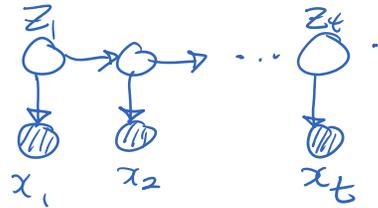


today: HMM & α - β recursion
ML for HMM (EM)

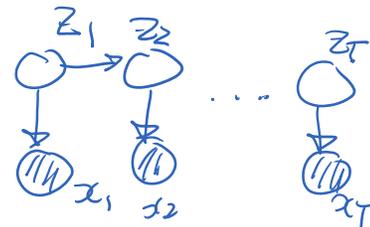
HMM (hidden Markov model)



$z_t \in \{1, \dots, k\}$ discrete
 x_t — cts. e.g. speech signal
 — discrete e.g. DNA sequence

(Often $z_t \sim$ Gaussian in the course \rightarrow Kalman filter)

HMM \rightarrow generalization of mixture model
 GMM \rightarrow add dependence on z_t



DGM:

$$p(x_{1:T}, z_{1:T}) = p(z_1) \prod_{t=1}^T \underbrace{p(x_t | z_t)}_{\text{emission probs.}} \prod_{t=2}^T \underbrace{p(z_t | z_{t-1})}_{\text{transition probs.}}$$

often, the emission probs & trans probs are homogeneous i.e. do not depend on t

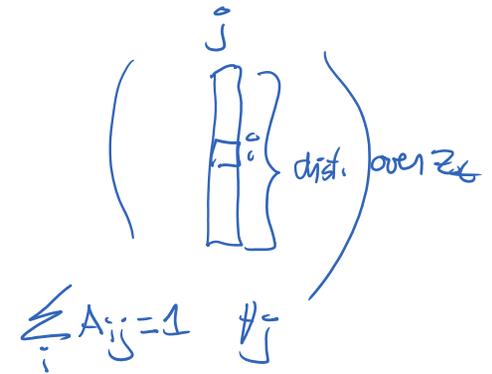
$$p_t(x_t | z_t) = f(x_t | z_t)$$

$$p_t(x_t | z_t) = f(x_t | z_t)$$

$$p_t(z_t = i | z_{t-1} = j) = A_{ij}$$

"stochastic matrix"

A



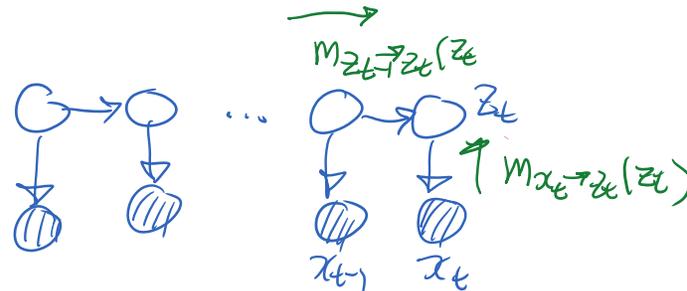
inference tasks:

prediction	$p(z_t x_{1:t-1})$	"where next?"
filtering	$p(z_t x_{1:t})$	"where now?"
smoothing	$p(z_t x_{1:T})$	"where in the past?"

$T > t$

α -recursion:

Let's run sum product here to derive recursions to compute prds.



to compute $p(z_t, \bar{x}_{1:t}) \propto p(z_t | \bar{x}_{1:t})$

$$p(z_t, \bar{x}_{1:t}) = \frac{1}{Z} \cdot m_{z_{t-1} \to z_t}(z_t) m_{x_t \to z_t}(z_t)$$

here $Z = 1$

$\alpha_t(z_t)$

$$M_{z_t \rightarrow z_t}(z_t) = \sum_{z_t} p(z_t | z_t) \delta(z_t, \bar{z}_t) = p(\bar{z}_t | z_t)$$

$$M_{z_{t-1} \rightarrow z_t}(z_t) = \sum_{z_{t-1}} p(z_t | z_{t-1}) \underbrace{M_{z_{t-2} \rightarrow z_{t-1}}(z_{t-1}) M_{z_{t-1} \rightarrow z_{t-1}}(z_{t-1})}_{p(z_{t-1}, \bar{x}_{1:t-1})}$$

define: $\alpha_t(z_t) \triangleq p(z_t | \bar{x}_{1:t})$

$\underbrace{p(z_{t-1}, \bar{x}_{1:t-1})}_{\alpha_{t-1}(z_{t-1})}$

$$\alpha_t(z_t) = \underbrace{p(\bar{x}_t | z_t)}_{\text{vector } (z_t)} \sum_{z_{t-1}} \underbrace{p(z_t | z_{t-1})}_{\text{matrix}} \underbrace{\alpha_{t-1}(z_{t-1})}_{\text{vector}}$$

α -recursion a.k.a "forward recursion" Like the "collect phase" in sum-product alg with z_t as the root

let $O_t(z_t) \triangleq p(\bar{x}_t | z_t)$

$\alpha_t = O_t \odot A \alpha_{t-1}$

↑
Klebsch product

initialization: $\alpha_1(z_1) = p(z_1, \bar{x}_1) = p(z_1) p(\bar{x}_1 | z_1)$

$\tilde{\alpha}_t \triangleq p(z_t | \bar{x}_{1:t})$ "filtering distribution"

time complexity:

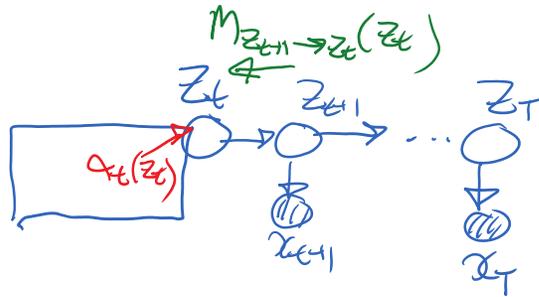
$O(t \cdot k^2)$

space complexity: $O(k)$ extra storage ($O(k^2)$ for string A)

Space complexity: $O(k)$ extra storage ($O(k^2)$ for string A)

$$\sum_{z_t} \underbrace{p(z_t, \bar{x}_{1:t})}_{\alpha_t(z_t)} = p(\bar{x}_{1:t}) \quad \text{"evidence probability"}$$

β recursion: smoothing



$$p(z_t, \bar{x}_{1:T}) = \frac{1}{\sum_{(=1)} \alpha_t(z_t)} \cdot \underbrace{M_{z_t \rightarrow z_t}(z_t)}_{\beta_t(z_t)}$$

turns out that $\beta_t(z_t) \triangleq p(\bar{x}_{t+1:T} | z_t)$

why?

$$p(z_t, \bar{x}) = p(\bar{x} | z_t) p(z_t) \stackrel{C.F.}{=} \underbrace{p(x_{t+1:T} | z_t)}_{\beta_t(z_t)} \underbrace{p(x_{1:t} | z_t)}_{\alpha_t(z_t)} p(z_t)$$

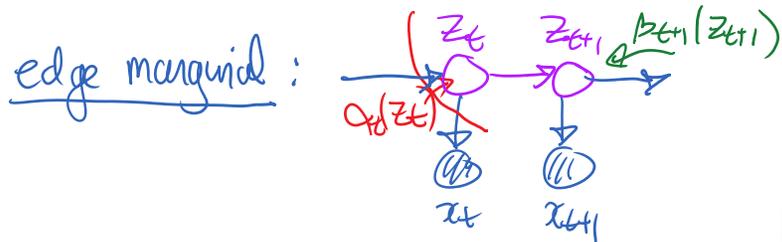
$$M_{z_{t+1} \rightarrow z_t}(z_t) = \sum_{z_{t+1}} p(z_{t+1} | z_t) p(\bar{x}_{t+1} | z_{t+1}) M_{z_{t+2} \rightarrow z_{t+1}}(z_{t+1})$$

$$\beta_t(z_t) = \sum_{z_{t+1}} p(z_{t+1} | z_t) p(\bar{x}_{t+1} | z_{t+1}) \beta_{t+1}(z_{t+1})$$

β -recursion (aka backward recursion)

initialization: $\beta_T(z_T) = 1 \quad \forall z_T$

initialization: $\beta_T(z_T) = 1 \quad \forall z_T$



$$p(z_t, z_{t+1}, \bar{x}_{1:T}) = \alpha_t(z_t) \beta_{t+1}(z_{t+1}) \cdot p(z_{t+1}|z_t) \cdot p(\bar{x}_{t+1}|z_{t+1})$$

15h28

Numerical stability trick:

issue: α_t & β_t can easily go to $1e-100$

two possibilities a) (general) use $\log(\alpha_t)$ instead

$$\log\left(\sum_i a_i\right) = \log\left(\tilde{a} \left(\sum_i \frac{a_i}{\tilde{a}}\right)\right)$$

call $\tilde{a} \triangleq \max_i a_i$ $= \log(\tilde{a}) + \log\left(1 + \sum_{j \neq i_{\max}} \exp(\log(a_j) - \log(\tilde{a}))\right)$

$i_{\max} \triangleq \text{argmax}_i a_i$

b) normalize the messages:

• α -recursion, use $\tilde{\alpha}_t(z_t) = p(z_t | \bar{x}_{1:t})$

before, $\alpha_t = \alpha_t \odot A \alpha_{t-1}$

$$\tilde{\alpha}_t = \alpha_t \odot A \tilde{\alpha}_{t-1}$$

$$\sum_{z_t} (\quad)$$

you can show that

$$\sum_{z_t} (\alpha_t \odot A \tilde{\alpha}_{t-1})(z_t) = p(\bar{x}_t | \bar{x}_{1:t-1})$$

$$p(\bar{x}_{1:T}) = \prod_{t=1}^T p(\bar{x}_t | \bar{x}_{1:t-1}) = \prod_{t=1}^T C_t$$

$$\sum_{z_t} (a_t \circ \tilde{a}_{t-1})(z_t) = p(\bar{x}_t | x_{1:t-1}) \triangleq C_t$$

• β -recursion:

define $\tilde{\beta}_t(z_t) \triangleq \frac{\beta_t(z_t)}{p(\bar{x}_{t+1:T} | x_{1:t})}$ $\prod_{u=t+1}^T C_u$

note $\sum_{z_t} \tilde{\beta}_t(z_t) \stackrel{\text{rec.}}{=} 1$

khco

exercise: derive $\tilde{\beta}$ -recursion

ML for HMM:

- suppose $p(x_t | z_t = k) = f(x_t | \eta_k)$ $\eta = (\eta_k)_{k=1}^K$ some parametric model e.g. Gaussian on x_t
- $p(z_{t+1} = i | z_t = j) = A_{ij}$
- $p(z_1 = i) = \pi_i$ $\Theta = (\eta, A, \pi)$

want to estimate $\hat{\eta}, \hat{A} \text{ \& } \hat{\pi}$ by ML from data $(x^{(i)})_{i=1}^N$
 $x^{(i)} = x_{1:T}^{(i)}$

\leadsto use EM

s^{th} iteration

E step: $Q_{s+1}(z) = p(z | x, \Theta^{[s]})$

M step: $\hat{\Theta}^{[s+1]} = \underset{\Theta \in \Theta}{\operatorname{argmax}} E_{q_{s+1}} [\log p(x, z)]$

complete log-likelihood:

$$\log p(x, z | \Theta) = \sum_{i=1}^N \left[\underbrace{\log p(z^{(i)})}_{\sum_k z_{t,k}^{(i)} \log \pi_k} + \sum_{t=1}^T \log p(\bar{x}_t^{(i)} | z_t^{(i)}) + \sum_{t=2}^T \log p(z_t^{(i)} | z_{t-1}^{(i)}) \right]$$

huge variables

$E_{q_{s+1}} [\log p(x, z | \Theta)] = \dots$

$E_{q_{s+1}} [z_{t,k}^{(i)}] = q_{s+1}(z_{t,k}^{(i)}=1) \triangleq \tau_{t,k}^{(i)}$

smoothing distribution
 $p(z_t | \bar{x}_{1:T})$

$q_{s+1}(z_{t,l}^{(i)}=1, z_{t,m}^{(i)}=1)$

$p(z_{t,l}^{(i)}=1, z_{t,m}^{(i)}=1 | \bar{x}_{1:T}^{(i)}, \Theta^{[s]})$

smoothing edge marginal

maximize with respect to Θ

$$\hat{\pi}_k^{[s+1]} = \sum_{i=1}^N \hat{\pi}_{i,k}^{(i)}$$

$\underbrace{\sum_{i=1}^N \sum_{k=1}^K \tau_{i,k}^{(i)}}_1 \} N$

$$\hat{A}_{l,m}^{[s+1]} = \frac{\sum_{i=1}^N \sum_{t=2}^T \tau_{t,l,m}^{(i)}}{\sum_u \sum_{i=1}^N \sum_{t=2}^T \tau_{t,u,m}^{(i)}}$$

$\hat{\pi}_k \rightarrow$ soft count ML
e.g. Gaussians
similar to GMM
"weighted empirical mean"

"Baum-Welch" alg. — forward-backward alg. (α-β recursion / sum product) + EM for GMM

Viterbi to compute $\operatorname{argmax}_{z_1, \dots, z_T} p(z_{1:T} | \bar{x}_{1:T})$
(max product)

https://en.wikipedia.org/wiki/Baum%E2%80%93Welch_algorithm

(62-1