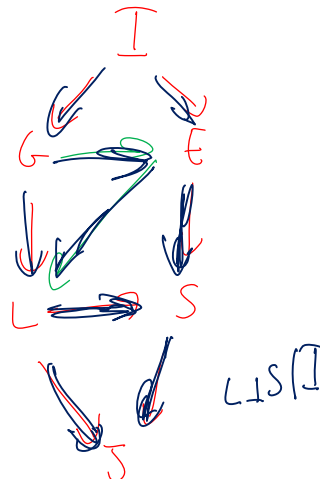
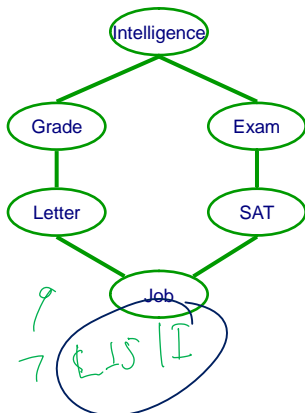
- **some** independencies that could be read from the BN graph become hidden



10-708 — Carlos Guestrin, 2006-2008

3

From Markov nets to Bayes nets



10-708 — Carlos Guestrin, 2006-2008

4

MNs \rightarrow BNs: Triangulation

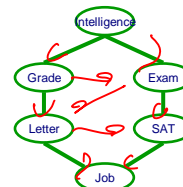
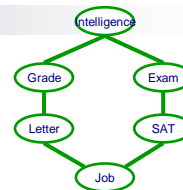
- **Theorem:** Given a MN H , let G be the Bayes net that is a *minimal I-map* for $I(H)$ then G must be **chordal** $I(G) \subseteq I(H) \subseteq I(P)$

- **Intuition:**

- v-structures in BN introduce immoralities
- these immoralities were not present in a Markov net
- the triangulation eliminates immoralities

- **Effect:**

- **many** independencies that could be read from the MN graph become hidden



10-708 — Carlos Guestrin, 2006-2008

5

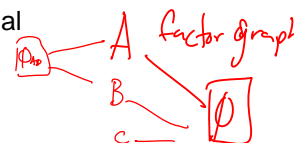
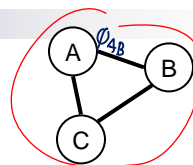
Markov nets v. Pairwise MNs

- Every Markov network can be transformed into a Pairwise Markov net

- introduce extra "variable" for each factor over three or more variables
- domain size of extra variable is exponential in number of vars in factor

- **Effect:**

- any local structure in factor is lost
- a chordal MN doesn't look chordal anymore



need to ensure that assignment of A & D agree:

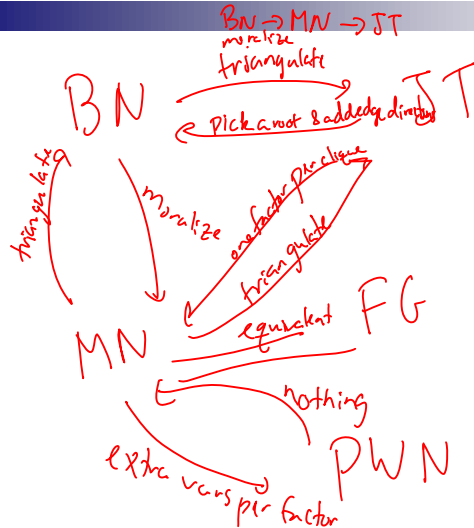
$$\psi(A, D) = \begin{cases} \phi & \text{when } a \text{ & } d \text{ are inconsistent} \\ 1 & \text{otherwise} \end{cases}$$

$\text{Dom}[D] \equiv \text{Dom}[A] \times \text{Dom}[B] \times \text{Dom}[C]$
node potential on D , $\psi(D)$
 $\psi(D) = \phi(a, b, c)$
a,b,c assignment of A, B, C consistent w. d

10-708 — Carlos Guestrin, 2006-2008

6

Overview of types of graphical models and transformations between them



10-708 — Carlos Guestrin, 2006-2008

7

Readings:
K&F: 10.1, 10.5

Mean Field and Variational Methods

First approximate inference

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 3rd, 2008

10-708 — Carlos Guestrin, 2006-2008

8

Approximate inference overview

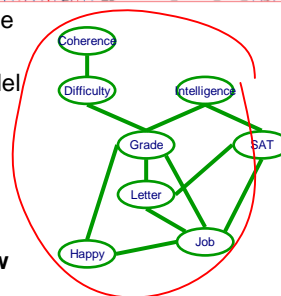
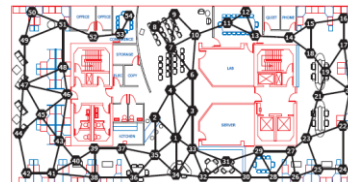
- So far: VE & junction trees
 - exact inference
 - exponential in tree-width
- There are many many many many approximate inference algorithms for PGMs
- We will focus on three representative ones:
 - sampling
 - variational inference
 - loopy belief propagation and generalized belief propagation

10-708 — Carlos Guestrin, 2006-2008

9

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, 2006-2008

10

KL divergence:

$$0 \log 0 = 0 \quad \lim_{x \rightarrow 0^+} x \log x = 0$$

Distance between distributions

- Given two distributions p and q KL divergence:

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

- $D(p||q) = 0$ iff $p=q$

- Not symmetric – p determines where difference is important

$$\exists \square p(x)=0 \text{ and } q(x) \neq 0 \quad p(x) \log \frac{p(x)}{q(x)} = 0 \log 0 = 0$$

$$\exists \square p(x) \neq 0 \text{ and } q(x)=0 \quad p(x) \log \frac{p(x)}{q(x)} = \epsilon \log \frac{\epsilon}{0} = +\infty$$

Find simple approximate distribution

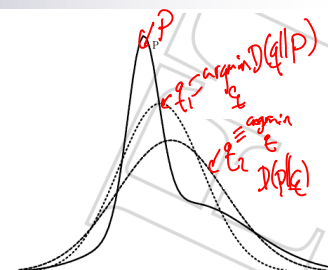
- Suppose p is intractable posterior
- Want to find simple q that approximates p
- KL divergence not symmetric

- $D(p||q)$

- ☐ 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



Back to graphical models

Inference in a graphical model:

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

drop e and assume e has been instantiated in every factor

What is the simplest q ?

- every variable is independent:
- mean field approximation
- can compute any prob. very efficiently

$$Q(\mathbf{x}) = \prod_i \phi_i(x_i)$$

easy to generalize to any $Q(\mathbf{x})$ where exact inference is feasible

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

■ $p: \frac{1}{Z} \prod_i \phi_i(c_i)$

■ $q: \prod_i \phi_i(x_i)$

■ $D(p||q) = \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} = \underbrace{\sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x})}_{-H_p(\mathbf{x})} - \underbrace{\sum_{\mathbf{x}} p(\mathbf{x}) \log q(\mathbf{x})}_{\text{cross entropy}}$

doesn't depend on q , I can ignore

$$- \sum_{\mathbf{x}} p(\mathbf{x}) \log \prod_i \phi_i(x_i) = - \sum_i \sum_{\mathbf{x}} p(\mathbf{x}) \log \phi_i(x_i)$$

$$\begin{aligned} & \sum_{x_1, \dots, x_n} p(x_1, \dots, x_n) \log \phi_i(x_i) \\ &= \sum_{x_i} p(x_i) \log \phi_i(x_i) \end{aligned}$$

the thing that was hard to compute in the first place!!

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

- p: $\frac{1}{Z} \prod_i \phi_i(c_i)$
- q: $\prod_j Q_j(x_j)$
- $D(q||p) = \sum_x q(x) \log \frac{q(x)}{p(x)}$
 $= \sum_x q(x) \log q(x) - \sum_x q(x) \log p(x)$

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

- p: $\frac{1}{Z} \prod_i \phi_i(c_i)$
 - q: $\prod_j Q_j(x_j)$
- first on entropy term*
- $$D(q||p) = \underbrace{\sum_x q(x) \log q(x)}_{-H_q(x)} - \sum_x q(x) \log p(x)$$
- $$-H_q(x) = \sum_x q(x) \log \prod_j Q_j(x_j) = \sum_j \sum_x q(x) \log Q_j(x_j)$$
- $$= \sum_j \sum_{x_j} q(x_j) \log Q_j(x_j) = \sum_j \sum_{x_j} Q_j(x_j) \log Q_j(x_j)$$
- easy to compute*
- $q(x_j)$ is marginal wrt x_j of $q(x) \equiv Q_j(x_j)$*

$D(q||p)$ for mean field – any proposed Q where exact inference is feasible allows us to compute $q(c_i)$
 KL the reverse direction: cross-entropy term

■ $p: \frac{1}{Z} \prod_i \phi_i(c_i)$

■ $q: \prod_j Q_j(x_j)$

$$D(q||p) = \sum_x q(x) \log q(x) - \sum_x q(x) \log p(x)$$

$\sum_x q(x) \log p(x) = \sum_x q(x) \log \frac{1}{Z} \prod_i \phi_i(c_i) = \sum_{i,x} q(x) \log \phi_i(c_i) - \sum_x q(x) \log Z$

$\sum_x q(x) \log \phi_i(c_i) = \sum_{c_i} q(c_i) \log \phi_i(c_i)$ easy to compute $= E_Q[\log \phi_i]$

for mean fields $q(c_i) = \prod_{x_j \in c_i} Q_j(x_j)$ (as long as c_i not too large)

What you need to know so far

■ Goal: $p(x|e) \approx \prod_j Q_j(x_j)$ $Q_j(x_j) \approx p(x_j|e)$

□ Find an efficient distribution that is close to posterior

■ 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

Reverse KL & The Partition Function

Back to the general case

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

□ p is Markov net P_F

$$p(x) = \frac{1}{Z} \prod_{\phi \in \mathcal{F}} \phi(c_\phi)$$

maximize

want to minimize

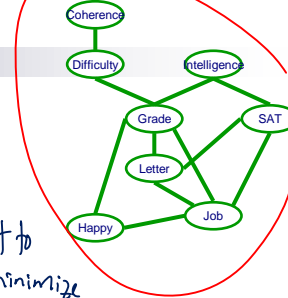
- Theorem:** $\ln Z = F[P_F, Q] + D(Q||P_F)$

- where energy functional:

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

$$Z = \sum_j q(c_j) \log \phi_j(c_j)$$

I know how to compute



10-708 — Carlos Guestrin, 2006-2008

19

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

$$D(q||p) \geq 0$$

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

$$F(P_F, Q) + D(Q||P_F) = \log Z$$

$$\log Z \geq F(P_F, Q) \leftarrow \text{what we maximize}$$

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

don't know how to compute Z , so we will try to find a lower bound

10-708 — Carlos Guestrin, 2006-2008

20

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., fully factorized (mean field) $q(x) = \prod_j Q_j(x_j)$

- Find Q that maximizes $F[P_{\mathcal{F}}, Q]$ $\leq \log Z$

- For mean field

$$\max_{Q_j} F[P_{\mathcal{F}}, \{Q_1, \dots, Q_n\}]$$

$$\text{subject to } Q_j(x_j) \geq 0$$

$$\sum_{x_j} Q_j(x_j) = 1$$

10-708 — Carlos Guestrin, 2006-2008

21

Optimization for mean field

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

$$\forall i, \sum_{x_i} Q_i(x_i) = 1$$

- Constrained optimization, solved via Lagrangian multiplier

□ 9 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 :

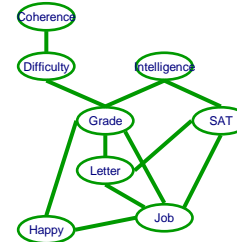
$$Q_i(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi \in \mathcal{F}} E_Q[\ln \phi | x_i] \right\}$$

10-708 — Carlos Guestrin, 2006-2008

22

Understanding fixed point equation

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



10-708 — Carlos Guestrin, 2006-2008

23

Q_i only needs to consider factors that intersect X_i

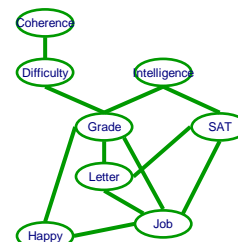
- **Theorem:** The fixed point:

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

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, 2006-2008

25

There are many stationary points!

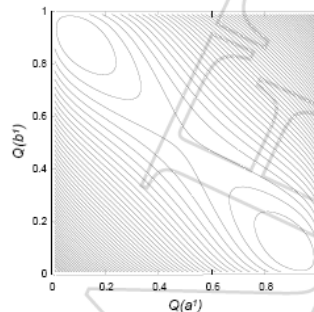


Figure 11.18 An example of a multi-modal mean field energy functional landscape. In this network, $P(a, b) = 0.25 - \epsilon$ if $a \neq b$ and ϵ if $a = b$. The axes correspond to the mean field marginal for A and B and the contours show equi-values of the energy functional.

10-708 — Carlos Guestrin, 2006-2008

26

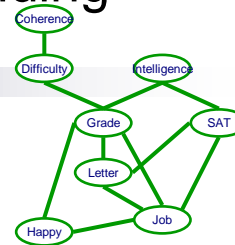
Very simple approach for finding one stationary point

- 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

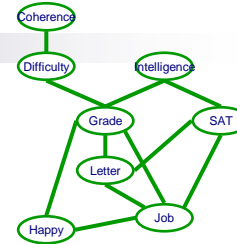


10-708 — Carlos Guestrin, 2006-2008

27

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:

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

- Very similar update rule

10-708 — Carlos Guestrin, 2006-2008

28

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, 2006-2008

31