

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Influence function of the error rate of generalized k-means

Joint work with G. Haesbroeck

Ch. Ruwet

UNIVERSITY OF LIÈGE

30 March 2009



Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Introduction

The k -means clustering method

Influence
function of
the error rate
of
generalized
 k -means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- **Aim of clustering** : Group similar observations in k clusters C_1, \dots, C_k .
- The k -means algorithm constructs clusters in order to minimize the within cluster sum of squared distances.
- Let us focus on $k = 2$ groups.

Classification based on clustering

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

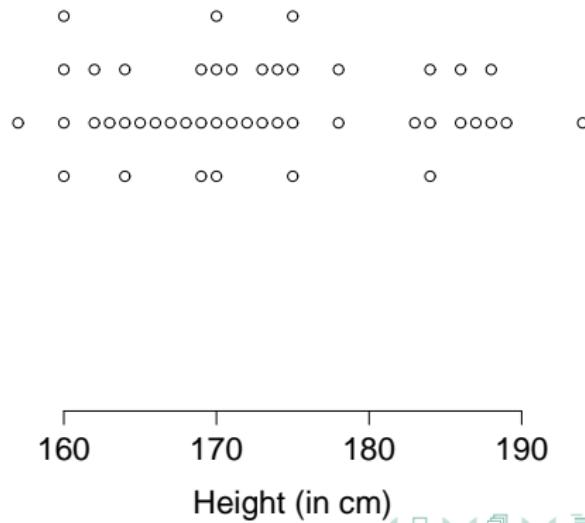
Future
research

In case of natural groups in the data, clustering may be used to find these groups via a **classification rule** :

$$x \in C_j \Leftrightarrow \|x - \bar{x}_j\| = \min_{1 \leq i \leq 2} \|x - \bar{x}_i\|$$

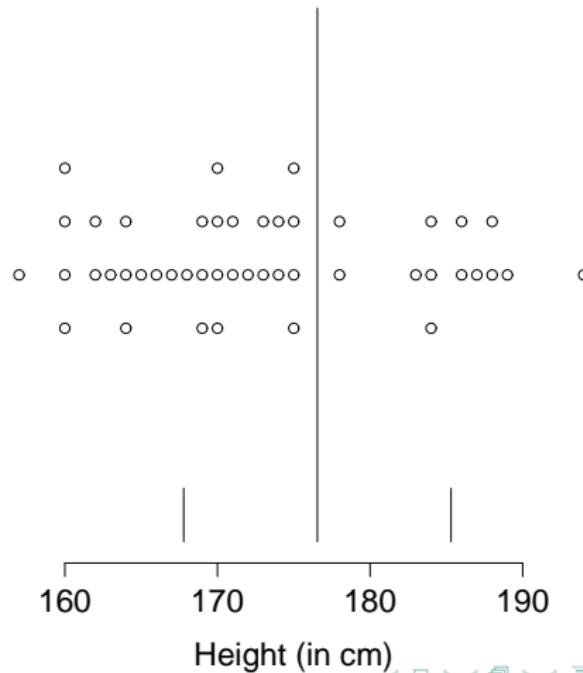
Example in 1D

Height of students in 1BM



Example in 1D

Height of students in 1BM



Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Example in 2D

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

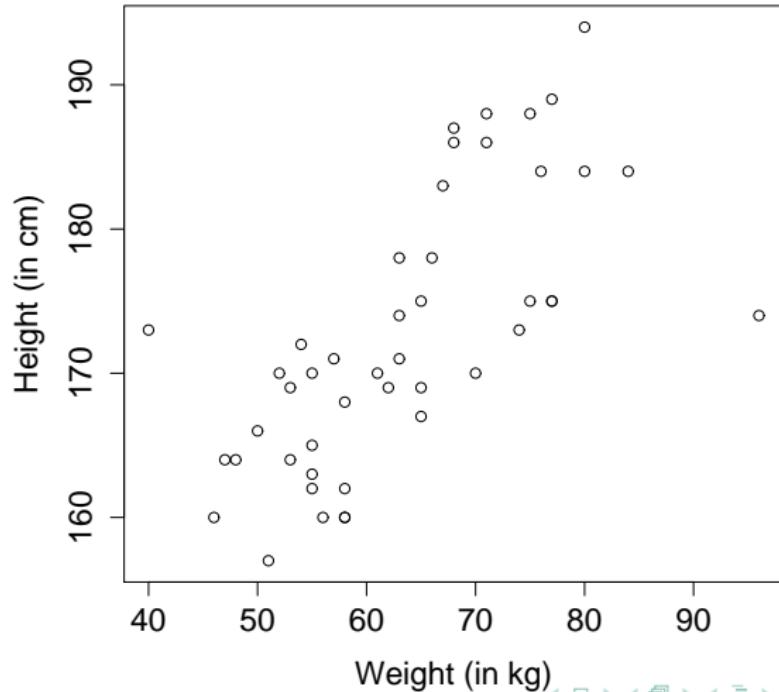
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Students of 1BM



Example in 2D

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

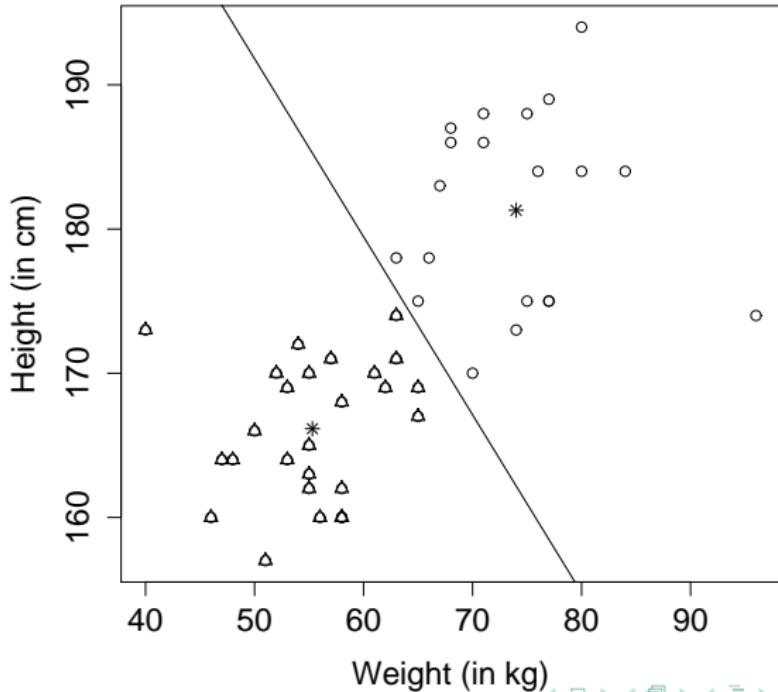
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Students of 1BM



Contaminated example in 1D

Height of students in 1BM

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

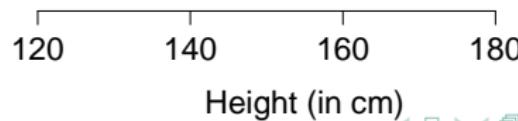
Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research



Contaminated example in 1D

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

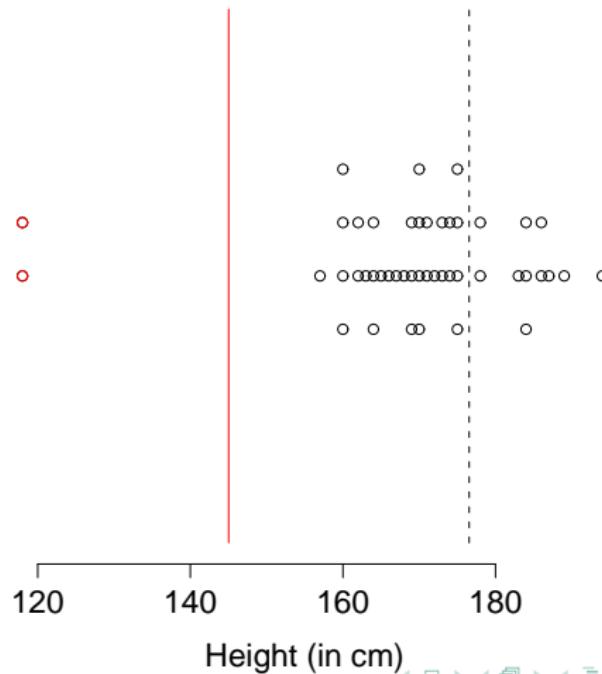
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Height of students in 1BM



The generalized 2-means clustering method

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- The clusters' centers (T_1, T_2) are solution of

$$\min_{\{t_1, t_2\} \subset \mathbb{R}^2} \sum_{i=1}^n \Omega \left(\inf_{1 \leq j \leq 2} \|x_i - t_j\| \right)$$

for a suitable nondecreasing penalty function Ω .

- Classical penalty functions :

$$\Omega(x) = x^2 \rightarrow \text{2-means method}$$

$$\Omega(x) = x \rightarrow \text{2-medoids method}$$

Contaminated example with the 2-medoids method

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

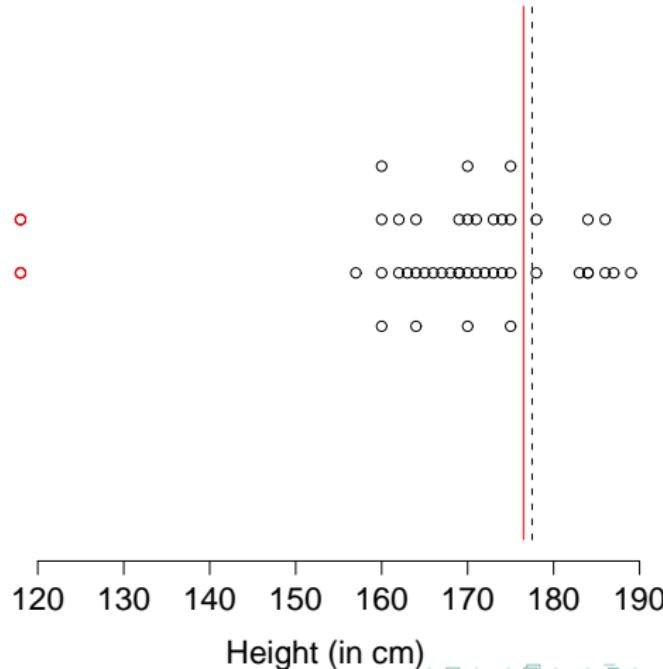
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Height of students in 1BM



The generalized 2-means clustering method

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

■ The classification rule :

$$x \in C_j \Leftrightarrow \Omega(\|x - T_j\|) = \min_{1 \leq i \leq 2} \Omega(\|x - T_i\|)$$

■ In one dimension, the estimated clusters are simply:

$$C_1 =]-\infty, C[$$

$$C_2 =]C, +\infty[$$

where $C = \frac{T_1 + T_2}{2}$ is the cut-off point.



**Influence
function of
the error rate
of
generalized
k-means**

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

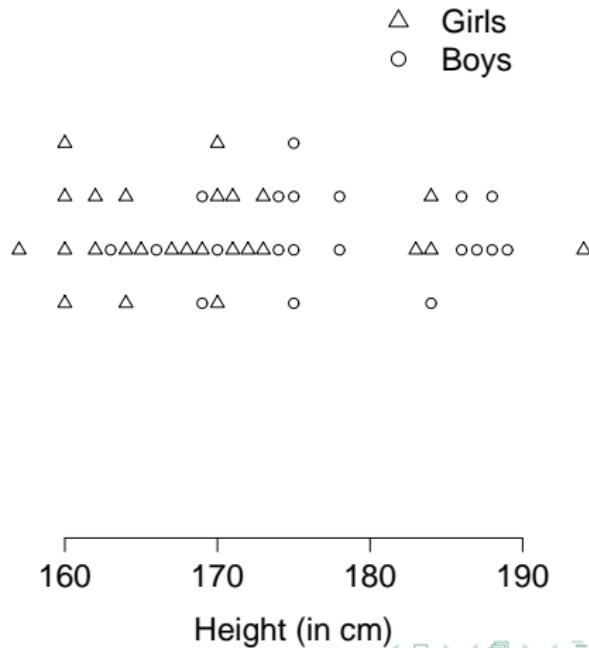
Error rate

Example in 1D

Influence function of the error rate of generalized k-means

Error rate

Height of students in 1BM



Example in 1D with the 2-means

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

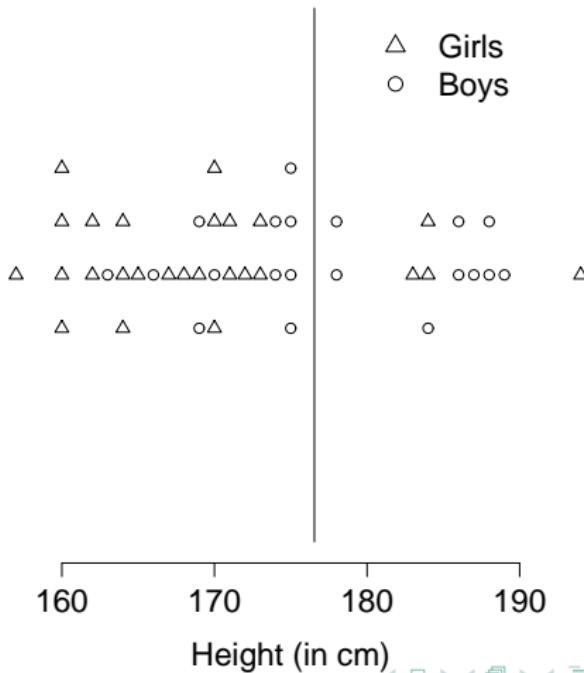
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

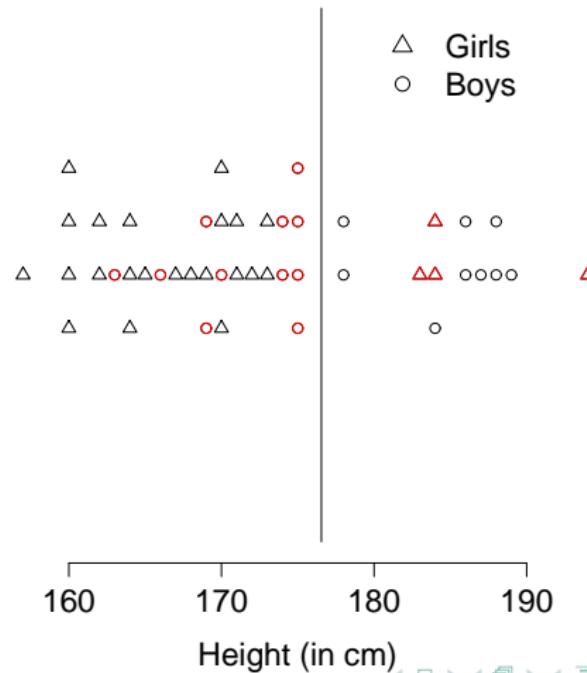
Future
research

Height of students in 1BM



Example in 1D with the 2-means

Height of students in 1BM



Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Example in 1D with the 2-means

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

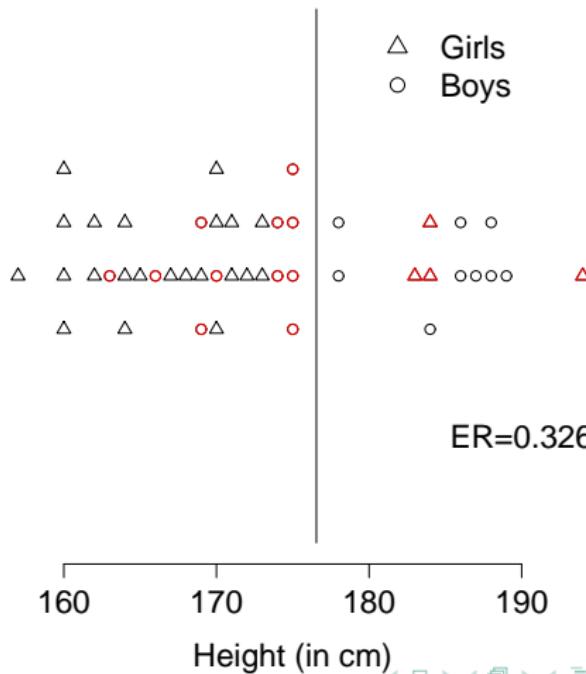
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Height of students in 1BM



Contaminated example in 1D

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Height of students in 1BM

△ Girls
○ Boys



Contaminated example in 1D with the 2-means

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

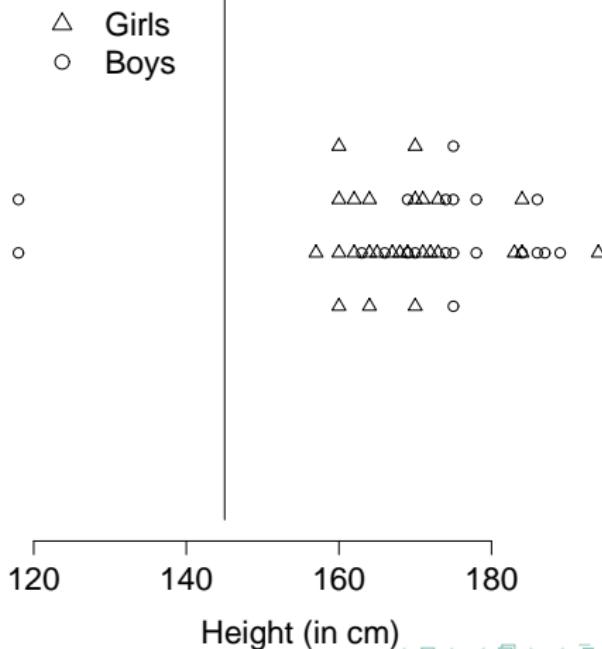
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

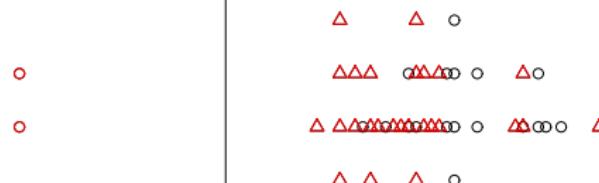
Height of students in 1BM



Contaminated example in 1D with the 2-means

Height of students in 1BM

△ Girls
○ Boys



120 140 160 180

Height (in cm)

Contaminated example in 1D with the 2-means

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

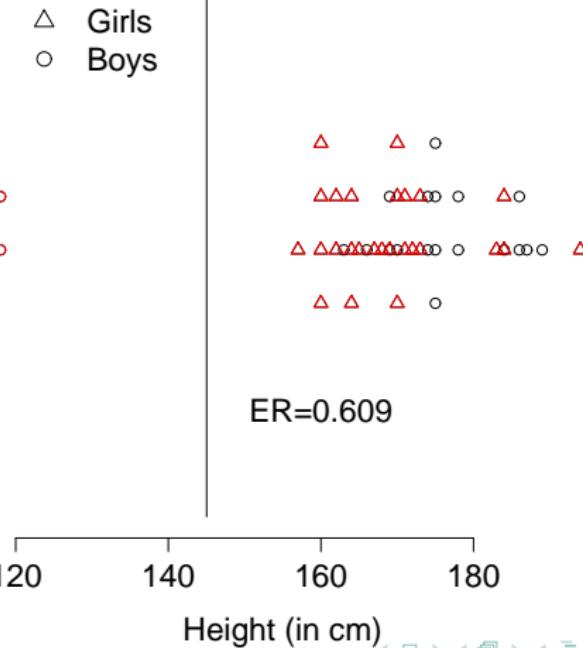
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Height of students in 1BM



Contaminated example in 1D with the 2-medoids

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Height of students in 1BM

△ Girls
○ Boys



ER=0.370

120 140 160 180

Height (in cm)

Classification setting

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Suppose

X arises from 2 groups G_1 and G_2 with $\pi_i(F) = \mathbb{P}_F[X \in G_i]$

then

F is a mixture of two distributions

$$F = \pi_1(F)F_1 + \pi_2(F)F_2$$

with densities f_1 and f_2 .

Additional assumption : one dimension !

The generalized 2-means as statistical functionals

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- The clusters' centers $(T_1(F), T_2(F))$ are solution of

$$\min_{\{t_1, t_2\} \subset \mathbb{R}^2} \int \Omega \left(\inf_{1 \leq j \leq 2} \|x - t_j\| \right) dF(x)$$

for a suitable nondecreasing penalty function Ω .

- The classification rule is

$$R_F(x) = C_j(F) \Leftrightarrow \Omega(\|x - T_j(F)\|) = \min_{1 \leq i \leq 2} \Omega(\|x - T_i(F)\|)$$

Optimality in classification

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- A classification rule is optimal if the corresponding error rate is minimal
- The optimal classification rule is the Bayes rule :

$$x \in C_1 \Leftrightarrow \pi_1(F)f_1(x) > \pi_2(F)f_2(x)$$

(Anderson, 1958)

- The 2-means procedure is optimal under the model

$$F_N = 0.5 N(\mu_1, \sigma^2) + 0.5 N(\mu_2, \sigma^2) \text{ with } \mu_1 < \mu_2$$

(Qiu and Tamhane, 2007)

Theoretical vs empirical error rate

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

■ Theoretical error rate :

- Training sample according to F : estimation of the rule
- Test sample according to F_m : evaluation of the rule
- In ideal circumstances : $F = F_m$

$$\text{ER}(F, F_m) = \sum_{j=1}^2 \pi_j(F_m) \mathbb{P}_{F_m} [R_F(X) \neq C_j(F) | G_j]$$

■ Empirical error rate :

- Training sample according to F : estimation and evaluation of the rule

$$\text{ER}(F, F) = \sum_{j=1}^2 \pi_j(F) \mathbb{P}_F [R_F(X) \neq C_j(F) | G_j]$$

Theoretical vs empirical error rate

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

■ Theoretical error rate :

- Training sample according to F : estimation of the rule
- Test sample according to F_m : evaluation of the rule
- In ideal circumstances : $F = F_m$

$$\text{ER}(F, F_m) = \sum_{j=1}^2 \pi_j(F_m) \mathbb{P}_{F_m} [R_F(X) \neq C_j(F) | G_j]$$

■ Empirical error rate :

- Training sample according to F : estimation and evaluation of the rule

$$\text{ER}(F, F) = \sum_{j=1}^2 \pi_j(F) \mathbb{P}_F [R_F(X) \neq C_j(F) | G_j]$$

Contaminated distribution

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

A contaminated distribution is defined by

$$F_\varepsilon \begin{cases} \downarrow \\ 1 - \varepsilon : F \\ \downarrow \\ \varepsilon : G \end{cases}$$

where G is an arbitrary distribution function.

Contaminated distribution

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

A contaminated distribution is defined by

$$F_\varepsilon \begin{cases} \downarrow \\ 1 - \varepsilon : F \\ \downarrow \\ \varepsilon : G \end{cases}$$

where G is an arbitrary distribution function.

To see the influence of one singular point x , $G = \Delta_x$ leading to

$$F_\varepsilon = (1 - \varepsilon)F + \varepsilon\Delta_x$$

Contaminated mixture

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

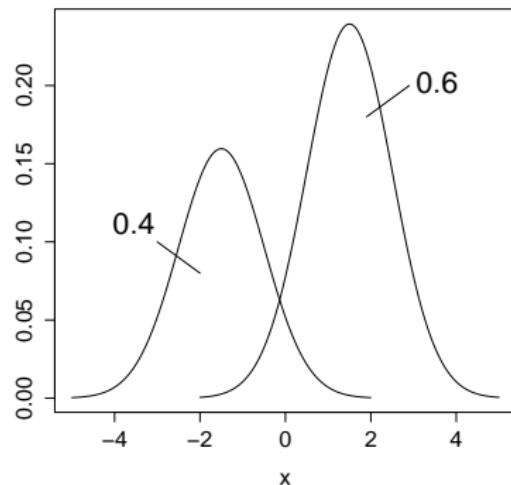
Influence
function of the
error rate

Bias of the
error rate

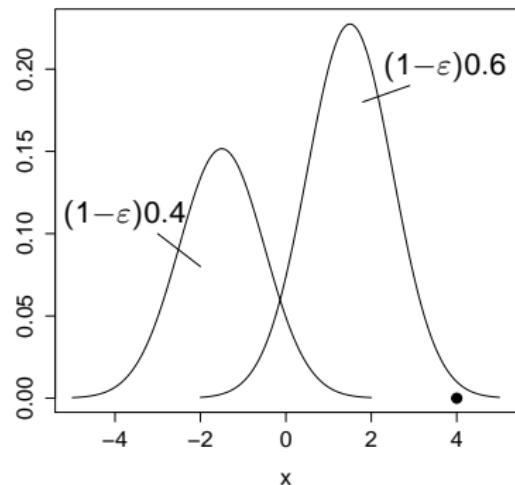
Simulation
study

Future
research

Mixture



Contaminated mixture



Theoretical vs empirical error rate under contamination

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Now, the training sample is distributed as F_ε which is a contaminated mixture.

■ Theoretical error rate :

$$\text{ER}(F_\varepsilon, F_m) = \sum_{j=1}^2 \pi_j(F_m) \mathbb{P}_{F_m} [R_{F_\varepsilon}(X) \neq C_j(F_\varepsilon) | G_j]$$

■ Empirical error rate :

$$\text{ER}(F_\varepsilon, F_\varepsilon) = \sum_{j=1}^2 \pi_j(F_\varepsilon) \mathbb{P}_{F_\varepsilon} [R_{F_\varepsilon}(X) \neq C_j(F_\varepsilon) | G_j]$$

Theoretical vs empirical error rate under contamination

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Now, the training sample is distributed as F_ε which is a contaminated mixture.

■ Theoretical error rate :

$$\text{ER}(F_\varepsilon, F_m) = \sum_{j=1}^2 \pi_j(F_m) \mathbb{P}_{F_m} [R_{F_\varepsilon}(X) \neq C_j(F_\varepsilon) | G_j]$$

■ Empirical error rate :

$$\text{ER}(F_\varepsilon, F_\varepsilon) = \sum_{j=1}^2 \pi_j(F_\varepsilon) \mathbb{P}_{F_\varepsilon} [R_{F_\varepsilon}(X) \neq C_j(F_\varepsilon) | G_j]$$

Theoretical vs empirical error rate under contamination

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Graphically :

- $F_m = F_N \equiv 0.5 N(-1, 1) + 0.5 N(1, 1)$ an optimal model
- $C(F_N) = \frac{-1+1}{2} = 0$ (Qiu and Tamhane, 2007)
- $F_\varepsilon = (1 - \varepsilon)F_m + \varepsilon\Delta_x$
- ε varying and $x = -0.5$
- $x \in G_1$ varying and $\varepsilon = 0.1$

Theoretical vs empirical error rate under contamination

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

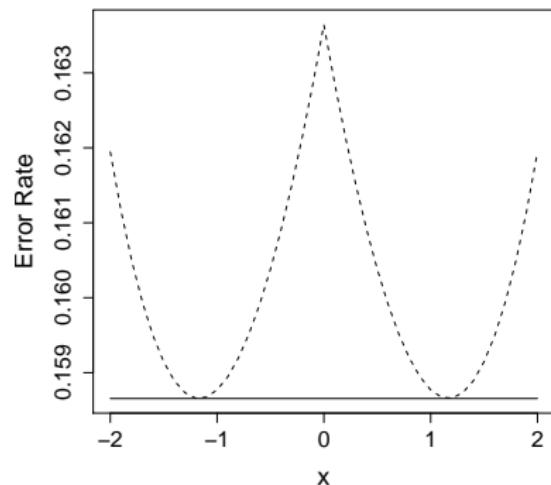
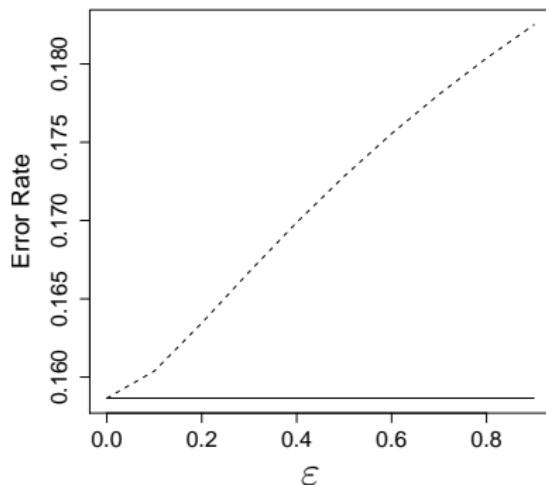
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- Theoretical error rate (with the 2-means) :



Theoretical vs empirical error rate under contamination

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

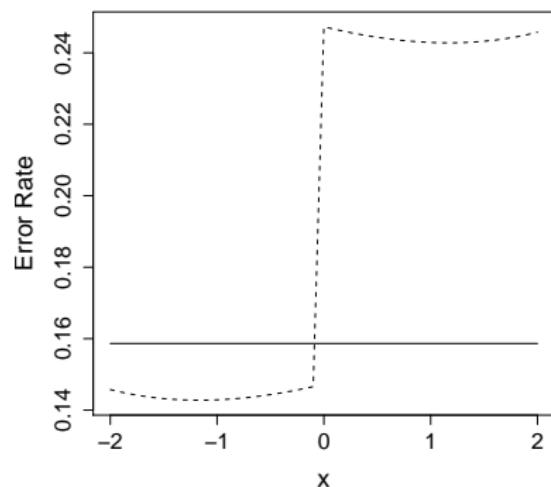
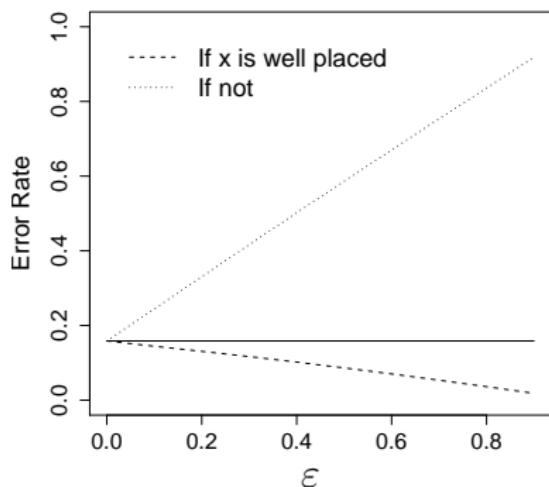
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- Empirical error rate (with the 2-means) :



Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Influence function of the error rate

Influence functions

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Hampel et al (1986) : For any statistical functional T and any distribution F ,

$$\blacksquare \text{ IF}(x; T, F) = \lim_{\varepsilon \rightarrow 0} \frac{T(F_\varepsilon) - T(F)}{\varepsilon} = \left. \frac{\partial}{\partial \varepsilon} T(F_\varepsilon) \right|_{\varepsilon=0} \text{ where } F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x;$$

- $E_F[\text{IF}(X; T, F)] = 0$;
- $T(F_\varepsilon) \approx T(F) + \varepsilon \text{IF}(x; T, F)$ for ε small enough
(First-order von Mises expansion of T at F).

Influence functions

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Hampel et al (1986) : For any statistical functional T and any distribution F ,

$$\blacksquare \text{ IF}(x; T, F) = \lim_{\varepsilon \rightarrow 0} \frac{T(F_\varepsilon) - T(F)}{\varepsilon} = \left. \frac{\partial}{\partial \varepsilon} T(F_\varepsilon) \right|_{\varepsilon=0} \text{ where } F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x;$$

- $E_F[\text{IF}(X; T, F)] = 0$;
- $T(F_\varepsilon) \approx T(F) + \varepsilon \text{IF}(x; T, F)$ for ε small enough
(First-order von Mises expansion of T at F).

Theoretical vs empirical error rate

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$\text{ER}(F_\varepsilon, F) \approx \text{ER}(F, F) + \varepsilon \text{IF}(x; \text{ER}, F)$$

$$\text{ER}(F_\varepsilon, F_\varepsilon) \approx \text{ER}(F, F) + \varepsilon \text{IF}(x; \text{ER}, F)$$

■ Theoretical error rate :

$$\text{ER}(F_\varepsilon, F_N) \geq \text{ER}(F_N, F_N) \Rightarrow \text{IF}(x; \text{ER}, F_N) \equiv 0$$

■ Empirical error rate : The IF does not vanish!
From now on, $\text{ER}(F) = \text{ER}(F, F)$.

Theoretical vs empirical error rate

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$\text{ER}(F_\varepsilon, F) \approx \text{ER}(F, F) + \varepsilon \text{IF}(x; \text{ER}, F)$$

$$\text{ER}(F_\varepsilon, F_\varepsilon) \approx \text{ER}(F, F) + \varepsilon \text{IF}(x; \text{ER}, F)$$

■ Theoretical error rate :

$$\text{ER}(F_\varepsilon, F_N) \geq \text{ER}(F_N, F_N) \Rightarrow \text{IF}(x; \text{ER}, F_N) \equiv 0$$

■ Empirical error rate : The IF does not vanish! From now on, $\text{ER}(F) = \text{ER}(F, F)$.

Theoretical vs empirical error rate

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$\text{ER}(F_\varepsilon, F) \approx \text{ER}(F, F) + \varepsilon \text{IF}(x; \text{ER}, F)$$

$$\text{ER}(F_\varepsilon, F_\varepsilon) \approx \text{ER}(F, F) + \varepsilon \text{IF}(x; \text{ER}, F)$$

■ Theoretical error rate :

$$\text{ER}(F_\varepsilon, F_N) \geq \text{ER}(F_N, F_N) \Rightarrow \text{IF}(x; \text{ER}, F_N) \equiv 0$$

■ Empirical error rate : The IF does not vanish! From now on, $\text{ER}(F) = \text{ER}(F, F)$.

ER(F_ε) = ?

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$\begin{aligned} \text{ER}(F) &= \sum_{j=1}^2 \pi_j(F) \mathbb{P}_F [R_F(X) \neq C_j(F) \mid G_j] \\ &= \pi_1(F) \{1 - F_1(C(F))\} + \pi_2(F) F_2(C(F)) \end{aligned}$$

ER(F_ε) = ?

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$\begin{aligned} \text{ER}(F) &= \sum_{j=1}^2 \pi_j(F) \mathbb{P}_F [R_F(X) \neq C_j(F) \mid G_j] \\ &= \pi_1(F) \{1 - F_1(C(F))\} + \pi_2(F) F_2(C(F)) \end{aligned}$$

Under $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_X$, one has

$$\text{ER}(F_\varepsilon) = \pi_1(F_\varepsilon) \{1 - F_{1,\varepsilon}(C(F_\varepsilon))\} + \pi_2(F_\varepsilon) F_{2,\varepsilon}(C(F_\varepsilon))$$

ER(F_ε) = ?

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$\begin{aligned} \text{ER}(F) &= \sum_{j=1}^2 \pi_j(F) \mathbb{P}_F [R_F(X) \neq C_j(F) \mid G_j] \\ &= \pi_1(F) \{1 - F_1(C(F))\} + \pi_2(F) F_2(C(F)) \end{aligned}$$

Under $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x$, one has

$$\text{ER}(F_\varepsilon) = \pi_1(F_\varepsilon) \{1 - F_{1,\varepsilon}(C(F_\varepsilon))\} + \pi_2(F_\varepsilon) F_{2,\varepsilon}(C(F_\varepsilon))$$

$\pi_i(F_\varepsilon) = ?$ and $F_{i,\varepsilon} = ?$

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$F_\varepsilon = (1 - \varepsilon)F + \varepsilon\Delta_x$$

■ $\pi_i(F_\varepsilon) = (1 - \varepsilon)\pi_i(F) + \varepsilon I\{x \in G_i\}$

■ $F_{i,\varepsilon} = \left(1 - \frac{\varepsilon I\{x \in G_i\}}{\pi_i(F_\varepsilon)}\right) F_i + \frac{\varepsilon I\{x \in G_i\}}{\pi_i(F_\varepsilon)} \Delta_x$

$$\Rightarrow F_\varepsilon = \pi_1(F_\varepsilon)F_{1,\varepsilon} + \pi_2(F_\varepsilon)F_{2,\varepsilon}$$

$\pi_i(F_\varepsilon) = ?$ and $F_{i,\varepsilon} = ?$

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x$$

■ $\pi_i(F_\varepsilon) = (1 - \varepsilon)\pi_i(F) + \varepsilon I\{x \in G_i\}$

■ $F_{i,\varepsilon} = \left(1 - \frac{\varepsilon I\{x \in G_i\}}{\pi_i(F_\varepsilon)}\right) F_i + \frac{\varepsilon I\{x \in G_i\}}{\pi_i(F_\varepsilon)} \Delta_x$

$$\Rightarrow F_\varepsilon = \pi_1(F_\varepsilon)F_{1,\varepsilon} + \pi_2(F_\varepsilon)F_{2,\varepsilon}$$

$\pi_i(F_\varepsilon) = ?$ and $F_{i,\varepsilon} = ?$

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$F_\varepsilon = (1 - \varepsilon)F + \varepsilon\Delta_x$$

■ $\pi_i(F_\varepsilon) = (1 - \varepsilon)\pi_i(F) + \varepsilon I\{x \in G_i\}$

■ $F_{i,\varepsilon} = \left(1 - \frac{\varepsilon I\{x \in G_i\}}{\pi_i(F_\varepsilon)}\right) F_i + \frac{\varepsilon I\{x \in G_i\}}{\pi_i(F_\varepsilon)} \Delta_x$

$$\Rightarrow F_\varepsilon = \pi_1(F_\varepsilon)F_{1,\varepsilon} + \pi_2(F_\varepsilon)F_{2,\varepsilon}$$

Influence function of the empirical error rate

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Proposition

For all $x \neq C(F)$,

$$\begin{aligned} \text{IF}(x; \text{ER}, F) = & -\text{ER}(F) \\ & + I\{x \leq C(F)\}(1 - 2I\{x \in G_1\}) + I\{x \in G_1\} \\ & + \frac{1}{2}(\text{IF}(x; T_1, F) + \text{IF}(x; T_2, F)) \\ & \{\pi_2(F)f_2(C(F)) - \pi_1(F)f_1(C(F))\}. \end{aligned}$$

Expressions of $\text{IF}(x; T_1, F)$ and $\text{IF}(x; T_2, F)$ have been computed by García-Escudero and Gordaliza (1999).

Graphics of $\text{IF}(x; \text{ER}, F)$

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

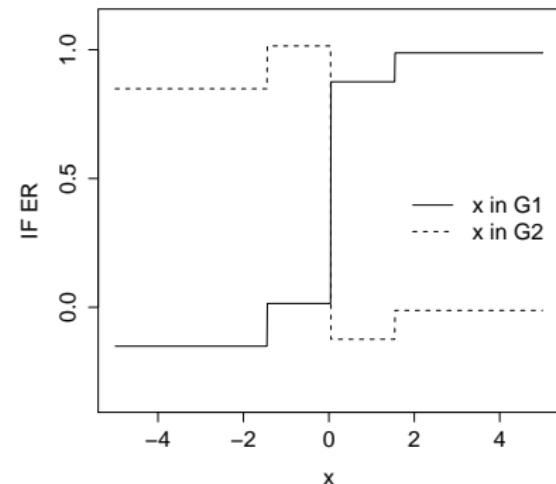
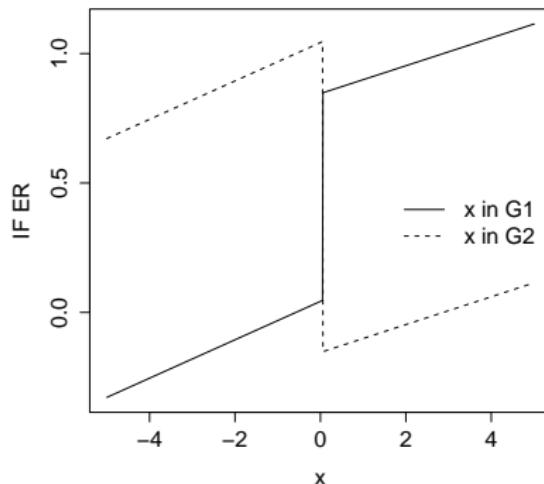
Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research



Graphics of $\text{IF}(x; \text{ER}, F)$

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

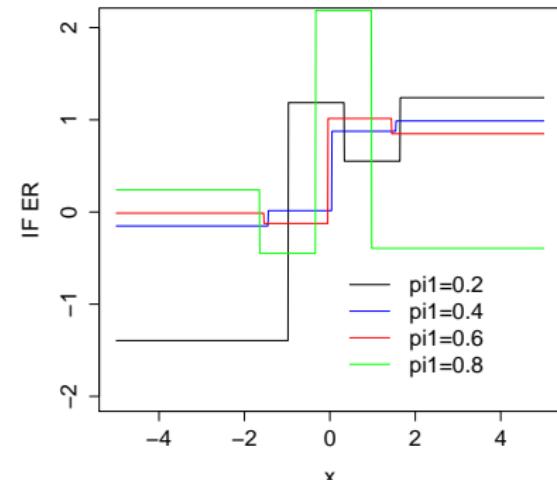
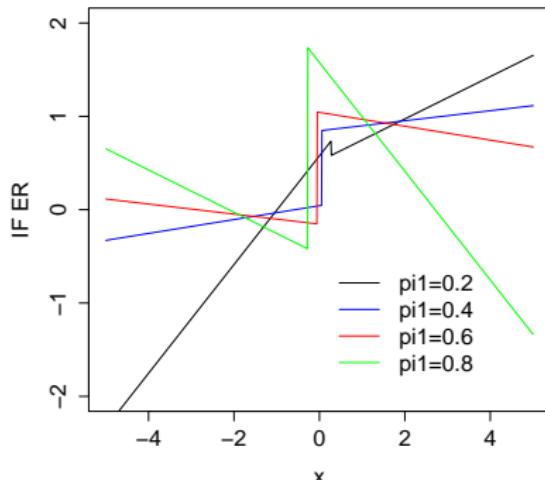
Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research



Graphics of $\text{IF}(x; \text{ER}, F)$

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

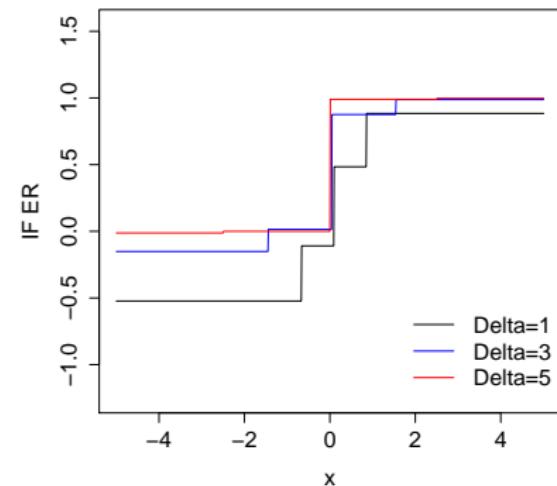
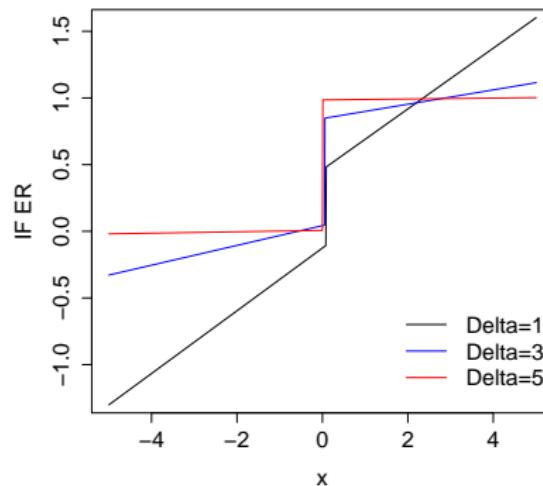
Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

$$\Delta = \mu_2 - \mu_1$$



Graphics of $\text{IF}(x; \text{ER}, F)$

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

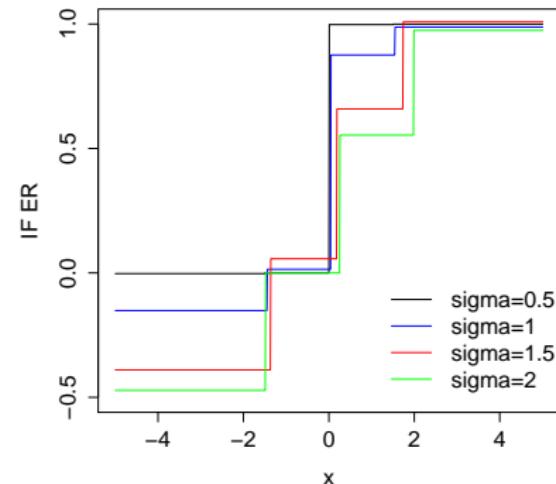
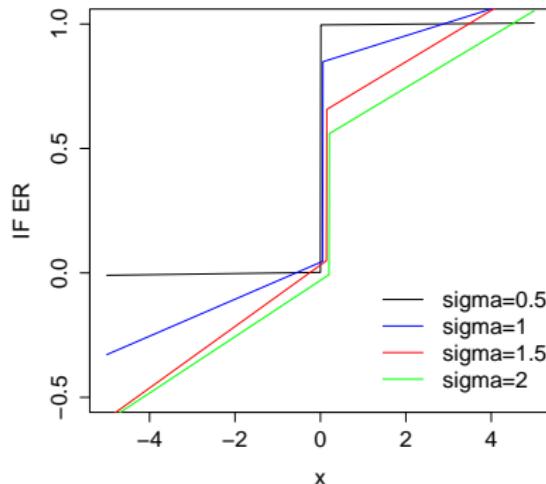
Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research





Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Bias of the error rate

Definition

- In practice, F is replaced by F_n , the empirical cdf
- The bias is defined by $B_n(\text{ER}) = E_F[\text{ER}(F_n) - \text{ER}(F)]$
- Fernholz (2001) :

$$B_n(\text{ER}) = \frac{1}{2n} E_F[\text{IF2}(X; \text{ER}, F)] + o(n^{-1})$$

where

$$\text{IF2}(x; \text{ER}, F) = \left. \frac{\partial^2}{\partial \varepsilon^2} \text{ER}((1 - \varepsilon)F + \varepsilon \Delta_x) \right|_{\varepsilon=0}.$$

This result is true if $\text{ER}(\cdot)$ is Hadamard or Fréchet differentiable.

Computation of the bias

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Proposition

Under asymptotic normality and consistency of T_1 and T_2 (Pollard, 1981 and 1982) :

$$\begin{aligned} B_n(\text{ER}) \approx & \frac{1}{4n} \{ \pi_2(F) f_2(C(F)) - \pi_1(F) f_1(C(F)) \} \\ & (E_F [\text{IF2}(X; T_1, F)] + E_F [\text{IF2}(X; T_2, F)]) \\ & + \frac{1}{8n} \{ \pi_2(F) f'_2(C(F)) - \pi_1(F) f'_1(C(F)) \} \\ & (\text{ASV}(T_1) + \text{ASV}(T_2) + 2 \text{ASC}(T_1, T_2)). \end{aligned}$$

Expressions of $\text{IF2}(X; T_1, F)$ and $\text{IF2}(X; T_2, F)$ have been computed.

Asymptotic difference

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

How much difference in error rate is to be expected by estimating a clustering rule from a finite sample ?

$$\text{A-Diff(ER)} = \lim_{n \rightarrow \infty} n B_n(\text{ER})$$

⇒ Graphical comparisons of the 2-means and 2-medoids methods :

- $F = \pi_1 N(-\Delta/2, 1) + (1 - \pi_1) N(\Delta/2, 1)$
 - π_1 varying and $\Delta = 3$
 - Δ varying and $\pi_1 = 0.4$
- $F_N = 0.5 N(-\Delta/2, 1) + 0.5 N(\Delta/2, 1)$

Graphics of A-Diff(ER) under F

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

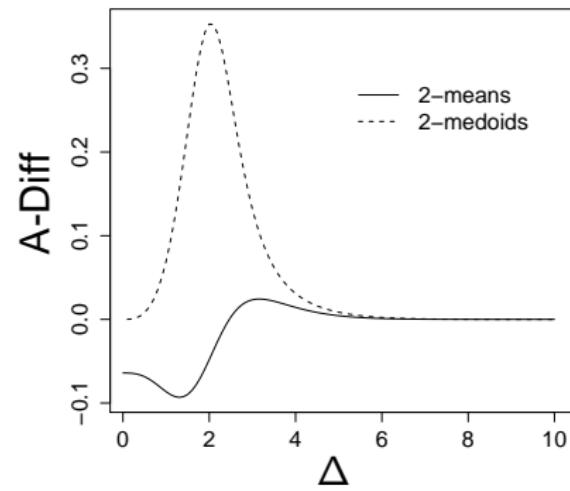
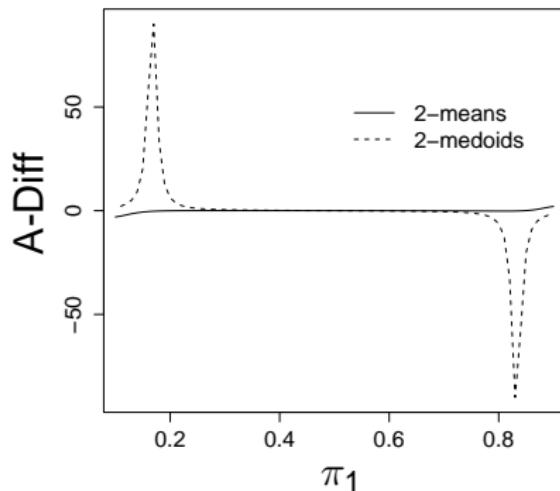
Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research



Graphic of A-Diff(ER) under F_N

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

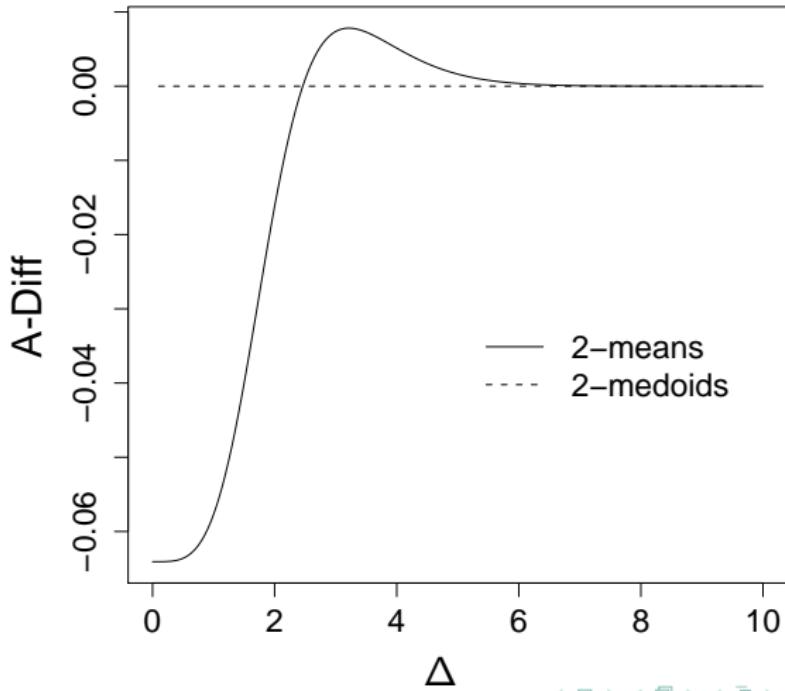
Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research





Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Simulation study

Simulation settings

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- $F = 0.4 N(-1.5, 1) + 0.6 N(1.5, 1)$
- $n = 100$
- $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x$ with
 - $\varepsilon = 0.01$ and $|x| = 5$
 - $\varepsilon = 0.05$ and $|x| = 5$
 - $\varepsilon = 0.01$ and $|x| = 50$
 - $\varepsilon = 0.05$ and $|x| = 50$
- $N = 1000$

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1				
from G_2				

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083		
from G_2	0.078	0.072		

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083	0.071	0.083
from G_2	0.078	0.072	0.081	0.073

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083	0.071	0.083
	0.064	0.139		
from G_2	0.078	0.072	0.081	0.073
	0.119	0.083		

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083	0.071	0.083
	0.064	0.139	0.066	0.127
from G_2	0.078	0.072	0.081	0.073
	0.119	0.083	0.116	0.072

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083	0.071	0.083
	0.064	0.139	0.066	0.127
	0.39	0.61		
from G_2	0.078	0.072	0.081	0.073
	0.119	0.083	0.116	0.072
	0.41	0.59		

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083	0.071	0.083
	0.064	0.139	0.066	0.127
	0.39	0.61	0.072	0.083
from G_2	0.078	0.072	0.081	0.073
	0.119	0.083	0.116	0.072
	0.41	0.59	0.081	0.073

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071)

2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083	0.071	0.083
	0.064	0.139	0.066	0.127
	0.39	0.61	0.072	0.083
	0.35	0.65		
from G_2	0.078	0.072	0.081	0.073
	0.119	0.083	0.116	0.072
	0.41	0.59	0.081	0.073
	0.45	0.55		

Simulation results

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

2-means
(0.071) 2-medoids
(0.072)

	in C_1	in C_2	in C_1	in C_2
from G_1	0.068	0.083	0.071	0.083
	0.064	0.139	0.066	0.127
	0.39	0.61	0.072	0.083
	0.35	0.65	0.35	0.65
from G_2	0.078	0.072	0.081	0.073
	0.119	0.083	0.116	0.072
	0.41	0.59	0.081	0.073
	0.45	0.55	0.45	0.55

Contaminated example with the 2-medoids

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

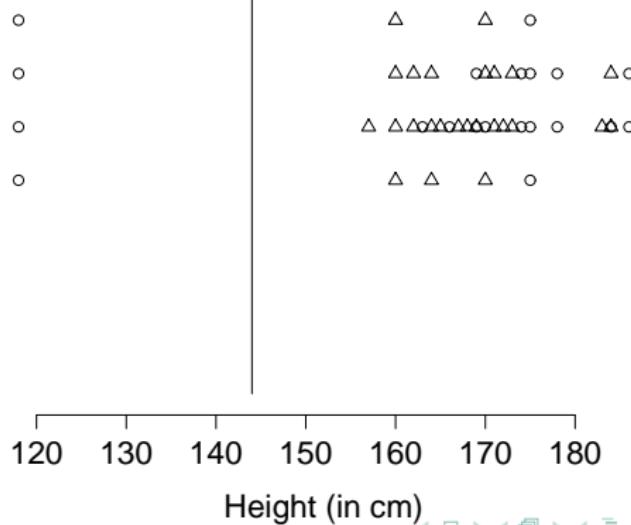
Bias of the
error rate

Simulation
study

Future
research

Height of students in 1BM

△ Girls
○ Boys



Conclusion of these simulations

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- A well-placed contamination in the smallest group makes the error rate decrease
- Too much contamination makes the error rate of the 2-means break down
- Unfortunately, the error rate of the 2-medoids also breaks down when there are too much and too far outliers !

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Future research

Future research

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- Generalized trimmed 2-means : for $\alpha \in [0, 1]$,
 $(T_1(F), T_2(F))$ are solution of

$$\min_{\{A: F(A)=1-\alpha\}} \min_{\{t_1, t_2\} \subset \mathbb{R}} \int_A \Omega \left(\inf_{1 \leq j \leq 2} \|x - t_j\| \right) dF(x).$$

- Theoretical error rate

$$\text{ER}(F_\varepsilon, F_m) = \sum_{j=1}^2 \pi_j(F_m) \mathbb{P}_{F_m} [R_{F_\varepsilon}(X) \neq C_j(F_\varepsilon) \mid G_j]$$

- More than 1 dimension or more than 2 groups

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

Thank you for your attention!

Bibliography

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- CROUX C., FILZMOSER P. and JOOSSENS K. (2008), Classification efficiencies for robust linear discriminant analysis, *Statistica Sinica* 18, pp. 581-599.
- CROUX C., HAESBROECK G. and JOOSSENS K. (2008), Logistic discrimination using robust estimators : an influence function approach, *The Canadian Journal of Statistics*, 36, pp. 157-174.
- FERNHOLZ L. T., On multivariate higher order von Mises expansions in *Metrika* 2001, vol. 53, pp. 123-140.

Bibliography

- GARCÍA-ESCUDERO L. A. and GORDALIZA A., Robustness Properties of k Means and Trimmed k Means, *Journal of the American Statistical Association*, September 1999, Vol. 94, n° 447, pp. 956-969.
- HAMPEL F.R., RONCHETTI E.M., ROUSSEEUW P.J., STAHEL W.A., *Robust Statistics : The Approach Based on Influence Functions*, John Wiley and Sons, New-York, 1986.
- ANDERSON T.W., *An Introduction to Multivariate Statistical Analysis*, Wiley, New-York, 1958, pp. 126-133.

Bibliography

Influence
function of
the error rate
of
generalized
k-means

Ch. Ruwet

Introduction

Error rate

Influence
function of the
error rate

Bias of the
error rate

Simulation
study

Future
research

- POLLARD D., Strong Consistency of k-Means Clustering, *The Annals of Probability*, 1981, Vol.9, n°4, pp.919-926.
- POLLARD D., A Central Limit Theorem for k-Means Clustering, *The Annals of Probability*, 1982, Vol.10, n°1, pp.135-140.
- QIU D. and TAMHANE A. C. (2007), A comparative study of the k-means algorithm and the normal mixture model for clustering : Univariate case, *Journal of Statistical Planning and Inference*, 137, pp. 3722-3740.