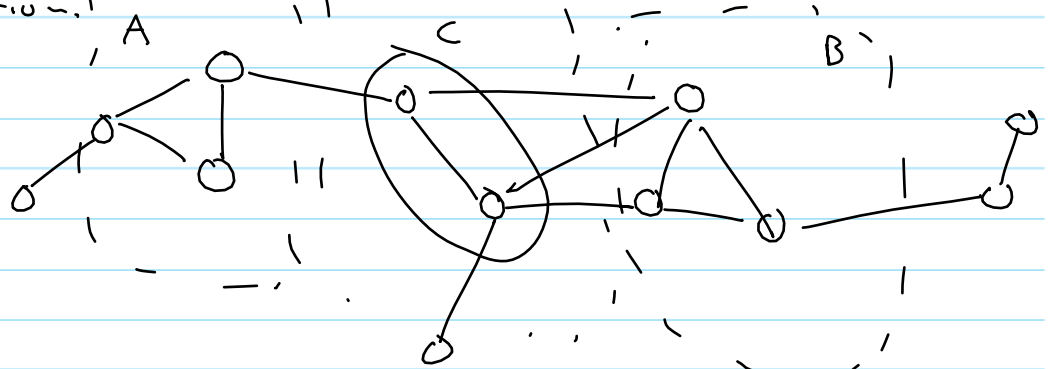


Markov Random Fields

- Described by undirected graphical models
- Major class of graphical models
- Also specify factorization & conditional independences
* though differently

MRF factorization - corresponding to graph separation



To test if $A \perp\!\!\!\perp B \mid C$ we look at all possible paths from all nodes in A to all nodes in B. If all pass through one or more nodes in C then such paths are blocked.

- No explaining away
- testing conditional independence simpler

Factorization properties

- if x_i and x_j are not connected by a link
 - $P(x_i, x_j \mid \vec{X}_{\setminus \{i, j\}}) = P(x_i \mid \vec{X}_{\setminus \{i, j\}}) P(x_j \mid \vec{X}_{\setminus \{i, j\}})$

because there is no direct path and all other nodes observed

How do we parameterize an MRF?

- Clique potentials, all other parameterizations would be redundant

Let C be a clique and the var's in that clique are \vec{x}_c , then the joint dist. is written as a product of pot. funcs $\psi_c(\vec{x}_c)$ over the max cliques of the graph.

$$p(\vec{x}) = \frac{1}{z} \prod_c \psi_c(\vec{x}_c)$$

where z , the "partition function" is a normalizing constant

$$z = \sum_x \prod_c \psi_c(\vec{x}_c)$$

we will consider only potential function $\psi_c(x_c) \geq 0$

- Note $\psi_c(x_c)$ need not be a prob. function!
- z is a problem, in model with M discrete K state nodes \sum_x involves K^M states
- partition function is usually needed for learning because it will be a function of whatever parameters the potential functions have.
- to compute local conditionals the partition func is not needed because the cond. is a ratio of terms that cancel.
- local marginals can be locally normalized

Conditional independence

Requires $\psi_c(x_c) > 0$
 \uparrow w/o se

Beyond scope for the class. the Hammersley-Clifford theorem states that the set of distributions that can be factorized into the product of maximal cliques is the same as the set of distributions whose conditional independence structure can be read off of the graph separability.

Often potentials are expressed as exponentials, i.e.

$$\psi_c(\vec{x}_c) = \exp(-E(x_c)) \quad \text{where}$$

$E(\vec{x}_c)$ is called an energy function.

Example application Image denoising

Binary images $y_i \in \{-1, +1\}$ pixels observed
 $i = 1 \dots N$

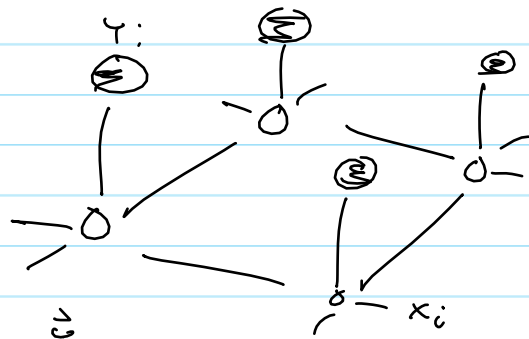
$x_i \in \{-1, +1\}$ noisy observed
 image

noise = flipping a bit w.p. 10%

assume strong correlation between x_i & y_i
 and between neighboring pixels x_i & x_j

$$E(\vec{x}, \vec{y}) = h \underbrace{\sum_i x_i}_{\text{bias}} - \beta \underbrace{\sum_{\{i,j\}} x_i x_j}_{\text{agreement between neighbors}} - \eta \underbrace{\sum_i x_i y_i}_{\text{agreement w/ observation}}$$

- explain signs



if we fix \vec{y}

this defines a conditional dist. $p(\vec{x} | \vec{y})$ over noise-free images. A goal is to find an \vec{x} that has high probability.

One approach, Iterated Conditional Modes

- basically coordinate-wise gradient ascent
- initialize all $x_i = \tau_i$, evaluate
- conditionals for single x_i and set x_i to one with lowest energy
- repeat

Relationship to directed graphs

Directed graph \rightarrow undirected

example

$$P(\vec{x}) = p(x_1) P(x_2 | x_1) \cdots P(x_N | x_{N-1})$$

to convert this to cliques simply write

$$P(\vec{x}) = \frac{1}{Z} \Psi_{1,2}(x_1, x_2) \Psi_{2,3}(x_2, x_3) \cdots \Psi_{x_{n-1}, x_n}(x_{n-1}, x_n)$$

where

$$\Psi_{1,2}(x_1, x_2) = P(x_1) P(x_2 | x_1)$$

$$\Psi_{i,j}(x_i, x_j) = P(x_j | x_i) \quad \forall \{i,j\} \neq \{1,2\}$$

and

$$Z = \int$$

Moralization

More generally this conversion requires "marrying the parents". In this ^{chain} example the moral graph is complete.

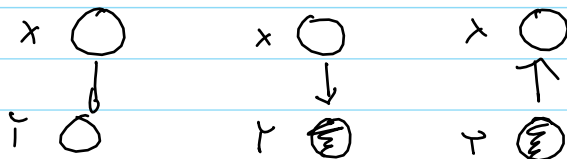
Recipe: directed G.M. \rightarrow undirected G.M.

- 1) Add links between all pairs of parents for all nodes in graph
- 2) Drop arrows
- 3) Initial clique potentials to 1
 - a) Multiply in all conditional dist's associated with each clique.
- 4) $z = 1$

* Inference in Graphical Models * \leftarrow KEY

Idea: exploit graphical structure in algorithms for inference.

First graphical Bayes Theorem



Joint $P(x, Y) = P(x) \cdot P(Y|x)$

If we obs. Y then $P(x)$ can be seen as a prior on x , and inferring the post. dist. of x can be the goal. To do this note

$$P(Y) = \sum_{x'} P(x', Y) = \sum_{x'} P(x') P(Y|x')$$

and

$$P(x|Y) = \frac{P(Y|x) P(x)}{P(Y)}, \text{ reversing the arrow}$$