

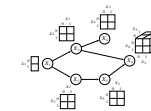
LECTURE 3:

UNDIRECTED GRAPHICAL MODELS

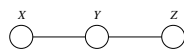
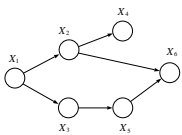
January 16, 2006

- Semantics:  $x \perp y | z$  if  $z$  d-separates  $x$  and  $y$
- d-separation:  $z$  d-separates  $x$  from  $y$  if along every undirected path between  $x$  and  $y$  there is a node  $w$  such that either:
  1.  $w$  has converging arrows along the path ( $\rightarrow w \leftarrow$ ) and neither  $w$  nor its descendants are in  $z$  or
  2.  $w$  does not have converging arrows along the path and  $w \in z$ .
- The “Bayes-Ball” algorithm can be used to check d-separation.
- It is always possible to find a distribution consistent with the graph. Most general such distribution is a product of parent-conditionals:

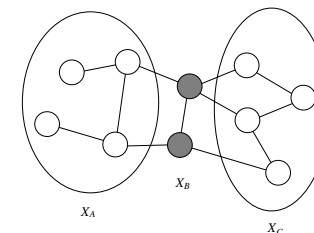
$$P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \prod_i P(\mathbf{x}_i | \mathbf{x}_{\pi_i})$$



- Graphical models aim to provide *compact factorizations* of large joint probability distributions.
- These factorizations are achieved using *local functions* which exploit *conditional independencies* in the models.
- The graph tells us a basic set of *conditional independencies* that must be true. From these we can derive more that also must be true. These independencies are crucial to developing efficient algorithms valid for *all* numerical settings of the local functions.
- Local functions tell us the quantitative details of the distribution.
- Certain numerical settings of the distribution may have more independencies present, but these do not come from the *graph*.

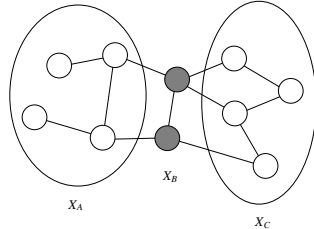


- Also graphs with one node per random variable and edges that connect pairs of nodes, but now the edges are undirected.
- Semantics: every node is conditionally independent from its non-neighbours given its neighbours, i.e.  $x_A \perp x_C | x_B$  if every path b/w  $x_A$  and  $x_C$  goes through  $x_B$



- Can model symmetric interactions that directed models cannot.
- aka Markov Random Fields, Markov Networks, Boltzmann Machines, Spin Glasses, Ising Models

- In undirected models, simple graph separation (as opposed to d-separation) tells us about conditional independencies.
- $\mathbf{x}_A \perp \mathbf{x}_C | \mathbf{x}_B$  if every path between  $\mathbf{x}_A$  and  $\mathbf{x}_C$  is blocked by some node in  $\mathbf{x}_B$ .



- “Markov Ball” algorithm:  
remove  $\mathbf{x}_B$  and see if there is any path from  $\mathbf{x}_A$  to  $\mathbf{x}_C$ .

- In directed models, we started with  $p(\mathbf{x}) = \prod_i p(\mathbf{x}_i | \mathbf{x}_{\pi_i})$  and we derived the d-separation semantics from that.
- Undirected models: have the semantics, need parametrization.
- What about this “conditional parameterization”?

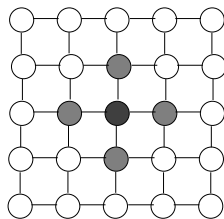
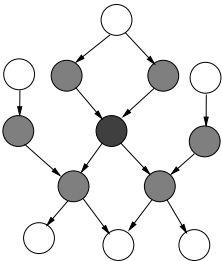
$$p(\mathbf{x}) = \prod_i p(\mathbf{x}_i | \mathbf{x}_{\text{neighbours}(i)})$$

- Good: product of local functions.  
Good: each one has a simple conditional interpretation.  
Bad: local functions cannot be arbitrary, but must agree properly in order to define a valid distribution.

- $\mathbf{b}$  is a “Markov blanket” for  $\mathbf{x}$  iff

$$\mathbf{x} \perp \mathbf{y} | \mathbf{b} \quad \forall \mathbf{y} \notin \mathbf{b}$$

- “Markov boundary”: minimal Markov blanket
- For undirected models, this is set of neighbours.
- Q: What is the Markov blanket (boundary) in a directed model?  
A: {parents+children+parents-of-children}



- OK, what about this “marginal parameterization”?

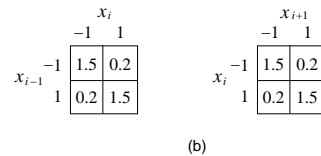
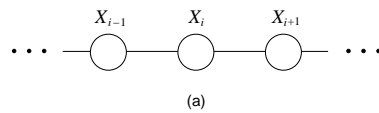
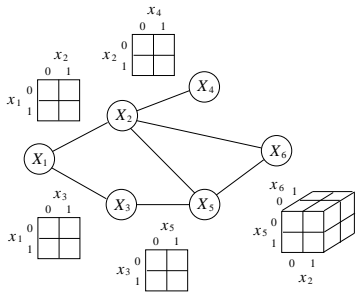
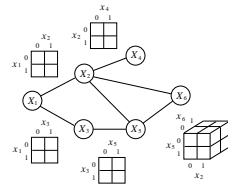
$$p(\mathbf{x}) = \prod_i p(\mathbf{x}_i, \mathbf{x}_{\text{neighbours}(i)})$$

- Good: product of local functions.  
Good: each one has a simple marginal interpretation.  
Bad: only very few pathological marginals on overlapping nodes can be multiplied to give a valid joint.

- Whatever factorization we pick, we know that only connected nodes can be arguments of a single local function.
- A *clique*  $\mathbf{x}_c$  is a fully connected subset of nodes.
- Thus, consider using a *product of positive clique potentials*:

$$P(\mathbf{x}) = \frac{1}{Z} \prod_{\text{cliques } c} \psi_c(\mathbf{x}_c) \quad Z = \sum_{\mathbf{x}} \prod_{\text{cliques } c} \psi_c(\mathbf{x}_c)$$

- The product of functions that don't need to agree with each other.
- Still factors in the way that the graph semantics demand.
- Without loss of generality we can restrict ourselves to *maximal cliques*. (Why?)



$$p(\mathbf{x}) = \psi_A(x_1, x_2)\psi_B(x_1, x_3) \psi_C(x_2, x_4)\psi_D(x_3, x_5) \psi_E(x_2, x_5, x_6)$$

$$p(\mathbf{x}) = \prod_i \psi(x_i, x_{i+1})$$

- We often represent the clique potentials using their logs:

$$\psi_C(\mathbf{x}_C) = \exp\{-H_C(\mathbf{x}_C)\}$$

for arbitrary real valued “energy” functions  $H_C(\mathbf{x}_C)$ .  
The negative sign is a standard convention.

- This gives the joint a nice additive structure:

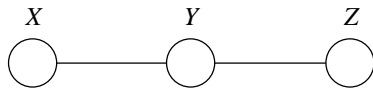
$$P(\mathbf{x}) = \frac{1}{Z} \exp\left\{-\sum_{\text{cliques } C} H_C(\mathbf{x}_C)\right\} = \frac{1}{Z} \exp\{-H(\mathbf{x})\}$$

where the sum in the exponent is called the “free energy”:

$$H(\mathbf{x}) = \sum_C H_C(\mathbf{x}_C)$$

- This way of defining a probability distribution based on energies is the “Boltzmann distribution” from statistical physics.

- Normalizer  $Z(\mathbf{x})$  above is called the “partition function”.
- Computing the normalizer and its derivatives can often be the hardest part of inference and learning in undirected models.
- Often the factored structure of the distribution makes it possible to efficiently do the sums/integrals required to compute  $Z$ .
- Don't *always* have to compute  $Z$ , e.g. for conditional probabilities.



- The model implies  $x \perp z \mid y$

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{y})p(\mathbf{x}|\mathbf{y})p(\mathbf{z}|\mathbf{y})$$

- We can write this as:

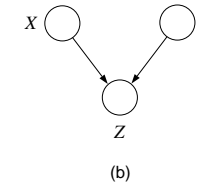
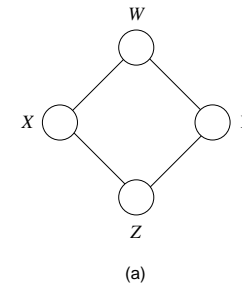
$$p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{x}, \mathbf{y})p(\mathbf{z}|\mathbf{y}) = \psi_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y})\psi_{\mathbf{y}\mathbf{z}}(\mathbf{y}, \mathbf{z})$$

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{z}, \mathbf{y}) = \psi_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y})\psi_{\mathbf{y}\mathbf{z}}(\mathbf{y}, \mathbf{z})$$

cannot have all potentials be marginals  
cannot have all potentials be conditionals

- The positive clique potentials can only be thought of as general “compatibility”, “goodness” or “happiness” functions over their variables, but not as probability distributions.

- Can we always convert directed  $\leftrightarrow$  undirected?
- No.

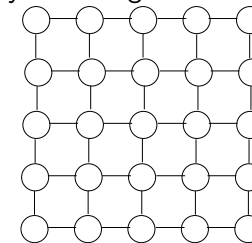


No directed model can represent these and only these independencies.  
 $x \perp y \mid \{w, z\}$   
 $w \perp z \mid \{x, y\}$

No undirected model can represent these and only these independencies.  
 $x \perp y$

- H-C theorem tells us that the family of distributions defined by the conditional independence semantics on the graph and the family defined by products of potential functions\* on maximal cliques are the same. (\* arbitrary real valued, but strictly positive)
- For directed models, there is a version of this theorem which tells us that the family of distributions defined by the conditional independencies semantics of the directed graph and the family defined by products of parent-conditionals are the same.
- Notice the crucial difference between *graphs*, which tells us independencies that are true *no matter what* local functions we choose, and *numerical functions* which could introduce some extra independencies, once we know them.

- Common model for binary nodes: spin-glass/ Ising lattice.
- Nodes are arranged in a regular topology (often a regular packing grid) and connected only to their geometric neighbours.



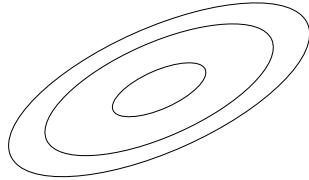
- For example, if we think of each node as a pixel, we might want to encourage nearby pixels to have similar intensities.
- Energy is of the form:

$$H(\mathbf{x}) = \sum_{ij} \beta_{ij} \mathbf{x}_i \mathbf{x}_j + \sum_i \alpha_i \mathbf{x}_i$$

- The most common and important undirected graphical model on a set of continuous valued nodes is the Gaussian (normal).
- It uses pairwise potentials between every pair of nodes to define an energy identical to the Ising model, but for continuous values:

$$H(\mathbf{x}) = \sum_{ij} (\mathbf{x}_i - \mu_i) V_{ij} (\mathbf{x}_j - \mu_j)$$

where  $\mu$  is the mean and  $V$  is the inverse covariance matrix.

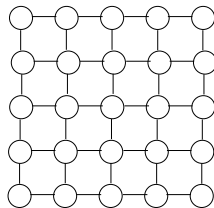


- Like a “fully connected” lattice.  
Also, the Gaussian is the maximum entropy distribution consistent with the mean and covariance defined by  $\mu$  and  $V$ .

- Fully observed Boltzmann machines are the binary equivalent of a Gaussian distribution: fully connected Ising models on a set of binary random variables. (Also maxent.)

Energy is the same:

$$H(\mathbf{x}) = \sum_{ij} \beta_{ij} \mathbf{x}_i \mathbf{x}_j + \sum_i \alpha_i \mathbf{x}_i$$



- 
- Boltzmann machines also add the possibility of having some units (random variables) which are never observed. These are called “hidden units” or “latent variables” and we will see much more about them later.
- For continuous variables, the equivalent of a Boltzmann machine with hidden units is called a “factor analysis” model.