# unsupervised learning

- → Rep. learning — PCA / kernel PCA
- → Clustering — Lloyd's / k-means
- → Estimation
  - → Max-likelihood
  - → Bayesian Modeling.

Today : Slightly complicated data    Estimation for.



$x_7 \; x_3 \; x_4$     $x_6 \; x_1 \; x_5$     $x_8 \; x_2 \; x_9$

- What could be a good "generative" story?



$N\left(\hat{\mu}_{ML}, \hat{\sigma}^2\right)$

$\hat{\mu}_{ML}$



Want a density like above to explain this data.

A NEW GENERATIVE MODEL

## MIXTURE OF GAUSSIANS

**STEP 1 :** Pick which mixture a data point comes from.
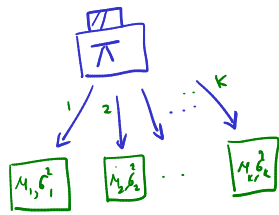
**STEP 2 :** Generate data point from that mixture.

---

**STEP 1 :** Generate a mixture component among $\{1, \cdots, k\}$    $z_i \in \{1, \cdots, k\}$

$$P\left(z_i = \ell\right) = \pi_\ell \qquad \left[\begin{array}{l} \sum_{i=1}^{K} \pi_i = 1 \\ 0 \leq \pi_i \leq 1 \quad +i \end{array}\right]$$

**STEP 2 :** Generate $x_i \sim N\left(\mu_{z_i}, \sigma^2_{z_i}\right)$

---



$\{x_1, \cdots, x_n\} \rightarrow$ OBSERVED

$\{z_1, \cdots, z_n\} \rightarrow$ UNOBSERVED/ LATENT.

Latent variable Models

**Parameters :** $\pi = [\pi_1 \; \pi_2 \cdots \pi_k]$     **Total :** $2K + K-1$

$+ \ell \left(\mu_\ell, \sigma^2_\ell\right)$             $3K - 1$

---

Max. Likelihood for GMM

$$L\left(\begin{array}{l} \mu_1, \cdots, \mu_K \\ \sigma^2_1, \cdots, \sigma^2_K, \\ \pi_1, \cdots, \pi_K \end{array} \; x_1, \cdots, x_n\right) = \prod_{i=1}^{n} f_{mix}\left(x_i \; ; \; \begin{array}{l} \mu_1, \cdots, \mu_K \\ \sigma^2_1, \cdots, \sigma^2_K \\ \pi_1, \cdots, \pi_K \end{array}\right)$$

$$= \prod_{i=1}^{n} \left[\sum_{\ell=1}^{K} \pi_\ell \cdot f\left(x_i ; \mu_\ell, \sigma^2_\ell\right)\right]$$

$\longrightarrow$ NORMAL/Gaussian Density.

$$X_1 = -15$$

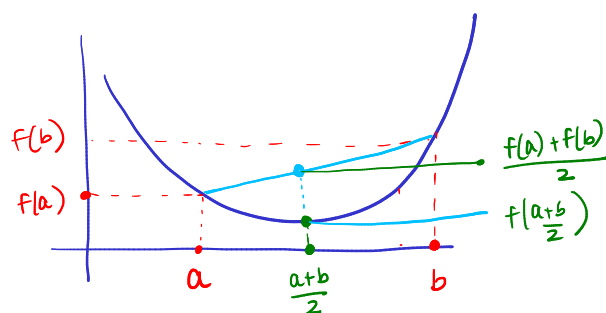$$\boxed{\pi_1 = 0.05} \qquad \pi_2 = 0.05 \qquad \boxed{\pi_3 = 0.9}$$

---

$$L(\theta) = \prod_{i=1}^{n} \left[ \sum_{k=1}^{K} \pi_k \frac{e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\sqrt{2\pi}\,\sigma_k} \right]$$

$\uparrow$
all parameters

$$\log L(\theta) = \sum_{i=1}^{n} \log \left( \sum_{k=1}^{K} \pi_k \frac{e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\sqrt{2\pi}\,\sigma_k} \right) \quad - \quad \circledast$$
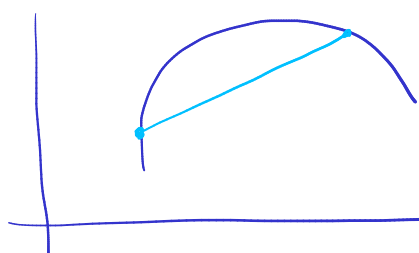
- Not possible to solve this analytically.

- Need an alternate way to solve this efficiently!

---

### Quick detour — Convex functions



$f(b)$

$f(a)$

$a \qquad \frac{a+b}{2} \qquad b$

$\frac{f(a) + f(b)}{2}$

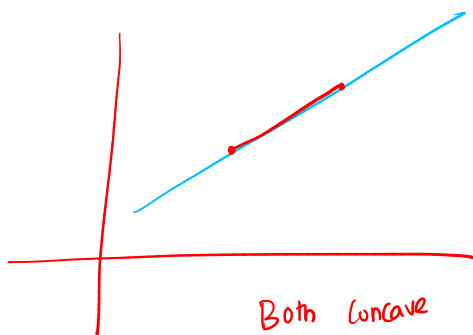$f\left(\frac{a+b}{2}\right)$

$\forall\, a, b$

$$f\left(\frac{a+b}{2}\right) \leq \frac{f(a) + f(b)}{2}$$
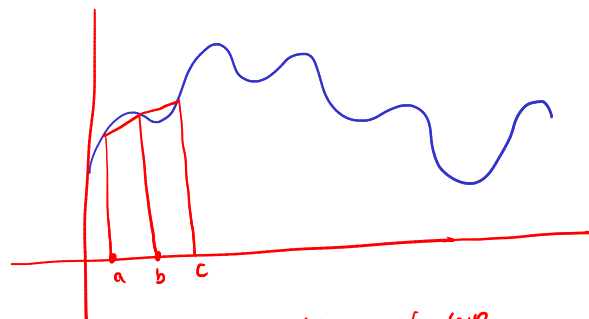
CONVEX FUNCTION

---



$\forall\, a, b$

$$f\left(\frac{a+b}{2}\right) \geq \frac{f(a) + f(b)}{2}$$
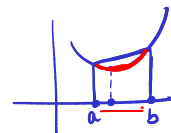
CONCAVE FUNCTION

Both Concave and Convex.

Neither Concave nor Convex.

$$f\left(\frac{1}{2}\cdot a + \frac{1}{2}b\right) \le \frac{1}{2}f(a) + \frac{1}{2}f(b)$$



$$\Rightarrow \quad f\left(\lambda a + (1-\lambda)b\right) \le \lambda f(a) + (1-\lambda)f(b) \qquad \underline{\lambda \in [0,1]}$$

<u>For Concave</u>

$$f\left(\lambda_1 a_1 + \lambda_2 a_2 + \cdots + \lambda_k a_k\right) \ge \lambda_1 f(a_1) + \cdots + \lambda_k f(a_k)$$

$$\left[\begin{array}{l}\sum_{i=1}^{k}\lambda_i = 1 \\ 0 \le \lambda_i \le 1\end{array}\right]$$



JENSEN'S INEQUALITY.

$$\boxed{f\left(\sum_{k=1}^{k}\lambda_k a_k\right) \ge \sum_{k=1}^{k}\lambda_k f(a_k)}$$

- Log is a concave function!  [ Why? Exercise ]

- How can we exploit <u>Jensen's</u> for performing maximum likelihood.

Recall

$$\circledast \quad \log L(\theta) = \sum_{i=1}^{n}\log\left(\sum_{k=1}^{k}\left(\pi_k \, e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \cdot \frac{1}{\sqrt{2\pi}\,\sigma_k}\right)\right)$$

- INTRODUCE for every data point $i$, two parameters

$$\{\lambda_1^i, \ldots, \lambda_k^i\} \quad s.t \quad \forall i \sum_{k=1}^{k} \lambda_k^i = 1, \quad 0 \le \lambda_k^i \le 1 \; \forall k$$

$$\log L(\theta) = \sum_{i=1}^{n} \log \left( \sum_{k=1}^{k} \lambda_k^i \left( \frac{\pi_k \, e^{\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}} \frac{1}{\sqrt{2\pi}\sigma_k}}{\lambda_k^i} \right) \right)$$

By Jensen's

$$\log L(\theta) \ge \text{modified\_} \log L(\theta, \lambda)$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{k} \lambda_k^i \log \left( \frac{\pi_k \, e^{\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}} \frac{1}{\sqrt{2\pi}\sigma_k}}{\lambda_k^i} \right)$$

- Note that the above modified log likelihood gives a
  lower bound for the true log likelihood at $\theta$
  for any choice of $\lambda$ $\searrow$

$$\begin{Bmatrix} \lambda_1^1, \ldots, \lambda_k^1 \\ \lambda_1^2, \ldots, \lambda_k^2 \\ \vdots \\ \lambda_1^n, \ldots, \lambda_k^n \end{Bmatrix}$$

$$\begin{Bmatrix} \mu_1, \ldots, \mu_k \\ \sigma_1^2, \ldots, \sigma_k^2 \\ \pi_1, \ldots, \pi_k \end{Bmatrix}$$

- But what are we gaining?

Key insight:

- If we fix $\lambda$, it is easy to maximize w.r.t $\theta$

- If we fix $\theta$, it is easy to maximize w.r.t $\lambda$.

Fix $\lambda$ and maximize over $\theta$

$$\max_{\theta} \sum_{i=1}^{n} \sum_{k=1}^{k} \lambda_k^i \left[ \log \left( \pi_k \, e^{-(x_i - \mu_k)^2/2\sigma_k^2} \frac{1}{\sqrt{2\pi}\sigma_k} \right) / \lambda_k^i \right]$$

$$= \max_{\theta} \sum_{i=1}^{n} \sum_{k=1}^{K} \left[ \lambda_k^i \log \pi_k - \lambda_k^i \frac{(x_i - \mu_k)^2}{2\sigma_k^2} - \lambda_k^i \log \sqrt{2\pi}\, \sigma_k \right]$$

Take derivative w.r.t $\mu, \sigma$ to get

$$\hat{\mu}_k^{MML} = \frac{\sum_{i=1}^{n} \lambda_k^i x_i}{\sum_{i=1}^{n} \lambda_k^i} \qquad \hat{\sigma}_k^{2\, MML} = \frac{\sum_{i=1}^{n} \lambda_k^i (x_i - \hat{\mu}_k^{MML})^2}{\sum_{i=1}^{n} \lambda_k^i}$$

$$\max_{\pi_1, \cdots, \pi_k} \sum_{i=1}^{n} \left( \sum_{k=1}^{K} \lambda_k^i \log \pi_k \right)$$

$$s.t \qquad \sum_{k} \pi_k = 1 \; ; \; \pi_k \geq 0$$

Can solve using method of Lagrange multipliers

$$\hat{\pi}_k^{MML} = \frac{\sum_{i=1}^{n} \lambda_k^i}{n}$$

Fixing $\lambda$, we get

$$\hat{\mu}_k^{MML} = \frac{\sum_{i=1}^{n} \lambda_k^i x_i}{\sum_{i=1}^{n} \lambda_k^i} \qquad\qquad \hat{\sigma}_k^{2\, MML} = \frac{\sum_{i=1}^{n} \lambda_k^i (x_i - \hat{\mu}_k^{MML})^2}{\sum_{i=1}^{n} \lambda_k^i}$$

$$\hat{\pi}_k^{MML} = \frac{\sum_{i=1}^{n} \lambda_k^i}{n}$$

- Fix $\theta$ and maximize $\lambda$

$$\sum_{i=1}^{n} \sum_{k=1}^{K} \lambda_k^i \log \left( \frac{\pi_k \, e^{\frac{-(x_i - \mu_k)^2}{2\sigma_k^2}} \frac{1}{\sqrt{2\pi}\, \sigma_k}}{\lambda_k^i} \right)$$

$$= \sum_{i=1}^{n} \left[ \sum_{k=1}^{K} \lambda_k^i \log(a_{ik}) - \lambda_k^i \log \lambda_k^i \right]$$

where $a_{ik} = \pi_k \, e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \frac{1}{\sqrt{2\pi}\sigma_k}$

Fix any $i$,

$$\max_{\lambda_1^i, \dots \lambda_k^i} \sum_{k=1}^{K} \left[ \lambda_k^i \log(a_k^i) - \lambda_k^i \log \lambda_k^i \right]$$

s.t $\sum_{k=1}^{K} \lambda_k^i = 1 \qquad 0 \leq \lambda_k^i \leq 1$

Can be solved analytically

$P(x_i | z_i = k)$

$P(z_i = k)$

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

$$\hat{\lambda}_k^{i \ MML} = \frac{\left( \frac{1}{\sqrt{2\pi}\sigma_k} \, e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \right) \cdot \pi_k}{\sum_{l=1}^{K} \left( \frac{1}{\sqrt{2\pi}\sigma_l} \, e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \cdot \pi_k \right)} \leftarrow P(x_i)$$

---

## ALGORITHM — E·M ALGORITHM (1970's Dempster et al)

iteration

→ Initialize $\theta^0 = \left\{ \begin{array}{l} \mathring{\mu}_1, \dots, \mathring{\mu}_K, \\ \mathring{\sigma}_1^2, \dots \mathring{\sigma}_L^2, \\ \mathring{\pi}_1, \dots \mathring{\pi}_K \end{array} \right\}$

usually comes from Llyod's

Tolerance parameter

→ until convergence $\left( \| \theta^{t+1} - \theta^t \| \leq \epsilon \right)$

- EM produces "soft clustering"

- EM takes variances into account.



- EM clusters need not be voronoi regions!

$$\lambda^{t+1} = \arg\max_{\lambda} \ modified\_\log L(\theta^t, \lambda)$$

$$\theta^{t+1} = \arg\max_{\theta} \ modified\_\log L(\theta, \lambda^{t+1})$$

→ end.

Maximization Step

Expectation Step

EM Converges to

a local-maximum

of log likelihood.

EM's converged Solution.

modified-$\log L(\theta, \lambda^{t+2})$

$\log L(\theta)$

modified-log lik$(\theta, \lambda^{t+1})$

$\theta^t$  $\theta^{t+1}$

$\theta^*$

$\theta$