# EXPECTATION MAXIMIZATION (EM) ①

① EM for GMM Training/Fitting
  k-Gaussians

Given: $D = \{\vec{x}_1, \ldots, \vec{x}_N\}$

Model: $Z \sim Cat(\vec{\pi})$
  $X \mid Z = c \sim N(\vec{\mu}_c, \Sigma_c)$

Goal: Estimate parameters:

$\Theta$ (parameters) $\begin{cases} \rightarrow \text{Means}: & \vec{\mu}_1, \vec{\mu}_2 - -, \vec{\mu}_k \quad \in \mathbb{R}^d \\ \rightarrow \text{Cov.}: & \Sigma_1 - - - \Sigma_k \quad \in \mathbb{R}^{d \times d} \quad \succeq 0 \; [P.S.D] \\ \rightarrow \text{Priors}: & \pi_1 - - - \pi_k \quad \geq 0 \quad \sum_{c=1}^k \pi_c = 1 \end{cases}$

Estimator: Maximum (Marginal) Likelihood

$$\hat{\Theta}_{MLE} = \underset{\Theta}{\arg\max} \sum_{i=1}^N \log P(\vec{X} = \vec{x}_i \mid \Theta)$$

$$\sum_{i=1}^N \log \sum_{c=1}^k P(\vec{X} = \vec{x}_i, Z = c \mid \Theta)$$
$\llcorner$ Problem!

Idea #1: What if someone told us the "GT" values of Z? Then it's easy!

$$\hat{\Theta}_{MLE} = \underset{\Theta}{\text{argmax}} \sum_{i=1}^{N} \left[ \log P(\vec{X}=\vec{x}_i, \underbrace{Z=z_i}_{\uparrow} \mid \Theta) \right]$$

GT provided

$$= \underset{\Theta}{\text{argmax}} \sum_{i=1}^{N} \left[ \log P(\vec{X}=\vec{x}_i \mid Z=z_i, \Theta) + \log P(Z=z_i \mid \Theta) \right]$$

all terms are "separable"

simple
algebra $\Longrightarrow$

$$\hat{\pi}_c = \frac{\# \text{ data-pts-assigned-to-Gaussian-in GT}}{N}$$

$$= \frac{\text{count}(z_i = c)}{N}$$

$$\hat{\mu}_c = \frac{\sum_{z_i = c} \vec{x}_i}{\sum_{z_i = c} 1} \qquad \hat{\Sigma}_c = \frac{\sum_{z_i = c} (\vec{x}_i - \vec{\mu}_c)(\vec{x}_i - \vec{\mu}_c)^T}{\sum_{z_i = c} 1}$$

Idea #2: Even if some one told us "soft-GT" or soft assignment of pt to Gaussian, we would be fine.

Why? We ~~do~~ could do "soft" versions of above. (See next idea)

Idea #3 (Final EM alg for GMMs): Alternating Minimization w/ soft-assignment

$\longrightarrow$ Initialize $\Theta^{(0)} = \{ \pi_c^{(0)}, \vec{\mu}_c^{(0)}, \Sigma_c^{(0)} \}$

At iteration $t$

$\rightarrow$ $\underline{E_{\text{(XPECTATION)}}}$ - Step: Fix $\Theta^{(t)}$; compute soft-assignment

$$a_{ic}^{(t)} = P(Z=c \mid \vec{X}=\vec{x}_i, \Theta^{(t)})$$

$$= \frac{P(Z=c, \vec{X}=\vec{x}_i \mid \Theta^{(t)})}{P(\vec{X}=\vec{x}_i \mid \Theta^{(t)})}$$

$$= \frac{P(\vec{X}=\vec{x}_i \mid Z=c, \Theta^{(t)}) \, P(Z=c \mid \Theta^{(t)})}{\sum_k P(\vec{X}=\vec{x}_i \mid Z=k, \Theta^{(t)}) \, P(Z=k \mid \Theta^{(t)})}$$

prior $\longrightarrow$ $= \dfrac{\pi_c^{(t)} \, N(\mu_c^{(t)}, \Sigma_c^{(t)})}{\underbrace{\sum_k \pi_k^{(t)} \, N(\mu_k^{(t)}, \Sigma_k^{(t)})}_{\text{Normalization}}}$ $\longleftarrow$ Likelihood

"Just compute prior & likelihood from each gaussian at time $(t)$ & renormalize"

$\rightarrow$ $\underline{M_{\text{(AXIMIZATION)}}}$ - Step: Fix $a_{ic}^{(t)}$; Compute $\Theta^{(t+1)}$

$$\Theta^{(t+1)} = \underset{\Theta}{\text{argmax}} \sum_{i=1}^{N} \sum_{c} a_{ic}^{(t)} \, \log P(\vec{X}=\vec{x}_i, Z=c \mid \Theta)$$

$\equiv$ learn from noisely / softly annotated dataset

| Dataset part 1 [all instances labelled class 1] | | Dataset part k [all instance labelled k] | |
|---|---|---|---|
| $(\vec{x}_1, 1)$ | $a_{11}$ | $(\vec{x}_1, k)$ | $a_{1k}$ |
| $(\vec{x}_2, 1)$ | $a_{21}$ | $(\vec{x}_2, k)$ | $a_{2k}$ |
| $\vdots$ | | $\vdots$ | |
| $(x_N, 1)$ | $a_{N1}$ | $(\vec{x}_N, k)$ | $a_{Nk}$ |

weights associated w/ each training pt

So we can easily derive estimators as:

$$\pi_c^{(t+1)} = \frac{\sum\limits_{i=1}^{N} a_{ic}}{\underbrace{\sum\limits_{c=1}^{K}\sum\limits_{i=1}^{N} a_{ic}}_{N}} \qquad \left\{ \begin{array}{l} \text{what } \overset{\text{fraction}}{\cancel{\text{percentage}}} \text{ of} \\ \text{"mass" is associated} \\ \text{with Gaussian } c \end{array} \right.$$

$$\mu_c^{(t+1)} = \frac{\sum\limits_{i} a_{ic}\, \vec{x}_i}{\sum\limits_{i} a_{ic}} \qquad \left\{ \text{Weighted mean} \right.$$

$$\Sigma_c^{(t+1)} = \frac{\sum\limits_{i} a_{ic}\, (\vec{x}_i - \vec{\mu}_c)(\vec{x}_i - \vec{\mu}_c)^{\top}}{\sum\limits_{i} a_{ic}} \qquad \left\{ \begin{array}{l} \text{Weighted} \\ \text{Co-variance} \end{array} \right.$$