

Embedding Privacy and Ethical Values in Big Data Technology

Michael Steinmann¹, Julia Shuster², Jeff Collmann³, Sorin Matei⁴, Rochelle Tractenberg⁵, Kevin FitzGerald⁶, Greg Morgan⁷, Douglas Richardson⁸

¹Stevens Institute of Technology Hoboken, New Jersey email: msteinma@stevens.edu

²Georgetown University Washington, DC email: js777@georgetown.edu

³Georgetown University Washington, DC email: collmanj@georgetown.edu

⁴Purdue University West Lafayette, Indiana email: smatei@purdue.edu

⁵Georgetown University Washington, DC email: ret7@georgetown.edu

⁶Georgetown University Washington, DC email: kf3@georgetown.edu

⁷Stevens Institute of Technology Hoboken, New Jersey email: gmorgan@stevens.edu

⁸Association of American Geographers Washington, DC email: drichardson@aag.org

Table of Contents

1. Introduction	3
2. Defining Big Data.....	4
3. Ethical Issues in the Use of Big Data	4
4.1. Defining Privacy	5
4.2. Four Normative Principles as Basis for the Ethical Analysis of Privacy	6
4.3. Remarks and Explanations.....	8
The Privacy Matrix: How to Think About Privacy in Big Data.....	13
3.1. Profiling Individuals with Big Data.....	16
3.2 Anonymity, manipulation and user consent in Online Communities.....	16
3.3. Sale of Big Data in Commercial contexts.....	Error! Bookmark not defined.
3.5 Protecting vulnerable populations in educational contexts	17
3.5. Unequal Access to Big Data in Scientific Research.....	Error! Bookmark not defined.
3.6 Big Data and government surveillance	18
5. Embedding Values in Big Data Technology.....	20
5.1 Values guide the use of Big Data Technology.....	20
5.2 Big Data tools enable realization of values	21
5.3 Basing Big Data Tool design on target values.....	23
6. Conclusion.....	24

1. Introduction

The phenomenon now commonly referred to as “Big Data” holds great promise and opportunity as a potential source of solutions to many societal ills ranging from cancer to terrorism; but it might also end up as “...a troubling manifestation of Big Brother, enabling invasions of privacy, decreased civil freedoms, (and) increased state and corporate control” (Boyd & Crawford, 2012: p 664). Discussions about the use of Big Data are widespread as “(d)iverse groups argue about the potential benefits and costs of analyzing genetic sequences, social media interactions, health records, phone logs, government records, and other digital traces left by people” (Boyd&Crawford, 2012: p 662). This chapter attempts to establish guidelines for the discussion and analysis of ethical issues related to Big Data in research, particularly with respect to privacy. In doing so, it does a new dimensions to the agenda setting goal of this volume. It is intended to help researchers in all fields, as well as policy-makers, to articulate their concerns in an organized way, and to specify relevant issues for discussion, policy-making and action with respect to the ethics of Big Data. On the basis of our review of scholarly literature and our own investigations with big and small data, we have come to recognize that privacy and the great potential for privacy violations constitute major concerns in the debate about Big Data. Furthermore, our approach and our recommendations are generalizable to other ethical considerations inherent in Big Data as we illustrate in the final section of the chapter.

To understand the ethical challenges that can arise from privacy concerns in Big Data, we first elucidate how privacy in Big Data can be analyzed using two dimensions: (1) different *contexts* in which privacy is relevant and (2) different *principles* that specify the ethical meaning of privacy. *Privacy contexts* refer to the various spheres of human existence and activity in which individuals might expect different forms and degrees of privacy, each of which requires a specifically targeted analysis (Nissenbaum, 2009). Simultaneously, privacy obtains normative meaning only with additional ethical *principles* that state what is permissible in each context, which means that a variety of principles pertains —and so must be considered and prioritized —in any given context. We refer to the alignment of these two dimensions as a *Privacy Matrix*. We hope that the Privacy Matrix stimulates fruitful discussions about the role of Big Data, and helps promote awareness of privacy-oriented ethical issues in Big Data, both those issues that are known and the possibility of additional issues arising in the future. Beyond the privacy of information, however, the ethical principles upon which we draw also help interpret the design, use, and evaluation of Big Data tools. Our analysis suggests as a whole that the process of building Big Data Technology (information and tools) implicitly or explicitly embeds values into its use. By highlighting these issues, we aim

to help scientists, engineers, and other Big Data Technology designers, builders and users better identify and explicitly reflect upon the ethical values their work entails.

2. Defining Big Data

The rise of personal computing, the Internet, inexpensive archiving, and advanced computational infrastructures have enabled a wide range of people such as scholars, marketers, governmental agencies, educational institutions and the general public to produce, share, and analyze vast amounts of readily available information, a phenomenon known as big data (Boyd&Crawford 2012). However, the notion of big data remains complex and multifaceted, and invites multiple definitions. Most definitions refer to, or highlight the “three V’s”: volume, velocity, and variety of data. For example, IBM describes big data as “being generated by everything around us at all times. Every digital process and social media exchange produces it. Systems, sensors and mobile devices transmit it. Big data is arriving from multiple sources at an alarming velocity, volume and variety” (IBM, 2014). Others emphasize the tendency of big data to exceed the management capabilities of conventional database tools (Einav&Levin 2013; Kaisler et al., 2013; Dumbill, 2014). For some people, big data potentially offers analytic power of mythological proportions with “the widespread belief that large data sets offer a higher form of intelligence” (Boyd&Crawford, 2012: p 663). For this chapter, we follow the National Science Foundation (2012) which defines big data as “large, diverse, complex, longitudinal, and/or distributed data sets generated from instruments, sensors, Internet transactions, email, video, click streams, and/or all other digital sources available today and in the future.”

3. Ethical Issues in the Use of Big Data

The challenge that is most frequently used in discussions of the ethical issues that arise from the research on, and use of, Big Data is *privacy*. But issues of privacy are often invoked without a proper specification of the defining qualities of the concept. It is not immediately clear what privacy is, exactly. Likewise, it is often not clear whether the term is used, or should be used, in a descriptive or a normative way. While certain kinds of data simply *are* private because of their content and origin, it is not clear that they also *ought* to be treated and respected as private matters. It is also not immediately clear where the normative implications of privacy lie. Exploring the definitions of privacy is a necessary first step for any fruitful discussion of the pertaining ethical issues. In the last chapter (5.) it becomes clear that appreciation for the complexity of the construct of privacy can also support growth in the awareness of other issues that arise in/from Big Data.

3.1. Defining Privacy

To start with some clarifications of the terms, privacy is a quality that can be attributed to actions, things, and pieces of information. Put differently, it applies to both to tangible and intangible things. Tangible and intangible things can be qualified as private insofar as they belong or relate to individuals or groups of individuals. This means that the privacy of things and data becomes thematic only insofar as persons can be concerned. Privacy is the right of the subjects (human individuals) to determine to what extent their thoughts, sentiments, emotions, or other personal and unique information is to be released to other individuals (Solove, 2008). Privacy is release of information without consent, where a legitimate expectation of non-inference is expected. It is not limited to dissemination. Privacy concerns collection and processing, as well as cases of invasions of privacy via forcible interrogations (Solove, 2008). Privacy concerns can emerge in various contexts, according to the types of human activities. In fact, according to Nissenbaum (2009), it is impossible to define privacy in the abstract or in the context of the pure individual person. Nissenbaum (2009) proposed that privacy need to be defined as control of flows of information in given socio-technical contexts. In this paper we add that for reasons of practicality and to reflect the historical evolution of privacy, which has in fact determined its contours in jurisdiction, we need to distinguish these contexts in terms of their distance from the ego and of the social-institutional agents that might control or distort it.

Given this clarification of “privacy”, we can see that privacy implies a relation among several persons or agents, both individual or corporate (organizational). Etymologically, privacy has a negative sense: the Latin *privo*, means to rob or spoil or to de-priv someone of something. Understood in this formal sense, privacy becomes thematic because one agent has an interest in something that another agent has an interest to withhold from him. Privacy can also be defined as the right of not being despoiled of something that by right belongs to the individual.

The number of the relata that are relevant for the determination of privacy is potentially unlimited. Privacy can become thematic between more than two agents. For example, Tom might not want to tell Joe anything because he thinks that Joe will tell it to Jack. Analogously, it might not be problematic if one party has access to personal data, but if that party gives them to another person, then it *might* be a problem; and so on. This character of privacy also means that privacy can be affected in indirect ways, and affections of privacy can have indirect consequences for persons. In addition, privacy not only has straightforwardly factual implications (for example, using or not using the property of others), it also has an emotional dimension, concerning one’s personal attitudes and feelings toward others (including towards ruling authorities).

From a methodological point of view privacy cannot be identified or defined in the way that other, more distinct and self-contained qualities are defined. For example, while the

term “yellow” designates the same color for everyone who is able to perceive light within a specified wavelength range, the term “private” acquires meaning depending on the context in which it is used, and based on the relations it involves.

Additionally, the relational character of privacy also entails that it is practically never valued or desired for its own sake, but always because of something else. This aspect will be further explained with respect to the different normative principles that have to be used in order to specify the ethical dimension of privacy.

Finally, the relational character of privacy also suggests that it is not a static concept, in the sense that certain actions or things, or certain information, are thought to be always and essentially private, while others are not. For example, it can have a liberating effect to make certain aspects of one’s private life public (e.g., “coming out,” “The personal is political.”). In turn, a liberating effect can occur if individuals can trust that aspects of their private life are not disclosed to, or used by, others. The concern for privacy is best understood as concern for the changing, and often fragile, demarcation between the private and the public (social, political, economic) sphere.

The term “privacy” is not unique in having a contextual meaning. Other notions that refer to social structures are used in an equally variable way. For example, privacy is similar to the notion of friendship. There is no primary and exclusive meaning of the term “friend”; individuals can be “friends” based on various roles and contexts, and with various degrees of intimacy.

The various meanings of privacy do not need to be unified. With Wittgenstein, we can assume that there are family resemblances between the various uses of the term, that is, similarities without one unifying and unchangeable core (Wittgenstein, *Philosophical Investigations*, paragraphs 65-67). In this sense, we can also assume that although there is no unified meaning, each meaning and use is definite. Therefore, as a concept, privacy is not “vague” or “difficult to define,” as some might think, but rather, has meaning that is determined by the context in which it is used. For example, the users of websites can make precise and legitimate claims with regard to the use of their private data by a specific Internet application, even if it is unclear how their privacy should be handled with respect to other applications.

3.2. Four Normative Principles as Basis for the Ethical Analysis of Privacy

Privacy can acquire a normative meaning insofar as it is possible to say that privacy “ought to be respected,” or is a “value” to which one should adhere. However, privacy does not have value that trumps all other values. That is, other values or principles *can* be prioritized over privacy, which can be seen in cases such as domestic violence and child abuse, which occur in the private sphere and have at times been treated as a private

matter, but are now no longer considered this way. Therefore, if privacy is taken as an ethical principle, it has to remain less fundamental than, for example, the respect for the dignity and integrity of persons.

On the other hand, “respect for privacy”, if taken as an ethical principle, is too vague to be meaningful in practical terms. Insofar as the meaning of privacy is contextual, one cannot “respect privacy” as such, but only in relation to specific conditions and agents. The level of privacy has to be spelled out more concretely in each context, for example as it can be in the patient-doctor confidentiality or informed consent.

This twofold limitation of privacy as an ethical principle – that it is either not fundamental enough or too general and vague – leads us to the conclusion that one cannot use privacy as ethical principle in an isolated way. If it is given a normative meaning, it has to be specified in relation to which principles or values this meaning is understood. This conclusion is also supported by the fact that privacy is desirable and obligatory not for its own sake but always because of something else, as stated above.

Four ethical principles can be used to specify the ethical meaning of privacy: Nonmaleficence; justice; autonomy; and trust. These are defined and their contributions to our specification of privacy are outlined below.

- Nonmaleficence: refers to the harm that can be experienced by an individual or a group of individuals. In the context of Big Data, “harm” has to be understood as direct harm (physical, psychological, social, economic, etc.) following the use of personal data. It might seem that all ethical concerns can be related to the principle of avoiding harm, so that only this one ethical principle would be necessary. To a certain degree this is true. However, there are cases in which the experience of harm is only an indirect consequence and the ethical concern relates more directly to practices and the values they incorporate. For example, one can be concerned about practices of democratic citizenship that are affected by the use of Big Data without having to prove that there are individuals who actually experience harm. If harm were the only ethical concern, one would overlook, for example, the role that intrinsic values play in human practice.
- Justice: refers to the distribution of opportunities, rights, and goods among the individuals or groups of individuals that are the target of data mining activities. For example, certain groups can suffer discrimination because statistical data show that they are less likely to succeed in facing specific challenges. Strong anecdotal evidence can be found that users of Big Data are looking for predictions about economically not lucrative segments of the population (Marwick 2014).
- Autonomy: refers to the decision-making capacities of individuals or groups of individuals. While the principle of nonmaleficence conceptualizes agents as objects, or targets, of data mining practices, the principle of autonomy

conceptualizes them as subjects. Autonomy can be understood in a narrow and practical sense with respect to Big Data, for example with respect to the question how much control the individual has about the use of his/her personal data. In a wider sense it concerns the question how much freedom is left to individual decision-making. The latter is relevant given the attempt to use Big Data for the prediction of individual behavior, which potentially eliminates the whole factor of individual decision-making.

- Trust: refers to the relation between the data sources and the agents who are interested in mining their data. It involves all relata, that is, agents in their commonly shared practices. Trust can be defined as collective attitude that reduces the burden of permanent mutual control, without, however, dispensing of it, as most social interactions involve a mixture of trust and control. Like “privacy” and “friend”, “trust” can be instantiated in various ways and to various degrees. Although specific institutional provisions can be implemented in order to establish trust in a given setting (e.g., “checks and balances”), it is often difficult to verify empirically whether trust really exists. . In the case of Big Data, trust is not a positive imperative but a category that can be used to assess critically the ways in which data are used. With respect to Big Data, agents have no obligation to trust any other agent, but they do have an obligation to be trustworthy. One can ask, for example, whether the government or companies are acting in a way to enable individuals to trust the ways in which data are used.

3.3. Remarks and Explanations

The four principles stated above are broad enough to cover other categories that are often used in order to conceptualize the ethical reach of privacy. For example, justice is a broad enough principle to cover concerns for discrimination and the access to data. It has been remarked that privacy has both *instrumental* and *intrinsic* value. Privacy can be both a means to an end and an end, desired because it is a good in itself (Moor, 1997). These two ways of conceiving of privacy are not mutually exclusive and covered by the principles mentioned here. The instrumental value of privacy can be captured by the use of non-maleficence and justice, the intrinsic value by autonomy and trust. In the case of the former, privacy is valued because it allows one to avoid harm and establish the fair treatment of all members of society, in the case of the latter, privacy is desired because being able to rely on the privacy of one’s data is a necessary part of the intrinsically valued practice of individual autonomy and the trust that enables social interaction.

Privacy has also been conceptualized through the distinction between the *restricted access theory* and the *control theory* (an overview is given by Tavani, 2008). The difference between these theories results from taking data sources (that is, individuals) as objects or subjects, respectively. While the restricted access theory sees data sources

passively as objects, the control theory involves them actively as subjects. It might be worth noting that there is nothing inherently wrong with taking individuals as objects in the course of an analysis. For example, a physician uses a patient's data without the latter having direct control. The patient has to have the confidence that the physician uses the data in a confidential way. Obviously, at one point at least the data sources have to be involved as subjects, for example by giving consent to the use of personal data, but the practice that ensues does not need to involve them actively. The principles stated above reflect this distinction and show again that the different approaches to privacy are not mutually exclusive: in some cases, privacy can be organized according to the idea of restricted theory while in other cases it requires active control.

It has been suggested to conceptualize privacy in terms of *rights*. Privacy would then be violated whenever the right to "life, liberty, and property" is being violated (Volkman, 2003). The advantage of this approach is that it makes it possible to specify concerns for privacy and relate them to well-known legal principles. As rights, privacy rights are indeed "derivative" from other, more fundamental rights. However, the focus on rights seems to obscure other normative concerns. For example, if one says: "The flow of information is not the problem. It is the illegitimate use of information that is of concern" (Volkman 2003, 209), then all cases in which the illegitimate use of information cannot be shown would raise no ethical concern, which is clearly not the case. A certain flow of information can, for example, erode trust even if no direct violation of privacy rights is proven. Analogously, if the commercial use of data is seen exclusively under the perspective of rights, one has to conclude that "prohibiting such capitalist acts between consenting adults is paternalistic and immoral" (Volkman 2003, 210). This perspective seems to shift the burden exclusively to the side of the providers of personal data, insofar as their active consent or refusal is needed in every case and they are given a certain degree of responsibility for the use (and misuse) of data. Also, no violation of privacy could be claimed as long as some consenting consumers were to be identified, which limits drastically the range of ethical analysis.

As an additional, fifth principle one could mention *beneficence*. It would refer to the possible goals and purposes of data mining. As a term, beneficence can be defined as active concern for the well-being of others. Like trust, it can be used as a critical category insofar as it is possible to ask whether data are used with the goal of improving the life of citizens. Critical concerns can raise the question whether the use of data is guided by genuine beneficence, and whether the principle is used in an inclusive and universal way. For example, cases where the access to Big Data is without from certain segments of the population because of a concern for profitability can be seen as a violation of beneficence. However, it is not clear yet whether the use of Big Data will be driven by

genuine beneficence as a concern for the well-being of others, and not rather by attempts at improving managing procedures and economic outcomes. The principle of beneficence seems therefore less relevant than the other four, although this situation might change in the future.

The principles used in the present paper are very close to the ones used in the **Menlo Report** (Dittrich & Kenneally, 2012). The Menlo Report is an important document written in 2011 to provide guidelines to researchers in the field of information and communication technologies (ICT). It is modeled on the paradigm of the Belmont Report which in 1979 established principles for biomedical research. The principles used in the Menlo Report are respect for persons (equivalent to the principle of autonomy used here), justice, and beneficence. The latter is defined as avoidance of harm and concern for public welfare, which means that it is wider than the present use and covers both what is distinguished here as nonmaleficence and beneficence. As a fourth principle, the Menlo Report mentions respect for law and public interest, which covers issues such as compliance, transparency, and accountability. In the present paper, these can be subsumed under the principle of trust. In general, the Menlo Report represents the attempt at establishing ethical principles and rules from within the community of researchers in the field of ICT, and it is important to note that the Privacy Matrix suggested here is in congruence to this attempt.

Autonomy can be understood in a twofold way. The first way relates to individual decision-making as a democratic practice. It has been noted that citizens in a constitutional democracy should be given the right to opacity so that they can legitimately refuse their lives “being read” by others (Hildebrandt, 2011). Autonomy must therefore not be reduced to procedures of informed consent, but concerns the roles that agents assume in social and professional interaction. One can claim, for example, that “meaningful autonomy requires a degree of freedom from monitoring, scrutiny, and categorization by others” (Cohen, 2000). This also calls for a positive attitude toward “semantic discontinuity” which entails more “contextually specific practices of self-definition” in the use and regulation of information systems (Cohen, 2012). Instead of the assumption that the Internet and data systems are to be covered by a single global regulation that assumes all individuals follow the same ideas of agency and privacy, it seems necessary to allow for more particular, either national or group-specific regulations that reflect the respective decision-making more accurately.

While such concerns can be qualified as soft insofar as no identifiable harm to individuals has to occur (harm is a “hard” category insofar as it has to be verifiable in each case) and the concerns are related to long-term changes in mentalities and practices which can only have an indirect effect on individuals, they show that an important part in the ethical reflection on big data is related to the evaluative attitudes with which it is received. From

an ethical point of view, there is no reason, or no possible justification, that would allow one to neglect such concerns. Even if Big Data do not necessarily have repercussions for specific individuals, they change the way in which society operates as a whole, which means that members of society can be legitimately concerned by it. Evaluative attitudes toward issues of trust and autonomy are necessary condition of individual agency and therefore need to be addressed. Otherwise, one would have to say that democratic practices are independent from the way agents experience their status vis-à-vis governing institutions, employers, and the like.

The second way in which autonomy can be understood concerns ontological conceptions of agency. Big data can be used for the “prediction, preemption, presumption” of individual behavior (Future of Privacy Forum 2013). Some see the risk of a reification of human cognitive processes (Hildebrandt, 2011). Pattern recognition entails a merely statistical conception of individual agency, which can have an impact on attitudes toward individual decision-making and the degree of freedom it is given in specific settings. If an inclination toward anticipatory or preemptive governance becomes an inherent part of policy-making, the autonomy of individuals or groups can be severely limited. From an ethical point of view, this means that individual decision-making has to be given an intrinsic value, especially in the light of statistical interpretations that can lead to qualifying particular decisions as arbitrary, detrimental, or defective. That is, ethically speaking individual decision-making has to be given the opportunity to define its own inherent standards, without being forced to resort to the “higher” vantage point of statistical data collection.

3.4 Contexts of privacy

Privacy has also been described as “*contextual integrity*” (Nissenbaum 2009). It has been remarked that approaches to privacy are often too general and do not take the “compatibility with presiding norms of information appropriateness and distribution” in given contexts into account (Nissenbaum 2004, 137). This insight is particularly relevant for the fine-tuning of privacy-related policies. It follows directly from the relational meaning of privacy explained above.

Contexts of privacy cannot, however, be defined arbitrarily. A possible approach is to project privacy onto the canvass of human experience. If privacy concerns the collection, processing, and ultimate dissemination of information from the individual to others (Solove, 2008), the trajectory of private information should start with the most intimate contexts of life: the bodily, interpersonal (family, friendship), and home based ones. This context is the most strongly defended by laws and customs. The 4th Amendment of the US constitutions refers mostly to it, especially in the residential context. Privacy in this context refers to information about the self that is deeply personal, including activities consumed in one’s home. It includes private diaries or other type of written records, oral

communication or interactions in one's own residence, and so on. Medical records, although recorded by various institutions, refer to one's body and also benefit from one of the highest level of protection. HIPAA regulations in the US came to strengthen this point of view.

As humans due to their social obligations participate in contexts outside these realms, contexts of privacy emerge at each turn of our social journey. The contexts need to be seen, however, as layers of sociability, increasingly distant from the most intimate context of privacy (bodily, interpersonal, residential). Thus, in layers of sociability that are increasingly distant from the self, claims to privacy become increasingly weak and legal protection correspondingly thinner. Such contexts would start with the ones that are the closest to our interests, choices, and control, such as voluntary participation in various social, religious, and civic organizations. These are the communitarian contexts. Here, the expectation is that our activities, to the degree to which they are not detrimental to others (such as participation in terrorist or criminal groups) should be protected from undue scrutiny. Of course, when the social participation in these organizations is public, such as an open religious mass or open civic event, the claim to privacy cannot be called in defense. Yet, confession, some donations to charitable organizations, use of public resources (e.g., libraries), or in kind community interactions that are by definition philanthropic entail a good degree of privacy that is recognized as such. For, example, the American Library Association has been staunchly and rightfully defended the right of library patrons not to have their reading records disclosed, not even in criminal cases, without a strongly determined due cause and a court order

<http://www.ala.org/advocacy/intfreedom/librarybill/interpretations/privacy>).

Closely connected and at times hard to differentiate from it is the educational layer of social interaction. In an educational context some information is strongly defended, while other less, according to the social implication of the data. Personally identifiable information that might put the person at a disadvantage or reduce his or her autonomy is strongly defended (grades, courses taken, etc.). Other type of information, especially if aggregated for assessment of educational policy impact is publicly available.

Following the track of human activities, privacy contexts escape more and more the control of the individual as he or she enters in transactions with organizations and institutions that have a legitimate claim to recording, preserving, and further disseminating the activities or information pertaining to the individual. A first context is that of our interactions with legal, political, governmental, or law enforcement organizations. Here, some types of information are legitimately public, while other ought to stay private. For example, individual contributions to political campaigns are public in an attempt to keep the political process transparent. Land records are also public, as are court records for most criminal cases. While we might consider the last two types of

information intimate and highly personal, the impact an individual's owning of a certain parcel of land or of their criminal activities is for the most part social and ought to be publicly accessible. On the other hand, while voter records are public, including party affiliation, voting behavior is not public, in a defense of our freedom of conscience and expression. Neither are tax records. The US Census bureau never asks questions related to religious affiliation.

More distant still from the most intimately private contexts are those pertaining to most commercial transactions, as in buying and selling or using commercial services. Such interactions demand transparency by definition, for the sake of enforcing contract laws in case of conflict. These are similar with most judicial transactions, which are to be open and subject to public scrutiny to prevent secret trials and abuse of power. At the same time, the definition of what is private and what is public is not rooted in abstract laws or principles, but in contractual obligations voluntarily accepted by the user. Most interactions and activities on social media enter in this category. What is and what is not private is subject to the contractual obligations that the users accepted when signing up and clicking the box for "accept terms of service."

Finally, privacy concerns may emerge in the context of participation as subject in scientific research. The situation is for the most part regulated by contractual terms, guaranteed by "informed consent." Yet, this is far from a clear cut situation, as some research contexts could intrude upon individual information of the most intimate kind (e.g. medical information).

In brief, our concept of concepts of privacy takes a layered approach. It orders contexts on a "distance" dimension, where some are closer, while other farther away from the most intimate and strongly defensible claims to privacy. In this respect, we follow the pragmatic approach of most legal literature, which aims to operationalize the contexts as areas of human activity with definite pragmatic implications and outcomes.

As we will explain below, the contexts become clearer and easier to comprehend when included in a Privacy Matrix that aligns different normative principles with a set of different levels and contexts of privacy.

4. 1 The Privacy Matrix: How to Think About Privacy in Big Data

This chapter suggests that privacy should be addressed according to two dimensions, referred to here jointly as the Privacy Matrix as shown in Table 1. The first dimension comprises the possible levels, or contexts, in which privacy can become relevant. Each level or context requires specific analysis from a privacy perspective. The second dimension is based on the ethical principles that specify the normative, or ethical meaning, that can be given to the idea of privacy. The combination of the two

dimensions, finally, is based on the idea that privacy, both on the descriptive and normative level, has to be further specified if an analysis is supposed to yield meaningful results. The Matrix is based on the assumption that practical concerns regarding the different levels of individual and social life can each be combined with different ethical principles.

Specifying Principles	Privacy Contexts					
	Individual	Community	Education	Governmental	Science	Commercial
Nonmaleficence						
Justice						
Autonomy						
Trust						

Table 1. The Privacy Matrix: the columns represent possible privacy contexts while the rows represents the ethical principles of privacy.

We suggest using the Privacy Matrix as a heuristic tool. The list of ethical principles, and the way they are understood, is not seen in opposition to existing approaches to the ethical analysis of privacy, but rather as an attempt at dealing with the necessary pluralism of principles in a more effective way. Very often, a variety of principles is suggested in a way that leaves it open which ones should be applied to the case at hand. Obviously, one would like to address all possible ethical concerns as one should not arbitrarily decide to leave some of them out if they are relevant, but not all of them can be applied to the same degree. The same can be said for the levels of privacy. With the Privacy Matrix, it is possible to start from the process of application. The question then becomes: on which level, or in which context, is privacy most relevant in the given case, and which normative concern is most relevant? It seems evident that no particular case can be limited to one combination of criteria only, but can always be conceptualized in various ways. However, one can assume, if only for heuristic purposes, that each case is relevant in one primary way, which then has to be taken as point of departure for an ethical analysis that is both specific and effective enough. Even if the search for the primary application of contextual and ethical criteria can seem arbitrary in certain cases, it supports at least one relevant issue being addressed, and it might help to identify others that have not yet been considered. The goal is to shift the focus of the analysis from the multitude of possible perspectives, which is often practically irrelevant, to the steps that are necessary to engage in a process of decision-making which is then, hopefully, practically relevant.

The privacy contexts included in Table 1 represent characteristic areas and fields for using big data. Electronic medical records, which the International Organization of

Standardization defines as “a repository of patient data in digital form, stored and exchanged securely, and accessible by multiple authorized users” represent a good example of big data related specifically to **individuals** (Hayrinen et al, 2008). The **community** context mainly consists of social media data, which investigators in many sectors are mining for useful information. For example, a recent study published in *Preventative Medicine* revealed the attempt at tracking real-time social media like Twitter for monitoring HIV exposure and drug-related behaviors with the intention of detecting and preventing future outbreaks (Stoove and Pedrana, 2014). **Educational** services and companies now use big data with the aim of improving teaching and learning. For example, Knewton, an education technology company, created digital courses in which students are tracked “as they play online games, watch videos, read books, take quizzes, and run laps in physical education” (Simon, 2014). The federal **government** employs big data sets from various programs for secondary purposes beyond the aim of their original collection. For example, law enforcement officials have attempted to develop predictive technologies using big data to anticipate, intervene, or prevent crime, including identification of terrorist networks, warning of impending attacks, and preventing the proliferation of weapons of mass destruction (Executive Office of the President, 2014). **Scientists** working on the Personal Genome Project at Harvard Medical School are investigating the utility of genetic data for enhancing health care in multiple ways such as increasing drug effectiveness, assessing predisposition to disease, and constructing microbiome profiles (Ball et al., 2014). Finally, **commercial** retailers and marketers analyze a wide range of customer activity, both on and offline, to provide, as they claim, more tailored recommendations and “optimal pricing”. For example, in April 2014 Verizon Wireless notified customers that it would begin gathering data about user activities and selling them to marketers (Lazarus, 2014).

These various examples of the contexts in which big data arise and get used illustrate that Big Data refers to the dynamic use of data for insight rather than static archives of massive amounts of information. Research using big data relies on massive, continuous, real-time data streams that might predict social behavior by those both creating and analyzing the data. Thus, the phenomenon of Big Data includes the potential for relatively continuous monitoring, control, and moderation of individual and societal behavior. “(O)ur ability to modify public behavior increases as the observed individuals are more exposed to our scrutiny and tracking.” (Matei Correspondence, 2014: PageXX). This suggests that we must acknowledge and make explicit tradeoffs between privacy and the utility of analyses based on Big Data. Our approach suggests, however, that the tradeoffs vary across the several contexts of Big Data collection, analysis and use. In the pages that follow, in which we explore the implications of this observation with respect to specific examples that serve as case studies in the contextual variation of privacy

concerns, we will show the implications of privacy for the realization or protection of specific ethical values such as autonomy and social justice (5.1-4.6.). This analysis illustrates use of our Privacy Matrix as a guide for inquiry into the relationship between privacy context and ethical principles in ethical reasoning about Big Data. In the final chapter, we will see how ethical values are directly related to, or even embedded into, the tools that use and research Big Data (5.)

4.1. Profiling Individuals with Big Data

Big data when is used to create profiles of individuals for various purposes brings risks. For example, the White House privacy report explains, “credit scores and other economic data could influence an individual’s opportunities to find housing, forecast their job security, or estimate their health outside of the protections of the Fair Credit Reporting Act. Individuals have little recourse to understand or contest the information that has been gathered about them or what that data, after analysis, suggests” (Executive Office of the President, 2014). The report further suggests that pricing and discrimination caused by big data could exacerbate existing socioeconomic disparities in education and the workforce setting. A more specific example of group profiling was revealed by The New York City Police, which reports stated that officers singled out mosques and collected attendees’ license plate numbers and plotted the locations on a map. “The Department of Homeland Security’s more recent plan to build a national license plate database— and the outcry it provoked — suggests that minorities may be especially vulnerable to what Americans would perceive as a violation of privacy” (Fung, 2014).

4.2 Anonymity, manipulation and user consent in Online Communities

In 2006, researchers at Harvard began gathering anonymized data on 1,700 college age Facebook users to study how interests and friendships changed over time (Boyd&Crawford, 2014). However, “these supposedly anonymous data were released to the world, allowing other researchers to explore and analyze them. What other researchers quickly discovered was that it was possible to de-anonymize parts of the data set: compromising the privacy of students, none of whom were aware their data were being collected” (Boyd&Crawford, 2014). Among many other studies conducted using Facebook data, in 2012, as described above, Facebook completed the “emotion contagion study” in which they skewed users newsfeed so they would see content happier or sadder than average “and when the week was over these manipulated users were more likely to post either especially positive or negative words themselves” (Meyer, 2014). The experiment was, technically speaking, legal, according to Facebook’s Terms of Service in which users relinquish their data by joining the social media site. Yet, the study was conducted prior to IRB approval. Experts and casual users alike have criticized the study saying “emotional well-being is sacred” and “research is different than marketing practices” (Boyd, 2014).

In addition to Facebook, Twitter has also been accumulating user data since 2006, in a rate of five hundred million tweets worldwide everyday, and announced that it is planning to release them all. The data is promising for scientists “looking to find patterns in human behaviors, tease out risk factors for health conditions and track the spread of infectious diseases” (Moyer, 2014). However, the question arises whether researchers may collect and use such data for research without the users’ consent and intention to be part of research. Other social media outlets have also used data without user consent. For example, Path was a social networking app for photo sharing and messaging. In 2012, the app was criticized for accessing and storing member phone contacts without their knowledge or permission. Path was fined \$800,000 by the Federal Trade Commission (Ramirez, 2014). Similarly, a Flashlight app failed to disclose to iPhone users that it was sharing their location data with advertising networks. Finally, Snapchat is an app that allows users to send pictures to friends that “self destruct” seconds after opening. There have been several simple ways identified that allow recipients to save the pictures indefinitely. Moreover, Snapchat was fined for security failures in which attackers compiled a database of 4.6 million Snapchat usernames and phone numbers (Ramirez, 2014).

4.3 Protecting vulnerable populations in educational contexts

Educational technologies firms are serving as third parties accumulating academic and behavioral data on students. The education company Knewton, as discussed above, observes students “monitoring every mouse click, every keystroke, every split-second hesitation as children work through digital textbooks, Knewton is able to find out not just what kids know, but how they think” (Simon, 2014). These companies are gathering up to 10 million unique data points on each child per day and despite extensive privacy policies and terms of service, an examination revealed there were “gaping holes in the protection of children’s privacy” (Simon, 2014). Another case is Learnboost, a third party that allows teachers to upload their notes and student attendance, test scores, behavior and more to a digital textbook. The teachers are then eligible to use, for example to email these grade books with no other regulation than “as they see fit”. A recent national study found that only 7% of contracts between schools and educational technology companies agreed not to sell the data for profit. Also, “few contracts required the companies to delete sensitive data when they were done with it. And just one in four clearly explained why the company needed personal student information in the first place” (Simon, 2014). For the company InBloom, privacy concerns resulted in school districts withdrawing from contracts and ultimately shutting the company down. InBloom a non-profit corporation was financed with \$100 million in seed money from the Bill and Melinda Gates Foundation as well as the Carnegie Corporation of New York to store and manage student data for public school districts across the country. However, once parents began to discover what kind of data was being collected, such as social security numbers, they

began to speak out causing multiple school districts to pull out and ultimately InBloom to close its doors (Singer, 2014).

4.4. Unequal Access to Big Data in Scientific Research

While the previous cases relate to the undue access to big data, it also has to be mentioned that not all potential users have equal access to data resources. “Top-tier, well-resourced universities will be able to buy access to data, and students from the top universities are the ones most likely to be invited to work within large social media companies,” resulting in a gap between researchers who have the potential to study big data and those who have not. Well-funded companies mostly likely will also have more access to data. This gap created by the difficulty and expenses associated with the access to big data results in a “restricted culture of research findings,” as large data companies have no requirement or responsibility to make their data available. In addition, “big data researchers with access to proprietary data sets are less likely to choose questions that are contentious to a social media company if they think it may result in their access being cut. The chilling effects on the kinds of research questions that can be asked – in public or private – are something we all need to consider when assessing the future of big data” (Boyd&Crawford, 2012: pg 674).

4.5 Big Data and government surveillance

Before the Snowden affair renewed vigorous debate about government surveillance, the Terrorist Information Awareness (TIA) program generated extensive controversy about balancing privacy with national security in contemporary America. (Cooper and Collmann 2005; Department of Defense, 2003). Begun in the wake of the 9/11 attacks, TIA mobilized enormous computer capability to search databases across the government in search of terrorists, an early application of Big Data before the term became popular. Critics argued that TIA posed multiple threats to the privacy of individual Americans (Safire, 2002; Washington Post, 2002; Crews, 2002; Simons and Spafford, 2003; Stanley and Steinhardt, 2003). TIA, they argue,

- Violates the Fourth Amendment of the Constitution by searching a data base containing detailed transaction information about all aspects of the lives of all Americans;
- Undermines existing privacy controls embodied in the Code of Fair Information Practices, such as improper reuse of personal data collected for a specific purpose;
- Overcomes “privacy by obscurity” including inappropriate coordination of commercial and government surveillance;
- Increases the risk of falsely identifying innocent people as terrorists;
- Increases the risk and cost of identity theft by collecting comprehensive archives of individually identifiable information in large, hard-to-protect archives;
- Accelerates development of the total surveillance society.

Critics also identified other potentially undesirable consequences in addition to invasion of privacy, including:

- Undermining the trust necessary for the successful development of the information economy and electronic commerce;
- Undesirably altering the ordinary behavior of the American population including quelling healthy civil disobedience, “normalizing” terrorist behavior, and inhibiting lawful behavior;
- Creating new, rich targets for cyberterrorism and other forms of individual malicious abuse of computerized personal information.

In addition to highlighting persistent concerns about privacy, civil liberties and government surveillance, the TIA controversy illustrates the need to reflect deeply on the ethical implications of any Big Data project during its design. Waiting until controversy erupts misses the opportunity to design a better application and sullies trust in scientific, political, educational and commercial leadership.

4.6. Sale of Big Data in Commercial contexts

Outside of social media, any activity an individual performs online can be tracked, resulting in information for commercial purposes such as marketing and behavior studies. For example, Disconnect is a program that lets users see who is tracking their visits to websites, revealing dozens of third parties observing and following their individual “click stream” (Kroft, 2014). These third parties, known as data brokers, “are collecting, analyzing and packaging some of our most personal information and selling it as a commodity...to each other, to advertisers even the government, often without our direct knowledge” (Kroft, 2014). In response, the White House report on Big Data and Privacy has highlighted the need for effective consumer privacy protections for the individuals (Executive Office of the President, 2014). But the report has also received criticism for stopping short of taking effective action to protect consumers, “such as requiring that private companies disclose to consumers what they know about them” (Lazarus, 2014).

Third parties also accumulate data without online sources as well. In Boston, an automated reader attached to the front of a “spotter car” takes a picture of every license plate it passes. These images, more than 8,000 per day, are then sent to Sousa, a company in Texas, that has over 1.8 billion plates from vehicles across the country. Typically, every license plate of a stolen or defaulted vehicle results in \$200 to \$400 for the company. In May 2014, a legislative committee was scheduled to hold a hearing on a bill that would ban most uses of license plate scanners. Jonathon Hecht, a Massachusetts representative said, “(w)e need to have a conversation about how to balance legitimate uses of this technology with protecting people’s legitimate expectation of privacy” (Musgrave, 2014). Kade Crockford of the American Civil Liberties Union of Massachusetts went on to explain, “it’s the wild west in terms of how companies can collect, process and sell this kind of data” (Musgrave, 2014).

5. Embedding Values in Big Data Technology

Although we may loosely refer to Big Data as if the data stand on their own free of any supporting technology, creating, analyzing and using Big Data depends on also creating complex computer infrastructures, applications, and devices. In this section, we will explore the ethical, or values-related, implications of this co-creation of Big Data and Big Data Technology. We suggest that, in the course of enabling Big Data, designers, users and analyzers embed and realize values in Big Data Tools. This process may occur in three, often interdependent ways, namely:

1. Values may guide the use of Big Data tools;
2. Big Data tools may enable realization of values, and;
3. Building tools to realize values may entail, and often requires basing their design on the target values themselves.

We will examine various chapters in the present book to elucidate each of these processes. In the course of our analysis, we will also distinguish between two types of value, technical and ethical values. For Big Data technology, computer scientists and engineers usually seek to design tools that effectively accomplish a technical purpose, such as enabling effective analysis and visualization of Twitter conversations in real time. Realizing such technical values through Big Data Technology, however, often helps realize ethical values, such as minimizing loss of life in a mass shooting or natural disaster. Performance requirements may link technical and ethical values through the effective design and functioning of Big Data Technology (Cooper and Collmann 2005). Finally, we should recognize two types of Big Data Technology, including:

1. Technology that produces, archives, protects, and displays Big Data or its constituent components, and;
2. Technology that enables description, analysis, interpretation and understanding of Big Data.

Distinguishing between these two types of technology reflects a primary consideration of this book: much more Big Data exist than we have the tools to exploit. From the perspective of Big Data Technology design and use, we observe how values condition tools and tools help realize values in a dynamic, interdependent embedding process (Collmann and Robinson 2010; Cooper, Collmann and Niedermeier 2008)

5.1 Values guide the use of Big Data Technology

Organizational meetings occur in a variety of formats in the 21st century, including face-to-face encounters, teleconferences, videoconferences and, quite commonly, mixed media meetings over the Internet. As Ahmed and Gavrilova note, meetings in all forms absorb much staff time and, thus, corporate money in the workday, including time that sometimes appears to have yielded little return. It makes good business sense to investigate the utility of employing advanced computerized technology to make best use of expensive meeting time. Ahmed and Gavrilova describe a multimodal physiological and behavioral biometric system (Microsoft Kinect v2) that records and analyses the overt participation of individuals in a meeting, including devices that record talk, gait, facial characteristics, and movement across the room. The authors emphasize the

effectiveness of their system in capturing data for analysis and, thus, its potential value as a tool for increasing the efficiency and productivity of individuals during a meeting. The tool's technical efficacy in rapidly capturing traits and identifying individuals relates directly to its avowed purposes of analyzing meeting workflow, characterizing and evaluating individual contribution level and analyzing group dynamics and behavior – all with the goal of improving the ability of meetings to achieve corporate objectives.

This Big Data producing technology poses several value-related problems for reflection. First, the Kinect v2 clearly places greater value on overt meeting behaviors such as talk and note taking that monitoring devices can detect and lesser value on covert meeting behaviors such as listening or thinking that remain undetectable. Second, the Kinect v2 as a surveillance tool gives expression to the concept of the Panopticon, a means for continuously documenting all behaviors of a target population with little regard for its own desires. The Kinect v2 bears comparison with the educational technologies described above which drew comment for inserting a “third party” in the educational process with few controls by the observed population over use of the information. Adult employees in an organization usually have greater control over themselves than children in a classroom; but, without guidelines for use and protection, such minutely documented, partial information poses relatively uncontrollable risks to their well-being by becoming a yardstick for job performance. From the perspective of safeguarding personal autonomy in the workplace, Kinect v2 directly challenges employee strategies of “stage management” that establish distance between an individual's private, backstage persona and their public, onstage performance. Kinect v2 constrains meeting members to participate overtly even if covert methods of contribution match their personal working style better, they prefer “off-line” contributions to project development, or specific meeting contexts favor reticence.

5.2 Big Data tools enable realization of values

Chapters in this book specifically address how Big Data technologies enable the realization of key values such as trust in information from diverse sources on the Internet (Ignjatovic et al) and self-organizing in task-based work groups (Matei) as well as protection against assaults on such values (Caverlee and Lee). Careful reading also gives evidence of the contributions of Big Data Tools to patterns of social injustice, for example exploitative crowdsourcing of problems from developed countries to low wage “knowledge workers” in developing countries. These examples set a precedent for making explicit reasoning about the ethical consequences of apparently value-free tools a typical dimension of their design and implementation.

In a deeply technical analysis, Ignjatovic and colleagues address an ethical value that lies at the heart of the instrumental effectiveness of Big Data technology by referring to trust in the reliability of data from disparate sources across cyberspace. We observe above that trust refers to the relation between the data sources and the agents who are interested in mining their data and involves all relata, that is, agents in their commonly shared practices. Thus, we define trust as a collective attitude that reduces the burden of permanent mutual control, without, however, dispensing of it, as most social interactions

involve a mixture of trust and control. Ignjatovic and colleagues note how the failure of typical social mechanisms of information credibility and trustworthiness in cyberspace necessitate automated methods, especially when drawing vast quantities of data from vast numbers of sources. One could reasonably argue that, without some basis for trust, Big Data as currently envisioned becomes impossible. From their perspective, trust in a flow of information derives from its provenance, or, the combined trustworthiness scores of sensor nodes and all the nodes through which data passes including terminal, intermediate and server nodes, and the emergent trustworthiness of data elements as they pass through the network. While offering detailed technical ideas about combatting collusion attacks in online rating systems, Ignjatovic and colleagues observe that achieving effective results depends, given the current state of the art, on compromising the privacy and anonymity of participants. As one can argue, this may be entirely appropriate given the public nature of many online communities. Yet, for instances in which this condition fails because participants do not expect all information to be public, Ignjatovic and colleagues offer ideas for further technical research on Big Data tools to help realize the twin values of trust in data from an unknowable social space and the privacy of its constituent members.

Caverlee and Lee address issues that emerge from a specifically malevolent corner of cyberspace, the world of weaponized crowdsourcing. From their perspective, crowdurfing “wherein masses of cheaply paid shills can be organized to spread malicious URLs in social media, form artificial grassroots campaigns (“astroturf”), spread rumor and misinformation and manipulate search engines (p.???)” poses clear threats to information quality and community trust of such systems. Their research focuses on developing automated means for detecting crowdurfing tasks and crowdurfing workers as well as other, off-line mechanisms such as increasing the cost of crowdurfing campaigns. In contrast to Ignjatovic and colleagues who worry about the integrity of data flowing through the network, Caverlee and Lee focus on how fraudulent use manipulates and potentially distorts the social perception of content. The phrase “cheaply paid shills” suggests a certain view of crowdurfing workers, however, that we may want to query. Caverlee and Lee do not cite any national or international laws prohibiting crowdurfing. Their analysis suggests that crowdurfing, at least in the forms they analyze, constitutes a form of cheating or false advertising, not a form of crime. Crowdurfing workers, thus, perform no criminal acts but only knowingly or unknowingly facilitate the misrepresentation of their subject matter. The data from Bangladesh suggests a social justice issue, however, in which a crowdurfing requester takes advantage of the international division of labor and wages to exploit crowdurfing workers. Crowdurfing constitutes a form of low wage piece work with no job security, no benefits, and no form of worker organization to prevent its worse abuses. Even if they earn more money doing crowdurfing than other workers in Bangladesh, they make less than crowdurfers who work from developed countries. Hence, crowdurfers may, indeed, be cheaply paid but the noun “shills” taints the workers not the requester, the victims of the international division of labor and wealth not its perpetrators and beneficiaries. No pure technical solution exists for crowdurfing as long as poor people with access to computers can fill its labor ranks.

Matei, Bruno, Fabiola and Morris explicitly identify the values they hope to realize through use of the tool Visible Effort (VE), self-guidance and self-actualization of collaborative online workgroups. In computer-mediated collaboration (CMC), VE enables groups to 1) measure and visualize the degree of collaborative unevenness, and the emergence of social structure, and 2) actively or passively steer the collaborative processes to attain specific goals (Matei et al, see above pg.4). Matei and colleagues set their discussion of VE in the context of the debate about social hierarchy and productivity in teams. In contrast to analysts who argue that flat, decentralized teams solve problems faster and more efficiently than hierarchical teams, they state that “CMC needs division of labor, coordination and clear goals” (pg 2)” They employ Shannon’s theory of social entropy to conceptualize their approach and design their application of the VE tool. Social entropy refers to varying levels of random individual participation and group structure. The greater the social entropy the more random individual participation and the less coordinated their activities. Matei and colleagues hypothesize that productive CMCs strike an effective balance between social hierarchy and social entropy and offer VE as a tool to help find the right balance for any specific project. As a tool for measuring and visualizing social entropy, VE offers a means for CMC to discover the “inflection point” between social hierarchy and social entropy that best suits their task requirements and modify their work processes to help sustain it.

5.3 Basing Big Data Tool design on target values

In their chapter on bottom-up decision making in urban infrastructure projects, Bakht and El-Diraby give a striking example of how enabling realization of a specific value (effective community participation in urban planning) and its beneficial consequences (creating an innovative and socially-savvy decision-making environment) affects the requirements and design of a Big Data Tool. They describe facing a specific problem: how to track citizen discussion over social media about urban infrastructure projects in order to incorporate relevant feedback into the planning process. Specifically, they wish to discover the semantic (ideas) and social (people) characteristics of “Infrastructure Discussion Networks (IDN)” as they emerge and evolve over Twitter. IDN constitute an example of “small world phenomena” with “relatively high clustering, comparatively small diameters, and short average path lengths” that “offer a good opportunity for information diffusion and viral marketing around the project” (Bakht and El-Diraby, p. 12). In addition to discovering IDN in the Twitter flow, analysts need to follow their evolution over time with respect to the ideas under circulation and the networks of people participating in the discussions through social media. IDNs mature; that is, grow in size and density of connections among members (triadic closure) as well as display changes in sentiments about specific subjects (dynamics of opinion). In the words of Bakht and El-Diraby, “by selecting a particular context for analysis of discussions over IDNs (such as sustainability), results of (Social Network Analysis) and lexical analysis can be aggregated to form the profile of online discussions for a particular project” (Bakht and El-Diraby, p. 16). Their tool produces a project discussion profile represented as a series of graphs over time for specific dimensions of specific issues in a specific project. From the noise of Twitter come community-based messages from influential community

members to aid planners in better meeting community needs. Without such a tool specifically designed to elucidate the dynamic social composition and meaning of ephemeral communications on Twitter, the opportunity to feed community discussion into the planning process within a realistic time horizon would not exist. From the perspective of values, the tool reflects the importance of autonomy in the twofold sense of allowing individual agency to have an impact on the planning process and integrating a plurality of viewpoints. However, the mere use of the tool realizes autonomy only in a partial and asymmetrical way, insofar as the information from the community is gathered without the possibility of interaction and dialogue, which is the reason why the authors indicate that the full realization of this value requires real-world interaction, such as meetings and public consultations. Big Data tools can stimulate, but not substitute democratic interaction in the traditional sense.

6. Conclusion

In this chapter we have developed an approach for systematically analyzing ethical values in the design, use and evaluation of both big data information and tools. We intend for the Privacy Matrix including the analysis of both privacy contexts and ethical principles to enhance explicit identification, discussion and reflection among scientists, engineers and other Big Data developers of the values they employ in their work. Building upon the concept of values-sensitive design (Friedman, Kahn, and Borning, nd), we argue that the domain of ethics in Big Data has little bearing if it remains a subspecialty discourse among professional philosophers, ethicists or social activists. We also argue, however, that, whether explicitly recognized or not, values always inform the design, development and use of Big Data Technology. As a community of Big Data technologists, we should attempt to make our values explicit and knowingly embed the values we seek to realize in the results of our work.

References:

- Ball, Madeleine P Bobe, J. R., Chou, M. F., Clegg, T., Estep, P. W., Lunshof, J. E., ... & Church, G. M. (2014) Harvard Personal Genome Project: lessons from participatory public research. *Genome medicine* 6(2): 10
- Boyd, Danah (2014, Jul 1) What does the Facebook Experiment teach us? *TheMedium*. Online.
- Boyd, Danah, and Kate Crawford (2012) Critical questions for big data: Provocations for a cultural, technological, and scholarly phenomenon. *Information, Communication & Society* 15(5): 662-679
- Cohen, J (2000) Examined Lives: Informational Privacy and the Subject as Object. *Stanford Law Review* 52: 1373-1438
- Cohen, J (2012) Configuring the Networked Citizen. In: Sarat, A, Douglas, L, Umphrey, MM (eds), *Imagining New Legalities: Privacy and Its Possibilities in the 21st Century*. Stanford University Press, Stanford, p129-53
- Collmann, J and A Robinson (2010) Designing Ethical Practice in Biosurveillance: The Project Argus Doctrine. In D Zeng, H Chen, C Castillo-Chavez, B Lober, M Thurmond (eds) *Infectious Disease Informatics and Biosurveillance: Research, Systems, and Case Studies*. Springer, New York
- Cooper Ted and Jeff Collmann (2005) Managing Information Security and Privacy in Health Care Data Mining. In Hsinchun Chen, Sherri Fuller, Carol Friedman, and William Hersh (eds) *Advances in Medical Informatics: Knowledge Management and Data Mining in Biomedicine*. Springer's Integrated Series in Information Systems, Vol 8. Springer, New York
- Ted Cooper, Collmann, J. and Henry Neidermeier (2008) Organizational repertoires and rites in health information security. *Cambridge Quarterly of Healthcare Ethics* 17(4):441-452
- Crews, Jr., C.W. (2002) The Pentagon's Total Information Awareness Project: Americans Under the Microscope? Cato Institute, Available via DIALOG <http://www.cato.org/publications/techknowledge/pentagons-total-information-awareness-project-americans-under-microscope> Accessed 4 Nov 2014
- Department of Defense, Office of the Inspector General, Information Technology Management (2003) Terrorist Information Awareness Program (D-2004-033). Arlington, VA
- Dittrich, D, Kenneally, E (2012) The Menlo Report: Ethical Principles Guiding Information and Communication Technology Research. US Department of Homeland Security. Available via DIALOG .

Matei, S. A., Russell, M., and Bertino, E eds. *Transparency on Social Media - Tools, Methods and Algorithms for Mediating Online Interactions*. New York: Springer Publishing House. 2015

<http://www.dhs.gov/sites/default/files/publications/CSD-MenloPrinciplesCORE-20120803.pdf> Accessed 4 Nov 2014

Dumbill, Edd (2012, Jan 12) What is Big Data? *O'Reilly Radar*. Online.

Einav, Liran, and Jonathan D. Levin (2013) The data revolution and economic analysis. *Innovation Policy and the Economy*. Doi: 10.3386/w19035

Executive Office of the President (2014) Big Data: Seizing Opportunities Preserving Values. Available via DIALOG.

http://www.whitehouse.gov/sites/default/files/docs/big_data_privacy_report_may_1_2014.pdf Accessed 4 Nov 2014

Friedman, B., Kahn, JR., P.H. and Borning, A., et al., (2001) Value Sensitive Design: Theory and Methods. UW CSE Technical Report 02-12-01. Available via DIALOG <http://www.urbansim.org/pub/Research/ResearchPapers/vsd-theory-methods-tr.pdf>. Accessed 2 Nov 2014

Fung, B. (2014, February 27) Why civil rights groups are warning against 'big data'. *The Washington Post*. Online.

Future of Privacy Forum (2013). Big Data and Privacy. Making Ends Meet. Available via DIALOG. <http://www.futureofprivacy.org/big-data-privacy-workshop-paper-collection/> Accessed 4 Nov 2014

Häyrinen, K., Saranto, K., & Nykänen, P. (2008). Definition, structure, content, use and impacts of electronic health records: a review of the research literature. *International journal of medical informatics* 77(5): 291-304

Hildebrandt, M (2011) Who Needs Stories if You Can Get the Data? ISPs in the Era of Big Data Crunching. *Philosophy of Technology* 24:371- 390.

IBM (2014) What is Big Data? <http://www.ibm.com/big-data/us/en/> Accessed 4 Nov 2014

Kaisler, Stephen, et al. (2013) Big data: Issues and challenges moving forward. *System Sciences (HICSS)*. 46th Hawaii International Conference, Maui, 2013

Knewton Inc. (2014) About Knewton. <http://www.knewton.com/about/>. Accessed 4 Nov 2014

Kroft, Steve (2014, Aug 24) The Data Brokers: Selling Your Personal Information. *60 Minutes*. Online

Lazarus, David (2014, April 24) Verizon Wireless sells out customers with creepy new tactic. *Los Angeles Times*. Online.

MacPhail, Theresa (2015) Data, Data, Everywhere... *Public Culture Forum* (in print)

Matei, Sorin, Email correspondence 2014

Matei, S. A., Russell, M., and Bertino, E eds. *Transparency on Social Media - Tools, Methods and Algorithms for Mediating Online Interactions*. New York: Springer Publishing House. 2015

- Marwick, A (2014) How Your Data Are Being Deeply Mined. New York Review of Books. Available via DIALOG. <http://www.nybooks.com/articles/archives/2014/jan/09/how-your-data-are-being-deeply-mined/?pagination=false> Accessed 4 Nov 2014
- Meyer, Robinson (2014, Jun 28) Everything We Know About Facebook's Secret Mood Manipulation Experiment. *The Atlantic*. Online.
- Moor, James (1997) Towards a theory of privacy in the information age. *Computers and Society* 27(3):27-32
- Moyer, Melinda (2014) Twitter to Release All Tweets to Scientists. *Scientific American* 310(6)
- Musgrave, Shaw (2014, May 3) A vast hidden surveillance network runs across America, powered by the repo industry. *The Boston Globe*. Online.
- National Science Foundation (2012). Directorate for Computer Science & Information Science & Engineering, Critical Techniques and Technologies for Advancing Big Data Science & Engineering (BIGDATA). Available via DIALOG. http://www.nsf.gov/funding/pgm_summ.jsp?pims_id=504767 Accessed 4 Nov 2014
- Nissenbaum, Helen (2004) Privacy as Contextual Integrity. *Washington Law Review* 79:101-139
- Nissenbaum, H. (2009). *Privacy in Context: Technology, Policy, and the Integrity of Social Life*. Stanford Law Books.
- Ramirez, Edith (2014) Protecting Consumer Privacy in the Big Data Age. Federal Trade Commission, Washington, DC, 2014
- Safire, W (2002, Nov 14) You are a Suspect. *New York Times*. Online.
- Simon, Stephanie (2014, May 15). Data mining your children. *Politico*. Online.
- Simons, B. Spafford, E.H., Co-chairs, US ACM Policy Committee, Association for Computing Machinery, Letter to Honorable John Warner, Chairman, Senate Committee on Armed Forces, January 23, 2003
- Singer, Natasha (2014, April 21). InBloom Student Data Repository to Close. *New York Times*. Online
- Solove, D. J. (2008). *Understanding Privacy*. Harvard University Press.
- Stanley, J., Steinhardt, B (2003) Bigger Monster, Weaker Chains: The Growth of an American Surveillance Society. American Civil Liberties Union, Technology and Liberty Program. Available via DIALOG. <https://www.aclu.org/technology-and-liberty/bigger-monster-weaker-chains-growth-american-surveillance-society> accessed 4 Nov 2014

Matei, S. A., Russell, M., and Bertino, E eds. *Transparency on Social Media - Tools, Methods and Algorithms for Mediating Online Interactions*. New York: Springer Publishing House. 2015

Stoové, Mark A., and Alisa E. Pedrana (2014) Making the most of a brave new world: Opportunities and considerations for using Twitter as a public health monitoring tool. *Preventive medicine* 63:109-111

Sullivan, Gail (2014, Jul 3) Sheryl Sandberg not sorry for Facebook mood manipulation study. *Washington Post*.Online.

Tavani, H. (2008) Informational Privacy: Concepts, Theories, and Controversies. In: Himma, KE, Tavani, H (eds), *The Handbook of Information and Computer Ethics*, Wiley, Hoboken, 131-164

Volkman, Richard (2003) Privacy as life, liberty, property. *Ethics and Information Technology*5:199–210