# Assignment 8

## Question 1. Identification risk in anonymized data

**(a) Pick two of the examples in Table 1 and describe in one or two paragraphs how the reidentification attack in both cases has a similar structure.**

In both the health insurance records case (Sweeney, 2002) and the Tastes, Ties, and Time case (Zimmer, 2010), the re-identification attack had a similar structure. The attackers used combinations of the personal characteristics to narrow down their guess. There is an anonymized dataset and a dataset with identifying information that both contain these attribute combinations. By linking these two sources, the attackers managed to uniquely identify the subject in question.

**(b) In one or two paragraphs, describe how the data could reveal sensitive information about the people in the dataset for each of your two examples in part (a).**

Both the datasets contained critical information about the subjects. In the first case, the health insurance records included not only medical records, but ZIP code, birth place, ethnicity and gender of the subject as well. (Sweeney, 2002, pp. 2) The subjects' names were removed, but there was another dataset, the voter registration list, that could match the combination of ZIP code, birth place, ethnicity and gender to individuals. In the second case, the Tastes, Ties and Time dataset included "each subjects' gender, race, ethnicity, hometown state, and major". (Zimmer, 2010, pp. 7) Although these attributes were not identifiable per se, their combinations could put together the portrait of a subject. Some majors and ethnicities were so rare that there was only one incoming student. As a matter of fact, the descriptive codebook alone revealed a lot of information to identify the "northeastern American university" (Zimmer, 2010, pp. 4).

## Question 2. Describing Ethical Thinking

**Researchers often struggle to describe their ethical thinking to each other and to the general public. After it was discovered that the Tastes, Ties, and Time study was reidentified, Jason Kauffman, the leader of the research team, made a few public comments about the ethics of the project. Read Zimmer (2010) and then rewrite Kauffman's comments using the principles and ethical frameworks that are described in Salganik (2018, Ch. 6). Specifically, Zimmer (2010) rewrite the following passages from Kauman (Sep. 30, 2008b) and Kauman (Sep. 30, 2008c).**

1. Kauffman believed that the Tastes, Ties, and Time (or "3T" henceforth) team's data collection process was in line with the principle of Beneficence, and fit in the consequentialism framework. They had a clear understanding of the risk and benefit of the data release. On the one hand, they took serious precautions to protect subject privacy, which suggests that they have put much emphasis on reducing the risk. On the other hand, as sociologists, they believed that the unique dataset could contribute a lot to research in social networks. Therefore, they concluded that the benefit of releasing the dataset outweighed the potential risk. [Zimmer (2010) citing Kauffman (Sep. 30, 2008b)]
2. Kauffman believed that the 3T team's data collection process agreed with the principle of Justice, and fit in the deontology framework. He argued that the hackers could not gain any additional benefits from their dataset other than those information already on the subjects' Facebook pages. Therefore, the project didn't add to the risk of the subjects' privacy being leaked. [Zimmer (2010) citing Kauffman (Sep. 30, 2008b)]

3. Kauffman believed that the 3T team's data collection process accorded with the principle of Respect for Persons, and fit in the deontology framework. They only collected the data that is accessible from the subject's Facebook page. They didn't approach any student for additional information. In other words, the students have consented others to browse their profiles and process these information. Therefore, it does not violate research ethics to collect data from their Facebook pages. [Zimmer (2010) citing Kauffman (Sep. 30, 2008b)]

## Question 3. Ethics of Encore

**Part (a)**

The authors provided an ethical analysis of the Encore study following the guidelines in the Menlo Report [1]. There is no easy answer to most of the questions, but the authors affirm the urgent need for ethical discussions regarding ICT research (research about or involving information and communication technologies, or "ICTR").

In ICT research, it is critical to identify the stakeholders. In the case of Encore, all the governments, censor authorities, computer users, web browses were involved. However, the scale of Encore is so wide that it is infeasible to figure out all the stakeholders. Furthermore, there was debate about whether the subjects are human. Technically speaking, Encore only collected data about the IP addresses and focused on the structure of the censorship systems, but the recorded IPs could be used to identify users. Therefore, even though the team didn't target human individuals, the internet users were indirectly impacted by the project.

The first principle for ICT research, "Beneficence", is inspired by the faith of "do no harm", and requires a systematically assessment of both the benefits and the risks. As for benefits, there is no doubt that the Encore helped to "illuminate censorship — both its motivations and the technologies behind it". (Narayanan and Zevenbergen, 2015, pp. 15) However, controversial voices exist as to whether censorship harms human rights and whether we should fight against censorship systems. Some scholars believed that the role of censorship ought to be discussed under certain political and cultural contexts. It might be implausible to analyze a censorship network with pure data science techniques.

The risks cast upon Internet users is another issue. Directors of Encore argued that "normal web browsing exposes users to the same risks that Encore does" (Narayanan and Zevenbergen, 2015, pp. 17). "The prevalence of malware and third-party trackers itself lends credibility to the argument that a user cannot reasonably control the traffic that their devices send". (Narayanan and Zevenbergen, (2015) citing Burnett and Feamster (2015)). Therefore, the directors believed their practices abided by the "minimal risk" criteria. As the authors pointed out, there are three caveats to their argument. First, neither Encore nor the third-party trackers respected "the users' expectations" (Narayanan and Zevenbergen, 2015, pp. 18). Second, the risk "may depend on the type of censored website" (Narayanan and Zevenbergen, 2015, pp. 18). Requesting the users to access less frequent websites may render them more harm. Third, the magnitude of potential harm may exceed individuals, and is beyond the researchers' control.

The principle of Beneficence is deeply rooted in consequentialist thinking (Salganik, 2018, Ch. 6). It calls for effective measures to "maximize probable benefits and minimize probable harms" (Salganik, 2018, Ch. 6). In the case of Encore, the team "limited the set of URLs that the script induced users to measure" (Narayanan and Zevenbergen, 2015, pp. 19) to mitigate the harm. Nevertheless, the authors criticized that a lot of concepts and mechanisms should have been more carefully defined and analyzed during research design.

The principle of Respect for Persons, Law, and the Public Interest is more consistent with a deontological thinking. It requires the researchers to "be transparent in methods and results", and "be accountable for actions" (Dittrich et al., 2012). Critiques for Encore focused on two aspects. First, they didn't seek informed consent from the subjects. Although they have put up notice on how Encore worked and offered options to opt out, Narayanan and Zevenbergen argued that the notice could have been strengthened. Second, their practices may have violated censorship laws under international jurisdictions.

As the authors noted, "Encore makes for a fascinating case study that presents a thick web of considerations and no easy answers" (Narayanan and Zevenbergen, 2015, pp. 22). While the Internet presents us with millions of interesting questions to explore, ICT researchers need to be more thoughtful on ethical issues. There may be more challenges awaiting in the future.

**Part (b)**

I will assess the ethics of the Encore study following the four principles: *Respect for Persons*, *Beneficence*, *Justice*, and *Respect for Law and Public Interest*. First of all, I believe the researchers didn't show enough respect for the stakeholders. They argued that they recorded nothing personal, but they were observing and recording browsing histories from a vantage point, which yielded them with more power than a normal passer-by. They should have respected the participants and sought informed consent. Second, I think the Encore team have satisfied the basic requirements of the Beneficence principle. The list of URLs they provided were common and less risky. The participants were unlikely to be exposed to harsh surveillance for trying to access these websites. Third, the principle of Justice was not well practiced. The Encore study involved people in very different cultural and political contexts. Also, more people could be impacted, not just their targets. Therefore, the consequences and risk of the study was beyond the researchers' control. Finally, the Encore study trespassed a grey zone in legal compliance. The project was based in the US, thus it was not clear whether they should abide by censorship laws in other countries. To conclude, although the study of Encore was fruitful, it alarmed future researchers to address ethical issues with more care.

# Reference

Burnett, Sam, and Nick Feamster. "Encore: Lightweight measurement of web censorship with cross-origin requests." *ACM SIGCOMM Computer Communication Review*. Vol. 45. No. 4. ACM, 2015.

Dittrich, David, and Erin Kenneally. "The Menlo Report: Ethical principles guiding information and communication technology research." *US Department of Homeland Security* (2012).

Kauman, Jason, I am the Principle Investigator...," Blog Comment, MichaelZimmer.org, http://www.michaelzimmer.org/2008/09/30/on-the-anonymity-of-the-facebook-dataset/, Sep. 30, 2008b.

—, \We did not consult...," Blog Comment, MichaelZimmer.org, http://www.michaelzimmer.org/2008/09/30/on-the-anonymity-of-the-facebook-dataset/, Sep. 30, 2008c.

Burnett, Sam, and Nick Feamster. "Encore: Lightweight measurement of web censorship with cross-origin requests." *ACM SIGCOMM Computer Communication Review*. Vol. 45. No. 4. ACM, 2015.

Salganik, Matthew J. *Bit by bit: social research in the digital age*. Princeton University Press, 2018.

Sweeney, Latanya. "k-anonymity: A model for protecting privacy." *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10.05 (2002): 557-570.

Zimmer, Michael. ""But the data is already public": on the ethics of research in Facebook." *Ethics and information technology* 12.4 (2010): 313-325.

---

1. "The Menlo Report: Ethical Principles Guiding Information and Communication Technology Research," U.S. Department of Homeland Security Science and Technology Directorate, Cyber Security Division, August 2012, accessed August 11, 2015, https://www.caida.org/publications/papers/2012/menlo_report_actual_formatted/menlo_report_actual_formatted.pdf.↵