

The Unexploited Mines of Academic and Official Statistics

Gilbert SAPORTA

Conservatoire National des Arts et Métiers, Chaire de Statistique Appliquée

292 rue Saint Martin, 75141 Paris Cedex 03, France

E-mail: saporta@cnam.fr

Abstract. There is often a large difference between methods proposed by academic researchers and those used in official statistics (numerical and non numerical multivariate methods, bayesian techniques, non-parametrics, etc.). Conversely, there is also some lack of interest of academic statisticians for official statistics. We will analyse the difference of scope between Official and Academic statisticians and try to propose ways to fill the gap..

Introduction

Official and academic (ie mathematical) statisticians are often evolving in two separate worlds and have few opportunities to meet. Despite the fact that official statisticians have been trained by academic statisticians, in many countries their paths and careers are separate, and as a result there are few communications between them.

This situation is unfortunate since both kinds of statisticians have an interest to cooperate, but often they do not know that! On one hand official statistics may provide stimulating research areas, and on the other hand mathematical statisticians have developed many methods which could be advantageously used by official statisticians.

1. Academic and Official Statistics: some misunderstandings

Academic and official statistics have a long and distinct history: official statistics was born several thousands of years ago with the objective of counting with accuracy the resources (human and material) of states. Its objective of quantitative assessment of the « state of the society and the economy » is still valid [1] but the users of official statistics are not only governments but have been enlarged to all economic agents.

Mathematical statistics began more recently in the 17th century with Pascal, De Moivre, Bernoulli and were associated to probability theory. Up to the end of the 19th century there were very few connections between official and mathematical statistics and the idea of using probability theory (ie samples instead of census) to get informations from a population was considered as heretic until the ISI session of 1925 [2] after years of debates and controversies. The problem then became: how to draw a sample instead of why.

This controversy is a very good example of the progress which can be obtained by using mathematical methodology, and of the misunderstanding between these two worlds, and also of the prominent role of an organisation gathering governmental and academic statisticians.

However, a kind of opposition between « blue-collar statisticians » (government statisticians) and « white-collar statisticians » (methodologists) is still vivid [3].

Mathematical statistics, or theoretical statistics is often far from the real world: excessive or undue formalisation gives too often to this discipline a rather negative feeling to users, not only in economics, but also in industry. They consider many results as useless refinements.

Academic statisticians are not always interested in official statistics for several reasons: they see in official statistics only the activity of collecting, controlling and editing data, a professional occupation which deserves respect, but which needs only know-how and no formalisation. As pointed in [3] « *those who works for Government behave very often as if no error or no uncertainty existed* ». Furthermore for some academics, official statistics may suffer from political pressures.

2. Official Statistics as a Mine of problems

Actually, there are many problems in official statistics which could stimulate the interest of academic statisticians: algorithms for ensuring confidentiality, small area estimation, visualisation techniques etc.

The problem of accuracy of data is of course a central one and due to the complexity of the processes of sampling and estimation, requires mathematical and/or computational skills. See again [3] who writes: « *The first and the most obvious is to develop means of measuring error where no such measurements exist and to ensure that new measurement tools are designed to a standard which in turn is made widely known to the user constituencies. The rigour that white collar statisticians bring to their task can at first appear to be misplaced particularly by government statisticians used as they are to the compromises that their applied discipline has asked them to adopt. But in the end the absence of rigour is what has led criticisms to be levied at the official numbers and in many cases has left official statisticians unable to provide an acceptable reply* ».

3. Existing Tools which are underused in official statistics

Among the enormous production of academic research, let us point out some methods which could be used with a high benefit by official statisticians and are not so widely used, according to my knowledge. I will classify these methods according to the field of applications.

3.1 sampling

The estimation of the size of a population is of crucial interest in industrial statistics, especially in countries where official information and is missing and where the informal sector is important. Biometricians used for many years capture-recapture methods for estimating the size of animal populations. This method may be transposed in official statistics without difficulty, see [4], and could provide confidence intervals for many ratios involving the population size as denominator.

3.2 exploratory analysis

Descriptive statistics are the first outputs published: graphics are generally limited to a few bar and pie charts or histograms, and Lorenz curves. The use of modern graphics such as box-plots, especially for comparing subpopulations, density estimation instead of histograms should be promoted.

Multidimensional data analysis (principal components, correspondence analysis) is now well established: mappings using principal axes are communications tools widely used in

industry and market research. They could also be advantageously used for regional comparisons in official statistics.

Methods for analysing textual data, see [5], are able to deal with open-ended questions by using stylometrics and multivariate techniques. They could be used for opinion surveys, leisure activity surveys etc.

3.3 Modelling

Econometricians are using generally sophisticated models based on strong probabilistic assumptions and with properties which are often valid only asymptotically. There could be advantages to use « soft modelling » based on weaker assumptions such as the PLS (partial least squares) approach initiated by H.Wold [6] which is successfully used in chemometrics for ill-conditioned models. PLS techniques are competitors to maximum likelihood [7].

Classical models (and PLS too) are explicit models in the sense that they provide equations where the main problem is to estimate parameters. Non parametric models have been developed for many years, but they do not seem to have been really applied to econometrics despite the efforts of their promoters (see for instance the web site <http://www.xplore-stat.de/>). One drawback of non explicit models was that until recently they were not conceived for prediction : now it is possible to use them not only for interpolation but also for extrapolation.

Computer intensive techniques such as neural networks for prediction are currently used in many fields of application : why not try them for economic forecasting? Of course they are black-boxes and do not provide deep understanding of phenomenons, but for short-term prediction it seems that the most important thing is to get good predictions more than having the most comprehensible model.

In official statistics a priori information is plentiful, and analysts use it implicitly. Why not do it explicitly and apply bayesian techniques? For a long time bayesian methodology was first a controversial topic, and also very difficult to apply due to the difficulty of programming posterior distributions, but now the status of bayesian analysis is well established and efficient methods for bayesian computing are available [8].

3.4 Data Mining

As D.Hand [9] writes « *Data Mining is the process of secondary analysis of large databases aimed at finding unsuspected relationships which are of interest for the data base owner* ». This new discipline is actually a merge of techniques coming from data-bases and statistics where some tools like cluster analysis, decision trees , bayesian networks are extensively used. The main field of application is for the moment customer and transactions data bases of large commercial companies. As far as NSIs manage large data bases on population, trade, agriculture, companies, they certainly would take a great profit of exploiting their mines of data.

In the same context, there are recent advances on data fusion (merging files from different sources with large scale estimation when questions are missing for some files) which could be exploited in official statistics [10].

Conclusion: some proposals to fill the gap

First of all it is necessary that official and academic statisticians meet in common organisations. At the international level, the role of ISI is very important, but at the national level both groups should participate to the same national statistical society. Common meetings where official statisticians expose their problems to academic ones are very useful .At the European level, Eurostat organizes several scientific meetings each year: they are

usually intended for official statisticians and academic researchers implied in research programs: the diffusion could be enlarged to the whole community of statisticians.

As we have pointed earlier, academics are not aware of the problems of official statistics, especially mathematical statisticians: publishing books about methodological issues of official statistics could be useful.

Developping institutional links between NSIs and Universities should be encouraged: NSI could give doctoral scholarships with a common supervision and help to introduce topics about official statistics in academic curriculums. One could also develop research contracts between NSI and Universities and foster exchanges between people: for instance the CREST (the research center of the French NSI) welcomes University professors for temporary positions (2 or 3 years) but the converse does not exist.

References

- [1] P.Cheung, Developments in official statistics and challenges for statistical education, ICOTS 5, Singapore, 1998
- [2] J.J.Droesbeke and P.Tassi, Histoire de la Statistique, , Que Sais-je, PUF Paris, 1990
- [3] J.Ryten, Blue and White collar statisticians, a gap revisited ,Conference of European Statisticians, Seminar on Official Statistics - Past and Future (Lisbon, Portugal, 25-27 September 1996)
- [4].Giommi et al., On the use of capture mark recapture methodology in estimating the size of an open population of firms. *Contributed paper, vol I, ISI Session , Pékin, 1995.*
- [5] .L.Lebart, A.Salem, L.Berry, Exploring textual data, Kluwer, 1998
- [6] .K.G. Jöreskog & H. Wold,; Contributions to Economic Analysis. Systems under indirect observation : causality, structure, prediction,., North-Holland, 1982.
- [7] M. Tenenhaus, La régression PLS : Théorie et Pratique. Editions Technip, 1998
- [8] C.Robert, The Bayesian Choice, Springer, 1994
- [9] D.Hand, Data Mining: Statistics and More?, The American Statistician, 52, 2 112-118, 1998
- [10] G. Saporta , V.Co, Data Fusion: A New Method Based on Homogeneity Analysis, *8th International Symposium Applied Stochastic Models and Data Analysis*, Contributed Papers, 395-399 .