

Assignment #8

MACS 30000, Dr. Evans

Ying Sun

1. Identification risk in anonymized data

Sweeney (2002) shows an example to illustrate that the re-identification can happen just by linking on some common fields of two datasets. More specifically, even though the medical records information looks anonymous because all the personal identification information was removed from the dataset, re-identification attacks can still be realized by linking the medical information and voter registration information on the common fields. Similarly, Zimmer (2010) illustrates the re-identification problem by citing another instance – the “Tastes, Ties and Time” project. The assumed anonymous Facebook data was quickly identified as coming from Harvard College just by using some public available data. The re-identifications in these two cases have a similar structure. Through some open available data, we can identify some sensitive information which is not originally included in the assumed anonymous dataset by matching some common fields in these datasets.

For the case in Sweeney (2002), although the patients’ names and other personally identifying information are not included in the medical records, we can link these medical records to specific people by combining with voter registration information. In the medical dataset, it contains sensitive medical information without personally identifying information. In the contrast, there are personally identifying information such as names, addresses in the voter registration dataset but no sensitive information. These two datasets share some attributes such as ZIP, birth date and gender. Through these common attributes, we can link these medical records to specific people.

For the case in Zimmer (2010), the researchers removed identifying information in the data retrieved from Facebook. However, “the uniqueness of the some of the data elements makes identifying the source of the data” (Zimmer, 2010, p.316). By using some “the freely available codebook and various public comments about the research” (Zimmer, 2010, p.316), the source of the data was identified. More specifically, based on some information such as private and co-educational institution, the number of students, the unique majors

and the method of choosing undergraduate housing, it is not difficult to confirm that the source of these data is Harvard College.

References:

Sweeney, Latanya, “K-Anonymity: A Model for Protecting Privacy,” *International Journal on Uncertainty Fuziness and Knowledge-Based Systems*, 2002, 10 (5), 557– 570.

Zimmer, Michael, “But the Data is Already Public: On the Ethics of Research in Facebook,” *Ethics and Information Technology*, 2010, 12 (4), 313–325.

2. Describing ethical thinking

Even though T3 research team had made a lot of efforts to protect the privacy of the subjects, the re-identification problem still existed in the project. After the problem was discovered, Jason Kauffman made a few public comments about the ethics of the project. However, these comments are not reasonable enough. We can use the principles and ethical frameworks to evaluate these comments thoroughly. According to Salganik (2018), “Four principles that can guide researchers facing ethical uncertainty are: Respect for Persons, Beneficence, Justice, and Respect for Law and Public Interest” (Salganik, 2018, p.294). These principles are “largely derived from two more abstract ethical frameworks: consequentialism and deontology” (Salganik, 2018, p.302). These principles and frameworks will be considered below.

First of all, Kaufman thought that the data collection process was in conformity with the principle of “Respect for Persons” (Salganik, 2018, p.295) and was reasonable in the deontology framework (Salganik, 2018, p.302). He emphasized that they “had not accessed any information not otherwise available on Facebook” (Kauffman, Sep. 30, 2008c). However, the problem was that the T3 research team didn’t receive informed consent from the actual participants – the students. The researchers only obtained the consent of the Harvard’s Committee on the Use of Human Subjects, but the actual participants didn’t get the relevant information about this T3 project in a comprehensive format and agree to contribute their personal Facebook information to this research.

Secondly, Kauffman thought that they obeyed the principle of “Beneficence” (Salganik, 2018, p.294) and their research process was proper in the consequentialism framework (Salganik, 2018, p.302). He argued that they “Were sociologists, not technologists” (Kauffman, Sep.30, 2008b). Kauffman and his colleagues hoped to know as much as possible about their research subjects. However, the researchers put the privacy of these students into a risk of being exposed. Kauffman looked at the problem only from a sociologist’s point of view and focused on the potential benefits of this research outcome, but he ignored the potential risks. In other words, he failed to strike an appropriate ethical balance between benefits and risks.

Thirdly, Encore failed to follow the principle of “Justice”. Kauffman believed that the Encore project would bring benefits to the society. However, he ignored the re-identification would impair students’ privacy. In other words, the students bore the cost of

this research while the sociologists or other research groups obtained their benefits. It is obvious that the risks and benefits are not distributed equally.

Finally, Kauffman thought that they followed the principle of “Respect for Law and Public Interest” (Salganik, 2018, p.294). He believed they had made a lot of efforts to protect the privacy of the subjects and their research followed the European Union’s guidance (Zimmer, 2010, p.319). Besides, he argued that hackers could also “crack the data and ‘see’ people’s Facebook information” (Kauffman, Sep.30, 2008b) from a consequentialism perspective of ethics. However, the re-identification problem still actually undermined the privacy of the students. Encore failed to show enough respect for law and public interest and there is no doubt that these researchers could do more in this aspect.

References:

Zimmer, Michael, “But the Data is Already Public: On the Ethics of Research in Facebook,” *Ethics and Information Technology*, 2010, 12 (4), 313–325.

Kauffman, 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.

Salganik, Matthew J., *Bit by Bit: Social Research in the Digital Age*, Princeton University Press, 2018.

3. Ethics of Encore

(a)

Narayanan and Zevenbergen (2015) firstly affirm the research value of the Encore study and further make an ethical analysis of this research project. Basically, the authors follow the structure and principles in “The Menlo Report” (Narayanan and Zevenbergen, 2015, p.11) in their ethical analysis.

First of all, identifying the stakeholders is of vital importance. In the Encore study, however, it is not an easy task because “analyzing all potential stakeholders individually in this project is infeasible” (Narayanan and Zevenbergen, 2015, p.12). What’s worse, this principle may undermine the goal of this study – scalability. Besides, it is necessary to make a judgement that whether Encore is a human-subjects research. In the eyes of the researchers, Encore is not a human-subjects research because they only collected the IP addresses without any further identification. However, “human subject research should now be considered as ‘human-harming research’” (Narayanan and Zevenbergen, 2015, p.13). More specifically, the research subjects – web users, may bear the “costs” in the Encore. In this sense, Encore should not be treated as a pure technical research because the inevitable interaction of technology and human.

Secondly, the authors conduct a systematic analysis of harms and benefits by citing the principle of Beneficence. Obviously, Encore provides an illumination of “motivations and technologies” (Narayanan and Zevenbergen, 2015, p.15) behind the web censorship, which is meaningful in the political science. However, this kind of censorship may impair the people’s freedom to express their opinions and “the freedom to seek, receive and impart information” (Narayanan and Zevenbergen, 2015, p.15).

Within the consequentialism framework, Encore researchers thought they achieved the “minimal risk” (Narayanan and Zevenbergen, 2015, p.17) compared with normal web browsing. However, the authors come up with several different opinions. Firstly, it is unreasonable to make such a comparison. An “ethical race to the bottom” (Narayanan and Zevenbergen, 2015, p.18) should be avoided. Secondly, “the probability and magnitude of harm may depend on the type of censored website” (Narayanan and Zevenbergen, 2015, p.18). What’s more, other unexpected types of harms may also occur.

In the deontological thinking, the principle of “Respect for Persons, Law and Public Interest” should be taken into consideration. Encore is not blameless in this aspect. The Encore researchers failed to seek informed consent from the study subjects even though the Encore website contains a related statement. Besides, the measurement actions of Encore might make “the issue of jurisdiction” (Narayanan and Zevenbergen, 2015, p.21) ambiguous.

(b)

According to Salganik (2018), “Four principles that can guide researchers facing ethical uncertainty are: Respect for Persons, Beneficence, Justice, and Respect for Law and Public Interest” (Salganik, 2018, p.294). I will assess the ethical quality of Encore following the four principles.

First of all, the Encore study violates the principle of “Respect for Persons” to some extent. Even though we could find an informed statement at the bottom of the Encore website, it is not a proper way to collect measurements before obtaining the informed consent. They definitely have several methods to strengthen notice. Second, it is hard to evaluate whether Encore derogates from the principles of “Beneficence”. There is no doubt that Encore will contribute a lot to the Internet censorship, but the harms and risks are difficult to evaluate. The risks and harms depend on the type of censored website and it is possible that other unexpected types of harms may also occur. Third, Encore project is not blameless in terms of “Justice”. More precisely, the benefits and risks are not distributed fairly. Because the Encore project involves different censorship systems in different cultural and political contexts. Finally, the Encore project doesn’t show enough respect for law and public interest. There are different privacy and data protection laws in the worldwide, it is impossible to “enumerate all possible legal risks to Encore users” (Narayanan and Zevenbergen, 2015, p.22). So Encore researchers cannot guarantee that this project would not violate any local laws.

References:

Burnett, Sam and Nick Feamster, “Encore: Lightweight Measurement of Web Censorship with Cross-Origin Requests,” 2015.

Narayanan, Arvind and Bendert Zevenbergen, “No Encore for Encore? Ethical QUestions for Web-based Censorship Measurement,” Technology Science, Decem- ber 15 2015.

Salganik, Matthew J., Bit by Bit: Social Research in the Digital Age, Princeton University Press, 2018.