

Gaussian mixture model (GMM)

If $\Sigma > 0$ then the density for $x \sim N_p(\mu, \Sigma)$ is

$$f(x; \mu, \Sigma) = \left(\frac{1}{2\pi}\right)^{p/2} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

GMM assumes that x_1, \dots, x_n are drawn i.i.d. from density:

$$p(x) = \sum_{k=1}^K \pi_k f(x; \mu_k, \Sigma_k), \quad \text{where}$$

mixing coefficients $\pi_k \geq 0$ and $\sum_{k=1}^K \pi_k = 1$

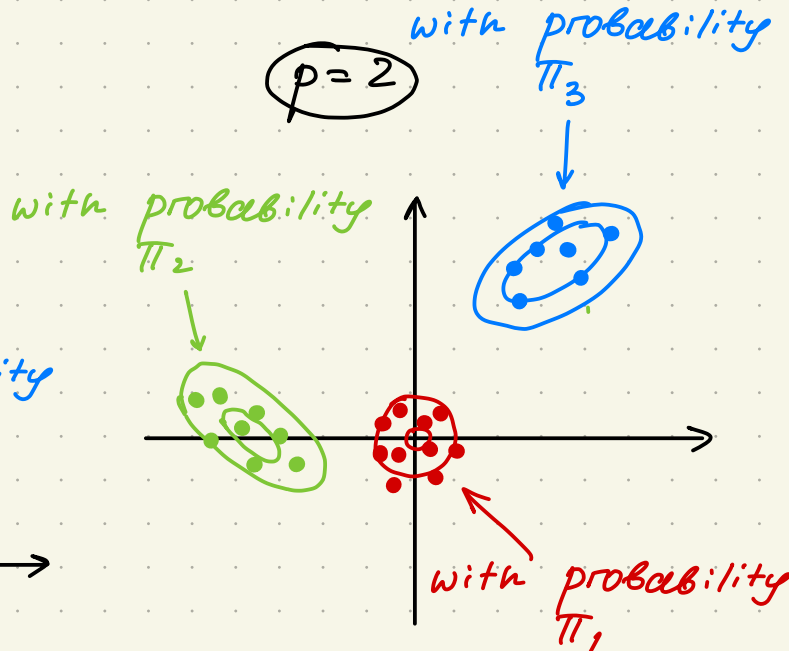
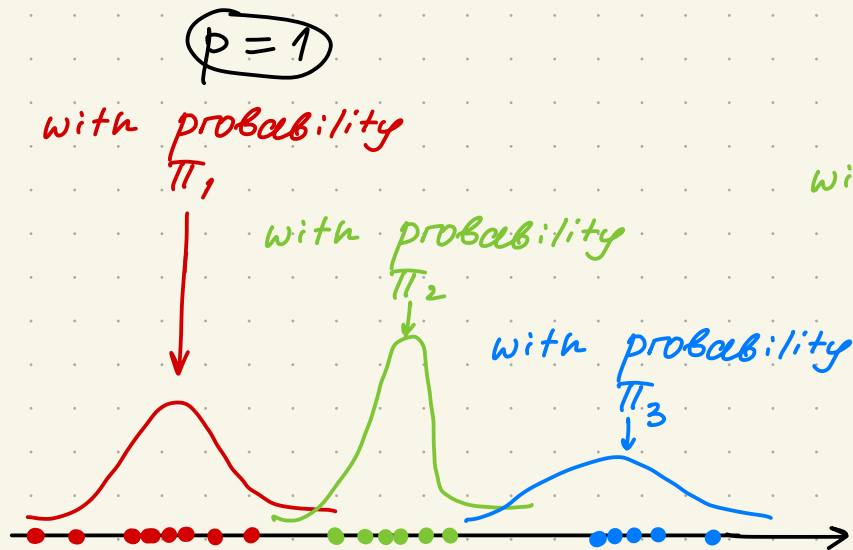
that is a **mixture of K multivariate Gaussian distributions.**

Unknown parameters: $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

generative model for GMM:

$$z_i = \begin{cases} 1 & \text{with probability } \pi_1 \\ \vdots & \vdots \\ k & \text{with probability } \pi_k \end{cases} \quad \text{latent variable}$$

$$x_i | z_i \sim \mathcal{N}(\mu_{z_i}, \Sigma_{z_i})$$



Maximize log-likelihood :

$$\sum_{i=1}^n \log p(x_i) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k f(x_i; \mu_k, \Sigma_k) \right)$$

• If $K=1$ the solution is easy :

$$\pi_1 = 1 \quad \mu_1 = \frac{1}{n} \sum_{i=1}^n x_i \quad \Sigma_1 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

• In general, $\log \left(\sum_{k=1}^K \dots \right)$ is the problem.

① If we know $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

we could compute $p(z_i = k | x_i)$

$$p(z_i = k | x_i) = \frac{p(z_i = k) \cdot p(x_i | z_i = k)}{p(x_i)} =$$

$$= \frac{\pi_k f(x_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j f(x_i; \mu_j, \Sigma_j)}$$

x_i belongs to $\begin{cases} \text{cluster 1 with } p(z_i = 1 | x_i) \\ \dots \\ \text{cluster } K \text{ with } p(z_i = K | x_i) \end{cases}$

We can pick $z_i = \underset{k=1 \dots K}{\operatorname{argmax}} p(z_i = k | x_i)$

② If we knew z_i then

- $\pi_k = \frac{1}{n} \sum_{i=1}^n I(z_i = k)$ - proportion of observations in cluster k

- $\mu_k = \frac{\sum_{i=1}^n I(z_i = k) x_i}{\sum_{i=1}^n I(z_i = k)}$ - sample mean of cluster k

- $\Sigma_k = \frac{\sum_{i=1}^n I(z_i = k) (x_i - \mu_k) (x_i - \mu_k)^T}{\sum_{i=1}^n I(z_i = k)}$ - sample co-variance of cluster k

$$\sum_{i=1}^n \log p(x_i, z_i) = \sum_{i=1}^n [\log p(x_i | z_i) + \log p(z_i)]$$

$$= \sum_{i=1}^n [\log f(x_i; \mu_{z_i}, \Sigma_{z_i}) + \log \pi_{z_i}] =$$

$$= \sum_{i=1}^n \sum_{k=1}^K I(z_i = k) (\log f(x_i; \mu_k, \Sigma_k) + \log \pi_k)$$

Estimating π_k :

$$\sum_{k=1}^K \log \pi_k \cdot \underbrace{\sum_{i=1}^n I(z_i = k)}_{n_k} = \sum_{k=1}^K n_k \log \pi_k \rightarrow \max$$

(*)

$$\sum_{k=2}^K n_k \log \pi_k + n_1 \log \left(1 - \sum_{k=2}^K \pi_k\right)$$

$$\nabla_{\pi_k} = \frac{n_k}{\pi_k} - \frac{n_1}{1 - \sum_{k=2}^K \pi_k} = 0 \Rightarrow n_k \underbrace{\left(1 - \sum_{k=2}^K \pi_k\right)}_{\pi_1} = n_1 \pi_k$$

$$n_k \pi_1 = n_1 \pi_k \Rightarrow \pi_k = n_k \cdot d$$

$$\pi_1 + \dots + \pi_K = 1 \Rightarrow d = \frac{1}{n_1 + \dots + n_K} = \frac{1}{n} \Rightarrow \pi_k = \frac{n_k}{n}$$

Estimating μ_k, Σ_k :

$$\sum_{i=1}^n I(z_i = k) \log f(x_i; \mu_k, \Sigma_k) = \sum_{i \in C_k} \log f(x_i; \mu_k, \Sigma_k)$$

$$\text{Thus } \mu_k = \frac{1}{n_k} \sum_{i \in C_k} x_i$$

$$\Sigma_k = \frac{1}{n_k} \sum_{i \in C_k} (x_i - \mu_k)(x_i - \mu_k)^T$$

① Given $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$, compute

$$p(z_i = k | x_i) = \frac{\pi_k f(x_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j f(x_i; \mu_j, \Sigma_j)}$$

$z_i = \underset{k=1 \dots K}{\operatorname{argmax}} p(z_i = k | x_i)$, $w_{ik} = I(z_i = k)$
 each x_i is assigned to cluster $\underset{\uparrow}{1} \dots \underset{\uparrow}{k}$
 $w_{i1} \dots w_{ik}$

② Given z_i compute

- $\pi_k = \frac{1}{n} \sum_{i=1}^n I(z_i = k) = \frac{\sum_{i=1}^n w_{ik}}{\sum_{i=1}^n \sum_{k=1}^K w_{ik}}$

- $\mu_k = \frac{\sum_{i=1}^n I(z_i = k) x_i}{\sum_{i=1}^n I(z_i = k)} = \frac{\sum_{i=1}^n w_{ik} x_i}{\sum_{i=1}^n w_{ik}}$

- $\Sigma_k = \frac{\sum_{i=1}^n I(z_i = k) (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^n I(z_i = k)} = \frac{\sum_{i=1}^n w_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^n w_{ik}}$

EM algorithm

Expectation (E) step: Given $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

$$w_{ik} = p(z_i = k \mid x_i) = \frac{\pi_k f(x_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j f(x_i; \mu_j, \Sigma_j)}$$

each x_i is *partially* assigned to cluster $1 \dots k$
 $\uparrow \quad \quad \quad \uparrow$
 $w_{i1} \dots w_{ik}$

Maximization (M) step: Given w_{ik}

$$\bullet \pi_k = \frac{\sum_{i=1}^n w_{ik}}{\sum_{i=1}^n \sum_{k=1}^K w_{ik}}$$

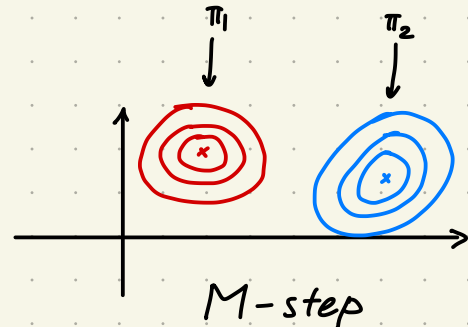
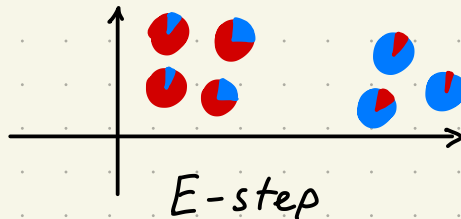
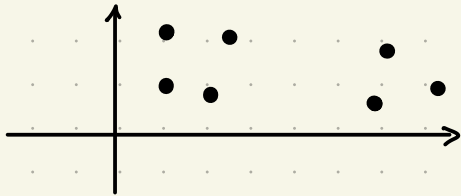
$$\bullet \mu_k = \frac{\sum_{i=1}^n w_{ik} x_i}{\sum_{i=1}^n w_{ik}}$$

$$\bullet \Sigma_k = \frac{\sum_{i=1}^n w_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^n w_{ik}}$$

E-step: which of $N_p(\mu_1, \Sigma_1) \dots N_p(\mu_k, \Sigma_k)$
generated each x_i ?

We are not sure, so we assign probabilities
 x_i came from $N(\mu_1, \Sigma_1)$ with probability w_{i1}
 \vdots
 $N(\mu_k, \Sigma_k)$ with probability w_{ik}

M-step: we want to estimate μ_1, \dots, μ_k and $\Sigma_1, \dots, \Sigma_k$
For μ_k and Σ_k each observation x_i will have
weight w_{ik} . We use weighted sample
mean and sample variance.



EM algorithm : motivation

Denote by Θ the set of parameters and
by x the set of observations.

$$\ell(x; \theta) = \sum_{i=1}^n \log p(x_i; \theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K p(x_i, z_i = k; \theta) \right)$$

• $\log p(\overset{\downarrow \text{any observation}}{x}; \theta) \geq \sum_{k=1}^K w_k \cdot \log \left(\frac{p(x, z = k; \theta)}{w_k} \right)$

for w_1, \dots, w_K such that $w_1 + \dots + w_K = 1$ and $w_k \geq 0$

| \log is concave function thus

$$\log \left(\sum_{k=1}^K a_k \right) = \log \left(\sum_{k=1}^K w_k \cdot \frac{a_k}{w_k} \right) \geq \sum_{k=1}^K w_k \log \left(\frac{a_k}{w_k} \right)$$

- If $w_k = p(z=k|x; \theta)$ then inequality becomes equality

$$\frac{p(x, z=k; \theta)}{w_k} = \frac{p(x, z=k; \theta)}{p(z=k|x; \theta)} = p(x; \theta)$$

$$\sum_{k=1}^K w_k \cdot \log\left(\frac{p(x, z=k; \theta)}{w_k}\right) = \sum_{k=1}^K w_k \cdot p(x; \theta) = p(x; \theta)$$

- Denote weights for observation x_i by $w_{i1} \dots w_{iK}$

E-step: at iteration t compute

$$w_{ik}^{(t)} = p(z_i=k|x_i; \theta^{(t-1)})$$

for each i , $w_{ik}^{(t)} \geq 0$ and $\sum_{k=1}^K w_{ik}^{(t)} = 1$

$$\ell(x; \theta) = \sum_{i=1}^n \log\left(\sum_{k=1}^K p(x_i, z_i=k; \theta)\right) \geq$$

$$\sum_{i=1}^n \sum_{k=1}^K w_{ik}^{(t)} \log\left(\frac{p(x_i, z_i=k; \theta)}{w_{ik}^{(t)}}\right)$$

- M-step: at iteration t maximize the lower-bound for $\ell(x; \theta)$

$$\sum_{i=1}^n \sum_{k=1}^K w_{ik}^{(t)} \log \left(\frac{p(x_i, z_i=k; \theta)}{w_{ik}^{(t)}} \right) =$$

$$\sum_{i=1}^n \sum_{k=1}^K w_{ik}^{(t)} [\log(p(x_i, z_i=k; \theta)) - \log w_{ik}^{(t)}] = Q(\theta) + \dots$$

↙ ignore

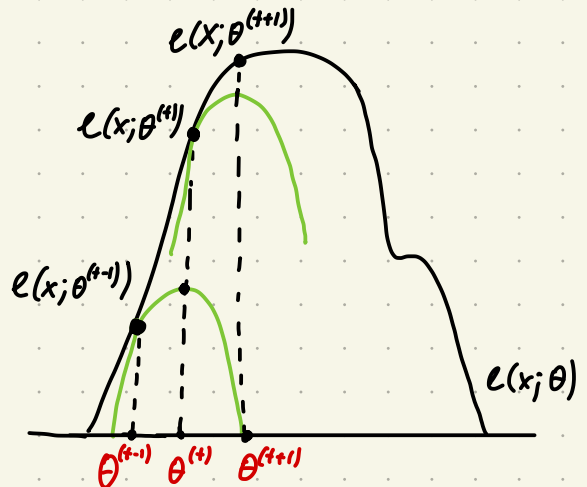
M-step: maximize $Q(\theta)$

- $\ell(x; \theta^{(t-1)}) \leq \ell(x; \theta^{(t)})$

$$\ell(x; \theta^{(t-1)}) = Q(\theta^{(t-1)}) + \dots \leq$$

$$\leq Q(\theta^{(t)}) + \dots = \ell(x; \theta^{(t)})$$

- log-likelihood converges to local maximum.



EM for GMM

E-step Given $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

$$w_{ik} = p(z_i = k | x_i) = \frac{\pi_k f(x_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j f(x_i; \mu_j, \Sigma_j)}$$

M-step maximize $Q(\theta) = \sum_{i=1}^n \sum_{k=1}^K w_{ik} \log(p(x_i, z_i = k; \theta)) =$
 $= \sum_{i=1}^n \sum_{k=1}^K w_{ik} (\log f(x_i; \mu_k, \Sigma_k) + \log \pi_k)$

Estimating π_k :

$$\sum_{i=1}^n \sum_{k=1}^K w_{ik} \log \pi_k = \sum_{i=1}^n \left(\sum_{k=2}^K w_{ik} \log \pi_k + w_{i1} \log \left(1 - \sum_{k=2}^K \pi_k \right) \right)$$

$$\nabla_{\pi_k} = \sum_{i=1}^n \left[w_{ik} / \pi_k - w_{i1} / \left(1 - \sum_{k=2}^K \pi_k \right) \right] = \frac{\sum_{i=1}^n w_{ik}}{\pi_k} - \frac{\sum_{i=1}^n w_{i1}}{\pi_1} = 0 \Rightarrow \pi_k = d \sum_{i=1}^n w_{ik}$$

as $\sum_{k=1}^K \pi_k = 1$ then $d = \frac{1}{\sum_{i=1}^n \sum_{k=1}^K w_{ik}} \Rightarrow \pi_k = \frac{\sum_{i=1}^n w_{ik}}{\sum_{i=1}^n \sum_{k=1}^K w_{ik}}$

Estimating μ_k :

$$\sum_{i=1}^n \sum_{k=1}^K w_{ik} (x_i - \mu_k)^T \sum_k (x_i - \mu_k)$$

$$\nabla_{\mu_k} = -2 \sum_{i=1}^n w_{ik} \sum_k (x_i - \mu_k) = -2 \sum_k \sum_{i=1}^n w_{ik} (x_i - \mu_k) = 0$$

$$\sum_{i=1}^n w_{ik} x_i = \mu_k \cdot \sum_{i=1}^n w_{ik} \Rightarrow \mu_k = \frac{\sum_{i=1}^n w_{ik} x_i}{\sum_{i=1}^n w_{ik}}$$