# Big Data PI Meeting Privacy and Ethics Breakout Session

Adam Smith
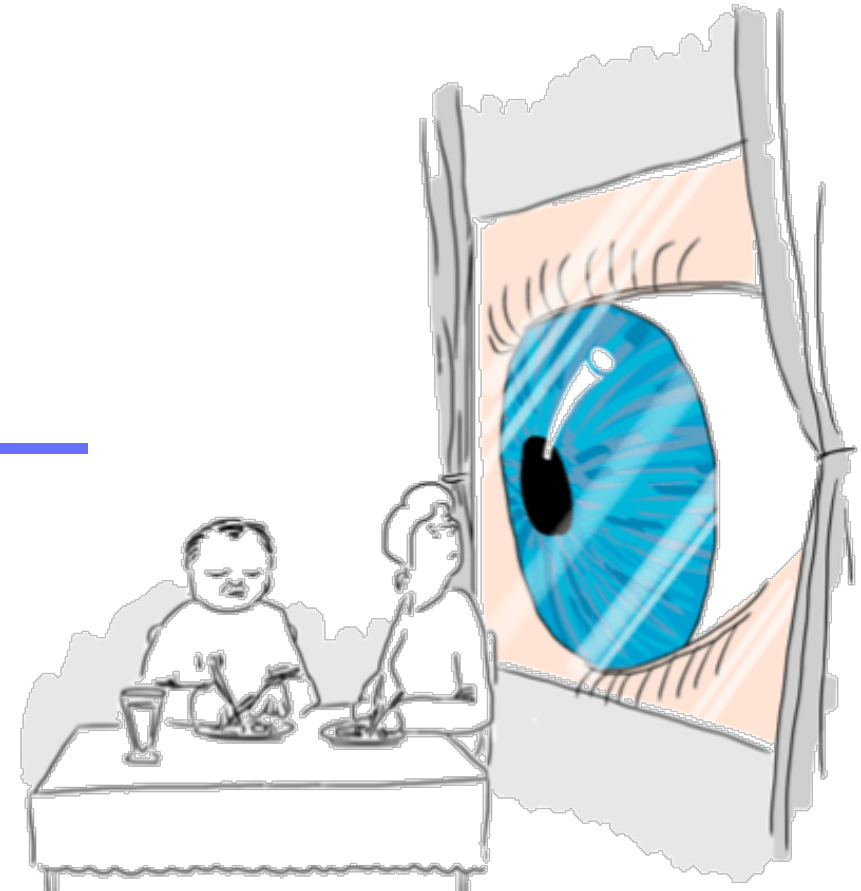Penn State EECS

Fang Liu
University of Notre Dame

& People attending the session

"Relax – it can only see metadata."

Cartoon: NOISE TO SIGNAL
RobCottingham.com

# *Topics*

- Overarching Themes in this Area

- Recent Successes (last 3 years)

- Major Obstacles impeding More Rapid Progress

- Areas of Neglect

- Strategic Priorities & Investments that will Advance Innovation

# *Topics*

- Overarching Themes in this Area

- Recent Successes (last 3 years)

- Major Obstacles impeding More Rapid Progress

- Areas of Neglect

- Strategic Priorities & Investments that will Advance Innovation

# *Overarching Themes*

What are the stakes?

- Surveillance / changing how enforcement works
  - ➢ Trust of the system is a big issue
- Asymmetry of information
  - ➢ trade secrets, market control
  - ➢ Data collection favors large organizations that have the resources to aggregate / use the data (e.g. farmers)
- Compliance
  - ➢ Infrastructure needs clear guidelines, especially when computation/storage are delegated
  - ➢ How can we make shared infrastructure (e.g. NSF supercomputer) with easy compliance certification?

# *Overarching Themes (cont'd)*

- Freedom of speech
  - ➢ Lack of privacy has a chilling effect on
    - Free speech
    - Free association
    - Free religion, etc
  - ➢ Made more acute by large-scale data aggregation across many sources
    - How can we understand / measure / control how disparate data sources are merged or aggregated?
  - ➢ If everything we say is recorded, how do we self-censor?
    - Protected conversations with lawyers, therapists, etc
    - Should other types of speech enjoy similar protection?

# *Overarching Themes (cont'd)*

- Fairness and disparate impact
  - ➢ Predictions and decisions should always come with measures of uncertainty (e.g. confidence / probabilities)
    - Understand how to interpret these measures
    - How do these measures correspond to legal standards?
  - ➢ Data-driven systems are large and complex
    - Consist of many institutions
    - Humans sometimes should be involved but actually not
- Accountability
  - ➢ How can we insist on transparency of algorithms
    - Requires algorithms in some human-interpretable form
    - How do we distinguish the algorithm from the data it relies on?
  - ➢ How can people correct/control the data about themselves that systems rely on?

# *Tools / successes*

- Encryption
  - ➢ Can we ensure end-to-end security with distributed data?
- Some commercial cloud providers do provide compliant services in some cases
- De-anonymization & "privacy-protected" data
  - ➢ Differential privacy
  - ➢ Synthetic data

# *Major Obstacles & Areas of Neglect*

- Data and technology as power
- New types of crime enabled by data and technology
  - ➤ "Cybercasing", e.g., figuring out which homes to rob based on YouTube vacation videos
- New types of harm
  - ➤ Information aggregated in new ways
  - ➤ E.g., Algorithmic discrimination
- New difficulties for regulation
  - ➤ Technology and data are not in a well-defined location
  - ➤ Fundamentally different value systems in different countries or regions
  - ➤ How can these heterogeneous constraints and commitments be resolved with in a small set of technological systems
  - ➤ Every system involves many collections of values and regulation

# *Strategic Priorities*

- Lots of work on learning and analytics but little on auditing
  - ➢ How can you "self-audit" to see if the algorithms being applied are a good fit to the data or setting we have
  - ➢ Formal verification of a large system for accountability
- When does learning "global" properties of a data set cause serious problems (e.g. does not prevent systemic discrimination)
- How can we mitigate unintended privacy/ethics consequences of big data?
- IRBs have successfully (?) forced articulation of human-subjects issues
  - ➢ Can we have similar structures for data privacy/ethics issues?
  - ➢ How do you create a culture of conversation around these issues?
  - ➢ Compliance?

# *Strategic Priorities*

- Education
  - ➢ Engaging technical / research community
  - ➢ Educating everyone else
  - ➢ Broad training in understanding these issues for everyone (lawyers, policy makers, everybody else)
- How we bridge gap between technical and nontechnical discussions?
  - ➢ How do you articulate natural language versions of technical tradeoffs?
  - ➢ Can we benefit from experience of debates in public health (surveillance vs public benefits)?

# *Strategic Priorities*

- Understanding what controls data-driven systems
  - ➢ Lessig's four forces, Law, Social norms, Technology, Market
- Technology researchers have a better understanding of technology's effect
  - ➢ How can we be encouraged to articulate and explicitly think about the risks of the technologies we work on?
  - ➢ How can we reward researchers who think carefully about these issues?
- How can education reflect this understanding?
  - ➢ For technical students?
  - ➢ More broadly?