

Mixture Models \Rightarrow clusters of Data

E.g. Spherical Gaussian MM

$$P(\underline{x}; \underline{\theta}, m) = \sum_{j=1}^m P(j) N(\underline{x}; \underline{m}_j, \sigma_j^2 I)$$

$$\underline{\theta} = \{ P(j), \underline{m}_j, \sigma_j^2 \}_{j=1}^m$$

Find where the clusters are, RATHER THAN $\underline{\theta}$.

\rightarrow To keep the model simpler, $P(j) = \frac{1}{m}, \forall j=1, 2, \dots, m$

$$P(\underline{x}; \underline{\theta}, m) = \frac{1}{m} \sum_{j=1}^m N(\underline{x}; \underline{m}_j, \sigma_j^2 I)$$

\rightarrow To simplify the model even more, let $\sigma_j^2 = \sigma^2 \quad \forall j=1, 2, \dots, m$

$$P(\underline{x}; \underline{\theta}, m) = \frac{1}{m} \sum_{j=1}^m N(\underline{x}; \underline{m}_j, \sigma^2 I)$$

$$\underline{\theta} = \{ \{\underline{m}_j\}_{j=1}^m, \sigma^2 \}$$

Q: How are points assigned to clusters in the EM alg.?

\rightarrow posterior assignment!

E-step: We compute $P(i|xt, \underline{\theta})$ under current pars $\underline{\theta}$

for every x_t & for every $i=1, 2, \dots, m$

Consider the points $\underline{x} \in \mathbb{R}^d$ s.t. $P(i|x_t, \theta) = P(j|x_t, \theta)$ for $i \neq j$

calculation

$$\frac{P(x_t|i, \theta)}{P(x_t|\theta)} = \frac{P(x_t|j, \theta)}{P(x_t|\theta)} \quad P(i) = P(j) = \frac{1}{m}$$

$$\Rightarrow P(x_t|i, \theta) = P(x_t|j, \theta)$$

$$\Rightarrow N(\underline{x}; \underline{\mu}_i, \sigma^2 I) = N(\underline{x}; \underline{\mu}_j, \sigma^2 I) \quad \text{常数项相等!}$$

$$\Rightarrow \exp\left(-\frac{1}{2\sigma^2} \|\underline{x} - \underline{\mu}_i\|^2\right) = \exp\left(-\frac{1}{2\sigma^2} \|\underline{x} - \underline{\mu}_j\|^2\right)$$

The set of points \underline{x} s.t. $P(i|\underline{x}, \theta) = P(j|\underline{x}, \theta)$ is precisely those \underline{x} s.t.

$$\|\underline{x} - \underline{\mu}_i\| = \|\underline{x} - \underline{\mu}_j\| \quad \text{二次項消失}$$

$$\Leftrightarrow 2\underline{x}^T (\underline{\mu}_i - \underline{\mu}_j) = \|\underline{\mu}_i\|^2 - \|\underline{\mu}_j\|^2$$

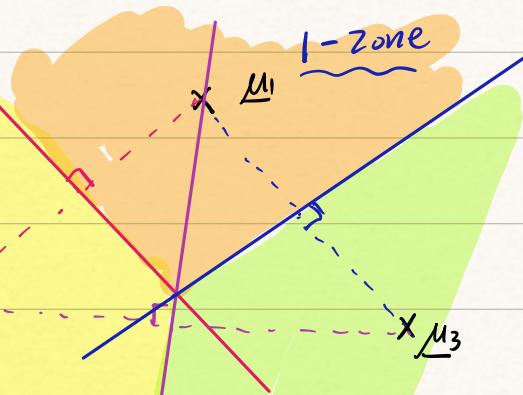


define a line!



Boundary is linear with \underline{x} !

Diagram:
equal mixture proportion
 $P(j) = \frac{1}{m} \quad \forall j = 1, 2, \dots, m$



$$x: P(x|1, \theta) = P(x|2, \theta)$$

$$x: P(x|1, \theta) = P(x|3, \theta)$$

$$x: P(x|2, \theta) = P(x|3, \theta)$$

2-zone

3-zone

(k-means)

Consider the following simpler (hard-assignment) version of EM.

different (which is hard assignment)

E-step: Assign each point x_t to the Voronoi region
with hard prob.

$$x_t \rightarrow j_t = \underset{1 \leq j \leq m}{\operatorname{argmin}} \|x_t - \mu_j\|^2$$

M-step: Recompute μ_j as the mean of assigned points

$$C_j = \{x_t : j_t = j\}$$

$$\mu_j = \frac{1}{|C_j|} \sum_{x \in C_j} x \quad \forall j=1, 2, \dots, m$$

Rmk: The M-step is equivalent to optimizing (maximizing) the
Complete Log-likelihood w.r.t the assignment j_1, \dots, j_t

(obtained in E-step)

over the unknown means μ_1, \dots, μ_j

K-means: Algorithm

Input: $D = \{x_t\}_{t=1}^n \subset \mathbb{R}^d$

1) step 0: Initialize centroids $\mu_j^{(0)} \quad j=1, 2, \dots, m$

SMARTLY

2) step $\ell \in \mathbb{N}$ (E-step): Assign each x_t to CLOSEST Mean

$$j_t^{(\ell)} = \underset{j=1, \dots, m}{\operatorname{argmin}} \|x_t - \mu_j\|^2 \quad t=1, 2, \dots, n$$

3) Step $\ell \in \mathbb{N}$ (M-step): Recompute mean:

$$\underline{\mu}_j^{(\ell+1)} = \frac{1}{|C_j^{(\ell)}|} \sum_{t: j_t^{(\ell)}=j} \underline{x}_t \quad \underline{C_j^{(\ell)} = \{ \underline{x}_t : j_t^{(\ell)} = j \}}$$

4) Terminate when cluster assignment $j_1^{(\ell)}, j_2^{(\ell)}, \dots, j_n^{(\ell)}$
do not change!



Natural terminate criterion

The cost Function for a Given partition $\mathcal{D} = \bigcup_{j=1}^m C_j$

$$\rightarrow C_j = \{ x \in \mathcal{D} : j = \operatorname{argmin}_{j'=1,2,\dots,m} \|x - \underline{\mu}_j\|^2 \}$$

Consider the cost function / objective function:

$$\Rightarrow J(C_1, \dots, C_m : \underline{\mu}_1, \underline{\mu}_2, \dots, \underline{\mu}_m) = \sum_{j=1}^m \sum_{x \in C_j} \|x - \underline{\mu}_j\|^2$$

\downarrow this part actually is fixed by certain partition
 \sum the distance² for each \underline{x} to their
own centroid $\underline{\mu}$

C_j : set of points closer to $\underline{\mu}_j$ than any other $\underline{\mu}_i$ ($i \neq j$)

Last Lecture teaches us!

[Proposition] : GMM \Rightarrow log-likelihood is non-decreasing ! (EM)

Let $J^{(e)} = J(C_1^{(e)}, \dots, C_m^{(e)}; \mu_1^{(e)}, \dots, \mu_m^{(e)})$

$\underbrace{C_1^{(e)}, \dots, C_m^{(e)}}_{\text{l}^{\text{th}} \text{ iteration of k-means}}, \underbrace{\mu_1^{(e)}, \dots, \mu_m^{(e)}}_{}$

①

The k-means algorithm results in :

$J^{(e+1)} \leq J^{(e)}$ for all $e \in \mathbb{N}$

this is not enough for termination!

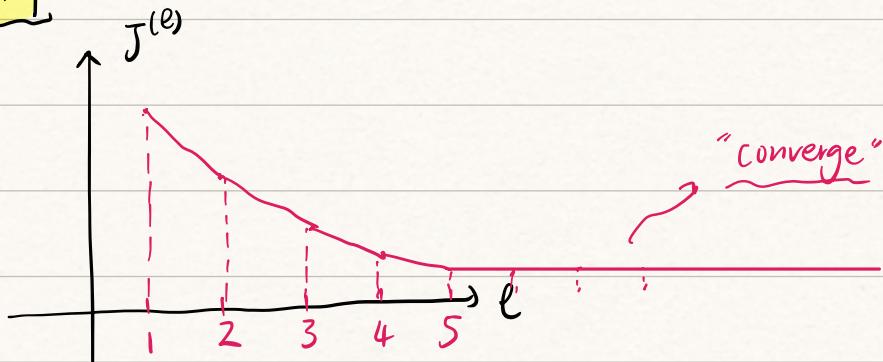
$J^{(e)}$ cannot increase!

②

The alg. terminates in a finite # of steps ($\leq m^n$)

↙ stop criterion is effective

diagram



Pf: Once \mathcal{D} are assigned to their initial clusters, we compute

$$\|x_t - \mu_{j_t}\|^2$$

Summing this up over all $t=1, 2, \dots, n$, yields :

$$J(C_1^{(1)}, \dots, C_m^{(1)}; \mu_1^{(1)}, \dots, \mu_m^{(1)})$$

$$j^{(1)} = \left\{ x_t \in \mathcal{D} : j = j_t^{(1)} = \arg \min_{j'=1,2,\dots,m} \|x_t - \mu_{j'}^{(1)}\|^2 \right\}.$$

For the M-step, we min the total dist. of each cluster's samples from the resp. centroid so the cost f is

$$J(C^{(1)}, \dots, C_m^{(1)}; \mu_1^{(2)}, \dots, \mu_m^{(2)})$$

The function $f(\mu) = \sum_{x \in G} \|x - \mu\|^2$ convex with μ .

$$\Rightarrow \nabla_\mu f(\mu) = 0 \Rightarrow \sum_{x \in G} x = \mu |G|$$

$$\Rightarrow \mu_j^{(2)} = \frac{1}{|G_j|} \sum_{x \in G_j} x$$

the M-step actually minimizes! (for fixed assignment)



the choice of update paras yields

$$\mu_j^{(l)} \rightarrow \mu_j^{(l+1)}$$

minimize for a fixed

assignment!

(M-step)

partition

$$J(C^{(1)}, \dots, C_m^{(1)}; \mu_1^{(2)}, \dots, \mu_m^{(2)})$$

$$\leq J(C^{(1)}, \dots, C_m^{(1)}; \mu_1^{(1)}, \dots, \mu_m^{(1)})$$

E-step: ① if a sample is assigned to a different centroid, it must be that the sample is closer to its new centroid than the previously assigned one!

② If a sample is not reassigned, then nothing happen!

From ① & ②

$$J(C^{(2)}, \dots, C_m^{(2)}; \mu_1^{(2)}, \dots, \mu_m^{(2)})$$

$$\leq J(C^{(1)}, \dots, C_m^{(1)}; \mu_1^{(1)}, \dots, \mu_m^{(1)})$$

Thus, $J(G^{(2)}, \dots, C_m^{(2)}; \mu_1^{(2)}, \dots, \mu_m^{(2)})$

$\leq J(G^{(1)}, \dots, C_m^{(1)}; \mu_1^{(1)}, \dots, \mu_m^{(1)})$

i.e., $\underbrace{J^{(2)}}_{\text{decrease!}} \leq J^{(1)} \rightarrow \text{decrease!}$

Pf^④: Procedure terminates in a finite # of steps

There is a finite # of samples n

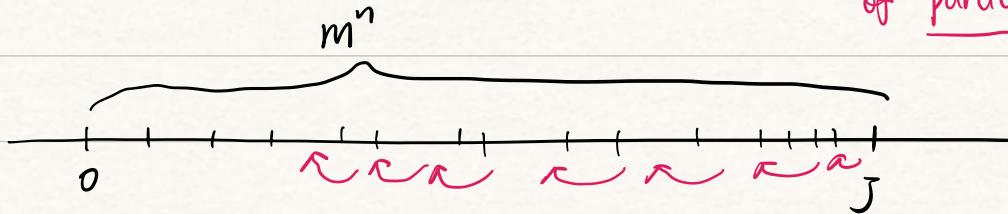
\Rightarrow There is a finite # of labellings $\rightarrow [m^n]$

Note that $J(G, \dots, C_m; \hat{\mu}_1, \dots, \hat{\mu}_m)$

where $\hat{\mu}_j = \frac{1}{|G_j|} \sum_{x \in G_j} x$

fixed

can only take on m^n values! \Rightarrow Actually J is a function of partition (m^n)



\Rightarrow Alg. will terminate at iteration $\ell^* \in \mathbb{N}$ with

$$J^{(\ell^*)} = J^{(\ell^*-1)} \Rightarrow \boxed{\text{then } J^\ell = J^{\ell^*} \text{ for all } \ell > \ell^*}$$

$0 \leq J^* \leq J^{(\ell^*)}$ \rightarrow the 'convergence' given by K-means

the true minimum of cost function

NP-Hard Problem to find J^*

↓ Reason

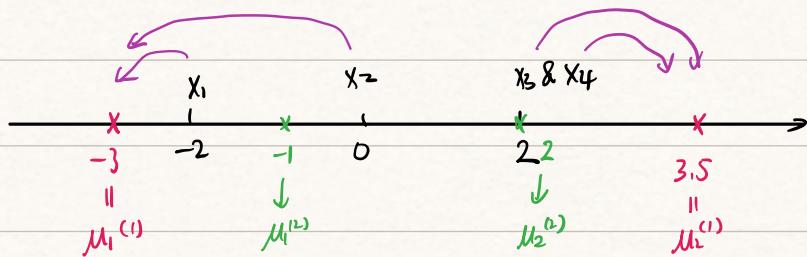
since $J^* = \min_{C_i} J(G_1, \dots, G_m; \hat{M}_1, \dots, \hat{M}_m)$

↓
find the optimal over all possible m

partition ($m^n \uparrow$) \Rightarrow NP-Hard

Ex: $x_1 = -2, x_2 = 0, x_3 = x_4 = 2$

a) $\mu_1^{(1)} = -3, \mu_2^{(1)} = 3.5$



Good initialization $J^{(1)} = 1^2 + 3^2 + 2 \times 1.5^2$

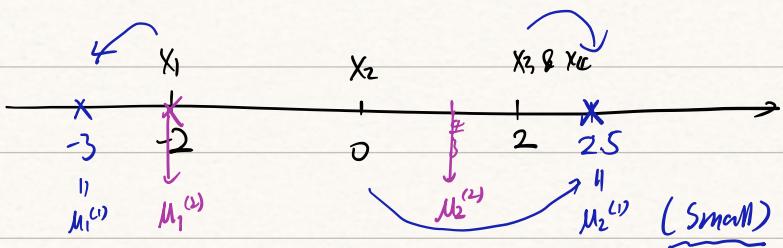
E-step: $j_1 = 1, j_2 = 1, j_3 = j_4 = 2$

M-step: $\mu_1^{(2)} = -1, \mu_2^{(2)} = 2$

$\Rightarrow J^{(2)} = 1^2 + 1^2 + 0^2 + 0^2 = 2$ shrink!

b) BAD Initialization

$\mu_1^{(1)} = -3, \mu_2^{(1)} = 2.5$



E-step: $j_1 = 1, j_2 = j_3 = j_4 = 2$

M-step: $\mu_1^{(2)} = -2, \mu_2^{(2)} = \frac{4}{3}$

$J^{(2)} = 0^2 + (\frac{4}{3})^2 + 2 \cdot (\frac{2}{3})^2 = \frac{16}{9} + \frac{8}{9} = \frac{8}{3} > 2$

Stuck in $J^{(2)} = \frac{8}{3}$ \rightarrow not global minimum of J

a BIG Cost! \Rightarrow 被困住了

Option Part:

K-means ++ Algorithm (2017)

(1) Take an center μ_1 uniformly chosen from $D = \{x_t\}_{t=1}^n$

(2) Take μ_i ($i > 1$) choosing $x \in D$ with prob. given by $\frac{D(x)^2}{\sum D(x')^2}$

$D(x)$: distance of $x \in D$ to the closest center we have chosen.
chosen from data points!

(3) Repeat Step (2) with K centroids are found.

Intuition: Select centroids using the idea that centroids should be as far apart as possible!



with BAD Initialization, we will have BAD $J^{(e*)}$

Given a set of K centers $C = \{\mu_1, \dots, \mu_K\}$

$$J = \sum_{x \in D} \min_{\mu \in C} \|x - \mu\|^2$$

$J_{\text{opt}} = \min_C J \rightarrow \text{NP-Hard Problem}$

Thm: [AV, 2007]

If C is constructed based on a K-means ++

let its cost function be \boxed{J} a random variable.

\hookrightarrow r.v since C is random

Then $\underbrace{\mathbb{E}[J]}_{\downarrow} \leq \delta(\ln K + 2) \underbrace{J_{opt}}_{\downarrow} = O(\log K) J_{opt}$

Global minimum

control by SMART Initialization!

\Rightarrow K-means ++ is $O(\log K)$ -competitive w.r.t

the optimal clustering (which is NP-HARD)

ICML (2020)

Informal

Using εK ($\varepsilon > 0$) local search steps, then with K-means ++,
we can attain $O(1)$ -competitive approx. the J_{opt} .

Some important point of K-means :

1. The decrease of cost function $J^{(l)}$

$$J^{(0)} = J(D_1^{(0)}, \dots, D_m^{(0)}; \mu_1^{(0)}, \dots, \mu_m^{(0)}).$$

partition according to $\mu_1^{(0)}, \dots, \mu_m^{(0)}$

\rightarrow decrease

decrease

$$J_1^{(e)} = J(D_1^{(e)}, \dots, D_m^{(e)}; \mu_1^{(e)}, \dots, \mu_m^{(e)})$$

reset of centroid

$$J^{(e+1)} = J(D_1^{(e+1)}, \dots, D_m^{(e+1)}; \mu_1^{(e+1)}, \dots, \mu_m^{(e+1)})$$

reassign (repartition)

① if $J^{(e)} = J^{(e)}$, \Rightarrow no reset of centroid $\Rightarrow J^{(e)} = J^{(e+1)} = \dots$

\Downarrow

no need of reassignment

② if $J_1^{(e)} = J^{(e+1)}$ \Rightarrow no reassignment

\Downarrow

no need of reset $\Rightarrow J^{(e+1)} = J^{(e+2)} = \dots$

③ if $J^{(e)} = J^{(e+1)}$, then $J^{(e)} \leq J_1^{(e)} \leq J^{(e+1)} \Rightarrow J^{(e)} = J^{(e+1)} = \dots$

no reset

no reassignment

$$\boxed{n \begin{pmatrix} m \\ n \end{pmatrix}}$$

2. local minimum of K-means

$$\boxed{J(D_1, \dots, D_m; \mu_1, \dots, \mu_m) = \sum_{i=1}^m \sum_{x \in D_i} \|x - \mu_i\|_2^2}$$

3. What kind of minimum (local) we can find?

