

---

# Retrospective correlation

---

Francesco Mulargia

*Dipartimento di Fisica, Settore di Geofisica, Università di Bologna, Italy*

## Abstract

One of the main tools in phenomenological studies is the identification of correlations among different processes. This is essentially effected in retrospective with the specific aim of finding a positive result, and that leads to a parameter optimization which introduces a bias, so far only marginally considered, in the significance level of the results. If the correlation can be validated in a forward study in which parameters are kept fixed, such a bias is irrelevant. Unfortunately, forward studies are often infeasible for either cost or intrinsic reasons. This is the case of geophysics, due to the comparatively long time scale of recurrence of the phenomena. Unbiased estimates can be obtained in retrospective if each of the optimal choices is properly identified and accounted for. An estimate of the bias is made in a specific case, which can be written in closed form. While simulation confirms its good performance, the latter shows that apparently highly significant retrospective correlations may be insignificant.

**Key words** *earthquake prediction – statistical analysis*

## 1. Introduction

Testing the validity of a hypothesized correlation deals primarily with establishing if it occurs more often than chance. The first problem is then how to define «chance». In geophysics, this is an often encountered problem, and a typical case is that of earthquake precursors. In the latter case, to gauge the correlation against «chance» the set of the hypothesized precursors is typically tested for performance against both the real catalog and a randomized seismic catalog. The problem is complicated due to the fact that randomized catalogs cannot be uniquely defined since the present knowledge of earthquake physics does not allow either a deterministic or a statistical satisfactory model to be derived (Kagan, 1994; Ben Menahem, 1995).

However, the lack of unique models is not the major cause of difficulties in establishing phenomenological correlations. There is a problem of greater importance, which has so far received only marginal consideration. This is the bias in the correlation estimates. In fact, the common procedure is to analyze data in retrospective and optimally select (rather than keeping *a priori* fixed or randomly choose) all free parameters in order to achieve the best apparent result, *i.e.* the correlation the least apparently probable by chance. The problem can be readily solved if the data are originated by repeatable experiments, because retrospective analysis can be used to cast an issue of correlation which can then be validated in an independent, typically forward, study (Mulargia and Gasperini, 1996). On the contrary, if only retrospective studies are feasible, a substantial distortion of reality in favour of the correlation can be expected. This is the case of some geophysical phenomena such as earthquakes and volcanic eruptions, since the comparatively long time scale of geologic processes would require waiting at least several decades for the forward occurrence of a sufficient number of events.

---

*Mailing address:* Prof. Francesco Mulargia, Dipartimento di Fisica, Settore di Geofisica, Università di Bologna, Viale C. Berti Pichat 8, 40127 Bologna, Italy; e-mail: mulargia@ibogfs.df.unibo.it

Unbiased retrospective estimates are possible if all the steps in parameter optimization are identified and properly taken into account in the number of favourable *vs.* number of total cases. This task is a difficult one because the tailoring of some parameters is often implicit, such as in the selection of the region in which the correlation is operative. In geophysics this selection is justified in terms of «tectonics», a term that implies more or less subjective *ad hoc* choices, which are not explicitly discussed (cf. Keilis-Borok *et al.*, 1988).

The present paper is devoted to estimating the importance of the bias induced by retrospective adjustment of parameters. This will be achieved by studying correlation in a particular simple case, regarding a single variable, as in the case in which two series are analyzed by considering only the times of occurrence of the events. In such a case correlation is commonly defined as *association*.

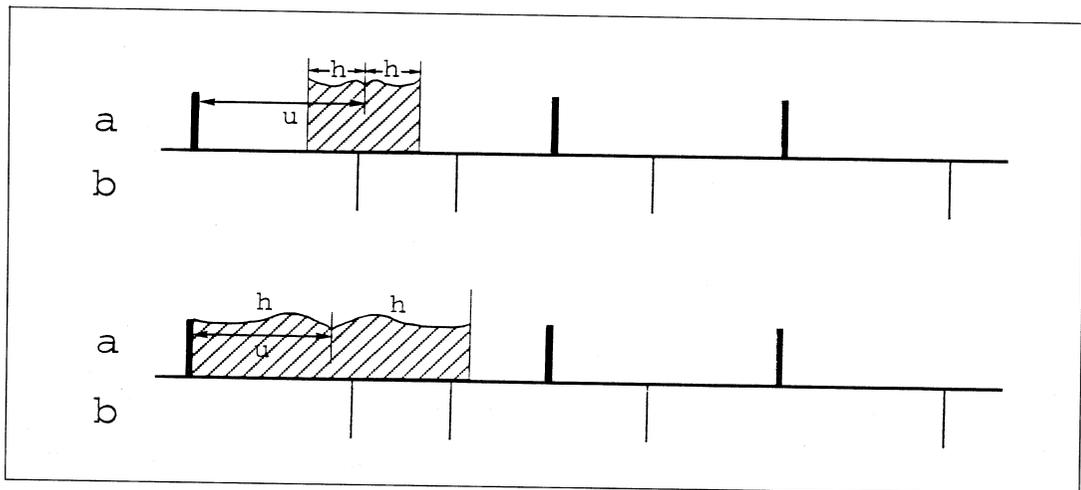
## 2. Retrospective association

In studying the association between two series of events, the typical record consists of a number  $m$  of type  $a$  events, and of a number  $n$

of type  $b$  events, in a time interval of length  $T$  time units (fig. 1). The series are said to be associated in time on the basis of the number of events  $b$  following the  $a$  events by a fixed lag  $u$ , *i.e.* when  $b$  events occur within windows of amplitude  $2h$  centered  $u$  instants after an event  $a$ . For simplicity, let us assume that the number of correlated events  $s$  is fixed *a priori*. This number is ideally equal to the minimum between  $m$  and  $n$ , but nobody would dismiss a correlation occurring on, say,  $s = 35$  among two series of  $m = 43$  and  $n = 52$  events. With this constraint, the only variable that can be optimally chosen is then the lag  $u$ , while the probability of finding  $s$  events associated in a window of width  $2h$  depends on the specific model assumed. To proceed further let us therefore first consider this.

### 2.1. Assumption of a specific model

A simple, yet widely applicable, model (Cox, 1955) puts no constraint on the events  $a$ , assumes that events  $b$  follow a simple (stationary) Poisson process, that  $m2h$  is small with respect to the total length  $T$  of the period analyzed, and that  $2h$  is small enough that the oc-



**Fig. 1.** The association parameters between a series of events  $a$  and a series of events  $b$  are the lag  $u$  and the window width  $2h$ . The general case is shown at the top, the case  $u = h$  is shown at the bottom.

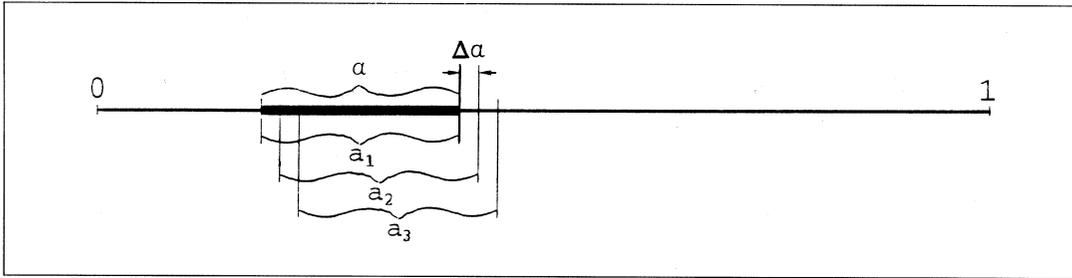


Fig. 2. Partition of the unit probability interval with equivalent association intervals.

currence of two  $a$  or two  $b$  events within  $2h$  can be neglected. Then, the probability of occurrence of one  $b$  event within  $(u-h, u+h)$  from an  $a$  event is approximately equal to  $m2h/T$ , and the probability that  $x$  events  $a$  and  $b$  are associated with parameters  $u$ ,  $h$  is binomial

$$p(x) = \binom{n}{x} (m2h/T)^x (1 - m2h/T)^{n-x}. \quad (2.1)$$

Note that if a lack of constraints were desirable on the  $b$  series, rather than on the  $a$  series, the reasoning could be reversed by simply using negative  $u$  values.

### 3. Bias introduced by the optimal choice of $u$

We have fixed  $s$ , and found the corresponding best  $h$  from the data once given the distribution. But we have not fixed the lag  $u$ , which is indeed the very target of any correlation analysis. This is chosen by retrospective optimal selection on a specific realization, *i.e.* the available data set. Should we have another realization, we would choose another optimal value of  $u$ . This optimal selection introduces a bias in the correlation estimates, which, according to eq. (2.1), are valid if  $u$  was kept fixed *a priori* at each realization. In other words, eq. (2.1) was valid if  $u$  had been determined through an independent set of measurements.

Let us now estimate the retrospective bias. Let us then keep  $h$  fixed and assume that a particular realization has its best association with a number  $s$  of  $a$  and  $b$  events associated in the window  $(\tilde{u}-\tilde{h}, \tilde{u}+\tilde{h})$ , corresponding to an apparent low probability of random occurrence equal to  $\tilde{p}$ . Actually, we would have been equally content with finding the same low probability of association with any other different value of  $u$ . To evaluate the importance of this effect it is convenient to transfer the reasoning on the unit probability interval  $(0,1)$ .

Consider first the case in which the number of associations  $s$  is equal to  $\min\{m, n\}$ , and let the corresponding probability be  $\tilde{p}$ . We can then draw on the unit probability interval  $(0,1)$  the set  $\alpha$  corresponding to the probability that one  $b$  event falls in the window  $(\tilde{u}-\tilde{h}, \tilde{u}+\tilde{h})$  which is  $(\tilde{p})^{1/s}$  (fig. 2). A set of other equivalent intervals, all with the same probability of containing one  $b$  event, and thus with the same  $\tilde{p}$  for  $s$  events, can then be drawn by shifting the set  $\alpha$  by fixed steps of amplitude  $\Delta\alpha$ . A total of  $\tilde{j} = 1 + (1 - \alpha)/\Delta\alpha$  are needed to cover the whole probability interval  $(0,1)$ . Obviously, the coverage will not be perfect, but the remainder can be controlled by choosing small values of  $\Delta\alpha$ , like  $\Delta\alpha = (1 - \alpha)/1000$ . Now consider the fact that one is not interested in the single event  $A_1$ ;  $\{s$  events in the interval  $a_1$  of length  $\alpha\}$ , but in the union of the events  $A_1 \cup A_2 \cup A_3 \cup \dots \cup A_{\tilde{j}}$ , *i.e.* in finding  $s$  events in any of the equivalent intervals  $a_1, a_2, \dots, a_{\tilde{j}}$



cal definition of probability), is the effect of the retrospective optimal choice of  $u$ . The blow-up factor calculated according to eq. (3.2) is also shown. It is immediately seen how the theoretical correction agrees well with the observed frequency, and provides a correct estimate of the blow-up factor. Note how the true significance levels are an order of magnitude or more larger than the uncorrected levels even with such small values of  $s$ ,  $m$ , and  $n$ .

## 5. Discussion and conclusions

The standard practice in correlation analysis is to retrospectively choose all free parameters in order to obtain the best apparent association. This optimal choice introduces a bias in the estimates. If the association can be validated in a forward study in which parameters are kept fixed, this bias is irrelevant. Unfortunately, forward studies are often infeasible for their cost, or intrinsically impossible for the long time scale of evolution of the phenomena, as mostly happens in astronomy and geophysics. In this case, unbiased estimates of the significance of the association may be obtained by identifying and accounting for the optimal choices implicit in the retrospective analysis. The explicit study of one particular case shows that the blow-up

factor can be written in a simple closed form. It also shows how this factor is very large and, unless corrected, may easily lead to conclude as highly significant correlations which are, indeed, insignificant.

## Acknowledgements

I am indebted to Professor Michele Caputo for many stimulating discussions. This work was performed with contributions CNR Gruppo Nazionale Difesa dai Terremoti and MPI 60% and 40%.

## REFERENCES

- BEN-MENACHEM, A. (1995): A concise history of mainstream seismology: origins, legacy, and perspectives, *Bull. Seism. Soc. Am.*, **85**, 1202-1205.
- COX, D.R. (1955): Some statistical methods connected with series of events, *J. R. Stat. Soc. B*, **17**, 129-164.
- KAGAN, Y.Y. (1994): Observational evidence for earthquakes as a nonlinear dynamic process, *Physica D*, **77**, 160-192.
- KEILIS-BOROK, V.I., L. KNOPOFF, I.M. ROTWAIN and C.R. ALLEN (1988): Intermediate-term prediction of occurrence times of strong earthquakes, *Nature*, **335**, 690-694.
- MULARGIA, F. and P. GASPERINI (1996): Precursor candidacy and validation: the VAN case so far, *Geophys. Res. Lett.*, **23**, 1323-1326.