Université
de Liège

**Influence
function of
the error rate
of
classification
based on
clustering**

Ch. Ruwet

Introduction

Error rate

IF of the
error rate

Conclusion

# Influence function of the error rate of classification based on clustering

*Joint work with G. Haesbroeck*

Ch. Ruwet

Department of Mathematics - University of Liège

19 May 2009

cruwet@ulg.ac.be

# Outline

Influence
function of
the error rate
of
classification
based on
clustering

Ch. Ruwet

Introduction

Error rate

IF of the
error rate

Conclusion

Influence
function of
the error rate
of
classification
based on
clustering

Ch. Ruwet

Introduction

Error rate

IF of the
error rate

Conclusion

Suppose

$X \sim F$ arises from $G_1$ 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 density $f = \pi_1(F)f_1 + \pi_2(F)f_2$.

Additional assumption : one dimension !

- Aim of clustering : Find estimations $C_1(F)$ and $C_2(F)$ (called clusters) of the two underlying groups.
- The clusters' centers $(T_1(F), T_2(F))$ are solutions of

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

for a suitable nondecreasing penalty function $\Omega : \mathbb{R}^+ \to \mathbb{R}^+$.
- Classical penalty functions :

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

- The classification rule is

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

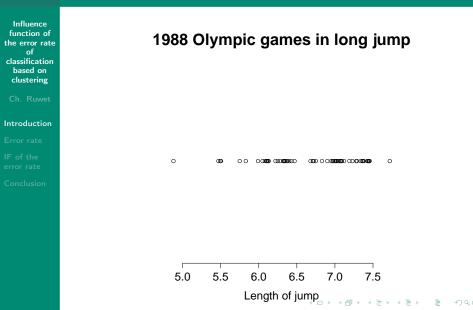- In one dimension, the clusters are simply :

$$C_1(F) = ]-\infty, C(F)[$$

$$C_2(F) = ]C(F), +\infty[$$

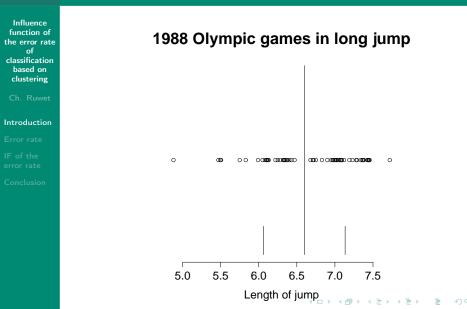where $C(F) = \dfrac{T_1(F) + T_2(F)}{2}$ is the cut-off point.

**1988 Olympic games in long jump**

1988 Olympic games in long jump

**1988 Olympic games in long jump**

Length of jump

**1988 Olympic games in long jump**

**1988 Olympic games in long jump**

Influence function of the error rate of classification based on clustering

Ch. Ruwet

Introduction

Error rate

IF of the error rate

Conclusion

**1988 Olympic games in long jump**

○   Men
△   Women

**1988 Olympic games in long jump**

○ Men
△ Women

Length of jump

**1988 Olympic games in long jump**

**1988 Olympic games in long jump**

# 1988 Olympic games in long jump

○ Men
△ Women



Length of jump

**1988 Olympic games in long jump**

Length of jump

**1988 Olympic games in long jump**

**1988 Olympic games in long jump**

Université
de Liège

**Influence
function of
the error rate
of
classification
based on
clustering**

Ch. Ruwet

Introduction

**Error rate**

IF of the
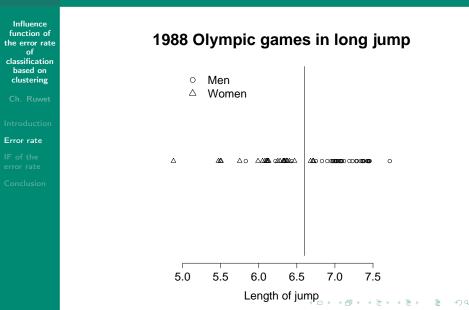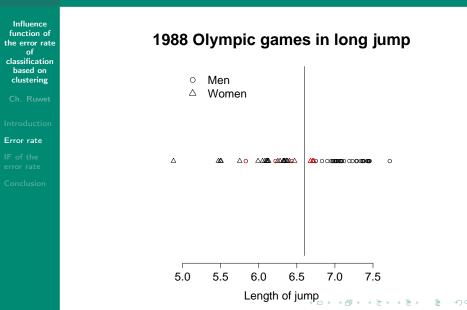error rate

Conclusion



**1988 Olympic games in long jump**

○   Men
△   Women

ER=0.12

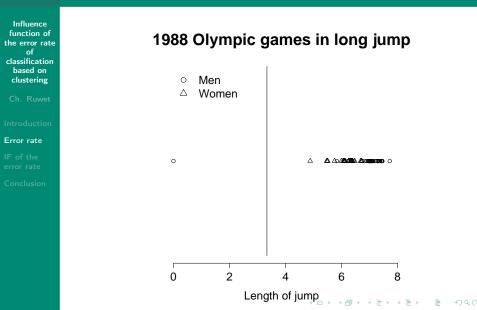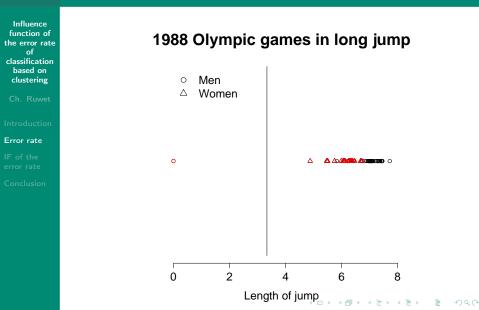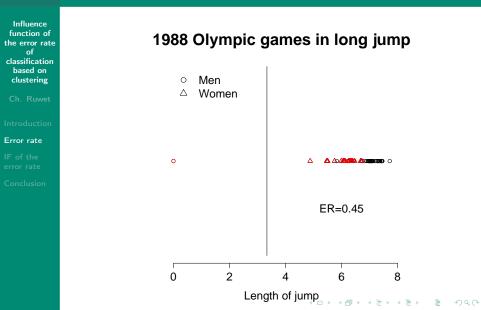Length of jump

0     2     4     6     8

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

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

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

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

(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)

# Outline

Influence
function of
the error rate
of
classification
based on
clustering

Ch. Ruwet

Introduction

Error rate

IF of the
error rate

First order IF
Second order
IF
ARCE

Conclusion

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

- $\mathsf{IF}(x; \mathsf{T}, F) = \lim_{\varepsilon \to 0} \dfrac{\mathsf{T}(F_\varepsilon) - \mathsf{T}(F)}{\varepsilon} = \left. \dfrac{\partial}{\partial \varepsilon} \mathsf{T}(F_\varepsilon) \right|_{\varepsilon=0}$ where
  $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x$
  (under condition of existence);

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

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

- $\mathsf{IF}(x; \mathsf{T}, F) = \lim_{\varepsilon \to 0} \dfrac{\mathsf{T}(F_\varepsilon) - \mathsf{T}(F)}{\varepsilon} = \left. \dfrac{\partial}{\partial \varepsilon} \mathsf{T}(F_\varepsilon) \right|_{\varepsilon = 0}$ where
  $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x$
  (under condition of existence);

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

Now, the training sample is distributed as $F_\varepsilon$ which is a contaminated mixture.

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

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

Now, the training sample is distributed as $F_\varepsilon$ which is a contaminated mixture.

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

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

- $\mathsf{ER}(F_\varepsilon, F_N) \approx \mathsf{ER}(F_N, F_N) + \varepsilon\mathsf{IF}(x; \mathsf{ER}, F_N)$
- $\mathsf{ER}(F_\varepsilon, F_N) \geq \mathsf{ER}(F_N, F_N)$

# First order influence function of the error rate

Now, the training sample is distributed as $F_\varepsilon$ which is a contaminated mixture.

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

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

- $\mathrm{ER}(F_\varepsilon, F_N) \approx \mathrm{ER}(F_N, F_N) + \varepsilon \mathrm{IF}(x; \mathrm{ER}, F_N)$
- $\mathrm{ER}(F_\varepsilon, F_N) \geq \mathrm{ER}(F_N, F_N)$

$$\Rightarrow \mathrm{IF}(x; \mathrm{ER}, F_N) \equiv 0$$

### Proposition

*The influence function of the error rate of the generalized 2-means classification procedure is given by*

$$IF(x; ER, F) = \frac{1}{2}\{IF(x; T_1, F) + IF(x; T_2, F)\}$$
$$\{\pi_2(F)f_2(C(F)) - \pi_1(F)f_1(C(F))\}$$

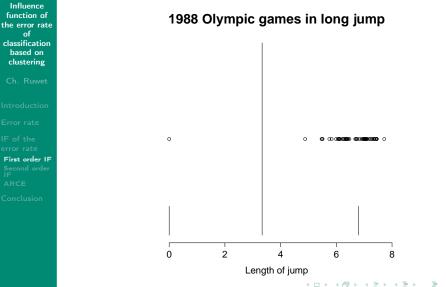*for all $x \neq C(F)$.*

Expressions of $IF(x; T_1, F)$ and $IF(x; T_2, F)$ were computed by García-Escudero and Gordaliza (1999).

**1988 Olympic games in long jump**



Length of jump

# Penalty functions

- 2-means : $\Omega(x) = x^2$
- 2-medoids : $\Omega(x) = x$

- 2-Tukey's : $\Omega(x) = \frac{b^2}{6} \begin{cases} 1 - \left[ 1 - \left( \frac{x}{b} \right)^2 \right]^3 & \text{if } |x| \leq b \\ 1 & \text{if } |x| > b \end{cases}$

  with $b = 2.795$

$$F = \pi_1 \, N(-\Delta/2, 1) + (1 - \pi_1) \, N(\Delta/2, 1)$$

with

- $\pi_1 = 0.4$ and $\Delta = 3$
- $\pi_1$ is varying and $\Delta = 3$
- $\pi_1 = 0.4$ and $\Delta$ is varying

**2–means**

2–medoids

Under the model $F_N$, $\mathrm{IF}(x; \mathrm{ER}, F_N) \equiv 0$

One needs to go a step further !

Under the model $F_N$, $\mathrm{IF}(x; \mathrm{ER}, F_N) \equiv 0$

One needs to go a step further !

For any statistical functional $T$ and any distribution $F$,

$$\mathrm{IF2}(x; T, F) = \left. \frac{\partial^2}{\partial \varepsilon^2} T(F_\varepsilon) \right|_{\varepsilon = 0}$$

where $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x$ (under condition of existence).

Under the model $F_N$, $\mathrm{IF}(x; \mathrm{ER}, F_N) \equiv 0$

One needs to go a step further !

For any statistical functional $T$ and any distribution $F$,

$$\mathrm{IF2}(x; \mathsf{T}, F) = \left. \frac{\partial^2}{\partial \varepsilon^2} \mathsf{T}(F_\varepsilon) \right|_{\varepsilon=0}$$

where $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \Delta_x$ (under condition of existence).
Second order von Mises expansion of $ER$ at $F_N$ :

$$\mathrm{ER}(F_\varepsilon, F_N) \approx \mathrm{ER}(F_N, F_N) + \frac{\varepsilon^2}{2} \mathrm{IF2}(x; \mathsf{T}, F_N)$$

for $\varepsilon$ small enough.

### Proposition

*Under the optimal model $F_N$, the second order influence function of the error rate of the generalized 2-means classification procedure is given by*

$$\mathrm{IF2}(x; \mathrm{ER}, F_N) = -\frac{1}{4}(f_{N1})' \left( \frac{\mu_1 + \mu_2}{2} \right)$$
$$\{IF(x; T_1, F_N) + IF(x; T_2, F_N)\}^2$$

*for all $x \neq \frac{\mu_1 + \mu_2}{2}$. This expression is always positive.*

**1** Introduction

**2** Error rate

**3** Influence functions of the error rate
  - First order influence function
  - Second order influence function
  - **Asymptotic relative classification efficiencies**

**4** Conclusions
  - Conclusions
  - Future research

A measure of the expected increase in error rate when estimating the optimal clustering rule from a finite sample with empirical cdf $F_n$ :

$$\text{A-Loss} = \lim_{n \to +\infty} n \, E_{F_N}[\text{ER}(F_n, F_N) - \text{ER}(F_N, F_N)].$$

As in Croux et al. (2008) :

### Proposition

*Under some regularity conditions of the clusters'centers estimators,*

$$A\text{-}Loss = \frac{1}{2} E_{F_N}[IF2(X; ER, F_N)]$$

A measure of the price one needs to pay in error rate for protection against the outliers when using a robust procedure instead of the classical one :

$$\text{ARCE(Robust,Classical)} = \frac{\text{A-Loss(Classical)}}{\text{A-Loss(Robust)}}.$$

$$F_N = 0.5\, N(-\Delta, 1) + 0.5\, N(\Delta, 1)$$

**1** Introduction

**2** Error rate

**3** Influence functions of the error rate
- First order influence function
- Second order influence function
- Asymptotic relative classification efficiencies

**4** Conclusions
- Conclusions
- Future research

**1** Introduction

**2** Error rate

**3** Influence functions of the error rate
- First order influence function
- Second order influence function
- Asymptotic relative classification efficiencies

**4** Conclusions
- **Conclusions**
- Future research

- The generalized 2-means procedure can give a more robust estimator of the error rate with a good choice of the penalty function;

- The price to pay is a loss in efficiency (depending also on the penalty function).

**1** Introduction

**2** Error rate

**3** Influence functions of the error rate
- First order influence function
- Second order influence function
- Asymptotic relative classification efficiencies

**4** Conclusions
- Conclusions
- **Future research**

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

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

(Cuesta-Albertos, Gordaliza, and Matrán, 1997);

- Other robustness properties of the generalized 2-means method defined with nondecreasing penalty function instead of strictly increasing;

- More than 1 dimension and/or more than 2 groups.

# Thank you for your attention!

- Anderson T.W., *An Introduction to Multivariate Statistical Analysis*, Wiley, New-York, 1958, pp. 126-133.

- Croux C., Filzmoser P., and Joossens K. (2008), Classification efficiences 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.

- 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., and Stahel W.A., Robust Statistics : The Approach Based on Influence Functions, John Wiley and Sons, New-York, 1986.

- Hand D.J., Daly F., Lunn A.D., McConway K.J., and Ostrowski E.,*A Handbook of Small Data Sets*, Chapman and Hall, London, 1991.

Influence
function of
the error rate
of
classification
based on
clustering

Ch. Ruwet

Introduction

Error rate

IF of the
error rate

Conclusion
Conclusion
Future
research

- Pollard D., Strong Consistency of k-Means Clustering, *The Annals of Probability*, 1981, Vol.9, nř4, pp.919-926.
- Pollard D., A Central Limit Theorem for k-Means Clustering, *The Annals of Probability*, 1982, Vol.10, nř1, 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.