

Big Data approaches in social and behavioral science: four key trade-offs and a call for integration

J Mahmoodi¹, M Leckelt², MWH van Zalk^{2,3}, K Geukes² and MD Back²



Big Data approaches have given rise to novel methodological tools to investigate human decisions and behaviors beyond what is possible with traditional forms of analysis. Like any other paradigm in the social and behavioral sciences, however, Big Data is not immune to a number of typical trade-offs: (1) Prediction versus explanation, pertaining to the overall research goals; (2) induction versus deduction, regarding the epistemological focus; (3) bigness versus representativeness in sampling approaches; and (4) data access versus scientific independence, addressing the forms of data usage. In this paper, we discuss these trade-offs and how Big Data and traditional approaches typically relate to them, and propose ways to overcome each trade-off by integrating advantages of different research approaches in the social and behavioral sciences with Big Data.

Addresses

¹ University of Geneva, Switzerland

² University of Münster, Germany

³ University of Oxford, United Kingdom

Corresponding author: Mahmoodi, J (jasmin.mahmoodi@unige.ch)

Current Opinion in Behavioral Sciences 2017, 18:57–62

This review comes from a themed issue on **Big data in the behavioural sciences**

Edited by **Michal Kosinski** and **Tara Behrend**

<http://dx.doi.org/10.1016/j.cobeha.2017.07.001>

2352-1546/© 2017 Elsevier Ltd. All rights reserved.

Through the exponential growth of the Internet and the large number of connected devices, people produce vast amounts of data and leave behind their digital footprints. The collection and analysis of these large amounts of information is, on the most basic level, referred to as *Big Data* (although this is a vaguely defined term; see e.g. [1]). Recently, Big Data has shown its potential to address longstanding questions in the social sciences — beyond what is possible with traditional forms of analysis [2,3] — and to open up great opportunities of investigating human experiences and behaviors [4,5]. Big Data has, for example, proven to successfully predict personality [6], movie box office success [7], stock market movements [8], and

economic welfare [9,10]. Thus, it provides wide access to unprecedented amounts of data, offers new insights into human emotions, cognitions, motivations, decisions, preferences, behaviors, and interactions, and facilitates the data-driven development of new conceptual ideas in the social sciences [11].

These are important empirical opportunities and methodological developments [12,13], and we share the enthusiasm that is contained regarding many of the newly available research options. Nonetheless, the Big Data approach is not an all-round, carefree package to the behavioral and social sciences (in the same way behaviorism, the cognitive ‘revolution’, neuroscience, or any other new paradigm that joined the mix of scientific approaches never was). Lazer *et al.* [14] even warn that ‘big data hubris’ leads to the assumption that data quantity is a substitute for knowledge-driven methodologies and theories (also see [15]). When integrating Big Data approaches in the social and behavioral sciences, one is, thus, well advised to carefully consider decades of social science and behavioral research and foundational issues of theorizing, measurement, and analysis.

Here, we discuss four broad and general trade-offs faced by researchers of all disciplines that we think are also crucial for successfully integrating Big Data approaches into social and behavioral science. These four trade-offs relate to (1) the overall goals that the research aims to achieve (prediction versus explanation), (2) the epistemological focus (induction versus deduction), (3) sampling approaches (bigness versus representativeness), and (4) the forms of data usage (data access versus scientific independence). We will, first, introduce these four trade-offs that are intrinsic to the social and behavioral sciences in general and further argue that Big Data approaches are no exception to this. That is, as any other approach, Big Data has to deal with the challenge that very often the optimization toward one end of the trade-off dimension comes at the cost of the optimization toward the other end (thus ‘trade-offs’).¹ We then, second, discuss how Big Data approaches typically relate to these trade-offs and,

¹ Of course, this is not meant in an absolute deterministic way. What we refer to with these trade-offs is the observation that in contemporary social and behavioral science (including Big Data approaches), the optimization of one of the ends of the spectrum typically reduces the probability that the other one is optimized. This does not exclude the possibility that one can develop approaches that optimize both ends (i.e., solve the trade-off). In fact, this possibility is exactly what fuels our call for integration between traditional and Big Data approaches.

third, highlight that more traditional and Big Data approaches are often complementary in that different trade-off aspects are optimized. As a consequence, strengths of one approach are often missing in the other, and the weaknesses of one approach are often (partially) solved by the other. Furthermore, given these inevitable trade-offs and the complementary nature of traditional and Big Data approaches, we argue that, fourth, more integrative approaches are necessary in optimizing research toward both ends of the spectrum (and thereby overcoming) each trade-off.

Prediction versus explanation

Prediction and *explanation* are the main goals of most scientific endeavor [16,17], which are qualitatively distinct in their approach, methodology, and statistical models, and, hence, answer different questions [18]. *Prediction* (or *engineering*) is mainly occupied with finding the best model to forecast future observations (e.g. behaviors). In contrast, *explanation* (or *science* [19^{••},20^{••}]) is about identifying and understanding causal structures and processes of what gives rise to observed phenomena (e.g. mental states or behaviors). Whereas models with high explanatory power do not necessarily possess the highest predictive power, predictive models cannot always explain the underlying (causal) mechanisms. Although prediction and explanation are not incompatible, they are hardly ever used in combination.

The Big Data approach is very successful in creating models, with which important phenomena in practically relevant contexts can be predicted; often more successful than by previous explanatory work (e.g. [21]). For example, predictive models based on Big Data approaches have forecasted commercial success [7,8] and have been used to predict psychological constructs such as personality [6,22], political orientation [23], and personal and sensitive information (e.g. age, sexual and political orientation, etc. [24]). The same models, however, do not always provide an explanation for the mechanisms that cause this prediction. As an illustrative example (see also [19^{••}]), variables that predict the performance of a hedge fund, even if theoretically implausible, can make a useful predictive model. But this same model does not shed light on the underlying factors that could explain this prediction. This example illustrates the distinction between prediction and explanation, showing that understanding underlying causal structures is not always a priority in prediction models as used by some data scientists, particularly not when the goal is to find a predictive model to maximize a company's performance and return.

An integration of Big Data and traditional approaches might help to optimize both the prediction and explanation of social and behavioral phenomena. On the one hand, “even in cases where causal explanation is the primary objective, machine learning concepts and

methods can still provide invaluable benefits when used instrumentally — by minimizing p-hacking, increasing research efficiency, facilitating evaluation of model performance, and increasing interpretability” ([20^{••}], p. 48). Specifically, prediction will foster a better understanding of human phenomena by identifying relevant antecedents that might then be targeted in process-oriented empirical work. Explanation, on the other hand, can contribute to better prediction models by specifying data-informed models and help to make them more robust and less dependent on continuous recalibration and arbitrary changes in the data systems they are based on (e.g. evolving algorithms). Hence, both outcome foci should be integrated, as the two approaches complement each other. Optimally, a productive cycle should be established that creates new theoretical insights based on prediction efforts and that uses explanatory insights to build better prediction models.

Induction (bottom-up, data-driven) versus deduction (top-down, theory-driven)

Similar to the explanation and prediction trade-off, historically, there have been extensive debates on inductive (i.e. data-driven) and deductive (i.e. theory-driven) scientific methods. Traditional deductive, that are knowledge-driven or theory-driven, approaches heavily rely on hypotheses testing (see also [15,25^{••}]). Writers such as Karl Popper [26] — one of the most influential philosophers of science — rejected inductivists' views on the scientific method and vindicated a top-down, theory-driven approach based on purely deductive logical reasoning to empirically test, criticize, and falsify theories. This deductivists' view relies on a priori hypotheses based on what (we think) we know and their critical, empirical test. Thus, it aims at circumventing loose and confirming interpretations and creation of non-robust post hoc theories that try to fit the pattern of results. A pure Popperian deductive approach is, however, also problematic since it is based on the unrealistic assumption that researchers can somehow — by pure thinking (but uninformed by empirical data, i.e., observations) — come up with relevant hypotheses. This might hinder a productive integration of inductive, data-driven findings. Ironically, it may also undermine the creation of robust theories, as it can force researchers to pretend that they already knew in advance what, in fact, was revealed by the data [27].

The accumulation of large quantities of data has brought forward a scientific practice that generates insights purely from data and stands in contrast with the more traditional deductive approaches in the social and behavioral sciences. This inductive, data-driven approach allows us to learn from actual observed and recorded (inter-)actions and behaviors in a bottom-up fashion, enabling researchers to derive theories from data. An advantage is that this avoids, for instance, that scientists fall victim to confirmation bias [28], which could result in the

Download English Version:

<https://daneshyari.com/en/article/5735720>

Download Persian Version:

<https://daneshyari.com/article/5735720>

[Daneshyari.com](https://daneshyari.com)