

Local independence assumptions for a Markov network Separation defines global independencies Pairwise Markov Independence: Pairwise Markov Independences: Total Tot

Equivalence of independencies in Markov networks

- **Soundness Theorem**: For all <u>positive</u> distributions *P*, the following three statements are equivalent:
 - \square *P* entails the global Markov assumptions $P \models \bot(H)$
 - \square *P* entails the pairwise Markov assumptions $P \models \mathcal{T}_{PW}(H)$
 - □ P entails the local Markov assumptions (Markov blanket)

10-708 – ©Carlos Guestrin 2006

Minimal I-maps and Markov Networks

- М
 - A fully connected graph is an I-map
 - Remember minimal I-maps?
 - \square A "simplest" I-map \rightarrow Deleting an edge makes it no longer an I-map
 - In a BN, there is no unique minimal I-map
- Theorem: In a Markov network, minimal I-map is unique!!
- Many ways to find minimal I-map, e.g.,
 - □ Take pairwise Markov assumption:
 - ☐ If P doesn't entail it, add edge:

10-708 - ©Carlos Guestrin 2006

How about a perfect map?



- Remember perfect maps?
 - \Box independencies in the graph are exactly the same as those in P
- For BNs, doesn't always exist
 - □ counter example: Swinging Couples
- How about for Markov networks?

0-708 – ©Carlos Guestrin 200

Δ

Unifying properties of BNs and MNs



BNs:

- ☐ give you: V-structures, CPTs are conditional probabilities, can directly compute probability of full instantiation
- but: require acyclicity, and thus no perfect map for swinging couples

MNs:

- □ give you: cycles, and perfect maps for swinging couples
- but: don't have V-structures, cannot interpret potentials as probabilities, requires partition function

Remember PDAGS???

- □ skeleton + immoralities
- □ provides a (somewhat) unified representation
- see book for details

10-708 - ©Carlos Guestrin 2006

What you need to know so far about Markov networks



Markov network representation:

- □ undirected graph
- □ potentials over cliques (or sub-cliques)
- □ normalize to obtain probabilities
- □ need partition function

Representation Theorem for Markov networks

- □ if P factorizes, then it's an I-map
- □ if P is an I-map, only factorizes for positive distributions

Independence in Markov nets:

- □ active paths and separation
- □ pairwise Markov and Markov blanket assumptions
- □ equivalence for positive distributions
- Minimal I-maps in MNs are unique
- Perfect maps don't always exist

10-708 - ©Carlos Guestrin 2006

Some common Markov networks and generalizations

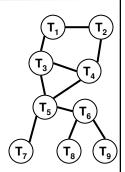
- Pairwise Markov networks
- A very simple application in computer vision
- Logarithmic representation
- Log-linear models
- Factor graphs

10-708 - ©Carlos Guestrin 2006

11

Pairwise Markov Networks

- All factors are over single variables or pairs of variables:
 - □ Node potentials
 - Edge potentials
- Factorization:



 Note that there may be bigger cliques in the graph, but only consider pairwise potentials

10-708 - ©Carlos Guestrin 2006

A very simple vision application



- Image segmentation: separate foreground from background
- Graph structure:
 - pairwise Markov net
 - □ grid with one node per pixel
- Node potential:
 - □ "background color" v. "foreground color"
- Edge potential:
 - □ neighbors like to be of the same class

10-708 - ©Carlos Guestrin 2006

13

Logarithmic representation



- Standard model: $P(X_1,...,X_n) = \frac{1}{Z} \prod_{i=1}^m \pi_i(\mathbf{D}_i)$
- Log representation of potential (assuming positive potential):
 - □ also called the energy function
- Log representation of Markov net:

10-708 - ©Carlos Guestrin 2006

Log-linear Markov network (most common representation)



- Feature is some function φ[D] for some subset of variables D
 - e.g., indicator function
- Log-linear model over a Markov network *H*:
 - \square a set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each **D**; is a subset of a clique in H
 - two ¢'s can be over the same variables
 - \square a set of weights $w_1, ..., w_k$
 - usually learned from data

$$\square P(X_1, ..., X_n) = \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i (\mathbf{D}_i) \right]$$

0-708 - ©Carlos Guestrin 2006

45

Structure in cliques



Possible potentials for this graph:



0-708 – ©Carlos Guestrin 2006

Factor graphs



- 1
- Very useful for approximate inference
 - □ Make factor dependency explicit
- Bipartite graph:
 - \square variable nodes (ovals) for $X_1,...,X_n$
 - \Box factor nodes (squares) for $\phi_1, ..., \phi_m$
 - \Box edge $X_i \phi_i$ if $X_i \in Scope[\phi_i]$
- More explicit representation, but exactly equivalent

10-708 - @Carlos Guestrin 2006

17

Exact inference in MNs and Factor Graphs



- Variable elimination algorithm presented in terms of factors → exactly the same VE algorithm can be applied to MNs & Factor Graphs
- Junction tree algorithms also applied directly here:
 - $\hfill\Box$ triangulate MN graph as we did with moralized graph
 - □ each factor belongs to a clique
 - $\hfill \square$ same message passing algorithms

10-708 - ©Carlos Guestrin 2006

Summary of types of Markov nets



- Pairwise Markov networks
 - □ very common
 - □ potentials over nodes and edges
- Log-linear models
 - □ log representation of potentials
 - □ linear coefficients learned from data
 - ☐ most common for learning MNs
- Factor graphs
 - □ explicit representation of factors
 - you know exactly what factors you have
 - □ very useful for approximate inference

10-708 - @Carlos Guestrin 2006

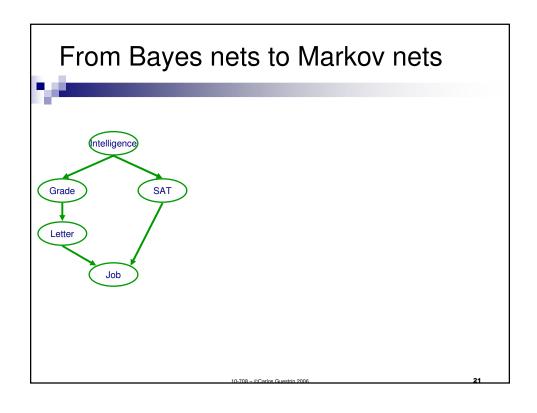
19

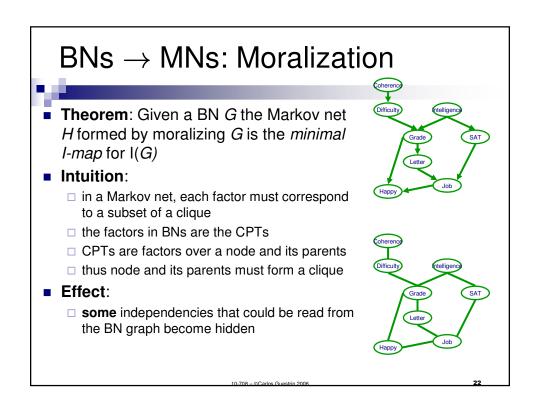
What you learned about so far

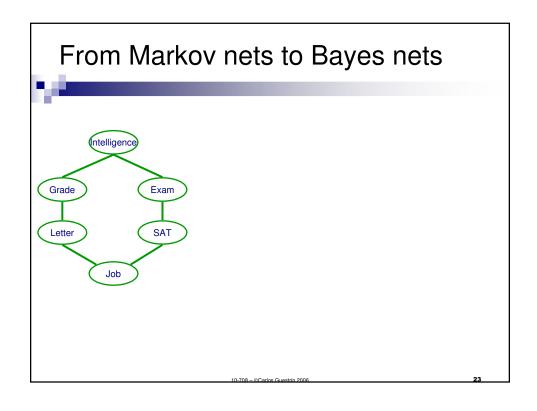


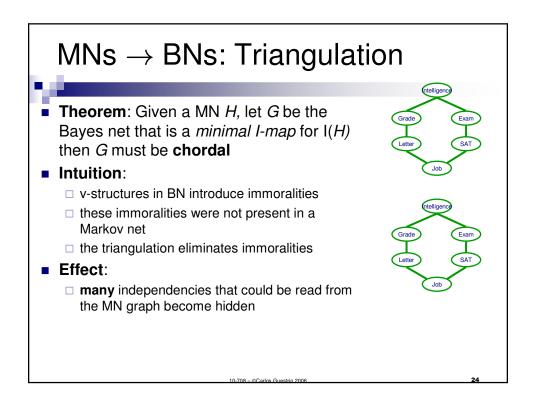
- Bayes nets
- Junction trees
- (General) Markov networks
- Pairwise Markov networks
- Factor graphs
- How do we transform between them?
- More formally:
 - □ I give you an graph in one representation, find an **I-map** in the other

10-708 - @Carlos Guestrin 2006









Markov nets v. Pairwise MNs

- Every Markov network can be transformed into a Pairwise Markov net
- A B
- □ 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

10-708 - ©Carlos Guestrin 2006

25

Overview of types of graphical models and transformations between them

0-708 – ©Carlos Guestrin 2006

Approximate inference overview

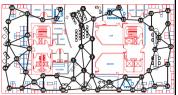
- ٧
 - So far: VE & junction trees
 - □ exact inference
 - □ exponential in tree-width
- There are many many many many approximate inference algorithms for PGMs
- We will focus on three representative ones:
 - sampling
 - □ variational inference
 - □ loopy belief propagation and generalized belief propagation
- There will be a special recitation by Pradeep Ravikumar on more advanced methods

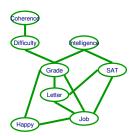
10-708 - ©Carlos Guestrin 2006

27

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

KL divergence: Distance between distributions

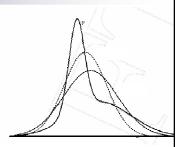
- Given two distributions p and q KL divergence:
- D(p||q) = 0 iff p=q
- Not symmetric p determines where difference is important
 - \Box p(x)=0 and q(x) \neq 0
 - \Box p(x) \neq 0 and q(x)=0

10-708 - @Carlos Guestrin 2006

29

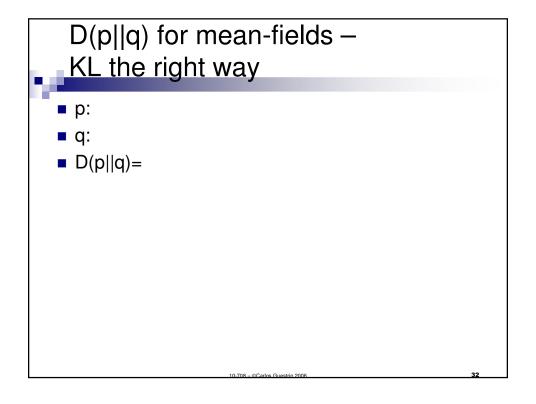
Find simple approximate distribution

- 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



10-708 - ©Carlos Guestrin 2006

Inference in a graphical model: P(x) = want to compute P(X_i|e) our p: What is the simplest q? every variable is independent: mean-fields approximation can compute any prob. very efficiently



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

- Ŋ
- p:
- **q**:
- D(p||q)=

10-708 - ©Carlos Guestrin 2006

33

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)
- The right KL is intractable, so we use the reverse KL

I0-708 – ©Carlos Guestrin 2006