



Resampling methods

MODULE DES130: COMPUTATIONAL STATISTICS

Dr. Erick A. Chacón Montalván (echacon@uni.edu.pe)

Escuela de Profesional de Ingeniería Estadística
Facultad de Ingeniería Económica, Estadística y Ciencias Sociales
Universidad Nacional de Ingeniería (UNI)
Lima – Perú

Introduction

Introduction

Basic concepts

Resampling for statistical inference

Motivation

Jackknife

Bootstrap

Permutation tests

Concept

Resampling is the procedure of generating new samples using observed data or a data generating mechanism.

Resampling is the same as:

- Given a set of observations x_1, \dots, x_n ; use the empirical cumulative density function \hat{F} as an estimate of the cumulative density function F
- Generate random values x_1^*, \dots, x_n^* coming from \hat{F} .

Common types

- Jackknife: Takes one-leave-out samples.
- Bootstrap: Take random samples with repetition.
- Permutations: Relabel the data.

Introduction

Basic concepts

Resampling for statistical inference

Motivation

Jackknife

Bootstrap

Permutation tests

Resampling is used with the goal of:

- Approximate the distribution of the statistics of interest.
- Estimate the standard error of an statistics.
- Compute confidence intervals.
- Perform hypothesis testing.

Resampling and monte carlo

Let consider a quantity of interest μ which can be written as the mean of $h(X)$ where X is a random variable with density function $f(x)$,

$$\mu = \int h(x)f(x)dx \approx \int h(x)\hat{f}(x)dx.$$

Considering the random sample $X_1^*, X_2^*, \dots, X_m^*$ of size m with density $\hat{f}(x)$, we define the resampling estimator

$$\hat{\mu}_{RE} = \frac{\sum_{i=1}^m h(X_i^*)}{m}.$$

The estimator improves as long as $m \rightarrow \infty$.

Introduction

Basic concepts

Resampling for statistical inference

Motivation

Jackknife

Bootstrap

Permutation tests

Motivation

- Classical statistical inference includes **asymptotic and distributional assumptions**.
- Asymptotic assumptions do not hold for **small datasets**.
- Uncertainty quantification is not obvious in **non-linear and hierarchical models**.
- Some statistics tests have **unknown analytical probability or density function**.
- Monte carlo methods **assumes a known cumulative density function**.

Jacknife

Introduction

Jackknife

 Estandar error

 Jackknife estimate

Bootstrap

Permutation tests

Let $\hat{\theta}$ be an estimator of θ . It is customary to work with the pivotal statistics

$$\frac{\hat{\theta} - \theta}{se(\hat{\theta})} \sim G(\cdot).$$

Using this pivotal statistics, we can obtain confidence interval or perform hypothesis testing, but we need a good computation of $se(\hat{\theta})$. Jackknife provides an alternative way to approximate $se(\hat{\theta})$.

Introduction

Jackknife

Estandar error

Jackknife estimate

Bootstrap

Permutation tests

Jackknife estimate

Let consider a random sample $\mathbf{X} = (X_1, X_2, \dots, X_n)$ with realization $\mathbf{x} = (x_1, x_2, \dots, x_n)$ such as $X_j \sim F$. Let $\hat{\theta}(\mathbf{X})$ be a statistic of and

$$\mathbf{x}_{(i)} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n),$$

such as

$$\hat{\theta}_{(i)} = \hat{\theta}(\mathbf{x}_{(i)}) \text{ for } i = 1, \dots, n.$$

Then the Jackknife estimation of the standard error is

$$\hat{s}e_{jack}(\hat{\theta}) = \left[\frac{n-1}{n} \sum_{i=1}^n (\hat{\theta}_{(i)} - \hat{\theta}_{(\cdot)})^2 \right]^{1/2}, \text{ with } \hat{\theta}_{(\cdot)} = \sum_{i=1}^n \hat{\theta}_{(i)} / n.$$

For $\hat{\theta} = \bar{X}$, it can be shown that

$$\widehat{se}_{jack}(\hat{\theta}) = \left[\sum_{i=1}^n (x_i - \bar{x})^2 / (n(n-1)) \right]^{1/2},$$

because $\hat{\theta}_{(i)} = (n\bar{x} - x_i)/(n-1)$.

Example for correlation

$$\hat{\theta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{[\sum_{i=1}^n (x_i - \bar{x}) \sum_{i=1}^n (y_i - \bar{y})]^{1/2}}$$

- The advantage of the jackknife estimate is that $\hat{\theta}$ could be any estimator or statistic of interest, and not only the sample average.
- It is nonparametric because nothing is assumed about F .
- Automatic.

Bootstrap

Introduction

Jackknife

Bootstrap

Introduction

Non parametric bootstrap

Parametric bootstrap

Bootstrap bias

Bootstrap regression

Bootstrap inference

Permutation tests

Example: sample average

TABLE 9.1 Possible bootstrap pseudo-datasets from $\{1, 2, 6\}$ (ignoring order), the resulting values of $\hat{\theta}^* = T(\hat{F}^*)$, the probability of each outcome in the bootstrapping experiment ($P^*[\hat{\theta}^*]$), and the observed relative frequency in 1000 bootstrap iterations.

λ^*	$\hat{\theta}^*$	$P^*[\hat{\theta}^*]$	Observed Frequency
1 1 1	3/3	1/27	36/1000
1 1 2	4/3	3/27	101/1000
1 2 2	5/3	3/27	123/1000
2 2 2	6/3	1/27	25/1000
1 1 6	8/3	3/27	104/1000
1 2 6	9/3	6/27	227/1000
2 2 6	10/3	3/27	131/1000
1 6 6	13/3	3/27	111/1000
2 6 6	14/3	3/27	102/1000
6 6 6	18/3	1/27	40/1000

A 93% confidence interval for θ is $[4/3, 14/3]$.

Theoretical concepts:

- F : Cumulative density function.
- X_i : A **random variable** with some associated density f and cumulative F functions.
- $\mathbf{X} = (X_1, \dots, X_n)^T$: A **random vector** representing a random sample such as $X_i \sim f$ with the cumulative density function F .
- $\theta = g(F)$: A parameter of interest of F . For example $\theta = \int x dF(x)$.
- $T(\mathbf{X}, F)$: An statistics of interest.

Empirical concepts:

- $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top$: Observed data, which is a realization of \mathbf{X} .
- \hat{F} : The empirical cumulative distribution function obtained using \mathbf{x} .
- $\mathbf{X}^* = (X_1^*, \dots, X_n^*)^\top$: A **random vector** representing a bootstrap sample such as X_i^* have the cumulative density function \hat{F} .
- $\hat{\theta} = g(\hat{F})$: The estimator of θ using \hat{F} . For example, $\hat{\theta} = \sum X_i/n$.
- $T(\mathbf{X}^*, \hat{F})$: An statistics of interest.

Bootstrap

The goal of bootstrap is to approximate the distribution of $T(\mathbf{X}, F)$ using the distribution of $T(\mathbf{X}^*, \hat{F})$.

Once we obtain the distribution of $T(\mathbf{X}^*, \hat{F})$, then we can estimate the standard error, provide confidence intervals or perform hypothesis testing.

Introduction

Jackknife

Bootstrap

Introduction

Non parametric bootstrap

Parametric bootstrap

Bootstrap bias

Bootstrap regression

Bootstrap inference

Permutation tests

In most applications, we will not be able to work with all the possible resamples. Lets define a number of resamples B such as $\mathbf{X}_i^* \sim \hat{F}$ for $i = 1, \dots, B$. Then the distribution of $T(\mathbf{X}, F)$ is approximated by the empirical distribution of $T(\mathbf{X}_i^*, \hat{F})$ for $i = 1, \dots, B$.

Introduction

Jackknife

Bootstrap

Introduction

Non parametric bootstrap

Parametric bootstrap

Bootstrap bias

Bootstrap regression

Bootstrap inference

Permutation tests

In case that data modelled comes from a parametric distribution $F(\mathbf{x}, \theta)$, then we can use $F(\mathbf{x}, \hat{\theta})$ to generate \mathbf{X}_i^* for $i = 1, \dots, B$.

Introduction

Jackknife

Bootstrap

Introduction

Non parametric bootstrap

Parametric bootstrap

Bootstrap bias

Bootstrap regression

Bootstrap inference

Permutation tests

A quantity of interest can be

$$T(\mathbf{X}, F) = g(\hat{F}) - g(F),$$

which has mean $\mathbb{E} [\hat{\theta}] - \theta$ representing the bias of $g(\hat{F})$. The bootstrap bias estimate is

$$\sum_{i=1}^B (\hat{\theta}(\mathbf{x}_i^*) - \hat{\theta}(\mathbf{x})) / B = \bar{\theta}^* - \hat{\theta}(\mathbf{x}).$$

Introduction

Jackknife

Bootstrap

Introduction

Non parametric bootstrap

Parametric bootstrap

Bootstrap bias

Bootstrap regression

Bootstrap inference

Permutation tests

$$z_i = (y_i, \mathbf{x}_i). \quad (1)$$

Introduction

Jackknife

Bootstrap

Introduction

Non parametric bootstrap

Parametric bootstrap

Bootstrap bias

Bootstrap regression

Bootstrap inference

Permutation tests

$1 - \alpha$ confidence interval

$$\Pr \left(Q_{\alpha/2}(\hat{\theta}(\mathbf{X}^*)) \leq \theta \leq Q_{1-\alpha/2}(\hat{\theta}(\mathbf{X}^*)) \right) = 1 - \alpha.$$

Justification of the percentile method

Let consider a strictly increasing transformation ϕ and a continuous and symmetric distribution function H such as:

$$\Pr \left(h_{\alpha/2} \leq \phi(\hat{\theta}) - \phi(\theta) \leq h_{1-\alpha/2} \right) = 1 - \alpha,$$
$$\Pr \left(\phi^{-1}(h_{\alpha/2} + \phi(\hat{\theta})) \leq \theta \leq \phi^{-1}(h_{1-\alpha/2} + \phi(\hat{\theta})) \right) = 1 - \alpha.$$

Similarly

$$\Pr \left(h_{\alpha/2} \leq \phi(\hat{\theta}) - \phi(\hat{\theta}(\mathbf{X}^*)) \leq h_{1-\alpha/2} \right) = 1 - \alpha,$$
$$\Pr \left(\phi^{-1}(h_{\alpha/2} + \phi(\hat{\theta})) \leq \hat{\theta}(\mathbf{X}^*) \leq \phi^{-1}(h_{1-\alpha/2} + \phi(\hat{\theta})) \right) = 1 - \alpha.$$

As an approximation we can obtain that $\phi^{-1}(h_{\alpha/2} + \phi(\hat{\theta}))$ is the $\alpha/2$ quantile of the bootstrap estimator. The approximation is good if $\phi(\hat{\theta})$ is unbiased and its variance does not depend on θ .

Bootstrap t (studentized)

Let consider $T(\mathbf{X}, F) = \frac{g(\hat{F}) - g(F)}{\sqrt{\mathbb{V}[g(\hat{F})]}}$ with cumulative distribution function G , it will be roughly pivotal.

$$\Pr \left(Q_{\alpha/2, G} \leq T(\mathbf{X}, F) \leq Q_{1-\alpha/2, G} \right) = 1 - \alpha,$$

$$\Pr \left(\hat{\theta} - \mathbb{V} [g(\hat{F})] Q_{1-\alpha/2, G} \leq \theta \leq \hat{\theta} - \mathbb{V} [g(\hat{F})] Q_{\alpha/2, G} \right) = 1 - \alpha,$$

Using bootstrap

$$\Pr \left(\hat{\theta}^* - \mathbb{V} [g(\hat{F})] Q_{1-\alpha/2, \hat{G}} \leq \theta \leq \hat{\theta}^* - \mathbb{V} [g(\hat{F})] Q_{\alpha/2, \hat{G}} \right) = 1 - \alpha.$$

Permutation tests

Hypothesis testing