

# Unsupervised learning: Anomaly detection

## Part I: Multivariate data

Pavlo Mozharovskyi

LTCI, Telecom Paris, Institut Polytechnique de Paris

Parcours Data Science BPCE

Paris, the 13th of June 2023

# Contents

## Introduction

## Non-parametric approaches

- One-class support vector machines

- Local outlier factor

- Isolation forest

## Systematic orderings: data depth

- The notion of data depth

- The Tukey depth function

- Central regions

- Further depth notions

## Practical session

# Contents

## Introduction

### Non-parametric approaches

- One-class support vector machines

- Local outlier factor

- Isolation forest

### Systematic orderings: data depth

- The notion of data depth

- The Tukey depth function

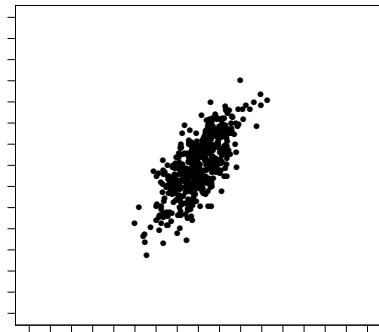
- Central regions

- Further depth notions

## Practical session

## A real task

Regard two measurements during a test in a production process:

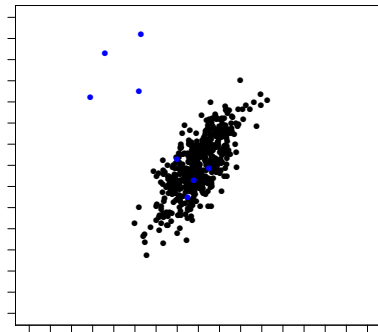


Given **training data**, polluted or not with anomalies:

- ▶ detect **anomalies** in the given data.

## A real task

Regard two measurements during a test in a production process:



Given **training data**, polluted or not with anomalies:

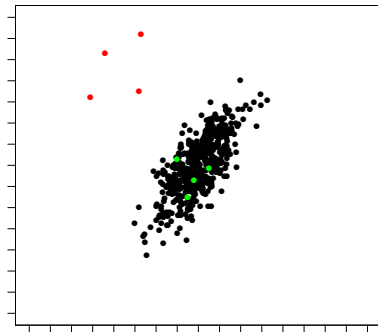
- ▶ detect **anomalies** in the given data.

For **new data**, determine:

- ▶ Whether new observations are **normal** data or **anomalies**?

## A real task

Regard two measurements during a test in a production process:



Given **training data**, polluted or not with anomalies:

- ▶ detect **anomalies** in the given data.

For **new data**, determine:

- ▶ Whether new observations are **normal** data or **anomalies**?

# Multivariate framework

- ▶ A training data set:

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$$

of observations in the  $d$ -dimensional Euclidean space.

# Multivariate framework

- ▶ A training data set:

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$$

of observations in the  $d$ -dimensional Euclidean space.

- ▶ Typical example: a table from a data base, with lines being observations (=individuals, items,...).



# Multivariate framework

- ▶ A training data set:

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$$

of observations in the  $d$ -dimensional Euclidean space.

- ▶ Typical example: a table from a data base, with lines being observations (=individuals, items,...).
- ▶ Construct a decision function:

$$\mathbb{R}^d \rightarrow \{-1, +1\} : \mathbf{x} \mapsto g(\mathbf{x}),$$

which attributes to any (possible)  $\mathbf{x} \in \mathbb{R}^d$  a label whether it is an anomaly (e.g., +1) or a normal observation (e.g., -1).

# Multivariate framework

- ▶ A training data set:

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$$

of observations in the  $d$ -dimensional Euclidean space.

- ▶ Typical example: a table from a data base, with lines being observations (=individuals, items,...).
- ▶ Construct a decision function:

$$\mathbb{R}^d \rightarrow \{-1, +1\} : \mathbf{x} \mapsto g(\mathbf{x}),$$

which attributes to any (possible)  $\mathbf{x} \in \mathbb{R}^d$  a label whether it is an anomaly (e.g., +1) or a normal observation (e.g., -1).

- ▶ It is more useful to provide an ordering on  $\mathbb{R}^d$ :

$$\mathbb{R}^d \rightarrow \mathbb{R} : \mathbf{x} \mapsto g(\mathbf{x}),$$

such that abnormal observations obtain higher anomaly score.

# Practical session (parts I and II)

## Notebooks:

- ▶ `anomdet_simulation1.Rmd`,
- ▶ `anomdet_hurricanes.Rmd`,
- ▶ `anomdet_cars.ipynb`,
- ▶ `anomdet_airbus.ipynb`.

## Data sets:

- ▶ `carsanom.csv`: Data set on anomaly detection for cars.
- ▶ `airbus_data.csv`: Data set from Airbus.
- ▶ `hurdat2-1851-2019-052520.txt`: Historical hurricane data.

## Supplementary scripts:

- ▶ `depth_routines.py`: Routines for data depth calculation.
- ▶ `FIF.py`: Implementation of the functional isolation forest.
- ▶ `depth_routines.R`: Routines for curves' parametrization.

# Contents

## Introduction

## Non-parametric approaches

- One-class support vector machines

- Local outlier factor

- Isolation forest

## Systematic orderings: data depth

- The notion of data depth

- The Tukey depth function

- Central regions

- Further depth notions

## Practical session

# Contents

## Introduction

## Non-parametric approaches

- One-class support vector machines

- Local outlier factor

- Isolation forest

## Systematic orderings: data depth

- The notion of data depth

- The Tukey depth function

- Central regions

- Further depth notions

## Practical session

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Generalized portrait:

- ▶ The method of the **generalized portrait** was introduced by Vapnik & Lerner (1963) and Vapnik & Chervonenkis (1974).

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Generalized portrait:

- ▶ The method of the **generalized portrait** was introduced by Vapnik & Lerner (1963) and Vapnik & Chervonenkis (1974).
- ▶ Generalized portrait is the vector:

$$\psi = \frac{\varphi}{\min_{\mathbf{x} \in \mathbf{X}} \langle \mathbf{x}, \varphi \rangle} \quad \text{with } \varphi \text{ from } \max_{\|\varphi\|=1} \min_{\mathbf{x} \in \mathbf{X}} \langle \mathbf{x}, \varphi \rangle.$$

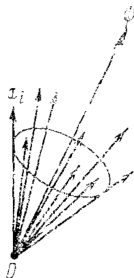


Рис. 24.

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Kernel trick** (Boser, Guyon, Vapnik; 1992):

- ▶ Let  $\Phi$  be a feature map:  $\mathbb{R}^d \mapsto \mathcal{H}$ .



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Kernel trick** (Boser, Guyon, Vapnik; 1992):

- ▶ Let  $\Phi$  be a feature map:  $\mathbb{R}^d \mapsto \mathcal{H}$ .
- ▶ Due to the **kernel trick**, the dot product in the image of  $\varphi$  can be computed by evaluation of a kernel  $K$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle .$$

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Kernel trick** (Boser, Guyon, Vapnik; 1992):

- ▶ Let  $\Phi$  be a feature map:  $\mathbb{R}^d \mapsto \mathcal{H}$ .
- ▶ Due to the **kernel trick**, the dot product in the image of  $\varphi$  can be computed by evaluation of a kernel  $K$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle .$$

- ▶ Example: **Gaussian kernel**

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}$$

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Kernel trick** (Boser, Guyon, Vapnik; 1992):

- ▶ Let  $\Phi$  be a feature map:  $\mathbb{R}^d \mapsto \mathcal{H}$ .
- ▶ Due to the **kernel trick**, the dot product in the image of  $\varphi$  can be computed by evaluation of a kernel  $K$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle .$$

- ▶ Example: **Gaussian kernel**

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{\gamma \|\mathbf{x}_i, \mathbf{x}_j\|}$$

**Soft margin** (Cortes, Vapnik; 1995):

- ▶ Allow for a portion of points from  $\mathbf{X}$  to be beyond the margin, label points far from the origin by “1”, those close by “-1”.

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Kernel trick** (Boser, Guyon, Vapnik; 1992):

- ▶ Let  $\Phi$  be a feature map:  $\mathbb{R}^d \mapsto \mathcal{H}$ .
- ▶ Due to the **kernel trick**, the dot product in the image of  $\varphi$  can be computed by evaluation of a kernel  $K$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle .$$

- ▶ Example: **Gaussian kernel**

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{\gamma \|\mathbf{x}_i, \mathbf{x}_j\|}$$

**Soft margin** (Cortes, Vapnik; 1995):

- ▶ Allow for a portion of points from  $\mathbf{X}$  to be beyond the margin, label points far from the origin by “1”, those close by “-1”.
- ▶ Controlled by a parameter  $\nu \in (0, 1)$   
(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999).

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Idea 1:** Separate points from the origin.

# One-class support vector machines (Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Idea 1:** Separate points from the origin.

This can be formulated as a quadratic programming problem

$$\begin{aligned} \min_{\psi \in \mathcal{H}, \xi \in \mathbb{R}^n, \rho \in \mathbb{R}} \quad & \frac{1}{2} \|\psi\|^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i - \rho \\ \text{subject to} \quad & \langle \psi, \Phi(\mathbf{x}_i) \rangle \geq \rho - \xi_i, \quad \xi_i \geq 0 \text{ for } i = 1, \dots, n, \end{aligned}$$

with  $\xi = (\xi_1, \dots, \xi_n)^\top$ .

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Idea 1:** Separate points from the origin.

This can be formulated as a quadratic programming problem

$$\begin{aligned} \min_{\psi \in \mathcal{H}, \xi \in \mathbb{R}^n, \rho \in \mathbb{R}} \quad & \frac{1}{2} \|\psi\|^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i - \rho \\ \text{subject to} \quad & \langle \psi, \Phi(\mathbf{x}_i) \rangle \geq \rho - \xi_i, \quad \xi_i \geq 0 \text{ for } i = 1, \dots, n, \end{aligned}$$

with  $\xi = (\xi_1, \dots, \xi_n)^\top$ .

The solution  $(\psi^*, \xi^*, \rho^*)$  yields the following **decision function**:

$$g_{\text{OC SVM}}(\mathbf{x}) = \text{sgn}(\langle \psi^*, \Phi(\mathbf{x}) \rangle - \rho^*).$$

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

**Idea 1:** Separate points from the origin.

This can be formulated as a quadratic programming problem

$$\begin{aligned} \min_{\psi \in \mathcal{H}, \xi \in \mathbb{R}^n, \rho \in \mathbb{R}} \quad & \frac{1}{2} \|\psi\|^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i - \rho \\ \text{subject to} \quad & \langle \psi, \Phi(\mathbf{x}_i) \rangle \geq \rho - \xi_i, \quad \xi_i \geq 0 \text{ for } i = 1, \dots, n, \end{aligned}$$

with  $\xi = (\xi_1, \dots, \xi_n)^\top$ .

The solution  $(\psi^*, \xi^*, \rho^*)$  yields the following **decision function**:

$$g_{\text{OCSVM}}(\mathbf{x}) = \text{sgn}(\langle \psi^*, \Phi(\mathbf{x}) \rangle - \rho^*).$$

One can reformulate the optimization problem to employ the **kernel trick**.



## One-class support vector machines (Schölkopf *et al.*, 1999)

In dual formulation, using the Lagrangian, one can restate the optimization problem as follows:

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{subject to} \quad 0 \leq \alpha_i \leq \frac{1}{\nu n} \text{ for } i = 1, \dots, n, \quad \sum_{i=1}^n \alpha_i = 1,$$

with  $\alpha = (\alpha_1, \dots, \alpha_n)^\top$ .

# One-class support vector machines (Schölkopf *et al.*, 1999)

In dual formulation, using the Lagrangian, one can restate the optimization problem as follows:

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{subject to} \quad 0 \leq \alpha_i \leq \frac{1}{\nu n} \text{ for } i = 1, \dots, n, \quad \sum_{i=1}^n \alpha_i = 1,$$

with  $\alpha = (\alpha_1, \dots, \alpha_n)^\top$ .

The [decision function](#) is then:

$$g_{\text{OC SVM}}(\mathbf{x}) = \text{sgn} \left( \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}) - \rho \right),$$

where  $\rho$  can be recovered from any  $\mathbf{x}_j$  such that  $0 < \alpha_j < \frac{1}{\nu n}$ :

$$\rho = \langle \psi, \Phi(\mathbf{x}_j) \rangle = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}_j).$$

# One-class support vector machines (Schölkopf *et al.*, 1999)

**Idea 2:** Put points into a small ball.

$$\begin{aligned} \min_{R \in \mathbb{R}, \xi \in \mathbb{R}^n, \mathbf{c} \in \mathcal{H},} \quad & R^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & \|\Phi(\mathbf{x}_i) - \mathbf{c}\| \leq R^2 + \xi_i, \quad \xi_i \geq 0 \text{ for } i = 1, \dots, n. \end{aligned}$$

# One-class support vector machines (Schölkopf *et al.*, 1999)

**Idea 2:** Put points into a small ball.

$$\begin{aligned} \min_{R \in \mathbb{R}, \xi \in \mathbb{R}^n, \mathbf{c} \in \mathcal{H},} \quad & R^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & \|\Phi(\mathbf{x}_i) - \mathbf{c}\| \leq R^2 + \xi_i, \quad \xi_i \geq 0 \text{ for } i = 1, \dots, n. \end{aligned}$$

This leads to the dual:

$$\begin{aligned} \min_{\alpha} \quad & \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}_i) \\ \text{subject to} \quad & 0 \leq \alpha_i \leq \frac{1}{\nu n}, \text{ for } i = 1, \dots, n, \quad \sum_{i=1}^n \alpha_i = 1. \end{aligned}$$

# One-class support vector machines (Schölkopf et al., 1999)

**Idea 2:** Put points into a small ball.

$$\begin{aligned} \min_{R \in \mathbb{R}, \xi \in \mathbb{R}^n, \mathbf{c} \in \mathcal{H},} \quad & R^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & \|\Phi(\mathbf{x}_i) - \mathbf{c}\| \leq R^2 + \xi_i, \quad \xi_i \geq 0 \text{ for } i = 1, \dots, n. \end{aligned}$$

This leads to the dual:

$$\begin{aligned} \min_{\alpha} \quad & \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}_i) \\ \text{subject to} \quad & 0 \leq \alpha_i \leq \frac{1}{\nu n}, \text{ for } i = 1, \dots, n, \quad \sum_{i=1}^n \alpha_i = 1. \end{aligned}$$

which leads to the **decision function**:

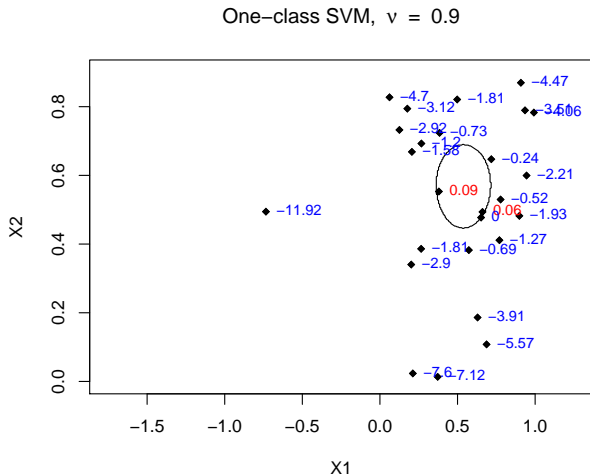
$$g_{\text{OCSVM}}(\mathbf{x}) = \left( R^2 - \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) + 2 \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}) - K(\mathbf{x}, \mathbf{x}) \right),$$

with  $R^2 = \sum_{i,j} \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) - 2 \sum_i \alpha_i K(\mathbf{x}_i, \mathbf{x}_k) + K(\mathbf{x}_k, \mathbf{x}_k)$  for any  $\mathbf{x}_k$  such that  $0 < \alpha_k < 1/(\nu n)$ .

# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

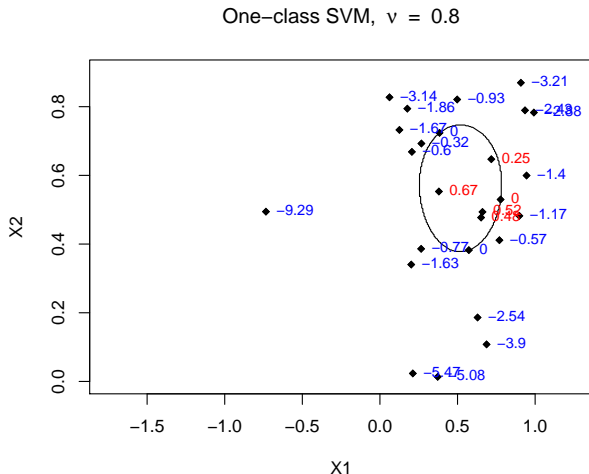
## Illustration: Case 1



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

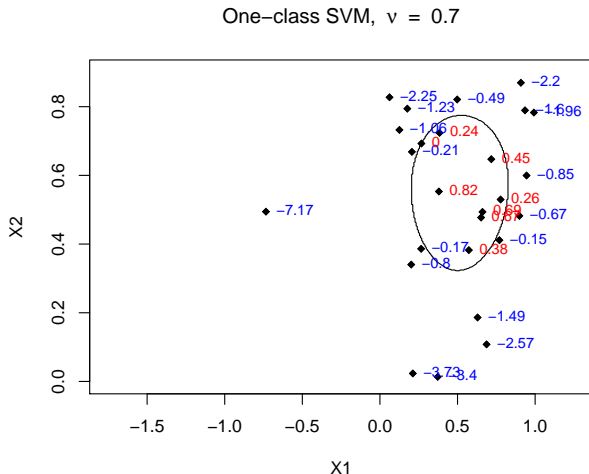
## Illustration: Case 1



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Illustration: Case 1



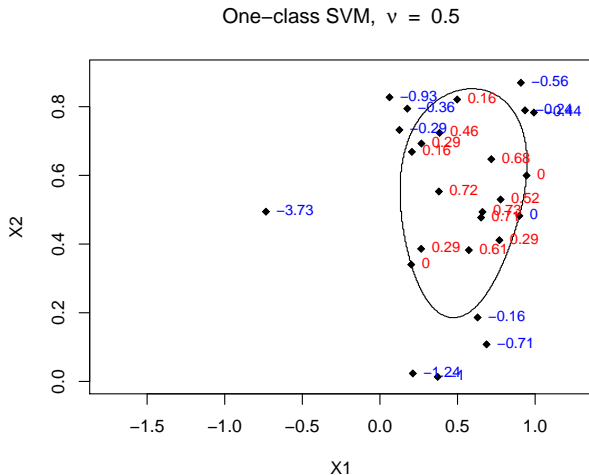




# One-class support vector machines

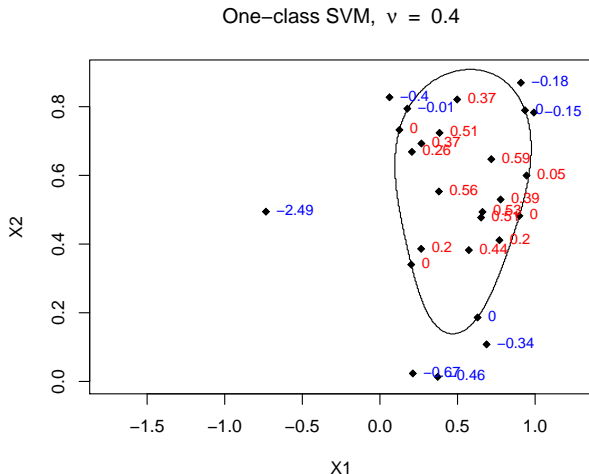
(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Illustration: Case 1



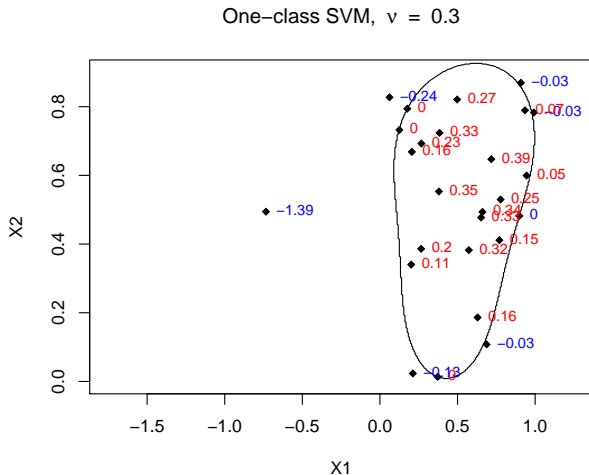
# One-class support vector machines (Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Illustration: Case 1



# One-class support vector machines (Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

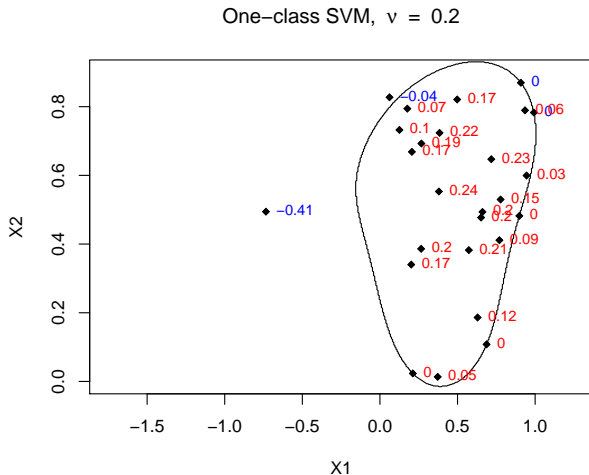
## Illustration: Case 1



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

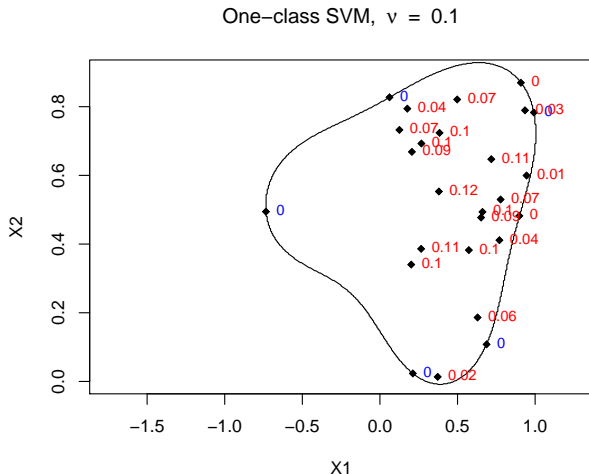
## Illustration: Case 1



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

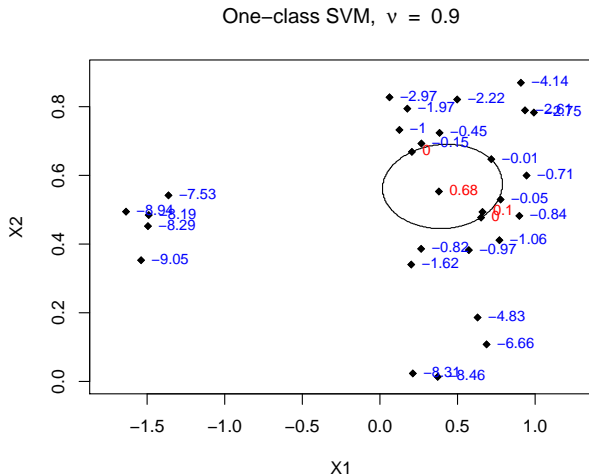
## Illustration: Case 1



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Illustration: Case 2



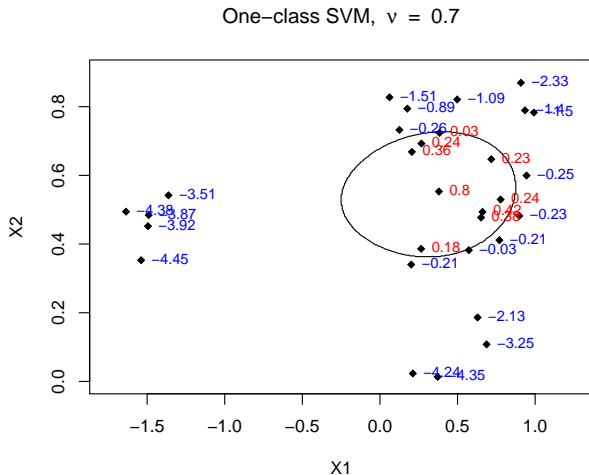




# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

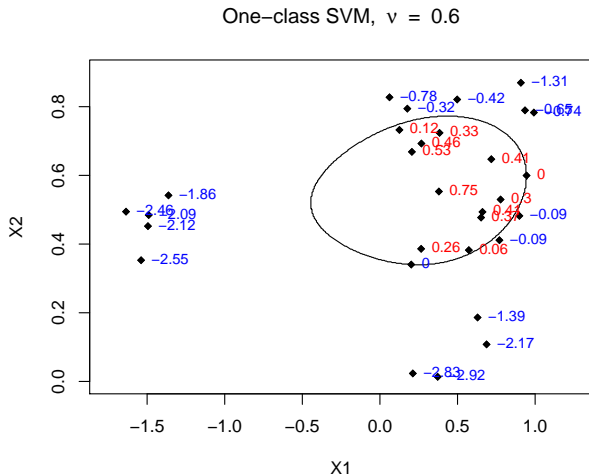
## Illustration: Case 2



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

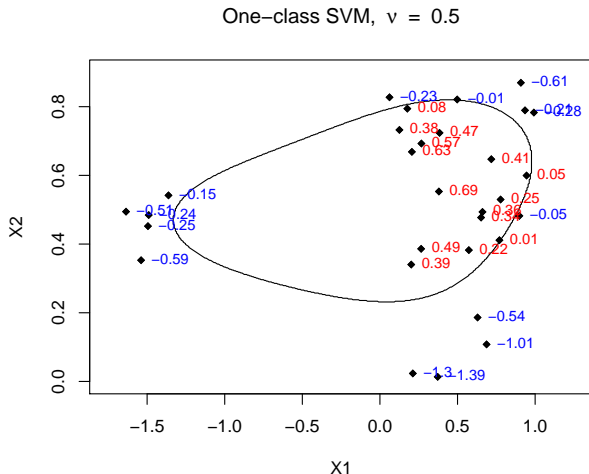
## Illustration: Case 2



# One-class support vector machines

(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

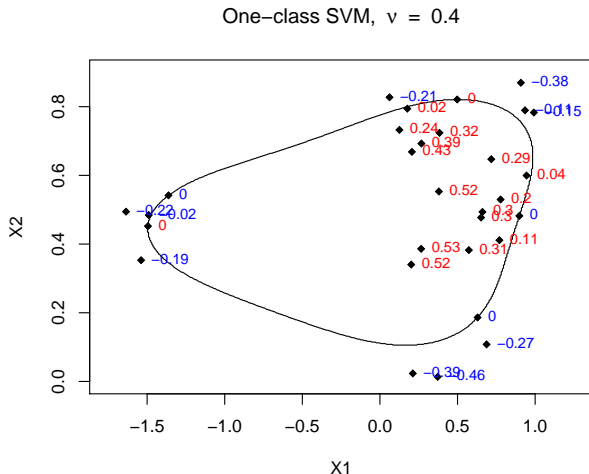
## Illustration: Case 2



# One-class support vector machines

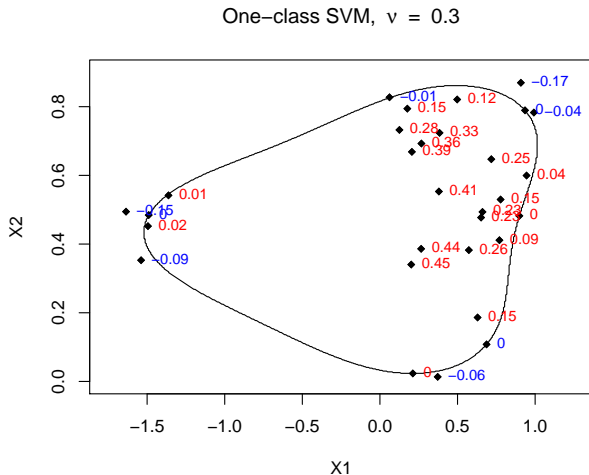
(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Illustration: Case 2



# One-class support vector machines (Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

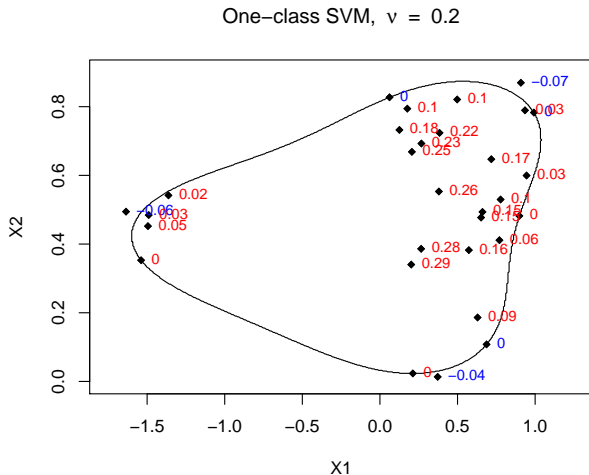
## Illustration: Case 2



# One-class support vector machines

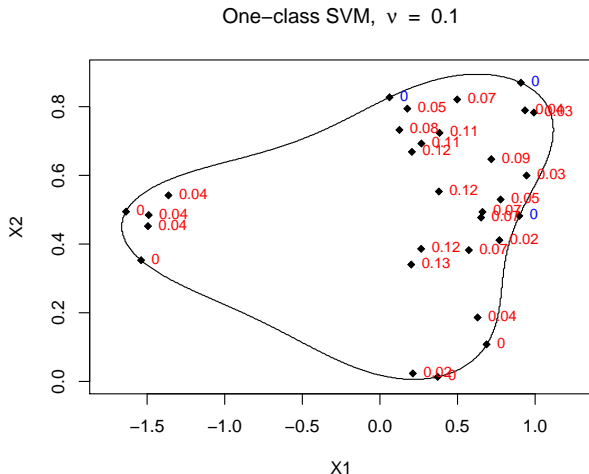
(Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Illustration: Case 2



# One-class support vector machines (Schölkopf, Platt, Shawe-Taylor, Smola, Williamson; 1999)

## Illustration: Case 2



# Contents

## Introduction

## Non-parametric approaches

One-class support vector machines

**Local outlier factor**

Isolation forest

## Systematic orderings: data depth

The notion of data depth

The Tukey depth function

Central regions

Further depth notions

## Practical session



# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**$k$ -distance** of a point  $\mathbf{x}$ :

For any integer  $k > 0$ , the  $k$ -distance of point  $\mathbf{x}$ , denoted as  $k\text{-dist}(\mathbf{x})$ , is defined as the distance  $d(\mathbf{x}, \mathbf{o})$  between  $\mathbf{x}$  and a point  $\mathbf{o} \in \mathbf{X}$  such that:

- ▶ for at least  $k$  points  $\mathbf{o}' \in \mathbf{X} \setminus \{\mathbf{x}\}$  it holds that  $d(\mathbf{x}, \mathbf{o}') \leq d(\mathbf{x}, \mathbf{o})$ , and
- ▶ for at most  $k - 1$  points  $\mathbf{o}' \in \mathbf{X} \setminus \{\mathbf{x}\}$  it holds that  $d(\mathbf{x}, \mathbf{o}') < d(\mathbf{x}, \mathbf{o})$ .

## Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**$k$ -distance** of a point  $\mathbf{x}$ :

For any integer  $k > 0$ , the  $k$ -distance of point  $\mathbf{x}$ , denoted as  $k\text{-dist}(\mathbf{x})$ , is defined as the distance  $d(\mathbf{x}, \mathbf{o})$  between  $\mathbf{x}$  and a point  $\mathbf{o} \in \mathbf{X}$  such that:

- ▶ for at least  $k$  points  $\mathbf{o}' \in \mathbf{X} \setminus \{\mathbf{x}\}$  it holds that  $d(\mathbf{x}, \mathbf{o}') \leq d(\mathbf{x}, \mathbf{o})$ , and
- ▶ for at most  $k - 1$  points  $\mathbf{o}' \in \mathbf{X} \setminus \{\mathbf{x}\}$  it holds that  $d(\mathbf{x}, \mathbf{o}') < d(\mathbf{x}, \mathbf{o})$ .

(=Distance from  $\mathbf{x}$  to its  $k$ th neighbor.)

# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**$k$ -distance** of a point  $\mathbf{x}$ :

For any integer  $k > 0$ , the  $k$ -distance of point  $\mathbf{x}$ , denoted as  $k\text{-dist}(\mathbf{x})$ , is defined as the distance  $d(\mathbf{x}, \mathbf{o})$  between  $\mathbf{x}$  and a point  $\mathbf{o} \in \mathbf{X}$  such that:

- ▶ for at least  $k$  points  $\mathbf{o}' \in \mathbf{X} \setminus \{\mathbf{x}\}$  it holds that  $d(\mathbf{x}, \mathbf{o}') \leq d(\mathbf{x}, \mathbf{o})$ , and
- ▶ for at most  $k - 1$  points  $\mathbf{o}' \in \mathbf{X} \setminus \{\mathbf{x}\}$  it holds that  $d(\mathbf{x}, \mathbf{o}') < d(\mathbf{x}, \mathbf{o})$ .

(=Distance from  $\mathbf{x}$  to its  $k$ th neighbor.)

**$k$ -neighborhood** of a point  $\mathbf{x}$ :

Given the  $k\text{-dist}(\mathbf{x})$ , the  **$k$ -neighborhood** of  $\mathbf{x}$ , denoted  $N_k(\mathbf{x})$ , contains every point whose distance from  $\mathbf{x}$  is not greater than the  $k\text{-dist}(\mathbf{x})$ , i.e.:

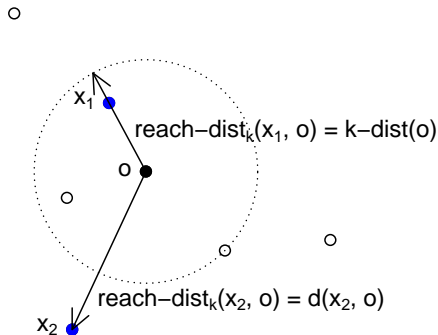
$$N_k(\mathbf{x}) = \{ \mathbf{q} \in \mathbf{X} \setminus \{\mathbf{x}\} \mid d(\mathbf{x}, \mathbf{q}) \leq k\text{-dist}(\mathbf{x}) \}.$$

## Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Reachability distance** of order  $k$  of point  $\mathbf{x}$  w.r.t. point  $\mathbf{o}$ :

For  $k \in \mathbb{N}$ , the **reachability distance** of order  $k$  of point  $\mathbf{x}$  with respect to point  $\mathbf{o}$  is defined as:

$$\text{reach-dist}_k(\mathbf{x}, \mathbf{o}) = \max\{k\text{-dist}(\mathbf{o}), d(\mathbf{x}, \mathbf{o})\}.$$



# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}.$$

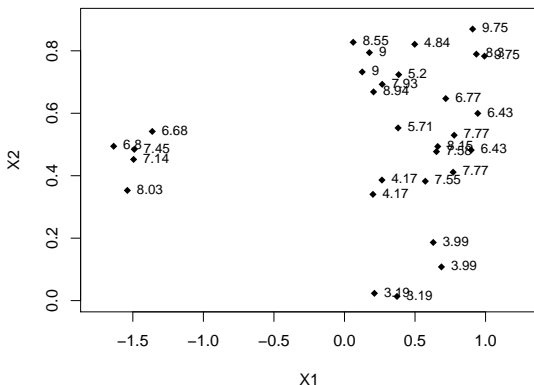
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

**Local reachability density,  $k = 2$**



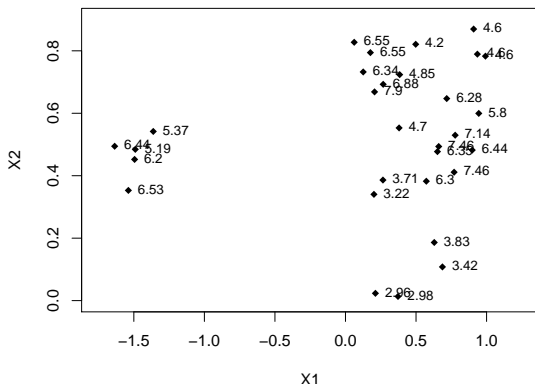
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

**Local reachability density,  $k = 3$**



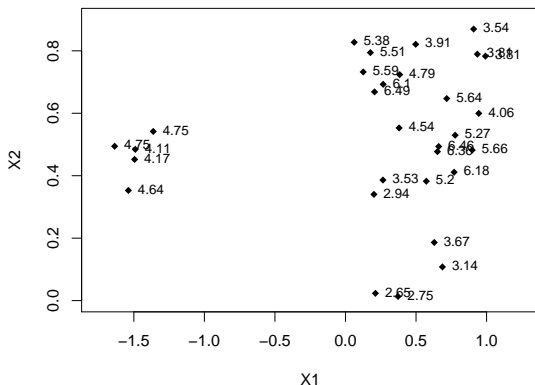
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}.$$

**Local reachability density,  $k = 4$**





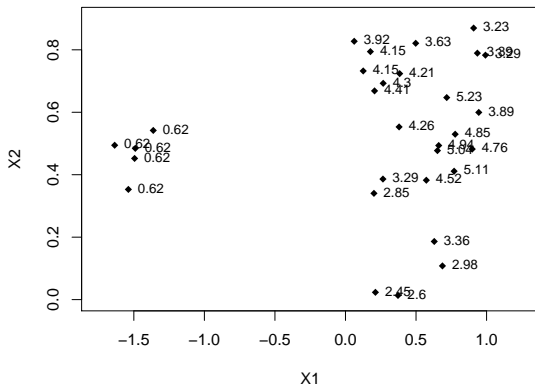
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}.$$

**Local reachability density,  $k = 5$**



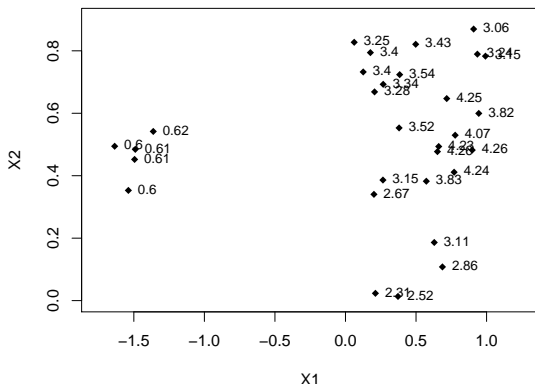
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

**Local reachability density, k = 6**



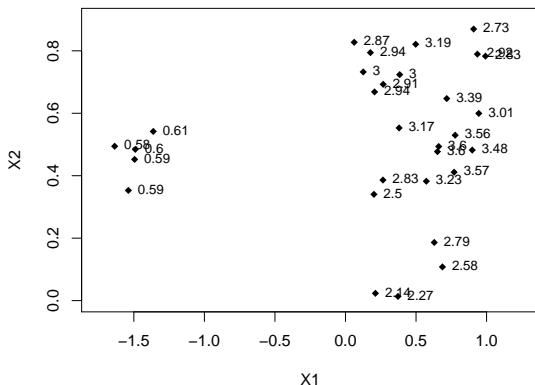
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

**Local reachability density,  $k = 7$**



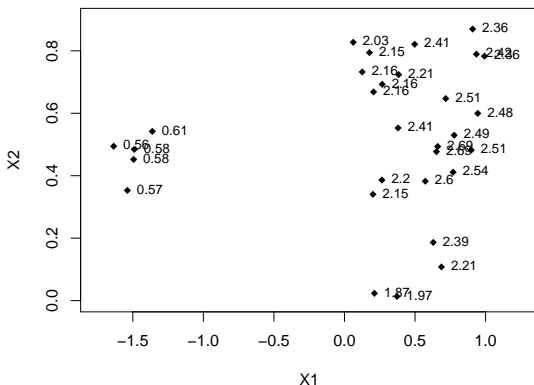
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}.$$

**Local reachability density,  $k = 10$**



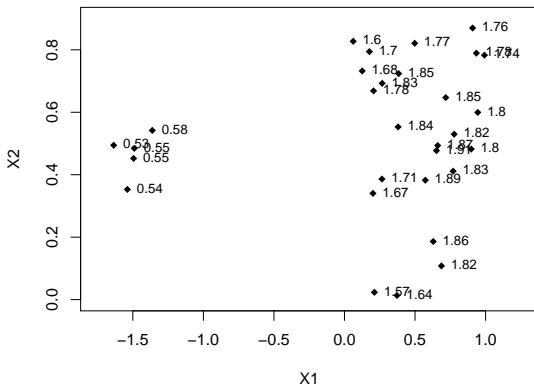
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

**Local reachability density, k = 15**



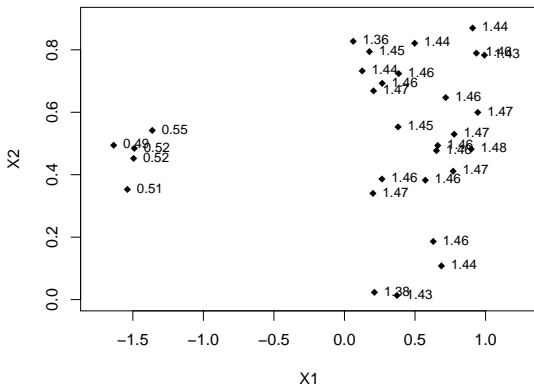
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

Local reachability density,  $k = 20$



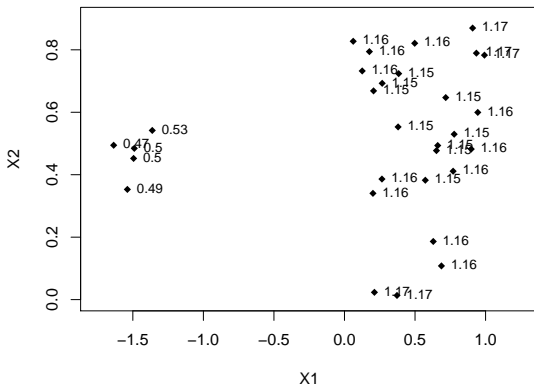
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

Local reachability density,  $k = 24$



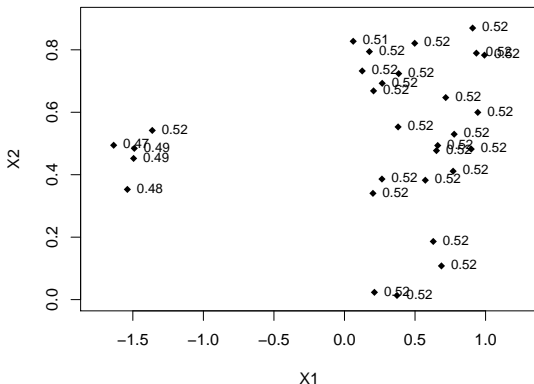
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}.$$

**Local reachability density,  $k = 25$**





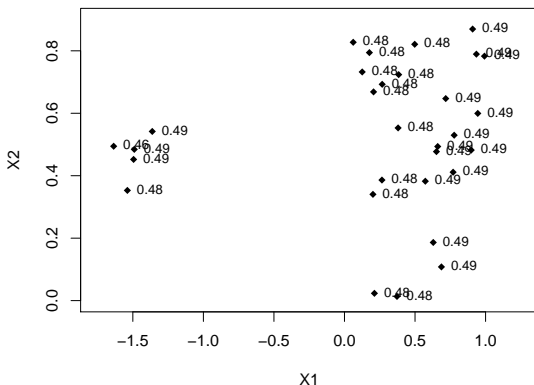
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}$$

Local reachability density,  $k = 26$



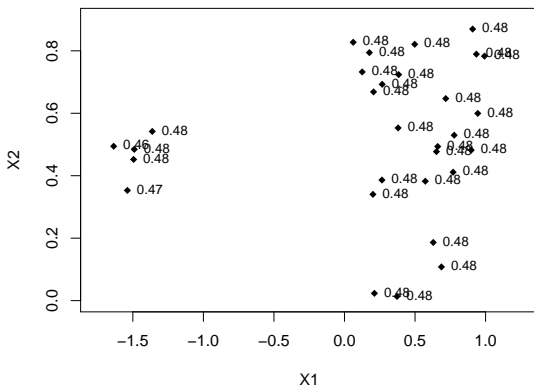
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local reachability density** of a point  $\mathbf{x}$ :

The **local reachability density** of  $\mathbf{x}$  is defined as:

$$lrd_k(\mathbf{x}) = \frac{|N_k(\mathbf{x})|}{\sum_{\mathbf{o} \in N_k(\mathbf{x})} reach-dist_k(\mathbf{x}, \mathbf{o})}.$$

Local reachability density,  $k = 27$



# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

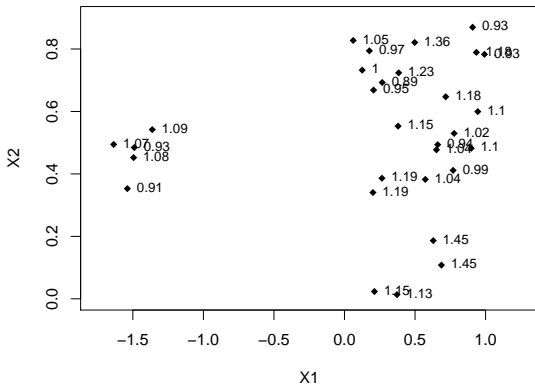
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor, k = 2**



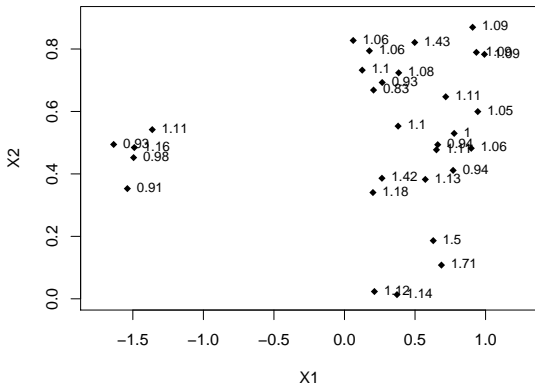
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor,  $k = 3$**



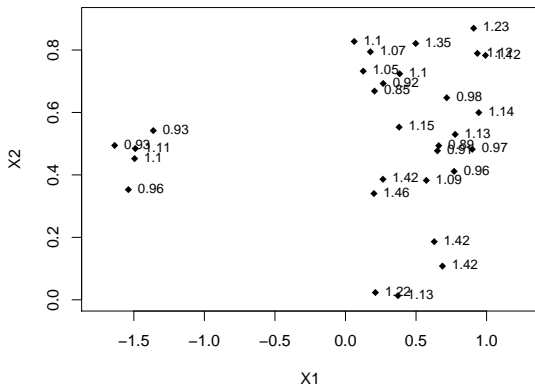
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor,  $k = 4$**



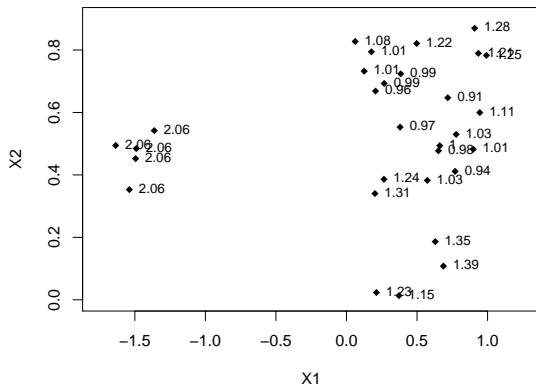
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor,  $k = 5$**



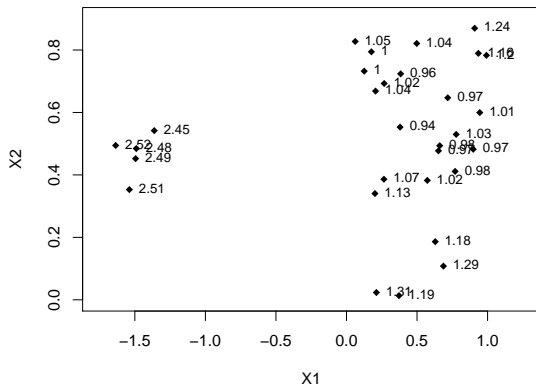
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor,  $k = 6$**





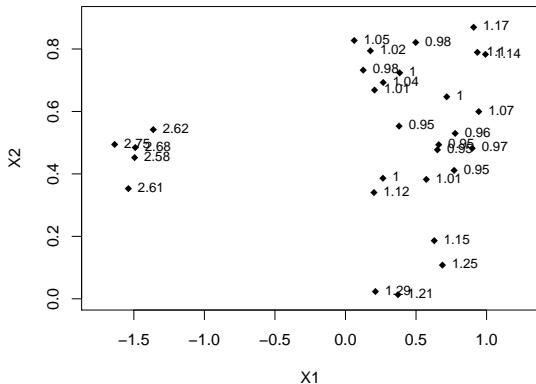
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor,  $k = 7$**



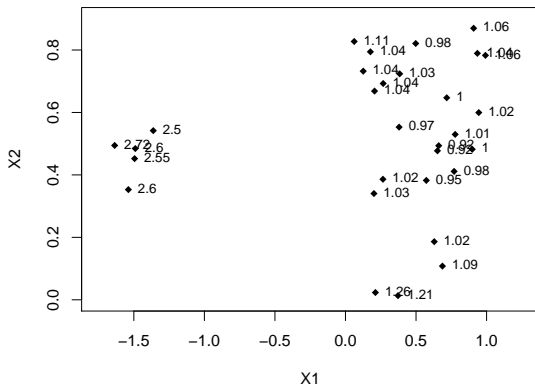
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

Local outlier factor,  $k = 10$



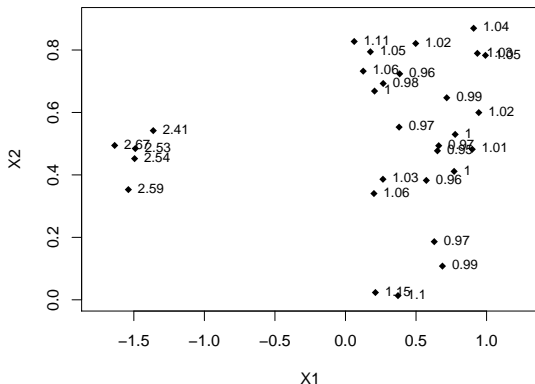
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor, k = 15**



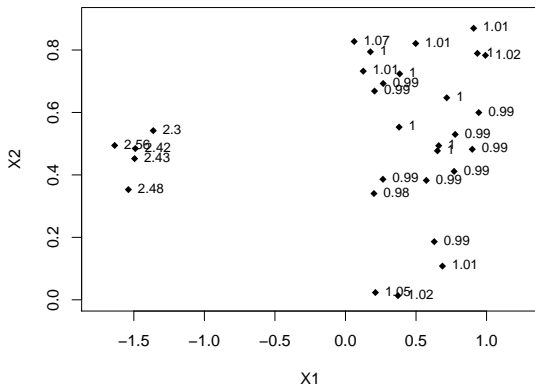
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

Local outlier factor,  $k = 20$



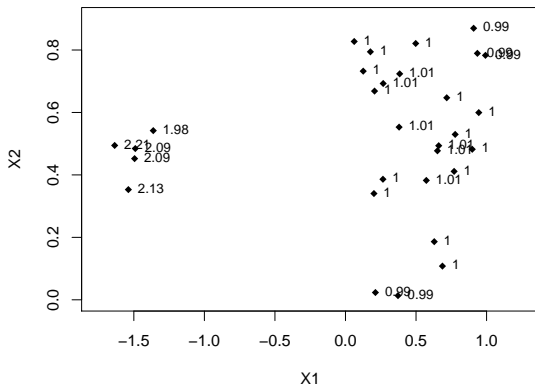
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor, k = 24**



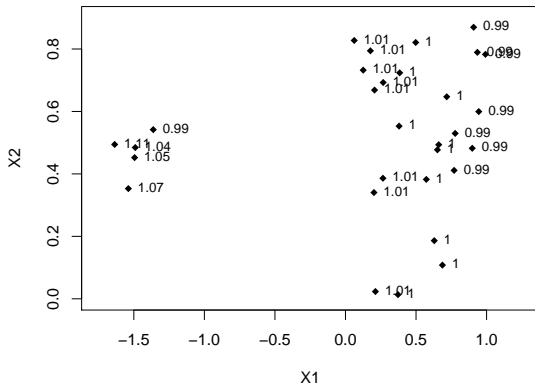
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

Local outlier factor,  $k = 25$



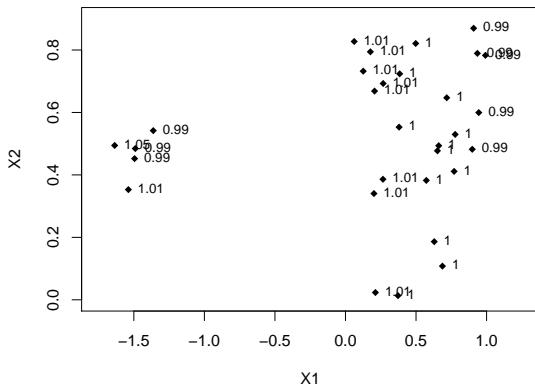
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor, k = 26**



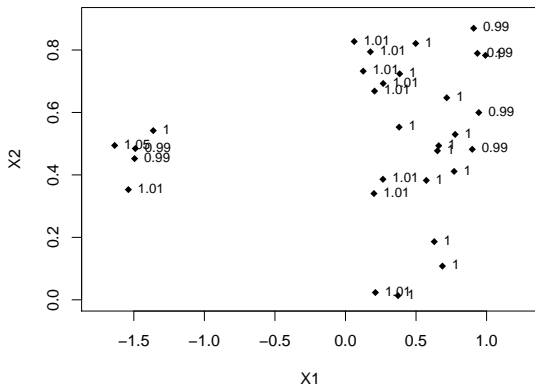
# Local outlier factor (Breunig, Kriegel, Ng, Sander; 2000)

**Local outlier factor** of a point  $\mathbf{x}$ :

The **local outlier factor** of  $\mathbf{x}$  is defined as:

$$LOF_k(\mathbf{x}) = \frac{\sum_{\mathbf{o} \in N_k(\mathbf{x})} \frac{lrd_k(\mathbf{o})}{lrd_k(\mathbf{x})}}{|N_k(\mathbf{x})|}.$$

**Local outlier factor, k = 27**





# Contents

## Introduction

## Non-parametric approaches

One-class support vector machines

Local outlier factor

Isolation forest

## Systematic orderings: data depth

The notion of data depth

The Tukey depth function

Central regions

Further depth notions

## Practical session

## Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ **Isolation forest** (Liu, Ting, Zhou; 2008) is an anomaly detection method inherited from the famous **random forest** algorithm (Breiman, 2001).
- ▶ Since no supervised feedback is given, isolation forest is based on **purely random** (uniform) variable-based partitioning.

## Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ **Isolation forest** (Liu, Ting, Zhou; 2008) is an anomaly detection method inherited from the famous **random forest** algorithm (Breiman, 2001).
- ▶ Since no supervised feedback is given, isolation forest is based on **purely random** (uniform) variable-based partitioning.
- ▶ **Main idea:** **Outlying observations are isolated faster.**

# Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ **Isolation forest** (Liu, Ting, Zhou; 2008) is an anomaly detection method inherited from the famous **random forest** algorithm (Breiman, 2001).
- ▶ Since no supervised feedback is given, isolation forest is based on **purely random** (uniform) variable-based partitioning.
- ▶ **Main idea:** **Outlying observations are isolated faster.**
- ▶ Tree-kind partitioning is done until “full isolation”: **outlying observations will have smaller depth** (on an average) in the **isolation tree**.
- ▶ A **monotone transform** is usually applied to the aggregated estimate.
- ▶ To reduce both **masking effect** and **computation cost**, small-size sub-sampling is used instead of bootstrap.

## Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ Each **isolation tree** is grown **recursively** using the described below **node-construction procedure**

## Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ Each **isolation tree** is grown **recursively** using the described below **node-construction procedure**

Non-terminal **node**  $(j, k)$ , **subspace**  $\mathcal{C}_{j,k}$ , **training subset**  $\mathcal{S}_{j,k}$ :

1. Choose a **split variable**  $l$  uniformly from  $\{1, \dots, d\}$ .

## Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ Each **isolation tree** is grown **recursively** using the described below **node-construction procedure**

Non-terminal **node**  $(j, k)$ , **subspace**  $\mathcal{C}_{j,k}$ , **training subset**  $\mathcal{S}_{j,k}$ :

1. Choose a **split variable**  $l$  uniformly from  $\{1, \dots, d\}$ .
2. Choose randomly and uniformly a **split value**  $\kappa$  in the interval

$$\left[ \min_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{e}_l \rangle, \max_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{e}_l \rangle \right].$$

## Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ Each **isolation tree** is grown **recursively** using the described below **node-construction procedure**

Non-terminal **node**  $(j, k)$ , **subspace**  $\mathcal{C}_{j,k}$ , **training subset**  $\mathcal{S}_{j,k}$ :

1. Choose a **split variable**  $l$  uniformly from  $\{1, \dots, d\}$ .
2. Choose randomly and uniformly a **split value**  $\kappa$  in the interval

$$\left[ \min_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{e}_l \rangle, \max_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{e}_l \rangle \right].$$

3. Form the children subsets

$$\begin{aligned} \mathcal{C}_{j+1,2k} &= \mathcal{C}_{j,k} \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{e}_l \rangle \leq \kappa\}, \\ \mathcal{C}_{j+1,2k+1} &= \mathcal{C}_{j,k} \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{e}_l \rangle > \kappa\}. \end{aligned}$$

as well as the children training datasets

$$\mathcal{S}_{j+1,2k} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k} \text{ and } \mathcal{S}_{j+1,2k+1} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k+1}.$$



## Isolation forest (Liu, Ting, Zhou; 2008)

- ▶ Each **isolation tree** is grown **recursively** using the described below **node-construction procedure**

Non-terminal **node**  $(j, k)$ , **subspace**  $\mathcal{C}_{j,k}$ , **training subset**  $\mathcal{S}_{j,k}$ :

1. Choose a **split variable**  $l$  uniformly from  $\{1, \dots, d\}$ .
2. Choose randomly and uniformly a **split value**  $\kappa$  in the interval

$$\left[ \min_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{e}_l \rangle, \max_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{e}_l \rangle \right].$$

3. Form the children subsets

$$\begin{aligned} \mathcal{C}_{j+1,2k} &= \mathcal{C}_{j,k} \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{e}_l \rangle \leq \kappa\}, \\ \mathcal{C}_{j+1,2k+1} &= \mathcal{C}_{j,k} \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{e}_l \rangle > \kappa\}. \end{aligned}$$

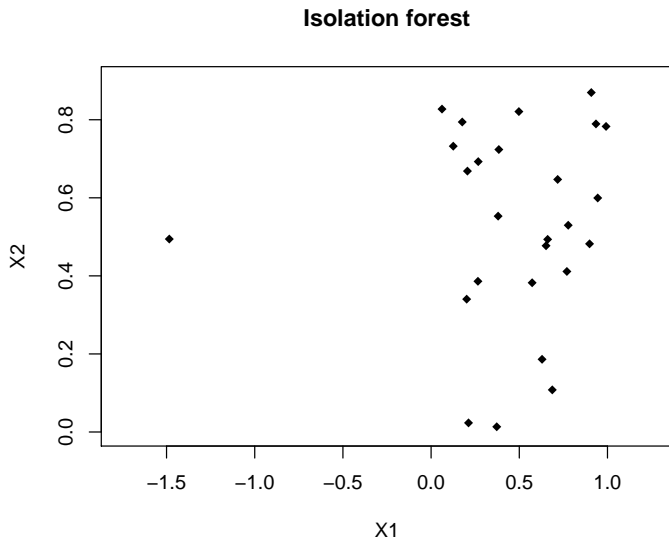
as well as the children training datasets

$$\mathcal{S}_{j+1,2k} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k} \text{ and } \mathcal{S}_{j+1,2k+1} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k+1}.$$

**Stop** when only one observation is in each node: **isolation**.

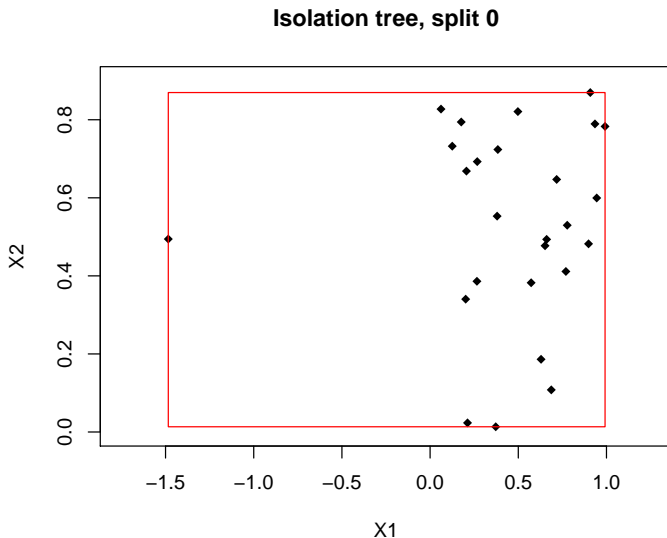
# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree



# Isolation forest (Liu, Ting, Zhou; 2008)

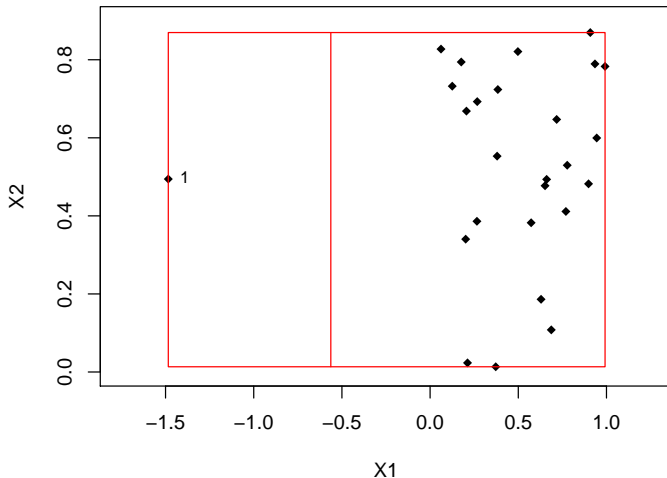
## Illustration: Isolation tree



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

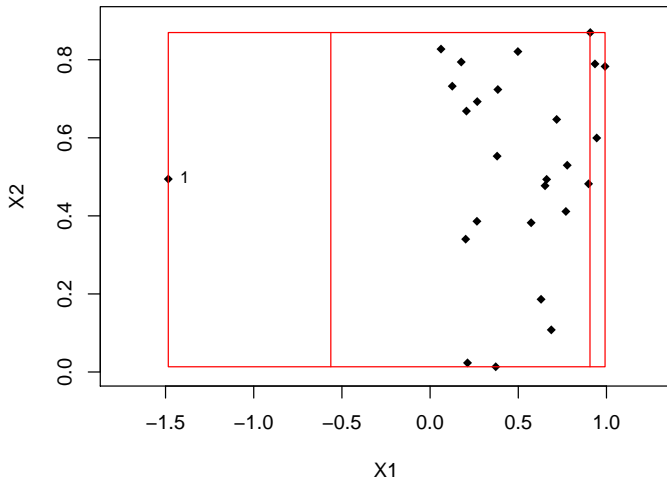
Isolation tree, split 1



# Isolation forest (Liu, Ting, Zhou; 2008)

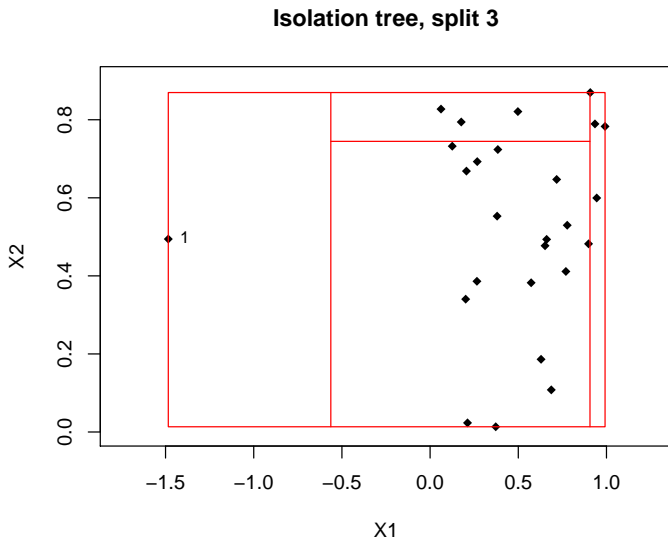
## Illustration: Isolation tree

Isolation tree, split 2



# Isolation forest (Liu, Ting, Zhou; 2008)

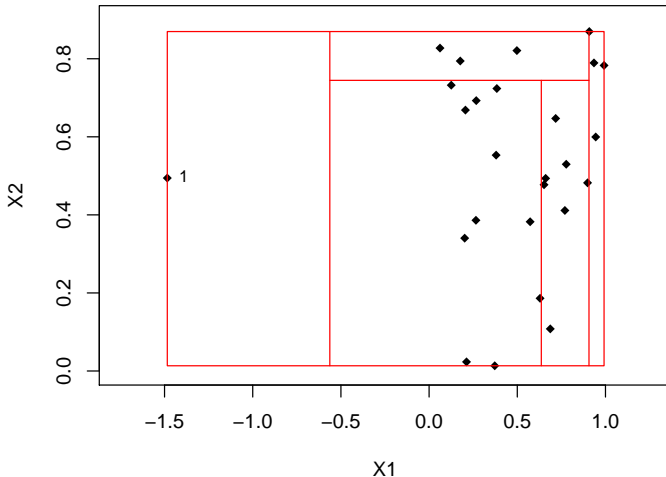
## Illustration: Isolation tree



# Isolation forest (Liu, Ting, Zhou; 2008)

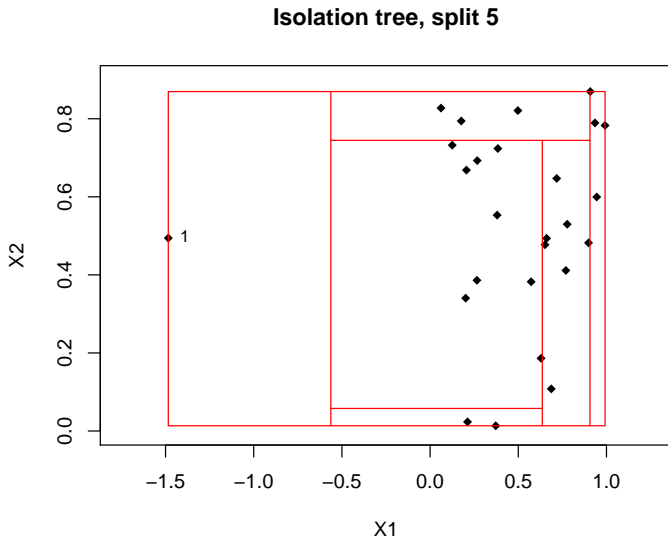
## Illustration: Isolation tree

Isolation tree, split 4



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

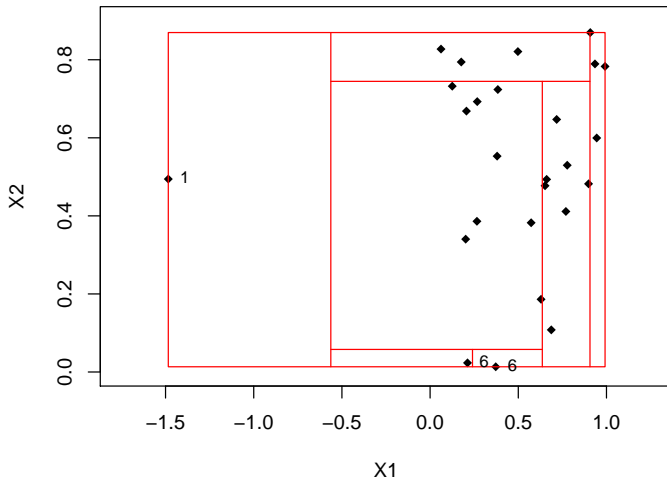




# Isolation forest (Liu, Ting, Zhou; 2008)

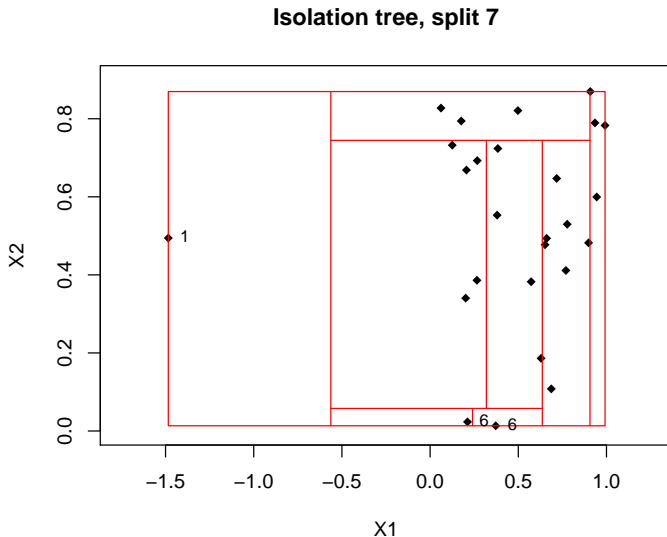
## Illustration: Isolation tree

Isolation tree, split 6



# Isolation forest (Liu, Ting, Zhou; 2008)

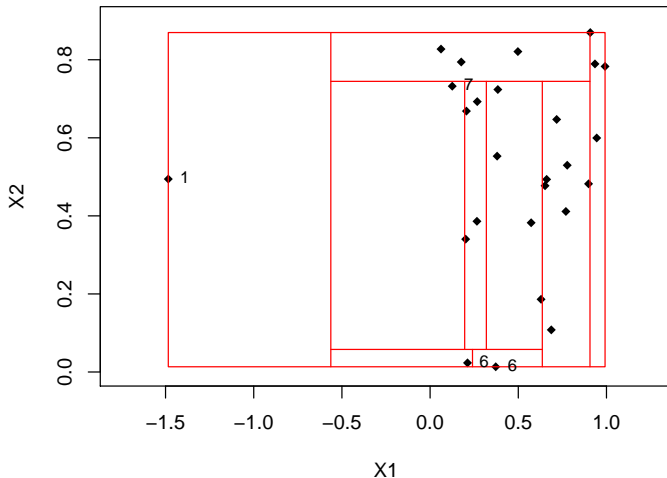
## Illustration: Isolation tree



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

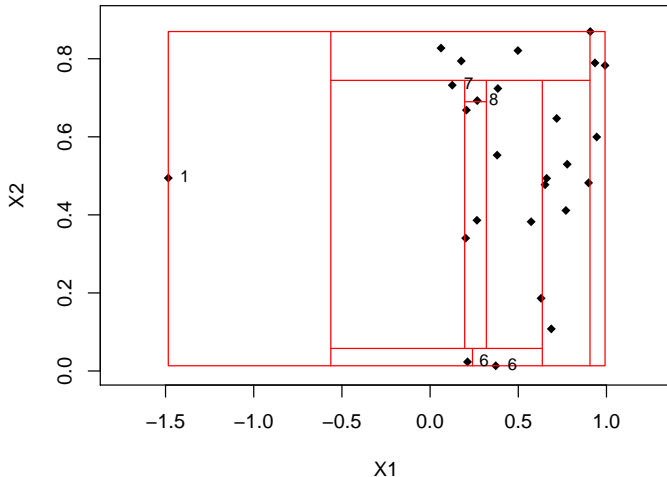
Isolation tree, split 8



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

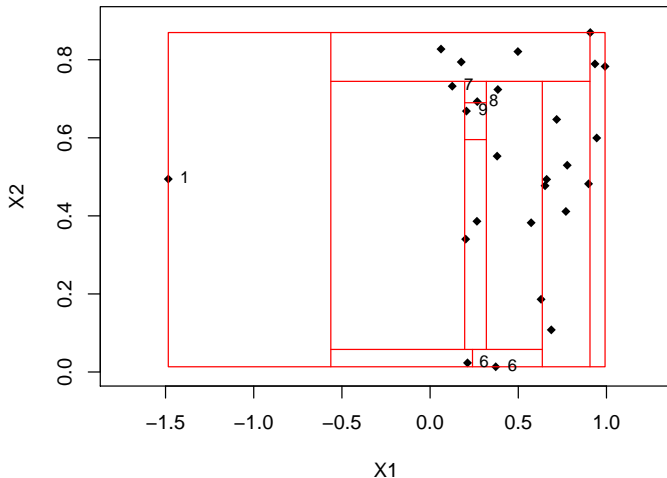
Isolation tree, split 9



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

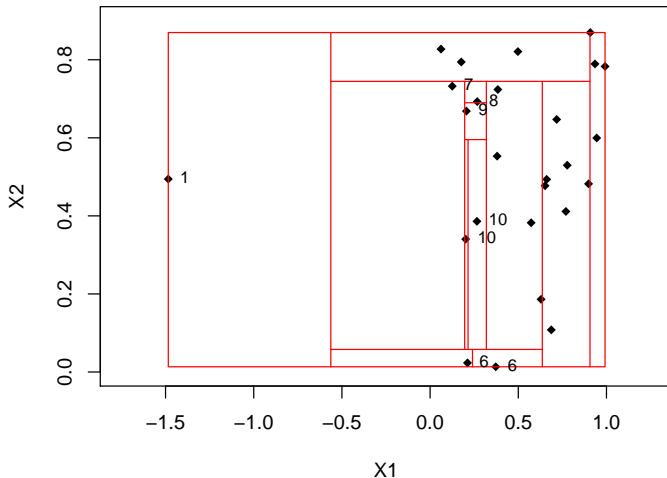
Isolation tree, split 10



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

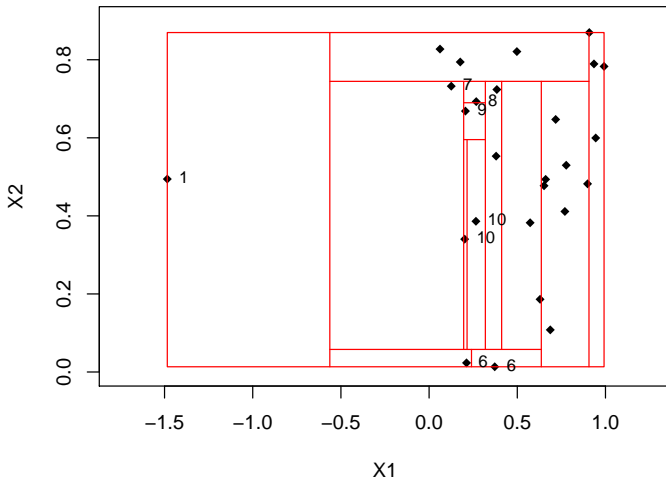
Isolation tree, split 11



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

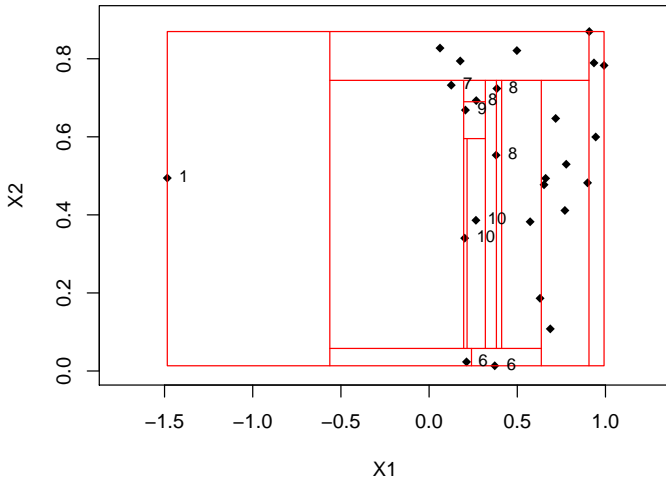
Isolation tree, split 12



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

Isolation tree, split 13





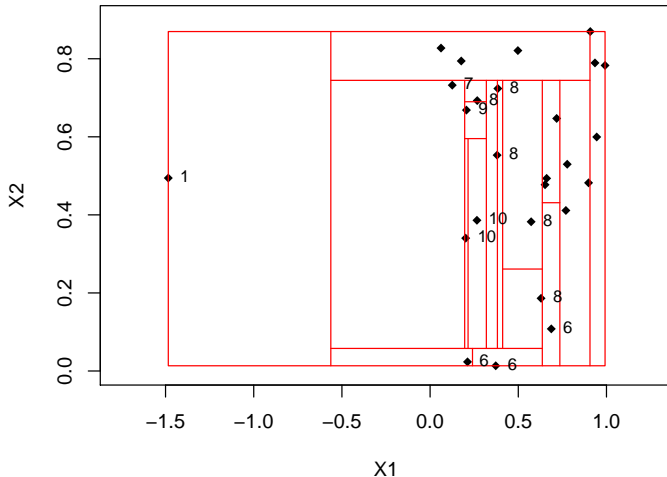




# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

Isolation tree, split 16



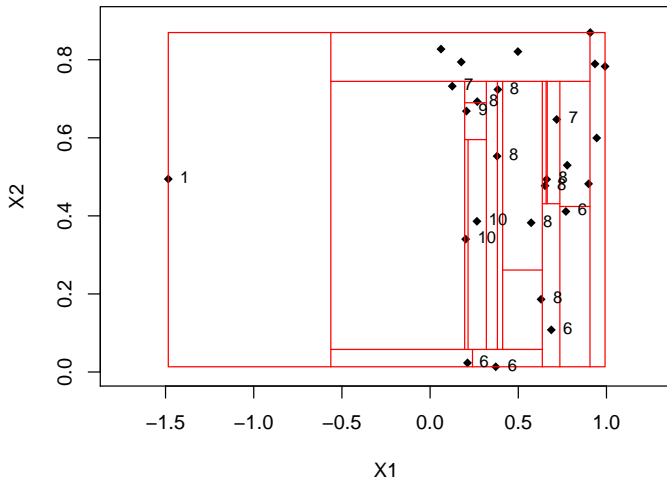




# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

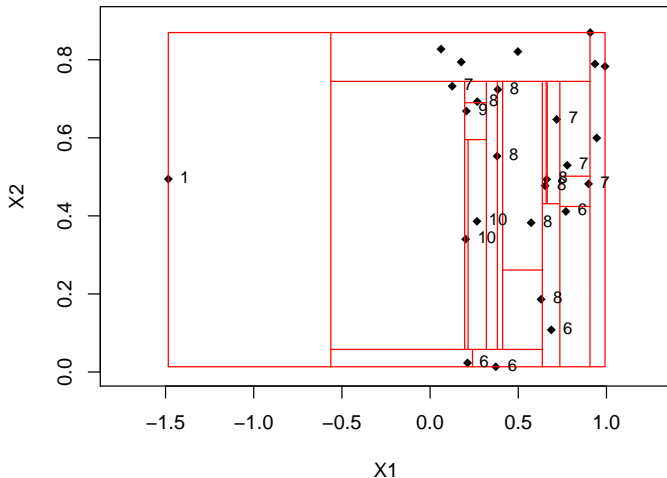
Isolation tree, split 19



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

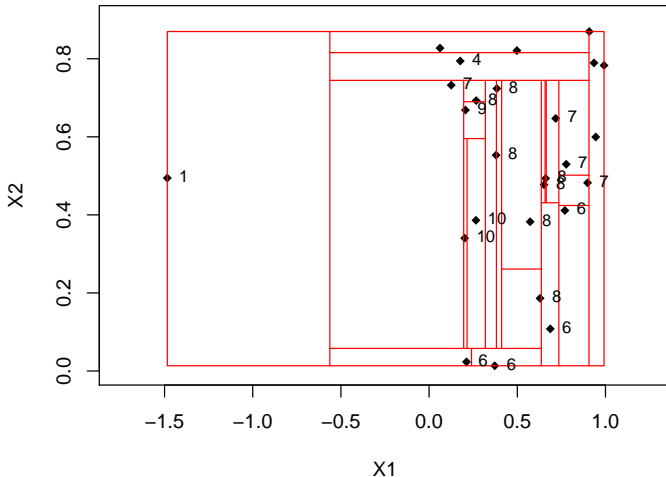
Isolation tree, split 20



# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

Isolation tree, split 21

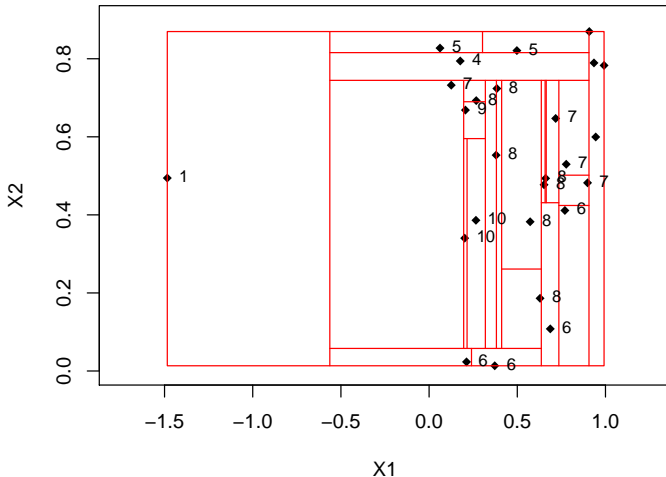




# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

Isolation tree, split 22

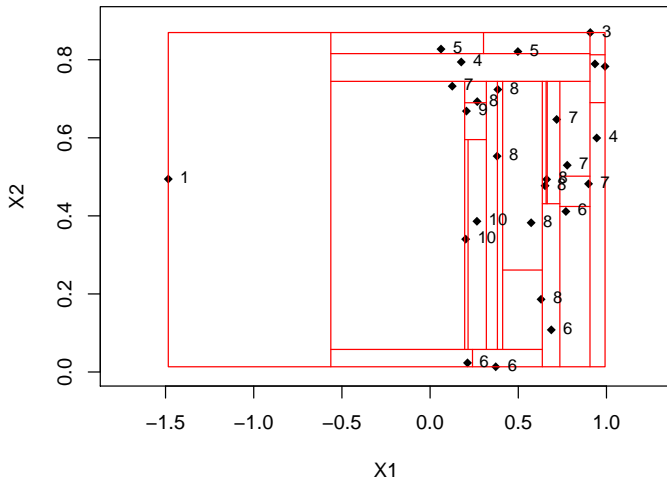




# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Isolation tree

Isolation tree, split 24





# Isolation forest (Liu, Ting, Zhou; 2008)

**Anomaly score calculation** for observation  $\mathbf{x}$ :

1. For each **isolation tree**  $i \in \{1, \dots, T\}$ , locate  $\mathbf{x}$  in a **terminal node** and calculate the **depth** of this node  $h_i(\mathbf{x})$ .
2. Attribute the **anomaly score**:

$$s(\mathbf{x}) = 2^{-\frac{\frac{1}{n} \sum_{i=1}^T h_i(\mathbf{x})}{c(n)}},$$

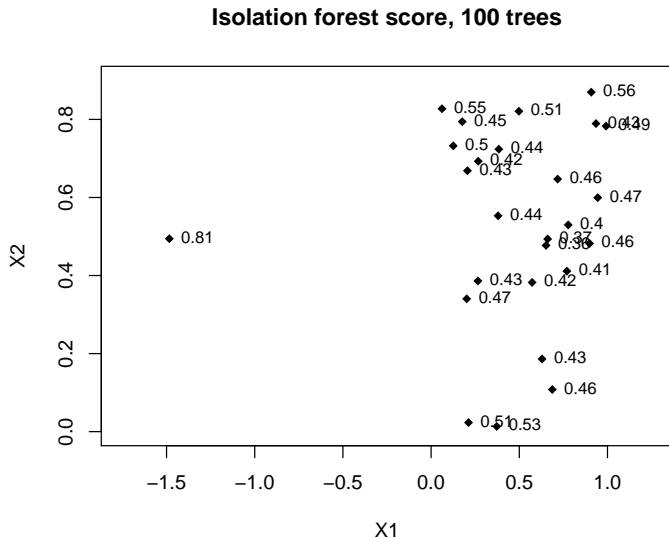
with  $c(n) = 2H(n-1) - \frac{2(n-1)}{n}$  where  $H(k)$  is the harmonic number and can be estimated by  $\ln(k) + 0.5772156649$ .

**Score behavior:**

- ▶ when  $\frac{1}{n} \sum_{i=1}^T h_i(\mathbf{x}) \rightarrow c(n)$ ,  $s(\mathbf{x}) \rightarrow 0.5$ ,
- ▶ when  $\frac{1}{n} \sum_{i=1}^T h_i(\mathbf{x}) \rightarrow 0$ ,  $s(\mathbf{x}) \rightarrow 1$ ,
- ▶ when  $\frac{1}{n} \sum_{i=1}^T h_i(\mathbf{x}) \rightarrow n-1$ ,  $s(\mathbf{x}) \rightarrow 0$ .

# Isolation forest (Liu, Ting, Zhou; 2008)

## Illustration: Anomaly score



# Contents

## Introduction

## Non-parametric approaches

- One-class support vector machines

- Local outlier factor

- Isolation forest

## Systematic orderings: data depth

- The notion of data depth

- The Tukey depth function

- Central regions

- Further depth notions

## Practical session

# Contents

## Introduction

## Non-parametric approaches

One-class support vector machines

Local outlier factor

Isolation forest

## Systematic orderings: data depth

The notion of data depth

The Tukey depth function

Central regions

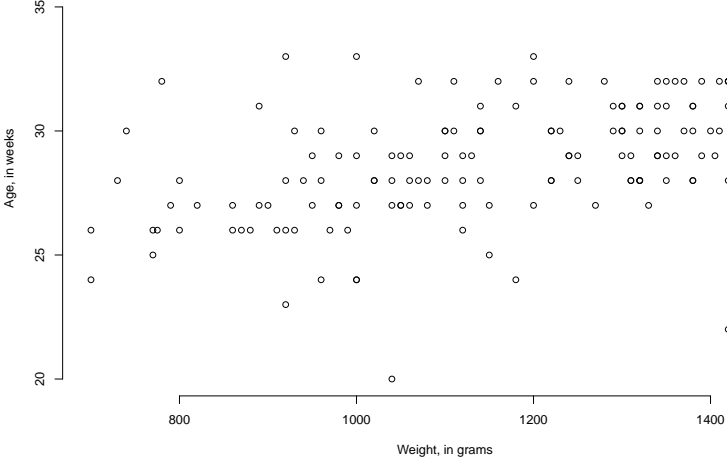
Further depth notions

## Practical session



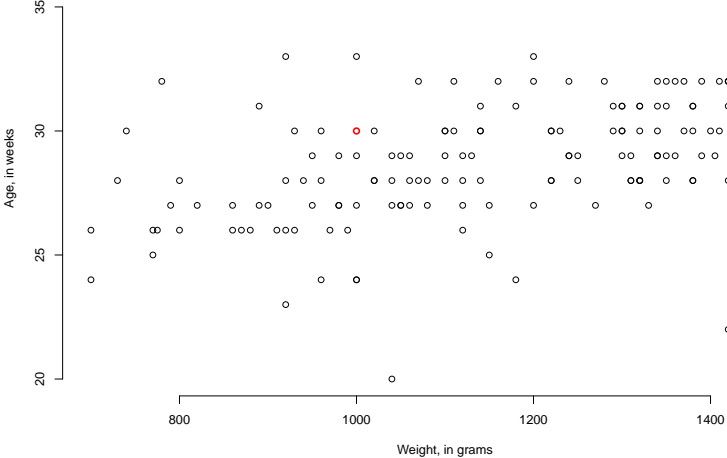
# Data depth

Babies with low birth weight



# Data depth

Babies with low birth weight



## Statistical data depth

A **data depth** measures how **close** a given point is located to the **center** of a distribution. For  $\mathbf{x} \in \mathbb{R}^p$  and a  $p$ -variate random vector  $X$  distributed as  $P \in \mathcal{P}$ , a data depth is a function

$$D : \mathbb{R}^p \times \mathcal{P} \rightarrow [0, 1], (\mathbf{x}, P) \mapsto D(\mathbf{x}|P)$$

## Statistical data depth

A **data depth** measures how **close** a given point is located to the **center** of a distribution. For  $\mathbf{x} \in \mathbb{R}^p$  and a  $p$ -variate random vector  $X$  distributed as  $P \in \mathcal{P}$ , a data depth is a function

$$D : \mathbb{R}^p \times \mathcal{P} \rightarrow [0, 1], (\mathbf{x}, P) \mapsto D(\mathbf{x}|P)$$

that is:

**D1 translation invariant:**  $D(\mathbf{x} + b|X + b) = D(\mathbf{x}|X)$  for any  $b \in \mathbb{R}^p$ ;

## Statistical data depth

A **data depth** measures how **close** a given point is located to the **center** of a distribution. For  $\mathbf{x} \in \mathbb{R}^p$  and a  $p$ -variate random vector  $X$  distributed as  $P \in \mathcal{P}$ , a data depth is a function

$$D : \mathbb{R}^p \times \mathcal{P} \rightarrow [0, 1], (\mathbf{x}, P) \mapsto D(\mathbf{x}|P)$$

that is:

- D1 translation invariant:**  $D(\mathbf{x} + b|X + b) = D(\mathbf{x}|X)$  for any  $b \in \mathbb{R}^p$ ;
- D2 linear invariant:**  $D(A\mathbf{x}|AX) = D(\mathbf{x}|X)$  for any  $p \times p$  non-singular matrix  $A$ ;

## Statistical data depth

A **data depth** measures how **close** a given point is located to the **center** of a distribution. For  $\mathbf{x} \in \mathbb{R}^p$  and a  $p$ -variate random vector  $X$  distributed as  $P \in \mathcal{P}$ , a data depth is a function

$$D : \mathbb{R}^p \times \mathcal{P} \rightarrow [0, 1], (\mathbf{x}, P) \mapsto D(\mathbf{x}|P)$$

that is:

- D1 translation invariant:**  $D(\mathbf{x} + b|X + b) = D(\mathbf{x}|X)$  for any  $b \in \mathbb{R}^p$ ;
- D2 linear invariant:**  $D(A\mathbf{x}|AX) = D(\mathbf{x}|X)$  for any  $p \times p$  non-singular matrix  $A$ ;
- D3 vanishing at infinity:**  $\lim_{\|\mathbf{x}\| \rightarrow \infty} D(\mathbf{x}|X) = 0$ ;

## Statistical data depth

A **data depth** measures how **close** a given point is located to the **center** of a distribution. For  $\mathbf{x} \in \mathbb{R}^p$  and a  $p$ -variate random vector  $X$  distributed as  $P \in \mathcal{P}$ , a data depth is a function

$$D : \mathbb{R}^p \times \mathcal{P} \rightarrow [0, 1], (\mathbf{x}, P) \mapsto D(\mathbf{x}|P)$$

that is:

- D1 translation invariant:**  $D(\mathbf{x} + b|X + b) = D(\mathbf{x}|X)$  for any  $b \in \mathbb{R}^p$ ;
- D2 linear invariant:**  $D(A\mathbf{x}|AX) = D(\mathbf{x}|X)$  for any  $p \times p$  non-singular matrix  $A$ ;
- D3 vanishing at infinity:**  $\lim_{\|\mathbf{x}\| \rightarrow \infty} D(\mathbf{x}|X) = 0$ ;
- D4 monotone on rays:** for any  $\mathbf{x}^* \in \operatorname{argmax}_{\mathbf{x} \in \mathbb{R}^p} D(\mathbf{x}|X)$ , any  $\mathbf{x} \in \mathbb{R}^p$ , and any  $0 \leq \alpha \leq 1$  it holds:  
 $D(\mathbf{x}|X) \leq D(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)|X)$ ;

## Statistical data depth

A **data depth** measures how **close** a given point is located to the **center** of a distribution. For  $\mathbf{x} \in \mathbb{R}^p$  and a  $p$ -variate random vector  $X$  distributed as  $P \in \mathcal{P}$ , a data depth is a function

$$D : \mathbb{R}^p \times \mathcal{P} \rightarrow [0, 1], (\mathbf{x}, P) \mapsto D(\mathbf{x}|P)$$

that is:

- D1 translation invariant:**  $D(\mathbf{x} + b|X + b) = D(\mathbf{x}|X)$  for any  $b \in \mathbb{R}^p$ ;
- D2 linear invariant:**  $D(A\mathbf{x}|AX) = D(\mathbf{x}|X)$  for any  $p \times p$  non-singular matrix  $A$ ;
- D3 vanishing at infinity:**  $\lim_{\|\mathbf{x}\| \rightarrow \infty} D(\mathbf{x}|X) = 0$ ;
- D4 monotone on rays:** for any  $\mathbf{x}^* \in \operatorname{argmax}_{\mathbf{x} \in \mathbb{R}^p} D(\mathbf{x}|X)$ , any  $\mathbf{x} \in \mathbb{R}^p$ , and any  $0 \leq \alpha \leq 1$  it holds:  
 $D(\mathbf{x}|X) \leq D(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)|X)$ ;
- D5 upper semicontinuous in  $\mathbf{x}$ :** the upper-level sets  
 $D_\alpha(X) = \{\mathbf{x} \in \mathbb{R}^p : D(\mathbf{x}|X) \geq \alpha\}$  are closed for all  $\alpha$ .



# Statistical data depth

Some remarks:

- ▶ **D4** implies star-shaped upper-level sets of  $D$ .

# Statistical data depth

Some remarks:

- ▶ **D4** implies star-shaped upper-level sets of  $D$ .

One can strengthen to:

- ▶ **D4con**:  $D(\cdot|X)$  is a **quasiconcave** function, *i.e.* the upper-level sets  $D_\alpha(X)$  are convex for all  $\alpha$ .

# Statistical data depth

Some remarks:

- ▶ **D4** implies star-shaped upper-level sets of  $D$ .

One can strengthen to:

- ▶ **D4con**:  $D(\cdot|X)$  is a **quasiconcave** function, *i.e.* the upper-level sets  $D_\alpha(X)$  are convex for all  $\alpha$ .
- ▶ **D1** and **D2** define **affine invariante depth**.

# Statistical data depth

Some remarks:

- ▶ **D4** implies star-shaped upper-level sets of  $D$ .

One can strengthen to:

- ▶ **D4con**:  $D(\cdot|X)$  is a **quasiconcave** function, *i.e.* the upper-level sets  $D_\alpha(X)$  are convex for all  $\alpha$ .
- ▶ **D1** and **D2** define **affine invariante depth**.

One can also weaken to:

- ▶ **D2iso**:  $D(Ax|AX) = D(x|X)$  for every isometric linear  $A$  to define **orthogonal invariant depth**;

# Statistical data depth

Some remarks:

- ▶ **D4** implies star-shaped upper-level sets of  $D$ .

One can strengthen to:

- ▶ **D4con**:  $D(\cdot|X)$  is a **quasiconcave** function, *i.e.* the upper-level sets  $D_\alpha(X)$  are convex for all  $\alpha$ .
- ▶ **D1** and **D2** define **affine invariante depth**.

One can also weaken to:

- ▶ **D2iso**:  $D(A\mathbf{x}|AX) = D(\mathbf{x}|X)$  for every isometric linear  $A$  to define **orthogonal invariant depth**;
- ▶ **D2sca**:  $D(\lambda\mathbf{x}|\lambda X) = D(\mathbf{x}|X)$  for any  $\lambda > 0$  to define **scale invariant depth**.

# Statistical data depth

Some remarks:

- ▶ **D4** implies star-shaped upper-level sets of  $D$ .

One can strengthen to:

- ▶ **D4con**:  $D(\cdot|X)$  is a **quasiconcave** function, *i.e.* the upper-level sets  $D_\alpha(X)$  are convex for all  $\alpha$ .
- ▶ **D1** and **D2** define **affine invariante depth**.

One can also weaken to:

- ▶ **D2iso**:  $D(Ax|AX) = D(x|X)$  for every isometric linear  $A$  to define **orthogonal invariant depth**;
- ▶ **D2sca**:  $D(\lambda x|\lambda X) = D(x|X)$  for any  $\lambda > 0$  to define **scale invariant depth**.

Depth notions: **Mahalanobis** ('36), **projection** (Stahel, '81; Donoho, '82), **simplicial volume** (Oja, '83), **simplicial** (Liu, '90), **zonoid** (Koshevoy, Mosler, '97), **spatial** (Vardi, Zhang, '00; Serfling, '02), **lens** (Liu, Modarres, '11), ... depth.

## Applications of data depth:

- ▶ **Multivariate data analysis** (Liu, Parelius, Singh '99);
- ▶ **Statistical quality control** (Liu, Singh '93);
- ▶ **Cluster analysis and classification** (Mosler, Hoberg '06; Li, Cuesta-Albertos, Liu '12; M., Mosler, Lange '15);
- ▶ **Tests for multivariate location, scale, symmetry** (Liu '92; Dyckerhoff '02; Dyckerhoff, Ley, Paindaveine '15);
- ▶ **Outlier detection** (Hubert, Rousseeuw, Segaert '15);
- ▶ **Multivariate risk measurement** (Casco, Mochalov '07);
- ▶ **Robust linear programming** (Bazovkin, Mosler '15);
- ▶ **Missing data imputation** (M., Josse, Husson '20);
- ▶ etc.

R-package **ddalpha** (Pokotylo, M., Dyckerhoff, Nagy):  
calculates a number of depths; performs depth-based classification  
of multivariate and functional data; contains 50 multivariate and 5  
functional data sets.

Python library **data-depth**: to be released soon.

# Contents

## Introduction

## Non-parametric approaches

One-class support vector machines

Local outlier factor

Isolation forest

## Systematic orderings: data depth

The notion of data depth

The Tukey depth function

Central regions

Further depth notions

## Practical session



## Tukey (=halfspace, location) depth

**Tukey (1975) — “Mathematics and the picturing of data”**

Tukey depth of  $\mathbf{x} \in \mathbb{R}^p$  w.r.t. a  $d$ -variate random vector  $X$  distributed as  $P$  is defined as the smallest probability mass of a closed halfspace containing  $\mathbf{x}$ :

$$D^T(\mathbf{x}|X) = \inf\{P(H) : H \text{ is a closed halfspace, } \mathbf{x} \in H\},$$

## Tukey (=halfspace, location) depth

**Tukey (1975) — “Mathematics and the picturing of data”**

Tukey depth of  $\mathbf{x} \in \mathbb{R}^p$  w.r.t. a  $d$ -variate random vector  $X$  distributed as  $P$  is defined as the smallest probability mass of a closed halfspace containing  $\mathbf{x}$ :

$$D^T(\mathbf{x}|X) = \inf\{P(H) : H \text{ is a closed halfspace, } \mathbf{x} \in H\},$$

and w.r.t. a data set  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^p$ :

$$D^{T(n)}(\mathbf{x}|\mathbf{X}) = \frac{1}{n} \min_{\mathbf{u} \in \mathbb{S}^{p-1}} \#\{i : \mathbf{u}'\mathbf{x}_i \geq \mathbf{u}'\mathbf{x}\}.$$

# Tukey (=halfspace, location) depth

**Tukey (1975) — “Mathematics and the picturing of data”**

Tukey depth of  $\mathbf{x} \in \mathbb{R}^p$  w.r.t. a  $d$ -variate random vector  $X$  distributed as  $P$  is defined as the smallest probability mass of a closed halfspace containing  $\mathbf{x}$ :

$$D^T(\mathbf{x}|X) = \inf\{P(H) : H \text{ is a closed halfspace, } \mathbf{x} \in H\},$$

and w.r.t. a data set  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^p$ :

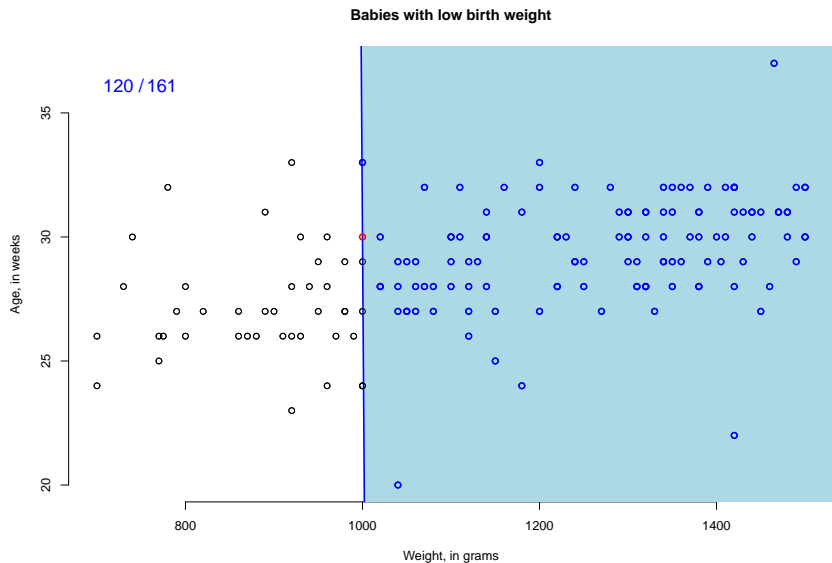
$$D^{T(n)}(\mathbf{x}|\mathbf{X}) = \frac{1}{n} \min_{\mathbf{u} \in \mathbb{S}^{p-1}} \#\{i : \mathbf{u}'\mathbf{x}_i \geq \mathbf{u}'\mathbf{x}\}.$$

## Tukey depth

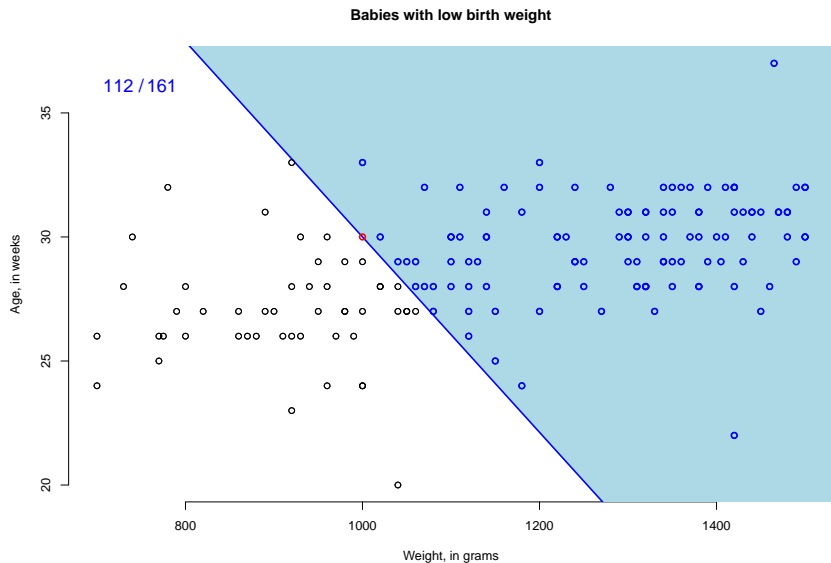
- ▶ satisfies all the above postulates,
- ▶ is purely non-parametric and robust,
- ▶ has direct connection to quantiles and many applications.



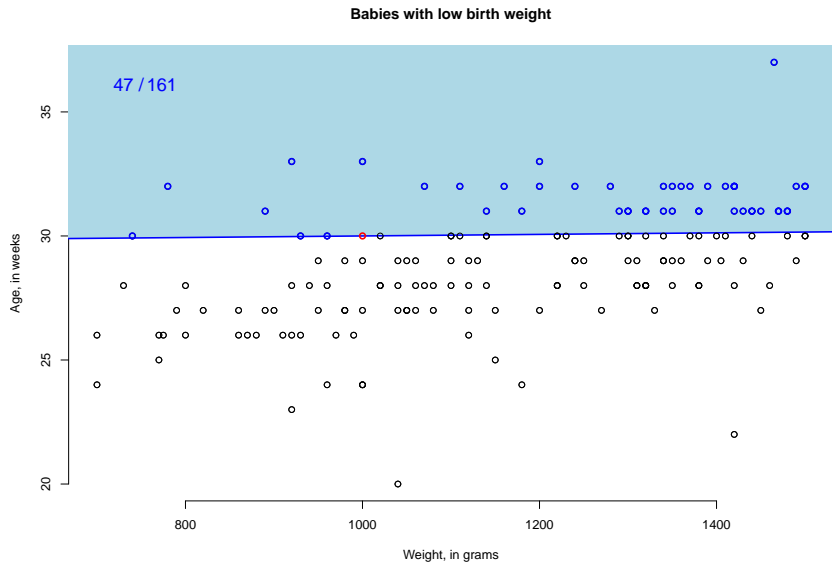
# Tukey (=halfspace, location) data depth



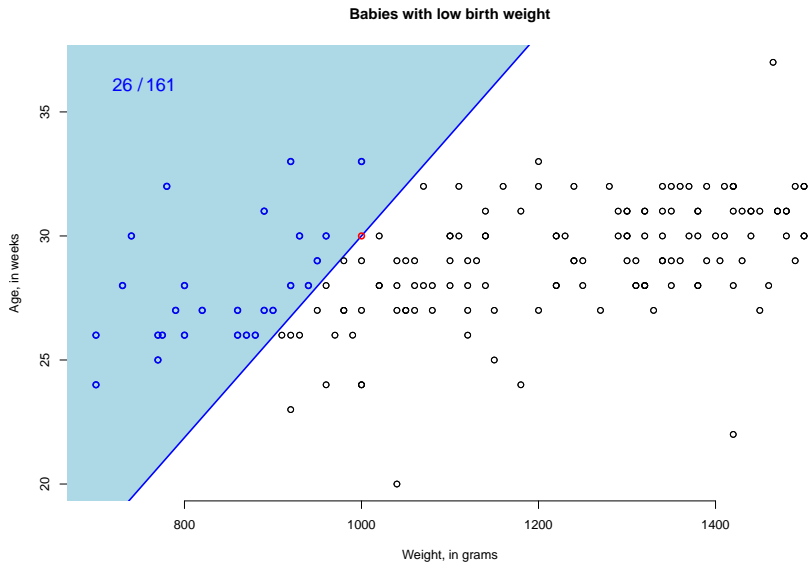
# Tukey (=halfspace, location) data depth



# Tukey (=halfspace, location) data depth



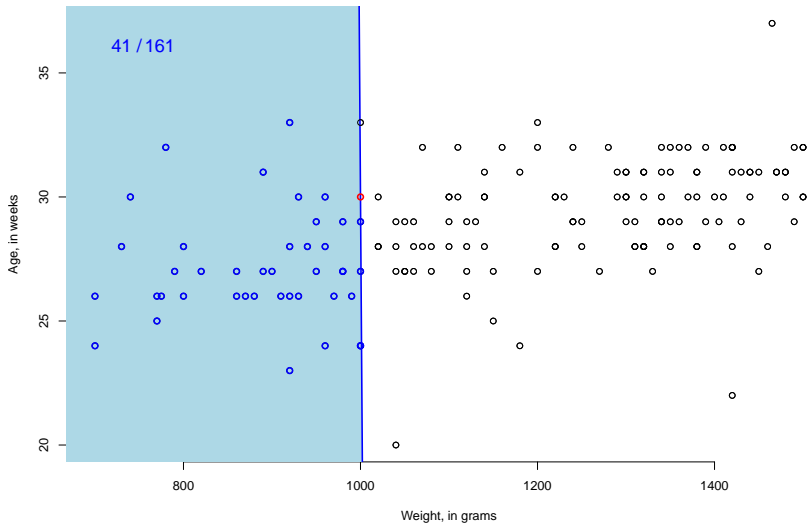
# Tukey (=halfspace, location) data depth





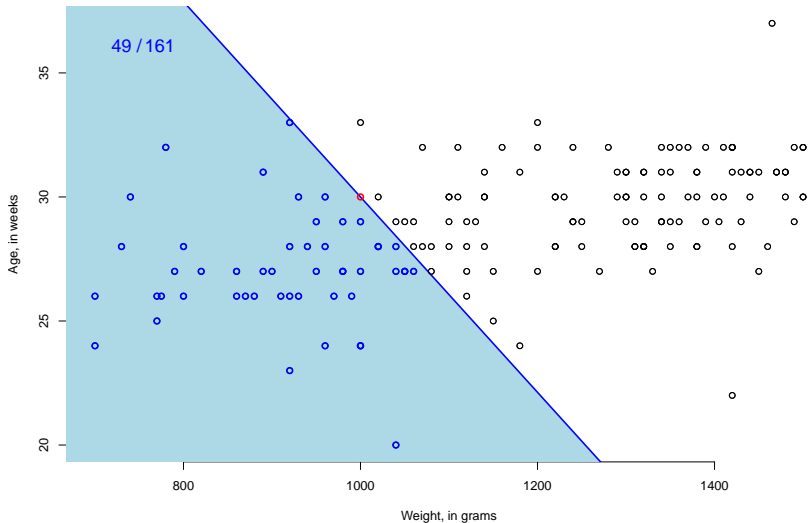
# Tukey (=halfspace, location) data depth

Babies with low birth weight



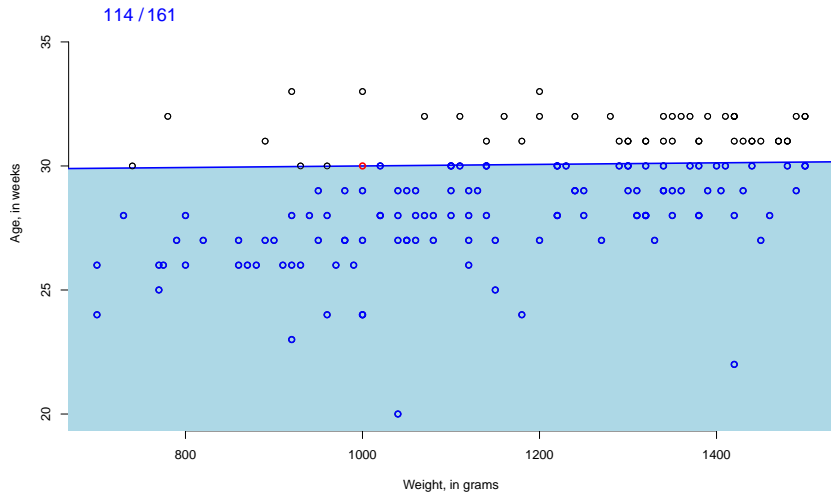
# Tukey (=halfspace, location) data depth

Babies with low birth weight

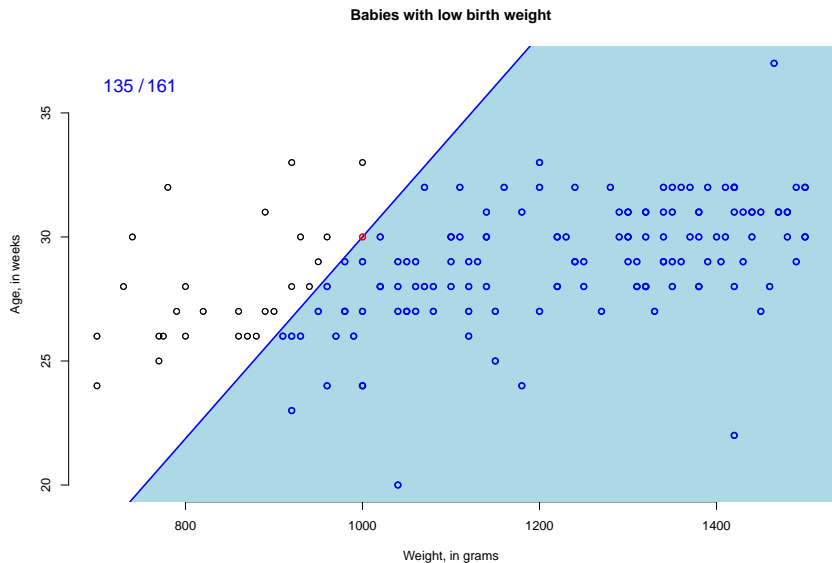


# Tukey (=halfspace, location) data depth

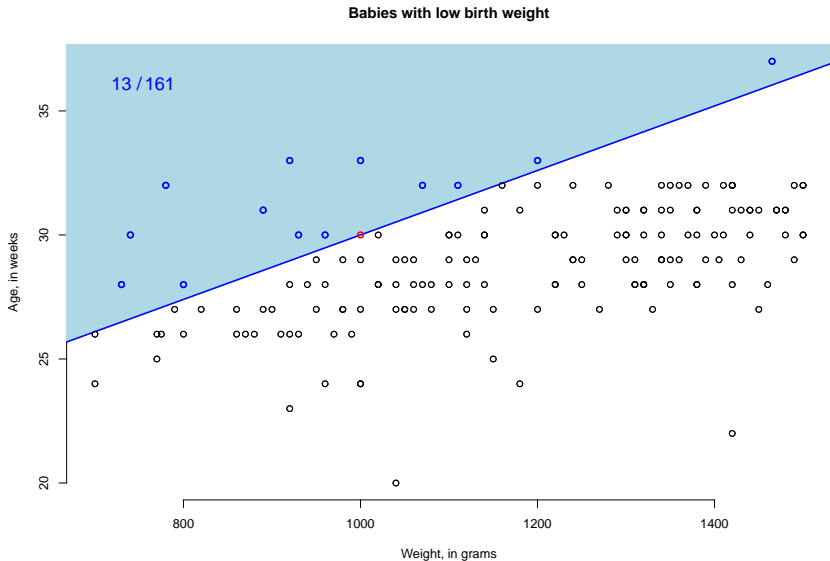
Babies with low birth weight



# Tukey (=halfspace, location) data depth

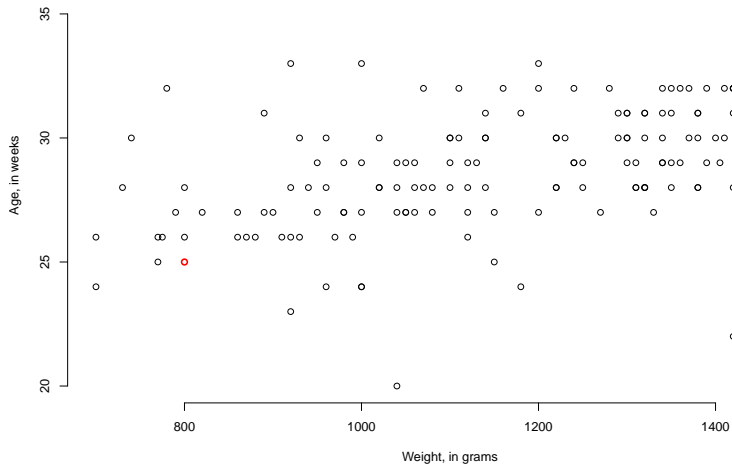


# Tukey (=halfspace, location) data depth



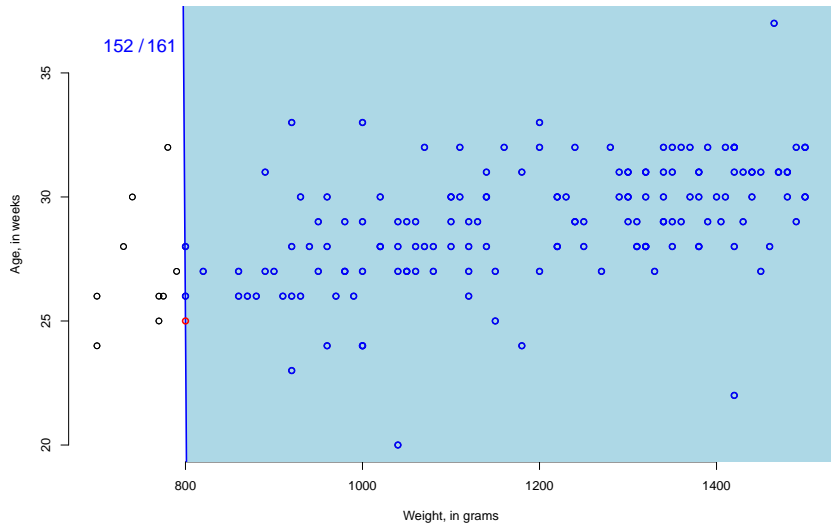
# Tukey (=halfspace, location) data depth

Babies with low birth weight

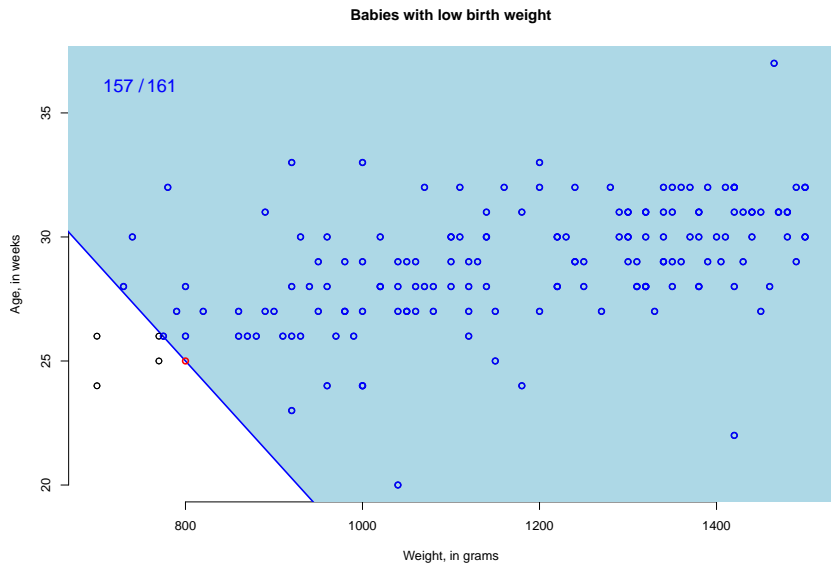


# Tukey (=halfspace, location) data depth

Babies with low birth weight

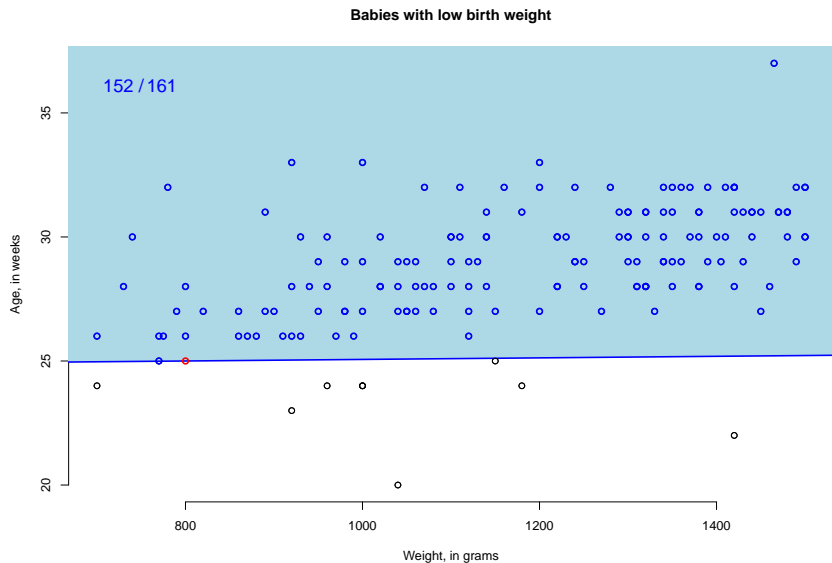


# Tukey (=halfspace, location) data depth

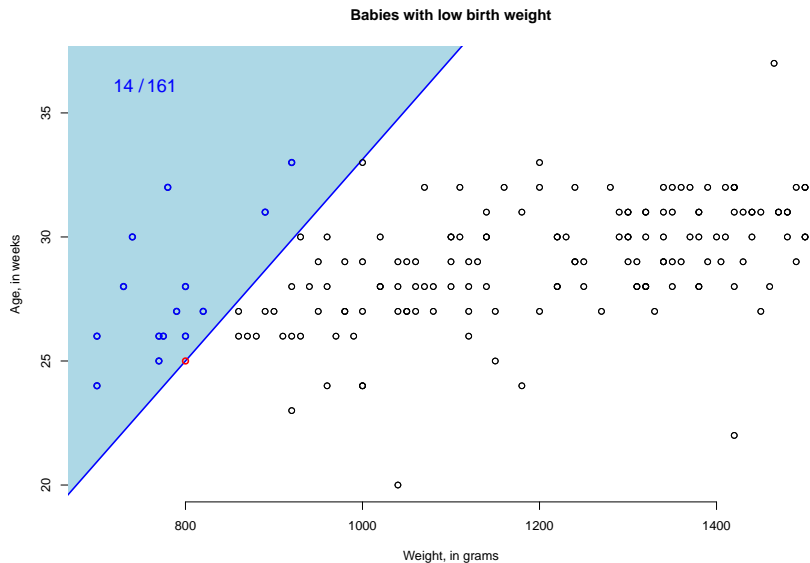




# Tukey (=halfspace, location) data depth

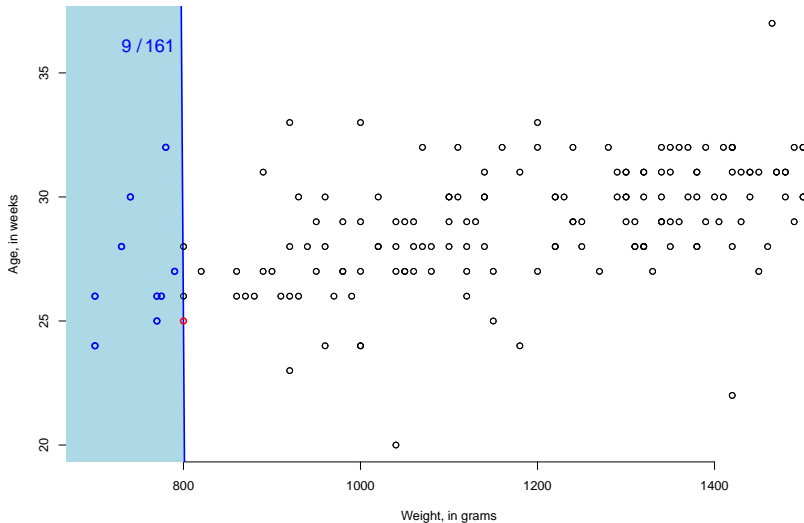


# Tukey (=halfspace, location) data depth



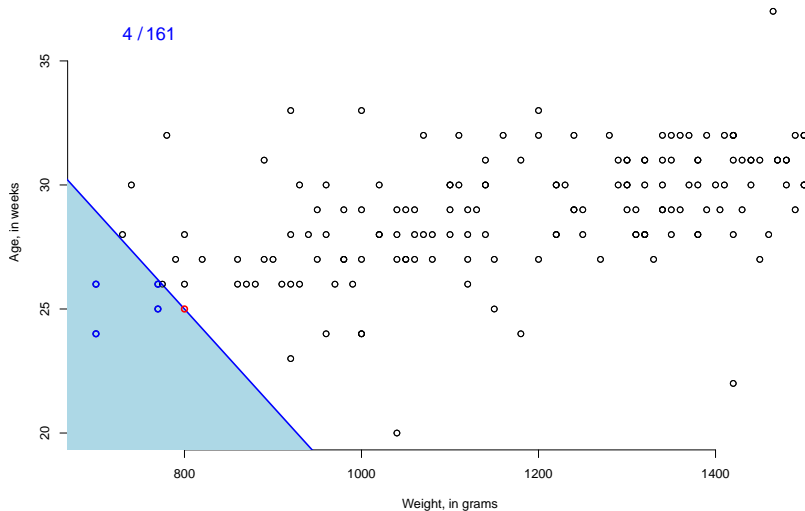
# Tukey (=halfspace, location) data depth

Babies with low birth weight

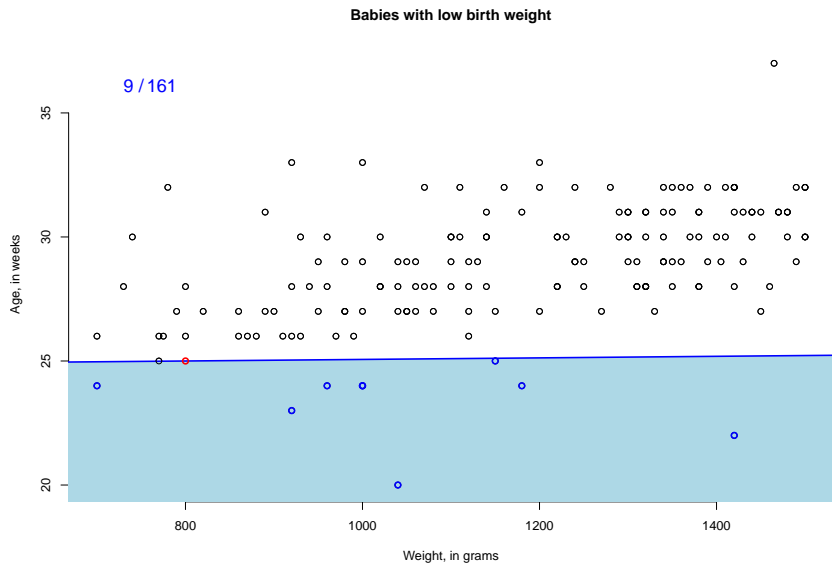


# Tukey (=halfspace, location) data depth

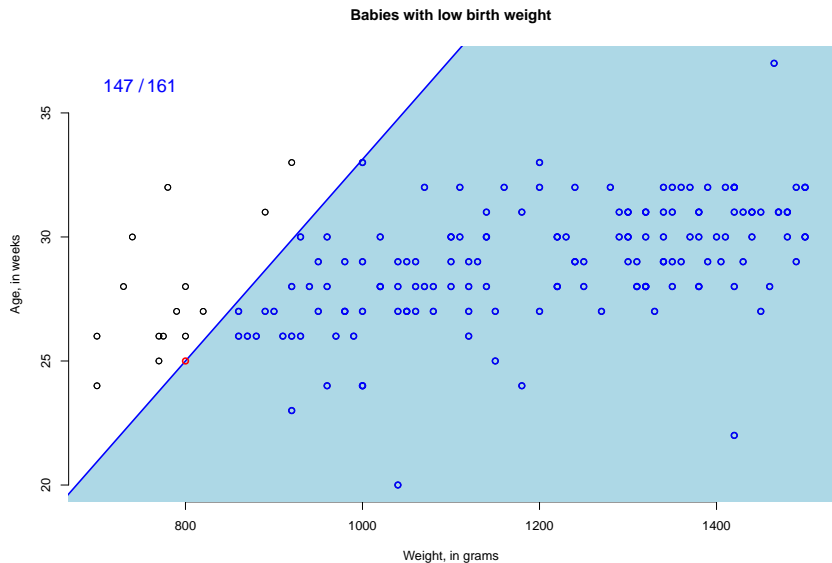
Babies with low birth weight



# Tukey (=halfspace, location) data depth

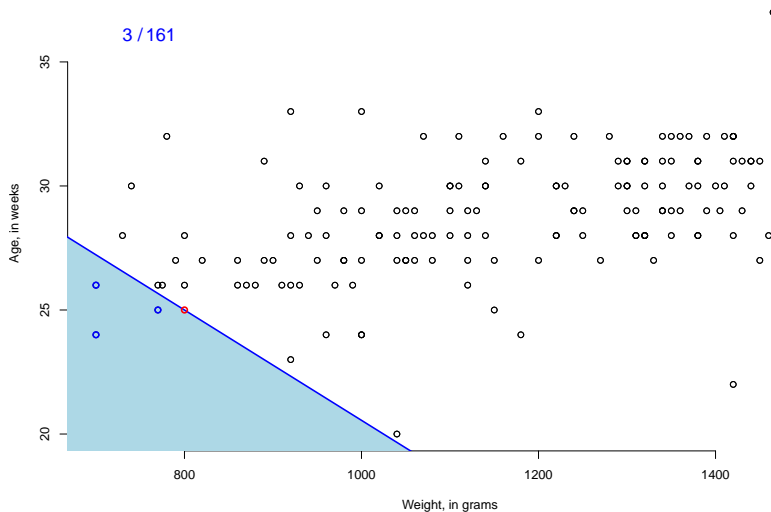


# Tukey (=halfspace, location) data depth

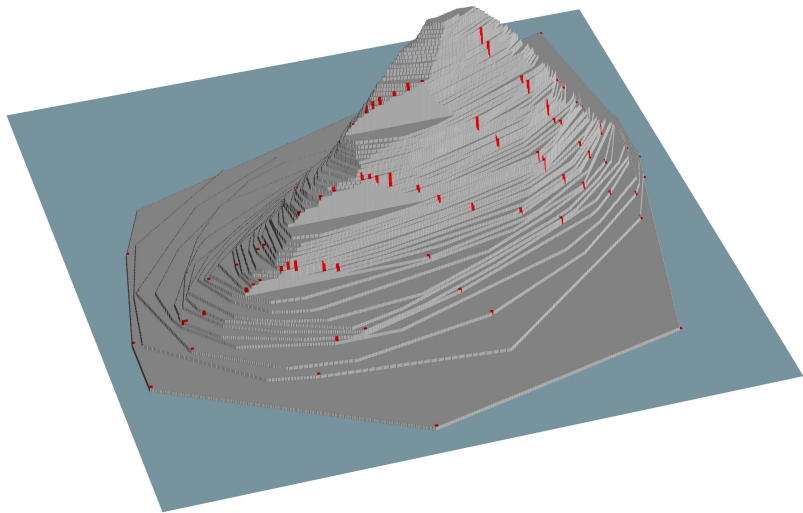


# Tukey (=halfspace, location) data depth

Babies with low birth weight



# Tukey (=halfspace, location) data depth





# Contents

## Introduction

## Non-parametric approaches

One-class support vector machines

Local outlier factor

Isolation forest

## Systematic orderings: data depth

The notion of data depth

The Tukey depth function

**Central regions**

Further depth notions

## Practical session

## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** of **central regions**.

## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** or **central regions**.
- ▶ The innermost region arises at some depth  $\alpha_{\max} \leq 1$ , which depends on the depth notion  $D$  and distribution  $P$ . Then  $D_\alpha(X)$  is the set of **deepest points**.

## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** of **central regions**.
- ▶ The innermost region arises at some depth  $\alpha_{\max} \leq 1$ , which depends on the depth notion  $D$  and distribution  $P$ . Then  $D_\alpha(X)$  is the set of **deepest points**.
- ▶ Central regions describe distribution w.r.t. location, dispersion, and shape.

## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** of **central regions**.
- ▶ The innermost region arises at some depth  $\alpha_{\max} \leq 1$ , which depends on the depth notion  $D$  and distribution  $P$ . Then  $D_\alpha(X)$  is the set of **deepest points**.
- ▶ Central regions describe distribution w.r.t. location, dispersion, and shape.
- ▶ Properties of central regions, for any  $\alpha$ :
  - ▶ Due to **D1** and **D2**  $D_\alpha(X)$  is **affine equivariant**:  
 $D_\alpha(AX + b) = AD_\alpha(X) + b$  for any  $p \times p$  non-singular matrix  $A$  and any  $b \in \mathbb{R}^p$ ;

## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** of **central regions**.
- ▶ The innermost region arises at some depth  $\alpha_{\max} \leq 1$ , which depends on the depth notion  $D$  and distribution  $P$ . Then  $D_\alpha(X)$  is the set of **deepest points**.
- ▶ Central regions describe distribution w.r.t. location, dispersion, and shape.
- ▶ Properties of central regions, for any  $\alpha$ :
  - ▶ Due to **D1** and **D2**  $D_\alpha(X)$  is **affine equivariant**:  
 $D_\alpha(AX + b) = AD_\alpha(X) + b$  for any  $p \times p$  non-singular matrix  $A$  and any  $b \in \mathbb{R}^p$ ;
  - ▶ Due to **D3**  $D_\alpha(X)$  is bounded;

## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** or **central regions**.
- ▶ The innermost region arises at some depth  $\alpha_{\max} \leq 1$ , which depends on the depth notion  $D$  and distribution  $P$ . Then  $D_\alpha(X)$  is the set of **deepest points**.
- ▶ Central regions describe distribution w.r.t. location, dispersion, and shape.
- ▶ Properties of central regions, for any  $\alpha$ :
  - ▶ Due to **D1** and **D2**  $D_\alpha(X)$  is **affine equivariant**:  
 $D_\alpha(AX + b) = AD_\alpha(X) + b$  for any  $p \times p$  non-singular matrix  $A$  and any  $b \in \mathbb{R}^p$ ;
  - ▶ Due to **D3**  $D_\alpha(X)$  is bounded;
  - ▶ Due to **D4**  $D_\alpha(X)$ -s are nested:  
if  $\alpha \geq \beta$ , then  $D_\alpha(X) \subseteq D_\beta(X)$ , and star-shaped;

## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** or **central regions**.
- ▶ The innermost region arises at some depth  $\alpha_{\max} \leq 1$ , which depends on the depth notion  $D$  and distribution  $P$ . Then  $D_\alpha(X)$  is the set of **deepest points**.
- ▶ Central regions describe distribution w.r.t. location, dispersion, and shape.
- ▶ Properties of central regions, for any  $\alpha$ :
  - ▶ Due to **D1** and **D2**  $D_\alpha(X)$  is **affine equivariant**:  
 $D_\alpha(AX + b) = AD_\alpha(X) + b$  for any  $p \times p$  non-singular matrix  $A$  and any  $b \in \mathbb{R}^p$ ;
  - ▶ Due to **D3**  $D_\alpha(X)$  is bounded;
  - ▶ Due to **D4**  $D_\alpha(X)$ -s are nested:  
if  $\alpha \geq \beta$ , then  $D_\alpha(X) \subseteq D_\beta(X)$ , and star-shaped;  
due to **D4con**  $D_\alpha(X)$  is in addition convex;



## Central regions

- ▶ For given distribution  $P$  and  $\alpha \in [0, 1]$ , the level sets  $D_\alpha(P)$  form a family of **depth-trimmed** or **central regions**.
- ▶ The innermost region arises at some depth  $\alpha_{\max} \leq 1$ , which depends on the depth notion  $D$  and distribution  $P$ . Then  $D_\alpha(X)$  is the set of **deepest points**.
- ▶ Central regions describe distribution w.r.t. location, dispersion, and shape.
- ▶ Properties of central regions, for any  $\alpha$ :
  - ▶ Due to **D1** and **D2**  $D_\alpha(X)$  is **affine equivariant**:  
 $D_\alpha(AX + b) = AD_\alpha(X) + b$  for any  $p \times p$  non-singular matrix  $A$  and any  $b \in \mathbb{R}^p$ ;
  - ▶ Due to **D3**  $D_\alpha(X)$  is bounded;
  - ▶ Due to **D4**  $D_\alpha(X)$ -s are nested:  
if  $\alpha \geq \beta$ , then  $D_\alpha(X) \subseteq D_\beta(X)$ , and star-shaped;  
due to **D4con**  $D_\alpha(X)$  is in addition convex;
  - ▶ Due to **D5**  $D_\alpha(X)$  is closed.

## Tukey-trimmed regions

Tukey depth defines a family of (depth-)trimmed (central) regions  $D_{\tau}^T(X)$ , the upper-level sets of the depth function:

$$D_{\tau}^T(X) = \{\mathbf{x} \in \mathbb{R}^p : D^T(\mathbf{x}|X) \geq \tau\}.$$

# Tukey-trimmed regions

Tukey depth defines a family of (depth-)trimmed (central) regions  $D_{\tau}^T(X)$ , the upper-level sets of the depth function:

$$D_{\tau}^T(X) = \{\mathbf{x} \in \mathbb{R}^p : D^T(\mathbf{x}|X) \geq \tau\}.$$

## Properties:

### Depth:

- ▶ Affine invariant;

### Regions:

Affine equivariant;

# Tukey-trimmed regions

Tukey depth defines a family of (depth-)trimmed (central) regions  $D_{\tau}^T(X)$ , the upper-level sets of the depth function:

$$D_{\tau}^T(X) = \{\mathbf{x} \in \mathbb{R}^p : D^T(\mathbf{x}|X) \geq \tau\}.$$

## Properties:

### Depth:

- ▶ Affine invariant;
- ▶ Vanishing at infinity;

### Regions:

- Affine equivariant;
- Bounded;

# Tukey-trimmed regions

Tukey depth defines a family of (depth-)trimmed (central) regions  $D_{\tau}^T(X)$ , the upper-level sets of the depth function:

$$D_{\tau}^T(X) = \{\mathbf{x} \in \mathbb{R}^p : D^T(\mathbf{x}|X) \geq \tau\}.$$

## Properties:

### Depth:

- ▶ Affine invariant;
- ▶ Vanishing at infinity;
- ▶ Monotone w.r.t. deepest point;

### Regions:

- Affine equivariant;
- Bounded;
- Nested;

# Tukey-trimmed regions

Tukey depth defines a family of (depth-)trimmed (central) regions  $D_{\tau}^T(X)$ , the upper-level sets of the depth function:

$$D_{\tau}^T(X) = \{\mathbf{x} \in \mathbb{R}^p : D^T(\mathbf{x}|X) \geq \tau\}.$$

## Properties:

### Depth:

- ▶ Affine invariant;
- ▶ Vanishing at infinity;
- ▶ Monotone w.r.t. deepest point;
- ▶ Upper-semicontinuous;

### Regions:

- Affine equivariant;
- Bounded;
- Nested;
- Closed;

# Tukey-trimmed regions

Tukey depth defines a family of (depth-)trimmed (central) regions  $D_\tau^T(X)$ , the upper-level sets of the depth function:

$$D_\tau^T(X) = \{\mathbf{x} \in \mathbb{R}^p : D^T(\mathbf{x}|X) \geq \tau\}.$$

## Properties:

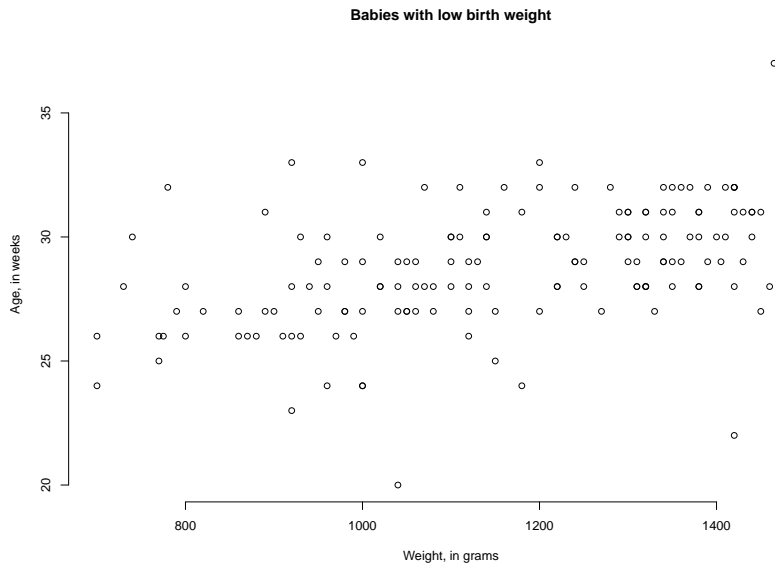
### Depth:

- ▶ Affine invariant;
- ▶ Vanishing at infinity;
- ▶ Monotone w.r.t. deepest point;
- ▶ Upper-semicontinuous;
- ▶ Quasiconcave.

### Regions:

- Affine equivariant;
- Bounded;
- Nested;
- Closed;
- Convex.

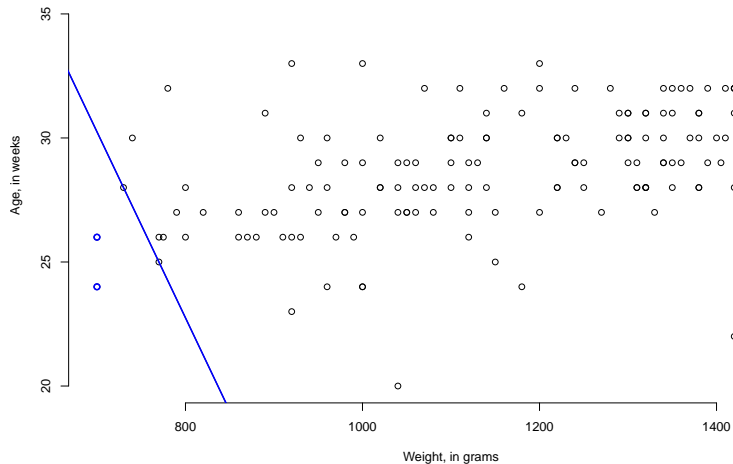
# Tukey (=halfspace, location) depth-trimmed regions





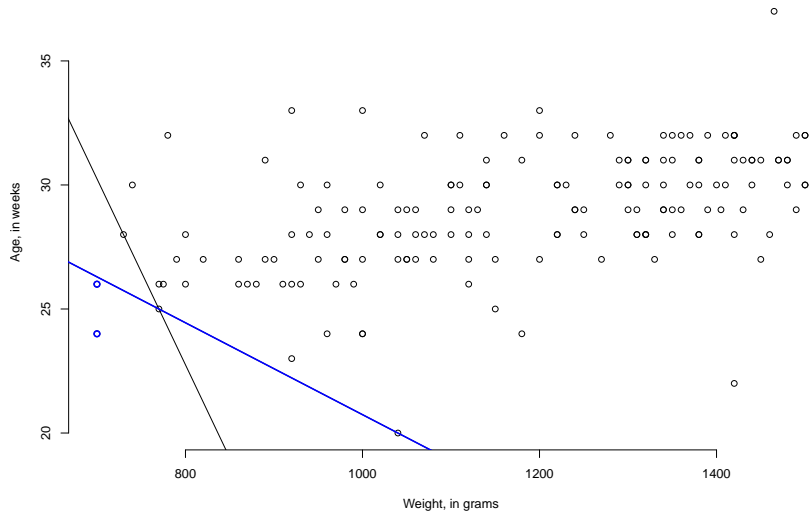
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



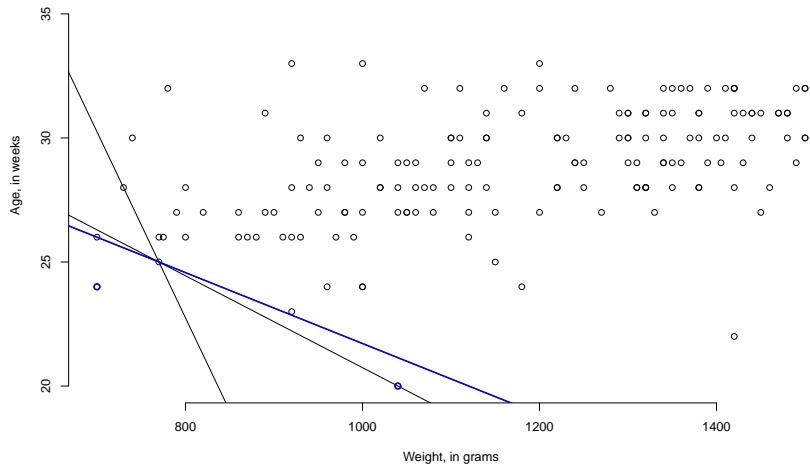
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



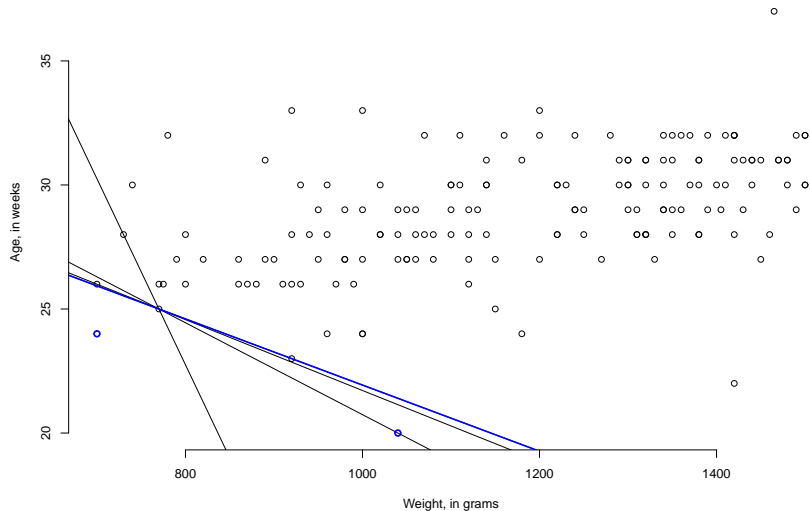
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



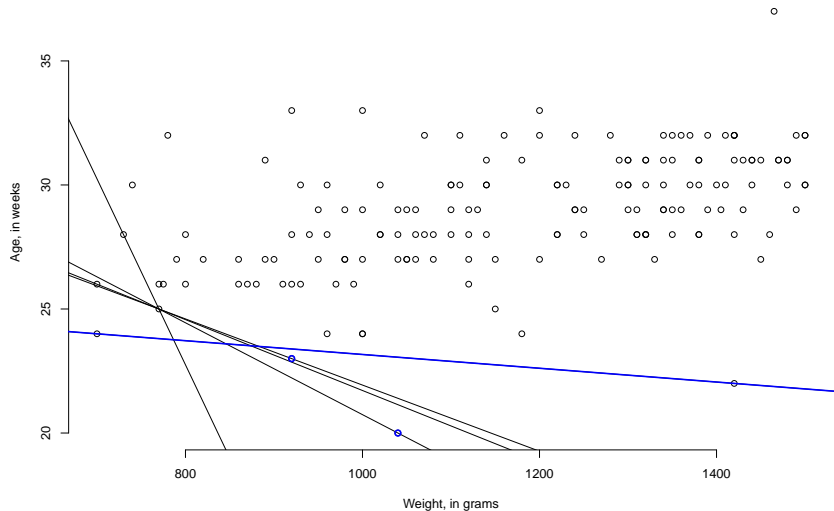
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



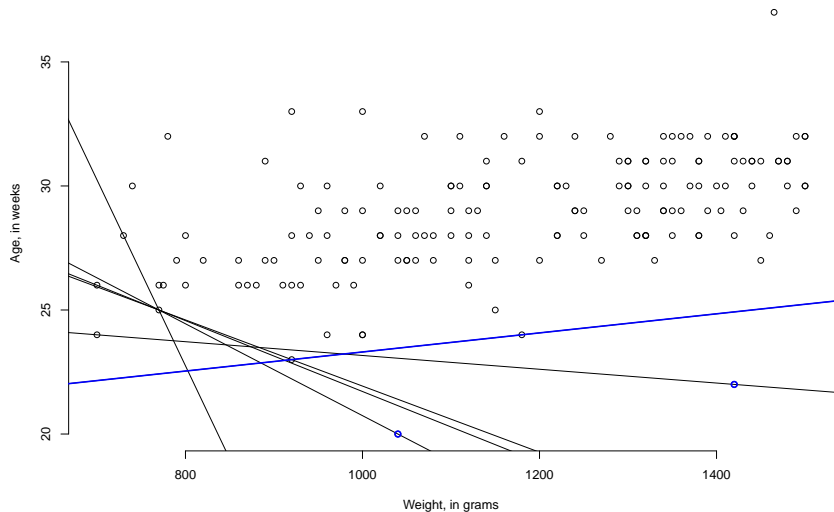
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



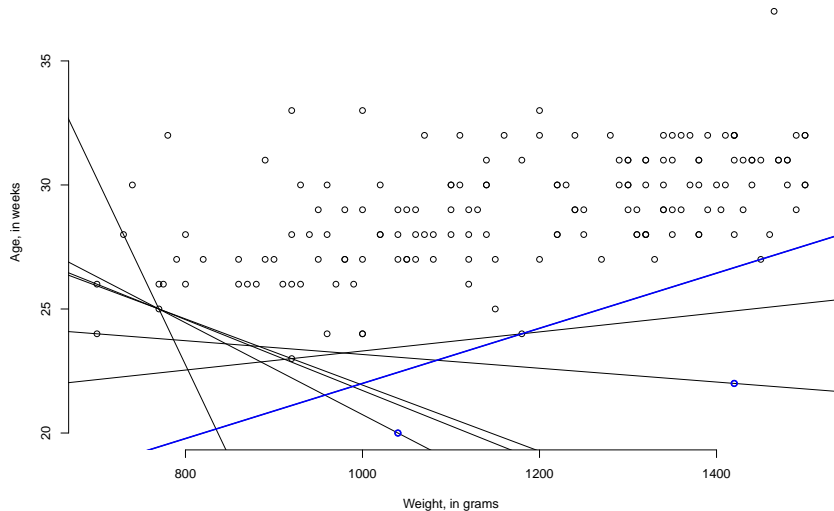
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



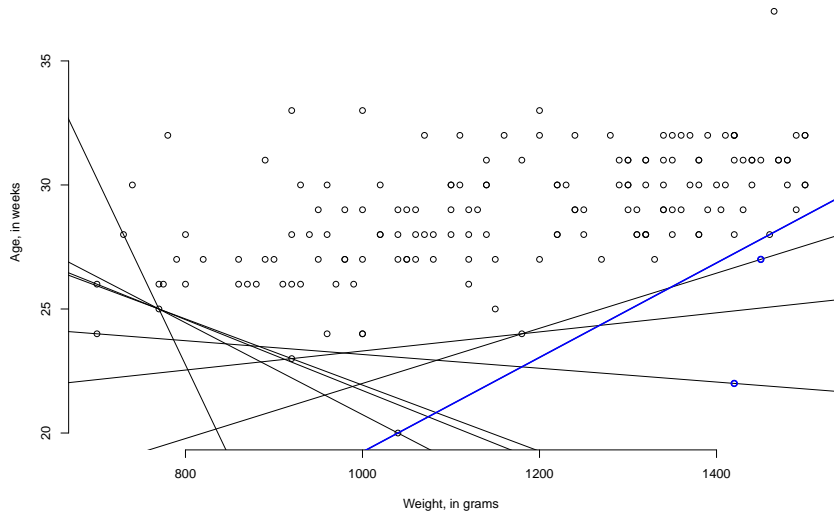
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



# Tukey (=halfspace, location) depth-trimmed regions

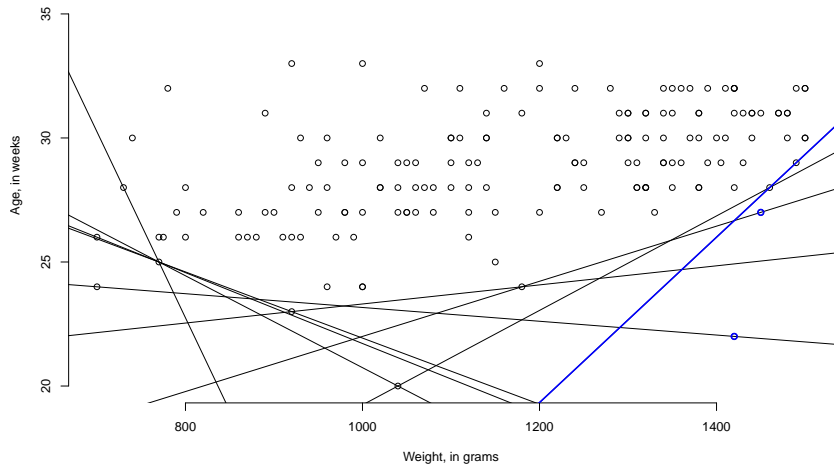
Babies with low birth weight





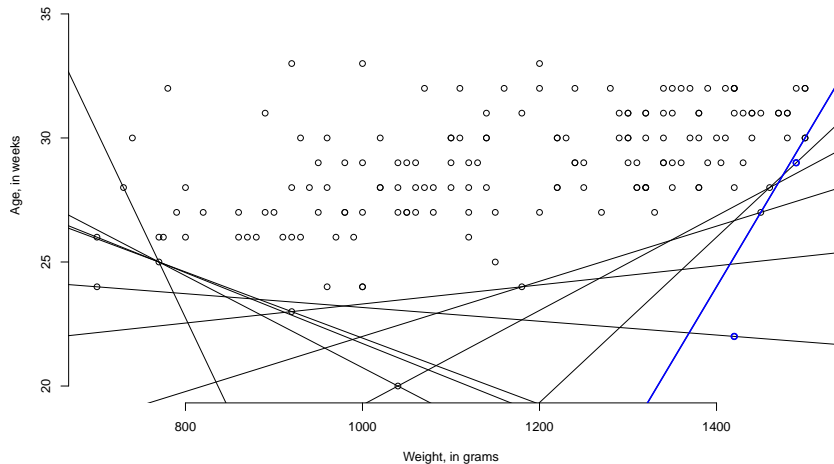
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight

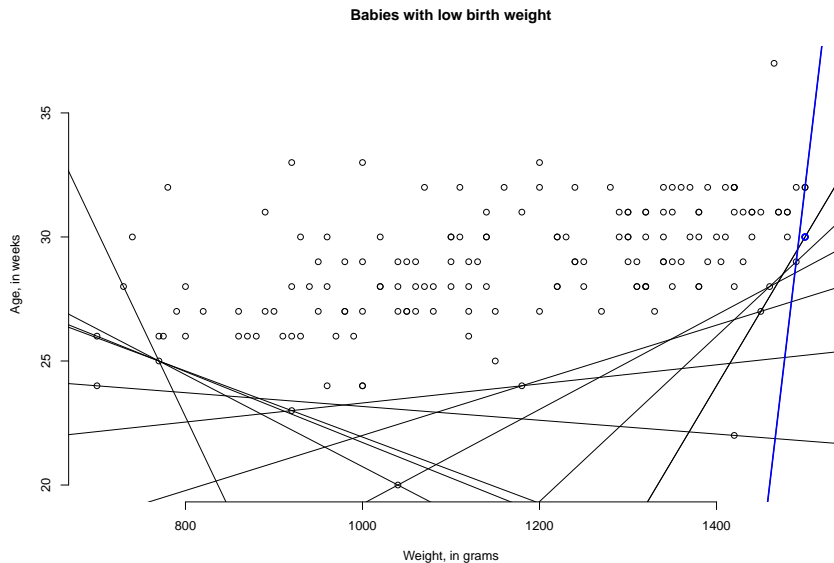


# Tukey (=halfspace, location) depth-trimmed regions

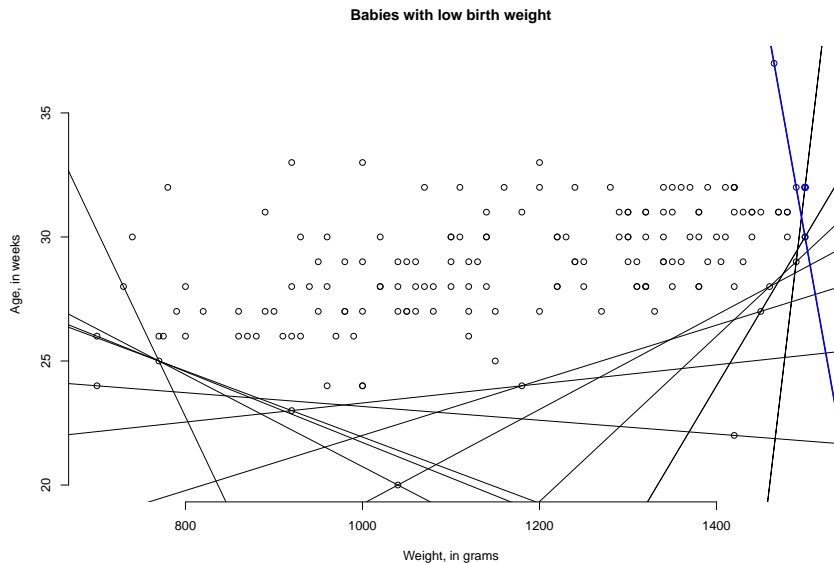
Babies with low birth weight



# Tukey (=halfspace, location) depth-trimmed regions

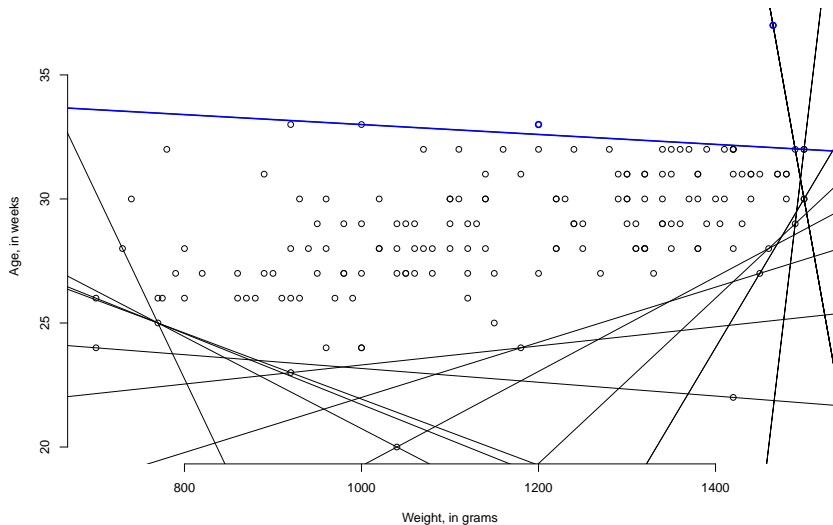


# Tukey (=halfspace, location) depth-trimmed regions

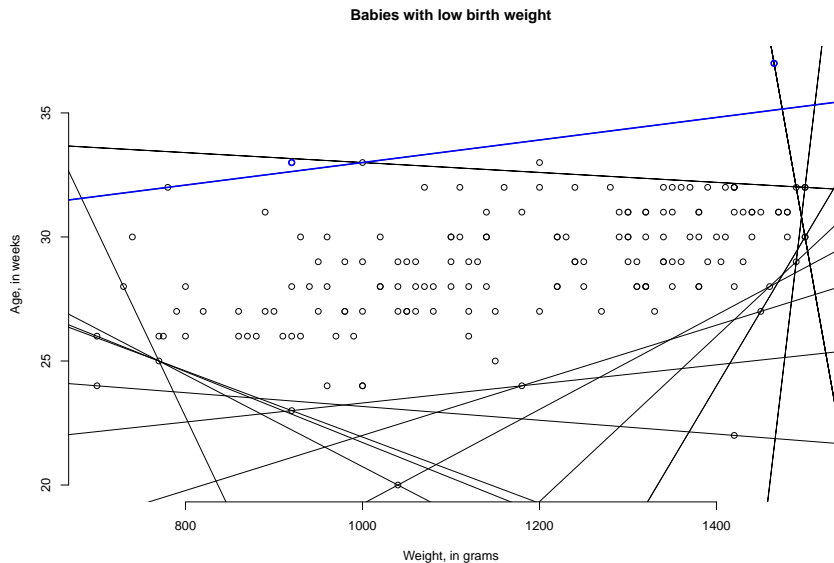


# Tukey (=halfspace, location) depth-trimmed regions

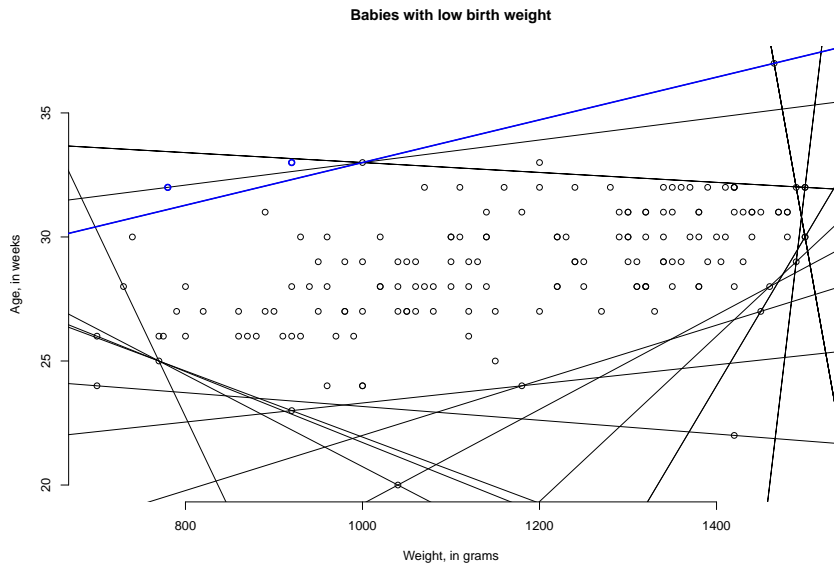
Babies with low birth weight



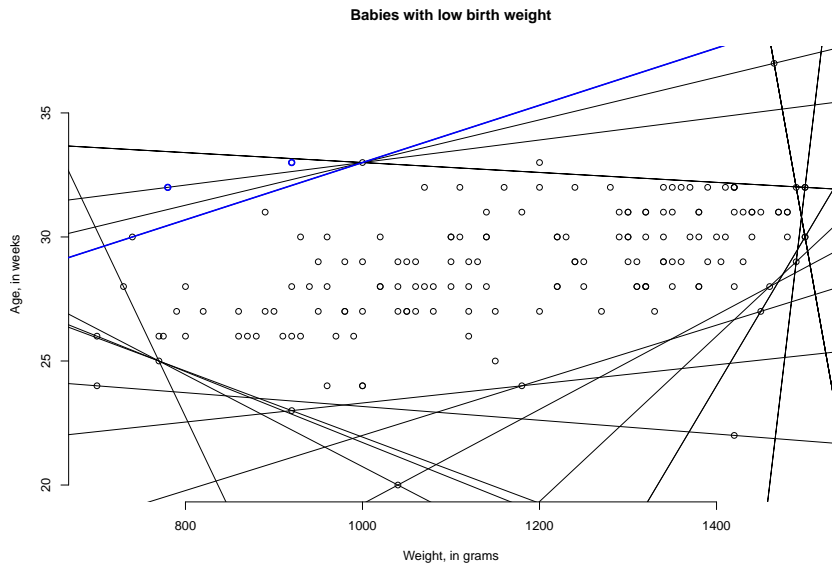
# Tukey (=halfspace, location) depth-trimmed regions



# Tukey (=halfspace, location) depth-trimmed regions

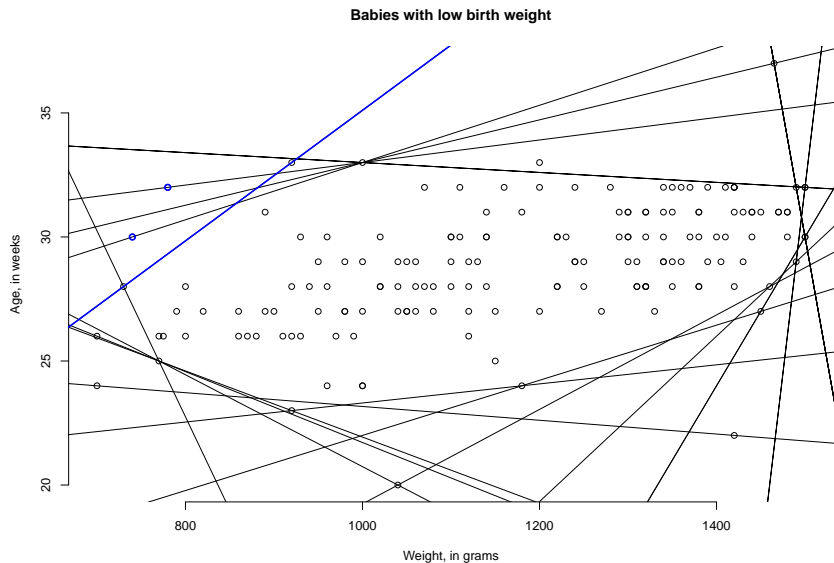


# Tukey (=halfspace, location) depth-trimmed regions

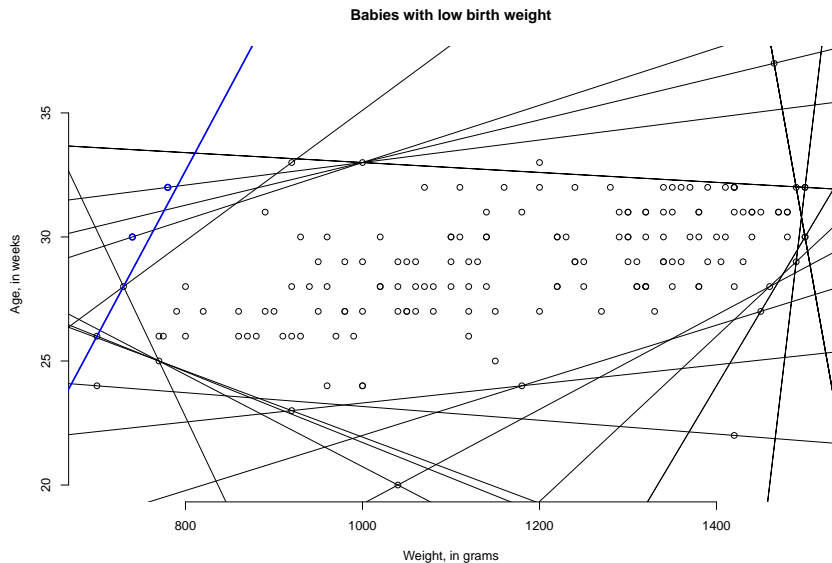




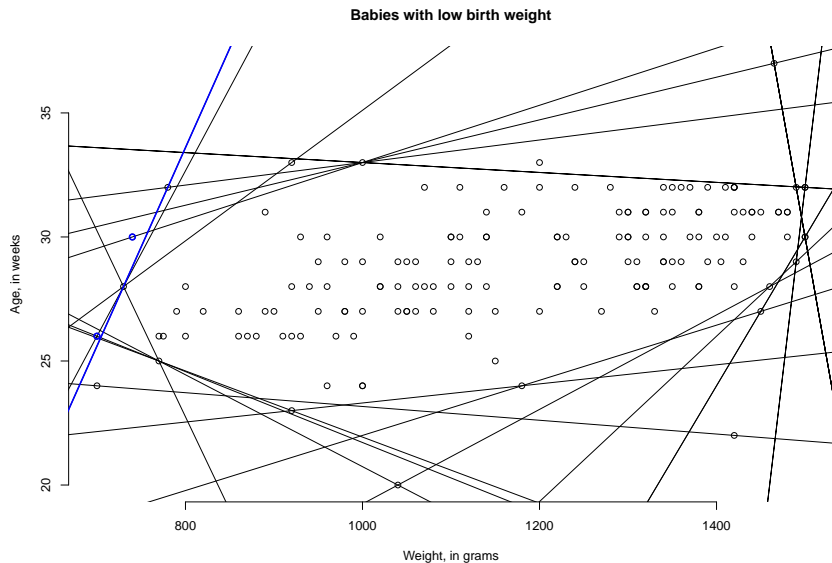
# Tukey (=halfspace, location) depth-trimmed regions



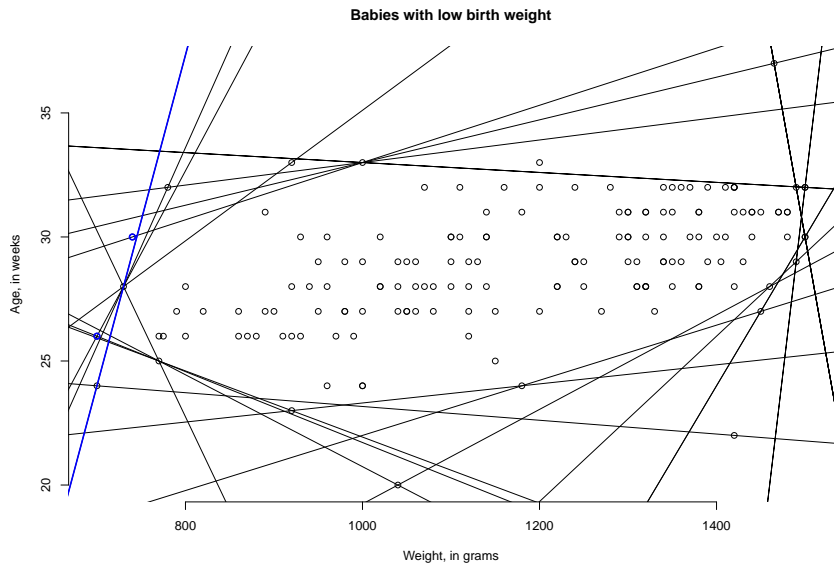
# Tukey (=halfspace, location) depth-trimmed regions



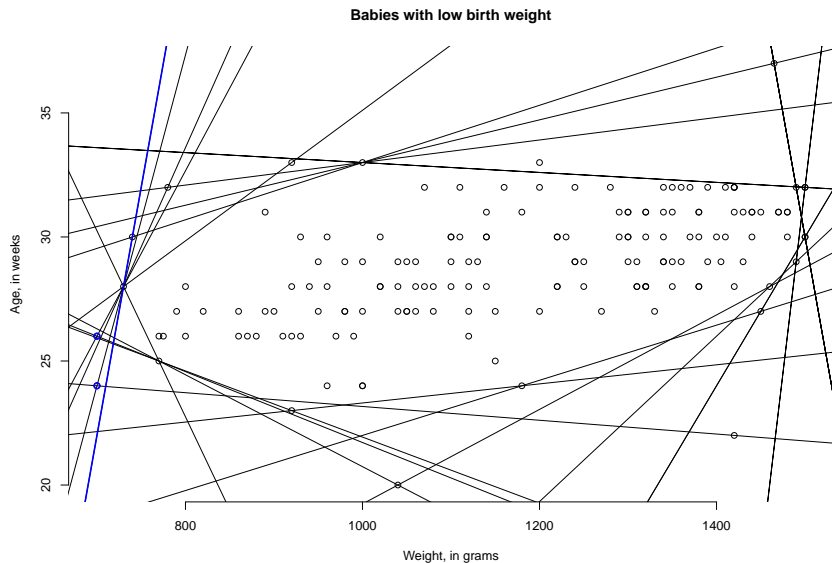
# Tukey (=halfspace, location) depth-trimmed regions



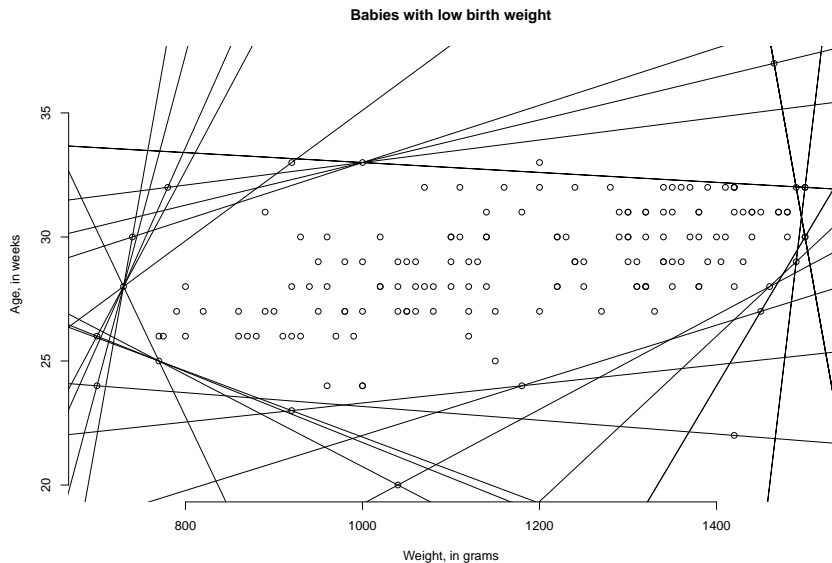
# Tukey (=halfspace, location) depth-trimmed regions



# Tukey (=halfspace, location) depth-trimmed regions

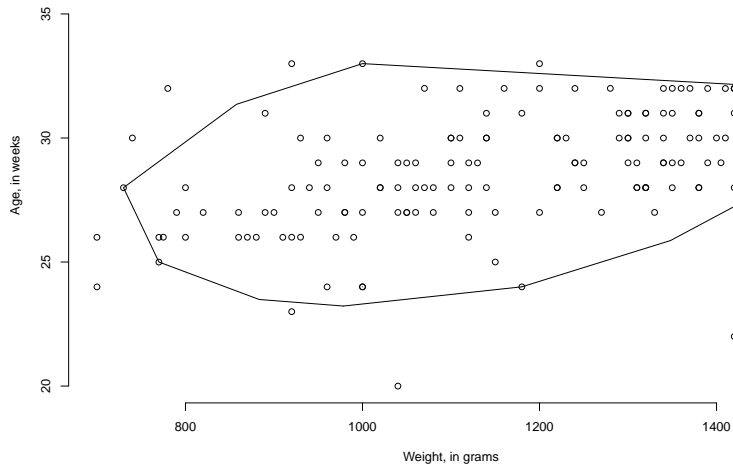


# Tukey (=halfspace, location) depth-trimmed regions



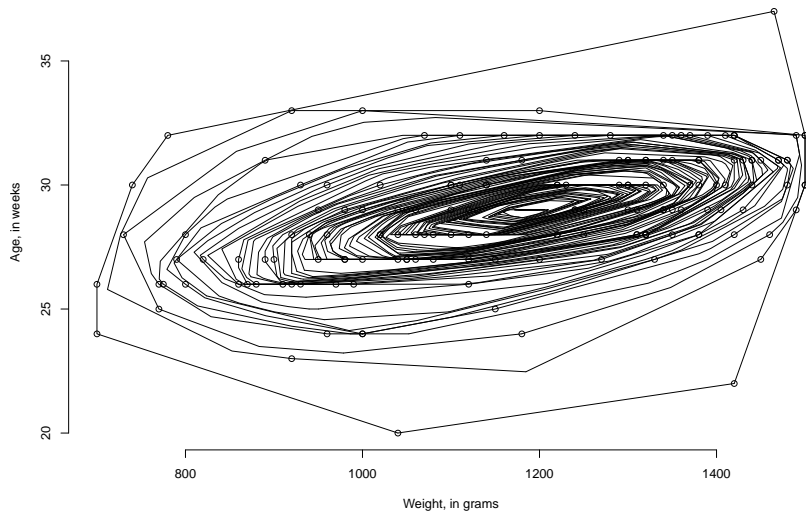
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight



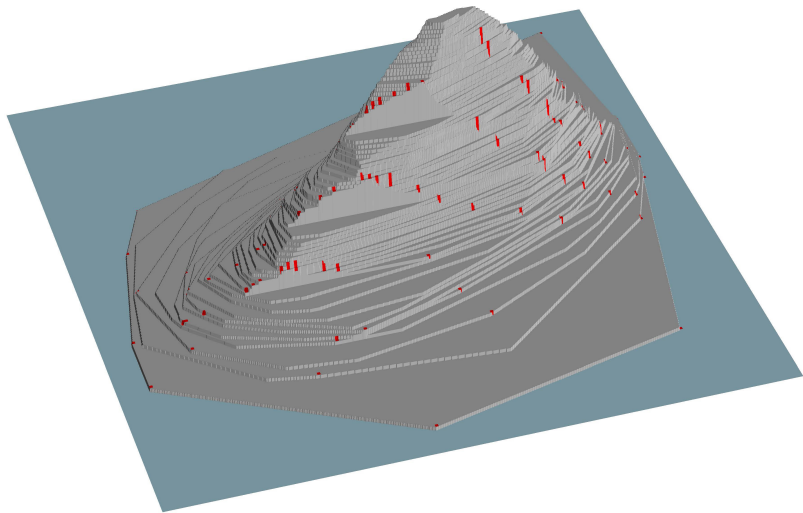
# Tukey (=halfspace, location) depth-trimmed regions

Babies with low birth weight

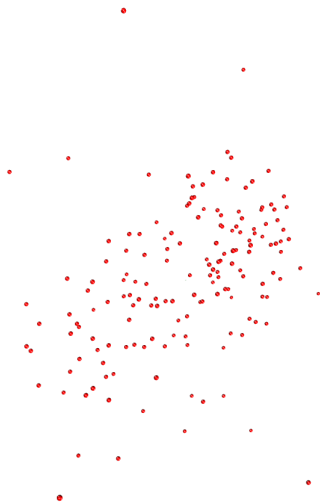




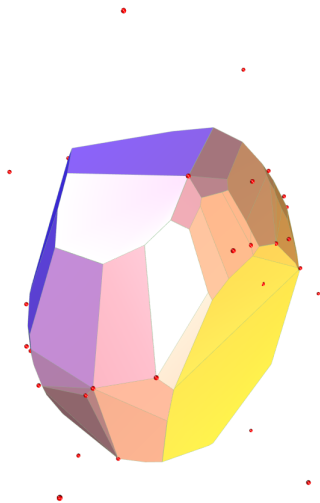
# Tukey (=halfspace, location) data depth



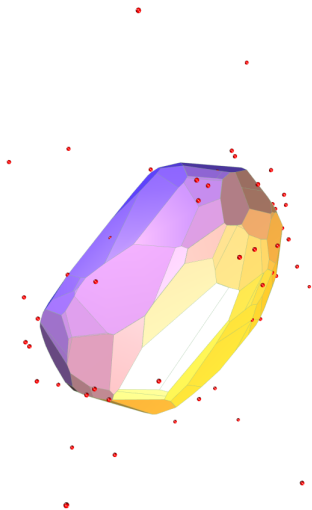
# Tukey (=halfspace, location) depth region



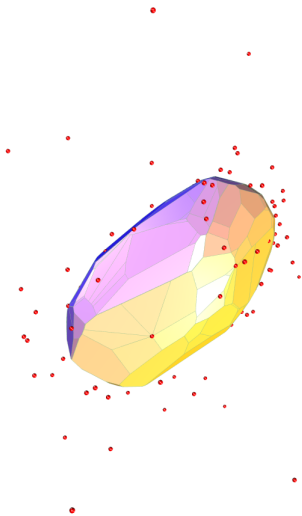
Tukey (=halfspace, location) depth region:  $\tau = 2/161$



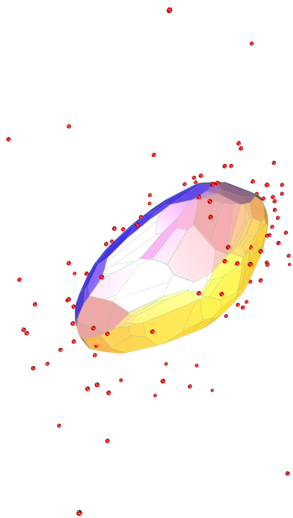
Tukey (=halfspace, location) depth region:  $\tau = 5/161$



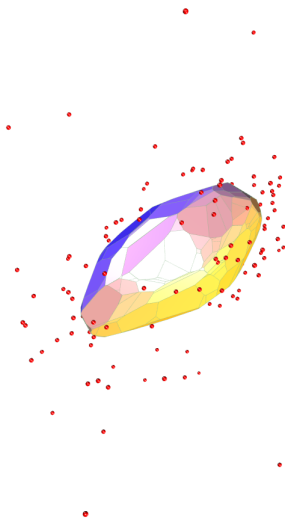
Tukey (=halfspace, location) depth region:  $\tau = 9/161$



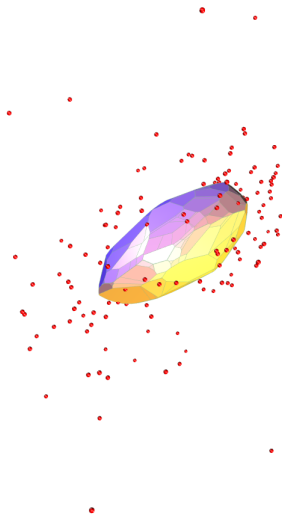
Tukey (=halfspace, location) depth region:  $\tau = 13/161$



Tukey (=halfspace, location) depth region:  $\tau = 17/161$

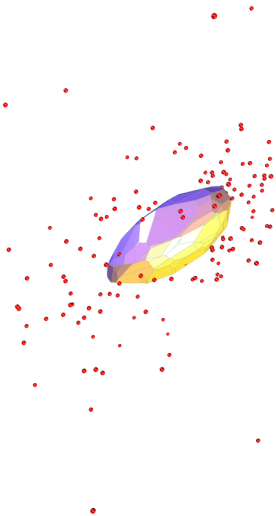


Tukey (=halfspace, location) depth region:  $\tau = 25/161$

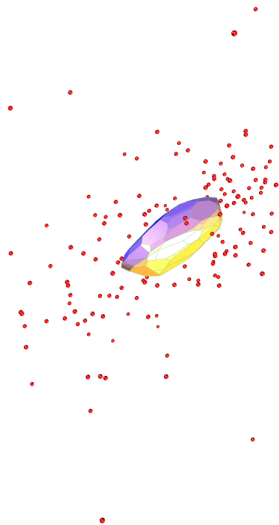




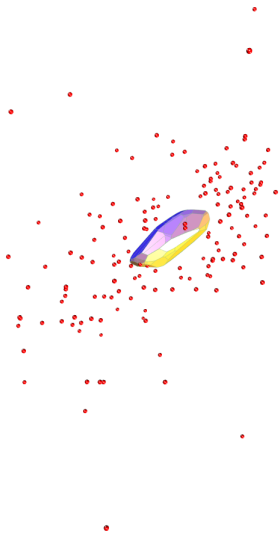
Tukey (=halfspace, location) depth region:  $\tau = 33/161$



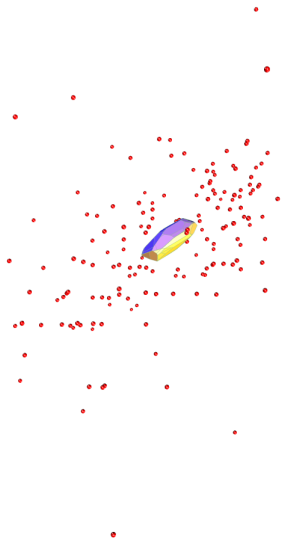
Tukey (=halfspace, location) depth region:  $\tau = 41/161$



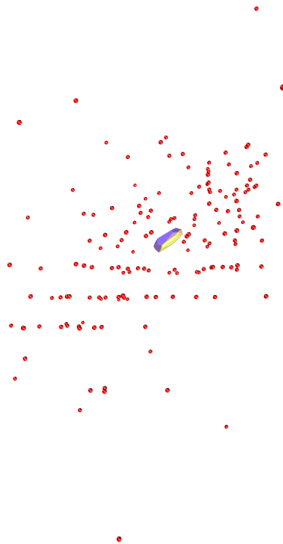
Tukey (=halfspace, location) depth region:  $\tau = 49/161$



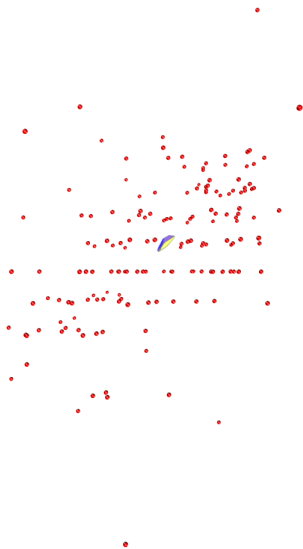
Tukey (=halfspace, location) depth region:  $\tau = 57/161$



Tukey (=halfspace, location) depth region:  $\tau = 65/161$



Tukey (=halfspace, location) depth region:  $\tau = 68/161$



# Contents

## Introduction

## Non-parametric approaches

One-class support vector machines

Local outlier factor

Isolation forest

## Systematic orderings: data depth

The notion of data depth

The Tukey depth function

Central regions

Further depth notions

## Practical session

## Mahalanobis depth (Mahalanobis, 1936)

- ▶ **Mahalanobis depth** is defined as:

$$D^{Mah}(\mathbf{x}|X) = \frac{1}{1 + (\delta^{Mah})^2(\mathbf{x}|X)},$$

based on Mahalanobis distance:

$$(\delta^{Mah})^2(\mathbf{x}|X) = (\mathbf{x} - \boldsymbol{\mu}_X)^T \boldsymbol{\Sigma}_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X).$$



## Mahalanobis depth (Mahalanobis, 1936)

- ▶ **Mahalanobis depth** is defined as:

$$D^{Mah}(\mathbf{x}|X) = \frac{1}{1 + (\delta^{Mah})^2(\mathbf{x}|X)},$$

based on Mahalanobis distance:

$$(\delta^{Mah})^2(\mathbf{x}|X) = (\mathbf{x} - \boldsymbol{\mu}_X)^T \boldsymbol{\Sigma}_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X).$$

- ▶ In the empirical version,  $\boldsymbol{\mu}_X$  and  $\boldsymbol{\Sigma}_X$  are substituted by suitable estimates:
  - ▶ moment estimates;
  - ▶ robust estimates such as **minimum volume ellipsoid** or **minimum covariance determinant (MCD)**.

## Mahalanobis depth (Mahalanobis, 1936)

- ▶ **Mahalanobis depth** is defined as:

$$D^{Mah}(\mathbf{x}|X) = \frac{1}{1 + (\delta^{Mah})^2(\mathbf{x}|X)},$$

based on Mahalanobis distance:

$$(\delta^{Mah})^2(\mathbf{x}|X) = (\mathbf{x} - \boldsymbol{\mu}_X)^T \boldsymbol{\Sigma}_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X).$$

- ▶ In the empirical version,  $\boldsymbol{\mu}_X$  and  $\boldsymbol{\Sigma}_X$  are substituted by suitable estimates:
  - ▶ moment estimates;
  - ▶ robust estimates such as **minimum volume ellipsoid** or **minimum covariance determinant (MCD)**.
- ▶ Properties:
  - ▶ satisfies **D1 – D5** and **D4con**, is continuous;

## Mahalanobis depth (Mahalanobis, 1936)

- ▶ **Mahalanobis depth** is defined as:

$$D^{Mah}(\mathbf{x}|X) = \frac{1}{1 + (\delta^{Mah})^2(\mathbf{x}|X)},$$

based on Mahalanobis distance:

$$(\delta^{Mah})^2(\mathbf{x}|X) = (\mathbf{x} - \boldsymbol{\mu}_X)^T \boldsymbol{\Sigma}_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X).$$

- ▶ In the empirical version,  $\boldsymbol{\mu}_X$  and  $\boldsymbol{\Sigma}_X$  are substituted by suitable estimates:
  - ▶ moment estimates;
  - ▶ robust estimates such as **minimum volume ellipsoid** or **minimum covariance determinant (MCD)**.
- ▶ Properties:
  - ▶ satisfies **D1 – D5** and **D4con**, is continuous;
  - ▶ being defined by  $d(d + 1)$  parameters, can be seen as a **parametric depth**;

## Mahalanobis depth (Mahalanobis, 1936)

- ▶ **Mahalanobis depth** is defined as:

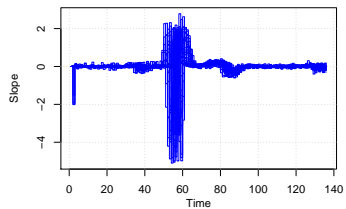
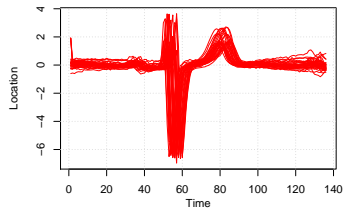
$$D^{Mah}(\mathbf{x}|X) = \frac{1}{1 + (\delta^{Mah})^2(\mathbf{x}|X)},$$

based on Mahalanobis distance:

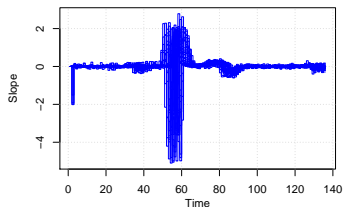
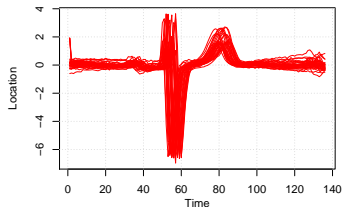
$$(\delta^{Mah})^2(\mathbf{x}|X) = (\mathbf{x} - \boldsymbol{\mu}_X)^T \boldsymbol{\Sigma}_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X).$$

- ▶ In the empirical version,  $\boldsymbol{\mu}_X$  and  $\boldsymbol{\Sigma}_X$  are substituted by suitable estimates:
  - ▶ moment estimates;
  - ▶ robust estimates such as **minimum volume ellipsoid** or **minimum covariance determinant (MCD)**.
- ▶ Properties:
  - ▶ satisfies **D1 – D5** and **D4con**, is continuous;
  - ▶ being defined by  $d(d + 1)$  parameters, can be seen as a **parametric depth**;
  - ▶ by a single elliptical contour characterizes a multivariate **normal distribution** or one within an affine **family of non-degenerate elliptical distributions**.

# ECG five days data



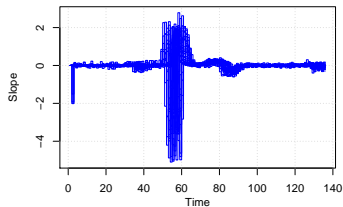
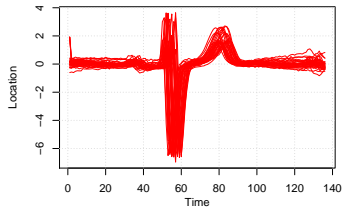
## ECG five days data



$$\hat{f}_i \mapsto \mathbf{x}_i = \left[ \int_0^T \hat{f}_i(t) dt, \int_0^T \hat{f}'_i(t) dt \right],$$

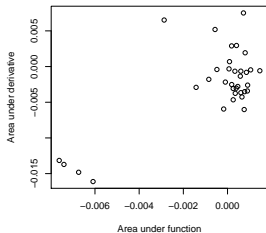
with  $\hat{f}_i(t)$  being the function obtained by connecting the points  $(t_{ij}, f_i(t_{ij}))$ ,  $j = 1, \dots, N_i$  with line segments,  $\hat{f}'_i(t)$  its derivative.

# ECG five days data

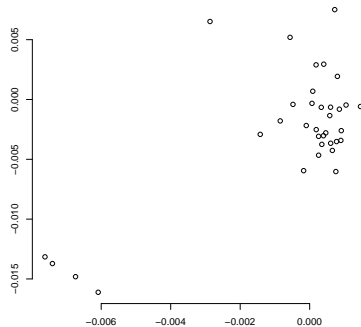


$$\hat{f}_i \mapsto \mathbf{x}_i = \left[ \int_0^T \hat{f}_i(t) dt, \int_0^T \hat{f}'_i(t) dt \right],$$

with  $\hat{f}_i(t)$  being the function obtained by connecting the points  $(t_{ij}, f_i(t_{ij})), j = 1, \dots, N_i$  with line segments,  $\hat{f}'_i(t)$  its derivative.

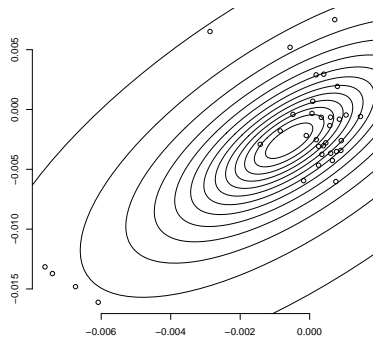
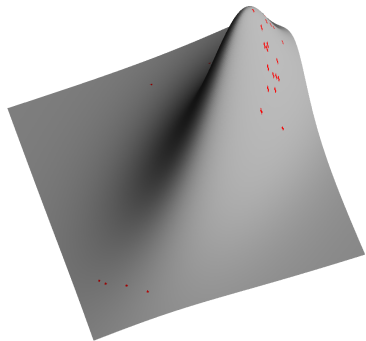


# Mahalanobis depth (Mahalanobis, 1936)

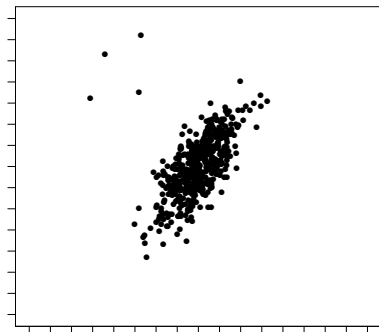




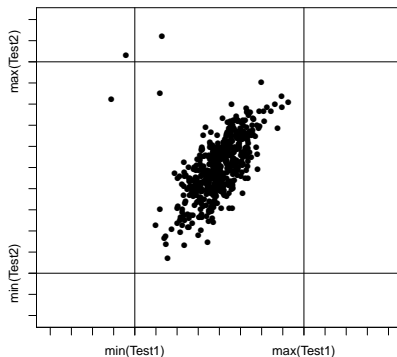
# Mahalanobis depth (Mahalanobis, 1936)



## Multivariate anomaly detection: an example

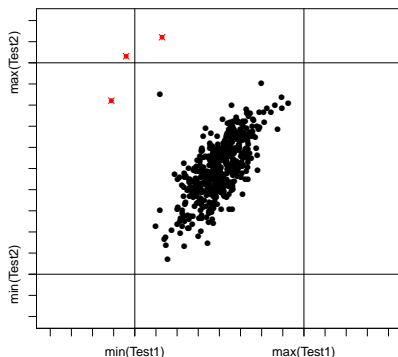


## Multivariate anomaly detection: an example



- ▶ Checking for **minimum** and **maximum** in each test result.

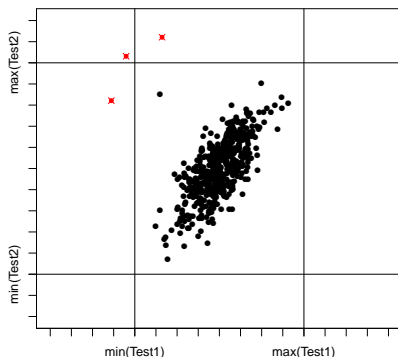
# Multivariate anomaly detection: an example



- ▶ Checking for **minimum** and **maximum** in each test result.
- ▶ Label observation  $x$  as **anomaly** if:

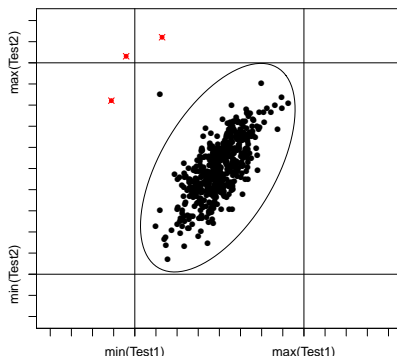
$$x \notin [\min(\text{Test1}), \max(\text{Test1})] \times [\min(\text{Test2}), \max(\text{Test2})].$$

## Multivariate anomaly detection: an example



- ▶ Checking for **minimum** and **maximum** in each test result.
- ▶ Label observation  $x$  as **anomaly** if:  
$$x \notin [\min(\text{Test1}), \max(\text{Test1})] \times [\min(\text{Test2}), \max(\text{Test2})].$$
- ▶ **Not all** anomalies can be detected.

## Multivariate anomaly detection: an example

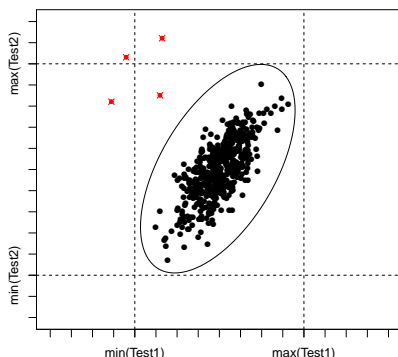


- ▶ **Mahalanobis distance** of an observation  $\mathbf{x} \in \mathbb{R}^2$  (from the mean) is defined as follows:

$$d_{Mah}(\mathbf{x}|\mathbf{X}) = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}),$$

where  $\boldsymbol{\mu}$  is the **mean** and  $\boldsymbol{\Sigma}$  is the **covariance** matrix.

## Multivariate anomaly detection: an example



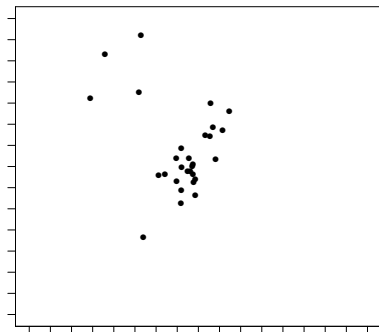
- ▶ **Mahalanobis distance** of an observation  $\mathbf{x} \in \mathbb{R}^2$  (from the mean) is defined as follows:

$$d_{Mah}(\mathbf{x}|\mathbf{X}) = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}),$$

where  $\boldsymbol{\mu}$  is the **mean** and  $\boldsymbol{\Sigma}$  is the **covariance** matrix.

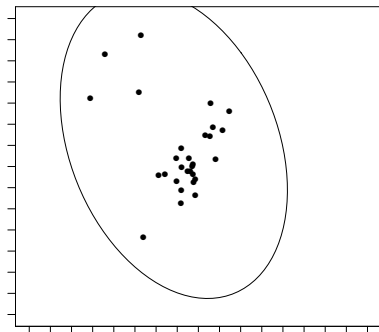
- ▶ Label  $\mathbf{x}$  as **anomaly**  $d_{Mah}(\mathbf{x}|\mathbf{X}) > \max(d_{Mah})$ .

## Multivariate anomaly detection: robustness



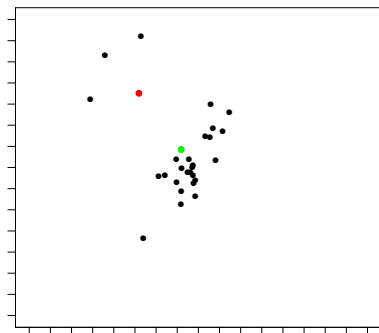


## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.

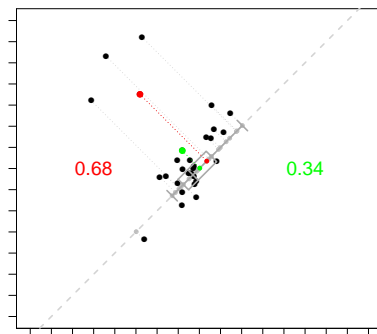
## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

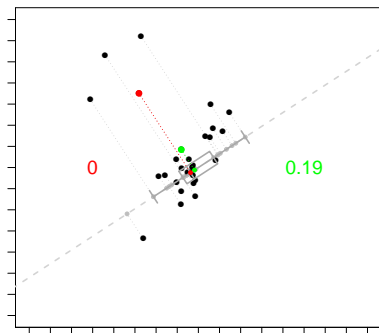
## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

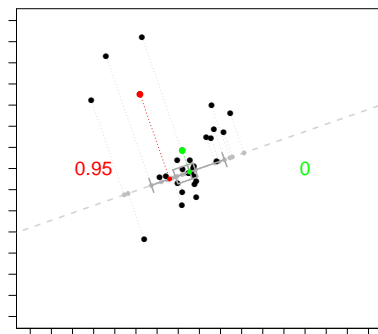
## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

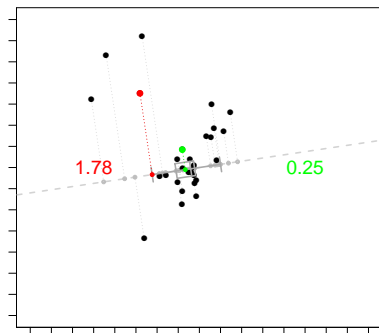
## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

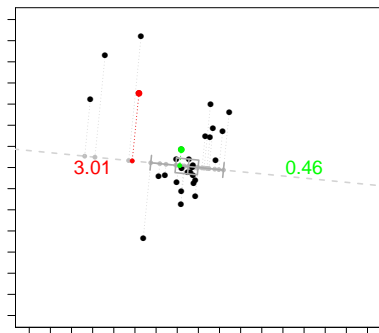
## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

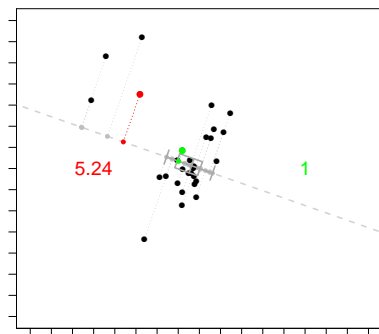
## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

## Multivariate anomaly detection: robustness

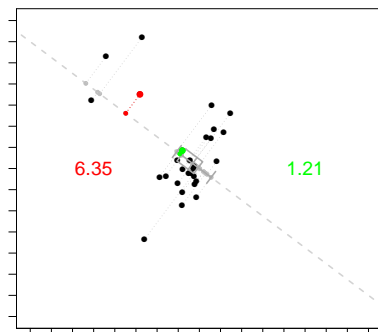


- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$



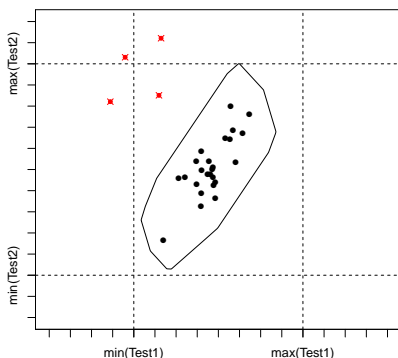
## Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

# Multivariate anomaly detection: robustness

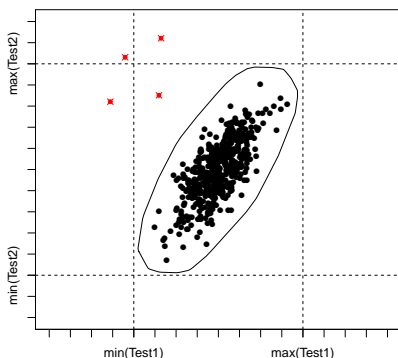


- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

- ▶ Label  $\mathbf{x}$  as **anomaly** if  $O_{SD}(\mathbf{x}|\mathbf{X}) > \max(O_{SD})$ .

# Multivariate anomaly detection: robustness



- ▶ Mahalanobis distance (moment estimators) **not robust**.
- ▶ **Stahel-Donoho outlyingness** of  $\mathbf{x}$  w.r.t.  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$O_{SD}(\mathbf{x}|\mathbf{X}) = \max_{\mathbf{u} \in \mathcal{S}^{d-1}} \frac{|\mathbf{x}^\top \mathbf{u} - \text{med}(\mathbf{X}\mathbf{u})|}{\text{MAD}(\mathbf{X}\mathbf{u})}.$$

- ▶ Label  $\mathbf{x}$  as **anomaly** if  $O_{SD}(\mathbf{x}|\mathbf{X}) > \max(O_{SD})$ .

## Projection depth (Zuo & Serfling, 2000)

According to Zuo & Serfling (2000), **projection depth** is defined as:

$$D^{prj}(\mathbf{x}|X) = \frac{1}{1 + O_{SD}(\mathbf{x}|X)},$$

## Projection depth (Zuo & Serfling, 2000)

According to Zuo & Serfling (2000), **projection depth** is defined as:

$$D^{prj}(\mathbf{x}|X) = \frac{1}{1 + O_{SD}(\mathbf{x}|X)},$$

where

$$O_{SD}(\mathbf{x}|X) = \sup_{\mathbf{r} \in S^{d-1}} \frac{|X^T \mathbf{r} - \text{med}(X^T \mathbf{r})|}{\text{MAD}(X^T \mathbf{r})}$$

is the **projected outlyingness** (Stahel, 1981; Donoho, 1982),  $\text{med}(Y)$  and  $\text{MAD}(Y) = \text{med}(|Y - \text{med}(Y)|)$  are the univariate median and median absolute deviation from the median, respectively.

## Projection depth (Zuo & Serfling, 2000)

According to Zuo & Serfling (2000), **projection depth** is defined as:

$$D^{prj}(\mathbf{x}|X) = \frac{1}{1 + O_{SD}(\mathbf{x}|X)},$$

where

$$O_{SD}(\mathbf{x}|X) = \sup_{\mathbf{r} \in S^{d-1}} \frac{|X^T \mathbf{r} - \text{med}(X^T \mathbf{r})|}{\text{MAD}(X^T \mathbf{r})}$$

is the **projected outlyingness** (Stahel, 1981; Donoho, 1982),  $\text{med}(Y)$  and  $\text{MAD}(Y) = \text{med}(|Y - \text{med}(Y)|)$  are the univariate median and median absolute deviation from the median, respectively.

Properties:

- ▶ Satisfies **D1 – D5** and **D4con**, is continuous;

## Projection depth (Zuo & Serfling, 2000)

According to Zuo & Serfling (2000), **projection depth** is defined as:

$$D^{prj}(\mathbf{x}|X) = \frac{1}{1 + O_{SD}(\mathbf{x}|X)},$$

where

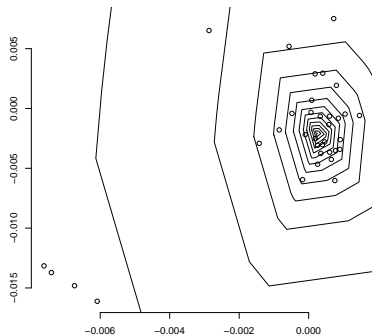
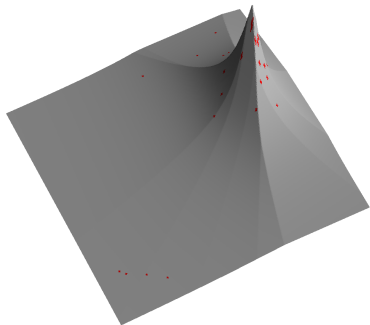
$$O_{SD}(\mathbf{x}|X) = \sup_{\mathbf{r} \in S^{d-1}} \frac{|X^T \mathbf{r} - \text{med}(X^T \mathbf{r})|}{\text{MAD}(X^T \mathbf{r})}$$

is the **projected outlyingness** (Stahel, 1981; Donoho, 1982),  $\text{med}(Y)$  and  $\text{MAD}(Y) = \text{med}(|Y - \text{med}(Y)|)$  are the univariate median and median absolute deviation from the median, respectively.

Properties:

- ▶ Satisfies **D1 – D5** and **D4con**, is continuous;
- ▶ its **median** has asymptotic breakdown point of 0.5.

# Projection depth (Zuo & Serfling, 2000)





## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

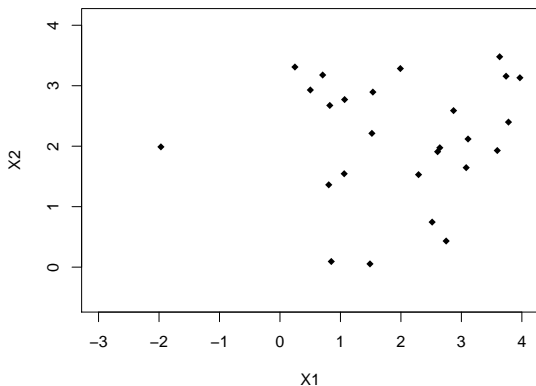
Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$

## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

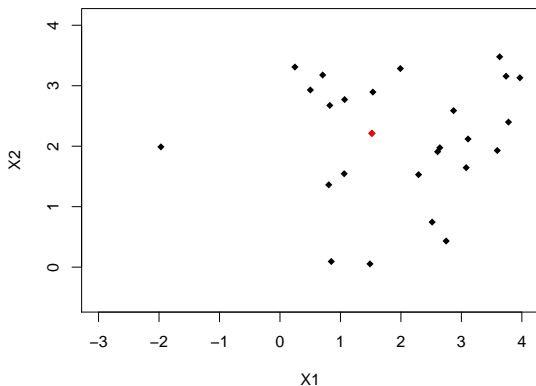
$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$



## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

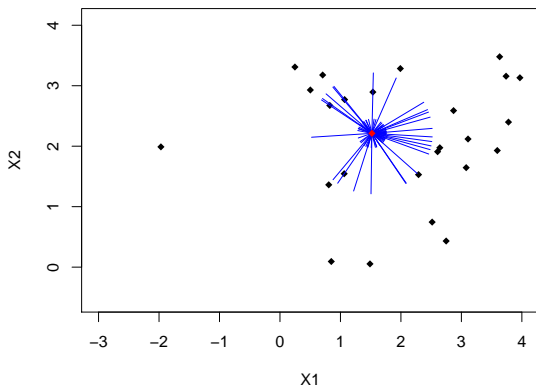
$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$



## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

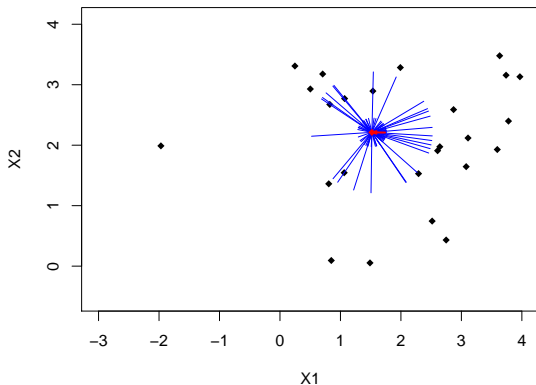
$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$



## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

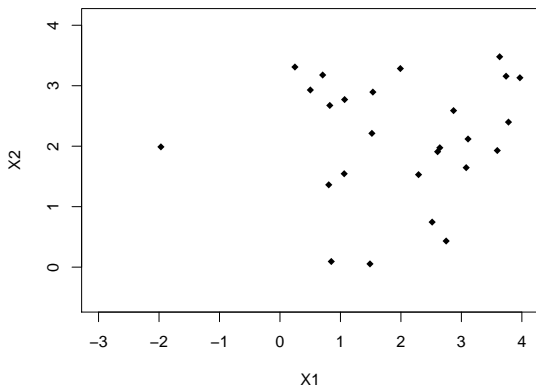
$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$



## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

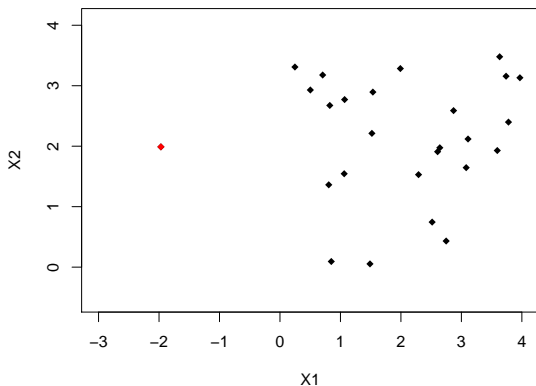
$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$



## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

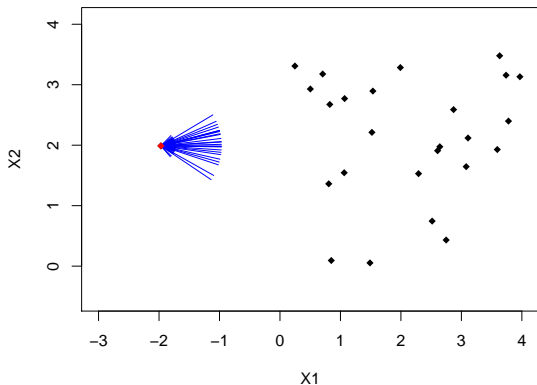
$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$



## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$

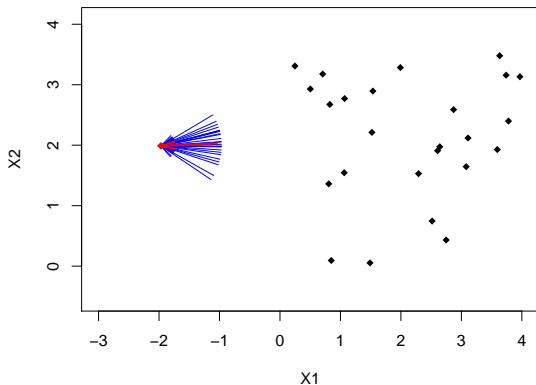




## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} \right] \right\| \quad \text{with} \quad \frac{\mathbf{x} - X}{\|\mathbf{x} - X\|} = 0 \quad \text{if} \quad \mathbf{x} - X = \mathbf{0}.$$



## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ v(\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{x} - X)) \right] \right\|,$$

with

$$v(\mathbf{y}) = \begin{cases} \frac{\mathbf{y}}{\|\mathbf{y}\|} & \text{if } \mathbf{y} \neq \mathbf{0}, \\ \mathbf{0} & \text{if } \mathbf{y} = \mathbf{0}. \end{cases}$$

## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ v(\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{x} - X)) \right] \right\|,$$

with

$$v(\mathbf{y}) = \begin{cases} \frac{\mathbf{y}}{\|\mathbf{y}\|} & \text{if } \mathbf{y} \neq \mathbf{0}, \\ \mathbf{0} & \text{if } \mathbf{y} = \mathbf{0}. \end{cases}$$

Properties:

- ▶ satisfies **D1** – **D5**, but not **D4con**, is continuous;

## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ v(\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{x} - X)) \right] \right\|,$$

with

$$v(\mathbf{y}) = \begin{cases} \frac{\mathbf{y}}{\|\mathbf{y}\|} & \text{if } \mathbf{y} \neq \mathbf{0}, \\ \mathbf{0} & \text{if } \mathbf{y} = \mathbf{0}. \end{cases}$$

Properties:

- ▶ satisfies **D1** – **D5**, but not **D4con**, is continuous;
- ▶ if  $\boldsymbol{\Sigma}$  is orthogonal, satisfies **D2iso** only;

## Spatial depth (Vardi & Zhang, 2000; Serfling 2002)

Exploiting the idea of spatial quantiles of Chaudhuri (1996) and Koltchinskii (1997), Vardi & Zhang (2000) and Serfling (2002) formulate the **spatial depth** (also  $L_1$ -depth) as:

$$D^{spt}(\mathbf{x}|X) = 1 - \left\| \mathbb{E} \left[ v(\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{x} - X)) \right] \right\|,$$

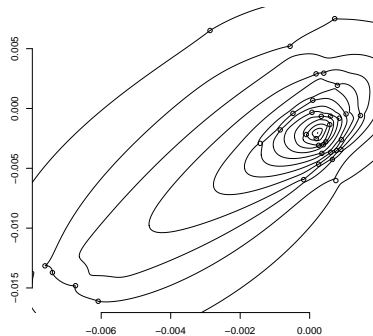
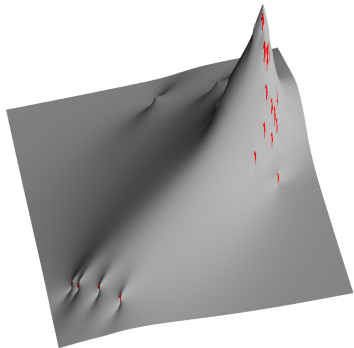
with

$$v(\mathbf{y}) = \begin{cases} \frac{\mathbf{y}}{\|\mathbf{y}\|} & \text{if } \mathbf{y} \neq \mathbf{0}, \\ \mathbf{0} & \text{if } \mathbf{y} = \mathbf{0}. \end{cases}$$

Properties:

- ▶ satisfies **D1** – **D5**, but not **D4con**, is continuous;
- ▶ if  $\boldsymbol{\Sigma}$  is orthogonal, satisfies **D2iso** only;
- ▶ with **D2iso** its maximum (say  $\mathbf{x}^*$ ) is referred to as **spatial median**, a multivariate location estimator having asymptotic breakdown point of 0.5.

# Spatial depth (Vardi & Zhang, 2000; Serfling 2002)



# Contents

## Introduction

## Non-parametric approaches

- One-class support vector machines

- Local outlier factor

- Isolation forest

## Systematic orderings: data depth

- The notion of data depth

- The Tukey depth function

- Central regions

- Further depth notions

## Practical session

# Thank you for attention! (and a short list of literature)

- ▶ Chandola, V., Banerjee, A., and Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(3):15, 1–58.
- ▶ Breunig, M.M., Kriegel, H.-P., Ng, R.T., and Sander, J. (2000). LOF: Identifying density-based local outliers. In: *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, 29, 93–104.
- ▶ Schölkopf, B., Platt, J.C., Shawe-Taylor, J., Smola, A., and Williamson, R. (2001). Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7), 1443–1471.
- ▶ Liu, F.T., Ting, K.M., and Zhou, Z. (2008). Isolation forest. In: *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, 413–422.
- ▶ Mosler, K. (2013). Depth statistics. In: *Robustness and Complex Data Structures: Festschrift in Honour of Ursula Gather*, 17–34.



# Practical session (part I)

## Notebooks:

- ▶ `anomdet_simulation1.Rmd`,
- ▶ `anomdet_hurricanes.Rmd`,
- ▶ `anomdet_cars.ipynb`,
- ▶ `anomdet_airbus.ipynb`.

## Data sets:

- ▶ `carsanom.csv`: Data set on anomaly detection for cars.
- ▶ `airbus_data.csv`: Data set from Airbus.
- ▶ `hurdat2-1851-2019-052520.txt`: Historical hurricane data.

## Supplementary scripts:

- ▶ `depth_routines.py`: Routines for data depth calculation.
- ▶ `FIF.py`: Implementation of the functional isolation forest.
- ▶ `depth_routines.R`: Routines for curves' parametrization.

## Literature (mentioned in the tutorial) (1)

- ▶ Boser, B.E., Guyon, I., and Vapnik, V.N. (1992). A training algorithm for optimal margin classifiers. In: *Proceedings of the Fifth Annual Workshop of Computational Learning Theory*, Pittsburgh, ACM, 5, 144–152.
- ▶ Breunig, M.M., Kriegel, H.-P., Ng, R.T., and Sander, J. (2000). LOF: Identifying density-based local outliers. In: *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, 29, 93–104.
- ▶ Chandola, V., Banerjee, A., and Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(3):15, 1–58.
- ▶ Chaudhuri P. (1996). On a geometric notion of quantiles for multivariate data. *Journal of the American Statistical Association*, 91, 862–872.
- ▶ Claeskens, G., Hubert, M., Slaets, L., and Vakili, K. (2014). Multivariate functional halfspace depth. *Journal of the American Statistical Association*, 109(505), 411—423.
- ▶ Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20, 273–297.
- ▶ Donoho D. (1982). *Breakdown Properties of Multivariate Location Estimators*. Ph.D. thesis, Harvard University.
- ▶ Donoho D.L., Gasko M. (1992). Breakdown properties of location estimates based on halfspace depth and projected outlyingness. *The Annals of Statistics*, 20, 1803–1827.

## Literature (mentioned in the tutorial) (2)

- ▶ Fraiman, R. and Muniz, G. (2001). Trimmed means for functional data. *TEST*, 10, 419—440.
- ▶ Hariri, S., Carrasco Kind, M., and Brunner, R.J. (2018). Extended isolation forest. [arXiv:1811.02141](https://arxiv.org/abs/1811.02141).
- ▶ Hubert, M., Rousseeuw, P.J., and Segaert, P. (2015). Multivariate functional outlier detection. *Statistical Methods & Applications*, 24(2), 177—202.
- ▶ Koltchinskii V. (1997). M-estimation, convexity and quantiles. *The Annals of Statistics*, 25, 435–477.
- ▶ Koshevoy G., Mosler K. (1997). Zonoid trimming for multivariate distributions. *The Annals of Statistics*, 25, 1998–2017.
- ▶ Liu R.Y. (1990). On a notion of data depth based on random simplices. *The Annals of Statistics*, 18, 405–414.
- ▶ Liu, Z. and Modarres, R. (2011). Lens data depth and median. *Journal of Nonparametric Statistics*, 23, 1063–1074.
- ▶ Liu, F.T., Ting, K.M., and Zhou, Z. (2008). Isolation forest. In: *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, 413–422.

## Literature (mentioned in the tutorial) (3)

- ▶ López-Pintado, S. and Romo, J. (2009). On the concept of depth for functional data. *Journal of the American Statistical Association*, 104(486), 718–734.
- ▶ Mahalanobis P.C. (1936). On the generalized distance in statistics. *Proceedings of the National Institute of Sciences of India*, 12, 49–55.
- ▶ Markou, M. and Singh, S. (2003). Novelty detection: a review - Part 1: Statistical approaches. *Signal Processing*, 83(12), 2481–2497.
- ▶ Markou, M. and Singh, S. (2003). Novelty detection: a review - Part 2: Neural network based approaches. *Signal Processing*, 83(12), 2499–2521.
- ▶ Miljković, D. (2010). Review of novelty detection methods. *The 33rd International Convention MIPRO*, Opatija, 593–598.
- ▶ Mosler, K. (2013). Depth statistics. In: *Robustness and Complex Data Structures: Festschrift in Honour of Ursula Gather*, 17–34.
- ▶ Oja, H. (1983). Descriptive statistics for multivariate distributions. *Statistics and Probability Letters*, 1, 327–332.
- ▶ Pimentel, M.A.F., Clifton, D.A., Clifton, L., and Tarassenko, L. (2014). A review of novelty detection. *Signal Processing*, 99, 215–249.
- ▶ Schölkopf, B., Platt, J.C., Shawe-Taylor, J., Smola, A., and Williamson, R. (2001). Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7), 1443–1471.

## Literature (mentioned in the tutorial) (4)

- ▶ Serfling, R. (2002). A depth function and a scale curve based on spatial quantiles. In: *Statistical Data Analysis Based on the  $L_1$ -Norm and Related Methods* Birkhäuser, Basel, 25—38.
- ▶ Stahel W. (1981). *Robust Estimation: Infinitesimal Optimality and Covariance Matrix Estimators (In German)*. Ph.D. thesis, Swiss Federal Institute of Technology in Zurich.
- ▶ Tukey J.W. (1975). Mathematics and the picturing of data. In: *Proceedings of the International Congress of Mathematicians*, volume 2, Canadian Mathematical Congress, 523–531.
- ▶ Vapnik, V. and Chervonenkis, A. (1974). *Theory of Pattern Recognition* (in Russian). Nauka, Moscow.
- ▶ Vapnik, V. and Lerner, A. (1963). Pattern recognition using generalized portraits. *Avtomatika i Telemekhanika*, 24, 774–780.
- ▶ Vardi Y., Zhang C. (2000). The multivariate  $L_1$ -median and associated data depth. *Proceedings of the National Academy of Sciences of the United States of America*, 97, 1423–1426.
- ▶ Zuo Y., Serfling R. (2000). General notions of statistical depth function. *The Annals of Statistics*, 28, 461–482.