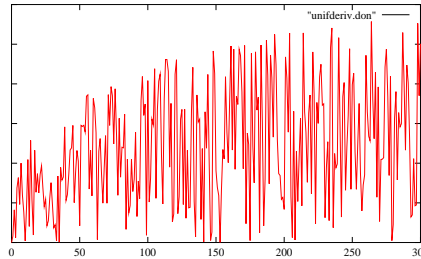# Sample Quality
**descriptive analysis of data**

Lucas Mello Schnorr, Jean-Marc Vincent
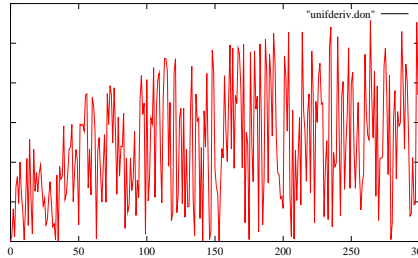
INF/UFRGS
Porto Alegre, Brazil – October 2018

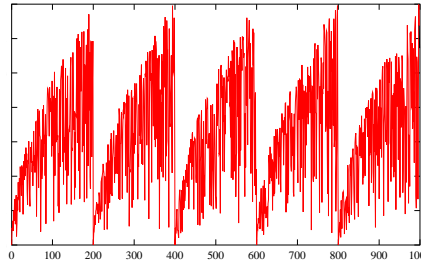# CONTROL OF EXPERIMENTS (1)

# CONTROL OF EXPERIMENTS (1)



**Tendency analysis**

**non homogeneous experiment**
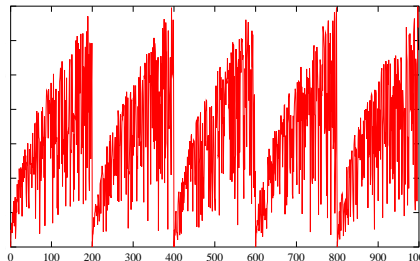$\Rightarrow$ model the evolution of experiment
estimate and compensate tendency
**explain why**

# CONTROL OF EXPERIMENTS (2)

# CONTROL OF EXPERIMENTS (2)



**Periodicity analysis**

**periodic evolution of the experimental environment ?**
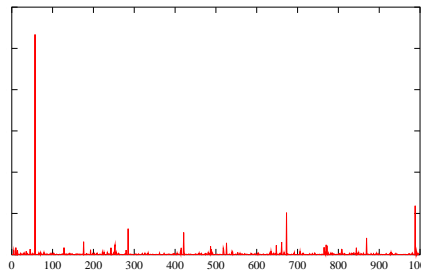$\Rightarrow$ model the evolution of experiment
Fourier analysis of the sample
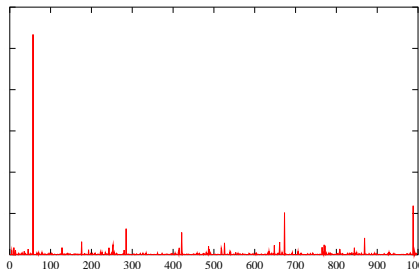Integration on time (sliding window analysis) Danger : size of the window
Wavelet analysis
**explain why**

UNIVERSITÉ
Grenoble
Alpes

# CONTROL OF EXPERIMENTS (3)

# CONTROL OF EXPERIMENTS (3)



**Non significant values**

**extraordinary behaviour of experimental environment**
rare events with different orders of magnitude
$\Rightarrow$ threshold by value
Danger : choice of the threshold : indicate the rejection rate
$\Rightarrow$ threshold by quantile
Danger : choice of the percentage : indicate the rejection value
**explain why**

UNIVERSITÉ
Grenoble
Alpes
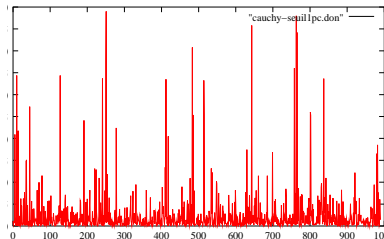
# CONTROL OF EXPERIMENTS (4)

Threshold value : 10



Threshold percentage : 1%

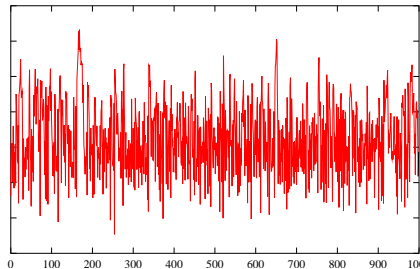# CONTROL OF EXPERIMENTS (5)

# CONTROL OF EXPERIMENTS (5)



**looks like correct experiments**
Statistically independent
Statistically homogeneous

## CONTROL OF EXPERIMENTS (5BIS)

Zooming

## CONTROL OF EXPERIMENTS (5BIS)

Zooming



**Autocorrelation**

Danger time correlation among samples
**experiments impact on experiments**
$\Rightarrow$ stationarity analysis
autocorrelation estimation (ARMA)

# EXPERIMENTAL RESULTS

- ► Deterministic (controlled error non significant (white noise))
- ► Statistic (the system is non deterministic)

## Sample analysis

- ► Identification of the response set
- ► Structure of the response set (measure)

UNIVERSITÉ
Grenoble
Alpes

# DISTRIBUTION ANALYSIS

Summarize data in a **histogram**



**Shape analysis**

- unimodal / multimodal
- variability
- symmetric / dissymmetric (skewness)
- flatness (kurtosis)

$\implies$ **Central tendency analysis**
$\implies$ **Variability analysis around the central tendency**

UNIVERSITÉ
Grenoble
Alpes

# MODE VALUE



## Mode

► **Categorical data**

► Most frequent value

► highly unstable value

► for continuous value distribution depends on the histogram step

► interpretation depends on the flatness of the histogram

$\Longrightarrow$ **Use it carefully**
$\Longrightarrow$ **Predictor function**

UNIVERSITÉ
Grenoble
Alpes

# MEDIAN VALUE

**Median**

- ► **Ordered data**
- ► Split the sample in two equal parts

$$\sum_{i \leqslant Median} f_i \leqslant \frac{1}{2} \leqslant \sum_{i \leqslant Median+1} f_i.$$

- ► more stable value
- ► does not depends on the histogram step
- ► difficult to combine (two samples)

$\Longrightarrow$ **Randomized algorithms**

UNIVERSITÉ
Grenoble
Alpes

# MEAN VALUE

**Mean**

- ▶ **Vector space**
- ▶ Average of values

$$Mean = \frac{1}{Sample\_Size} \sum x_i = \sum_x x.f_x.$$

- ▶ stable value
- ▶ does not depends on the histogram step
- ▶ easy to combine (two samples $\Rightarrow$ weighted mean)

$\Longrightarrow$ **Additive problems (cost, durations, length,...)**

UNIVERSITÉ
Grenoble
Alpes

# CENTRAL TENDENCY



**Complementarity**

- Valid if the sample is "Well-formed"
- **Semantic of the observation**
- Goal of analysis

$\Longrightarrow$ **Additive problems (cost, durations, length,...)**

UNIVERSITÉ
Grenoble
Alpes

# CENTRAL TENDENCY (2)

**Summary of Means**

- ▶ Avoid means if possible
  Loses information

- ▶ Arithmetic mean
  When sum of raw values has physical meaning
  Use for summarizing times (not rates)

- ▶ Harmonic mean
  Use for summarizing rates (not times)

- ▶ Geometric mean
  Not useful when time is best measure of perf
  Useful when multiplicative effects are in play

UNIVERSITÉ
Grenoble
Alpes

# VARIABILITY

**Categorical data (finite set)**

$f_i$ : empirical frequency of element $i$
Empirical entropy

$$H(f) = -\sum_i f_i \log f_i.$$

Measure the empirical distance with the uniform distribution

- $H(f) \geqslant 0$
- $H(f) = 0$ iff the observations are reduced to a unique value
- $H(f)$ is maximal for the uniform distribution

UNIVERSITÉ
Grenoble
Alpes

# **VARIABILITY (2)**

**Ordered data**

Quantiles : quartiles, deciles, etc
Sort the sample :

$$(x_1, x_2, \cdots, x_n) \longrightarrow (x_{(1)}, x_{(2)}, \cdots, x_{(n)});$$

$$Q_1 = x_{(n/4)}; \ Q_2 = x_{(n/2)} = \text{Median}; \ Q_3 = x_{(3n/4)}.$$

For deciles

$$d_i = \text{argmax}_i \{\sum_{j \leqslant i} f_j \leqslant \frac{i}{10}\}.$$

Utilization as quantile/quantile plots to compare distributions

UNIVERSITÉ
Grenoble
Alpes

# VARIABILITY (3)

**Vectorial data**

Quadratic error for the mean

$$Var(X) = \frac{1}{n} \sum_{1}^{n} (x_i - \bar{x}_n)^2.$$

**Properties :**

$$
\begin{aligned}
Var(X) &\geqslant 0; \\
Var(X) &= \overline{x^2} - (\bar{x})^2, \ \text{ where } \ \overline{x^2} = \frac{1}{n} \sum_{i=1}^{n} x_i^2. \\
Var(X + cste) &= Var(X); \\
Var(\lambda X) &= \lambda^2 Var(X).
\end{aligned}
$$

UNIVERSITÉ
Grenoble
Alpes