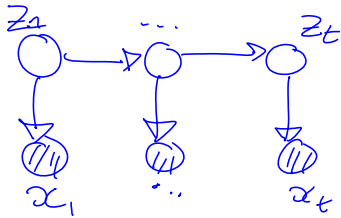


Lecture 14 - scribbles

Tuesday, October 18, 2016  
14:50

today: HMM  
 • max entropy, exp family  
 and beautiful duality

HMM: hidden Markov model



$z_t \in \{1, \dots, K\}$  discrete

$x_t$   $\begin{cases} \text{cts} \rightarrow \text{eg speech signal} \\ \text{discrete} \rightarrow \text{DNA sequence} \end{cases}$

speech recognition  $x_t \rightarrow$  speech signal  
 $z_t \rightarrow$  phonemes

HMM  $\rightarrow$  generalization of mixture model



DGM:

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

often: emission probs and transition probs are homogeneous i.e. do not depend on t

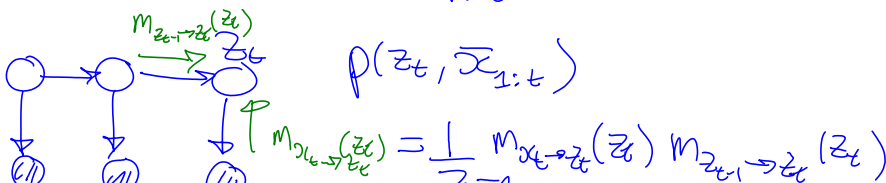
$$p_t(x_t | z_t) = f(x_t | z_t)$$

$$p_t(z_t = i | z_{t-1} = j) = A_{ij} \quad A \left( \begin{matrix} j \\ \text{dist. over } z \end{matrix} \right)$$

$$\sum_i A_{ij} = 1$$

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?"



$$M_{z_t \rightarrow z_t}(z_t) = \prod_{z=1} M_{z_t \rightarrow z_t}(z_t)$$

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

$$M_{z_{t-1} \rightarrow z_t}(z_t) = \sum_{z_{t-1}} p(z_t | z_{t-1}) \underbrace{M_{z_{t-1} \rightarrow z_{t-1}}(z_{t-1}) M_{z_{t-2} \rightarrow z_{t-1}}(z_{t-1})}_{\text{matrix}}$$

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

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

$\alpha$ -recursion, forward recursion like the "collect phase" in sum-product with  $z_t$  the root

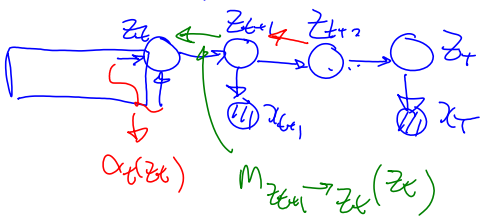
$$\alpha_1(z_1) = p(z_1, \bar{x}_1) = p(z_1) p(\bar{x}_1 | z_1)$$

time complexity:  $O(t \cdot k^2)$

space complexity  $O(k)$  extra ( $O(k^2)$  count storing A)

$$\sum_{z_t} \alpha_t(z_t) = p(\bar{x}_{1:t}) \triangleq \beta_t(z_t)$$

smoothing:  $p(z_t, \bar{x}_{1:T}) = \prod_{z=1} \alpha_t(z_t) M_{z_{t+1} \rightarrow z_t}(z_t)$



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

$$\beta_t(z_t) = \sum \beta_{t+1}(z_{t+1})$$

initialization:  $\rightarrow \beta$ -recursion "backward recursion"

$$\beta_T(z_T) = M_{z_{T+1} \rightarrow z_T}(z_T) = 1$$

finally, can get marginals on two nodes



$$p(z_t, z_{t+1}, \bar{x}_{1:T})$$

$$\alpha_t(z_t) \beta_{t+1}(z_{t+1}) p(\tilde{x}_t | z_t) p(\tilde{x}_{t+1} | z_{t+1}) \cdot p(z_{t+1} | z_t)$$

ML in HMM:

- suppose  $p(x|z=k) = f(x|\theta_k)$
- $p(z_{t+1}=i | z_t=j) = A_{ij}$

want to estimate  $\hat{A}, \hat{\Theta}$  by ML from data

$\{x^{(i)}\}_{i=1}^N$  where  $x^{(i)} = x_{1:T}^{(i)}$

EM write complete log-likelihood  $\log p(x, z)$

$q_{s+1}(z) \rightarrow p(z|x, \theta^{[s+1]})$   
E-step  
parameter of s-th iteration

M-step:  $\hat{\theta}^{[s+1]} = \arg \max_{\theta \in \Theta} \mathbb{E}_{q_{s+1}}[\log p(x, z)]$

$$\log p(x, z) = \sum_{i=1}^N \left[ \underbrace{\log p(z_1^{(i)})}_{\substack{\text{huge} \\ \text{variables}}} + \sum_{t=1}^T \log p(\tilde{x}_t^{(i)} | z_t^{(i)}) + \sum_{t=2}^T \log p(z_t^{(i)} | z_{t-1}^{(i)}) \right]$$

$$\sum_k z_{1,k}^{(i)} \log \pi_k \quad \sum_k z_{t,k}^{(i)} \log f(\tilde{x}_t^{(i)} | \theta_k) \quad \sum_{l,m} z_{t,l}^{(i)} z_{t-1,m}^{(i)} \log A_{l,m}$$

$\mathbb{E}_{q_{s+1}}[\log p(x, z)] = \dots$

$q_{s+1}(z_{t,l}^{(i)}=1, z_{t-1,m}^{(i)}=1)$   
 smoothing marginals  $p(z_t^{(i)}=l, z_{t-1}^{(i)}=m | x_{1:T}^{(i)}, \theta^{[s+1]})$

maximizing for  $\Theta$ :

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

$\tau_{t,k}^{(i)}$  soft-counts

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

$$e^{-\log(\cdot)}$$

KL divergence and density estimation

statistical decision theory

recall loss:

$$L(P, a)$$

if action is estimate a distribution say  $\hat{P}$

standard (ML) loss is log-loss  $L(P, \hat{P}_a) = \mathbb{E}_{X \sim P} [-\log \hat{P}_a(X)]$

cross-entropy

suppose  $P = P_{\theta}$

then this loss becomes:  $\sum_{x \in \Omega_X} -p(x) \log p(x) \stackrel{\Delta}{=} H(p)$  entropy of  $P$

excess loss for action  $a = P_{\theta}$ ,

$$\begin{aligned} L(P, P_{\theta}) - \min_{a \in \mathcal{A}} L(P, a) &= - \sum_{x \in \Omega_X} p(x) \log \frac{P_{\theta}(x)}{p(x)} \\ &\stackrel{\Delta}{=} L(A, P) \\ &= KL(P \parallel P_{\theta}) \end{aligned}$$