

Statistical challenges of administrative and transaction data

David J. Hand

Imperial College London and Winton Capital Management, London, UK

[Read before The Royal Statistical Society at a meeting organized by the Official Section on Wednesday, November 15th, 2011, Mr M. Baxter in the Chair]

Summary. Administrative data are becoming increasingly important. They are typically the side effect of some operational exercise and are often seen as having significant advantages over alternative sources of data. Although it is true that such data have merits, statisticians should approach the analysis of such data with the same cautious and critical eye as they approach the analysis of data from any other source. The paper identifies some statistical challenges, with the aim of stimulating debate about and improving the analysis of administrative data, and encouraging methodology researchers to explore some of the important statistical problems which arise with such data.

Keywords: 'Big data'; Data quality; Management data; Operational data; Repurposed data

1. Introduction

Administrative data are data generated during the course of some operation, and then retained in a database. They are becoming increasingly important as the potential for discovery from such sources of data is being recognized and as alternative sources of data become more costly or difficult to use (e.g. because of declining response rates in surveys). In the main, this means that the analysis of administrative data is *secondary*—the data are being *repurposed*—although, as explained below, this is not always so. The existence of large, often administrative, data sets, offering potential for secondary analysis, was one of the primary drivers behind the development of data mining technology (Hand *et al.*, 2000) as well as the modern rise of interest in 'big data'. But the analysis of administrative data presents new statistical challenges. This can be seen by a cursory examination of the examples in most basic statistics texts, which will almost all involve 'random samples': administrative data are, by definition, typically not random samples. The aim of this paper is to explore these statistical challenges and to stimulate discussion. The hope is that it will help to focus attention on what is needed for valid and accurate analysis of administrative data. The need is illustrated by the comment made by Wallgren and Wallgren (2014), page 3, on the closely related topic of analysing data from statistical registers:

'Although register-based statistics are a common form of statistics used for official statistics and business reports, no well-established theory in the field exists. There are no recognized terms or principles, which makes the development of register-based statistics and register-statistical methodology all the more difficult. As a consequence, *ad hoc methods are used instead of methods based on a generally accepted theory*'.

Address for correspondence: David J. Hand, Department of Mathematics, Huxley Building, Imperial College London, South Kensington Campus, London, SW7 2AZ, UK.
E-mail: d.j.hand@imperial.ac.uk

1 It is hoped that this paper will serve as a framework to stimulate discussion about what ‘generally
2 accepted theory’ might be taught for the analysis of administrative data.

3 There are many definitions of statistics. This is because the discipline has various aspects,
4 including the study of methods for collecting, presenting, interpreting and analysing data, but
5 also because it involves expertise in coping with uncertainty and chance. My own definition
6 (Hand, 2008) tries to capture this diversity: *statistics is the technology of extracting meaning*
7 *from data and of handling uncertainty*.

8 There are fewer definitions of administrative data. The Organisation for Economic Co-
9 operation and Development (Organisation for Economic Co-operation and Development, 2016)
10 defined administrative data as having the following features:

- 11 (a) the agent that supplies the data to the statistical agency and the unit to which the data
12 relate are *usually* different, in contrast with most statistical surveys;
- 13 (b) the data were originally collected for a definite non-statistical purpose that might affect
14 the treatment of the source unit;
- 15 (c) complete coverage of the target population is the aim;
- 16 (d) control of the methods by which the administrative data are collected and processed rests
17 with the administrative agency.

18 The definition continues by saying that

19 ‘In most cases it is normal to accept (and expect) that the administrative agency will be a government
20 unit that is responsible for implementing an administrative regulation’.

21 That leads to a rather narrower definition than is taken in this paper. For example, it excludes
22 corporate use of administrative data, describing the workforce, products, processes, and so on,
23 as well as narrowly restricting the uses to which such data are put. Instead, although accepting
24 that the features described above do characterize administrative data, I shall follow Nordbotten
25 (2010) and simply distinguish between statistical data and administrative data. Statistical data
26 are collected primarily for statistical purposes—e.g. to summarize in order to shed light on
27 the system generating the data, or to make predictions. In contrast, administrative data are
28 initially collected for some administrative purpose—to run an organization, such as a company,
29 government, charity, school, hospital, and so on. Running the organization might require on-
30 going operational analysis of the data but, once collected and stored, the data can later be
31 analysed to shed light on what has happened, to help to predict what might happen in the future,
32 and to evaluate systems and their performance, i.e. the data can later be subjected to statistical
33 analysis. Often statistical data consist of mere samples from the universe of possible values
34 which could have been obtained, and these will have been collected by surveys or experiments
35 for example. In contrast, administrative data will ideally consist of data on all of the cases, records
36 or transactions in some population. This leads to something of a conceptual distinction: sample
37 data are used to obtain *estimates* of a population *parameter*. In contrast, administrative data
38 are *summarized* to obtain a descriptive *feature* of the population.

39 Transaction data are an important kind of administrative data concerned with *events*, typically
40 with sequences of events. Usually the prime operational purpose of collecting the data is to
41 inform the transaction (e.g. to decide how much to charge a supermarket customer or to decide
42 how much tax someone should pay), but once collected the data can be retained in a database
43 and analysed to improve understanding of the organization’s operations.

44 I used the word ‘operational’ above. Occasionally we see the terms ‘operational data’, ‘man-
45 agement data’ or ‘management information’ used to describe data collected and analysed to
46 guide the operation of a system. It is clear from the above discussion that the data, once
47
48

collected and placed in a database, are no different from administrative data. What differs is the way in which the data are being used—from immediate decisions to more considered analysis with longer-term implications. Operational data become administrative data when they are stored and used for some purpose beyond the day-to-day operations of the organization. In a sense, then, administrative data are *data exhaust*: that which is left over after the organizational machinery has used the data to drive itself forward.

Incidentally, in this paper, I shall use the term ‘survey’ to refer to data collection by sample survey, so that it is contrasted with administrative data, collected as a side effect of an operational activity. This is different from, for example, Statistics Canada (2009), page 11, which used the term survey ‘generically to cover any activity that collects or acquires statistical data’, including the collection of data from administrative records.

At first glance, though we shall see that appearances can be deceptive, administrative data appear to have several advantages compared with statistical data.

- (a) Since the data have already been collected, no additional cost appears to be incurred in collecting them.
- (b) In a sense, we might reasonably expect that ‘all’ the data are available. After all, a company will certainly process and can retain details of all its transactions.
- (c) The data might be of high quality, since the effectiveness of the operation of the organization depends on this.
- (d) The stored data will certainly be timely and might be regarded as as up to date as it is possible to achieve, since they describe the organization as it is, or at least as it was when the last change was made. This advantage is strikingly illustrated in the use of administrative data to derive estimates of population attributes at times that are intermediate between decadal censuses, and in essentially real-time estimates of price inflation.
- (e) In a real sense administrative data often tell us what people *are* and what they *do*, not what they *say* they are and what they *claim* to do. We might thus argue that such data get us closer to social reality than survey data.
- (f) Administrative data may provide tighter definitions than alternative sources of data. Wallgren and Wallgren (2014), page 33, gave examples of data about income and children in families. Where the time restrictions on eliciting responses to a survey might mean one must simply ask ‘what is your yearly income before tax?’, administrative data might, depending on the source of the data, specify whether this means ‘disposable income, taxable income, earned income or income including unearned income...’.

Unfortunately, although all of those advantages of administrative data might apply in an ideal world, in practice things are typically not so straightforward. Regarding (a), effort will normally be required to extract the data, to clean them and possibly to link them to other data sets. Moreover, although data may be free for the organization which collected them, other organizations which wish to use these data may have to pay—and the cost must be balanced against that of data from alternative sources, such as surveys or administrative data collected by other organizations. Regarding (b), data will usually enter a database via a complex social process—the sample of records in a database may not be representative of the population to which one wishes to make an inference. An operational database might not have a form which is convenient for statistical analysis exercises. In particular, different parts of an organization might use different database systems—indeed, there is a large amount of current activity as organizations seek to put all of their data into a single data repository (a data warehouse, for example). The notion of ‘data = all’ is discussed in Section 3. Regarding (c), although sampling variation issues may not apply, administrative data will have other sources of uncertainty, and

1 unfortunately these may be various and diverse, and not susceptible to resolution machinery
2 as mathematically elegant or unified as sampling theory. More generally, although in principle
3 one might expect administrative data to be of high quality, in practice all data sets, perhaps
4 especially those involving human beings, are susceptible to quality issues. A particular issue
5 with administrative data sets arises from the very fact that they were not deliberately collected
6 to answer the later statistical question being addressed. This means that the data may not
7 be ideally suited to answer the question; there is often a compromise between cost and rele-
8 vance. For example, a costly survey could be designed to answer the specific question whereas
9 ‘free’ administrative data might be only roughly suitable (but see principle 8, clause 8.1, of the
10 European statistics code of practice (Eurostat, 2011), which says ‘When European Statistics are
11 based on administrative data, the definitions and concepts used for administrative purposes
12 are a good approximation to those required for statistical purposes’). Moreover, the definitions
13 that are involved in administrative data are subject to changes for operational purposes which
14 might impact the research questions that can reasonably be asked. It follows that time series of
15 administrative data might exhibit discontinuities: an administrative database containing details
16 of unemployment benefits might appear to be ideal to address questions of unemployment rates
17 but, if the definition that is used in assessing benefits changed over time, then it might limit
18 what can be done. Regarding (d), although administrative data may be instantly available to the
19 organization collecting them, it may not be so to other organizations which wish to use them.
20 Regarding (e), there are important kinds of administrative data which are not automatically
21 accrued through some transaction but are specifically sought and reported—e.g. income tax
22 data. And, finally, (f) is not universally true: credit card transaction data contain considerable
23 detail of the nature of the item purchased, but not necessarily to a level that is adequate for all
24 potential analyses.

25 An example of the relative merits of administrative and statistical data is given by crime
26 statistics in England and Wales. There are two main sources of such data, the crime survey
27 for England and Wales and police-recorded crime (PRC) (Office for National Statistics, 2016a).
28 These can sometimes show trends going in opposite directions. The reason is that definitions—of
29 crimes, of victims, of what data are collected—differ. With the administrative data (the police-
30 recorded crime data) we must make use of the data that we have, whereas with survey data (the
31 crime survey for England and Wales) we can decide what data are best collected to answer the
32 questions we want to ask. Indeed, there has been extensive research on how best to formulate
33 survey questions to elicit the information that is sought (see, for example, Presser *et al.*, (2004)
34 and Fowler (1995)).

35 There are also other issues which arise with administrative data which are slightly different
36 from those arising with survey data. An obvious problem relates to privacy and confidentiality—
37 discussed in more detail below. Since survey data are collected purely for statistical analysis, data
38 released for analysis would not normally retain identifying information: apart from its use in
39 ensuring consistency and representativeness, the identifying information is not relevant to the
40 use of the data. In contrast such information is central to the initial (operational) purpose for
41 which administrative data are collected. The (hoped for) comprehensiveness of administrative
42 data increases the risk of reidentification—and perhaps the public concern.

43 I could go on, but it is clear that, although administrative data have merits, the statistician
44 should approach such data sets with the same critical eye that they approach any other data
45 set.

46 It is worth noting that it is sometimes useful to distinguish between two kinds of administrative
47 data. The first kind is that which is *necessarily* collected during the course of some operation.
48 Credit card transaction data, for example, necessarily involve the recording of the amount spent,

1 the currency, and the business where the transaction occurs, since these items of information are
2 needed to run the credit card operation. The second kind is additional information which is not
3 needed for an operation, but which is helpful for other reasons, and which is collected during the
4 administrative process. The age and gender of a customer might fall into this category: a product
5 might be bought by anyone, but it could be useful to analyse the customer's details later to inform
6 new marketing strategies. In some sense this second kind of data lies between administrative
7 and statistical: they are collected for statistical rather than operational purposes, but they are
8 collected during and as part of the administrative process. There is an important lesson to be
9 taken from this: benefits can be gained from involving the statisticians and data analysts in the
10 data collection stage. This is not a new lesson: we recall Ronald Fisher's comment that

11 'to call in the statistician after the experiment is done may be no more than asking him to perform a
12 post-mortem examination: he may be able to say what the experiment died of.'

13
14 This last point leads us into the modern world of so-called 'big data'. The term has no
15 universally accepted definition, but we might define it as the result of some automatic data
16 collection system. Indeed, I have argued elsewhere that the data revolution is not so much a
17 consequence of the size of modern data sets and the ability to store them (big data) but rather
18 of the fact that data are nowadays largely collected automatically without requiring explicit
19 human effort. Examples of automatically collected data are everywhere and include personal
20 health data collected by wrist monitors, automated monitoring of tickets as people travel through
21 a rail network, telemetry of engine functioning and recording of metadata of phone calls. Data
22 arising from the so-called '*Internet of things*' would clearly be of this type.

23 Given that so many official and economic statistics are based on administrative data, or on a
24 combination of administrative and survey data, we might have expected there to be a substantial
25 literature in the leading methodological statistical journals describing the statistical challenges
26 and how to overcome them. This appears not to be so, with such journals carrying relatively
27 few papers on the statistical challenges of administrative data (being mostly focused on the con-
28 sequences of sampling theory). For example, a search of the *Journal of the American Statistical*
29 *Association* for occurrences of the phrase 'administrative data' yielded 44 results. Obviously the
30 *Journal of Official Statistics* is an exception, although even there most of the papers including
31 the phrase 'administrative data' in the title are concerned with particular applications. More
32 generally, papers on the topic seem to be widely scattered and often appear in the proceedings of
33 conferences and workshops, or perhaps as reports of official exercises (e.g. from official statistics
34 offices). Given this wide scattering, it is certain that important contributions have been omitted
35 from this paper, and I welcome attention being drawn to them in the discussion.

36 One of the best discussions is the excellent and comprehensive introduction to register-based
37 statistics by Wallgren and Wallgren (2014). In one sense, this has a much wider scope than
38 this paper, including discussion of register structures and the creation of registers, but in an-
39 other sense it is narrower, being focused on official statistics and not including, for example,
40 commercial or engineering applications of administrative data.

41 The series of conferences on 'New techniques and technologies for statistics, (e.g. New Tech-
42 niques and Technologies for Statistics (2013, 2015)), organized by the European Commission,
43 often have items touching on administrative data challenges: the phrase 'administrative data'
44 appeared in the 2013 conference proceedings 220 times. To give the flavour of the breadth of
45 topics that were covered, these proceedings included papers on state space models (Horn and
46 Czaplewski, 2013), structural equation models (Scholtus and Bakker, 2013), business statistics
47 and other topics. Romanov and Gubman (2013) explored regression to the mean in survey re-
48 sponses to questions on income by using administrative (tax) data, and Lewis and Woods (2013)

described some of the issues which must be tackled when using administrative data in the form of value-added tax and company accounts data, as the basis for business statistics. As well as problems of matching and cleaning administrative and survey data, they also discussed differences of timeliness and periodicity. Kloek and Vâju (2013) discussed integration of administrative data with other kinds of data. They characterized five different kinds of use: direct use at microlevel, use as auxiliary information at microlevel, use as auxiliary information at aggregate level, use as a source for the population frame and use as circumstantial evidence. They also explored the distinction between administrative data for business and those for households. This, of course, reflects the general point that data describing different kinds of entities might have different characteristics (e.g. more pronounced skewness for some variables for business data compared with household data). Četković *et al.* (2013) have provided an elaborate example: the Austrian register-based census, involving seven base registers and several comparison registers which are provided with data from 35 data holders. They characterized data quality in terms of several ‘hyperdimensions’ described by Berka *et al.* (2012) (see below). Antoni (2013) linked survey and administrative employment data.

Other sources which have relevant materials include the following:

- (a) the United Nations Economic Commission for Europe data collection workshops (see, for example, United Nations Economic Commission for Europe (2012), on new frontiers in data collection);
- (b) the ‘ESS vision 2020’, which includes discussion of administrative data sources and challenges (European Statistical System, 2020) (their ‘Administrative data sources’ project is exploring how administrative data may be used to increase data availability and to reduce costs (European Statistical System Admin, 2015));
- (c) ESSNet has current and previous projects on administrative data topics (ESSNet, 2017) (see, for example, ESSNet Admin Data Workshop (2013));
- (d) the US *Review of Administrative Data Sources* (Ruggles, 2015);
- (e) the Statistics New Zealand ‘Guide to reporting on administrative data quality’ (Statistics New Zealand, 2016);
- (f) the use of administrative data at Statistics Canada (Statistics Canada, 2015);
- (g) the administrative data quality assurance documents produced by the UK Statistics Authority (UK Statistics Authority, 2014, 2015);
- (h) the checklist of quality of statistical outputs in van Nederpelt (2009);
- (i) ‘Pros and cons for using administrative records in statistical bureaus’, from the Israel Central Bureau of Statistics (2007);
- (j) the Organisation for Economic Co-operation and Development compilation ‘Short-term economic statistics (STES) administrative data: two frameworks of papers’ (Organisation for Economic Co-operation and Development, 2016) is a particularly valuable source of examples of the use and challenges of administrative data, albeit focused mainly on economic uses.

The structure of this paper is as follows. Section 2 describes a fundamental problem that is relevant to all data analysis, no matter what the source of the data, namely data quality. But the challenges—and even the recognition that there are challenges—that are presented by administrative data differ from those presented by other sources. We look at some of these challenges and how they differ from those of other types of data.

Section 3 addresses the notion that one might have ‘all’ of the data. This is typically regarded as one of the particular merits of administrative data but, as we show, it is all too often an unjustified assumption.

Section 4 explores the fact that administrative data are collected for operational purposes, and not with specific research questions in mind. The consequence is that the data may be far from ideal for addressing those questions.

Sections 5, 6 and 7 look at deeper issues where the nature of administrative data impacts other aspects of analysis, including efforts to identify causation, merging data from multiple sources and the thorny issues of confidentiality, privacy and anonymization of records.

Section 8 draws some conclusions.

2. Data quality

The value of administrative data for producing official statistics has attracted increasing attention recently. In large part this is in the hope that they can replace more conventional survey data, motivated on the one hand by a worldwide decrease in survey response rates, and on the other by a perceived lower cost in using administrative data, since they have already been collected. However, as the UK Statistics Authority put it,

‘we have been surprised by the general assumption made by many statistical producers that administrative data can be relied upon with little challenge, and, unlike survey-based data, are not subject to any uncertainties’

(UK Statistics Authority, 2014). Because of this, the UK Statistics Authority has produced a report on quality issues in administrative data, summarizing the lessons learnt from a review of users of administrative data for statistical purposes and describing a toolkit to monitor data quality in this context (UK Statistics Authority, 2014, 2015).

Other explorations of the quality aspect of administrative data include the model that Daas *et al.* (2008) have developed for Statistics Netherlands. Noting that a key issue with administrative data is that the source of the data is typically some other body, Daas *et al.* (2008) pointed out that the collection and maintenance are not within the control of the analyst: when data are collected by bodies other than those undertaking the analysis, issues of data provenance and curation are critical. In their review of earlier work on administrative data quality, Daas *et al.* (2008) observed that different researchers have identified

‘a remarkable difference between the number and types of quality groups or dimensions identified for the statistical quality aspects of administrative data’.

They attributed this partly to the complexity of the problem and partly to the fact that different researchers had different perspectives on the topic. Their paper is then an attempt to integrate the various views into a single framework. Their conclusions include the observation that administrative data quality is a multi-dimensional issue, with a hierarchy of dimensions (Karr *et al.*, 2006).

Other work (e.g. Eurostat (2003)) has explored the potential uses of administrative data. This is an important point when attempting to evaluate data quality, as data may be ‘good’ for one purpose but ‘bad’ for another: quality is not a property of the data set itself, but of the interaction between the data set and the use to which it is put. And yet other, more general, work on data quality, especially in the context of official statistics, inevitably touches on administrative sources (e.g. van Nederpelt (2009), Memobust Handbook (2014) and Statistics Netherlands (2014)). Once again, we stress that work on the quality of administrative data has appeared in diverse publications, from a wide range of sources.

The common misconception that quality issues are less important for administrative data than they are for survey data seems to be based chiefly on the belief that data that have initially been acquired for operational purposes must necessarily be both complete and error free, whereas

1 survey data will be based on a mere sample from the population being studied, so the results will
2 vary between possible samples. The fact is, however, that administrative data may be neither
3 complete nor error free. As far as ‘complete’ is concerned, incompleteness can manifest either in
4 the form of partial records—records in which some of the fields are missing—or in the form of
5 entire records missing, so that the data set does not in fact cover the entire population. And, as
6 far as ‘error free’ is concerned, errors can arise in an unlimited number of ways. To paraphrase
7 Leo Tolstoy: ‘A perfect data set is perfect in only one way; each imperfect data set is imperfect
8 in its own way’. This means that we can never be sure that all the errors have been detected. The
9 problem is analogous to that of testing random-number generators: we can look for particular
10 kinds of departures from randomness, but there will always be kinds that we have not thought
11 of. Unfortunately, one of the lessons that we have learnt from data mining practice over the past
12 20 years is that most of the unusual structures in large data sets arise from data errors, rather
13 than anything of intrinsic interest. We should be suspicious of any data set (large or small)
14 which appears perfect. A standard check that I carry out is to ask those providing the data what
15 they have done about missing values. Often this has resulted in surprising responses which the
16 researchers would not have thought to mention if the question had not been explicitly asked. For
17 example, it is not uncommon for researchers to have removed any incomplete records from the
18 data set, introducing unknown selection bias. Caruana *et al.* (2015) described the development
19 of a machine learning diagnostic system based on hospital administrative data which classified
20 high risk asthma patients as low risk because such cases had been excluded from the training
21 data.

22 Although the technology of data editing and imputation has been substantially developed,
23 with entire books being written about it (e.g. de Waal *et al.* (2011)), it is not the case that detected
24 errors can necessarily usefully be corrected. This means that commercial tools for detecting and
25 correcting data errors are unlikely to be 100% effective, whatever they may claim.

26 Statisticians know very well that it is common for the major part of their time on a project to
27 be spent cleaning data before actual statistical analysis. This is all very well when the data set is of
28 moderate size, but it becomes more of a problem when the data set is massive—as is increasingly
29 the case in the ‘big data’ world, and is particularly the case with administrative data and data
30 which are captured automatically. Especially in such contexts, the computer is a necessary
31 intermediary between the analyst and the data, with consequent risks of missing important
32 shortcomings of the data—and indeed, even creating extra errors during an automatic data
33 cleaning process. For example, rule-based correction mechanisms can distort perfectly good,
34 though unusual, data values, and an unfortunately all-too-common strategy for coping with
35 missing values is to substitute the mean of the observed values (so leading to an underestimate
36 of variance).

37 Familiarity with the fact that data are often not of the highest quality has led to the devel-
38 opment of relevant statistical methods and tools, such as detection methods based on integrity
39 checks and on statistical properties (e.g. comparing distributions with expected distributions in
40 electoral data, or using the Benford distribution for leading digits); see, for example, Hellerstein
41 (2008) and de Jonge and van der Loo (2013). However, this emphasis has often not been matched
42 within the realm of machine learning, which places more emphasis on the final modelling stage
43 of data analysis. This can be unfortunate: feed data into an algorithm and a number will emerge,
44 whether or not it makes sense. However, even within the statistical community, most teaching
45 implicitly assumes perfect data. This is entirely reasonable: if one is aiming to teach the basic
46 concepts of regression, one does not want to spend time pointing out the consequences of miss-
47 ing data, digit heaping or digit transposition. Nonetheless, students do need to understand the
48 reality of data analysis. This leads to our first challenge.

1 *Challenge 1.* Statistics teaching should cover data quality issues.

2
3 Even if data may depart from perfect quality in an unlimited number of ways, it is important to
4 characterize as many ways as possible, and Kim *et al.* (2003) have produced a general ‘taxonomy
5 of dirty data’. They characterized data as dirty ‘if the user or application ends up with a wrong
6 result or is not able to derive a result due to certain inherent problems with the data’, and they
7 identified various possible causes of the problem, including data entry errors data update errors,
8 data transmission errors and also bugs in a data processing system. Particular applications are
9 likely to have their own characteristic types of error, and it seems likely that an 80/20 rule
10 will often apply, with a large proportion of errors being of just a few types, so that relatively
11 little effort will lead to substantial initial improvement in overall quality. An illustration of
12 this was given by Lewis and Woods (2013), who identified the main causes of error in value-
13 added tax data to be just four types: scanning errors, unit errors, incorrect quarterly data, and
14 errors in individual responses. De Veaux and Hand (2005) gave examples of data errors and
15 their consequences, and national statistical institutes often define several dimensions of quality,
16 including accuracy, relevance, timeliness, existence, coherence, completeness, accessibility and
17 security (see, for example, Eurostat (2000) and the archives of other national statistical institutes;
18 Meader and Tily (2008) and Biemer *et al.* (2014)), though these will affect administrative data
19 in varying degrees.

20 The ‘relevance’ aspect in the national statistical institutes list is more subtle than simply find-
21 ing a mistake in the data. Even perfectly accurate data may be useless for answering a particular
22 research question if the data have not been collected with the research question in mind—as is
23 typical with administrative data. Clearly we can try to ease that difficulty if we know beforehand
24 what questions are likely to occur, but even then difficulties can arise. For example, in a project
25 aimed at constructing a scorecard to predict likely default on bank loans, one of the (relatively
26 highly predictive) variables was ‘is the applicant a home owner or renter?’. This was adminis-
27 trative data of the second kind mentioned above—the question was not relevant to everyday
28 operations. But as a consequence of this the people tasked with recording the data failed to see its
29 importance, with the result that they initially recorded it for only a small percentage of customers.

30 If administrative data are subject to restrictions arising from operational imperatives, they
31 are also subject to possible constraints from the opposite direction: administrative data are
32 often communicated, compared and aggregated across bodies collecting the data. For example,
33 national statistics for US states and countries within the European Union will be aggregated
34 to produce Federal statistics and European Union statistics respectively. The need to do this
35 imposes constraints on what must be collected and on its format, with particular standards
36 requiring particular structures, formats and protocols, as well as content.

37 As mentioned above, administrative data are also susceptible to changes of definition, which
38 can adversely affect things like time series, rendering them non-comparable over time. Since
39 much administrative data, especially those concerned with government and public policy, are
40 subject to regulation and legislation, changes in laws can have an unfortunate effect, at least
41 from the perspective of the statistician hoping to use the data to make inferences. Changes in
42 what data can be stored, or the characteristics which are allowed to be used in statistical models,
43 can mean that earlier models become unusable.

44 Data can be incorrectly entered, even for operational purposes. We have all heard of ‘fat
45 finger’ errors leading to mistaken financial transactions. Other classic examples include things
46 like weights of 1 lb being miscoded as 11 lb, data being entered in incorrect columns, abbrevi-
47 ations leading to confusion (e.g. MS for Microsoft or Morgan Stanley), incorrect time stamps
48 due to clocks being misset, mistakes in the use of measurement units, simple misspellings and

instrument failures not being detected (leading to, for example, an unnoticed stream of 0-values). The list of examples is endless. Kruskal (1981) observed that

‘A reasonably perceptive person, with some common sense and a head for figures, can sit down with almost any structured and substantial data set or statistical compilation and find strange-looking numbers in less than an hour’.

However, even data which are entered correctly and unambiguously for operational purposes can lead to errors when subjected to statistical analysis. Alternative, equally legitimate spellings or identifiers (e.g. David Hand, David J. Hand and D. J. Hand) may not be recognized as equivalent in a subsequent analysis unless they have been explicitly characterized as so. Conversely, identical entries might refer to different objects (e.g. father and son with the same name). Missing values for age coded as 999 can be analysed as legitimate ages, with obvious adverse consequences. Although clearly this should be flagged in the metadata, we note, again, that large data sets necessarily involve an opacity that does not affect small data sets, in that the computer is a necessary intermediary between the data and the analyst. Mistakes and ambiguities can slip through.

The argument has been made that errors in data will often affect only a very small part of the data, and so will, for example, have no significant effect on large-scale conclusions. Although this may be true, large-scale conclusions are typically not the only ones which will be drawn from administrative data. One of the particular strengths of such data is that they are also used for small-scale investigations—to explore subgroups or for small area statistics, for example. In such cases, errors in only a few records can have important consequences.

Quality issues may also arise when data sets, of adequate quality in themselves, are merged. Take, for example, time series which are out of phase, or have different frequencies of publication, or publish on different dates or, even worse, are irregular.

These considerations lead to several challenges.

Challenge 2. Develop detectors for particular quality issues.

Challenge 3. Construct quality metrics and quality scorecards for data sets

Challenge 4. Audit data sources for quality.

Challenge 5. Be aware of time series discontinuities arising from changing definitions.

Challenge 6. Evaluate the impact of data quality on statistical conclusions.

3. ‘Data=all’?

The phrase ‘data = all’ is sometimes encountered in the context of administrative data. This is intended to convey the notion that the data are not merely a sample from the population of objects but are its entirety: all credit card transactions, all supermarket purchases, all tax records, and so on. The implication is that having data describing the entire population means that we need not worry about sampling errors or errors arising from non-representativeness. This, however, is misleading. Administrative data tell us what happened with a particular group of people, but this group of people may or may not be the group about which we wish to make statements or from which we wish to generalize. Very often, for example, the selection process which results in their being chosen will include an aspect of self-selection.

A few examples will illustrate some of the difficulties.

Retail banks and other financial institutions construct *scorecards* to predict likely customer behaviour with financial products. For example, such models are used to predict who is likely to

1 default on a loan, and hence whom to give loans to. Administrative data are then collected on the
 2 customers who are awarded loans as they make their repayments. In particular, outcome data—
 3 whether they defaulted or not—are collected. Such data, the outcomes along with potential
 4 predictor variables (from application forms or behaviour on other financial products), can then
 5 be used to construct models to make loan decisions on future applicants. Unfortunately, this
 6 data set will not be representative of the population of applicants. It will only include people who
 7 were previously thought to be good risks. This means that models based on it could give seriously
 8 distorted predictions for people who are drawn from the entire population of applicants (Hand
 9 and Henley, 1993; Hand, 2001). ‘All’ of the data are there, but they are not all of the data that
 10 one needs.

11 An example that is currently attracting a huge amount of attention is publication bias and
 12 associated phenomena in scientific literature. We can obtain data on all papers that are pub-
 13 lished, but they arrive at publication through a complex sociological selection process: papers
 14 reporting positive results are more likely to be submitted, editors are more likely to publish
 15 them, anomalous results may be regarded as errors so the work is not written up, and so on.
 16 So what we see is a distorted view of the work that is done and the results that are obtained, so
 17 much so, in fact, that John Ioannidis could publish a paper with the title ‘Why most published
 18 research findings are false’ (Ioannidis, 2005), stimulating much interest and subsequent work.
 19 The notion that the published scientific literature represents ‘all’ the relevant material is simply
 20 false.

21 The Crimemaps system provides another example. Originally developed in Chicago, based on
 22 police-recorded crimes, this gives (approximate) locations of crimes, displayed on maps so that
 23 people can see which areas are dangerous. However, research from the Direct Line insurance
 24 company in the UK suggests that large numbers of people are not reporting crimes because of the
 25 potentially adverse effect it will have on house prices and hence their ability to sell or rent their
 26 house (Direct Line, 2011). The data are purely administrative—from the police databases—but
 27 can become progressively more distorted. This may mean not only that the data are of limited
 28 value for determining which areas are risky, but also that they become increasingly less valuable
 29 for their original purposes. This is a straightforward illustration of Campbell’s law:

30
 31 ‘The more any quantitative social indicator is used for social decision-making, the more subject it will
 32 be to corruption pressures and the more apt it will be to distort and corrupt the social processes it is
 33 intended to monitor’.

34 This law applies just as much to administrative data as to survey or other data.

35 Even something as apparently automatic and complete as data from financial markets exhibit
 36 numerous errors and omissions. These can arise from the use of different time stamps on trans-
 37 actions, ambiguity over time resolution of transactions, the extent and method of aggregation of
 38 data, whether or not certain types of transactions are present (e.g. so-called ‘dark pool trades’),
 39 changes to symbols identifying corporations perhaps leading to mismatches or failed matches
 40 when data are linked, confusion arising from stock splits or mergers, and so on. The financial data
 41 sources described in the Caltech Quantitative Finance Group guide to market data at [http://](http://quant.caltech.edu/historical-stock-data.html)
 42 quant.caltech.edu/historical-stock-data.html illustrate the problems.

43 The bulk of traditional survey work is based on known probabilities of including each person
 44 in the study, with sampling theory permitting solid inferential conclusions, but sometimes non-
 45 probability samples are chosen, e.g. convenience sampling, matched sampling, network sampling
 46 and situations in which people are allowed to opt in or opt out. These kinds of non-probability
 47 samples have stimulated research on drawing valid population conclusions (see, for example,
 48 Baker *et al.* (2013) and Bethlehem (2010)) although it may not be straightforward to apply the

1 methods to administrative data. In general, the data distortion will have different consequences
 2 in different contexts and with different problems. For example, for one particular credit scoring
 3 data set, Crook and Banasik (2004) drew the conclusion that

4 ‘even where a very large proportion of applicants are rejected, the scope for improving on a model
 5 parameterised only on those accepted appears modest. Where the rejection rate is not so large, that
 6 scope appears to be very small indeed’.

7
 8 So, for their example at least, the problem seems not to be too bad, but it would be unwise to
 9 assume that this is generally so.

10 On top of all this, there are more complicated questions of what is meant by ‘complete’. For
 11 example, many systems are dynamic and constantly changing. A database of all the people in a
 12 company today will provide at best only a snapshot—it is likely that some employees will have
 13 moved on or been recruited by next year. Indeed, it is *certain* that the individual employees
 14 will have changed by next year—if only because they will have aged, let alone possible name
 15 changes due to marriage, address changes due to moving house, and so on. This *population drift*
 16 poses interesting statistical challenges, and it again points out the weakness of the assertion that
 17 administrative data represent ‘all’ the data one needs.

18 Statistical methods have been developed for correcting for sample distortion (e.g. Heckman
 19 (1976) and Copas and Li (1997)), but they depend on making assumptions about the form
 20 of the distortion. Statisticians can do amazing things, but they cannot perform miracles and
 21 if the data have been chosen in an arbitrary and unspecified way there is little that can be
 22 done. If this were not so we could always draw accurate conclusions from the most limited of
 23 data. And this is precisely why survey sampling and experimental design have grown into such
 24 elaborate disciplines: they specify and constrain how the data must be collected so that valid
 25 conclusions can be drawn from a statistical analysis. Administrative data, in contrast, without
 26 this underlying statistical imperative, may not be so useful for drawing statistical conclusions.
 27 They may be selected in precisely those ‘arbitrary and unspecified ways’.

28 In short, the fact that administrative data arise through an administrative process does not
 29 mean that they represent the entire population of interest. Some major successes in the world
 30 of ‘big data’ have been achieved by simply analysing the data as they present—but some major
 31 failures—such as the initial Google flu trends projections (Hodson, 2014)—have also arisen
 32 from taking the data at face value.

33 These points lead to the following challenges:

34 *Challenge 7.* Explore potential sources of non-representativeness in the data.

35
 36 *Challenge 8.* Develop and adopt tools for adjusting conclusions in the light of the data selection
 37 processes.

38 39 **4. Answering the right question**

40
 41 As the previous sections have illustrated, there can be difficulties in using administrative data to
 42 answer specific research questions. This might be because the data were not collected with those
 43 questions in mind, because of quality issues that are irrelevant to operations but highly relevant
 44 to subsequent statistical analysis, because of changes in definitions of the recorded data items
 45 or for other reasons. This brings us back to a point that was made earlier: it can be useful, if it
 46 is possible, to have statisticians involved in the data collection process. They might be able to
 47 think ahead and to expand the range of data collected so that they will be more able to answer
 48 future questions.

1 Statistical analysis methods are often divided into *descriptive* and *inferential*. Descriptive
 2 methods are used to summarize a body of data so that the important messages within it can
 3 be readily grasped. We might summarize a distribution of values by their mean and standard
 4 deviation, or the results of a census by using a series of counts organized as cross-tabulations. It
 5 goes without saying that the summary statistics that are appropriate will depend on the subject
 6 matter and on the questions to which answers are sought. Administrative data are often used
 7 for purely descriptive purposes—perhaps especially so in official statistics contexts, where we
 8 might want to establish the characteristics of some population.

9 In contrast, inferential methods are used to make a statement about unobserved values or
 10 underlying mechanisms. We might be trying to infer the disease of a new patient, on the basis
 11 of analysis of patients with similar symptoms diagnosed in the past. We might be trying to
 12 forecast whether inflation will go up or down next month. We might be trying to elucidate an
 13 underlying mechanism, so that we can understand how the data were generated, and perhaps
 14 influence things in the future. Much of the statistical theory of inference is based on the notion
 15 of random sampling from a (possibly infinite) population of values. Because the sampling is
 16 random, solid mathematics (such as the law of large numbers and the central limit theorem)
 17 means that sound statements can be made about the characteristics of the population from
 18 summary statistics obtained from the sample. Moreover, error bounds can be put on the con-
 19 clusions. We can say things such as ‘on average, 99 out of 100 of our intervals will cover the true
 20 population mean’, so we can be confident of our results (always subject to data quality issues,
 21 of course).

22 But administrative data are not collected by such a random sampling process. We can certainly
 23 calculate descriptive statistics, summarizing the data before us and, if we are willing to assume
 24 that the data are perfect, with no missing or distorted values (a brave assumption, as the above
 25 discussion has illustrated), then this will accurately summarize the population which led to our
 26 data. We can make a statement such as ‘*this* is the true population mean’.

27 But, if our aim is not really to summarize the data at hand but to make an inference to another
 28 (often future) population, then additional, unknown and quite possibly unquantified, sources
 29 of uncertainty may be relevant—possible incompleteness of the data set, discussed above, is
 30 only one example. This has consequences for inferential statements.

31 Of course, these unknown and possibly unquantified sources of uncertainty beyond sampling
 32 variation also affect the sampling approach. Indeed, it is likely that they have been underappre-
 33 ciated in many contexts, where the elegance of the sampling theory mathematics has distracted
 34 attention from the fact that there are other sources of uncertainty. This could be one of the
 35 drivers behind the phenomena to which John Ioannidis has drawn attention, mentioned above.
 36 Hand (2006) has given examples and implications of this oversight in a different context.

37 Administrative data are often highly complex. For example, a single credit card transaction
 38 leads to 70–80 items of data being recorded, whereas Web search and social media data have
 39 an elaborate graph structure. And this leads to the observation that data capture technology
 40 changes rapidly. In the context of the first of these examples, a recent change is a shift towards
 41 mobile phone banking, leading to new and additional transaction characteristics being available.
 42 In the context of the second, changes in social media platforms mean that there is a very real
 43 risk that any particular kind of model based on data recorded from Web transactions may be
 44 impossible to build just a couple of years in the future, as social interaction media change and
 45 evolve. And similar problems apply in other areas—in medicine, for example, with different
 46 kinds of bioinformatic measurement methods, in the financial sector with short time series
 47 because of changing regulations, and so on. At a substantive level, this has clear implications
 48 for studying how society is developing. At a statistical analysis level, it drives home the point

1 that was made above, that administrative data are collected for operational reasons, and may
2 have serious weaknesses for subsequent analysis.

3 Economic and social measures such as gross domestic product, the consumer price index
4 CPI and national wellbeing are what are called *pragmatic* measures (see, for example, Hand
5 (2004)): the definition of the concept and the way that it is measured are two sides of the
6 same coin. Change the measurement procedure and you change the thing being measured, with
7 different measures being suitable for different purposes. This is why, in the UK, we have CPI,
8 CPIH, RPI, RPIJ, and so on. It is not a question of any of these being more ‘right’ than the
9 others, but simply that they measure slightly different things. This means that they have different
10 properties and are suited to answering different questions. Increasingly, interest is turning to the
11 possibility of using administrative data for measuring productivity and price inflation. Instead
12 of conducting surveys of businesses to obtain data, the data can be automatically transmitted
13 from the transaction to the database. Scanner data, such as retail purchase data obtained directly
14 from the point-of-sale machine, provide an example, yielding data that are ideal for use in price
15 index calculation. Moreover, such data also give information on the volume of different goods
16 purchased, so that weights can be chosen. But issues of selection bias still apply: not all purchases
17 are made through such routes, and we cannot assume that those purchases which are made in
18 this way represent a random or representative sample of all purchases.

19 A variant of this uses Web-scraped price collection, being explored by various national sta-
20 tistical institutes. For example, the billion prices project (Cavallo and Rigobon, 2016) seeks to
21 collect massive amounts of price data from Web sites. Apart from its vast coverage (‘big data’)
22 this means that much more timely estimates can be obtained much more cheaply than by tra-
23 ditional methods. Note, however, that Cavallo and Rigobon (2016) did not describe this as an
24 alternative to traditional methods, but as a complement.

25 This approach has had some notable successes—for example, Cavallo (2013), using such on-
26 line data collection, estimated Argentina’s annual inflation rate between 2007 and 2011 to be
27 over 20%, which was a striking contrast with the 8% claimed by the Argentinian Government.
28 Moreover the on-line estimates were reported daily.

29 The apparent simplicity of this approach risks concealing various complications. It is still nec-
30 essary to decide which Web sites to collect data from—and this might be biased towards larger
31 retailers. In fact, Cavallo and Rigobon (2016) said ‘we... focus almost exclusively on large mul-
32 tichannel retailers and tend to ignore online-only retailers (such as Amazon.com)’. The basket
33 of goods to be included still must be chosen. Cavallo and Rigobon noted that on-line prices
34 cover a much smaller set of retailers and product categories than is covered by the traditional
35 approach. Also, on-line prices are one thing, but they say nothing about the quantity sold.

36 A key question, perhaps obvious in view of the earlier sections, is how representative the
37 on-line prices are of prices in general: what about prices of goods or services that are not bought
38 on line? Moreover, in certain contexts, such as airline tickets, dynamic pricing systems operate,
39 which introduces not only changes over time, but also the effect of gaming strategies.

40 One of the problems with Web-based tools is the rate of change of that technology. Companies
41 appear, grow to a massive size and vanish at a dramatic pace. Bebo, for example, launched in
42 January 2005, and sold to America On Line just 3 years later for \$850 million. But, then, in
43 May 2013 it voluntarily filed for Chapter 11 bankruptcy protection. Worse still, the algorithms
44 that are used by the companies can change arbitrarily. Google’s search algorithm is constantly
45 being redeveloped. As we have noted above, this means that administrative data may have a
46 short shelf life, in the sense that comparative data and time series may not merely have disconti-
47 nuities as definitions change but may also experience changes in ill-defined, even ill-understood,
48 ways.

1 Since survey data will be collected with a view to answering specific questions, the variables
 2 will be relevant by definition. Variables from administrative data may be less relevant. This
 3 means that *derived* variables will be more important for the analysis of administrative data.
 4 These are variables that are created by combining other variables. For example, whereas a
 5 survey question might ask about disposable income directly, to obtain the corresponding value
 6 from administrative data might require adding earned income, interest from bank and other
 7 deposits as well as other sources of income, subtracting tax paid, and so on.

8 As a final example of definitional difficulties, the media were recently exercised by an ap-
 9 parent discrepancy between the number of long-term migrants to the UK, estimated by the
 10 International Passenger Survey, and the number of national insurance number registrations
 11 (administrative data). Close examination (Office for National Statistics, 2016b) revealed that
 12 the discrepancy was due to differences in definitions. They commented that

13 “it is not possible to provide an accounting type reconciliation that simply “adds” and “subtracts”
 14 different elements of the NINo registrations to match the LTIM definitions’.

15 All of this leads to the next challenge.

16
 17 *Challenge 9.* Explore how suitable the administrative data are for answering the questions.
 18 Identify their limitations, and be wary of changes of definitions and data capture methods over
 19 time.
 20

21 Administrative data, typically being observational data, permit hypothesis-generating exer-
 22 cises. Whereas survey data have the advantage that they will be tuned to answer the survey
 23 questions, and administrative data may not be well suited to answer those questions, the con-
 24 verse also applies: administrative data, often being much richer than survey data, can be used
 25 to explore other questions and to generate hypotheses based on relationships that are observed
 26 in the data.

27 Here is one example illustrating both the complexity of human behaviour and the use of
 28 administrative data in detecting unsuspected patterns in that behaviour.

29 Hand and Blunt (2001) sought to model the distribution of sizes of credit card transactions
 30 at petrol stations in the UK, on the basis of administrative data recorded at petrol stations.
 31 Superficially, the distribution was as expected—roughly normal, but with some right skewness
 32 since it could take only positive values. However, closer investigation revealed some anomalous
 33 spikes. The size of the data set meant that these spikes could not be attributed to random variation
 34 arising from the particular period being studied but must have represented an underlying reality.
 35 (Note, the data were *all* of the transactions but were being analysed as a sample to make
 36 inferences about underlying mechanisms, and hence what might be expected to happen at other
 37 times). Closer investigation led to the observation that there are two different types of behaviour
 38 pattern: some people simply fill the petrol tank at each purchase, whereas others seek to hit a
 39 convenient whole number of pounds cost, such as £20 or £30. Noting this, and digging deeper, led
 40 us to recognize further patterns of behaviour: there was more overshoot than undershoot; people
 41 preferred to hit whole numbers of pounds of *any* magnitude than numbers ending in a non-
 42 zero number of pennies (though especially those which were a multiple of £10); subject to that,
 43 they particularly favoured numbers ending in 50p, and then 25p, and so on. Things were further
 44 complicated by the fact that a significant proportion of spend in such situations is in the forecourt
 45 shop, where goods have particular prices, often with special values of their own (e.g. ending in
 46 99p). And, as if all that was not enough; there were further features in the data which arose as a
 47 consequence of marketing initiatives run by the forecourt operations. In the end, we constructed
 48 an elaborate mixture model which tried to take all these phenomena into account. This model

1 was purely descriptive, though inference was needed to decide whether effects were sufficiently
2 large (in the context of the particular data set being supposedly drawn from a superpopulation
3 of possible such data sets) to be included. However, the aim was not merely to describe what
4 customers had done in the past, but to use the model to inform future pricing strategies.

5 A particular merit of administrative data, and especially of transaction data, is that it is
6 recorded as time progresses. Unlike data that are recorded at a particular time, or a discrete
7 sequence of times as in repeated surveys, it is essentially continuous. This means that admin-
8 istrative data can be very useful for early detection of changes in populations. Indeed, often
9 one of the operational reasons for collecting the data in the first place will be for monitoring
10 processes. But an assertion that a time series (of gross domestic product or unemployment,
11 say) has changed will typically not be intended merely as a face value assertion that the raw
12 numbers differ, but rather as an assertion that some underlying reality has changed. And this
13 should not be based on a simple comparison of the numbers, but rather on a comparison of
14 the difference between the numbers with the inaccuracy of measurement. It should answer the
15 question: is this difference larger than we would expect, given the intrinsic uncertainty in how
16 the measured numbers represent the reality, or is it well within the scope of what we might
17 expect with no change in the underlying reality? And the crucial point is that this intrinsic
18 uncertainty should include all sources of uncertainty, not merely sampling variation (if that is
19 indeed relevant).

20 This leads to the following challenge.

21 *Challenge 10.* Report changes and time series with appropriate measures of uncertainty, so that
22 both the statistical and the substantive significance of changes can be evaluated. The measures
23 of uncertainty should include all sources of uncertainty which can be identified.
24

25 5. Causality and intervention

26 As is well known, observational data present challenges in establishing causality. If we observe
27 a difference in some outcome measure (e.g. income) between two groups, and we note that
28 the groups differ in various properties (e.g. education) we cannot be sure that the observed
29 differences in their properties explain or cause the difference in outcome. To establish causality,
30 we need to intervene to break all possible causal links except the link that we wish to test (but
31 see also Pearl *et al.* (2016)). The most common way to do this is via a properly controlled
32 experiment involving randomization. Usually this is difficult with administrative data, not least
33 because it requires modifying the standard operation of the organization, although occasionally
34 experimental designs are built into on-going operations, enabling comparisons to be made by
35 using administrative data. In such situations the designs will typically be fairly simple, such as
36 merely comparing two groups.
37

38 This notion of modifying an operation so that we can learn from it, as well as simply using it
39 to carry out its normal function, can manifest in other ways. One mail order organization that
40 we worked with enrolled a ‘gold sample’—a small set of people regarded as poor risks (who
41 would normally be rejected), just so that they could collect data across the entire population
42 distribution, and hence enhance their models and improve future predictions. Scholtus *et al.*
43 (2015) have also explored the use of such a gold sample, in their case to yield estimates of
44 intercept bias in a model.

45 We see from this that the needs of planned subsequent statistical analysis can sometimes
46 influence what administrative data are to be collected. Occasionally, such intentions can lead to
47 data being collected during the operations which are not required to run the organization but
48 which can be used subsequently—the second kind of administrative data mentioned in Section 1.

1 *Challenge 11.* Be aware that administrative data are observational data, and exercise due
 2 caution about claiming causal links.
 3

4 **6. Combining data from different sources**

5 Combining data and evidence from different sources is increasingly important in statistics and
 6 elsewhere. This can be for statistical purposes, such as to yield an improved or more comprehen-
 7 sive estimate (e.g. Ashley *et al.* (2005) and Cunningham and Jeffery (2007)), or simply because
 8 information is needed for a higher level organization (e.g. combining statistics from several
 9 countries to give European-Union-wide statistics). But it can also be at the individual level, e.g.
 10 in detecting fraud or adverse drug reactions, or tracking terrorist activity.

11 Even if the data are, at least in principle, of the same type, such as combining economic
 12 statistics from different countries, they may have been collected by using different methods or
 13 definitions, so producing combinations or aggregates is not necessarily straightforward. Vâju
 14 *et al.* (2015) describe

15 ‘a huge number of possible sources of lack of comparability, given by combinations of (i) national legal
 16 and institutional environments, (ii) acceptable trade-off between quality dimensions at national level,
 17 (iii) appropriate trade-off between costs and benefits in terms of output data quality at national level,
 18 (iv) methodological choices to integrate the several data sources’.

19 At a lower level, problems might arise because a particular characteristic might be grouped
 20 in different ways in two data sets (e.g. age classified into 10-year bands or into young *versus*
 21 old), or observations might be taken or recorded with different periodicities. Possible strategies
 22 for overcoming such problems include the latent variable perspective, with the observed data
 23 being regarded as a coarsened or grouped version, or state space models (Horn and Czaplewski,
 24 2013).
 25

26 The situation is further complicated because the data are often of different types—survey
 27 data, administrative data, Web-scraped data, social network data, data collected from wrist
 28 health and activity monitors, and even non-numerical forms of data such as speech and image
 29 data. This is perhaps where the real opportunities, and statistical challenges, arise. Medicine,
 30 in particular, is making extensive use of such approaches, combining medical images, clinical
 31 trial reports, epidemiological data and health registry data. Credit bureaus combine credit card
 32 transaction records from several operators to build a single database from which they can
 33 construct a generally applicable credit scorecard. An example from official statistics in the UK
 34 is the estimation of income within small geographical areas, based on linking data from the
 35 Family Resources Survey and administrative data from benefit claimant counts, council tax
 36 bandings and tax credit claims. Vâju *et al.* (2015) pointed out that, even if the accuracy of
 37 the separate sources of data can be measured, assessing the sensitivity of the accuracy of the
 38 final combined data set to the source-specific errors and the integration methods can be very
 39 difficult.
 40

41 As far as merging data from different sources goes, reasons include the following.

- 42 (a) *Complement:* different sources of data and different types of data, can each serve as a
 43 complement to each other by providing different types of information. This is perhaps
 44 particularly true for administrative and survey data. Some types of variables—attitudes
 45 and opinions, for example—do not normally naturally arrive in administrative data but
 46 must be collected by surveys (or panels, or some other purposive data collection strategy).
 47 Surveys can be designed so that they shed light on tightly focused research questions,
 48 whereas with administrative data we may have to be satisfied with questions which are a

1 little different from those we would ideally like to ask, perhaps because they are based
2 on slightly different definitions of the concepts involved. In contrast, administrative data
3 sets are likely to be larger, with better population coverage (though possibly vulnerable
4 to the other data quality issues that were mentioned above).

- 5 (b) *Supplement*: although administrative data are often thought of as an alternative to survey
6 data, they are at least as valuable when used in conjunction with survey data. Survey
7 data can be used to pinpoint particular research questions, but cost necessarily limits
8 coverage. However, relationships that are found from survey data can be extrapolated
9 to yield estimates from overall populations and smaller groups by using such tools as
10 regression estimation applied to an administrative data population base. This can be
11 useful to yield small area and regional estimates. Indeed, such statistical tools can be used
12 to improve estimates from survey data. A further point is that surveys require sampling
13 frames, and administrative data are central to their construction.
- 14 (c) *Accuracy*: we have stressed issues of data quality above. Triangulation and imputation
15 from multiple sources of data and reconciliation between sources of data are good ways to
16 tackle these issues. Berka *et al.* (2012) gave an example, exploring accuracy in the Austrian
17 register-based census of 2011. They noted the use of surveys to check register data but
18 pointed out that this is resource intensive. They evaluated the quality of data at the raw data
19 level in terms of three ‘hyperdimensions’, assessing documentation (e.g. plausibility and
20 legal aspects), preprocessing (formal methods for testing for errors and inconsistencies)
21 and comparison with an external source. The results are three measures, each scored in
22 the interval 0–1. A weighted average is taken to yield an overall quality indicator for
23 each register and attribute. The fundamental challenge here is that of combining quality
24 indicators from different sources, and Berka *et al.* (2012) explored the use of Dempster–
25 Shafer theory to do this.

27 Another example was given by Romanov and Gubman (2013), who used administrative data
28 to explore bias in answers to survey questions about income. Discrepancies pinpoint potential
29 errors and issues to be resolved. Of course, there are complications. Errors can propagate and
30 perhaps not all of them can be resolved. Worse, especially in the context of administrative data,
31 this jigsaw solution is vulnerable to one of the pieces disappearing as the operational imperatives
32 generating the administrative data change. Moreover, as we have repeatedly stressed, one must
33 be alert to different sources of data using different definitions.

34 A special case of merging data from different sources is matching data from different adminis-
35 trative databases. For example, we may have identified data on individuals, collected for different
36 reasons and stored in two distinct databases, and we may want to combine them. But, of course,
37 the problem is not restricted to data on individuals: Lewis and Woods (2013) described a prob-
38 lem of incompatible business registers, with different identifiers in the two databases. Because
39 of its importance the matching of corresponding records from different databases has been
40 the focus of much research effort—see, for example, Christen (2012), D’Orazio *et al.* (2006)
41 and Rässler (2002). It faces various data analytic challenges, including deciding when to match
42 two records given that they do not have unique and identical identifiers, detection of duplicate
43 records (again, because slightly different identifiers may refer to the same individual person or
44 object) and merging of duplicate records into a single entity (or *deduplication*).

45 A traditional, and still widely used, method, at least for small data sets, is manual matching.
46 This has some obvious shortcomings, including a scalability cost (in various measures), subject-
47 ivity arising from human biases, variation between people, variation within any one person as
48 they become tired or bored, and the difficulty of objectively improving performance. A modern

variant of manual matching, for contexts where confidentiality is not important or where the data to be matched can be effectively encrypted, is crowdsourcing, enlisting the help of large numbers of people.

Computational methods can be divided into two classes: deterministic and probabilistic or statistical.

Deterministic methods simply see whether two records agree on all of a specified set of identifiers. This is clearly very quick. It can be a single-step procedure or can proceed through sequential steps, beginning with stringent matching criteria and progressively relaxing them.

Probabilistic methods relax the requirement of an exact match and instead calculate a dissimilarity measure for each field in the pair of records being compared. The choice of dissimilarity measure will depend on the context (e.g. approximate string matches for some text fields, matches that allow different date formats and matches that allow the given name and surname to occur in the reverse order). The separate field dissimilarity measures are then combined (e.g. added or used to maximize the likelihood of a match, given a probability model) to yield an overall dissimilarity measure for the record pair. In the simplest approaches, these dissimilarity measures can then be compared with a threshold to yield a match–non-match classification. More sophisticated approaches (e.g. the classic work of Fellegi and Sunter (1969)) follow the ‘reject option’ and define three types of decisions: match, non-match or possible-match. The third class is then subjected to a second stage of investigation, which often a manual comparison. Winkler (2006) has reviewed linkage methods.

Clearly methods which are based on pair-by-pair comparisons run the risks of intransitivity, of several records from one database being matched to a single record in the other and of computational intractability if all possible pairs are compared. The first two problems, at least, can be eased if a higher level view of the matching process is taken, in which constrained groups of records are compared. To take a simple example, suppose that we wanted to match a collection of left shoes with a collection of right shoes, to find which shoes belonged in a pair. One strategy would be simply to calculate similarities between shoes, one from each collection, and to choose the pairs which had the greatest similarity—but this would be susceptible to the first two of the problems just listed. An alternative approach would be (computation allowing) to look at all possible pairings of shoes, one from each collection, and to choose the set of pairings which maximized the likelihood. Exactly this sort of approach has been used in chromosome matching.

These considerations lead to the following challenges.

Challenge 12. Be aware of the risks that are associated with linked data sets and the potential effect on the accuracy and validity of any conclusions. Recognize that quality issues of individual databases may propagate and amplify in linked data. Develop better measures of overall combined data quality.

Challenge 13. Continue to develop statistically principled and sound methods for record linkage and evidence assimilation, especially from non-structured data and data of different modes.

Challenge 14. Develop improved methods for data triangulation, combining different sources and types of data to yield improved estimates.

7. Confidentiality, privacy and anonymization

A common challenge with all data describing human beings is the need to preserve confidentiality and privacy, but this often seems to be a particularly sensitive issue with administrative data.

This may be because, unlike with surveys, there may be no choice about being included (at least, if one wants access to the service or product) or perhaps because it is obvious that the identifier must be retained in the data (since it is needed for operational reasons—one cannot run a credit card operation without being able to match transactions to customers). There seems to be growing concern about the *data shadows* that we all inevitably leave as we access administrative services, whether corporate or public.

Anonymization and deidentification tools do exist—e.g. based on aggregating data, perturbing data or randomly generating data with statistical properties the same as the raw data—but they all have shortcomings. An overview of such methods is given in Duncan *et al.* (2011) and see also Karr *et al.* (2006), Reiter (2005), Matthews and Harel (2011) and McClure and Reiter (2016). One of the most challenging—and probably intractable—problems is that it is often possible to combine a data set with other publicly available data to identify an individual and to reveal something about them. There have been several well-known public incidents of this kind, such as the identification of individual subscribers from the Netflix prize data set (Narayanan and Shmatikov, 2008) and the identification of the medical records of Massachusetts Governor William Weld (Anderson, 2009).

From the perspective of statistical challenges, work continues to develop statistical methods of disclosure control—such as the development of differential privacy (Dwork and Roth, 2014). More generally, statistical tools are being developed to permit analysis without divulging the identity of individuals. For example, multiparty computation is a strategy to calculate aggregate statistics for a collection of individuals without requiring any individual to give away their value (Cramer *et al.*, 2015).

Challenge 15. Continue to explore anonymization and de-identification methods.

8. Conclusion

In the paper, I have sought to identify and characterize what I thought were the main statistical challenges arising from administrative data. There are other challenges, including the following three:

- (a) *The communication of uncertainty:* as statisticians, we are familiar with uncertainty arising from sampling variation, and with methods of communicating that uncertainty, such as confidence intervals. However, since the sources of uncertainty in administrative data are many and diverse, and may not include sampling variation, we need to find other ways to communicate (and indeed perhaps even to define) such uncertainties. In some contexts this is already done. For example, the Bank of England's August 2016 inflation report (Bank of England, 2016), chart 5.1, shows a fan chart with

‘To the left of the vertical dashed line, the distribution reflects the likelihood of revisions to the data over the past; to the right, it reflects uncertainty over the evolution of GDP growth in the future’. More, however, remains to be done. Manski (2014) has a good discussion of the issues.

- (b) *Statistical education:* challenge 1 above was about statistical education, although limited to the context of data quality. Administrative data are becoming so important, and so widely used (as a consequence of automatic data capture), that one can argue a case for more specialized teaching of specific methods related to administrative data.
- (c) *Legal environment:* the growth of awareness of modern data analysis technology has stimulated considerable legal and regulatory thought, much a consequence of the privacy and confidentiality issues discussed above. On April 14th, 2016, the European Union's General

1 Data Protection Regulation (European Union, 2017) was adopted by the European Par-
 2 liament, and on April 27th, 2017, the UK's Digital Economy Act received Royal assent
 3 (Her Majesty's Government, 2017). These changes will certainly impact how personal
 4 data are stored and are likely to impact statistical analyses of administrative data.
 5

6 As a final comment, applied statisticians often emphasize the importance of being familiar
 7 with the data generation process. Understanding where the data come from and how they are
 8 collected can lead to the avoidance of many misunderstandings and mistakes. At first glance it
 9 might seem as if this is less critical for administrative data. This, however, is not so. Issues of
 10 data quality, changes over time, changing regulatory and legal environments, advances in data
 11 capture and access technology and a host of other factors are likely to impact administrative
 12 data, and their analysis. In fact, because the data will have been primarily collected for some
 13 operational purpose, these changes will almost certainly have been made without any subsequent
 14 statistical analysis in mind. They may not even be reported to the statistician who is later
 15 analysing the data. It is thus even more important—perhaps essential—that the statistician
 16 understands the data collection process. But note that this is a two-way communication. If they
 17 are aware of the analyses to be undertaken later, the data producers will be able to adjust their
 18 data collection and recording processes to facilitate the subsequent analyses.

19 The aim of this paper is to raise awareness and to stimulate discussion among statisticians of
 20 the need for methodological statistical work on administrative data. Such data are being used
 21 increasingly more widely—partly a consequence of the 'big data' revolution. But drawing valid
 22 conclusions from such data encounters problems that are distinct from the more familiar and
 23 well-trodden paths of sampling theory inference. The problems are diverse and heterogeneous,
 24 so it is doubtful that a unifying theory as elegant as that of sampling theory can be developed.
 25 But nevertheless some principles apply. These include the need to cope with rather different
 26 kinds of data quality issues, the recognition that, despite superficial appearances, we typically
 27 do not have 'all' the data, possible mismatches between the question we want to answer and
 28 the information in the available data, challenges arising from the fact that the data are (usually)
 29 merely observational, so elucidation of causality is difficult, the need to combine data from
 30 multiple rather different sources, and issues of confidentiality, privacy, and anonymization
 31 which might be rather different from those of survey data.
 32

33 Acknowledgements

34 The first draft of this paper was written as part of the Isaac Newton Institute programme on 'Data
 35 Linkage and Anonymization', July–December 2016. I would like to express my appreciation to
 36 the three referees and the Associate Editor for their detailed and helpful comments, which led
 37 to substantial improvement of the paper. The opinions expressed in this paper are the personal
 38 opinions of the author and do not necessarily reflect those of any organization with which the
 39 author is associated.
 40

41 References

- 42
 43 Anderson, N. (2009) "Anonymized" data really isn't—and here's why not. (Available from <https://arstechnica.com/tech-policy/2009/09/your-secrets-live-online-in-databases-of-ruin/>.)
 44
 45 Antoni, M. (2013) Linking survey data with administrative employment data: the case of the German ALWA
 46 survey. *NTTS 2013*, 279–289.
 47
 48 Ashley, J., Driver, R., Hayes, S. and Jeffery, C. (2005) Dealing with data uncertainty. *Bnk Engl. Q. Bull.*,
 spring. (Available from <http://www.bankofengland.co.uk/publications/Documents/quarterlybulletin/qb050101.pdf>.)

- Baker, R., Brick, J. M., Bates, N. A., Battaglia, M., Couper, M. P., Dever, J. A., Gile, K. J. and Tourangeau, R. (2013). Summary report of the AAPOR task force on non-probability sampling. *J. Surv. Statist. Methodol.*, **1**, 90–143.
- Banasik, J. and Crook, J. (2004) Does reject inference really improve the performance of application scoring models? *J. Bnkng Finan.*, **28**, 857–874.
- Berka, C., Humer, S. and Moser, M. (2012) Combination of evidence from multiple administrative data sources: quality assessment of the Austrian register-based Census 2011. *Statist. Neerland.*, **66**, 18–33.
- Bethlehem, J. G. (2010) Selection bias in web surveys. *Int. Statist. Rev.*, **78**, 161–188.
- Biemer, P., Trewin, D., Bergdahl, H. and Japac, L. (2014) A system for managing the quality of official statistics. *J. Off. Statist.*, **30**, 381–415.
- Cavallo, A. (2013) Online and official price indexes: measuring Argentina's inflation. *J. Monet. Econ.*, **60**, 152–165.
- Cavallo, A. and Rigobon, R. (2016) The billion prices project: using online prices for measurement and research. *J. Econ. Perspect.*, **30**, 151–178.
- Ćetković, P., Humer, S., Kausl, A., Lenk, M., Moser, M., Rechta, H. and Schnetzer, M. (2013) Quality measurement in administrative statistics with a special focus on quality assessment of imputations. *NTTS 2013*, 247–256.
- Christen, P. (2012) *Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*. Heidelberg: Springer.
- Copas, J. B. and Li, H. G. (1997) Inference for non-random samples (with discussion). *J. R. Statist. Soc. B*, **59**, 55–95.
- Cramer, R., Damgård, I. B. and Nielsen, J. B. (2015) *Secure Multiparty Computation and Secret Sharing*. New York: Cambridge University Press.
- Cunningham, A. and Jeffery, C. (2007) Extracting a better signal from uncertain data. *Q. Bull. Bnk Engl.*, quarter3. (Available from <http://www.bankofengland.co.uk/publications/Documents/quarterlybulletin/qb070301.pdf>.)
- Daas, P. J. T., Arends-Tóth, J., Schouten, B. and Kuijvenhoven, L. (2008) Proposal for a quality framework for the evaluation of administrative and survey data. (Available from <http://www.pietdaas.nl/beta/pubs/pubs/ESSnet.Vienna.paper.pdf>.)
- De Veaux, R.D. and Hand, D. J. (2005) How to lie with bad data. *Statist. Sci.*, **20**, 231–238.
- Direct Line (2011) Direct Line, Leeds. (Available from <https://www.directline.com/media/archive-2011/news-11072011>.)
- D'Orazio, M., Di Zio, M. and Scanu, M. (2006) *Statistical Matching: Theory and Practice*. Chichester: Wiley.
- Duncan, G. T., Elliott, M. and Salazar-González, J.-J. (2011) *Statistical Confidentiality: Principles and Practice*. New York: Springer.
- Dwork, C. and Roth, A. (2014) The algorithmic foundations of differential privacy. *Foundns Trends Theoret. Comput. Sci.*, **9**, 211–407.
- ESSNet (2017) (Available from https://ec.europa.eu/eurostat/cros/page/essnet_en.)
- ESSNet Admin Data Workshop (2013) (Available from <https://ec.europa.eu/eurostat/cros/content/essnet-admin-data-workshop-using-administrative-data-sts-evaluation-questionnaire-main.en>.)
- European Statistical System (2020) European Statistical System Vision 2020. European Statistical System (Available from <http://ec.europa.eu/eurostat/web/ess/about-us/ess-vision-2020>.)
- European Statistical System Admin (2015) Administrative data sources business project. European Statistical System Admin. (Available from <http://ec.europa.eu/eurostat/documents/7330775/7339647/ADMIN+fact+sheet.pdf/cbb590b2-9d6f-439c-af2d-ca8b5e9cf1f7>.)
- European Union (2017) (Available from <http://www.eugdpr.org/>.)
- Eurostat (2000) Assessment of the quality in statistics. Eurostat, Luxembourg. (Available from <http://citeseerx.ist.psu.edu/viewdoc/download;jsessionid=FDD8499ECE9F43685CAFDD2A10EBE317?doi=10.1.1.5.8718&rep=rep1&type=pdf>.)
- Eurostat (2003) Eurostat, Luxembourg. (Available from http://ec.europa.eu/eurostat/documents/64157/4374310/36-QUALITY-ASSESSMENT-ADMINISTRATIVE-DATA-STATISTICAL-PURPOSES_2003.pdf/37373e67-d69c-4215-b727-5b036393b80f.)
- Eurostat (2011) European statistics code of practice. Eurostat, Luxembourg. (Available from <http://ec.europa.eu/eurostat/documents/3859598/5921861/KS-32-11-955-EN.PDF/5fa1ebc6-90bb-43fa-888f-dde032471e15>.)
- Fellegi, I. P. and Sunter, A. B. (1969) A theory for record linkage. *J. Am. Statist. Ass.*, **64**, 1183–1210.
- Fowler, F. J. (1995) *Improving Survey Questions: Design and Evaluation*. Thousand Oaks: Sage.
- Hand, D. J. (2001) Reject inference in credit operations. In *Handbook of Credit Scoring* (ed. E. May), pp. 225–240. Chicago: Glenlake.
- Hand, D. J. (2004) *Measurement Theory and Practice: the World Through Quantification*. Chichester: Wiley.
- Hand, D. J. (2006) Classifier technology and the illusion of progress (with discussion). *Statist. Sci.*, **21**, 1–34.
- Hand, D. J. (2008) *Statistics: a Very Short Introduction*. Oxford: Oxford University Press.

- 1 Hand, D. J. and Blunt, G. (2001) Prospecting for gems in credit card data. *IMA J. Mangmnt Math.*, **12**, 173–200.
- 2 Hand, D. J., Blunt, G., Kelly, M. G. and Adams, N. M. (2000) Data mining for fun and profit. *Statist. Sci.*, **15**,
3 111–126.
- 4 Hand, D. J. and Henley, W. E. (1993) Can reject inference ever work? *IMA J. Math. Appl. Bus. Indstry*, **5**, 45–
5 55.
- 6 Heckman, J. J. (1976) The common structure of statistical models of truncation, sample selection and limited
7 dependent variables, and a simple estimator for such models. *Ann. Econ. Soci Measmnt*, **5**, 475–492.
- 8 Hellerstein, J. (2008) Quantitative data cleaning for large databases. (Available from [http://db.cs.
9 berkeley.edu/jmh/papers/cleaning-unece.pdf](http://db.cs.berkeley.edu/jmh/papers/cleaning-unece.pdf).)
- 10 Her Majesty's Government (2017) Digital Economy Act 2017. London: Stationery Office. (Available from
11 <http://www.legislation.gov.uk/ukpga/2017/30/contents/enacted>.)
- 12 Hodson, H. (2014) Google Flu Trends gets it wrong three years running. *New Scient.*, Mar. 13th.
- 13 Horn, S. and Czaplewski, R. (2013) Combining survey and administrative data using state space models. *NTTS
14 2013*, 174–183.
- 15 Ioannidis, J. (2005) Why most published research findings are false. *PLoS Med.*, **2**, 696–701.
- 16 Israel Central Bureau of Statistics (2007) Pros and cons for using administrative records in Statistical Bureaus.
17 Statistical Commission, Economic Commission for Europe, Brussels. Central Bureau of Statistics (Available
18 from <http://www.oecd.org/std/41143741.pdf>.)
- 19 de Jonge, E. and van der Loo, M. (2013) *An Introduction to Data Cleaning with R*. The Hague: Statistics Nether-
20 lands.
- 21 Karr, A. F., Kohnen, C. N., Oganian, A., Reiter, J. P. and Sanil, A. P. (2006) A framework for evaluating the utility
22 of data altered to protect confidentiality. *Am. Statistn*, **60**, 224–232.
- 23 Karr, A. F., Sanil, A. P. and Banks, D. L. (2006) Data quality: a statistical perspective. *Statist. Methodol.*, **3**,
24 137–173.
- 25 Kim, W., Choi, B.-Y., Hong, E.-K., Kim, S.-K. and Lee, D. (2003) A taxonomy of dirty data. *Data Minng Knowl.
26 Discov.*, **7**, 81–99.
- 27 Kloek, W. and Váju, S. (2013) The use of administrative data in integrated statistics. *NTTS 2013*, 128–138.
- 28 Lewis, D. and Woods, J. (2013) Issues to consider when turning to the use of administrative data: the UK experience.
29 *NTTS 2013*, 549–557.
- 30 Manski, C. (2014) Communicating uncertainty in official economic statistics. *Working Paper 20098*.
31 National Bureau of Economic Research, Cambridge. (Available from [http://papers.ssrn.
32 com/sol3/papers.cfm?abstract_id=2432840](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2432840).)
- 33 Matthews, G. J. and Harel, O. (2011) Data confidentiality: a review of methods for statistical disclosure limitation
34 and methods for assessing privacy. *Surv. Statist.*, **5**, 1–29.
- 35 McClure, D. and Reiter, J. P. (2016) Assessing disclosure risks for synthetic data with arbitrary intruder knowledge.
36 *Statist. J. Int. Ass. Off. Statist.*, **32**, 109–126.
- 37 Meader, R. and Tily, G. (2008) Monitoring the quality of national accounts. *Econ. Lab. Markt Rev.*, **2**, 24–33.
- 38 Memobust Handbook (2014) Quality of statistics module. (Available from [https://ec.europa.eu/
39 eurostat/cros/system/files/Quality%20Aspects-01-T-Quality%20of%20Statistics%20
40 v1.0.pdf](https://ec.europa.eu/eurostat/cros/system/files/Quality%20Aspects-01-T-Quality%20of%20Statistics%20v1.0.pdf).)
- 41 Narayanan, A. and Shmatikov, V. (2008) Robust de-anonymization of large sparse datasets. (Available from
42 https://www.cs.cornell.edu/~shmat/shmat_oak08netflix.pdf.)
- 43 van Nderpelt, P. W. M. (2009) Checklist quality of statistical output. (Available from [https://www.cbs.nl/
44 NR/rdonlyres/4119715F-7437-4379-9A70-90A0893F949E/0/2009ChecklistQualityofSta
45 tisticalOutput.pdf](https://www.cbs.nl/NR/rdonlyres/4119715F-7437-4379-9A70-90A0893F949E/0/2009ChecklistQualityofStatisticalOutput.pdf).)
- 46 New Techniques and Technologies for Statistics (2013) New Techniques and Technologies for Statistics: the meet-
47 ing place for research in official statistics. (Available from [https://ec.europa.eu/eurostat/cros/
48 system/files/NTTS2013%20Proceedings.0.pdf](https://ec.europa.eu/eurostat/cros/system/files/NTTS2013%20Proceedings.0.pdf).)
- 49 New Techniques and Technologies for Statistics (2015) New Techniques and Technologies for Statistics: re-
50 liable evidence for a society in transition. (Available from [https://ec.europa.eu/eurostat/cros/
51 system/files/NTTS2015%20proceedings.pdf](https://ec.europa.eu/eurostat/cros/system/files/NTTS2015%20proceedings.pdf).)
- 52 Nordbotten, S. (2010) The use of administrative data in official statistics—past, present and future: with special
53 reference to the Nordic countries. In *Official statistics: Methodology and Applications in Honour of Daniel
54 Thorburn* (eds M. Carlson, H. Nyquist and M. Villani), pp. 205–223, Stockholm, Statistics Sweden.
- 55 Office for National Statistics (2016a) Crime in England and Wales, year ending Mar 2016. *Statistical Bul-
56 letin*. Office for National Statistics, Newport. (Available from [http://www.ons.gov.uk/people
57 populationandcommunity/crimeandjustice/bulletins/crimeinenglandandwales/yeare
58 ndingmar2016](http://www.ons.gov.uk/peoplepopulationandcommunity/crimeandjustice/bulletins/crimeinenglandandwales/yearendingmar2016).)
- 59 Office for National Statistics (2016b) Note on the difference between National Insurance registrations and
60 the estimate of long-term international migration: 2016. Crime in England and Wales, year ending Mar
61 2016. *Statistical Bulletin*. Office for National Statistics, Newport. (Available from [https://www.ons.
62 gov.uk/peoplepopulationandcommunity/populationandmigration/internationalmigra
63 tion/articles/noteonthedifferencebetweennationalinsurancenumberregistrationsa
64 ndtheestimateoflongterminternationalmigration/2016](https://www.ons.gov.uk/peoplepopulationandcommunity/populationandmigration/internationalmigration/articles/noteonthedifferencebetweennationalinsurancenumberregistrationsandtheestimateoflongterminternationalmigration/2016).)

- 1 Organisation for Economic Co-operation and Development (2016) Short-term economic statistics (STES)
2 administrative data: two frameworks of papers. Organisation for Economic Co-operation and Development,
3 (Available from <http://www.oecd.org/std/short-term-economic-statistics/stes-administrative-data-two-frameworks-of-papers.htm#Definition>.)
- 4 Pearl, J., Glymour, M. and Jewell, N. P. (2016) *Causal Inference in Statistics: a Primer*. Chichester: Wiley.
- 5 Presser, S., Rothgeb, J. M., Couper, M. P., Lessler, J. T., Martin, E., Martin, J. and Singer, E. (eds) (2004) *Methods
6 for Testing and Evaluating Survey Questionnaires*. New York: Wiley.
- 7 Rässler, S. (2002) *Statistical Matching: a Frequentist Theory, Practical Applications, and Alternative Bayesian
8 Approaches*. New York: Springer.
- 9 Reiter, J. P. (2005) Estimating risks of identification disclosure for microdata. *J. Am. Statist. Ass.*, **100**, 1103–
10 1113.
- 11 Romanov, D. and Gubman, Y. (2013) Estimation of measurement error in categorical income survey data. *NTTS
12 2013*, 78–87.
- 13 Ruggles, P. (2015) Review of administrative data sources. (Available from <https://www.census.gov/content/dam/Census/library/working-papers/2015/acs/2015Ruggles.01.pdf>.)
- 14 Scholtus, S. and Bakker, B. F. M. (2013) Estimating the validity of administrative and survey variables by means
15 of structural equation models. *NTTS 2013*, 290–299.
- 16 Scholtus, S., Bakker, B. F. M. and van Delden, A. (2015) Modelling measurement error to estimate bias in
17 administrative and survey variables. *NTTS 2015*, 451–455.
- 18 Statistics Canada (2015) Statistics Canada, Ottawa. (Available from <http://www.statcan.gc.ca/pub/12-539-x/2009001/administrative-administratives-eng.htm>.)
- 19 Statistics Netherlands (2014) Quality guidelines 2014: Statistics Netherlands quality assurance framework
20 at process level. Statistics Netherlands, Voorburg. (Available from <http://ec.europa.eu/eurostat/documents/64157/4374310/01-Quality-Guidelines-2014-Statistics-Netherlands-Quality.pdf/292b18bc-9bfd-426d-9282-785aabc43126>.)
- 21 Statistics New Zealand (2016) Guide to reporting on administrative data quality. Statistics New Zealand. (Available
22 from <http://www.stats.govt.nz/methods/data-integration/guide-to-reporting-on-admin-data-quality/explaining-framework.aspx#>.)
- 23 Trépanier, J., Pignal, J. and Royce, D. (2013) Administrative data initiatives at Statistics Canada. (Available from
24 <http://www.copafs.org/UserFiles/file/fcsm/G1.Trepanier.2013FCSM.pdf>.)
- 25 UK Statistics Authority (2014) *Quality Assurance and Audit Arrangements for Administrative Data*. London: UK
26 Statistics Authority.
- 27 UK Statistics Authority (2015) *Administrative Data Quality Assurance Toolkit*. London: UK Statistics Authority.
- 28 United Nations Economic Commission for Europe (2016) (Available from <http://www1.unece.org/stat/platform/display/Collection/2012+Data+Collection+Seminar%3A+Documents>.)
- 29 Văju, S., Agafiței, M., Gras, F., Kliek, W. and Reis, F. (2015) Measuring the quality of multisource statistics.
30 *NTTS 2015*, 456–459.
- 31 de Waal, T., Pannekoek, J. and Scholtus, S. (2011) *Handbook of Statistical Data Editing and Imputation*. Chichester:
32 Wiley.
- 33 Wallgren, A. and Wallgren, B. (2014) *Register-based Statistics: Statistical Methods for Administrative Data*, 2nd
34 edn. Chichester: Wiley.
- 35 Winkler, W. E. (2006) Overview of record linkage and current research directions. Statistical Research Division,
36 US Bureau of the Census, Washington DC.
- 37
38
39
40
41
42
43
44
45
46
47
48