

today: • EM for HMM
 • info. theory & KL
 man, entropy

continue α -recursion:

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

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

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

initialization:

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

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

Hadamard product

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

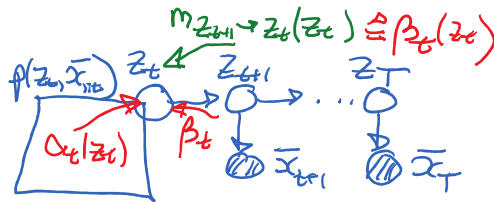
space complexity: $O(k)$ extra storage

time complexity: $O(tk^2)$

$$\left. \begin{aligned} &= \alpha_t(z_t) \\ &\sum_{z_t} \end{aligned} \right\} = \sum_{z_t} p(z_t, \bar{x}_{1:t}) = p(\bar{x}_{1:t})$$

"evidence probability"

β -recursion (smoothing)



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

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

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

β -recursion (aka backward recursion)

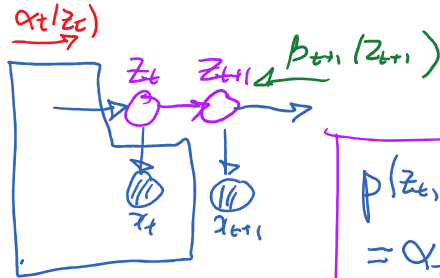
turns out that $\beta_t(z_t) \triangleq p(\bar{x}_{t+1:T} | z_t)$ why?

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

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

$$\Rightarrow \beta_t(z_t) \stackrel{\text{CI}}{=} p(\bar{x}_{t+1:T} | z_t) \cdot \underbrace{p(\bar{x}_{1:t} | z_t) p(z_t)}_{\alpha_t(z_t)} //$$

edge marginal



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

Numerical stability trick:

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

two possibilities:

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

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

call $\tilde{a} \triangleq \max_i a_i$
 $i_{\max} \triangleq \text{argmax}_i a_i$

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

b) normalize the message

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

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

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

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

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

• β -recursion:

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

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

exercice: derive β -recursion

14h24

ML & HMM

- suppose $p(x_t | z_t = k) = f(x_t | \mu_k)$
- $p(z_{t+1} = i | z_t = j) = A_{ij}$
- $p(z_1 = i) = \pi_i$

some parametric model for dist. on x_t
e.g. Gaussian on x_t

$$\mu = (\mu_k)_{k=1}^K$$

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

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

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

→ use EM

s^{th} iteration

E step:

$$q_{s+1}(z) = p(z | x, \theta^{[s]})$$

M step:

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

Complete log-likelihood:

$$\log p(x, z | \theta) = \sum_{i=1}^N \left(\log p(z_i^{(i)}) + \sum_{t=1}^{T_i} \log p(\tilde{x}_t^{(i)} | z_t^{(i)}) + \sum_{t=2}^{T_i} \log p(z_t^{(i)} | z_{t-1}^{(i)}) \right)$$

\nearrow $\sum_k z_{1:k}^{(i)} \log \pi_k$
 \nearrow $\sum_k z_{t:k}^{(i)} (\log f(\tilde{x}_t^{(i)} | \mu_k))$
 \nearrow $\sum_{l,m} z_{t:l}^{(i)} z_{t+1:m}^{(i)} \log A_{lm}$

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

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

smoothing dist $p(z_{t,k} = 1 | \tilde{x}_{1:T}^{(i)}; \theta^{[s]})$

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

$$= p(z_{t,l} = 1, z_{t+1,m} = 1 | \tilde{x}_{1:T}^{(i)}; \theta^{[s]}) \triangleq \tau_{t,l,m}^{(i)}$$

A- β recursion

$$= p(z_{t,l}^{(i)}=1, z_{t,l}^{(i)}=m=1 | \bar{x}_{1:T}^{(i)}; \theta^{(i)}) \triangleq \tau_{t,l,m}^{(i)}$$

↑ smoothing edge marginal $(l \rightarrow \uparrow)$

maximize with respect to θ :

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

$$\left. \sum_{i=1}^N \sum_{l=1}^L \tau_{l,e}^{(i)} \right\} N$$

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

$\hat{\pi}_k \rightarrow$ soft-counts ML

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

$$\hat{\mu}_k = \frac{\sum_{i=1}^N \sum_{t=2}^{T_i} \tau_{t,k}^{(i)} x_t^{(i)}}{\sum_{i=1}^N \sum_{t=2}^{T_i} \tau_{t,k}^{(i)}}$$

correction
2020/11/19
from original
scribbles

Viterbi to compute $\arg \max_{z_{1:T}} p(z_{1:T} | \bar{x}_{1:T})$
(max product)

Information theory

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

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

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

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

[if $\exists x$ st.
 $q(x) = 0$
but $p(x) \neq 0$

$$-p(x) \log q(x) = +\infty$$

if support of $p \not\subseteq$ support of $q \Rightarrow KL(p \parallel q) = +\infty$

motivation from density estimation

recall statistical decision theory would here, estimation of dist., say \hat{q}
(statistical) loss $L(p_\theta, \hat{q})$

standard (MSE) loss is log-loss $L(p, \hat{q}) = \mathbb{E}_{x \sim p} [-\log \hat{q}(x)]$

if use $\hat{q} = p$, then get $\sum_{x \in \mathcal{X}} -p(x) \log p(x) \triangleq H(p)$ (cross-entropy)
 entropy of p

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

$$L(p, \hat{q}) - \min_q L(p, q) = -\sum_x p(x) \log \frac{\hat{q}(x)}{p(x)}$$

$\log \triangleq \log_2 \rightarrow$ "bits" $\log_e \rightarrow$ "nats"

Coding theory:

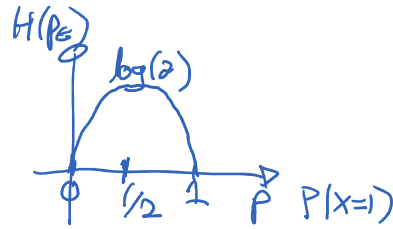
use length of code $\propto -\log p(x)$

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

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

Example:

entropy for 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:

• $KL(p||q) \geq 0$ \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. $KL(p||\cdot) : \mathbb{R}^k \rightarrow \mathbb{R}$
 $KL(\cdot||q)$

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

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

strictly convex

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