$(k = 3)$
K-means

$(\varepsilon, n)$
DBSCAN

$\chi$

hard

Clustering

Eval.

Soft

HDBSCAN

Gaussian
Mixture
Model

DATA

# Clustering algs:

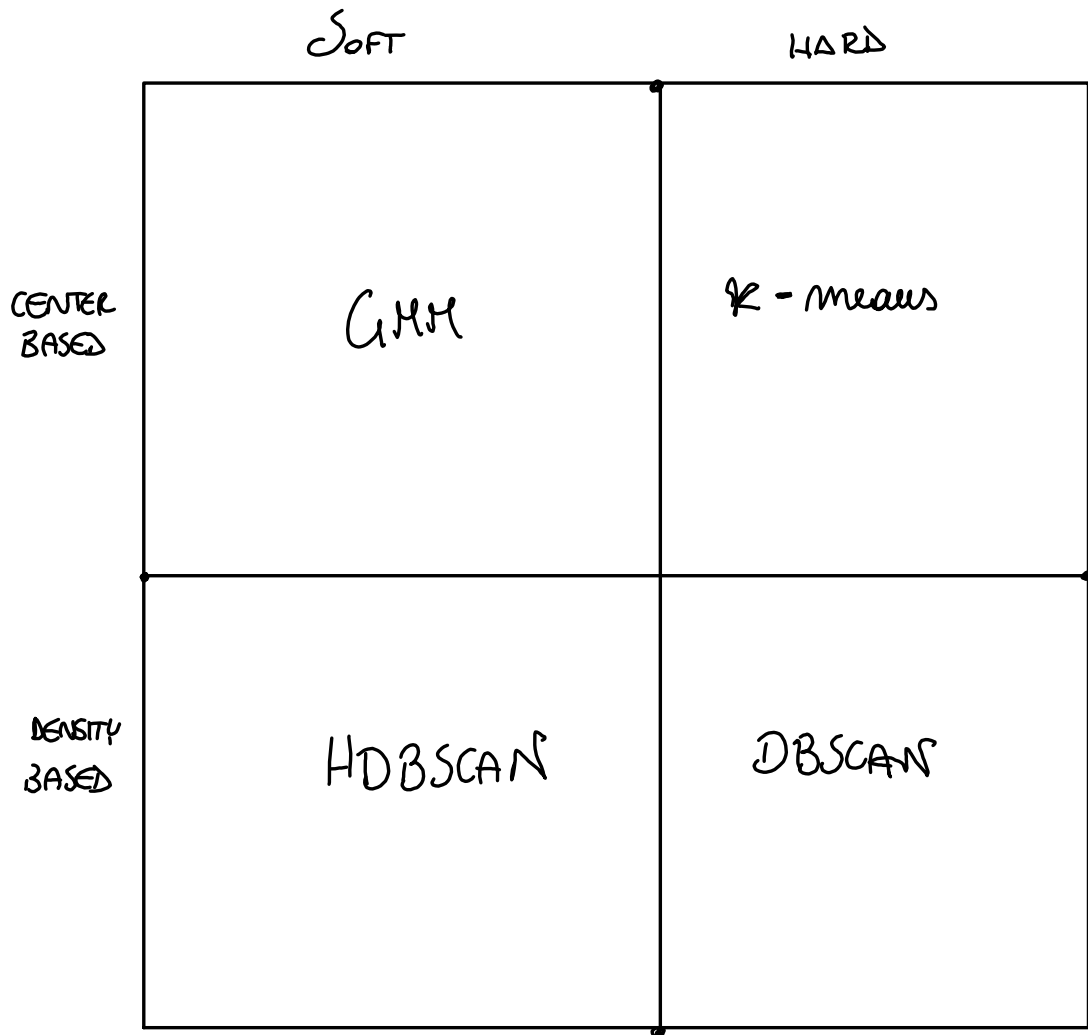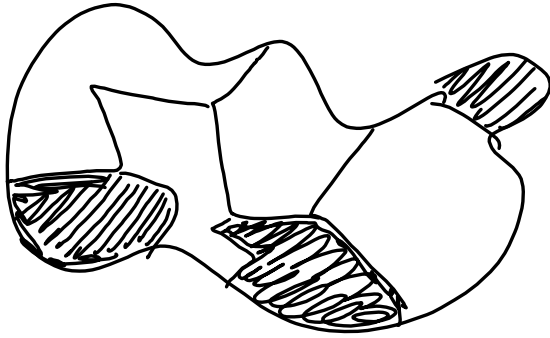|  | Soft | Hard |
|---|---|---|
| **Center Based** | GMM | K-means |
| **Density Based** | HDBSCAN | DBSCAN |

# K-Means:
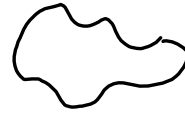
find K clusters within the data. K must be set beforehand.

## Idea:

hyperparameter: K



$C_1 \cup C_2 \cup \dots \cup C_k$

and $C_i \cap C_j = \emptyset$

$(i \neq j)$

$1 \dots \dots \dots \dots n$

Goal function: with Cluster variation

$$\min_{C_1,\dots,C_k} \left\{ \sum_{k=1}^{K} W(C_k) \right\}$$

The definition of $W(C_k)$ depends on the distance definition we use:

ex Euclidean distance

$$W(C_k) := \frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^{p} (x_{ij} - x\, x_{i'j})^2$$

— n° features

# K-means: algorithm

1) Randomly choose K clusters among the data
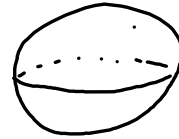
2) Repeat

   2.a) For each cluster, compute the cluster centroid.

   $$x_k = \bar{x} \mid x \in C_k = (\bar{x}_1, \ldots, \bar{x}_p)$$

   2.b) Reassign points to the cluster whose centroid is closer
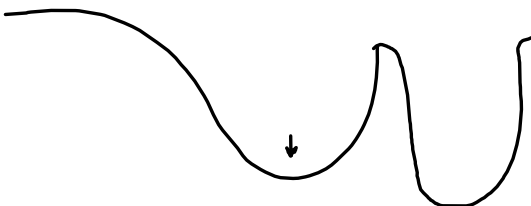
Good for: spheric data clusters found

Idea:

$$W(C_k) = 2 \sum_{i \in C_k} \sum_{j=1}^{P} (x_{ij} - \bar{x}_{kj})^2$$

mean of feature $j$ in cluster $C_k$

2.a) minimizes the sum of square deviations

2.b) minimize at each iteration $W(C_k)$

Danger:

Find local minima.
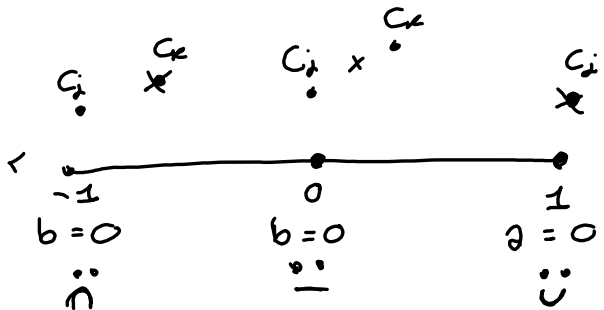
Multiple runs are recommended.

# KMEANS:
## performance measures

1. Overall $d^2$ from centroids:

$$\sum_{i=1}^{N} \sum_{\substack{x_i \in C_j \\ j=1}}^{K} d^2(x_i, C_k)$$

2. Silhouette coefficient: $S = \dfrac{1}{N} \sum_{i=1}^{N} S(x_i)$

$$S(x_t) = \frac{(b-a)}{\max(b;a)}$$



where:

$$a := \frac{1}{|C_j|} \sum_{x \in C_j} d^2(x_t, x) \quad \curvearrowright \text{ mean distance point-point in the same cluster}$$

$$b := \frac{1}{k} \sum_{j=1}^{k} d^2(x_t; c_j) \quad \curvearrowright \text{ mean distant point-other centroids}$$

# Hierarchical Clustering.

A class of Clust. algs. that do not need to choose
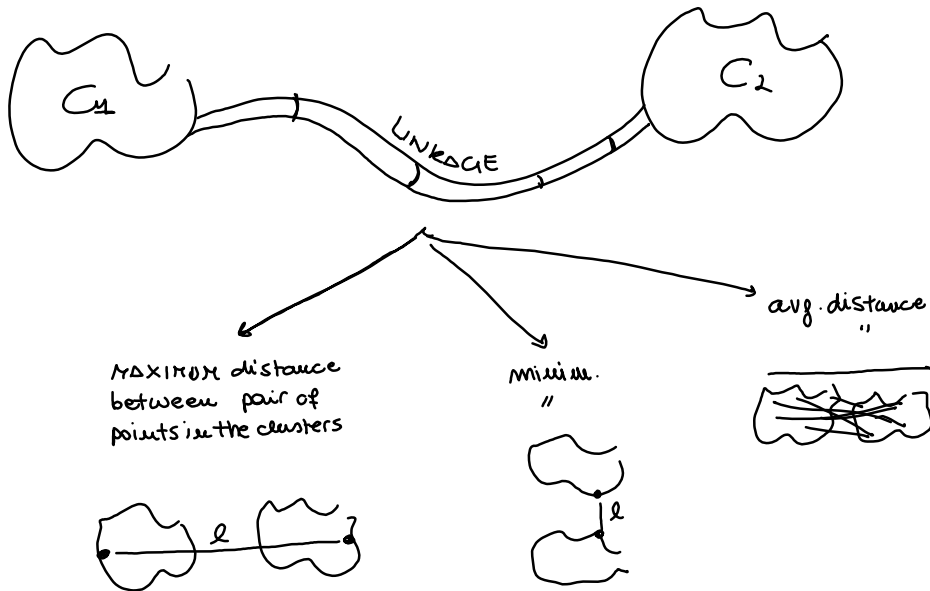k a priori.

BOTTOM-UP  : :  → ⌒  (AGGLOMERATIVE)

TOP-DOWN: ⌒ → : :  (DIVISIVE)

A dendrogram
can be associated

you can choose
the degree of clust.
precision arbitrarily.

In agglomerative clustering, in the beginning every point
is treated like a cluster. Then most similar points get
joined in the same cluster. Clusters are joined following
a chosen LINKAGE (infra-cluster similarity).

$C_1$     $C_2$     LINKAGE

avg. distance
"

MAXIMUM distance
between pair of
points in the clusters

minim.
"

$\ell$

$\ell$

hyperparams)  → similarity meas./distance (ex: euclidean dist $d^\ell$)  ⎤
            → linkage (min, max, avg...)                                ⎬ GRIDSEARCH
            → feature scaling (YES; NO)                                 ⎦

# DBSCAN algorithm

1) Define a radius $\varepsilon$ and a threshold $m$:
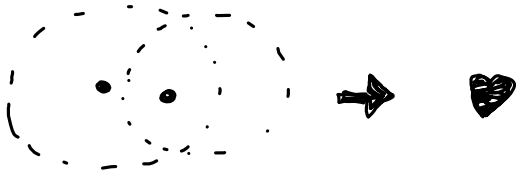


**NOTE!**
One point is a core if it has at least $m$ point within its radius

2) Label points as core, border, noise ↙

3) Remove noise points

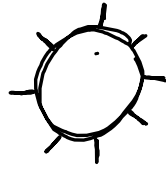4) Make clusters with nearest core points:



considering also core clusters
⌄

5) Assign border points to the nearest core point:

# DBSCAN: Setting



<u>VOTE</u> How do we set the hyperparams? ($\varepsilon$ and $m$)



distance from k-th neighbor

$\varepsilon$

Points are ordered for ascending distance from the k-th nearest point.

$\boxed{K = m}$

$\underline{\underline{ex}}$ $K = 4$ frequently

points

NON-NOISE pts.

NOISE points

# HDBSCAN : fast density based clustering.

Hierarchical density based clusting. ($\varepsilon$ is not needed)

hyperparameter : m (minimum n° of points needed to have a cluster).

IDEA :

underlying PDF
$f_{\sim}(x)$

DATA

Level map

Level sets

Visual approx. of PDF

level

$C_1$

$C_2$
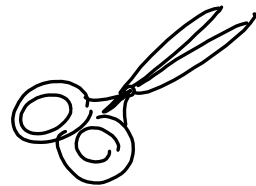
obtain clusters

PROBLEM:

• do not know PDF

?

complexity:
$O(n^2)$

• not comput. efficient

• hard to know cutting level a priori. The cutting level $\ell$ corresponds to the threshold density to form a cluster
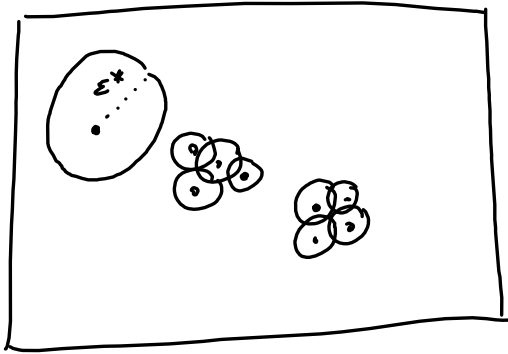
# HDBSCAN : how it works

How do we compute $f_{\tilde{x}}(x)$ ?

1. locally approximate the density with $\varepsilon^*$ :

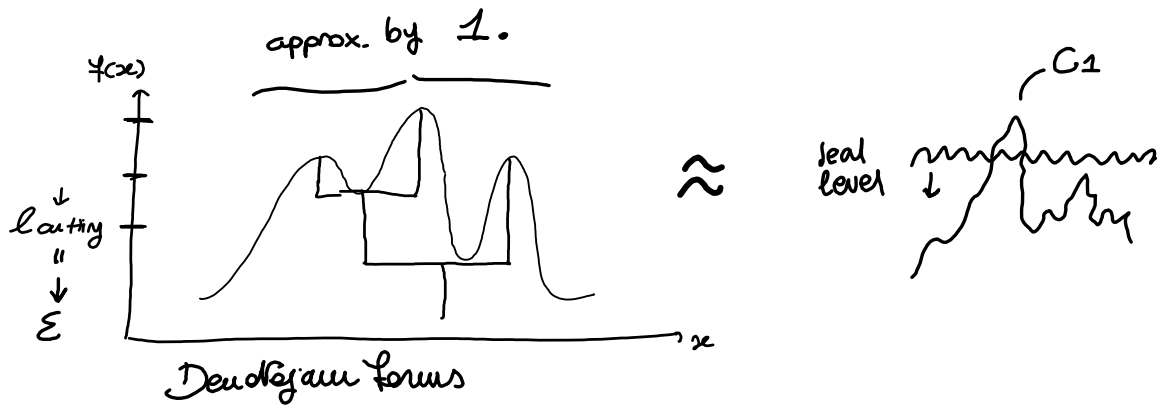$$\varepsilon^* = \min \varepsilon \;\Big|\; \text{Point is a core point}$$



Denser regions automatically detected.

This draws a sketch (approx.) of $f_{\tilde{x}}(x)$.

How do choose the $\ell$ cutting ?

2. decreasing $\ell$ cutting form a tree ( short. dendrogram)



Dendrogram forms

NOTE : in choosing $\varepsilon$ in DBSAN we implicitly choose the cutting level $\ell$.

# HDBSCAN : how it works

3. Cutting the tree according to mutual Reachibility distance :

$$d_{mreac}(x_i, x_j) = \begin{cases} \max\{ K(x_i); K(x_j); d(x_i; x_j)\} & \\ 0 & x_i = x_j \end{cases}$$

where $K(x)$ = distance $(x, K-th \text{ near. neighbour } x)$



4. Actual cut :



- vetical cut with the constraint thats if you choose one vertical cluster you cannot choose its descendants.

If you choose vertical clusters that exist for long enough maybe you are choosing the NATURAL clusters you data shows

# GMM : Gaussian Mixture Models

Assumption: data were drawn from a population: having:

$$f_{\tilde{x}}(x) = \sum_{j=1}^{m} \phi_j \, f_{\mu_j, \Sigma_j}^{(x)}$$

multivariate normal distribution.

$$P(\tilde{X} \in C_j)$$

weight

$$P(x_i \in C_j)$$

$\approx$ importance of the ith (more data were likely drawn from it.)

$\phi_j, \mu_j, \Sigma_{i,j}$ are obtained through expectance maximization



$$f_{\mu_j, \Sigma_j}^{(x)}$$

data

$$\downarrow$$

data :   $+ x x^+_x$   $+$   $- \not{\vdash}$
          $+$          $+$
          $+^+ _+ +$

output · estimates for $f_{\mu_j, \Sigma_j}(x)$ withdrawn from data
( $j = 1, \dots, m$ ).

# GMM: expectation maximization algorithm

FOR $j = 1, \ldots, m$

0) Guess initial values for all $\varphi_j, \mu_j, C_j$  $\left(\varphi_j = \frac{1}{m}\right)$

1) FOR ALL $i = 1, \ldots, N$ compute:

$$f_{\mu_j, C_j}(x_i)$$

2) Use Bayes Rule

$$P(x_i \in C_j \mid x_i) = \frac{P(x_i \mid x_i \in C_j) P(x_i \in C_j)}{P(x_i)}$$

$$= \frac{f_{\mu_j, C_j}(x_i)\, \varphi_j}{\sum_{k=1}^{m} f_{\mu_k, C_k}(x_i)\, \varphi_k}$$

3) compute $\mu_j, \sigma_j^2$

$$\mu_j = \frac{\sum_{i=1}^{N} P(x_i \in C_j \mid x_i)\, x_i}{\sum_{i=1}^{N} P(x_i \in C_j \mid x_i)}$$

$$\sigma_j^2 = \frac{\sum_{i=1}^{N} P(x_i \in C_j \mid x_i)(x_i - \mu_j)^2}{\sum_{i=1}^{N} P(x_i \in C_j \mid x_i)}$$
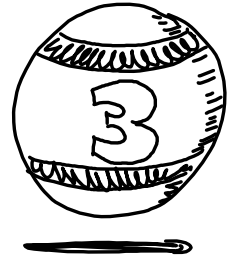
Since we do NOT know if an element belongs certainely to a cluster, the best thing we can do is updating the MVN params $(\mu_j; \sigma_j^2)$ as a prob. weigh. term over the data. ✳

4) compute $\varphi_j$:

$$\varphi_j \leftarrow \frac{1}{N} \sum_{i=1}^{N} P(x_i \in C_j \mid x_i) \quad \left.\right\} \begin{array}{l} \text{Average:} \\ \overline{P(x_i \in C_j)} \end{array}$$

STOP: $\mu_j, \sigma_j^2$ don't change much.

✳ As a result $\mu_j, \sigma_j^2$ are modified heavily just by the instances that are very likely to $\in C_j$

# GMM:
## Scores for hyperparameter tuning ($n°$ clusters)

$$AIC = 2p - 2\log L$$
$$BIC = p\log m - 2\log L$$
$\left.\right\}$ penalized likelihood scores

$$\max L \iff P(\vartheta = \vartheta_\alpha \mid \underline{x}) \text{ is high} \iff \text{the estimations made by the GMM are good}$$

$\updownarrow$

the right number of cluster has been chosen

low AIC and BIG $\longmapsto$ good $n°$ clust $\cdots\cdots$

# Bayesian GMM: auto-evaluating $n°$ clusters

Another approach would be using Bayesian GMM with the advantages of:

$\rightarrow$ $n°$ clusters > real $n°$ clust $\Rightarrow$ the alg. automatically detects the relevant $n°$ of clust. looking at the data

$(\varphi)$

$\rightarrow$ prior knowledge could be used about the weights of each $\pi\text{VN}$.

**NOTE**
the more data we have the less te prior matters.

$$\ll 1 \qquad\qquad \gg 1$$
$$\bullet\overline{\phantom{xxxxxxxxxxxxxxxx}}\bullet \quad \alpha$$
$$(\varphi)\approx 1 \qquad\qquad (\varphi\approx 0)$$
few clusters found $\qquad$ more clusters found

# Evaluation: PREDICTION SCORE



1) Run a Clustering algorithm on both:

$$C(S_{tr}, k) = A \ , \quad C(S_{te}, k)$$

considering $k$ clusters

2) Build a co-membership matrix:

$$M \in N^{|S_{te}| \times |S_{te}|} \ ; \quad m_{ij} := \begin{cases} 1 & x_i, x_j \in \text{same cluster according to } A \\ 0 & \text{otherwise} \end{cases}$$

3) Compute the prediction strength

$$PS(k) := \min_{j=1,\ldots,k} \frac{1}{|A_j| (|A_j| - 1)} \sum_{i,j \in A_j} m_{ij}$$



$A \neq C(S_{te}, k)$                                   $A \approx C(S_{te}, k)$

low                           FLAG                      high

bad $k$                         value                    good $k$