

# Exploratory Data Analysis Tools

Stelian Ion\*

Technical Reports

## Abstract

In this paper we review some mathematical tools to analyze ecological data. We focus on cluster analysis and linear regression analysis. Cluster analysis is the organization of a collection of objects into cluster based on similarity. The clustering process involves the following steps: object representation, object proximity, grouping step, validation of the clusters. As any clustering algorithm produces a cluster structure, is fundamental to know if the structure is valid or not and how many clusters are there. We address this problems by considering the cophenetic correlation coefficient, statistical significance of the cophenetic index and the stable association with respect to the number of the species in the ecological data matrix. A mathematical issue on multiple linear regression is to estimate the parameters in the linear regression disposing on a set of the measurements of the variables in the model. We discuss the bootstrap and jackknife resampling techniques in view to estimate the parameters, covariance matrix and confidence interval.

## Contents

<b>1</b>	<b>Cluster analysis</b>	<b>2</b>
1.1	Object representation and feature selection . . . . .	2
1.2	Object proximity . . . . .	2
1.3	Grouping step . . . . .	2
1.4	Validation of the clusters . . . . .	3
1.5	Internal validity index . . . . .	3
1.5.1	Silhouette coefficients . . . . .	3
1.5.2	Cophenetic correlation coefficient . . . . .	4
1.6	Stable associations . . . . .	4
<b>2</b>	<b>Multiple regression</b>	<b>5</b>
2.1	Model equation . . . . .	5
2.2	Ordinary least squares estimator . . . . .	6
2.3	Confidence region of regression coefficients. Normal case . . . . .	9
2.4	Bootstrap Confidence Intervals . . . . .	16
2.5	BC <sub>a</sub> method . . . . .	17

---

\*Institute of Statistical Mathematics and Applied Mathematics, Calea 13 Septembrie 13, Bucharest, Romania.

## 1 Cluster analysis

Cluster analysis consists in organizing a collection of objects into clusters based on similarity. Intuitively, objects belonging to same cluster are more similar to each other than they are to an object from a different cluster. The clustering process involves the following steps:

1. object representation,
2. object proximity,
3. grouping step,
4. validation of the clusters.

### 1.1 Object representation and feature selection

The feature selection is a process of identifying the most relevant features to be used in the cluster analysis. Let  $N$  be the number of objects and  $L$  be the number of selected features. Let us denote by  $\mathbf{X}^\alpha$ ,  $\alpha = \overline{1, N}$  the objects, by  $\mathcal{O}$  the set of objects and by  $\mathcal{Y}$  the data matrix whose elements  $y_i^\alpha$  quantify the values of  $i$ -feature of the  $\alpha$ -object. One can think each object  $\mathbf{X}^\alpha$  as a point in a mathematical  $L$ -dimensional space.

### 1.2 Object proximity

The similarity is measured by a function defined on pairs of objects

$$S : \mathcal{O} \times \mathcal{O} \rightarrow \mathbb{R},$$

$s_{\alpha, \beta}$  define the index of similarity of the two objects  $\mathbf{X}^\alpha$  and  $\mathbf{X}^\beta$ . The dissimilarity of two objects is given by  $d_{\alpha, \beta} = 1 - s_{\alpha, \beta}$  and the proximity matrix contains the dissimilarity between each pair of objects.

The data matrix  $\mathcal{Y}$  and the proximity matrix  $D$  are the base of the clustering process.

### 1.3 Grouping step

Next step in cluster analysis is to select a method to partition the set of the object  $\mathcal{O}$  into subsets. Agglomerative hierarchical clustering algorithms produce a nested series of partitions, called dendrogram, based on a criterion for merging cluster by using similarity. There are three main hierarchical clustering algorithms: single-link [19], complete-link [7], and minimum-variance [21] and [15]. The first two differ in the way they characterize the similarity between two clusters. If in the single-link method, the distance between two clusters is given by the *minimum* of the distances between all pairs of objects belonging the two different clusters in the complete-link method the distance between the two clusters is given by the *maximum* of distance of all pairs of objects.

The results of the hierarchical cluster algorithm is a tree of clusters and the similarity levels at which grouping change. Each dendrogram can be translated into a *cophenetic* matrix [11]. For any two objects  $\mathbf{X}^\alpha$  and  $\mathbf{X}^\beta$  the cophenetic distance  $h_{\alpha\beta}$  equals the level at which the two objects belong to same cluster for the first time. The cophenetic matrix is the matrix with elements  $h_{\alpha,\beta}$ ,  $1 \leq \alpha, \beta \leq N$ .

#### 1.4 Validation of the clusters

Any clustering algorithm has a cluster structure as output for any data set. Consequently, is fundamental to know if such structure is valid or not. Cluster validity is a procedure of evaluating the results of a cluster algorithm. The quality of the structure is quantified by a cluster validity index that gives an information about the possibility that the structure have occurred by chance or is a relevant result of the clustering algorithm. There are three main methods to investigate cluster validity:

1. *external criteria methods*,
2. *internal criteria methods*, and
3. *relative criteria methods*.

The internal examination of validity evaluates the result of the clustering process using only the quantities and features existent in the data matrix.

#### 1.5 Internal validity index

A lot of internal validity indexes were defined (see [8], [3], for example).

Let  $\mathcal{P} = \{C_a\}_{a=\overline{1,k}}$  be a partition of the set  $\mathcal{O}$ , each subset  $C_a$  represents a cluster and one denotes by  $|C_a|$  the number of elements of cluster  $C_a$ . Corresponding to partition  $\mathcal{P}$ , the membership function  $\mathcal{U}_{\mathcal{P}}$ :

$$\mathcal{U}_{\mathcal{P}} : \mathcal{O} \rightarrow \{0, 1\}^k$$

given by

$$u_{\mathcal{P}}^{\alpha a} = \begin{cases} 1, & \text{if } \mathbf{X}^\alpha \in C_a \\ 0, & \text{otherwise} \end{cases}$$

is defined, together with a function that indicates the cluster to which an object belongs

$$v_{\mathcal{P}} : \mathcal{O} \rightarrow \{0, 1, \dots, .k\}$$

$$v_{\mathcal{P}}(\mathbf{X}^\alpha) = b \text{ iff } \mathbf{X}^\alpha \in C_b$$

##### 1.5.1 Silhouette coefficients

The *average dissimilarity of a object and a cluster* is defined by

$$\bar{d}(\mathbf{X}^\alpha, C_a) = \frac{1}{|C_a|} \sum_{\beta} u_{\mathcal{P}}^{\beta a} d_{\alpha\beta}. \quad (1)$$

The quantity  $a_{\mathcal{P}}(\mathbf{X}^\alpha) = \bar{d}(\mathbf{X}^\alpha, C_{v_{\mathcal{P}}(\mathbf{X}^\alpha)})$  can be viewed as a measure of the cohesion of the cluster  $C_{v_{\mathcal{P}}(\mathbf{X}^\alpha)}$  with respect to  $\mathbf{X}^\alpha$ . The separation of the object  $\mathbf{X}^\alpha$  from the clusters that do not contain it can be evaluated by

$$b_{\mathcal{P}}(\mathbf{X}^\alpha) = \min_{a \neq v_{\mathcal{P}}(\mathbf{X}^\alpha)} \bar{d}(\mathbf{X}^\alpha, C_a)$$

and the the silhouette of an object  $\mathbf{X}^\alpha$  is defined by:

$$s_{\mathcal{P}}(\mathbf{X}^\alpha) = \frac{b_{\mathcal{P}}(\mathbf{X}^\alpha) - a_{\mathcal{P}}(\mathbf{X}^\alpha)}{\max\{a_{\mathcal{P}}(\mathbf{X}^\alpha), b_{\mathcal{P}}(\mathbf{X}^\alpha)\}}. \quad (2)$$

An overall quality measure for a given partition is its average silhouette coefficient:

$$s_{\mathcal{P}} = \frac{1}{N} \sum_{\alpha} s_{\mathcal{P}}(\mathbf{X}^\alpha). \quad (3)$$

### 1.5.2 Cophenetic correlation coefficient

Let  $(h_{\alpha,\beta})_{\alpha,\beta}$  be the cophenetic matrix associated to a dendrogram and  $(d_{\alpha,\beta})_{\alpha,\beta}$  the dissimilarity matrix. The cophenetic index is the Mantel normalized statistic of the two matrices

$$Z(h, d) = \frac{\sum_{\alpha,\beta,\alpha<\beta} (h_{\alpha,\beta} - \bar{h}) (d_{\alpha,\beta} - \bar{d})}{\sqrt{\left(\sum_{\alpha,\beta,\alpha<\beta} (h_{\alpha,\beta} - \bar{h})^2\right) \left(\sum_{\alpha,\beta,\alpha<\beta} (d_{\alpha,\beta} - \bar{d})^2\right)}}$$

where the summation is over  $\alpha$  and  $\beta$  with  $\alpha$  ranging from one to  $N - 1$ ,  $\bar{h}$  and  $\bar{d}$  are the means of the matrices  $h$  and  $d$  respectively.

We accept a dendrogram as valid if the cophenetic index is close to one and its value is statistical significant at a given level  $p_0$ .

We test the statistical significance of the cophenetic index  $Z(h, d)$  by using the Monte Carlo method. We generate a number,  $n$ , of *random dendrograms* and for each such dendrogram we calculate the cophenetic matrix  $h_i$  and the cophenetic index  $Z(h_i, d)$ . Then we calculate the fraction of the cophenetic indices greater than  $Z(h, d)$

$$p = \frac{|\{i | Z(h_i, d) > Z(h, d), i = \overline{1, n}\}|}{n}$$

where for any finite set  $A$ ,  $|A|$  means the number of the elements of the set  $A$ .  $p$  measures the probability to obtain a value of cophenetic index greater then  $Z(h, d)$  by chance.

### 1.6 Stable associations

A subset  $\mathcal{M} \subset \mathcal{O}$  is named a *stable association* with respect to some perturbed process if for any perturbation the set is still a cluster of the hierarchical structure.

We assume here that we study a series of data taken at different moment of time. We investigate the stable associations by considering variable number of species (hence the perturbation of the set will consists in the occurrence or absence of some species). Let  $m$  be an integer number greater than 1 and smaller than the number of time observations.

Let us denote

$$\mathcal{O}_m = \{\mathbf{X}^\alpha | y_i^\alpha > 0; \text{for at least } m \text{ values of } i\},$$

and by  $\mathcal{Y}_m$  denote the corresponding data matrix.

It is obvious that the  $\mathcal{O}_m$  is a descending sequence with respect to  $m$ ,  $\mathcal{O}_m \subset \mathcal{O}_{m-1} \subset \dots \subset \mathcal{O}_1 = \mathcal{O}$ .

## 2 Multiple regression

The mathematical problem of multiple linear regression is to estimate the parameters in the linear regression disposing on a set of measurements of the variables in the model. We restrict our presentation to the case of one dependent variable  $\mathbf{y}$  and  $p$  independent variables  $\{\mathbf{x}^a\}_{a=\overline{1,p}}$ .

**Set of measures.** To find the parameters one makes a number of observations on the *state of the system*. Each observation records the measured value of the dependent variable  $\mathbf{y}$  and the set of the measured values of the independent variables  $\{\mathbf{x}^a\}_{a=\overline{1,p}}$ . One assumes that the values of the independent variables are measured without errors but there can exist errors in the measured values of the dependent variables. Thus the dependent variable is a random variable while the independent variables are not random variables. We define  $\mathbf{y} | \{\mathbf{x}^a\}_{a=\overline{1,p}}$  to be the random variable  $\mathbf{y}$  corresponding to the fixed values  $\{\mathbf{x}^a\}_{a=\overline{1,p}}$  and denote by  $E(\mathbf{y} | \{\mathbf{x}^a\}_{a=\overline{1,p}})$  its mean.

Let us denote by  $N$  the number of the observations, by  $\tilde{y}_i$  the measured value of the dependent variable corresponding to the  $i$ -th observation, by  $x_i^a$  the values of the  $a$ -variable,  $a = \overline{1,p}$  corresponding to the  $i$ -th observation,  $i = \overline{1,N}$ .

### 2.1 Model equation

In a linear regression model one assumes that the mean of dependent variable is linearly related to the independent variable by *population regression equation*:

$$E(\mathbf{y} | \{\mathbf{x}^a\}_{a=\overline{1,p}}) = \beta_0 + \sum_{a=1}^p \beta_a \mathbf{x}^a l y$$

A random variable  $\mathbf{y} | \{\mathbf{x}^a\}_{a=\overline{1,p}}$  may be described by the multiple linear regression model

$$\mathbf{y} = E(\mathbf{y} | \{\mathbf{x}^a\}_{a=\overline{1,p}}) + \epsilon = \beta_0 + \sum_{a=1}^p \beta_a \mathbf{x}^a + \epsilon, \quad (4)$$

where the random errors  $\epsilon$  must have zero mean.

Each observation  $(\tilde{y}_i, x_i^1, x_i^2, \dots, x_i^p)$  satisfies the equation

$$\tilde{y}_i = \beta_0 + \sum_{\alpha=1}^p \beta_\alpha x_i^\alpha + \epsilon_i, \quad (5)$$

for some values  $\epsilon_i$  of the random error variable. We assume that the random errors  $\epsilon_i$  are independent distributed and each of them has zero mean and variance  $\sigma_i^2$

$$E(\epsilon_i) = 0, \text{Var}(\epsilon_i) = \sigma_i^2. \quad (6)$$

## 2.2 Ordinary least squares estimator

Given the set of observations  $\{\tilde{y}_i, \{x_i^a\}_{a=1, \dots, p}\}_{i=1, \dots, N}$  the ordinary least squares method estimates the parameters  $\boldsymbol{\beta}$  as the point minimizer of the *merit function*  $\phi$  defined by

$$\phi(\boldsymbol{\eta}) = \sum_{i=1}^N \left( \tilde{y}_i - \left( \eta_0 + \sum_{a=1}^p \eta_a x_i^a \right) \right)^2. \quad (7)$$

Consequently the point minimizer  $\hat{\boldsymbol{\beta}}$  satisfies the equations

$$\frac{\partial \phi(\boldsymbol{\eta})}{\partial \eta_a} = 0, a = \overline{0, p}$$

which is the *normal equation*

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\eta} = \mathbf{X}^T \tilde{\mathbf{y}}, \quad (8)$$

where the matrix  $\mathbf{X}$ , the *design matrix*, and the vector  $\tilde{\mathbf{y}}$  are given by

$$\mathbf{X} = \begin{pmatrix} 1 & x_1^1 & \cdots & x_1^p \\ 1 & x_2^1 & \cdots & x_2^p \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N^1 & \cdots & x_N^p \end{pmatrix}, \quad \tilde{\mathbf{y}} = \begin{pmatrix} \tilde{y}_1 \\ \tilde{y}_2 \\ \vdots \\ \tilde{y}_N \end{pmatrix}$$

respectively. If  $\mathbf{X}^T \mathbf{X}$  is an invertible matrix then the OLS estimator is given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tilde{\mathbf{y}}. \quad (9)$$

### SVD method for OLS estimation

The *singular value decomposition* (SVD) method allows one to find an estimate  $\hat{\boldsymbol{\beta}}$  of true parameter  $\boldsymbol{\beta}$  in population regression equation by using the SVD of design matrix. Theoretically the method gives the same estimate as that found by solving the normal equation.

For convenience let us introduce some notations. We introduce the scalar product on  $\mathbb{R}^N$  by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^N x_i y_i \quad (10)$$

and on  $\mathbb{R}^{p+1}$  by

$$(\boldsymbol{\eta}, \boldsymbol{\gamma}) = \sum_{a=0}^p \eta_a \gamma_a \quad (11)$$

We denote by  $X_i$  the element of the space  $\mathbb{R}^{p+1}$  that has the components  $X_i^0 = 1, X_i^a = x_i^a, a = \overline{1, p}$ . Instead of resolving the normal equation one tries to solve the system of equations

$$\eta_0 + \sum_{a=1}^p \eta_a x_i^a = \tilde{y}_i, \quad i = \overline{1, N}$$

with  $\boldsymbol{\eta}$  as unknown. In matrix notation

$$\mathbf{X}\boldsymbol{\eta} = \tilde{\mathbf{y}}.$$

The system is resolved by using the singular values decomposition of the design matrix  $\mathbf{X}$ .

$$\mathbf{X} = \mathbf{U}\mathbf{W}\mathbf{V}^T; \quad X_i^a = \sum_{b,c=0}^p u_i^b w_{bc} v^{ca}, \quad i = \overline{1, N}, a = \overline{0, p}. \quad (12)$$

The matrix  $\mathbf{U}, \mathbf{W}$  and  $\mathbf{V}$  have the following properties

$$\langle U^a, U^b \rangle = \delta^{ab}, \quad \langle V^a, V^b \rangle = \delta^{ab}, \quad W_{ab} = w_{bb} \delta_{ab}; \quad a, b = \overline{0, p}.$$

The solution is given by

$$\hat{\boldsymbol{\beta}} = \mathbf{V}\mathbf{W}^{-1}\mathbf{U}^T\tilde{\mathbf{y}}, \quad (13)$$

or, componentwise

$$\hat{\beta}_a = \sum_{b=0}^p \frac{v^{ab}}{w_{bb}} \langle U^b, \tilde{\mathbf{y}} \rangle \quad (14)$$

For brevity we introduce the matrix

$$\mathbf{R} = \mathbf{V}\mathbf{W}^{-1}\mathbf{U}^T$$

and note that

$$\mathbf{R} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad \text{and} \quad \mathbf{R}\mathbf{X} = \mathbf{I}_{p+1}.$$

The the fitted multiple linear regression model reads as

$$\hat{\mathbf{y}}(\mathbf{x}) = \hat{\beta}_0 + \sum_{a=1}^p \hat{\beta}_a \mathbf{x}^a \quad (15)$$

Each observation  $(\tilde{y}_i, x_i^1, x_i^2, \dots, x_i^p)$  satisfies the equation

$$\tilde{y}_i = \hat{\beta}_0 + \sum_{\alpha=1}^p \hat{\beta}_\alpha x_i^\alpha + e_i, \quad (16)$$

where  $e_i$  is the fitted error at the  $i$  observation.

*Remark.* In the regression analysis one deals with two kinds of errors, the fitted errors  $e_i$  that measure how well the hyperplan  $\mathcal{H}_{\hat{\boldsymbol{\beta}}}$  interpolates the  $N$  points of observations  $(\tilde{y}_i, \mathbf{x}_i), i = \overline{1, N}$ , and model errors  $\epsilon_i$  which, for a given state  $\mathbf{x}$ , measure the departure of the dependent variable  $y$  from the linear model.

## Exploratory Data Analysis Tools

---

There exists an explicit relation between estimated parameters  $\hat{\boldsymbol{\beta}}$  and true parameters  $\boldsymbol{\beta}$  in the linear model that involves the random errors

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = \mathbf{R}\boldsymbol{\epsilon}. \quad (17)$$

Theoretically if one knows the errors, then one can precisely determine the true parameters! Another interesting relation is between the errors in the model and fitted errors, namely

$$\mathbf{e} = (\mathbf{I} - \mathbf{U}\mathbf{U}^T) \boldsymbol{\epsilon}. \quad (18)$$

Unfortunately, relation (18) does not allow us to know the errors.

### Prediction error, variance and covariance of parameters

Note that the OLS estimates  $\{\hat{\beta}_a\}_{a=0,p}$  are dependent on the set of observations. If we take another set of  $N$  observations keeping the values of the independent variables, same design matrix  $\mathbf{X}$ , we measure different values of dependent variables  $\mathbf{y}$  and consequently we obtain different values of the parameters. One can think the different estimates as the values of the random variables  $\{\mathcal{B}_a\}_{a=0,p}$  defined by

$$\mathcal{B} = \mathbf{R}\mathbf{y}. \quad (19)$$

In the sequel the media and the covariance matrix of  $\mathcal{B}$  must be read as conditional quantities in the sense "for a given state  $\mathbf{X}$ "

$$E(\mathcal{B}) = \boldsymbol{\beta} + \mathbf{R}E(\boldsymbol{\epsilon}) = \boldsymbol{\beta},$$

which shows that  $\mathcal{B}$  is an *unbiased* estimator of  $\boldsymbol{\beta}$ . The covariance of the random parameters is given by

$$\begin{aligned} \text{Cov}(\mathcal{B}_a, \mathcal{B}_b) &= \sum_{i,j=1}^N R_{ai}R_{bj} \text{Cov}(\varepsilon_i, \varepsilon_j) \\ &= \sum_i^N R_{ai}R_{bi}\sigma_i^2. \end{aligned} \quad (20)$$

In the case of equal variance  $\sigma_i = \sigma$

$$\text{Cov}(|\beta_a, |\beta_b) = \sigma^2 \sum_{c=0}^p \frac{v^{ac}v^{bc}}{w_{cc}^2}, \quad (21)$$

or

$$\text{Cov}(|\beta_a, |\beta_b) = \sigma^2 (\mathbf{X}^T \mathbf{X})_{ab}^{-1}. \quad (22)$$

We introduce the *prediction error* as the total sum of square of fitted errors:

$$\mathcal{E} := \sum_i (\tilde{y}_i - \hat{y}_i)^2.$$

Taking into account the relation (18) and using

$$(\mathbf{I} - \mathbf{U}\mathbf{U}^T) (\mathbf{I} - \mathbf{U}\mathbf{U}^T)^T = \mathbf{I} - \mathbf{U}\mathbf{U}^T,$$



one gets

$$\mathcal{E} = \langle \mathbf{e}, \mathbf{e} \rangle = \langle (\mathbf{I} - \mathbf{U}\mathbf{U}^T) \boldsymbol{\varepsilon}, (\mathbf{I} - \mathbf{U}\mathbf{U}^T) \boldsymbol{\varepsilon} \rangle = \langle \boldsymbol{\varepsilon}, (\mathbf{I} - \mathbf{U}\mathbf{U}^T) \boldsymbol{\varepsilon} \rangle.$$

The predicted errors  $\mathcal{E}$  can be interpreted as random variables whose values depend on the random errors

$$E(\mathcal{E}) = \sum_{ij} \left( \delta_{ij} - (UU^T)_{ij} \right) E(\varepsilon_i \varepsilon_j) = \sum_i \sigma_i^2 - \sum_i (UU^T)_{ii} \sigma_i^2.$$

In the case of equal variance  $\sigma_i = \sigma$  one obtains

$$E(\mathcal{E}) = (N - (p + 1))\sigma^2. \quad (23)$$

Instead of using  $\mathcal{E}$ , it is possible to rescale it as

$$s^2 = \frac{\mathcal{E}}{N - (p + 1)}. \quad (24)$$

### 2.3 Confidence region of regression coefficients. Normal case

Let us assume that the random errors  $\varepsilon_i$  are independent and normally distributed with zero mean and variance  $\sigma_i^2$ . The assumption of the normal distribution of the random errors in regression model allows us to draw the confidence regions for all parameters in the  $\mathbb{R}^{p+1}$  space and also for each individual parameter.

**THEOREM 1** *Assume that the errors  $\varepsilon_i$  are independent, normal distributed, with zero mean, and square variance  $\sigma_i^2$ . Then:*

(a)  $\{\mathcal{B}_a\}_{a=0,p}$  has a multivariate normal distribution with probability density function given by

$$f_{\mathcal{B}}(\boldsymbol{\eta}) = \sqrt{\frac{\det \mathbf{M}}{(2\pi)^{p+1}}} \exp \left( -\frac{1}{2} \sum_{a,b} M^{ab} (\eta_a - \beta_a)(\eta_b - \beta_b) \right), \boldsymbol{\eta} \in \mathbb{R}^{p+1} \quad (25)$$

where the matrix  $\mathbf{M}$  is the inverse of the covariant matrix of the parameters. If  $\sigma_i = \sigma$  for all  $i$  then  $\mathbf{M}$  is given by

$$M^{ab} = \sum_c \frac{w_{cc}^2}{\sigma^2} v^{ac} v^{bc}$$

and

$$\det \mathbf{M} = \frac{\prod w_{aa}^2}{\sigma^{2(p+1)}}.$$

If in addition  $\sigma_i = \sigma$  for all  $i$  then (b) The marginal probability density function of  $\mathcal{B}_a$  is given by

$$f_a(t) = \frac{1}{\sigma \sqrt{2\pi} \sqrt{\sum_c (w_{cc})^{-2} v^{ac} v^{ac}}} \exp \left( -\frac{(t - \beta_a)^2}{2\sigma^2 \sum_c (w_{cc})^{-2} v^{ac} v^{ac}} \right); \quad (26)$$

- (c)  $\left( \mathbf{WV}^T \frac{\mathbf{B} - \boldsymbol{\beta}}{\sigma}, \mathbf{WV}^T \frac{\mathbf{B} - \boldsymbol{\beta}}{\sigma} \right)$  is chi-squared distributed with  $p + 1$  degrees of freedom;  
 (d)  $\mathcal{E}/\sigma^2$  is chi-squared distributed with  $N - p - 1$  degrees of freedom;  
 (e)  $\frac{1}{p+1} \frac{(\mathbf{WV}^T (\mathbf{B} - \boldsymbol{\beta}), \mathbf{WV}^T (\mathbf{B} - \boldsymbol{\beta}))}{s^2}$  is  $F_{p+1, N-p-1}$  distributed.

*Proof* (a) To prove the theorem we will use the characteristic function of a random variable, see Appendix.

Let  $\mathbf{t} \in \mathbb{R}^{p+1}$  and let

$$\alpha = i \sum_a t_a \mathcal{B}$$

that can be rewritten as

$$\alpha = i(\mathbf{t}, \boldsymbol{\beta}) + i\mathbf{t}^T \mathbf{L}\boldsymbol{\varepsilon}.$$

One has

$$\varphi_{\mathbf{B}}(\mathbf{t}) = E(e^\alpha) = \frac{\exp i(\mathbf{t}, \boldsymbol{\beta})}{(\sigma\sqrt{2\pi})^N} \int_{\mathbb{R}^N} \exp \left[ i \sum_j m_j y_j - \frac{1}{2} \sum_j \frac{y_j^2}{\sigma^2} \right] d\mathbf{y} =$$

where  $m_j = \sum_a t_a l_{aj}$ . By using the formula ??, see Appendix, we have

$$\varphi_{\mathbf{B}}(\mathbf{t}) = \exp i(\mathbf{t}, \boldsymbol{\beta}) \exp \left( -\frac{1}{2} \sum_j m_j^2 \sigma_j^2 \right).$$

By using

$$\sum_j m_j^2 \sigma_j^2 = \sum_j \sum_{a,b} l_{aj} l_{bj} t_a t_b \sigma_j^2 = \sum_{a,b} \text{Cov}(\mathcal{B}_a, \mathcal{B}_b) t_a t_b,$$

one has

$$\varphi_{\mathbf{B}}(\mathbf{t}) = \exp \left( i(\mathbf{t}, \boldsymbol{\beta}) - \frac{1}{2} \sum_{a,b} \text{Cov}(\mathcal{B}_a, \mathcal{B}_b) t_a t_b \right). \quad (27)$$

By using Theorem ?? the affirmation (a) is proved.

The affirmation (b) results from

$$\varphi_{|\beta_a}(t) = \varphi_{|\boldsymbol{\beta}}(0, \dots, t, \dots, 0) = \exp \left( it\beta_a - \frac{t^2}{2} \text{Cov}(|\beta_a, |\beta_a) \right).$$

(c) First we note that

$$\mathbf{WV}^T (\mathbf{B} - \boldsymbol{\beta}) = \mathbf{U}^T \boldsymbol{\varepsilon}.$$

Then

$$\xi := \left( \mathbf{WV}^T \frac{\mathbf{B} - \boldsymbol{\beta}}{\sigma}, \mathbf{WV}^T \frac{\mathbf{B} - \boldsymbol{\beta}}{\sigma} \right) = (\mathbf{U}^T, \mathbf{U}^T) / \sigma^2.$$

The characteristic function of  $\xi$  is given by

$$\begin{aligned}\varphi_\xi(t) &= \int_{\mathbb{R}^N} \exp \left[ i \sum_{jk} (UU^T)_{jk} y_j y_k - \sum_j \frac{1}{2} j^j \right] dy = \\ &= \int_{\mathbb{R}^N} \exp \left( -\frac{1}{2} \sum_{jk} Q_{jk} y_j y_k \right) dy\end{aligned}$$

where the matrix  $\mathbf{Q}$  is given by

$$Q_{ij} = \delta_{ij} - 2it \sum_a U_i^a U_j^b.$$

Concerning the matrix  $\mathbf{Q}$ , its  $N$  eigenvalues are

$$\lambda_1 = \lambda_2 = \cdots = \lambda_{p+1} = 1 - 2it, \quad \lambda_{p+2} = \cdots = \lambda_N = 1$$

and there exists an orthogonal matrix  $\mathbf{R}$  such that  $\mathbf{R}^T \mathbf{Q} \mathbf{R}$  is a diagonal matrix with diagonal entries  $\lambda_1, \lambda_2, \dots, \lambda_N$ .

The set of the vectors  $\{\mathbf{U}_a\}_{a=0, \overline{p}}$  spans a linear subspace of  $\mathbb{R}^N$ , say  $L$ , of dimension  $p+1$  and it is a orthonormal basis of that subspace. Let  $H$  be the orthogonal complement of  $L$  and let  $\{\mathbf{S}_a\}_{a=\overline{p+1}, \overline{N-1}}$  be an orthonormal basis of  $H$  that is

$$\langle \mathbf{S}_a, \mathbf{U}_b \rangle = 0; \forall a, b, \quad \text{and} \quad \langle \mathbf{S}_a, \mathbf{S}_b \rangle = \delta_{ab}, \quad a, b = \overline{p+1}, \overline{N-1}.$$

The vectors  $\{\mathbf{U}_a\}_{a=0, \overline{p}}$  and  $\{\mathbf{S}_a\}_{a=\overline{p+1}, \overline{N-1}}$  are eigenvectors of  $\mathbf{Q}$ , indeed

$$\sum_j Q_{ij} U_{cj} = \sum_j \left( \delta_{ij} - 2it \sum_a U_{ai} U_{bj} \right) U_{cj} = (1 - 2it) U_{ci},$$

and

$$\sum_j Q_{ij} S_{cj} = \sum_j \left( \delta_{ij} - 2it \sum_a U_{ai} U_{aj} \right) S_{cj} = S_{ci}.$$

Denote by  $\mathbf{R}$  the matrix whose columns are the vectors  $\mathbf{U}_0, \dots, \mathbf{U}_p, \mathbf{S}_{p+1}, \dots, \mathbf{S}_{N-1}$  and by  $\mathbf{\Lambda}$  the diagonal matrix with entries  $\lambda_1, \dots, \lambda_N$ . One has

$$\mathbf{Q} = \mathbf{R} \mathbf{\Lambda} \mathbf{R}^T$$

and by the transformation of variable

$$z^a = \sum_i R_i^a y_i,$$

we obtain

$$\varphi_\xi(t) = \int_{\mathbb{R}^N} \exp \left( -\frac{1}{2} \sum_i \lambda_i (z_i)^2 \right) dz$$

## Exploratory Data Analysis Tools

---

and

$$\varphi_\xi(t) = \left( \prod \lambda_i \right)^{-1/2} = (1 - 2it)^{-\frac{(p+1)}{2}}$$

which is the characteristic function of  $\chi^2$  with  $N - (p + 1)$  degrees of freedom, see Appendix.

(d) We will also calculate the characteristic function of random variable  $|\mathcal{E}/\sigma^2$ . By some algebraic calculation we can rewrite the expression (??) as

$$\mathcal{E}/\sigma^2 = \left\langle \frac{\mathbf{Y} - E(\mathbf{Y})}{\sigma}, \frac{\mathbf{Y} - E(\mathbf{Y})}{\sigma} \right\rangle - \left( \mathbf{U}^T \left( \frac{\mathbf{Y} - E(\mathbf{Y})}{\sigma} \right), \mathbf{U}^T \left( \frac{\mathbf{Y} - E(\mathbf{Y})}{\sigma} \right) \right) \quad (28)$$

The characteristic function of  $|\mathcal{E}/\sigma^2$  is given by

$$\begin{aligned} \varphi_{\mathcal{E}/\sigma^2} &= E \left( e^{i\mathcal{E}/\sigma^2} \right) = \frac{1}{\sqrt{2\pi}^N} \int_{\mathbb{R}^N} \exp \left[ it \sum_{i,j} \delta_{i,j} y_i y_j - it \sum_{i,j} \sum_a U_i^a U_j^a y_i y_j - \frac{1}{2} \sum_{i,j} \delta_{i,j} y_i y_j \right] d\mathbf{y} \\ &= \frac{1}{\sqrt{2\pi}^N} \int_{\mathbb{R}^N} \exp \left( -\frac{1}{2} \sum_{i,j} Q_{i,j} y_i y_j \right) d\mathbf{y}, \end{aligned}$$

where the matrix  $\mathbf{Q}$  is given by

$$Q_{ij} = (1 - 2it)\delta_{ij} + 2it \sum_a U_i^a U_j^b.$$

The  $N$  eigenvalues of the matrix  $\mathbf{Q}$  are

$$\lambda_1 = \lambda_2 = \dots = \lambda_{p+1} = 1, \lambda_{p+2} = \dots = \lambda_N = 1 - 2it$$

and there is an orthogonal matrix  $\mathbf{R}$  such that  $\mathbf{R}^T \mathbf{Q} \mathbf{R}$  is a diagonal matrix with diagonal entries  $\lambda_1, \lambda_2, \dots, \lambda_N$ .

The set of the vectors  $\{\mathbf{U}^a\}_{a=\overline{0,p}}$  spans a linear subspace of  $\mathbb{R}^N$ , say  $L$ , of dimension  $p + 1$  and it is a orthonormal basis of that subspace. Let  $H$  be the orthogonal complement of  $L$  and let  $\{\mathbf{S}^a\}_{a=\overline{p+1,N-1}}$  be an orthonormal basis of  $H$ , that is

$$\langle \mathbf{S}^a, \mathbf{U}^b \rangle = 0; \forall a, b, \text{ and } \langle \mathbf{S}^a, \mathbf{S}^b \rangle = \delta_{ab}, a, b = \overline{p+1, N-1}.$$

The vectors  $\{\mathbf{U}^a\}_{a=\overline{0,p}}$  and  $\{\mathbf{S}^a\}_{a=\overline{p+1,N-1}}$  are eigenvectors of  $\mathbf{Q}$ . Indeed

$$\sum_j Q_{ij} U_j^c = \sum_j \left( (1 - 2it)\delta_{ij} + 2it \sum_a U_i^a U_j^b \right) U_j^c = U_i^c,$$

and

$$\sum_j Q_{ij} S_j^c = \sum_j \left( (1 - 2it)\delta_{ij} + 2it \sum_a U_i^a U_j^b \right) S_j^c = (1 - 2it)S_i^c.$$

Denote by  $\mathbf{R}$  the matrix whose columns are the vectors  $\mathbf{U}^0, \dots, \mathbf{U}^p, \mathbf{S}^{p+1}, \dots, \mathbf{S}^{N-1}$  and by  $\mathbf{\Lambda}$  the diagonal matrix with entries  $\lambda_1, \dots, \lambda_N$  one has

$$\mathbf{Q} = \mathbf{R}\mathbf{\Lambda}\mathbf{R}^T$$

and by the transformation of the variable

$$z^a = \sum_i R_i^a y_i$$

we obtain

$$\varphi_{\mathcal{E}/\sigma^2}(t) = \frac{1}{\sqrt{2\pi}^N} \int_{\mathbb{R}^N} \exp\left(-\frac{1}{2} \sum_i \lambda_i (z_i)^2\right) dz$$

and

$$\varphi_{\mathcal{E}/\sigma^2}(t) = \left(\prod \lambda_i\right)^{-1/2} = (1 - 2it)^{-\frac{N-(p+1)}{2}}$$

which is the characteristic function of  $\chi^2$  with  $N - (p + 1)$  degrees of freedom, see Appendix.

(e) Is a consequence of (c), (d) and proposition 1.

To draw the confidence regions for the true parameter  $\beta$  one uses the estimated parameters  $\hat{\beta}$  and one of the distributions that appears in Theorem 1.

Consider an  $p$ -ellipsoid in  $\mathbb{R}^{p+1}$  with principal axes given by the unitary vectors  $v^a \in \mathbb{R}^{p+1}$ ,  $a = \overline{0, p}$  and the length  $\rho/w_{aa}$ , respectively,

$$\Delta(\rho) = \left\{ \eta \in \mathbb{R}^{p+1} \mid \sum_c (w_{cc})^2 \sum_{a,b} v^{ac} v^{bc} \eta_a \eta_b \leq \rho^2 \right\} \quad (29)$$

Let  $\hat{\beta}$  be a set of estimated parameters and suppose that  $\sigma$  is a known quantity. The point  $\frac{\hat{\beta} - \beta}{\sigma}$  belongs to  $\Delta(\rho)$  if

$$\sum_c (w_{cc})^2 \sum_{a,b} v^{ac} v^{bc} \frac{\hat{\beta}_a - \beta_a}{\sigma} \frac{\hat{\beta}_b - \beta_b}{\sigma} \leq \rho^2$$

but the left hand term is chi-square distributed so that we obtain

$$P\left(\frac{\hat{\beta} - \beta}{\sigma} \in \Delta(\rho)\right) = \chi_{p+1}(\rho^2). \quad (30)$$

If  $\sigma$  is unknown, as usual, we use the estimation  $s$ , formula (24), of  $\sigma$  and we have

$$P\left(\frac{\hat{\beta} - \beta}{s\sqrt{p+1}} \in \Delta(\rho)\right) = F_{p+1; N-p-1}(\rho^2). \quad (31)$$

## Exploratory Data Analysis Tools

---

For an individual parameter  $\beta_a$  one can use the random variable

$$\mathbf{T} = \frac{|\beta_a - \hat{\beta}_a|}{\sigma \sqrt{(\mathbf{X}^T \mathbf{X})_{aa}^{-1}}} \frac{1}{\sqrt{\frac{|\mathcal{E}|}{\sigma^2(N-p-1)}}} = \frac{|\beta_a - \hat{\beta}_a|}{\sqrt{(\mathbf{X}^T \mathbf{X})_{aa}^{-1}}} \frac{1}{\sqrt{\frac{|\mathcal{E}|}{N-p-1}}} = \frac{|\beta_a - \hat{\beta}_a|}{s \sqrt{\sum_c (w_{cc})^{-2} v^{ac} v^{bc}}} \quad (32)$$

which has  $t$ -distribution with  $N - p - 1$  degrees of freedom.

The table 2.3 summarizes the essential facts concerning the OLS estimation of the parameters.

Table 1: OLS estimation of the parameters in linear regression. Random errors are independent and identical normal distributed with zero mean and  $\sigma$  variance.

Predicted value	$\hat{y}(\mathbf{x}_i) = \hat{\beta}_0 + \sum_{a=1}^p x_i^a \hat{\beta}_a$	
Estimation of the variance	$s^2 = \frac{\sum_i (\hat{y}_i - \hat{y}(\mathbf{x}_i))^2}{N-(p+1)}$	
Covariance matrix	$\text{Cov}( \beta_a,  \beta_b) = \sigma^2 \sum_c (w_{cc})^{-2} v^{ac} v^{bc}$ $V( \beta_a) = \text{Cov}( \beta_a,  \beta_b) / \sigma^2$	
Confidence region	$\Delta(\rho) = \left\{ \eta \in \mathbb{R}^{p+1} \mid \sum_c (w_{cc})^2 \sum_{a,b} v^{ac} v^{bc} \eta_a \eta_b \leq \rho^2 \right\}$	
	$ \beta_a$ marginal distribution	joint distribution
$\sigma$ known	$P \left( \frac{ \hat{\beta}_a - \beta_a }{\sigma \sqrt{V( \beta_a)}} \leq \sqrt{2}\xi \right) = \text{erf}(\xi)$	$P \left( \frac{\hat{\beta} - \beta}{\sigma} \in \Delta(\rho) \right) = \chi_{p+1}(\rho^2)$
$\sigma$ unknown	$P \left( \frac{ \hat{\beta}_a - \beta_a }{s \sqrt{V( \beta_a)}} \leq \xi \right) = t(\xi)$	$P \left( \frac{\hat{\beta} - \beta}{s \sqrt{p+1}} \in \Delta(\rho) \right) = F_{p+1; N-p-1}(\rho^2)$

**Goodness-of-fit**

A lot of criteria to evaluate the quality of the estimate linear model ?? exist. The most used criterion is  $R^2$  which is defined as:

$$R^2 = 1 - \frac{\sum_i e_i^2}{\sum_i (\tilde{y}_i - \bar{\tilde{y}})^2}. \quad (33)$$

**Hypothesis Testing** Let us consider the following type of statistical test: *null hypothesis*  $H_0 : \beta \in S_m$  and the *alternative*  $H_1 : \beta \in \mathcal{CS}_m$ . The set  $S_m$  is a given m-plane

$$\eta_a = \tilde{\beta}_a + \sum_{i=1}^m \tau_a^i z_i \quad z_i \in \mathbb{R}. \quad (34)$$

For a given  $\alpha$  let  $\rho_\alpha$  be chosen such that  $F_{p+1;N-p-1}(\rho^2) = \alpha$ . It follows that

$$P\left(\frac{\hat{\beta} - \beta}{s\sqrt{p+1}} \in \Delta(\rho)\right) = \alpha$$

In the space of the parameters consider the elliptical domain

$$\Delta_\alpha = \left\{ \eta \in \mathbb{R}^{p+1} \mid \sum_{a,b} C^{ab} (\hat{\beta}_a - \eta_a) (\hat{\beta}_b - \eta_b) \leq \rho_{s\alpha}^2 \right\}$$

where  $C^{ab} = \sum_c (w_{cc})^2 v^{ac} v^{bc}$  and  $\rho_{s\alpha}^2 = \rho_\alpha^2 s^2 (p+1)$ .

We reject  $H_0$  if m-plane  $S_m$  do not intersect  $\Delta_\alpha$  and accept it if contrary.

We introduce the function

$$f(z_1, z_2, \dots, z_m) = \sum_{a,b} C^{ab} \left( \hat{\beta}_a - \tilde{\beta}_a - \sum_i \tau_a^i z_i \right) \left( \hat{\beta}_b - \tilde{\beta}_b - \sum_i \tau_b^i z_i \right) - \rho_{s\alpha}^2$$

which measures the "distance" from a point  $z$  of  $S_m$  to the center  $\hat{\beta}$  of the elliptical domain  $\Delta_\alpha$ . Firstly we determine the point of the m-plane  $S_m$  which minimize the "distance"  $f$ .

We have

$$\frac{\partial f}{\partial z_i} = 2 \left( \sum_{a,b} C^{ab} \tau_a^i \sum_j \tau_b^j z_j - \sum_{a,b} C^{ab} \tau_a^i (\hat{\beta}_b - \tilde{\beta}_b) \right)$$

and

$$\frac{\partial^2 f}{\partial z_i \partial z_j} = \sum_{a,b} C^{ab} \tau_a^i \tau_b^j.$$

It follows that the solution of the equations

$$\frac{\partial f}{\partial z_i} = 0$$

minimizes  $f$ .

Denote

$$B^{ij} = \sum_{a,b} C^{ab} \tau_a^i \tau_b^j, \quad \gamma^i = \sum_{a,b} C^{ab} \tau_a^i (\hat{\beta}_b - \tilde{\beta}_b).$$

The solution is given by

$$\hat{\mathbf{z}} = \mathbf{B}^{-1} \boldsymbol{\gamma} \tag{35}$$

and

$$f(\hat{\mathbf{z}}) = \sum_{a,b} C^{ab} (\hat{\beta}_a - \tilde{\beta}_a)(\hat{\beta}_b - \tilde{\beta}_b) - \sum_{i,j} B_{ij}^{-1} \gamma^i \gamma^j - \rho_{s\alpha}^2. \tag{36}$$

To test the null hypothesis one has to compare the  $f(\hat{\mathbf{z}})$  with 0.

If  $f(\hat{\mathbf{z}}) > 0$ , since  $f(\mathbf{z}) \geq f(\hat{\mathbf{z}})$  for any  $\mathbf{z}$ , the  $m$ -plane  $S_m$  does not intersect the  $\Delta_\alpha$  which implies that the null hypothesis is false and we reject it at significance level  $1 - \alpha$ . If  $f(\hat{\mathbf{z}}) \leq 0$  the  $m$ -plane  $S_m$  intersects the  $\Delta_\alpha$  which implies that the null hypothesis is true.

*Remark.* The "distance"  $f(\mathbf{z})$  does not depend of the point  $\tilde{\boldsymbol{\beta}}$  nor on  $S_m$  and neither of the choice of the basis vectors  $\boldsymbol{\tau}^i$  of  $S_m$ !

## 2.4 Bootstrap Confidence Intervals

In many practical situations the distribution of the errors violates the normality conditions.

In such situations we can yet use the OLS estimator to estimate the parameters in the regression model but we can not use the formulae (30), (31) and (32) to estimate the confidence region or the formula (35) to test hypothesis. The bootstrap and jackknife are resampling techniques that are meant to supply informations about the distribution of the errors in the model by using only the measured data. The general framework is that of [20].

Let us denote by  $\mathcal{D}_0$  the set of measured values  $\{y_i, \mathbf{x}_i\}_{i=1, \dots, N}$  and let  $\hat{\boldsymbol{\beta}}$  be the OLS estimate corresponding to it.

A resampling procedure produces a synthetic data set  $\mathcal{D}_*$  from  $\mathcal{D}_0$  and by using new data  $\mathcal{D}_*$  one can define the OLS estimate  $\hat{\boldsymbol{\beta}}^*$ . The distribution of  $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}$  is approximated by the distribution on  $\hat{\boldsymbol{\beta}}^* - \hat{\boldsymbol{\beta}}$ ! Starting with the pioneering paper [?] several techniques were proposed to be used in regression analysis. Here we attempt to review some of the most relevant techniques in the field of numerical ecology and to explain how these are working when applied to estimate the parameters, covariance matrix and confidence interval.

### Resampling procedures

*The Generalized Residual Bootstrap* (GBS), [1]. Basically, the techniques classified as (GBS) build up new dependent data by resampling the fitted errors. Denote

$$e_i = \tilde{y}_i - \hat{y}_i \text{ and } e_i^s = e_i - \bar{e}.$$

Let  $\mathbf{W}$  be a random matrix and define

$$\mathbf{y}^* = \mathbf{X} \hat{\boldsymbol{\beta}} + \mathbf{W} \mathbf{r}.$$

where  $\mathbf{r}$  is either  $\mathbf{e}$  or  $\mathbf{e}^s$ . The resampled data set  $\mathcal{D}_W$  consists in  $\{y_i^*, \mathbf{x}_i\}_{i=1, \dots, N}$ . The class of GBS includes:



(a) the *classical bootstrap* (GBSa), [6], the rows  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N$  of random matrix  $\mathbf{W}$  are i.i.d samples from *Multinomial*  $(1; 1/N, 1/N, \dots, 1/N)$  and  $\mathbf{r} = \mathbf{e}^s$ ;

(b) the *weighted bootstrap* (GBSb), [13],  $\mathbf{r} = \mathbf{e}$  and  $\mathbf{W} = (\mathbf{I} - \mathbf{J})\mathbf{W}_1$ , where  $\mathbf{W}_1$  is diagonal matrix with entries that are i.i.d observations from a random variable with zero mean and unit variance and  $\mathbf{J}_{ij} = 1/N, i, j = \overline{1, N}$ ;

(c) the *external or wild bootstrap* (GBSc), [22] and [14],  $\mathbf{r} = \mathbf{e}$  and  $\mathbf{W}$  is diagonal matrix with entries that are i.i.d observations from a random variable with zero mean and unit variance.

*Block Bootstrap Methods*, [9] Let  $l$  be an integer  $1 \leq l < N$ ,  $l$  denotes the block length in the method. One forms the series  $\mathcal{S} = \{\tilde{y}_{0;i}; \mathbf{x}_{0;i}\}_{i>1}$  by periodic extension,  $(\tilde{y}_{0;i}; \mathbf{x}_{0;i}) = (\tilde{y}_j; \mathbf{x}_j)$  if  $i = m * N + j$  for some integer  $m$  and  $1 \leq j \leq N$ . Define the blocks of length  $k \geq 1$  based on the series  $\mathcal{S}$  by  $\mathcal{B}(i, k) = ((y_{0;i}; \mathbf{x}_{0;i}), \dots, (y_{0;i+k-1}; \mathbf{x}_{0;i+k-1}))$ . Different methods in this class are defined by resampling from a specific subcollection of the blocks  $\{\mathcal{B}(i, k) : i, k \geq 1\}$ .

(MBB). This method resamples blocks randomly, with replacement from  $\{\mathcal{B}(i, l) : i = 1, \dots, N-l+1\}$ . The synthetic data  $\mathcal{D}_{MBB} = \{\mathcal{B}(I_{11}, l), \mathcal{B}(I_{12}, l), \dots, \mathcal{B}(I_{1b}, l)\}$ , where  $b = \lfloor n/l \rfloor$  (the integer part of  $n/l$ ) and  $I_{1i}$  are i.i.d samples from *Multinomial*  $(n-l+1; 1/(n-l+1), \dots, 1/(n-l+1))$ .

(NBB). The same as MBB but this method uses a subcollection of disjoint blocks  $\{\mathcal{B}((i-1)l+1, l) : 1 \leq i \leq b\}$ .

*The Jackknife resampling procedures*. [6] A synthetic data set in delete-1 Jackknife method consists in all data  $\mathcal{D}_0$  except one,  $\mathcal{D}_{(-i)} = \{(\tilde{y}_j; \mathbf{x}_j) : j = \overline{1, N}, j \neq i\}$ .

As we remarked before, the synthetic data sets are used to approximate the distribution of the parameters. By using the sample distribution one can estimate the confidence interval and test the hypothesis.

## 2.5 BC<sub>a</sub> method

The BC<sub>a</sub> ("bias-corrected and accelerated") procedure is a method of setting approximate confidence intervals for parameters from percentiles of bootstrap histogram [5]. It involves the cumulative distribution function of the resampling replicants and other two parameters: the bias correction  $z_0$  and the acceleration  $a$ . Let  $\hat{\beta}(l)$  be the OLS estimate from a synthetic data set  $\mathcal{D}_l$  and let  $B$  be the number of synthetic data sets used in the resampling procedure. Let  $\hat{F}_b(z)$  be the cumulative distribution function of  $B$  resampling replications  $\hat{\beta}_b(l)$

$$\hat{F}_b(z) = \#\{\hat{\beta}_b(l) < z\}/B. \quad (37)$$

A confidence interval is set by BC<sub>a</sub> end points. By definition BC<sub>a</sub> is

$$\hat{\beta}_{b,BC_a}(\alpha) = \hat{F}_b^{-1} \left( \Phi \left( z_0 + \frac{z_0 + z^{(\alpha)}}{1 - a(z_0 + z^{(\alpha)})} \right) \right), \quad (38)$$

where  $\Phi$  is the standard distribution function with  $z^{(\alpha)} = \Phi(\alpha)$ . The central 0.90 BC<sub>a</sub> interval is given by

$$\left( \hat{\beta}_{b,BC_a}(0.05), \hat{\beta}_{b,BC_a}(-.95) \right). \quad (39)$$

## Exploratory Data Analysis Tools

---

The several variants of  $BC_a$  differ by the method used to estimate  $z_0$  and  $a$ . Here we consider, as in [5],

$$\hat{z}_0 = \Phi^{-1} \left( \hat{F}_b(\hat{\beta}_b) \right), \quad (40)$$

$$\hat{a} = \frac{1}{6} \frac{\sum_{i=1}^N A_i^3}{\left( \sum_{i=1}^N A_i^2 \right)^{3/2}} \quad (41)$$

and

$$A_i = (N - 1)(\hat{\beta}_b - \hat{\beta}_b^{(-i)}) \quad (42)$$

where  $\hat{\beta}_b^{(-i)}$  is replicant of the jackknife synthetic data  $\mathcal{D}_i$ .

### 3 Appendix

Useful formula. Let  $\xi$  be a real random variable. We denote by  $P_\xi$  the *probability distribution* of it, by  $F_\xi$  its *distribution function* i.e.

$$P_\xi(\xi \leq x) = F_\xi(x), \tag{43}$$

and by  $f_\xi$  the density of distribution function, if it exists., i.e

$$F_\xi(x) = \int_{-\infty}^x f_\xi(t) dt. \tag{44}$$

For any real numbers  $a \leq b$  one has

$$P_\xi((a, b]) = F_\xi(b) - F_\xi(a). \tag{45}$$

Table 2: Important density distributions and their parameters.

Distribution	Density	Parameters
Normal, $N(\mu, \sigma)$	$\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), x \in \mathbb{R},$	$\mu \in \mathbb{R}, \sigma > 0$
Gamma	$\frac{x^{\alpha-1} \exp(-x/\beta)}{\Gamma(\alpha)\beta^\alpha}, x \geq 0,$	$\alpha > 0, \beta > 0$
Beta	$\frac{x^{r-1} (1-x)^{s-1}}{B(r, s)}, 0 \leq x \leq 1,$	$r > 0, s > 0$
Chi-squared, $\chi^2$	$\frac{x^{\frac{n}{2}-1} \exp(-x/2)}{2^{n/2}\Gamma(\frac{n}{2})}, x \geq 0,$	$n = 1, 2 \dots$
$F_{m;n}$	$\frac{(m/n)^{m/2}}{B(m/2, n/2)} \frac{x^{m/2-1}}{(1 + \frac{mx}{nx})^{(m+n)/2}}, x > 0,$	$m = 1, 2 \dots, n = 1, 2, \dots$

If  $\xi$  and  $\eta$  are two random variable with joint distribution  $F_{\xi\eta}(x, y)$  and  $\phi(x, y)$  is a Borel function, if we put  $\psi = \phi(\xi, \eta)$  we have

$$F_\psi(z) = \int \int_{\{(x,y)|\phi(x,y)\leq z\}} dF_{\xi\eta}. \tag{46}$$

If  $\xi$  and  $\eta$  are independent then  $F_{\xi\eta}(x, y) = F_\xi(x)F_\eta(y)$ .

**PROPOSITION 1** *If  $\xi$  and  $\eta$  are independent and chi-squared distributed with  $m$  and  $n$  degrees of freedom respectively, then  $(\xi/m)/(\eta/n)$  is  $F(m;n)$  distributed.*

$$f_{(\xi/m)/(\eta/n)} = F_{m;n} \tag{47}$$

## Exploratory Data Analysis Tools

---

Table 3: Distribution functions and probability density functions of the arithmetic operations of two independent random variable.

$\phi(x, y)$	$F_\phi(z) =$	$f_\phi(z) =$
$\phi(x, y) = x + y$	$\int_{-\infty}^{\infty} F_\xi(z - y) dF_\eta(y)$	$\int_{-\infty}^{\infty} f_\xi(z - y) f_\eta(y) dy$
$\phi(x, y) = xy$	$\int_{-\infty}^0 (1 - F_\eta\left(\frac{z}{x}\right)) dF_\xi(x) + \int_{-\infty}^0 (1 - F_\eta\left(\frac{z}{x}\right)) dF_\xi(x)$	$\int_{-\infty}^{\infty} f_\eta\left(\frac{z}{x}\right) f_\xi(x) \frac{dx}{ x }$
$\phi(x, y) = \frac{x}{y}$	$\int_{-\infty}^0 (1 - F_\xi(zy)) dF_\eta(y) + \int_0^{\infty} F_\xi(zy) dF_\eta(y)$	$\int_{-\infty}^{\infty} f_\xi(zy) f_\eta(y) \frac{dy}{ y }$

*Proof.* By applying the ratio formula from table 3 and by taking into account that  $f_{\xi/m}(x) = m f_\xi(mx)$ , (47) is obtained.  $\square$

Let  $F = (F_1, \dots, F_n)$  be a  $n$ -dimensional distribution function in  $\mathbb{R}^n$ . Its *characteristic function* is

$$\varphi(t) = \int_{\mathbb{R}^n} e^{i\langle t, x \rangle} dF(x), t \in \mathbb{R}^n. \quad (48)$$

If  $\xi = (\xi_1, \dots, \xi_n)$  is a random vector with values in  $\mathbb{R}^n$ , its *characteristic function* is

$$\varphi_\xi(t) = \int_{\mathbb{R}^n} e^{i\langle t, x \rangle} dF_\xi(x), t \in \mathbb{R}^n; \quad (49)$$

if  $\xi$  has probability density function then

$$\varphi_\xi(t) = \int_{\mathbb{R}^n} e^{i\langle t, x \rangle} f_\xi(x) dx, \quad (50)$$

which is the Fourier transform of  $f_\xi$ . The characteristic function of a random vector can be also defined by

$$\varphi_\xi(t) = E\left(e^{i\langle t, \xi \rangle}\right). \quad (51)$$

**PROPOSITION 2** (a) *The characteristic function of chi-squared distribution with  $n$  degree of freedom is given by*

$$\varphi_{\chi_n^2}(t) = (1 - 2it)^{-\frac{n}{2}} \quad (52)$$

(b) *The characteristic function of gaussian distribution  $N(\mu, \sigma)$  is given by*

$$\varphi_{N(\mu, \sigma)}(t) = \exp\left(it\mu - \frac{t^2\sigma^2}{2}\right) \quad (53)$$

*Proof.* The analytic function

$$f(z) = z^{\frac{n}{2}-1} \exp\left(-\frac{z}{2}\right)$$

is considered.

In the complex plane we take the point  $a + ib$  with  $a > 0$  and the path

$$C_\varepsilon = \{(x, 0) | 0 \leq x < \varepsilon\} \cup \{(\varepsilon, x) | 0 \leq x \leq \frac{b\varepsilon}{a}\} \cup \{(ax, bx) | 0 \leq x < \frac{\varepsilon}{a}\}.$$

Since

$$\int_{C_\varepsilon} f(z) dz = 0,$$

we have

$$\int_0^\varepsilon x^{\frac{n}{2}-1} \exp\left(-\frac{x}{2}\right) dx + i \int_0^\varepsilon (\varepsilon + ix)^{\frac{n}{2}-1} \exp\left(-\frac{\varepsilon + ibx}{2}\right) dx = (a + ib)^{\frac{n}{2}} \int_0^\varepsilon x^{\frac{n}{2}-1} \exp\left(-\frac{a + ib}{2}x\right) dx.$$

By using the fact that the second integral on the left hand side becomes zero as  $\varepsilon \rightarrow \infty$ , one obtains

$$\int_0^\infty x^{\frac{n}{2}-1} \exp\left(-\frac{x}{2}\right) dx = (a + ib)^{\frac{n}{2}} \int_0^\infty x^{\frac{n}{2}-1} \exp\left(-\frac{a + ib}{2}x\right) dx. \square$$

**THEOREM 2** Let  $\mathbf{Q}$  be a  $n \times n$  real symmetric and positive definite matrix. The nonnegative function

$$f(\mathbf{x}) = \sqrt{\frac{\det(\mathbf{Q})}{(2\pi)^n}} \exp\left(-\frac{(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q} (\mathbf{x} - \boldsymbol{\mu})}{2}\right) \quad (54)$$

defines the probability density function of some random vector, normal distributed with vector mean  $\boldsymbol{\mu}$  and  $\mathbf{Q}^{-1}$  covariant matrix,  $\boldsymbol{\xi} = (\xi_1, \xi_2, \dots, \xi_n)$ .

The characteristic function of  $\boldsymbol{\xi}$  is given by

$$\varphi_{\boldsymbol{\xi}}(\mathbf{t}) = \exp\left(i \langle \mathbf{t}, \boldsymbol{\mu} \rangle - \frac{\mathbf{t}^T \mathbf{Q}^{-1} \mathbf{t}}{2}\right) \quad (55)$$

**COROLLARY 1** If the random vector  $(\xi_1, \xi_2, \dots, \xi_n)$  is normal distributed with vector mean  $\boldsymbol{\mu}$  and covariance matrix  $\mathbf{Q}^{-1}$ , then the  $(n-1)$ -dimensional vector  $(\xi_1, \xi_2, \dots, \xi_{n-1})$  is also normal distributed with vector mean  $(\mu_1, \mu_2, \dots, \mu_{n-1})$  and covariance matrix  $\mathbf{Q}^{-1}$ . The matrix  $\mathbf{Q}^{(n-1)}$  is given by

$$Q_{ij}^{(n-1)} = Q_{ij} - \frac{Q_{ni} Q_{nj}}{Q_{nn}}, \quad i, j = \overline{1, n-1}$$

### References

- [1] Arup Bose and Snigdhanu Chatterjee, Comparison of Bootstrap and Jackknife Variance Estimators in Linear Regression: Second Order Results, *Statistica Sinica*, **12**(2002), pp. 575–598.
- [2] Raffaele Giancarlo, David Scaturo, and Filippo Utra, A Tutorial on Computational Cluster Analysis with Application to Pattern Discovery in Microarray Data, *Math.comput.sci.*, **1** (2008), pp. 655–672.
- [3] Frederic Cao, Julie Delon, Agnes Desolneux, Pablo Muse, Frederic Sur, An *a contrario* approach to hierarchical clustering validity assessment, *CMLA report*, 2004-16, November 2004.
- [4] James Carpenter and John Bithell, Bootstrap confidence intervals: when, which, what? A practical guide for medical statisticians, *Statist. Med.*, **19**(2000), pp. 1141–1164.
- [5] Thomas J. DiCiccio and Bradley Efron, Bootstrap Confidence Intervals, *Statistical Science*, **11**(1996), pp.189–228.
- [6] Efron, B., Bootstrap methods: Another look at the jackknife, *Ann. Statist.* , **7**(1979), pp. 1–26.
- [7] King, B., Step-wise clustering procedures, *J. Am. Stat. Assoc.*, **69**(1967), pp. 86–101.
- [8] Maria Halkidi, Yannis Batistakis and Michalis Vazirgiannis, On Clustering Validation Techniques, *Journal of Intelligent Information System*, **17**:2/3 (2001), pp. 107–145.
- [9] S.N. Lahiri, Theoretical Comparison of Block Bootstrap Methods, *Ann. Statist.*, **27**(1999), pp. 386–404.
- [10] S.N. Lahiri and Juan Zhu, Resampling Methods for Spatial Regression Models under a Class of Stochastic Design, *Ann. Statist.*, **34**(2006), pp. 1774–1813.
- [11] Francois-Joseph Lapoint and Pierr Legendre, A Statistical Framework to Test the Consensus of Two Nested Clasifications, *Syst. Zool.*, **39**(1)(1990), pp. 1–13.
- [12] Regina Y. Liu and Kesar Singh, Efficiency and Robustness in Resampling, *Ann. Statist.* , **20**(1992), pp. 370–384.
- [13] Regina Y. Liu, Bootstrap procedure under some non i.i.d models, *Ann. Statist.* , **16**(1988), pp. 1697–1708.
- [14] Mammen, Bootstrap and wild bootstrap for high dimensional linear models, *Ann. Statist.* , **21**(1993), pp. 255–285.
- [15] Murtagh, F., A survey of recent advances in hierarchical clustering algorithms which use cluster centers, *Comput. J.*, **26**(1984), pp. 354–359.
- [16] Real, R. and Vargas, J. M., The probabilistic basis of Jaccard’s index of similarity. *Syst. Biol.*, **45**(1996),pp. 380–385.

- [17] Real, R., Tables of significant values of Jaccard's index of similarity. *Misc. Zool.*, **22.1**(1996),pp. 29–40.
- [18] V. K. Rohatgi, *An Introduction to probability Theory and Mathematical Statistics*, John Wiley and Sons, New York, 1976.
- [19] Sneath, P.H.A. and Sokal, R.R., *Numerical Taxonomy*, Freeman, London, UK, 1973.
- [20] William H. Press, Saul A. Teukolsky, William T. Vetterling, Brian P. Flannery, *Numerical Recipes in C*, Cambridge University Press, 1992.
- [21] Ward, J. H. Jr., Hierarchical grouping to optimize an objective function, *J. Am. Stat. assoc.*, **58**(1963), pp. 236–244.
- [22] Wu, C.F.J, Jackknife, bootstrap and other resampling methods in regression analysis, *Ann. Statist.* , **14**(1986), pp. 1261–1350.