

EM: can we use latent variables to devise algo?

- will consider this during lecture to be in an MLE setting, altho would be fun as MAP (just would have prior terms floating around)

We start  $\theta^{(0)}$  randomly. Run til "convergence"  
i = iteration count.

E step: we have cluster assignments  $c_t$  ( $c_t \in 1, \dots, k$ )  
define a distribution wrt each  $t$ :

$$Q_t^{(i)}(c_t) = P(c_t | x_t, \theta^{(i-1)})$$

$$= \begin{bmatrix} .3 \\ .7 \\ 0 \end{bmatrix} \begin{matrix} \text{lemons} \\ \text{oranges} \\ \dots \end{matrix}$$

This gives you a way to have ~~an~~ a distribution to take expectation wrt.

maximization step:  
for  $Q_t^{(i)}(c_t)$  fixed, maximize a "weighted" likelihood:

$$\theta^{(i)} = \underset{\theta}{\operatorname{argmax}} \sum_t \sum_{c_t} \underbrace{Q_t^{(i)}(c_t)}_{\text{weight}} \log P(x_t, c_t | \theta)$$

looks like a likelihood

< may have to manually ~~adjust~~ adjust the  $\theta$ 's &  $c_t$ 's >

< e.g. repeat what 1st step is doing >

example for GMM: for every  $k \in [k]$ ,

$$Q_t^{(i)}(c_t = k) = P(c_t = k | x_t, \theta^{(i-1)})$$

$$= \frac{P(x_t | c_t = k, \theta^{(i-1)}) P(c_t = k | \theta^{(i-1)})}{P(x_t, \theta^{(i-1)})}$$

< typo on slide >

$$\text{Gaussian}(x_t; \mu_k^{(i-1)}, \Sigma_k^{(i-1)})$$

this is  $\pi_k$ , our mixture guess (marginal prob)

you could expand this out.

this is easy (- in theory).

- now, for the M step:  $\theta$  not similar to last time:

$$\theta^{(i)} = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_t \sum_k Q_t^{(i)}(k) \log P(x_t, c_t = k | \theta)$$

(note: if you take derivatives, the  $\frac{\partial}{\partial \cdot}$  will go inside the sum.)

question: are the  $Q$ 's like what we saw last time??  
Yes, by making the  $Q_t^{(i)}$ 's to be "indicator fns".

for GMM, similar also doing:  $\underset{\pi, \mu, \Sigma}{\operatorname{argmax}} \sum_c \sum_{c=1}^k Q_t^{(i)}(k) (\log \text{Gauss}(x_t | \mu_k, \Sigma_k) + \log \pi_k)$

and the sol'n will consist to: for each  $k$ :

$$\mu_k^{(i)} = \frac{\sum_t Q_t^{(i)}(k) x_i}{\sum_t Q_t^{(i)}(k)}$$

$$\sigma_k^{(i)} = \frac{\sum_t Q_t^{(i)}(k) (x_i - \mu_k^{(i)}) (x_i - \mu_k^{(i)})^T}{\sum_t Q_t^{(i)}(k)}$$

weight on outer product that is your computation of the empirical covariance.

$$\pi_k^{(i)} = \frac{\sum_t Q_t^{(i)}(k)}{n}$$

will post the calc's for  $\sum_i \pi_i$ , but let's try  $\hat{\theta}_k$ :

→ actually, let's talk about "why it makes sense", then go back to derivative of  $\mu_k$ .

Why does this work? Intuition:

Each  $\theta$  never decreases log-likelihood, and neither does M-step.

→ to show:  $\log \text{Lik}(\theta^{(i+1)}) \geq \log \text{Lik}(\theta^{(i)})$ .

steps: insert latent vars, use Jensen's inequality, massage.

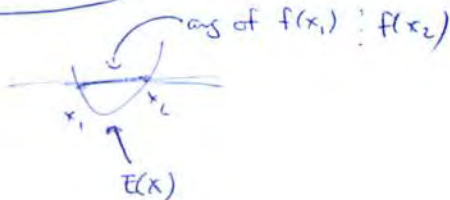
→ cornerstone of convex (did he say analysis) (line is above the curve).  
log is concave, which goes the other way.

$$\sum_t \log P(x_t | \theta^{(i+1)}) = \sum_t \log \sum_{k=1}^K P(x_t, c_t = k | \theta^{(i+1)})$$

$$= \sum_t \log \sum_{k=1}^K \frac{Q_t^{(i+1)}(k)}{Q_t^{(i+1)}(k)} P(x_t, c_t = k | \theta^{(i+1)})$$

← intro of variable you will marginalize out

Jensen's in = : for  $f$  convex,  $f(E(x)) \leq E(f(x))$   
↑  
avg value



(it's practically a defn of convexity)

consider: and let's call it ~~Q~~  $H(k)$ , so to get  $\sum_t \log \sum_{k=1}^K Q_t^{(i+1)}(k)$

note: fn of  $k$

$$\sum_t \log E_{k \sim Q_t^{(i+1)}} [H(k)]$$

→ by Jensen on concave fns, such as the log.

$$\geq \sum_t \sum_k \underbrace{Q_t^{(i+1)}(k)}_{\downarrow} \cdot \log H(k)$$

$$= \sum_t \sum_k Q_t^{(i+1)}(k) \left[ \log P(x_t, c_t = k | \theta^{(i+1)}) \right] = \frac{\sum_t \sum_k Q_t^{(i+1)}(k) \log Q_t^{(i+1)}(k)}{\log}$$

↳ no need to depend on  $\theta$

↓ no depend on  $\theta$

for the  $i+1$ th iteration, we're picking the thing  $\theta$  that maximizes 1<sup>st</sup> term.

→ 1<sup>st</sup> term must be bigger than its value @  $\theta = \theta^{(i)}$ .

So, then, why did we do the expected log-likelihood?

this comes in handy here:

$$\sum_n \sum_{\ell} Q_t^{(in)}(\ell) \log \left( \frac{P(x_t, c_t = \ell | \theta^{(i)})}{Q_t^{(in)}(\ell)} \right) = \sum_n \sum_{\ell} Q_t^{(in)}(\ell) \log \frac{P(c_t = \ell | x_t, \theta^{(i)})}{Q_t^{(in)}(\ell)} P(x_t | \theta^{(i)})$$

no dependence on  $\ell$ , so

$$= \sum_n \sum_{\ell} Q_t^{(in)}(\ell) \log P(x_t | \theta^{(i)}) = \sum_n \log P(x_t | \theta^{(i)})$$

can hit a stationary point

an intuitive: hallucinate the hidden variables' values.

this isn't specific to cluster assignments



k-means: a hard version (hard assignment) on G's. (and assume spherical covariance)

soft k-means: still assume spherical Gaussians.