

Performance evaluation

How to deal with experiments

Jean-Marc Vincent

MESCAL-INRIA Project
Laboratoire d'Informatique de Grenoble
Universities of Grenoble, France
{Jean-Marc.Vincent}@imag.fr

This work was partially supported by CAPES/COFECUB and STIC-AMSUD



Outline

- 1 Typical software engineering problem
- 2 Analysis of experimental data
- 3 Estimation
- 4 References

Introduction

Aim of this course

Give basic concepts of experimental design

- Statistical analysis of experimental data
- Modelling and parameters estimation
- Measurement in distributed systems and large scale systems
- Design of experiments

Interactive course : discussion about your own experiments

Outline

- 1 Typical software engineering problem
- 2 Analysis of experimental data
- 3 Estimation
- 4 References

Application problems

Procedure TRC

```
int TRC(int *T,int p, int r)
{
int cpt=0;
int i, j,x;
if (r-p>0)
{
i=p; j=r;
while (j-i>0)
{
if (T[i] >= T[i+1]) { x=T[i]; T[i]=T[i+1]; T[i+1]=x; i++;}
else {x=T[j]; T[j]=T[i+1]; T[i+1]=x; j-;}
cpt++;
}
cpt+=TRC(T,p,i-1); cpt+=TRC(T,i+1,r);
}
return cpt;
}
```

First step

Specification

- The procedure $\text{TRC}(T,i,j)$ sorts in place elements of an array T from index i to j included;

How to check it ?

Methods

- Exhaustive checking : enumerate all arrays and check them one by one
- Subset checking : use a representative subset of all arrays
- Statistical testing : generate uniformly an arbitrary subset size of array (confidence)

Example

Second step

Specification

- The procedure $\text{TRC}(T,i,j)$ sorts in place elements of an array T from index i to j included;
- The procedure $\text{TRC}(T,i,j)$ sorts n elements in $\mathcal{O}(n \log n)$ comparisons

How to evaluate it ?

Methods

- Exhaustive evaluation : enumerate all arrays and compute cost for each one
- Statistical evaluation : generate uniformly an arbitrary subset size of array (confidence) compute the empirical distribution of the cost, test the model $n \log n$

Example



Third step

Specification

- The procedure $\text{TRC}(T,i,j)$ sorts in place elements of an array T from index i to j included;
- The procedure $\text{TRC}(T,i,j)$ sorts n elements in $\mathcal{O}(n \log n)$ comparisons
- The procedure $\text{TRC}(T,i,j)$ sorts n elements efficiently

How to measure it ?

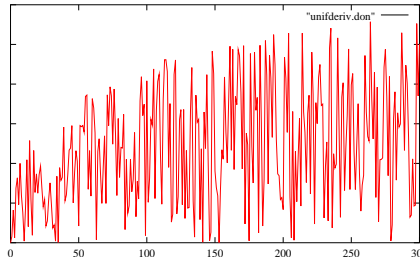
Methods

- Put probes inside the program/system/processor...
- Performance evaluation : generate uniformly an arbitrary subset size of array (confidence) measure the execution time, compute the empirical distribution of the execution time, fit with some models

Example



Sample analysis



Tendency analysis

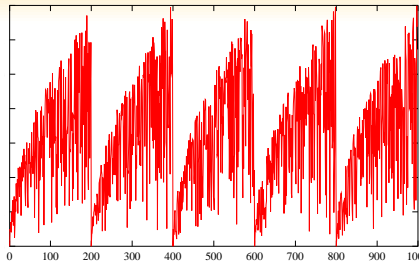
non homogeneous experiment

⇒ model the evolution of experiment

estimate and compensate tendency

explain why

Sample analysis (2)



Periodicity analysis

periodic evolution of the experimental environment ?

⇒ model the evolution of experiment

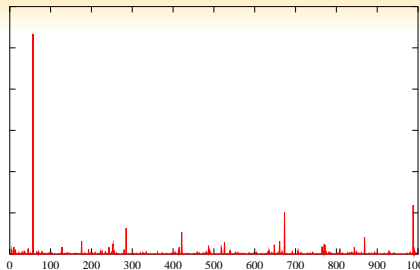
Fourier analysis of the sample

Integration on time (sliding window analysis) Danger : size of the window

Wavelet analysis

explain why

Sample analysis (3)



Non significant values

extraordinary behaviour of experimental environment

rare events with different orders of magnitude

⇒ threshold by value

Danger : choice of the threshold : indicate the rejection rate

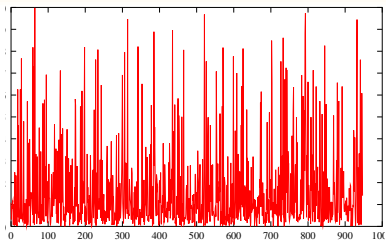
⇒ threshold by quantile

Danger : choice of the percentage : indicate the rejection value

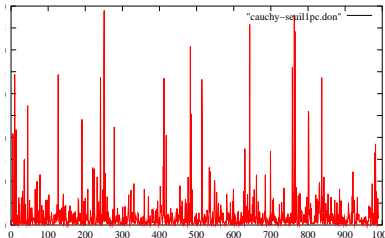
explain why

Sample analysis (4)

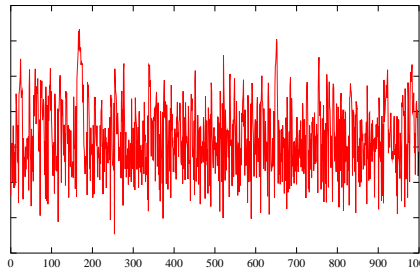
Threshold value : 10



Threshold percentage : 1%



Sample analysis (5)



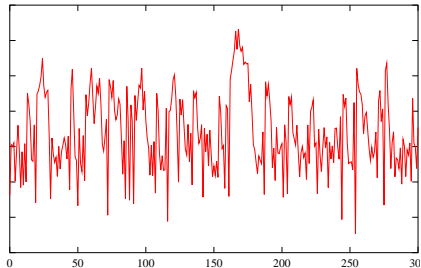
looks like correct experiments

Statistically independent

Statistically homogeneous

Sample analysis (5bis)

Zooming



Autocorrelation

Danger time correlation among samples

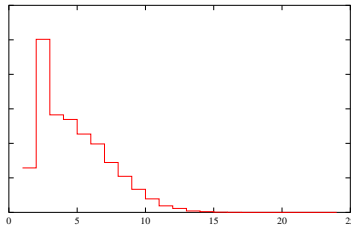
experiments impact on experiments

⇒ stationarity analysis

autocorrelation estimation (ARMA)

Distribution analysis

Summarize data in a **histogram**



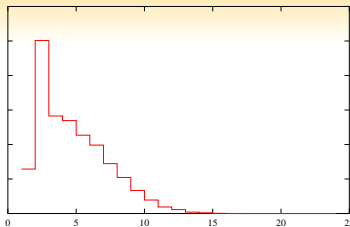
Shape analysis

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

⇒ **Central tendency analysis**

⇒ **Variability analysis around the central tendency**

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

⇒ **Use it carefully**

⇒ **Predictor function**

Median value

Median

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

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

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

⇒ **Randomized algorithms**

Mean value

Mean

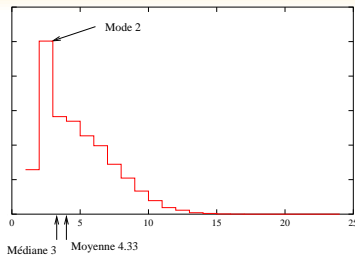
- **Vector space**
- Average of values

$$\text{Mean} = \frac{1}{\text{Sample_Size}} \sum x_i = \sum_x x \cdot f_x.$$

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

\Rightarrow **Additive problems (cost, durations, length,...)**

Central tendency



Complementarity

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

⇒ **Additive problems (cost, durations, length,...)**

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

Computational aspects

- Mode : computation of the histogram steps, then computation of max $O(n)$ “off-line”
- Median : sort the sample $O(n\log(n))$ or $O(n)$ (subtile algorithm) “off-line”
- Mean : sum values $O(n)$ “on-line” computation

Is the central tendency significant ?
⇒ Explain variability.

Computational aspects

- Mode : computation of the histogram steps, then computation of max $O(n)$ “off-line”
- Median : sort the sample $O(n \log(n))$ or $O(n)$ (subtile algorithm) “off-line”
- Mean : sum values $O(n)$ “on-line” computation

Is the central tendency significant ?
⇒ Explain variability.

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

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

Variability (2)

Ordered data

Quantiles : quartiles, deciles, etc

Sort the sample :

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

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

For deciles

$$d_i = \operatorname{argmax}_i \left\{ \sum_{j \leq i} f_j \leq \frac{i}{10} \right\}.$$

Utilization as quantile/quantile plots to compare distributions

Variability (3)

Vectorial data

Quadratic error for the mean

$$\text{Var}(X) = \frac{1}{n} \sum_1^n (x_i - \bar{x}_n)^2.$$

Properties:

$$\text{Var}(X) \geq 0;$$

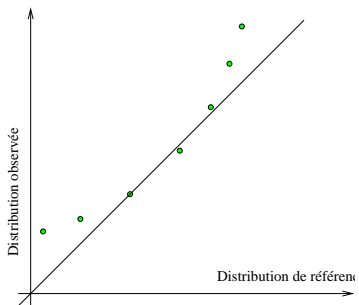
$$\text{Var}(X) = \overline{x^2} - (\bar{x})^2, \text{ where } \overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2.$$

$$\text{Var}(X + \text{cste}) = \text{Var}(X);$$

$$\text{Var}(\lambda X) = \lambda^2 \text{Var}(X).$$

Comparison with other distributions

Describe how an empirical distribution fits a theoretical one



Modelling

Basic assumptions :

- Data are considered as random variables
- Mutually independent
- Same probability distribution

Check Check Check

The distribution is given by

- Probability density function (pdf) (asymptotic histogram)

$$f_X(x) = \mathbb{P}(x \leq X \leq x + dx)/dx = F'_X(x).$$

- Cumulative distribution function

$$F_X(x) = \mathbb{P}(X \leq x);$$

- Moments : $M_n = \mathbb{E}X^n$, Variance

Average convergence

Law of large numbers

Let $\{X_n\}_{n \in \mathbb{N}}$ be a iid random sequence with finite variance, then

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{i=1}^n X_i = \mathbb{E}X, \quad \text{almost surely in } L^1.$$

- convergence of empirical frequencies
- for any experience we get the same result
- fundamental theorem of probability theory

$$\text{Notation : } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Law of errors

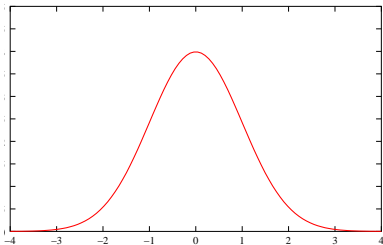
Central limit theorem (CLT)

Let $\{X_n\}_{n \in \mathbb{N}}$ be a iid random sequence with finite variance σ^2 , then

$$\lim_{n \rightarrow +\infty} \frac{\sqrt{n}}{\sigma} (\bar{X}_n - \mathbb{E}X) \stackrel{\mathcal{L}}{=} \mathcal{N}(0, 1).$$

→ error law (Gaussian law, Normal distribution, Bell curve,...)

→ Normalized mean = 0, variance = 1



Distribution

$$\mathbb{P}(X \in [-1, 1]) = 0.68;$$

$$\mathbb{P}(X \in [-2, 2]) = 0.95;$$

$$\mathbb{P}(X \in [-3, 3]) \geq 0.99.$$

Confidence intervals

Confidence level α compute ϕ_α

$$\mathbb{P}(X \in [-\phi_\alpha, \phi_\alpha]) = \alpha$$

For n sufficiently large ($n > 50$)

$$\mathbb{P}\left(\left[\bar{X}_n - \frac{\phi_\alpha \sigma}{\sqrt{n}}, \bar{X}_n + \frac{\phi_\alpha \sigma}{\sqrt{n}}\right] \ni \mathbb{E}X\right) = 1 - \alpha.$$



Confidence intervals (2)

Need an estimator of the variance

$$\hat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Danger n too small \rightarrow with a normal hypothesis take Student statistic
Three step method

- 1 In a first set of experiments check that the hypothesis is valid
- 2 Estimate roughly the variance
- 3 Estimate the mean and control the number of experiment by a confidence interval

Bibliography

- **The Art of Computer Systems Performance Analysis : Techniques for Experimental Design, Measurement, Simulation and Modeling.** Raj Jain *Wiley 1991*
- **Measuring Computer Performance: A Practitioner's Guide** David J. Lilja Cambridge University Press, 2000.