

# Resampling techniques for statistical modeling

Gianluca Bontempi

Département d'Informatique  
Boulevard de Triomphe - CP 212  
<http://www.ulb.ac.be/di>

# Rationale

- All the methods presented so far are strongly based on two assumptions:
  1. the form of distribution underlying the data is parametric and known. The only uncertainty concerns the values of some parameters.
  2. we focus on parameters  $\theta$  (e.g. the mean) and corresponding estimators  $\hat{\theta}$  (e.g. the average) for which the analysis is simple.

# Rationale(II)

- For a generic parameter in a parametric setting, the distribution of the statistic  $\hat{\theta}$  cannot be easily defined.
- For a generic problem, when only a set of observations are available, the parametric assumption cannot be made without strongly biasing the procedure.
- Two alternatives:
  1. make simpler the complex problem, with the risk of oversimplifying.
  2. use computer intensive techniques.
- Note that, also when the configuration is parametric, computer intensive methods could still be useful to assess the robustness of conclusions drawn from a parametric analysis.

# Estimation of arbitrary statistics

- Consider a set  $D_N$  of  $N$  data points sampled from a one-dimensional distribution of a r.v.  $\mathbf{x}$ .
- Suppose that our quantity of interest is the mean of  $\mathbf{x}$ . It is straightforward to compute the estimate  $\hat{\mu}$  of the mean, the bias and the variance of the estimator:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \text{Bias}[\hat{\mu}] = 0, \quad \text{Var}[\hat{\mu}] = \frac{1}{N(N-1)} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

- Consider now another quantity of interest, for example the median or a mode of the distribution. While it is easy to have an estimate of these quantities, their accuracy is difficult to be computed.
- In other terms, given an arbitrary estimator  $\hat{\theta}$ , the analytical form of the variance  $\text{Var}[\hat{\theta}]$  and the bias  $\text{Bias}[\hat{\theta}]$  is typically not available.

# Example: the patch data

- Eight subjects wore medical patches designed to infuse a certain hormone in the blood.
- Each subject had his hormone levels measured after wearing three different patches: a placebo, an “old” patch and a “new” patch.
- The goal is to show bioequivalence. In other terms the Food and Drug Administration (FDA) will approve the new patch for sale only if the new patch is bioequivalent to the old one.
- The FDA criterion is

$$\theta = \frac{|E(\text{new}) - E(\text{old})|}{E(\text{old}) - E(\text{placebo})} \leq 0.2$$

# Example: the patch data (II)

subj	plac	old	new	z=old-plac	y=new-old	
1	9243	17649	16449	8406	-1200	
2	9671	12013	14614	2342	2601	
3	11792	19979	17274	8187	-2705	
...	...	...	...	...	...	
8	18806	29044	26325	10238	-2719	
mean:				6342	-452.3	Estimate

$$\hat{\theta} = t(\hat{F}) = \frac{|\bar{y}|}{\bar{z}} = \frac{452.3}{6342} = 0.07$$

Does it satisfy the FDA criterion?

What about its accuracy, bias, variance?

# Jackknife

- The **jackknife** (or **leave-one-out**) resampling technique aims at providing a computational procedure to estimate the variance and the bias of a generic estimator  $\hat{\theta}$ .
- The technique was first proposed by Quenouille in 1949.
- The technique is based on removing samples from the available dataset and recalculating the estimator.
- It is a general-purpose tool which is easy to implement and solves a number of problems.

# Jackknife for the mean

In order to show the theoretical foundation of the jackknife, we apply first this technique to the estimator  $\hat{\mu}$  of the mean .

Suppose we have available a dataset  $D_N$ . Let us remove the  $i$ th sample from  $D_N$  and let us calculate the **leave-one-out (l-o-o) mean** estimate from the  $N - 1$  remaining samples

$$\hat{\mu}_{(i)} = \frac{1}{N-1} \sum_{j \neq i}^N x_j = \frac{N\hat{\mu} - x_i}{N-1}$$

Observe that the following relation holds

$$x_i = N\hat{\mu} - (N-1)\hat{\mu}_{(i)}$$

that is, we can calculate  $x_i$  if we know both  $\hat{\mu}$  and  $\hat{\mu}_{(i)}$ .

# Jackknife for an arbitrary statistic

Suppose that we wish to estimate some parameter  $\theta$  as a very complex statistic of the  $N$  data points

$$\hat{\theta} = f(D_N) = f(x_1, x_2, \dots, x_N)$$

We first compute

$$\hat{\theta}_{(i)} = f(x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_N)$$

which is called the  $i$ th **jackknife replication** of  $\hat{\theta}$ . Then by analogy with the relation holding for the mean we define the  $i$ -th **pseudovalue** by

$$\eta_{(i)} = N\hat{\theta} - (N - 1)\hat{\theta}_{(i)}$$

These pseudovalues assume the same role as the  $x_i$  in calculating the sampled mean.

# The jackknife estimate of bias

Hence the **jackknife estimate** of  $\theta$  is given by

$$\hat{\theta}_{\text{jk}} = \frac{1}{N} \sum_{i=1}^N \eta_{(i)} = \frac{1}{N} \sum_{i=1}^N \left( N\hat{\theta} - (N-1)\hat{\theta}_{(i)} \right) = N\hat{\theta} - (N-1)\hat{\theta}_{(\cdot)}$$

where

$$\hat{\theta}_{(\cdot)} = \frac{\sum_{i=1}^N \hat{\theta}_{(i)}}{N}$$

Since  $\text{Bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$ , if we put  $\text{Bias}_{\text{jk}}[\hat{\theta}] = \hat{\theta} - \hat{\theta}_{\text{jk}}$ , we obtain the **jackknife estimate of the bias of  $\hat{\theta}$**

$$\text{Bias}_{\text{jk}}[\hat{\theta}] = (N-1)(\hat{\theta}_{(\cdot)} - \hat{\theta})$$

Note that in the particular case of a mean estimator (i.e.  $\hat{\theta} = \hat{\mu}$ ), we see that we obtain, as expected,  $\text{Bias}_{\text{jk}}[\hat{\mu}] = 0$ .

# The jackknife estimate of variance

An estimate of the variance of  $\hat{\theta}$  can be obtained from the sample variance of the pseudo-values

$$\begin{aligned}\text{Var}_{\text{jk}}[\hat{\theta}] &= \text{Var} [\hat{\theta}_{\text{jk}}] = \\ &= \frac{\text{Var} [\eta_{(i)}]}{N} = \frac{\sum_{i=1}^N (\eta_{(i)} - \hat{\theta}_{\text{jk}})^2}{N(N-1)} = \left( \frac{N-1}{N} \sum_{i=1}^N (\hat{\theta}_{(i)} - \hat{\theta}_{(\cdot)})^2 \right)\end{aligned}$$

Note that in the particular case of a mean estimator (i.e.  $\hat{\theta} = \hat{\mu}$ ), it can be shown that

$$\text{Var}_{\text{jk}}[\hat{\theta}] = \text{Var} [\hat{\mu}] = \frac{\sum_{i=1}^N (x_i - \hat{\mu})^2}{N(N-1)}$$

# Consideration on jackknife

- The major motivation for jackknife estimates is that they reduce bias.
- However, the jackknife can fail miserably if the statistic  $\hat{\theta}$  is not smooth (i.e. small changes in data cause small changes in the statistic).
- An example of non-smooth statistic for which the jackknife works badly is the median.

# Bootstrap: why this name?

- Bootstrap is a data-based simulation method for statistical inference.
- The use of the term *bootstrap* derives from the phrase *to pull oneself up by one's bootstrap*, widely thought to be based on one of the eighteenth century *Adventures of Baron Munchausen*, by R.E. Raspe.
- The Baron had fallen to the bottom of a deep lake. Just when it looked like all was lost, he thought to pick himself up by his own bootstraps.
- It is not the same as the term “bootstrap” used in computer science to “boot” a computer from a set of core instructions (though the derivation is similar).

# The bootstrap

- The bootstrap was introduced in 1979 as a computer-based method for estimating the standard error of an estimator  $\hat{\theta}$ .
- It enjoys the advantage of being completely automatic.
- The bootstrap estimate of standard errors requires no theoretical calculations and is available no matter how mathematically complicated the estimator  $\hat{\theta}$  may be.
- The bootstrap is extensively based on the plug-in principle.

# The setting of the problem

- Consider a dataset  $D_N$ , whose sample values are the outcomes of iid r.v.s  $z_i$  whose probability density is  $p_{\mathbf{z}}$ .
- We want to make inference about  $\theta = t(F)$  which is a characteristic of the population.
- We have chosen an estimator  $\hat{\theta}$ .
- Our attention is focused on questions concerning the probability distribution of  $\hat{\theta}$ . For example we would like to know about
  - its bias,
  - its standard error
  - its quantiles
  - its confidence values and so on.

# The parametric case

- Suppose that  $z \sim F_{\mathbf{z}}(z, \psi)$  and that we are interested in a generic parameter  $\theta = t(F)$ .
- Note that we have distinguished  $\theta$  from  $\psi$  in order to have a more general setting. For example if  $z \sim \mathcal{N}(\mu, \sigma^2)$ , it could be that  $\psi = \mu$  and  $\theta = \log \mu$ .
- Suppose we have an estimate  $\hat{\psi}$  and a corresponding fitted distribution  $F(z, \hat{\psi})$ . We will define  $\mathbf{z}^* \sim F(z, \hat{\psi})$ .
- Suppose we want to calculate the distribution (e.g. bias or variance) of  $\hat{\theta}$ . For a generic  $\hat{\theta}$  this operation is arduous.
- The idea is to estimate  $\hat{\theta}$  by using a simulated dataset.

# The parametric case (II)

- Let us consider the dataset

$$D_{(b)} = \{z_1^*, z_2^*, \dots, z_N^*\}$$

obtained by iid sampling the **known distribution**  $\mathbf{z}^* \sim \hat{F}(z, \hat{\psi})$ .  $D_{(b)}$  is called the **bootstrap sample**.

- Consider  $B$  repetitions of the  $D_{(b)}$  sampling.
- Let us define by

$$\hat{\theta}_{(b)} = t(D_{(b)}) \quad b = 1, \dots, B$$

the statistic computed by  $D_{(b)}$ ,  $b = 1, \dots, B$ .  $\hat{\theta}_{(b)}$  is called **bootstrap replication** of  $\hat{\theta}$ .

- The statistical properties of  $\hat{\theta}$  can now be calculated on the basis of the distribution of  $\hat{\theta}_{(b)}$ .

# Bootstrap estimate of the variance

The bootstrap estimate of the variance of the estimator  $\hat{\theta}$ , is the variance of the set  $\hat{\theta}_{(b)}$ ,  $b = 1, \dots, B$ .

$$\text{Var}_{\text{bs}}[\hat{\theta}] = \text{Var} \left[ \hat{\theta}_{(b)} \right] = \frac{\sum_{b=1}^B (\hat{\theta}_{(b)} - \hat{\theta}_{(\cdot)})^2}{(B - 1)}$$

where

$$\hat{\theta}_{(\cdot)} = \frac{\sum_{b=1}^B \hat{\theta}_{(b)}}{B}$$

If  $\hat{\theta} = \hat{\mu}$ , for  $B \rightarrow \infty$  the bootstrap estimate  $\text{Var}_{\text{bs}}[\hat{\theta}]$  converges to the variance  $\text{Var}[\hat{\mu}]$ .

# Bootstrap estimate of bias

Let  $\hat{\theta}$  be the estimator based on the original sample  $D_N$  and

$$\hat{\theta}_{(\cdot)} = \frac{\sum_{b=1}^B \hat{\theta}_{(b)}}{B}$$

Since  $\text{Bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$ , the bootstrap estimate of bias is

$$\text{Bias}_{\text{bs}}[\hat{\theta}] = \hat{\theta}_{(\cdot)} - \hat{\theta}$$

Then, since

$$\theta = E[\hat{\theta}] - \text{Bias}[\hat{\theta}]$$

the **bootstrap bias corrected** estimate is

$$\hat{\theta}_{\text{bs}} = \hat{\theta} - (\hat{\theta}_{(\cdot)} - \hat{\theta}) = 2\hat{\theta} - \hat{\theta}_{(\cdot)}$$

# Bootstrap confidence interval

- **Standard bootstrap confidence limits** are based on the assumption that the estimator  $\hat{\theta}$  is normally distributed with mean  $\theta$  and variance  $\sigma^2$ .
- Taking the bootstrap estimate of variance, an approximate  $100(1 - \alpha)\%$  confidence interval is given by

$$\hat{\theta} \pm z_{\alpha/2} \sqrt{\text{Var}_{\text{bs}}[\hat{\theta}]} = \hat{\theta} \pm z_{\alpha/2} \sqrt{\frac{\sum_{b=1}^B (\hat{\theta}_{(b)} - \hat{\theta}_{(\cdot)})^2}{(B - 1)}}$$

- An improved interval is given by using the bootstrap bias corrected estimate of  $\theta$ . The interval becomes

$$\hat{\theta}_{\text{bs}} \pm z_{\alpha/2} \sqrt{\text{Var}_{\text{bs}}[\hat{\theta}]} = 2\hat{\theta} - \hat{\theta}_{(\cdot)} \pm z_{\alpha/2} \sqrt{\frac{\sum_{b=1}^B (\hat{\theta}_{(b)} - \hat{\theta}_{(\cdot)})^2}{(B - 1)}}$$

# Bootstrap percentile confidence interval

- A more direct approach for constructing a  $100(1 - \alpha)\%$  confidence interval is to use the upper and lower  $\alpha/2$  values of the bootstrap distribution.
- The approaches using the full bootstrap distribution are often referred to as **percentile confidence limits**.
- If  $\hat{\theta}_{L,\alpha/2}$  denotes the value such that only a fraction  $\alpha/2$  of all bootstrap estimates are inferior to it, and likewise  $\hat{\theta}_{H,\alpha/2}$  is the value exceeded by only  $\alpha/2$  of all bootstrap estimates, then an approximate confidence interval is given by

$$[\hat{\theta}_{L,\alpha/2}, \hat{\theta}_{H,\alpha/2}]$$

also called the **Efron's percentile confidence limit**.

# Studentized bootstrap statistic

Recall that for  $\mathbf{z} \sim \mathcal{N}(\mu, \sigma^2)$  with  $\mu$  and  $\sigma^2$  unknown, the  $(1 - \alpha)$  interval of confidence of  $\mu$  is

$$[\hat{\mu} - t_{(N-1, \alpha/2)} \hat{\sigma} / \sqrt{N}, \hat{\mu} + t_{(N-1, \alpha/2)} \hat{\sigma} / \sqrt{N}]$$

For a generic distribution of  $\mathbf{z}$  the Student  $\mathcal{T}$  distribution is no more appropriate. At its place, the bootstrap method can be used to approximate the distribution of the random variable

$$\frac{\hat{\theta} - \theta}{\hat{\sigma} / \sqrt{N}}$$

The idea is to estimate the above variable through the **Studentized bootstrap statistic**

$$\mathbf{s}_{(b)} = \frac{\hat{\theta}_{(b)} - \hat{\theta}}{\sqrt{\text{Var} [\hat{\theta}_{(b)}]}}$$

# Studentized bootstrap confidence interval

Given  $B$  bootstrap replications we obtain the empirical distribution of  $\mathbf{s}_{(b)}$ .

Hence the studentized bootstrap confidence interval for  $\theta$  is

$$[\hat{\theta} - s_{(\alpha/2)} \hat{\sigma} / \sqrt{N}, \hat{\mu} + s_{(\alpha/2)} \hat{\sigma} / \sqrt{N}]$$

where  $s_{(\alpha/2)}$  is the  $\alpha/2$  upper point of the empirical distribution of  $\mathbf{s}_{(b)}$ .

# The nonparametric case

- Suppose that we have no parametric model but that we can sample in an iid manner a set  $D_N$  from a **completely unknown** distribution function.
- We will use the empirical distribution function  $\hat{F}$  just as we would with a parametric model.
- Since the empirical  $\hat{F}$  puts equal probabilities on each original data value, each simulated dataset  $D_{(b)}$  is a random sample taken **with replacement** from  $D_N$ .
- This resampling procedure is called the **nonparametric bootstrap**.
- It differs from the parametric case only for the generation of the  $B$  datasets  $D_{(b)}$ .

# The nonparametric bootstrap algorithm

Suppose we want to compute the standard error of an estimator  $\hat{\theta}$ .

The algorithm has the following steps:

1. Select  $B$  independent bootstrap samples  $D_{(b)} = \{z_1^*, z_2^*, \dots, z_N^*\}$ ,  $b = 1, \dots, B$ , each consisting of  $N$  data values drawn with replacement from  $D_N$ .
2. Evaluate the bootstrap replication corresponding to each bootstrap sample

$$\hat{\theta}_{(b)} = t(D_{(b)}) \quad b = 1, \dots, B$$

3. Estimate the standard error  $\sqrt{\text{Var} [\hat{\theta}]}$  by the sample standard deviation  $\sqrt{\text{Var} [\hat{\theta}_{(b)}]}$  of the  $B$  replications.

# Practical example

- Consider a dataset containing the service hours between failures of the air-conditioning equipment in a Boeing aircraft

$$D_N = \{3, 5, 7, 18, 43, 85, 91, 98, 100, 130, 230, 487\}$$

- We wish to estimate the mean or its reciprocal, the failure rate. Then in this practical example  $\hat{\theta} = \hat{\mu}$ .
- We could adopt a parametric approach and suppose that the probabilistic distribution underlying the data is exponential
- Note that in this case it is known that

$$\text{Bias}[\hat{\mu}] = 0, \quad \text{Var}[\hat{\mu}] = \hat{\mu}^2 / N$$

- Let us see the relation between empirical and exponential distribution function ( $\hat{\lambda} = 1/\hat{\mu}$ ).
- R script `airco.R`.

# Exponential distribution

A continuous random variable  $\mathbf{z}$  is said to be **exponentially distributed** with rate  $\lambda > 0$  (also  $\mathbf{z} \sim \mathcal{E}(\lambda)$ ) if its probability density function is given by

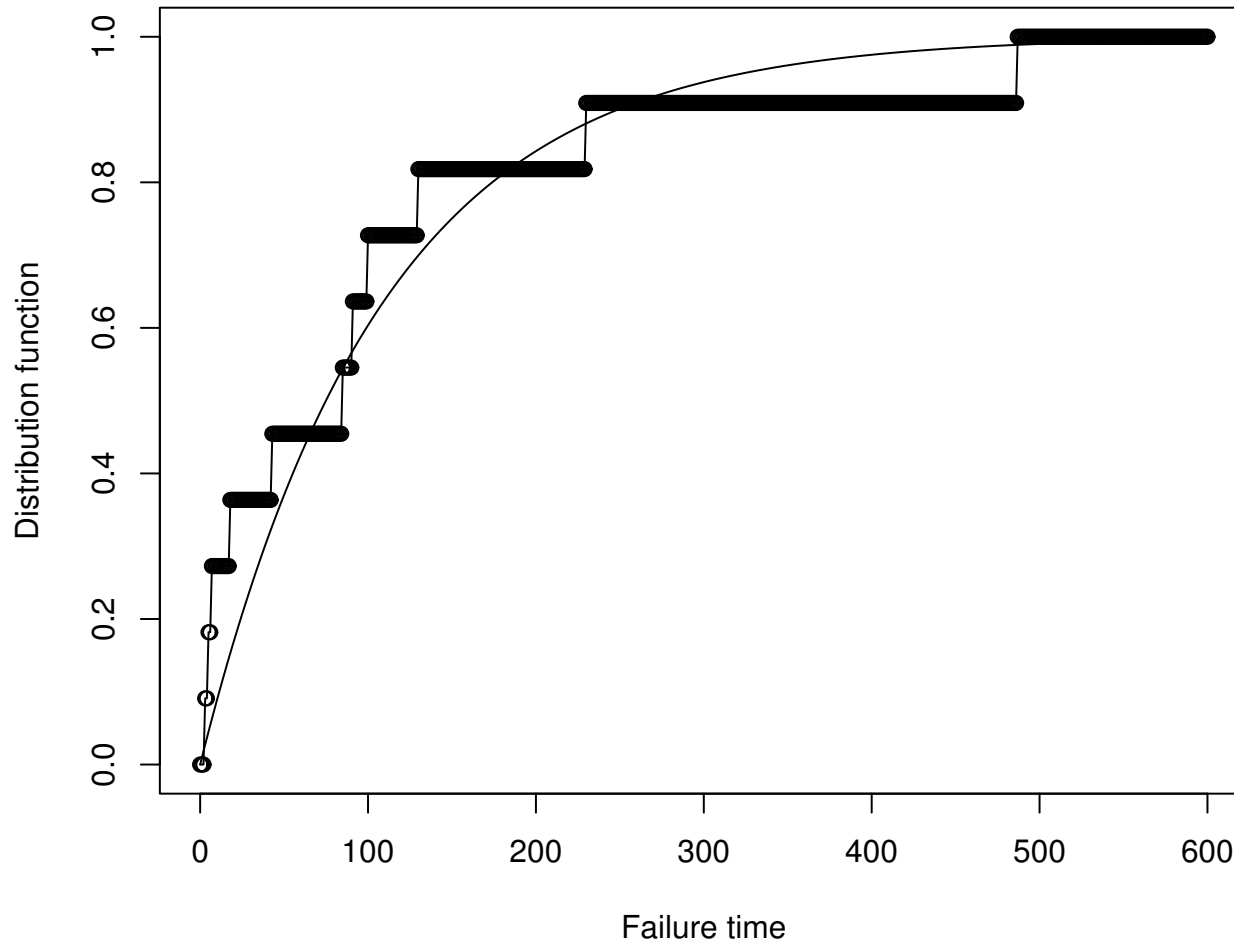
$$p_{\mathbf{z}}(z) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

- The mean of  $\mathbf{z}$  is  $1/\lambda$ .
- The variance of  $\mathbf{z}$  is  $1/\lambda^2$ .
- It can be considered as continuous approximation to the geometric distribution.
- Like the geometric distribution it satisfies the **memoryless property**

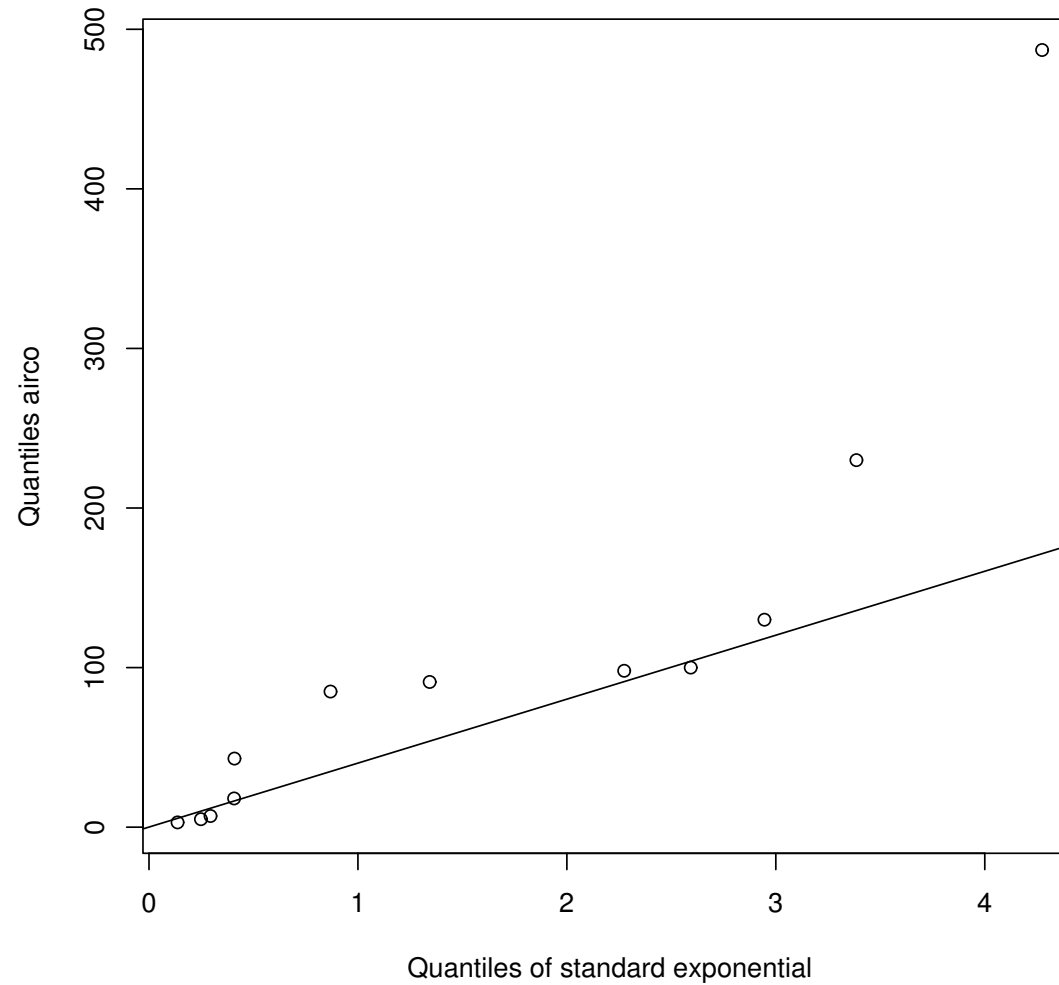
$$\text{Prob} \{ \mathbf{z} \geq z_1 + z_2 | \mathbf{z} = z_1 \} = \text{Prob} \{ \mathbf{z} \geq z_2 \}$$

- It is used to describe physical phenomena (e.g. radioactive decay time or failure time).

# Ex: empirical and exponential distribution



# Ex: empirical and exponential quantiles



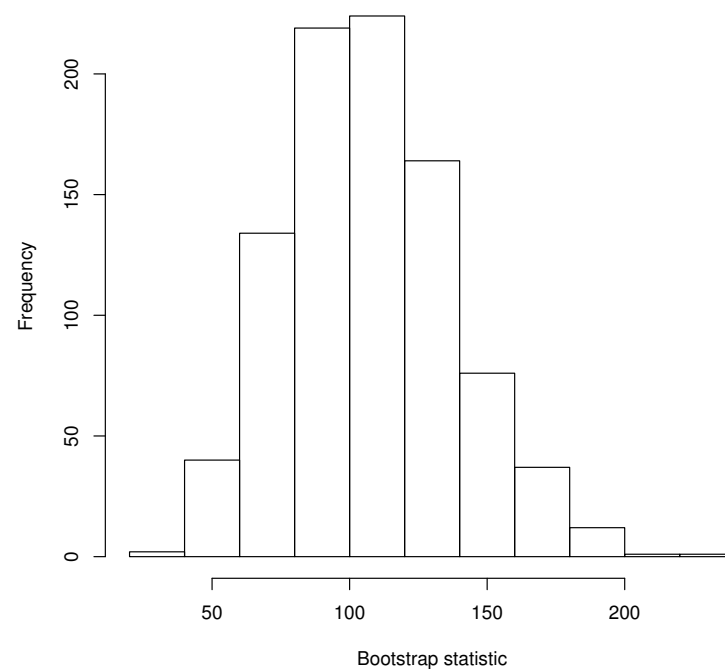
# Ex: parametric bootstrap

- Consider as statistic  $\hat{\theta}$  the sample average  $\hat{\mu}$ .
- Assume we know that the distribution is exponential but that we have not access to the analytical expressions of bias and variance of  $\hat{\mu}$ .
- We perform parametric bootstrap to estimate these quantities.
- We sample  $B$  times the exponential distribution  $\mathcal{E}(1/\hat{\mu})$ .
- In the following we report the outcome of a set of bootstrap experiments.

# Ex: distribution of $\hat{\theta}_{(b)}$

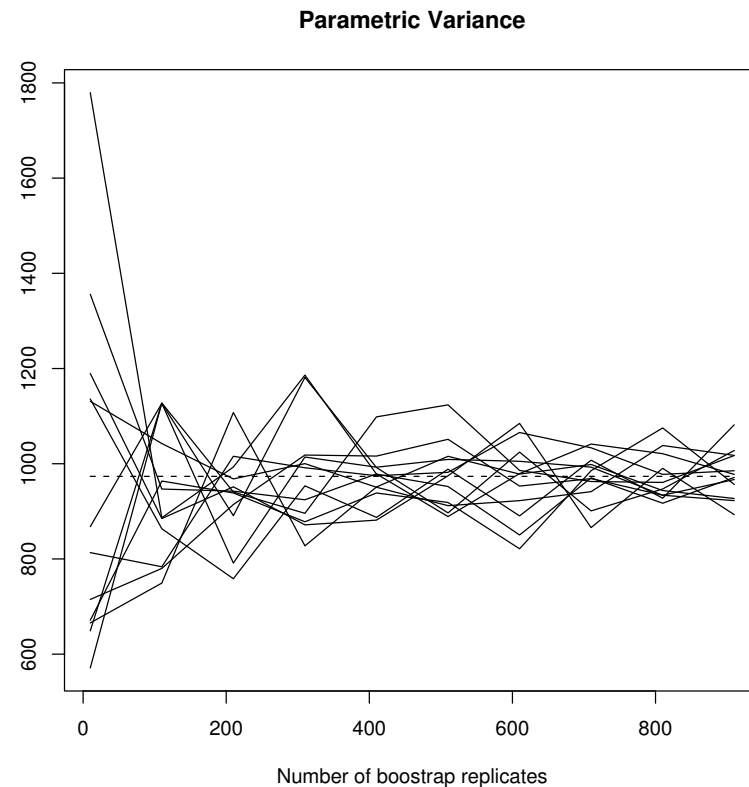
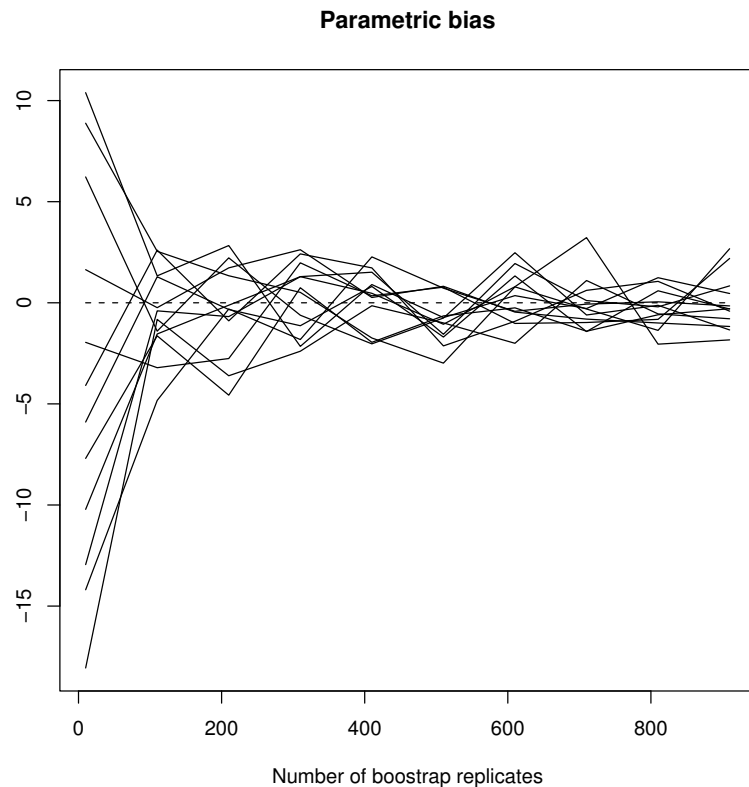
Consider the statistic  $\hat{\theta} = \hat{\mu}$ .

The distribution of  $\hat{\theta}_{(b)}$  on the basis of  $B = 100$  of bootstrap parametric repetitions is represented by the following histogram



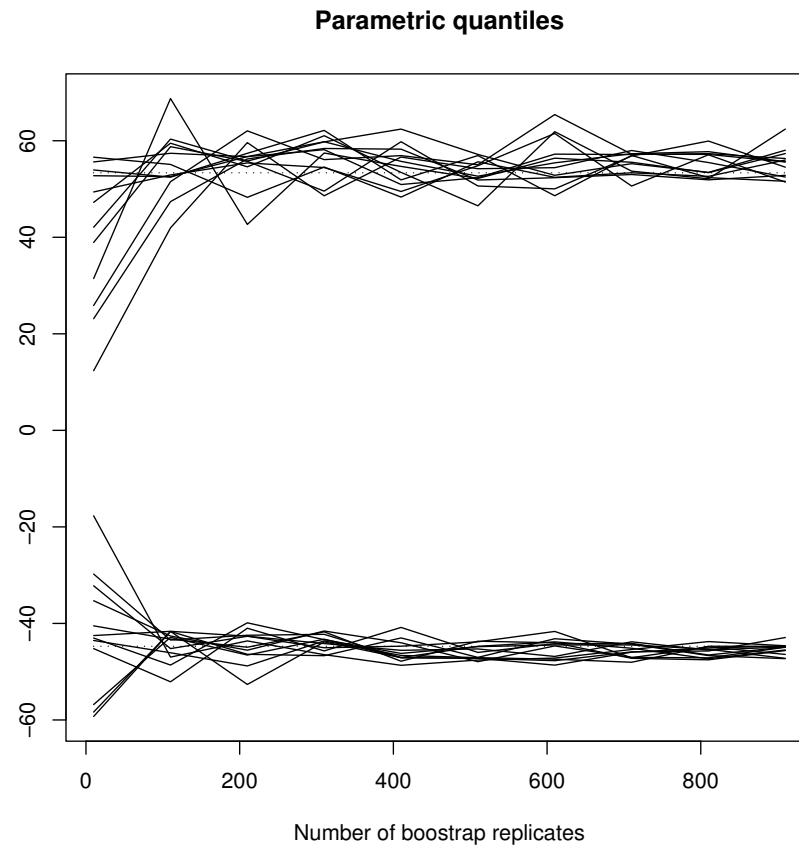
# Ex: parametric bootstrap

We report the bias and the variance of  $\hat{\theta}$  obtained in different random experiments. Each experiment is based on a set of bootstrap repetitions obtained by varying the number  $B$ .



# Ex: parametric bootstrap

We report the quantiles of  $\hat{\theta}$  obtained in different random experiments. Each experiment is based on a set of bootstrap repetitions obtained by varying the number  $B$ .



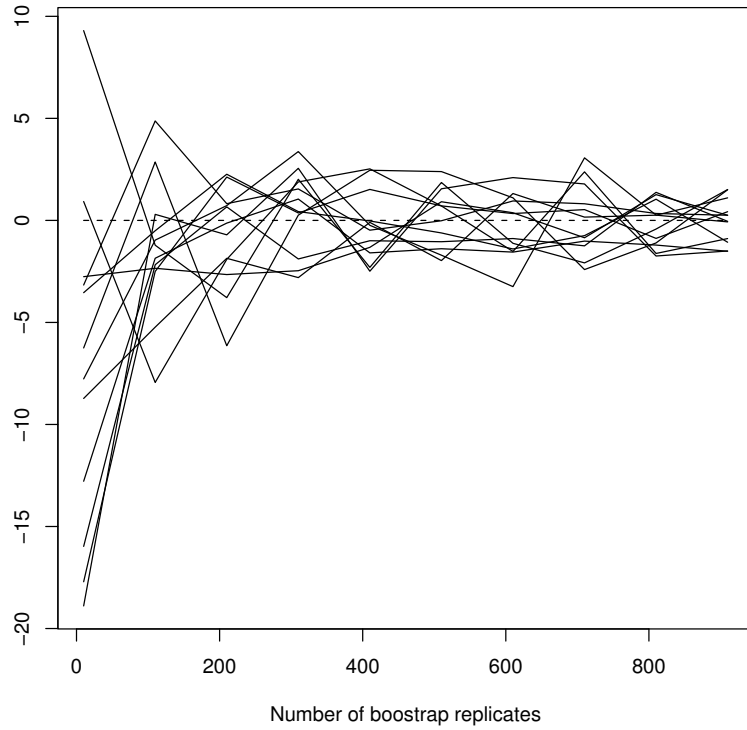
Dotted lines: theoretical quantiles under the hypothesis of exponential model.

# Ex: nonparametric bootstrap

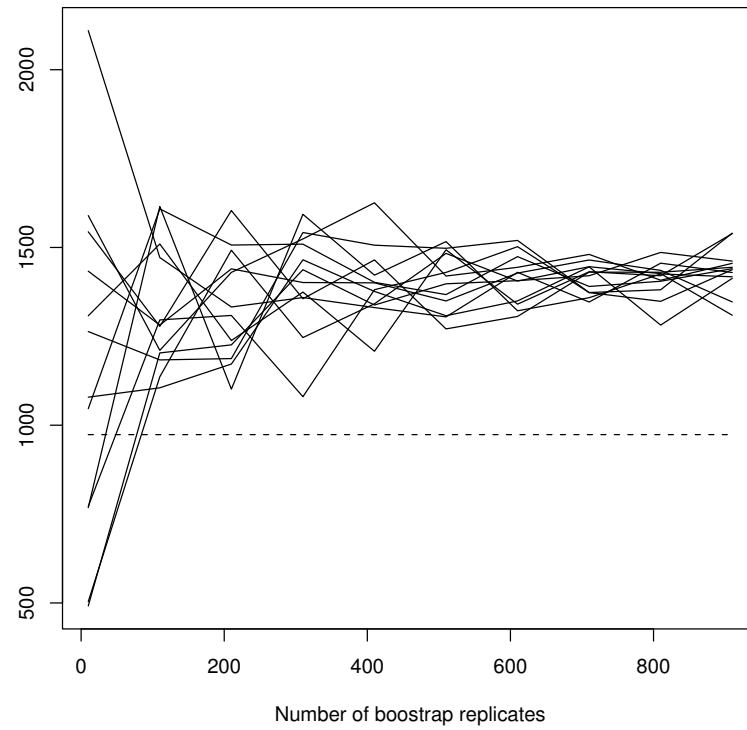
- Assume we do not know that the distribution is quasi exponential.
- We perform nonparametric bootstrap to estimate bias and variance of the estimator  $\hat{\mu}$ .
- We sample  $B$  times with replacement the dataset  $D_N$ .
- We report the bias, the variance and the quantiles for different repetitions each having a set of  $B$  nonparametric trials.

# Ex: nonparametric bootstrap

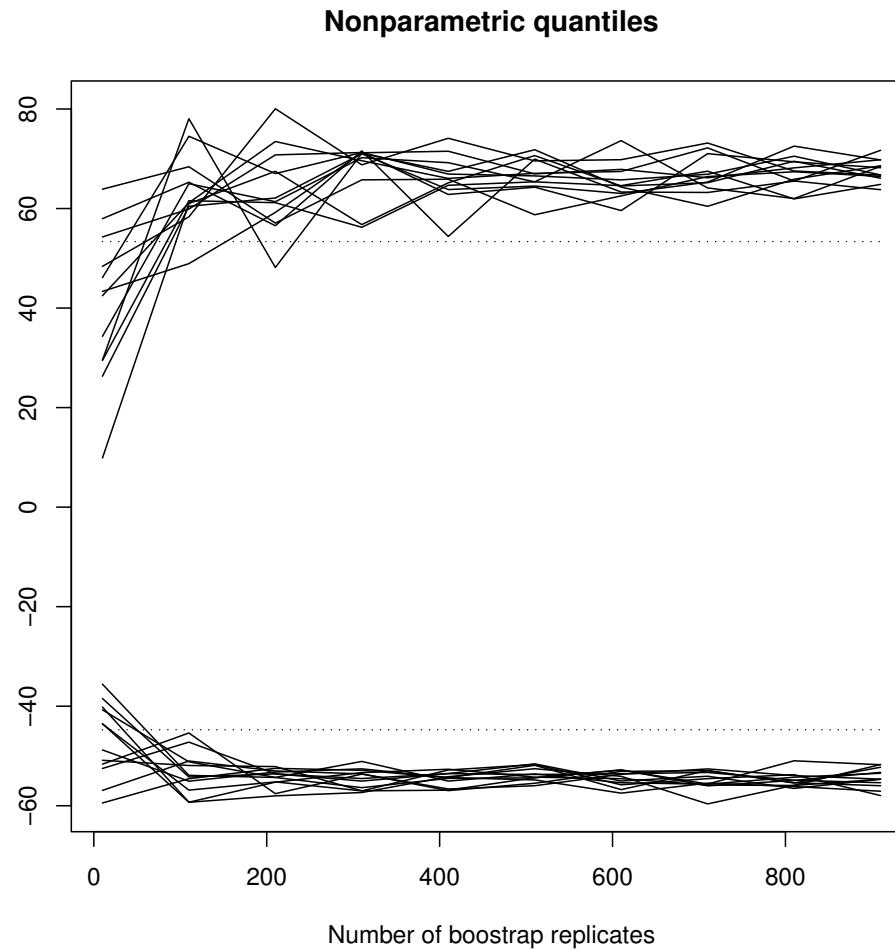
Nonparametric bias



Nonparametric Variance



# Ex: nonparametric bootstrap



Note that the nonparametric quantiles do not converge to the theoretical quantiles (dotted lines) based on the hypothesis of exponential model.

# Parametric and nonparametric bootstrap

- For continuous data, a major difference between parametric and nonparametric resampling lies in the discreteness of the latter. Under nonparametric resampling the  $\hat{\theta}_{(b)}$  will have always a discrete distribution.
- However, given  $N$  distinct samples there are up to

$$B_{\max} = \binom{2N - 1}{N - 1} = \frac{(2N - 1)!}{N!(N - 1)!}$$

different bootstrap samples and consequently up to  $B_{\max}$  possible values of  $\hat{\theta}_{(b)}$  depending on the smoothness of  $t$ . For example for  $N = 11$ ,  $B_{\max} = 352716$ .

- For many practical applications, the effects of the discreteness are likely to be fairly minimal.

# The bootstrap principle

- What we would like to know in an estimation problem is the distribution of  $\hat{\theta} - \theta$ .
- What we have in bootstrap is a Monte Carlo approximation to the distribution  $\hat{\theta}_{(b)} - \hat{\theta}$ .
- The key idea of the bootstrap is that for  $N$  sufficiently large we expect the two distributions to be nearly the same.
- In other terms the variability of  $\hat{\theta}_{(b)}$  (based on the empirical distribution) around  $\hat{\theta}$  is expected to be similar (or mimic) the variability of  $\hat{\theta}$  (based on the true distribution) around  $\theta$ .
- There is good reason to believe this will be true for large  $N$ , since as  $N$  gets larger and larger, the empirical  $\hat{F}(\cdot)$  converge to  $F(\cdot)$  (see the Glivenko-Cantelli theorem for iid samples).
- This idea is sometimes referred to as the *bootstrap principle*.

# Error in resampling methods

- The error in resampling methods is generally a combination of **statistical error** and **simulation error**.
- Statistical error is due to the difference between the underlying distribution  $F$  and the empirical distribution  $\hat{F}$ . The magnitude of this error depends on the choice of  $t(F)$ .
- The use of rough statistics  $t(F)$  (e.g. unsmooth or unstable) can make the resampling approach behave wildly. Example of nonsmooth statistics are sample quantiles and the median.
- The simulation error is due to the use of empirical (Monte Carlo) properties of  $t(F)$  rather than exact properties.
- Simulation error decreases by increasing the number  $B$  of bootstrap replications.

# Convergence of bootstrap estimate

In general terms for iid observations, we require

1. the convergence of  $\hat{F}$  to  $F$  (satisfied by the Glivenko-Cantelli theorem) for  $N \rightarrow \infty$ ;
2. an estimator such that the estimate  $\hat{\theta}$  is the corresponding functional of the empirical distribution.

$$\theta = t(F) \rightarrow \hat{\theta} = t(\hat{F})$$

This is satisfied for sample means, standard deviations, variances, medians and other sample quantiles.

3. a smoothness condition on the functional. This is not true for extreme order statistics such as the minimum and the maximum values.

# When might the bootstrap fail?

So far we have assumed that the dataset  $D_N$  is iid sampled from a distribution  $F$ .

In some non conventional configurations, bootstrap might fail. For example

- Incomplete data (survival data, missing data).
- Dependent data (e.g. variance of a correlated time series).
- Dirty data (outliers)

For a critical view on bootstrap, see the publication *Exploring the limits of bootstrap* edited by Le Page and Billard which is a compilation of the papers presented at a special conference of the Institute of Mathematical Statistics held in Ann Arbor, Michigan, 1990.