# TOWARDS A LICENSING MODEL AND ECOSYSTEM FOR DATA SHARING

SAM MADDEN (MIT)

W/CARSTEN BINNIG (BROWN), JANE GREENBERG (DREXEL), TIM KRASKA (BROWN)

NSF Big Data Hub Spoke Grant

## WHY IS DATA SHARING IMPORTANT!



## **Different Reasons**

- Combining Data
- Sharing with Experts
- •

# **Promise:** Better Insights into "Big Data"

# COMBINING DATA: COLLABORATIVE CANCER CLOUD



# How open data on agriculture & nutrition can solve world hunger

**07 SEPTEMBER 2016** 





#### DESIGN / TRANSPORTATION / ENVIRONMENT / EQUITY / LIFE Q

**Guai** susta busin Value busin



-range

tect

Noi Jaitang, interviewed as part of the World Resources Institute report, waters his garden in Thailand // Laura Villadiego

#### How to Solve the Environmental Information Divide TERESA MATHEW SEP 5, 2017

SundayReview | OPINION

#### Give Up Your Data to Cure Disease

By DAVID B. AGUS FEB. 6, 2016

# The New York Times

February 2016



MARK WARREN NATIONAL FRONTIERS SCIENCE 10.19.16 6:55 AM

# THE CURE FOR CANCER IS DATA— MOUNTAINS OF DATA

# ※ WIRED

October 2016



#### December 2013

# Yes, Big Data Can Solve Real World Problems



Greg Satell, CONTRIBUTOR

Opinions expressed by Forbes Contributors are their own.





Forbes, Working with IBM, the Memphis Police Dept. managed to reduce crime by 30% using big data analytics

# SHARING DATA WITH RESEARCHERS: FINANCIAL DATA

**Challenge:** Consumer credit risk analysis and forecasting

**Approach:** Machine learning



Machine-learning detects potential defaults more accurately than FICO scores!

## BUT DATA SHARING IS HARD!



## OPEN DATA



# DataSNE

#### DFC DataNet FEDERATION CONSORTIUM



# CLOSED DATA



#### Intel-Collaborative Cancer Cloud

(CCC) (Dana-Farber, OHSU, Ontario Institute for Cancer Research (OICR))



FICO

### Collaborative Genomics Cloud

(CGC )colocalizing massive genomics datasets)

**FICO** score (Fair Isaac Corporation)

## WHY NOT OPEN DATA?



## **BARRIERS TO DATA SHARING**

Must go beyond "creative commons"

Incentives – why would someone go to all the effort to share their valuable data?

**Concerns over sensitive information (e.g., PII)** 

**Regulations governing use of data in different domains** 

Not just "throwing it over the wall"!

- Do not want to loose control over data
- Can I get my data back?
- Has to be updated, requires training, redacted etc.

## SHARING DATA TODAY

## No data sharing without a legal agreement



Involve lawyers to create individual agreements  $\rightarrow$  often prevents sharing!

## DATA SHARING SPOKE: COMPONENTS

- 1. Data-sharing Licensing Framework / Generator
- 2. Data-Sharing Platform (Enforce Licenses)

**Principle: Solve the 80% case!** 

# **GOAL: LICENSING FRAMEWORK**

# Standard terms that researchers, lawyers, and compliance teams conform with



Tracking of access

Usage rights (e.g., publication, copying)

Duration of use

Warrantees of correctness/completeness/availability

## LICENSES: FIRST RESULTS

### Data-Sharing Workshop 2016 (Metadata Research Center @ Drexel):

- Approx. 60 participants form industry + academia
- Hear from the trenches
- What works? What doesn't? What are the biggest barriers? (What are the non-barriers?)
- Brainstorm solutions: would standardized licenses, use-cases/best practices help? Would better technologies help?
- Forge a path forward, together

**Agenda and Report:** <u>http://cci.drexel.edu/mrc/news/</u> 2016-11-bigdatahubworkshop/

## LICENSES: FIRST RESULTS

### **Collected sharing agreements from academic institutions**

### **Compile list of standard terms for**

- General (Time period, Use of data, ...)
- Privacy & Protection (PII, Security, Training)
- Access (Who?, How?)
- Responsibility (Indemnity clause, Ownership, Rights)
- Compliance (Background checks, Right to audit, ...)
- Data Handling (Allowed Methods of Data Transfer, ...)
   Initial analysis suggests there is much commonality
   Send us your (anonymized) licenses: <a href="mailto:smg383@drexel.edu">smg383@drexel.edu</a>

## **CONTENT ANALYSIS**

## 1. Data collection

26 data sharing agreements, industry, academia, government

## 2. Content analysis

- Confirm data sharing in closed environment
- Focused, language parsed for higher-level general categories; mid, lower-level to → specifications to data handling

## 3. Concept clustering

• Classes, sub-classes, attributes organized on a spreadsheet in a classified, hierarchical arrangement.

## 4. Metadata labeling

Language of the categories and attributes was refined

#### **Licenses: First Results**

Categori

gh-level

(Sam Grabus: <u>smg383@drexel.edu)</u>

General: attributes relating to the project and the agreement itself

Privacy & Protection: the protection of sensitive information and security

Access: who and how contact may be made with the data e.g., Who has access, Method of access (approved hardware or software)

e.g., Description of the data,

**Definition of terms** 

e.g., Individual identifiers removed

prior to transfer,

Encryption

Responsibility: legal, financial, ownership, and rights management pertaining to the data

Compliance: ensuring fulfilment of agreement terms

Data Handling: specifics of permissible interactions with the data Establishment of data ownership

e.g., Indemnity clause,

e.g., Third party compliance with contract, Background checks for personnel

> e.g., Publication of data, Conditions for Termination

Privacy & Protection										
	Sensitive Information									
Regulations	Preparing data	Access								
<ul> <li>Regulation used to define sensitive data (e.g., HIPAA, FERPA, etc.)</li> <li>Compliance with federal/state/international data protection laws and regulations</li> </ul>	<ul> <li>Identification of confidential/special categories of information (e.g., pii, proprietary)</li> <li>Individual identifiers removed/anonymized prior to transfer</li> </ul>	<ul> <li>Who has access to pii/confidential data</li> <li>Who has access to proprietary information</li> </ul>								
Privacy	Avoiding re-identification	Exceptions								
<ul> <li>Anonymization of data</li> <li>Confidentiality and safeguarding of PII/sensitive data</li> <li>Removal/nondisclosure of company/personnel identification in materials and publications</li> <li>No contact with data subjects</li> </ul>	<ul> <li>No direct/indirect re- identification</li> <li>Statistical cell size (how many people, in aggregated form, can be released in groups)</li> <li>Merging data with other sets (e.g., allowed with aggregated data—not in any way that will re-identify</li> </ul>	<ul> <li>Exceptions to confidentiality</li> <li>Conditions of proprietary information disclosure</li> <li>Conditions of pii disclosure (who, what, and for what purpose?)</li> <li>Limitations on obligations if data becomes public</li> <li>Limitations on obligations if data is already known prior to agreement</li> <li>Limitations on obligations if data given by 3<sup>rd</sup> party without restriction</li> </ul>								
	Security									
<ul> <li>Sharing non-confidential data</li> <li>Password protection/authentica</li> <li>Encryption</li> </ul>	<ul> <li>Security training</li> <li>tion of files</li> <li>Establishing introduction</li> <li>confidential data</li> </ul>	ng for involved personnel frastructure to safeguard ata								

Data Handling										
Us	se la	Physical								
<ul> <li>Each data field/elements to be accessed</li> <li>Use of data: only for project-specific/research, or analytical use</li> <li>Documenting all projects using the data</li> </ul>	<ul> <li>Modification of data</li> <li>Compliance with data updates (changes, removal, corrections)</li> <li>Sharing data</li> </ul>	<ul> <li>Copy/reproduction of data</li> <li>Storage of data</li> <li>Transfer of data (e.g., allowed methods)</li> </ul>								
Res	Personal Gain									
<ul> <li>Presentation of data</li> <li>Publication of data (e.g., prior approval needed or right to publically disclose publication)</li> </ul>	<ul> <li>Results/reports and associated documents (e.g., must be provided copies)</li> <li>Right to remove/delete confidential data from proposed publications</li> </ul>	<ul> <li>Sale of/profit from data (e.g., noncommercial use only)</li> <li>Licensing of data</li> <li>No reverse engineering</li> </ul>								
	Termination									
<ul> <li>Conditions for termination</li> <li>Destruction or return of data after</li> <li>3<sup>rd</sup> party destruction or return of</li> <li>Confirmation of data destruction</li> </ul>	or used for period of time after nd obligations remain in effect ion									

6,~40,90+

### **Privacy & Protection**

#### Security

- Sharing non-confidential data  $\rightarrow$  Sharing non-confidential data
- Password protection/authentication of files  $\rightarrow$  Password protection
- Encryption  $\rightarrow$  Encryption
- Security training for involved personnel → Personnel Security Training
- Establishing infrastructure to safeguard confidential data → Establishing Infrastructure

## Data Handling

#### 🛛 Use

- Each data field/elements to be accessed  $\rightarrow$  Fields Accessed
- Use of data: only for project-specific/research, or analytical use  $\rightarrow$  Research Use Only
- Documenting all projects using the data  $\rightarrow$  Projects involved
- Modification of data → Modification
- Compliance with data updates (e.g., changes, removal, corrections) → Data Updates
- Sharing data  $\rightarrow$  Data Sharing

## NLTK – PARSING TERMS

Set maximum keywords length: 5 List top 1/5 of all the keywords **Result:** Keyword: research studies involving human subjects, score: 20.4583333333 Keyword: district assigned student identification numbers, score: 18.8387650086 Keyword: includes personally identifiable student information, score: 17.6168132942 Keyword: district initiated data research projects, score: 14.8577044025 Keyword: support effective instructional practices, score: 13.0 Keyword: personally identifiable information shared, score: 11.3440860215 Keyword: disclose personally identifiable information, score: 11.1440860215 Keyword: policy initiatives focused, score: 9.0 Keyword: informing education policies, score: 9.0

#### Sample 30 agreements

-5	-4	-3	-2	-1	0 1 2		3	4	5	
			educational	right	privacy	act	health	insurance	portability	accountability
applicable	federal	law	regulation	protecting	privacy	citizen	including	family		
	license	agreement	authorized	protect	privacy	individual	subject	nd	study	
				applicable	privacy	law				
consistent	federal	family	educational	right	privacy	act	department	designates	education	alliance
subject	federal	family	educational	right	privacy	act	authorized			
education	record	covered	family	educational	privacy	act	amended			
recipient	agent	subcontractor	violation	agreement	privacy	rule	security	rule	implementing	regulation
comply	applicable	state	local	security	privacy	law	extent	protective	individual	privacy
		data	security	protection	privacy					
information	identified	family	educational	right	privacy	act				
		de	identified	applicable	privacy	law				
				applicable	privacy	law	permit	data	provider	provide
				federal	privacy	act	requirement	apply	agreement	entered
shared	state	subjected	applicable	requirement	privacy	confidentiality				
resolved	permit	covered	entity	comply	privacy	rule				
time	covered	entity	comply	requirement	privacy	rule	hipaa			
		reference	agreement	section	privacy	rule	mean	section	amended	renumbered
					privacy	rule	extent	information	created	received
					privacy	rule	standard	privacy	individually	identifiable
					privacy	rule	include	person	qualifies	personal
tern	defined	agreement	meaning	term	privacy	rule				
set	accordance	term	agreement	hipaa	privacy	security	rule			
hipaa	regulation	promulgated	thereunder	governing	privacy	security	health	information		

## MOST FREQUENT TERMS

privacy	protection	set						
applicable	privacy	law						
privacy	rule	standard	privacy	individually	identifiable			
definition	set	privacy	rule					
data	security	protection	privacy	Frequency from the				
				most to	the least:			

## GOAL: HOSTED DATA-SHARING PLATFORM



### <u>Tech</u>

- Access control & rights management
- Expiration
- Logging & auditing
- Provenance/Finger printing
- **De-identification**
- "Noising"
- Aggregation

#### Agreement Clauses

- Controlled access (who & where)
- Tracking of access
- Usage rights (e.g., publication, copying)
- Duration of use
- Warrantees of correctness/completeness/ava ilability
- Other requirements and regulations

# IS THIS POSSIBLE: TECHNOLOGY 🖂 SHARING AGREEMENTS

### <u>Tech</u>

# Access control & rights management

Expiration

Logging & auditing

Provenance/Finger printing

**De-identification** 

"Noising"

Aggregation

#### Agreement Clauses

# Controlled access (who & where)

Tracking of access

Usage rights (e.g., publication, copying)

#### **Duration of use**

Warrantees of correctness/completeness/avail ability

### <u>Tech</u>

Access control & rights management

Expiration

# Logging & auditing

Provenance/Finger printing

**De-identification** 

"Noising"

Aggregation

#### Agreement Clauses

Controlled access (who & where)

#### **Tracking of access**

Usage rights (e.g., publication, **copying**)

Duration of use

Warrantees of correctness/completeness/avail ability

### <u>Tech</u>

- Access control & rights management
- Expiration
- Logging & auditing

# **Provenance/Finger**<br/>printing

- **De-identification**
- "Noising"

Aggregation

#### Agreement Clauses

Controlled access (who & where)

Tracking of access

**Usage rights** (e.g., **publication**, **copying**)

#### **Duration of use**

Warrantees of correctness/completeness/avail ability

### <u>Tech</u>

- Access control & rights management
- Expiration
- Logging & auditing
- Provenance/Finger printing
- **De-identification**
- "Noising"
- Aggregation

<u>Agreement Clauses</u>

Controlled access (who & where)

Tracking of access

Usage rights (e.g., publication, copying)

Duration of use

Warrantees of correctness/completeness/availabi lity



DataHub - Chromium						tu ⊠ ⊲× 9:29 PM ⊀‡
* DataHub ×						0
← → C (i) localhost/licenses/create?						☆ :
Dat	taHub		脅 Home	e 🔇 Public Data 🚬	SQL Console	johndoe <del>-</del>
Cre	eate New License					
Ger	neral					
Owner	r.	organization				
Licens	se Name:	license name				
Priv	acy and Protection					
Regu	lations					
	PAA					
E FE	RPA					
Priva	асу					
🗆 PII	Anonymized or Removed					
I PII	Anonymized					
	ntions					
Beid	entification					
	e K-Anonymity					
L oo	ize					
	BUCKET SIZE TOF K					
Crea	ate					

About Documentation GitHub Repo API

DataHub - Chromium				ti, ⊠ ≪ 9:30 PM ∜ ● ↔
DataHul	)	# Hor	me 🛛 Public Data 🚬 SQL Console 🛔	johndoe ▼
Create	New License			
General				
Owner:	health data res	earch org		
License Name:	new ferpa rem	oved		
Privacy a	nd Protection			
Regulations				
III HIPAA				
I FERPA				
Privacy				
PII Anonymi	zed or Removed			
PII Anonymi.	zea 1			
Exceptions	a			
Reidentificat	tion			
🔲 Use K-Anon	ymity			
K-size	Bucket Size for K			
Create				
	About	Documentation GitHub Repo API		

DataHub - Chromium								tų ⊠ ∉× 9:31 PM 🔱
☆ DataHub ×								Θ
← → C () localhost/licenses/johndoe/test_repo								☆ :
۵	DataHub			🖀 Home	Public Data	>_ SQL Console	å johndoe 🗸	,
jo	ohndoe / test_repo /							
Ν	Manage Repository	y Licenses						
	Create New License							
	Add Existing License							
1	License Name	Applied To Tables			Manage Tab	le Application		
f	ferpa removed		License Applied To All Tables		Ма	anage		
1	hipaa removed		License Applied To All Tables		Ма	anage		
		Abou						

DataHub - Cl	romium										t∎ 🖾 ∥× 9:32 PM 🔱
👫 DataHi	1p × di										θ
← → C	(i) localhost/licenses/johndoe/te	est_repo									☆ :
Chasse							di Usura	() Dublic Data	> SOL Consolo	& johndoo	
Choose Li	Cense					^	- Aome	e Fusiic Dala	- oge console	E johnuoe	
License: ferp	a removed ferpa License: hipaa re	emoved hipaa L	license: new ferpa rer	moved ferpa							
License ID	Name	PII Def	PII Removed	PII Anonymized	Link To Repo						
1	ferpa removed	ferpa	True	False	Link						
2	hipaa removed	hipaa	True	False	Link						
3	new ferpa removed	ferpa	True	False	Link						
								Manage Ta	ble Application		
							oplied To All Tables				
					Ē	nter			lanage		
			hipaa removed			icense A	oplied To All Tables		anage		
					About	Dod	umentation GitHub Repo API				

DataHub - Chromium								<b>↑</b> 🐼 🖅 9:45 PM 🐴
* DataHub ×								•
← → C () localhost/licenses/johndoe/test_repo								☆:
	DataHub			😤 Home	<b>Q</b> Public Data	≥ SQL Console	🚨 johndoe 👻	
	johndoe / test_repo /							
	Manage Reposito	y Licenses						
	Create New License							
	Add Existing License							
	License Name	Applied To Tables			Manage Tal	ble Application		
	ferpa removed		License Applied To All Tables		м	lanage		
	hipaa removed		License Applied To All Tables		м	lanage		
	test hipaa 3		License Not Applied To All Tables		м	lanage		
		Ab	bout Documentation GitHub Repo	API				

aHub - Chromium License backend - d 🗙 🖽 F	Protect the Open In 🗙 M Screenshots of Dat. X M Screenshots of Dat. X M demo_1_13.png - fa X	🗖 Data Sharing Work 🗙 🐼 Cog	gnitive Complex 🗙 💥	DataHub ×	tµ :	∦ 🖂 ब× 12:
→ C 🛈 localhost/licenses/jo	hndoe/test_repo/8/manage					(
I	DataHub	<b>谷</b> Home	Public Data	≥_SQL Console	🛓 johndoe 🔻	
j	phndoe / test_repo /					
٦	Manage License					
t	est hipaa 3					
1	View Details					
	Base Tables +					
	test		License not a	pplied 🗙 App	ly To Table	
(	Collaborators					
	X user1					
	¥ user2					
ŀ	dd Collaborators					
	Username					
F	ermissions for repo database tables:					
	) select ) update					
	) insert					

delete

DataHub - Chromium		M demo 1 13 ppg - [	a X 🗖 Data Sharin	a Work - X	e Complex X 💥 DataHub	×	ttt 🖇 🖂 ब× 12:06 AM 🖑 ອ
← → C ③ localhost/licenses/johndoe/test_repo/8/manage		(Fracino_r_isiping i					@☆:
DataHub	ſ			🛪 Home 🕠	<mark>@ Publi</mark> c Data >_ SQI	. Console 🔒 johndoo	3 ▼
johndoe / test_repo /	Remove Column				×		
Manage License test hipaa 3	Column: state Remove column						
View Details				CI	lose		
I Base Tables +	daniel NY jane CA	25         20000           20         100000	food server counselor	0			
					nse not applied 🗙	Apply To Table	
				En	nter		
Collaborators				· · ·			
X user1							
× user2							
Add Collaborators							
⊌ select							
✓ update							

DataHub - Chromium	Protect the Open In X M Screenshots of Data	× M Screen	shots of Data	K M demo 1	13.png - fa 🗙 🗖 Data Sł	aring Work 🗙 🔽 Cognitive Co	mplex × 💥 DataHub	×	‡∎ ∦ 🖂 ╡× 12:06 AM 🔱 ⊖
$\leftrightarrow$ $\rightarrow$ C (i) localhost/licens	ses/johndoe/test_repo/8/manage	\ <u>.</u>		(		5			@☆:
	DataHub					🏶 Home 🛛 🚱 Pu	ublic Data 🚬 SQL	Console 🔒 johndo	e <del>-</del>
	johndoe / test_repo /	Data Pre Please ren	eview nove all PII ad						
	Manage License	name edit	age edit	income edit	job edit	num_publications edit			
	test hipaa 3	john	54	22000	physicist	500			
	View Details	sarah	23	58000	economist	20			
		peter	24	120000	fincancial advisor	0			
		daniel	25	20000	food server	0			
	⊞ Base Tables     +	jane	20	100000	counselor	10			
	test						nse not applied 🗙	Apply To Table	
	Collaborators			_	_	Enter			
	X user2								
	Add Collaborators								
	Username								
	Permissions for repo database tables:								
	✓ select								
	<ul> <li>✓ update</li> <li>✓ insert</li> </ul>								

DataHub - Chromium	uborg/da × M Screenshots of Data	× M demo_1_13.png-fa ×	📮 Data Sharing Work 🗵 🗙	Cogn	itive Complex 🗙	🗧 DataHub	×		<b>1</b> ∎ ∦ 🖂 ब× 12:20
← → C ③ localhost/licenses/johndoe/test_repo/8/manage									0
DataHub			*	Home	Public Data	≻_ SQL Co	onsole	🛔 johndoe 🔻	
johndoe / test_repo /									
Manage License									
test hipaa 3									
View Details									
Base Tables +									
test					License a	applied 🗸	Apply	To Table	
test_license_view_8						<u>ا</u>			
Collaborators									
× user1									
× user2									
Add Collaborators									
Username									
Permissions for repo database tab	es:								

select

update

ataHub - Chromium					îtµ 🖾 ╡× 9:47 PM 🎎 [€]
C localhost/browse/johndoe/test_repo/table/test_license_	view_8				☆ :
DataHu	DataHub			ublic Data 🚬 SQL Console 🛔 johndo	€ ▼
johndoe / te	st_repo / table / test_lice	ense_view_8		Open in App 🗸	
					-
Query Builde	Type SQL query here			Run Query	
Tables & Vie	ws Files Cards				
test_lice No description	test_license_view_8 ← No description yet ℃				
Run Sentim	Run Sentiment Analysis -				
state		job	num_publications	income	
МА		physicist	500	22000	
CA		economist	20	58000	
NY		fincancial advisor	0	120000	
NY		food server	0	20000	
CA		counselor	10	100000	
state		job	num_publications	income	

About Documentation GitHub Repo API

## PLATFORM: DEIDENTIFICATION

**De-identification is a major obstacle for data sharing (e.g., HIPAA, FERPA, ...)** 



## HIPAA: INTERACTIVE DE-IDENTIFICATION



## HIPAA: INTERACTIVE DE-IDENTIFICATION

Id	ame	Street	City	State	P-Code	Age
I	J Smith	123 University Ave	Seattle	Washington	98106	42
2	Mary Jones	245 3rd St	Redmond	WA	98052-1234	30
3	Bob Wilson	345 Broadway	Seattle	Washington	98101	19
4	M Jones	245 Third Street	Redmond	NULL	98052	299
5	Robert Wilson	345 Broadway St	Seattle	WA	98101	19
6	James Smith	123 Univ Ave	Seatle	WA	NULL	41
7	JWidom	123 University Ave	Palo Alto	CA	94305 🛕	NULL
•••						



data owner

## HIPAA: INTERACTIVE DE-IDENTIFICATION



## ADDITIONAL TECHNICAL DIRECTIONS

- Support support k-anonymous vies
  - Upload or choose an existing generalization hierarchy
- Support for differentially private views
  - Currently support histogram-based views employing Laplacian noise
  - Adding support for automatically selecting the best differentially private mechanism given a target error rate
  - Based on analysis of data properties to determine best mechanism

## NEXT STEPS

#### Next Data Sharing Spoke Workshop (Summer 2018)

**Collect more agreements and create license framework** 0.1

Extend tooling support, integrated into MIT datahub:

- Watermarking
- Auditing
- Time-based fine-grained access control

- ...

Metadata support

