

# HW 1

1 a) WE HAVE TO SHOW THAT  $\sum_{x \in C_J} d(x; m_J)^2 = \frac{1}{2} \frac{1}{|C_J|} \sum_{x \in C_J} \sum_{x' \in C_J} d(x; x')^2$

NOW:  $\frac{1}{|C_J|} \sum_{x \in C_J} \sum_{x' \in C_J} \underbrace{\|x\|^2 + \|x'\|^2 - 2|x \cdot x'|}_{\substack{\text{SAME TERM} \\ \text{GIVEN THE} \\ \text{DOUBLE SUM}}} = \frac{2}{|C_J|} \sum_{x \in C_J} \sum_{x' \in C_J} (|x|^2 - x \cdot x') \quad (1)$

ON THE OTHER HAND

$$\begin{aligned} \sum_{x \in C_J} d(x; \mu)^2 &= \sum_{x \in C_J} \left( |x|^2 + \left( \frac{1}{|C_J|} \sum_{x' \in C_J} x' \right)^2 - 2x \cdot \frac{1}{|C_J|} \sum_{x' \in C_J} x' \right) = \\ &= \sum_{x \in C_J} \left( |x|^2 + \frac{1}{|C_J|^2} \sum_{x' \in C_J} \sum_{x'' \in C_J} x' \cdot x'' - 2x \cdot \frac{1}{|C_J|} \sum_{x' \in C_J} x' \right) \quad \begin{array}{l} \text{BROKE SUM} \\ \sum (a_i + b_i) = \\ \sum a_i + \sum b_i \end{array} \\ &= \left( \sum_{x \in C_J} |x|^2 \right) + \frac{1}{|C_J|^2} \sum_{x' \in C_J} \sum_{x'' \in C_J} x' \cdot x'' - \left( \frac{2}{|C_J|} \sum_{x \in C_J} \sum_{x' \in C_J} x \cdot x' \right) \end{aligned}$$

$\frac{\sum_{i=1}^{C_J} A}{C_J} = A$   
IF A CONST

$$\begin{aligned} &= \sum_{x \in C_J} |x|^2 - \frac{1}{C_J} \sum_{x \in C_J} \sum_{x' \in C_J} x' \cdot x'' \\ &= \sum_{x \in C_J} \sum_{x' \in C_J} \frac{1}{C_J} (|x|^2 - x' \cdot x) \quad (2) \end{aligned}$$

COMPARING (1) AND (2)  $\rightarrow J_{IC} = 2 J_{AVG}^2$

$$1(b) \quad J_{avg}^2 = \sum_{i=1}^m d(x_i; m_{\delta_i})^2$$

• BY DEFINITION WHEN  $(m_1, \dots, m_k)$  IS FIXED  $\delta_i \leftarrow \arg \min_{j \in \{1, \dots, k\}} d(x_i; m_j)$

MINIMIZE  $\sum d(x_i; m_{\delta_i})$ ; BUT SINCE  $( )^2$  IS MONOTONE IT ALSO

MINIMIZE  $\sum d(x_i; m_{\delta_i})^2$

• TO MINIMIZE  $\sum_{i=1}^m d(x_i; m_{\delta_i})$  WITH  $\delta_i$  FIXED  $\frac{\partial}{\partial m} \sum_{i=1}^m |x_i|^2 + |m_{\delta_i}|^2 - 2 x_i m_{\delta_i} =$

$$\sum_{i=1}^m 2|m_{\delta_i}| - 2x_i = 0 \longrightarrow \sum_{i: \delta_i = j} x_i = |C_j| m_j \quad \text{WHICH CORRESPONDS}$$

TO THE UPDATING RULE; NOTE THAT THIS IS BASICALLY THE DEFINITION OF THE MEAN.

c) THE K MEANS IS BASED ON AN ITERATIVE STEP WITH THE UPDATE RULES IN EXERCISE 1b; SINCE BOTH MINIMIZE  $J$  KEEPING A SET OF ITS ARGUMENT FIXED THE FULL STEP MINIMIZE  $J_{avg}^2(\delta_i; m_i)$

d) THE WORST CASE SCENARIO WOULD BE WHEN THE ALGORITHM HAS TO TRY ALL THE POSSIBLE ASSIGNMENTS OF EACH DATA TO A CLUSTER BEFORE IT FINDS THE RIGHT ONE SO WE HAVE  $O(K^m)$ ; THE ALGORITHM WON'T VISIT THE SAME COMBINATION TWICE BECAUSE IT IS DESIGNED TO MINIMIZE THE COST FUNCTION MONOTONICALLY; THERE ARE STRONGER LIMITS IN THE LITERATURE:

$\times O(n^{kd})$  in  $d$  dimension FROM VORONOI TASSILLATION CONSIDERATION INABA et AL

$\times O(m \Delta^2)$  WITH  $\Delta$  SPREAD OF POINT  $\Delta = \left( \frac{\text{LARGEST PAIRWISE DISTANCE}}{\text{SMALLER " "}} \right)$

3) AS A FUNCTION OF THE PARAMETER THE LIKELIHOOD CAN BE SEEN AS  $P(X|\theta)$ ;  $P(X|\theta) = \prod_{i=1}^m \pi_{z_i} N(x_i | \mu_i; \Sigma_i)$ ;

$$\text{SO } \ln P(X|\theta) = \sum_{i=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k) \right\};$$

THIS IS THE LIKELIHOOD FOR THE COMPLETE DATA SET  $P(X|\pi, \mu, \Sigma)$  WHERE  $X$  ALSO CONTAINS THE LATENT VARIABLE  $z$   $X = (X, \bar{z})$

b) THIS CAN BE OBTAINED USING THE BAYES THEOREM

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$\text{WE WANT } P(z_{i,j}=1 | x_i) = \frac{P(x_i | z_{i,j}=1) P(z_{i,j}=1)}{P(x_i)}$$

NOW,  $P(z_k) = \pi_k$  AND  $P(x_i | z_{i,j}=1)$  WAS DEFINED IN 3a

$$P(z_{i,j}=1 | x_i) := p_{i,j} = \frac{\pi_j N(x_i; \mu_j, \Sigma_j)}{\sum_{j=1}^J \pi_j N(x_i; \mu_j, \Sigma_j)}$$

c) NOW GIVEN THE LATENT VARIABLE  $z$  WE HAVE

$$P(z) = \prod_{k=1}^K \pi_k^{z_k} \quad \text{AND} \quad P(X|z) = \prod_{n=1}^N N(x_n; \mu_{z_n}, \Sigma_{z_n})^{z_{z_n}}$$

$$\text{SO } \ln P(X, z | \mu, \Sigma, \pi) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln N(x_n | \mu_k, \Sigma_k) \}$$

WITH  $z_{nk}$   $k$  COMPONENT OF  $\bar{z}$ ; NOW WITH BAYES THEOREM

$$P(z | X, \mu, \Sigma, \pi) \propto \prod_{n=1}^N \prod_{k=1}^K [\pi_k N(x_n | \mu_k, \Sigma_k)]^{z_{nk}} \quad \text{SO EASILY } E[z_{nk}] = p_{nk}$$

$$\text{THIS LEADS TO } E_z[\ln P(X, z | \mu, \Sigma, \pi)] = \sum_{n=1}^N \sum_{k=1}^K p_{nk} \{ \ln \pi_k + \ln N(x_n | \mu_k, \Sigma_k) \}$$

3d) WE ARE LOOKING FOR  $\frac{\partial}{\partial \pi_j} \bar{e}_{out}(\theta) = 0$ ; TO IMPOSE THE CONSTRAINT  $\sum_{j=1}^K \pi_j = 1$  WE CAN USE THE

LAGRANGE MULTIPLIERS

$$\sum_{i=1}^m \sum_{j=1}^K P_{i,j} (\log \pi_j) + \lambda \left( \sum_{j=1}^K \pi_j - 1 \right)$$

$$\downarrow \frac{\partial}{\partial \pi_j}$$

$$\sum_{i=1}^m \frac{P_{i,j}}{\pi_j} + \lambda = 0 \rightarrow \pi_j = - \frac{\sum_i P_{i,j}}{\lambda} \quad (1)$$

$$\frac{\partial}{\partial \lambda} = 0$$

$$\sum_{j=1}^K \pi_j = 1 \rightarrow - \sum_{j=1}^K \sum_i \frac{P_{i,j}}{\lambda} \rightarrow \lambda = -m \quad (2)$$

$$\left[ \pi_j = \frac{1}{m} \sum_i P_{i,j} \right] \Delta \text{ COMBINING } (1) \ (2)$$

3e)  $\boxed{\mu}$  WE WANT  $\frac{\partial}{\partial \mu} = 0$  OF THE LOG LIKELIHOOD

$$\frac{\partial}{\partial \mu_n} \sum_{i=1}^N \sum_{k=1}^K P_{i,k} \log N(x_i | \mu_n \Sigma_n) = 0 \quad \text{SINCE } \pi_k \text{ AND}$$

$P_{i,j}$  DOES NOT DEPEND ON  $\mu$ .

$$\sum_{i=1}^N \sum_{k=1}^K P_{i,k} \frac{\partial}{\partial \mu_n} \left[ (\bar{x} - \mu_n)^T \Sigma_n^{-1} (\bar{x} - \mu_n) \right] = \sum_{i=1}^N P_{i,k} - \frac{1}{2} (-2 \Sigma_n^{-1} (\bar{x} - \mu_n))$$

$$= \sum_{i=1}^N P_{i,k} \Sigma_n^{-1} (\bar{x} - \mu_n)$$

$$\text{NOW } \frac{\partial}{\partial \mu} ( ) = 0 \quad \mu_n = \frac{\sum_{i=1}^m P_{i,k} x_i}{\sum_{i=1}^m P_{i,j}}$$

use  
 $\partial(\bar{A} \cdot \bar{B})$   
 $= \partial(\bar{A}) \bar{B} +$   
 $\bar{A} \partial(\bar{B})$

$$3e) \boxed{\Sigma} \quad \frac{\partial}{\partial \Sigma} ( ) = \sum_{i=1}^m P_{i,n} \frac{\partial}{\partial \Sigma_n} \left\{ -\frac{1}{2} \ln(|\Sigma|) - \frac{1}{2} (x_n - \mu)^T \Sigma_n^{-1} (x_n - \mu) \right\}$$

$$\text{Now } \frac{\partial}{\partial \Sigma} \ln(|\Sigma|) = (\Sigma^{-1})^T$$

$$\text{AND } \frac{\partial}{\partial \Sigma} (x - \mu)^T \Sigma^{-1} (x - \mu) = (\Sigma^{-1} (x_i - \mu)(x_i - \mu)^T \Sigma^{-1})^T$$

THE LATTER CAN BE SEEN USING  $\text{Tr}[AB] = \text{Tr}[BA]$

$$\sum_{n=1}^N (x - \mu)^T \Sigma^{-1} (x - \mu) =$$

$$\begin{aligned} \text{SO } \frac{\partial}{\partial \Sigma} ( ) &\propto \sum_{i=1}^m P_{i,n} \left\{ (\Sigma_n^{-1})^T - (\Sigma_n^{-1} (x_i - \mu_n)(x_i - \mu_n)^T \Sigma_n^{-1})^T \right\} \\ &\propto (\Sigma_n^{-1})^T \sum_{i=1}^m P_{i,n} \left\{ 1 - (\Sigma_n^{-1} (x_i - \mu_n)(x_i - \mu_n)^T \Sigma_n^{-1})^T \right\} \end{aligned}$$

WHERE WE NEGLECT CONSTANT BECAUSE WE ARE INTERESTED IN  $\frac{\partial}{\partial \Sigma} ( ) = 0$

$$\begin{aligned} \text{NOW } \frac{\partial}{\partial \Sigma} ( ) = 0 \quad \sum_i P_{i,n} &= \Sigma_n^{-1} \sum_i P_{i,n} (x_i - \mu_n)(x_i - \mu_n)^T \\ \downarrow \\ \boxed{\Sigma_n} &= \frac{\sum_i P_{i,n} (x_i - \mu_n)(x_i - \mu_n)^T}{\sum_i P_{i,n}} \end{aligned}$$

3(4)

## GAUSSIAN MIXTURE (GM)

• START WITH GUESS FOR  $\pi_j, \mu_j, \sigma_j$

• ASSIGN A WEIGHT TO EACH POINT TO BELONG TO SAME CLUSTER WITH  $P_{ij} = \pi_j N(x_i | \mu_j, \sigma_j)$

• NOW MAXIMIZE GIVEN THE PARAMETERS  $P_{ij}$ ;

$$\pi_j \leftarrow \frac{1}{n} \sum_i P_{ij}; \quad \mu_j = \frac{\sum_i P_{ij} x_i}{\sum_i P_{ij}}$$

$$\sigma_j^2 \leftarrow \frac{\sum_i P_{ij} (x_i - \mu_j)(x_i - \mu_j)^T}{\sum_i P_{ij}}$$

THEY ARE SIMILAR BUT 1) IN GM YOU HAVE A MODEL; SO YOU HAVE TO MAXIMIZE 3 PARAMS  $\pi, \sigma, \mu$  INSTEAD OF THE ONLY  $m_j$  YOU HAVE IN KMEANS

2) IN GM ALL THE POINTS BELONG TO ALL THE GAUSSIAN INSTEAD OF THE BINARY ASSIGNATION YOU HAVE IN KMEANS

## K-MEANS

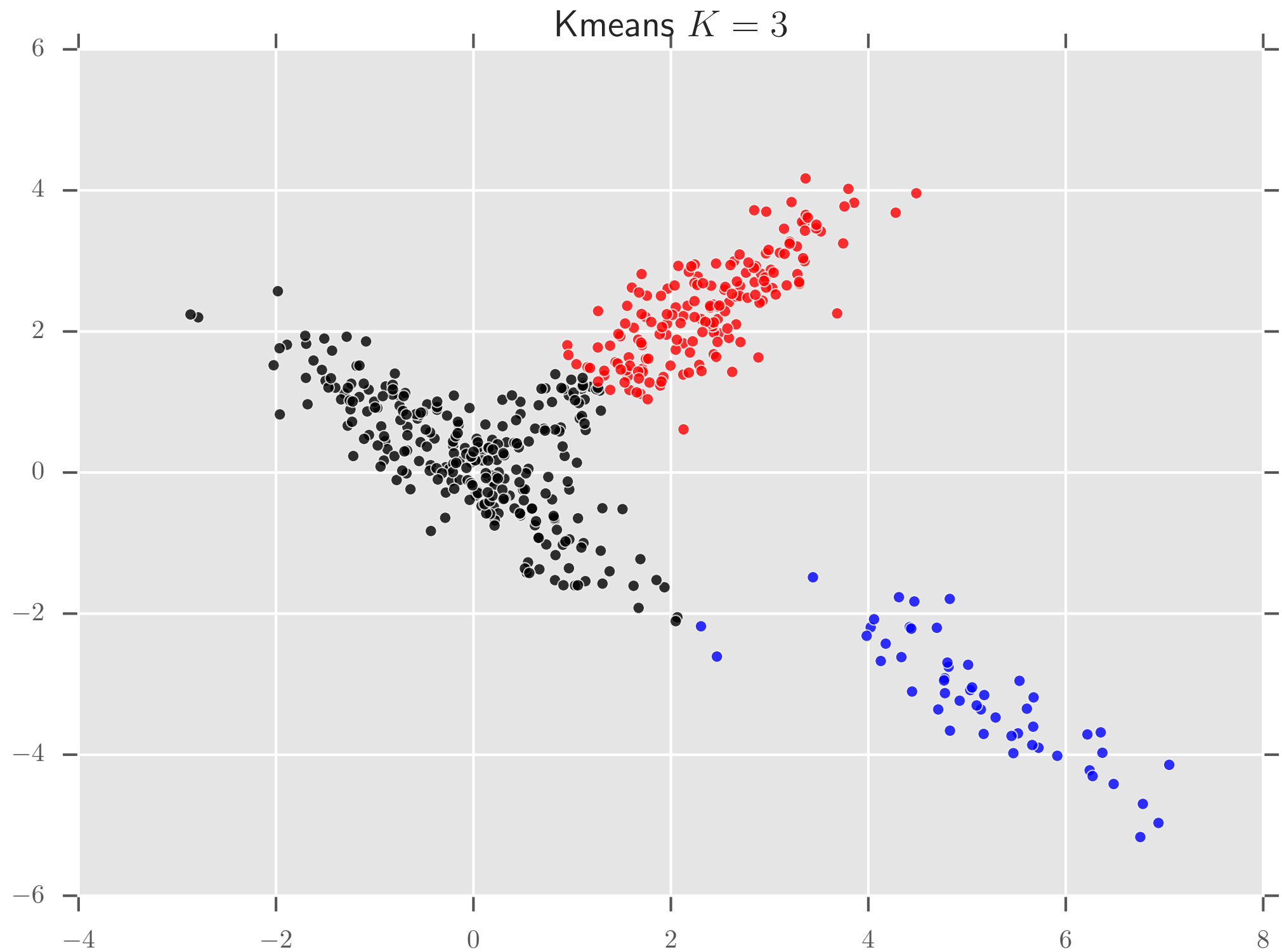
• INITIALIZE RANDOMLY THE  $K$  CENTROID  $\mu$  SAMPLE OR USING K++

• ASSIGN EACH POINT TO A CLUSTER BY MINIMIZING  $\sum_i d(x_i, m_{\delta_i})^2$  WITH  $\bar{m}_i$  FIXED

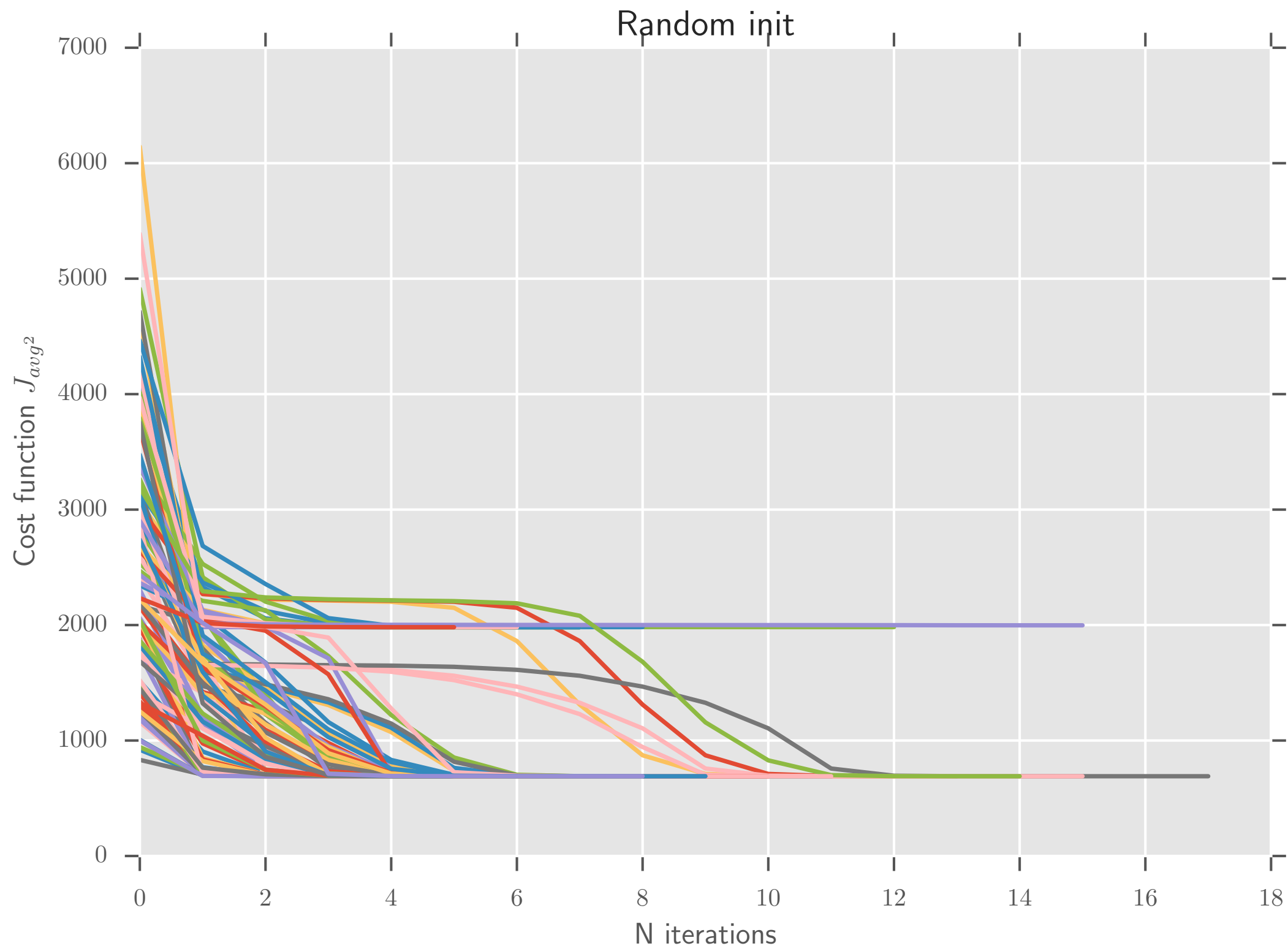
↳ THIS IS A ONE VS ALL ASSIGNMENT; IN GM YOU ASSIGN EACH POINT TO ALL OF THE GAUSSIAN WITH DIFFERENT WEIGHTS

• MAXIMIZE GIVEN THE ASSIGNMENT  $\delta_i$

$$m_j = \frac{1}{|C_j|} \sum_{i: \delta_i = j} x_i$$

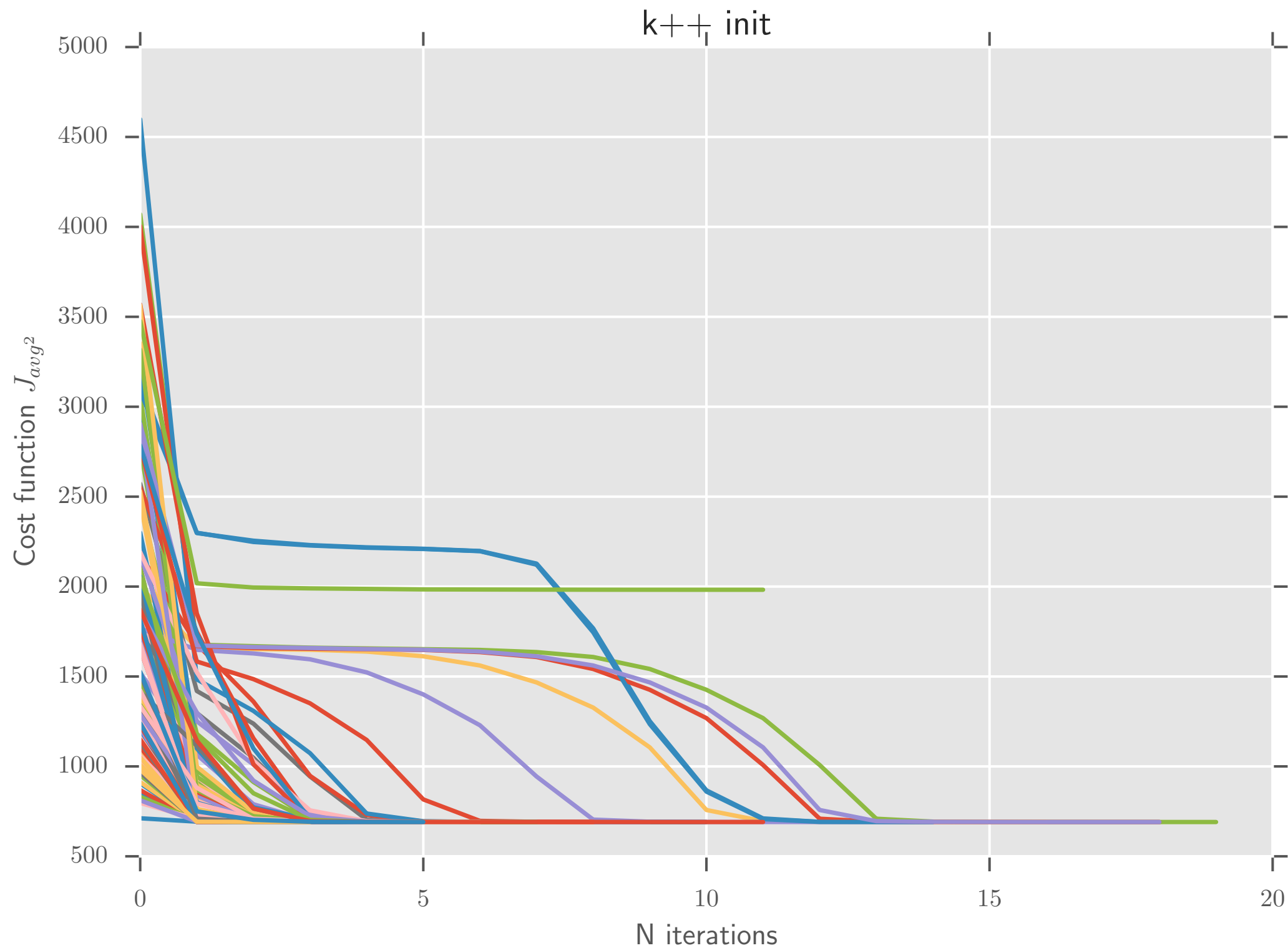


Results of the K-means algorithm. Each color represents a different cluster

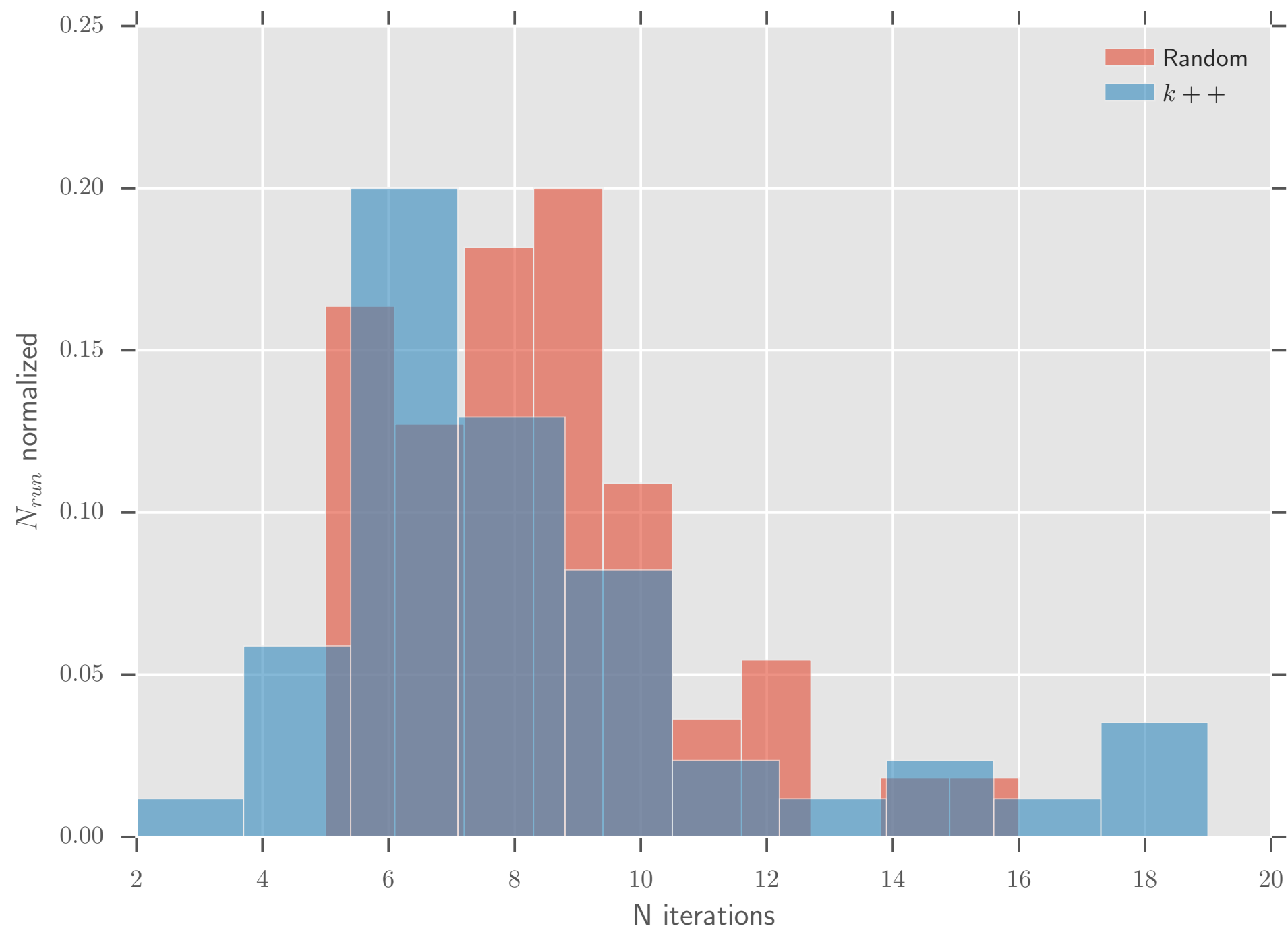


Cost function as a function of the iteration for the Kmeans algorithm (100 runs). As expected the function is monotonically decreasing at each step and most of the curves converge on the right minima of  $J \sim 680$ . Some of the curves does not converge to the minima because they are stuck in local minima. This is why one should always run the K-mean algorithm multiple time and pick the configuration that minimize  $J$  among them

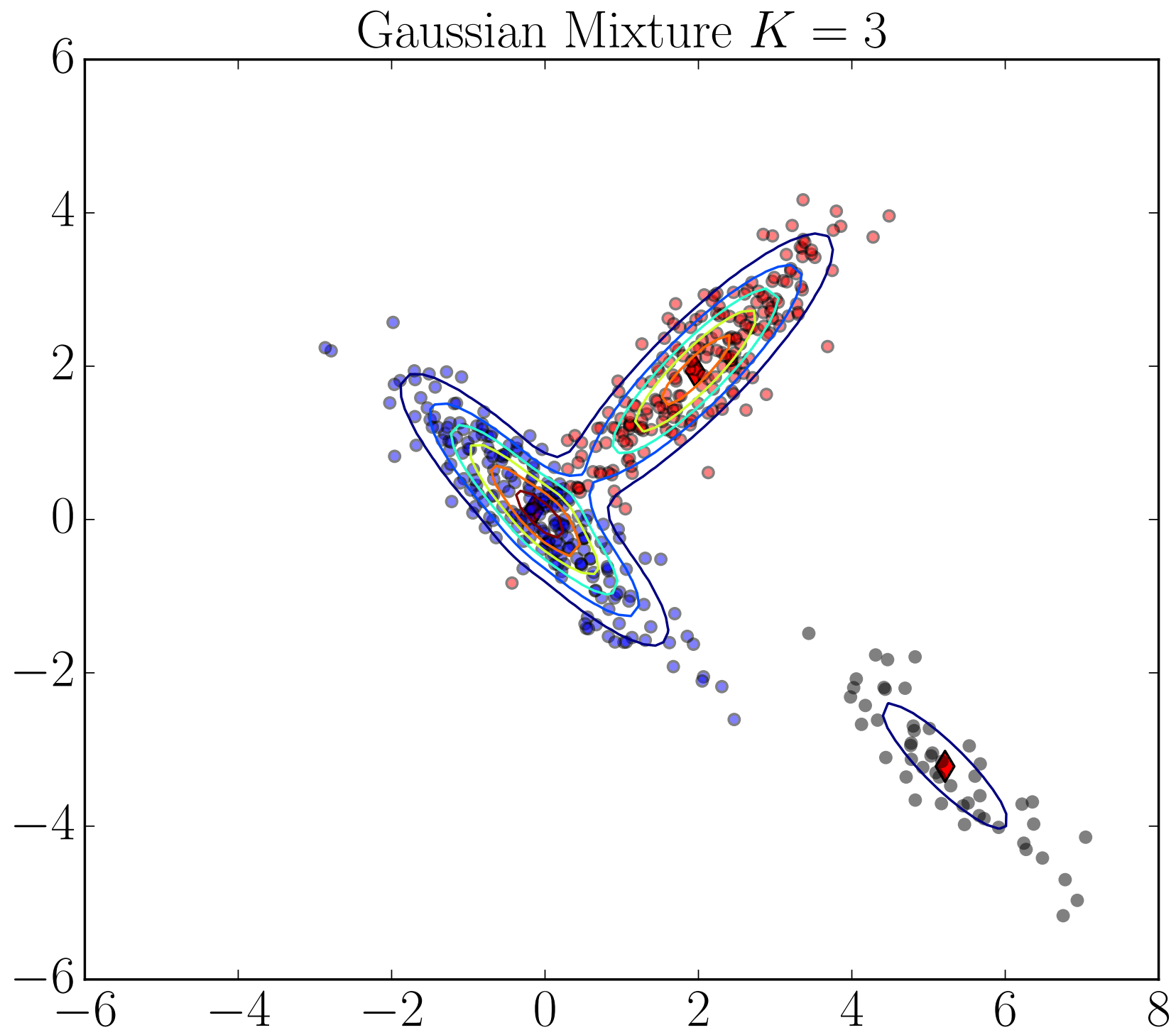




Cost function as a function of the iteration for the Kmeans algorithm (100 runs) with a k++inicialization. The same comments of the previous figures apply here. Note that it can be see by eye that the cost function start at lower values because of the clever initialization and that the curves converge a little bit faster than in the random init case. See next figure on this.



Even if the number of runs of the algorithm is not enough to have a statistical sample it can be seen with K++ it converges a little bit faster than with a random init. Indeed the distribution is peaked a little bit on the left. Also the problem here is not hard enough to see the full power of the k++ initialization



Results of the gaussian mixture algorithm. The red diamonds correspond to the mean of the three gaussian and the contour plot allows to have an intuition of the different weights and the direction of the covariance matrices.

The weights are 0.44321029, 0.10000401, 0.4567857