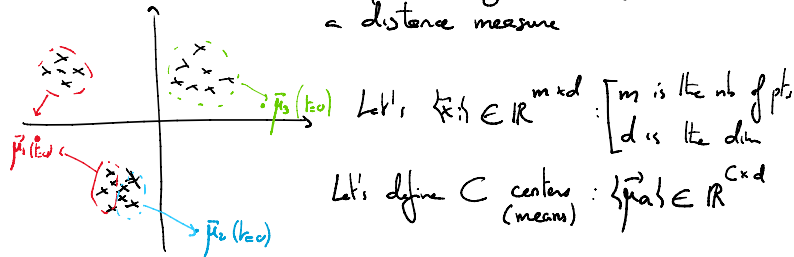


## II Gaussian Mixture model

A) First approach: k-means → intuitive & simple approach of clustering

Based on the following idea:
 

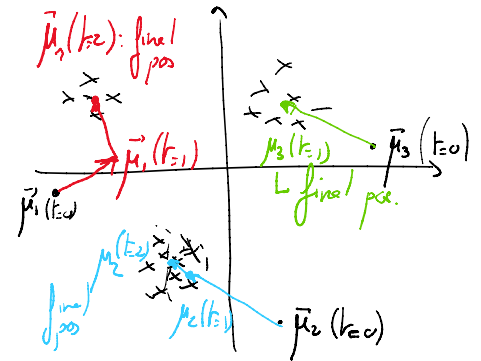
- a cluster of points is a set of points nearby each other
- we put them together according to a distance measure



Algo: 1) assign each pt to the closest centre

$$r_a^i = \begin{cases} 1 & \text{if the centre } a \text{ is the closest of the data } i \\ 0 & \text{otherwise} \end{cases}$$

2) move the centres to the center of mass of the datapoints that have been assigned:

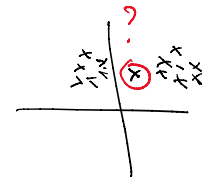


$$\vec{\mu}^a(t+1) = \frac{\sum_{i=1}^m \vec{x}_i \cdot r_a^i}{\sum_{i=1}^m r_a^i}$$

→ very intuitive

Some drawbacks:

- hard assignment
- no scale

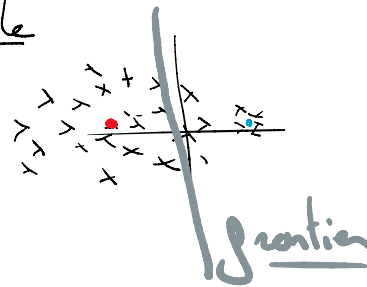


• no geometry

• no quality indicator

• hard to choose initial cdt.

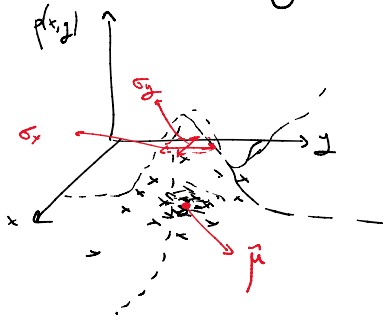
• \_\_\_\_\_ the number  $C$  of clusters



## B) A bayesian approach: GMM

we need a probabilistic model of our clusters

A simple choice: describe a cluster by both- $\mu$ 's  
a gaussian distribution



$$p(\vec{x}|\theta) = \sum_{k=1}^C \rho_k \prod_{i=1}^d \left[ \frac{1}{\sqrt{2\pi\sigma_{ki}^2}} \exp\left(-\frac{(x_i - \mu_{ki})^2}{2\sigma_{ki}^2}\right) \right]$$

$$\Theta = \begin{cases} \cdot \rho_k: \text{density of cluster } k \left( \approx \frac{\# \text{ of pts in class } k}{\text{Total \# of pts}} \right) \\ \cdot \vec{\mu}_k: \text{position of the center (means)} \\ \cdot \vec{\sigma}_k: \text{variances in all directions} \\ \quad \text{of class } k. \end{cases}$$

### ⚠ How the variances is parametrized:

- here we have a variance for each direction along the axis

Other cases: • having the same variance in all directions

• having a full covariance matrix

⚠ huge # of parameters:  $\frac{d(d-1)}{2} \cdot K \sim O(d^2 K)$

Bayes Thm :

$$p(A|B) = \frac{p(B|A) \cdot p(A)}{p(B)}$$

posterior dist.

$$p(\theta | \{\vec{x}\})$$

likelihood

a prior

$$= \frac{p(\{\vec{x}\}|\theta) \cdot p(\theta)}{p(\{\vec{x}\})}$$

prob of my param  $\theta$   
give the data.

$p(\{\vec{x}\})$  (=  $\int d\theta p(\{\vec{x}\}|\theta) p(\theta)$ )  
evidence  $\rightarrow$  we ignore it

Goal: to maximize  $p(\theta | \{\vec{x}\})$  :  $\theta^* = \arg \max p(\theta | \{\vec{x}\})$

when using a uniform prior for  $p(\theta)$

$$\rightarrow \boxed{\theta^* = \arg \max p(\{\vec{x}\}|\theta)}$$

)

- \_\_\_\_\_ the number C of clusters

Let's do it in 1D : impossible to obtain directly the global max!

→ gradient ascent

log-likelihood  $\mathcal{L} = \log p(x|\theta)$

$$\frac{\partial \mathcal{L}}{\partial \mu_k} = \sum_{i=1}^m \left[ \frac{\frac{\rho_k}{\sqrt{2\pi\sigma_k^2}} \frac{(x_i - \mu_k)}{\sigma_k^2} \cdot e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sum_{k'=1}^C \rho_{k'} \frac{1}{\sqrt{2\pi\sigma_{k'}^2}} \cdot \exp\left[-\frac{(x_i - \mu_{k'})^2}{2\sigma_{k'}^2}\right]} \right] = 0$$

no explicit form for the gradient

To find the max: using  $\left[ \begin{array}{l} \text{Expectation} \\ \text{Maximization} \end{array} \right]$  ~ similar to the way k-means work

Update eqs with EM

1) we assign (in proba) each data to a cluster

the responsibility

$$r_a^i = \frac{\frac{\rho_a}{\sqrt{2\pi\sigma_a^2}} \exp\left(-\frac{(x_i - \mu_a)^2}{2\sigma_a^2}\right)}{\sum_k \frac{\rho_k}{\sqrt{2\pi\sigma_k^2}} \exp\left[-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right]} : \underline{\underline{\text{Expectation}}}$$

ii) update of the parameters :

$\sum_i : R \cdot x$   
 $R^{c \times m} R^{m \times d}$

$$\left[ \begin{array}{l} \mu_a = \frac{\sum_{i=1}^m r_a^i x_i}{\sum_{i=1}^m r_a^i} \quad (R_a := \sum_i r_a^i) \\ \sigma_a = \frac{\sum_{i=1}^m r_a^i (x_i - \mu_a)^2}{R_a} \\ \rho_a = \frac{R_a}{\sum_a R_a} \end{array} \right.$$

Maximization  
step



