

Sheet

(Module 3)

Science IFT3700 data /

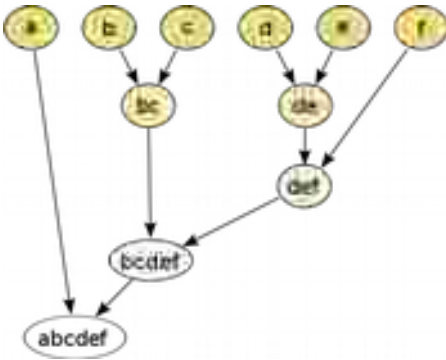
IFT6758 Fall 2018 ©

Alain Tapp

Contents

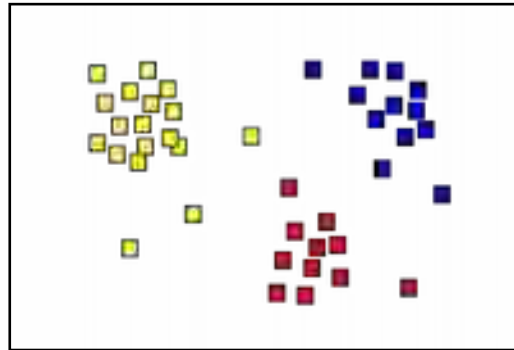
- Partitionnement de données
- Distance et similarité
- Hierarchical clustering
- Expectation maximization algorithm
 - - average
 - GMM
- DBSCAN

grouping data partitioning methods



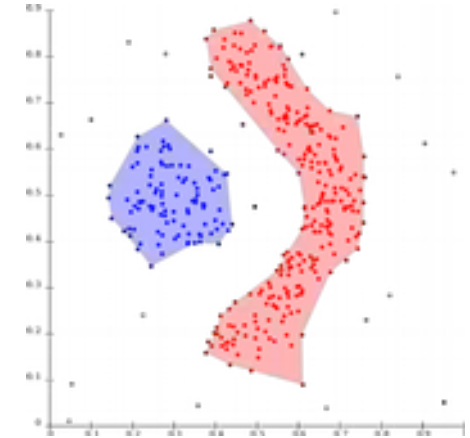
hierarchical

- Similarity
- by division
- by grouping



centroid

- Similarity or distance
- Average
- Representative



neighborhood

- Similarity distance
- neighborhood graph
- Density

Partitioning data

Given a set of points, with a notion of distance between points, group points to a number of groups, so that:

- The group members are close / similar to each other.
- Members of different groups are different.
- Points are scattered in a large space.
- The proximity is defined with the aid of a measure of similarity or distance.

Partitioning data

Clustering in two dimensions look easy.
Le regroupement en deux dimensions semble facile.

Many applications involve a large number of dimensions (MNIST.)
De nombreuses applications impliquent un grand nombre de dimensions (MNIST, $d = 784$)

In large areas, it is often noted that almost all pairs of points are about the same distance.
Dans les espaces de grande dimension, on remarque souvent que presque toutes les paires de points sont à peu près à la même distance.

Examples

galaxies

A catalog of 2 billion "celestial objects" represents objects by their radiation in 7 dimensions (frequency bands). How to group similar objects, such as galaxies, nearby stars, quasars, etc.

Music

The music is divided into categories and customers prefer certain categories. What is the best ranking. Represent parts of a set of customers who bought it.

Texts and articles

Representing the documents by all the significant words they contain.

distances

Distance

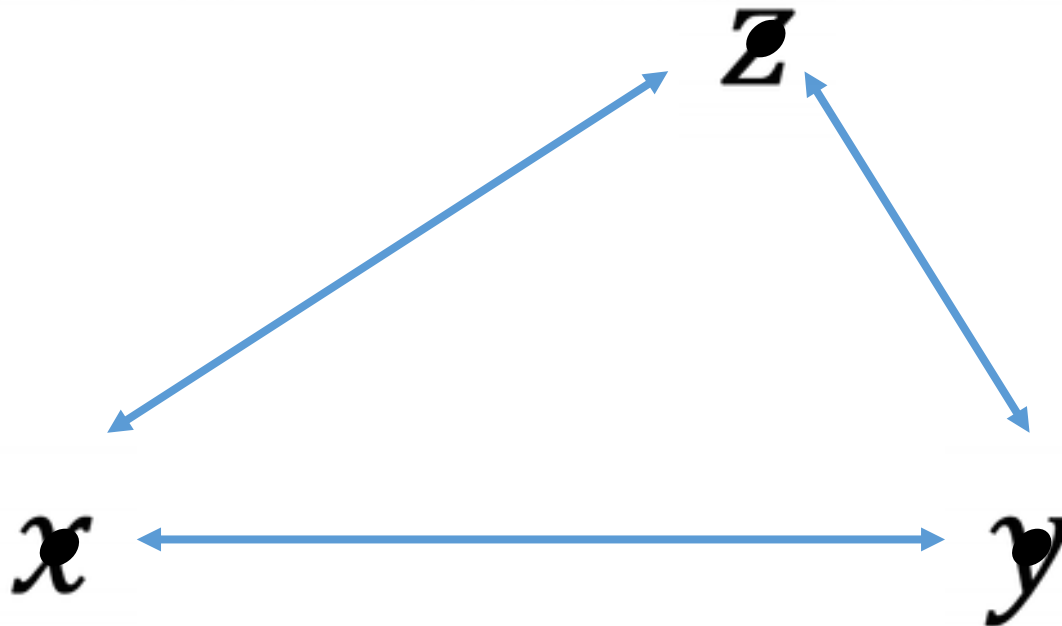
$$D: E \times E \rightarrow \mathbb{R}^+$$

Positivité: $D(x, y) \geq 0$

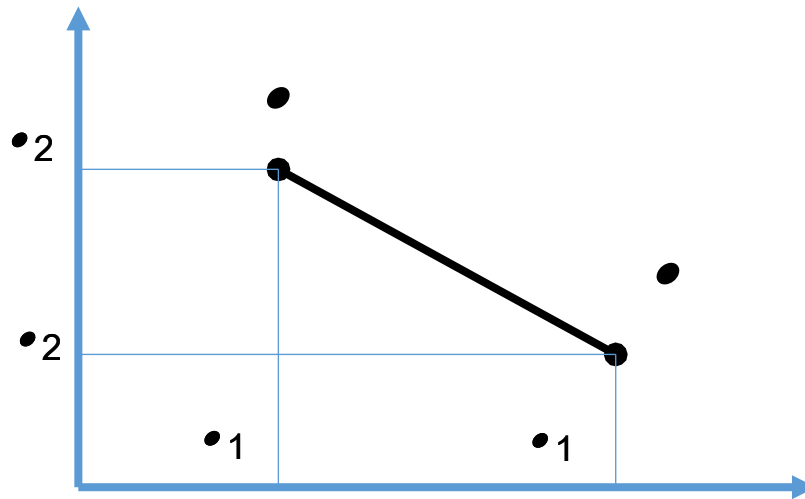
Symétrie: $D(x, y) = D(y, x)$

Séparation: $D(x, y) = 0 \rightarrow x = y$

Inégalité triangulaire: $D(x, z) \leq D(x, y) + D(y, z)$

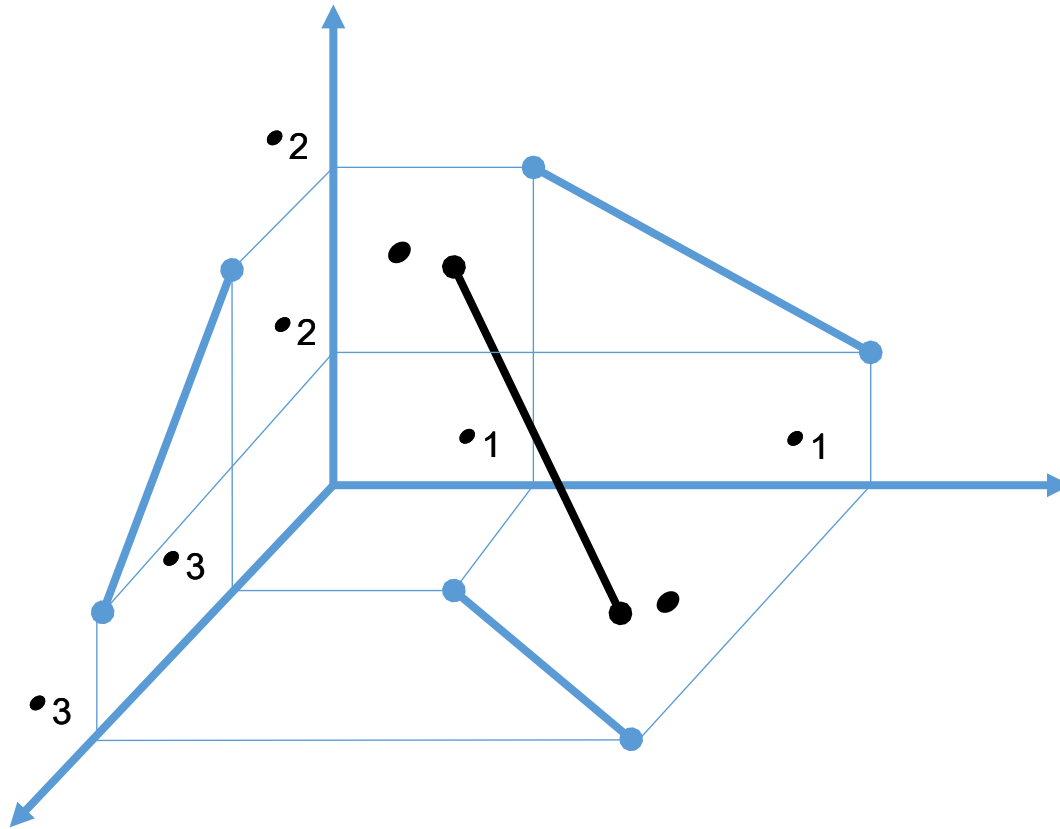


Euclidean Distance



.

Euclidean Distance



.

$$\bullet = (\bullet_1 \bullet_2 \bullet_3)$$

$$\bullet = (\bullet_1 \bullet_2 \bullet_3)$$

$$\bullet(\bullet, \bullet) = \sqrt{(\bullet_1 - \bullet_1)^2 + (\bullet_2 - \bullet_2)^2 + (\bullet_3 - \bullet_3)^2}$$

$$\bullet = (\bullet_1 \bullet_2 \bullet_3 \bullet_4)$$

$$\bullet = (\bullet_1 \bullet_2 \bullet_3 \bullet_4)$$

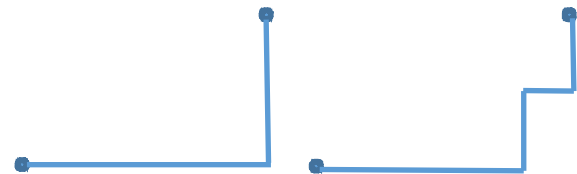
$$\bullet(\bullet, \bullet) = \sqrt{(\bullet_1 - \bullet_1)^2 + (\bullet_2 - \bullet_2)^2 + (\bullet_3 - \bullet_3)^2 + (\bullet_4 - \bullet_4)^2}$$

$$\bullet = (\bullet_1 \bullet_2 \bullet_3 \cdots, \bullet_{784}) \quad \bullet = (\bullet_1 \bullet_2 \bullet_3 \cdots, \bullet_{784})$$

$$\bullet(\bullet, \bullet) = \sqrt{(\bullet_1 - \bullet_1)^2 + (\bullet_2 - \bullet_2)^2 + \cdots + (\bullet_{784} - \bullet_{784})^2}$$

$$L_1(x, y) = \sum_{i=1}^n |x_i - y_i|$$

Manhattan

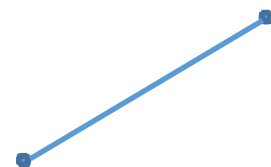


$$L_2(x, y) = \sum_{i=1}^n |x_i - y_i|^2$$

square error

$$L_2(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^2 \right)^{1/2}$$

Euclidean



$$L_\infty(x, y) = \max_i |x_i - y_i|$$

Max



cosine pseudodistance

$$x \cdot y = \sum_{i=1}^n x_i y_i \quad \dots$$

$$\|x\| = \sqrt{x \cdot x}$$

Cosine similarity with the angle θ between the vectors x and y .

$$\cos(\theta) = S(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

$$-1 \leq S(x, y) \leq 1$$

Pseudorange cosine nonnegative vectors

$$D(x, y) = 1 - S(x, y)$$

Mahalanobis Distance

Each coordinate is normalized according to the variance in that dimension.
On normalise chaque coordonnée en fonction de la variance dans cette dimension.

$$d(x, y) = \sqrt{\sum_{i=1}^n \left(\frac{x_i - y_i}{\sigma_i} \right)^2}$$

In general, if the covariance matrix of the distribution is obtained:

En général, si la matrice de covariance de la distribution est obtenue:

$$D(x, y) = \sqrt{(x - y)^T S^{-1} (x - y)}$$

Pseudodistance Jaccard

For sets can use the index of Jaccard.
 Pour des ensembles on peut utiliser l'index de Jaccard.

For an ensemble x , $|x|$ est sa cardinalité.

$$J(x, y) = \frac{|x \cap y|}{|x \cup y|} = \frac{|x \cap y|}{|x| + |y| - |x \cap y|}$$

$$0 \leq J(x, y) \leq 1$$

$$x = \{\spadesuit, \clubsuit, \heartsuit, \diamondsuit, \circ, \dagger, \nabla, \infty, \div, \otimes, \oslash, \square\}$$

$$|x| = 12$$

$$y = \{\spadesuit, \clubsuit, \heartsuit, \diamondsuit, \nabla, \infty, \div, \otimes, \pitchfork, \oslash, \triangledown\}$$

$$|y| = 11$$

$$x \cap y = \{\spadesuit, \clubsuit, \heartsuit, \diamondsuit, \nabla, \infty, \div, \otimes, \oslash\}$$

$$|x \cap y| = 9,$$

$$J(x, x) = J(y, y) = 100\%$$

$$J(x, y) = \frac{9}{12 + 11 - 9} = 64\%$$

Hamming distance and edit distance

For character sequences, bit or byte, of the same length, can be used Hamming distance. utiliser la distance de Hamming.

$$H(x, y) = |\{i | x_i \neq y_i\}|$$

The edit distance counts the number of replacement, addition or deletion necessary to turn x into y .

- $O_1: x_i \rightarrow a$
- $O_2: x_i \rightarrow x_i \cdot a$
- $O_3: x_{i-1} \cdot x_i \cdot x_{i+1} \rightarrow x_{i-1} \cdot x_{i+1}$

The edit distance can be calculated using dynamic programming time dynamique en temps $O(n^2)$

$$x = 01100101$$

$$y = 01101010$$

$$H(x, y) = 4$$

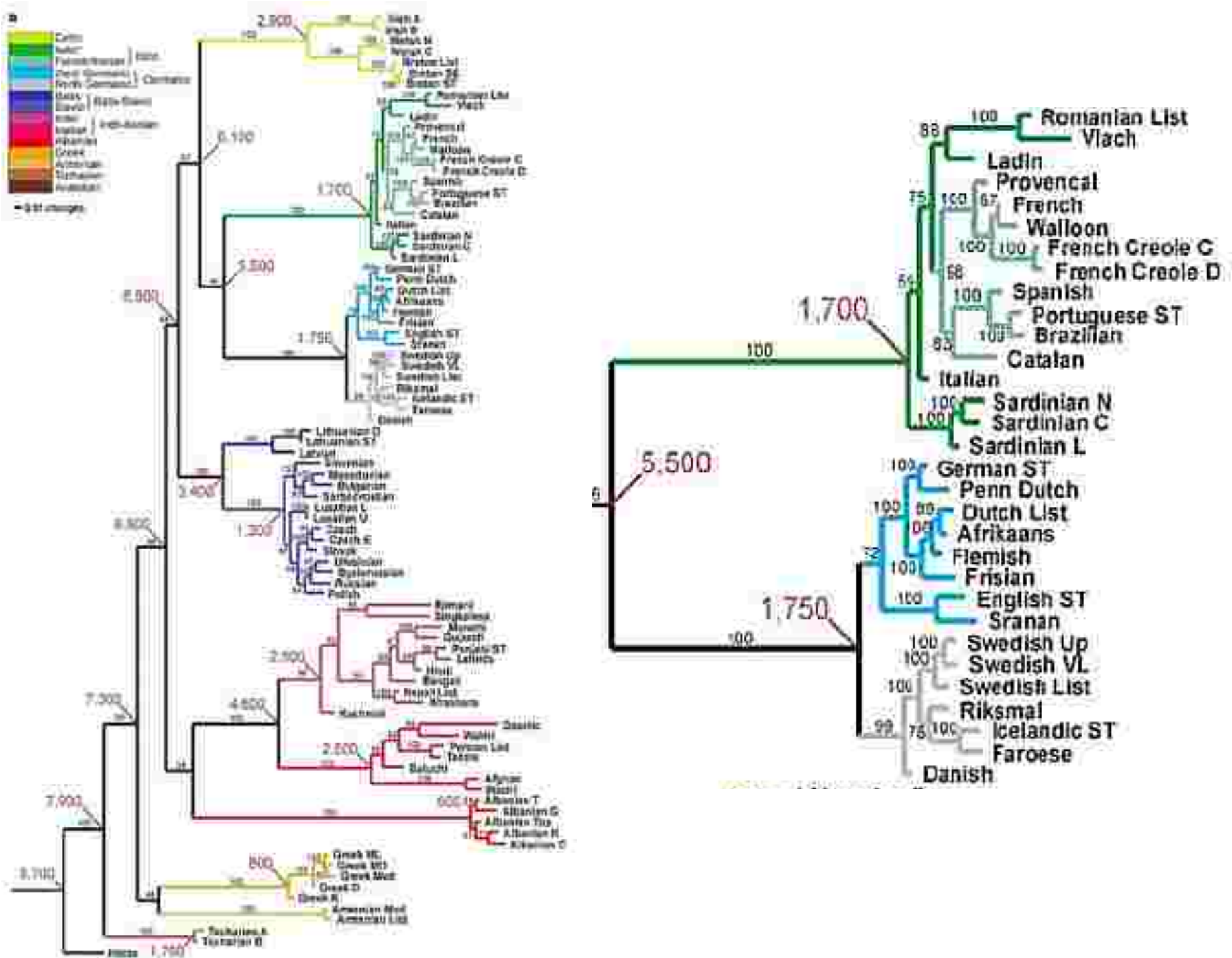
$$01100101 = x$$

$$0110101$$

$$01101010 = y$$

$$E(x, y) = 2$$

hierarchical Sheet

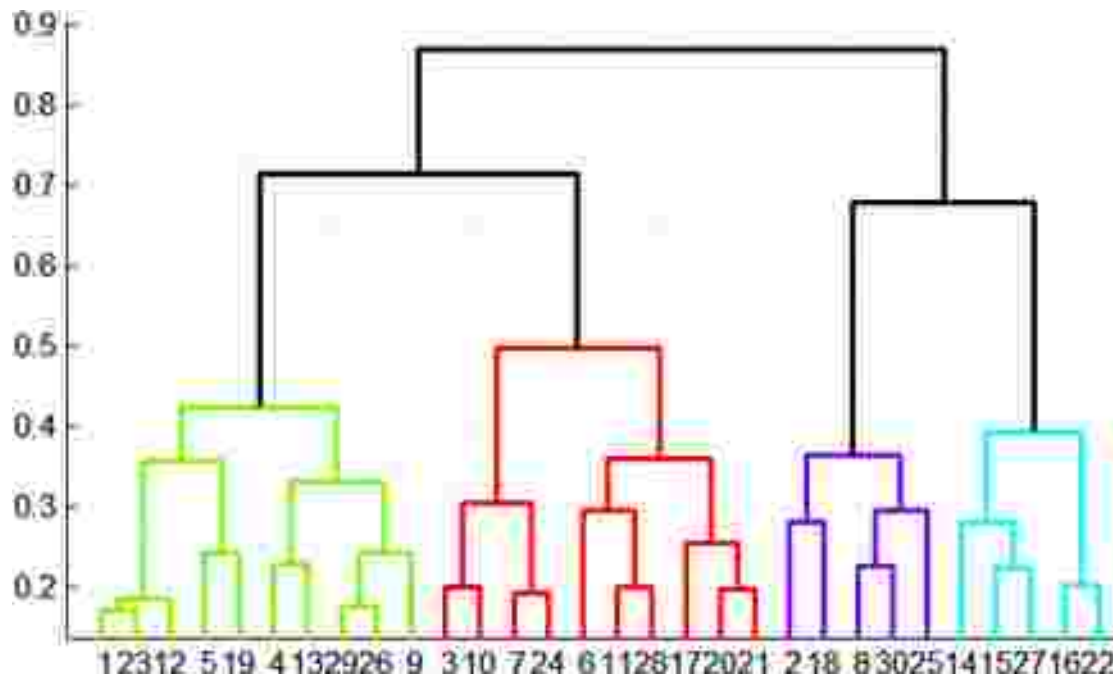


binary hierarchical partition

key operation: **Combine** the two closest groups.

Two important questions:

- How to represent a group?
- How do you determine the "closeness" of the groups?



hierarchical binary partition algorithm

$X = \{x_1, x_2, \dots, x_n\}$	Elements For	Les éléments
Pour $x_i \in U, I(x_i) = s_i \in S$	initialization	Initialisation
For $s_1, s_2 \in S, J(s_1, s_2) = s \in S$	Attach two groups	Joindre deux groupes
For $s_1, s_2 \in S, D(s_1, s_2) = d \in \mathbb{R}$	The distance between two groups	La distance entre deux groupes

CEBH(X)

$G = \{I(x_1), \dots, I(x_n)\}$

Finement: — 1 fois:

Find minimally Show
Trouver $s_i, s_j \in G$ avec $D(s_i, s_j)$ minimale

Afficher (i, j)

$s_i = J(s_i, s_j)$

Take back from
Retirer s_j de G

hierarchical partition algorithm

binary

The naive implementation is in $O(n^2)$ and an implementation using a priority queue can reduce the time but it is impractical for very large data sets.
L'implémentation naïve est dans $O(n^2)$ et une implémentation utilisant une file d'attente prioritaire peut réduire le temps à $O(n^2 \log n)$ mais cela reste impraticable pour les très gros jeux de données.

Example

$X = \{x_1, x_2, \dots, x_n\}$ avec les $x_i \in U$

Les éléments

Pour $x_i \in U, I(x_i) = s_i \in S$

Initialisation

Pour $s_1, s_2 \in S, J(s_1, s_2) = s \in S$

Joindre deux groupes

Pour $s_1, s_2 \in S, D(s_1, s_2) = d \in \mathbb{R}$

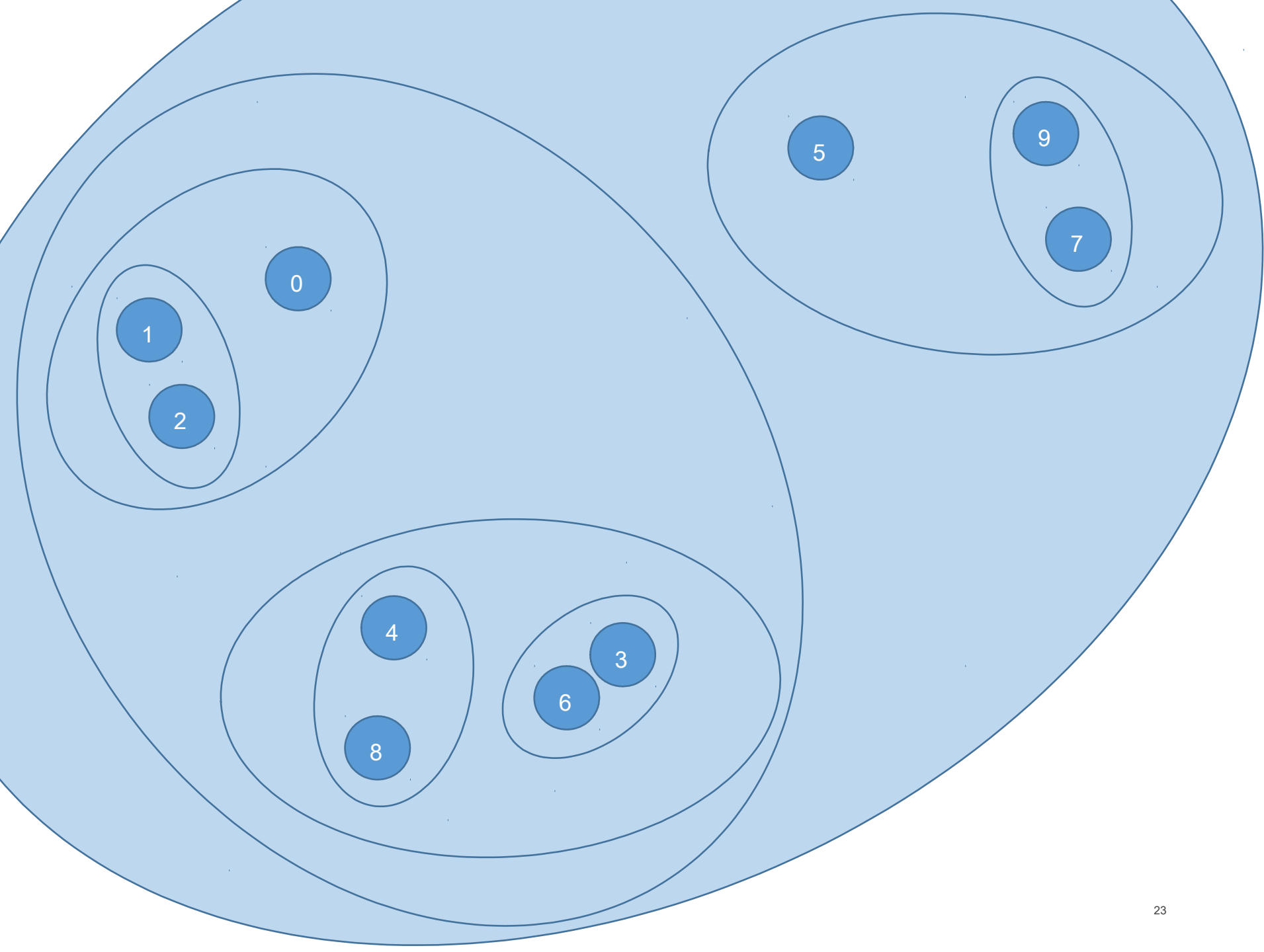
La distance entre deux groupes

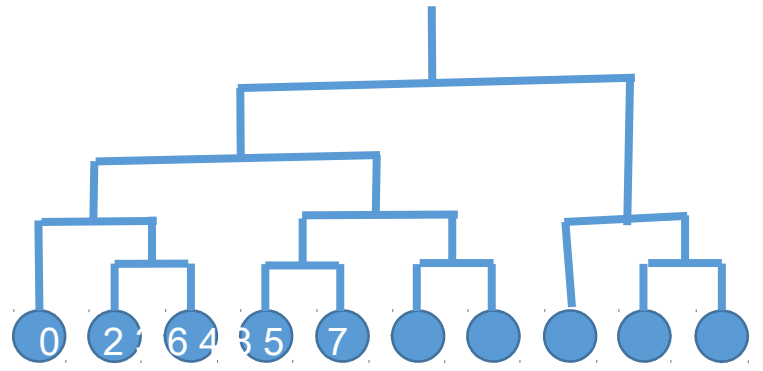
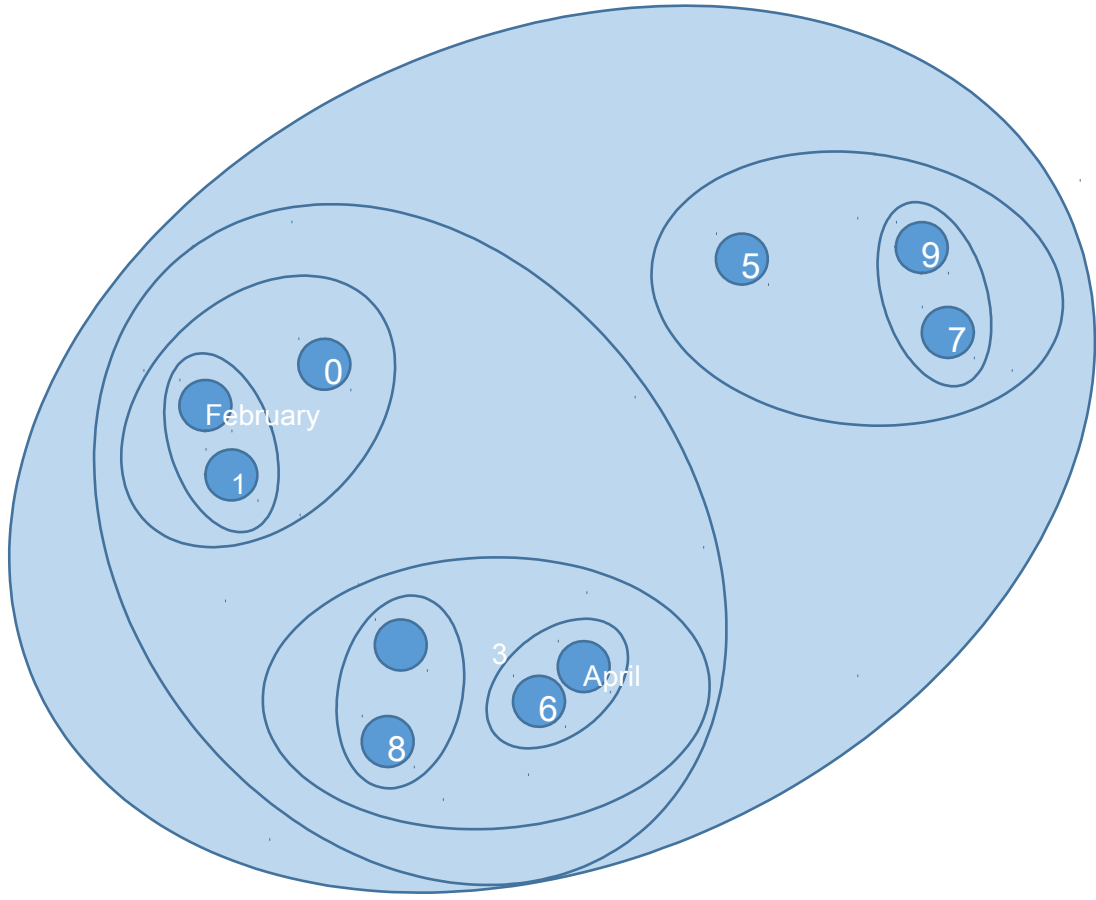
In the Euclidean case, a group can be represented by its cardinality and its mean.
Dans le cas euclidien, un groupe peut être représenté par sa cardinalité et sa moyenne.

$$I(x_i) = [1, x_i]$$

$$J([n_1, m_1], [n_2, m_2]) = \left[n_1 + n_2, \frac{n_1 m_1 + n_2 m_2}{(n_1 + n_2)} \right]$$

$$D([n_1, m_1], [n_2, m_2]) = L_2(m_1, m_2)$$





non-Euclidean case

We can not calculate average, but still has a notion of similarity between the elements.

approach 1

Choose a representative for each grouping. For example the one that minimizes the distance with the other members of the grouping. The similarity groups is obtained by calculating the similarity of Representatives.

approach 2

Use the mean similarity between peer element from each group.

approach 3

Select the combination that maximizes cohesion. Use the diameter of the merged group. The diameter is the maximum distance between the points of the group. Using the average distance between the points of the group. Density can also promote dividing by the number of points.

Partitions

Sheet

The set of elements is partitioned.
L'ensemble d'éléments a partitionné.

$$X = \{x_1, x_2, \dots, x_n\}$$

The set is partitioned into k groups.
L'ensemble est partitionné en k groupes.

$$\{G_1, G_2, \dots, G_k\}$$

$$X = \bigcup_{i=1}^k G_i$$

$$\forall i \neq j, G_i \cap G_j = \emptyset$$

Requires notion of similarity (or distance) between the elements.

Nécessite une notion de similarité (ou de distance) entre les éléments.

$$S: X \times X \rightarrow R$$

Useful to understand the data and that schematize and to make decisions.

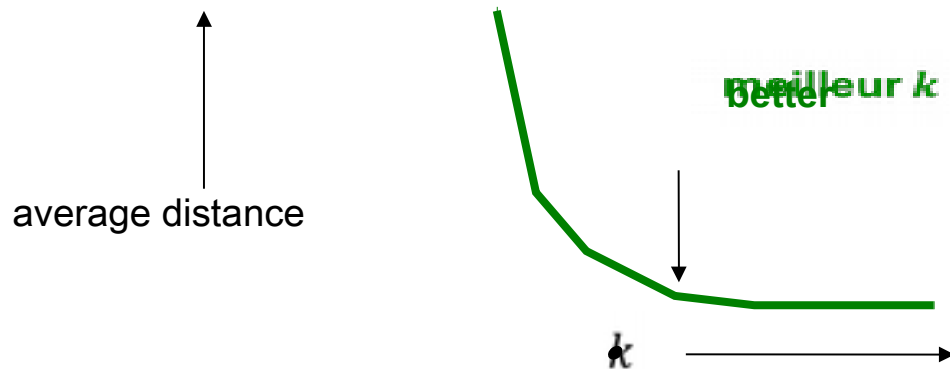
Utile pour comprendre les données et les schématiser ainsi que pour prendre des décisions. Utile pour appliquer une solution personnalité en fonction des groupes obtenue.

Choose the number of group (k)

Use group cohesion.

Use the diameter of the merged group.

Using the average distance between the points of the groups. Using the average distance between the points of the group and their centers.



expectation-algorithm maximization

The expectation-maximization algorithm is an iterative algorithm to find the parameters of maximum likelihood of a probabilistic model when it depends on unobserved latent variables.

The expectation-maximization algorithm includes:

- an evaluation expectancy step (E), which calculates the likelihood expectancy taking into account the last observed variables,
- a maximization step (M), wherein the maximum likelihood parameter is estimated by maximizing the likelihood found in step E.

It then uses the settings found in M as the starting point of a new phase of assessment of hope, and it iterates well.

***K*-moyenne
-average**

- ~~average~~ K-moyenne

When the data belong to a Euclidean space, k-mean algorithm allows a set divided into groups for a given k .
Lorsque les données appartiennent à un espace euclidien, l'algorithme k-moyenne permet de diviser un ensemble en k groupes pour un k donné.

$X = \{x_1, x_2, \dots, x_n\}$ avec les $x_i \in \mathbb{R}^d$

$D(x, y)$ la distance euclidienne sur \mathbb{R}^d

KMEAN(X, k)

1. Choisir $C = \{m_1, \dots, m_k\}$ avec les $c_i \in \mathbb{R}^d$ les centres pour les k groupes.
To choose centers for groups.

111 Pour i de 1 à k : $G_i = \{x \mid \forall j \neq i, D(x, m_i) \leq D(x, m_j)\}$
to:

2. Return 1 to 1 if the **stop condition** is not satisfied. Return
2. Pour i de 1 à k : $m_i = \frac{1}{|G_i|} \sum_{x \in G_i} x$

Retourner a 1 si la **condition d'arrêt** n'est pas satisfaite.

Retourner $\{G_1, \dots, G_k\}$

- ~~average~~ K moyenne

Choisir l'initialisation des centres
To choose initializing centers

Choisir k élément uniformément au hasard et affecter les k centres avec ces k valeurs.
Choose items randomly and evenly assigned centers with these values.

Choisir un élément au hasard, pour les $k - 1$ élément suivant, choisir à maximise la distance qu'ils ont avec les points déjà choisis.
Select an item at random, for next item, choose each time point that maximizes the distance to the points already selected.

Stop condition of the algorithm.

The new centers are identical to those of the previous iteration précédente.

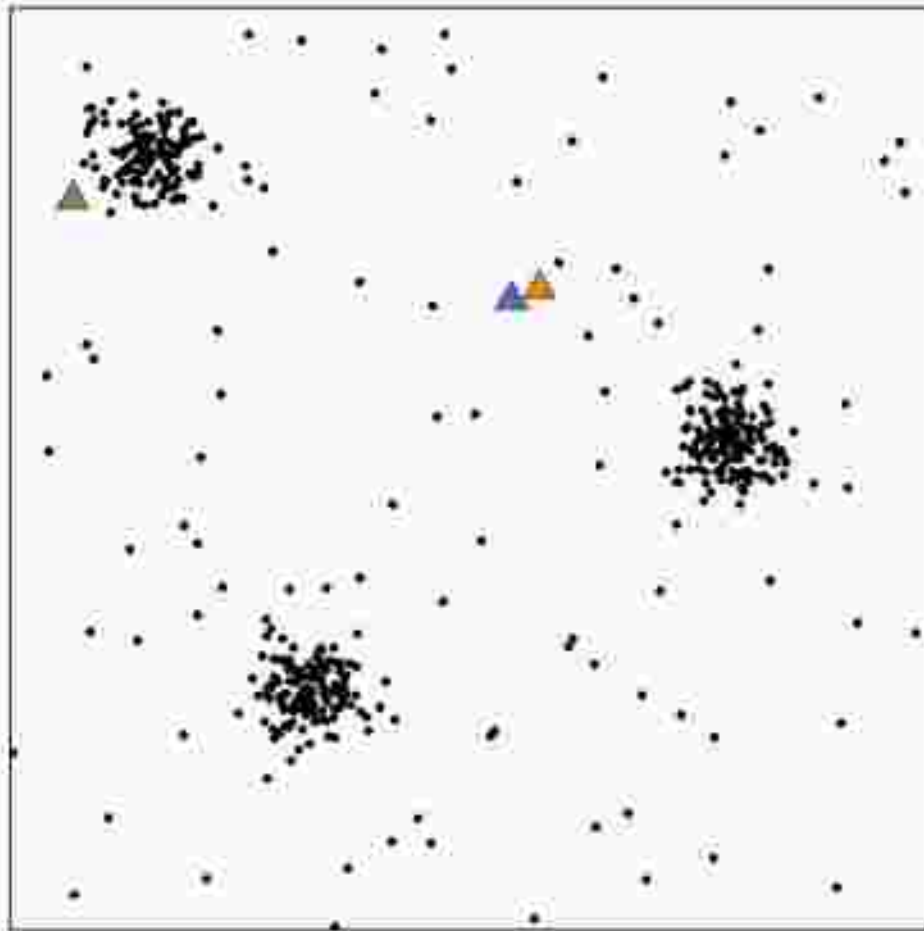


Bishma Stornelli

Published on March 22, 2014

A simple example of a real-time simulation of the K-Means Clustering Algorithm using different values for n and k .

Visualizing K-Means Clustering



Mean square point-centroid distance: not yet calculated

The k-means algorithm is an iterative method for clustering a set of N points (vectors) into k groups or clusters of points.

Algorithm

Repeat until convergence:

Find closest centroid

Find the closest centroid to each point, and group points that share the same closest centroid.

Update centroid

Update each centroid to be the mean of the points in its group.

Find closest centroid

Data

Clustered points ☐ Random

Number of clusters:

Number of centroids:

New points

New centroids

256 images and $k = 22$ groups

3	7	7	0	1	3	4	0	3	1	1	4	8	2	9	6
6	0	3	9	5	1	5	5	4	8	4	4	0	7	8	0
1	7	0	7	1	0	0	3	6	4	2	8	7	3	1	7
7	1	4	5	8	8	0	4	4	6	0	5	6	8	4	1
5	2	1	0	3	9	5	3	2	3	9	7	3	7	4	7
9	9	3	6	3	3	9	7	0	8	7	3	7	7	4	9
3	0	9	5	5	6	2	0	4	8	0	8	1	0	3	3
2	1	8	6	3	8	3	3	1	8	7	5	7	7	2	6
8	9	1	6	1	2	5	5	3	5	2	8	7	5	4	0
0	0	7	2	9	0	5	1	7	1	2	8	3	5	3	1
6	3	0	4	3	2	9	5	5	8	9	1	8	5	1	2
8	7	9	8	6	4	8	7	6	6	6	6	5	6	7	1
6	2	9	0	0	4	7	6	8	7	1	1	0	7	7	6
7	6	3	0	7	4	3	6	8	9	8	0	6	8	0	5
4	3	7	4	0	1	4	0	2	8	7	6	7	8	4	7
3	1	5	4	9	6	0	4	0	4	1	2	0	3	2	0

iteration 1

8 8 8 8 8 7 9 5 8

2 2 2 2 4

6 6 2

4 4 9 4 7 7 4 4 7 7 4 4 7 9 3 7 4 6 9 6 7 7 7 8 7 2 4 7 7 7 7 8 8

8 8 5 3 5

0 0 7 0 0 5 9 7 5 2 9

1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 7 7 4 6 8 5 1 1 1 1 1 1 8 1 2 3 9 4 5

2 2 5

0 0 0 0 0 0 3 0 0 0

6 6 6 6 6 6 4 4 6 2 4 4 4 2

3 3 3 3 3 3 3 3 8 3 3 3 8 5 3 2 3 2 5 9 8 4 8 5 8

7 7 7 7 7 7 9

6 6 6 6 6 6 6 6 6 6 6

3 3 3 3 5 5 6 5 3 8 5 3 3 6 3 1 5 5 5 8

3 3 3 3

6 6 4 0

0 0 0 0 0 0 0 0 0 0 0 0 0 5 0 0

3 3 1 8 3 2 8 3 2 8 2 2 2 0

9 9 9 7 8 4 7 7 7 0 9 7 7 7 7 4 4 4 4

7 7 9 9 9 9 7 4 9 4 8

5 5 5 5 5 1 1 1 1 1 7 5 5 2 0

0 0 0 0 0

iteration 2

8 8 8 8 8 8 8 8 7 8
 2 2 2 2 9 2
 6 6 2
 4 4 9 7 7 7 4 7 7 7 4 4 7 9 9 4 9 7 3 4 7 7 4 4 7 7 9 7 7 9
 5 8 5 5 3 5
 0 0 7 9 9 5
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 4 1 6 7 1 5 1 7 6 6 7
 2 2 2
 0 0 0 0 0 0 0 0 3 0 0 0
 4 4 6 6 6 4 6 6 4 6 4
 3 3 3 3 3 3 3 3 3 3 8 3 3 5 3 8 2 9 5 9 8 2 3 5 8 2 8
 7 7 7 7 7 7 9 7
 6 6 6 6 6 6 6 6 6 6 6 6
 3 3 3 3 5 5 3 6 3 8 5 5 3 1 3 5 6 3 5 5 8
 3 3 3 3
 0 0 4 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 5 0
 3 3 8 3 8 2 8 2 1 2 2 2 2
 9 9 9 9 0 9 7 4 7 7 9 7 7 2 4
 9 9 9 9 4 9 7 9 4 4 4 7 4 4 4 4 8
 7 1 1 3 1 5 1 1 5 7 5 3 5 5 7 2
 0 0 0 0 0 0

iteration 5

8	8	8	8	8	8	8	8	7	8	8	8																		
2	2	2	2	2	2																								
6	6	2																											
7	7	9	7	9	7	7	4	7	4	7	4	7	9	7	4	7	4	7	4	4	3	7	7	9	7	7	7	7	9
5	5	5	5	8	5																								
9	9	7	9	9	5																								
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	7	1	1	4	9	6	1	7	5	1	6	6	1	
2	5	2																											
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0														
6	6	6	6	6	6	6	6	6	6	4	4																		
3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	3	3	8	8	8	5	8	6	8	5	8	1	3	8	2
7	7	7	7	7	7	9	7																						
6	6	6	6	6	6	6	6	6	6	6	6	6	6																
8	3	3	5	5	3	3	6	3	5	5	5	5	5	3	8	6	3	3	8										
3	3	3	3																										
10	4	4	0																										
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0												
2	2	8	8	2	8	2	3	2	8	3	2	2	2																
7	7	9	9	7	0	9	7	7	4	7	9	7	7	4															
9	4	9	9	9	9	4	9	4	4	4	4	4	7	4	4	7	4	4	8										
8	1	1	1	5	1	5	1	7	5	5	3	5	5	2	0														
0	0	0	0	0	0																								

iteration 8

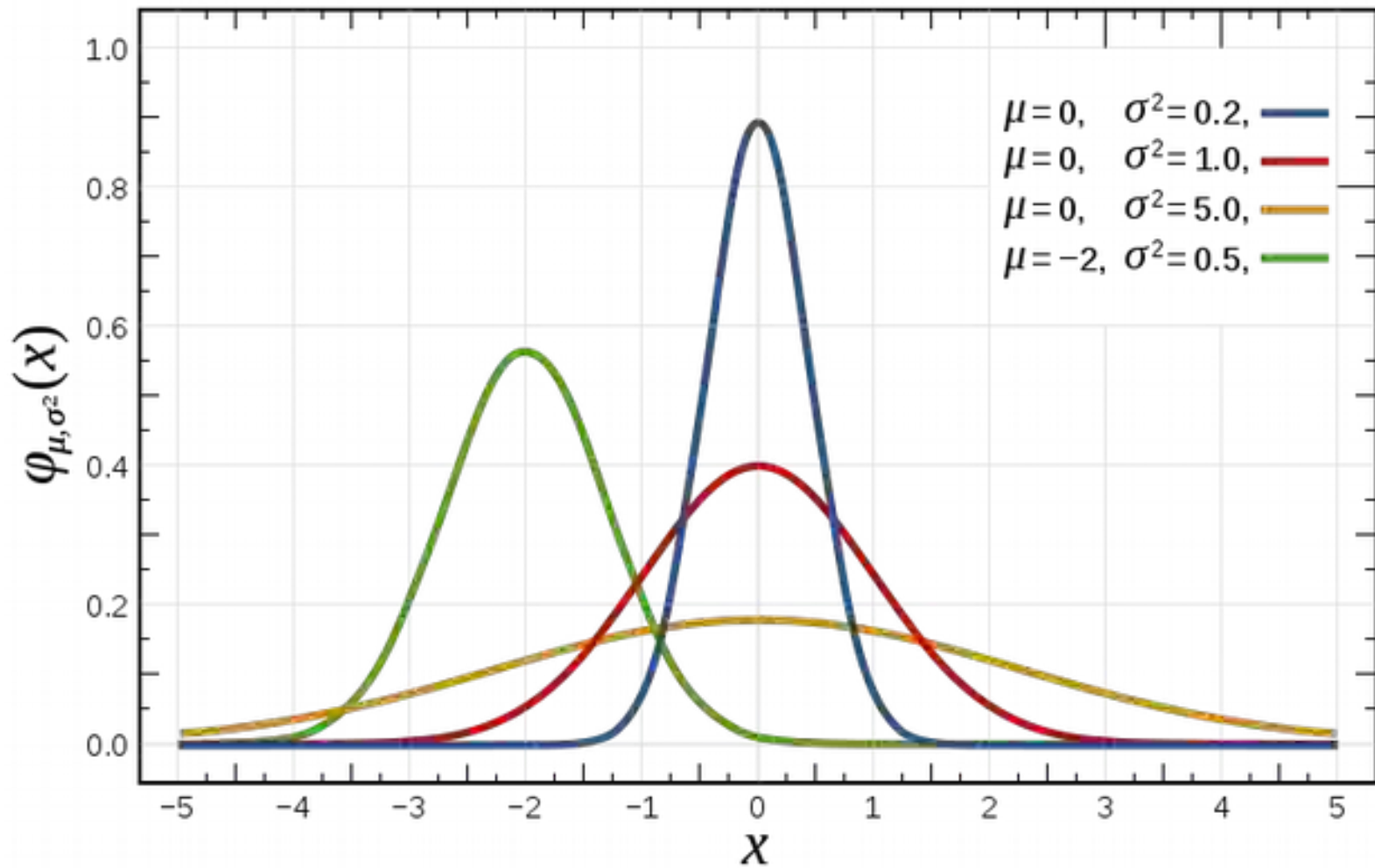
8	8	8	8	8	8	8	7	8	8	8																					
2	2	2	2	2	2																										
6	6	2																													
7	7	9	7	4	7	7	4	7	4	7	4	7	9	7	4	7	4	4	3	7	7	9	7	7	7	9					
5	5	5	5	8	5																										
9	0	7	9	9	5																										
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	7	1	1	4	8	0	1	7	5	1	6	6	2			
2	8	2																													
0	0	0	0	0	0	0	0	0	0	0	3	0																			
6	6	6	6	6	6	6	6	6	4	4																					
3	3	3	3	3	3	3	3	5	5	3	3	3	3	3	3	3	8	3	5	9	8	8	8	8	8	5	9	1	3	8	2
7	7	7	7	7	7	7	9	7																							
6	6	6	6	6	6	6	6	6	6	6	6	6																			
5	5	3	3	3	5	5	3	6	5	8	5	3	6	3	3	8															
3	3	3	3																												
10	10	4	0																												
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0																
2	3	8	8	2	8	2	3	2	8	3	2	2																			
9	7	9	9	7	0	9	7	7	4	7	9	7	7	4																	
4	4	9	9	9	9	4	9	4	4	4	4	7	4	4	7	4	4	8													
8	1	1	1	5	1	5	1	7	5	5	7	5	5	2	0																
0	0	0	0	0	0																										

GMM

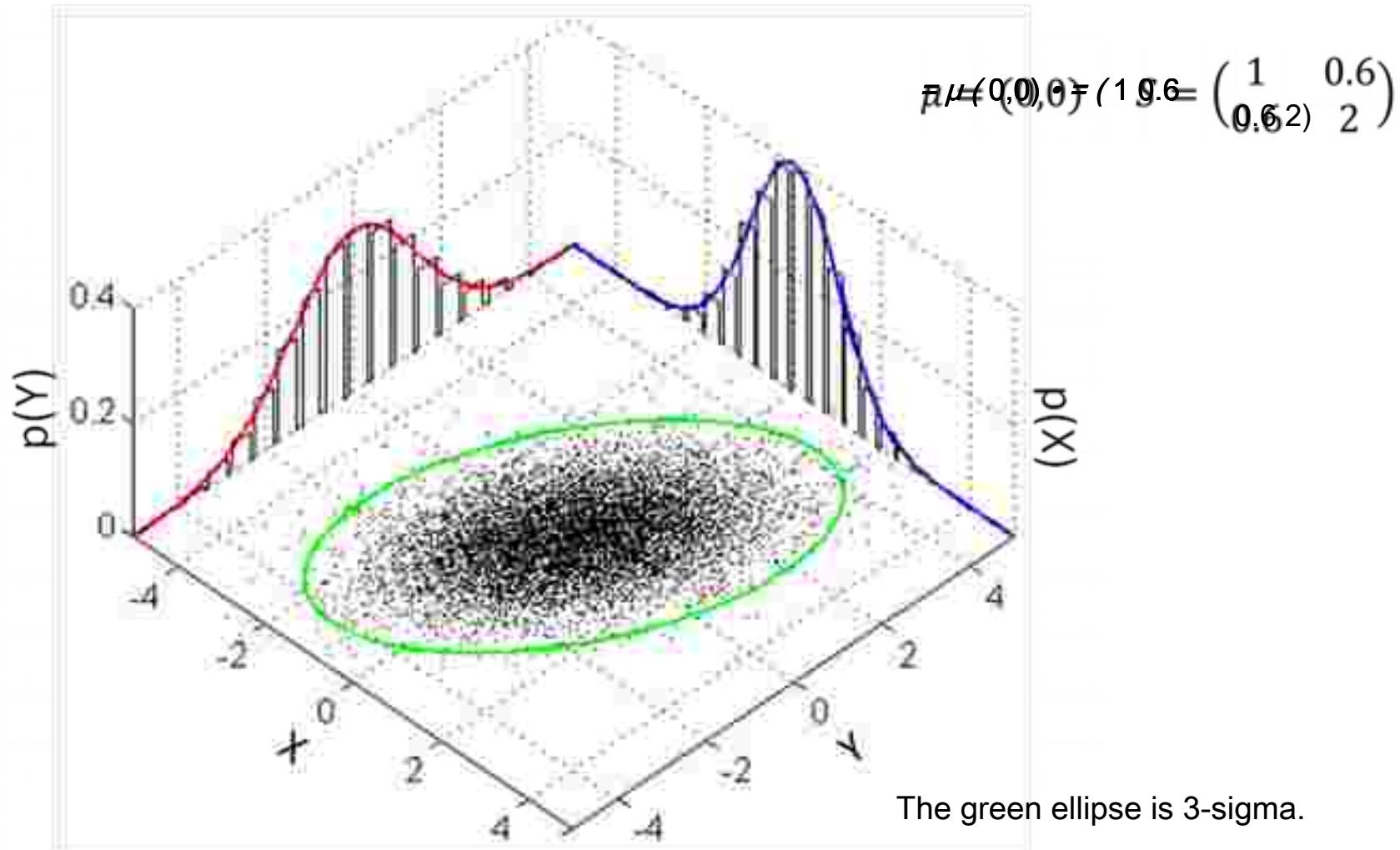
Gaussian Mixture Model

normal mixture

Normal law



multivariate normal distribution



statistical Reminder

the data

$$X = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^n$$
$$Y = \{y_1, y_2, \dots, y_n\} \in \mathbb{R}^n$$

The average

$$\mu_X = \frac{1}{n} \sum_{l=1}^n x_l$$

The variance

$$\sigma_X = \frac{1}{n} \sum_{l=1}^n (x_l - \mu_X)^2$$

covariance

$$\sigma_{X,Y} = \frac{1}{n} \sum_{l=1}^n (x_l - \mu_X)(y_l - \mu_Y)$$

Correlation coefficient

$$r_{X,Y} = \frac{\sigma_{X,Y}}{\sigma_X \sigma_Y}$$

covariance matrix

$X = \{x_1, x_2, \dots, x_n\}$ avec les $x_i \in \mathbb{R}^d$

The Cov $\text{Cov}(X)_{ij} = \sigma_{i,j}$

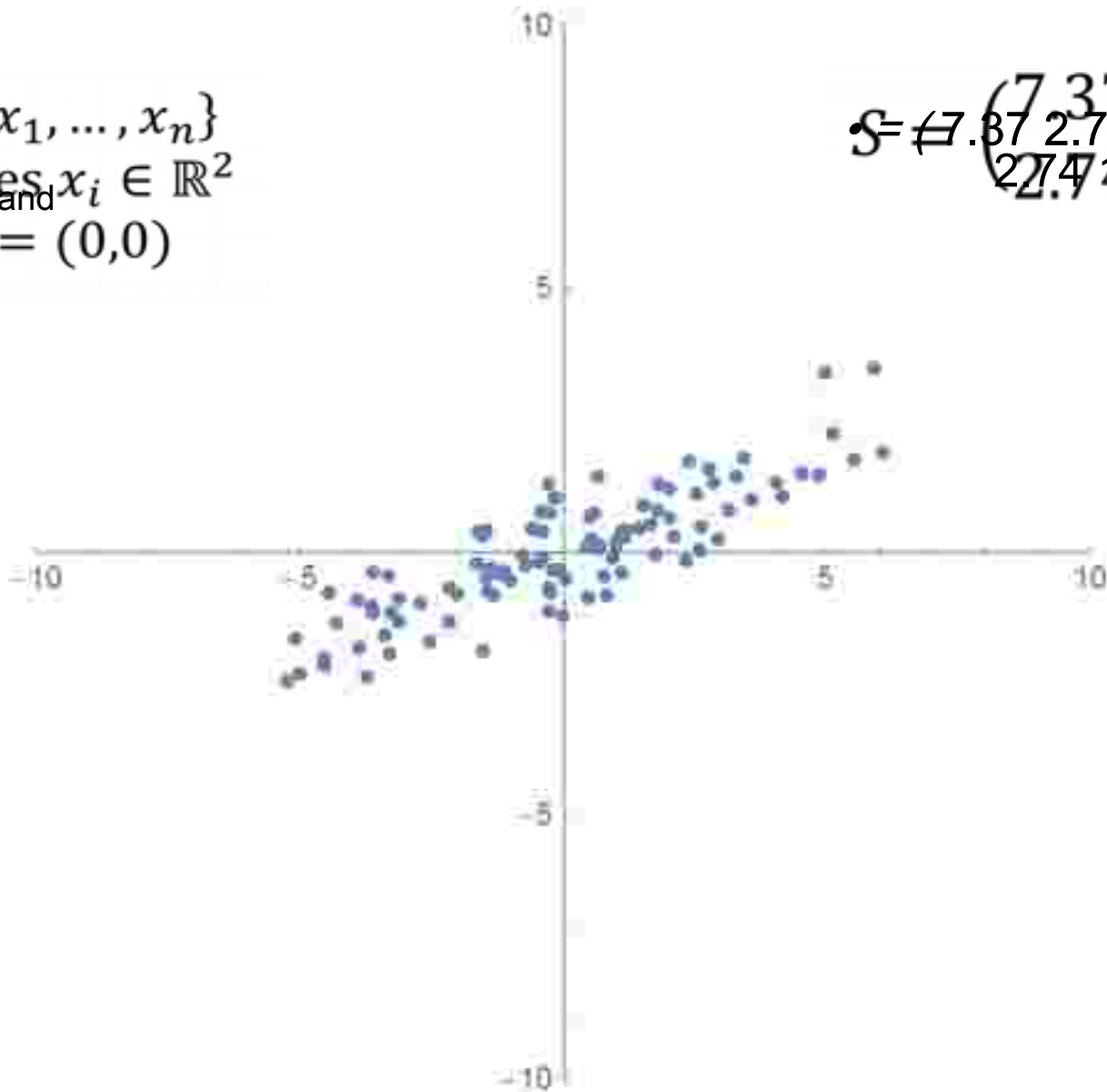
$$d = 3$$

$$\begin{pmatrix} \sigma_{1,1} & \sigma_{1,2} & \sigma_{1,3} \\ \sigma_{2,1} & \sigma_{2,2} & \sigma_{2,3} \\ \sigma_{3,1} & \sigma_{3,2} & \sigma_{3,3} \end{pmatrix} = \begin{pmatrix} \sigma_1 & \sigma_{1,2} & \sigma_{1,3} \\ \sigma_{1,2} & \sigma_2 & \sigma_{2,3} \\ \sigma_{1,3} & \sigma_{2,3} & \sigma_3 \end{pmatrix} \quad \begin{matrix} \sigma_{1,2} & \sigma_{1,3} \\ \sigma_{2,3} \\ \sigma_3 \end{matrix}$$

covariance

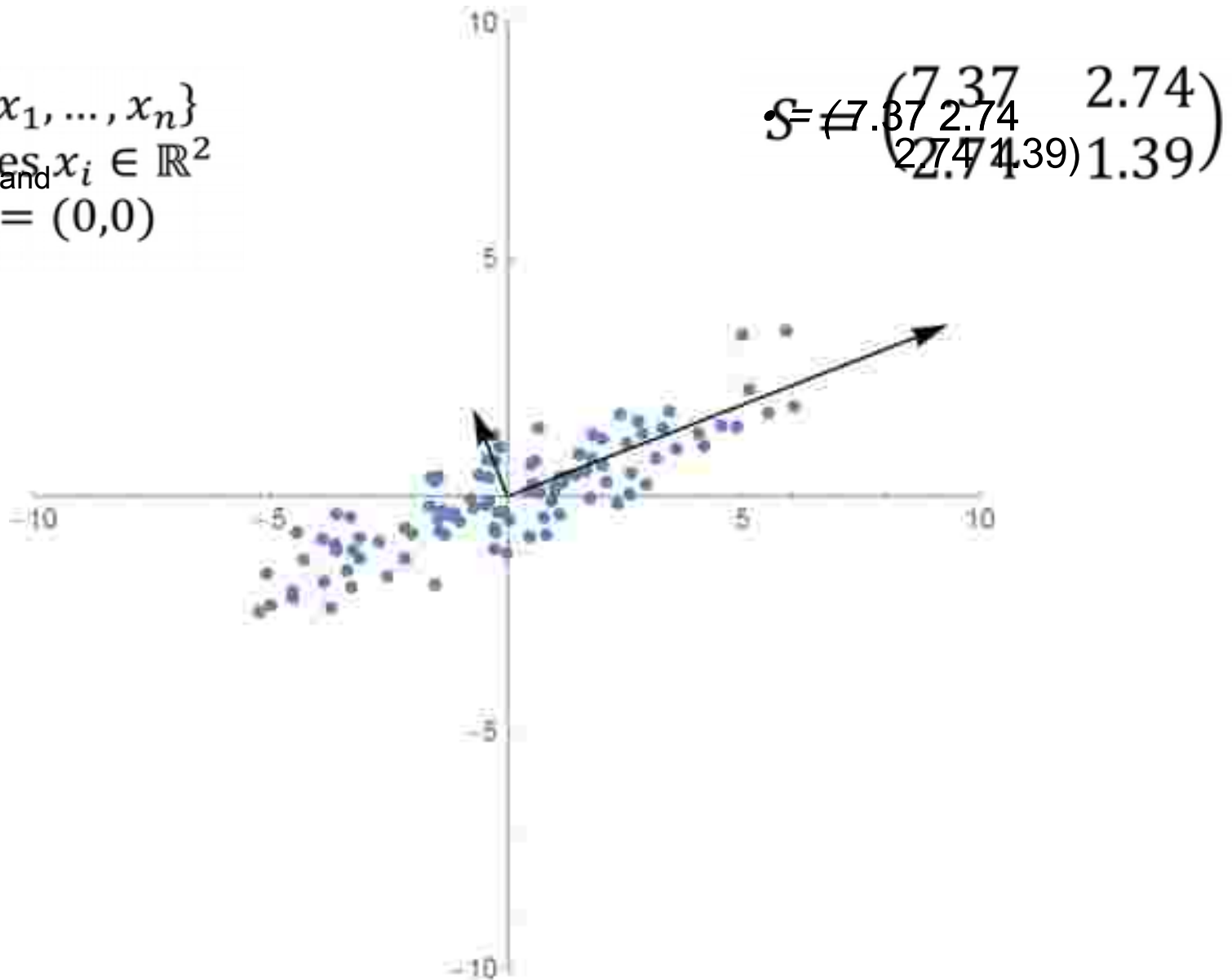
$X = \{x_1, \dots, x_n\}$
 avec les $x_i \in \mathbb{R}^2$
 and
 et $\mu_X = (0,0)$

$$S = \begin{pmatrix} 7.37 & 2.74 \\ 2.74 & 1.39 \end{pmatrix}$$



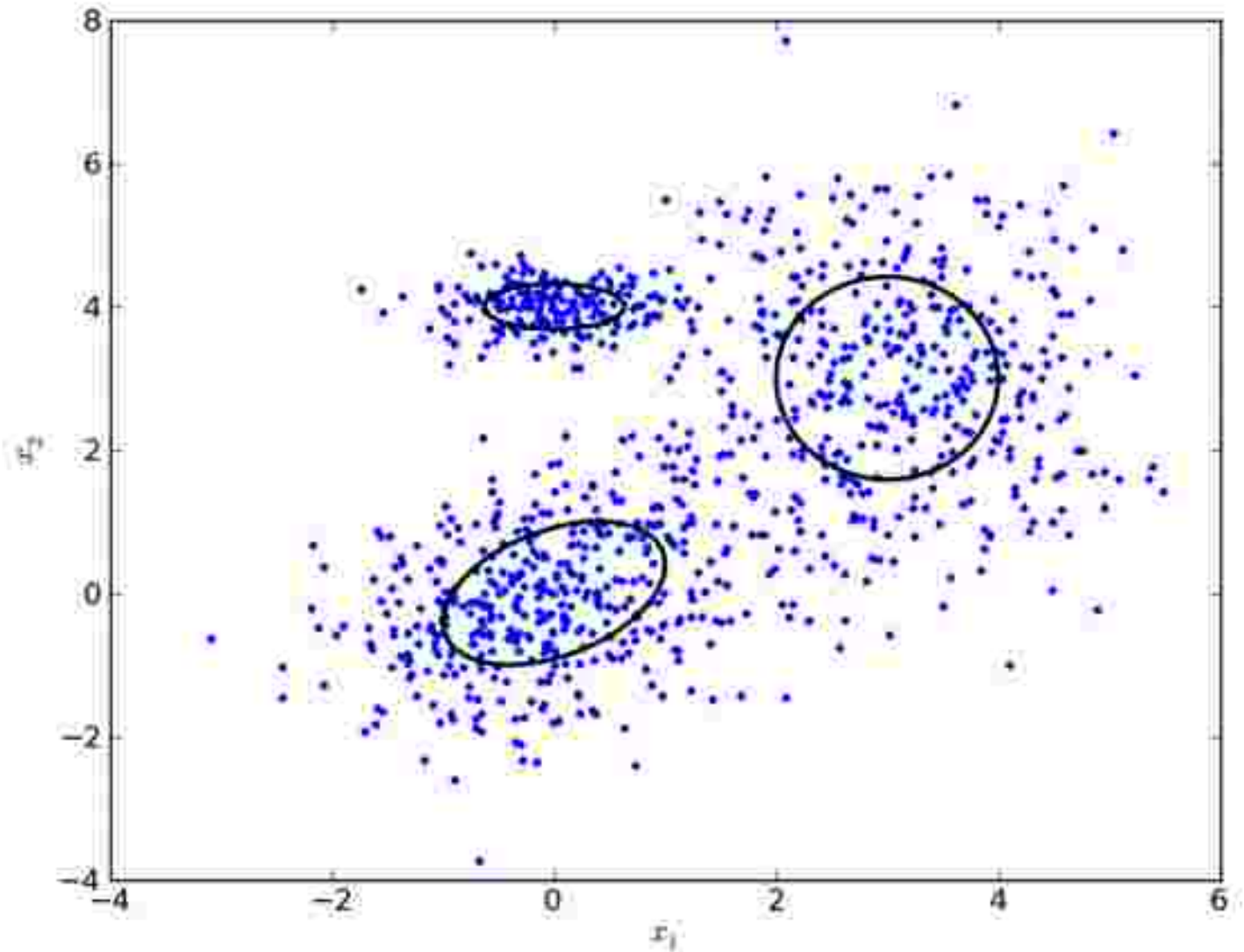
The proper vectors

$X = \{x_1, \dots, x_n\}$
 avec les $x_i \in \mathbb{R}^2$
 and
 et $\mu_X = (0,0)$

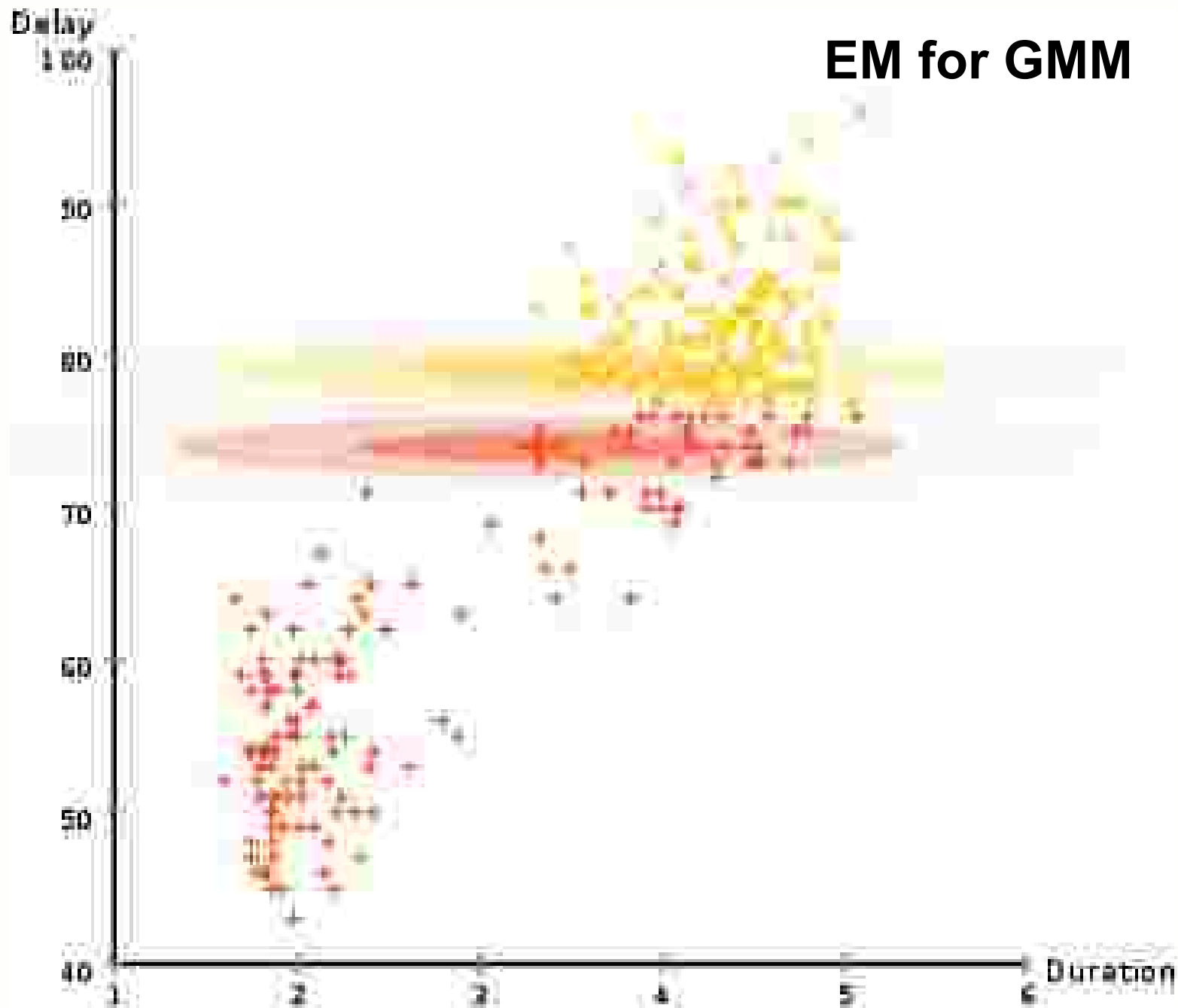


Gaussian Mixture Model

normal mixture



EM for GMM



GMM

$$X = \{x_1, x_2, \dots, x_n\} \text{ avec les } x_i \in \mathbb{R}^d$$

$$\mathcal{G} = \{(p_1, \mu_1, S_1), (p_2, \mu_2, S_2), \dots, (p_k, \mu_k, S_k)\}$$

$$\sum_{i=1}^k p_i = 1$$

$$\mathcal{N}(x|\mu, S) = \frac{1}{\sqrt{(2\pi)^d |S|}} \exp\left(-\frac{1}{2}(x - \mu)^T S^{-1}(x - \mu)\right)$$

Vraisemblance d'un élément

Likelihood of an element

$$v(x) = \sum_{i=1}^k p_i \mathcal{N}(x|\mu_i, S_i)$$

Trouver \mathcal{G} qui maximise la vraisemblance des donnés

Find that maximizes the likelihood given $\sum_{i=1}^n \ln \mathcal{N}(x_i)$

La probabilité pour un point x_i d'appartenir au groupe g est donné par

The probability for a point to belong to the group is given by

$$p_{i,g} = \frac{p_g \mathcal{N}(x_i|\mu_g, S_g)}{v(x_i)}$$

EM for GMM

$$X = \{x_1, x_2, \dots, x_n\} \text{ avec les } x_i \in \mathbb{R}^d$$

$$\mathcal{G} = \{(p_1, \mu_1, S_1), (p_2, \mu_2, S_2), \dots, (p_k, \mu_k, S_k)\}$$

GMM(X, k)

(GMM) Choisir uniformément $(\mu_1, \mu_2, \dots, \mu_k)$, un sous-ensemble de X
 Initially choosing a subset of the data to initialise and

Espérance

Hope Pour tout x_i et g calculer $\gamma_{i,g}$ la probabilité que x_i soit dans g .

For and calculate the probability that either.

$$\gamma_{i,g} = \frac{p_g v(x_i | \mu_g, S_g)}{v(x_i)}$$

Maximisation

Maximization

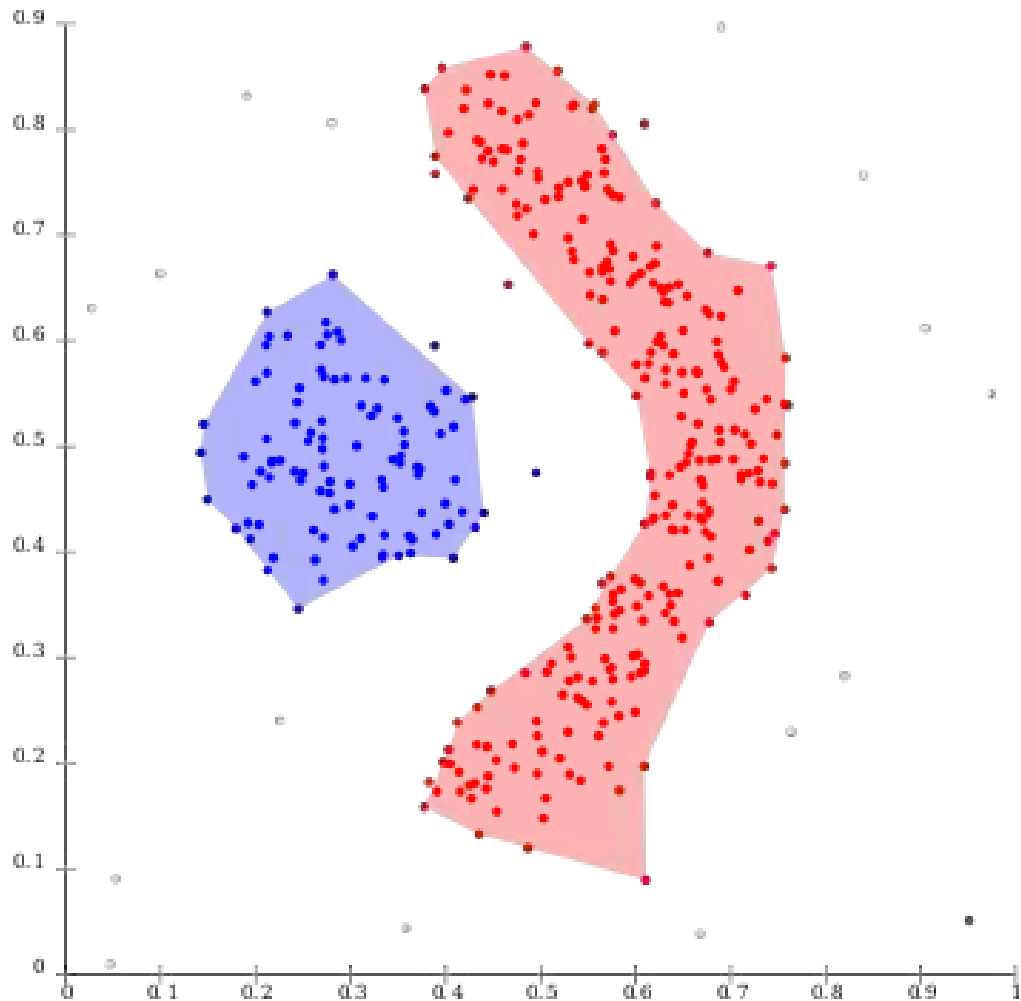
Réévaluer les paramètres de \mathcal{G}
 Reassess parameters For For For

Pour tout $g, p_g = \frac{1}{n} \sum_{i=1}^n \gamma_{i,g}$

Pour tout $g, \mu_g = \frac{\sum_{i=1}^n \gamma_{i,g} x_i}{\sum_{i=1}^n \gamma_{i,g}}$

Pour tout $g, S_g = \frac{\sum_{i=1}^n \gamma_{i,g} (x_i - \mu_g)(x_i - \mu_g)^T}{\sum_{i=1}^n \gamma_{i,g}}$

dbscan algorithm



Dbscan can find separable groups with disparate forms. This data set can not be partitioned appropriately with Kmean GMM.

Scores based on density

Dbscan is an approach to the partition based on the density. Main characteristics:

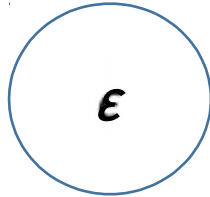
- Discover arbitrary form groups
- Manage Noise
- Effective
- Depends density parameters

dbscan

- ϵ is the radius of the neighborhoods
• ϵ est le rayon des voisinages
- Density is the number of points in a neighborhood.
• La densité est le nombre de points dans un voisinage.
- Δ is the critical density
• Δ est la densité critique
- Points **internal** is a point with the critical density in its vicinity.
• Points **internes**, est un point avec la densité critique dans son voisinage.
- A point **border** is not an internal point, but has an internal point in its vicinity.
• Un point de **bordure** n'est pas un point interne, mais possède un point interne dans son voisinage.
- The other points are points **noise**.
• Les autres points sont des points de **bruit**.

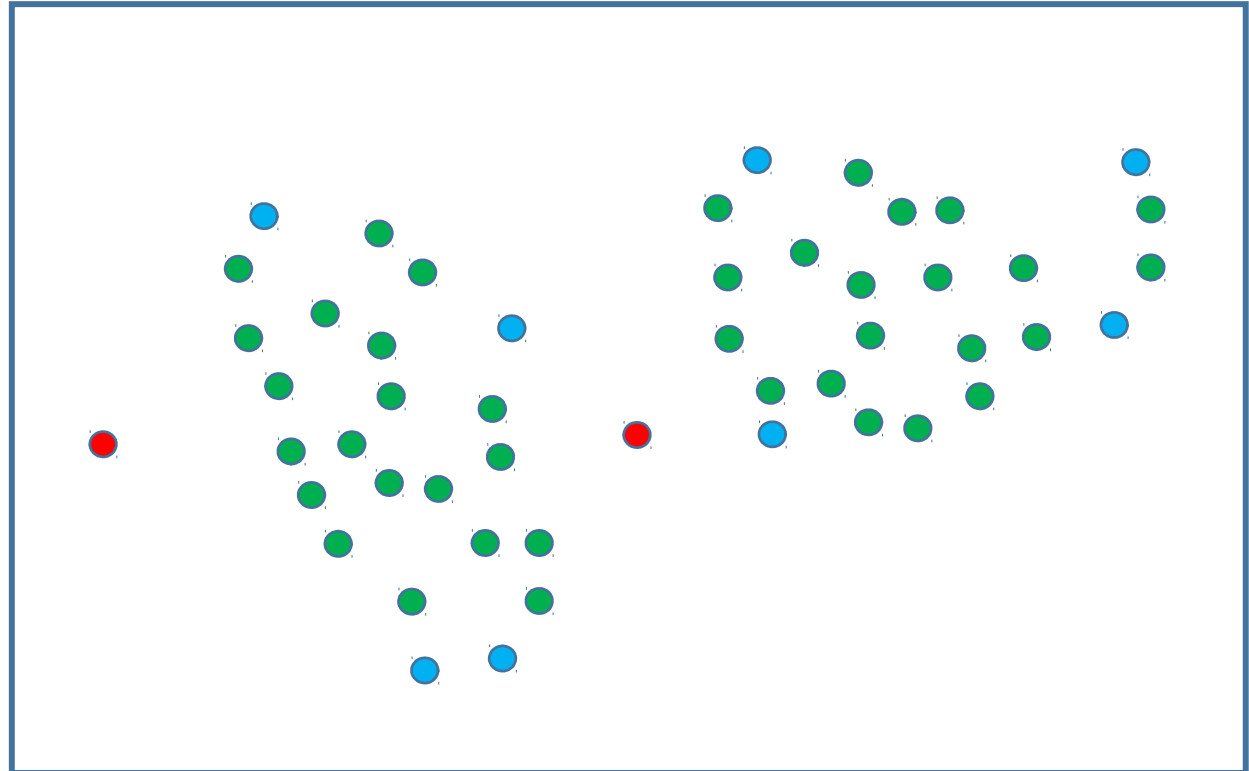
Dbscan: internal, border, and noise

$$\Delta = 4$$



2 groups

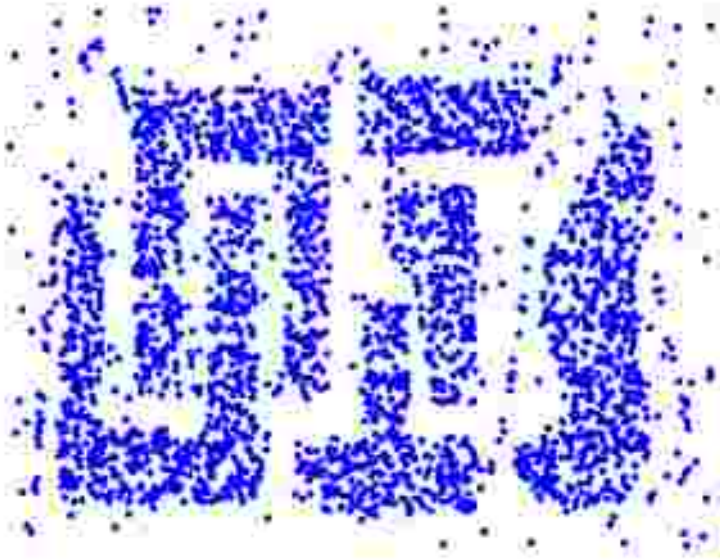
- noise
- edge
- inner



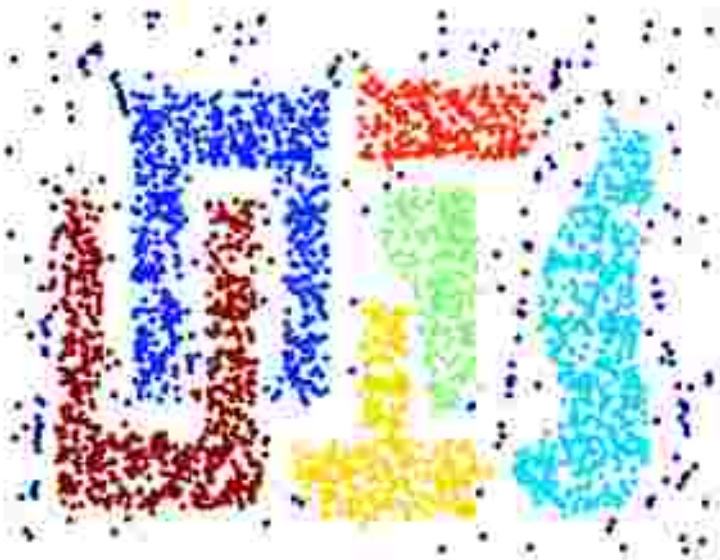
Dbscan algorithm (simplified)

1. Select and. Δ et ϵ .
2. Create a graph whose nodes are the points to regroup.
créer un graphe dont les nœuds sont les points à regrouper.
3. For each item, make a stop at each element in the vicinity.
Pour chaque élément x , créez une arête de x à chaque élément dans le voisinage.
4. Until all elements
tant qu'il reste des éléments
 - Select an item, if type noise is removed. Or create a group with all the elements in its vicinity and remove the group overall.
choisir un élément, s'il est de type bruit on l'enlève. Sinon créer un groupe avec tous les éléments dans son voisinage et retirer le groupe de l'ensemble.

Example where dbSCAN works well

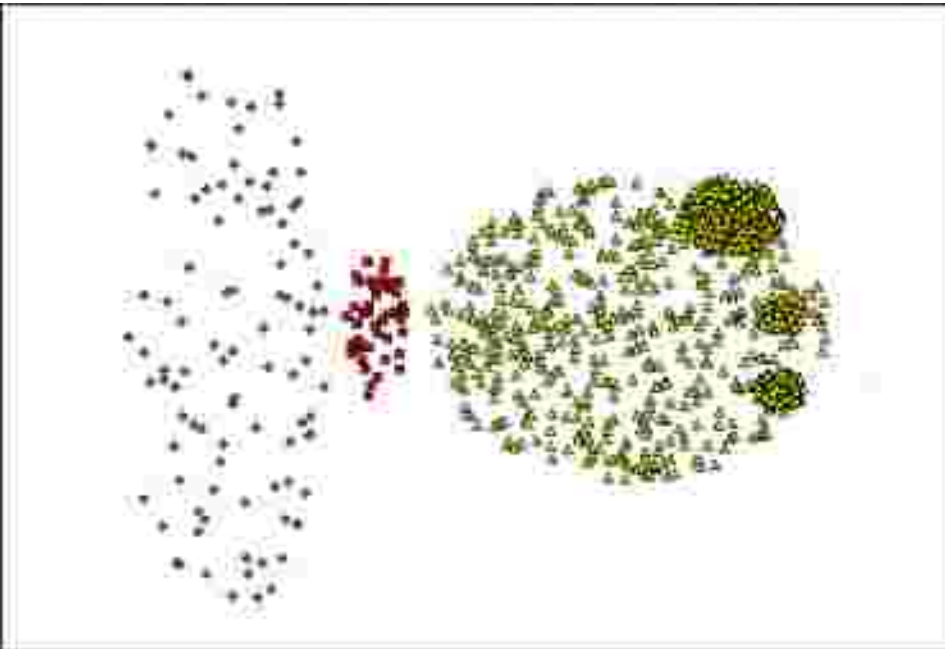


internal , border and noise

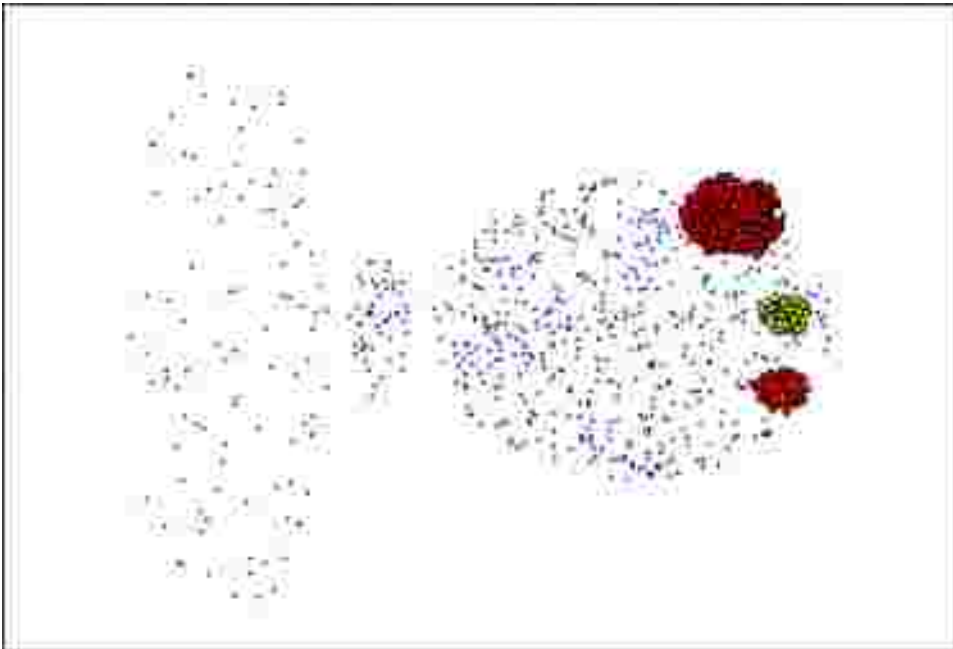


6 groups and noise

Example where dbscan works not good



$$\Delta = 4 - \epsilon = 9.75$$



$$\Delta = 4 - \epsilon = 9.92$$