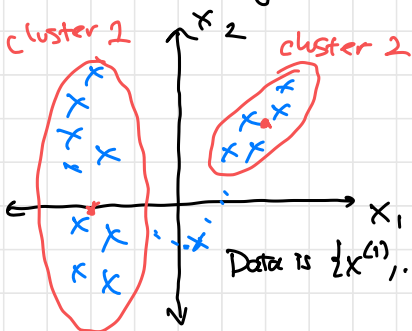


3/23/2023: Gaussian Mixture Models (GMM)

1. What is a GMM?

2. Inference — Assign datapoint to a cluster

3. Learning — Decide mean & covariance of each cluster
location *shape*



How was the data generated?

1. Randomly choose cluster 1 or 2

2. Sample from a Gaussian distribution for the chosen cluster

↑ Latent variable

Data is $\{x^{(1)}, \dots, x^{(n)}\}$

For this dataset: #clusters = $k = 2$

$$\pi_1 = 2/3, \quad \pi_2 = 1/3$$

$$\mu_1 = \begin{bmatrix} -2 \\ 0 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 0.9 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$$

— π_i = Probability of choosing cluster i in Step 1

Reasonable parameters we want to learn given this dataset

Inference: Infer probability distribution of a latent random variable
unobserved in the data

(In contrast, learning means fitting parameters)

Terminology: For each $i = 1, \dots, n$

Z_i is latent variable denoting cluster choice $\in \{1, \dots, k\}$

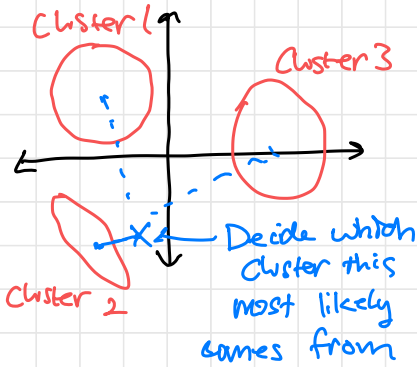
X_i is random variable we do observe as $x^{(i)}$

Inference problem: Given:

$X_i = x^{(i)}$ and

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

Compute $P(Z_i | X_i = x^{(i)}; \pi_{1:k}, N_{1:k}, \Sigma_{1:k})$



$P(Z_i = c | X_i = x^{(i)})$

$$= \frac{P(Z_i = c) P(X_i = x^{(i)} | Z_i = c)}{\sum_{b=1}^K P(Z_i = b) P(X_i = x^{(i)} | Z_i = b)}$$

① $P(Z_i = c) = \pi_c$

Gaussian w/
mean N_c , covariance Σ_c

② $P(X_i = x^{(i)} | Z_i = c)$

Conditioned on being in Cluster c , what is probability of observing $x^{(i)}$

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

d = dimension of data

Comparison:
Univariate Gaussian

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

Announcements

Learning GMMs:

Comparison with k-Means

- Assignments ("hard assignments") \approx Latent variables ("soft assignments")
- Centroids \approx parameters

k-Means: Alternate between updating assignments & updating centroids

Today: Expectation-Maximization (EM)

Generally used when you have both

- Latent variables AND
- Unknown parameters

① E-step: Infer latent variable's distributions given current guess of parameters

② M-step: Choose best parameters that fit the data based on the inferred distribution of latent variables

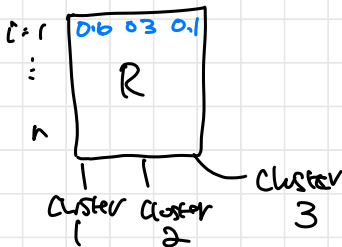
\approx make assignments based on current centroids

\approx choose new centroids based on current assignment

EM for GMMs:

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

Call $r_{ic} = P(Z_i = c \mid X_i = x^{(i)}; \text{current parameter guess})$



Comes from inference procedure above

M-step We have:

Actual value of all the X_i 's
Distribution for each Z_i

Can't directly do MLE, but can do something similar
we will maximize Expected Complete Loglikelihood (ECLL)

$$ECLL(\pi_{1:k}, N_{1:k}, \Sigma_{1:k}) = \sum_{i=1}^n \underbrace{\sum_{c=1}^k r_{ic}}_{\text{"Expected"}} \underbrace{\log P(X_i = x^{(i)}, Z_i = c; \pi, N, \Sigma)}_{\text{"Complete" b/c compute likelihood of both } X_i \text{ \& } Z_i}$$

What choices of $\pi_{1:k}, N_{1:k}, \Sigma_{1:k}$ maximize ECLL?
Start with N_1 (N_2, \dots, N_k are symmetric)

Taking gradient wrt N_1 & set it = 0

$$\begin{aligned} \nabla_{N_1} ECLL &= \sum_{i=1}^n r_{i1} \nabla \log P(X_i = x^{(i)}, Z_i = 1) \\ &= \sum_{i=1}^n r_{i1} \nabla \left[\cancel{\log P(Z_i = 1)} + \log P(X_i = x^{(i)} | Z_i = 1) \right] \\ &\quad \text{doesn't depend on } N_1 \\ &= \sum_{i=1}^n r_{i1} \nabla \log P(X_i = x^{(i)} | Z_i = 1) \end{aligned}$$

$$= \frac{1}{(2\pi)^{D/2}} \cdot \frac{1}{\sqrt{\det(\Sigma_c)}} \cdot \exp\left(-\frac{1}{2} \cdot (x^{(i)} - N_c)^T \Sigma_c^{-1} (x^{(i)} - N_c)\right)$$

Constant

Constant
wrt N_1

$$= \sum_{i=1}^n r_{ii} \nabla_{N_i} \left[-\frac{1}{2} (x^{(i)} - N_i)^T \Sigma_i^{-1} (x^{(i)} - N_i) \right]$$

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

$$= \sum_{i=1}^n r_{ii} \cdot 2 \Sigma_i^{-1} (x^{(i)} - N_i) = 0$$

$$= \sum_{i=1}^n r_{ii} \Sigma_i^{-1} (x^{(i)} - N_i) = 0$$

$$\Sigma_x \left(\sum_{i=1}^n \Sigma_i^{-1} \right) \cdot \sum_{i=1}^n r_{ii} (x^{(i)} - N_i) = 0$$

$$\sum_{i=1}^n r_{ii} (x^{(i)} - N_i) = 0 \iff$$

$$\sum_{i=1}^n r_{ii} x^{(i)} = N_i \sum_{i=1}^n r_{ii}$$

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

Weighted average of $x^{(i)}$'s where weights are how likely each example was in cluster 1

$\left. \begin{array}{l} p(\text{example 1 in cluster 1}) \\ + p(\text{example 2 in cluster 1}) \\ + \dots \\ + n \text{ in cluster 1} \end{array} \right\} = \text{Expected number of examples in cluster 1}$

$$\pi_c = \frac{\sum_{i=1}^n r_{ic}}{n} \quad \left. \vphantom{\sum_{i=1}^n r_{ic}} \right\} \text{Soft version of counting \# points in cluster } c \text{ / total \# points}$$

$$\Sigma_c = \frac{\sum_{i=1}^n r_{ic} (x^{(i)} - N_c)(x^{(i)} - N_c)^T}{\sum_{i=1}^n r_{ic}} \quad \left. \vphantom{\sum_{i=1}^n r_{ic}} \right\} \text{Expectation of } (x - N_c)(x - N_c)^T \text{ using } r_{ic} \text{ as weights}$$

↑
This is one definition of covariance