

# Lecture 15 - scribbles

Tuesday, October 31, 2017

14:30

today: finish HMM + EM  
information theory & KL

HMM continued:

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

$\alpha$ -recursion:  $\alpha_t(z_t) = p(\bar{x}_t | z_t) \sum_{z_{t-1}} p(z_t | z_{t-1}) \alpha_{t-1}(z_{t-1})$

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

$\downarrow$  observation prob.  $p(\bar{x}_t | \cdot)$   
vec for  $z_t$

numerical trick & issue:  $\alpha_t$  can easily be  $1e-100$

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

$$a+b = \exp(\log(a)) + \exp(\log(b))$$

say  $a \geq b$

$$= a (1 + \exp(\log(b) - \log(a)))$$
$$\log(a+b) = \log(a) + \log(1 + \exp(\log(b) - \log(a)))$$

in general  $\sum_i a_i$   
first pick  $\max_i a_i$

b) normalize the messages

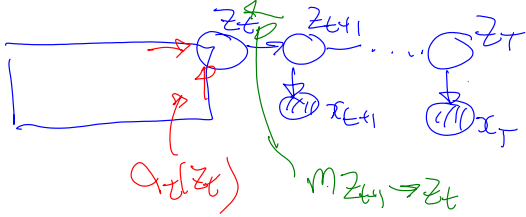
ie.  $\tilde{\alpha}_t(z_t) = p(z_t | \bar{x}_{1:t})$

before  $\alpha_t = O_t \odot A \alpha_{t-1}$       $\tilde{\alpha}_t = O_t \odot A \tilde{\alpha}_{t-1}$      you can show that...

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

$\tilde{\alpha}_t = \frac{O_t \odot A \tilde{\alpha}_{t-1}}{\sum_{z_t} [ \text{ " } ]}$  you can show that  
 $\tilde{\alpha}_t$  is  $\propto p(\tilde{x}_t | x_{1:t-1})$   
normalization constant

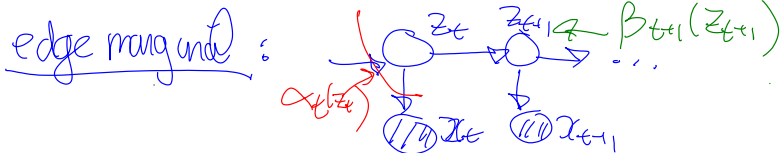
smoothing:  $p(z_t, \tilde{x}_{1:T}) = \prod_{z=1}^T \alpha_t(z_t) \underbrace{M_{z_{t+1} \rightarrow z_t}(z_t)}_{\triangleq \beta_t(z_t)}$



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

$\beta_t(z_t) = \sum_{z_{t+1}} p(z_{t+1} | z_t) p(\tilde{x}_{t+1} | z_{t+1}) \beta_{t+1}(z_{t+1})$   
 $\beta$ -recursion (backward recursion)

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



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

ML for HMM:

- suppose  $p(x_t | z_t = k) = f(x_t | n_k)$        $n = (n_k)_{k=1}^K$
- $p(z_{t+1} = i | z_t = j) = A_{ij}$
- $p(z_1 = i) = \pi_i$        $\Theta = (n, A, \pi)$

want to estimate  $\hat{n}, \hat{A}, \hat{\pi}$  by ML from data  $\{x^{(i)}\}_{i=1}^N$ ;  $x^{(i)} = x_{1:T}^{(i)}$   
 → use EM      sth iteration

E step:  $Q_{S_{t+1}}(z) = p(z | x, \hat{\theta}^{[S_t]})$

M step:  $\hat{\theta}^{[S_{t+1}]} = \arg \max_{\theta \in \Theta} E_{Q_{S_{t+1}}} [\log p(z, z)]$

complete log-likelihood:

$$\log p(x, z | \theta) = \sum_{i=1}^N \left[ \underbrace{\log p(z_1^{(i)})}_{\sum_k z_{1,k}^{(i)} \log \pi_k} + \sum_{t=1}^T \underbrace{\log p(\bar{x}_t^{(i)} | z_t^{(i)})}_{\sum_k z_{t,k}^{(i)} \log \xi(x_t^{(i)} | m_k)} + \sum_{t=2}^T \underbrace{\log p(z_t^{(i)} | z_{t-1}^{(i)})}_{\sum_{l,m} z_{t,l}^{(i)} z_{t-1,m}^{(i)} \log A_{lm}} \right]$$

$E_{Q_{S_{t+1}}} [\log p(z, z)] = \dots$

$E_{Q_{S_{t+1}}} [z_{t,k}^{(i)}] = Q_{S_{t+1}}(z_{t,k}^{(i)} = 1) \triangleq \hat{\pi}_{t,k}^{(i)}$

$Q_{S_{t+1}}(z_{t,l}^{(i)} = 1, z_{t-1,m}^{(i)} = 1) = p(z_{t,l}^{(i)} = 1, z_{t-1,m}^{(i)} = 1 | x_{1:t}^{(i)}, \hat{\theta}^{[S_t]})$   
 smoothing marginal

$\hat{\pi}_{t,l,m}^{(i)}$

maximizing with respect to  $\theta$ :

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

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



"Baum-Welch" alg.

- forward-backward alg.  
 (alpha-beta recursion / sum-product)  
 + EM for HMM

[https://en.wikipedia.org/wiki/Baum%E2%80%93Welch\\_algorithm](https://en.wikipedia.org/wiki/Baum%E2%80%93Welch_algorithm)

note:

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

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

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

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

$$= p(\bar{x}_{t+1:T} | z_t) \text{ by cond. indep.}$$

Information theory:

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

$$KL(p||q) \triangleq \sum_x p(x) \log \frac{p(x)}{q(x)}$$

$0 \cdot \log 0 = 0$

motivation from density estimation:

recall statistical decision theory

(statistical) loss  $L(p, a)$

action is estimating a distribution say  $\hat{q}$

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

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

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

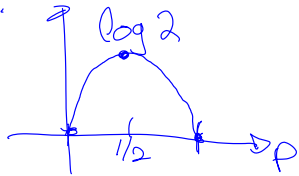
$$L(p, \hat{q}) - \min_q L(p, q) = \frac{L(p, \hat{q})}{L(p, p)} = -\sum_{x \in \Omega_x} p(x) \log \frac{\hat{q}(x)}{p(x)} = KL(p||\hat{q})$$

coding theory:

use length of code  $\alpha_q - \log p(x)$   $\log_2 \rightarrow$  "bits"  $\log \rightarrow$  "nats" [minimal expected length]

expected length of code:  $\sum_x p(x) (-\log p(x))$  (like the entropy)  
 KL divergence  $\rightarrow$  excess cost (in terms of <sup>expected</sup> length of code)  
 to use dist.  $q$  for coding instead of  $p$

entropy a Bernoulli:



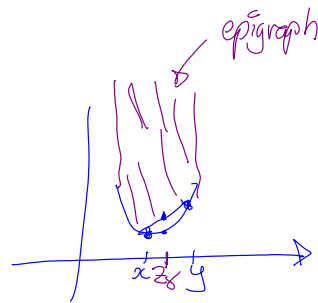
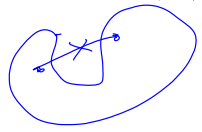
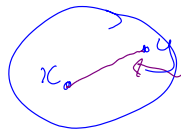
entropy of uniform dist. on  $k$  states:  $-\sum_x \frac{1}{k} \log(\frac{1}{k}) = \log k$

properties of KL:

- $KL(p||q) \geq 0$
- is strictly convex in each argument i.e.  $KL(p||\cdot)$   
 $KL(\cdot||q)$

convexity review:

set  $A$  is convex  $\Leftrightarrow \forall x, y \in A$   
 $(1-\gamma)x + \gamma y \in A$   
 for  $\gamma \in [0,1]$



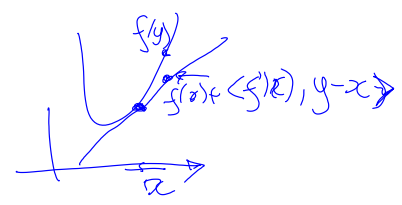
function  $f$  is convex  $\Leftrightarrow$  its epigraph is convex

$$\hookrightarrow \{ (x,t) : x \in \text{dom}(f), t \geq f(x) \}$$

$$\Leftrightarrow f(\underbrace{(1-\delta)x + \delta y}_{z_\delta}) \stackrel{\text{convex}}{\leq} (1-\delta)f(x) + \delta f(y) \quad \forall \delta \in [0,1] \quad \forall x,y \in \text{dom}(f)$$

$<$   
"strictly convex"

If  $f$  is differentiable,  $f(y) \geq f(x) + \langle f'(x), y-x \rangle \quad \forall y, x$



recall Jensen:  $f(\mathbb{E}X) \leq \mathbb{E}f(X) \quad (f \text{ convex})$

ML and KL minimization:

$\{P_\theta\}_{\theta \in \Theta}$  parametric family empirical dist.

then ML for  $\Theta \Leftrightarrow \min_{\theta \in \Theta} \text{KL}(\hat{P}_n \| P_\theta)$

proof:  $\text{KL}(\hat{P}_n \| P_\theta) = \sum_x \hat{p}_n(x) \log \frac{\hat{p}_n(x)}{P_\theta(x)}$

$$= -H(\hat{P}_n) - \sum_x \hat{p}_n(x) \log P_\theta(x)$$

$\hookrightarrow \sum_{i=1}^n \delta(x, x^{(i)})$  (Kronecker-delta)

$$= \underbrace{-H(\hat{P}_n)}_{\text{constant}} - \frac{1}{n} \sum_{i=1}^n \log P_\theta(x^{(i)})$$

$\log \prod_{i=1}^n P_\theta(x^{(i)})$

next: max entropy

next: max entropy

//