

Small data in the era of big data

Rob Kitchin · Tracey P. Lauriault

Published online: 11 October 2014
© Springer Science+Business Media Dordrecht 2014

Abstract Academic knowledge building has progressed for the past few centuries using small data studies characterized by sampled data generated to answer specific questions. It is a strategy that has been remarkably successful, enabling the sciences, social sciences and humanities to advance in leaps and bounds. This approach is presently being challenged by the development of big data. Small data studies will however, we argue, continue to be popular and valuable in the future because of their utility in answering targeted queries. Importantly, however, small data will increasingly be made more big data-like through the development of new data infrastructures that pool, scale and link small data in order to create larger datasets, encourage sharing and reuse, and open them up to combination with big data and analysis using big data analytics. This paper examines the logic and value of small data studies, their relationship to emerging big data and data science, and the implications of scaling small data into data infrastructures, with a focus on spatial data examples.

Keywords Big data · Small data · Data infrastructures · Cyber-infrastructures · Ontology · Epistemology

Introduction

Until recently, academic knowledge building was conducted through what, in the context of emerging big data, might now be termed small data studies: that is, studies underpinned by data produced in tightly controlled ways using sampling techniques that limited their scope, temporality, size and variety, and which tried to capture and define their levels of error, bias, uncertainty and provenance (Miller 2010). Small data are thus characterized by their generally limited volume, non-continuous collection, narrow variety, and are usually generated to answer specific questions. In contrast, new forms of big data produced predominantly through new information and communication technologies (ICTs) are characterised as being large in volume, produced continuously, and varied in nature, although they are often a by-product of systems rather than being designed to investigate particular phenomena or processes (Laney 2001; Mayer-Schonberger and Cukier 2013). The rapid growth and impact of big data has led some to ponder whether big data might lead to the demise of small data, or whether the stature of studies based on small data might be diminished, due to their limitations in size, temporality and relative

R. Kitchin (✉) · T. P. Lauriault
NIRSA, National University of Ireland Maynooth,
County Kildare, Ireland
e-mail: Rob.Kitchin@nuim.ie

T. P. Lauriault
e-mail: Tracey.Lauriault@nuim.ie

cost. Indeed, Sawyer (2008) notes that funding agencies are evermore pushing their limited funding resources to data-rich areas and big data analytics at the expense of small data studies, a trend that has continued in recent years (Kitchin 2013).

This paper scrutinizes such concerns by considering the value of small data in an emerging era of big data and how they are being reconceived in the context of new data archiving and sharing infrastructures. We examine how small data are increasingly being pooled, linked and scaled into data infrastructures that make them more big data-like—that is, amenable to combination with big data and open to analysis using big data analytics, though the data themselves do not hold the inherent ontological characteristics of big data. As such, our focus is not big data per se, though we do discuss big data in order to help make sense of the changes occurring with respect to small data.

The principal arguments we develop are three fold. First, despite the rapid growth of big data and associated analytics, small data studies will continue to flourish because they have a proven track record of answering specific questions. Second, the data from these studies will more and more be pooled, linked, and scaled through new data infrastructures, with an associated drive to try to harmonize small data with respect to data standards, formats, metadata, and documentation, in order to increase their value through combination and sharing. Third, scaling small data exposes them to the new epistemologies of data science and to incorporation within new multi-billion data markets being developed by data brokers, thus potentially enrolling them in pernicious practices such as dataveillance, social sorting, control creep, and anticipatory governance, for which they were never intended. Small data studies might continue to be a vital component of the research landscape, but their position and role within it are thus changing.

Small data versus big data

The distinction between small and big data is a recent one. Prior to 2008, data were rarely considered in terms of being ‘small’ or ‘big’. All data were, in effect, what is now sometimes referred to as ‘small data’ regardless of their volume. Due to factors such as cost, resourcing, and the difficulties of generating, processing, analyzing and storing data, limited volumes of

Table 1 Comparing small and big data

Characteristic	Small data	Big data
Volume	Limited to large	Very large
Exhaustivity	Samples	Entire populations
Resolution and indexicality	Coarse and weak to tight and strong	Tight and strong
Relationality	Weak to strong	Strong
Velocity	Slow, freeze-framed	Fast
Variety	Limited to wide	Wide
Flexible and scalable	Low to middling	High

high quality data were produced through carefully designed studies using sampling frameworks designed to ensure representativeness. In the last decade or so, small data have been complemented by what has been termed ‘big data’, which have very different ontological characteristics (see Table 1).

As detailed in Kitchin (2013: 262), big data are:

- huge in *volume*, consisting of terabytes or petabytes of data;
- high in *velocity*, being created in or near real-time;
- diverse in *variety* in type, being structured and unstructured in nature, and often temporally and spatially referenced;
- *exhaustive* in scope, striving to capture entire populations or systems ($n = \text{all}$);
- fine-grained in *resolution*, aiming to be as detailed as possible, and uniquely *indexical* in identification;
- *relational* in nature, containing common fields that enable the conjoining of different data sets;
- *flexible*, holding the traits of extensionality (can add new fields easily) and scalability (can expand in size rapidly).

(Boyd and Crawford 2012; Dodge and Kitchin 2005; Marz and Warren 2012; Mayer-Schonberger and Cukier 2013).

The term ‘big’ then is somewhat misleading as big data are characterized by much more than volume. Indeed, some ‘small’ datasets can be very large in size, such as national censuses that also seek to be exhaustive and have strong resolution and relationality. However, census datasets lack velocity (usually conducted once every 10 years), variety (usually c.30 structured questions), and flexibility (once a census is set and is being administered it is all but impossible to

tweak the questions or add new questions or remove others and generally the fields are fixed, typically across censuses, to enable time-series analysis; Kitchin 2014a). Other small datasets also consist of a limited combination of big data's characteristics. For example, a qualitative dataset such as interview transcripts are usually relatively small in size (perhaps a couple of dozen respondents), have a non-continuous temporality (one-off interviews or a sequence over a number of months), possess weak relationality, and are limited in variety (text transcripts), though they have strong resolution and flexibility.

In contrast, big data have all these characteristics, or nearly all depending on their form (for example, sensor data lack variety but have the other characteristics), with the crucial qualities being velocity and exhaustivity. The rapid growth of big data has arisen due to the simultaneous development of a number of enabling technologies, infrastructures, techniques and processes, and their rapid embedding into everyday business and social practices and spaces, such as fixed and mobile internet, the embedding of computation into all kinds of objects, machines and systems that are networked together, advances in database design (especially the creation of NoSQL databases), new forms of social media and online interactions and transactions, and new kinds of data analytics designed to cope with data abundance as opposed to data scarcity (Kitchin 2013). Indeed, the practices of everyday life and the places in which we live are now augmented, monitored and regulated by dense assemblages of data-enabled and data-producing infrastructures and technologies, such as traffic and building management systems, surveillance and policing systems, government databases, customer management and logistic chains, financial and payment systems, and locative and social media (Kitchin and Dodge 2011). Within these socio-technical systems much of the data generation is automated through algorithmically-controlled cameras, sensors, scanners, digital devices such as smart phones, clickstreams, or are the by-product of networked interactions (such as the records of online transactions), or are volunteered by users through social media or crowd sourcing initiatives.

Collectively, such systems produce massive, exhaustive, dynamic, varied, detailed, indexical, inter-related, low cost per data point datasets that are flexible and scalable. To take just two examples as

way of illustration. In 2011, Facebook's active users spent more than 9.3 billion hours a month on the site (Manyika et al. 2011), and by 2012 Facebook reported that it was processing 2.5 billion pieces of content (links, stores, photos, news, etc.) and 500 + terabytes of data, 2.7 billion 'Like' actions and 300 million photo uploads *per day* (Constine 2012), each accompanied by associated metadata. Walmart was generating more than 2.5 petabytes of data relating to more than 1 million customer transactions *every hour* in 2012. These data are very different to traditional small data, consisting of a rapid, continuous torrent of highly resolute, indexical, relational and scalable data. Whereas small datasets were largely oases of data within data deserts, big data produce a veritable data deluge that seemingly enable research to shift from: "data-scarce to data-rich; static snapshots to dynamic unfoldings; coarse aggregation to high resolution; relatively simple hypotheses and models to more complex, sophisticated simulations and theories" (Kitchin 2013: 263).

These promises of big data potentially threaten the status of small data studies by positioning big data as being of more value and utility to the academy and business. However, such a framing misunderstands both the nature of big data and the value of small data. Big data may seek to be exhaustive, but as with all data they are both a representation and a sample. What data are captured is shaped by:

- the field of view/sampling frame (where data capture devices are deployed and what their settings/parameters are; who uses a space or media, e.g., who belongs to Facebook or shops in Walmart);
- the technology and platform used (different surveys, sensors, lens, textual prompts, layout, etc. all produce variances and biases in what data are generated);
- the context in which data are generated (unfolding events mean data are always situated with respect to circumstance);
- the data ontology employed (how the data are calibrated and classified), and;
- the regulatory environment with respect to privacy, data protection and security (Kitchin 2013, 2014b).

Indeed, all data provide oligoptic views of the world: views from certain vantage points, using

particular tools, rather than an all-seeing, infallible god's eye view (Haraway 1991; Amin and Thrift 2002). As such, big data constitute a 'series of partial orders, localised totalities, with their ability to gaze in some directions and not others' (Latour cited in Amin and Thrift 2002: 92). Big data undoubtedly strive to be more exhaustive and provide dynamic, fine-grained insight but, nonetheless, their promise can never be fully fulfilled. Big data generally capture what is easy to ensnare—data that are openly expressed (what is typed, swiped, scanned, sensed, etc.; people's actions and behaviours; the movement of things)—as well as data that are the 'exhaust', a by-product, of the primary task/output. Tackling a question through big data often means repurposing data that were not designed to reveal insights into a particular phenomenon, with all the attendant issues of such a maneuver, for example creating ecological fallacies.

In contrast, small data may be limited in volume and velocity, but they have a long history of development across science, state agencies, non-governmental organizations and businesses, with established methodologies and modes of analysis, and a record of producing meaningful answers. Small data studies can be much more finely tailored to answer specific research questions and to explore in detail and in-depth the varied, contextual, rational and irrational ways in which people interact and make sense of the world, and how processes work. Small data can focus on specific cases and tell individual, nuanced and contextual stories. Small data studies thus seek to mine gold from working a narrow seam, whereas big data studies seek to extract nuggets through open-pit mining, scooping up and sieving huge tracts of land.

These two approaches of narrow versus open mining have consequences with respect to data quality, fidelity and lineage. Given the limited sample sizes of small data, data quality—how clean (error and gap free), objective (bias free) and consistent (few discrepancies) the data are; veracity—the authenticity of the data and the extent to which they accurately (precision) and faithfully (fidelity, reliability) represent what they are meant to; and lineage—documentation that establishes provenance and fit for use; are of paramount importance (Lauriault 2012). Much work is expended on limiting sampling and methodological biases as well as ensuring that data are as rigorous and robust as possible before they are analyzed or shared. In contrast, it has been argued by some that big data

studies do not need the same standards of data quality, veracity and lineage because the exhaustive nature of the dataset removes sampling biases and more than compensates for any errors or gaps or inconsistencies in the data or weakness in fidelity (Mayer-Schonberger and Cukier 2013). The argument for such a view is that "with less error from sampling we can accept more measurement error" (p.13) and "tolerate inexactitude" (p. 16). Viewed in this way, Mayer-Schonberger and Cukier (2013: 13) thus argue "more trumps better." Of course, this presumes that all uses of big data will tolerate inexactitude, when in fact many big data applications do require precision (e.g., finance data), or at least data with measurable error parameters.

Moreover, the warning "garbage in, garbage out" still holds. Big datasets that generate dirty, gamed or biased data, or data with poor fidelity, are going to produce analysis and conclusions that have weakened validity and deliver fewer benefits to those that analyze and seek to exploit them. And by dint of their method of production big data can suffer from all of these ills. The data can be dirty through instrument error or biased due to the demographic being sampled (e.g., not everybody uses Twitter) or the data might be gamed or faked through false accounts or hacking (e.g., there are hundreds of thousands of fake Twitter accounts seeking to influence trending and direct clickstream trails; Bollier 2010; Crampton et al. 2012). With respect to fidelity there are question marks as to the extent to which social media posts really represent peoples' views and the faith that should be placed on them. Manovich (2011: 6) warns that "[p]eoples' posts, tweets, uploaded photographs, comments, and other types of online participation are not transparent windows into their selves; instead, they are often carefully curated and systematically managed."

There are issues of access to both small and big data. Small data produced by academia, public institutions, non-governmental organizations and private entities can be restricted in access, limited in use to defined personnel or available for a fee or under license. Increasingly, however, public institution and academic data are becoming more open. Big data are, with a few exceptions such as satellite imagery and national security and policing, mainly produced by the private sector. Access is usually restricted behind pay walls and proprietary licensing, limited to ensure competitive advantage and to leverage income

through their sale or licensing (CIPPIC 2006). Indeed, it is somewhat of a paradox that only a handful of entities are drowning in the data deluge (boyd and Crawford 2012) and companies such as mobile phone operators, app developers, social media providers, financial institutions, retail chains, and surveillance and security firms are under no obligations to share freely the data they collect through their operations. In some cases, a limited amount of the data might be made available to researchers or the public through application programming interfaces (APIs). For example, Twitter allows a few companies to access its firehose (stream of data) for a fee for commercial purposes (and have the latitude to dictate terms with respect to what can be done with such data), but researchers are restricted to a ‘gardenhose’ (c. 10 % of public tweets), a ‘spritzer’ (c. 1 % of public tweets), or to different subsets of content (‘white-listed’ accounts), with private and protected tweets excluded in all cases (boyd and Crawford 2012). The worry is that the insights that privately owned and commercially sold big data can provide will be limited to the business sector, or maybe only opened to a privileged set of academic researchers whose findings cannot be replicated or validated (Lazer et al. 2009).

Given these limitations of big data and the strengths of small data, small data studies will continue to be an important element of the research landscape. However, such data will increasingly come under pressure to utilize the new archiving technologies, being scaled-up within digital data infrastructures in order that they are preserved for future generations, become accessible to reuse and combination with other small and big data, and more value and insight can be extracted from them through the application of big data analytics. Considerable resources have already been invested in creating such data infrastructures. In the remainder of this paper we examine the scaling of small data into data infrastructures and the implications of such a scaling with respect to exposing small data to new big data epistemologies and repurposing, focusing on spatial data examples.

Pooling, scaling, preserving, sharing and reusing small data: creating data infrastructures

Data have been collected together and stored for much of recorded history. Such practices have been both

informal and formal in nature. The former consists simply of gathering data and storing them, whereas the latter consists of a set of curatorial practices and institutional structures designed to ensure that data are preserved for future generations. The former might best be described as data holdings, or backups, whereas the latter are data archives. Archives are formal collections of data that are actively structured, curated and documented, are accompanied by appropriate metadata, and where preservation, access and discoverability are integrated into technological systems and institutions designed to last the test of time (Lauriault et al. 2013). Archives explicitly seek to be long term endeavours, preserving the full record set—data, metadata and associated documentation—for future reuse.

The ability to store data digitally and to structure them within databases has radically transformed the volume of data that can be stored and efficiently and effectively handled and queried and has enabled the creation of extensive digital holdings and archives. Such digital data can be easily shared and reused for a low marginal cost, although the cost of both the soft (institutional, policies, standards, human resources) and hard (technology, servers, software, delivery mechanisms, portals) infrastructures are not in the least bit inexpensive. Moreover, these data can be manipulated and analyzed by exposing them to computational algorithms. As such, procedures and calculations that would be difficult to undertake by hand or using analogue technologies become possible in just a few microseconds, enabling more and more complex analysis to be undertaken or the replication of objects (i.e., an atlas) and results. Further the data can also be relatively easily linked together and scaled into other forms of data infrastructure.

A data infrastructure is a digital means for storing, sharing and consuming data across networked technologies. Over the past two decades in particular, considerable effort has been expended on developing and promoting data access and discovery infrastructures, which take a number of forms: catalogues, directories, portals, clearinghouses and repositories (Lauriault et al. 2007). These terms are often used interchangeably and are confused for one another, though they are slightly different types of entities. Catalogues, directories and portals are centralized resources that may detail and link to individual data archives (e.g., Earth Observation Data Management

Service of the Canada Centre for Remote Sensing) or data collections held by individual institutions (e.g. Australian National Data Service) or are federated infrastructures which provide the means to access the collections held by many (e.g., US National Sea Ice Data Center). They might provide fairly detailed inventories of the datasets held, and may act as metadata aggregators but do not necessarily host the data (e.g., GeoConnections Discovery Portal; European; O’Carroll et al. 2013). Single site repositories host all the data sets in a single site, accessible through a web interface, though they may maintain backup or mirror sites in multiple locations (e.g., The UK Data Archive). A federated data repository or clearing house can be a shared place for storing and accessing data [e.g. US National Database for Autism Research (NDAR), NASA’s Global Change Master Directory]. It might provide some data services in terms of search and retrieval, and data management and processing, but each holding or archive has been produced independently and may not share data formats, standards, metadata, and policies. Nevertheless, the repository seeks to ensure that each archive meets a set of requirement specifications and uses audit and certification to ensure data integrity and trust amongst users (Dasish 2012).

A cyber-infrastructure is more than a collection of digital archives and repositories. It consists of a suite of dedicated networked technologies, shared services (relating to data management and processing), analysis tools such as data visualizations (e.g., graphing and mapping apps), and shared policies (concerning access, use, IPR, etc.) which enable data to be distributed, linked together and analyzed (e.g. a spatial data infrastructure; Cyberinfrastructure Council 2007). Whilst it is sometimes used to denote the infrastructure that enables a federated repository to function, here we use it to denote a data infrastructure in which data share common technical specifications relating to formats, standards, and protocols. In other words, there are strong rules relating to data standardization and compliance within the infrastructure. Such cyber-infrastructures include those implemented by national statistical agencies and national spatial data infrastructures (SDIs) that require all data stored and shared to comply with defined parameters in order to maximize data interoperability and ensure data quality, fidelity and integrity that promotes trust. The objectives of SDIs are to ensure that users from

multiple sectors and jurisdictions can seamlessly re-use these data and link them into their systems. A cross border natural disaster, for instance, would require multiple agencies, in different countries along with sub-national entities, under severe time constraints and pressures, to access, model and visualize spatial data in near real time while also inputting newly acquired data to respond to and inform an emergency response arena. In less stressful environments, SDIs enable the management of cross border shared services and natural resources (e.g. EU Water Framework Directives).

The rationale for scaling small data into data infrastructures

The arguments for the storing, sharing and scaling of data within repositories and across data infrastructures centre on the promises of new discoveries and innovations through the combination of datasets and the crowdsourcing of minds. Individual datasets are valuable in their own right, but when combined with other datasets or examined in new ways fresh insights can potentially be discerned and new questions answered (Borgman 2007). By combining datasets, it is contended that the cumulative nature and pace of knowledge building is accelerated (Lauriault et al. 2007). Moreover, by preserving data over time it becomes possible to track trends and patterns, and the longer the record, the greater the ability to build models and simulations and have confidence in the conclusions drawn (Lauriault et al. 2007). Over time then, the cumulative value of data infrastructures increases as the data become more readily and broadly available, both in scope and temporality. Such a sharing strategy is also more likely to spark new interdisciplinary collaborations between researchers and teams and to foster enhanced skill through having access to new kinds of data (Borgman 2007). Moreover, the sharing of data and the adoption of infrastructure standards, protocols and policies increases data quality and enables third party data and study verification, thus increasing data integrity (Lauriault et al. 2007).

The financial benefits of data infrastructures centre on the scales of economy created by sharing resources and avoiding replication, the leveraging effects of re-using costly data, the generation of wealth through new discoveries, and producing more efficient

societies. Research and the production of administrative, statistical and geomatics data are typically costly undertakings, with various funding agencies collectively spending billions of dollars every year to fund research activities. Rather than creating a plethora of ad hoc archives, it makes more sense to establish a smaller number of dedicated institutional repositories or infrastructures which undertake basic data standardization and produce significant efficiencies in effort, as well as enable broader access to data for individual researchers/institutions where entry costs to a field would normally be prohibitive (Fry et al. 2008). As well as reducing wastage, preserving and sharing the fruits of such endeavors is more likely to maximize the return on investment by enabling as much value as possible to be extracted from the data (Lauriault et al. 2007). That said, the sustainability of these research data infrastructures are often an issue as these are funded through a mix of mechanisms such as state and research funds, community based organization infrastructures rely on small grants and membership fees, while the open data infrastructures run by civil society organizations are built by volunteers. SDIs, alternatively are funded by national and sub-national governments to ensure that all sectors and jurisdictions can seamlessly interoperate and build upon and access the same framework datasets. This allows for a decentralized and distributed data infrastructure that enables the linking of thematic datasets from multiple sources, and ensuring that these are managed by their producers, but done so in such a way that they can be combined when necessary.

Given the anticipated gains from sharing data, over the past three decades supranational bodies such as the European Union, national governments, research agencies, philanthropic and civil society organizations, have invested extensively in funding a wide variety of data and cyber-infrastructure initiatives.

Some example data infrastructures

Spatial data infrastructures (SDIs) are the archetype cyber-infrastructure. National scale SDIs are normally institutionally located in national mapping organizations, national surveys, or the departments that manage natural resources. They are an assemblage of institutions (e.g., government, geomatics), policies (e.g., data sharing protocols), laws (e.g. licenses, legislation, regulation), technologies (e.g., data

portals, storage, software), processes (e.g., web mapping, metadata aggregation), standards (e.g., metadata, file transfer, data quality) and specifications (e.g., interoperability), scientific and computing knowledge, skilled human resources, discovery and access portals, framework data (e.g., common datasets upon which others can build such as road networks) and mapping services that direct the who, how, what and why geospatial data are collected, stored, manipulated, analyzed, transformed and shared. They are intersectoral, cross-domain, trans-disciplinary, interdepartmental, and require much consensus building. Supranational SDIs such as the Infrastructure for Spatial Information in the European Community (INSPIRE), are very similar to national SDIs, however in the case of INSPIRE, it governs how nations are to construct their infrastructures via rules, directives, and policies that will lead to data, geomatic systems and services being seamlessly interoperable across 27 member states. Each SDI, irrespective of its scale and jurisdiction, is therefore unique, but by adhering to a shared set of standards, policies and technologies they can be joined up. In addition, INSPIRE includes a GeoPortal which is a federated catalog that aggregates the metadata of member state SDIs thus providing users with a single point to discover and view EU geospatial data.

On a smaller scale, and in a different domain, the UK Data Archive, is an example of a research data infrastructure that acquires, curates and provides access to social science and humanities data. Data are discovered via the UK Data Service which is a catalogue that provides access to hosted national and international survey data collections, international databanks, census data and qualitative data. Secure data services for access and use of more sensitive research data are also provided. Data are described with standard metadata, and a number of educational resources are provided for users to work with the data once they have been downloaded. Although not a certified trusted digital repository, the UK Data Archive aims to maintain its large collections of data for long-term reuse, and provides a number of capacity building resources to enable researchers to manage and deposit their data.

There are not many examples of data infrastructures in the non-profit and charitable sector. The Canadian Council on Social Development, Community Data Program (CDP), is however an example of a small data

infrastructure created for the specific purpose of enabling small area, evidence based decision making in the social sector. It is funded by its members through a consortia model. Members are city based networks of municipal administrators, school boards, community health centres, social planning councils and a number of charitable and non-profit organizations. The CDP acquires and disseminates mostly public sector data and custom ordered cross tabulated data aggregated into neighbourhood, city ward, small area census geographies and postal codes. These are stored into a database and delivered to members via an online catalog. In this instance, members not only benefit from the data, but also from services where experts negotiate data acquisition based on community needs and specifications, and a knowledge sharing network between super users and novices.

Finally, since 2009 open data infrastructures have been created by national governments, sub national governments such as cities, provinces, counties and states, and civil society organizations such as the UK-based Open Knowledge Foundation (OKF), and to a lesser extent research and private sector entities. The objectives of these data infrastructures are to unlock access to public sector datasets and make them accessible via a discovery and access portals for free and under open licences. The OKF is an open data supranational organization which provides direction to governments and civil society groups and helps build capacity in terms of the deployment of catalogs (e.g., CKAN), and has created a set of open data principles and open license specifications. Open data portals have not yet matured into cyber-infrastructure, although government funded open data portals do manifest some of their qualities. Unlike SDIs, these are not grounded in a domain, discipline or the sciences, and often open data infrastructures are administered in information management/technology departments and championed by chief technology officers, or are created and supported by volunteer groups composed of new media enthusiasts and app developers.

These four cases are but a small sample of the innumerable data and cyber-infrastructure currently in operation. In all four cases, the data found in their portals are small data, SDIs being the exception as remote sensing data and many environmental sensors produce data that have the qualities of big data. Alternatively, geodemographic data infrastructures,

discussed later, exemplify the scaling of small data with big data.

Making small data more big data-like

Whilst the scaling of small data into data infrastructures does not create big data, in the sense that the data still lack velocity and exhaustivity, it does make them more big data-like by making them more extensive, relational and interconnected, varied, and flexible. This enables two effects to occur. First, it opens scaled small data to new epistemologies and, in particular, to new forms of big data analytics (Kitchin 2014a). Second, it facilitates small data being conjoined with big data to produce more complex, inter-related and wide-ranging data infrastructures that are presently driving the rapid growth of commercial data brokers, including the burgeoning geodemographics industry (also known as locational targeted niche marketing tools). Both have consequences with respect to how small data are being used and raise normative questions concerning the creation and use of data infrastructures.

Exposing small data to new epistemologies

Traditional small data methods of analysis have primarily been designed to extract insights from scarce, static, clean and weak relational data sets that have been sampled and adhere to strict assumptions (such as independence, stationarity, and normality), and were generated and analyzed with a specific question in mind (Miller 2010). The challenge with big data is to cope with abundance and exhaustivity (including sizable amounts of data with low utility and value), timeliness and dynamism, messiness and uncertainty, high relationality, semi-structured or unstructured content, and the fact that much of them are generated with no specific question in mind or are a by-product of another activity. The solution has been new data analytics that utilize the power of algorithms and computation to process and provide insight into datasets that would simply be too costly, difficult and time-consuming to analyze otherwise. Such analytics scale-up existing statistical methods, such as regression, model building, data visualization and mapping, as well as employing new machine learning and visual analytics techniques that computationally mine meaning from data and detect, classify and segment

meaningful patterns, relationships, associations and trends between variables, and build predictive, simulation and optimization models (Han et al. 2011; Hastie et al. 2009). These data analytics can equally be applied to scaled small data to extract and model insights.

Data analytics are reflective of a particular way of making sense of the world; they are the manifestation of a particular epistemology. Some envisage them as a new form of empiricism that enables data to speak for themselves free of theory. For example, Anderson (2008) argues that “the data deluge makes the scientific method obsolete”. He continues, “We can analyze the data without hypotheses about what it might show. We can throw the numbers into the biggest computing clusters the world has ever seen and let statistical algorithms find patterns where science cannot... Correlation supersedes causation, and science can advance even without coherent models, unified theories, or really any mechanistic explanation at all.” In other words, rather than testing whether certain hypothesized patterns or relationships exist within a dataset, algorithms are set to work on big data to discover meaningful associations between data without being guided by hypotheses. In this epistemological vision, scaled small data are made sense of through a purely inductive approach.

In contrast, data-driven science seeks to hold to the tenets of the scientific method, but uses a combination of abductive, inductive and deductive approaches to advance the understanding of a phenomenon (Kitchin 2014a). It differs from the traditional deductive approach in that it seeks to generate hypotheses and insights ‘born from the data’ rather than ‘born from the theory’ (Kelling et al. 2009: 613). It thus seeks to incorporate induction into the initial stages of the research design guided by abduction (logical inference and reasoning based on established theory), though explanation through induction is not the intended endpoint. Here, the patterns, associations and trends identified through initial data analytics are used to identify potential hypotheses worthy of further examination and testing. As such, the epistemological strategy adopted within data-driven science is to use guided knowledge discovery techniques to identify valuable insights that traditional ‘knowledge-driven science’ might fail to spot and then to investigate these further (Kelling et al. 2009; Miller 2010; Loukides 2010).

With respect to the social sciences and humanities, data infrastructures, new data analytics and associated epistemologies offer the potential to transform the research landscape (Kitchin 2013, 2014a; Ruppert 2013). As noted, data infrastructures provide access to large collections of data for reuse and analysis. These data can be conjoined in new ways and the relationships and associations between them explored using data analytics. With respect to structured data, it becomes possible to produce more refined and sophisticated models and to test the veracity of these models across a multitude of groups, settings and situations (Lazer et al. 2009). This includes the production of more elaborate and robust spatial models (Batty 2013). The volume of unstructured digital data is multiplying rapidly, including access to new sources of information (e.g., social media) and sources which have heretofore been difficult to access (e.g., millions of books, documents, newspapers, photographs, art works, and material objects; Cohen 2008). These data are opened up to the power of computation, including sophisticated tools for handling, searching, linking, sharing and analyzing data that seek to complement and augment existing humanities methods and traditional forms of interpretation and theory building (Berry 2011; Manovich 2011), as well utilizing new data analytics that provide new means to make sense of such data (Moretti 2005). Typically humanities research has progressed by providing a close reading of a handful of sources, however new machine learning techniques mean that thousands of sources can be mined, graphed and mapped, finding patterns and insights that an individual would find difficult to spot without the help of ‘reading machines’ (Ramsay 2010).

Such approaches are not without critique, with detractors arguing that data analytics are mechanistic, reductionist, functionalist, and parochial, reducing diverse individuals and complex, multidimensional social structures to mere data points (Wyly 2014), thus fostering weak, surface analysis, rather than deep, penetrating insight; that they sacrifice specificity, context and depth for scale, automation and breadth. Indeed, Brooks (2013) contends that data analytics: struggle with the social (people are not rational and do not behave in predictable ways; human systems are incredibly complex, having contradictory and paradoxical relations); and with context (data are largely shorn of the social, political and economic and

historical context); create bigger haystacks (consisting of many more spurious correlations making it difficult to identify needles); have trouble addressing big problems (especially social and economic ones); favor memes over masterpieces (identifies trends but not necessarily significant features that may become a trend); and obscure values (of the data producers and those that analyze them and their objectives). Such debates over the value and appropriateness of new analytics and epistemologies, and their application to scaled small data, seem set to continue for the foreseeable future (Kitchin 2014a).

Normative concerns related to scaled small data

Scaled small data also gain in value as a commodity, especially when they can be conjoined with big data. In contrast to academic, research-orientated or governmental data infrastructures, data brokers (sometimes called data aggregators, consolidators or resellers) gather together data into privately held infrastructures for resale on a for-profit basis. They source data from both public and private sources. For example, from public sector sources they gather data relating to individuals and aggregates (e.g., groups, places) concerning health, education, crime, property, travel, environment, etc., matching these with private sector data related to or captured within retail, financial, logistics, business intelligence, real estate, private security, political polling, transportation, media, and so on. The potential to link data across domains is high. For example, the Dutch Data Protection Authority estimates that the average Dutch citizen is included in 250–500 databases, with more socially active people included in up to 1,000 databases (Koops 2011). More recently, data brokers have been combining these data with the metadata and content from locative (e.g., smart phone apps) and social media (e.g. Twitter and Facebook). For example, Facebook is partnering with large data brokers and marketers in order to merge together the profiles, networks and uploaded content of its billion users (their likes, comments, photos, videos, etc.) with non-Facebook purchasing and behaviour data (Edwards 2013).

These interconnected data infrastructures bind together a vast array of personal data and are used to construct a suite of derived data products, wherein value is added through integration and data analytics,

creating profiles of individuals, groups and places, and predictions as to what people might do under different circumstances. In the main, profiles are used to micro-target advertising and niche marketing campaigns, assess how such targets might behave and be nudged into a particular response (e.g., selecting and purchasing a particular item), assess credit worthiness and socially sort individuals (determine whether one might receive a service or set personalized pricing), and provide detailed business analytics, whilst reducing their overheads in terms of wastage and loss through risky investments (Lyon 2002; Graham 2005; Siegel 2013). Acxiom, for instance, seeks to mesh offline, online and mobile data in order to create a ‘360-degree view’ of consumers, using these data to create detailed profiles and robust predictive models which it sells to interested parties (Singer 2012).

Geodemographic segmentation is a data analytical process which can combine both small and big data in order to create quantitatively based classification systems of groups of people at a particular geographic unit of analysis, often at postal code geographies. Once classification systems are developed, primarily with small data inputs, big data such as purchasing histories, which use postal codes as unique identifiers, can be matched to these classifications to assess consumption patterns and to refine the groupings. These data infrastructures, while they can be used to better understand population dynamics in cities, are mostly developed by the private sector to geo-target marketing. As an illustration, the Environics Analytics PRiZMC2 segmentation tool classifies Canadians into 66 lifestyle types such as ‘cosmopolitan elite’ or ‘Les Chics and Lunch at Tim’s’ (short for Tim Horton Donuts) “based on their demographics, marketplace preferences and psychographic Social Values”. This company also produces a product called WealthScapes Dollar and Sense which provides marketers with a similar service (Environics Analytics 2013a, b). The algorithms, methodological assumptions and the mix of datasets used to produce the geodemographic profiles are proprietary and protected by intellectual property regimes and are not subject to public scrutiny. Irrespective, by using such products companies seek to become more effective and efficient in their operations with respect to targeting customers and siting stores.

The scaling of small data, mashing them with big data, and subjecting them to data analytics, can have

profound implications for citizens and the services and opportunities extended to them. The worry for some is that a form of ‘data determinism’ is being practiced in which individuals are not profiled and judged just on the basis of what they have done, but on the prediction of what they might do in the future (Ramirez 2013). A new probability market is emerging—although gambling industry odds compilers and security markets have been around for some time—which constitutes a new phase in the era of probabilistic thinking (Hacking 1975, 1990), one that is making up new kinds of people (Hacking 2007) and new kinds of places (Lauriault 2012), led by the private sector and surveillance institutions, mostly for the purpose of marketing products and security. Moreover, there are concerns regarding the extent to which scaled small data and data infrastructures facilitate dataveillance (surveillance enacted through the processing and analyzing of data records), infringe on privacy and other human rights, affect access to private health insurance and its rates, stigmatize and redline areas, pose significant data security concerns with regards to data being stolen and exploited criminally, and enable control creep wherein data generated for one purpose is used for another (Clarke 1988; Innes 2001; Solove 2006; CIPPIC 2006). Citizens may have not agreed with the entities producing the data as to how data about themselves are used (CIPPIC 2006). As such, whilst scaling small data does offer a number of benefits they also can have differential and negative consequences. There are thus a number of fundamental normative questions that need urgent reflexive consideration concerning the production of data infrastructures if we are to maximize their benefits whilst minimizing their more pernicious effects.

Conclusion

We are presently witnessing a fast changing landscape with respect to data. Not only are we witnessing the roll-out of a new form of data in the guise of big data, but traditional small data are evolving through new data infrastructures that enable them to be scaled and analyzed in new ways. In this paper we have compared small and big data before going on to examine how small data are being scaled, combined with big data, and being made amenable to big data analytics. Our argument has been three fold.

First, despite the rapid growth of big data and associated new analytics, small data will continue to be a vital part of the research landscape. There will not be a paradigm shift in the near future in which studies using big data replace those employing small data, rather small and big data will complement one another; mining narrow seams of high quality data will continue alongside open pit mining because it enables much more control of the research design and to answer specific, targeted questions. As such, rather than directing research funding to projects that have access to vast quantities of data in the hope that they will inherently produce useful insights, funding needs to be focused on answering critical questions, whether they are tackled using small or big data (Sawyer 2008).

Second, the small data landscape is changing through the development of data infrastructures. Small data gain value and utility when made accessible for reuse and are combined with other datasets. As a consequence, much effort is being directed at building such infrastructures and in trying to harmonize small data, with respect to data standards, formats, metadata, and documentation, to ensure their compatibility with systems, maximize discoverability, and facilitate the linking together of datasets. The pressure to harmonize, share and reuse small data will continue to grow as research funders seek to gain the maximum return on their investment through new knowledge and innovations.

Third, the scaling of small data into data infrastructures has three consequences. One: by pooling and linking small data to create larger, interconnected datasets, small data are opened up to analysis by big data analytics. Small data are thus exposed to the new epistemologies of data science, fostering the growth of new approaches such as the digital humanities and computational social sciences. Two: small data are more easily conjoined with big data to produce more diverse derived data that enables more wide-ranging and extensive analysis. This reconfiguration of the data landscape is facilitating the rapid growth of data brokers and new data products, including detailed profiling. Three: the scaling of small data, and their combination with big data and exposure to big data analytics, produces a set of potential pernicious effects such as dataveillance, social sorting, control creep, and anticipatory governance that impinge on privacy, social freedoms and have structural consequences for individual lives. As such, the scaling of small data

raises normative questions concerning how data should be managed and utilized. We have barely begun to examine these consequences, with developments running ahead of critical and normative reflection and political, policy and legal reaction.

Small data are set to continue being an important component of research endeavors. However, they are in the process of taking on new forms that have consequences for how we think about and utilize such data. We have made an initial attempt to detail some of these transformations, but further critical reflection and normative thinking is required to make sense of the changes taking place and their implications.

Acknowledgments The research conducted for this paper was made possible with funding from the European Research Council (ERC-2012-AdG-323636) and Science Foundation Ireland.

References

- Amin, A., & Thrift, N. (2002). *Cities: Reimagining the urban*. London: Polity.
- Anderson, C. (2008). The end of theory: The data deluge makes the scientific method obsolete. *Wired*, June 23, 2008, http://www.wired.com/science/discoveries/magazine/16-07/pb_theo-ry. Accessed 12 Oct 2012.
- Batty, M. (2013). *The new science of cities*. Cambridge, MA: MIT Press.
- Berry, D. (2011). The computational turn: Thinking about the digital humanities. *Culture Machine* 12. <http://www.culturemachine.net/index.php/cm/article/view/440/470>. Accessed 3 Dec 2012.
- Bollier, D. (2010). *The promise and peril of big data*. The Aspen Institute. http://www.aspeninstitute.org/sites/default/files/content/docs/pubs/The_Promise_and_Peril_of_Big_Data.pdf. Accessed 1 Oct 2012.
- Borgman, C. L. (2007). *Scholarship in the digital age*. Cambridge, MA: MIT Press.
- boyd, D., & Crawford, K. (2012). Critical questions for big data. *Information, Communication and Society*, 15(5), 662–679.
- Brooks, D. (2013). What data can't do. *New York Times*, <http://www.nytimes.com/2013/02/19/opinion/brooks-what-data-cant-do.html>. Accessed 18 Feb 2013.
- Canadian Internet Public Policy Interest Clinic (CIPPIC). (2006). *On the data trail: How detailed information about you gets into the hands of organizations with whom you have no relationship*. Ottawa: A Report on the Canadian Data Brokerage Industry. <https://www.cippic.ca/sites/default/files/May1-06/DatabrokerReport.pdf>.
- Clarke, R. (1988). Information technology and dataveillance. *Communications of ACM*, 31(5 May 1988), 498–512.
- Cohen, D. (2008). Contribution to: The promise of digital history (roundtable discussion). *Journal of American History*, 95(2), 452–491.
- Constine, J. (2012). *How big is facebook's data? 2.5 billion pieces of content and 500 + terabytes ingested every day*, 22 August 2012, <http://techcrunch.com/2012/08/22/how-big-is-facebooks-data-2-5-billion-pieces-of-content-and-500-terabytes-ingested-every-day/>. Accessed 28 Jan 2013.
- Crampton, J., Graham, M., Poorthuis, A., Shelton, T., Stephens, M., Wilson, M. W., et al. (2012). *Beyond the Geotag? Deconstructing "big data" and leveraging the potential of the geoweb*. http://www.uky.edu/~tmute2/geography_methods/readingPDFs/2012-Beyond-the-Geotag-2012.10.01.pdf. Accessed 21 Feb 2013.
- Cyberinfrastructure Council. (2007). *Cyberinfrastructure vision for 21st century discovery*. <http://www.nsf.gov/pubs/2007/nsf0728/index.jsp?org=ECC> Washington, DC: National Science Foundation. Accessed 17 Jan 2014.
- Dasish. (2012). *Roadmap for preservation and curation in the social sciences and humanities*. http://dasish.eu/publications/projectreports/D4.1_-_Roadmap_for_Preservation_and_Curation_in_the_SSH.pdf. Accessed 15 Oct 2013.
- Dodge, M., & Kitchin, R. (2005). Codes of life: Identification codes and the machine-readable world. *Environment and Planning D: Society and Space*, 23(6), 851–881.
- Edwards, J. (2013). Facebook is about to launch a huge play in 'big data' analytics. *Business insider*, May 10th <http://www.businessinsider.com/facebook-is-about-to-launch-a-huge-play-in-big-data-analytics-2013-5>. Accessed 18 Sept 2013.
- Environics Analytics. (2013a). *WealthScapes: Dollars and sense*, <http://www.environicsanalytics.ca/environics-analytics/data/financial-data/wealthscapes>. Accessed 26 Nov 2013.
- Environics Analytics. (2013b). *PRiZMc2 segmentation lifestyle lookup tool*, <http://www.environicsanalytics.ca/prizm-c2-cluster-lookup>. Accessed 26 Nov 2013.
- Fry, J., Lockyer, S., Oppenheim, C., Houghton, J. W., & Rasmussen, B. (2008). *Identifying benefits arising from the curation and open sharing of research data produced by UK higher education and research institutes*. London and Bristol: JISC. <http://repository.jisc.ac.uk/279/>. Accessed 8 Oct 2014.
- Graham, S. (2005). Software-sorted geographies. *Progress in Human Geography*, 29(5), 562–580.
- Hacking, I. (1975). *The emergence of probability*. Cambridge: Cambridge University Press.
- Hacking, I. (1990). *The taming of chance*. Cambridge: Cambridge University Press.
- Hacking, I. (2007). Kinds of people, moving targets. In *Proceedings of the British Academy* (Vol. 151, pp. 285–318), 2006 Lectures. British Academy Lecture, Read at the Academy 11 April 2006.
- Han, J., Kamber, M., & Pei, (2011). *Data mining: Concepts and techniques* (3rd ed.). Waltham: Morgan Kaufmann.
- Haraway, D. (1991). *Simians, cyborgs and women: The reinvention of nature*. New York: Routledge.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd edition ed.). Berlin: Springer.
- Innes, M. (2001). Control creep. *Sociological Research Online*, 6(3), <http://www.socresonline.org.uk/6/3/innes.html>. Accessed 8 Oct 2014.
- Kelling, S., Hochachka, W., Fink, D., Riedewald, M., Caruana, R., Ballard, G., et al. (2009). Data-intensive science: A new

- paradigm for biodiversity studies. *BioScience*, 59(7), 613–620.
- Kitchin, R. (2013). Big data and human geography: Opportunities, challenges and risks. *Dialogues in Human Geography*, 79(1), 1–14.
- Kitchin, R. (2014a). Big data, new epistemologies and paradigm shifts. *Big Data and Society*, 1(1), 1–12.
- Kitchin, R. (2014b). The real-time city? Big data and smart urbanism. *GeoJournal*, 3(3), 262–267.
- Kitchin, R., & Dodge, M. (2011). *Code/space: Software and everyday life*. Cambridge, MA: MIT Press.
- Koops, B. J. (2011). Forgetting footprints, shunning shadows: A critical analysis of the ‘right to be forgotten’ in big data practice. *SCRIPTed*, 8(3), 229–256.
- Laney, D. (2001). 3D data management: Controlling data volume, velocity and variety. *Meta Group*. <http://blogs.gartner.com/doug-laney/files/2012/01/ad949-3D-Data-Management-Controlling-Data-Volume-Velocity-and-Variety.pdf>. Accessed 16 Jan 2013.
- Lauriault, T. P. (2012). *Data, infrastructures and geographical imaginations: Mapping data access discourses in Canada*. PhD Thesis, Ottawa: Carleton University.
- Lauriault, T. P., Craig, B. L., Taylor, D. R. F., & Pulsifier, P. L. (2007). Today’s data are part of tomorrow’s research: Archival issues in the sciences. *Archivaria*, 64, 123–179.
- Lauriault, T. P., Hackett, Y., & Kennedy, E. (2013). *Geospatial data preservation primer*. Arthurs and Low: Hickling.
- Lazer, D., Pentland, A., Adamic, L., Aral, S., Barabási, A.-L., Brewer, D., et al. (2009). Computational social science. *Science*, 323, 721–733.
- Loukides, M. (2010). What is data science? *O’Reilly Radar*, 2 June 2010, <http://radar.oreilly.com/2010/06/what-is-data-science.html>. Accessed 28 Jan 2013.
- Lyon, D. (2002). Everyday surveillance: Personal data and social classifications. *Information, Communication and Society*, 5, 242–257.
- Manovich, L. (2011). *Trending: The promises and the challenges of big social data*. http://www.manovich.net/DOCS/Manovich_trending_paper.pdf. Accessed 9 Nov 2012.
- Manyika, J., Chiu, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., et al. (2011). *Big data: The next frontier for innovation, competition, and productivity*. McKinsey Global Institute.
- Marz, N., & Warren, J. (2012). *Big data: Principles and best practices of scalable realtime data systems*. Manning: MEAP edition.
- Mayer-Schonberger, V., & Cukier, K. (2013). *Big data: A revolution that will change how we live*. John Murray: Work and Think.
- Miller, H. J. (2010). The data avalanche is here. Shouldn’t we be digging? *Journal of Regional Science*, 50(1), 181–201.
- Moretti, F. (2005). *Graphs, maps, trees: Abstract models for a literary history*. London: Verso.
- O’Carroll, A., Collins, S., Gallagher, D., Tang, J., & Webb, S. (2013). *Caring for digital content, mapping international approaches Nui Maynooth*. Dublin: Trinity College Dublin, Royal Irish Academy and Digital Repository of Ireland.
- Rameriz, E. (2013). The privacy challenges of big data: A view from the lifeguard’s chair. *Technology Policy Institute Aspen Forum*, August 19th. <http://ftc.gov/speeches/ramirez/130819bigdataaspen.pdf>. Accessed 11 Oct 2013.
- Ramsay, S. (2010). *Reading machines: Towards an algorithmic criticism*. Champaign, IL: University of Illinois Press.
- Ruppert, E. (2013). Rethinking empirical social sciences. *Dialogues in Human Geography*, 3(3), 268–273.
- Sawyer, S. (2008). Data wealth, data poverty, science and cyberinfrastructure. *Prometheus: Critical Studies in Innovation*, 26(4), 355–371.
- Siegel, E. (2013). *Predictive analytics*. Hoboken, NJ: Wiley.
- Singer, N. (2012). You for sale: Mapping, and sharing, the consumer genome. *New York Times*, 17th June, www.nytimes.com/2012/06/17/technology/acxiom-the-quiet-giant-of-consumer-database-marketing.html. Accessed 11 Oct 2013.
- Solove, D. J. (2006). A taxonomy of privacy. *University of Pennsylvania Law Review*, 154(3), 477–560.
- Wyly, E. (2014). Automated (post) positivism. *Urban Geography*, 35(5), 669–690.