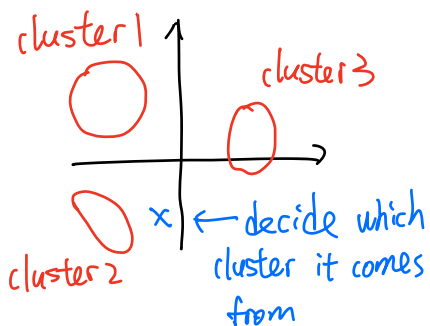


Inference problem: Given

- $x_i = x^{(i)}$

- knowing $\pi_{1:k}, \mu_{1:k}, \Sigma_{1:k}$

Compute $P(z_i | x_i = x^{(i)}; \pi_{1:k}, \mu_{1:k}, \Sigma_{1:k})$



$$P(z_i = c | x_i = x^{(i)})$$

$$= \frac{P(z_i = c) P(x_i = x^{(i)} | z_i = c)}{\sum_{b=1}^k P(z_i = b) P(x_i = x^{(i)} | z_i = b)}$$

① $P(z_i = c) = \pi_c$

② $P(x_i = x^{(i)} | z_i = c)$ Gaussian w/ mean μ_c , covariance Σ_c
conditioned on being in cluster c , what is prob of observing $x^{(i)}$

$$= \frac{1}{(2\pi)^{d/2}} \cdot \frac{1}{|\det(\Sigma_c)|} \cdot \exp\left(-\frac{1}{2} \cdot (x^{(i)} - \mu_c)^T \Sigma_c^{-1} (x^{(i)} - \mu_c)\right)$$

d : dimension of data

compare with univariable Gaussian:

$$\frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \cdot \exp\left(-\frac{1}{2} \cdot \frac{(x - \mu)^2}{\sigma^2}\right)$$

Learning GMMs:

Comparison with k-means

k-means

- Assignments ('hard assignment')
- Centroids
- alternate between update assignments / centroids

Gaussian-mixture

- Latent variables ('soft ...')
- parameters
- expectation-maximization

Expectation - Maximization (EM)

generally used when have both

- Latent variables
- unknown parameters

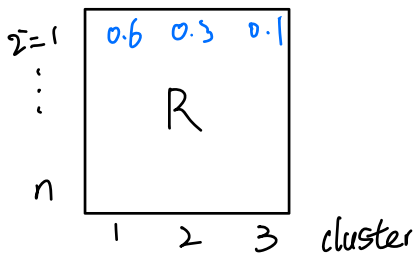
① E-step: infer latent variables distributions given current guess of params } \approx make assignments

② M-step: choose best params that fit the data based on the inferred dist. of latent variables } choose centroids

EM for GMM's

E-step For each $i=1, \dots, n$, infer dist. for z_i

call $r_{ic} = P(z_i = c | x_i = x^{(i)}; \text{current params guess})$



M-step: we have $\left\{ \begin{array}{l} \text{actual value of all the } x_i \text{'s} \\ \text{dist. for each } z_i \end{array} \right.$

\Rightarrow can't do MLE

we will maximize Expected Complete Loglikelihood (ECLL)

$$ECLL(\pi_{i:|k}, \mu_{i:|k}, \Sigma_{i:|k}) =$$

$$\sum_{i=1}^n \underbrace{\sum_{c=1}^k \Gamma_{ic}}_{\text{"expected"}} \underbrace{\log P(x_i = x^{(i)}, z_i = c; \pi, \mu, \Sigma)}_{\text{"complete" b/c compute likelihood of both } x_i \text{ and } z_i}$$

What choice of $\pi_{i:|k}, \mu_{i:|k}, \Sigma_{i:|k}$ maximize ECLL?

start with μ_i ($\mu_2 \dots \mu_k$)

Taking gradient with μ_1 , set to 0

$$\begin{aligned} \nabla_{\mu_1} ECLL &= \sum_{i=1}^n \Gamma_{i1} \nabla \log P(x_i = x^{(i)}, z_i = 1) \\ &= \sum_{i=1}^n \Gamma_{i1} \nabla \left[\cancel{\log P(z_i = 1)} + \log P(x_i = x^{(i)} | z_i = 1) \right] \\ &\quad \text{doesn't depend on } \mu_1 \\ &= \sum_{i=1}^n \Gamma_{i1} \nabla \log P(x_i = x^{(i)} | z_i = 1) \end{aligned}$$

$$\frac{1}{(2\pi)^{d/2}} \cdot \frac{1}{\det(\Sigma_c)} \cdot \exp\left(-\frac{1}{2} \cdot (x^{(i)} - \mu_c)^T \Sigma_c^{-1} (x^{(i)} - \mu_c)\right)$$

const const
w.r.t. μ_1

$$= \sum_{i=1}^n r_{i1} \nabla \left[-\frac{1}{2} \cdot (x^{(i)} - \mu_1)^T \Sigma_1^{-1} (x^{(i)} - \mu_1) \right]$$

quadratic form

$$\nabla_x x^T A x = 2Ax$$

$$= \cancel{\frac{1}{2}} \sum_{i=1}^n r_{i1} \cdot \cancel{2} \sum_1 (x^{(i)} - \mu_1) \cdot \cancel{(-1)} = 0$$

$$= \sum_{i=1}^n r_{i1} \sum_1^{-1} (x^{(i)} - \mu_1) = 0$$

$$\sum_1 (\sum_1^{-1}) \cdot \sum_{i=1}^n r_{i1} (x^{(i)} - \mu_1) = \sum_1 0$$

$$\sum_{i=1}^n r_{i1} (x^{(i)} - \mu_1) = 0$$

$$\sum_{i=1}^n r_{i1} x^{(i)} = \mu_1 \sum_{i=1}^n r_{i1}$$

$$\Rightarrow \mu_1 = \frac{\sum_{i=1}^n r_{i1} x^{(i)}}{\sum_{i=1}^n r_{i1}}$$

} weighted coverage of $x^{(i)}$,
where weights are how likely
example is in cluster 1

$$= P(\text{example 1 in cluster 1})$$

+ ...

$$+ P(\text{example } n \text{ in cluster 1})$$

$$= \mathbb{E} \text{ number of examples in cluster 1}$$

$$\bar{\Pi}_c = \frac{\sum_{i=1}^n r_{ic}}{n} \left. \vphantom{\frac{\sum_{i=1}^n r_{ic}}{n}} \right] \text{ "soft version" of counting } \frac{\# \text{ points in cluster } c}{\text{total } \# \text{ points}}$$

$$\Sigma_c = \frac{\sum_{i=1}^n r_{ic} (x^{(i)} - \mu_c)(x^{(i)} - \mu_c)^T}{\sum_{i=1}^n r_{ic}} \left. \vphantom{\frac{\sum_{i=1}^n r_{ic} (x^{(i)} - \mu_c)(x^{(i)} - \mu_c)^T}{\sum_{i=1}^n r_{ic}}} \right] \begin{array}{l} \text{expectation of} \\ (x - \mu_c)(x - \mu_c)^T \\ \text{using } r_{ic} \text{ as weights} \\ \uparrow \end{array}$$

definition of covariance