Data Privacy Decision Spaces

Why Privacy is not Dead!

Bradley Malin, Ph.D.

Assoc. Prof. & Vice Chair of Biomedical Informatics, School of Medicine

Assoc. Prof. of Computer Science, School of Engineering

Affiliated Faculty, Center for Biomedical Ethics & Society

Vanderbilt University

22/9/2015

Acknowledgements

- Edoardo Airoldi, Ph.D. (Harvard U.)
- Kathleen Benitez
- Mustafa Canim, Ph.D. (IBM)
- Ellen Wright Clayton, M.D., J.D.
- Josh Denny, M.D.
- Xiaofeng Ding, Ph.D. (U. So. Australia)
- Khaled El Emam, Ph.D. (U. Ottawa)
- Aris Gkoulalas-Divanis, Ph.D. (IBM)
- Jonathan Haines, Ph.D. (Case Western)
- Raymond Healtherly, Ph.D. (Shyft Analytics)

eMERGE Teams

- Boston Children's Hospital
- Children's Hospital of Philadelphia
- Cincinnati Children's Hospital
- Geisinger Health System
- Group Health Research Institute / U. Washington
- Marshfield Clinic
- Mayo Clinic
- Mt. Sinai Medical Center / Columbia University
- Northwestern University
- Vanderbilt University

- Murat Kantarcioglu, Ph.D. (U. of Texas)
- Jiuyong Li, Ph.D. (U. So. Australia)
- Grigorios Loukides, Ph.D. (Cardiff U)
- Dan Masys, M.D. (U. Washington)
- Dan Roden, M.D.
- Laura Rodriguez, Ph.D. (NHGRI / NIH)
- Latanya Sweeney, Ph.D. (Harvard U.)
- Acar Tamersoy, M.S. (Georgia Tech)
- Weiyi Xia, M.S.
- Eugene Vorobeychik, Ph.D.
- Zhiyu Wan
- Funding
- NHGRI @ NIH
 - U01 HG006385 (eMERGE)
 - U01 HG006378 (VGER)
- NLM @ NIH
 - R01 LM009989
- Trustworthy Computing @ NSF
 - CCF-0424422 (TRUST)

U.S. <u>National Institutes of H</u>ealth Data Sharing "Policy"

• 2003 Final Data Sharing Policy:

- Receive $\$500k \rightarrow$ must have data sharing plan (or say why not possible)
- Recommended sharing data devoid of identifiers

• 2014 Genome Data Sharing Policy

• Studies involving > \$0

Identifiable?

A Legal View Privacy

EU Data Protection Directive:

"principles of protection shall not apply to data rendered anonymous in such a way that the data subject is no longer identifiable"

US Medical Regulation:

"information that does not identify an individual and ... no reasonable basis ... information can be used to identify an individual"



US Medical Privacy Rule

Identified Patient Data	 Waiver of consent: data is "on the shelf" Consent is impractical to obtain
Limited Data Set	 Removal of 16 designated attributes Recipient signs data use contract
De-identified	 See previous page

Recipes for Privacy

Field	Detail	
Names	Related to patient (not provider)	
Unique Numbers	Phone, Social Security Number,	
Internet	Email, URL, IP addresses,	
Biometrics	Finger, voice,	
Dates	Less specific than year Ages > 89	Limited Dataset
Geocodes	Town, County, Less specific than Zip-3 (assuming > 20,000 people in zone)	

The Concern

Ethnicity Visit date Diagnosis Procedure Medication

Total charge

Hospital

Discharge Data

ZIP Code Birthdate

Gender

Address

Name

Date registered

Party affiliation

Date last voted

Voter List

High Profile Re-identification



Sweeney. Journal of Law, Medicine, & Ethics. 1997

5-Digit ZIP

+ Birthdate

+ Gender

63-87% of USA estimated to be unique

Sweeney Tech Report 2000; Golle WPES 2006; Benitez & Malin JAMIA 2010

Set the World





The AOL >> Search Log Case (2006)

Pseudo	Name	Query	Date	Time
1		Books	1/2/05	16:52
2		Payscale	1/4/05	23:41
1		Porn	1/8/05	03:15

Goal: Support web information retrieval research

- 650 K customers, 20 M queries, 3 MONTH period
- Names replaced with persistent pseudonyms

Barbaro & Zeller. A face exposed for AOL searcher no. 4417749. <u>New York Times</u>. Aug 9, 2006.



Thelma Arnold

& Dudley



[Your Favorite Feature] Distinguishes You!!

- Demographics (Sweeney '97; Bacher '02; Golle '06; El Emam '08; Koot '10; Li '11)
- Diagnosis Codes (Loukides '10; Tamersoy '10, '12)
- Laboratory Tests (Cimino '12, Atreya '13)
- DNA (Malin '00, Lin '04; Malin '05; Homer '08; Gymrek '13, Ayday'14, Huttenhower '15)
- Health Survey Responses (Solomon '12)
- Location Visits (Malin '04; Golle '09; El Emam '11)
- Pedigree Structure (Malin '06, Ayday '13)
- Movie Reviews (Narayanan '08)
- Social Network Structure (Backstrom '07; Narayanan '09; Yang '12)
- Search Queries (Barbaro '06)
- Internet Browsing (Malin '05; Eckersley '10; Banse '11; Herrmann '12, Olejnik '12)
- Smart Utility Meter Usage (Buchmann et al '12)

HIPAA Expert Determination (abridged)

Certify via "generally accepted statistical and scientific principles & methods, that the **rigk i** very mall that the information could be used, alone or in combination with other **reasonably available info**rmation, by the anticipated recipient to identify the subject of the information."

A Brief History of Data Protection Models



k-Based Models

Age	Sex	Zip	Age	Sex	Zip	Age	Sex	Zip
30	Μ	15213	30	Μ	15213	30	Μ	15213
33	Μ	15217	33	*	1521*	33	*	1521*
33	F	15213	3*	*	15213	33	*	1521*
30	Μ	15213	30	Μ	15213	30	Μ	15213

Private Records

2-Abiguous

2-Anonymous

Thanks to Vitaly Shmatikov

Differential Privacy (informal)

Output is similar whether any is included in the database or

If there is already <u>some risk</u> of revealing a secret of C by combining auxiliary information and something learned from DB



C is no worse off because her record is included in the computation

Thanks to Vitaly Shmatikov

Achieving DP with Laplace Noise

Theorem

 $\overline{If A(x)} = f(x) + \mathsf{Lap}\left(\frac{\mathsf{GS}_f}{\varepsilon}\right) \text{ then } A \text{ is } \varepsilon \text{-indistinguishable.}$

Laplace distribution $Lap(\lambda)$ has density $h(y) \propto e^{-\frac{\|y\|_1}{\lambda}}$



Sliding property of $Lap\left(\frac{GS_f}{\varepsilon}\right)$: $\frac{h(y)}{h(y+\delta)} \le e^{\varepsilon \cdot \frac{\|\delta\|}{GS_f}}$ for all y, δ *Proof idea:* A(x): blue curve A(x'): red curve $\delta = f(x) - f(x') \le GS_f$





Communications of the ACM

DOI:10.1145/1743546.1743558

Arvind Narayanan and Vitaly Shmatikov

Privacy and Security Myths and Fallacies of "Personally Identifiable Information"

Developing effective privacy protection technologies is a critical challenge for security and privacy research as the amount and variety of data collecte individuals increase exponentially.

HE DIGITAL ECONOMY relies on the collection of personal data on an ever-increasing scale. Information about our searches, browsing history,

social relationships, medical history, and so forth is collected and shared with advertisers, researchers, and government agencies. This raises a number of interesting privacy issues. In today's data protection practices, both in the U.S. and internationally, "personally identifiable information" (PII)-or, as the U.S. Health Insurance Portability and Accountability Act (HIPAA) refers to it, "individually identifiable" information-has become the lapis philosophorum of privacy. Just as medieval alchemists were convinced a (mythical) philosopher's stone can transmute lead into gold, today's privacy practitioners believe that records containing sensitive individual data can be "de-identified" by removing or modifying PII.

What is PII?

For a concept that is so pervasive in enable identity theft. Therefore, they both legal and technological discourse focus solely on the types of data that

Any information that distinguishes one person from another can be used for re-identifying data.

on data privacy, PII is surprisingly difficult to define. One legal context is provided by breach-notification laws. California Senate Bill 1386 is a representative example: its definition of personal information includes Social Security numbers, driver's license numbers, financial accounts, but not, for example, email addresses or telephone numbers. These laws were enacted in response to security breaches involving customer data that could enable identity theft. Therefore, they focus solely on the types of data that

24 COMMUNICATIONS OF THE ACM | JUNE 2010 | VOL. 53 | NO. 6

NIH Public A Author Manuscript

Am J Bioeth. Author manuscript; available in PMC 2011 February 2

Published in final edited form as:

Am J Bioeth. 2010 September; 10(9): 3-11. doi:10.1080/15265161.2010.494215.

Is Deidentification Sufficient to Protect Health Privacy in

n laws. a repition of s Social

Mark A. Rothstein

University of Louisville School of Medicine

The revolution in health information technology has enabled the compilation and use of large data sets of health records for genomic and other research. Extensive collections of health records, especially those linked with biological specimens, are also extremely **20** valuable for outcomes research, quality assurance, public health surveillance, and other

American Journal of Bioethics







De-identification May Be Safe

- Reviewed all <u>actual</u> re-identification attempts
- Attacks on health data
 - 14 published re-identification attacks on any type of data
 - 11 of 14 were conducted by researchers as demo attacks
 - Only 2 of 14 attacks followed any standard
 - Only case with health data subject to "Safe Harbor" had a success rate of 0.00013

A Case Study on Demographics

• Details at

http://www.ehcca.com/presentations/HIPAAWest4/lafky_2.pdf

- Challenge issued by U.S. Dept. Health & Human Services
- Gave 15,000 Safe Harbor records to investigative team at University of Chicago
- Team purchased public records from commercial broker
- Correctly identified <u>2 people</u>

Considerable Knowledge

• Identifiability is proportional to

Uniqueness (must distinguishable)xReplicability (must be reproducible)xAvailability (must be accessible)

• A drug dose may be unique, but may not be accessible to the public in any known resource

Only 1 in 650,000 people identified



Thelma Arnold & Dudley

What's Going On?

A Very Simplified View on Risk



- Uniqueness
- Replicability
- Availability

Data



An Augmented View of Data Privacy



Central Dogma of Re-identification



Malin, Benitez, Loukides, Clayton. Human Genetics. 2011.

...We've Been Looking at Worst Case Scenarios...

- How would you <u>use</u> demographics?
- Could link to registries
 - Birth Marriage
 - Death Professional (Physicians, Lawyers)
- What's in vogue?

Back to voter registration databases

Going to the Source



The Availability of Demographics Varies...

	IL	MN	TN	WA	WI
WHO	Registered Political Committees (ANYONE – In Person)	MN Voters	Anyone	Anyone	Anyone
Format	Disk	Disk	Disk	Disk	Disk
Cost	\$500	\$46; "use ONLY for elections, political activities, or law enforcement"	\$2500	\$30	\$12,500
Name	•	•			\bullet
Address	•	•			
Date of Birth	•	0			
Sex	•				
Race					
Phone Number	•	•			

Who is Like You?

Time for Estimation (Golle 2006)

- Disclose sample with {*dob, gender, zip*}, but don't know the population's values
- Don't always have exact knowledge of what a data recipient has access to

	Birthdate					
	1/1/80	•••	12/31/80			
zip1						
zip2						
•••						
zip m						

Time for Estimation (Golle 2006)

- May know population counts, such as
 - Census aggregates for {*year of birth, gender, county*}

	Birth Year			Birthdate		
	1980			1/1/80		12/31/80
zip1	12000		zip1			
 zip2	50000	סוק	zip2			
•••			•••			
zip m	10000		zip m			
Time for Estimation (Golle 2006)

- May know population counts, such as
 - U.S. Census aggregates for {*year of birth, gender, county*}
- Use a disaggregation model to estimate

Birth Year					Birthdate			SUM	
	1980		_			1/1/80		12/31/80	3011
	zip1	12000			zip1	random		random	12000
	zip2	50000		סוק	zip2	random		random	50000
				ZIP	•••				
	zip m	10000			zip m	random		random	10000

It's an Occupancy Problem (Golle 2006)



- *n* people in aggregated bin
- *b* disaggregated bins
- the expected # of bins with exactly *i* people
- Total number of people in a group of size less than *k*

$$f_i(n) = \binom{n}{i} b^{1-n} b^{n-i}$$

$$r_k(n) = \sum_{i=1}^{k-1} f_i(n)$$

All U.S. States

Safe Harbor





Benitez & Malin, JAMIA. 2010.

Risk...

- ... means something different to everyone
- Can be modeled in various ways (Dankar & El Emam, 2010)
 - Prosecutor $\leftarrow 1/\min_i(f_i)$
 - Journalist $\leftarrow 1/\min_i(F_i)$
 - Marketer $\leftarrow n^{-1} \sum_{i} \frac{f_i}{F_i}$

 f_i = size of group in sample F_i = size of group in population n = sample size



The Availability of Demographics Varies...

	IL	MN	TN	WA	WI
WHO	Registered Political Committees (ANYONE – In Person)	MN Voters	Anyone	Anyone	Anyone
Format	Disk	Disk	Disk	Disk	Disk
Cost	\$500	\$46; "use ONLY for elections, political activities, or law enforcement"	\$2500	\$30	\$12,500
Name	•	•			\bullet
Address	•	•			
Date of Birth	•	0			
Sex	•				
Race					
Phone Number	•	•			

Identifiability Changes!

Limited Data Set

Limited Data Set $\leftarrow \rightarrow$ Voter Reg.



What About Cost?

(Consider Marketer Risk)

IIS State	Limi	ted Dataset	Safe Harbor		
U.S. State	At Risk	Cost per Re-id	At Risk	Cost per Re-id	
Virginia	3159764	\$0	221	\$0	
South Carolina	2231973	\$0	1386	\$0	

We Need Policy Alternatives

How Can we Find Policies?

 Model acceptable data abstractions as a lattice and "search" for low risk



Simple Discovery Decision



Benitez, Loukides, and Malin. ACM IHI. 2010. Malin, Benitez, and Masys. JAMIA. 2011. Xia, et. Al. ACM CODASPY. 2013

Searching the Lattice

- Risk is monotonic on the graph
- Search space is huge and there are multiple "optimal" solutions
- We can search for "good" solutions using a ILP



• Faster - bisecting strategy

Vandy ECG Case Study

Who	State	State Population Size (2000 Census)	Cohort Size	Patients >89 years old
Vanderbilt	TN	5,689,283	2,983	12

Doligy		Dick			
Ροπογ	Gender	Race	Age	NISK	
Safe Harbor	Ø	Ø	[90 - 120]	0.909	
Alternative 1	[M or F]	Ø	Ø	0.476	
Alternative 2	Ø	[Asian or Other]	Ø	0.857	
Alternative 3	Ø	Ø	[52 - 53]	0.875	

Evaluation of De-id Model

 Cohorts from the Electronic Medical Records and Genomics Consortia (http://www.gwas.net)

Pheno.	Cohort	Who	State	State Population Size (2000 Census)	Clinical Finding of Interest	Cohort Size	Patients >89 years old
	G_{Dem}	GHC	WA	5,894,121	Dementia	3,616	1,483
	R _{Cat}	Marshfield	WI	5,363,675	Cataracts	2,646	269
Primary	Y _{PAD}	Mayo	MN	4,919,479	Peripheral Arterial Disease	3,412	29
	N _{T2D}	Northwestern	IL	1,2519,293	Type-II Diabetes	3,383	6
	V _{ORS}	Vanderbilt	TN	5,689,283	QRS Duration	2,983	12
Quality	N _{ORS}	Northwestern	IL	1,2519,293	QRS Duration	149	0
Control	V_{T2D}	Vanderbilt	TN	5,689,283	Type-II Diabetes	2,015	18

Analysis for eMERGE

Risk Model: Uniques

Are the number of uniques expected to be greater than Safe Harbor?

Disclosure		Acceptable?						
Policy	G _{DEM}	R_{CAT}	Y _{PAD}	N _{T2D}	V _{ORS}	N _{ORS}	V _{T2D}	
Generalized Ethnicity (Black, White, Other)					~	~		
Age at 5 Year Bins								
Generalized Ethnicity AND Age at 5 year bins								
Age at 10 Year Bins								

Red = more risk than Safe Harbor

Green = risk no worse than Safe Harbor

Malin, Benitez, & Masys. JAMIA. 2011.

Dual Optimization Extension (Risk-Utility Frontier)



Xia, Heatherly, Ding, Li, and Malin; 2013; 2015

Enhancement: Risk-Utility Frontiers

(Xia, et al. ACM CODASPY. 2013; JAMIA. 2015)



Example: State of Hawaii

- Simulation of 30,000 records from the Adult Census Database
 - {Age, Race, Gender}
- Appended 5-digit ZIP proportional to public use statistics
- Space of 2⁷⁰⁰ policies
- Frontier of ~400 policies discover by evaluating ~20,000



We are Driven By Incentives

(under rational assumptions)

Sharing Strategy 1 Utility 1 Risk ???

Strategies:

- Generalize Demographics
- Perturb Statistics
- Apply Data Use Agreement
- •••
- Charge for Access

Publisher

Attack Strategy A

Utility A

Risk A







Attack Strategy A Utility A Risk A



Publisher

56



57

Sharing Strategy 1

Utility 1

Risk B

Sharing Strategy 2 Utility 2 Risk ???

Attack Strategy A

Utility A

Risk A



Risk C

Recipient

Publisher



Sharing Strategy 1 Utility 1 Risk B



Publisher

Sharing Strategy 1 Utility 1 Risk B



Choose Strategy that maximizes overall benefit

- Optimal Utility / Risk Tradeoff

Payoffs!

Gain	No Attack	Attack
Publisher	Vg	
Attacker	0	

- g : Generalization level
- $v_{\rm g}$: Value of record at g

Payoffs!

Gain	No Attack	Attack
Publisher	Vg	v_g - L π_g
Attacker	0	

- g : Generalization level
- v_g : Value of record at g
- π_{g} : Probability of successful attack at g
- L: Loss to publisher for successful attack

Payoffs!

Gain	No Attack	Attack
Publisher	Vg	v_g - L π_g
Attacker	0	$L\pi_g - c$

- g : Generalization level
- v_g : Value of record at g
- π_{g} : Probability of successful attack at g
- L: Loss to publisher for successful attack
- c : Cost to run attack

Wan et al, PLoS One. 2015

Game Variations

- Safe Harbor (SH) Game
 - Defender shares data according to federal policy
- Basic Game
 - Defender shares data to maximize overall payoff
- SH-Friendly
 - Defender constrains strategy space to disclose no greater detail than SH
- No Attack
 - Defender constrains strategy space to disclose no greater detail than SH

Solving the Game?

- The sublattice search will work... but it's not optimal
- Alternatives
 - Backward Induction Search
 - For each generalization level
 - choose the one that maximizes publisher's utility
 - Exhaustive search over combinatorial space of data representation
 - ILP or something else.

Intelius - Find People with ×	
← → C 🗋 www.intelius.com	දූ 🔁 🕄 🔳
🏥 Apps 🕥 Barclays 🍥 ScholarOne Manusc 👿 HRPP/IRB Vanderb 🥰 V	BA Express : Excel 😕 Home - PubMed » 🎦 Other bookmarks
	F Like {28k G+1 {17k Help Sign ★ Bookmark this 5
People Search Background Check Criminal Records Reverse Look	Intelius Premier Identity Protection Employee Screening
People Search Email Lookup Social Network Search Property Records 2	24-Hour People Search Pass
People Search - Updated	Daily, Accurate and Fast!
People Search	
First Name M.I. Last Name required	City and/or State
	Search
Reverse Phone Lookup	
Phone Number	More ways to get info you need:
Search	Perform a Background Check Run a Background Check by SSN
	Perform an Address Lookup
	Do a Reverse Phone Lookup
What is People Search?	What is Reverse Phone Lookun?
It's a confidential way to find people so you	What is reverse Phone Lookup?
can reconnect or just get more info on a	phone number belongs to. Reverse phone
person. People Search reports can include phone numbers, address history, age &	search works for landline, unlisted & non- published numbers, and cell phone lookups.

person you're curious about - search

date of birth, relatives, and more. Find a



4



umbers, and e lookup Reports can include phone type, owner name, address & more. Curious? Do a phone number lookup!



⊩

Entity	Fine	#Records	Fine/record	Date
New York and Presbyterian Hospital	\$4,800,000	6,800	\$705.9	May 7, 2014
QCA Health Plan, Inc.	\$250,000	148	\$1689.2	Apr 22, 2014
Skagit County, Washington	\$215,000	118,000	\$1.8	Mar 7, 2014
Adult and Pediatric Dermatology	\$150,000	2,200	\$68.2	Dec 26, 2013
Affinity Health Plan, Inc.	\$1,215,780	344,579	\$3.5	Aug 14, 2013
WellPoint Inc.	\$1,700,000	612,402	\$2.8	Jul 11, 2013
Idaho State University	\$400,000	17,500	\$22.9	May 21, 2013
The Hospice of North Idaho	\$50,000	441	\$113.4	Jan 2, 2013

individuals. These breaches are now posted in a new, more accessible format that allows users to search and sort the posted breaches. Additionally, this new format includes brief summaries of the breach cases that OCR has investigated and closed, as well as the names of private practice providers who have reported breaches of unsecured protected health information to the Secretary. The following breaches have been reported to the Secretary:

Show Advanced Options

Breach Report Results 📎 🎽 📥 📥							
	Name of Covered Entity ≎	State \$	Covered Entity Type ≎	Individuals Affected ≎	Breach Submission Date ≎	Type of Breach	Location of Breached Information
0	Brooke Army Medical Center	ТХ	Healthcare Provider	1000	10/21/2009	Theft	Paper/Films
0	Mid America Kidney Stone Association, LLC	MO	Healthcare Provider	1000	10/28/2009	Theft	Network Server
0	Alaska Department of Health and Social Services	AK	Healthcare Provider	501	10/30/2009	Theft	Other, Other Portable Electronic Device
0	Health Services for Children with Special Needs, Inc.	DC	Health Plan	3800	11/17/2009	Loss	Laptop
0	L. Douglas Carlson, M.D.	CA	Healthcare Provider	5257	11/20/2009	Theft	Desktop Computer
0	David I. Cohen, MD	CA	Healthcare Provider	857	11/20/2009	Theft	Desktop Computer
0	Michele Del Vicario, MD	CA	Healthcare Provider	6145	11/20/2009	Theft	Desktop Computer
0	Joseph F. Lopez, MD	CA	Healthcare Provider	952	11/20/2009	Theft	Desktop Computer
0	Mark D. Lurie, MD	CA	Healthcare Provider	5166	11/20/2009	Theft	Desktop Computer

- \$1200: Benefit per record
- \$300: Cost per violation
- Average Payoff Per Record

- \$4: Access cost per record
- ~30,000 Census records



- \$1200: Benefit per record
- \$300: Cost per violation
- Average Payoff Per Record

- \$4: Access cost per record
- ~30,000 Census records



- \$1200: Benefit per record
- \$300: Cost per violation
- Average Payoff Per Record

- \$4: Access cost per record
- ~30,000 Census records



- \$1200: Benefit per record
- \$300: Cost per violation
- Average Payoff Per Record

- \$4: Access cost per record
- ~30,000 Census records


Case Study



Wan et al, PLoS One. 2015

Sensitivity Analysis

Publisher Payout

Probability of Attack



V = Value of Record

L = Loss due to re-identification

We Must Account for Process

Adversaries Make Sequential Decisions!

(Xia et al, CIKM 2015)



link the record to the external dataset



Intelius \$3.95/record USSearch \$1.45/record NY voter \$0 registration broker ? t2.micro t2.medium m3.large

\$0.013/hour \$0.052/hour \$0.133/hour





exploit an individual that matches the record

phone \$0. marketing/con firmation email marketing Publish the information Penalty \$?

<u>\$0.013/mi</u>

Outcome of Adversary's Actions Are Stochastic

 Before accessing the external dataset, the adversary may not know how many individuals can be linked to the record (unknown equivalence group size)



Adversary Makes a <u>Series of Decisions</u> to Complete an Attack

- The state of the attack can be represented by a set of variables
- Given state and action, attacker is granted a reward or pays a penalty
- Next state of attack depends on the current state & decision
 - Before accessing the external dataset, the attacker is uncertain about the content of the dataset
 - Before exploiting an individual, the attacker is not sure if the outcome will be success or fail

Factored Markov Decision Process (FMDP) to Represent Adversarial Behavior

- Assumption: Adversary always wants to maximize expected payoff
- MDP: models decision making in situations where outcomes are partly random and partly under the control of decision maker
- An optimal policy policy(x): action that maximizes the expected payoff across states
- Methods
 - Linear Programming, Value iteration, Policy iteration
- Challenges
 - State explosion of the FMDP
 - Dependency explosion of the Dynamic Bayesian Network

Experimental Setting

- Linking Variables: [Age, ZIP-5, Gender, Race]
- De-identified data: Adult Census + North Carolina ZIP-5
 - ~30K people
- Identified dataset: NC Voter's registration
 - ~6M people
- US census 2010 data from NC used to estimate the equivalence group size distribution

Costs & Gains

- \$100 Cost to access identified dataset
- \$10 Cost to conduct exploit
- \$8000 Gain for successful adversary
- \$10000 Penalty for adversary when exploit detected

(Sensitivity analysis in CIKM paper)

	Sequential	Baseline*
Exploit (per de-id record)		1 individual
Decision Making		

	Sequential	Baseline*
Exploit (per de-id record)	\geq 1 individual	1 individual
Decision Making		

	Sequential	Baseline*
Exploit (per de-id record)	\geq 1 individual	1 individual
Decision Making		Single attack decision

$$payoff = G \times \left(\frac{prior_{r,D_e}}{G_{r,D_e}}\right) - p_{det} \times C_p - C_d - C_l - C_e$$

*Wan, et al. PLoS One 2015.

	Sequential	Baseline*
Exploit (per de-id record)	\geq 1 individual	1 individual
Decision Making	Sequence of decisions	Single attack decision

$$payoff = G \times \left(\frac{prior_{r,D_e}}{G_{r,D_e}}\right) - p_{det} \times C_p - C_d - C_l - C_e$$

*Wan, et al. PLoS One 2015.

Adversary May Exploit > 1 Person!



Baseline Can Underestimate Risk!



Uncertainty in Equivalence Group Size can Lead to Higher Risk!



So Where Are We Now?

- Privacy is NOT dead
- It is a much more complex landscape than has been suggested
- The space of options for data and social manipulation can be modeled... but it's huge!
- Challenges
 - Reliable cost estimates
 - Must Beware of "Baiting"!
 - Non-monotonic privacy and utility functions
 - State explosion in process models
 - Multiple publisher / adversary scenarios

Questions?

b.malin@vanderbilt.edu

Health Information Privacy Laboratory http://www.hiplab.org/