

## Protocol Validation

- typical random access protocol to a common channel (CSMA family)

Protocol- Send (M)

While Message is not sent

Send (N)

If Collision

W-Random (1/n)

$n = n + 1$

$I_{n+1} = g(n, I_n)$

Wait (W)

What should be the amount of time?

- Protocol dimensioning

Waiting time:

- random

- uniformly distributed on an interval  $[0, I_0]$

...

## Samples

The observation: a sequence of  $n$  waiting times

$$\{x_1, \dots, x_n\}$$

The stochastic model: the observations are considered to be realizations of independent random variables with the same probability law  $F$

$$\{X_1, \dots, X_n\}$$

Question: What could be said on  $F$  from the observations  $\{x_1, \dots, x_n\}$ ?

### ▷ A priori Knowledge:

- the shape of the law is known and depends on some parameters unknown: parametric estimation

$$F_0(x) = \begin{cases} 0 & x \leq 0 \\ \frac{x}{I_0} & 0 < x \leq I_0 \\ 1 & x > I_0 \end{cases}$$

- the shape of the law is unknown and some priors are under study: non-parametric estimation

## BASIC CONCEPTS

- A statistic is a function of observations :  $t_n(x_1, \dots, x_n)$ ; usually summarizes some param. of the distribution
- An estimator is a random variable  $T_n = t(X_1, \dots, X_n)$  (model of statistic)

Example:

$$t(x_1, \dots, x_n) = \max_{1 \leq i \leq n} x_i \quad \text{is a statistic on the samples}$$

and  $T_n = \max_{1 \leq i \leq n} X_i$  the corresponding estimator

Law of  $T_n$  under the hypothesis  $X_i$  uniformly distributed on  $[0, \theta]$ :

$$F_n^{\theta}(x) = P(\max_{1 \leq i \leq n} X_i \leq x) = \left(\frac{x}{\theta}\right)^n \quad \text{using independence and uniformity law of } X_i;$$

$$\text{and density } f_n^{\theta}(x) = \frac{1}{\theta^n} n \cdot x^{n-1}$$

## BIAS

- An estimator  $T_n$  of some parameter  $\theta$  is UNBIASED if  $E(T_n) = \theta$

Example :

$$E(T_n) = \int x b_n^{\theta}(x) dx = \dots = \frac{n}{n+1} \theta \neq \theta !$$

$T_n$  is biased estimator, on average it underestimate  $\theta$   
But, for large samples, the bias decreases to 0

$\lim_{n \rightarrow \infty} E(T_n) = \theta$ ,  $T_n$  is asymptotically unbiased

to compensate the bias,

$$T_n' = \frac{n+1}{n} T_n \quad \text{which is } \underline{\text{unbiased}}$$

## RISK

The quality  $T_n$  of an estimator is evaluated by RISK function:

$$R(T_n) = E(T_n - \theta)^2 = E(T_n - E(T_n))^2 + (E(T_n) - \theta)^2 = \text{Var} T_n + (\text{bias})^2$$

- $\text{Var} T_n$ : the concentration of the distribution
- $(E(T_n) - \theta)^2$ : impact of the bias

For an unbiased estimator  $R(T_n) = \text{Var}(T_n)$

$$\text{Var} T_n = \mathbb{E} (T_n - \mathbb{E} T_n)^2 = \mathbb{E} T_n^2 - (\mathbb{E} T_n)^2 = \int_0^\infty x^2 f_{T_n}(x) dx - \left(\frac{\theta}{n+2}\right)^2 = \left(\frac{n}{n+2} - \frac{n^2}{(n+1)^2}\right) \theta^2$$

$$\text{Var} T_n' = \text{Var} \frac{u^2}{n} T_n = \frac{(n+1)^2}{n^2} \text{Var} T_n = \frac{(n+1)^2}{n^2} \left(\frac{n}{n+2} - \frac{n^2}{(n+1)^2}\right) \theta^2 = \frac{1}{n(n+2)} \theta^2$$

Another estimator

$$U_n = \frac{2}{n} \sum_{i=1}^n X_i$$

$$\mathbb{E} U_n = \frac{2}{n} \sum_{i=1}^n \mathbb{E} X_i = \frac{2}{n} \sum_{i=1}^n \frac{\theta}{2} = \theta \quad \Rightarrow U_n \text{ is unbiased estimator of } \theta$$

Risk of  $U_n$ :

$$\text{Var} U_n = \text{Var} \left( \frac{2}{n} \sum_{i=1}^n X_i \right) = \left(\frac{2}{n}\right)^2 \sum \text{Var} X_i \quad \text{because of the independence of } X_i$$

$$= \frac{4}{n^2} \cdot \frac{\theta^2}{n^2} = \frac{\theta^2}{3n}$$

the risk is much larger than the risk of  $T_n'$   
 so we prefer  $T_n'$

### Synthesis

#### Parametric Estimation

- Assumptions: the results of experiments are modeled by a sequence of independent identically distributed with a **KNOWN** probability law
- An estimator
- The bias of an estimator is expected difference between the estimator & the real value  $\mathbb{E}_\theta(T_n - \theta)$
- The risk is the expectation of the error around the real value  $\mathbb{E}_\theta(T_n - \theta)^2$

#### Non-parametric Estimation - estimation of the mean

- Assumptions: the results of experiments are modeled by a sequence of independent identically distributed with a **UNKNOWN** probability law, but with some properties mean  $m$  & variance  $b^2$
- The average  $\frac{1}{n} \sum X_i$  is an unbiased estimator of  $m$
- The error  $\frac{1}{n} \sum X_i - m$  converges almost surely to 0 (strong law of large numbers)
- The law of errors satisfies the CLT (ci and confidence intervals):

$$\frac{1}{\sqrt{n}} \sum \frac{(X_i - m)}{b^2} \xrightarrow{d} N(0, 1)$$

# Data Statistics Introduction

\* 3 parts:

- ① The Problem : The Data Set - Building a DataSet (DATA)
- ② The Problem : The Data Set - How data was produced
- ③ Explorations

(1)

## DATA PRODUCTION

- Why this data has been produced? (purpose)
- Which approach has been used? (method)
- How this dataset has been practically produced? (observations)

## ANALYSIS OF THE SET OF VARIABLES

- Identification of Variable types
- Identification of the Variable role
- Identification of the Variables semantics

## ANALYSIS OF THE TYPE OF VARIABLES

- Nominal Variables : classification, membership
- Ordinal Variables : comparison! level
- Quantitative Variables : quantities

## USAGE OF VARIABLES

- Response Variables → response time, duration ...
- Explanatory Variables → size, load (affect response)
- Univariate → 1 variable is used  
Multivariate → several --

TO BE CONTINUED