Privacy Protections as an Incentive for Collaborative Research on Human Health

Anand D. Sarwate

Department of Electrical and Computer Engineering Rutgers, the State University of New Jersey

April 24, 2017





Human health research



There are many data sharing challenges in human health research

- Secondary use of clinical data for research
- Multi-site studies on QA or comparative effectiveness



• Joint (secondary) analyses on aggregated research data

Institutions often want to share data





Sarwate

DIMACS > Human health research

Institutions often want to share data



• Different research groups using the same type of measurements want to do a joint analysis.





3 / 23

Institutions often want to share data



- Different research groups using the same type of measurements want to do a joint analysis.
- Sharing requires lawyers at each institution to generate Data Use Agreements.



Institutions often want to share data



- Different research groups using the same type of measurements want to do a joint analysis.
- Sharing requires lawyers at each institution to generate Data Use Agreements.
- Resulting months of negotiation makes even small-scale collaboration too complicated.



Collaborative research systems



Research consortia are common in many research areas involving human health:







Research consortia are common in many research areas involving human health:

• Foster collaborative research about a particular condition (Alzheimer's, autism, breast cancer, etc.)







Research consortia are common in many research areas involving human health:

- Foster collaborative research about a particular condition (Alzheimer's, autism, breast cancer, etc.)
- Automated sharing is challenging, but this is changing.







Research consortia are common in many research areas involving human health:

- Foster collaborative research about a particular condition (Alzheimer's, autism, breast cancer, etc.)
- Automated sharing is challenging, but this is changing.

Goal: use privacy protections to encourage consortium growth.

COllaborative Informatics Neuroimaging Suite

#COINS

- End-to-end system for managing data for studies on the brain
- Current usage: 37,903 participants in 42,961 scan sessions from 612 studies for a total of 486,955 clinical assessments.
- Data from 34 states, 38 countries
- Partners with research consortia such as the Autism Brain Imaging Data Exchange (ABIDE)



Rutgers









• Goal: build a system that can identify schizophrenia.







- Goal: build a system that can identify schizophrenia.
- **Data:** MRIs from multiple studies (healthy controls and schizophrenics).







- Goal: build a system that can identify schizophrenia.
- **Data:** MRIs from multiple studies (healthy controls and schizophrenics).
- **Algorithm:** classification using machine learning (e.g. support vector machine).







- Goal: build a system that can identify schizophrenia.
- **Data:** MRIs from multiple studies (healthy controls and schizophrenics).
- Algorithm: classification using machine learning (e.g. support vector machine).
- **Privacy risk:** each study has to allow access to sensitive subject data.



Rutgers

State of the art: ENIGMA



http://enigma.ini.usc.edu

"The ENIGMA Network brings together researchers in imaging genomics to understand brain structure, function, and disease, based on brain imaging and genetic data."

- MA = meta analysis : focused on
- Goals: improve reproducibility, sample sizes
- Validation: found genetic variations associated with neurophysiological characteristics (e.g. hippocampal/intercranial volumes)



Workflows in ENIGMA



http://enigma.ini.usc.edu

ENIGMA has 30+ working groups on diseases, genomics, population variation, and methods. To do a study:

- Study proposal is approved by ENIGMA managers.
- Analyses performed on local sites and emailed to ENIGMA manager as Excel spreadsheets.
- Manager has to perform "manual" meta-analysis.





COINSTAC works in a different way: data is registered in the system and analyses are performed/aggregated automatically through message passing.







COINSTAC works in a different way: data is registered in the system and analyses are performed/aggregated automatically through message passing.

• Study is proposed specifying data needed.





Sarwate



COINSTAC works in a different way: data is registered in the system and analyses are performed/aggregated automatically through message passing.

- Study is proposed specifying data needed.
- Local sites approve access to data.







COINSTAC works in a different way: data is registered in the system and analyses are performed/aggregated automatically through message passing.

- Study is proposed specifying data needed.
- Local sites approve access to data.
- Analyses are run and aggregated automatically.







COINSTAC works in a different way: data is registered in the system and analyses are performed/aggregated automatically through message passing.

- Study is proposed specifying data needed.
- Local sites approve access to data.
- Analyses are run and aggregated automatically.

minimize this can be significantly faster than the ENIGMA approach.

The COINSTAC workflow

	Conso	rtia		
E Consortia	Schizophrenia		Adolescent Brain	
Cristaple Damper despedancy	A study of stru disease Ø View Tage:	Users: Rator-osciocal Snapor-osc Harmory2.local	Tracking Chan the Growing B Tags:	pes and Anomalies in an Usens: Raptor-osk.loca Snapper-osk Harmony@.local
	Imaging-Ger Multimodal An Tage:	netics alyris at a Large Scale Users: Reptor-osx/local Snapper-osx Harmory2.local	Substance U Functional Effe without Border @ View Tage:	BB Lts of Substance Abuse Users: Raptor-osx.local Srapper-osx Harmony2.local

In COINSTAC, research groups install the software and register their data in the system:

- Form ongoing and ad-hoc "consortia" (slow, requires approval)
- Once established, consortium members can initiate a joint analysis
- Computation is performed locally and messages passed between sites



DIMACS > COINSTAC

What's in the medium term



COINSTAC prototype is currently "demo-able" but not up and running.

- Compute more than summary statistics, ridge regression, etc.
- Improve user interface and usability for practitioners, including visualization tools.
- Initial subject focus for new results: addiction studies.
- Incorporate/test differentially private methods for machine



learning.

Rutgers

Focusing on "old" algorithms



Because the focus is on usability, we are working on methods popular in neuroimaging:

- Feature discovery: ICA, IVA, NMF, deep learning, etc.
- Regression and classification: ridge regression, LASSO, SVM, etc.
- Visualization: t-SNE, network visualization, etc.



COINSTAC vs. other health data systems



COINSTAC is a solution that works for typical neuroimaging research initiatives.





13 / 23

COINSTAC vs. other health data systems



COINSTAC is a solution that works for typical neuroimaging research initiatives.

1 Data is "big" from the *perspective of the domain area*.







COINSTAC is a solution that works for typical neuroimaging research initiatives.

- 1 Data is "big" from the *perspective of the domain area*.
- **2** Methods with asymptotic guarantees may not be ideal.





COINSTAC is a solution that works for typical neuroimaging research initiatives.

- 1 Data is "big" from the *perspective of the domain area*.
- 2 Methods with asymptotic guarantees may not be ideal.
- **8** Strong formal privacy may be trumped by utility requirements.





Sad news: no privacy is enough?



From the perspective of IRBs and other regulatory bodies, decentralized/distributed algorithms may be "good enough."

- Getting them to work on the computing infrastructure is itself challenging.
- Threat models and surface are different than "typical" data sharing scenarios.
- Provides a useful test case for "newer" privacy technologies: differential privacy, multiparty computation.



Rutgers

Making formal privacy guarantees

Currently working on making *differentially private* versions of existing algorithms. Differential privacy involves introducing randomization (e.g. noise) in computations.

- Small number of subjects \rightarrow larger noise \rightarrow more error.
- Neuroimaging data is high-dimensional: need some dimension reduction.
- Preference for stronger $(\varepsilon,0)$ guarantees, but improved analyses give $(\varepsilon,\delta).$



DIMACS > Privacy challenges

Dealing with federated infrastructure



Uploading all the data to EC2 or Azure is not an acceptable (yet?)





16 / 23

Dealing with federated infrastructure



Uploading all the data to EC2 or Azure is not an acceptable (yet?)

• Local storage overhead can be challenging.





Dealing with federated infrastructure



Uploading all the data to EC2 or Azure is not an acceptable (yet?)

- Local storage overhead can be challenging.
- Local processing costs are heterogeneous.



Dealing with federated infrastructure



Uploading all the data to EC2 or Azure is not an acceptable (yet?)

- Local storage overhead can be challenging.
- Local processing costs are heterogeneous.
- Communication can act as a real bottleneck.



Compromises, compromises



At the moment we are making many compromises:

- Utility first: practical values of ε for differential privacy are large.
- Low communication: focus on one-shot aggregation over iterative methods.
- Simple tasks: stick with developing distributed methods for known algorithms.



Policy and privacy and systems, oh my!

Data sharing in health research may be different than open sharing or industry/academia sharing.

- Different regulations around human subjects for experimental data or for PHI in clinical data.
- Informed consent model allows *subject-level* and *study-level* privacy preferences.
- Data sharing is contingent and possibly transient: revert to access only.





Shared-access models with "privacy protections" (formal or not) can encourage researchers to join consortia.

- Benefits/risks align with the desires of data holders/researchers.
- Data holders retain control over access and allowed computations.
- Data users can use automated computations for hypothesis generation.













• start small: variability of problem types is large







- start small: variability of problem types is large
- challenging to bridge gaps between algorithmist and developer





- start small: variability of problem types is large
- challenging to bridge gaps between algorithmist and developer
- communication requirements are nontrivial





The iDASH center at UCSD is working on larger-scale human health research involving clinical records.





21 / 23



The iDASH center at UCSD is working on larger-scale human health research involving clinical records.

• Goal: to make clinical data warehouse more useful to researchers







The iDASH center at UCSD is working on larger-scale human health research involving clinical records.

- Goal: to make clinical data warehouse more useful to researchers
- Diverse range of problems in compression, genomics, NLP, etc. with privacy







The iDASH center at UCSD is working on larger-scale human health research involving clinical records.

- Goal: to make clinical data warehouse more useful to researchers
- Diverse range of problems in compression, genomics, NLP, etc. with privacy
- Spurring research transition through data challenges, internships, etc.







The iDASH center at UCSD is working on larger-scale human health research involving clinical records.

- Goal: to make clinical data warehouse more useful to researchers
- Diverse range of problems in compression, genomics, NLP, etc. with privacy
- Spurring research transition through data challenges, internships, etc.



The features of these problems are very different!

🧠 ≠ 🌈 ≠ 🏷 ≠ 🌹





 $\neq \swarrow \neq \checkmark \neq \checkmark \neq$ 5

22 / 23

Sarwate

• Recognize that "medical data" is at best a placeholder and at worst semantically void.



 $\neq \swarrow \neq \checkmark$ Ş

- Recognize that "medical data" is at best a placeholder and at worst semantically void.
- Spend some time delineating the problem space and domain-specific challenges.





22 / 23

- Recognize that "medical data" is at best a placeholder and at worst semantically void.
- Spend some time delineating the problem space and domain-specific challenges.
- For theorists: can we get out of asymptopia?





 $\neq \swarrow \neq$

- Recognize that "medical data" is at best a placeholder and at worst semantically void.
- Spend some time delineating the problem space and domain-specific challenges.
- For theorists: can we get out of asymptopia?
- For practitioners: what do you want to *do* versus *how* do you want to do it?

