

Some Thoughts on Official Statistics and its Future

October 19, 2021

Abstract

In this article, we share some reflections on the state of statistical science and its evolution in the production systems of official statistics. We first try to make a synthesis of the evolution of statistical thinking. We then examine the evolution of practices in official statistics, which had to face very early on a diversification of sources: first with the use of censuses, then sample surveys and finally administrative files. At each stage, a profound revision of methods was necessary. We show that since the middle of the 20th century, one of the major challenges of statistics has been to produce estimates from a variety of sources. To do this, a large number of methods have been proposed which are based on very different foundations. The term “big data” encompasses a set of sources and new statistical methods. We first examine the potential of valorization of big data in official statistics. Some applications such as image analysis for agricultural prediction are very old and will be further developed. However, we report our skepticism towards web-scraping methods. Then we examine the use of new deep learning methods. These methods are promising but raise new epistemological questions. With access to more and more sources, the great challenge will remain the valorization and harmonization of these sources.

Keywords: deduction, foundations, induction, Lasso, p -value, registers, sampling, statistical learning

1 Introduction

Official statistics is a somewhat special field of statistics. The methods used there have been developed to deal with particular problems. In official statistics, the main concern is not decision-making but the quality of the proposed estimates. In this article, we try to glimpse the future of statistical methodology. Currently, statistics is confronted with many epistemological questions, the most emblematic of which is the crisis in the p -value. We describe the specificities of official statistics by tracing the history of sources and the history of controversies around methodology. Then, we analyze the impact of new data sources and new statistical methods. An evolution of the methodology will be necessary and must be supported by quality fundamental research.

2 The Statistical Science

Statistics is a science that aims to study a reality through the processing, analysis, modeling and interpretation of data. Traditionally, there are several statistical approaches.

1. *Descriptive or exploratory statistics*, which consists of presenting data in a more condensed way by means of tables and graphs. The aim consists of reducing the complexity of the data by means of dimensionality reduction or cluster analysis. Descriptive statistics is above all linked to an interpretive and exploratory process. John Wilder Tukey was a strong supporter of descriptive statistics. His book *Exploratory Data Analysis* (Tukey, 1977) remains a reference of the approach. Similarly, in France, the school of data analysis initiated by Jean-Paul Benzécri promoted descriptive and exploratory statistics (Benzécri, 1973a,b; Bastin et al., 1980).
2. *Analytical or inferential* statistic which aims to deduce properties on a population from data. This population can be notional and be, for example, a probability distribution or a model. Inferential statistics rely entirely on probability calculations. It includes the theory of statistical hypothesis testing and decision theory. The main founders of this theory are Karl Pearson, William Sealy Gosset, Ronald Fisher and Jerzy Neyman.

3. *Modeling* consists in describing reality by a general model described by one or more equations. A model is necessarily a simpler approximation than reality, as Box and Draper (2007) wrote: “Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful”. Models can be used either to describe relationships between variables or to make predictions.

3 Interpretation of the Statistical Approaches

One usually distinguishes two scientific approaches: The *inductive approach* is based on observations and data and aims to build a general theory. The *deductive approach* (or hypothetico-deductive) is based on a theory. It aims to deduce particular results from a general theory in order to check whether these particular results can be confirmed by observations.

Exploratory statistics can be associated with an inductive approach. In this approach, we start from the data to try to empirically find a global explanation. Inferential statistics and more particularly the theory of statistical hypothesis testing can be seen as a hypothetico-deductive approach. One formulates a hypothesis and then we decide on that hypothesis by confronting it with data. The two approaches are qualified as complementary. The exploratory analysis allows hypotheses to be formulated. Then the inferential statistic can possibly confirm them. A return to exploratory analysis finally makes it possible to formulate new hypotheses and so on. This is referred to as the induction-deduction cycle which would be the engine of the production of knowledge. Often this cycle has been completely perverted, for example, by using the same data to formulate and to test hypotheses which makes the application of test theory completely lapsed. This approach is ironically called HARKing by Norbert L. Kerr (1998) whose acronym comes from *Hypothesizing After the Results are Known*.

However, this vision of the induction/deduction opposition is far from being unanimously shared. The role of induction has been the subject of controversy between Jerzy Neyman and Ronald Fisher. Fisher (1935) promoted the inductive reasoning. Jerzy Neyman (1957) had a more decision-making approach in the hypothesis testing theory but then he was reluctant to extend this to

inductive reasoning. Inferential statistics can be considered as an inductive reasoning because it aims to extrapolate the results of a sample to a population or to a model. It is therefore a generalization process which can be seen as fundamentally inductive (see among others Lehmann, 1993; Capel et al., 1996).

Costantini and Galavotti (1986) assert that estimation methods, such as the maximum likelihood method, correspond to an inductive approach and theory of hypothesis testing typically corresponds to a deductive approach. Other interpretations exist. For example, Andrew Gelman (2011) discusses a philosophy of statistics that sees frequentist statistics as deductive and Bayesian statistics as inductive.

Faced with these different interpretations, we can only subscribe to the assertion of Capel et al. (1996) who wrote: “In any event, a real understanding of the role of inductive reasoning in the humanities that is common to the statistician, the researcher and the practitioner clearly does not exist, and in particular, as we have seen, with regard to the use of hypothesis tests”¹.

Moreover, the distinction between induction and deduction is sometimes considered outdated. Karl Popper (2005) rejects the description of the scientific process as an induction/deduction cycle. He outright rejects the interest of the inductive approach in science: “Yet even supposing this were the case – for after all, ‘the whole of science’ might err – I should still contend that a principle of induction is superfluous, and that it must lead to logical inconsistencies.” Karl Popper maintains that the scientific method consists in formulating scientific propositions which must be able to be falsified by an experiment. Indeed, no theory can be proven by an experiment. The fact that a model is compatible with data never proves that the model is true. Indeed, another model could also be compatible with the same dataset.

The induction/deduction cycle also completely loses its meaning in certain statistical applications. In official statistics, this distinction can hardly be applied. Indeed, the objective often consists in estimating certain characteristics of a population using scattered sources (surveys, administrative files, censuses).

¹Translated from French: “En tout état de cause, une compréhension réelle du rôle du raisonnement inductif dans les sciences humaines qui soit commune au statisticien, au chercheur et au praticien n’existe manifestement pas, et en particulier, comme nous l’avons vu, en ce qui concerne l’usage des tests d’hypothèse.”

The objective is therefore not to establish a scientific theory but simply to overcome the impossibility of obtaining a complete and correct measurement on all the units of the population of interest.

This induction/deduction cycle also loses its meaning with the advent of new methods known as *statistical learning* (support vector machine, neural network, nearest neighbor, sparse methods, random forests) which make it possible to predict without actually modeling. One could see these methods as typically inductive because they do not imply any a priori formalization of reality (Harman and Kulkarni, 2012). The word *statistical learning* is however a bit misleading because these methods do not really lead to a theorization which would be an automatic modeling of data. The so-called induction/deduction cycle is somehow bypassed. For example, if we do not know the income of an individual in a database, we can assign the income of the statistical unit that most closely resembles her/him. This is a forecast by the nearest neighbor method. Can we however say that this approach consists in modeling the income? At most, we can discuss what it means to “resemble”, but we are often forced to use only the variables available in the files to define a distance between the units. However, this approach can be very effective. It is not clear how one could say that this approach would be inductive or deductive. For the statistician, the main question concerning this type of method is to evaluate and estimate its precision. Neural network methods also do not allow us to finely understand the relationship between variables. However, the family of Classification And Regression Tree (CART) allows for post-interpretation. With random forests, it is also possible to have an idea of the importance of the variables in the predictions.

The new methods of *statistical learning* therefore shake up a basic principle of the scientific process which would like all knowledge to be able to be transmitted through discourse. While forecasts by *neural network*, *nearest neighbors* or *random forest* can work very well to make a prediction, they do not allow us to establish a general theory or a principle that would allow us to understand the relationships between variables.

4 Crisis of the p -value

In statistical hypothesis testing, the p -value is the probability that under a hypothesis (conventionally called the null hypothesis), we obtain the same value or an even more extreme value than that obtained with observed data. If the p -value is low, then we reject this hypothesis. The error of type 1 is defined as the probability of rejecting the null hypothesis given that it is true. If we perform the test with an error of type 1 of 5% for example, we reject the hypothesis if the p -value is less than 5%.

The p -value has become a decisive argument in many sciences: humanities, economics, finance or biology. However, in scientific journals of statistics, one may be surprised at the very small number of p -values used in published articles. The p -value is often misinterpreted as the probability that the null hypothesis is true, which is obviously not the case.

In many publications several p -values (sometimes tens) appear without any reflection on the probability of having at least one p -value less than 5% if all the null hypotheses were true. Very often, researchers perform a large number of hypothesis tests and only publish results for p -values below 5%. This approach is similar to HARKing. Another methodological error consists in identifying a model by choosing the variables having p -values less than 5% among a very large number of variables. This procedure inevitably leads to over-specification of the model. However, new solutions have been developed to select variables using for example the Lasso method (Tibshirani, 1996, 2011).

These uses of p -values were denounced in a provocative article: *Why Most Published Research Findings Are False* (Ioannidis, 2005). The subject has become controversial to such an extent that one can say that there is a crisis in statistics (Gelman and Loken, 2014; Fraser and Reid, 2016) and that the *American Statistical Association* felt compelled to issue a statement on the p -value (Wasserstein and Lazar, 2016).

5 What about Official Statistics?

The Swiss Official Statistics Charter is available on the website of the Ethics Council for Official Statistics (Conseil d'éthique de la statistique, 2012). In

this charter, the mission of official statistics is defined: “The mission of official statistics is to meet the needs for statistical information of general interest of society as well as those relating to the conduct of public policies”².

The mission therefore does not consist of interpreting, modeling, establishing knowledge, or deciding. The induction-deduction cycle for the production of knowledge is thus not an appropriate framework in official statistics. In fact, in official statistics, the approach is neither exploratory nor decisional. The statistical methods developed are specific, because it is not a question of carrying out scientific research to establish new knowledge, but of meeting information needs by providing reliable, lasting and high-quality statistics.

There exists however a cycle between civil and political society and official statistics in order to determine the needs. Official statistics must ensure continuity but also renew their statistical production by opening up to new themes such as gender inequalities or environmental concerns. The statistical methods used in official statistics are therefore neither exploratory nor decision-making. Official statisticians have focused on quality. An exemplary document is *The Statistics Canada’s Quality Assurance Framework* whose basic principles are relevance, accuracy, timeliness, accessibility, coherence, interpretability (Statistics Canada, 2017).

6 Data Sources in Official Statistics and their Integration

Official statistics have therefore developed somewhat specific methods to meet their objectives. The advantage is that official statistics are immune to the p -value crisis because its use is relatively limited. The history of official statistics methods can be summed up in a few eras (see among others Hansen and Madow, 1974; Kruskal and Mosteller, 1980; Hansen, 1987; Bellhouse, 1988; Bethlehem, 2009; Tillé, 2020). First, there is the era of censuses which covers the entire 19th century. During this period, only comprehensive data compilation was

²Translated from French: “La statistique publique a pour mission de répondre aux besoins d’informations statistiques d’intérêt général de la société ainsi qu’à ceux relatifs à la conduite des politiques publiques.”

considered scientific. This doctrine is clearly stated by the statistician Adolphe Quételet (1846). The rupture is initiated by the director of the Norwegian Institute of Statistics: Anders Nicolai Kiær (1896, 1899, 1903, 1905) who proposed to use partial data and therefore a sample. After a long controversy, the idea of using samples was finally accepted by the International Statistical Institute (Jensen, 1926). Ken Brewer (2013) interprets this debate as the first controversy in survey sampling.

This opens the era of sampling, which from the start has been the subject of assiduous and fruitful research, (notably on the part of Jerzy Neyman, 1934, 1938, 1952). He shows in particular that it is necessary to over-represent in the sample the categories where the dispersion is larger. This result renders the still too often used concept of representativeness obsolete. With the development of information technologies, then came the era of registers and administrative files in the 1970s. Administrative data are often presented as the new source of data. However, it should not be forgotten that some countries like Finland have had a population register for more than 50 years. The technical means for creating a register have been around for a very long time. Obstacles to the use of registers for statistical purposes are primarily organizational, political and legal. They are often the result of data protection rules or communication difficulties between different administrations. For a long time, the problems have obviously not been technical.

The development of new data sources has not put an end to previous practices. Many National Institutes of Statistics still carry out censuses and almost all of them carry out sample surveys. The integration of administrative files into official statistics is far from obvious. The registers depend on the administrative specificities of the countries, which does not facilitate statistical harmonization. Administrative files often contain many errors, as evidenced by the Serafe crisis in Switzerland.

Serafe has been a private company responsible for collecting audiovisual license fees in Switzerland since 2019. During the first invoices sent, a large number of errors were noted. Serafe had sent the invoices on the basis of the “control of inhabitants” registers updated by the municipalities. Many errors appeared to be related to household composition in multi-unit buildings or to

obsolescence of data. These population control registers are also used for administrative and statistical purposes, but the impact of these errors is relatively small for these applications. In administrative files, it is often observed that only the variables which are necessary for immediate administrative operation are correctly updated.

We cannot therefore expect that all official statistics will be produced from the same type of source. Each type of source contains specific errors, whether they are data produced by statistical offices or not. Problems to be addressed can be sampling errors, errors due to non-response, under-coverage, over-coverage, duplicates, measurement errors, errors due to fraud. Official statistics are a perpetual struggle against these errors.

In each source, one can find a particular reliability. One type of reliability is the correct identification of statistical units, which can often be obtained through an administrative file or register. However, this register may contain poor quality or obsolete variables. A sample survey, even one with non-response, may contain more reliable measures for the variables of interest. The problem then comes down to combining what is most reliable in the different sources. The crucial issue in official statistics is thus the integration of data from different sources in order to best enhance the reliability of each source.

From the beginning, research in official statistics has been concerned with the problem of data integration. Almost a century ago, the controversy between, on the one hand, Corrado Gini and Luigi Galvani and, on the other hand, Jerzy Neyman (Gini and Galvani, 1929) revolves around balanced sampling and random sampling. Gini and Galvani had selected a sample of 29 districts (*circondari*) out of 214 so as to return the same means as those of the census for several known variables. It is therefore a question of using a source (the census) to improve the collection of a sample. Neyman criticizes this way of doing things, because the selection of the sample is not random and therefore it is not possible to make an inference. We obviously now know that it is possible to select a sample that is both balanced and random (Deville and Tillé, 2004). Ken Brewer (2013) interprets this discussion as the second controversy in survey sampling.

Another issue is the adjustment of survey data to census data. On this subject, the article by W. Edwards Deming and Frederick F. Stephan of 1940 is considered as a founding text in official statistics (Deming and Stephan, 1940). The problem dealt with is the adjustment of a table obtained by sampling on marginal totals known by a census. The integration of different sources is therefore present from the start of the use of sampling methods. For the record, the article is however mathematically false since it maintains that the *raking ratio* method is obtained by minimizing the chi-square distance under the constraints given by known marginal totals. This is obviously not true since the Kullback-Leibler divergence must be minimized. The fact remains that this article and all those that follow seek to optimize the integration of survey and census data.

The article of Deville and Särndal (1992) which defines the general method of calibration is the outcome of this research. These authors provide a methodology for calibrating survey data on census data while freeing itself from modeling. The originality is that the processing is carried out in order to obtain a weighting system applicable to any variable, which makes its application extremely practical. Deville and Särndal also provide a method for measuring the accuracy of the estimates obtained. The method has become essential and typical of official statistics. It is now applied by all survey statisticians (see also Särndal, 2007; Devaud and Tillé, 2019a,b). The use of calibration has become widespread, especially since calibration is currently used not only to correct the sampling error but also to correct the nonresponse and measurement error (see among others Dupont, 1994; Fuller et al., 1994; Lundström and Särndal, 1999; Deville, 2000; Särndal and Lundström, 2005; Kott, 2006; Brick, 2013; Valliant et al., 2013; Haziza and Lesage, 2016; Devaud and Tillé, 2019a).

A developing problem is the integration not only of two sources but of a multitude of distinct sources. The calibration methods can be generalized to harmonize several sources. These sources can be two samples or several samples and a census (see among others Guandalini and Tillé, 2017). Yang and Kim (2020) provide an overview of modern methods for integrating data from different sources.

The use of sources may change quite rapidly. Many European countries have abandoned large population censuses. Statistics based on these censuses are now produced by using registers and new surveys complete these registers. National Institutes of Statistics still carry out a lot of sample surveys. Sampling methods are increasingly used to check or improve the quality of censuses or registers. Thus, sampling may become more of a quality control tool than a direct method of producing data.

The evolution of practices does not go without raising new questions. Most emblematic has been the resignation of Martha Farnsworth Riche from the head of the United States Census Bureau. According to the Wikipedia contributors (2020) site: “Although she cited only personal reasons in her resignation, it was seen as a sign that Congressional Republicans were winning in their fight to prevent the Census Bureau from using sampling techniques to correct for persistent undercounting of minorities and other underrepresented groups.” This political intervention in the methodology is obviously extremely worrying and goes against a principle of methodological independence of the National Institutes of Statistics.

7 Modeling in Official Statistics

The debate on the place of the model in survey sampling theory is interpreted by Ken Brewer (2013) as the third controversy in survey sampling. The idea of introducing a model to exploit the auxiliary information resulting from a census in a survey was initially proposed by this same Ken Brewer (1963). However, this idea was mainly developed by Royall (1970, 1971, 1976) (see also Valliant et al., 2000; Chambers and Clark, 2012). The model-based approach involves building a model that is estimated using the sample, and then predicting population data that is not in the sample.

The model approach is opposed to the so-called design-based approach which consists in weighting the sample units by the inverse of their probabilities of being selected. The initial survey weights are the inverse of the inclusion probabilities. These weights are then slightly modified by means of the calibration

technique of Deville and Särndal in order to return exactly the totals of the variables known by a register or by a census.

Official statisticians are inherently reluctant to model. The reason is this: Modeling is somehow expressing an opinion on observations. Modeling can therefore be interpreted as contrary to the principle of impartiality of official statistics. However, model-based and design-based approaches are not necessarily contradictory. Särndal et al. (1992) advocate a model-assisted and nonetheless valid design-based approach. It is therefore possible to construct the concept of double robustness in the sense of an estimate that would be valid either when the model is correct or when the inclusion and response probabilities are correctly identified.

However, there are issues that cannot be solved without some modeling. Research concerning estimation for small areas has been very active during the last decades. These methods consist in producing estimates at very low levels (districts, municipalities) from a sample survey and a register or census. Once again, the problem is to make the best use of the two sources of information. The sample contains the variable of interest. The register contains the list of all units of the population and auxiliary variables. Composite estimators are often used, which are mixtures of direct estimates computed from the sample and estimates obtained from a model linking the auxiliary variables. These composite estimators can be obtained by means of mixed models for which the domains are random effects (see the very complete book of Rao and Molina, 2015).

The National Institutes of Statistics, however, have been relatively cautious and have rarely published estimates obtained by these methods. One of the reasons is that there may be peculiarities well known at the level of local communities that would not be taken into account by a model. For example, Molina and Strzalkowska-Kominiak (2020) proposed labor force estimates in Swiss districts using the Swiss Structural Survey. A Swiss district bordering on Liechtenstein was peculiar because a large number of Swiss border workers work in Liechtenstein. This district has therefore undergone special treatment. However, if it is not possible to identify the singularities before modelling, the estimates may turn out to be very far from reality.

Another area where it is difficult to work without modeling is non-response. There are two ways to do this. Either we can predict the missing values (imputation), or we can estimate the probability that these values are missing in order to weight the responding observations. The models can be very simple. One can impute by a closest neighbor. One can also impute by a simple ratio or by a prediction by regression. One can also predict by taking at random an individual belonging to a small homogeneous stratum. In order to avoid a bad specification of the model, one often seeks to treat the non-response in a doubly robust way. Two models are used: the first one allows us to predict the missing value and the second one allows us to predict the probability of being missing. Double robustness means that the estimation is then approximately unbiased if at least one of the two models is well specified (see on this subject Kang and Schafer, 2007; Han and Wang, 2013; Kim and Haziza, 2014; Boistard et al., 2016; Chen and Haziza, 2017).

8 Big Data and Official Statistics

Tim Harford (2014) gives a strong opinion on big data: “As with so many buzzwords, “big data” is a vague term, often thrown around by people with something to sell. Some emphasise the sheer scale of the data sets that now exist – the Large Hadron Collider’s computers, for example, store 15 petabytes a year of data, equivalent to about 15,000 years’ worth of your favorite music.” Indeed, we can often attend presentations of “experts” announcing the era of “big data” whose slides only contain lists of words starting with “V”, numbers in yottabytes and potatoes linked by arrows without ever giving a real and concrete application. The slides are embellished with pontificating quotes and cartoons. Making lists of words beginning with a letter is the opposite of a serious scientific approach.

The fascination with big data can be irritating in several ways. Talking about “big data” in official statistics is not appropriate. Large administrative files such as population or business registers cannot be qualified as big data. They can contain millions of records but are processable over any desktop computer. These large files have been around in some countries for over fifty years and

these countries have always found the computer means to manage them. They are not big data.

I am convinced that the use of web scraping will remain very limited in producing official statistics. For example, inflation cannot be calculated exclusively by automatically fetching prices from the Internet. ten Bosch et al. (2018) give a list of experiences and lay the foundations of a methodology for web scraping in official statistics. They show that, in most cases, projects rarely end up going into production. The main reasons are that official statistics must be sustainable and that the Internet is not stable over time. Second, because official statistics must have a documented methodology and sources. Because it is not enough to look at a price on the Internet, it must also be possible to verify them. And above all because the National Institutes of Statistics have not waited for the fashion for large data to use directly the price files of the main distributors. Often the data collected by scraping the web can be more simply obtained by another, more stable means.

The number of applications of web scraping actually put into practice for the production of official statistics is therefore very limited. One can mention some applications in price indices for products sold almost exclusively via the Internet. For instance the US Bureau of Labor Statistics use the web to track the price changes of airline fares (US Bureau of Labor Statistics, 2021). However watching prices on the web cannot be really called “web scraping”. Social network analysis can however be interesting for social science analysis, as noted by Connelly et al. (2016). Nevertheless, official statistics needs a stable and reproducible production system that allows temporal comparisons, which is difficult to obtain from the web.

An anonymous referee drew our attention to the billion price project that plans to calculate a price index using web scraping (see Cavallo and Rigobon, 2016). This initiative does not come from the world of official statistics but from two professors of MIT Sloan and Harvard Business School. This approach can be interesting when official statistics are deficient, as was the case in Argentina, but it has not been integrated into official statistics until now. I am also convinced that the methodology is more important than the sample size and that observing

more than a billion prices is not necessarily a guarantee of the quality of the estimates.

The huge data streams of social networks are not valuable either. They belong to the giants of the web. This data has immense commercial value for performing advertising profiling. It would also be extremely dangerous to become dependent on these giants. Above all, these data are not reliable in terms of identifying statistical units and potentially usable variables.

However, there are areas where the data is really massive, such as image analysis for territorial statistics. In this case, the data can be produced (aerial or satellite photos, cartography) by the country which uses them. For instance, the Joint Research Centre of the European Union in Ispra has a long tradition of estimation of crop production using satellite images (see for instance Gallego et al., 1993; Taylor et al., 1997; Gallego, 2004; Carfagna and Gallego, 2005; Kussul et al., 2016). However, this practice was developed before the popularization of the word “big data”. Indeed, image analysis is an area of research that emerged with the early days of computer science. Significant progress has been made, notably through the use of neural networks, which allows us to foresee a major development of applications in official statistics.

9 Statistical Learning in Official Statistics

We often talk about new statistical methods. Most of these methods are actually not so new and were all developed in the 20th century. Several methods (random forests, support vector machine, neural network, nearest neighbors) make it possible to make predictions without having to think about the relationships between the dependent variables and the variable of interest. Can we follow the famous sentence of Deng Xiaoping who said “it doesn’t matter whether a cat is black or white, if it catches mice it is a good cat”? Can we use in official statistics methods that make it possible to predict without understanding? Indeed, the Statistical Charter of the Swiss Ethics Council for Official Statistics specifies that “Statistical information is documented in order

to facilitate its understanding and allow its correct use”³. Is it compatible with statistical learning methods? Is it sufficient to specify the method used?

We believe that *statistical learning* methods can be used, but with certain precautions. The generalized regression estimator (see Särndal et al., 1992) allows a prediction to be incorporated into the estimate while remaining approximately unbiased under the sampling design. This estimator thus avoids slippage due to poor specification of the model. So, official statisticians use fairly simple models that can be integrated in a method that remains valid under the design.

One might think that public statistics are not affected by the p -value crisis because they mainly produce descriptive statistics. However, we have seen that underneath this production, there is a very technical machinery. Surveys are imputed, weighted by non-response models, and then calibrated. Thus, it is necessary to choose calibration variables, design nonresponse models and imputation models. These operations require modeling and this modeling is often based on hypothesis testing, i.e. p -values. The use of parametric statistics is therefore a hidden task. In this field, the use of new deep learning methods can be very promising.

Official statisticians have also long been familiar with light non-parametric methods for dealing with non-response, such as imputation by the nearest neighbor or by an individual selected from a homogeneous stratum. Methods like *support vector machine* or *random forests* ultimately consist in splitting the space of explanatory variables to define one or more neighborhoods of the unit in order to make a forecast.

The generalized regression estimator has enabled Breidt and Opsomer (2000), for example, to construct estimates assisted by a model whose predictions are made by the local polynomial method. By following this same approach, all forecasting methods can be used without introducing disproportionate risks into the estimates. Thus, the recent work of Beaumont and Bocci (2008); Goga and Shehzad (2014); Breidt and Opsomer (2017); McConville et al. (2017); Mayor-Gallego et al. (2019); Chen et al. (2019); Tan (2020); Dagdoug et al. (2020a,b)

³Translated from French: “Les informations statistiques sont documentées afin d’en faciliter leur compréhension et leur utilisation correcte.”

which integrates the methods of *shrinkage* and *statistical learning* in calibration and the treatment of non-response, probably shows the way forward on future research in official statistics.

The so-called *shrinkage* methods like the Lasso make it possible to choose variables in a modeling or a calibration. Survey statisticians often tend to over-calibrate surveys because more and more variables are available in the registers. However, calibration is above all an estimation technique which aims to reduce the variances of the estimates. In general, both calibration and modeling must be subject to a principle of parsimony. It consists of finding the most efficient model as possible while being as simple as possible. The Lasso method makes it possible to reduce the number of calibration variables. It has already been applied to survey data by McConville et al. (2017).

10 Conclusions

The term big data should maybe be dropped. When one asks what “big data” is, one often gets an answer worthy of Borges’ classification of animals in his new “The Analytical Language of John Wilkins” (Borges, 2012). The term “big data” can include both certain types of data or a set of methods. For example, image analysis, social network stream analysis, large administrative files, smart meter data, neural network methods, random forest methods, support-vector machine methods, sparse statistical methods (such as Lasso). The word “big data” thus gathers a heterogeneous set of data and methods in which it is necessary to make a distinction. It is important to keep things in perspective. Image analysis is a very old problem which has been used in public statistics for at least 40 years. The “new” statistical methods were almost all developed in the last century before the big data fad. Obviously all the “new” statistical methods are interesting and are increasingly applied in public statistics sometimes to “classical” problems such as the treatment of non-response where the data sets are not necessarily large. The valorisation of administrative data is a very important challenge. However, I am particularly sceptical about statistical production directly from the web.

The multiplication of sources in official statistics can be misleading because it does not necessarily imply an improvement in quality (see among others Deville, 1997). As Tim Harford (2014) reminds us, abundance of data is not and never will be synonymous with quality. Administrative records often contain a large number of errors because they were not designed to be used for statistical purposes. These new sources must be combined with all other available sources. The integration of sources is an issue that researchers have tackled from the very beginning of official statistics. These methods must be further developed, because more and more sources will be available. It would be useful to have a general theoretical framework for data integration.

New statistical methods must also be integrated or combined with existing methods. We believe that there will certainly not be a clean slate of methodology but an evolution towards practices that are more varied and more focused on the types of data to be processed. More than ever, official statistics must promote and develop methodological research around its issues.

An anonymous arbitrator rightly pointed out to us that an important issue is to maintain and increase the independence of national statistical institutes. A number of autocrats have questioned the legitimacy of official statistical offices, which can bring their output into disrepute. This is obviously a very important issue that we cannot unfortunately develop here. We are convinced that international institutions have a role to play in regulating these practices, but this is well beyond the scope of this article.

References

- Bastin, C., J. P. Benzécri, C. Bougarit, and P. Cazès (1980). *Pratique de l'Analyse des Données*. Paris: Dunod.
- Beaumont, J.-F. and C. Bocci (2008). Another look at ridge calibration. *Metron* 66(1), 5–20.
- Bellhouse, D. R. (1988). A brief history of random sampling methods. In P. R. Krishnaiah and C. R. Rao (Eds.), *Handbook of Statistics Volume 6: Sampling*, New York, Amsterdam, pp. 1–14. Elsevier/North-Holland.
- Benzécri, J.-P. (1973a). *L'analyse des données : tome 1 : La taxinomie*. L'analyse des données. Paris: Bordas.
- Benzécri, J.-P. (1973b). *L'analyse des données : tome 2 : L'analyse des correspondances*. L'analyse des données. Paris: Bordas.
- Bethlehem, J. G. (2009). *The rise of survey sampling*. The Hague, Statistics Netherlands.

- Boistard, H., G. Chauvet, and D. Haziza (2016). Doubly robust inference for the distribution function in the presence of missing survey data. *Scandinavian Journal of Statistics* 43(3), 683–699.
- Borges, J. (2012). *Inquisiciones — Otras inquisiciones*. Penguin Random House Grupo Editorial España.
- Box, G. E. P. and N. R. Draper (2007). *Response Surfaces, Mixtures, and Ridge Analyses*, Volume 649. Hoboken: John Wiley & Sons.
- Breidt, F. J. and J. D. Opsomer (2000). Local polynomial regression estimators in survey sampling. *The Annals of Statistics* 28(4), 1026–1053.
- Breidt, F. J. and J. D. Opsomer (2017). Model-assisted survey estimation with modern prediction techniques. *Statistical Science* 32(2), 190–205.
- Brewer, K. R. W. (1963). A model of systematic sampling with unequal probabilities. *Australian Journal of Statistics* 5, 5–13.
- Brewer, K. R. W. (2013). Three controversies in the history of survey sampling. *Survey Methodology* 39(2), 249–262.
- Brick, M. J. (2013). Unit nonresponse and weighting adjustments: A critical review. *Journal of Official Statistics* 29(3), 329–353. cited By 32.
- Capel, R., D. Monod, and J.-P. Müller (1996). Essai sur le rôle des tests d’hypothèse en sciences humaines. *Actualités Pédagogiques* 1, 1–51.
- Carfagna, E. and F. J. Gallego (2005). Using remote sensing for agricultural statistics. *International statistical review* 73(3), 389–404.
- Cavallo, A. and R. Rigobon (2016). The billion prices project: Using online prices for measurement and research. *Journal of Economic Perspectives* 30(2), 151–78.
- Chambers, R. L. and R. G. Clark (2012). *An Introduction to Model-Based Survey Sampling with Applications*. Oxford: Oxford University Press.
- Chen, J. K. T., R. L. Valliant, and M. R. Elliott (2019). Calibrating non-probability surveys to estimated control totals using LASSO, with an application to political polling. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 68(3), 657–681.
- Chen, S. and D. Haziza (2017). Multiply robust imputation procedures for the treatment of item nonresponse in surveys. *Biometrika* 104(2), 439–453.
- Connelly, R., C. J. Playford, V. Gayle, and C. Dibben (2016). The role of administrative data in the big data revolution in social science research. *Social science research* 59, 1–12.
- Conseil d’éthique de la statistique (2012). Charte de la statistique publique de la Suisse. Office fédéral de la statistique, Neuchâtel.
- Costantini, D. and M. C. Galavotti (1986). Induction and deduction in statistical analysis. *Erkenntnis* 24, 73–94.
- Dagdoug, M., C. Goga, and D. Haziza (2020a). Imputation procedures in surveys using nonparametric and machine learning methods: an empirical comparison. arXiv 2007.06298.
- Dagdoug, M., C. Goga, and D. Haziza (2020b). Model-assisted estimation through random forests in finite population sampling. arXiv 2002.09736.
- Deming, W. E. and F. F. Stephan (1940). On a least square adjustment of sampled frequency table when the expected marginal totals are known. *Annals of Mathematical Statistics* 11,

427–444.

- Devaud, D. and Y. Tillé (2019a). Deville and Särndal’s calibration: revisiting a 25 years old successful optimization problem. *TEST* 4, 1033–1065.
- Devaud, D. and Y. Tillé (2019b). Rejoinder on: Deville and Särndal’s calibration: revisiting a 25 years old successful optimization problem. *TEST* 28, 1087–1091.
- Deville, J.-C. (1997). *Une bonne petite enquête vaut-elle mieux qu’un mauvais recensement ?* Document de travail – Institut national de la statistique et des études économiques. Insee.
- Deville, J.-C. (2000). Generalized calibration and application to weighting for non-response. In *Compstat – Proceedings in Computational Statistics: 14th Symposium Held in Utrecht, The Netherlands* (Softcover ed.), New York, pp. 65–76. Springer.
- Deville, J.-C. and C.-E. Särndal (1992). Calibration estimators in survey sampling. *Journal of the American Statistical Association* 87, 376–382.
- Deville, J.-C. and Y. Tillé (2004). Efficient balanced sampling: The cube method. *Biometrika* 91, 893–912.
- Dupont, F. (1994). Calibration used as a nonresponse adjustment, studies in classification, data analysis, and knowledge organization. In E. Diday (Ed.), *New Approaches in Classification and Data Analysis*, pp. 539–548. Springer-Verlag.
- Fisher, R. A. (1935). The logic of inductive inference. *Journal of the royal statistical society* 98(1), 39–82.
- Fraser, D. A. S. and N. Reid (2016). Crisis in science? or crisis in statistics! mixed messages in statistics with impact on science. *Journal of Statistical Research* 48(1), 1–9.
- Fuller, W. A., M. M. Loughin, and H. D. Baker (1994). Regression weighting in the presence of nonresponse with application to the 1987/1988 nationwide food consumption survey. *Survey Methodology* 20, 75–85.
- Gallego, F. J. (2004). Remote sensing and land cover area estimation. *International Journal of Remote Sensing* 25(15), 3019–3047.
- Gallego, F. J., J. Delincé, and C. Rueda (1993). Crop area estimates through remote sensing: stability of the regression correction. *International Journal of Remote Sensing* 14(18), 3433–3445.
- Gelman, A. (2011). Induction and deduction in Bayesian data analysis. *Rationality, Markets and Morals* 2(67-78), 1999.
- Gelman, A. and E. Loken (2014). The statistical crisis in science: data-dependent analysis—a “garden of forking paths”—explains why many statistically significant comparisons don’t hold up. *American scientist* 102(6), 460–466.
- Gini, C. and L. Galvani (1929). Di una applicazione del metodo rappresentativo al censimento italiano della popolazione (1. dicembre 1921). *Annali di Statistica Series 6, 4*, 1–107.
- Goga, C. and M. A. Shehzad (2014). A note on partially penalized calibration. *Pakistan Journal of Statistics* 30(4), 429–438.
- Guandalini, A. and Y. Tillé (2017). Design-based estimators calibrated on estimated totals from multiple surveys. *International Statistical Review* 85, 250–269.
- Han, P. and L. Wang (2013). Estimation with missing data: beyond double robustness. *Biometrika* 100(2), 417–430.

- Hansen, M. H. (1987). Some history and reminiscences on survey sampling. *Statistical Science* 2, 180–190.
- Hansen, M. H. and W. G. Madow (1974). Some important events in the historical development of sample survey. In D. B. Owen (Ed.), *On the History of Statistics and Probability*, pp. 75–102. New York: Marcel Dekker.
- Harford, T. (2014). Big data: A big mistake? *Significance* 11(5), 14–19.
- Harman, G. and S. Kulkarni (2012). *Reliable Reasoning: Induction and Statistical Learning Theory*. Cambridge, Massachusetts: MIT Press.
- Haziza, D. and É. Lesage (2016). A discussion of weighting procedures for unit nonresponse. *Journal of Official Statistics* 32(1), 129–145.
- Ioannidis, J. P. (2005). Why most published research findings are false. *PLoS medicine* 2(8), e124.
- Jensen, A. (1926). Report on the representative method in statistics. *Bulletin of the International Statistical Institute* 22, 359–380.
- Kang, J. D. Y. and J. L. Schafer (2007). Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data. *Statistical Science* 22(4), 523–539.
- Kerr, N. L. (1998). Harking: Hypothesizing after the results are known. *Personality and Social Psychology Review* 2(3), 196–217.
- Kiær, A. N. (1896). Observations et expériences concernant des dénombrements représentatifs. *Bulletin de l'Institut International de Statistique* 9, 176–183.
- Kiær, A. N. (1899). Sur les méthodes représentatives ou typologiques appliquées à la statistique. *Bulletin de l'Institut International de Statistique* 11, 180–185.
- Kiær, A. N. (1903). Sur les méthodes représentatives ou typologiques. *Bulletin de l'Institut International de Statistique* 13, 66–78.
- Kiær, A. N. (1905). Discours sans intitulé sur la méthode représentative. *Bulletin de l'Institut International de Statistique* 14, 119–134.
- Kim, J. K. and D. Haziza (2014). Doubly robust inference with missing data in survey sampling. *Statistica Sinica* 24(1), 375–394.
- Kott, P. S. (2006). Using calibration weighting to adjust for nonresponse and coverage errors. *Survey Methodology* 32, 133–142.
- Kruskal, W. and F. Mosteller (1980). Representative sampling, IV: The history of the concept in statistics, 1895-1939. *International Statistical Review* 48, 169–195.
- Kussul, N., G. Lemoine, F. J. Gallego, S. V. Skakun, M. Lavreniuk, and A. Y. Shelestov (2016). Parcel-based crop classification in ukraine using landsat-8 data and sentinel-1a data. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 9(6), 2500–2508.
- Lehmann, E. L. (1993). The Fisher, Neyman-Pearson theories of testing hypotheses: one theory or two? *Journal of the American Statistical Association* 88(424), 1242–1249.
- Lundström, S. and C.-E. Särndal (1999). Calibration as a standard method for treatment of nonresponse. *Journal of Official Statistics* 15, 305–327.

- Mayor-Gallego, J., J. Moreno-Rebollo, and M. Jiménez-Gamero (2019). Estimation of the finite population distribution function using a global penalized calibration method. *AStA Advances in Statistical Analysis* 103(1), 1–35.
- McConville, K. S., F. J. Breidt, T. C. M. Lee, and G. G. Moisen (2017). Model-assisted survey regression estimation with the lasso. *Journal of Survey Statistics and Methodology* 5(2), 131–158.
- Molina, I. and E. Strzalkowska-Kominiak (2020). Estimation of proportions in small areas: application to the labour force using the swiss census structural survey. *Journal of the Royal Statistical Society A* 183(1), 281–310.
- Neyman, J. (1934). On the two different aspects of the representative method: The method of stratified sampling and the method of purposive selection. *Journal of the Royal Statistical Society* 97, 558–606.
- Neyman, J. (1938). Contribution to the theory of sampling human population. *Journal of the American Statistical Association* 33, 101–116.
- Neyman, J. (1952). *Lectures and Conferences on Mathematical Statistics and Probability*. Washington: Graduate School, U. S. Department of Agriculture.
- Neyman, J. (1957). “Inductive Behavior” as a basic concept of philosophy of science. *Revue de l’Institut International de Statistique*, 7–22.
- Popper, K. (2005). *The logic of scientific discovery*. London: Routledge.
- Quételet, A. (1846). *Lettres à S. A. R. le Duc régnant de Saxe-Cobourg et Gotha sur la théorie des probabilités appliquées aux sciences morales et politiques*. Bruxelles: M. Hayez.
- Rao, J. N. K. and I. Molina (2015). *Small Area Estimation*. New York: Wiley.
- Royall, R. M. (1970). On finite population sampling theory under certain linear regression models. *Biometrika* 57, 377–387.
- Royall, R. M. (1971). Linear regression models in finite population sampling theory. In V. P. Godambe and D. A. Sprott (Eds.), *Foundations of Statistical Inference*, Toronto, Montréal, pp. 259–279. Holt, Rinehart et Winston.
- Royall, R. M. (1976). The linear least squares prediction approach to two-stage sampling. *Journal of the American Statistical Association* 71, 657–664.
- Särndal, C.-E. (2007). The calibration approach un survey theory and practice. *Survey Methodology* 33, 99–119.
- Särndal, C.-E. and S. Lundström (2005). *Estimation in Surveys with Nonresponse*. New York: Wiley.
- Särndal, C.-E., B. Swensson, and J. H. Wretman (1992). *Model Assisted Survey Sampling*. New York: Springer.
- Statistics Canada (2017). Statistics Canada’s Quality Assurance Framework. Documentation of the internet site of Statistics Canada, Statistics Canada, Ottawa.
- Tan, Z. (2020). Regularized calibrated estimation of propensity scores with model misspecification and high-dimensional data. *Biometrika* 107(1), 137–158.
- Taylor, J., C. Sannier, J. Delincé, and F. J. Gallego (1997). Regional crop inventories in europe assisted by remote sensing. *Synthesis Report, Office for Publications of the European Commission*.

- ten Bosch, O., D. Windmeijer, A. van Delden, and G. van den Heuvel (2018). Web scraping meets survey design: Combining forces. In *Big Data Meets Survey Science Conference, Barcelona, Spain*.
- Tibshirani, R. J. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* 58(1), 267–288.
- Tibshirani, R. J. (2011). Regression shrinkage and selection via the lasso: a retrospective. *Journal of the Royal Statistical Society B73*(3), 273–282.
- Tillé, Y. (2020). *Sampling and Estimation From Finite Populations*. Hoboken: Wiley.
- Tukey, J. W. (1977). *Exploratory Data Analysis*, Volume 2. Reading, MA: Addison-Wesley.
- US Bureau of Labor Statistics (2021). Consumer price index <https://www.bls.gov/cpi/factsheets/airline-fares.htm>. web site visited on 2020-01-08.
- Valliant, R., J. A. Dever, and F. Kreuter (2013). *Practical Tools for Designing and Weighting Survey Samples*. New York: Springer.
- Valliant, R., A. H. Dorfman, and R. M. Royall (2000). *Finite Population Sampling and Inference: A Prediction Approach*. New York: Wiley.
- Wasserstein, R. L. and N. A. Lazar (2016). The ASA statement on p -values: Context, process, and purpose. *The American Statistician* 70(2), 129–133.
- Wikipedia contributors (2020). Martha Farnsworth Riche – Wikipedia, the free encyclopedia. [Online; accessed 24-August-2020].
- Yang, S. and J. K. Kim (2020). Statistical data integration in survey sampling: A review. *Japanese Journal of Statistics and Data Science* 3, 625–650.