Free the Data, Manage the Risks

Sponsored by California HealthCare Foundation

Moderator: Andy Krackov Senior Program Officer, Market and Policy Monitor Program, CHCF Khaled El Emam, CEO, Privacy Analytics Daniel Barth-Jones, Assistant Professor, Columbia University Barbara Rudolph, Consultant, NAHDO

CALIFORNIA HEALTHCARE FOUNDATION

29th Annual NAHDO Conference



PRIVACYANALYTICS Data Anonymization Solutions

Free the Data Manage the Risks

Khaled El Emam (PhD)

NAHDO 29th Annual Conference 9th October 2014

www.privacyanalytics.ca | 855.686.4781 info@privacyanalytics.ca

251 Laurier Avenue W, Suite 200 Ottawa, ON Canada K1P 5J6



- It is possible to share data and:
 —Protect privacy
 —Meet regulatory and legal
 - requirements
 - -Get high quality data



There are some standards and more detailed (operational) standards are in development for deidentification











- There are different levels of data release, and the de-identification scheme needs to match that:
 - -Public
 - -Quasi-public
 - -Non-public



De-identification is a risk management exercise that takes into account the context of the data release



© 2014 Privacy Analytics, Inc.



There are many precedents on what is acceptable risk – this is not something we should be debating (there are bigger issues)



More sophisticated methods (computational and statistical methods) are needed to be able to generate high quality data and meet the risk thresholds



There is a real need for multidisciplinary education and certification around deidentification practices to start building a community of practice



The discourse on reidentification attacks is of deep concern from a scientific, ethical, and integrity perspective



Our biggest challenge is transitioning good de-identification methodologies into practice



10

The existing legal framework is fine as it is and is quite robust – we can do a lot within it to free the data







PRIVACYANALYTICS Data Anonymization Solutions

NAHDO 29th Annual Conference Oct 8-9 2014

Free the Data, Manage the Risks: Reality Check for Data Privacy Risks

Daniel C. Barth-Jones, M.P.H., Ph.D Assistant Professor of Clinical Epidemiology, Mailman School of Public Health Columbia University

E-mail: db2431@columbia.edu

On Twitter : <u>@dbarthjones</u>

A Historic and Important Societal Debate is underway...



Public Policy Collision Course

The Societal Value of De-identified Data

- Properly de-identified health data is an *invaluable "public good"*. The broad availability of de-identified data is an essential tool for society supporting scientific innovation and health system improvement and efficiency.
- De-identified data does and can serve as the engine driving forward innumerable essential health systems improvements: quality improvement, health systems planning, healthcare fraud, waste and abuse detection, and medical/public health research (e.g. comparative effectiveness research, adverse drug event monitoring, patient safety improvements and reducing health disparities).
- De-identified health data greatly benefits our society and provides strong privacy protections for the individuals. As the promise of EHRs and Health IT yields richer de-identified clinical data, the progress of our nation's healthcare reform will likely be built on a foundation of such de-identified health data.

General Health Status among US Adults*, by Race or Ethnicity



Counting and Tabulating is Essential to Public Health and Population Science...

- -The foundational acts of counting and tallying individual characteristics underlie our ability to analyze the population distributions and determinants of disease—which is essential to medical and population health science.
- -But some risk of re-identification exists with every characteristic that we collect and report.
- -Thus, the important ongoing debate about health data de-identification and the ethical and public policy implications for research conducted with de-identified data.

Unfortunately, de-identification public policy has often been driven by anecdotal and limited evidence, privacy folklore, and targeted reidentification demonstration attacks which fail to provide reliable evidence about real world re-identification risks



Bloomberg Our Company Professional Anywhere



Your Health Data for Sale: Who's Selling, Buying?

Misconceptions about HIPAA De-identified Data:

"It doesn't work..." "easy, cheap, powerful re-identification" (Ohm, 2009 "Broken Promises of Privacy")

*Pre-HIPAA Re-identification Risks {Zip5, Birth date, Gender} able to identify 87%?, 63%?, 27%? of US Population (Sweeney, 2000, Golle, 2006, Sweeney, 2013)

Reality: HIPAA compliant de-identification provides important privacy protections

- Safe harbor re-identification risks have been more recently estimated at 0.04% (4 in 10,000) (Sweeney, NCVHS Testimony, 2007)
- Safe Harbor de-identification provides protections that have been estimated to be a minimum of 400 to 1000 times more protective of privacy than permitting direct PHI access.
 (Benitez & Malin, JAMIA, 2010)

Reality: Under HIPAA de-identification requirements, reidentification is expensive and time-consuming to conduct, requires serious computer/mathematical skills, is rarely successful, and uncertain as to whether it has actually succeeded Misconceptions about HIPAA De-identified Data:

"It works perfectly and permanently..."

Reality:

- Perfect de-identification is not possible
- De-identifying does not free data from all possible subsequent privacy concerns
- Data is never permanently "de-identif<u>ied"</u>... (There is no guarantee that de-identified data will remain de-identified regardless of what you do to it after it is de-identified.)
- Simply collapsing your coding categories until the data is "k-anonymous" can make the data unsuitable for many statistical analyses

Myth of the "Perfect Population Register"

- The critical part of many re-identification efforts that is often assumed by disclosure scientists is the assumption of a perfect population register.
- All Population registers will have data errors and be incomplete to some extent. (e.g. Nationwide voter registration levels typically are about 70%)
 - -However, some types of data errors are more critical than others.
 - -Persons who are not included in population registers will not have identifiers which can be linked to identify them.
 - Persons who are not in a population register can not reidentified, but they also indirectly reduce the probability of correct re-identification for others.
 - If only one person within a quasi-identifier set is missing from the population register, then the probability of correct reidentification drops to 50%; if two persons are missing, then the probability of correct re-identification is 33%, and so on.

Importance of "Data Divergence"

- Errors and inconsistencies in the linking data between the sample and the population create "data divergence":
 - -Time dynamics in the variables (e.g. changing Zip Codes when individuals move, Change in Martial Status, Income Levels, etc.),
 - -Missing and Incomplete data and
 - -Keystroke or other coding errors in either dataset,
- But even probabilistic record linkage methods, which can help address such challenges, are subject to uncertainty. The data intruder is never really certain that the correct persons have been re-identified.
- The recent Personal Genome Project re-identification attack using {Zip5, Gender and DoB} was able to achieve only a 27% re-identification rate (not 87%) due to these issues.



information loss which may limit the usefulness of the resulting Complete *health information*"(p. 8, Guidance) **Protection Protection Bad Decisions / Bad Science Ideal Situation** (Perfect **Information &** Trade-Off between Perfect Information **Protection**) isclosure Quality Unfortunately, and not achievable **Privacy Protection** due to Poor mathematical **Privacy** constraints **Protection** No Protection Information **Optimal Precision**, No Lack of Bias Information

Balancing Disclosure Risk/Statistical Accuracy

- Balancing disclosure risks and statistical accuracy is essential because some popular de-identification methods (e.g. k-anonymity) can unnecessarily, and often undetectably, degrade the accuracy of deidentified data for multivariate statistical analyses or data mining (distorting variance-covariance matrixes, masking heterogeneous sub-groups which have been collapsed in generalization protections)
 - This problem is well-understood by statisticians, but not as well recognized and integrated within public policy.
 - Poorly conducted de-identification can lead to "bad science" and "bad decisions".

Reference: C. Aggarwal <u>http://www.vldb2005.org/program/paper/fri/p901-aggarwal.pdf</u>



Statistical methods can help reveal the true signal; But... **Kernel Density Estimation** (2) 0 0 O O _0



and this problem becomes more severe with with higher multi-dimensional space...





and... K-anonymity Hides Heterogeneities



2 Percent Sample from Population

- Not Pop Unique (56%)
 Sample Unique, but not Pop Unique (40.6%)
- Pop Unique (3.5%)









State Specific Re-identification Risks: Population Uniqueness



Challenge: Subtraction Geography (i.e., Geographical Differencing)

- Challenge: Data recipients often request reporting on more than one geography (e.g., both State and 3 digit Zip code).
- Subtraction Geography creates disclosure risk problems when more than one geography is reported for the same area and the geographies overlap.
- Also called geographical differencing, this problem occurs when the multiple overlapping geographies are used to reveal smaller areas for re-identification searches.

Example: OHIO Core-based Statistical Areas



Tennessee - ZCTA5 Populations





Tennessee - County Populations





Tennessee - ZCTA5 X County Populations











Challenge: "Geoproxy" Attacks

- Challenge: Data intruders can use Geographic Information Systems (GIS) to determine the likely locations of patients from the locations of their healthcare providers
 - Retail Pharmacy Locations
 - Physician or Healthcare Provider Locations
 - Hospital Locations

Geoproxy attacks have become much easier to conduct using newly available tools (e.g., Web mapping & "Mash-up" technologies) on the internet. So, How Do We Move Beyond Anecdotes to a Rigorous, Scientific, Evidence-Based Risk Management Approach for Dealing with Re-identification Risks?

Quantitative Policy Analyses have been used for decades by many government agencies (EPA, Energy Dept.) to help address challenging policy decisions regarding difficult risk management questions where considerable uncertainty exists for important risk management questions.

Quantitative Policy Analyses for De-identification Policy:

De-identification policy is the subject of considerable controversy because it must balance important risks and benefits to individuals and societies and both sides of this question are subject to important uncertainties and competing values.

Essential to recognize that complex social, psychological, economic and political motivations can underlie whether reidentification attempts are made.

Three Main Data Intrusion Scenarios:

- Specific-Target (aka "Nosy Neighbor") Attacks (Have specific target individuals in mind: acquaintances or celebrities)
- Marketing Attacks (Want as many re-identifications as possible in order to market to these individuals, may tolerate a high proportion of incorrect reidentifications, but this can come at the risk of being caught re-identifying)
- Demonstration Attacks (Want to demonstrate reidentification is possible to discredit the practice or to harm the data holder; Doesn't matter who is reidentified so unverified re-identifications may also achieve intended goals)

Data Intrusion Scenarios:

Prob(Re-identification) = Prob(Re-ident|Attempt)*Prob(Attempt)

- Note that Prob(Attempt) & Prob(Reident|Attempt) are actually not likely to be independent - higher reidentification probabilities are likely to increase reidentification attempts.
- Some very useful frameworks exist for characterizing Data Intrusion Scenarios:
 - Elliot & Dale, 1999, Duncan & Elliot Chapter 2, 2011
- We can frame the Prob(Attempt) in terms of: Motivation, Resources, Data Access, Attack Methods, Quasi-identifier Properties and Sets, Data Divergence Issues, and Probability of Success, Consequences and Alternatives for Goal Achievement

Quantitative Policy Science

Conducting systematic quantitative costbenefit policy analyses using state-of-theart uncertainty and sensitivity analysis methods (e.g. with Latin-Hypergrid exploration of uncertain parameters) allows us to properly deal with the many important unknowns which could impact whether re-identification attempts under various data intrusion scenarios are likely be economically viable and realistic.



Uncertainty Analyses



Intrusion Scenarios



The 'Re-Identification' of Governor William Weld's Medical Information: A Critical Re-Examination of Health Data Identification Risks and Privacy Protections, Then and Now

http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2076397



Bill of Health

Examining the intersection of law and health care, biotech & bioethics A blog by the Petrie-Flom Center and friends



Online Symposium on the Law, Ethics & Science of Re-identification Demonstrations

- <u>http://blogs.law.harvard.edu/billofhealth/2013/05/29/public-policy-</u> <u>considerations-for-recent-re-identification-demonstration-attacks-on-</u> <u>genomic-data-sets-part-1-re-identification-symposium/</u>
- <u>https://blogs.law.harvard.edu/billofhealth/2013/10/01/press-and-</u> reporting-considerations-for-recent-re-identification-demonstrationattacks-part-2-re-identification-symposium/
- <u>http://blogs.law.harvard.edu/billofhealth/2013/10/02/ethical-</u> <u>concerns-conduct-and-public-policy-for-re-identification-and-de-</u> <u>identification-practice-part-3-re-identification-symposium/</u>

Reserve Slides for Questions

Measuring Disclosure Risks





Sample

Records

Records that are unique in the sample but which aren't unique in the population, would match with more than one record in the population, and only have a probability of being identified Only records that are unique in the sample and the population are at clear risk of being identified with exact linkage

> Population Uniques

Records that are not unique in the sample cannot be unique in the population and, thus, aren't at definitive risk of being identified

Sample

Uniques

Links

Records that are not in the sample also aren't at risk of being identified 59

Population

Records

Re-identification Failure and Success Conditions

HOSPITAL DATA SET (Found In Data Set)	VOTER DATA SET (Found in Data Set)	NON-VOTERS (in Population)
1 Not in Hospital Data	Male 1/1/1945 02138 Can't Re-identify (No Match)	
2 Male 1/2/1945 02138	Not in Voter Data	Male 1/2/1945 02138 Can't Re-identify (No Match)
3 Male 1/3/1945 02138 3 Male Male 1/3/1945 02138	Male 1/3/1945 02138 Can't Re-identify (>1 Match)	
4 Male 1/4/1945 02138 Can't Re-identify (>1 Match)	Male Male 1/4/1945 1/4/1945 02138 02138	
Male 1/5/1945 02138 9 Presumed Re-identification (Has Only 50% Chance of Being a Correct Match)	Male 1/5/1945 02138	Male 1/5/1945 02138 Directly Protected From Re-identification
6 Male 1/6/1945 02138 Correct Re-identification	Male 1/6/1945 02138	

lote: Figure illustrates only those imited cases where only one or two persons vith shared quasi-identifier" characteristics exist in either he healthcare lata set or in the voter registration ist.

Myth of the "Perfect Population Register"

Note that in Row 5 on previous slide:

- Every person not within the voter list is directly protected from re-identification.
- Furthermore, their absence from the population register also reduces the probability that others who share their quasi-identifier set would be correctly reidentified.

This is an extremely important limitation on re-identification when imperfect population registers are used.

Challenge: Geoproxy Attacks



Example: Patient location as revealed within data set, but further narrowed to probable "hotspots" by using healthcare provider location data





Directional (Standard Deviation Ellipse) distributions and "Hot Spot" analysis (Z-score color coding zip codes for Getis-Ord Gi* statistics)





»Free the Data, Manage the Risks

Barbara Rudolph PHD, MSSW UW-Madison, Center for Health Systems Research and Analysis NAHDO Consultant

History of NAHDO Efforts in Privacy of Data

- » Testimony at NCVHS
- White Papers and Contract
 Project Reports
- » Published articles
- Technical Tools—Inventory for
 Screening (Data) Release
- » Administrative Data Committee leading to PHDSC _____
- » Guidance Documents and Technical Assistance (NY)_
- » Privacy Workgroups

NAHDO Documents Covering Data Management and Release Policies

Rudolph, B. Davis, R. Administrative Data and Disease Surveillance: An Integration Toolkit. NAHDO-CDC Cooperative Agreement Project, CDC Assessment Initiative

Rudolph B., Shah, G., Love D. Small numbers, disclosure risk, security, and reliability issues in web-based data query systems. J Public Health Management Practice, 2006.

Person-level Data: An Inventory for Screening Release, NAHDO, 2008.

PHDSC, PRISM--A Privacy Toolkit for Public Health Professionals.; Glossary.

Guidance Document on Creating and Releasing Hospital and Facility Discharge Data Public Use Files, NAHDO, 2012.

Key Privacy Points From NAHDO Documents

>Define policies for release keeping in mind important societal needs for data and privacy for the individual patient. The utility of data must be preserved as well as the privacy of patient data or collection will cease.

>Establish management controls such as: data use agreements, limited access to files, limited access to data locations, limited data use beyond premises, encryption, etc.

>Promulgate administrative rules on: data release; exceptions to open record statutes; IRBs and/or privacy boards; direct and indirect or sensitive data elements and their release; limits for re-release of data elements by users; limits on attempt to re-identify patients; penalties for misuse and penalty enforcement

> Establish statistical processes, define specific tests, or establish rules for reducing risk of re-identification from release of sensitive data elements, such as detailed geographic identifiers, service or other dates, rare diagnoses, race and ethnicity.

Key Privacy Points from NAHDO Documents

>Establish limits for cell sizes in tables or reports; web-query systems, tables and maps and determine if you must aggregate years of data or geography to address small cell sizes

>Define a responsible person(s) within your organization for each of these tasks in job descriptions to assure internal compliance with the processes

>Provide annual staff training on the processes for maintaining privacy of the patient