

## Lecture 9.

## Dynamic Programming: Sampling &amp; Expectation

Note Title

4/17/2011

$$\pi(\underline{x}) = \frac{1}{Z} e^{-E[\underline{x}]}, \quad E[\underline{x}] = \sum_{i=1}^N \phi_i(x_{i-1}, x_i)$$

$\begin{matrix} x_0 & x_1 & & & x_{N-1} & x_N \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{matrix}$

Previous lecture showed how to use DP to estimate  $\hat{\underline{x}} = \underset{\underline{x}}{\text{ARG MIN}} E[\underline{x}]$ , and compute  $E[\hat{\underline{x}}]$   
 Forward Pass computes  $E[\hat{\underline{x}}]$   
 Backward Pass gives  $\hat{\underline{x}}$ .

This lecture will show how to use DP to obtain samples from  $\pi(\underline{x})$ , to compute  $Z$ , and to  $\sum_{\underline{x}} h(\underline{x}) \pi(\underline{x})$  (sometimes by sampling, sometimes exactly).

Note: this is a different algorithm than the one used to compute  $\hat{\underline{x}} = \underset{\underline{x}}{\text{ARG MIN}} E[\underline{x}]$  — but it is closely related and is based on the same principles.

(1) Sampling from  $\pi(\underline{x})$ 

Claim — we can use DP to convert  $\pi(\underline{x}) = \frac{1}{Z} e^{-E(\underline{x})}$

into an alternative form  $\pi_0(x_0) \pi_1(x_1|x_0) \dots \pi_N(x_N|x_{N-1})$   
 or, equivalently,  $\pi_0(x_0|x_1) \pi_1(x_1|x_2) \dots \pi_{N-1}(x_{N-1}|x_N) \pi_N(x_N)$   
 (not the same  $\pi$ 's in these two cases).

If we can express  $\pi(\underline{x})$  in one of these two forms then we can obtain samples  $\underline{x}^1, \underline{x}^2, \dots, \underline{x}^M$  as follows

$x_0^1$	~ sample from $\pi_0(x_0)$	$x_0^2$	~ from $\pi_0(x_0)$
$x_1^1$	~ sample from $\pi_1(x_1 x_0^1)$	$x_1^2$	~ from $\pi_1(x_1 x_0^2)$
$x_2^1$	~ sample from $\pi_2(x_2 x_1^1)$	$x_2^2$	~ from $\pi_2(x_2 x_1^2)$
$\vdots$		$\vdots$	
$x_N^1$	~ sample from $\pi_N(x_N x_{N-1}^1)$	$x_N^2$	~ from $\pi_N(x_N x_{N-1}^2)$

Or  $x_N^1$  ~ from  $\pi_N(x_N)$ ,  $x_{N-1}^1$  ~ from  $\pi_{N-1}(x_{N-1}|x_N^1)$ , etc

So only need to sample from  $\pi(x_i|x_{i+1})$  or  $\pi(x_{i+1}|x_i)$   
 ~ use techniques from earlier lectures.

How to convert from  $\pi(x) = \frac{e^{-E(x)}}{Z}$

to these forms - eg.  $\pi_N(x_N) \pi_{N-1}(x_{N-1}|x_N) \dots \pi_1(x_0|x_1)$ ?

Special case: if  $N=1$   $\pi(x) = \pi(x_0, x_1)$

$$\begin{aligned} \text{then it is easy, } \pi(x_0, x_1) &= \pi(x_0|x_1) \pi(x_1) \\ &= \pi(x_1|x_0) \pi(x_0) \end{aligned}$$

Note: Markov Condition  $\rightarrow$  the graph structure of this model (due to potentials  $\phi_i(x_{i-1}, x_i)$ ) means that the model has a local Markov structure  $\rightarrow$  eg.  $\pi(x_i | x_{i-1}, x_{i-2}, \dots, x_0) = \pi(x_i | x_{i-1})$   
If we know  $x_{i-1}$ , then knowing  $x_{i-2}, \dots, x_0$  also gives no more information about  $x_i$ .

DP Algorithm: To compute  $\pi_0(x_0|x_1) \dots \pi_{N-1}(x_{N-1}|x_N) \pi_N(x_N)$  from  $\pi(x)$

- Define  $V_1(x) = \sum_{s \in S} e^{-\phi_1(s, x)}$

- Recursively compute for  $i=2, \dots, N$   
$$V_i(x_i) = \sum_{y \in S} V_{i-1}(y) e^{-\phi_i(y, x_i)}$$

Then we compute efficiently (in  $O(k^2 N)$ ):

(A) The normalization constant:  $Z = \sum_{x_N \in S} V_N(x_N)$

(B) The marginal  $\pi_N(x_N) = \frac{V_N(x_N)}{Z}$

(C) The conditionals  $\pi_i(x_i | x_{i+1}) = \frac{V_i(x_i) e^{-\phi_{i+1}(x_i, x_{i+1})}}{\sum_y V_i(y) e^{-\phi_{i+1}(y, x_{i+1})}}$

(Note: easy to modify algorithm to compute  $\pi_0(x_0) \pi_1(x_1|x_0) \dots \pi_N(x_N|x_{N-1})$ )

Page 3

To estimate  $I = \sum_x h(x) \pi(x)$ ,

where  $h(x)$  is a quantity you want to estimate, like previous lectures.

Obtain samples  $x^1, x^2, \dots, x^m$  from  $\pi(x)$

Estimator:  $I_m = \frac{1}{m} \sum_{i=1}^m h(x^i)$ . (as before)

Unbiased  $\langle I_m \rangle = E_{\pi} I_m = I$

efficiency  $\sim \frac{1}{m} \text{Var}_{\pi} h(x)$ .

But, if  $h(x)$  takes certain special forms, then we can use DP to compute  $\sum_x h(x) \pi(x)$  exactly.

Observe that:  $V_1(x_1) = \sum_{x_0} e^{-\phi_1(x_0, x_1)}$ ,  $V_2(x_2) = \sum_{x_0, x_1} e^{-\phi_1(x_0, x_1) - \phi_2(x_1, x_2)}$   
 $V_n(x_n) = \sum_{x_0, x_1, \dots, x_{n-1}} e^{-\phi_1(x_0, x_1) - \dots - \phi_n(x_{n-1}, x_n)}$   $\rightarrow$  computed by DP (on page 2)

Suppose  $h(x) = h_0(x_0) + \dots + h_n(x_n)$

Then to compute  $\sum_x h(x) \pi(x)$  requires:

(i) computing  $Z$ , done by DP on page 2

(ii) for each  $i$ , compute  $\sum_{x_0, x_1, \dots, x_n} h_i(x_i) e^{-\phi_1(x_0, x_1) - \dots - \phi_n(x_{n-1}, x_n)}$

$\rightarrow$  done by modifying DP on page 2.

e.g.  $\tilde{V}_i(x_i) = \sum_{y \text{ yes}} \tilde{V}_{i-1}(y) e^{-\phi_i(y, x_i)} h_i(y, x_i)$

other updates as before.

new

This requires more computation -  $O(k^2 n^2)$

$O(k^2 n)$  for each  $h_i(\cdot)$

We can extend this in the obvious way if

$h(x) = h_1(x_0, x_1) + \dots + h_n(x_{n-1}, x_n)$

e.g.  $\tilde{V}_i(x_i) = \sum_{y \text{ yes}} \tilde{V}_{i-1}(y) e^{-\phi_i(y, x_i)} h_i(y, x_i)$

for computing  $\sum_{x_0, x_1, \dots, x_n} h_i(x_{i-1}, x_i) e^{-\phi_1(x_0, x_1) - \dots - \phi_n(x_{n-1}, x_n)}$ .

Note: this can be extended to other cases - e.g. terms like  $h_i(x_{i-1}, x_{i-2}, x_i)$ , but gets more complicated.

Page 4

### Special Case:

Ising Spin Model

$$\pi(\underline{x}) = \frac{1}{Z} e^{\beta(x_0 x_1 + \dots + x_{N-1} x_N)}$$

Note  $(-1)^2 = 1 = 1^2$   
 so most probable states  
 are  $\underline{x} = (1, 1, \dots, 1)$  and  
 $\underline{x} = (-1, -1, \dots, -1)$

$x_i \in \{-1, +1\}$

$$V_1(x_1) = e^{\beta x_1} + e^{-\beta x_1} = e^{\beta} + e^{-\beta}$$

(note: this is a very special case, usually  $V_1(x_1)$  is a function of  $x_1$ )

$$V_2(x_2) = \sum_{y \in S} V_1(y) e^{\beta y x_2}$$

$$= (e^{\beta} + e^{-\beta}) \sum_{y \in S} e^{\beta y x_2} = (e^{\beta} + e^{-\beta})^2$$

In general,  $V_t(x_t) = (e^{\beta} + e^{-\beta})^t$  //  $V_d(x_d) = (e^{\beta} + e^{-\beta})^d$

$$\text{Hence } Z = \sum_{\underline{x} \in S} V_N(x_N) = 2 (e^{\beta} + e^{-\beta})^N$$

$S = \{-1, +1\}$

Marginal Density

$$\pi(x_N) = \frac{V_N(x_N)}{Z} = \frac{1}{2}$$

Conditional Distribution:

$$\pi_t(x_t | x_{t+1}) = \frac{(e^{\beta} + e^{-\beta})^t e^{\beta x_t x_{t+1}}}{\sum_{y \in S} (e^{\beta} + e^{-\beta})^t e^{\beta y x_{t+1}}}$$

$$\pi_t(x_t | x_{t+1}) = \frac{e^{\beta x_t x_{t+1}}}{e^{\beta} + e^{-\beta}} \quad //$$

page 5

The Ising Model was invented by Physicists to study phase transitions  $\rightarrow$  eg. how does ice become water at a critical temperature  $T = 0^\circ$  centigrade.  $\beta = 1/T$ .

For small  $\beta$  (high temperature  $T$ ), the Ising model is disordered, samples from Ising will be like

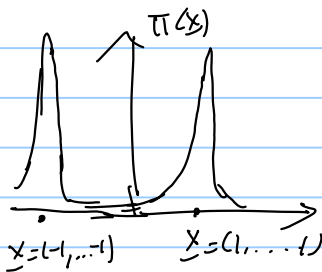
1 -1 -1 1 1 1 1 -1 1 -1 1 1

For large  $\beta$  (small temperature  $T$ ), the Ising model is ordered, samples will tend to be either

1 1 1 1 1 1 1 1 1 ie.  $x_0 = x_1 = \dots = x_n$   
or -1 -1 -1 -1 -1

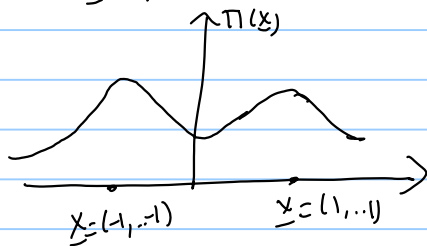
At the critical temperature  $T_c = 1/\beta_c$ , the model will change from order to disorder (e.g. from ice to water)

At low T.



Sharp peaks in  $\pi(x)$  at  $\underline{x} = (1, \dots, 1)$  and  $\underline{x} = (-1, \dots, -1)$

At high T.



peaks in  $\pi(x)$  at the same places  $\underline{x} = (1, \dots, 1)$  &  $\underline{x} = (-1, \dots, -1)$ .  
but  $\pi(x)$  is also quite large at other places.

Samples from  $\pi(x)$  at low temperature will be almost all at  $(1, 1, \dots, 1)$  or  $(-1, \dots, -1)$  the ordered states

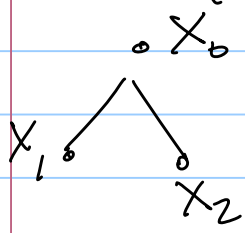
Samples at high temperature will also occur at the disordered states - eg  $(1, -1, 1, 1, -1, -1, \dots)$ .

They will still be most probable at  $(1, 1, \dots, 1)$  or  $(-1, \dots, -1)$  but there are many more disordered states (exponentially more). So we expect that a sample will probably be a disordered state.

page 6 DP can be extended to apply to any graph without closed loops.

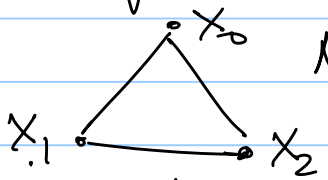
(Both the  $m$  computations in previous lecture and the  $V$  computations in this lecture.)

E.G.



yes for DP.

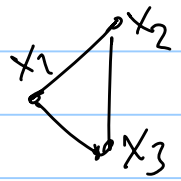
But not if there are closed loops.



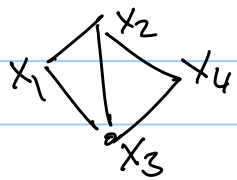
No for DP

But there are ways — the junction tree algorithm — which allows us to convert a prob. model on a graph with closed loops into a new model without closed loops — by augmenting the variables

E.G.



Define new variable  $z_1 = (x_1, x_2, x_3)$



etc.  $z_1, z_2$   $z_1 = (x_1, x_2, x_3)$   $z_2 = x_4$

But augmenting variables may make DP impractical. It may require adding extra nodes (i.e. making  $N$  very large) or make the number of state values  $k$  too large.

E.G. If  $x_i$  has  $k$  possible values  $(s_1, s_2, \dots, s_n)$  then  $z_1 = (x_1, x_2, x_3)$  has  $k^3$  possible values. Beyond scope of course.

Also, the  $m$ -update and  $V$ -update algorithms can be applied to graphs with closed loops as approximation. This gives belief propagation (sum-product & sum-max).