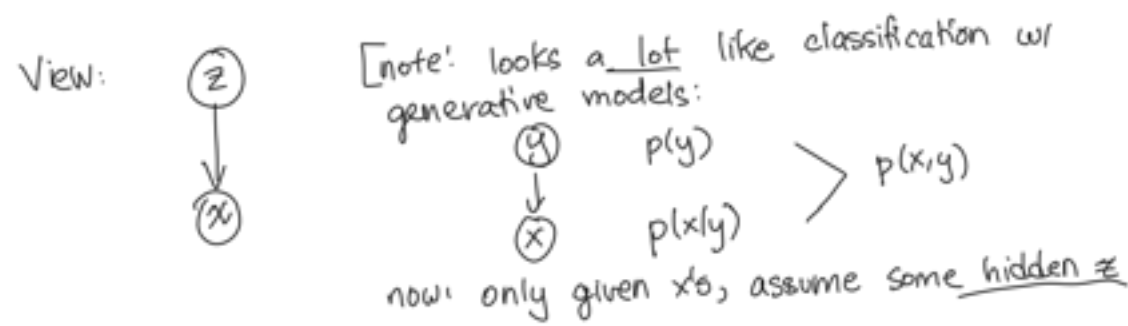


Last time: clustering: intuitive.

- how do you measure a "good" clustering?
- how to handle: unclear points?
missing data / missing dims? $d(x, x')$

Today: probabilistic clustering \Rightarrow mixture models



Gen. Process:

$$z_n \sim \pi$$

$$x_n \sim p(x_n | z_n, w)$$

$$p(x_n, z_n) \rightarrow \underline{p(x)} = \int_z p(x, z) dz$$

Notes:

1) This is density estimation

2) We can predict which mixture is likely given x, w

$$p(z|x, w) \propto p(x|z, w)p(z)$$

3) e.g. masses of car



"source separation"

Let's make this specific: Gaussian mixture model:

$$p(z_n = k) = \pi_k$$

$$p(x | z = k) = \mathcal{N}(\mu_k, \Sigma_k)$$

Complete data log likelihood: $p(x, z | w)$

$$\log \mathcal{L}(w) = \sum_n \log p(x_n, z_n | w)$$

$$= \sum_n \log p(x_n | z_n, w) + \log p(z_n | w)$$

$$= \sum_n \log \mathcal{N}(x_n | \mu_{z_n}, \Sigma_{z_n}) + \log \pi_{z_n}$$

$$= \sum_n \sum_k z_{nk} \log \mathcal{N}(x_n | \mu_k, \Sigma_k) + z_{nk} \log \pi_k$$

We can solve for the MLE analytically:

set mixture μ, Σ to

$$\pi_k = \frac{N_k}{N} \quad \mu_k = \text{mean}(x_n) \quad \Sigma_k = \text{var}(x_n) \quad \text{empirical values.}$$

Problem: we don't have z ! [aside: search over the z 's is combinatorial]

• Let's write down $p(x|W)$

$$\begin{aligned} \mathcal{L}_{\text{marg}}(W) &= \sum_n \log p(x_n|W) \\ &= \sum_n \log \left[\sum_k \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) \right] \end{aligned} \quad \left. \vphantom{\sum_n} \right\} \text{model p(x)}$$

no analytic solution
gradients are messy...

• Can we bound the marginal? [New idea]

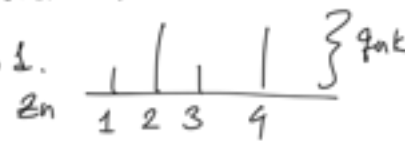
↳ lower bound.

Turns out, the expected complete data log likelihood is a lower bound for the marginal likelihood $\max_z \int p(x, z) dz \rightarrow$ we subbed in π_k 's.

$$\begin{aligned} \mathcal{L}_{\text{marg}} &= \log \int p(x, z) dz \\ E[\mathcal{L}_{\text{complete}}] &= \int \log p(x, z) dz \end{aligned}$$

$$\begin{aligned} E_z[\mathcal{L}_{\text{complete}}(W)] &= E_z[\log p(z|x_n, W)] \\ &= E_z[\sum_{n,k} z_{nk} \log \pi_k + z_{nk} \log \mathcal{N}(x_n | \mu_k, \Sigma_k)] \end{aligned}$$

Let q_{nk} be ~~the~~ a distribution over z_n
 q_n has K elements, sum to 1.



- Apply block coordinate ascent $[q] \cdot [\pi, \mu, \Sigma]$
- This is nice for optimizing w.r.t. $\pi, \mu, \Sigma \rightarrow$ corresponds to a weighted ~~feature~~ cluster assignment.

$$\pi_k = \frac{\sum_n q_{nk}}{N}$$

$$\mu_k : \frac{\partial \mathcal{L}}{\partial \mu_k} = q_{nk} \left(-\frac{1}{2} \right) (2 \Sigma_k^{-1} \mu_k - 2 \Sigma_k^{-1} x_n)$$

$$\hookrightarrow \text{gives us } \hat{\mu}_k = \frac{\sum_n q_{nk} x_n}{\sum_n q_{nk}}$$

$$\Sigma_k = \frac{1}{\sum_n q_{nk}} \sum_n q_{nk} (x_n - \hat{\mu}_k) (x_n - \hat{\mu}_k)^T$$

↳ nice! [often the case w/ exponential families]

Name: M-step [Maximizing expected complete D.L. w.r.t. global params].

• 2nd part of the block update: q_{nk}

best choice of $q_{nk} = p(z_n | x_n, W)$
 $\propto \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$

E-step: Expectations over z [locals].

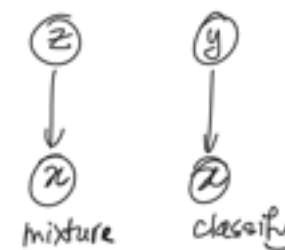
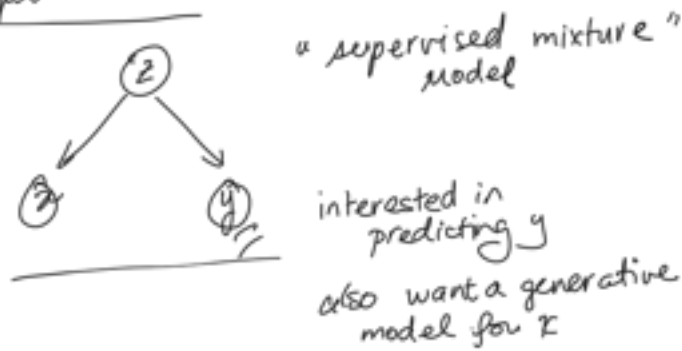


Back to properties of EM:

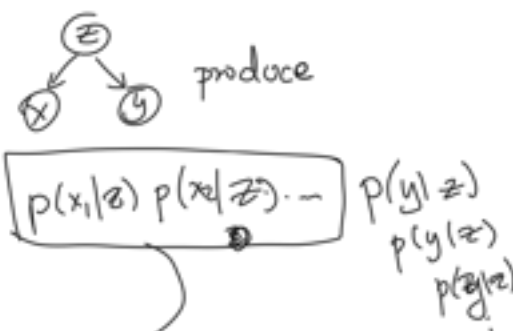
- general: E-step: local hidden
M-step: opt. global
- \Rightarrow monotonically improves lower bound [coding check]
- converges to local optima (random restarts)
- somewhat sensitive to init: $\pi_k = 1/K$ & all μ_k, Σ_k equal (saddle)
- variations MoG: Σ_k be full, diagonal, isotropic... shared Σ_k 's. (versatile)
- MoG w/ small variance \rightarrow K-means.



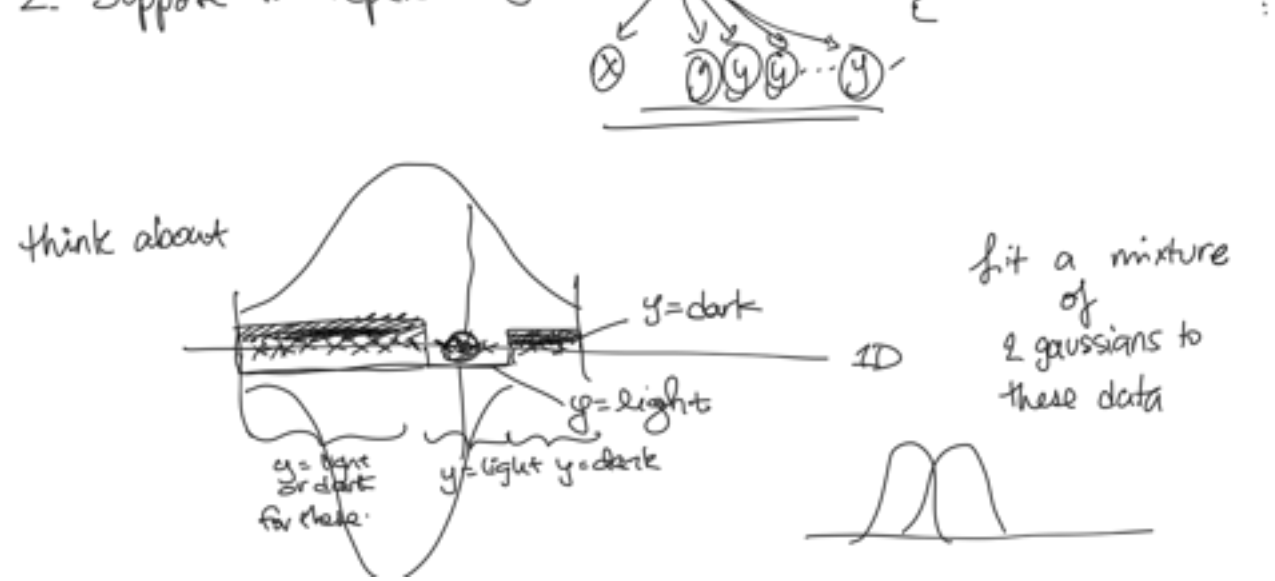
Concept Check:



1. Suppose $D_x \gg D_y$. Will mixture model different $\{\mu_k, \Sigma_k\}$ than $\{x\}$?



2. Suppose we replicate y :



Mixture Models