



STATISTICS IN TRANSITION

new series

An International Journal of the Polish Statistical Association and Statistics Poland

CELEBRATING 100TH ISSUE AND THE 30TH ANNIVERSARY

Okrasa W., Rozkrut D., Preface	I
Invited Paper	
Kalton G., Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day	1
Pfeffermann D., Comments	23
Lehtonen R., Comments	27
Gershunskaya J., Lahiri P., Discussion	31
Münnich R., C., Discussion	39
Kalton G., Rejoinder	43

Preface

This issue is the hundredth in 30 years of publishing *Statistics in Transition*. The first issue appeared in July 1993, and for the next fifteen years it was a semi-annual publication. In 2007 the title of the journal was slightly changed to *Statistics in Transition new series* and it became a quarterly publication. To celebrate the historical significance of these milestones, we dedicate the first part of this issue to them, opening it with a specially prepared Invitation Paper, along with four discussion pieces of the issues raised in that paper.

With a sense of deep gratitude and the highest appreciation we would like to thank, both personally and on behalf of all the editorial bodies, Professor Graham Kalton for preparing his Invited Paper entitled *Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day*. Dr. Kalton is a long-time friend of our journal, and he serves as a member of our Editorial Board. The issues discussed in Dr. Kalton's paper are particularly appropriate at this time as major changes are taking place in survey research methods and in sources of official statistics. The paper and the discussion pieces should therefore be of interest to members of the international statistician community and to members of national statistical offices.

Despite the relatively short time for reactions, we are grateful to five eminent experts, four of whom are associated with *SiTns*, for preparing four discussion pieces related to the paper. The authors of the four discussions are Professor Danny Pfeffermann, Dr. Julie Gershunskaya and Professor Partha Lahiri, Professor Risto Lehtonen, and Professor Ralf Münnich. Each of the discussions provides insightful observations supplementing some of the issues picked out from those discussed by Graham Kalton. They share concerns about the current challenges to probability sampling and design-based inference primarily caused by the serious declines in response rates, especially in high-income countries. They point to the possibilities of using alternative modalities (administrative data, big data, internet data, scientific data, etc.) for data collection that can supplement or replace probability samples. They describe the considerable body of research that is in progress to enable these alternative data sources to produce valid population estimates from the nonprobability samples associated with the modalities, and to the data integration methods that are being developed to combine the data obtained from different sources.

An *addendum* to this section contains a paper by Professor Jacek Wesolowski entitled *Rotation schemes and Chebyshev polynomials*, as being inspired in a way by the Invited Paper, and as an indication of other types of effects that it may have as well.

It is noteworthy that as our journal celebrates its 30th anniversary, the journal's name *Statistics in Transition* well reflects the radical changes in the methodology of survey statistics and official statistics that are currently underway, as indicated in the Invited Paper and the discussions in this section.

Włodzimierz Okrasa
Editor, *Statistics in Transition new series*

Dominik Rozkrut
President, Statistics Poland

Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day

Graham Kalton¹

Abstract

At the beginning of the 20th century, there was an active debate about random selection of units versus purposive selection of groups of units for survey samples. Neyman's (1934) paper tilted the balance strongly towards varieties of probability sampling combined with design-based inference, and most national statistical offices have adopted this method for their major surveys. However, nonprobability sampling has remained in widespread use in many areas of application, and over time there have been challenges to the Neyman paradigm. In recent years, the balance has tilted towards greater use of nonprobability sampling for several reasons, including: the growing imperfections and costs in applying probability sample designs; the emergence of the internet and other sources for obtaining survey data from very large samples at low cost and at high speed; and the current ability to apply advanced methods for calibrating nonprobability samples to conform to external population controls. This paper presents an overview of the history of the use of probability and nonprobability sampling from the birth of survey sampling at the time of A. N. Kær (1895) to the present day.

Key words: Anders Kær, Jerzy Neyman, representative sampling, quota sampling, hard-to-survey populations, model-dependent inference, internet surveys, big data, administrative records.

1. Introduction

This paper presents a selection of the major developments that have taken place over the years since social surveys were first introduced in the late 19th century. I restrict my coverage to surveys of households and persons and my focus is on the sampling methods used to conduct such surveys. Major changes have also taken place in modes of data collection, in questionnaire design, and in other aspects of survey research over the years, but these topics are outside the scope of this paper. My paper on the more general theme of survey research over the past 60 years overlaps with this paper and gives greater coverage on some topics (Kalton, 2019).

The changes that have occurred in methods of survey sampling have arisen for many reasons, including developments in sampling theory, the continuing growth in computer power (that was non-existent for the first fifty years of survey research), new sampling frames, and the problems created by a broader and more challenging range of applications of social surveys that has occurred as the potential for survey research has been more fully recognized. While acknowledging these changes, it is noteworthy that many aspects of the sampling methods that have been superseded over time have remained relevant. Indeed, much of the current discussion of the use of nonprobability sampling and big data sources has roots in the early days of survey research.

Without attempting to date the origins of survey research, early applications of survey research for studying the social conditions of populations took off in the late 1800's. English examples include Charles Booth's large-scale survey of the social conditions of the population of London that was started in 1886, Seebohm Rowntree's survey of working-class poverty in York that was conducted a decade later, and Bowley's survey of working-class conditions in Reading in 1912, which he followed up with surveys in four other English towns (three of which were conducted by Burnett-Hurst under Bowley's direction). See Caradog Jones (1949) for the early surveys in England, Converse (2017) for an account of the history of survey research in the United States from its beginnings at the turn of the century through until 1960, and Stephan (1948) for a history of the use of sampling procedures dating back from earlier times through until the 1940's, primarily in the United States.

The London and York surveys were complete censuses of the surveys' target populations. As complete censuses, they were deemed statistically acceptable at the time; they were known as 'monographs' of their local communities. For the London survey, the target population was households with school-aged children, while for the York survey it was households that did not have servants (conducted only in streets that were likely to contain households without servants). Bowley had long argued for the use of sampling for such surveys, and he played a major role in its adoption (Aldrich, 2008). He used sampling for the first time in the five towns surveys, where systematic sampling was employed (Bowley, 1913), and he introduced the idea of measuring sampling errors for survey estimates.

As Kish (1995) notes, the emergence of the field of survey sampling can be dated from work led by the Norwegian statistician Anders Kær, the first Director of Statistics Norway. Kær developed a sampling method that he termed "representative sampling". Kær's method of purposive sampling is worth reviewing both for the procedures he devised to make a sample nationally 'representative' and for the reactions to the method from statisticians attending meetings of the International Statistical Institute (ISI) at the time. The next section provides a brief overview of these issues.

¹ Joint Program in Survey Methodology, University of Maryland, College Park, MD, USA.
E-mail: gkalton@gmail.com. ORCID: <https://orcid.org/0000-0002-9685-2616>.

2. Kiær's Representative Method of Statistical Surveys

Kiær's sampling methodology is described in detail in his monograph *The Representative Method of Statistical Surveys*, first published in Norwegian in 1897 and republished in 1976 with an English translation (Kiær, 1976). The monograph provides a good deal of detail on the sample designs Kiær developed for two large-scale surveys—one on personal income and property (PIP) and the other on living conditions (LC)—as well as reporting the objections to his methods that he received when he presented them at ISI meetings. As distinct from the surveys of English towns cited above, Kiær aimed to produce survey estimates for the whole of Norway. For this purpose, he developed two-stage area sample designs for his surveys: at the first stage, he selected a “representative” sample of administrative districts (rural districts or counties, towns, and cities); at the second stage, he drew samples of people for each survey. The choice of the sampled first-stage units was carefully fashioned to give geographical spread and to achieve a good representation of the Norwegian population in terms of characteristics collected in the 1891 Population Census (e.g., age, marital status, occupation, urbanicity).

The sample for the PIP survey was defined as men aged 17, 22, 27, etc. who had names starting with certain letters, selected from 1891 census records that were being processed at the time, with a total sample size of around 11,400 men. The sample size for the LC survey was around 80,000 adults. The sample size to be obtained in each selected rural county was specified based on calculations from census data; within selected counties, the enumerators were instructed to follow certain routes and to select different types of houses, but otherwise they were left to make the selections. In the smaller towns, every 9th, 5th, or 3rd house was selected. An extra sampling stage was introduced in the largest towns. For example, the sample of houses in Oslo was selected within a sample of streets. Moreover, a higher proportion of the streets with larger populations was included in the sample, but this feature was counterbalanced by the selection of houses at a lower rate in the large streets.

The driving objective with Kiær's approach was to produce a representative sample that would constitute a microcosm of the Norwegian population. He invented some intricate methods to attempt to achieve this objective. His purposive selection of first stage administrative units sometimes incorporated ideas of probability proportional to size sampling and subsampling at different rates in compensation, thereby avoiding an excessive sample concentration in a few large districts. Similarly, his street sample in Oslo has the same feature. He also employed a simple 2:1 weighting adjustment to compensate for the smaller proportion of members of the rural population in the PIP survey. (Before the advent of computers, anything other than simple integer weighting adjustments would have been extremely difficult to routinely apply.)

Despite his thoughtful approach, Kiær encountered a great deal of criticism of his methods when he presented them to the ISI in 1895. The dominant criticism, however, was not of the representative method, *per se*, but rather of a sample-based enquiry rather than a complete enumeration. In the words of one strong critic, von Mayr: “We remain firm and say: no calculations when observations can be made”. Kiær also made presentations on the representative method at the 1897, 1901, and 1903 ISI sessions, at which they were subjected to similar criticisms, together with another one. At the 1903 session, von Bortkiewicz reported the results of a significance test he had conducted that found that Kiær's representative samples were not truly representative. See Kruskal and Mosteller (1980) for a detailed account of the ISI sessions.

At the same time, Kiær expertise was under attack at home for the LC survey, which was conducted on behalf of a parliamentary labor commission to inform a very contentious social security act that would provide highly expensive disability insurance. A three-person “critique committee” was established to review the commission's major recommendation and its statistical basis. One committee member, the actuary Jens Hjorth, was extremely critical of Kiær's statistics, including the survey design, the representative sample design, and the analysis. The attacks on the statistics that Kiær's produced for the commission were forceful, extensive, and widely debated. In the end, based on the results of some new surveys, Kiær admitted that he had initially seriously underestimated the extent of disability. After that time, representative sampling for large-scale surveys disappeared in Norway. Lie (2002) provides an informative account of the rise and fall of Kiær's representative sampling method.

The ISI discussion of survey sampling fell into abeyance until 1924 when the ISI appointed a commission for studying the application of the representative method in statistics. By that time, the idea of a “partial investigation” was widely accepted. In its 1926 report (Jensen, 1926), the Commission concluded that a sample was acceptable if it was sufficiently representative of the whole. To satisfy this condition the sample could be produced either by random selection with equal probability or by purposive selection of groups with a representative overall sample. The report also recommended that the survey results should, wherever possible, be accompanied by an indication of the errors to which they are liable.

3. Neyman's Seminal Paper

In 1934, Neyman presented his classic paper comparing the methods of random and purposive selection to the Royal Statistical Society (Neyman, 1934). Covering more than the comparison, the paper contained a detailed discussion of a methodology for making inferences from random—or, more generally, probability—samples of finite populations, including providing a definition of a confidence interval in this context.

He also critically examined the assumptions made when using data from a purposive sample to produce an accurate estimate of a population parameter.

He discussed the sample design of purposive selection of groups used by Gini and Galvani in selecting a sample of records from the already-processed Italian General Census of 1921 that was to be used as the basis for later analysis. For their sample, Gini and Galvani (1929) selected a sample of twenty-nine of the 214 districts in Italy, balanced on seven covariables (note that departs from Kiær's stipulation that a large wide-spread sample of areas is needed). While the sample worked well for the averages of the control variables, it often failed to adequately represent the national population for other characteristics, and for the distributions of the control variables. These findings led them to raise questions about representative sampling.

Neyman's paper was a watershed for survey sampling, leading to widespread adoption of probability sampling, particularly by national statistical offices. It also led to the development of an extensive range of sampling methods and the associated theory applicable to a variety of practical survey problems, as described in the several texts on survey sampling that appeared in the 1950's. The many contributions of statisticians at the U.S. Census Bureau led by Morris Hansen are particularly noteworthy; see, for example, the two-volume text by Hansen, Hurwitz, and Madow (1953). Statisticians active in research on sample designs for agricultural surveys, such as Yates in England and Mahalanobis in India, also made important contributions to the advancement of the subject. The sampling text by Yates (1949) was among the first books on survey sampling methods. In 1950, Mahalanobis went on to establish and lead the famous socio-economic National Sample Survey (NSS) of India. An interesting feature of the NSS sample design was that the sample was composed of four replicate samples. The survey results were presented for each replicate separately as well as for the full sample, with the aim of communicating to readers an indication of the amount of sampling error in the survey estimates (see, for example, Mahalanobis, 1946). This was thus a forerunner of variance estimation using replication methods.

Note that perfect application of Neyman's design-based inference for probability sampling depends on:

- The availability of a sampling frame that provides complete coverage of the finite target population;
- A sample design that assigns known and non-zero selection probabilities to every element in the target population;
- Survey responses from every sampled unit; and
- The use of survey weights in the analysis to compensate for unequal selection probabilities.

Under these conditions (and assuming no response errors), survey estimates can be computed that are design-consistent estimates of the population parameters without the need to make any assumptions about the characteristics of the survey population. Model assumptions made about the population structure may be used to make the sample design more efficient or in the computation of the survey estimates, but the consistency of the survey estimates remains irrespective of the validity of the model. What the model assumptions do affect is the precision of the survey estimates. For example, in a stratified sample, if the sampling fraction in a stratum is set at a higher rate because the elements in a stratum are incorrectly modeled to be more variable, the (weighted) sample mean will still be unbiased, but it will be less precise than if the stratum element variance has been correctly modeled. Similarly, if a set of auxiliary variables \mathbf{X} is available for all population elements, and a function of the \mathbf{X} 's, $f(\mathbf{X})$, is used as a working model to predict the survey variable y , then the finite population total may be estimated by

$$\hat{Y}_d = \sum_U \hat{f}(\mathbf{X}_i) + \sum_S w_i e_i, \quad (1)$$

where \sum_U and \sum_S denote summations over the population and sample respectively, $\hat{f}(\mathbf{X}_i)$ denotes the model estimate of y_i using the sample estimates of the unknown model parameters, $e_i = y_i - \hat{f}(\mathbf{X}_i)$, and the weight w_i is the inverse of element i 's selection probability. By including the weighted estimate of the population total of the e_i 's in this estimate, \hat{Y}_d is a consistent estimator of the population total Y irrespective of the suitability of the working model; the choice of working model affects only the precision of the estimate \hat{Y}_d . This estimator is model-assisted, using the terminology coined by Särndal, Swensson, and Wretman (1992), but it is not model-dependent. For simple random sampling, Cochran (1953) gave an early example of a model-assisted estimator with the ratio estimator $\hat{Y} = (\bar{y}/\bar{x})X$, where X denotes the population total for the auxiliary variable x . An additional, important, feature of design-based inference is that estimates of the variances of sample estimates can be computed from the sample itself.

While the lack of dependence of design-based inference on model assumptions is the major attraction of probability sampling, it needs to be acknowledged that probability sampling is rarely perfectly executed in practice. There are two main sources of imperfection: noncoverage and nonresponse. Noncoverage, which arises because the sampling frame fails to include some elements of the target population, is widespread and its magnitude is often underrated. Area sampling is widely used in social surveys, selecting a probability sample of geographical areas, listing the households or dwelling units in the sampled areas, selecting a probability sample of households, and selecting either all or a probability sample of persons in those households. Even when the sample

of areas provides complete geographical coverage, noncoverage arises often from incomplete listing of households or dwelling units within sampled areas, and from incomplete listing of persons within sampled households. Nonresponse occurs when a sampled element fails to provide acceptable responses to some or all the survey questions. In the early years of probability sampling, response rates were high, and these two sources of imperfection were treated as minor blemishes that received little attention. They were either ignored or treated by simple weighting adjustments (simple, in part because more complex adjustments were computationally infeasible at the time).

Probability sampling has two main drawbacks to be balanced against the theoretical attractions of design-based inference: cost and timeliness. The extra costs of probability sampling include the costs of tracking down sampled individuals, including repeat calls when the individual is not initially available. When area sampling is used, the sampling costs also include the costs of listing units within sampled areas. For similar reasons, collecting survey data from a probability sample takes longer, making the production of the survey estimates less timely. Timeliness is important for all surveys, but particularly for surveys where the results are highly time-dependent, such as political polls, surveys of outbreaks of certain infections, and surveys of areas that have experienced a recent disaster.

A variety of less rigorous sampling methods are used in an attempt to apply a probability sampling approach to address these drawbacks. However, since all these methods require modeling assumptions, none of them can be classified as probability sampling. For convenience, they are called ‘pseudo-probability’ methods in what follows. In the early days of design-based inference, the quasi-probability sampling method known as quota sampling was widely used in market research and in other applications. That method is described in Section 4. Three other quasi-probability sampling methods are described briefly in Section 5.

4. Quota Sampling

To set the scene for the need for imposing quota controls on a sample of the general population, consider the infamous Literary Digest Poll of 1936. To forecast the outcome of the 1936 U.S. Presidential Election, the Literacy Digest mailed a questionnaire to a sample of ten million individuals selected from telephone directories, lists of automobile owners, and registered voters. The results obtained from the two million respondents indicated a clear-cut victory for Alf Landon with 57 percent of the vote, whereas in fact Franklin Roosevelt won with 61 percent of the vote. The upper-class bias of the sample, and of the respondents within the sample, is a major part of the explanation of the discrepancy between these percentages. No weighting adjustments

were employed to attempt to address the bias at the time. (Lohr and Brick, 2017, reweighted the sample using respondents’ reports of their voting in the 1932 election, and these adjustments led to a correct prediction of the outcome, but the estimate of the vote for Roosevelt still fell far short of the actual vote.) This study serves to demonstrate that a large sample size does not necessarily yield good estimates. See Converse (2017) for more details.

Market researchers and pollsters developed the methods of quota sampling separately from the developments in probability sampling, with the aim of addressing the biases from uncontrolled sampling. There are various forms of quota sampling, with the essence of all of them being to control the types of persons to be interviewed. Interviewers are instructed to make their samples of respondents conform to specified quota controls by such characteristics as sex, age group, and employment status. The controls could be independent (e.g., so many men and so many women, so many persons over 35 and so many persons 35 years of age or less) or the numbers to be interviewed could be interrelated (e.g., so many men over 35, so many women over 35). Sudman (1966) describes a method of quota sampling for national face-to-face interview surveys that he termed “probability sampling with quotas”. He employed the four quota control groups of men under 35, men 35 and older, employed women and unemployed women, with the control groups chosen to give appropriate representation to young men and employed women. See also Stephenson (1979). The interviewing field force would generally be distributed across the country in a balanced way, either in areas selected to be representative, along the lines employed by Kiær, or in areas selected by a probability sample design. Sometimes additional controls are imposed, for example specifying the routes the interviewers were to follow, with no more than one person sampled in any household. Quota controls can also be applied in telephone surveys, mall intercept surveys, internet surveys (see Section 6), and other types of survey.

Quota sampling has two main advantages over probability sampling: cost and timeliness. Quota sampling is less costly because interviewers do not need to chase up elusive sampled units and because it avoids the costs of sampling specific households or persons (often including the associated listing costs). For the same reasons, a quota sample can be speedily fielded, and the data collected more rapidly than with a probability sample.

Quota sampling is a form of nonprobability sampling that assumes that the respondents in a quota group are an equal probability sample of the population in that group. Note that this assumption also assumes that nonrespondents in the group are missing at random; nonresponse occurs with quota sampling, in essence with respondents substituted for the nonrespondents. Studies that have been conducted to evaluate quota sampling have found that the results are often similar to those produced

by probability sampling, but this is not always the case (see Moser and Stuart, 1953, also Moser and Kalton, 1971; Stephan and McCarthy, 1958). For further references on quota sampling, see Kruskal and Mosteller (1980).

Random Route Sampling. Random route, or random walk, sampling is another quasi-probability sampling method that avoids the cost of, and associated time involved with, the listing operation. There are various versions of this method, but each starts with a random selection of a starting household and the interviewers then follow specified rules for walking patterns to follow and selection methods to use for serially identifying the subsequent households. The method has often been used in Europe and it is used in the Expanded Programme of Immunization (EPI) sampling method described in Section 5. Bauer (2014, 2016) discusses the selection errors that can occur with random route sampling and demonstrates that the method does not produce an equal probability sample, as its users generally assume.

5. Pseudo-Probability Sample Designs for “Hard-to-Survey Populations”

Recent years have seen a major increase in the use of social survey methods to study the characteristics of “hard-to-survey populations” (Tourangeau, Edwards, Johnson, Wolter, Bates, 2014). Such populations are of various types, but all comprise only a small proportion of the general population and a population for which there is no separate sampling frame. This section presents three examples of sample designs for such populations. The first example is an inexpensive method that has been very widely used for vaccination surveys of the extremely rare population of 1-year-old children. The other two examples describe methods for sampling rare populations where membership of that population is a sensitive characteristic.

a. *The EPI sampling method.*

For almost 50 years, the World Health Organization’s Expanded Programme on Immunization (EPI) has used simple, inexpensive, sample designs in developing countries for measuring childhood immunization at the district level. Many thousands of EPI surveys have been conducted over this period, and the sample design has evolved over time. The sample design is a two-stage sample of clusters of communities (e.g., villages, towns, health service districts) that are sampled with outdated measures of estimated population sizes, with samples of eligible children selected within selected communities. The standard overall sample size is small, with the selection of 30 clusters and 7 children in each cluster. The design is often known as 30×7 design. Except in smaller communities, no household listings are made. Instead, the interviewer goes to the center of the village, chooses a random direction by spinning a bottle on the ground, and counts the number of households in that direction to the edge of the

community. The interviewer then chooses a random number (for instance, from the numbers on a banknote) to identify the first sampled household. The second sampled household is then the one closest to the first, and so on, sequentially until survey data are collected on seven eligible children. Levy and Lemeshow (2008, pp. 427–428) describe the EPI sampling methods and Bennett (1993) describes some of the modifications to the original method.

The US Centers for Disease Control and Prevention (CDC) recommends a probability 30×7 sample design for its rapid needs assessment tool, the Community Assessment for Public Health Emergency Response (CASPER) program. In this case, the clusters are generally census blocks with counts of households obtained from the U.S. Census Bureau or by using a GIS program for use in the PPES selection of thirty clusters. The fieldworker counts or estimates the number of households in a sampled cluster, divides that number by seven to give the sampling interval for systematic sampling, proceeds to select the sample from a random starting point, selecting subsequent households using a serpentine walking procedure. A crude weighting adjustment is proposed for use in the data analysis. Details are provided by CDC (2019).

b. *Venue-Based Sampling*

Venue-based sampling (also known as location sampling, time-space sampling, center sampling, and intercept sampling) is used for sampling members of a rare population at places that they frequent. It is applicable for rare populations that visit certain locations. It can be used to survey nomadic populations and for sampling hidden rare populations where the membership of that population is a sensitive matter. The method requires the construction of a frame of locations and a decision on the overall time period for the survey, selecting a sample of location/time periods for data collection, and selecting all or a sample of members of the survey population visiting each sampled location in the sampled data collection time period (Kalton, 1991). Two issues of concern arise when sampling hidden populations. One relates to the population coverage provided by the frame of locations and the overall time period: What proportion of the population will fail to visit any of the locations in that time period? Another issue relates to the multiplicity problem: How to account for the variability in the numbers of visits made to any of the locations by different sample members during the overall time period? These numbers are needed for use in weighting to compensate for unequal selection probabilities, but they are unknown. At best, they can be estimated by asking respondents questions about their general frequencies of visiting the locations. See MacKellar, Gallagher, Findlayson, Lansky, and Sullivan (2007) for a description of the sampling methods used for surveying men who have sex with men (MSM) in a number of metropolitan areas in the United States.

c. Respondent Driven Sampling

Respondent driven sampling (RDS) is a form of link-trace sampling that selects the sample based on the social networks that exist for some populations. RDS has become a popular method for sampling rare hidden populations that have this feature, such as injection drug users and sex workers. The method starts by identifying a small set of members of the population of interest, who serve as *seeds* for the subsequent sample. The seeds respond to the survey, including responding to a question asking how many members of the survey population they know. They are then asked to recruit a set number of members of that population for the survey, the *alters*. The alters then go through the same process, recruiting further sample members. Under idealized circumstances, Heckathorn (1997) has shown that RDS produces a probability sample. However, the many conditions required for this to apply will not hold in practice (Gile and Hancock, 2010).

6. Internet Surveys

Recruiting the sample via the internet is a relatively recent approach for conducting social research. This approach has become extremely popular and has led to several alternative methods. See, for example, Baker, Blumberg, Brick *et al.* (2010) for a review of these methods. Surveys based on internet sampling have the great attractions of obtaining responses from large samples at low cost and high speed. However, their nonprobability sampling methods raise concerns about potential biases in the survey estimates. Those without, or with limited, access to the internet are excluded from these surveys and the survey respondents are clearly not a representative sample of the general population.

One form of internet sampling, known as river sampling, attaches invitations to participate in a survey on a number of internet sites, usually with offers of some form of compensation. The biases in the sample selection process make the representativeness of the sample highly questionable. Questions also need to be raised about the honesty and thoughtfulness of the responses.

Another form of internet sampling employs an opt-in internet panel. (An opt-in internet panel is distinct from an internet panel that selects a household panel by probability sampling and then conducts many data collections from the panel over time, albeit typically with low response rates). Extremely large numbers of people are recruited for opt-in internet panels to be available to be approached to respond to surveys over time, sometimes as one of a range of services they may be asked to provide, in exchange for a payment for their services. The panel members can then be selected for invitation to respond to a given survey based on their responses to the screening instrument used in their recruitment.

In some ways, these large-scale nonprobability internet surveys bring to mind the abysmal results obtained from the 1936 Literacy Digest Poll referred to early. However, there are two major differences from the uncontrolled sample in the Digest Poll. One is the attempt to select a representative quota sample in design with internet panels. The other is the use of weighting adjustments in the analysis to achieve the same purpose. Before around 1970, lacking today's computers, complex calibration weighting adjustments were infeasible, but now advanced adjustment methods have been developed and are readily employed for both probability samples (particularly those with low response rates) and for nonprobability samples. With river sampling, a limited number of variables can be collected as part of the data collection for use in calibrating the sample to known or estimated population characteristics. The data collected in the screening instrument for an on-line panel can provide a much greater range of variables that can be used in sample selection and in the application of complex calibration adjustments to make the weighted sample correspond to a wide range of external controls. Nevertheless, serious doubts will persist about whether external data are available for the key auxiliary calibration variables at the population level or for a probability sample of that population, and whether the responses to the on-line survey can be treated as equal to the responses from the external source. Thus, for any given survey estimate, there must be concerns about how representative the nonprobability sample members are of the general population within the controls imposed in design or weighting. There will inevitably remain some residual biases of unknown magnitude and, with large samples, these biases can have a dominant influence on the level of accuracy of the survey estimates (Meng, 2018; Kalton, 2021, pp. 136–137).

7. Model-Dependent Inference

In 1976, Fred Smith—my late friend and colleague at the University of Southampton at that time—wrote a paper reviewing the foundations of survey sampling in which he raised the question of why finite population inference should be so different from inference in the rest of statistics. His view at the time was that 'survey statisticians should accept their responsibility for providing stochastic models for finite populations in the same way as statisticians in the experimental sciences' (Smith, 1976); he moderated his position in a subsequent paper (Smith, 1994). Smith (1976) and papers by Brewer (1963), Royall (e.g., 1970, 1976) and others led to a spirited and longstanding debate about the choice between design-based (model-assisted) inference and model-dependent (or model-based) inference. I was a discussant of Fred's 1976 paper and I subsequently published two papers on the role of models in survey sampling inference, with a defense of design-based inference in most circumstances applicable in large-scale social surveys (Kalton, 1983, 2002). However, models are needed to deal

with the sampling imperfections of noncoverage and nonresponse, and they are needed for subgroup analyses in which the sample sizes are not adequate to provide design-based estimators of adequate precision. With the large decline in response rates that has occurred since the 1970's, it is no longer possible for survey statisticians to treat nonresponse as a minor blemish that can be brushed under the carpet in using design-based inference. I will return to this point later.

The model-dependent approach has led to the development of the prediction approach to survey inference. With this approach, an estimate of the population total Y is given by

$$\hat{Y}_m = \sum_{i \in S} y_i + \sum_{i \notin S} \hat{f}(X_i) \quad (2)$$

where the first summation is over the observed values in the sample S of size n and the second summation is over the model predictions of the y values for the nonsampled elements in the population. For comparison with the model-assisted design-based estimator \hat{Y}_d in (1), the model-dependent estimator may be expressed as $\hat{Y}_m = \sum_S e_i + \sum_U \hat{f}(X_i)$. In practice, greater care is used to develop the model for \hat{Y}_m than is the case in developing the working model for \hat{Y}_d . If the same model is used, \hat{Y}_m likely has lower variance than \hat{Y}_d . However, \hat{Y}_m has a design bias if the model is mis-specified, as is always the case to some extent, and the magnitude of the bias is unknown. The texts by Valliant, Dorfman, and Royall (2000) and Chambers and Clark (2012) describe the prediction approach in detail. The first chapter of Valliant et al. (2000) provides a useful review of design-based and model-based inference and includes further references. Note that the equation for \hat{Y}_m does not include selection probabilities (except possibly for estimating the model parameters) and does not require a probability sample. However, as Valliant, Dorfman, and Royall (2000, pp. 19–22) argue, randomization has the benefit of giving some protection against imbalance in factors uncontrolled in the design.

In my experience, until recently the prediction approach has had limited utility for large-scale social surveys of households and persons for the following reasons:

1. As distinct from surveys of establishments, there are generally little, if any, data available from the sampling frame about every member of the target population for use in the prediction models. Although some countries maintain up-to-date population registers that contain a selection of individual characteristics, in many countries area sampling is used, with frame construction for individuals or households being performed only in selected areas. In these latter countries, no frame data is available for all members of the target population.
2. Social surveys are multipurpose in nature. They collect survey data on many variables, often numbering in the hundreds, and these data are analyzed in many ways, producing thousands of estimates. As a rule, these surveys are

primarily conducted to produce descriptive estimates of parameters of the survey's finite population. These estimates need to be produced rapidly and to be consistent with each other. (These days, analytic estimates are also often produced, mostly through secondary analyses—see section 7).

3. A large proportion of the variables collected in social surveys are categorical in nature. They often cannot be as well predicted from auxiliary data as is the case with some of the continuous variables collected in business surveys.

However, even with large-scale social surveys, model-dependent estimation has a role to play in the production of descriptive estimates for small subclasses for which the sample sizes are too small to yield design-based estimates of adequate precision. This situation occurs particularly when the subclasses are geographical-defined administrative areas. The growth of interest by policy makers and others in separate estimates for administrative districts of all sizes has led to the development of the subject known as *small area estimation*. For many years, small area estimates, which are obtained using model-dependent prediction methods, were viewed with considerable skepticism by design-based statisticians but they have now become widely accepted in many fields of application. Ghosh (2020) gives a history of the development of small area estimation over five decades and Rao and Molina (2015) give a detailed description of this large and growing field.

The theoretical developments in model-based inference have now become increasingly relevant for social surveys to address the sampling imperfections and limitations with probability samples, and for the analyses of nonprobability samples; the use of nonprobability sampling for social research has grown rapidly in recent years, in particular for internet surveys.

8. Analytic Uses of Survey Data

As computing power and software came into widespread use in the 1970's, survey data collected using complex sample designs were used, mostly in secondary analyses, to produce analytic statistics that studied the relationships between variables, often looking for causal connections. Initially, multiple regression was the main form of analysis, with interest directed to the magnitude of the regression coefficients. Many analysts argued that their interest in the results of these analyses was not for the specific finite population surveyed, but rather as estimates of superpopulation parameters of universal generality, and that, with the "correct" model, aspects of the sample design were irrelevant. From this perspective, probability sampling of the finite population becomes irrelevant and, unless survey weights and clustering were important as predictor variables, their inclusion in the analysis in a standard design-based way serves only to lower the precision of the estimated regression coefficients. The counter

position was that no model is totally correct and that the estimation of the population regression coefficients, often termed census parameters, using the survey weights provides a safer approach. There is extensive literature on this topic. See, for example, DuMouchel and Duncan (1983).

Over time, the use of regression methods with survey data has been extended to include a wide range of regression models and other multivariate analysis techniques such as categorical data analysis, multilevel modeling, and longitudinal analyses. It is outside the scope of this paper to describe the application of these methods with complex survey data. See Skinner, Holt, and Smith (1989), Chambers and Skinner (2003). Applications of a range of multivariate methods with complex survey data are well described in the texts by Korn and Graubard (1999) and Heeringa, West, and Berglund (2017).

9. Administrative Records and Big Data

A great deal of attention has been paid recently to the use of administrative records as an alternative source of research data. There are obvious serious issues of privacy and confidentiality to be addressed when government-maintained administrative data are used in this way. For this reason, this approach is particularly suited to researchers in government agencies. The approach has notable potential attractions in terms of cost and sample size, but it needs to be recognized that it has its limitations. For instance, what is the coverage of the frame of the records, especially regarding program enrollment versus eligibility? Do the records contain the data needed to measure the concepts as the researcher would like to define them? Are the record data measured consistently across the population, or are there differences in the procedures used in different administrative areas? Are the data measured consistently over time to enable time series data to be validly analyzed? How might changes in program rules affect temporal comparisons? How long is the period between data collection and the researcher's access to an analyzable dataset? Do the records contain the full set of variables needed for the analyses? In many cases, a single set of administrative records does not contain all the variables needed for the analyses. In this situation, it may be possible to link two or more sets of records, but record linkage problems need to be overcome and greater issues of confidentiality must be addressed.

How accurate are the data recorded in the records? Survey researchers have devoted a great deal of effort to training a relatively small number of interviewers to ask and record respondents' answers in a standard way. The situation is different with administrative records. Charlie Cannell, my late friend and colleague at the University of Michigan's Survey Research Center, had the following quotation from Josiah Stamp (1880–1941) in a plaque on his office wall:

“The government are very keen on amassing statistics. They collect them, add them, raise them to the n th power, take the cube root and prepare wonderful diagrams. But you must never forget that every one of these figures comes in the first instance from the village watchman, who just puts down what he damn pleases.”

While not claiming that current administrative records are as bad as this quotation might suggest, those who use such records for statistical purposes should carefully assess their quality and the distortions to which they may be subjected. See the paper by Hand (2018) and the ensuing discussion for a detailed discussion of the advantages and limitations of administrative records for research purposes.

In addition to government-maintained administrative records, there are other sources of social research data. In some cases, nongovernment records, such as those maintained by private organizations, may contain relevant information. However, they are subject to similar quality concerns, and access to the records may be hard to obtain. There are also sources of big data that occur on a flow basis, such as from linking cell phones to their GPS locations. The cell phone locations can be used to provide information about commuter times and even about long-distance travel trips if the home location is identified. Another source of big data is from scrapings on the web. Google Flu Trends (GFT) is a well-known and cautionary example. By analyzing extremely large numbers of flu-related searches on the web, Google developed models to predict local flu outbreaks in real time, avoiding the inevitable delay with other data sources. However, the models have since been found to fail (Lazer, Kennedy, King, and Vespignani, 2014), which serves as a warning that the apparent attraction of very big datasets can be illusory. For another example, see Bradley, Kuriwaki, Isakov, Sejdinovic, Meng, and Flaxman (2021).

10. Concluding Remarks

As illustrated in previous sections, the choice between purposive selection and probability sampling was a subject of debate in the early period of survey research. It was not until after Neyman's (1934) paper that probability sampling and design-based inference were established as the gold standard for large-scale surveys conducted by national statistical offices. With a perfectly executed probability sample and no response error, the analyst has the security of being able to report the survey findings as being subject only to a measurable degree of sampling error, whereas with nonprobability sampling the analyst can always be challenged that a purposive sample is not representative of the population with respect to the variables of analytic interest.

The preeminence of probability sampling for government surveys in the years from 1940 to, say, 2010 was not universal. There are costs incurred with probability sampling

and a probability sample takes more time to draw and data collection takes longer. As illustrated in earlier sections, failures to devise probability sampling methods that can be applied with acceptable cost and timeliness for certain populations has given rise to the development of shortcut methods that depart in varying degrees from rigorous probability sampling.

In the early days, the idea of a “representative sample” was restricted to a sample that was representative in its design, as was the case with Kiær’s designs. The use of weighting adjustments in the analysis to achieve representativeness was seldom considered. The failure of the Literacy Digest poll in predicting the result of the U.S. Presidential election made clear that an extremely large unrepresentative sample could, without weighting adjustments, yield bad results.

Over the years, the implementation of probability sampling in social surveys has been increasingly challenged in many—but not all—countries by a steady decline in the willingness of the public to participate in surveys. Despite greater efforts to encourage response, response rates have declined dramatically in recent years. In reaction, greater efforts have been made to compensate for nonresponse, with major advances in the techniques employed. While replication methods of variance estimation can be applied to reflect the effect of the use of these techniques on the precision of the survey estimates, their use results in lower precision. Furthermore, the nonresponse adjustment model cannot be assumed to be “correct,” and the extent of any remaining nonresponse bias cannot be assessed. With its current heavy reliance on nonresponse models, in many countries probability sampling with design-based inference no longer retains its status as the undisputed gold standard. Moreover, the current levels of nonresponse have led to a marked increase in the costs of conducting a survey with probability sampling, both because of the increase in the initial sample size needed to produce the required sample size and because of the increased efforts to counteract nonresponse. For example, in the U.S. random digit dialing (RDD) was widely used with telephone surveying in the later part of the last century and the early part of this one because of the cost-efficiency of this modality (particularly for surveying rare populations). However, response rates for RDD surveys have plummeted to a level as low as 10 to 20 percent, largely ruling out this form of sampling.

With the security of model-free probability sampling with design-based inference now a thing of the past, model-dependent methods appear to be taking on a major role in social statistics. Research on making valid inferences from nonprobability samples is ongoing (see, for example, Valliant, 2020). Models are increasingly used to analyze data from a combination of data sources, including survey data from probability and nonprobability samples, administrative records, and other sources of big data. Thus, there is much research currently underway on making inferences from combinations

of probability and nonprobability samples and from probability samples and other data sources (Kim and Wang, 2019; Beaumont and Rao, 2021; Rao, 2021),

In summary, after a long period in which probability sampling methods have dominated, the current situation is in a state of flux. New methods involving nonprobability sampling, internet sampling, administrative records, and big data are under constant modification and development. Brackstone (1999) lists six aspects of data quality for a statistical agency that remain applicable: relevance (how well the data meet the needs of the clients); accuracy (including both bias and variance); timeliness (time between the reference point and the time of data availability); interpretability (availability of relevant metadata); and coherence (ability to bring the data into a broader framework, including over time). The new data collection methods need to be assessed against these measures and, furthermore, the extensive research on response errors that has been conducted in the past now needs to be applied with the new methods of data collection. This is an exciting and challenging time for survey methodologists.

References

- Aldrich, J., (2008). Professor A. L. Bowley’s theory of the representative method. (Discussion Papers in Economics and Econometrics, 801) University of Southampton. <https://eprints.soton.ac.uk/150493>.
- Baker, R., Blumberg, S. J., Brick, J. M., Couper, M. P., Courtright, M., Dennis, J. M., Dillman, D., Frankel, M. R., Garland, G., Groves, R. M., Kennedy, C., Krosnick, J., Lavrakas, P. J., (2010). AAPOR Report on Online Panels. *Public Opinion Quarterly*, 74(4), pp. 711–781.
- Bauer J. J., (2014). Selection errors of random route samples. *Sociological Methods and Research*, 43(3), pp. 519–544.
- Bauer J. J., (2016). Biases in random route surveys. *Journal of Survey Statistics and Methodology*, 4(2), pp. 263–287.
- Beaumont J-F., Rao, J. N. K., (2021). Pitfalls of making inferences from non-probability samples: Can data integration through probability samples provide remedies? *The Survey Statistician*, 83, pp. 11–22.
- Bennett, S., (1993). Cluster sampling to assess immunization: a critical appraisal. *Bulletin of the International Statistical Institute*, 49th Session, 55(2), pp. 21–35.
- Bowley, A. L., (1913). Working-class households in Reading. *Journal of the Royal Statistical Society*, 76(7), pp. 672–701.

- Bradley, V. C., Kuriwaki, S., Isakov, M., Sejdinovic, D., Meng, X.-L., Flaxman, S., (2021). Unrepresentative big surveys significantly overestimated US vaccine uptake. *Nature*, 600, pp. 695–700.
- Brewer, K. R. W., (1963). Ratio estimation in finite populations: some results deducible from the assumption of an underlying stochastic process, *Australian Journal of Statistics*, 5, pp. 93–105.
- Caradog Jones, D., (1949). *Social Surveys*. Hutchinson's University Library, London.
- CDC, (2019). Community Assessment for Public Health Emergency Response (CASPER) Toolkit. 3rd ed., CDC, Atlanta. <https://www.cdc.gov/nceh/casper/>.
- Chambers, R., Clark, R., (2012). *An Introduction to Model-Based Survey Sampling with Applications*. Oxford University Press, Oxford.
- Chambers, R. L., Skinner, C. J., Eds., (2003). *Analysis of Survey Data*. Wiley, Chichester.
- Cochran, W. G., (1953). *Sampling Techniques*. Wiley, New York.
- Converse, J. M., (2017). *Survey Research in the United States: Roots and Emergence 1890-1960*. Routledge, New York.
- DuMouchel, W. H., Duncan, G. J., (1983). Using sample survey weights in multiple regression analyses of stratified samples. *Journal of the American Statistical Association*, 78, pp. 535–543.
- Ghosh, M., (2020). Small area estimation: its evolution in five decades (with discussion). *Statistics in Transition*, 21(4), pp. 1–67.
- Gile, K. J., Hancock, M. S., (2010). Respondent-driven sampling: an assessment of current methodology. *Sociological Methodology*, 40(1), pp. 285–327.
- Gini, C., Galvani, L., (1929). Di una applicazione del metodo rappresentativo. *Annali di Statistica*, 6(4), pp. 1–107.
- Hand, D. J., (2018). Statistical challenges of administrative and transaction data (with discussion). *Journal of the Royal Statistical Society, A*, 181(3), pp. 555–605.
- Hansen, M. H., Hurwitz, W. N. and Madow, W. G. (1953). *Sample Survey Methods and Theory. Volume I: Methods and Applications. Volume II: Theory*. Wiley, New York.
- Heckathorn, D. D., (1997). Respondent-driven sampling: a new approach to the study of hidden populations. *Social Problems*, 44(2), pp. 174–199.
- Heeringa, S. G., West, B. T., Berglund, P. A., (2017). *Applied Survey Data Analysis*. Chapman & Hall/ CRC, Boca Raton, FL.

- Jensen, A., (1926) The report on the representative method in statistics. *Bulletin of the International Statistical Institute*, 22, pp. 355–376.
- Kalton, G., (1983). Models in the practice of survey sampling. *International Statistical Review*, 51, pp. 175–188.
- Kalton, G., (1991). Sampling flows of mobile human populations. *Survey Methodology*, 17(2), pp. 183–194.
- Kalton, G., (2002). Models in the practice of survey sampling (revisited). *Journal of Official Statistics*, 18, pp. 129–154.
- Kalton, G., (2019). Developments in survey research over the past 60 years: A personal perspective. *International Statistical Review*, 87 (S1), pp. S10–S30.
- Kalton, G., (2021). *Introduction to Survey Sampling*. 2nd ed. SAGE Publications, Thousand Oaks, California.
- Kiær, A. N., (1976). *The Representative Method of Statistical Surveys*. English translation, Statistisk Centralbyro, Oslo.
- Kim, J. K., Wang, Z., (2019). Sampling techniques for big data analysis. *International Statistical Review*, 87 (S1), pp. S177–S191.
- Kish, L., (1995). The hundred years' war of survey sampling. *Statistics in Transition*, 2(5), pp. 813–830.
- Korn, E. L., Graubard, B. I., (1999). *Analysis of Health Surveys*. Wiley, New York.
- Kruskal, W., Mosteller, F., (1980). Representative sampling, IV: The history of the concept in statistics, 1895-1939. *International Statistical Review*, 48(2), pp. 169–195.
- Lazer, D., Kennedy, R., King, G., Vespignani, A., (2014). The parable of Google flu: Traps in big data analysis. *Science*, 343, pp. 1203–1205.
- Levy, P. S., Lemeshow, S., (2008). *Sampling of Populations. Methods and Applications*. 4th ed. Wiley, Hoboken, NJ.
- Lie, E., (2002). The rise and fall of sampling surveys in Norway, 1875–1906. *Science in Context*, 15(3), pp. 385–409.
- Lohr, S. L., Brick, J. M., (2017). Roosevelt predicted to win: Revisiting the 1936 Literary Digest Poll. *Statistics, Politics, and Policy*, 8(1), pp. 65–84.
- MacKellar, D. A., Gallagher, K. M., Findlayson, T., Sanchez, T., Lansky, A., Sullivan, P. S., (2007). Surveillance of HIV risk and prevention behaviors of men who have sex with men—a national application of venue-based, time-space sampling. *Public Health Reports*, 122 (1), Supplement 1, pp. 39–47.

- Mahalanobis, P. C., (1946). Recent experiments in statistical sampling in the Indian Statistical Institute (with discussion). *Journal of the Royal Statistical Society*, 109, pp. 325–378.
- Meng, X-L., (2018). Statistical paradises and paradoxes in big data (I): Law of large populations, big data paradox, and the 2016 US presidential election. *Annals of Applied Statistics*, 12(2), pp. 685–726.
- Moser, C. A., Kalton, G., (1971). *Surveys Methods in Social Investigation*. 2nd ed. Heinemann, London.
- Moser, C. A., Stuart, A., (1953). An experimental study of quota sampling. *Journal of the Royal Statistical Society, A*, 116, pp. 349–405.
- Neyman, J., (1934). On 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, pp. 558–625.
- Rao, J. N. K., (2021). On making valid inferences by integrating data from surveys and other sources. *Sankhya B*, 83, pp. 242–272.
- Rao, J. N. K., Molina, I., (2015). *Small Area Estimation*. 2nd ed. Wiley, Hoboken, N. J.
- Royall, R. M., (1970). On finite population sampling theory under certain regression models. *Biometrika*, 57, pp. 377–387.
- Royall, R. M., (1976). The linear least squares prediction approach to two-stage sampling. *Journal of the American Statistical Association*, 71, pp. 657–664.
- Särndal, C.E., Swensson, B., Wretman, J., (1992). *Model Assisted Survey Sampling*. Springer-Verlag, New York.
- Skinner, C. J., Holt, D., Smith, T. M. F., Eds., (1989). *Analysis of Complex Surveys*. Wiley, Chichester.
- Smith, T. M. F., (1976). The foundations of survey sampling: a review. *Journal of the Royal Statistical Society, A*, 139, pp. 183–204.
- Smith, T. M. F., (1994). Sample surveys 1975-90; an age of reconciliation? *International Statistical Review*, 62, pp. 5–34.
- Stephan, F. F., (1948). History of the uses of modern sampling procedures. *Journal of the American Statistical Association*, 43(241), pp. 12–39.
- Stephan, F. F., McCarthy P. J., (1958). *Sampling Opinions. An Analysis of Survey Procedures*. Wiley, New

- Stephenson, C. B., (1979). Probability sampling with quotas: An experiment. *Public Opinion Quarterly*, 43(4), pp. 477–497.
- Sudman, S., (1966). Probability sampling with quotas. *Journal of the American Statistical Association*, 61, pp. 749–771.
- Tourangeau, R., Edwards, B., Johnson, T. P., Wolter, K.M., Bates, N., Eds., (2014). *Hard-to-Survey Populations*. Cambridge University Press, Cambridge, U. K.
- Valliant, R. (2020). Comparing alternatives for estimation from nonprobability samples. *Journal of Survey Statistics and Methodology*, 8(2), pp. 231–263.
- Valliant, R., Dorfman, A. H., Royall, R. M., (2000). *Finite Population Sampling and Inference. A Prediction Approach*. Wiley, New York.
- Yates, F., (1949). *Sampling Methods for Censuses and Surveys*. Griffen, London.

Comments on „Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day” by Graham Kalton

Danny Pfeffermann¹

I like to congratulate Professor Kalton for writing this very constructive article on probability versus nonprobability sampling. I learned a lot from reading it. In what follows, I add a few comments on this topic.

1- Professor Kalton emphasizes the issue of representative samples. In my view, probability samples and obviously nonprobability samples are practically never representative, even if balanced in advance on certain control (covariate) variables. A major reason for this is nonresponse, which might be “not missing at random” (NMAR), in which case the response probabilities depend on the target study variable, even after conditioning on known covariates. However, even in the case of simple random sampling and complete response, the actual sample may not be representative with respect to the unknown study variables, simply because of the randomness of the sample selection, unless the sample size is sufficiently large. Clearly, this problem worsens when sampling with unequal probabilities. Classical design-based theory overcomes this problem by restricting the inference to the randomization distribution over all possible sample selections. Thus, an estimator of a population mean is unbiased if its average over all possible samples that could have been drawn equals the true population mean, but in practice, we only have one sample. The use of models does not solve this problem either. A good model has to account for the sampling probabilities and the model assumed for the population values, and the inference need to account for both stochastic processes. As illustrated in many articles, ignoring the sampling process when fitting models to the sample data results with biased estimators of the model parameters in the case of informative sampling, by which the sampling probabilities are correlated with the outcome variables, again after conditioning on the model covariates. See, e.g. Pfeffermann and Sverchkov (1999) for empirical illustrations. In the case of NMAR nonresponse, the model has to account also for the unknown response probabilities.

¹ Department of Statistics, Hebrew University, Jerusalem, Israel & Southampton Statistical Sciences Research Institute, University of Southampton, UK. E-mail: msdanny@mail.huji.ac.il; msdanny@soton.ac.uk. ORCID: <https://orcid.org/0000-0001-7573-2829>.

2- The problem of nonresponse is indeed troubling and requires the use of models in the case of NMAR nonresponse, even in the case of design-based inference. The use of a response model enables to adjust the base sampling weights by the inverse of the estimated response probabilities, viewed as a second stage of the sampling process. I should say though that unlike a common perception, the response model can be tested, by testing the model of the study variable holding for the responding units, which accounts for the sampling design and the response. See, e.g. Pfeffermann and Sikov (2011).

3- Professor Kalton discusses the pros and cons of internet surveys “standing on their own”. I like to add that internet surveys are often used as one, out of several possible modes of response. For example, a questionnaire is sent to all the sampled units. It encourages them to respond via the internet. Those who do not respond are approached by telephone. When no response is obtained, an interviewer is sent for a face-to-face interview.

A well-known problem with this procedure is of mode effects; different estimates obtained from the respondents to the different modes, either because of differences between the characteristics of respondents responding with the different modes, (selection effect), or because of responding differently by the same sampled unit, depending on the mode of response (measurement effect). Several approaches to deal with this problem have been proposed in the literature. See, e.g. De Leeuw et al. (2018) for a comprehensive review.

My last 2 comments refer to inference from nonprobability samples:

4- Denote by S_{NP} the nonprobability sample. Rivers (2007) proposes to deal with the possible non-representativeness of S_{NP} by the use of sample matching. (Rivers considers a Web sample as the nonprobability sample but here I extend the idea to a more general nonprobability sample.) The approach consists of using a probability (reference) sample S_R from the target population, drawn with probabilities $\pi_k = \Pr(k \in S_R)$, and matching to every unit $i \in S_R$ an element $k \in S_{NP}$, based on known auxiliary (matching) variables x . Denote by S_M the matched sample. Suppose that it is desired to estimate a population total of a study variable Y , based on measurements $\{\tilde{y}_j, j \in S_{NP}\}$. Estimate, $\hat{Y}_T = \sum_{j \in S_M} w_j \tilde{y}_j$; $w_j = (1 / \pi_j)$. Clearly, the base sampling weights can be modified to account for nonresponse. This is an intriguing approach, but its success depends on the existence of a reference probability sample S_R , which allows sufficiently close matching, and ignorability of membership in the nonprobability sample S_{NP} , conditional upon

the matching variables. I do not know whether this approach is used in practice, but I think that it deserves further investigation, with proper modifications.

- 5- The last two decades have witnessed the rapid growing of data science. One of the facets of this growth is that some people are agitating that the existence of all sorts of “big data” and the new advanced technologies that have been developed to handle these data, will soon replace the use of sample surveys. In an article I published in 2015, I overviewed some of the problems with the use of big data for the production of official statistics but clearly, when such data sources are available, accessible and timely, they cannot and should not be ignored. Big data can be viewed as a big, nonprobability sample, which for all kinds of reasons is not representative of the target population, and relying just on them can yield biased inference. Integrating big data with surveys is a major issue for research. See, e.g. Kim and Zhonglei (2018) and Rao (2021) for possible approaches, with references to other studies.

I conclude my discussion by congratulating Statistics in Transition for its 30th anniversary and the publication of its 100th issue. This is one of the best journals of its kind and I wish it to continue prospering in the coming years.

References

- De Leeuw, E. D., Suzer-Gurtekin, Z. and Hox, J., (2018). The Design and Implementation of Mixed Mode Surveys. In *Advances in Comparative Survey Methodology*. Wiley, New York.
- Kim, J. K. and Zhonglei Wang, (2008). Sampling Techniques for Big Data. *International Statistical Review*, 87, pp. 177–191.
- Pfeffermann, D., (2015). Methodological Issues and Challenges in the Production of Official Statistics. The *Journal of Survey Statistics and Methodology* (JSSAM), 3, pp. 425–483.
- Pfeffermann, D. and Sverckov, M., (1999). Parametric and Semi-Parametric Estimation of Regression Models Fitted to Survey Data. *Sankhya*, 61, pp. 166–186.
- Pfeffermann, D. and Sikov, A., (2011). Imputation and Estimation under Nonignorable Nonresponse in Household Surveys with Missing Covariate Information. *Journal of Official Statistics*, 27, pp. 181–209.
- Rao, J. N. K., (2021). On making valid inferences by integrating data from surveys and other sources. *Sankhya B*, 83, pp. 242–272.
- Rivers, D., (2007). Sampling for Web Surveys. Joint Statistical Meetings, Proceedings of the Section on Survey Research Methods.

Comments on „Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day” by Graham Kalton

Risto Lehtonen¹

I would like to congratulate Professor Graham Kalton for his significant and inspiring article entitled as "Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day". The article provides an elegant overview of the history of survey sampling, covering the purposive approaches that dominated the sampling field in the early days but from the 1940s, at least in official statistics, were gradually replaced entirely by probability-based approaches. Today we may be facing a paradigm shift again, but the direction is the opposite. Non-probability-based approaches are becoming viable, if not the only option, in fields that are moving towards big data and other new data sources and new methodological approaches.

The country's data infrastructure forms the basis of official statistics and opens up for me an important perspective on Kalton's presentation. Both probability and non-probability sampling and inference can benefit from statistical data infrastructures that contain a rich selection of micro-level covariates drawn from a variety of administrative and other registers. Perhaps the best options are in countries where population data from register sources and sample data are linked for combined micro-level databases. However, the utility of model-based (prediction) approaches for large-scale social surveys of households and persons will be limited if unit-level data for population members is missing from the sampling frames, as pointed out by Prof. Kalton. This is an important point and I think it can be extended to design-based model-assisted approaches that use mixed models in particular.

Countries differ much in terms of infrastructures based on administrative data. For example, Constance Citro calls for a move to multiple data sources that include administrative records and, increasingly, transaction and Internet-based data (Citro 2014). Eric Rancourt argues that Statistics Canada is facing the new data world by modernizing itself and embracing an admin-first (in the broadest sense) paradigm as a statistical paradigm for the agency (Rancourt 2018). According to the United

Nations Economic Commission for Europe (UNECE) report on register-based statistics in the Nordic countries, Central Population Registers of Denmark, Finland, Norway and Sweden were established in the sixties, and for example a totally register-based census was first implemented in Denmark (1981) and next in Finland (1990) (UNECE 2007). The number of national statistical institutes that have adopted or are developing administrative data infrastructures is increasing, as also described in the UNECE report on the use of registers and administrative data for population and housing censuses (UNECE 2018). This development can enhance the use of methods that utilize modeling and individual-level population frame data for model-assisted or prediction-based estimation with probability-based or non-probability-based sample data sets and their combinations.

The situation is different in countries that do not have similar high-quality population registers as for example in the Nordic countries. A recent contribution by Dunne and Zhang (2023) provides one important methodological approach for such countries. The authors present an innovative system (the PECADO application) for population estimates compiled from administrative data only.

Today, in the Nordic countries, as Finland, a majority of official statistics are based on administrative register combinations. In Finland, official statistics are produced by 13 expert organisations in the field of public administration and is coordinated by Statistics Finland. Probability samples are mainly used for regular social surveys such as labour force surveys and special surveys, e.g. Time Use survey. In these surveys, the sample elements can be uniquely linked with the elements in the register databases that often contain a lot of important background data including demographic, regional, socio-economic, income, educational, labour force status, and other variables. Thus these data need not to be collected by direct data collection methods from the respondents, and measurement errors are avoided. In addition, these variables are also used for calibration and model-assisted estimation procedures.

As an example, let me describe briefly the sampling and estimation design of the Labour Force Survey (LFS) of Finland. According to the quality description, in most European countries the LFS is based on a sample of households, and all members of a sample household living at the same address are interviewed. Finland is one of the Nordic countries where LFS is based on sampling of individual persons. The sample of about 12,500 persons is drawn by stratified probability sampling from Statistics Finland's population database, which is based on the Central Population Register. Auxiliary information from registers include gender, age, region and language and selected register variables on employment, completed education and degrees, and income from the Employment Service Statistics of the Ministry of Economic Affairs and Employment, Statistics Finland's Register of Completed Education and Degrees, and the Tax Administration's Incomes Register (Quality Description: Labour Force

¹ University of Helsinki, Finland. E-mail: risto.lehtonen@helsinki.fi.

Survey, Statistics Finland 2022). Sample data are linked to data from the registry using unique ID keys that exist across all data sources and are used in estimation procedures, including nonresponse adjustments. My experience is that this type of data infrastructure can also provide an excellent sampling and auxiliary data platform for e.g. methodological research in survey statistics; see for example Lehtonen, Särndal and Veijanen (2003, 2005).

Data infrastructures based on integrated administrative and other registers should be based on appropriate statistical theory and methodology for quality assessment and control and quality improvement. Recent sources in the field are for example Zhang (2012), Zhang and Haraldsen (2022) and the book on register-based statistics by Anders Wallgren and Britt Wallgren (2014). Research in statistical data integration and data science methods relevant for official statistics also is extending. A recent source is Yang and Kim (2020).

Experiences show that data infrastructures for official statistic containing a wealth of micro-level information on the population and an option for integration of the various register and sample data sources provide a flexible and efficient framework for survey estimation with probability-based samples. For non-probability samples, the variables of interest are typically in the non-probability data source. Most current methods for valid inference require an auxiliary data source containing the same covariates as the non-probability sample. These data can be obtained from the statistical population register or, more commonly, from a probability sample from it (e.g. Kim, Park, Chen and Wu 2021; Wu 2022). It can be foreseen that although the golden age of probability sampling may be over, probability sampling and non-probability sampling are not in conflict, but can complement each other.

References

- Citro, C. F., (2014). From multiple modes for surveys to multiple data sources for estimates. *Survey Methodology*, 40(2), pp. 137–161.
- Dunne, J. and Zhang, L.-C., (2023). A system of population estimates compiled from administrative data only. *Journal of the Royal Statistical Society Series A: Statistics in Society*. <https://doi.org/10.1093/jrssa/qnad065>.
- Kim, J.-K., Park, S., Chen, Y. and Wu, C., (2021). Combining non-probability and probability survey samples through mass imputation. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 184, pp. 941–963.
- Lehtonen, R., Särndal, C.-E. and Veijanen, A., (2003). The effect of model choice in estimation for domains, including small domains. *Survey Methodology*, 29(1), pp. 33–44.

- Lehtonen, R., Särndal, C.-E. and Veijanen, A., (2005). Does the model matter? Comparing model-assisted and model-dependent estimators of class frequencies for domains. *Statistics in Transition*, 7(3), pp. 649–673.
- Quality Description: Labour force survey, Statistics Finland 2022, (2022). https://www.tilastokeskus.fi/til/tyti/2022/01/tyti_2022_01_2022-02-22_laa_001_en.html
- Rancourt, E., (2018). *Admin-First as a statistical paradigm for Canadian official statistics: Meaning, challenges and opportunities*. Proceedings of Statistics Canada Symposium 2018.
- United Nations Economic Commission for Europe, (2007). *Register-based statistics in the Nordic countries: Review of best practices with focus on population and social statistics*. United Nations, New York. <https://digitallibrary.un.org/record/609979?ln=en>
- UNECE, (2018). *Guidelines on the use of registers and administrative data for population and housing censuses*. United Nations, New York and Geneva. <https://unece.org/guidelines-use-registers-and-administrative-data-population-and-housing-censuses-0>
- Yang, S. and Kim, J. K., (2020). Statistical data integration in survey sampling: a review. *Jpn J Stat Data Sci*, 3, pp. 625–650.
- Zhang, L.-C., (2012). Topics of statistical theory for register-based statistics and data integration. *Statistica Neerlandica*, 66(1), pp. 41–63.
- Zhang, L.-C. and Haraldsen, G., (2022). Secure big data collection and processing: framework, means and opportunities. *Journal of the Royal Statistical Society: Series A, Statistics in Society*, (In Press).
- Wallgren, A. and Wallgren, B., (2014). *Register-Based Statistics: Administrative Data for Statistical Purposes*. Second edition. Wiley.
- Wu, C., (2022). Statistical inference with non-probability survey samples. *Survey Methodology*, 48(2), pp. 283–311.

STATISTICS IN TRANSITION new series, June 2023

Vol. 24, No. 3, pp. 31–37, <https://doi.org/10.59170/stattrans-2023-032>

Received – 25.05.2023; accepted – 31.05.2023

Discussion of “Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day” by Graham Kalton

Julie Gershunskaya¹, Partha Lahiri²

In this excellent overview of the history of probability and nonprobability sampling from the end of the nineteenth century to the present day, Professor Graham Kalton outlines the essence of past endeavors that helped to define philosophical approaches and stimulate the development of survey sampling methodologies. From the beginning, there was an understanding that a sample should, in some ways, resemble the population under study. In Kær’s ideas of “representative sampling” and Neyman’s invention of probability-based approach, the prime concern of survey sampling has been to properly plan for representing characteristics of the finite population. Poststratification and other calibration methods were developed for the same important goal of better representation.

Professor Kalton’s paper underscores growing interest in the use of nonprobability surveys. With recent proliferation of computers and the internet, wealth of data becomes available to researchers. However, “opportunistic” information collected with present-day capabilities usually is not purposely planned or controlled by survey statisticians. No matter how big such a nonprobability sample could be, it may inaccurately reflect the finite population of interest, thus presenting a substantial risk of an estimation bias.

Below, we discuss several recent papers that propose ways to incorporate nonprobability surveys to produce estimates for both large and small areas. Specifically, we will consider two situations often encountered in practice. In the first situation, a nonprobability sample contains the outcome variable of interest, and the main task is to reduce the selection bias with the help of a reference probability sample that does not contain the outcome variable of interest. In the second situation, a probability sample contains the outcome variable of interest, but there is little or no sample available to produce granular level estimates. For such a small area estimation problem, we consider a case when we have access to a large nonprobability sample that does not contain the outcome variable but contains some related auxiliary variables also present in the probability sample. In both situations, researchers have discussed statistical data integration techniques in which a reference probability sample is combined with a nonprobability sample in an effort to overcome deficiencies associated with both probability and nonprobability samples.

¹U.S. Bureau of Labor Statistics, 2 Massachusetts Ave NE Washington, DC 20212, USA, E-mail: Gershunskaya.Julie@bls.gov. ORCID: <https://orcid.org/0000-0002-0096-186X>.

²University of Maryland, College Park, MD 20742. USA. E-mail: plahiri@umd.edu. ORCID: <https://orcid.org/0000-0002-7103-545X>.

© Julie Gershunskaya, Partha Lahiri. Article available under the CC BY-SA 4.0 licence



One way to account for the selection bias of a nonprobability sample is by estimating the sample inclusion probabilities, given available covariates. Then, the inverse values of estimated inclusion probabilities are used, in a similar manner as the usual probability sample selection weights, to obtain estimates of target quantities. Several approaches to estimation of nonprobability sample inclusion probabilities (or propensity scores) have been considered in the literature. Recent papers by Chen et al. (2020), Wang et al. (2021), and Savitsky et al. (2022) propose ways to estimate these probabilities based on combining nonprobability and probability samples. Kim J. and K. Morikawa (2023) propose an empirical likelihood based approach under a different setting. To save space, we will not discuss their approach. We now review three statistical data integration methods.

The approaches concern with the estimation of probabilities $\pi_{ci}(x_i) = P\{c_i = 1|x_i\}$ to be included into the nonprobability sample S_c , for units $i = 1, \dots, n_c$, where c_i is the inclusion indicator of unit i taking on the value of 1 if unit i is included into the nonprobability sample, and 0 otherwise; x_i is a vector of known covariates for unit i ; n_c is the total number of units in sample S_c . The problem, of course, is that we cannot estimate π_{ci} based on the set of units in nonprobability sample S_c alone, because $c_i = 1$ for all i in S_c . The probabilities are estimated by combining set S_c with a probability sample S_r . Due to its role in this approach, the probability sample here is also called “reference sample”.

Assuming both nonprobability and probability samples are selected from the same finite population P , Chen et al. (2020) write a log-likelihood, over units in P , for the Bernoulli variable c_i :

$$\ell_1(\theta) = \sum_{i \in P} \{c_i \log[\pi_{ci}(x_i, \theta)] + (1 - c_i) \log[1 - \pi_{ci}(x_i, \theta)]\}, \quad (1)$$

where θ is the parameter vector in a logistic regression model for π_{ci} .

Since finite population units are not observed, Chen et al. (2020) employ a clever trick and re-group the sum in (1) by presenting it as a sum of two parts: part 1 involves the sum over the nonprobability sample units and part 2 is the sum over the whole finite population:

$$\ell_1(\theta) = \sum_{i \in S_c} \log \left[\frac{\pi_{ci}(x_i, \theta)}{1 - \pi_{ci}(x_i, \theta)} \right] + \sum_{i \in P} \log [1 - \pi_{ci}(x_i, \theta)]. \quad (2)$$

Units in part 1 of the log-likelihood in (2) are observed; for part 2, Chen et al. (2020) employ the pseudo-likelihood approach by replacing the sum over the finite population with its probability sample based estimate:

$$\hat{\ell}_1(\theta) = \sum_{i \in S_c} \log \left[\frac{\pi_{ci}(x_i, \theta)}{1 - \pi_{ci}(x_i, \theta)} \right] + \sum_{i \in S_r} w_{ri} \log [1 - \pi_{ci}(x_i, \theta)], \quad (3)$$

where weights $w_{ri} = 1/\pi_{ri}$ are inverse values of the reference sample inclusion probabilities π_{ri} . Estimates are obtained by solving respective pseudo-likelihood based estimating equations.

One shortcoming of the Chen et al. (2020) approach is that their Bernoulli likelihood is formulated with respect to an unobserved indicator variable. Although the regrouping

employed in (2) helps to find a solution, results obtained by Wang et al. (2021) indicate that it is relatively inefficient, especially when the nonprobability sample size is much larger than the probability sample size.

Wang et al. (2021) formulate their likelihood for an *observed* indicator variable and thus their method is different from the approach of Chen et al. (2020). To elaborate, Wang et al. (2021) introduce an imaginary construct consisting of two parts: they *stack* together non-probability sample S_c (part 1) and finite population P (part 2). Since nonprobability sample units belong to the finite population, they appear in the stacked set twice. Let indicator variable $\delta_i = 1$ if unit i belongs to part 1, and $\delta_i = 0$ if i belongs to part 2 of the stacked set; the probabilities of being in part 1 of the stacked set are denoted by $\pi_{\delta_i}(x_i) = P\{\delta_i = 1|x_i\}$. Wang et al. (2021) assume the following Bernoulli likelihood for observed variable δ_i :

$$\ell_2(\tilde{\theta}) = \sum_{i \in S_c} \log [\pi_{\delta_i}(x_i, \tilde{\theta})] + \sum_{i \in P} \log [1 - \pi_{\delta_i}(x_i, \tilde{\theta})], \quad (4)$$

where $\tilde{\theta}$ is the parameter vector in a logistic regression model for π_{δ_i} . Since the finite population is not available, they apply the following pseudo-likelihood approach:

$$\hat{\ell}_2(\tilde{\theta}) = \sum_{i \in S_c} \log [\pi_{\delta_i}(x_i, \tilde{\theta})] + \sum_{i \in S_r} w_{ri} \log [1 - \pi_{\delta_i}(x_i, \tilde{\theta})]. \quad (5)$$

Existing ready-to-use software can be used to obtain estimates of π_{δ_i} . However, the actual goal is to find probabilities π_{ci} rather than probabilities π_{δ_i} . Wang et al. (2021) propose a two-step approach, where at the second step, they find π_{ci} by employing the following identity:

$$\pi_{\delta_i} = \frac{\pi_{ci}}{1 + \pi_{ci}}. \quad (6)$$

Savitsky et al. (2022) use an exact likelihood for the estimation of inclusion probabilities π_{ci} , rather than a pseudo-likelihood based estimation. They propose to stack together nonprobability, S_c , and probability, S_r , samples. In this stacked set, S , indicator variable z_i takes the value of 1 if unit i belongs to the nonprobability sample (part 1), and 0 if unit i belongs to the probability sample (part 2). In this construction, if there is an overlap between the two samples, S_c and S_r , then the overlapping units are included into stacked set S twice: once as a part of the nonprobability sample (with $z_i = 1$) and once as a part of the reference probability sample (with $z_i = 0$). We do not need to know which units overlap or whether there are any overlapping units. The authors use first principles to prove the following relationship between probabilities $\pi_{z_i}(x_i) = P\{z_i = 1|x_i\}$ of being in part 1 of the stacked set and the sample inclusion probabilities π_{ci} and π_{ri} :

$$\pi_{z_i} = \frac{\pi_{ci}}{\pi_{ri} + \pi_{ci}}. \quad (7)$$

A similar expression (7) was derived by Elliott (2009) and Elliott and Valliant (2017) under the assumption of non-overlapping nonprobability and probability samples. The derivation given in Savitsky et al. (2022) does not require this assumption.

To obtain estimates of π_{ci} from the combined sample, Beresovsky (2019) proposed to parameterize probabilities $\pi_{ci} = \pi_{ci}(x_i, \theta)$, as in Chen et al. (2020), and employ identity (7) to present π_{z_i} as a composite function of θ ; that is, $\pi_{z_i} = \pi_{z_i}(\pi_{ci}(x_i, \theta)) = \pi_{ci}(x_i, \theta) / (\pi_{ri} + \pi_{ci}(x_i, \theta))$.

The log-likelihood for observed Bernoulli variable z_i is given by

$$\ell_3(\theta) = \sum_{i \in S_c} \log [\pi_{z_i}(\pi_{ci}(x_i, \theta))] + \sum_{i \in S_r} \log [1 - \pi_{z_i}(\pi_{ci}(x_i, \theta))]. \quad (8)$$

Since the log-likelihood *implicitly* includes a logistic regression model formulation for probabilities π_{ci} , Beresovsky (2019) labeled the proposed approach Implicit Logistic Regression (ILR). For the maximum likelihood estimation (MLE), the score equations are obtained from (8) by taking the derivatives, with respect to θ , of the composite function $\pi_{z_i} = \pi_{z_i}(\pi_{ci}(\theta))$. This way, the estimates of π_{ci} are obtained directly from (8) in a single step. Savitsky et al. (2022) parameterized the likelihood, as in (8), and used the Bayesian estimation technique to fit the model.

Note that to implement the ILR approach, the reference sample inclusion probabilities π_{ri} have to be known for all units in the combined set. This is not a limitation for many probability surveys. As discussed in Elliott and Valliant (2017), if probabilities π_{ri} cannot be determined exactly for units in the nonprobability sample, they can be estimated using a regression model. Savitsky et al. (2022) used Bayesian computations to simultaneously estimate π_{ri} and π_{ci} for nonprobability sample units, given available covariates x_i .

It must be noted that the estimation method of Wang et al. (2021) can be similarly modified to avoid the two-step estimation procedure: a logistic regression model could be formulated for inclusion probabilities π_{ci} , while probabilities π_{δ_i} in (6) could be viewed as a composite function, $\pi_{\delta_i} = \pi_{\delta_i}(\pi_{ci}(x_i, \theta)) = \pi_{ci}(x_i, \theta) / (1 + \pi_{ci}(x_i, \theta))$. This approach is expected to be more efficient. Moreover, it avoids π_{ci} estimates greater than 1 that could occur when the estimation is performed in two steps. Once modified this way, preliminary simulations indicate that Wang et al. (2021) formulation would produce more efficient estimates than the Chen et al. (2020) counterpart, unless in a rare situation where the whole finite population rather than only a reference sample is available.

Simulations show that the exact likelihood method based on formulation of Savitsky et al. (2022) and Beresovsky (2019) performs better than the pseudo-likelihood based alternatives. In the usual situation where the reference probability sample fraction is small, the relative benefits of the exact likelihood approach are even more pronounced.

The existence of a well-designed probability reference sample plays a crucial role in the efforts to reduce the selection bias of a nonprobability sample. Importantly, an ongoing research indicates that the quality of estimates of the nonprobability sample inclusion probabilities is better if there is a good overlap in domains constructed using covariates from both samples. This observation harks back to problems appearing in traditional poststratification methods and to the notion of "representative sampling." Since survey practitioners usually do not have control over the planning or collection of the emerging multitude of nonrandom opportunistic samples, efforts should be directed to developing and maintaining comprehensive probability samples that include sets of good quality covariates. Beaumont et al. (2023)

proposed several model selection methods in application of the modeling nonprobability sample inclusion probabilities.

We now turn our attention to the second data integration situation involving small area estimation, a topic Professor Kalton touched on. This is a problem of great interest for making public policies, fund allocation and regional planning. Small area estimation programs already exist in some national statistical organizations such as the Small Area Income and Poverty Estimates (SAIPE) program of the US Census Bureau (Bell et al., 2016) and Chilean government system (Casas-Cordero Valencia et al., 2016.) The importance placed in the United Nations Sustainable Development Goals (SDG) for disaggregated level statistics is expected to increase the demand for such programs in various national statistical offices worldwide. Standard small area estimation methods generally use statistical models (e.g., mixed models) that combine probability sample data with administrative or census data containing auxiliary variables correlated with the outcome variable of interest. For a review of different small area models and methods, see Jiang and Lahiri (2006), Rao and Molina (2015), Ghosh (2020), and others.

A key to success in small area estimation is to find relevant auxiliary variables not only in the probability sample survey but also in the supplementary big databases. Use of a big probability or nonprobability sample survey could be useful here as surveys typically contain a large number of auxiliary variables that are also available in the probability sample survey. In the context of small area estimation, Sen and Lahiri (2023) considered a statistical data integration technique in which a small probability survey containing the outcome variable of interest is statistically linked with a much bigger probability sample, which does not contain the outcome variable but contains many auxiliary variables also present in the smaller sample. They essentially fitted a mixed model to the smaller probability sample that connects the outcome variable to a set of auxiliary variables and then imputed the outcome variable for all units of the bigger probability sample using the fitted model and auxiliary variables. Finally, they suggested to produce small area estimates using survey weights and imputed values of the outcome variable contained in the bigger probability sample survey. As discussed in their paper, such a method can be used even if the bigger sample is a nonprobability survey using weights constructed by methods such as the ones described earlier.

The development of new approaches demonstrates how the methods of survey estimation continue to evolve by taking into the future the best from fruitful theoretical and methodological developments of the past. As Professor Kalton highlights, we will increasingly encounter data sources that are not produced by standard probability sample designs. Statisticians will find ways to respond to new challenges, as is reflected in the following amusing quote:

...D.J. Finney once wrote about the statistician whose client comes in and says, "Here is my mountain of trash. Find the gems that lie therein." Finney's advice was to not throw him out of the office but to attempt to find out what he considers "gems". After all, if the trained statistician does not help, he will find someone who will....(source: David Salsburg, ASA Connect Discussion)

Of course, nonprobability samples should not be viewed as a "mountain of trash." Indeed, they can contain a lot of relevant information for producing necessary estimates. It is just that one needs to explore different innovative ways to use information contained in nonprobability samples. In the United States federal statistical system, the need to innovate for combining information from multiple sources has been emphasized in the National Academies of Sciences and Medicine (2017) report on Innovations in Federal Statistics. As discussed, statisticians have been already engaged in suggesting new ideas, such as statistical data integration, to extract information out of multiple non-traditional databases. In coming years, statisticians will be increasingly occupied with finding solutions for obtaining useful information from non-traditional data sources. This is indeed an exciting time for survey statisticians.

References

- Beaumont, J.-F., K. Bosa, A. Brennan, J. Charlebois, and K. Chu (2023). Handling non-probability samples through inverse probability weighting with an application to statistics canada's crowdsourcing data. *Survey Methodology* (accepted in 2023 and expected to appear in 2024).
- Bell, W. R., W. W. Basel, and J. J. Maples (2016). *An overview of the US Census Bureau's small area income and poverty estimates program*, pp. 349–378. Wiley Online Library.
- Beresovsky, V. (2019). On application of a response propensity model to estimation from web samples. In ResearchGate.
- Casas-Cordero Valencia, C., J. Encina, and P. Lahiri (2016). *Poverty mapping for the Chilean Comunas*, pp. 379–404. Wiley Online Library.
- Chen, Y., P. Li, and C. Wu (2020). Doubly robust inference with nonprobability survey samples. *Journal of the American Statistical Association* 115(532), 2011–2021.
- Elliott, M. R. (2009). Combining data from probability and non-probability samples using pseudo-weights. *Survey Practice* 2, 813–845.
- Elliott, M. R. and R. Valliant (2017). Inference for Nonprobability Samples. *Statistical Science* 32(2), 249 – 264.
- Ghosh, M. (2020). Small area estimation: Its evolution in five decades. *Statistics in Transition New Series, Special Issue on Statistical Data Integration*, 1–67.
- Jiang, J. and P. Lahiri (2006). Mixed model prediction and small area estimation, editor's invited discussion paper. *Test* 15, 1–96.
- Kim J. and K. Morikawa (2023). An empirical likelihood approach to reduce selection bias in voluntary samples. *Calcutta Statistical Association Bulletin* 35 (to appear).

National Academies of Sciences, E. and Medicine (2017). *Innovations in Federal Statistics: Combining Data Sources While Protecting Privacy*. Washington, DC: The National Academies Press.

Rao, J. N. K. and I. Molina (2015). *Small Area Estimation, 2nd Edition*. Wiley.

Savitsky, T. D., M. R. Williams, J. Gershunskaya, V. Beresovsky, and N. G. Johnson (2022). Methods for combining probability and nonprobability samples under unknown overlaps. <https://doi.org/10.48550/arXiv.2208.14541>.

Sen, A. and P. Lahiri (2023). Estimation of finite population proportions for small areas: a statistical data integration approach. <https://doi.org/10.48550/arXiv.2305.12336>.

Wang, L., R. Valliant, and Y. Li (2021). Adjusted logistic propensity weighting methods for population inference using nonprobability volunteer-based epidemiologic cohorts. *Stat Med.* 40(4), 5237–5250.

Discussion of "Probability vs. Nonprobability Sampling: From the Birth of Survey Sampling to the Present Day" by Graham Kalton

Ralf Münnich¹

Let me first thank Dr. Kalton for his amazing historical review of the development of survey sampling from its origin, contrasting purposive sampling, until now, where some elements of purposive sampling in terms of web or big data seem to supersede the well-elaborated theory of survey statistics. Shall the message be that we do not need any sampling courses at universities anymore, that official statistics should turn to modelling using data with unknown data generating processes, or actually even be substituted by (commercial) *data krakens*? Hardly so! Graham Kalton emphasises a modern thinking about the use of these new data sources which may also have some advantages and he urges future research on data integration methods using (very) different kinds of data while strongly taking quality aspects into account.

Within the last decade, we could observe many new uses of classical data like administrative data and new types of data stemming from internet sources or technical measurement processes such as satellite, mobile phone or scanner data. Already the availability of these new data leads to a huge increase in developing new methodologies and uses. Indeed, official statistics also forced research on new data types, such as scanner data or web-scraped data and others. In Europe, these statistics are often called experimental statistics to emphasise that these statistics cannot (yet) be evaluated using the classical quality concepts, as, e.g. proposed within the European Statistics Code of Practice (<https://ec.europa.eu/eurostat/web/quality/european-quality-standards/european-statistics-code-of-practice>). Some examples can be drawn from https://www.destatis.de/EN/Service/EXDAT/_node.html or <https://ec.europa.eu/eurostat/web/experimental-statistics>.

During the Covid crisis, and especially in light of the political discussion in Germany, however, one could observe little understanding of data quality and statistics. Timeliness – with its urge of getting data and producing statistics immediately – often lead to the use of available (infection) data, which certainly were influenced by unknown biases. The impact of statistics on these available data in terms of evidence-based policy could hardly be understood at the time, but still legal processes like

lockdowns were initiated. To state this message more strongly: whenever a legislation process is involved, and especially so if a direct impact on society is the outcome, we must make sure that high quality requirements on data gathering and statistical methodology are set as well as met. High quality typically cannot be achieved with low costs. England was one of the few very good examples during the pandemic, since they were setting up a special Covid survey to better understand the pandemic and to provide adequate and reliable information.

Certainly, this example already shows some critical aspects in data gathering and data quality. Dr. Kalton was emphasising timeliness and accuracy as very important goals of data quality. For sure, these are of utmost importance! However, in practice, both quality principles suffer from budget constraints and cost controls. This directly leads to two questions: Do modern data help to provide more timely and accurate statistics at lower costs? Is there, in case of conflicts, an *ultimate* quality principle?

The first question is already answered by Dr. Kalton. Of course, modern web or big data can help to gather information quickly. Interesting approaches are of course the use of satellite or scanner data. With electronic cash systems, price changes could be tracked much faster than via the use of survey data. However, one always has to understand the advantages as well as the disadvantages of these data generation processes, and one must be able to measure the quality of the output.

Let me briefly sketch one current German debate which, in my view, perfectly fits into this discussion. In the past years, more and more internet surveys were preferred to data from traditional market and opinion research. This immediately led to a discussion on the quality of the outcomes. And certainly, timeliness, accuracy, and costs played an important role within this discussion. The two major arguments were the following: internet surveys suffer from unknown biases. Classical surveys, in the meantime, have to consider response rates considerably below 20%. Under these conditions, most likely both areas have to consider statistical models with strong assumptions to at least reduce possible biases induced by either web surveys or non-response. In my view, one important question has not been raised yet. What is the aim of the survey?

The ultimate aim that necessitates data collection in the first place is of crucial importance for evaluating the importance of the different quality principles. In case one is interested in getting information on current public opinion, probably timeliness and costs are more important than high accuracy. However, in evidence-based policy making, and especially when information for legislative action is needed, I must stress that accuracy must always be considered to be the major principle. This is even more important when large budgets or financial equalization schemes are involved. Additionally, in these cases one must also be able to measure the quality of the outcome of the statistics. This is still a major drawback of using web or big data. And to stress this point, in legislation processes, I strongly urge to involve independent official statistics with its transparent data production process.

¹ Economics, Economic and Social Statistics Department, Trier University, Germany.

E-mail: muennich@uni-trier.de. ORCID: <https://orcid.org/0000-0001-8285-5667>.

With this discussion, I do not want to be misunderstood. Modern data and modern statistical methods are important. And the direction of research, as Dr. Kalton pointed out, will be complex modelling and data integration. Also administrative, register, and related data are important and can provide very good information. However, with all these data, we always have to understand their quality and we should be able to measure the quality of the resulting statistics. Especially in the context of big data, quality measurement may have to be enhanced (cf. Münnich and Articus, 2022, and the citations therein).

Sampling itself may also follow new directions. Classical sampling optimization may be adequately applied in more special cases that allow focusing on specific goals, e.g. the design optimization in the German Censuses 2011 and 2022 (see Münnich et al., 2012, and Burgard, Münnich, and Rupp, 2020). However, likely robustness of methods against assumptions has to be incorporated in design optimization. On the other hand, data integration, multi-source environments, geo-spatial modelling, small area estimation and other modern methods may yield new ideas and directions in sampling theory and application. One example may be sampling from big data sources to reduce complexity.

Despite the mentioned new directions, many ideas have been well-known for a long time. In data analytics, we differentiate between descriptive, predictive, and prescriptive aims. Data that were gathered to describe a state of a system cannot be used to analyse interventions on the system. Indeed, we need the right data and not just merely available data. In conclusion, the exact purpose of the statistics under consideration plays an extremely important role for the selection of data and the priority of the different quality principles.

References

- Burgard, J. P., Münnich, R., & Rupp, M., (2020). Qualitätszielfunktionen für stark variierende Gemeindegrößen im Zensus 2021. *AStA Wirtschafts- und Sozialstatistisches Archiv*, 14(1), pp. 5–65. With discussion.
- Münnich, R., Articus, C., (2022): Big Data und Qualität – ist viel gleich gut? Pp 85–101. In: Wawrzyniak, B., Herter, M. (Ed.): *Neue Dimensionen in Data Science*. Wichmann.
- Münnich, R., Gabler, S., Ganninger, M., Burgard, J. P.; Kolb, J.-P., (2012): Stichprobenoptimierung und Schätzung im Zensus 2011. *Destatis: Wiesbaden, Statistik und Wissenschaft*, Vol. 21, https://www.zensus2011.de/SharedDocs/Downloads/DE/Publikationen/Monografien_Archiv/2012_07_Destatis_Stochprobenoptimierung_und_Schaetzung_im_Zensus.pdf?__blob=publicationFile&v=12

Rejoinder

Graham Kalton¹

I should like to thank the discussants for their kind remarks, for their valuable comments on the present state and future directions of the field, and for the many references they cite. Since I have no disagreements with them, I will confine my rejoinder to a few issues that their contributions have surfaced for me.

I will start by rectifying an oversight in my treatment of the early history of survey research and survey sampling: Carl-Erik Särndal has reminded me of the major developments that occurred in Russia during the early years. The impetus for these developments was the need for local self-government units known as *zemstva* to collect data about their populations for administrative purposes. Initially such data were collected with 100% enumerations, but around 1875 sample surveys were introduced for cost savings. The survey procedures were coordinated across *zemstva* and a number of sampling methods were evaluated with input from theoretical statisticians. These statisticians made a number of important contributions, including an impressive early text (1924) entitled *The Foundations of the Theory of the Sampling Method* by A. G. Kowalsky. Although Russian statisticians were at the frontiers of developments in survey sampling until the late 1920's, their contributions were not fully recognized outside Russia. For example, Tschuprow (1923) and Kowalsky in his 1924 text both derived the optimum allocation formula for stratified sampling a decade before Neyman did so in his famous 1934 paper (after learning of Tschuprow's paper, Neyman (1952) recognized Tschuprow's priority for the results). Mespoulet (2002), Zarkovic (1956), Zarkovic (1962), and Seneta (1985) provide further details about early survey research and research on survey sampling in Russia.

Danny Pfeffermann has pointed out that probability samples are almost never representative because of nonresponse—and I would add noncoverage—that is not missing completely at random (NMAR or MCAR). Moreover, I do not think the nonresponse should be viewed as missing at random (MAR), that is MCAR after conditioning on known covariates. Using standard weighting adjustments based on known covariates will not make the sample representative. My favorite quotation from George Box is “Essentially, all models are wrong, but some are useful.” Nonresponse adjustments should be viewed from this perspective as useful but not perfect. Another George Box quotation: “Statisticians, like artists, have the bad habit of falling in love with

their models.” But there is a difference: artists have artistic license to paint over a model's blemishes whereas statisticians should attempt to identify and repair the blemishes.

Risto Lehtonen points out the considerable attractions of population registers, as have existed for some time in several Scandinavian countries and are in development elsewhere. Such registers can be viewed as surveys with 100% samples, and the quality of their data should be assessed accordingly: What is their actual coverage? How up-to date are they? How accurate are the data they contain?

Risto's discussion of population registers also reminded me of a point that I should have addressed more fully: there is a wide variation in the data infrastructure for social research across countries. For example, most developing countries are not in a position to use administrative records or the internet. They rely on probability sample surveys to satisfy their data needs. Fortunately, they have not yet experienced the severe declines in response rates that are so harmful to surveys in most high-income countries.

Julie Gershunskaya and Partha Lahiri address two important current areas of research. One is the research on how to employ a probability sample to reduce the bias in estimates from a nonprobability sample, making use of auxiliary variables collected in both samples. The auxiliary variables aim to capture the key variables that are predictors of membership in the nonprobability sample. Challenges to be addressed with this approach include identifying the key variables; dealing with the fact that some response categories that occur frequently in the probability sample are very sparsely represented in the nonprobability sample; and concerns about the equivalence of the responses to the key variables obtained in the two samples that use different modes of collection. The results from this approach should be viewed with caution. However, recalling George Box's quotation above, imperfect models can be useful. Julie and Partha rightly say that the aim of these models is to reduce, not eliminate, bias. The question to be asked is how to assess whether the models have reduced bias to an acceptable level.

The second area that Julie and Partha address is small area estimation. I should have written more about this methodology whose use has now become so widespread. My first practical exposure to small area estimation occurred in the late 1990's, when I chaired a National Academy of Sciences' panel that was asked to advise about the quality of the small area estimates of the numbers of poor school-aged children that were being developed in the U.S. Census Bureau's Small Area Income and Poverty Estimates (SAIPE) program. The central issue was whether the estimates, which were produced for 3,000 counties and 14,000 school districts, were appropriate and sufficiently reliable to be used in allocating very large sums of money directly to school districts. At that time, this was a novel application of small area estimates, and subject to considerable questioning. After extensive evaluation of the area level models by both the Panel and the Census Bureau (Citro and Kalton, 2000), the Panel concluded that the small area estimates were “fit for use” for the purpose of this fund allocation, despite a recognition of substantial errors in the individual estimates. The Panel was influenced by the fact that the legislation stipulated that the funds should be distributed directly to the school districts and that, even though the small area estimates were not ideal, they were the best available. I was persuaded by my experience on the Panel that, with strong predictors and careful model

¹ Joint Program in Survey Methodology, University of Maryland, College Park, MD, USA.

E-mail: gkalton@gmail.com. ORCID: <https://orcid.org/0000-0002-9685-2616>.

development and testing, small area estimation methods have an important role to play in responding to policy makers' increasing demands for local area estimates.

Ralf Münnich emphasizes the importance of assessing the overall quality of statistical estimates in the light of the uses of the estimates. As he notes, timeliness is often in conflict with accuracy. In some situations, timeliness may be paramount, and accuracy may suffer. However, one must guard against the risk that accuracy is so low that the resulting estimates are misleading. Estimates based on big data sources or even large surveys conducted with an overriding emphasis on speed may, because of their sample sizes, appear to be well-grounded but that may well be illusory.

It is often argued that although individual estimates may be subject to serious biases, these biases will cancel out for differences between estimates, either between subgroups of the sample or across time. While the underlying model for that argument often appears reasonable, the assumptions underpinning it need to be carefully assessed in each case.

Ralf also points out the importance of cost constraints. When the cost constraints severely limit a study to a very small sample size, it may be preferable to forego the extra costs involved in selecting and fielding a probability sample, in favor of a quasi-probability sample or a nonprobability sample design. As Kish (1965, p. 29) notes: "Probability sampling is not a dogma, but a strategy, especially for large numbers."

Finally, Ralf and other discussants have pointed out the attractions of data integration. I also see these attractions, but I think that the challenges of mode effects arising from different data sources should not be underestimated.

In conclusion, I congratulate Statistics in Transition on celebrating its 30th anniversary. It plays a distinct and important role among statistics journals. With the major changes in statistical methodology taking place in official statistics and in social research, it has a bright future for the contributions it can make.

References

- Citro, C. F., Kalton, G. Eds., (2000). *Small-Area Estimates of School-Age Children in Poverty: Evaluation of Current Methodology*. National Academy Press, Washington D.C.
- Kish, L., (1965). *Survey Sampling*. Wiley, New York.
- Mespoulet, M., (2002). From typical areas to random sampling: sampling methods in Russia from 1875 to 1930. *Science in Context*, 15(3), pp. 411–425.
- Neyman, J., (1952). Recognition of priority. *Journal of the Royal Statistical Society, A*, 115(4), 602.
- Seneta, E., (1985). A sketch of the history of survey sampling in Russia. *Journal of the Royal Statistical Society, A*, 148(2), pp. 118–125.
- Tschuprow, A. A., (1923). On the mathematical expectation of the moments of frequency distributions in the case of correlated observations. *Metron*, 2(4), pp. 646–683.
- Zarkovic, S. S., (1956). Note on the history of sampling methods in Russia. *Journal of the Royal Statistical Society, A*, 119(3), pp. 336–338.
- Zarkovic, S. S., (1962). A supplement to "Note on the history of sampling methods in Russia". *Journal of the Royal Statistical Society, A*, 125(4), pp. 580–582.