

# A spatial classification : application to ecological data

L. Bel<sup>1</sup>, A. Bar-Hen<sup>3</sup>, D. Allard<sup>4</sup>J-M., Laurent<sup>2</sup>, R. Cheddadi<sup>2</sup>

<sup>1</sup>Probabilités, Statistique et Modélisation, Université Paris-Sud, Orsay, France

<sup>2</sup>Institut des Sciences de l'Evolution, CNRS and Université Montpellier II, France

<sup>3</sup>Institut National d'Agronomie, Paris, France

<sup>4</sup>Institut National de la Recherche Agronomique, Unité de Biométrie, Avignon, France

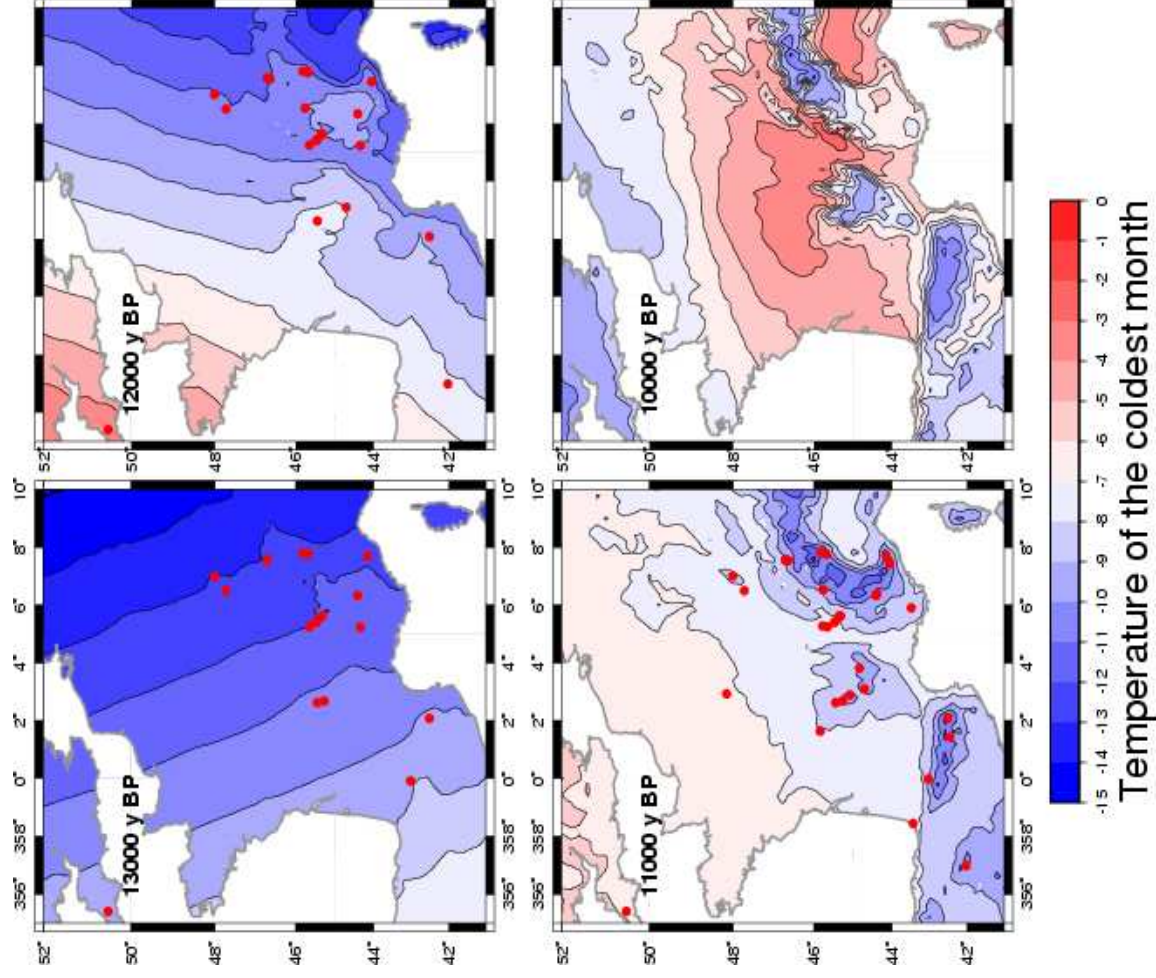
# Paleoecological Questions

**Paleoecology** : reconstructing past environments and their evolution.  
Pollen frequencies reproduce plant ranges.

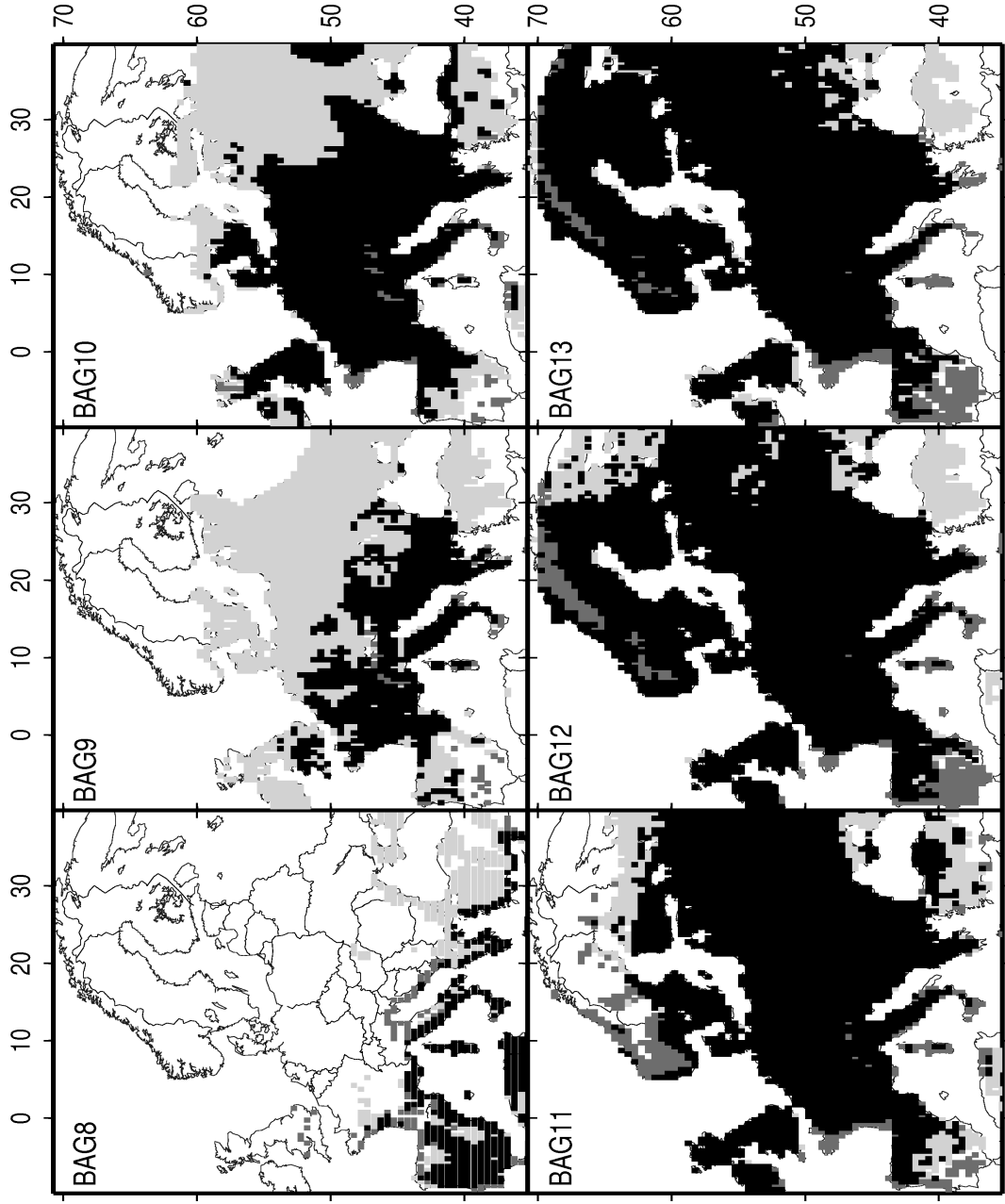
**Goal** : discriminate between absence, presence and abundance of plant using pollen frequencies.

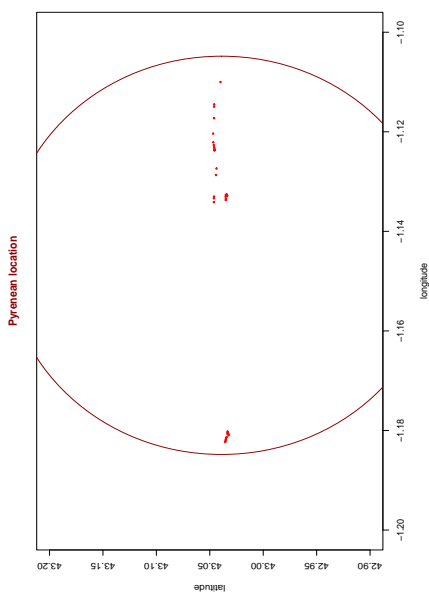
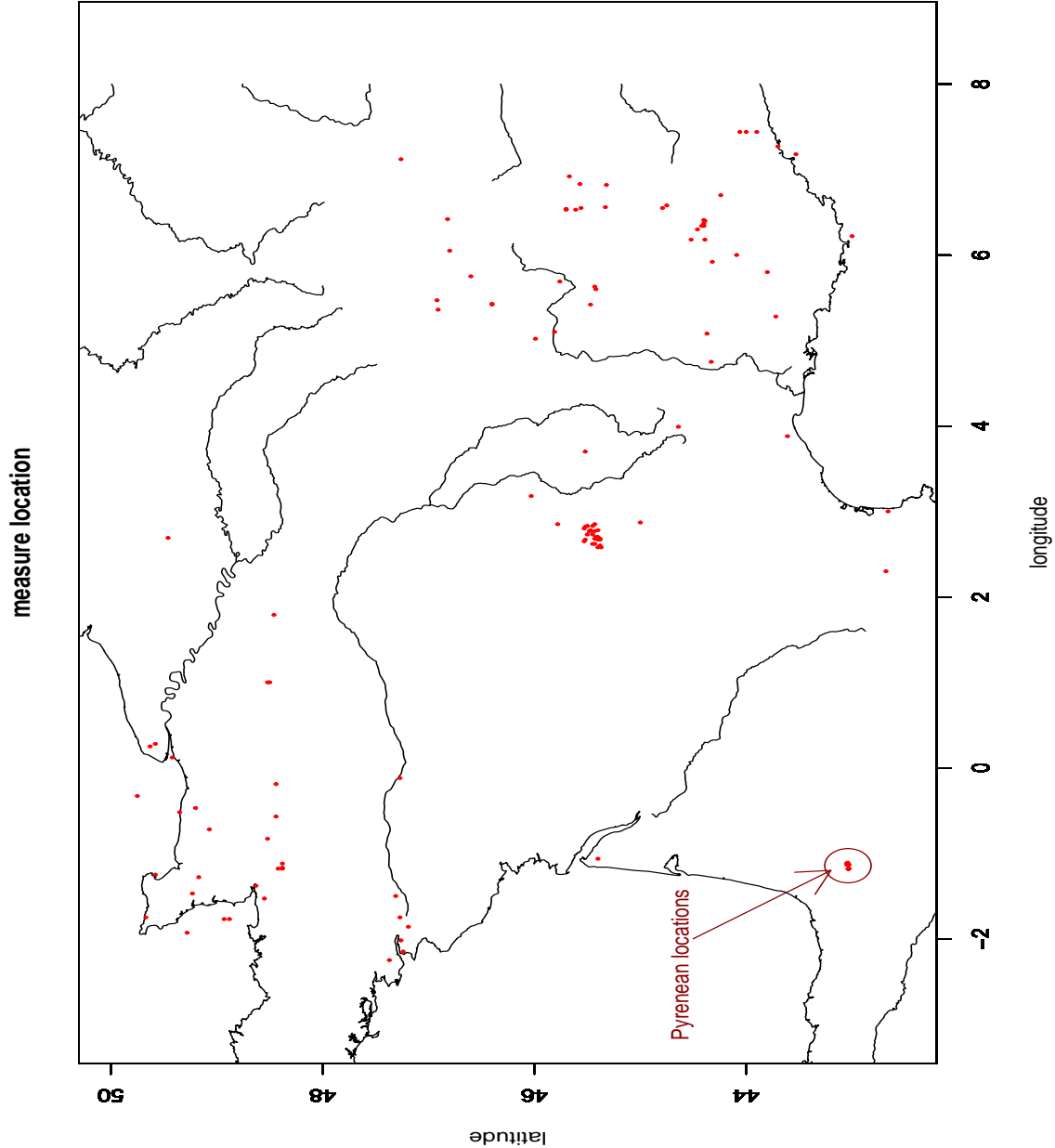
**Supervised Classification** with actual data as training set.  
Sampling spatial scheme of pollen data very irregular.

# Example of climate reconstruction



# European distribution of BAGs





# CART

## Classification And Regression Tree

$X^1, \dots, X^p$ : explanatory variables,  $Y$ : categorical explained variable  
Binary tree through binary recursive partitioning

Example: predict Ozone pollution level from climatic variables

### 3 variables :

temperature (temp)

wind (vent)

O<sub>3</sub> previous day (pollj)

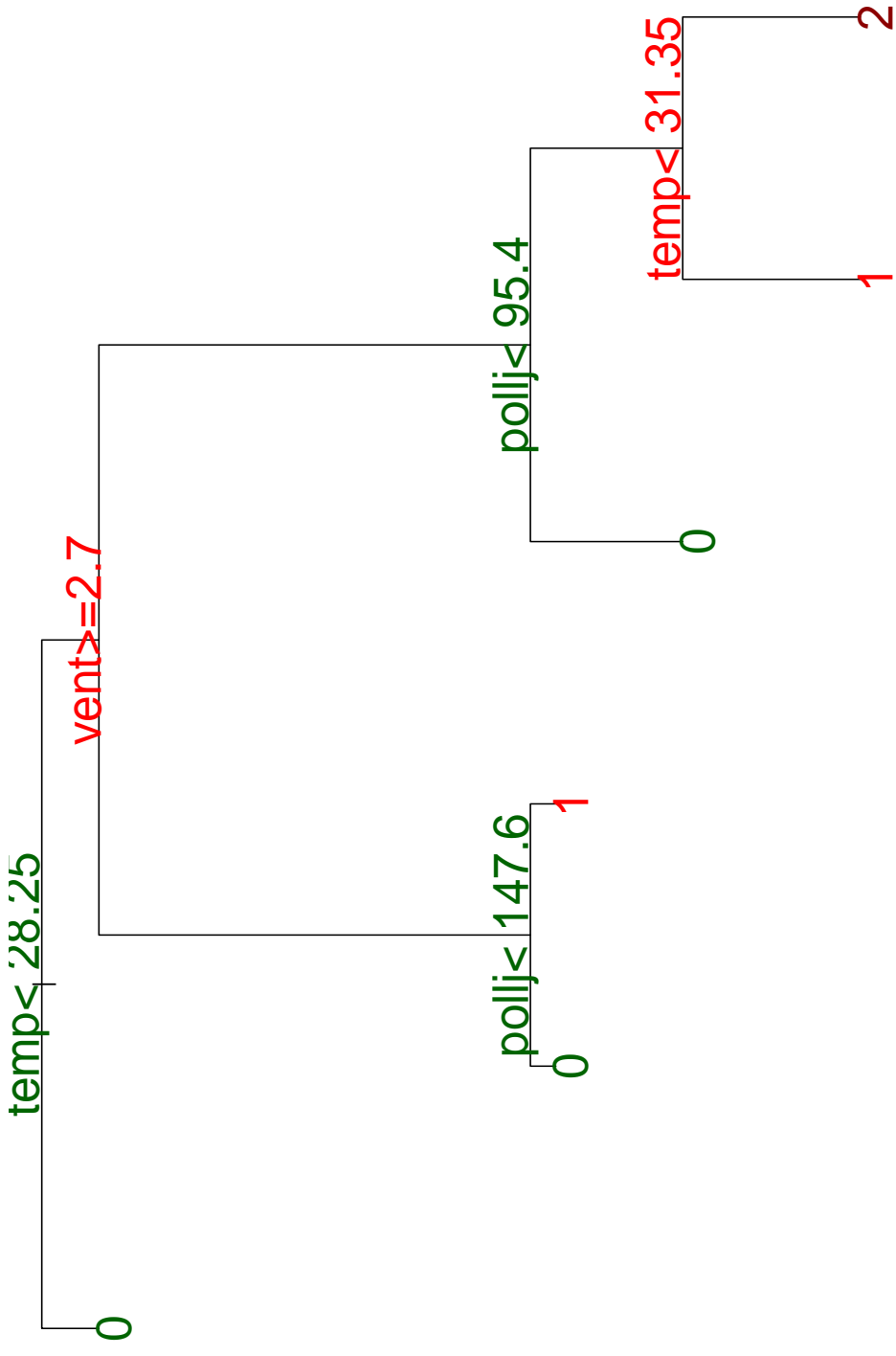
### 3 classes :

0 : low pollution

1 : medium pollution

2 : high pollution

## CART EXAMPLE : pollution level



# CART : algorithm

- Heterogeneity criterion

$$\text{Gini Index : } D_l = \sum_{i,j} p_{il}p_{jl} = 1 - \sum_i p_{il}^2$$

- Splitting rule

$$D_l - p_{l+}D_{l+} - p_{l-}D_{l-} > 0$$

- Criterion error prediction  $R(T) = P[T(X) \neq Y]$
- Penalized criterion  $S(T) = \text{argmin}\{\hat{R}(T) + \mu \cdot \text{size}(T)\}$



# Weighting CAR I

- Classically assume independence of the data

$$p_{il} = \frac{n_{il}}{n+l} \quad \hat{R}(T) = \frac{1}{n} \sum_{\alpha=1}^n \mathbb{I}(T(X_\alpha) \neq Y_\alpha)$$

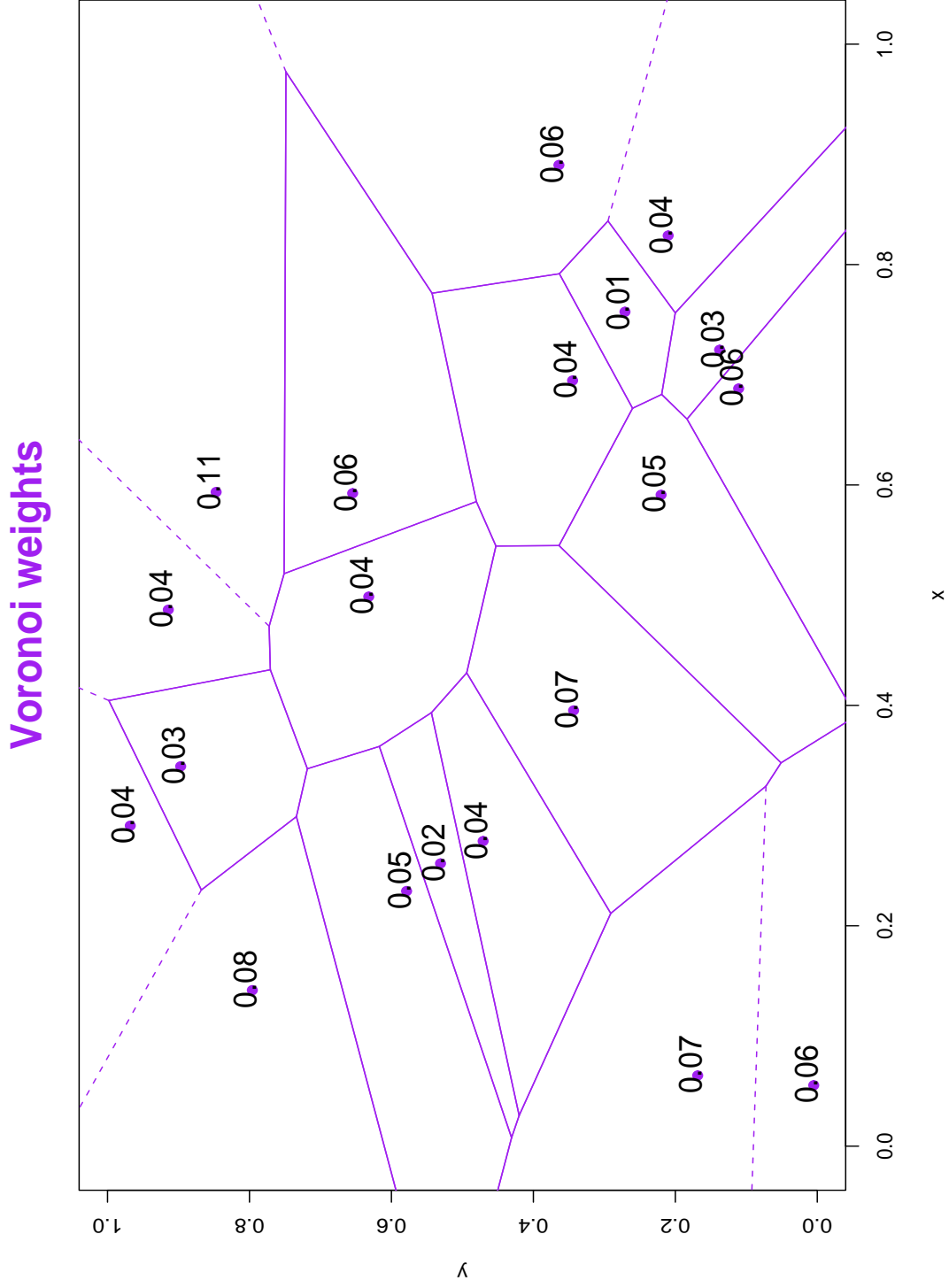
- Dependence through weights for individuals

$$p_{il} = \sum_{\alpha \in l} w_\alpha \mathbb{I}(Y_\alpha = i) \quad \hat{R}(T) = \sum_{\alpha=1}^n w_\alpha \mathbb{I}(T(X_\alpha) \neq Y_\alpha)$$

## Possible weights:

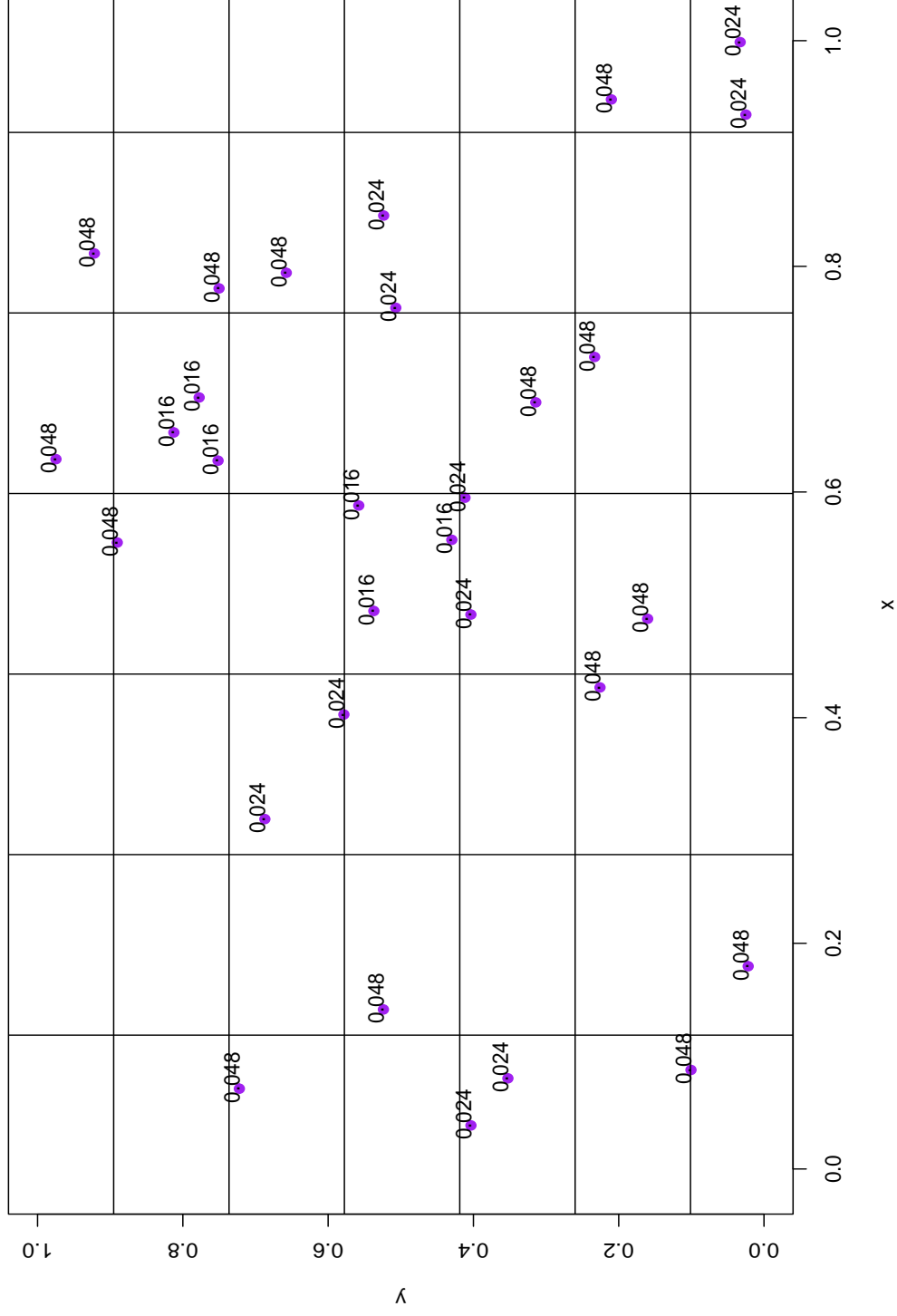
- Voronoï tessellation
- constant for each square of a grid
- kriging weights

# Voronoi tessellation



# Grid method

## Grid weights



# Kriging Weights

Estimating the mean over a domain  $\mathcal{D}$  of dependent data  $(X_\alpha)_{\alpha=1,n}$   
Best (minimal variance) Linear Unbiased Estimator

$$\hat{\mu} = \sum_{\alpha=1}^n \lambda_\alpha X_\alpha \quad \Lambda = \begin{pmatrix} (\lambda_\alpha)_{\alpha=1,n} \\ m \end{pmatrix}$$

with  $\tilde{C}\Lambda = \tilde{c}$  and

$$\tilde{C} = \begin{pmatrix} C & 1 \\ 1 & 0 \end{pmatrix} \quad \tilde{c} = \begin{pmatrix} c(x_\alpha, \mathcal{D}) \\ 1 \end{pmatrix}$$

$$C_{\alpha,\beta} = \text{Cov}(X_\alpha, X_\beta)$$

$$c(x_\alpha, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \int_{\mathcal{D}} \text{Cov}(x_\alpha, y) dy.$$

Kriging of the mean :  $(\lambda_\alpha)_{\alpha=1,n}$  Kriging Weights  
with the constraint  $\lambda_\alpha > 0$

“natural declustering”

# Spatial CART

$p_{il}$  estimation of  $P(Y = i \mid X \in B_l)$

$$P(Y = i \mid X \in B_l) = E(\mathbb{I}_{Y=i} \mid X \in B_l)$$

estimated by kriging the mean of  $\mathbb{I}_{Y=i} \mid X \in B_l$ .

$$\hat{p}_{il} = \sum_{\alpha: X_\alpha \in B_l} \lambda_\alpha \mathbb{I}(Y_\alpha = i)$$

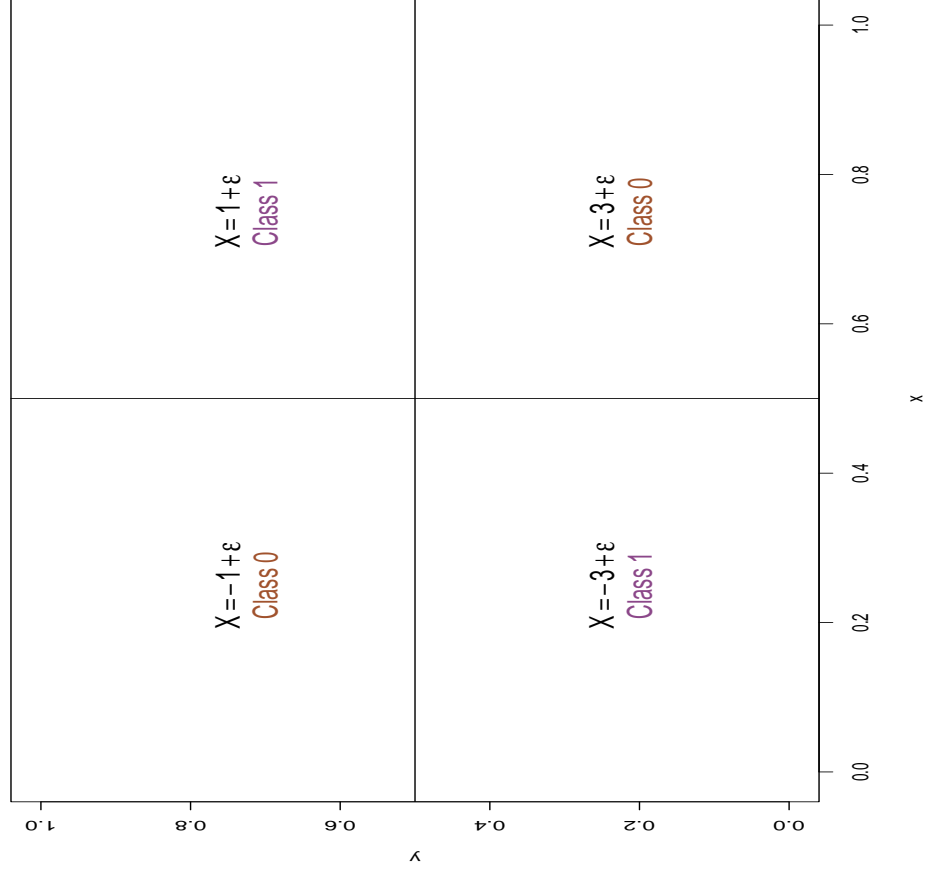
$R(T) = E[\mathbb{I}_{T(X) \neq Y}]$  kriging the mean of  $\mathbb{I}_{T(X) \neq Y}$ .

$$\hat{R}(T) = \sum_{\alpha} \lambda_\alpha \mathbb{I}(T(X_\alpha) \neq Y_\alpha)$$

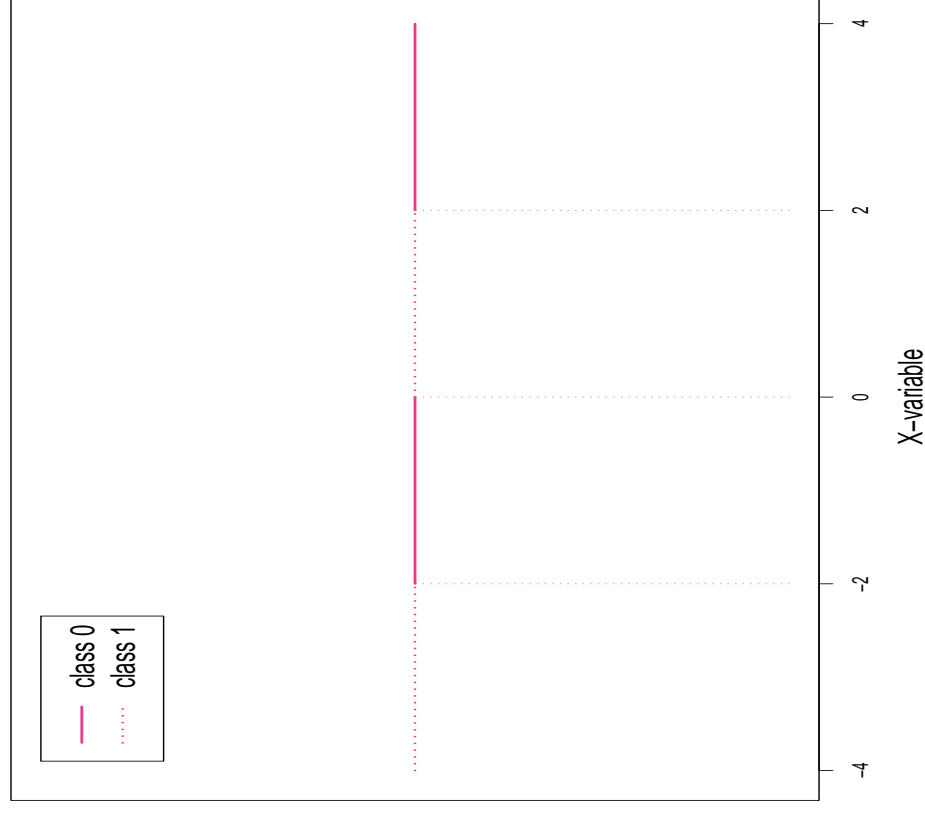
# Simulations

## Simulation design

simulation design



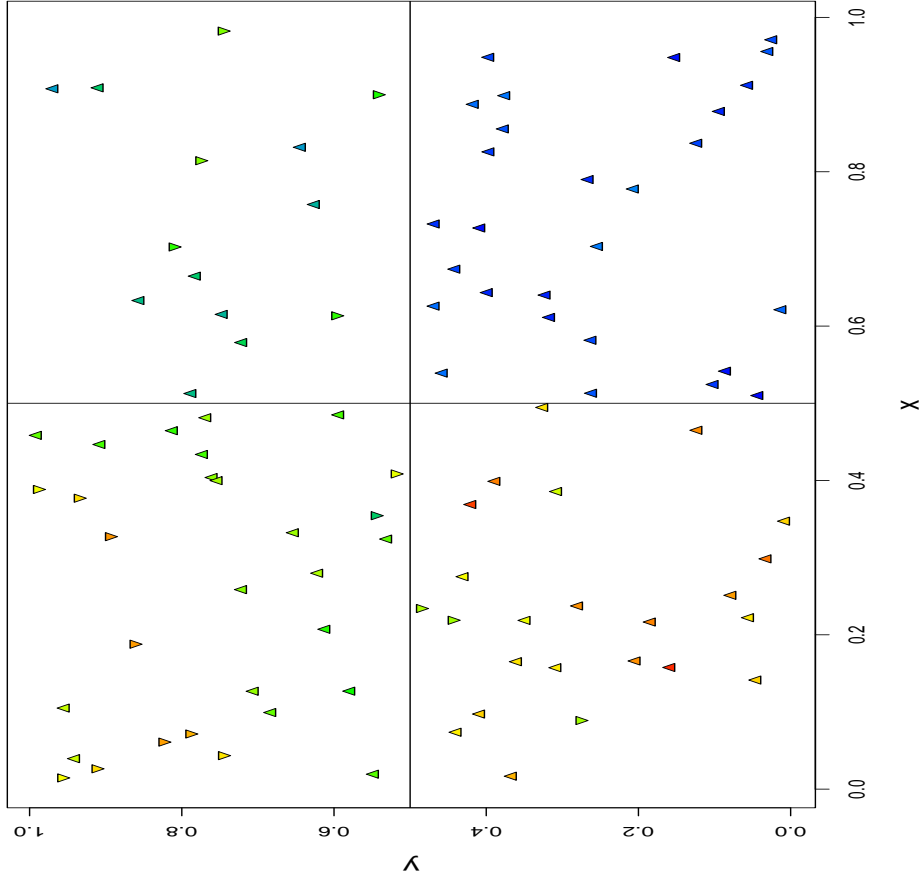
Expected classification rule



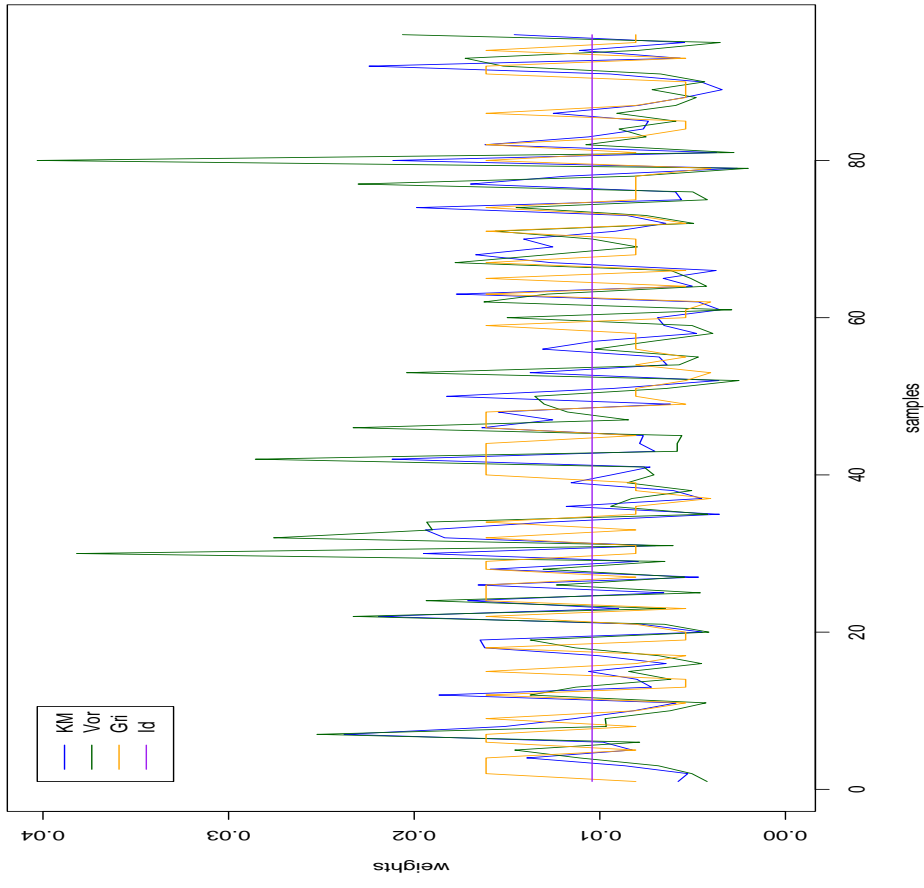
# Simulations

## Simulation 0 clusters

a) simulation with 0 clusters



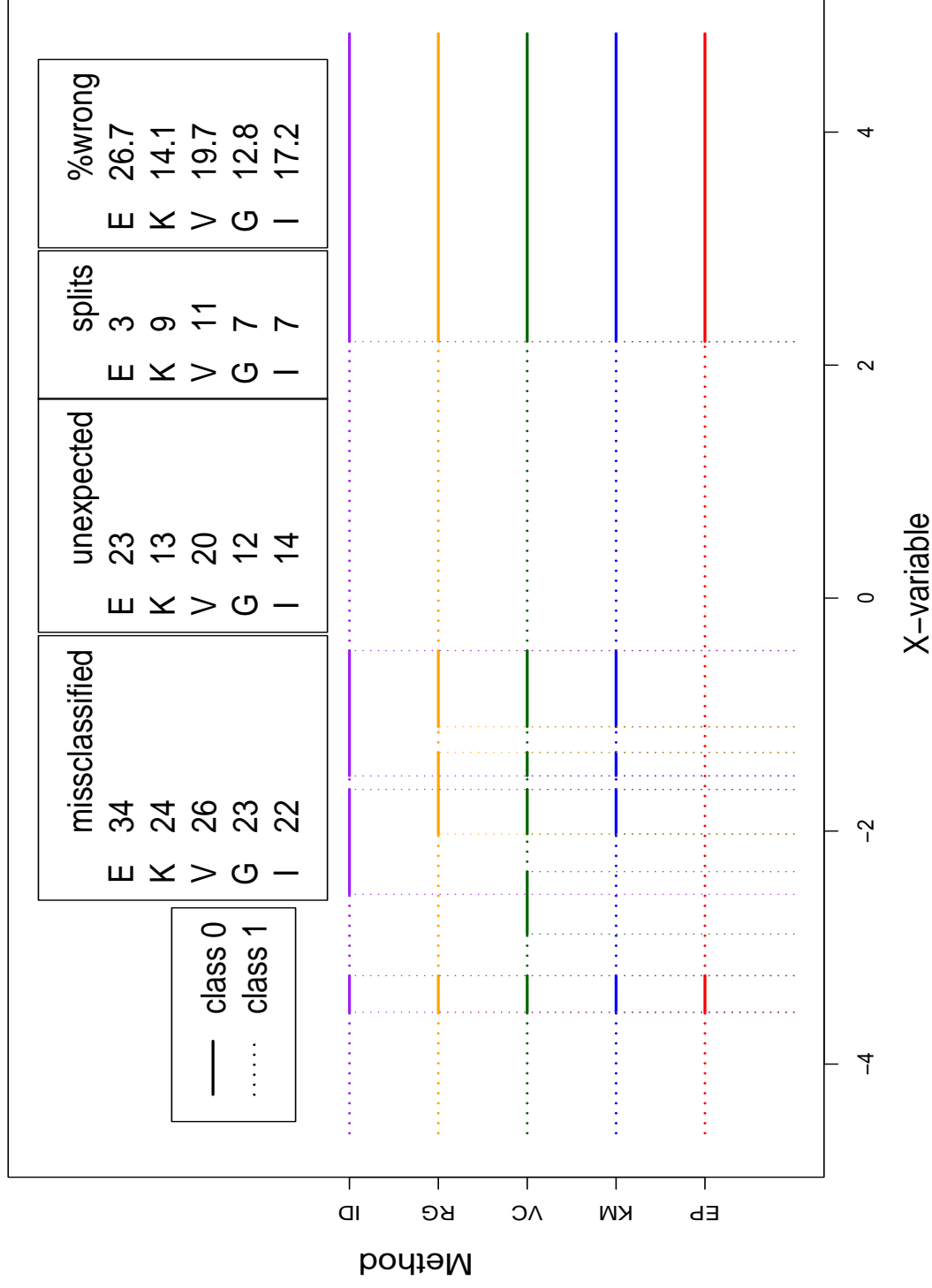
weights comparison, simulation with 0 clusters



# Simulations

## Simulation 0 clusters

### b) Comparison of the 5 methods, simulation with 0 clusters

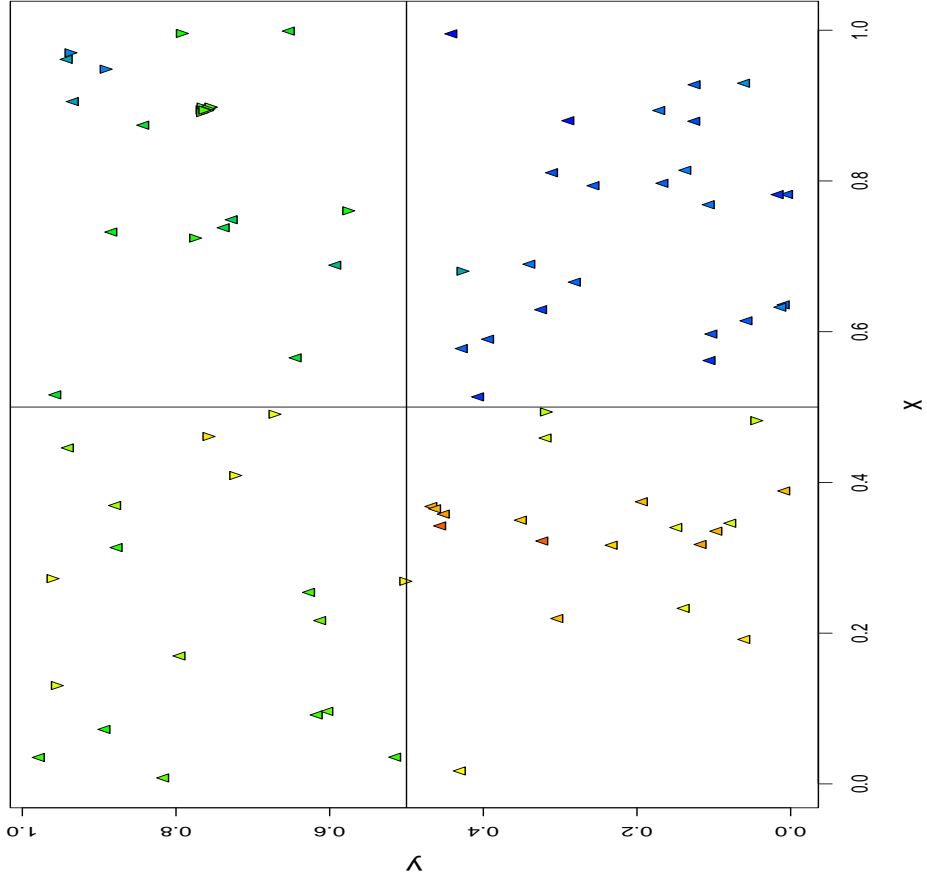




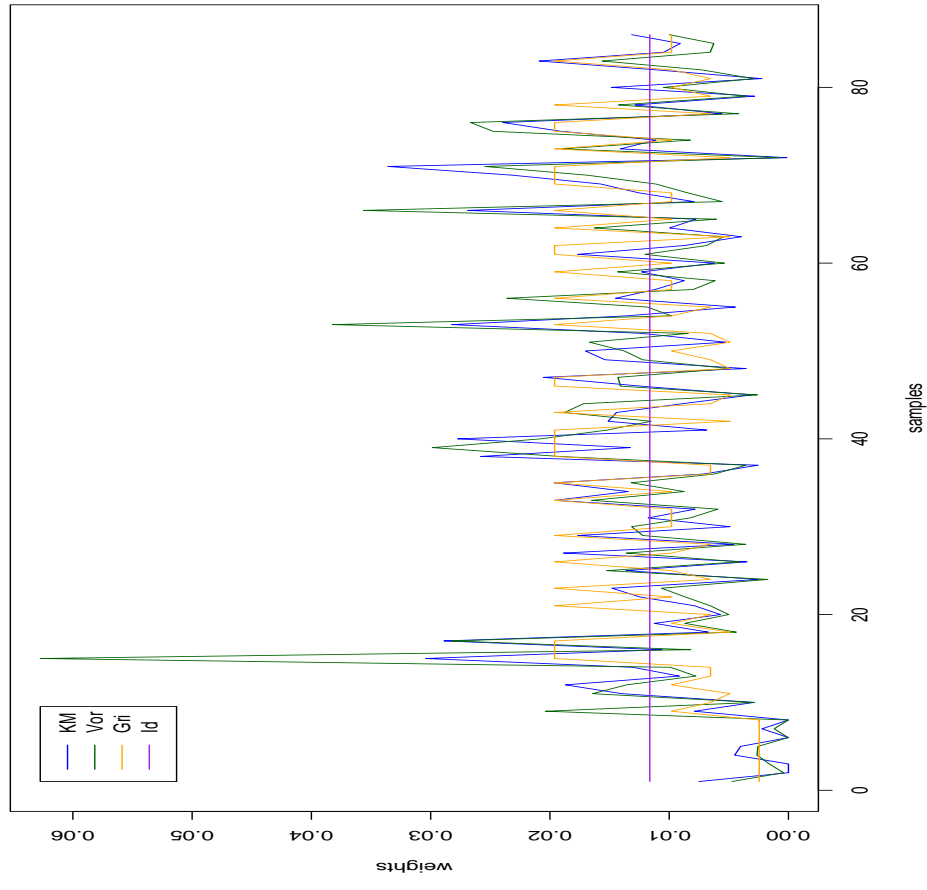
# Simulations

## Simulation 1 clusters

a) simulation with 1 clusters



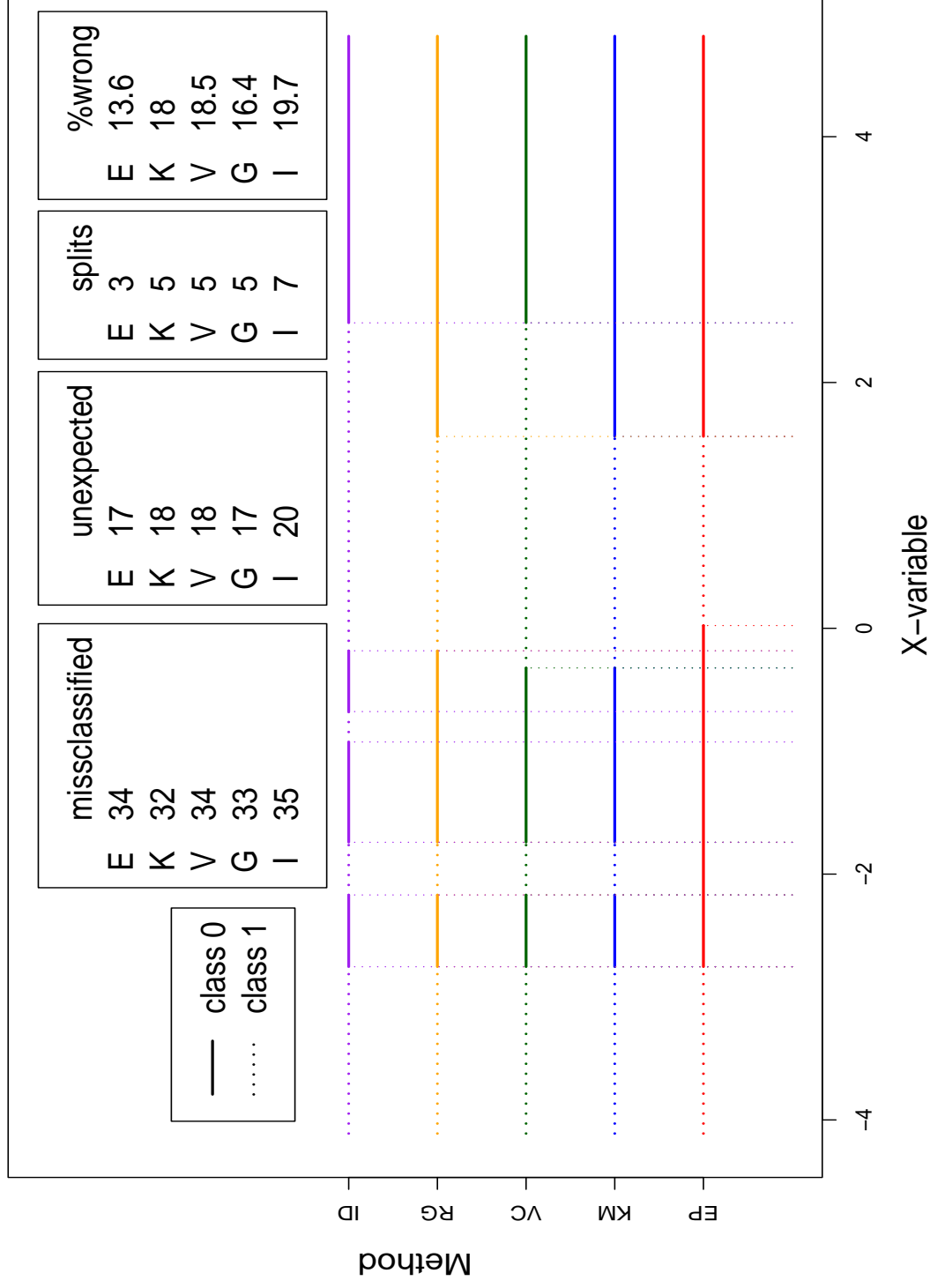
weights comparison, simulation with 1 clusters



# Simulations

## Simulation 1 clusters

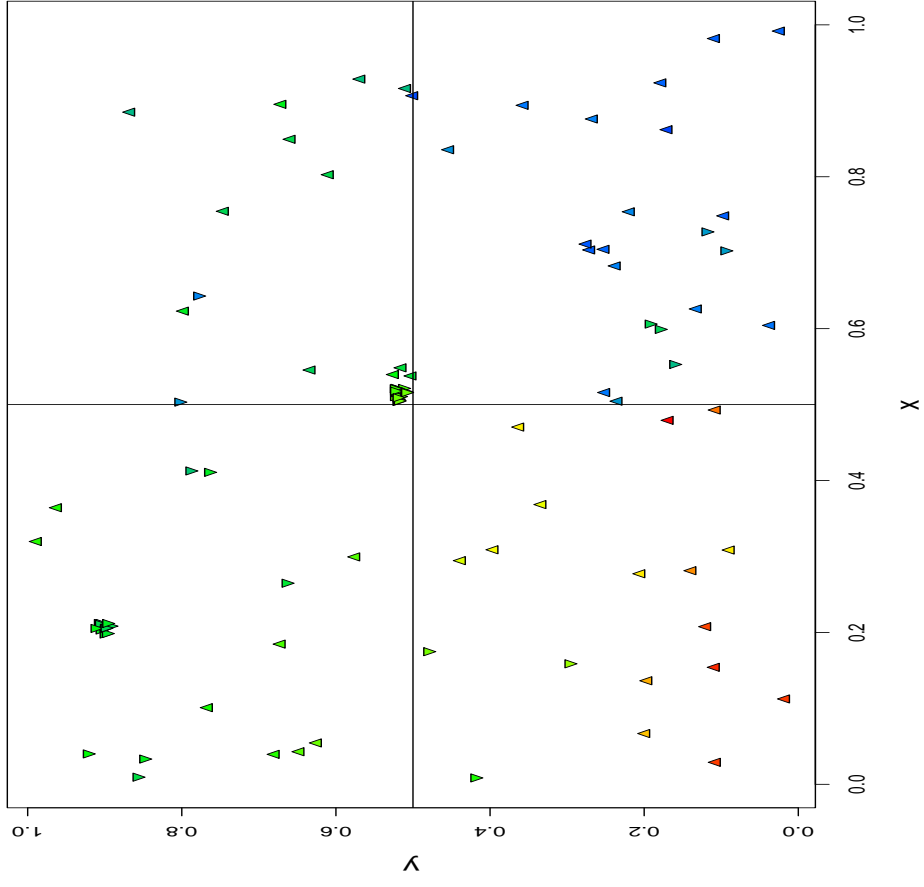
### b) Comparison of the 5 methods, simulation with 1 clusters



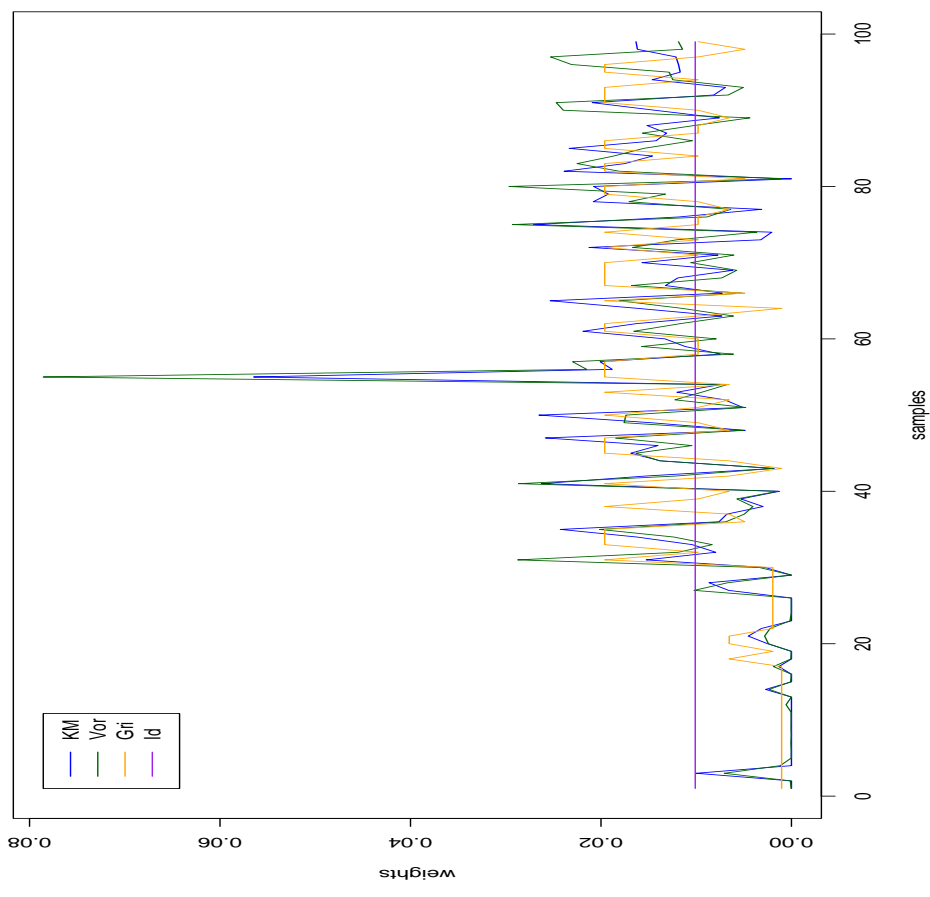
# Simulations

## Simulation 2 clusters

a) simulation with 2 clusters



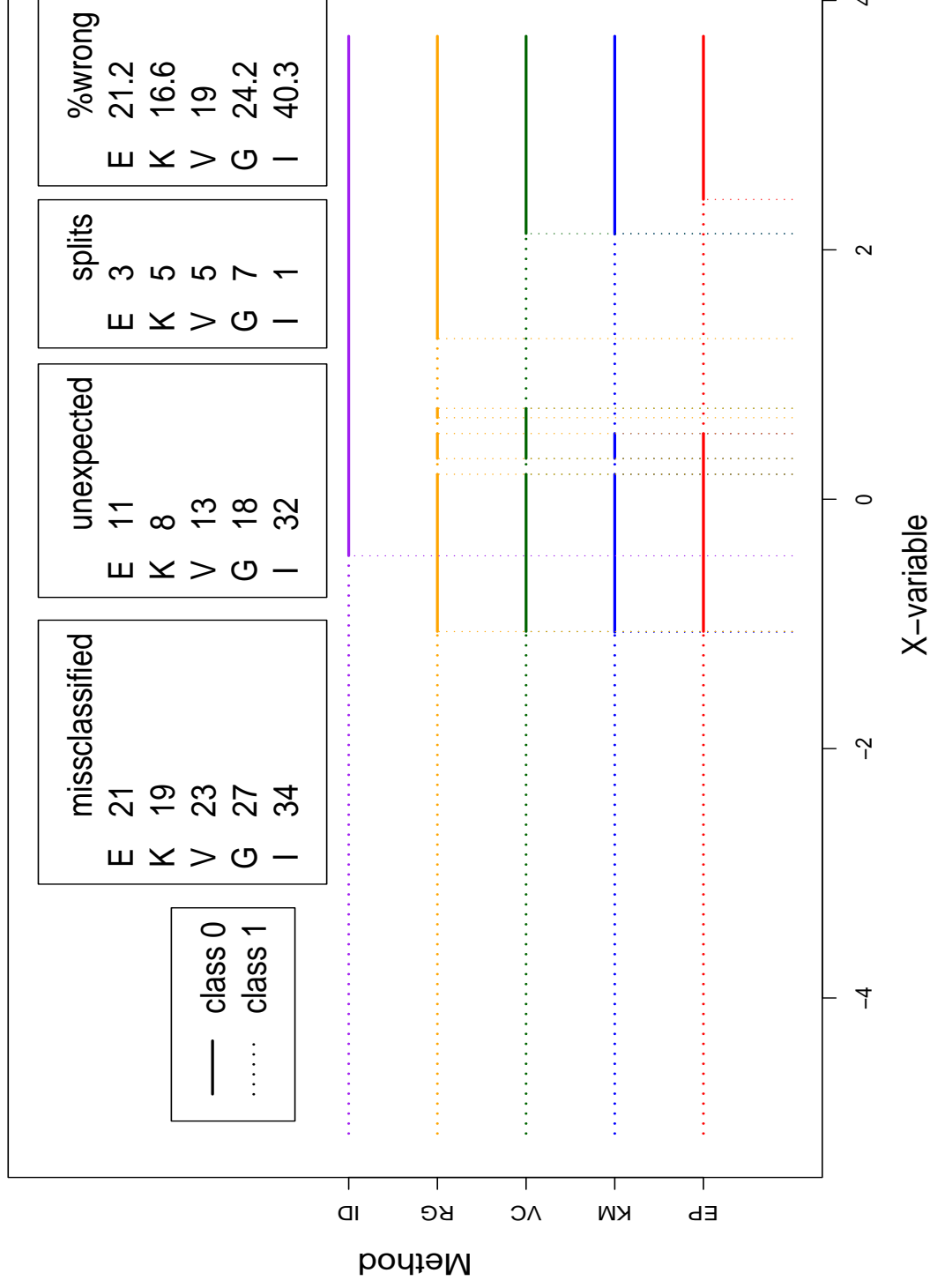
weights comparison, simulation with 2 clusters



# Simulations

## Simulation 2 clusters

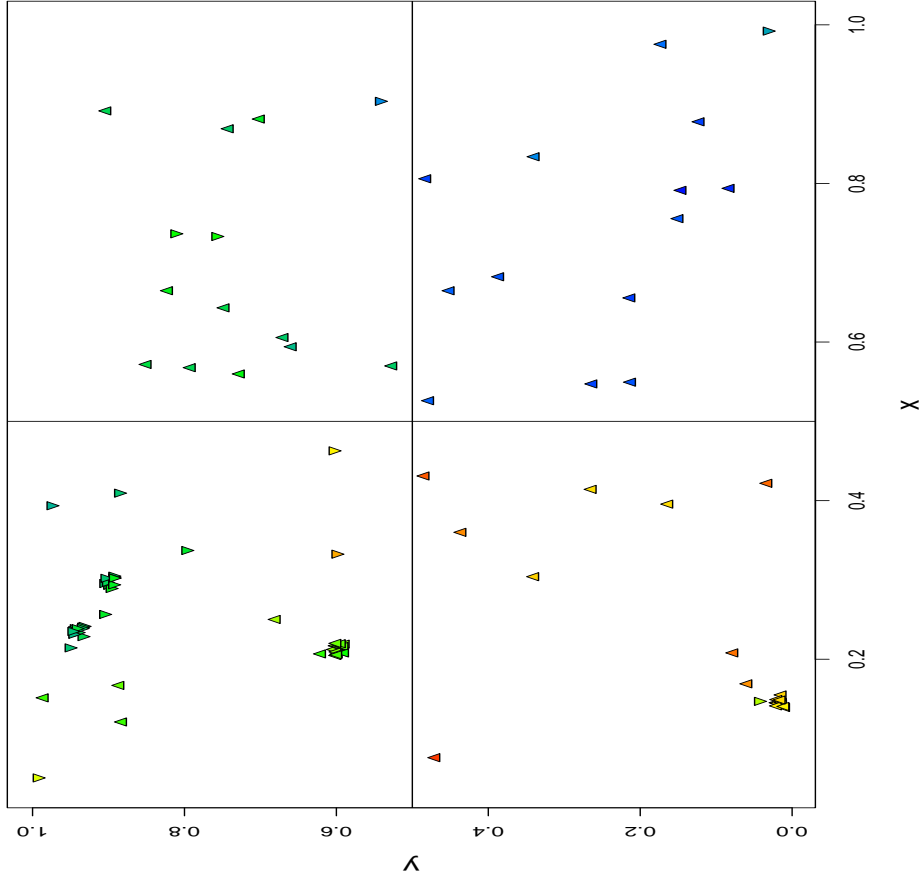
### b) Comparison of the 5 methods, simulation with 2 clusters



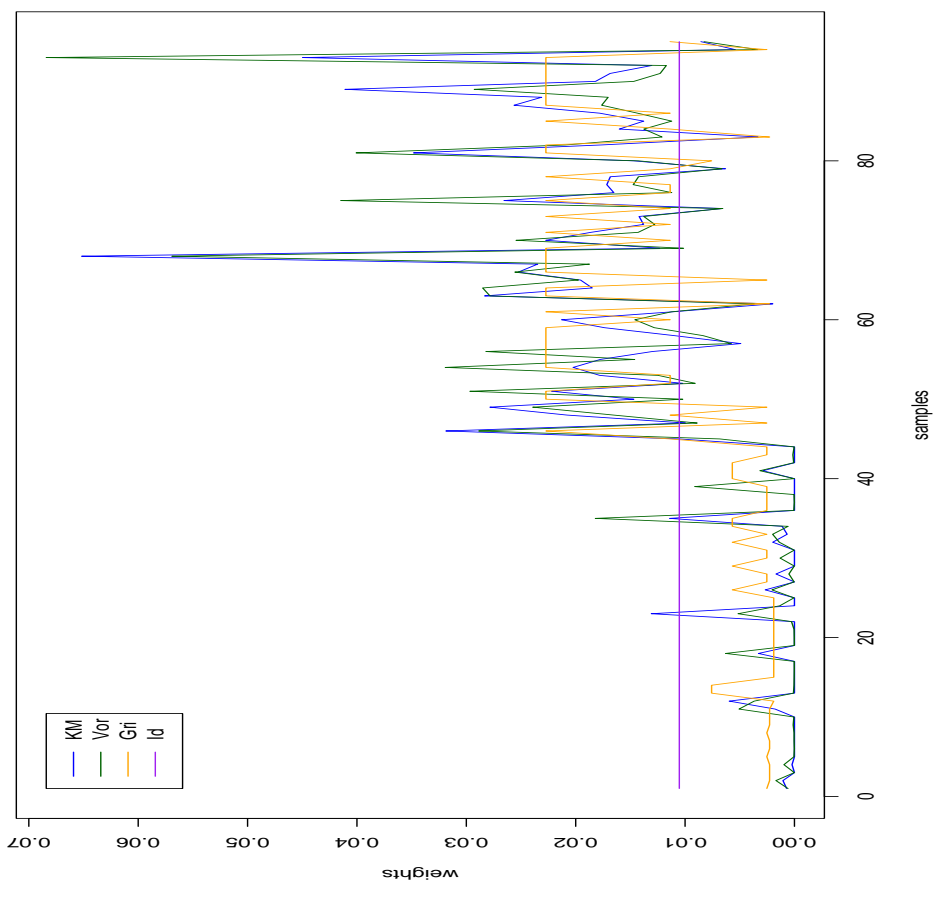
# Simulations

## Simulation 4 clusters

a) simulation with 4 clusters



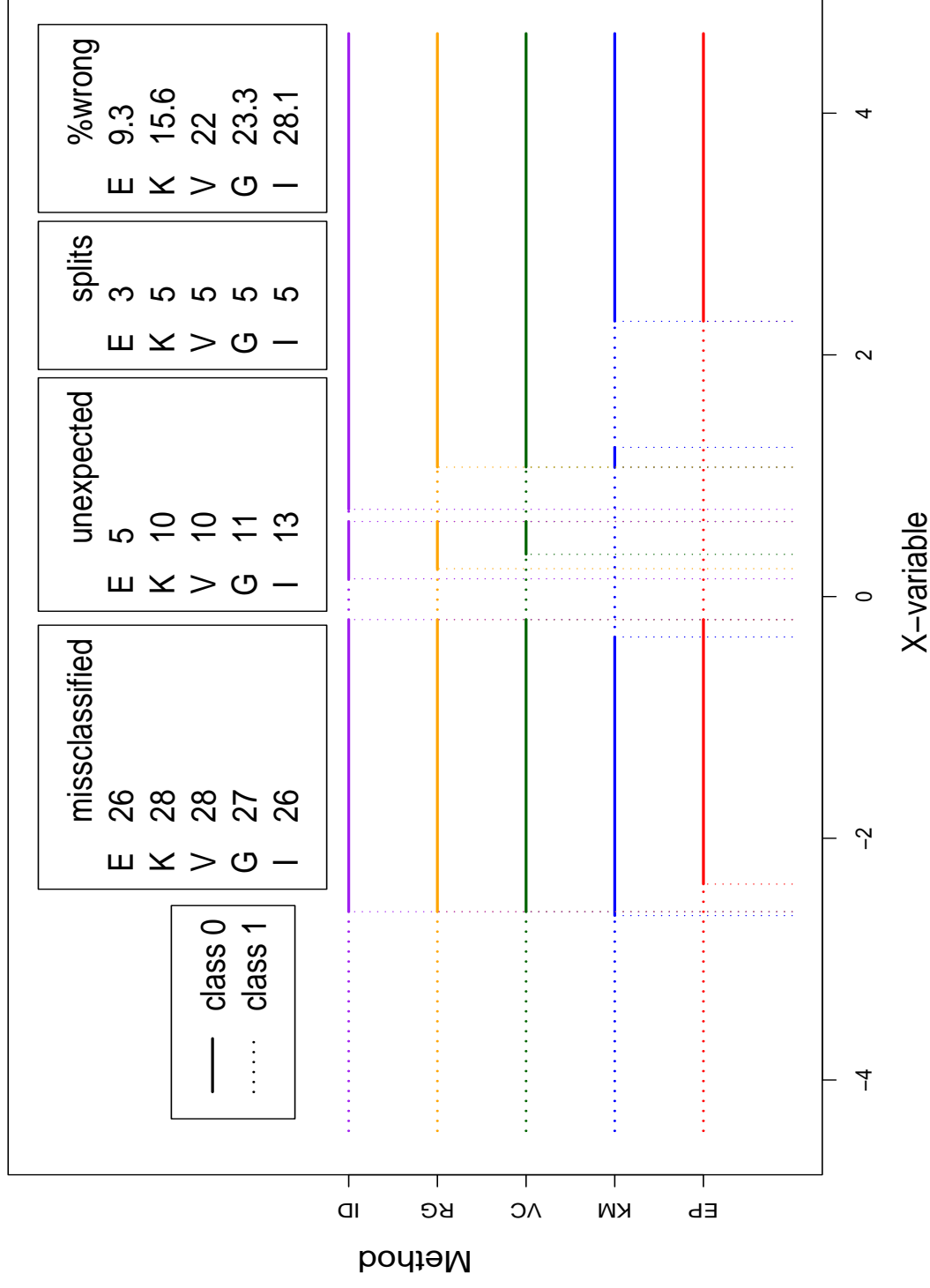
weights comparison, simulation with 4 clusters



# Simulations

## Simulation 4 clusters

### b) Comparison of the 5 methods, simulation with 4 clusters



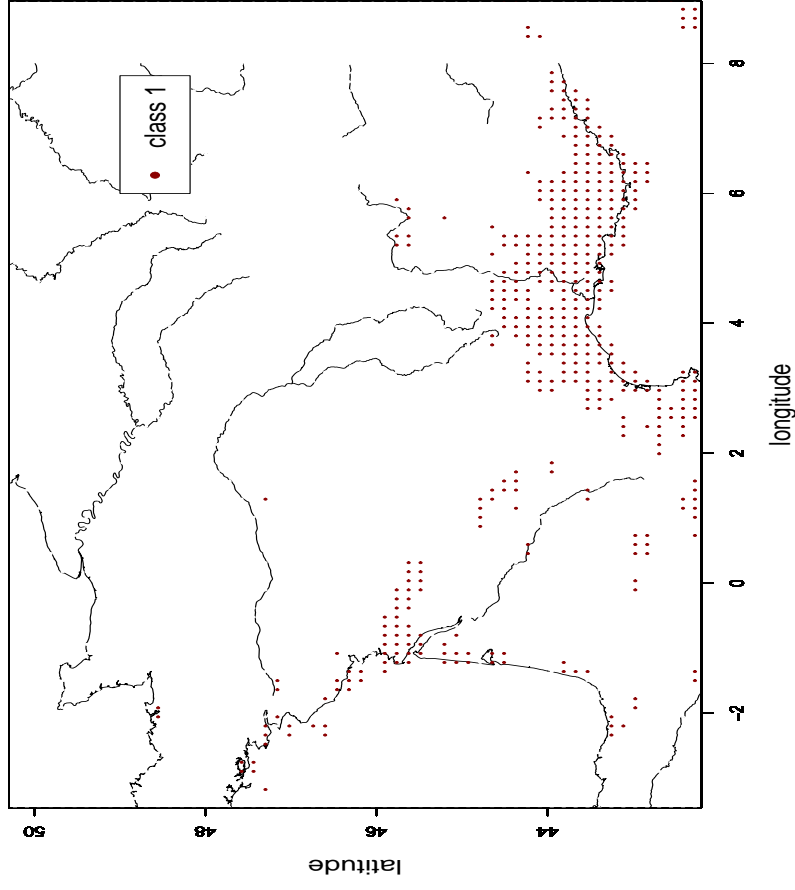
# Simulations

Statistics on missclassified points on test samples.

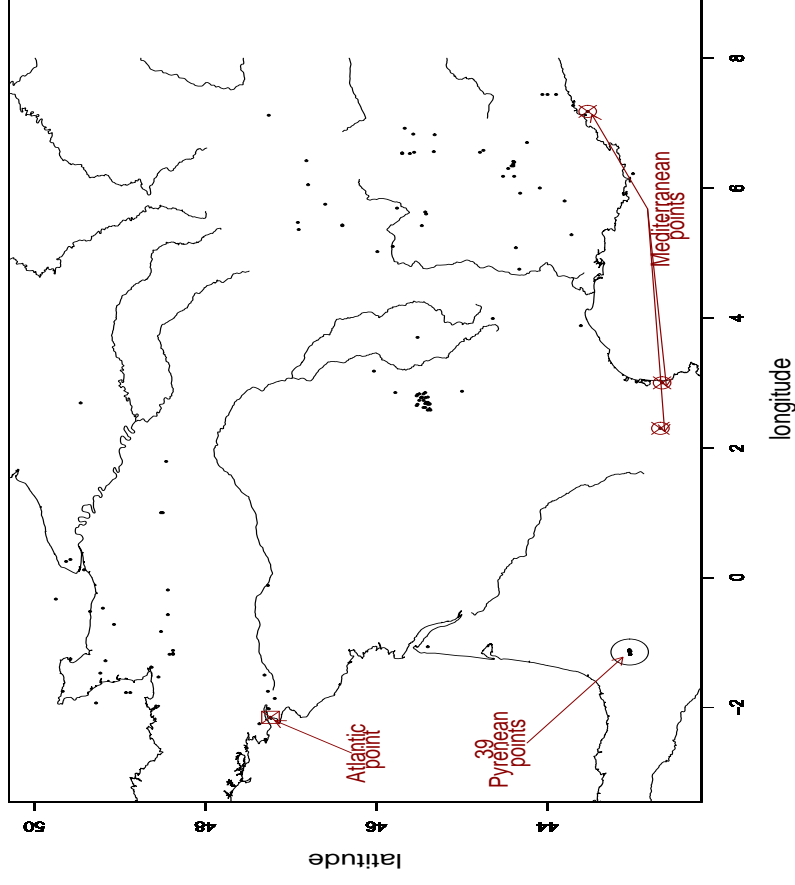
| clusters | EP          |             | KM          |             | Vor         |             | Gri         |             | Id          |             |
|----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|          | total       | rule        | total       | rule        | total       | rule        | total       | rule        | total       | rule        |
| 0        | <b>32.8</b> | <b>17.7</b> | 29.9        | 12.7        | 30.4        | 13.3        | 30.3        | 11.9        | <b>29.7</b> | <b>10.8</b> |
| 1        | 28.3        | 14.5        | 27.7        | 12.1        | <b>26.5</b> | <b>10</b>   | 27.8        | 12.8        | <b>29.8</b> | <b>16.5</b> |
| 2        | 28.3        | 14.1        | 25.9        | <b>10.8</b> | 26.9        | 11.7        | <b>26.3</b> | <b>10.8</b> | <b>28.9</b> | <b>15.6</b> |
| $\geq 3$ | 29.3        | 14.4        | <b>27.6</b> | 11.3        | 27.9        | <b>10.7</b> | 28.5        | 11.3        | <b>31.9</b> | <b>17.6</b> |

# Ecological Data

a) presence/absence SOPHY8



b) pollen sample sites



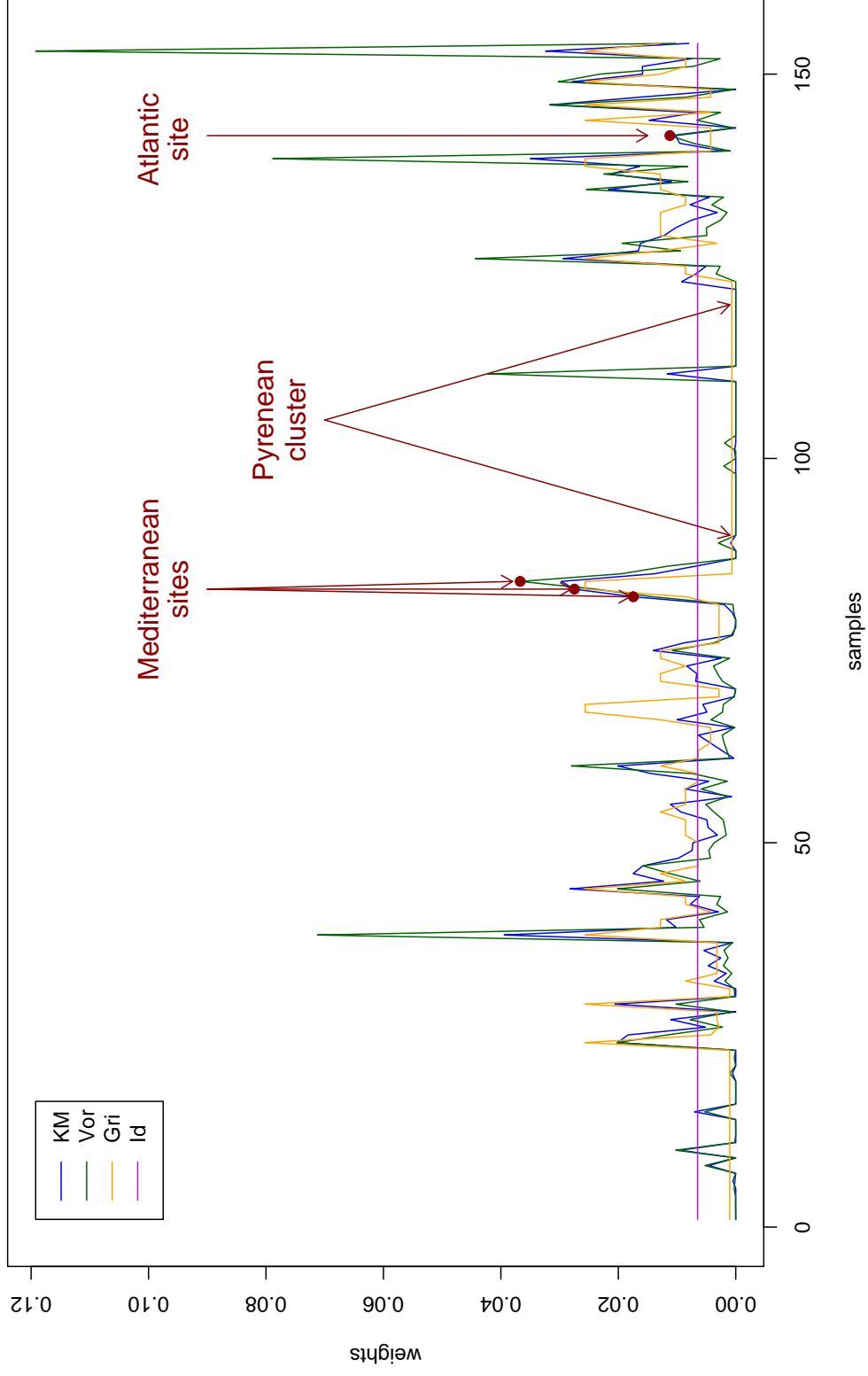
a) Presence and absence of specie from BAG8;

b) boxplots of pollen frequencies according to presence or absence



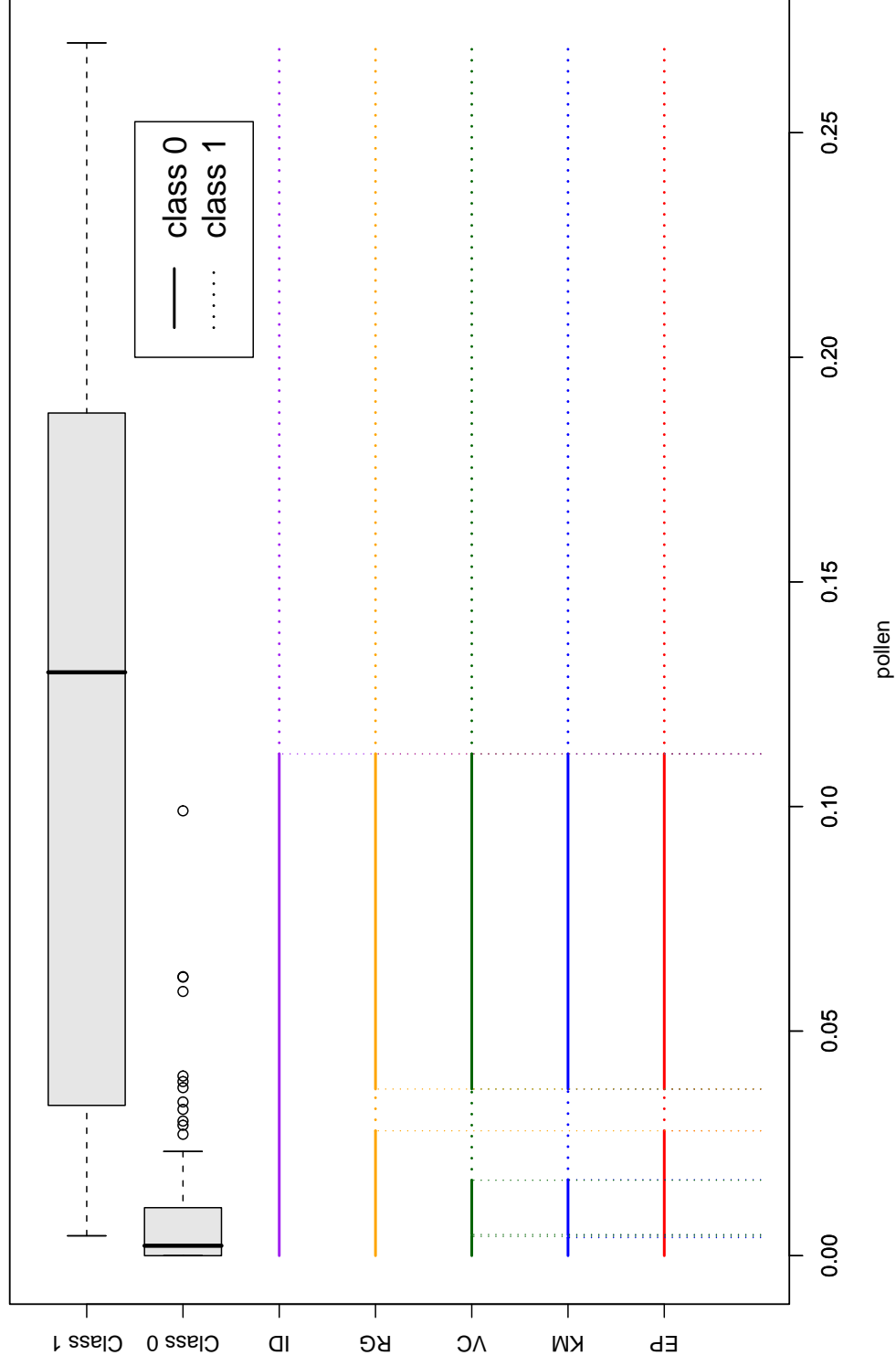
# Weights

## BAG8 weights comparison



# Partitionning

BAG8 : comparaison of the five methods



# Interpretation

- pollen frequency  $> 0.11$ , classified as presence for all methods.
- presence sites with low pollen frequency:
  - one at the border: large weight and well classified for Voronoi and Kriging the Mean.
  - 3 have similar pollen frequencies as the Pyrenean locations, but are isolated
    - \* Standard CART forget them
    - \* Spatial methods create a new cut

# Conclusion

## Weighted and Spatial CART

- discard clustered points of the analysis based on proximity or correlation properties
- provide weights easily interpretable
- highlight the effect of disturbance on the clustered points
- give a better understanding of the sampling scheme for the ecological data

Mostly useful when clustered data perturb the classification rule.