

ADVANCED BIOSTATISTICS

ABSTAT17

Monte Carlo and Bootstrap methods

Carina Silva

(*carina.silva@estesl.ipl.pt*)

Higher School of Technologies and Health of Lisbon &
Center of Statistics and Applications, University of Lisbon

IGC, April, 10–13, 2017

Introduction

In statistical models construction there are some steps to consider:

- ▶ Model hypothesis.
- ▶ Parameter estimation of the model.
- ▶ Validation of the model (goodness-of-fit, interpretation)
- ▶ Model comparisons (parsimony).

Monte Carlo (MC) simulation to the rescue

- ▶ An estimator or test statistic has a true sampling distribution under a particular set of conditions (finite sample size, true distribution of the data, etc.).

Monte Carlo (MC) simulation to the rescue

- ▶ An estimator or test statistic has a true sampling distribution under a particular set of conditions (finite sample size, true distribution of the data, etc.).
- ▶ Ideally, we would want to know this true sampling distribution.

Monte Carlo (MC) simulation to the rescue

- ▶ An estimator or test statistic has a true sampling distribution under a particular set of conditions (finite sample size, true distribution of the data, etc.).
- ▶ Ideally, we would want to know this true sampling distribution.
- ▶ But sometimes derivation of the true sampling distribution is not tractable or impossible.

Monte Carlo (MC) simulation to the rescue

- ▶ An estimator or test statistic has a true sampling distribution under a particular set of conditions (finite sample size, true distribution of the data, etc.).
- ▶ Ideally, we would want to know this true sampling distribution.
- ▶ But sometimes derivation of the true sampling distribution is not tractable or impossible.
- ▶ MC simulation allows **approximate** the **sampling distribution** of an estimator or test statistic **under a particular set of conditions**.

Monte Carlo simulation

A typical Monte Carlo simulation involves the following:

- ▶ Generate S independent data sets under the conditions of interest.

Monte Carlo simulation

A typical Monte Carlo simulation involves the following:

- ▶ Generate S independent data sets under the conditions of interest.
- ▶ Compute numerical value of the estimator/test statistic, for each data set $\implies t_1^*, \dots, t_S^*$

Monte Carlo simulation

A typical Monte Carlo simulation involves the following:

- ▶ Generate S independent data sets under the conditions of interest.
- ▶ Compute numerical value of the estimator/test statistic, for each data set $\implies t_1^*, \dots, t_S^*$
- ▶ If S is large enough, summary statistics across t_1^*, \dots, t_S^* should be good **approximations** to the true sampling properties of the estimator/test statistic under the conditions of interest.

Example

For an estimator for a parameter θ : T_S is the value of T from the s th data set, $s = 1, \dots, S$

- ▶ The sample mean over S data sets is an estimate of the true mean of the sampling distribution of the estimator

Simple example

Compare three estimators for the mean μ of a distribution based on i.i.d. draws Y_1, \dots, Y_n

- ▶ Sample mean: $T^{(1)}$
- ▶ Sample 20% trimmed mean: $T^{(2)}$
- ▶ Sample median: $T^{(3)}$

Remarks:

- ▶ If the distribution of the data is symmetric, all three estimators indeed estimate the mean.
- ▶ If the distribution is skewed, they do not.

Simulation procedure

For a particular choice of μ , n and true underlying distribution:

- ▶ Generate independent draws Y_1, \dots, Y_n from the distribution
- ▶ Compute $T^{(1)}$, $T^{(2)}$ and $T^{(3)}$
- ▶ Repeat S times:
 $T_1^{(1)}, \dots, T_S^{(1)}; T_1^{(2)}, \dots, T_S^{(2)}; T_1^{(3)}, \dots, T_S^{(3)}$

Mean Squared Error

MSE: the mean squared error (MSE) of an estimator measures the average of the squares of the “errors”, that is, the difference between the estimator and what is estimated.

An MSE of zero, meaning that the estimator $\hat{\theta}$ predicts observations of the parameter θ with perfect accuracy, is the ideal, but is practically never possible.

Relative efficiency

For any estimators for which:

$$E(T^{(1)}) = E(T^{(2)}) = \mu \implies RE = \frac{\text{var}(T^{(1)})}{\text{var}(T^{(2)})}$$

is the *relative efficiency* of estimator 2 to estimator 1.

- ▶ When the estimators are not unbiased it is standard to compute

$$RE = \frac{MSE(T^{(1)})}{MSE(T^{(2)})}$$

- ▶ In either case $RE < 1$ means estimator 1 is preferred (estimator 2 is inefficient relative to estimator 1 in this sense)

Simulation procedure

- Compute for $k = 1, 2, 3$:

$$\widehat{mean}_{MC}^{(k)} = \frac{\sum_{s=1}^S T_s^{(k)}}{S} = \bar{T}^{(K)};$$

$$\widehat{bias}_{MC}^{(k)} = \bar{T}^{(K)} - \theta;$$

$$\widehat{SD}_{MC}^{(k)} = \sqrt{\frac{\sum_{i=1}^S (T_s^{(k)} - \bar{T}^{(K)})^2}{S-1}};$$

$$\widehat{MSE}_{MC}^{(K)} = \frac{\sum_{s=1}^S (T_s^{(k)} - \mu)^2}{S} \approx (\widehat{SD}^{(K)})^2 + (\widehat{bias}^{(K)})^2$$

Exercise 1

1. Generate 1000 MC replicates from samples of size 15 of a normal distribution with $\mu=1$ and $\sigma=1.7$.
2. Get the following MC estimates: MC mean, MC standard deviation, MC bias and MC MSE, of the true parameter μ , considering the sample mean, 20% trimmed mean and median estimators.
3. Compare the results and choose the best estimator.

Suppose that we can feasibly take another N samples from our population. Then from each sample we can generate another estimator:

$$x_1, \dots, x_M \xrightarrow{\text{sample}} \begin{cases} X_1^{(1)}, \dots, X_n^{(1)} & \rightarrow \hat{\theta}^{(1)} \\ \vdots & \\ X_1^{(N)}, \dots, X_n^{(N)} & \rightarrow \hat{\theta}^{(N)} \end{cases}$$

we can estimate properties of $\hat{\theta}$ using the estimates $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(N)}$. For example:

$$\widehat{E(\hat{\theta})} = \frac{1}{N} \sum_{i=1}^N \hat{\theta}^{(i)}$$

In most cases, however, it is not feasible to take new samples from our population. How else can we generate samples to estimate properties of $\hat{\theta}$?

Suppose that we **know the distribution** of the population (e.g. $F_\theta = \text{Exponential}(4)$). Then we can simulate N IID samples from F_θ and calculate the estimators from each sample

$$F_\theta \xrightarrow{\text{simulate}} \begin{cases} X_1^{(1)}, \dots, X_n^{(1)} & \rightarrow \hat{\theta}^{(1)} \\ \vdots & \\ X_1^{(N)}, \dots, X_n^{(N)} & \rightarrow \hat{\theta}^{(N)} \end{cases}$$

we can estimate properties of $\hat{\theta}$ using the estimates $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(N)}$.
For example:

$$\widehat{\text{Var}(\hat{\theta})} = \frac{1}{N} \sum_{i=1}^N \left(\hat{\theta}^{(i)} - \frac{1}{N} \sum_{j=1}^N \hat{\theta}^{(j)} \right)^2$$

In most cases, however, we don't know what F_θ is.

Introduction

- ▶ Bootstrapping is a method which uses random sampling techniques to estimate properties (such as bias, variance, confidence intervals, etc.) of an estimator, $\hat{\theta}$, when we don't know the true distribution, F_{θ} , of our data (and we cannot feasibly draw new samples from our population).

Introduction

The term “bootstrap” comes from literature. In “The Adventures of Baron Munchausen”, by Rudolph Erich Raspe, the Baron had fallen to the bottom of a deep lake, and he thought to get out by pulling himself up by his own bootstraps.

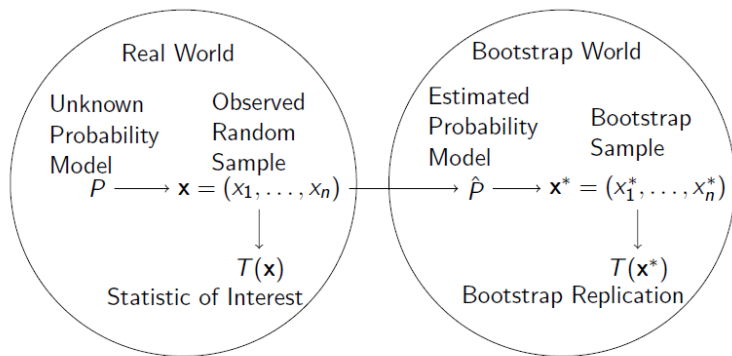


How does it work?

- ▶ *Step 1: Resampling.* A sampling distribution is based on many random samples from the population. In place of many samples from the population, create many resamples by repeatedly sampling with replacement from this one random sample. Each resample is the same size as the original random sample.

How does it work?

- ▶ *Step 2: Bootstrap distribution.* The sampling distribution of a statistic collects the values of the statistic from many samples. The bootstrap distribution of a statistic collects its values from many resamples.



Bootstrap

We will focus on two approaches to generating bootstrap samples:

- ▶ **Parametric bootstrap:** Estimate $\hat{\theta}$ from our original sample, X_1, \dots, X_n and generate samples X_1^*, \dots, X_n^* from $F_{\hat{\theta}}$, which approximates F_{θ} .
- ▶ **Non-parametric bootstrap:** Take samples X_1^*, \dots, X_n^* with replacement from our original sample X_1, \dots, X_n .

Once we have B bootstrap samples, we can generate B estimates of $\hat{\theta}$:

$$\begin{array}{ccc} X_1^{*(1)}, \dots, X_n^{*(1)} & \rightarrow & \hat{\theta}^{*(1)} \\ & \vdots & \\ X_1^{*(B)}, \dots, X_n^{*(B)} & \rightarrow & \hat{\theta}^{*(B)} \end{array}$$

We can use these bootstrapped estimators, $\hat{\theta}^{*(1)}, \dots, \hat{\theta}^{*(B)}$, to estimate properties of $\hat{\theta}$

Introduction

The number of bootstrap replications B

- ▶ A small number of replications $B=25$, is usually informative according to Efron.
- ▶ Fifty replications is often enough to give good estimate of standard error.
- ▶ Much bigger values of B are required for bootstrap confidence intervals.
- ▶ 1000 replications is recommended by Harrel and others for stable confidence intervals.

Parametric Bootstrap

Recall that we are using bootstrapping because we don't know the true parameter θ , but we want to identify the distribution of our estimator $\hat{\theta}$ (e.g. the MLE).

Suppose that we have observed a sample X_1, \dots, X_n , $X_i \sim F_\theta$.

Suppose further that we know F , but θ is unknown (For example, we know that our sample is $N(\theta, 2)$ distributed, but we don't know θ).

We can calculate $\hat{\theta} = T(X_1, \dots, X_n)$ for our sample.

- To identify the distribution of our estimator, $\hat{\theta}$, we want to obtain more observations of $\hat{\theta}$, but we only have one sample!

How can we get more (approximated) values of $\hat{\theta}$?

Parametric Bootstrap

Why not approximate the distribution, F_θ , using our estimate, $\hat{\theta}$?

We can then draw samples from the approximated distribution $F_{\hat{\theta}}$:

$$F_{\hat{\theta}} \xrightarrow{\text{simulate}} \begin{cases} X_1^{*(1)}, \dots, X_n^{*(1)} & \rightarrow \hat{\theta}^{*(1)} \\ \vdots & \\ X_1^{*(N)}, \dots, X_n^{*(N)} & \rightarrow \hat{\theta}^{*(N)} \end{cases}$$

We can now estimate properties of $\hat{\theta}$ using these bootstrapped estimates, for example:

$$\text{Var}_{boot}(\hat{\theta}) = \frac{1}{B} \sum_{i=1}^B \left(\hat{\theta}_i^* - \frac{1}{B} \sum_{j=1}^B \hat{\theta}_j^* \right)^2$$

Parametric bootstrap CI

Ingredients to use parametric bootstrap approach to estimate a confidence interval for a parameter:

- ▶ Data x_1, x_2, \dots, x_n drawn from a distribution $F(\theta)$ with unknown parameter θ .
- ▶ A statistic $\hat{\theta}$ that estimates θ .
- ▶ B Bootstrap samples drawn from $F(\hat{\theta})$.

Bootstrap CI

Percentile Method to a $(1-\alpha)*100\%$ CI for θ .

- ▶ For each bootstrap sample, $x_1^*, x_2^*, \dots, x_n^*$, we compute $\hat{\theta}^*$.
- ▶ Get the empirical percentiles $\hat{\theta}_{(\alpha/2)}^*$ and $\hat{\theta}_{(1-\alpha/2)}^*$.
- ▶ The bootstrap CI for θ will be given by $(\hat{\theta}_{(\alpha/2)}^*, \hat{\theta}_{(1-\alpha/2)}^*)$.

Bootstrap CI

Three types of confidence intervals:

► **Normal interval:**

$$\hat{\theta} \pm q_{(1-\alpha/2)} \sqrt{\text{Var}_{boot}(\hat{\theta})}$$

where q is the $1 - \alpha/2$ quantile of the standard normal distribution.

► **Pivotal interval:**

$$(2\hat{\theta} - \hat{\theta}_{(1-\alpha/2)}^*, 2\hat{\theta} - \hat{\theta}_{(\alpha/2)}^*)$$

Exercise 2

Suppose it was drawn a sample of 300 observations from an $\exp(\lambda)$ distribution. Estimate λ using a 95% parametric bootstrap confidence interval for λ .

Exercise 3

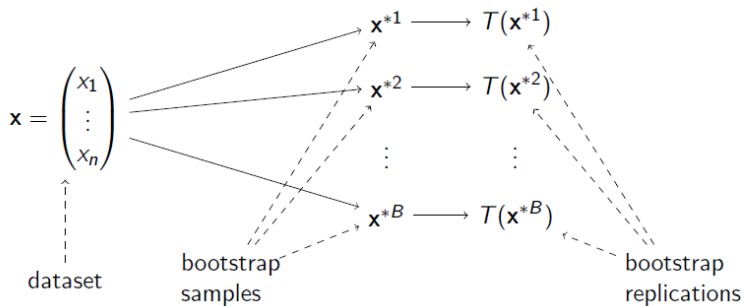
The following measurements were given for weights (Kg) of 11 children with ages between 8 and 10 years old with renal disfunction: 38.43, 38.43, 38.39, 38.83, 38.45, 38.35, 38.43, 38.31, 38.32, 38.48, 38.50.

a) Find the 95% parametric bootstrap confidence interval for μ assuming the normal distribution for the observations and $\sigma = 0.57$. Compare with the classical analytic approach based on the t -distribution.

N.B: Use $B=1000$ bootstrap samples (each sample hence consisting of 11 measurements).

Non-parametric Bootstrap

- ▶ Consider the situation:
 - ▶ We have a sample $\mathbf{x} = (x_1, \dots, x_n)$ from an unknown distribution function F .
 - ▶ We wish to make inferences about a parameter $\theta = t(F)$ based on \mathbf{x} .
- ▶ Let \hat{F} be the empirical distribution function, which assigns the probability $1/n$ to each observed value $x_i, i = 1, \dots, n$.
- ▶ A bootstrap sample, $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$ is a random sample of size n (the size of the original sample) from \hat{F} , ie, is a sample of size n obtained, **with replacement**, from a population of n objects x_1, \dots, x_n .



Non-parametric Bootstrap: bootstrap estimation

- ▶ Based on data \mathbf{x} we obtain an estimate of θ , say $\hat{\theta} = t(\mathbf{x})$.
- ▶ Based on the bootstrap sample \mathbf{x}^* we obtain a new estimate for θ , say $\hat{\theta}^* = t(\mathbf{x}^*)$
- ▶ Repeating the procedure m times, the bootstrap estimate of θ is

$$\hat{\theta}_B = \frac{1}{m} \sum_{i=1}^m t(\mathbf{x}_i^*) .$$

Non-parametric Bootstrap: estimation of the standard error

- ▶ Bootstrap can be used to estimate the standard error of $\hat{\theta}$, $se(\hat{\theta})$. It is given by the standard error of the bootstrap estimate of $\hat{\theta}$, ie,

$$\hat{se}_B(\hat{\theta}) = \sqrt{\frac{1}{m-1} \sum_{i=1}^m [\hat{\theta}_i^* - \hat{\theta}_B]^2}.$$

Exercise 4

To illustrate the bootstrap procedure, let's bootstrap a small random sample:

3.12 0.00 1.57 19.67 0.22 2.20

1. Create 1000 replicates of size 6 with replacement.
2. Calculate the sample mean for each of the replicates.
3. Make a histogram and a normal quantile plot of the 1000 means. Make the density plot of the 1000 replicates. This is the bootstrap distribution.
4. Calculate the bootstrap estimates of the mean and the standard error.

Confidence interval for a correlation coefficient **Exercise 5**

(Adapted from Applied Statistics for Bioinformatics using R, Wim P. Krijnen)

Consider two sets of expression values of the MCM3 gene of the Golub *et al.* (1999) data. This data set is a gene expression data (3051 genes and 38 tumor mRNA samples) from the leukemia microarray study. This gene encodes for highly conserved mini-chromosome maintenance proteins (MCM) which are involved in the initiation of eukaryotic genome replication.

```
source("https://bioconductor.org/biocLite.R")
biocLite("multtest")
library(multtest)
data(golub)
x <- golub[2289,]; y <- golub[2430,]
cor(x,y)
[1] 0.6376217
```

- 1 Obtain a bootstrap sample from (x, y) , and compute the correlation coefficient for the bootstrap sample.
- 2 Repeat the procedure [1] $B=1000$ times.
- 3 From the sample of size $n = 38$ of the bootstrapped correlation coefficients obtain the 0.025 and 0.975 percentiles.
- 4 This pair is a bootstrap 95% confidence interval for the correlation coefficient ρ .

Compare with the interval obtained using normality assumption for the sampling distribution of the empirical correlation coefficient.

Non-parametric Bootstrap: test of hypothesis

- Suppose X follows an unknown distribution F , with expected value μ . We want to perform the test

$$H_0 : \mu = \mu_0 \quad \text{v.s.} \quad H_1 : \mu \neq \mu_0.$$

Non-parametric Bootstrap: test of hypothesis

- ▶ Suppose X follows an unknown distribution F , with expected value μ . We want to perform the test

$$H_0 : \mu = \mu_0 \quad \text{v.s.} \quad H_1 : \mu \neq \mu_0.$$

- ▶ Let X_1, \dots, X_n be a random sample from X . Let $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$, be the test statistic where \bar{X} , is the sample mean and S the sample standard deviation.

Non-parametric Bootstrap: test of hypothesis

- ▶ Suppose X follows an unknown distribution F , with expected value μ . We want to perform the test

$$H_0 : \mu = \mu_0 \quad \text{v.s.} \quad H_1 : \mu \neq \mu_0.$$

- ▶ Let X_1, \dots, X_n be a random sample from X . Let $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$, be the test statistic where \bar{X} , is the sample mean and S the sample standard deviation.
- ▶ Let $\mathbf{x} = (x_1, \dots, x_n)$ an observed sample and t_{obs} the observed value of the test statistic.

Non-parametric Bootstrap: test of hypothesis

- ▶ Suppose X follows an unknown distribution F , with expected value μ . We want to perform the test

$$H_0 : \mu = \mu_0 \quad \text{v.s.} \quad H_1 : \mu \neq \mu_0.$$

- ▶ Let X_1, \dots, X_n be a random sample from X . Let $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$, be the test statistic where \bar{X} , is the sample mean and S the sample standard deviation.
- ▶ Let $\mathbf{x} = (x_1, \dots, x_n)$ an observed sample and t_{obs} the observed value of the test statistic.
- ▶ Since we need an estimate of the sampling distribution of T under the null hypothesis, we cannot use the empirical distribution function \hat{F} since it does not obey H_0 .

Non-parametric Bootstrap: test of hypothesis

- ▶ One way of overcoming this problem is to apply a transformation to F such that the expected value is μ_0 . This can be done if we consider the empirical distribution function $\mathbf{z} = (z_1, \dots, z_n)$, where $z_i = x_i - \bar{x} + \mu_0, i = 1, \dots, n$.

Non-parametric Bootstrap: test of hypothesis

- ▶ One way of overcoming this problem is to apply a transformation to F such that the expected value is μ_0 . This can be done if we consider the empirical distribution function $\mathbf{z} = (z_1, \dots, z_n)$, where $z_i = x_i - \bar{x} + \mu_0, i = 1, \dots, n$.
- ▶ We obtain then m bootstrap samples from \mathbf{z} , $\mathbf{z}_1^*, \dots, \mathbf{z}_m^*$, their averages and standard deviations and with them we build a sequence of statistics $t_i^*, i = 1, \dots, m$ with

$$t_i^* = \frac{\bar{z}_i^* - \mu_0}{s_i^* / \sqrt{n}},$$

Non-parametric Bootstrap: test of hypothesis

- ▶ One way of overcoming this problem is to apply a transformation to F such that the expected value is μ_0 . This can be done if we consider the empirical distribution function $\mathbf{z} = (z_1, \dots, z_n)$, where $z_i = x_i - \bar{x} + \mu_0, i = 1, \dots, n$.
- ▶ We obtain then m bootstrap samples from \mathbf{z} , $\mathbf{z}_1^*, \dots, \mathbf{z}_m^*$, their averages and standard deviations and with them we build a sequence of statistics $t_i^*, i = 1, \dots, m$ with

$$t_i^* = \frac{\bar{z}_i^* - \mu_0}{s_i^* / \sqrt{n}},$$

- ▶ The p -value, $p = 2P_{H_0}(T > |t_{obs}|)$ is then approximated by

$$p \approx 2 \frac{\#\{i : |t_i^*| > |t_{obs}|\}}{m}.$$

Exercise 6: Gdf5 gene from the Golub et al. (1999) data.

(Adapted from Applied Statistics for Bioinformatics using R, Wim P. Krijnen)

- ▶ The corresponding expression values are contained in row 2058.
- ▶ A quick search through the NCBI site makes it likely that this gene is not directly related to leukemia.
- ▶ Hence, we may hypothesize that the population mean of the ALL expression values equals zero.
- ▶ Accordingly, we test $H_0 : \mu = 0$ v.s. $H_0 : \mu > 0$.
- ▶ a t test gives a p – value = 0.499 and clearly H_0 is not rejected.
- ▶ How can we use bootstrap to test the present hypothesis?

Non-parametric Bootstrap: test of hypothesis - two samples

- ▶ Suppose we have $X_1 \sim F_1$, with expected value μ_1 and $X_2 \sim F_2$, with expected value μ_2 , where X_1 is independent of X_2 and F_1 e F_2 are unknown.
- ▶ The objective is to test

$$H_0 : \mu_1 = \mu_2 \quad v.s. \quad H_1 : \mu_1 \neq \mu_2.$$

based on the test statistic

$$T = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S_1^2/n_1 + S_2^2/n_2}},$$

where (\bar{X}_k, S_k^2, n_k) , $k = 1, 2$ are, respectively, the sample means, sample variances and sample sizes of the samples \mathbf{X}_k , $k = 1, 2$, from the two populations.

Non-parametric Bootstrap: test of hypothesis - two samples

- ▶ Since we need the distribution of T under H_0 and this hypothesis assumes only the equality of the mean values we have to resample from samples which obey this condition. For that we do as follows:

1. Compute the combined mean of the two samples

$$\bar{x} = \frac{\sum_{i=1}^{n_1} x_{1i} + \sum_{j=1}^{n_2} x_{2j}}{n_1 + n_2},$$

2. Consider $z_{ki} = x_{ki} - \bar{x}_k + \bar{x}$, $i = 1, \dots, n_k$, $k = 1, 2$.

3. Obtain m bootstrap samples from $(z_{1i}, i = 1, \dots, n_1)$ and m bootstrap samples from $(z_{2i}, i = 1, \dots, n_2)$, and compute the respective sample means and sample variances.

$$(\bar{z}_{ki}^*, s_{ki}^{2*}), i = 1, \dots, m, k = 1, 2$$

4. Obtain the test statistic for each bootstrap samples

$$t_i^* = \frac{\bar{z}_1^* - \bar{z}_2^*}{\sqrt{s_1^{2*}/n_1 + s_2^{2*}/n_2}},$$

Non-parametric Bootstrap: test of hypothesis - two samples

- An approximation for the p -value $p = 2P_{H_0}(T > |t_{obs}|)$ is

$$p \approx 2 \frac{\#\{i : |t_i^*| > |t_{obs}|\}}{m}.$$

Exercise 7: gene CCND3 Cyclin D3

- ▶ Golub et al. (1999) argue that gene CCND3 Cyclin D3 plays an important role with respect to discriminating ALL from AML patients.
- ▶ We are then interested in testing the null hypothesis of equal means, ie, we want to test

$$H_0 : \mu_{ALL} = \mu_{AML} \quad \text{vs} H_1 : \mu_{ALL} \neq \mu_{AML}.$$

- ▶ A t test for the equality of means would give strong support for the rejection of H_0 .
- ▶ How would you implement a bootstrap test for this problem?

Non-parametric Bootstrap: permutation tests

- ▶ Another way of resampling in a manner that is consistent with the null hypothesis is to use permutation resampling.

Non-parametric Bootstrap: permutation tests

- ▶ Another way of resampling in a manner that is consistent with the null hypothesis is to use permutation resampling.
- ▶ Permutation resample from a sample is done by sampling without replacement.

Non-parametric Bootstrap: permutation tests

- ▶ Another way of resampling in a manner that is consistent with the null hypothesis is to use permutation resampling.
- ▶ Permutation resample from a sample is done by sampling **without replacement**.
- ▶ Choose from the combined sample n_1 elements without replacement as a sample of the first population; the other n_2 constitute the sample from the second population.

Non-parametric Bootstrap: permutation tests

- ▶ Another way of resampling in a manner that is consistent with the null hypothesis is to use permutation resampling.
- ▶ Permutation resample from a sample is done by sampling **without replacement**.
- ▶ Choose from the combined sample n_1 elements without replacement as a sample of the first population; the other n_2 constitute the sample from the second population.
- ▶ Repeat the procedure m times.

Non-parametric Bootstrap: permutation tests

- ▶ Another way of resampling in a manner that is consistent with the null hypothesis is to use permutation resampling.
- ▶ Permutation resample from a sample is done by sampling **without replacement**.
- ▶ Choose from the combined sample n_1 elements without replacement as a sample of the first population; the other n_2 constitute the sample from the second population.
- ▶ Repeat the procedure m times.
- ▶ The distribution of the statistic from these resamples estimates the sampling distribution under H_0 . Its is called a permutation distribution.

Non-parametric Bootstrap: permutation tests

- ▶ Another way of resampling in a manner that is consistent with the null hypothesis is to use permutation resampling.
- ▶ Permutation resample from a sample is done by sampling **without replacement**.
- ▶ Choose from the combined sample n_1 elements without replacement as a sample of the first population; the other n_2 constitute the sample from the second population.
- ▶ Repeat the procedure m times.
- ▶ The distribution of the statistic from these resamples estimates the sampling distribution under H_0 . Its is called a permutation distribution.
- ▶ Locate the observed value of the test statistic under consideration in this sampling distribution to get the p -value.

Non-parametric Bootstrap: General procedure for permutation tests

To carry out a permutation test based on a statistic that measures the size of an effect of interest do as follows

- ▶ Compute the statistic for the original data.

Non-parametric Bootstrap: General procedure for permutation tests

To carry out a permutation test based on a statistic that measures the size of an effect of interest do as follows

- ▶ Compute the statistic for the original data.
- ▶ Choose permutation resamples from the data without replacement in a way that is consistent with the null hypothesis of the test and with the study design.

Non-parametric Bootstrap: General procedure for permutation tests

To carry out a permutation test based on a statistic that measures the size of an effect of interest do as follows

- ▶ Compute the statistic for the original data.
- ▶ Choose permutation resamples from the data without replacement in a way that is consistent with the null hypothesis of the test and with the study design.
- ▶ Construct the permutation distribution of the statistic from its values in a large number of resamples.

Non-parametric Bootstrap: General procedure for permutation tests

To carry out a permutation test based on a statistic that measures the size of an effect of interest do as follows

- ▶ Compute the statistic for the original data.
- ▶ Choose permutation resamples from the data without replacement in a way that is consistent with the null hypothesis of the test and with the study design.
- ▶ Construct the permutation distribution of the statistic from its values in a large number of resamples.
- ▶ Find the p -value by locating the observed value of the statistic on the permutation distribution.

When Can We Use Permutation Tests?

- ▶ We can use a permutation test only when we can see how to resample in a way that is consistent with the study design and with the null hypothesis.

We now know how to do this for the following types of problems:

- ▶ Two-sample problems when the null hypothesis says that the two populations are identical.
- ▶ Matched pairs designs when the null hypothesis says that there are only random differences within pairs.
- ▶ Relationships between two quantitative variables when the null hypothesis says that the variables are not related. The correlation is the most common measure.

Exercise 8: permutation tests

Use the data in exercise 7 and use a permutation test instead.

Parametric vs non-parametric

- ▶ When we sample with replacement from the data to the non-parametric bootstrap, we are simply taking a random sample from the empirical distribution.
- ▶ The empirical distribution is the maximum likelihood estimate of the population distribution.
- ▶ The parametric bootstrap is simply sampling from the fitted model.

Parametric vs non-parametric

- ▶ Simulation from empirical distribution is computationally cheap and easier to implement;
- ▶ simulation from the fitted model may be computationally expensive and/or difficult to implement;

Anyway, one of the main motivations of the bootstrap was a desire to shake off the restrictions of conventional parametric modelling.

- ▶ Acknowledgements:
Antónia Turkman (DEIO-FCUL), for allowing the use of some material produced by us in previous courses.
- ▶ A good book for introducing bootstrap methods: Efron, B and Tibshirani (1993) An Introduction to the Bootstrap, New York: Chapman and Hall.