

# Sparsity 101: Statistical estimators

## Central location and dispersion

Laurent Duval

IFP Energies nouvelles

2013

# Introduction

What is the trend? Where is the outlier?

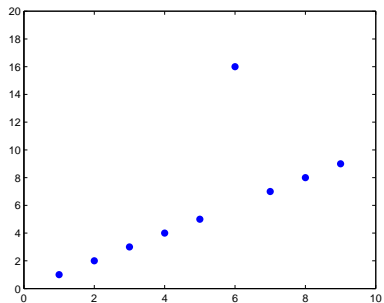


Figure : Toy noise problem

# Introduction

What is the trend? Where is the outlier?

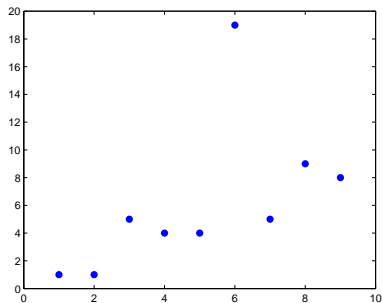


Figure : Toy noise problem

# Introduction

What is the trend? Where is the outlier?

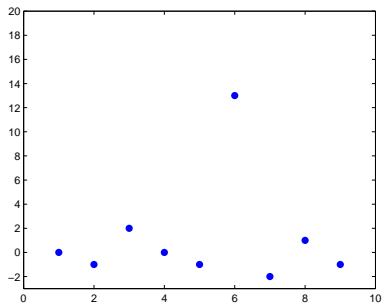


Figure : Toy noise problem

# Introduction

What is the trend? Where is the outlier?

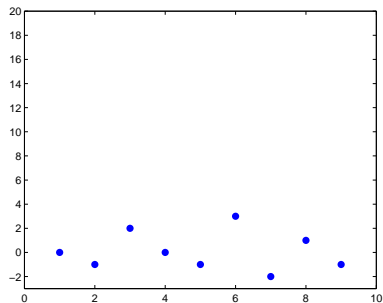


Figure : Toy noise problem

# Introduction

What is the trend? Where is the outlier?

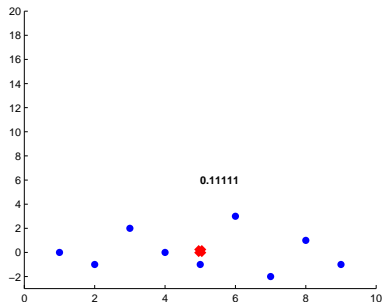


Figure : Toy noise problem

# Introduction

What is the trend? Where is the outlier?

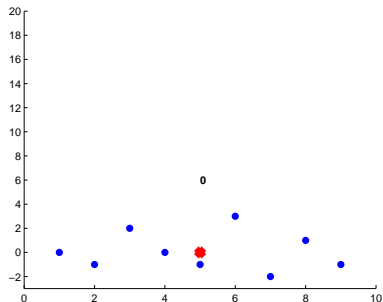


Figure : Toy noise problem

# Introduction

What is the trend? Where is the outlier?

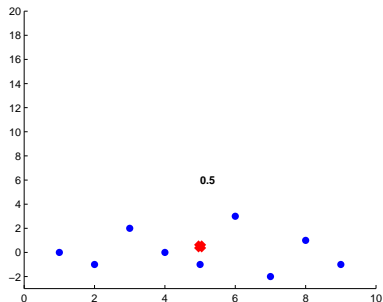


Figure : Toy noise problem



# Introduction

What is the trend? Where is the outlier?

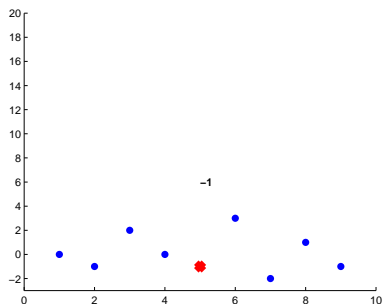


Figure : Toy noise problem

# Standard estimators

Compute a representative **central location  $\bar{m}$**  ?

► mean:  $\bar{m} = \frac{1}{N} \sum x_i$

## Standard estimators

Compute a representative **central location  $\bar{m}$**  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$  ( $\frac{N+1}{2}$  position)

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$

## Standard estimators

Compute a representative **central location  $\bar{m}$**  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$

## Standard estimators

Compute a representative **central location  $\bar{m}$**  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$

## Standard estimators

Compute a representative **central location  $\bar{m}$**  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$  (and anything else in between?)

## Standard estimators

Compute a representative **central location  $\bar{m}$**  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?



## Standard estimators

Compute a representative **central location  $\bar{m}$**  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \leftarrow \arg \min \sum (x_i - m)^2$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \leftarrow \arg \min \sum w_i (x_i - m)^2$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \leftarrow \arg \min \sum (x_i - m)^2$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \leftarrow \arg \min \sum |x_i - m|$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \leftarrow \arg \min \sum (x_i - m)^2$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \leftarrow \arg \min \sum w_i |x_i - m|$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i)$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \leftarrow \arg \min \sum (x_i - m)^2$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \leftarrow \arg \min \sum |x_i - m|$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i) \leftarrow \arg \min \max |x_i - m|$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i))$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \leftarrow \arg \min \sum (x_i - m)^2$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \leftarrow \arg \min \sum |x_i - m|$
- ▶ *mid-range*:  $\frac{1}{2}(\min x_i + \max x_i) \leftarrow \arg \min \max |x_i - m|$
- ▶ *mid-hinge*:  $\frac{1}{2}(Q1(x_i) + Q3(x_i)) \leftarrow \arg \min \text{med}|x_i - m|$
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**

## Standard estimators

Compute a representative **central location**  $\bar{m}$  ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*
- ▶ *mid-hinge*
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**
- ▶ with a **natural spread**



## Standard estimators

Compute a representative **dispersion/spread** ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*
- ▶ *mid-hinge*
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**
- ▶ with a **natural spread**

## Standard estimators

Compute a representative **dispersion/spread** ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \rightarrow \sigma = \sqrt{\frac{1}{N} \sum (x_i - \bar{m})^2}$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m})$
- ▶ *mid-range*
- ▶ *mid-hinge*
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**
- ▶ with a **natural spread**

## Standard estimators

Compute a representative **dispersion/spread** ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \rightarrow \sigma = \sqrt{\frac{1}{N} \sum (x_i - \bar{m})^2}$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \rightarrow \text{MAD} = \text{med}|x_i - \bar{m}|$
- ▶ *mid-range*
- ▶ *mid-hinge*
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**
- ▶ with a **natural spread**

## Standard estimators

Compute a representative **dispersion/spread** ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \rightarrow \sigma = \sqrt{\frac{1}{N} \sum (x_i - \bar{m})^2}$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \rightarrow \text{MAD} = \text{med}|x_i - \bar{m}|$
- ▶ *mid-range*  $\rightarrow \frac{1}{2}(\max x_i - \min x_i)$  (half-range, mid-span)
- ▶ *mid-hinge*
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**
- ▶ with a **natural spread**

## Standard estimators

Compute a representative **dispersion/spread** ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \rightarrow \sigma = \sqrt{\frac{1}{N} \sum (x_i - m)^2}$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \rightarrow \text{MAD} = \text{med}|x_i - m|$
- ▶ *mid-range*  $\rightarrow \frac{1}{2}(\max x_i - \min x_i)$  (half-range, mid-span)
- ▶ *mid-hinge*  $\rightarrow \frac{1}{2}(Q3(x_i) - Q1(x_i))$  (IQR/mid-spread)
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**
- ▶ with a **natural spread**

## Standard estimators

Compute a representative **dispersion/spread** ?

- ▶ mean:  $\bar{m} = \frac{1}{N} \sum x_i \rightarrow \sigma = \sqrt{\frac{1}{N} \sum (x_i - m)^2}$
- ▶ median:  $\#(x_i < \bar{m}) = \#(x_i > \bar{m}) \rightarrow \text{MAD} = \text{med}|x_i - m|$
- ▶ *mid-range*  $\rightarrow \frac{1}{2}(\max x_i - \min x_i)$  (half-range, mid-span)
- ▶ *mid-hinge*  $\rightarrow \frac{1}{2}(Q3(x_i) - Q1(x_i))$  (IQR/mid-spread)
- ▶ mode:  $\arg \max p(x)$
- ▶ even:  $\min x_i, \max x_i$
- ▶ arbitrary algorithmic choice?
- ▶ no, answer to an **optimization problem**
- ▶ with a **natural spread**

**Bonuses:** add shape factor (skewness, kurtosis);  $\frac{Q1+2 \times \text{med}+Q3}{4} \dots$

## Standard estimators: $\ell_2$

Least-squares estimator of weighted central location:

$$f(m) = \sum w_i (x_i - m)^2$$

$$\frac{df}{dm} = \sum -2w_i (x_i - m)$$

$$\frac{df}{dm} = 0 \Leftrightarrow \sum -2w_i x_i = \sum -2w_i m$$

$$\sum w_i x_i = \sum w_i m$$

$$\bar{m} = \frac{\sum w_i x_i}{\sum w_i}$$

Weighted sum location  $\rightarrow$  all linear filters; Gaussian noise

## Standard estimators: $\ell_1$

Least-magnitude estimator of (weighted) central location:

$$f(m) = \sum w_i |x_i - m|$$

$$\frac{df}{dm} \approx \sum -w_i \text{sign}(x_i - m)$$

$$\frac{df}{dm} = 0 \Leftrightarrow \#_{w_i}(x_i < \bar{m}) = \#_{w_i}(x_i > \bar{m})$$

Equality reached when  $\bar{m}$  stands “in between”

$$\bar{m} = \text{median}_{w_i} x_i$$

Sorted location  $\rightarrow$  gen. weighted median filters; Laplace noise



## Other (less standard) estimators

### Importance of measurement units

- ▶ *harmonic*, geometric, arithmetico-geometric means  
100 km at 150 km h<sup>-1</sup>, 100 km at 100 km h<sup>-1</sup> → 120 km h<sup>-1</sup>
- ▶ M-estimators, L-estimators, robust statistics
- ▶ no natural dispersion in general
- ▶ time vs individuals; 1D/2D; representative scale; transforms

### A robust Gaussian noise dispersion estimator (details)

$$\hat{\sigma} \simeq \frac{\text{median}|c_j|}{0.6745}$$

Use with wavelet shrinkage operators (soft , hard, garrote, etc.)

$$\hat{c}_j = S(c_j, \Lambda(\sigma))$$

## Noise level estimation and wavelets



Figure : Noise estimation with standard 1D wavelets

# Noise level estimation and wavelets

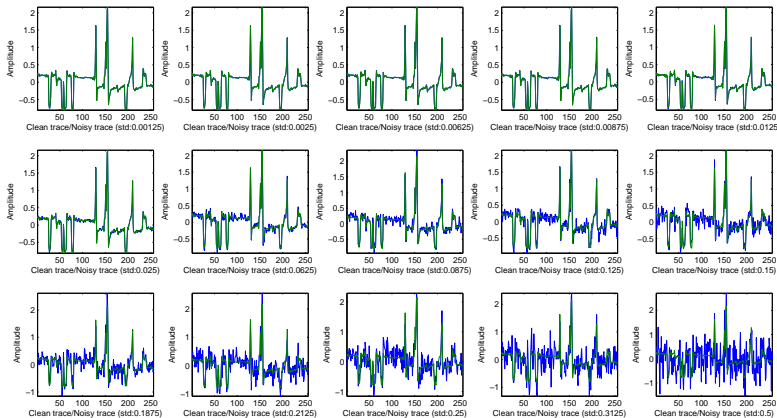


Figure : Noise estimation with standard 1D wavelets

# Noise level estimation and wavelets

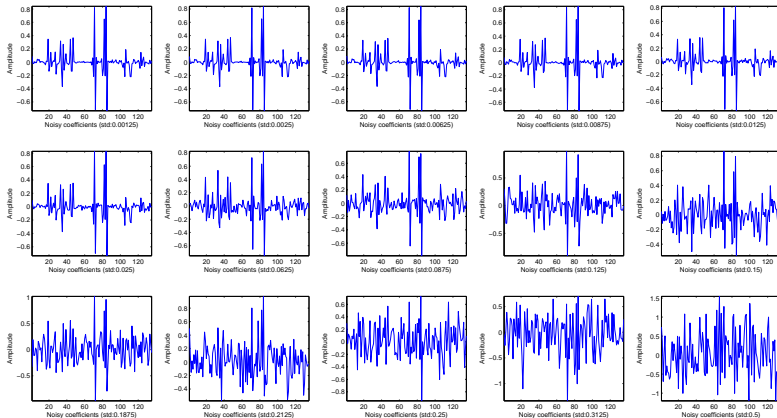


Figure : Noise estimation with standard 1D wavelets

## Noise level estimation and wavelets

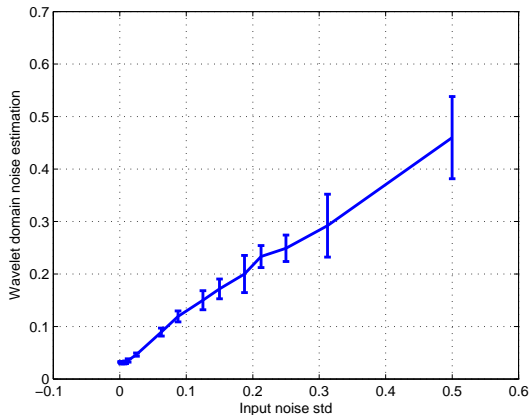


Figure : Noise estimation with standard 1D wavelets