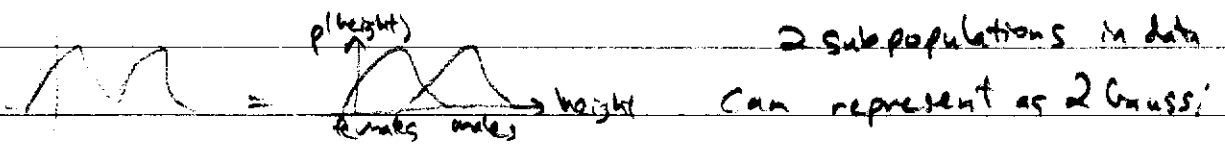


Mixture Models

In e-fam: $\log \left(\prod_{n=1}^N p(x_n | \theta) \right)$
 $\sum_{n=1}^N \log p(x_n | \theta)$ if $p(x|\theta) = \exp(\eta(x) \cdot \theta - A(\theta))$
 log & exp annihilate

Easy MLE, nice unique max. step. $\frac{\partial}{\partial \theta} := 0$

But what if want more general / powerful distrib's?



2 subpopulations in data can represent as 2 Gauss.

or non-linear complicated distrib 3 Gauss.

consider mixtures of Gaussians (multinomial e-families)

$$p(x) = \sum_i \alpha_i \frac{1}{(2\pi)^{D/2} |\Sigma_i|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\right)$$

$N(x | \mu_i, \Sigma_i)$

mixture \rightarrow hidden variable / latent / unobserved class label

- male / female
- or conceptual hidden variable



$$p(x|\theta) = \sum_i \pi_i N(x | \mu_i, \Sigma_i) \quad \text{or} \quad \sum_i \pi_i \exp(\eta(x) + x^T \theta_i - A(\theta_i)) \quad \begin{matrix} \pi_i \geq 0 \\ \sum \pi_i = 1 \end{matrix}$$

- subpopulations in data
- hidden variable selects between them (coin flip)
- clustering: interpret the sub-pops. (i.e. data contains 2 doc types)
- density est: more complicated distrib
- classification: we observed label, here it is hidden

in mixture model, label is unobserved

hidden vars: $z = \begin{bmatrix} z^1 \\ \vdots \\ z^k \end{bmatrix}$ $z^i = 1$ or not or $z = 1 \dots k$

$z \rightarrow x$
 $p(z|x) = p(z)p(x|z)$

$p(x|\theta) = \sum_z p(x,z|\theta) = \sum_z p(z|\theta) p(x|z, \theta)$

or

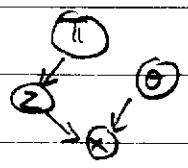
$p(x|\theta) = \sum_i p(z^i=1|\pi_i) p(x|z^i=1, \theta)$

$\pi_i = p(z^i=1) = p(z^i=1|\pi_i)$ = prior probs = mixing proportions, $\sum \pi_i = 1$

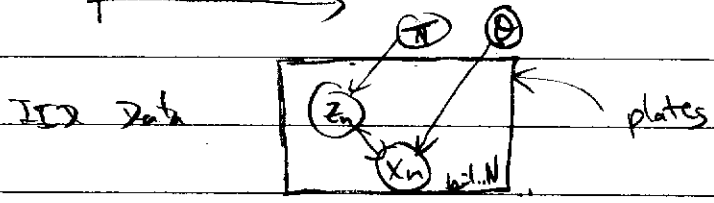
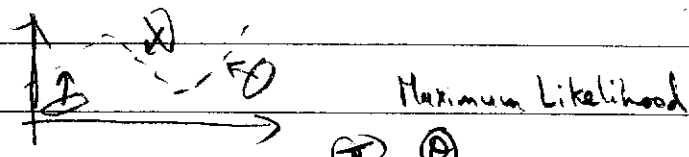
$p(x|z^i=1, \theta)$ = mixture components

$\gamma_i = p(z^i=1|x, \theta)$ = posterior probs = responsibilities, $\sum \gamma_i = 1$

$= \frac{p(x|z^i=1, \theta) p(z^i=1|\theta)}{p(x|\theta)}$



Learning with Mixtures? Given Data how to fit model?



$p(\theta|x) \propto p(x|\theta)$ if ML, uniform prior on model

$P(x|\theta) = \prod_{n=1}^N p(x_n|\theta)$

$l(\theta) = \log P(x|\theta) = \sum_{n=1}^N \log p(x_n|\theta) = \sum_{n=1}^N \log \sum_{i=1}^K \pi_i N(x_n | \mu_i, \Sigma_i)$

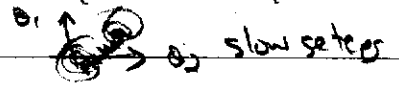
maximization step is ugly

or expo-fam

- parameters are coupled via log-sum

- separate or decouple? Reuse E-Family Machinery?

- gradient descent / newton raphson



Problem with log-sum:

$$\sum e\text{-family} \neq e\text{-family}$$

$$\prod e\text{-family} = e\text{-family}$$

So mixture models break clean ML for mix models

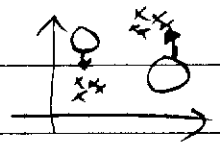
- We know how to estimate 1 single distrib in mix with ML

Can we break down several mixture model estimates into multiple e-family estimates?

- Heuristic Divide & Conquer ... K-means ... formal ... EM

K-Means - Old heuristic algorithm

- "gobble up" data with divide & conquer



Timeline:

Year	Contributors
1977	Dempster, ... , Baum & Elgin
1980's	CS. Szar & Tusnady
1993	Neal & Hinton
1995	em, Amari, information geometry

Chicken & Egg: If we know classes, easy to get model (e-family)
If we know model, easy to get classes.

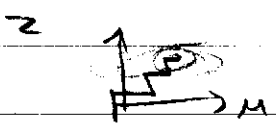
- for each x_n , define z_n indicator vector which classifies it, $\sum_i z_n^i = 1$

K-Means:
$$\mu_i = \frac{\sum_n z_n^i x_n}{\sum_n z_n^i}$$

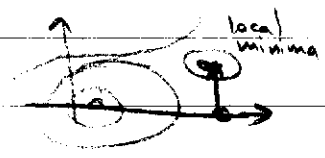
$$z_n^i = \begin{cases} 1 & \text{if } i = \text{argmin}_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

↑ closest point gets "gobbled up"

K-means Cost fn $\min_{\mu} \min_z J(\mu_1, \dots, \mu_k, z_1, \dots, z_n) = \sum_{k=1}^K \sum_{n: z_n=k} \|x_n - \mu_k\|^2$



coordinate descent
axis parallel optimization
alternating minimization



EM - Maximum Likelihood for Mix Models

outline: - EM for mix of Gaussians

- EM as bound maximization
- EM as alternating maximization
- EM as shading nodes
- EM as Min KL & Info-Geometry...

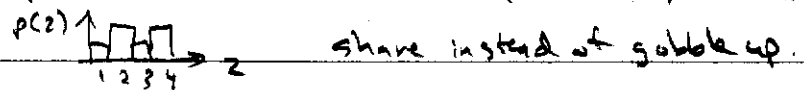
EM mix of Gaussians - "fuzzify" - Softening K-means, instead of winner takes all

slower than
K-means but
less greedy

K-means: $z = \arg \max_i p(x|z^i=1, \theta)$

Soft assignment, $\propto \pi_i \frac{1}{\sigma_i} \exp(-\frac{1}{2} \|x_n - \mu_i\|^2)$

look at $\tau_n = p(z|x_n)$ as shared responsibility for a point



$$\mu_i = \frac{\sum_n \tau_n^i x_n}{\sum_n \tau_n^i}$$

← use expectation of z, instead of hard select

$$\tau_n^i = E[z_n^i | x_n] = p(z^i=1 | x_n)$$

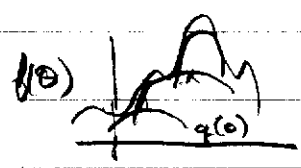
E Step: $\tau_n^i = \frac{\pi_i \mathcal{N}(x_n | \mu_i^{(t)}, \sigma_i^{(t)})}{\sum_j \pi_j \mathcal{N}(x_n | \mu_j^{(t)}, \sigma_j^{(t)})}$

M Step: $\mu_i^{(t+1)} = \frac{\sum_n \tau_n^i x_n}{\sum_n \tau_n^i}$ $\sigma_i^{(t+1)} = \frac{\sum_n \tau_n^i (x_n - \mu_i^{(t+1)}) (x_n - \mu_i^{(t+1)})^T}{\sum_n \tau_n^i}$

$$\pi_i^{(t+1)} = \frac{1}{N} \sum_n \tau_n^i$$

Not yet formal. Not coord ascent.
Why converge? Why divide & conquer?
divide & conquer need not always work...

EM as Bound Maximization



How to Bound? Jensen
for concave f:
 $f(E\{x\}) \geq E\{f(x)\}$
i.e. if $\sum p_i = 1$ & $p_i \geq 0$
 $f(\sum p_i x_i) \geq \sum p_i f(x_i)$

 $f(\sum p_i x_i) \geq \sum p_i f(x_i)$

$$f(\theta) \geq q_t(\theta) \quad \forall \theta \quad \& \quad \forall t$$

$$f(\theta^t) = q_t(\theta^t) \quad q_t(\theta^{t+1}) > q_t(\theta^t)$$

$$f(\theta^{t+1}) \geq q_t(\theta^{t+1}) > q_t(\theta^t) = f(\theta^t)$$

$\therefore f(\theta^{t+1}) > f(\theta^t)$ monotonic increase

$$\max_{\theta} \ell(\theta) = \max_{\theta} \sum_{n=1}^N \log \sum_z p(x_n, z | \theta)$$

constant

$$\equiv \max_{\theta} \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) - \ell(\theta^t) \equiv \max_{\theta} \Delta \ell(\theta)$$

Same max locus yet $\Delta \ell(\theta^t) = 0$

$$\sum_{n=1}^N \log \sum_z p(x_n, z | \theta) - \sum_{n=1}^N \log \sum_z p(x_n, z | \theta^t)$$

$$\sum_{n=1}^N \log \frac{\sum_z p(x_n, z | \theta)}{\sum_z p(x_n, z | \theta^t)} = \sum_n \log \sum_z \frac{p(x_n, z | \theta)}{\sum_z p(x_n, z | \theta^t)} \times 1$$

$$= \sum_n \log \sum_z \frac{p(x_n, z | \theta)}{\sum_z p(x_n, z | \theta^t)} \frac{p(z | x_n, \theta^t)}{p(z | x_n, \theta^t)}$$

$$= \sum_n \log \sum_z p(z | x_n, \theta^t) \frac{p(x_n, z | \theta)}{\left(\sum_z p(x_n, z | \theta^t)\right) p(z | x_n, \theta^t)} = \sum_n \log \frac{p(x_n, z | \theta)}{p(x_n, z | \theta^t) p(z | x_n, \theta^t)}$$

$$\geq \sum_n \sum_z p(z | x_n, \theta^t) \log \frac{p(x_n, z | \theta)}{p(x_n, z | \theta^t) p(z | x_n, \theta^t)}$$

"Pulled" log inside sum

$Q(\theta | \theta^t)$ auxiliary function

$$\theta^{t+1} = \arg \max_{\theta} Q(\theta | \theta^t)$$

$$= \arg \max_{\theta} \sum_n \sum_z \tau_n^z \log p(x_n, z | \theta)$$

straight forward to derive w.r.t θ & set to 0

$$\frac{\partial}{\partial \mu_k} \left(\sum_n \tau_n^k (\mu_k - x_n) \right) = 0$$

EM as Expected Likelihood

$Q(\theta|\theta^t) = \text{expected "complete" likelihood} + \text{constant}$

- ② → ⊙ incomplete data since z unobserved (only x)
- ⊙ → ⊙ complete data if z is observed (like classification)

$l(\theta) = \text{incomplete log-likelihood} = \sum_n \log \sum_z p(x_n, z|\theta)$

$l^c(\theta) = \text{complete log-likelihood} = \sum_n \log p(x_n, z_n|\theta)$ ← no log-sum

but we don't know z , so use expected values of z assuming current model θ^t , i.e. $p(z_n|x_n, \theta^t)$

$E_{\text{distrib. over } z} [l^c(\theta)] = \sum_{z_1=1}^K \dots \sum_{z_N=1}^K \prod_{i=1}^N p(z_i|x_i, \theta^t) l^c(\theta)$

rewrite: $l^c(\theta) = \sum_n \log p(x_n, z_n|\theta)$

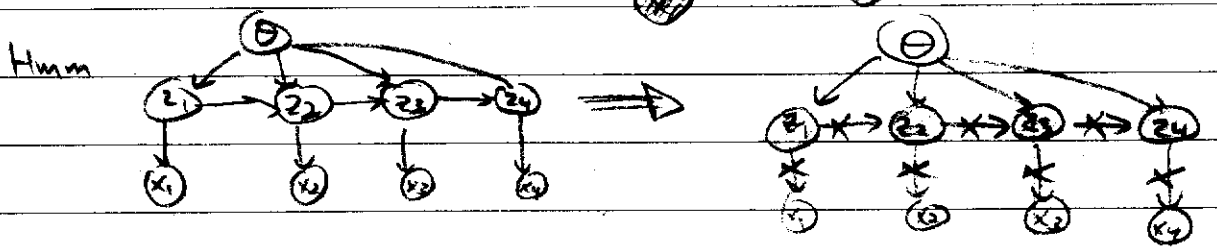
$E[l^c] = \sum_n \dots \sum_{z_n=1}^K \prod_{i=1}^N p(z_i|x_i, \theta^t) \sum_n \log p(x_n, z_n|\theta)$

$E[l^c] = \sum_n \sum_{z_n} p(z_n|x_n, \theta^t) \log p(x_n, z_n|\theta) \sum_{z_1=1}^K \dots \sum_{z_N=1}^K \prod_{i=1}^N p(z_i|x_i, \theta^t)$

$E[l^c] = \sum_{n=1}^N \sum_{j=1}^K p(j|x_n, \theta^t) \log p(x_n, j|\theta)$

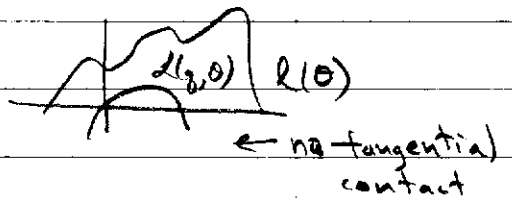
↑ like $Q(\theta|\theta^t)$ function + constant

EM as "filling in" Nodes



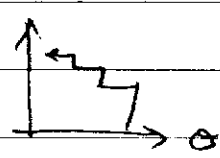
EM as Alternating Maximization ← more general (Neal & Hinton)

$$\begin{aligned}
 \ell(\theta) &= \sum_{n=1}^N \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) \\
 &= \sum_{n=1}^N \sum_z \log \sum_q q_n(z | x_n) \frac{p(x_n, z | \theta)}{q_n(z | x_n)} \quad \text{any } q \geq 0 \ \& \ \sum_z q_n(z | x_n) = 1 \\
 &\geq \sum_{n=1}^N \sum_z q_n(z | x_n) \log \frac{p(x_n, z | \theta)}{q_n(z | x_n)} \\
 &\geq \ell(q, \theta)
 \end{aligned}$$



E: $q^{t+1} = \arg \max_q \ell(q, \theta^t)$

M: $\theta^{t+1} = \arg \max_{\theta} \ell(q^{t+1}, \theta)$



$$\ell(q, \theta) = \sum_{n=1}^N \sum_z q_n(z | x_n) \log p(x_n, z | \theta) - \sum_{n=1}^N \sum_z q_n(z | x_n) \log q_n(z | x_n)$$

↑ Expected q [0, 1]
↑ constant with θ

different "q" though

MStep $\max_{\theta} \ell$ or $\frac{\partial}{\partial \theta} \ell = 0$ is easy, for e-family m-step is standard else "GEM" is partial m-step

EStep maximize over q , closes gap with current tangential contact

~~graph~~ $\max_q \ell(q, \theta) \Rightarrow q_n^* = p(z | x_n, \theta^t)$

check: $\ell(\theta) \geq \ell(\theta, q)$

$$\ell(\theta^t) \geq \ell(\theta^t, q^*) = \sum_{n=1}^N \sum_z p(z | x_n, \theta^t) \log \frac{p(x_n, z | \theta^t)}{p(z | x_n, \theta^t)} \Big|_{\theta = \theta^t}$$

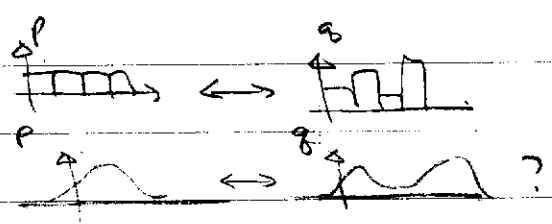
$$\ell(\theta^t) \geq \sum_{n=1}^N \sum_z \underbrace{p(z | x_n, \theta^t)}_{=1} \log p(x_n | \theta^t)$$

$$\ell(\theta^t) \geq \sum_{n=1}^N \log p(x_n | \theta^t) = \ell(\theta^t)$$

can't increase q any more since R.H.S. is equal to L.H.S. upper bound

else "Incremental EM" is partial E-step

Alternative Deriv. of ML



distance between two distrib

Dist $(p, q) = ?$ Euclidean? $\|p - q\|^2$ no, since $\sum p = 1$ & $p \geq 0$

Generalized Divergences: Bregman Divs.

$$\text{Euclidean, Mahalanobis, KL, } \kappa(p) - \kappa(q) - \kappa'(q)^T (p - q)$$

for any convex κ function on the space of p

for $\kappa = \text{neg entropy}$, $\kappa(p) = -\sum p \log p$, get KL (p, q)
 entropy \uparrow , disorder, more uniform

Shannon, Kullback-Leibler
 1940s
 information theory

natural metric on probability manifolds



KL-divergence: $KL(p, q) = \int p \log \frac{p}{q} = \sum_i p_i \log \frac{p_i}{q_i}$

$D(p \| q) \neq D(q \| p)$ asymmetric

abstract non-Euclidean prob. space

Empirical Distrib of a data set: $\hat{p}(x) = \frac{1}{N} \sum_{n=1}^N \delta(x, x_n)$

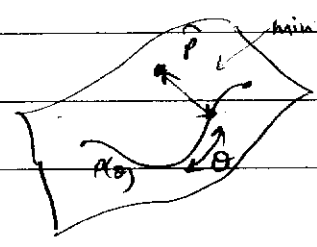
not a very good model, ∞ if we observe a point, 0 otherwise
 no generalization power.

$$\begin{aligned} D(\hat{p}(x) \| p(x|\theta)) &= \sum_x \hat{p}(x) \log \frac{\hat{p}(x)}{p(x|\theta)} = \sum_x \hat{p}(x) \log \hat{p}(x) - \sum_x \hat{p}(x) \log p(x|\theta) \\ &= \sum_x \hat{p}(x) \log \hat{p}(x) - \sum_x \frac{1}{N} \sum_{n=1}^N \delta(x, x_n) \log p(x|\theta) \\ &= \sum_x \hat{p}(x) \log \hat{p}(x) - \frac{1}{N} \sum_{n=1}^N \log p(x_n|\theta) \end{aligned}$$

constant (neg. entropy) $\leftarrow \frac{1}{N} \times \log\text{-likelihood}$

$\min D(\hat{p} \| p(x|\theta)) \equiv \text{max. likelihood}$

min. KL, geometric picture.



Information Geometry of EM (Amari)

consider single data point (easy to extend to more)

$$l(\theta) = \log p(x|\theta) \geq \alpha(q, \theta) \text{ from before but i.e.} \dots$$

$$\min_{\theta} D(\tilde{p}(x) \parallel p(x|\theta)) = \sum_x \tilde{p}(x) \log \tilde{p}(x) - \sum_x \tilde{p}(x) \log p(x|\theta)$$

"constant" "l(\theta)"

$$= \sum_x \tilde{p}(x) \log \tilde{p}(x) - \sum_x \tilde{p}(x) l(\theta)$$

$$\leq \sum_x \tilde{p}(x) \log \tilde{p}(x) - \sum_x \tilde{p}(x) \alpha(q, \theta)$$

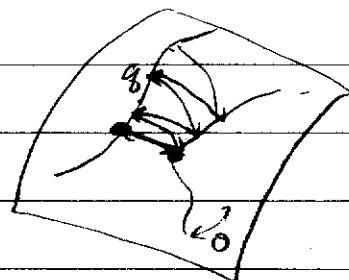
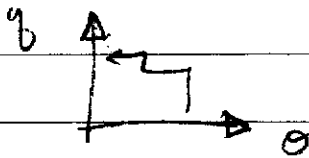
$$\leq \sum_x \tilde{p}(x) \log \tilde{p}(x) - \sum_x \tilde{p}(x) \sum_z q(z|x) \log \frac{p(x,z|\theta)}{q(z|x)}$$

$$\leq \sum_x \tilde{p}(x) \sum_z q(z|x) \left(\log \frac{\tilde{p}(x) q(z|x)}{p(x,z|\theta)} \right)$$

$$D(\tilde{p}(x) \parallel p(x|\theta)) \leq D(\tilde{p}(x)q(z|x) \parallel p(x,z|\theta))$$

incomplete divergence

"complete" divergence



iterated projections (i.e. min KL's)

e-step $q^{t+1}(z|x) = \operatorname{argmin}_q D(\tilde{p}q \parallel p(x,z|\theta^t))$

m-step $\theta^{t+1} = \operatorname{argmin}_{\theta} D(\tilde{p}q^t \parallel p(x,z|\theta))$

little em yields slightly different results than EM