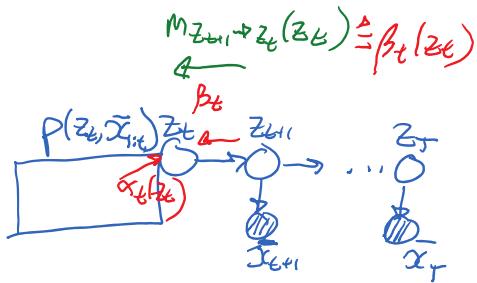


today:

- EM for HMM
- info. theory & KL

$\beta$ -recursion (smoothing)



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

$$\begin{aligned} M_{z_{t+1} \rightarrow z_t}(z_t) &= \sum_{z_{t+1}} p(z_{t+1}|z_t) p(\bar{x}_{t+1}|z_{t+1}) M_{z_{t+2} \rightarrow z_{t+1}}(z_{t+1}) \\ \beta_t(z_t) &= \sum_{z_{t+1}} p(z_{t+1}|z_t) p(\bar{x}_{t+1}|z_{t+1}) \beta_{t+1}(z_{t+1}) \end{aligned}$$

$\beta$ -recursion (aka. backward recursion)

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

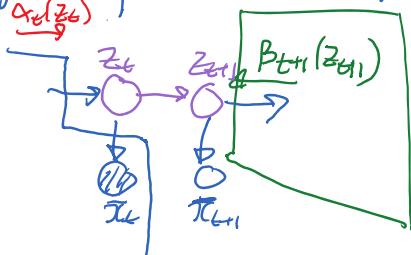
why?  $p(z_t, \bar{x}_{1:T}) = p(\bar{x}|z_t) p(z_t)$

$$\stackrel{\text{c.i.}}{=} p(\bar{x}_{t+1:T}|z_t) p(\bar{x}_{1:t}|z_t) p(z_t)$$

$$\Rightarrow \beta_t(z_t)$$

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

edge marginal



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

numerical stability trick:

issue:  $\alpha_t$  &  $\beta_t$  can easily go to  $1e-100$

two possibilities:

a) (general) store  $\log(\alpha_t)$  instead

$$\log\left(\frac{\alpha_t}{\tilde{\alpha}_t}\right) = \log\left(\frac{\tilde{\alpha}_t}{\alpha_t}\right) \quad (\alpha_t > 0)$$

$$\log \left( \sum_i a_i \right) = \log \left( \tilde{a} \left( \sum_i a_i \right) \right) \quad (a_i > 0)$$

call  $\tilde{a} \triangleq \max_i a_i$

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

$\log(\tilde{a}) \triangleq \max_i \log(a_i)$

$i_{\max} \triangleq \arg \max_i a_i$

b) normalize the message

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

before  $\alpha_t = \alpha_0 \odot A \alpha_{t-1}$

$$\tilde{\alpha}_t = \frac{\alpha_0 \odot A \tilde{\alpha}_{t-1}}{\sum_{z_t} (\quad)} \quad \} \triangleq c_t$$

You can show that  $c_t = \sum_{z_t} (\alpha_0 \odot A \tilde{\alpha}_{t-1})(z_t) = p(\bar{x}_t | \bar{x}_{1:t})$

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

- $\beta$ -recursion:

$$\text{define } \tilde{\beta}_t(z_t) = \frac{p(\bar{x}_{t+1:T} | z_t)}{p(\bar{x}_{t+1:T} | \bar{x}_{1:t})} \quad \} \frac{1}{\prod_{u=t+1}^T c_u} \quad \text{note: } \sum_{z_t} \tilde{\beta}_t(z_t) \neq 1$$

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

15h50

### ML for HMM

some parametric model for dist. on  $z_t$

- suppose  $p(x_t | z_t=k) = f(x_t | \eta_k)$  e.g. Gaussian on  $x_t$
  - $p(z_{t+1}=i | z_t=j) = A_{ij}$
  - $p(z_1=i) = \pi_i$
- want to estimate  $\hat{\Pi}, \hat{A}, \hat{\pi}$  by ML from data  $\mathcal{X} = (x^{(i)})_{i=1}^N$
- $$x^{(i)} = x_{1:T_i}^{(i)}$$

$$\Theta = (\eta, A, \pi)$$

$$x^{(i)} = x_{1:T_i}^{(i)}$$

→ use EM at  $s$ th iteration

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

$$\text{M step: } \hat{\Theta}^{[s+1]} = \underset{\Theta \in \Theta}{\operatorname{arg\,max}} \mathbb{E}_{q_{st+1}} [\log p(x, z | \Theta)]$$

Complete log-likelihood

$$\log D(x, z | \Theta) = \sum_{i=1}^N \left( \log D(z_i^{(i)}) + \sum_{t=1}^{T_i} \log p(\bar{x}_{t+1}^{(i)} | z_t^{(i)}) + \sum_{t=1}^{T_i} \log p(z_t^{(i)} | z_{t-1}^{(i)}) \right)$$

Complete Log-Likelihood

$$\log p(z|z_t|G) = \sum_{i=1}^N \left( \underbrace{\log p(z_i^{(i)})}_{\text{huge variable}} + \sum_{t=1}^{T_i} \log p(\bar{z}_{t,i}^{(i)} | z_t^{(i)}) + \sum_{t=2}^{T_i} \log p(z_t^{(i)} | z_{t-1}^{(i)}) \right)$$

$$\mathbb{E}_{q_{S+1}} [\log p(z|z_t|G)] = \sum_K \bar{z}_{t,K}^{(i)} \log \pi_K$$

$$\dots$$

$$\mathbb{E}_{q_{S+1}} [z_{t,K}^{(i)}] = q_{S+1}(z_{t,K}^{(i)} = 1) \triangleq \gamma_{t,K}^{(i)}$$

smoothing dist  $p(z_{t,K=1}^{(i)} | \bar{x}_{1:T_i}^{(i)}; G^{(i)})$

$$q_{S+1}(z_{t,l,m}^{(i)} = 1, z_{t-1,m}^{(i)} = 1) = p(z_{t,l}^{(i)} = 1, z_{t-1,m}^{(i)} = 1 | \bar{x}_{1:T_i}^{(i)}; G^{(i)})$$

$$\Delta \text{marginal } m$$

$$\gamma_{t,l,m}^{(i)} \quad \text{note: you should have: } \sum_m \gamma_{t,l,m}^{(i)} = \gamma_{t,m}^{(i)}$$

smoothing edge  
marginal on HM  
( $\alpha$ - $\beta$  recursion)

maximize with respect to  $G$ :

$$\hat{\pi}_K^{[S+1]} = \sum_{i=1}^N \gamma_{1,K}^{(i)}$$

$$\hat{\pi}_{l,m}^{[S+1]} = \sum_{i=1}^N \sum_{t=2}^{T_i} \gamma_{t,l,m}^{(i)}$$

$\hat{\pi}_K^{[S+1]} \rightarrow \text{soft counts}$

$$\sum_u \left( \sum_{i=1}^N \sum_{t=2}^{T_i} \gamma_{t,u,m}^{(i)} \right)$$

e.g. (Gaussians similar to GMM "weighted empirical mean" with weights  $\gamma_{t,K}^{(i)}$ )

$$\hat{\mu}_K = \frac{\sum_{i=1}^N \sum_{t=1}^{T_i} \gamma_{t,K}^{(i)} \bar{z}_{t,i}^{(i)}}{\sum_{i=1}^N \sum_{t=1}^{T_i} \gamma_{t,K}^{(i)}}$$

Viterbi to compute  $\arg\max_{\bar{x}_{1:T_i}^{(i)}} p(z_{1:T_i}^{(i)} | \bar{x}_{1:T_i}^{(i)}; G^{(i)})$

(max product)

Information theory

KL divergence: for discrete dist.  $p \neq q$

$$|\text{KL}(p || q) \triangleq \sum_{x \in \Omega} p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_p [\log \frac{p(x)}{q(x)}]$$

$$0 \cdot \log 0 = 0$$

$$(\lim_{x \rightarrow 0^+} x \log x = 0)$$

$$[\text{if } \exists x \text{ s.t. } q(x)=0$$

$$\text{but } p(x) \neq 0$$

$$-\underbrace{p(x) \log q(x)}_{\text{not zero}} = +\infty$$

If support of  $p \not\subseteq$  support of  $q$   
 $\Rightarrow KL(p||q) = +\infty$

motivation from density estimation

recall statistical decision theory

(statistical) loss  $L(p_\theta, a)$

world

here, estimation of dist., say  $\hat{q}$

Standard (MLE) loss is log-loss  $L(p_\theta, \hat{q}) = \mathbb{E}_{x \sim p_\theta} [-\log \hat{q}(x)]$

If we  $\hat{q} = p_\theta$ , then get  $\sum_{x \in \Omega_X} -p_\theta(x) \log p_\theta(x) \stackrel{\text{(cross-entropy)}}{=} H(p_\theta)$   
 entropy of  $p_\theta$

excess loss for action  $a = \hat{q}$

$$L(p, \hat{q}) - \min_q L(p, q) = \sum_{x \in \Omega_X} -p(x) \log \frac{\hat{q}(x)}{p(x)} = KL(p||\hat{q})$$

coding theory :

use length of code  $\log -\log p(x)$

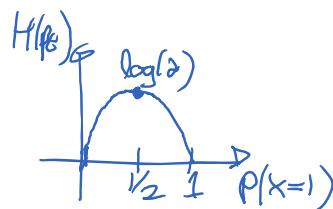
expected length of code :  $\sum_x p(x) (-\log p(x))$  (entropy measured in bits)

$KL$  divergence  $\rightarrow$  interpreted as excess length of cost (in terms of length to use dist. of to design code of code vs. optimal dist.  $p_{true}$ )

example :

entropy of a Bernoulli:

$$-p \log p - (1-p) \log(1-p)$$



entropy for a uniform dist. on  $k$  states

$$\sum_{x=1}^k -\frac{1}{k} \log \left(\frac{1}{k}\right) = \log k$$

(max-entropy dist. over  $k$  states)

properties of KL

•  $KL(p||q) \geq 0$  & to show this, use Jensen's inequality  $f(\mathbb{E}x) \leq \mathbb{E}f(x)$

when  $f$  is convex

•  $KL$  is strictly convex in each of its arguments i.e.  $KL(p||\cdot) : \mathbb{R}^k \rightarrow \mathbb{R}$

INFO  
THEORY  
BLOCK

$KL(\cdot \| q)$ Strictly  
convex  
sf.

- not symmetric:  $KL(p\|q) \neq KL(q\|p)$  in general  $KL(p\|q) = 0 \iff p \in \Delta_K$

symmetrized  
version  $\frac{1}{2} (KL(p\|q) + KL(q\|p))$