# MACS 30000 Assignment 8

Shanglun Li

December 3, 2018

## Problem 1.

(a) The two examples sharing the similar structure of re-identification attacks are the health insurance records (Sweeney, 2002) and the Netflix movie rating data (Narayanan and Shmatikov, 2008). In order to retrieve the sensitive information of a specific individual, the authors always have two datasets. One of them contains sensitive information of individuals but no personally identifying information, the other one contains personally identifying information but no sensitive information. In addition, the two datasets share some common variables so that one can merge these two datasets according to the same variables and combine personally identifying information with sensitive information.

(b) Sweeney (2012) managed to identify William Weld (68th Governor of Massachusetts)'s medical records to illustrate the threat of healthcare information leakage. Thus, the author retrieves the two datasets in the database. One of them is a medical dataset containing individuals' potentially sensitive information like the medical visit date, diagnosis, procedure, medication, and total medical charge. In this dataset, it doesn't include identifying information, but contains individuals' ZIP code, sex, and date of birth. The other dataset is "the voter registration" data (Sweeney, 2002, p2) including people's names together with ZIP code, sex, and date of birth. In order to retrieve William Weld's medical records, the author can merge these two datasets by the common fields (ZIP code, sex, and date of birth).

In Narayanan and Shmatikov (2008), the Netflix movie rating dataset contains people's rating to movies. In addition, the personally identifying information has been removed. However, if someone has the information about when some individual rated a specific movie and the corresponding rating, even it is not perfectly accurate, there will be a large chance that the individual's record in the Netflix movie rating dataset can be identified. Then, this individual's sensitive information can also be accessed.

# Problem 2.

Indeed, Researchers should make every effort on their research to get significant findings. However, they should also prevent potential harms to research subjects according to the principle of beneficence (Salganik, 2018, p296). An increasing number of researchers is finding more detailed data to get more knowledge. Kauffman and his colleagues wanted the dataset to include more details to create more new knowledge by studying the dataset. However, when more detailed personal informations are retrieved by the researchers, the more dangerous the research subjects can be, since there will be a higher probability that their personal informations will be identified and their privacy will be revealed. Thus, Kauffman and his colleagues should have better balanced the gain from the data and the potential harm to the students.

Kauffman comment that hackers could also get people's private information directly from Facebook. He focused on the expected outcome of researches instead of the potential risk of revealed sensitive information. He also defended the validity of their aggregated dataset as it "contains almost no information that isnt on Facebook" (Kauffman, Sep.30, 2008b). Although the procedure of building this Facebook dataset might be ethically questionable, the consequence was just that the students' information that was already available on Facebook and included in the dataset. His consequentialism view demonstrates that it doesn't make any difference since those who want to crack the data can get the same information from Facebook.

Kauffman was arguing from a deontology view of ethics by mentioning that the research only included the information that can be found on Facebook and did not disclose any subject's privacy. They just included public available informations in their research and respected every research subject. In addition, they protected the content of the dataset well. He also claimed that if someone cracked the dataset and made use of the private informations to do harm to the subjects, people should blame the hacker instead of the researchers, since he thought he finished his duties well to protect the database.

# Problem 3.

(a)     In Burnett and Feamster (2015), the controversial censoring measurement project Encore is discussed. Encore study injects "an invisible element into the page that tells the browser to download and execute a code piece" (Burnett and Feamster, 2015, p2). This code will allow the browser to send requests to certain websites without informing the user about their success and sending the results to the researchers. Although Encore is an efficient, scalable strategy and has collected valuable large-scale censorship data around the world, an ethical dilemma remains.

Some intrinsic features of the Internet have enabled computer science research to track the behavior of Internet users "without the consent of the user" (Burnett and Feamster, 2015, p4). Some of them are extremely disputed as a result of they create use of web security holes, whereas Encore, probes devices on the net while not "exploitation of any security holes" (Burnett and Feamster, 2015, p4). Alternative options like avoiding hiring volunteer or government intervention, and manufacturing global-wide fine-grained knowledge during a scalable and automatic manner, create Encore even a lot of technically enticing. However, because the committee of "ACM SIGCOMM 2015"(Burnett and Feamster, 2015, p7) realized, there are 3 moral issues in Encore. First, being a "third-party requests used for ad tracking" (Burnett and Feamster, 2015, p7), Encore ought to have au fait users. Second, it probably exposes the users within

the risk of being chastised by the regime if their on-line activities are monitored; Third, if the users were truly au fait, they "would be unlikely to consent" (Burnett and Feamster, 2015, p7) since they live underneath censorship.

The author consider the question "who are the stakeholders"(Burnett and Feamster, 2015: p.9) of Encore to be difficult because when the "users browser sends a request to a potentially censored website", "the users IP address may be recorded by the server hosting that website" (Burnett and Feamster, 2015, p9). Then, the author also commented that the question whether Encore is a "human-subjects research" (Burnett and Feamster, 2015, p10) is also hard to answer. "Neither the Princeton nor the Georgia Tech IRB considered Encore to be human-subjects research" (Burnett and Feamster, 2015, p10), since it neither involves "intervention or interaction with individual", nor includes "identifiable private information" (Burnett and Feamster, 2015, p10). Moreover, the author considered the potential harm of Encore to individual internet users is difficult to assess "due to the complex, dynamic, and innovative nature of the Internet" (Burnett and Feamster, 2015, p11). From a consequentialism view, the authors of Encore argued that even without Encore, people are exposed to the risk of third-party tracks, while others argue that researchers "should not participate in and facilitate a race to the bottom" (Burnett and Feamster, 2015, p13).

Burnett and Feamster(2015) also discussed about how to reduce the harms of Encore in terms of "informed consent, transparency and accountability" (Burnett and Feamster, 2015, p14). The difficulty for informed consent is that it is hard to get it. It can potentially reduce the benefit of Encore, and put higher risks on the users. For transparency, Encore provides website operators with an instruction on how to inform visitors of Encore and provide the option of disabling it. However, it is optional for them. In term of accountability, it is complicated because machines are around the world. Thus, researchers cannot be completely sure whether they violate any law of any region or not.

(b) In my opinion, the Encore study failed to obey the principle of Respect for Persons since it

ask for the users' consent without informing them any detail. I think it will be dangerous for users to make a consent to connect to certain blocked websites. Nowadays, people save all their personal information in the computer and internet, like credit card information, personal address, and password. Thus, if someone take the advantage of such a connection to the blocked website, the damage to the users will be irreversible.

# References

Barbaro, Michael and Tom Jr. Zeller, "A Face Is Exposed for AOL Searcher No. 4417749," New York Times, August 9, 2006.

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

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.

Kauffman, Jason, "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.

Mayer, Jonathan, Patrick Mutchler, and John C. Mitchell, "Evaluating the Privacy Properties of Telephone Metadata," Proceedings of the National Academy of Sciences of the USA, 2016, 113 (20), 55365541.

Montjoye, Yves-Alexandre de, Laura Radaelli, Vivek Kumar Singh, and Alex Sandy Pentland, "Unique in the Shopping Mall: On the Reidentifiability of Credit Card Metadata," Science, 2015, 347 (6221), 536539.

Narayanan, Arvind and Bendert Zevenbergen, "No Encore for Encore? Ethical QUestions for Web-based Censorship Measurement," Technology Science, December 15 2015.

Narayanan and Vitaly Shmatikov, Robust De-Anonymization of Large Sparse Datasets, 2008.

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 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), 313325.