# Event Summary: The Social, Cultural, & Ethical Dimensions of "Big Data"

March 17, 2014 - New York, NY
http://www.datasociety.net/initiatives/2014-0317/

*This event summary attempts to capture the broad issues and concerns raised by "The Social, Cultural, and Ethical Dimensions of 'Big Data'" conference. Not all attendees were involved in every part of the conversation, nor does this document necessarily reflect the views and beliefs of individual attendees. To learn more about the conference, watch the videos, or read the other documents produced for the event, please visit http://www.datasociety.net/initiatives/2014-0317/*

On March 17, 2014, the Data & Society Research Institute, the White House Office of Science & Technology Policy, and New York University's Information Law Institute co-hosted a public event entitled "The Social, Cultural, & Ethical Dimensions of 'Big Data'". The purpose of this event was to convene key stakeholders and thought leaders from across academia, government, industry, and civil society to examine the social, cultural, and ethical implications of "big data," with an eye to both the challenges and opportunities presented by the phenomenon.

As attendees noted, the big data phenomenon affects many facets of contemporary life and has the potential to alter governance, the economy, and the very structures of society. The battles over data bridge public and private sectors in fundamental ways, raising key questions about the role of civil society and the practice of journalism. Technology and data are poised to increasingly alter many different sectors, including healthcare, education, finance, security, marketing, and transportation. Furthermore, what's unfolding reveals the ways in which these seemingly discrete sectors are intertwined; data increasingly flows through and is used by people and organizations across sectors and across domains.

Questions, concerns, and opportunities surround the big data phenomenon. Some of this is grounded in actual practices, while other hopes and fears step from mythical understandings of how technology operates. The practices that are driving the phenomenon are often invisible to the public or difficult to interpret. It is not always clear who has access to what data, let alone who is using it for what purposes or how they're accounting for potential abuses. This is particularly challenging when the data at play is transformed through big data practices, such that ownership becomes murky.

What's at stake is not simply data, but modification, aggregation, usage, and interpretation. Algorithmic accountability is an uncertain enterprise, forcing people to try to reverse engineer what's happening when something seems amiss. These dynamics leave many rightfully skittish, especially given that there's a long history of discrimination in the United States.

We generally lack the necessary frames to tease out what's at stake socially. The driving goal of this conference was to bring together diverse constituencies in an effort to examine potential approaches for framing the big data phenomenon. We examined socio-technical practices such as predicting human behavior and developing inferences alongside metaphors such as data supply chains, structural issues like inequalities and asymmetries alongside unwanted outcomes such as those that occur when interpretation goes wrong. These conceptual models provided anchors for debates, while also highlighting how much more work is needed to develop strong operational frames.

Throughout the conference, participants questioned how big data is different from other related phenomena. Some of the recurring themes included: the relationship between data analytics and power; the emergence of new forms of discrimination along with the reinforcement of existing forms; the accountability of data caretakers in unregulated spaces; and trade-offs between increased surveillance of vulnerable populations and gaps in data that might be used to empower or assist groups rendered invisible by targeted data collection. When discussing interventions, participants raised the need for updated ethical and legal frameworks alongside the complex dynamics of consent and transparency and the power of designers and algorithms.

Although there was no attempt by the organizers to formally define big data or the surrounding algorithmic phenomena, participants kept returning to questions about what is different with these practices compared to other forms of data collection and analysis. What, specifically, do algorithms allow people and groups to do, for better or worse? Are algorithms permitted to make calculations that individuals are prohibited from making explicitly? For instance, algorithms can easily use proxies for protected variables, such as religion or race, in a decision-tree that can effect discrimination on prohibited grounds, without being directly accountable for doing so. However, this is not particularly different than other forms of discrete judgment calls that occur in routine social interaction. Is there something more insidious about decision making around developments of big data that makes people uncomfortable? Inquiries and juxtapositions like these raised the larger issue of, what does big data replace? What are its alternatives? Are our concerns about big data different from our concerns about quantitative thinking?

One of the most salient themes of the day was the relationship between big data uses and power, notably the issue of power differentials in relation to data analytics. This issue is most visible in the conversations on privacy and surveillance, but it even plays out when well-intended organizations are seeking to help people. The role of power was also discussed in terms of different kinds of relationships between organizations and data subjects as well as the existence of power in the design, production, and use of algorithmic technologies. For example, what power does an algorithm implemented by a search engine or a social media platform have to alter information flows? How accurately can data that is removed from its original context and reproduced in a bigger composite picture represent an individual, and how accurate are the subsequent inferences and connections made about them? Who or what is reading inferences or making connections about us based on our data, and when it is appropriate or inappropriate to make those connections? Similarly, what are the harms or benefits to an individual or user who is reading an inference into the products presented to them as a result of data-driven associations in an online or offline, and private or public marketplace of products, services, or information?

Discussions of power inevitably led to concerns about the entrenchment or reinforcement of existing forms of inequality and discrimination. Marginal populations may be subjected to increased surveillance by both public and private actors. If predictive algorithms deem them to be "at-risk," they may be labeled as such and further marginalized. The problem with algorithmic discrimination is that it may be harder to detect. How can an individual be sure that the ads she is being shown are similar to what other users see? What if algorithms determine that she only sees advertisements or opportunities based on her race and gender? How might her agency then be limited and what does it mean if algorithms can affect educational, housing, and employment choices? Discrimination can occur when individuals volunteer their personal data, but this data is then applied to family members or others who did not choose to release this information (as with genomic data, for example). Individual choices can affect entire groups of people and the big data phenomenon makes it difficult to draw barriers between the personal and the collective.

As participants grappled with systemic concerns, they turned to a question of accountability. What would it mean to hold an algorithm, or its implementers, or designers, or any other agent or organization associated with it, accountable for some of its negative outcomes? When an algorithmic process or system results in a negative social outcome, who can be sought out to rectify it? Invisibility was a problem that emerged not only in relation to marginalized groups who may be left out of big data analytics, but also when it comes to data brokers and other actors who are less visible and thus not held accountable. The discussion on accountability intersected with a

discussion on data interpretation, or misinterpretation. Who is equipped to interrogate algorithmic systems? Peer reviews, accountability systems, credibility, and generally, a guarantee that the information being disseminated has some integrity to it is fundamental to the public trust in data and data-driven decision making processes. Several participants voiced the thought that if you offload data-interpretation to individuals, the organization or experts responsible for data use and interpretation evade a responsibility that should be theirs.

Discussion about collective responsibility circled around the fact that those in marginalized or precarious positions may be hurt by increased surveillance, but that opting out or losing access to data is just as pressing of a concern. Which people or communities come under closer surveillance when policing authorities adopt predictive technologies that use data analytics to anticipate criminal geographies or people? Who is not represented in the data collected for analysis? For instance, those who don't possess smartphones or who don't participate in particular data environments may be rendered yet more invisible when policies are developed around data-driven outcomes. Will big data cleave greater divides between the haves and have-nots? If there are benefits to using data tracking and prediction methods, how can individuals gain access to this information and use it for their own ends? Those made invisible by the big data phenomenon may be left out of civic improvement projects, while those who are tracked by data analytics may lose access to their information or may be unaware of how it is being used. How can researchers attempt to respectfully integrate marginal communities into their studies? How can individuals learn more information about what is being done with their data, tracing it as it move across platforms and goes to different data caretakers? The need for individuals to be granted more agency and room to make choices was emphasized throughout the day.

Discussions of accountability forced the group to reflect on the ethical and legal frameworks needed to manage the collection, use, and maintenance of data. As data moves from public to private sectors and changes hands, and as the lines between government and corporate power are blurred, it is unclear who or what should be regulating this movement. Private companies like ancestry DNA tests or online fitness tracking applications, for example, may not be covered under the jurisdiction of Health Insurance Portability and Accountability Act (HIPAA) and are thus not be subject to the same kinds of regulation. How can entities be held accountable for their actions in this undefined space? For entities that exist outside of regulatory jurisdiction, how can these unregulated spaces be monitored? One conceptual frame that was offered was to think through data due process, or finding a way of subjecting new technologies to ethical and legal protocols.

Legal frames are not always the best way to grapple with ethical concerns. Indeed, much of what's at stake isn't running afoul of laws *per se*, even if it is still making people uncomfortable. What is that individuals find unsettling about data analytics? At various times, participants noted the problems associated with data uncovering information we'd rather have hidden or making inaccurate assumptions about us. As data moves across different spaces and is shared or combined with other datasets, it loses its contextual meaning. How do we find ethical and legal frameworks for handling data's predictive capabilities, its potential persistence, and its ability to move across various sectors and platforms?

Participants in the event also discussed the potential power that algorithms and their designers have, as well as the ways this power is mitigated by others groups or organizations' interests. Particular individuals create algorithms and they may have different goals than the companies they work for. Developers may not properly account for the effects of their own metrics or of users' reactions and hence may produce erroneous results. Managers may override technical specialists' judgments and push for the deployment of algorithms whose use is not technically warranted or recommended by technical specialists. On the other hand, developers may not know how to translate ethical guidelines into code. Regardless of whether they are aware and their actions are intentional, they have the means to hide value judgments undetectably in algorithms. Engineers thus have a great deal of power and the fact that marginal groups and women are less likely to be trained as engineers, for a variety of structural reasons, may lead to even more subtle forms of discrimination and exacerbate existing inequalities.

Although consent is often identified as a key intervention, how this plays out in practice is often complicated. Consent often gets boiled down to obtuse contracts. Moreover, as different companies take ownership of data, terms of service and terms of use may become obsolete. Data is portable so the same regulations do not apply to all actors. While an individual may agree to give her or his information to one application, for example, she or he might be upset to learn that this data has been sold or, if the original company merges with another entity, has changed hands without giving proper notice. Even when individuals consent to provide seemingly innocuous information like a zip code or birth date, all data has the capacity to become personally identifying information (PII) through combination and re-identification. Even if people do give consent and allow companies or public entities to use their data, do they know what they are consenting to?

Discussion about consent led to further discussion of transparency, and other ways of arriving at certain knowledge about what happens to our data and how it is interpreted. Some participants argued that the algorithms that produce data should be made public. Others pointed out that if people lack the education or knowledge to fully

understand the implications of what they are consenting to, public access might not really offer transparency, just as terms of service don't really address issues of consent How can we ensure that the general population is educated about data analytics and can understand the ways in which their data might be used? Because of data's mobility, transparency at one stage does not guarantee transparency at the next one. In order for transparency to be realized, data collection and use would be to be visible and fully explained at every step of the data supply chain. Does transparency refer to the data, the actors, or the information flow? Whom is transparency for, and to what end? If individuals are able to see what is happening to their data, does this place the burden of responsibility on end users rather than on those actually in power? Can individuals not only know what is being done with their data but also access it and use it themselves?

Through discussions and debates during the whole day, it became evident that existing approaches to managing the potential harms of big data are ineffective. Existing approaches include privacy notices, informed consent, and data use and sale agreements. These are ineffective for many reasons, including technical problems, regulatory issues, social implications, organizational complications, and market-related concerns. What the discussions surrounding the big data phenomenon make clear is that divisions between public and private, individual and collective, and opportunity and problem are murky, requiring a sophisticated legal and ethical framework for handling these emerging problems in a world where various types of data intersect and are moved across multiple platforms by many parties, mostly without regulatory oversight. Beyond a narrow regulatory vision, how might experts from a variety of sources, including academia, civil rights groups, industry, and others help to bridge the gap in the wider public understanding of what the big data phenomenon is about, and how it affects people in varied ways? In sum, there is a wealth of expertise that can be applied to creating a comprehensive framework for understanding the social, cultural, and ethical dimensions of big data, but there are still a lot of unanswered questions about the best ways to address the benefits and risks associated with the big data phenomenon.

As is evident from this write-up, the day prompted more questions than answers. Given the importance of these issues, it is clear that much more work is needed to systematically untangle different aspects of this puzzle in order to build a coherent framework that will allow all involved to move forward in an ethical fashion. Achieving this will be no small feat.