

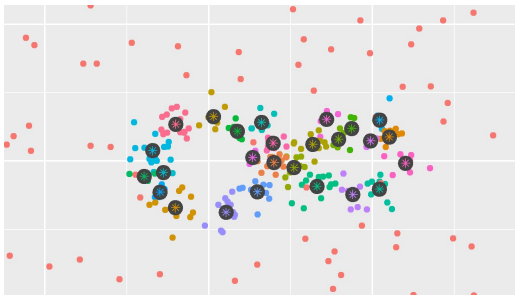
# Robust approximation of compact sets with unions of ellipsoïds. Application to data clustering.

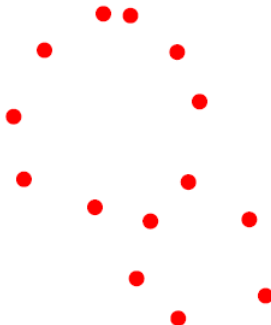
Claire Brécheteau

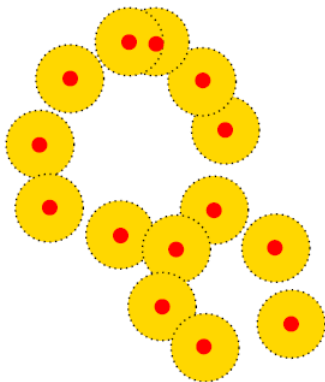
Université Rennes 2, IRMAR - with Clément Levrard, Université de Paris, LPSM

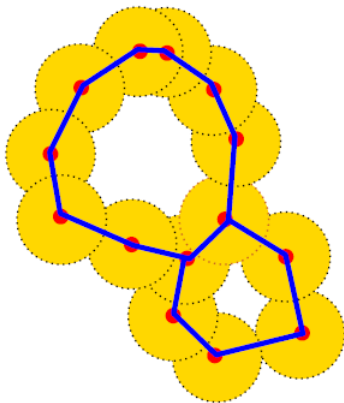
January, 29th 2021

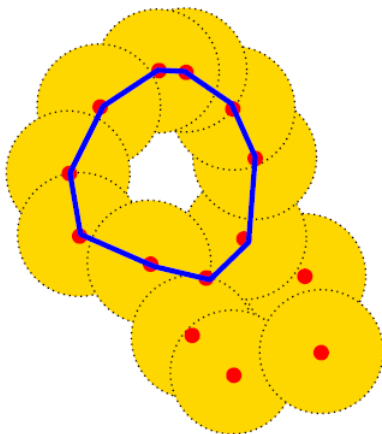
How to approximate a manifold with a set of  $k$  points,  
from a noisy sample?

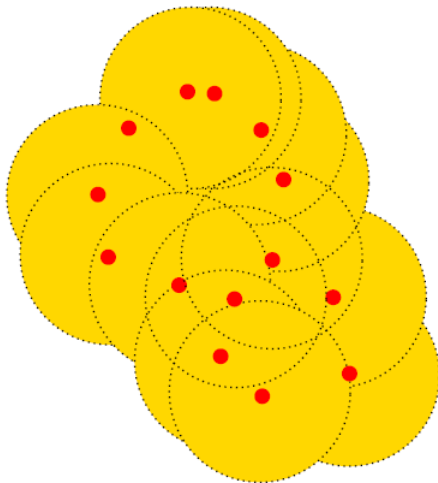












## The $k$ -means method

$P$  distribution on  $\mathbb{R}^d$

$\mathbf{c} = (c_1, c_2, \dots, c_k) \in (\mathbb{R}^d)^k$  codebook

### Definition

The optimal codebook  $\mathbf{c}^*$  minimizes the  $k$ -means loss function

$$R : \mathbf{c} \mapsto P \min_{i=1..k} \|\cdot - c_i\|^2.$$



FIGURE – Lloyd's algorithm method

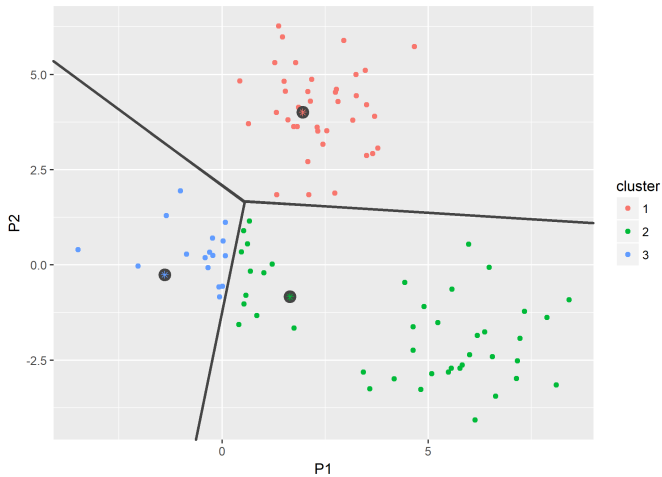


FIGURE – Lloyd's algorithm method

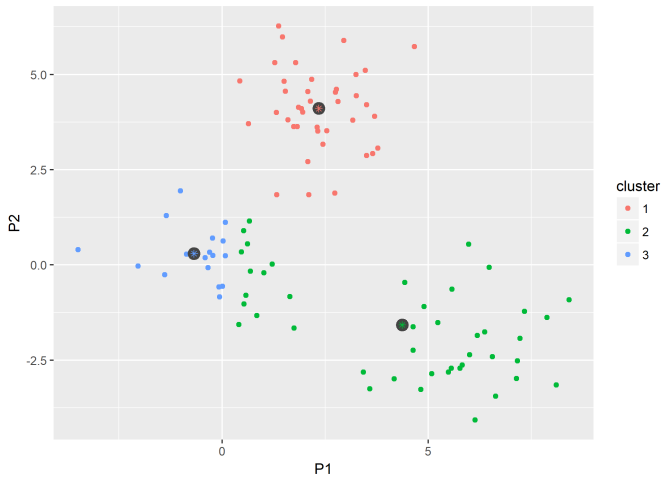


FIGURE – Lloyd's algorithm method

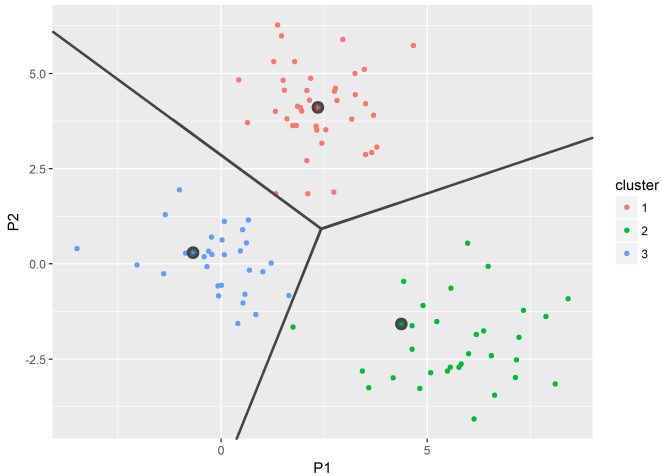


FIGURE – Lloyd's algorithm method

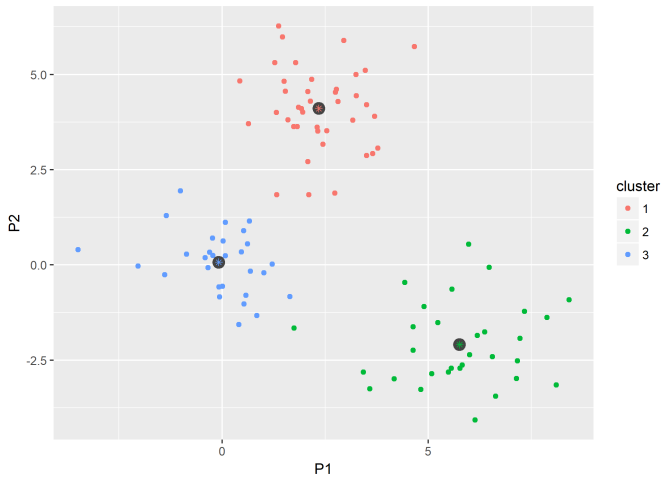
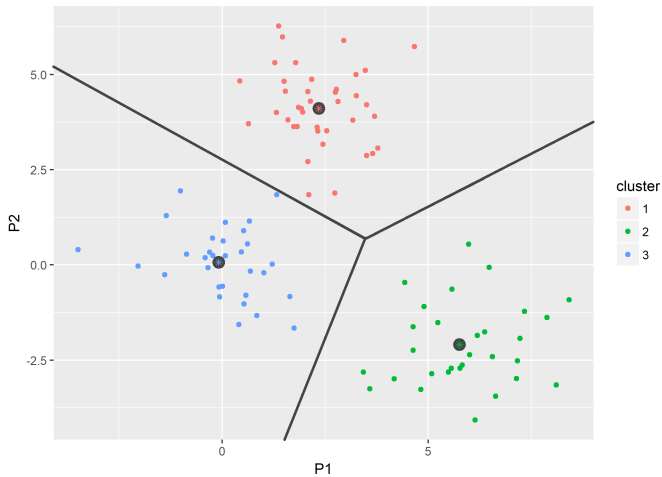


FIGURE – Lloyd's algorithm method

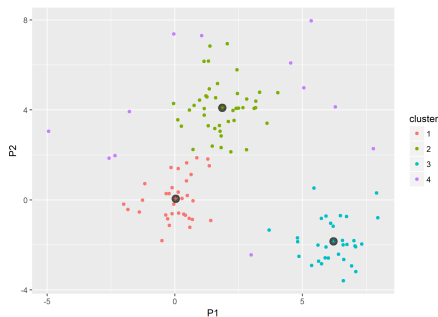


## Definition

The optimal trimmed codebook  $\mathbf{c}_h^*$  minimizes the trimmed  $k$ -means loss function

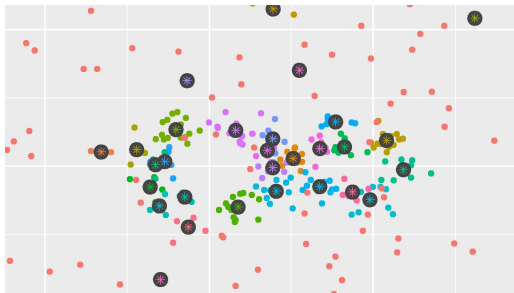
$$R_h : \mathbf{c} \mapsto \inf_{\tilde{P}, h\tilde{P} \leq P} \tilde{P} \min_{i=1..k} \|\cdot - c_i\|^2.$$

FIGURE – Trimmed  $k$ -means

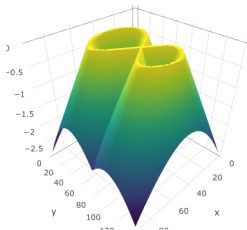


# (Trimmed) $k$ -means for the distance-approximation problem ?

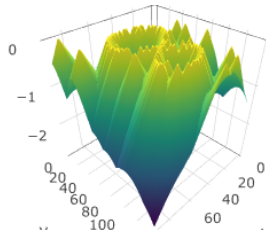
Trimmed  $k$ -means :



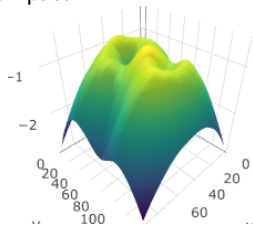
⇒ Not working...



Distance to the compact



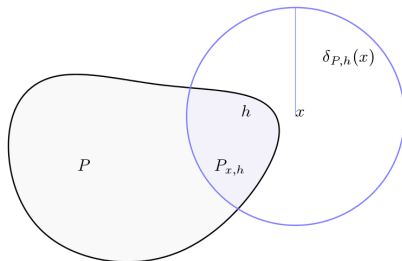
Distance to the sample



Distance to the empirical  
measure



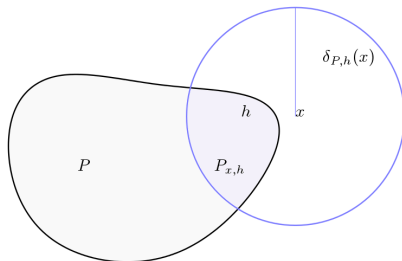
## Distance-to-measure (DTM)



$$\begin{aligned} d_{P,h}^2(x) &= P_{x,h} \|\cdot - x\|^2 \\ &= \inf_{t \in \mathbb{R}^d} P_{t,h} \|\cdot - x\|^2 \\ &= \|m(P_{x,h}) - x\|^2 + v(P_{x,h}) \\ &= \inf_{t \in \mathbb{R}^d} \|m(P_{t,h}) - x\|^2 + v(P_{t,h}) \end{aligned}$$

Notation : Mean  $m(P)$ , Variance  $v(P)$ .

## Distance-to-measure (DTM)



$$\begin{aligned} d_{P,h}^2(x) &= P_{x,h} \|\cdot - x\|^2 \\ &= \inf_{t \in \mathbb{R}^d} P_{t,h} \|\cdot - x\|^2 \\ &= \|m(P_{x,h}) - x\|^2 + v(P_{x,h}) \\ &= \inf_{t \in \mathbb{R}^d} \|m(P_{t,h}) - x\|^2 + v(P_{t,h}) \end{aligned}$$

Notation : Mean  $m(P)$ , Variance  $v(P)$ .

$\rightsquigarrow$  Sublevel sets of the DTM : union of balls.

$$\mathbf{t}^* \in \arg \min_{\mathbf{t}} P \min_{j=1..k} \|\cdot - m(P_{t_j, h})\|^2 + v(P_{t_j, h}).$$

### Definition

The  $k$ -power distance-to-measure ( $k$ -PDTM)  $d_{P, h, k}$  is defined for  $x \in \mathbb{R}^d$  by :

$$d_{P, h, k}^2(x) = \min_{j=1..k} \|x - m(P_{t_j^*, h})\|^2 + v(P_{t_j^*, h})$$

$$P \left| d_{Q_n, h, k}^2(\cdot) - d_{P, h}^2(\cdot) \right| \leq P \left| d_{Q_n, h, k}^2(\cdot) - d_{Q, h, k}^2(\cdot) \right| + P \left| d_{Q, h, k}^2(\cdot) - d_{P, k}^2(\cdot) \right|$$

## Proposition

If  $\text{Supp}(P) \subset B(0, K)$ , and  $Q\|\cdot\| < \infty$ ,

then  $P\left|d_{Q,h,k}^2(\cdot) - d_{P,h}^2(\cdot)\right|$  is bounded from above by

$$3\|d_{Q,h}^2 - d_{P,h}^2\|_{\infty, B(0,K)} + P\left(d_{P,h,k}^2(\cdot) - d_{P,h}^2(\cdot)\right) + 4W_1(P, Q) \sup_{s \in \mathbb{R}^d} \|m(P_{s,h})\|$$

with  $P\left(d_{P,h,k}^2(\cdot) - d_{P,h}^2(\cdot)\right)$  of order  $k^{-\frac{2}{d'}}$  for a “ $d'$ -dimensional distribution”.

$$\text{Supp}(P) = \mathcal{X} \subset B(0, K)$$

$X_i = Y_i + Z_i$ ,  $Y_i$  and  $Z_i$  all independent,  $Y_i \sim P$ ,  $Z_i$  sub-Gaussian with variance  $\sigma^2 \leq K^2$

$$Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}.$$

Theorem (B. - Levrard 2020)

For every  $p > 0$ , with probability larger than  $1 - 10n^{-p}$ , we have

$$\left| Pd_{Q_n, h, k}^2(\cdot) - d_{Q, h, k}^2(\cdot) \right| \leq C \sqrt{k \log(k) d} \frac{K^2 ((p+1) \log(n))^{\frac{3}{2}}}{h \sqrt{n}} + C \frac{K \sigma}{\sqrt{h}}.$$

$$\text{Supp}(P) = \mathcal{X} \subset B(0, K)$$

$X_i = Y_i + Z_i$ ,  $Y_i$  and  $Z_i$  all independent,  $Y_i \sim P$ ,  $Z_i$  sub-Gaussian with variance  $\sigma^2 \leq K^2$

$$Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}.$$

Theorem (B. - Levrard 2020)

For every  $p > 0$ , with probability larger than  $1 - 10n^{-p}$ , we have

$$\left| Pd_{Q_n, h, k}^2(\cdot) - d_{Q, h, k}^2(\cdot) \right| \leq C \sqrt{k \log(k) d} \frac{K^2 ((p+1) \log(n))^{\frac{3}{2}}}{h \sqrt{n}} + C \frac{K \sigma}{\sqrt{h}}.$$

$\rightsquigarrow$  optimize in  $k$  the quantity

$$\frac{C \sqrt{k \log(k) d} K^2 ((p+1) \log(n))^{\frac{3}{2}}}{h \sqrt{n}} + C_{P, h} k^{-\frac{2}{d'}}.$$

$$\text{Supp}(P) = \mathcal{X} \subset B(0, K)$$

$X_i = Y_i + Z_i$ ,  $Y_i$  and  $Z_i$  all independent,  $Y_i \sim P$ ,  $Z_i$  sub-Gaussian with variance  $\sigma^2 \leq K^2$

$$Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}.$$

Theorem (B. - Levrard 2020)

For every  $p > 0$ , with probability larger than  $1 - 10n^{-p}$ , we have

$$\left| Pd_{Q_n, h, k}^2(\cdot) - d_{Q, h, k}^2(\cdot) \right| \leq C \sqrt{k \log(k) d} \frac{K^2 ((p+1) \log(n))^{\frac{3}{2}}}{h \sqrt{n}} + C \frac{K \sigma}{\sqrt{h}}.$$

$\rightsquigarrow$  optimize in  $k$  the quantity

$$\frac{C \sqrt{k \log(k) d} K^2 ((p+1) \log(n))^{\frac{3}{2}}}{h \sqrt{n}} + C_{P, h} k^{-\frac{2}{d'}}.$$

Optimal choice  $k \sim n^{\frac{d'}{d'+4}}$ .



---

**Algorithm** Approximation of the  $k$ -PDTM centers
 

---

```

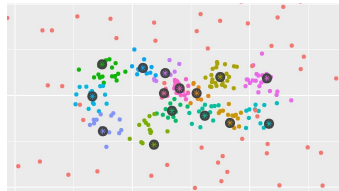
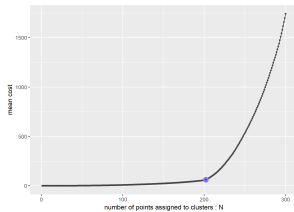
1: Input  $\mathbb{X}_n$  a  $n$ -sample from  $P$ ,  $q$  and  $k$ 
2: Sample  $t_1, t_2, \dots, t_k$  from  $\mathbb{X}_n$  without replacement.
3: while the  $t_i$ 's vary do
4:   for  $j$  in  $1..n$  do
5:     Add  $X_j$  to some  $\mathcal{C}(t_i)$  satisfying
        $\|X_j - m(t_i)\|^2 + v(t_i) \leq \|X_j - m(t_l)\|^2 + v(t_l) \forall l \neq i$ 
6:   end for
7:   for  $i$  in  $1..k$  do
8:      $t_i = \frac{1}{|\mathcal{C}(t_i)|} \sum_{X \in \mathcal{C}(t_i)} X$ 
9:   end for
10: end while
11: Output  $(t_1, t_2, \dots, t_k)$ .
    
```

---

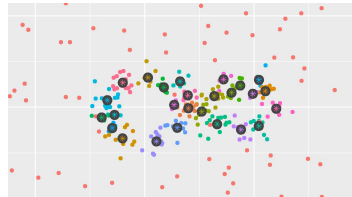
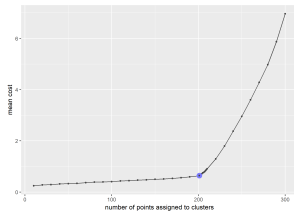
Let  $h = \frac{q}{n}$ ,  $q \in \mathbb{N}^*$ . For  $t \in \mathbb{R}^d$ ,  $m(t) = \frac{1}{q} \sum_{i=1}^q X_i(t)$ ,  $v(t) = \frac{1}{q} \sum_{i=1}^q (X_i(t) - m(t))^2$  with  $X_i(t)$  an  $i$ -th nearest neighbor of  $t$  and  $\mathcal{C}(t)$  the weighted Voronoï cell of  $t$ .

B. - Levrard, 2020

Convergence to a local minimum of  $\mathbf{t} \mapsto P_n \min_{i=1..k} \|\cdot - m(P_{n t_i, h})\|^2 + v(P_{n t_i, h})$ .



$k$ -PDTM



Trimmed  $k$ -PDTM

# Comparison to other methods



FIGURE – Comparison of the basic methods

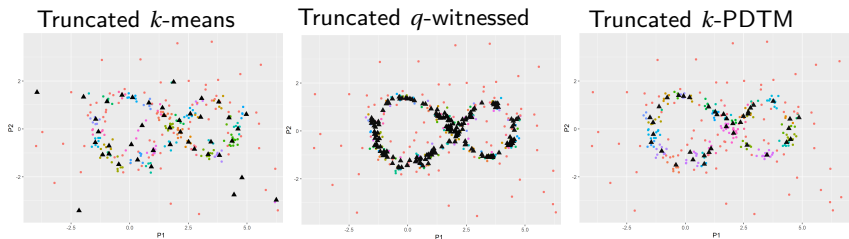
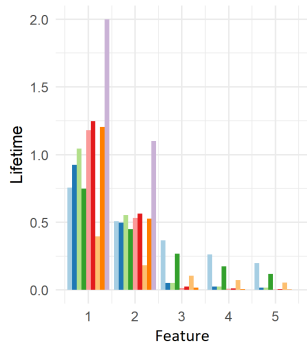
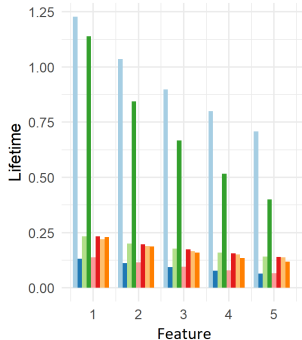


FIGURE – Comparison of the methods after thresholding

# Features lifetimes comparison



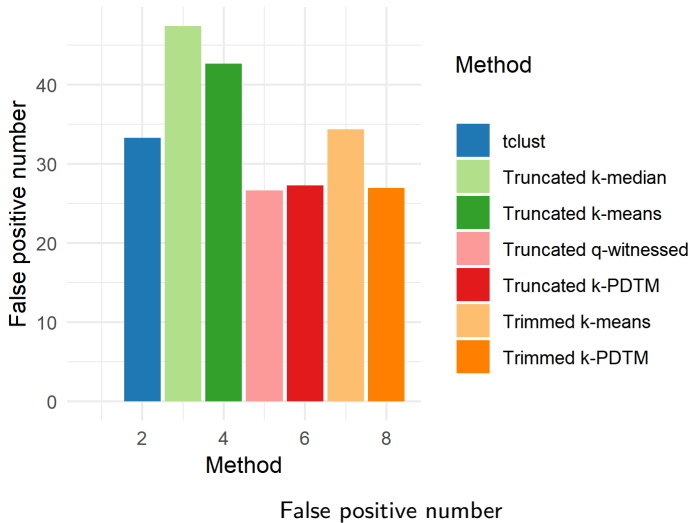
Holes



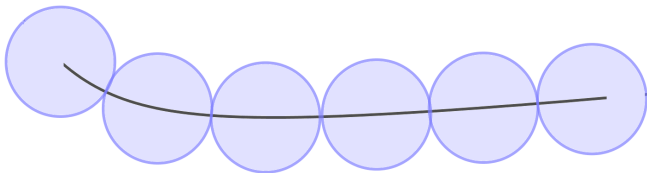
Connected components

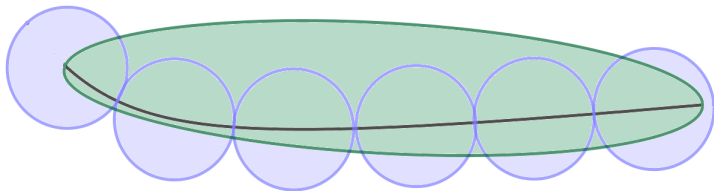
## Method

- k-means
- q-witnessed distance
- k-PDTM
- Truncated k-means
- Truncated q-witnessed
- Truncated k-PDTM
- Trimmed k-means
- Trimmed k-PDTM
- Ground truth











How to approximate a manifold with a set of  $k$  ellipsoids,  
from a noisy sample?

How to approximate a manifold with a set of  $k$  ellipsoids,  
from a noisy sample?

By modifying the criterion, with Mahalanobis norms...

Squared  $\Sigma$ -Mahalanobis norm :  $\|y\|_{\Sigma^{-1}}^2 = y^T \Sigma^{-1} y$ .

### Definition

The optimal codebook  $\theta^* = (\theta_i^*)_{i=1..k} = ((t_i^*, \Sigma_i^*))_{i=1..k}$  : a minimizer of

$$\theta \mapsto P \left( \min_{i=1..k} \|\cdot - m_{\theta_i, h}\|_{\Sigma_i^{-1}}^2 + v_{\theta_i, h}^{\Sigma_i} + \log(\det(\Sigma_i)) \right),$$

with  $m_{\theta_i, h}$  the expectation of  $P_{\theta_i, h}$  and  $v_{\theta_i, h}^{\Sigma_i}$  its variance (for  $\|\cdot\|_{\Sigma_i^{-1}}$ ).

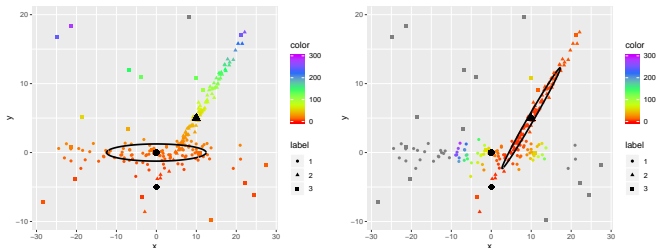


FIGURE – Illustration of the assignment phase

---

## Algorithm Approximation of the $k$ -PLM centers

---

- 1: **Input**  $\mathbb{X}_n$  a  $n$ -sample from  $P$ ,  $q = hn$  and  $k$
  - 2: Sample  $t_1, t_2, \dots, t_k$  from  $\mathbb{X}_n$  without replacement. Set  $\Sigma_i = I_d$  for  $i$  in  $1..k$ .
  - 3: **while** the  $\theta_i = (t_i, \Sigma_i)$ s vary **do**
  - 4:   Set  $\mathcal{C}(\theta_i) = \{\}$  for  $i$  in  $1..k$ .
  - 5:   **for**  $j$  in  $1..n$  **do**
  - 6:     Add  $X_j$  to some  $\mathcal{C}(\theta_i)$  satisfying  

$$\|X_j - m(\theta_i)\|_{\Sigma_i^{-1}}^2 + w(\theta_i) \leq \|X_j - m(\theta_l)\|_{\Sigma_l^{-1}}^2 + w(\theta_l) \forall l \neq i$$
  - 7:   **end for**
  - 8:   **for**  $i$  in  $1..k$  **do**
  - 9:      $t_i = \frac{1}{|\mathcal{C}(\theta_i)|} \sum_{X \in \mathcal{C}(\theta_i)} X$ ;  $\Sigma_i = \Sigma(t_i, \Sigma_i, \mathcal{C}(\theta_i))$
  - 10:     $\theta_i = (t_i, \Sigma_i)$
  - 11:   **end for**
  - 12: **end while**
  - 13: **Output**  $(\theta_1, \theta_2, \dots, \theta_k)$ .
- 

For  $\theta = (t, \Sigma)$ ,  $m(\theta) = \frac{1}{q} \sum_{i=1}^q X_i(\theta)$ ,  $w(\theta) = \frac{1}{q} \sum_{i=1}^q \|X_i(t) - m(\theta)\|_{\Sigma^{-1}}^2 + \log(\det(\Sigma))$  with  $X_i(\theta)$  an  $i$ -th  $\|\cdot\|_{\Sigma}$ -nearest neighbor of  $t$  and  $\mathcal{C}(\theta)$  the cell of  $\theta$ . Also, for  $l, m = 1..d$ ,  $\Sigma(t, \Sigma, \mathcal{C})_{l,m} = \frac{1}{|\mathcal{C}|} \sum_{X \in \mathcal{C}} P_{n,(t,\Sigma),h}(X_{(l)} - \cdot_{(l)})(X_{(m)} - \cdot_{(m)})$ .

$f$  : density on  $\mathbb{R}$

$\Sigma$  : scatter matrix

$\mu \in \mathbb{R}^d$  : location parameter

$P$  on  $\mathbb{R}^d$ , with density :

$$f_{\mu, \Sigma, f} : x \mapsto \frac{C_{d,f}}{\sqrt{\det(\Sigma)}} f(\|x - \mu\|_{\Sigma^{-1}}).$$

Theorem (B. - Levrard - Michel, 2020)

If  $f$  is non-increasing, then,  $t^* = \mu$  and  $\Sigma^* = \left(1 + \frac{1}{h} \frac{M_{d+1}^{r_h}(f)}{M_{d+1}(f)}\right) \text{Cov}(P)$ .

$$M_{d+1}^r(f) = \int_{u=0}^r u^{d+1} f(u) du$$

$$r_h \in \mathbb{R}_+ \text{ such that } h = \frac{M_{d-1}^{r_h}(f)}{M_{d-1}(f)}.$$

$$\mathcal{S}_{d,d',\sigma_{\min}^2}^* = \{\Sigma = PDP^T \mid PP^T = I_d,$$

$$D = \text{diag}(\lambda_1, \dots, \lambda_{d'}, \sigma^2 \dots \sigma^2),$$

$$\lambda_1 \geq \dots \geq \lambda_{d'} \geq \sigma^2 \geq \sigma_{\min}^2\}.$$

$P$  sub-Gaussian with variance  $V^2$  :  $\left( \forall r > V, P(B(0, r)^c) \leq \exp\left(-\frac{r^2}{2V^2}\right) \right)$

$$R_h(\theta) = P\|\cdot - m_{\theta,h}\|_{\Sigma^{-1}}^2 + v_{\theta,h}^\Sigma + \log(\det(\Sigma))$$

$\theta_{d'}^*$  (resp.  $\hat{\theta}_{d'}$ ) an  $R_h$ -minimizer (resp.  $\hat{R}_h$ ) in  $\mathbb{R}^d \times \mathcal{S}_{d,d',\sigma_{\min}^2}^*$ .

Theorem (B. - Levrard - Michel, 2020)

$$\mathbb{E} \left[ R_h(\hat{\theta}_{d'}) - R_h(\theta_{d'}^*) \right] \leq \frac{CV^2}{h\sigma_{\min}^2} \sqrt{\mathcal{D}_{d'}} \frac{\log n}{\sqrt{n}}$$

for some absolute constant  $C > 0$ , with

$$\mathcal{D}_{d'} = \max \left( d' \left( d - \frac{d'+1}{2} \right) \log(d'), (d'+1) \left( d+1 - \frac{d'}{2} \right) \right).$$

$$\mathcal{C}_{d,d'}^* = \{\{x \mid \|x - c\|_{\Sigma^{-1}}^2 \leq r\}, \Sigma \in \mathcal{S}_{d,d',0}^*, c \in \mathbb{R}^d, r > 0\}.$$

Lemma (B. - Levrard - Michel, 2020)

*The Vapnik-Chervonenkis dimension of  $\mathcal{C}_{d,d'}^*$  is bounded above by*

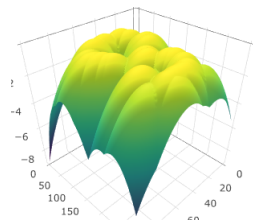
$$(12.416...) \left( d + 1 - \frac{d'}{2} \right) (d' + 1).$$

(Dudley, 1979)

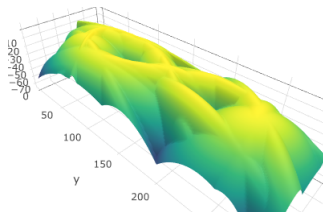
$$VC(\mathcal{C}_{d,0}^*) = d + 2$$

(Akamaa, Irie, 2011)

$$VC(\mathcal{C}_{d,d}^*) = \frac{d^2 + 3d}{2}$$



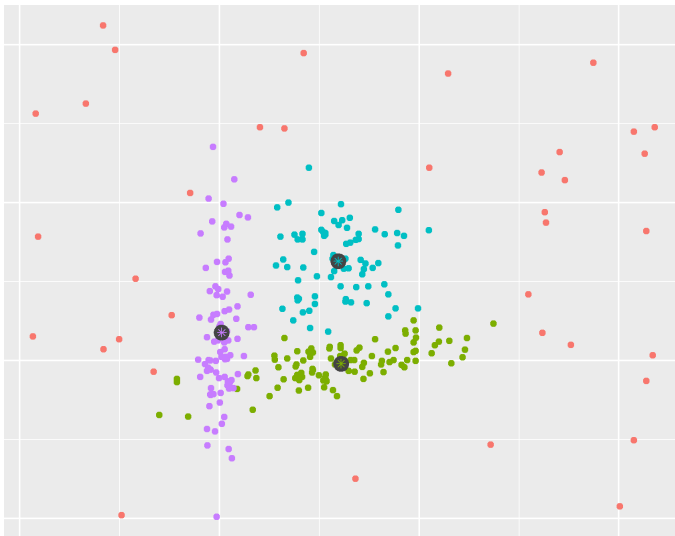
$k$ -PDTM



$k$ -PLM



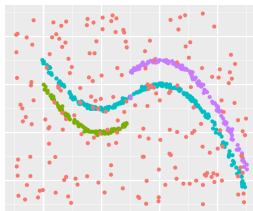
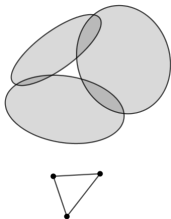
# Clustering with the $k$ -PLM



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

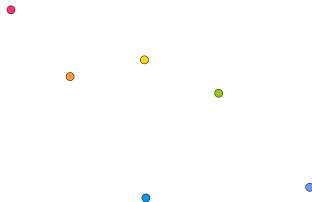
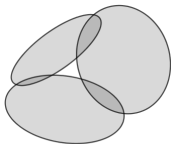
with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

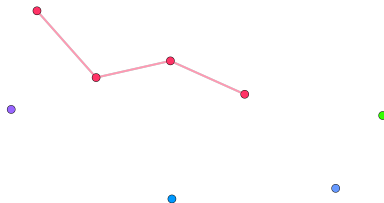
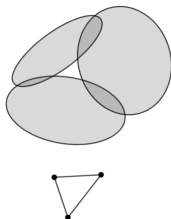
with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

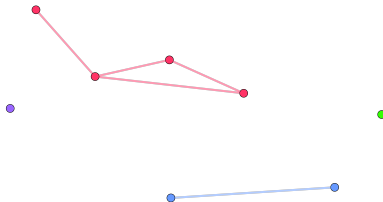
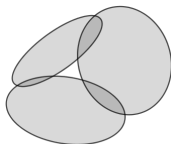
with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

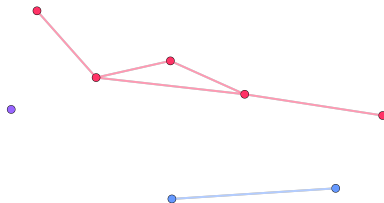
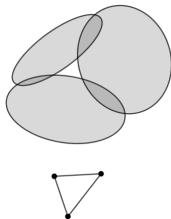
with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

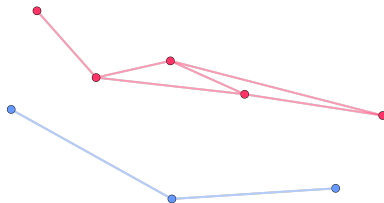
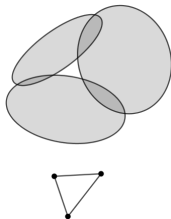
with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

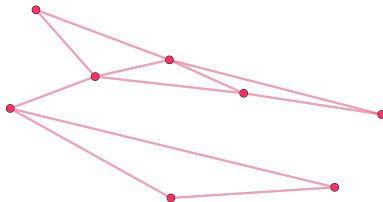
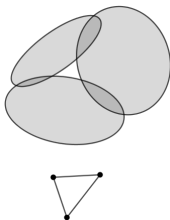
with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .

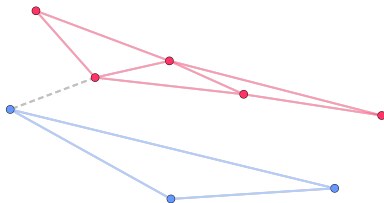
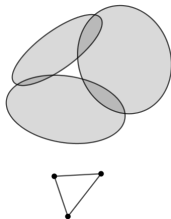




We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

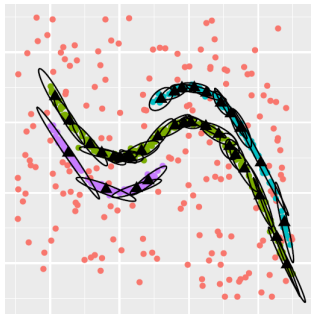
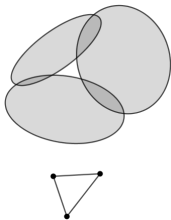
with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



We consider unions of ellipsoids indexed by  $\alpha \in \mathbb{R}$  :

$$\bigcup_{i=1..k} \{x \mid \|x - c_i\|_{\Sigma_i^{-1}}^2 + \omega_i \leq \alpha\},$$

with  $\omega_i \in \mathbb{R}$ ,  $c_i \in \mathbb{R}^d$  and covariance matrices  $\Sigma_i$ .



Partitionnement (B. 2020)

Thank you !