

STATISTICAL ASPECTS OF GOWER'S  
INTERPOLATION: ERROR TERM AND ITS  
INFLUENCE ON PREDICTION

*Eva Boj, Maria Mercè Claramunt,  
Aurea Grané and Josep Fortiana*

No. 378  
October 2005

**Mathematics Preprint Series**



**U**  
UNIVERSITAT DE BARCELONA



# Statistical Aspects of Gower's Interpolation: Error Term and its Influence on Prediction

Eva Boj<sup>1</sup>      Maria Mercè Claramunt<sup>1</sup>  
Aurea Grané<sup>2</sup>    Josep Fortiana<sup>3</sup>

July 22, 2005

New points can be superimposed on a Euclidean configuration obtained as a result of a metric Multidimensional Scaling at coordinates given by Gower's interpolation formula. The procedure amounts to discarding a, possibly non-null, coordinate along an additional dimension. We compute this error term, assessing its influence on distance-based predictions.

**Keywords:** Distance-based prediction, Metric multidimensional scaling, Gower interpolation, Typicality in discrimination.

**AMS Subject Classification:** 62H25, 62H30, 62G99, 62H99.

## 1. Introduction

Gower's *interpolation formula* (Gower 1968, see also Gower and Hand 1996, pp. 246ff) is a well-known device for inserting new points in a given Euclidean map obtained from a set of interdistances by means of metric Multidimensional Scaling (MMDS). General context references are Borg and Groenen (1997), Cox and Cox (1994), and Krzanowski and Marriott (1994), as well as Gower and Hand (1996). In general, a faithful representation of the augmented set should require a new dimension for each new point. For a quick illustration,

---

<sup>1</sup>Departament de Matemàtica Econòmica, Financera i Actuarial, Facultat de Ciències Econòmiques i Empresarials. Universitat de Barcelona, Avinguda Diagonal, 690, 08034 Barcelona, Spain.

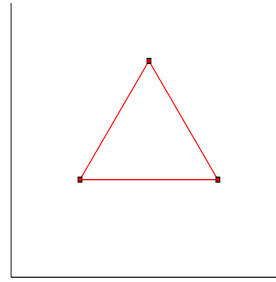
Tel: +34 934 035 744. Fax: +34 934 034 892. Email: evaboj@ub.edu

<sup>2</sup>Departamento de Estadística, Universidad Carlos III, Calle Madrid, 126, 28903 Getafe (Madrid), Spain.

<sup>3</sup>Departament de Probabilitat, Lògica i Estadística, Facultat de Matemàtiques, Universitat de Barcelona, Gran Via de les Corts Catalanes, 585, 08007 Barcelona, Spain.

Figure 1: Interdistances matrix for a set of three objects admitting a two-dimensional Euclidean representation.

$$D_3 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}$$



let us assume a set of three objects, with interdistances shown in Figure 1, to which a fourth object is added, equidistant at distance  $\delta$  from the former three, resulting:

$$D_4 = \begin{pmatrix} 0 & 1 & 1 & \delta \\ 1 & 0 & 1 & \delta \\ 1 & 1 & 0 & \delta \\ \delta & \delta & \delta & 0 \end{pmatrix}.$$

Depending on  $\delta$ , there are four possible situations.

1. If  $\delta < \frac{1}{2}$ , the dissimilarity defined by  $D_4$  is nonmetric, as it fails to satisfy the triangular inequality.
2. If  $\frac{1}{2} \leq \delta < \frac{\sqrt{3}}{3}$ ,  $D_4$  defines a non Euclidean metric (Figure 2, left hand side).
3. If  $\delta = \frac{\sqrt{3}}{3}$ , the dissimilarity defined by  $D_4$  is a Euclidean metric admitting a two-dimensional representation (Figure 2, right hand side).
4. If  $\delta > \frac{\sqrt{3}}{3}$ , the dissimilarity defined by  $D_4$  is a Euclidean metric admitting a three-dimensional representation (Figure 3).

Figure 2: A non Euclidean metric configuration (left hand side) and a two-dimensional Euclidean metric configuration (right hand side).

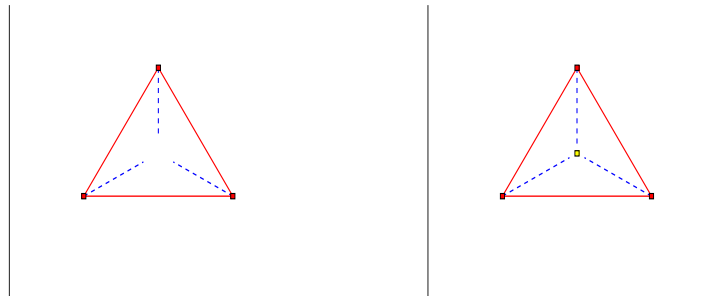


Figure 3: A three-dimensional Euclidean metric configuration.

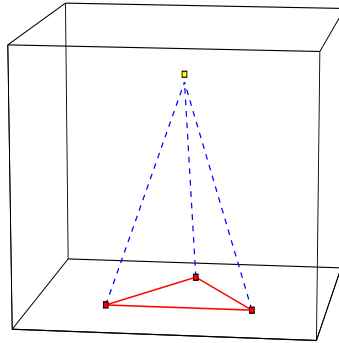


Figure 4: A new observation can be assigned to one of the two known groups (left hand side) or it is far from both (right hand side).

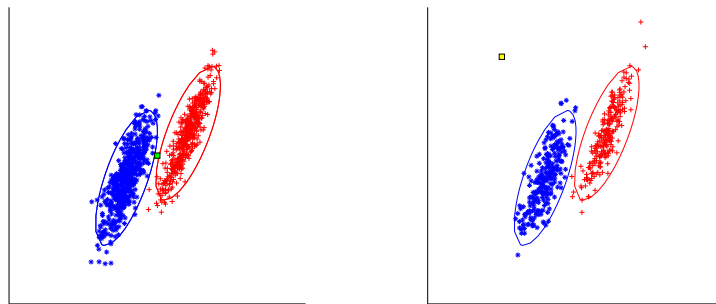


Figure 5: Existence of a new group may be inferred.

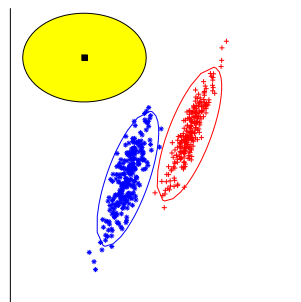
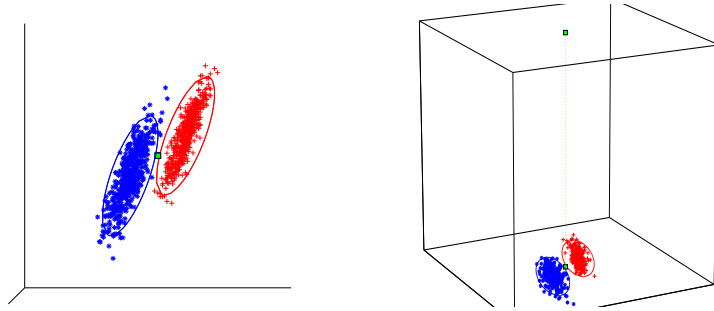


Figure 6: A new observation appears close to the known groups due to the projection effect.



If the exact representation is projected onto the preexisting lower dimensional subspace, which is what Gower's interpolation does, this additional coordinate is lost. The *squared error term*, the square of the additional coordinate, lost in projection, is a measure of interpolation quality. Its significance may lead to incorrect results in distance-based statistical prediction methods. For instance, the *typicality problem* in Discriminant Analysis (Rao 1973) amounts to deciding whether a new observation to be classified belongs to one of the already known groups or, instead, it falls far apart from all of them (Figure 4), so much so as to lead us to postulate the existence of a third, so far unobserved, group (Figure 5). In (Cuadras and Fortiana 1998) the problem of typicality is considered in the context of Distance-Based Discrimination (Cuadras et al. 1997). A shortcoming of the proposed method is its implicit use of Gower's interpolation, which may cause that a new observation appears as assignable to the known groups, when it actually lies very far along the additional dimension which is lost under projection (Figure 6). Similarly, predictions in Distance Based Regression (see, e.g., Cuadras 1989, Cuadras et al. 1996) implicitly rely upon Gower's interpolation and will be systematically biased for observations with a significant projection error.

This paper is organized as follows: In Section 2 we establish notations and assumptions. In Section 3 we evaluate this squared error term, rederiving Gower's own formula. In Section 4 we address the problem of learning the probability distribution of the observed error term in a given statistical context, showing a particular instance thereof for an ideal gaussian configuration. For actual data, in Section 5 we suggest a method for estimating the squared error term distribution under the null hypothesis that the underlying geometrical quantity is zero. Concluding remarks are in Section 6 and Appendix A contains accessory derivations.

## 2. Notations and assumptions

We will denote by  $\Omega_n = \{\omega_1, \dots, \omega_n\}$  a set of  $n$  points or observational units, and by  $D_n = (d_{ij})$  an  $n \times n$  matrix of squared interdistances between them. We assume throughout that the Euclidean property is satisfied, thus for some  $k > 0$  there exists a set of  $n$  vectors,  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , in a  $k$ -dimensional Euclidean space such that

$$\|\mathbf{x}_i - \mathbf{x}_j\|^2 = d_{ij}, \quad 1 \leq i, j \leq n.$$

Such a set of vectors is called a *Euclidean configuration of  $D_n$*  or of  $\Omega_n$ , and the ambient Euclidean space will be denoted by  $\mathcal{F}_n \subset \mathbb{R}^k$ . The usual assumption that such configurations are centered entails no loss of generality. Let

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{pmatrix}$$

be an  $n \times k$  matrix containing the  $\mathbf{x}_i$  vectors as rows. The matrix equation  $\mathbf{1}' \cdot \mathbf{X} = \mathbf{0}$ , where  $\mathbf{1}_n$  is the  $n \times 1$  column of ones, expresses the centering condition. Clearly there are many centered Euclidean configurations of a given  $D_n$ . More precisely, a centered matrix with  $n$  rows,  $\mathbf{X}$ , is a centered Euclidean configuration of  $D_n$  if its *inner products matrix*  $\mathbf{G}_n = \mathbf{X} \cdot \mathbf{X}'$  satisfies the equality

$$\mathbf{G}_n = -\frac{1}{2} \mathbf{J}_n \cdot \mathbf{D}_n \cdot \mathbf{J}_n, \quad (1)$$

where  $\mathbf{J}_n = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \cdot \mathbf{1}_n'$  is the *centering matrix*. This same equation (1), interpreted now as a definition of  $\mathbf{G}_n$ , leads to the statement of Schoenberg's theorem (Schoenberg 1935), that  $D_n$  is a Euclidean distance matrix if  $\mathbf{G}_n$  is positive semidefinite. Expansion of (1) gives:

$$\mathbf{G}_n = -\frac{1}{2} \left( \mathbf{D}_n - \bar{D}_n \cdot \mathbf{1}_n' - \mathbf{1}_n \cdot \bar{D}_n' + \bar{\bar{D}}_n \mathbf{1}_n \cdot \mathbf{1}_n' \right),$$

where

$$\bar{D}_n = \frac{1}{n} \mathbf{D}_n \cdot \mathbf{1}_n, \quad \bar{\bar{D}}_n = \frac{1}{n} \mathbf{1}_n' \cdot \bar{D}_n = \frac{1}{n^2} \mathbf{1}_n' \cdot \mathbf{D}_n \cdot \mathbf{1}_n.$$

Additional notations, also used in (Cuadras et al. 1997), are the *geometric variability of  $D_n$* , defined as:

$$v_n = \frac{1}{2} \bar{\bar{D}}_n, \quad (2)$$

and the column  $\mathbf{g}_n$ , containing the diagonal entries in  $\mathbf{G}_n$ , the  *$n$  proximity functions* of the elements of  $\Omega_n$ , i.e., the squared norms of their Euclidean representatives in  $\mathcal{F}_n$ ,

$$\mathbf{g}_n = \text{diag}(\mathbf{G}_n) = \bar{D}_n - \frac{1}{2} \bar{\bar{D}}_n \mathbf{1}_n = \bar{D}_n - v_n \mathbf{1}_n = \begin{pmatrix} \|\mathbf{x}_1\|^2 \\ \vdots \\ \|\mathbf{x}_n\|^2 \end{pmatrix}. \quad (3)$$

Note that these squared norms depend only on the distances themselves; in particular they are invariant across all centered Euclidean configurations of  $D_n$ . The mean squared norm is equal to the geometric variability,  $\frac{1}{n} \mathbf{1}_n' \cdot \mathbf{g}_n = \frac{1}{n} \text{tr}(\mathbf{G}_n) = v_n$ .

### 3. The error term

By adding a new  $\omega_{n+1}$  to  $\Omega_n$ , we have to consider an  $(n+1) \times (n+1)$  matrix of squared interdistances of the augmented set  $\Omega_{n+1} = \Omega_n \cup \{\omega_{n+1}\}$ , which may be represented as a

$2 \times 2$  block matrix,

$$\mathbf{D}_{n+1} = \left( \begin{array}{c|c} \mathbf{D}_n & \mathbf{d}_n \\ \hline \mathbf{d}_n' & 0 \end{array} \right), \quad (4)$$

where the column  $\mathbf{d}_n$  contains the  $n$  squared distances from  $\omega_{n+1}$  to  $\omega_i \in \Omega_n$ . As in (1),

$$\begin{aligned} \mathbf{G}_{n+1} &= -\frac{1}{2} \mathbf{J}_{n+1} \cdot \mathbf{D}_{n+1} \cdot \mathbf{J}_{n+1} \\ &= -\frac{1}{2} (\mathbf{D}_{n+1} - \bar{\mathbf{D}}_{n+1} \cdot \mathbf{1}_{n+1}' - \mathbf{1}_{n+1} \cdot \bar{\mathbf{D}}_{n+1}' + \bar{\bar{\mathbf{D}}}_{n+1} \mathbf{1}_{n+1} \cdot \mathbf{1}_{n+1}'). \end{aligned} \quad (5)$$

Computing in terms of previously defined quantities and of

$$\bar{\mathbf{d}}_n = \frac{1}{n} \mathbf{1}_n' \cdot \mathbf{d}_n, \quad (6)$$

$$\phi_{n+1}^2 = \bar{\mathbf{d}}_n - v_n,$$

we obtain:

$$\bar{\mathbf{D}}_{n+1} = \frac{1}{n+1} \mathbf{D}_{n+1} \cdot \begin{pmatrix} \mathbf{1}_n \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{n}{n+1} \bar{\mathbf{D}}_n + \frac{1}{n+1} \mathbf{d}_n \\ \frac{n}{n+1} \bar{\mathbf{d}}_n \end{pmatrix} \quad (7)$$

and

$$\bar{\bar{\mathbf{D}}}_{n+1} = \frac{1}{n+1} (\mathbf{1}_n', 1) \cdot \bar{\mathbf{D}}_{n+1} = \frac{1}{n+1} \left( \frac{n^2}{n+1} \bar{\bar{\mathbf{D}}}_n + 2 \frac{n}{n+1} \bar{\mathbf{d}}_n \right). \quad (8)$$

In particular, the geometric variability of the augmented set is:

$$v_{n+1} = \left( \frac{n}{n+1} \right)^2 \left( v_n + \frac{1}{n} \bar{\mathbf{d}}_n \right). \quad (9)$$

Similarly, the vector of proximities  $\mathbf{g}_{n+1} = \text{diag}(\mathbf{G}_{n+1})$  is:

$$\begin{aligned} \mathbf{g}_{n+1} &= \bar{\mathbf{D}}_{n+1} - v_{n+1} \mathbf{1}_{n+1} = \begin{pmatrix} \frac{n}{n+1} \bar{\mathbf{D}}_n + \frac{1}{n+1} \mathbf{d}_n \\ \frac{n}{n+1} \bar{\mathbf{d}}_n \end{pmatrix} - \left( \frac{n}{n+1} \right)^2 \left( v_n + \frac{1}{n} \bar{\mathbf{d}}_n \right) \begin{pmatrix} \mathbf{1}_n \\ 1 \end{pmatrix} \\ &= \begin{pmatrix} \frac{n}{n+1} \mathbf{g}_n + \frac{1}{n+1} \mathbf{d}_n - \frac{n}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n \\ \left( \frac{n}{n+1} \right)^2 \phi_{n+1}^2 \end{pmatrix}, \end{aligned}$$

thus, denoting by  $\mathbf{x}_{n+1}$  the vector representing  $\omega_{n+1}$  in a centered Euclidean configuration of  $\Omega_{n+1}$ , we have

$$\|\mathbf{x}_{n+1}\|^2 = \left( \frac{n}{n+1} \right)^2 \phi_{n+1}^2. \quad (10)$$

Note that the form of (10) makes it clear that the equality holds for *any* centered Euclidean configuration of  $\Omega_{n+1}$ , a result which will become useful below.

On the other hand, we have the *superimposed coordinates* of  $\omega_{n+1}$ ,  $\hat{\mathbf{x}}_{n+1} \in \mathcal{F}_n$ , result of Gower's interpolation, or *add-a-point* formula (Gower 1968, see also Gower and Hand 1996, pp. 246ff). In order to derive it, we assume a Euclidean configuration for  $\Omega_{n+1}$ ,

$$\mathbf{Y} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{n+1} \end{pmatrix}, \quad (11)$$

in a subspace  $\mathcal{F}_{n+1} \subset \mathbb{R}^{k+1}$ .  $\mathbf{Y}$  consists of  $n$  vectors from  $\mathbf{X}$ ,

$$\mathbf{y}_i = (\mathbf{x}_i, 0) \in \mathbb{R}^{k+1}, \quad \text{for } 1 \leq i \leq n,$$

plus another  $\mathbf{y}_{n+1} = (\hat{\mathbf{x}}_{n+1}, r_{n+1}) \in \mathbb{R}^{k+1}$  with a possibly nonnull  $(k+1)$ -th coordinate. Below it will become clear that such a configuration is always possible provided that  $\mathbf{D}_{n+1}$  satisfies the Euclidean condition. Since  $\mathbf{Y}$  is obviously noncentered,  $\mathbf{y}_{n+1} \neq \mathbf{x}_{n+1}$  in (10). From the  $n$  entries in  $\mathbf{d}_n = (d_1, \dots, d_n)'$ , defined in (4) as the squared distances between  $\omega_{n+1}$  and the elements of  $\Omega_n$ , we have the following set of equations,

$$\begin{aligned} d_i &= (\mathbf{y}_{n+1} - \mathbf{y}_i) \cdot (\mathbf{y}_{n+1} - \mathbf{y}_i)' \\ &= \|\mathbf{y}_{n+1}\|^2 + \|\mathbf{y}_i\|^2 - 2\mathbf{y}_{n+1} \cdot \mathbf{y}_i' \\ &= \|\mathbf{y}_{n+1}\|^2 + \|\mathbf{x}_i\|^2 - 2\hat{\mathbf{x}}_{n+1} \cdot \mathbf{x}_i', \end{aligned}$$

which, taking into account (3), can be written as a single matricial equation as

$$(\mathbf{d}_n - \mathbf{g}_n)' = \|\mathbf{y}_{n+1}\|^2 \mathbf{1}_n - 2\hat{\mathbf{x}}_{n+1} \cdot \mathbf{X}',$$

Posmultiplication by  $\mathbf{X}$  gives

$$(\mathbf{d}_n - \mathbf{g}_n)' \cdot \mathbf{X} = -2\hat{\mathbf{x}}_{n+1} \cdot (\mathbf{X}' \cdot \mathbf{X}),$$

since  $\mathbf{X}$  is centered. If  $\mathbf{X}' \cdot \mathbf{X}$  is nonsingular,  $\hat{\mathbf{x}}_{n+1}$  is given by the usual *add-a-point* formula

$$\hat{\mathbf{x}}_{n+1} = \frac{1}{2} (\mathbf{g}_n - \mathbf{d}_n)' \cdot \mathbf{X} \cdot (\mathbf{X}' \cdot \mathbf{X})^{-1}, \quad (12)$$

otherwise the inverse must be replaced by a  $g$ -inverse, resulting in a nonunique  $\hat{\mathbf{x}}_{n+1}$ . Nonetheless, its squared norm

$$\|\hat{\mathbf{x}}_{n+1}\|^2 = \frac{1}{4} (\mathbf{g}_n - \mathbf{d}_n)' \cdot \mathbf{X} \cdot (\mathbf{X}' \cdot \mathbf{X})^{-2} \cdot \mathbf{X}' \cdot (\mathbf{g}_n - \mathbf{d}_n),$$

is *invariant* under different choices of a  $g$ -inverse, since  $\mathbf{X} \cdot (\mathbf{X}' \cdot \mathbf{X})^{-2} \cdot \mathbf{X}' = \mathbf{G}_n^+$ , the Moore-Penrose pseudo-inverse of  $\mathbf{G}_n$ . Hence

$$\|\hat{\mathbf{x}}_{n+1}\|^2 = \frac{1}{4} (\mathbf{g}_n - \mathbf{d}_n)' \cdot \mathbf{G}_n^+ \cdot (\mathbf{g}_n - \mathbf{d}_n). \quad (13)$$



The Euclidean configuration  $\mathbf{Y}$ , defined in (11), can easily be centered by subtracting from each row their overall average, giving, in an obvious block matrix notation,

$$\mathbf{Y}_0 = \left( \begin{array}{c|c} \mathbf{X} - \frac{1}{n+1} \mathbf{1}_n \cdot \hat{\mathbf{x}}_{n+1} & -\frac{1}{n+1} \mathbf{1}_n r_{n+1} \\ \hline \frac{n}{n+1} \hat{\mathbf{x}}_{n+1} & \frac{n}{n+1} r_{n+1} \end{array} \right). \quad (14)$$

Equating the squared norm of its  $(n+1)$ -th vector,

$$\left( \frac{n}{n+1} \right)^2 \left( \|\hat{\mathbf{x}}_{n+1}\|^2 + r_{n+1}^2 \right)$$

with  $\|\mathbf{x}_{n+1}\|^2$  in (10), we have an expression for the squared error term, the squared norm of the residual of the orthogonal projection  $\mathcal{F}_{n+1} \rightarrow \mathcal{F}_n$ ,

$$r_{n+1}^2 = \phi_{n+1}^2 - \frac{1}{4} (\mathbf{g}_n - \mathbf{d}_n)' \cdot \mathbf{G}_n^+ \cdot (\mathbf{g}_n - \mathbf{d}_n). \quad (15)$$

An alternative derivation of Gower's interpolation formula, and a further insight, can be obtained by explicitly writing the relationship between  $\mathbf{G}_n$  and  $\mathbf{G}_{n+1}$  (see Appendix A).

## 4. Distribution of the error term

In a statistical setting, for instance in the problem of typicality referred to in the Introduction, we would like to decide whether the residual term is significantly non-null or else it can be safely ignored. To this end we need to know the probability distribution of the (observed) residual under the null hypothesis that the geometrical residual is equal to zero. Analytical derivation of such a distribution is an exceedingly complex task for any realistic situation, hence for actual data resampling methods provide sensible estimations —see Section 5. Meanwhile, and in order to gain some theoretical insight on the qualitative properties of error term distributions we pose an ideal scenario, in which we assume that normality holds *for the Euclidean configuration*: A random  $(n+1)$ -th point is added to  $n$  points in  $\mathcal{F}_n = \mathbb{R}^p$ , realized as the hyperplane of all vectors in  $\mathcal{F}_{n+1} = \mathbb{R}^{p+1}$  with a null  $(p+1)$ -th coordinate. Let us denote them by

$$\boldsymbol{\mu}_i = (\boldsymbol{\theta}_i, 0) \in \mathbb{R}^{p+1}, \quad \boldsymbol{\theta}_i \in \mathbb{R}^p, \quad 1 \leq i \leq n.$$

In order to evaluate the projection we do not need the actual coordinates  $\boldsymbol{\theta}_i$ , only the subspace they generate. For  $n \geq p$  we can, with probability 1, take them to be a set of  $p$  constant vectors generating  $\mathcal{F}_n$ . The additional random point will be

$$\mathbf{y}_{p+1} = (\mathbf{x}_{p+1}, t) \sim \mathbf{N}(\boldsymbol{\mu}_{p+1}, \boldsymbol{\Sigma}_{p+1}),$$

where  $\boldsymbol{\mu}_{p+1} = (\boldsymbol{\theta}_{p+1}, \tau) \in \mathbb{R}^{p+1}$  and  $\boldsymbol{\Sigma}_{p+1}$ , which can be assumed nonsingular, has the block structure:

$$\boldsymbol{\Sigma}_{p+1} = \left( \begin{array}{c|c} \boldsymbol{\Sigma}_p & \boldsymbol{\Sigma}_{xt} \\ \hline \boldsymbol{\Sigma}'_{xt} & \sigma_t^2 \end{array} \right), \quad (16)$$

where  $\boldsymbol{\Sigma}_p$  is  $p \times p$ ,  $\boldsymbol{\Sigma}_{xt}$  is  $p \times 1$  and  $\sigma_t^2$  is a scalar. The projection of  $\mathbf{y}_{p+1}$  on the subspace generated by the rows of

$$\mathbf{M} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \vdots \\ \boldsymbol{\mu}_p \end{pmatrix}$$

can be expressed as:

$$\langle \mathbf{y}_{p+1}, \mathbf{M} \rangle \langle \mathbf{M}, \mathbf{M} \rangle^{-1} \mathbf{M},$$

where the inner products are to be computed with the Mahalanobis metric:

$$\langle \mathbf{y}_{p+1}, \mathbf{M} \rangle = \mathbf{y}_{p+1} \cdot \boldsymbol{\Sigma}_{p+1}^{-1} \cdot \mathbf{M}', \quad \langle \mathbf{M}, \mathbf{M} \rangle = \mathbf{M} \cdot \boldsymbol{\Sigma}_{p+1}^{-1} \cdot \mathbf{M}'.$$

Using a well-known formula to invert a block nonsingular symmetric matrix (see, e.g., Seber 1984, Appendix),

$$\boldsymbol{\Sigma}_{p+1}^{-1} = \left( \begin{array}{c|c} \mathbf{A} & \mathbf{b} \\ \hline \mathbf{b}' & a \end{array} \right),$$

where

$$\mathbf{A} = \left[ \boldsymbol{\Sigma}_p - \frac{1}{\sigma_t^2} \boldsymbol{\Sigma}_{xt} \cdot \boldsymbol{\Sigma}'_{xt} \right]^{-1}, \quad \mathbf{b} = -\frac{1}{\sigma_t^2} \mathbf{A} \cdot \boldsymbol{\Sigma}_{xt}, \quad a = \frac{1}{\sigma_t^2} + \frac{1}{\sigma_t^4} \boldsymbol{\Sigma}_{xt} \cdot \mathbf{A} \cdot \boldsymbol{\Sigma}_{xt}.$$

If we assume, as we are entitled to do, that  $\{\theta_i\}$  are the  $\mathbb{R}^p$  standard basis vectors,

$$\theta_i = \mathbf{e}_i, \quad 1 \leq i \leq n = p,$$

then  $\mathbf{M} = (\mathbf{I}_p, \mathbf{0})$ ,  $\langle \mathbf{M}, \mathbf{M} \rangle^{-1} = \mathbf{A}^{-1}$  and

$$\langle \mathbf{y}_{p+1}, \mathbf{M} \rangle = (\mathbf{x}_{p+1}, t) \cdot \begin{pmatrix} \mathbf{A} \\ \mathbf{b}' \end{pmatrix} = \mathbf{x}_{p+1} \cdot \mathbf{A} + t \mathbf{b}' = \left( \mathbf{x}_{p+1} - \frac{t}{\sigma_t^2} \boldsymbol{\Sigma}'_{xt} \right) \cdot \mathbf{A}.$$

Thus, the projected vector and its residual are

$$\begin{aligned} \hat{\mathbf{y}}_{p+1} &= (\mathbf{x}_{p+1} - \frac{t}{\sigma_t^2} \boldsymbol{\Sigma}'_{xt}, 0), \\ \tilde{\mathbf{y}}_{p+1} &= \mathbf{y}_{p+1} - \hat{\mathbf{y}}_{p+1} = \left( \frac{t}{\sigma_t^2} \boldsymbol{\Sigma}'_{xt}, t \right). \end{aligned}$$

The error term is the squared norm

$$Q = \|\tilde{\mathbf{y}}_{p+1}\|^2 = t^2 \left( 1 + \frac{\|\boldsymbol{\Sigma}_{xt}\|^2}{\sigma_t^4} \right).$$

Since  $t$ , the  $(p + 1)$ -th coordinate of  $y_{p+1}$ , is normally distributed with mean  $\tau$  and variance  $\sigma_t^2$ , under the null hypothesis that  $\tau$  is null we have  $t = Z \sigma_t$ , where  $Z \sim N(0,1)$ . Hence  $Q$  is a multiple of a  $\chi_1^2$  variate, with the proportionality factor

$$\kappa = \sigma_t^2 + \frac{\|\Sigma_{xt}\|^2}{\sigma_t^2}.$$

Under  $H_1 : \tau \neq 0$  we can write

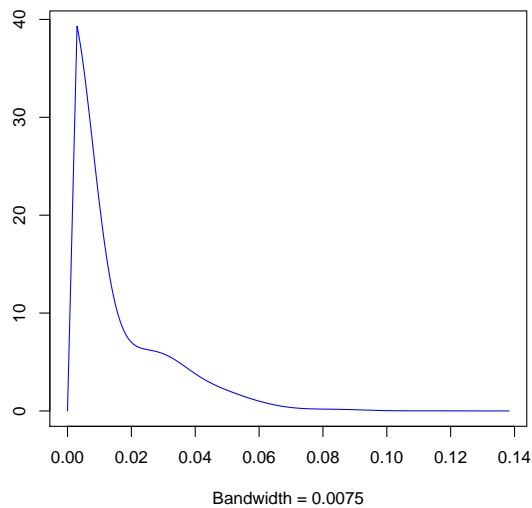
$$t = \tau + \sigma_t Z = \sigma_t (\delta + Z), \quad \text{where } \delta = \tau/\sigma_t,$$

then  $t^2 = \sigma_t^2 W$  and  $Q = \kappa W$ , where  $W = (\delta + Z)^2$  is a noncentral  $\chi_1^2$  with  $\delta^2$  as its noncentrality parameter.

## 5. Computations with real data

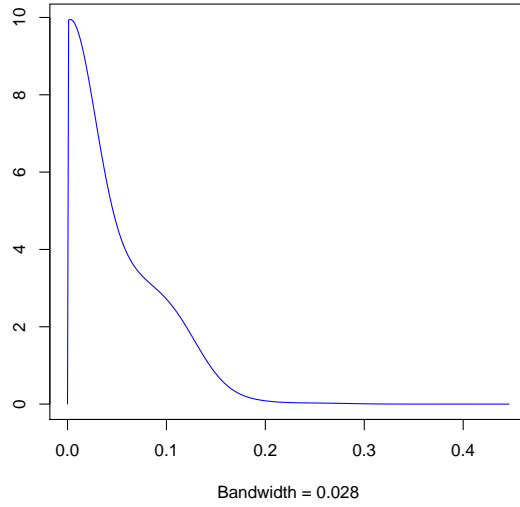
A dataset consisting of observations on an  $n$ -set  $\Omega_n = \{\omega_1, \dots, \omega_n\}$  is understood as an  $n$ -sample drawn from a population  $\Omega$ . Within this framework the Euclidean space  $\mathcal{F}_n$  derived from  $\Omega_n$  as explained above is an “estimation” of the Euclidean space  $\mathcal{F}$  likewise associated with  $\Omega$ . By hypothesis each  $\omega_i$  in  $\Omega_n$  satisfies the null hypothesis that it can be exactly represented by a vector in  $\mathcal{F}$  or, equivalently, that when projected on  $\mathcal{F}$  the error

Figure 7: Density estimation for Data set 1



term is zero. Thus we can implement a cross-validatory scheme by randomly splitting  $\Omega_n$  in two disjoint subsets,  $\Omega_A$  and  $\Omega_B$ , with  $n_A$  and  $n_B$  observations, respectively. We take

Figure 8: Density estimation for Data set 2



$\mathcal{F}_A$ , the space generated by the Euclidean configuration of the  $n_A$  observations in  $\Omega_A$ , as an estimation of  $\mathcal{F}$ , and the  $n_B$  residuals obtained by projecting on  $\mathcal{F}_A$  the observations in  $\Omega_B$ , i.e., by applying (15), as samples of the observed squared error term under  $H_0$ . We perform a suitable number  $\nu \leq \binom{n}{n_A}$  of such splittings so to yield  $N = \nu n_B$  values of  $r^2$  to estimate its probability density. In particular, an expedient choice is the jackknife-like, or leave-one-out, computation, with  $n_A = n - 1$  and  $N = n$ .

To illustrate this procedure we provide two examples, in the contexts of classification (Data set 1) and regression (Data set 2). Data set 1, taken from Krzanowski (1975), has observations of six continuous and three qualitative variables on 186 subjects from two groups, of 99 and 87 cases, respectively. We shall not need to delve into specifics on meaning of variables and groups. Let us only point out that for such data, i.e., intended for a classification problem, the subsets  $\Omega_A$  and  $\Omega_B$  should be formed with proportional numbers of cases in each group, as representing the whole sample. Data set 2 is Gasoline Mileage Data, from Henderson and Velleman (1981), where a continuous numerical response is regressed on 8 continuous numerical and two categorical predictors, observed for 32 cars. In both examples the distance matrix is derived from the mixed predictor variables using the widely used metric derived from Gower's omnibus similarity coefficient (Gower 1971). Figures 7 and 8 show kernel-smoothed  $r^2$  density estimations for both data sets, derived from  $\nu = 2500$  resamplings with  $n_B \approx 30\%$  of the total sample size  $n$ . Each resampling is realized as a boolean  $n$ -vector, randomly drawn from a Bernoulli distribution with parameter 0.30.

## 6. Concluding remarks

In the first place, the present work provides some additional insight on the mechanics of distance-based prediction methods, especially on the role played by Gower's interpolation and the dangers of implicitly and uncritically assuming the equal Euclidean dimension hypothesis. Figure 6 shows an example, with an artificial data-set, where consideration of the error term is highly relevant to avoid misleading results. As a tangible result of our investigation we propose an auxiliary statistical tool, useful as a cautionary device to verify the validity of the aforementioned hypothesis.

---

## Appendices

### A. Expressing $G_{n+1}$ in terms of $G_n$ and Gower's interpolation

Direct computation from the equality

$$-2 G_{n+1} = D_{n+1} - \mathbf{g}_{n+1} \cdot (\mathbf{1}_n', 1) - \begin{pmatrix} \mathbf{1}_n \\ 1 \end{pmatrix} \cdot \mathbf{g}_{n+1}',$$

gives

$$G_{n+1} = \left( \frac{\mathbf{G}_n - \frac{1}{2(n+1)} [(\mathbf{g}_n - \mathbf{d}_n) \cdot \mathbf{1}_n' + \mathbf{1}_n \cdot (\mathbf{g}_n - \mathbf{d}_n)'] - \frac{n}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n \cdot \mathbf{1}_n'}{\frac{1}{2} \left(\frac{n}{n+1}\right) \left[ (\mathbf{g}_n - \mathbf{d}_n)' + \frac{(n-1)}{(n+1)} \phi_{n+1}^2 \mathbf{1}_n' \right]} \right) \left( \frac{\frac{1}{2} \left(\frac{n}{n+1}\right) \left[ (\mathbf{g}_n - \mathbf{d}_n) + \frac{(n-1)}{(n+1)} \phi_{n+1}^2 \mathbf{1}_n \right]}{\left(\frac{n}{n+1}\right)^2 \phi_{n+1}^2} \right), \quad (17)$$

where notations are as in Section 3. A somewhat clearer expression is obtained by defining

$$\mathbf{k}_n = \frac{1}{2} \left[ (\mathbf{g}_n - \mathbf{d}_n) + \phi_{n+1}^2 \mathbf{1}_n \right]. \quad (18)$$

Note that  $\mathbf{1}_n \cdot \mathbf{k}_n = 0$ , taking into account (3) and (6). Thus,

$$G_{n+1} = \left( \frac{\mathbf{G}_n - \frac{1}{n+1} (\mathbf{k}_n \cdot \mathbf{1}_n' + \mathbf{1}_n \cdot \mathbf{k}_n') + \frac{1}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n \cdot \mathbf{1}_n'}{\frac{n}{n+1} \mathbf{k}_n' - \frac{n}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n'} \right) \left( \frac{\frac{n}{n+1} \mathbf{k}_n - \frac{n}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n}{\left(\frac{n}{n+1}\right)^2 \phi_{n+1}^2} \right). \quad (19)$$

Since also  $\mathbf{G}_{n+1} = \mathbf{Y}_0 \cdot \mathbf{Y}_0'$ , comparing with

$$\mathbf{G}_{n+1} = \left( \frac{\mathbf{X} \cdot \mathbf{X}' - \frac{1}{n+1} [\mathbf{X} \cdot \hat{\mathbf{x}}_{n+1}' \cdot \mathbf{1}_n' + \mathbf{1}_n \cdot \hat{\mathbf{x}}_{n+1} \cdot \mathbf{X}'] + \frac{1}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n \cdot \mathbf{1}_n'}{\frac{n}{n+1} \hat{\mathbf{x}}_{n+1} \cdot \mathbf{X}' - \frac{n}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n'} \right) \left( \frac{\frac{n}{n+1} \mathbf{X} \cdot \hat{\mathbf{x}}_{n+1}' - \frac{n}{(n+1)^2} \phi_{n+1}^2 \mathbf{1}_n}{\left(\frac{n}{n+1}\right)^2 \phi_{n+1}^2} \right), \quad (20)$$

we recover the interpolation formula (12) equating  $\mathbf{k}_n$  in (18) with  $\mathbf{X} \cdot \hat{\mathbf{x}}_{n+1}$  in (20):

$$\mathbf{X} \cdot \hat{\mathbf{x}}_{n+1} = \frac{1}{2} \left[ (\mathbf{g}_n - \mathbf{d}_n) + \phi_{n+1}^2 \mathbf{1}_n \right].$$

## References

- Borg, I. and P. Groenen (1997). *Modern Multidimensional Scaling*. New York: Springer-Verlag.
- Cox, T. F. and M. A. A. Cox (1994). *Multidimensional Scaling*. London: Chapman & Hall.
- Cuadras, C. M., C. Arenas, and J. Fortiana (1996). Some computational aspects of a distance-based model for prediction. *Communications in Statistics B. Simulation and Computation* 25, 593–609.
- Cuadras, C. M. (1989). Distance analysis in discrimination and classification using both continuous and categorical variables. In Y. Dodge (Ed.), *Statistical Data Analysis and Inference*, Amsterdam, The Netherlands, pp. 459–473. North-Holland Publishing Co.
- Cuadras, C. M. and J. Fortiana (1998). Typicality in discriminant analysis with mixed variables. In A. Rizzi, M. Vichi, and H. H. Bock (Eds.), *Advances in Data Science and Classification. Proceedings of the 6th Conference of the International Federation of Classification Societies (IFCS-98)*. Rome, Italy, July 21-24, 1998, Heidelberg, pp. 82–85. Springer-Verlag.
- Cuadras, C. M., J. Fortiana, and F. Oliva (1997). The proximity of an individual to a population with applications in discriminant analysis. *Journal of Classification* 14, 117–136.
- Gower, J. C. (1968). Adding a point to vector diagrams in multivariate analysis. *Biometrika* 55, 582–585.
- Gower, J. C. (1971). A general coefficient of similarity and some of its properties. *Biometrics* 27, 857–874.
- Gower, J. C. and D. J. Hand (1996). *Biplots*. London, UK: Chapman & Hall.

- Henderson, H. V. and P. F. Velleman (1981). Building multiple regression models interactively. *Biometrics* 37, 391–411.
- Krzanowski, W. J. (1975). Discrimination and classification using both binary and continuous variables. *Journal of the American Statistical Association* 70, 782–790.
- Krzanowski, W. J. and F. H. C. Marriott (1994). *Multivariate Analysis. Part 1, Distributions, ordination and inference*. London: Edward Arnold.
- Rao, C. R. (1973). *Linear Statistical Inference and its Applications*. New York, NY, USA: John Wiley & Sons.
- Schoenberg, I. J. (1935). Remarks to Maurice Fréchet's article "Sur la Définition Axiomatique d'une Classe d'Espaces Distanciés Vectoriellement Applicables sur l'Espace de Hilbert". *The Annals of Mathematics* 36, 724–732.
- Seber, G. A. F. (1984). *Multivariate observations*. New York, NY, USA: John Wiley & Sons.



## Relació dels últims Preprints publicats:

- **361** *Strict modules and homotopy modules in stable homotopy.* Javier J. Gutiérrez. MS Subject Classification: 55P43, 18E30, 55P42. October 2004.
- **362** *Lattice approximation for a stochastic wave equation.* Lluís Quer-Sardanyons and Marta Sanz-Solé. MS Subject Classification: 60H35, 60H15. October 2004.
- **363** *An adaptive goodness-of-fit test.* Aurea Grané and Josep Fortiana. AMS Subject Classification: 62G10, 62G30, 62E15. November 2004.
- **364** *A stabilization phenomenon for a class of stochastic partial differential equations.* David Nualart and Pierre-A. Vuillermot. AMS Subject Classification: 60H15, 35R60. November 2004.
- **365** *Homotopy localization of groupoids.* Carles Casacuberta, Marek Golasinski, and Andrew Tonks. AMS Subject Classification (2000): Primary 18B40, 55P60; Secondary 18E35, 18A40. November 2004.
- **366** *Bifurcations of three-dimensional diffeomorphisms with non-simple quadratic homoclinic tangencies and generalized Hénon maps.* S.V. Gonchenko, V.S. Gonchenko, and J.C. Tatjer. AMS Subject Classification: 37G25, 37C29, 37G05. December 2004.
- **367** *Bootstrapping repeated measures data in a nonlinear mixed-models context.* Jordi Ocaña, Rachid El Halimi, M. Carme Ruiz de Villa, and Josep A. Sánchez. AMS Subject Classification: 62F35, 62F40, 62G09, 62P10. February 2005.
- **368** *The orthogonal subcategory problem in homotopy theory.* Carles Casacuberta and Boris Chorny. AMS Subject Classification: 55U35, 55P60, 18G55. February 2005.
- **369** *On two fragments with negation and without implication of the logic of residuated lattices.* Félix Bou, Àngel García-Cerdaña, and Ventura Verdú. AMS Subject Classification: 03B47, 03B50, 03F99, 06D15, 06B99. February 2005.
- **370** *Variational solutions for partial differential equations driven by a fractional noise.* David Nualart and Pierre-A. Vuillermot. AMS Subject Classification: 60H15, 35R60. March 2005.
- **371** *Malliavin calculus for stochastic differential equations driven by a fractional Brownian motion.* David Nualart and Bruno Saussereau. AMS Subject Classification: 60H05, 60H07. March 2005.
- **372** *Glivenko's Theorem in algebraizable logics.* Antoni Torrens. AMS Subject Classification: 03G, 08C15. March 2005.
- **373** *Power variation of some integral long-memory processes.* José Manuel Corcuera, David Nualart, and Jeannette H.C. Woerner. AMS Subject Classification: 60F05, 60G15, 60G18, 62M99. April 2005.
- **374** *Intersection local time for two independent fractional Brownian motions.* David Nualart and Salvador Ortiz. AMS Subject Classification (2000): 60G15, 60F05, 60F25, 60G18, 60J55. July 2005.
- **375** *On the infinite-valued Lukasiewicz logic that preserves degrees of truth.* Josep Maria Font, Àngel J. Gil, Antoni Torrens, and Ventura Verdú. AMS Subject Classification (2000): 03B50, 03G20, 06D35, 03B22. July 2005.
- **376** *A lattice scheme for stochastic partial differential equations of elliptic type in dimension  $d \geq 4$ .* Teresa Martínez and Marta Sanz-Solé. AMS Subject Classification (2000): 60H15, 60H35, 35J05. September 2005.
- **377** *Approximation of rough paths of fractional Brownian motion.* Annie Millet and Marta Sanz-Solé. AMS Subject Classification (2000): Primary 60G15; Secondary 60H05, 60H07. September 2005.