

parameter learning - the case of HMMs

parameters  $\theta =$

So far: given a Bayesian Network (including the conditional distributions), how can we do inference

ex: compute  $P(\text{state } S \text{ has value } j, \text{ state } S' \text{ has value } l)$

Today: learning the parameters from data

Focus: EM on HMMs

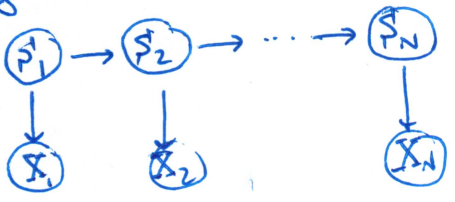
- from last time: HMMs are trees, so inference is easy and "half" of their nodes are leaves  $\Rightarrow$  inference is "easy" (compared to arbitrary BNs)

prep for A3.

Setting and notation

- also: HMMs are nice models for time series (sequential data)

Setting and notation



} unobserved. Each  $S_t$  takes values from  $1 \dots K$  state @ time  $t$

} observed

our usual notation for "hidden compressed dimensionality"

example: (instead of modeling a student as in the Koller/Friedman BN) no going to model a lecturer

lecture state values: "announcement" = 1, "important content" = 2, "optional content" = 3, "joke" = 4. (let's not say "unimportant")

observed: loudness relative, modeled as a Gaussian with (hidden) mean; variance

- we do relative loudness so that  $\theta$ 's are OK.

parameters

parameters  $\theta =$

Trans( $j \leftarrow i$ ) = prob, given in state  $i$ , that next state is  $j$

Out( $x \uparrow i$ ) = prob, given in state  $i$ , that volume is  $x \in \mathbb{R}$

controlled by  $\mu_i; \sigma_i^2$ , mean; variance for state  $i$  (very loud  $\uparrow$  joke)

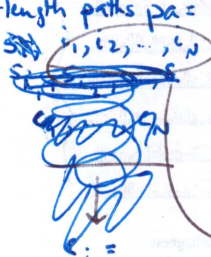
Start( $i$ ) = prob the first state is  $i$ .

Task: (for all of these, the same) you weren't listening to the words, just her volume. Can you infer her model?

to be it would be convenient to do so.  
 Once you've got a model, you can infer things like ~~you can~~ for future lectures like:  
~~given~~ I've been kind of zoned out in lecture until know, ~~what's the~~ but  
 attuned to the volume;  
 what's the prob that she's about to say sth important : I should  
 tune in?  
 or, what's the prob that she just made a joke so I should  
 humor her : laugh?

1st attempt: find  $\theta^*$  maximizing log-likelihood of observed  $x_1, \dots, x_N = \vec{x}$

$$\theta^* = \arg \max_{\substack{\text{Trans}(\dots) \\ \theta_1, \dots, \theta_k \\ \sigma_1^2, \dots, \sigma_k^2 \\ \text{start}(1), \dots, \text{start}(k)}}} \log \left( \sum_{\substack{\text{N-length paths } pa = \\ (i_1, i_2, \dots, i_N)}} \text{start}(i_1) \text{Out}(x_1, i_1) \prod_{t=2}^N \text{Trans}(i_{t-1} \leftarrow i_t) \text{Out}(x_t, i_t) \right)$$



- Lagrange constraints  
 i.e.,  $\sum_j \text{Trans}(j \leftarrow i) = 1$   
 $\forall i$

can the students work this out themselves? clicker-ize?   
 (maybe in words) instead of the whole formula!

Most students had significant trouble: they could not get started.  
 Karthik afterwards suggested: tell them:  $\prod P(n | \text{parents}(n))$   
 <that seems like a good idea>

So, what would we do?

- Take the derivative w.r.t., say,  $\text{Trans}(4 \leftarrow 0)$ , set to 0, and solve,   
 (partial)

<he also noted that maybe the notation shift was hard>

hope  
 1  
 (giant sum w/ all of  $\theta$  in it)

$$\frac{\partial}{\partial \text{Trans}(4 \leftarrow 0)} \sum_{\substack{\text{N-length paths} \\ pa}} \dots$$

(monomial in  $\text{Trans}(4 \leftarrow 0)$ )

$$\frac{\partial}{\partial \text{Trans}(4 \leftarrow 0)} (\text{const}) \text{Trans}(4 \leftarrow 0)^m$$

$$= (\text{const}) m \text{Trans}(4 \leftarrow 0)^{m-1}$$

so that's easy.

$$= \frac{(\text{const}) m \text{Trans}^m}{\text{Trans}}$$

that's bad.

this is impossible: we want to solve for Trans independent of all the other variables.

•  $-\lambda_i$

But thinking back over techniques that have helped us before:

EM: when guessing the hidden structure would make life easier.

~~take~~  $\log \cancel{P(\vec{x}_1, \dots, \vec{x}_N, i_{pa} | \theta)}$

add in the hidden structure  
normalize

b/c that would  
sure make things  
easier.

$\sum_{pa} P(\vec{x}_1, \dots, \vec{x}_N, i_{pa} | \theta)$

since I don't know the hidden structure  
I'll take an expectation over all the possibilities  
But how am I supposed to know what this distribution  
is?

$\theta^{old}$  - ~~guess~~ use a previous guess!

Is taking the <sup>partial</sup> derivative of this going to be easier?

this portion is constant in  $\text{Trans}(2 \leftarrow 1)$ .  
(it has terms like  $\text{Trans}^{old}(2 \leftarrow 1)$ , but so  
what :))

So!  $\frac{\partial}{\partial \text{Trans}(4 \leftarrow 3)} \left[ \sum_{pa} P(pa | \vec{x}, \theta^{old}) \log P(pa, \vec{x} | \theta) - \lambda_1 \left( \sum_j \text{Trans}(j \leftarrow 3) - 1 \right) \right]$   
- (stuff irrelevant to  $\text{Trans}(4 \leftarrow 3)$ )

$= \left[ \sum_{pa} P(pa | \vec{x}, \theta^{old}) \cdot \frac{\partial}{\partial \text{Trans}(4 \leftarrow 3)} P(pa, \vec{x} | \theta) \right] - \lambda_1$   
this we saw

= the great thing about monomials, as seen above.

$= \left[ \sum_{pa} P(pa | \vec{x}, \theta^{old}) \cdot \frac{1}{P(pa, \vec{x} | \theta)} \cdot P(pa, \vec{x} | \theta) \cdot \# \left( \begin{matrix} 4 \\ \leftarrow \\ 3 \end{matrix} \text{ in } pa \right) \cdot \frac{1}{\text{Trans}(4 \leftarrow 3)} \right] - \lambda_1$

set to 0, solve:

$\text{Trans}(4 \leftarrow 3) = \frac{1}{\lambda_1} \sum_{pa} P(pa, \vec{x} | \theta^{old}) \cdot \# \left( \begin{matrix} 4 \\ \leftarrow \\ 3 \end{matrix} \text{ in } pa \right)$

normalizing constant, which is expected # of  
times you'd see  $j \leftarrow i$  ...  
which makes sense!

this is, though not impossible, still bad.

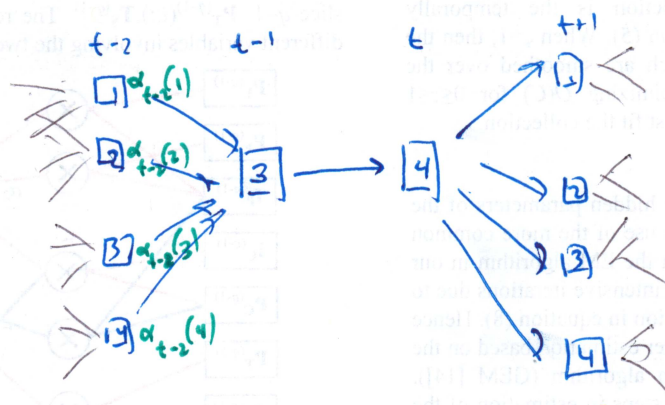
= expectation over all paths of  $S_{t+1} \rightarrow 3, S_t = 4$ . Now, b/c these are HMMs:

which we can ~~reconcept~~ reconceptualize as this:

$P(\bar{x}, \theta^{old}) \sum_{t=2}^N P(x_1, \dots, x_{t-1}, S_{t-1} = 3 | \theta^{old}) \text{Trans}^{old}(4 \leftarrow 3) \text{Out}^{old}(x_t | 4)$

called this "Sridharan's  $\alpha$ "

called this "Sridharan's  $\beta$ "



Now this certainly looks better, b/c we syntactically got rid of the sum over an exponential # of things.

But not so fast, you say - it's just hidden in the notation, right?

Ah, but here's the saving grace...

We've showed how to compute something like this in previous lectures!

Or, you can look @ the trellis graph.

notice that if we call  $P(x_1, \dots, x_{t-1}, S_{t-1} = 3 | \theta^{old}) = \alpha_{t-1}(3)$

then  $\alpha_{t+1}(3) = \sum_{i=1}^k P(x_1, \dots, x_{t-2}, S_{t-2} = i) \text{Trans}^{old}(3 \leftarrow i) \cdot \text{Out}^{old}(x_{t+1} | 3)$

alternately present as more wishful thinking

if we know  $P(x_1, \dots, x_{t-2}, S_{t-2} = i)$  for each  $i \dots$

$P(x_1, \dots, x_{t-2}, S_{t-2} = i)$

$\alpha_{t-2}(i)$

dynamic-program your way: to compute left-to-right to compute an  $\alpha$  per node.

polytime computation (poly # of nodes, linear work/node).

- and that's why we care about those  $\alpha$ 's.

the last lecture

what about this? You also saw how to do this in Prof Sridharan did

marginal

However, does <sup>the interior</sup> it look familiar...?  
that's like our inference stuff!

So, if only we could get rid of the ~~exponential~~ need to sum over an exponential # of things...

Flashback: we've seen how to, in HMMs, to magically get an ~~exponential~~ sum over exponential # of things to disappear.

- Prof Sridharan's lecture:

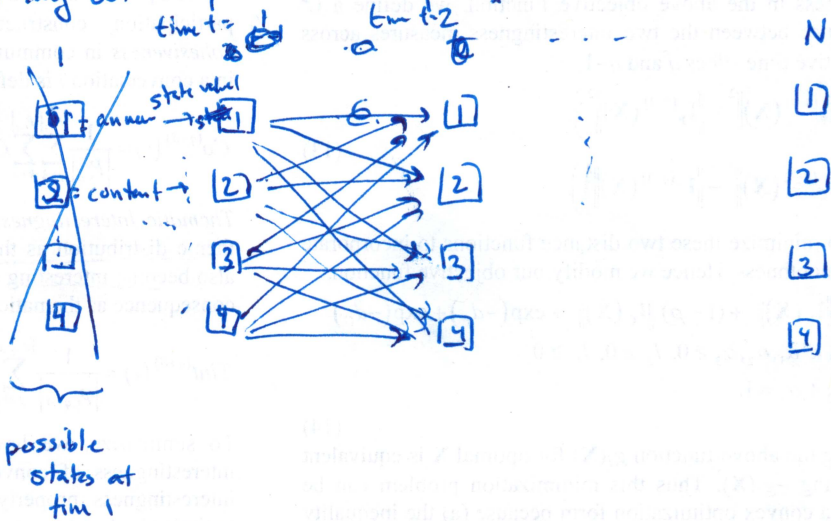
$P(s_t = i | \vec{x}, \theta^{\text{train}}) =$ , mathematically,  $\sum_{pa} P(pa | \vec{x}, \theta^{\text{train}})$  sum over all paths for  $\vec{x}$  of prob that ~~that~~  $pa$  has  $i$  as  $t$ th state. looks like a problem.

But we can reorganize as a dynamic program; given coord index, as follows: stepping through time indices (those  $\alpha$  &  $\beta$  quantities)

So, here:

$$\sum_{pa} P(pa | \vec{x}, \theta^{\text{old}}) \# (\overset{4}{\circ} \leftarrow \overset{3}{\circ} \text{ in } pa)$$

let's lay out all paths in compact trellis diagram:



again, reorganize by time steps:

$$\sum_{t=2}^N \sum_{pa} P(pa | \vec{x}, \theta^{\text{old}}) \cdot \mathbb{1}(\text{state at time } t-1 \text{ is } 3, \text{ state at time } t \text{ is } 4)$$

↑ remember we like marginals.

$$= \frac{1}{P(\vec{x}, \theta^{\text{old}})} \sum_{t=2}^N \sum_{pa} P(pa, \vec{x} | \theta^{\text{old}}) \mathbb{1}(s_{t-1}^{\text{pa}} = 3, s_t^{\text{pa}} = 4)$$

} so, we are summing up the probs of all paths that have  $3 \rightarrow 4$  @ time  $t-1 \rightarrow t$

You can similarly solve for  $\mu_3$  and  $\sigma_3$ .

↳ in the log-likelihood w/ hidden structure, you get terms ~~that are~~ linear in  $\mu_3$ .