

COMPSCI 688: Probabilistic Graphical Models

Lecture 8: Undirected Graphical Models: Inference

Dan Sheldon

Manning College of Information and Computer Sciences
University of Massachusetts Amherst

Partially based on materials by Benjamin M. Marlin (marlin@cs.umass.edu) and Justin Domke (domke@cs.umass.edu)

1 / 33

Review

2 / 33

Markov Random Fields

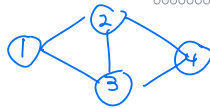
- ▶ Markov random field

x_1, \dots, x_N

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \phi_c(\mathbf{x}_c)$$

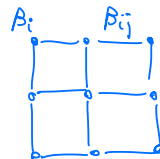
$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} \cdot \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4)$$

$$x_1 \perp x_4 \mid x_2, x_3$$



- ▶ *Dependence graph* \mathcal{G} : where nodes i and j are connected by an edge if they appear together in some factor
- ▶ *Ising Model*: grid-structured graph, unary/pairwise potentials express local preferences for values of x_i or (x_i, x_j) pairs

$$p(\mathbf{x}) = \frac{1}{Z} \prod_i \beta_i(x_i) \prod_{(i,j) \in E} \beta_{ij}(x_i, x_j)$$



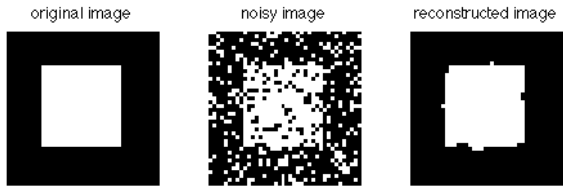
3 / 33

Examples

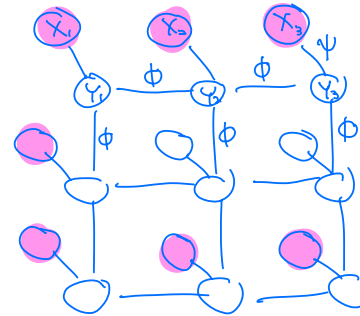
4 / 33

Example: Statistical Image Models

The Ising model with pairwise potentials encourages smoothness and can be used as a model for images for denoising:



Example: Image Denoising



$$p(\vec{Y} | \vec{X})$$

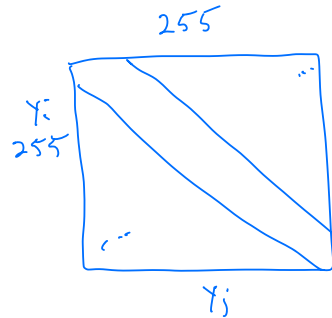
$Y_i \in \{0, \dots, 255\}$
unknown true pixel values

$X_i \in \{0, \dots, 255\}$
observed noisy pixel values

$$p(\vec{X}, \vec{Y}) = \frac{1}{Z} \cdot \prod_{(i,j) \in E} \Phi(y_i, y_j) \cdot \prod_i \Psi(x_i, y_i)$$

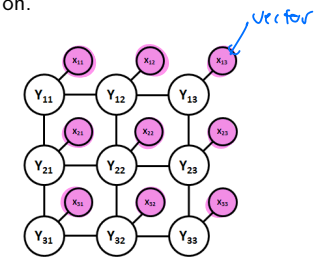
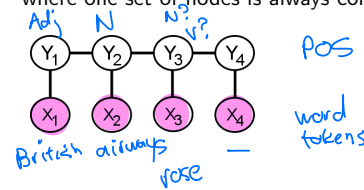
$$\Psi(y_i, y_j) = \exp\left(\frac{1}{T} y_i y_j\right)$$

↑ ↑
 $\{0, \dots, 255\}$



Conditional Random Fields

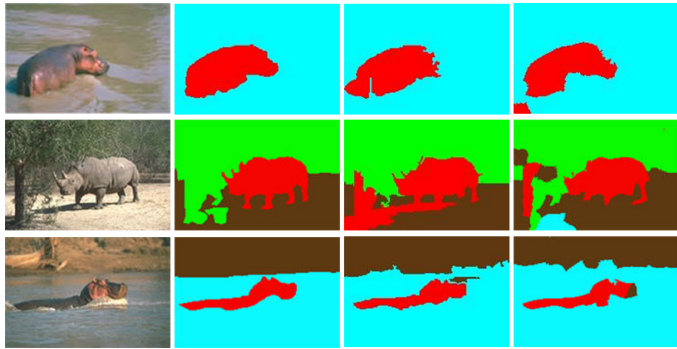
The image denoising model is one example of a **conditional random fields (CRFs)**, a very important model class in machine learning. A CRF is essentially a Markov network where one set of nodes is always conditioned on.



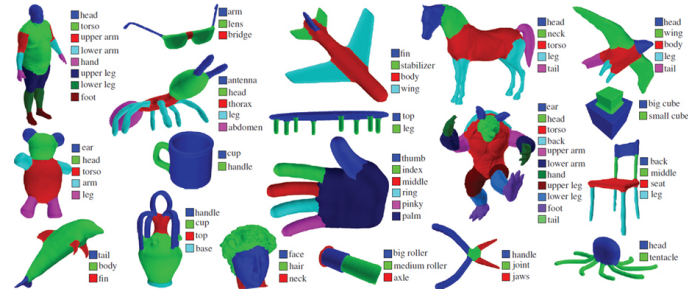
The y nodes are *labels*, and the x nodes are *features*.

$$\Psi(\vec{X}_i, y_i)$$

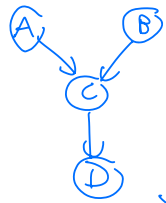
Example: Image Segmentation



Example: 3D Mesh Segmentation

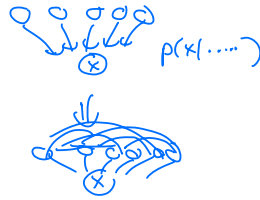
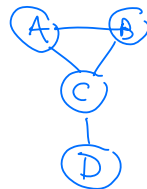


Example: Bayes Nets as MRFs



$$p(a, b, c, d) = p(a) p(b) p(c|a, b) p(d|c)$$

$$= \phi_1(a) \phi_2(b) \phi_3(a, b, c) \phi_4(c, d)$$

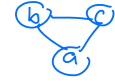


Example: Bayes Nets as MRFs

Some structure is lost in this transformation. When we replace $p(a|b, c)$ by $\phi(a, b, c)$, we “forget” that a Bayes net is **locally normalized**



$$\sum_a \phi(a, b, c) = 1 \quad \forall b, c.$$



This is a special property of Bayes nets and is central to V-structures, explaining away, and D-separation. It occurs “internally” to the factor $\phi(a, b, c)$ and is not represented in the MRF graph structure.

Similarly, when we replace $\prod_i p(x_i | \mathbf{x}_{pa(i)})$ by $\frac{1}{Z} \prod_{c \in C} \phi_c(\mathbf{x}_c)$, we “forget” that a Bayes net is **globally normalized**:

$$\sum_x \prod_{c \in C} \phi_c(\mathbf{x}_c) = 1 \implies Z = 1.$$

$$\sum_a \sum_b p(a) \phi(b|a) = 1$$

This is another special property of Bayes nets that makes learning easy.

Inference: Conditioning

Inference in Markov Networks

$$p(x) = p(x_Q, x_U, x_E) = \frac{1}{Z} \prod_c \phi_c(x_c)$$

- ▶ Given a Markov network, the main task is *probabilistic inference*, which means answering probability queries of the form

$$p(x_Q | x_E) = \sum_{x_U} p(x_Q, x_U | x_E)$$

- ▶ condition on *evidence variables* x_E
- ▶ marginalize *unobserved variables* x_U
- ▶ compute the joint distribution over *query variables* x_Q
- ▶ ... often by transforming Markov network into one with fewer or simpler factors
- ▶ Conditioning is easy
- ▶ Marginalization is hard! (formally NP-hard, easy in some cases)

Conditioning: Single Factor

Suppose we have a single-factor MRF $p(x_1, x_2) = \frac{1}{Z} \phi(x_1, x_2)$ for two binary variables. We are given a fixed value for x_2 , and want an MRF for $p(x_1 | x_2)$, i.e.:

$$p(x_1 | x_2) = \frac{1}{Z'} \phi'(x_1) \quad \phi'(x) := \phi(x, x_2)$$

Observe

$$p(x_1 | x_2) = \frac{p(x_1, x_2)}{p(x_2)} = \frac{1}{p(x_2)} \cdot \frac{1}{Z} \phi(x_1, x_2) = \frac{1}{Z'} \cdot \phi'(x_1)$$

For fixed x_2 , the conditional $p(x_1 | x_2)$ is *proportional* to the joint $p(x_1, x_2)$. We can use the same factor, but hard-code x_2 so that only x_1 is a free variable:

$$\phi'(x_1) = \phi(x_1, x_2), \quad Z' = p(x_2)Z$$

$$p(x_1, x_2) = \frac{1}{Z} \phi(x_1, x_2)$$

x_1	x_2	ϕ	$p(x_1, x_2)$
0	0	1	.1
0	1	4	.4
1	0	3	.3
1	1	2	.2
			$Z=10$

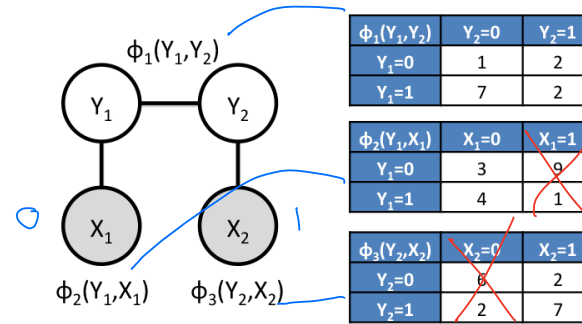
observe
 $x_2=0$

x_1	ϕ'	$p(x_1 x_2)$
0	1	.25
1	3	.75
		$Z'=4$

Conditioning: General Case

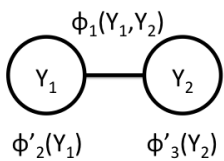
For a general MRF, we can apply the same reasoning to *reduce* every factor by hard-coding the evidence variables

Factor Reduction: Example



Query: $P(Y_1, Y_2 \mid X_1=0, X_2=1)$

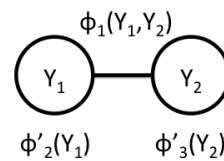
Factor Reduction: Step 1



$\phi_1(Y_1, Y_2)$	$Y_2=0$	$Y_2=1$
$Y_1=0$	1	2
$Y_1=1$	7	2
$\phi_2(Y_1, X_1)$	$X_1=0$	$X_1=1$
$Y_1=0$	3	9
$Y_1=1$	4	1
$\phi_3(Y_2, X_2)$	$X_2=0$	$X_2=1$
$Y_2=0$	6	2
$Y_2=1$	2	7

Query: $P(Y_1, Y_2 \mid X_1=0, X_2=1)$

Factor Reduction: Step 2



$\phi_1(Y_1, Y_2)$	$Y_2=0$	$Y_2=1$
$Y_1=0$	1	2
$Y_1=1$	7	2
$\phi_2(Y_1, X_1)$	$X_1=0$	$X_1=1$
$Y_1=0$	3	9
$Y_1=1$	4	1
$\phi_3(Y_2, X_2)$	$X_2=0$	$X_2=1$
$Y_2=0$	6	2
$Y_2=1$	2	7

Query: $P(Y_1, Y_2 \mid X_1=0, X_2=1)$

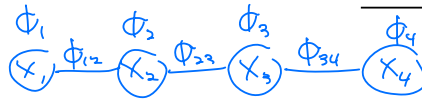
Example: Variable Elimination on a Chain

Consider the following MRF on a four-node "chain" graph:

$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} \phi_1(x_1) \phi_2(x_2) \phi_3(x_3) \phi_4(x_4) \phi_{12}(x_1, x_2) \phi_{23}(x_2, x_3) \phi_{34}(x_3, x_4)$$

x_i	$\phi_i(x_i)$
0	1
1	2

x_i	x_j	$\phi_{ij}(x_i, x_j)$
0	0	2
0	1	1
1	0	1
1	1	2



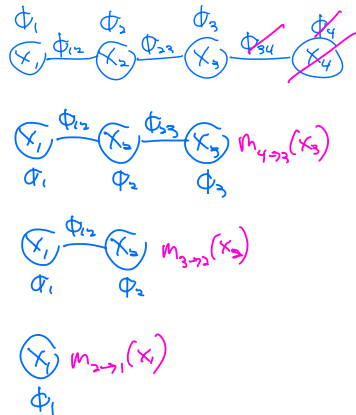
Let's compute Z :

$$\begin{aligned} &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \phi_1(x_1) \phi_2(x_2) \phi_3(x_3) \phi_4(x_4) \phi_{12}(x_1, x_2) \phi_{23}(x_2, x_3) \phi_{34}(x_3, x_4) \\ &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \phi_1(x_1) \phi_2(x_2) \phi_3(x_3) \phi_{12}(x_1, x_2) \phi_{23}(x_2, x_3) \sum_{x_4} \phi_4(x_4) \phi_{34}(x_3, x_4) \\ &= \sum_{x_1} \sum_{x_2} \phi_1(x_1) \phi_2(x_2) \phi_{12}(x_1, x_2) \sum_{x_3} \phi_3(x_3) \phi_{23}(x_2, x_3) m_{4 \rightarrow 3}(x_3) \\ &= \sum_{x_1} \phi_1(x_1) \sum_{x_2} \phi_2(x_2) \phi_{12}(x_1, x_2) m_{3 \rightarrow 2}(x_2) \\ &= \sum_{x_1} \phi_1(x_1) \cdot m_{2 \rightarrow 1}(x_1) \end{aligned}$$

We eliminated x_4, x_3, x_2, x_1

Pictorially, this is how we changed the MRF

x_3	x_4	$\phi_3 \phi_4$	$m_{4 \rightarrow 3}(x_3)$
0	0	1 · 2 = 2	4
0	1	2 · 1 = 2	1
1	0	1 · 1 = 1	1
1	1	2 · 2 = 4	5



What if we want to compute the *unnormalized* marginal $\hat{p}(x_1)$?

$$\begin{aligned} \hat{p}(x_1) &= \sum_{x_2} \sum_{x_3} \sum_{x_4} \phi_1(x_1) \phi_2(x_2) \phi_3(x_3) \phi_4(x_4) \phi_{12}(x_1, x_2) \phi_{23}(x_2, x_3) \phi_{34}(x_3, x_4) \\ &= \dots \\ &= \dots \\ &= \phi_1(x_1) \cdot m_{2 \rightarrow 1}(x_1) \end{aligned}$$

What if we want to compute the *actual* marginal $p(x_1)$?

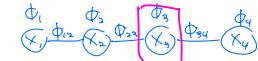
Take $\hat{p}(x_1)$ and normalize it

$$Z = \sum_{x_1} \hat{p}(x_1), \quad p(x_1) = \frac{1}{Z} \hat{p}(x_1)$$

Lesson: always normalize at the end

x_1	$\hat{p}(x_1)$	$p(x_1)$
0	3	3/20
1	17	17/20
	20	

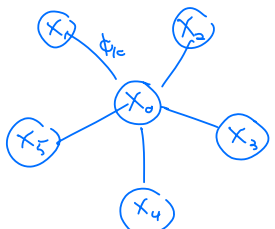
What if we eliminate x_3 first?



$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \phi_1(x_1) \phi_2(x_2) \phi_3(x_3) \phi_4(x_4) \phi_{12}(x_1, x_2) \phi_{23}(x_2, x_3) \phi_{34}(x_3, x_4) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_4} \phi_1(x_1) \phi_2(x_2) \phi_4(x_4) \phi_{12}(x_1, x_2) \underbrace{\sum_{x_3} \phi_3(x_3) \phi_{23}(x_2, x_3) \phi_{34}(x_3, x_4)}_{\tau_{24}(x_2, x_4)} \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_4} \phi_1(x_1) \phi_2(x_2) \phi_4(x_4) \phi_{12}(x_1, x_2) \tau_{24}(x_2, x_4) \\
 &= \dots
 \end{aligned}$$

Correct, but less efficient due to larger intermediate factor

What if our graph is a star graph?



$$\begin{aligned}
 Z &= \sum_{x_0} \sum_{x_1} \dots \sum_{x_5} \phi_0(x_0, x_1) \dots \phi_5(x_0, x_5) \\
 &= \sum_{x_0} \dots \sum_{x_5} \underbrace{\sum_{x_0} \phi_0(x_0, x_1) \dots \phi_5(x_0, x_5)}_{\tau(x_1, \dots, x_5)}
 \end{aligned}$$

elim x_0

elim x_1

The Variable Elimination Algorithm

Variable elimination is an algorithm to compute any marginal distribution in any MRF

In words: pick a variable x_i to eliminate, multiply together all factors containing x_i to get an intermediate factor, then sum out x_i

- ▶ Let $F = \{\phi_c : c \in C\}$ be the set of factors
- ▶ For each variable i in **some elimination order** (may not include all variables)
 - ▶ Let $A = \{\phi_c \in F : i \in c\}$ be the set of factors whose scope contains i
 - ▶ Let $\phi_a(\mathbf{x}_a) = \prod_{\phi_c \in A} \phi_c(\mathbf{x}_c)$ be the product of factors in A , with scope a equal to the union of the scopes of the individual factors
 - ▶ Let $\psi_i(\mathbf{x}_{a \setminus i}) = \sum_{x_i} \phi_a(\mathbf{x}_a)$ be the result of summing out x_i
 - ▶ Let $F = F \setminus A \cup \{\psi_i\}$

The final set of factors forms an MRF for the marginal distribution of the variables that were not eliminated.

Variable Elimination Discussion

- ▶ The efficiency of variable elimination depends on the maximum size of the intermediate factors created, which depends on the elimination ordering
 - ▶ Inference in MRFs is NP-hard, so we can't always find a good elimination ordering.
 - ▶ Finding the best elimination ordering for a given MRF is also NP-hard!
- ▶ It's always efficient to eliminate leaves if present (intermediate factors are no larger than original ones)
 - ▶ \implies for trees, we can find an efficient elimination ordering
 - ▶ In fact, because the elimination ordering is predictable in trees, we can realize extra efficiencies when answering multiple queries through a dynamic programming approach known as **message passing**