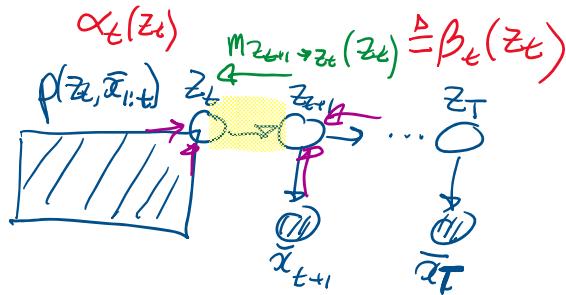


today:
 - EM for HMM
 - info theory { KL }

β -recursion: (smoothing)



$$p(z_t, \bar{x}_{1:T}) = \sum_z \alpha_t(z_t) \underbrace{m_{z_{t+1} \rightarrow z_t}(z_t)}_{\hat{=} \beta_t(z_t)}$$

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

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

β -recursion (aka. backward recursion)

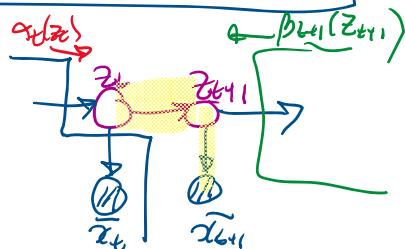
turns out that $\boxed{\beta_t(z_t) \hat{=} p(\bar{x}_{t+1:T}|z_t)}$

why? $p(z_t, \bar{x}_{1:T}) = p(\bar{x}_{1:T}|z_t) p(z_t)$
 $\stackrel{\text{def}}{=} p(\bar{x}_{t+1:T}|z_t) p(x_{1:t}|z_t) p(z_t)$

$$\Rightarrow \beta_t(z_t)$$

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

edge marginal:



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

Numerical stability trick:

$\dots \approx \dots$

Numerical stability trick:

Issue: $\alpha_t \& \beta_t$ can easily go to $1e-1000$

two possibilities:

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

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

call $\tilde{\alpha} \triangleq \max_i \alpha_i$

$$= \log \tilde{\alpha} + \log\left(\sum_i \frac{\alpha_i}{\tilde{\alpha}}\right)$$

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

$$\log\left(\sum_i \exp(\log \alpha_i - \log(\tilde{\alpha}))\right)$$

$$\log(\tilde{\alpha}) = \log(\alpha_{i_{\max}}) = \max_i \log(\alpha_i)$$

$$\log\left(\sum_i \alpha_i\right) = \log(\tilde{\alpha}) + \log\left(1 + \sum_{i \neq i_{\max}} \exp(\overbrace{\log(\alpha_i)}^{\log(\alpha_i) - \log(\tilde{\alpha})})\right)$$

b) normalize the message

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

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

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

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

• β -recursion:

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

note: $\sum_{z_t} \tilde{\beta}_t(z_t) \neq 1$

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

15h21

ML for HMM

for some parametric model for dist. on z_t

ML for HMM

for some parametric model for dist. on x_t

- suppose $p(x_t | z_t=k) = f(x_t | \eta_k)$ e.g. Gaussian on x_t

$$\eta = (\eta_k)_{k=1}^K$$

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

- $p(z_1=i) = \pi_i$

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

want to estimate $\hat{\pi}, \hat{A} \notin \hat{\pi}$ by ML from data $x = (x^{(i)})_{i=1}^N$

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

→ use EM at s^{th} iteration

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

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

complete log-likelihood

$$\log p(x, z | \Theta) = \sum_{i=1}^N \left(\underbrace{\log p(z_i^{(i)})}_{\substack{\text{huge variable} \\ \sum_k z_{i,k} \log \pi_k}} + \sum_{t=1}^{T_i} \underbrace{\log f(\bar{x}_t^{(i)} | z_t^{(i)})}_{\sum_k z_{t,k}^{(i)} \log f(\bar{x}_t^{(i)} | \eta_k)} + \sum_{t=2}^{T_i} \underbrace{\log p(z_t^{(i)} | z_{t-1}^{(i)})}_{\sum_{l,m} z_{t,l}^{(i)} z_{t-1,m}^{(i)}} \right)$$

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

$$\mathbb{E}_{q_{s+1}} [z_{t,R}^{(i)}] = q_{s+1}(z_{t,R}^{(i)}=1) \triangleq \gamma_{t,R}^{(i)} \quad \begin{array}{l} \text{calculated using} \\ \text{eq recursion} \\ \text{on } i^{\text{th}} \text{ time series} \end{array}$$

$$\text{smoothing def: } p(z_{t,R}^{(i)}=1 | x_{1:T_i}^{(i)}; \Theta^{[s]})$$

$$q_{s+1}(z_{t,R}^{(i)}=1, z_{t-1,m}=1) = p(z_{t,R}^{(i)}=1, z_{t-1,m}=1 | x_{1:T_i}^{(i)}; \Theta^{[s]}) \triangleq \gamma_{t,R,m}^{(i)}$$

A matrix
 $\begin{pmatrix} & m \\ l & \otimes \end{pmatrix}$

note: you should have

$$\sum_l \gamma_{t,R,m}^{(i)} = \gamma_{t-1,m}^{(i)}$$

smoothing edge marginal
prob. in HMM
(α - β recursion)

max. with respect to Θ :

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

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

$\hat{\pi}_k \rightarrow$ soft count
ML

$$\pi_{t,K}^{(i)} = \sum_{j=1}^K \gamma_{t,j,K}$$

$\sum_{j=1}^K \sum_{i=1}^N \gamma_{t,j,K}^{(i)}$

1

$$A_{l,m} = \sum_{i=1}^N \sum_{t=2}^T \gamma_{t,i,l,m}$$

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

$M_K \rightarrow$ soft count
ML
 $(\hat{M}_K, \hat{\Sigma}_K)$
for Gaussian

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

$$\hat{M}_K = \frac{\sum_{i=1}^N \sum_{t=1}^T \gamma_{t,i,K}^{(i)} x_t^{(i)}}{\sum_{i=1}^N \sum_{t=1}^T \gamma_{t,i,K}^{(i)}}$$

Viterbi to compute argmax $p(z_{1:T} | \tilde{x}_{1:T})$
(max product)

Information Theory

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

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

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

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

[If $\exists x$ s.t. $q(x) = 0$
but $p(x) \neq 0$
 $-p(x) \log q(x)$ not zero $\rightarrow +\infty$]

$$KL = \infty \text{ if } \text{sup}(p) \not\subseteq \text{sup}(q)$$

motivation from density estimation

recall statistical decision theory

(statistical) loss

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

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

cross-entropy

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

If use $\hat{q} = p_\theta$, then get

$$\sum_{x \in \Omega_X} -p_\theta(x) \log p_\theta(x) \triangleq H(p_\theta)$$

entropy of p_θ

$$\text{If } p(x) = q(x), \text{ then } \sum_{x \in \Omega_X} p(x) \log \frac{q(x)}{p(x)} = 0$$

entropy of p

excess loss for action $a = q$

$$L(p, \hat{q}) - \min_q L(p, q) = \sum_{x \in \Omega} -p(x) \log \frac{\hat{q}(x)}{p(x)} = kL(p, \hat{q})$$

Coding theory:

use length of code $\mathcal{L} = -\lg p(x)$

$\lg \equiv \log_2 \rightarrow \text{"bits"}$
 $\log_e \rightarrow \text{"nats"}$

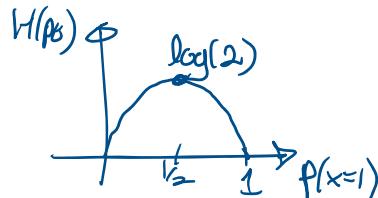
expected length of code : $\sum_x p(x) (-\lg p(x))$ (entropy measured in bits)

KL divergence \rightarrow interpreted as excess length cost (in terms of length of code)
 to use dist q to design code vs optimal dist (the p)

Example:

entropy of a Bernoulli

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



entropy of 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

• $KL(p || q) \geq 0$ \leftarrow 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. $KL(p || \cdot)$: $\Delta_k \subseteq \mathbb{R}^k \rightarrow \mathbb{R}$
 ↳ strictly convex

$$KL(\cdot || q)$$

• not symmetric $KL(p || q) \neq KL(q || p)$
 in general

$$KL(p || p) = 0 \quad \forall p \in \Delta_k$$

commutative

$$(f_1 \circ f_2)(x) = f_2(f_1(x))$$

symmetrized
version

"in general"

$$\frac{1}{2}(k_L(p||q) + k_C(q||p))$$