

Lecture 5: Undirected graphical models

- Semantics of probabilistic models over undirected graphs
- Parameters of undirected models
- Example applications

Undirected graphical models

- So far we have used directed graphs as the underlying structure of a Bayes net
- Why not use *undirected* graphs as well?
E.g., variables might not be in a “causality” relation, but they can still be correlated, like the pixels in a neighborhood in an image
- An undirected graph over a set of random variables $\{X_1, \dots, X_n\}$ is called a **undirected graphical model** or **Markov random field (MRF)** or **Markov network**

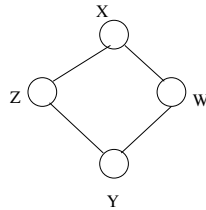
Conditional independence

- We need to be able to specify, for a given graph, if $X \perp\!\!\!\perp Z|Y$, for any disjoint subsets of nodes X, Y, Z .
- In directed graphs, we did this using the Bayes Ball algorithm
- In undirected graphs, independence can be established simply by graph separation: if every path from a node in X to a node in Z goes through a node in Y , we conclude that $X \perp\!\!\!\perp Z|Y$
- Hence, independence can be established by removing the nodes in the conditioning set then doing reachability analysis on the remaining graph.
- What is the Markov blanket of a node in an undirected model?

How expressive are undirected models?

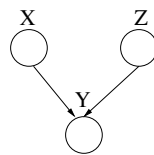
- Are undirected models more expressive than directed models?
I.e. for any directed model, can we find an undirected model that satisfies exactly the same conditional independence relations?
- Are undirected models less expressive?
I.e. for any undirected model, can we find a directed model that satisfies exactly the same conditional independencies?

Example: An undirected graph



Can we find a directed graph that satisfies the same independence relations?

Example: A directed graph



Can we find an undirected graph that satisfies the same independence relations?

Expressiveness of undirected models

- Undirected models are neither more nor less expressive than directed models; they are simply different
- The semantics of an undirected model naturally capture *correlation* of r.v.s, not causation
- If you ever want, in an application, to write a Bayes net with cycles, it is a sign that the right model is undirected.

Local parameterization

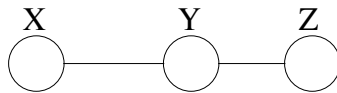
- In directed models, we had local probability models (CPDs) attached to every node, giving the conditional probability of the corresponding random variable given its parents
- The joint probability distribution expressed by a directed model factorizes over the graph
- This means that the joint can be written as a product of “local” factors, which depend on subsets of the variables.
- We want a similar property for directed models.
- But what should the local factors be?

What about local marginal parameterizations?

- Suppose we express the joint as:

$$p(X_1, \dots, X_n) = \prod_i p(X_i, \text{Neighbors}(X_i))$$

- It is local and has a nice interpretation
- So consider using it for an example:



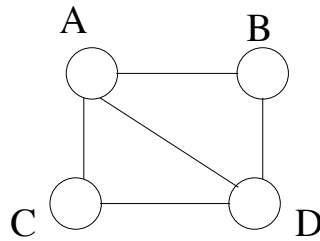
Local parameterizations: Try 2

- Consider a pair of nodes X and Y that are not directly connected through an arc
- According to the conditional independence interpretation, X and Y are independent given all the other nodes in the graph

$$X \perp\!\!\!\perp Y \mid \{X_1, \dots, X_n\} - X - Y$$

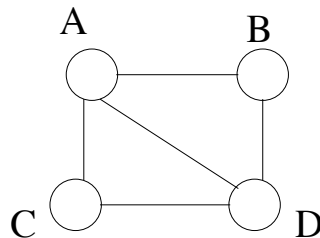
- Hence, there must be a factorization in which they do not appear in the same factor
- This suggests that we should define factors on **cliques**
Recall that a clique is a fully connected subset of nodes (i.e., there is an arc between every pair of nodes)

Example: what are the cliques?



Defining parameters on cliques

The main idea is that if variables do not have an arc between them, they are conditionally independent given the rest of the graph, and hence should not be in the same local model.



Clique potentials

- We will represent the joint distribution as a **product of clique potentials**:

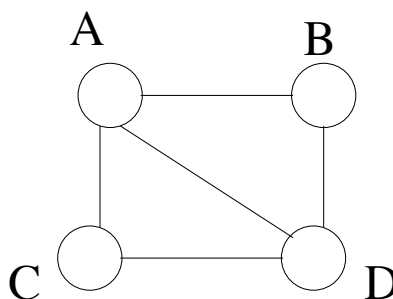
$$p(X_1 = x_1, \dots, X_n = x_n) = \frac{1}{Z} \prod_{\text{cliques } C} \psi_C(\mathbf{x}_C)$$

where \mathbf{x}_C are the values for the variables that participate in clique C and Z is a normalization constant, to make probabilities sum to 1:

$$Z = \sum_{\mathbf{x}} \prod_{\text{cliques } C} \psi_C(\mathbf{x}_C)$$

- Without loss of generality, we can consider only maximal cliques
These are the cliques that cannot be extended with other nodes without losing the fully connected property

Example

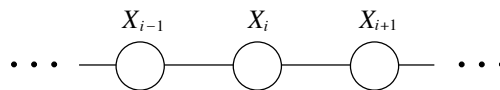


Normalizing constant

- The normalizing constant Z can be ugly to compute, since we have to sum over all possible assignments of values to variables
- Depending on the shape of the graph, the summation could be done efficiently
- However, if we are interested in conditional probabilities, we do not even need to compute it! (why?)

Interpretation of clique potentials

- Potentials are *NOT* probabilities (conditional or marginal)
- But they do have a natural interpretation as “agreement” or “energy”
- Example: spin glass model



(a)

		x_i			x_{i+1}		
		-1	1		-1	1	
	-1	1.5	0.2		-1	1.5	0.2
	1	0.2	1.5		1	0.2	1.5

(b)

More on spin glasses

- In general, a spin glass is a collection of magnetic moment (spins) whose low temperature state is disordered.
- They have been studied a lot in statistical physics and they can model many practical materials
- These models have two important features:
 - There is competition among the interactions between moments, so there is no configuration of spins that is favored by all interactions; this is called **frustration**
 - Interactions are at least partially random
- There are many states whose energy is locally optimal (low)
- Finding such a state can be done by probabilistic inference.

Boltzmann (Gibbs) distribution

- The fact that potentials must be non-negative is annoying
- We can escape from that by using the exponential function, which is non-negative:

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

- Now we have to define $H_C(\mathbf{x}_C)$, which can be anything!
- Moreover, the joint also has a nice form:

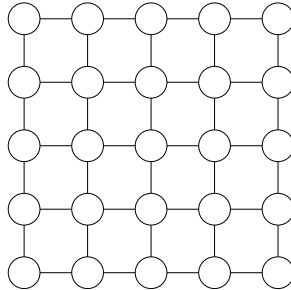
$$p(\mathbf{x}) = \frac{1}{Z} \prod_C e^{-H_C(\mathbf{x}_C)} = \frac{1}{Z} e^{-\sum_C H_C(\mathbf{x}_C)} = \frac{1}{Z} e^{-H(\mathbf{x})}$$

where $H(\mathbf{x}) = \sum_C H_C(\mathbf{x}_C)$ is the “free energy”

- Hence, p is represented using a Boltzmann distribution

Special case: Ising Model

- All r.v.s are binary and nodes are arranged in a regular fashion and connected only to geometric neighbors.
- E.g., Spin glass in 2D:



- Energy has the form:

$$H(\mathbf{x}) = \sum_{i,j} \beta_{ij} x_i x_j + \sum_i \alpha_i x_i$$

Applications of the Ising model

- Very popular for explaining the effect of “society” or “environment” on a “component” or “individual”
 - Flocking behavior
 - Behavior of neural networks
 - Sociology studies
- In all these cases, the effort is both to find, from the data, what the model should be, as well as to use inference in order to determine what will be the next state of minimum energy to which the model “settles”

Choosing the parameters of a Markov network

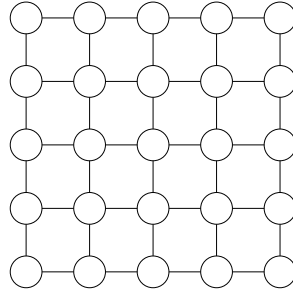
- These nets often have a regular structures, and parameters may be similar (or identical) in all cliques of a given type
- As a result, optimization or learning are often the preferred ways of coming up with parameters

Important result (for strictly positive distributions)

- Consider the family of probability distributions that respect all the conditional independencies implied by an undirected graph G . These are the distributions that satisfy the *global Markov properties* of the graph
- Consider the family of probability distributions defined by ranging over all allowed maximal clique potential functions. These are the distributions that *factorize* on the graph G .
- The **Hammersley-Clifford theorem** shows that these two families are identical.
- This is a similar result to the “soundness and completeness” of d-separation which we discussed for directed models.

A real example: Texture synthesis

- You are given a small patch of texture and want to produce a “similar” larger patch
- We can define a Markov random field over pixels, e.g:



- The “potentials” favor certain configurations of pixels over others
- We get the texture by doing inference (and sometimes learning) for this model

More applications

- MRFs are used extensively in computer vision, e.g for labeling tasks
- Labeling can be low-level, like labeling edges or other pixel configurations, or high-level, like labeling objects in an image
- Images often obey constraints (e.g. smoothness of surfaces, texture) which can be captured easily as a MRF structure
- Labeling then becomes a search for a pattern of minimum energy, which is often solved by optimization
- Learning can help establish the parameters of the model.