Solon Barocas and Helen Nissenbaum

# Computing Ethics
# Big Data's End Run Around Procedural Privacy Protections

*Recognizing the inherent limitations of consent and anonymity.*

**P**RIVACY PROTECTIONS FOR the past 40 years have concentrated on two types of procedural mitigations: informed consent and anonymization. Informed consent attempts to turn the collection, handling, and processing of data into matters of individual choice, while anonymization promises to render privacy concerns irrelevant by decoupling data from identifiable subjects. This familiar pairing dates to the historic 1973 report to the Department of Health, Education & Welfare, *Records, Computers and the Rights of Citizens*, that articulated a set of principles in which informed consent played a pivotal role (what have come to be known as the Fair Information Practice Principles (FIPPs)) and proposed distinct standards for the treatment of statistical records (that is to say, records not identifiable with specific individuals).

In the years since, as threats to privacy have expanded and evolved, researchers have documented serious cracks in the protections afforded by informed consent and anonymity.[a] Nevertheless, many continue to see them as the best and only workable solutions for coping with privacy hazards.[b] They do not deny the practical challenges,

---

a  For concise explanations for the growing lack of confidence in anonymization and informed consent, see Narayanan and Shmatikov[7] and Solove, D.J.,[11] respectively.
b  As evidenced by their continued centrality in the policy documents under discussion in both the U.S. and Europe and the entrenched role that they play in commercial privacy policies and codes of conduct.

but their solution is to try harder—to develop more sophisticated mathematical and statistical techniques and new ways of furnishing notice tuned to the cognitive and motivational contours of users. Although we applaud these important efforts, the problem we see with informed consent and anonymization is not only that they are difficult to achieve; it is that, even if they were achievable, they would be ineffective against the novel threats to privacy posed by big data. The cracks become impassable chasms because, against *these* threats, anonymity and consent are largely irrelevant.[1]

### Informed Consent

Long-standing operational challenges to informed consent ("notice and choice") have come to a head with online behavioral advertising. Companies eager to exploit readily available transactional data, data captured through customer tracking, or data explicitly provided by users have crafted notices comprising a mish-mash of practices and purposes. Even before big data entered common parlance, authors of privacy policies faced profound challenges in trying to explain complex flows, assemblages, and uses of data. Anxious to provide easily digestible accounts of information practices, they have confronted something we have called the *transparency paradox*[8]: simplicity and fidelity cannot both be achieved because details necessary to convey properly the impact of the information practices in question would confound even sophisticated users, let alone the rest of us.

Big data extinguishes what little hope remains for the notice and choice regime. Stated simply, upfront notice is not possible because new classes of goods and services reside in future and unanticipated uses.[4] Two decades ago, experts were already warning that data mining posed insurmountable challenges to the foundations of emerging privacy law;[5,9,10] the situation now is worse than they had feared. Even if it were possible, as a theoretical matter, to achieve meaningful notice and render informed, rational decisions concerning our *own* privacy, these decisions would nevertheless affect what companies can then infer about *others*, whether or not these others

have consented. The willingness of a few individuals to disclose information about themselves may implicate others who happen to share the more easily observable traits that correlate with the traits disclosed. We call this the *tyranny of the minority* because it is a choice forced upon the majority by a consenting minority.

How might this happen? Consider the familiar trope about "the company you keep." What your friends say or do (on social networking sites, for example) can affect what others infer about you.[c] Information about social ties, however, is unnecessary when drawing such inferences, as we learned from Target's infamous pregnancy prediction score. To discover the telltale signs of pregnancy, Target looked over the purchase histories of those few customers who also made use of the company's baby shower registry. Analysts then employed data mining techniques to isolate the distinctive shopping habits of these women and then searched for similar purchases in other customers' records to identify those who were likely to be pregnant. Target was thus able to induce a rule about the relationship between certain purchases and pregnancy from what must have been a tiny proportion of all its customers.

When analysts can draw rules from the data of a small cohort of consenting individuals that generalize to an entire population, consent loses its practical import. In fact, the value of a particular individual's withheld con-

---

c   See, for example, Mislove, M. et al.[6]

---

## Long-standing operational challenges to informed consent ("notice and choice") have come to a head with online behavioral advertising.

sent diminishes the more effectively a company can draw inferences from the set of people that do consent as it approaches a representative sample. Once a dataset reaches this threshold, the company can rely on readily observable data to draw probabilistic inferences about an individual, rather than seeking consent to obtain these details. This possibility also reveals why the increasingly common practice of vacuuming up innocuous bits of data may not be quite so innocent: who knows what inferences might be drawn on the basis of which bits?

### Anonymity

Most online outfits make a serious show about anonymity.[d] But when they claim they only maintain anonymous records,[12] they rarely mean they have no way to distinguish a specific person—or his browser, computer, network equipment, or phone—from others. Nor do they mean they have no way to recognize him as the same person with whom they have interacted previously, to associate observed behaviors with the record assigned to him, or to tailor their content and services accordingly. They simply mean they do not rely on the limited set of information commonly conceived as "personally identifiable" (for example, name, Social Security number, date of birth, and so forth), while still employing unique *persistent identifiers*. Hence the oxymoronic notion of an "anonymous identifier"[2]—more accurately labeled a pseudonym. These identifiers are anonymous only insofar as they do not depend on traditional categories of identity while still serving the function of persistent identification.

Such practices may make it more difficult for someone to show up on a person's doorstep with a folder full of embarrassing, discrediting, or incriminating facts, but they do nothing to limit the ability of these companies to draw upon this information in shaping a growing share of everyday experiences that take place on these companies' platforms. Stated differently,

---

d   That this type of anonymity bears little resemblance to the rigorous specifications of anonymity developed by computer scientists is not our concern here; ours is a discussion of the *value* of anonymity evinced by these techniques.

---

while anonymous identifiers can make it more difficult to use information about a specific user outside an organization's universe, they do nothing to alleviate worries individuals might have about their fates within it—the information they are presented, the opportunities they are offered, or the way they are treated in the marketplace.

Whatever protections this arrangement offers are further undermined by the kinds of inferences companies can draw having discovered patterns in large assemblages of diverse datasets. A company that may have been unable to learn about individuals' medical conditions without matching records across datasets using personally identifiable information may be able to infer these conditions from the more easily observable or accessible qualities that happen to correlate with them.[13] If organizations become sufficiently confident to act on these uncertain inferences, the ability to draw these inferences will pose as serious a threat to privacy as the increasingly well-recognized risk of de-anonymization. But rather than going to the trouble of attempting to re-associate "anonymized" medical files with specific individuals, companies might instead discover patterns that allow them to estimate the likelihood someone has a particular medical condition. That certain institutions could meaningfully affect a person's experiences and prospects in the absence of identifying information or without violating record-keepers' promises of anonymity defies the most basic intuitions about the *value* of anonymity.

## We Are Not Saying…

*There is no role for consent and anonymity in privacy protection.* Consent and anonymity should not bear, and should never have borne, the entire burden of protecting privacy. Recognizing their limits allows us to assess better where and under what conditions they may perform the work for which they are well suited.

*Privacy loses the trade-off with big data.* This tired argument misunderstands the nature and value of privacy and mistakes means for ends. Weaknesses in existing *procedures* for protecting privacy do not undercut the viability of privacy itself.

## Big data extinguishes what little hope remains for the notice and choice regime.

*We need to try even harder to achieve fail-safe anonymization and effectively operationalize notice and consent.* Though worthy goals, the practices described here bypass not only weak mechanisms but also defeat the ideal.

### What to Do?

Mathematicians and computer scientists will continue to worry about re-identification. Policymakers will continue down the rabbit hole of defining personally identifiable information and informed consent. Social scientists and designers will continue to worry about refining notice and choice. In the meantime, miners of big data are making end runs around informed consent and anonymity.

A lesson may be drawn from biomedicine where informed consent and anonymity function against a rich ethical backdrop. They are important but not the only protective mechanisms in play. Patients and research subjects poised to sign consent forms know there are limits to what may be asked of them. Treatment or research protocols that lie outside the norm or involve a higher than normal risk must have passed the tests of justice and beneficence. In other words, clinicians and researchers must already have proven to their expert peers and institutional review boards that the protocols being administered or studied are of such great potential value to the individual subject or to society that the reasonable risks are worthwhile. Consent forms have undergone ethical scrutiny and come at the end of a process in which the values at stake have been thoroughly debated. The individual's signature is not the sole gatekeeper of welfare.

By contrast, informed consent and anonymity have served as the sole gatekeepers of informational privacy.

When consent is given (or not withheld) or the data is anonymized, virtually any information practice becomes permissible. These procedural mitigations have long relieved decision-makers of the burden of rendering judgment on the substantive legitimacy of specific information practices and the ends that such practices serve. It is time to recognize the limits of purely procedural approaches to protecting privacy. It is time to confront the substantive values at stake in these information practices and to decide what choices can and cannot legitimately be placed before us—for our consent. **ⓒ**

### References
1. Barocas, S. and Nissenbaum, H. Big data's end run around anonymity and consent. In *Privacy, Big Data, and the Public Good: Frameworks for Engagement.* J. Lane, V. Stodden, S. Bender, and H. Nissenbaum, Eds. Cambridge University Press, NY, 2014.
2. Barr, A. Google may ditch 'cookies' as online ad tracker. *USA Today* (Sept. 17, 2013).
3. Cate, F.H. The failure of fair information practice principles. In *Consumer Protection in the Age of the "Information Economy."* J.K. Winn, Ed. Ashgate, Burlington, VT, 2006, 341–378.
4. Cate, F.H. and Mayer-Schonberger, V. Notice and consent in a world of big data. *International Data Privacy Law 3,* 2 (May 20, 2013), 67–73.
5. Klösgen, W. KDD: Public and private concerns. *IEEE Expert: Intelligent Systems and Their Applications 10,* 2 (Feb. 1995), 55–57.
6. Mislove, M. et al. You are who you know: Inferring user profiles in online social networks. In *WSDM '10 Proceedings of the Third ACM International Conference on Web Search and Data Mining.* ACM, NY, 2010, 251–60; DOI: 10.1145/1718487.1718519.
7. Narayanan, A. and Shmatikov, V. Myths and fallacies of 'personally identifiable information.' *Commun. ACM 53,* 6 (June 2010), 24; DOI: 10.1145/1743546.1743558.
8. Nissenbaum, H. A contextual approach to privacy online. *Daedalus 140,* 4 (Oct. 2011), 32–48; DOI: 10.1162/DAED_a_00113.
9. O'Leary, D.E. Some privacy issues in knowledge discovery: The OECD personal privacy guidelines. *IEEE Expert: Intelligent Systems and Their Applications 10,* 2 (Apr. 1995), 48–59.
10. Piatetsky-Shapiro, G. Knowledge discovery in personal data vs. privacy: A mini-symposium. *IEEE Expert: Intelligent Systems and Their Applications 10,* 2 (Apr. 1995), 46–47.
11. Solove, D.J. Privacy self-management and the consent dilemma. *Harvard Law Review 126,* 7 (May 2013), 1880.
12. Steel, E. and Angwin, J. On the Web's cutting edge, anonymity in name only. *The Wall Street Journal* (Aug. 4, 2010).
13. Walker, J. Data mining to recruit sick people. *The Wall Street Journal* (Dec. 17, 2013).

**Solon Barocas** (sib237@nyu.edu) is a postdoctoral research associate at the Center for Information Technology Policy at Princeton University.

**Helen Nissenbaum** (helen.nissenbaum@nyu.edu) is a professor of media, culture, and communication at New York University.