# Towards an Axiomatization of Privacy and Utility
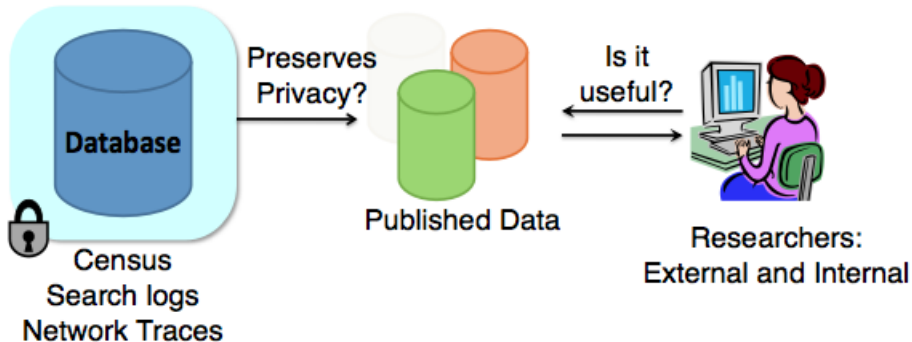
Daniel Kifer
Bing-Rong Lin

Department of Computer Science & Engineering
Penn State University

# Guiding Principles?

| SSN | Gender | Age | Zip Code | Disease |
|-----|--------|-----|----------|---------|
| 111111111 | M | 25 | 90210 | AIDS |
| 222222222 | F | 43 | 90211 | AIDS |
| 333333333 | M | 29 | 90212 | Cancer |
| 456456456 | M | 41 | 90213 | AIDS |
| 567867867 | F | 41 | 07620 | Cancer |
| 654321566 | F | 40 | 33109 | Cancer |
| 799999999 | F | 40 | 07620 | Flu |
| 800000000 | F | 24 | 33109 | None |
| 934587938 | M | 48 | 07620 | None |
| 109494949 | F | 40 | 07620 | Flu |
| 112525252 | M | 48 | 33109 | Flu |
| 121111111 | M | 49 | 33109 | None |

# Guiding Principles?

- We know this is not enough

| SSN | Gender | Age | Zip Code | Disease |
|---|---|---|---|---|
| ~~111111111~~ | M | 25 | 90210 | AIDS |
| ~~222222222~~ | F | 43 | 90211 | AIDS |
| ~~333333333~~ | M | 29 | 90212 | Cancer |
| ~~456456456~~ | M | 41 | 90213 | AIDS |
| ~~567867867~~ | F | 41 | 07620 | Cancer |
| ~~654321566~~ | F | 40 | 33109 | Cancer |
| ~~799999999~~ | F | 40 | 07620 | Flu |
| ~~800000000~~ | F | 24 | 33109 | None |
| ~~934587938~~ | M | 48 | 07620 | None |
| ~~109494949~~ | F | 40 | 07620 | Flu |
| ~~112525252~~ | M | 48 | 33109 | Flu |
| ~~121111111~~ | M | 49 | 33109 | None |

# So what happens?

- Aug 6, 2006 - AOL releases data
  - 20 Million Search Queries from 3 months
  - 650,000 users
- How is data protected: Change AOL id to a number.
- What happened?
  - NYT identified user # 4417749
    - People search for names of friends/relatives/self
    - People search for locations "What to do in State College"
    - Age-related searches
  - Many people got fired.

# Outline

# Statistical Privacy

- Art of turning sensitive data into nonsensitive data suitable for public release.
- Sensitive data:
    - Cannot release sensitive data directly.
    - Detailed information about individuals (search logs, health records, census/tax data, etc.)
    - Proprietary secrets (search logs, network traces, machine debug info)
- Want to release useful but non-private information from this data.
    - Typical user web search behavior
    - Demographics
    - Information that can be used to build models
    - Information that can be used to design & evaluate algorithms
- **Mechanism**: a (randomized) algorithm that converts sensitive into nonsensitive data.
- Goal: Design a mechanism that protects privacy and provides utility

# Privacy & Utility

- What does privacy mean?
    - Many, many privacy definitions in the literature.
    - How do I compare them?
    - How do I identify strengths and weaknesses?
    - How do I customize them (for an application)?
    - How do I design one?
    - Does it really do what I want it to do?
    - What statements are/aren't privacy definitions?
- What does utility mean?
    - Many, many measures of utility in the literature:
        - KL-divergence.
        - Expected (Bayesian) utility.
        - Minimax estimation error.
        - Task-specific measures.
    - Which one should I choose?
    - Does it do what I want it to do?
    - How do I design one?
    - Does it make sense in statistical privacy?

# A Common Approach

1. Start with a privacy mechanism.
   - Generalization (e.g. coarsen "state college" → "Pennsylvania")
   - Suppression (remove parts of data items)
   - Add random noise
2. Create privacy definition that feels most natural with this privacy mechanism.
3. Create utility measure that feels most natural for this mechanism.
   - # of generalizations
   - # of suppressions
   - variance of noise
   - anything we can borrow from statistics
   - often can't compare utility across mechanisms
4. (Usually) Find flaws, revise steps 2 and 3.

# The Axiomatic Approach

- What if we did this in reverse? For a given application:
  1. Identify properties we think a privacy definition should satisfy.
  2. Identify properties we think a utility metric should satisfy.
  3. Find a privacy mechanism that satisfies those properties.
- Benefits of axiomatization:
  - Apples to apples comparison of properties of privacy definitions.
  - Small set of axioms easier to study than large set of privacy definitions.
  - Abstract approaches yield general results and insights (e.g. group theory, vector spaces, etc.)
  - Can study relationships between axioms.
  - Easier to identify weaknesses.
  - Design mechanisms by picking axioms depending on application.
  - Can study consequences of omitting axioms.
- Is it really necessary for privacy and utility?
  - Let's look at some illustrative results.

# Outline

# Axioms for Privacy

- Hard to create a good privacy definition.
- Simple things usually don't work.
- Different applications have different privacy requirements.
- Instead of starting from a privacy definition:
  - Identify axioms you want it to support.
  - Determine the privacy definition implied by axioms
  - Let axioms be the building blocks.
- It is easier to reason about axioms that about entire privacy definitions.
- Efficiency: insights into 1 axiom lead to insights into many privacy definitions.
- Example: how to relax differential privacy.

# Some definitions

- Abstract input space $\mathfrak{I}$ (all possible data).
    - Semantics (e.g. neighboring databases in differential privacy) should be given by axioms.
- Abstract output space $\mathfrak{O}$.
    - Semantics (e.g. query answers, synthetic data, utility) should be given by axioms.

### Definition (Randomized Algorithm)

A randomized algorithm $\mathcal{A}$ is a regular conditional probability distribution $P(O \mid I)$ with $O \subset \mathfrak{O}$ and $I \subset \mathfrak{I}$

- Privacy definition: intentionally undefined (all parameters must be instantiated).

### Definition (Privacy Mechanism for $D$)

A privacy mechanism $\mathfrak{M}$ is a randomized algorithm that satisfies privacy definition $D$.

# Two Simple Privacy Axiom

- Intuition: postprocessing the output of a privacy mechanism should still maintain privacy.

### Axiom (Transformation Invariance)

*Given a privacy mechanism $\mathfrak{M}$ and a randomized algorithm $\mathcal{A}$ (independent of the data and $\mathfrak{M}$), the composition $\mathcal{A} \circ \mathfrak{M}$ is a privacy mechanism.*

- Intuition: it does not matter which privacy mechanism I choose.

### Axiom (choice)

*If $\mathfrak{M}_1$ and $\mathfrak{M}_2$ are privacy mechanisms for $D$, then the process of choosing $\mathfrak{M}_1$ with probability $c$ and $\mathfrak{M}_2$ with probability $1 - c$ (with randomness independent of the data, $\mathfrak{M}_1$, and $\mathfrak{M}_2$) results in a privacy mechanism for $D$.*

# Two Simple Privacy Axiom

### Axiom (Transformation Invariance)

*Given a privacy mechanism $\mathfrak{M}$ and a randomized algorithm $\mathcal{A}$ (independent of the data and $\mathfrak{M}$), the composition $\mathcal{A} \circ \mathfrak{M}$ is a privacy mechanism.*

### Axiom (choice)

*If $\mathfrak{M}_1$ and $\mathfrak{M}_2$ are privacy mechanisms for $D$, then the process of choosing $\mathfrak{M}_1$ with probability $c$ and $\mathfrak{M}_2$ with probability $1 - c$ (with randomness independent of the data, $\mathfrak{M}_1$, and $\mathfrak{M}_2$) results in a privacy mechanism for $D$.*

- Consistency conditions for privacy definitions
- Thus privacy definitions should discuss how they are affected by postprocessing.
- Privacy definitions cannot focus only on deterministic mechanisms.
- Many privacy definitions do not satisfy these axioms!

# Applications Differential Privacy

### Definition (Differential Privacy [Dwo06, DMNS06])

$\mathfrak{M}$ satisfies $\epsilon$-differential privacy if $P(\mathfrak{M}(i_1) \in S) \leq e^{\epsilon} P(\mathfrak{M}(i_2) \in S)$ for all measurable $S \subset \mathfrak{O}$ and all neighboring input databases $i_1, i_2 \in \mathfrak{I}$.
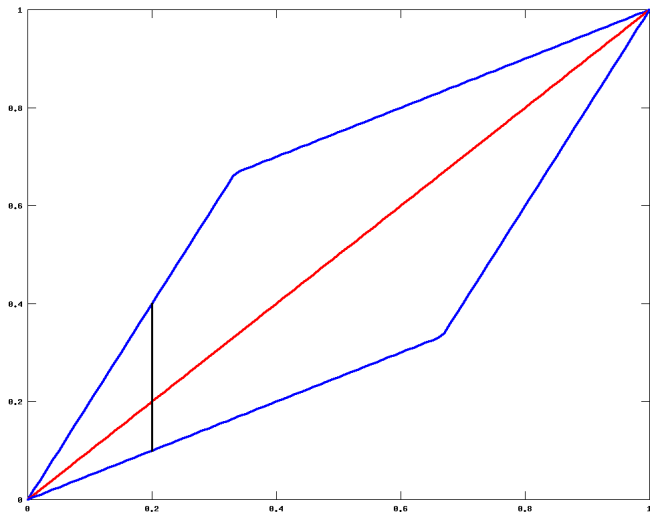
- There has been interest in relaxing differential privacy. For example: For example:

$$P(\mathfrak{M}(i_1) \in S) \leq e^{\epsilon} P(\mathfrak{M}(i_2) \in S) + \delta$$

# Example

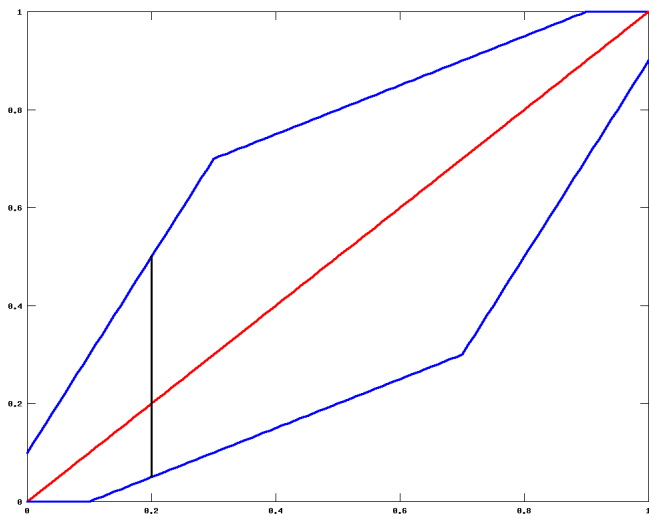$$a = P(\mathfrak{M}(i_1) \in S) \qquad b = P(\mathfrak{M}(i_2) \in S) \qquad a \leq 2b$$

# Example

$a = P(\mathfrak{M}(i_1) \in S)$    $b = P(\mathfrak{M}(i_2) \in S)$    $a \leq 2b + .1$

# Applications Differential Privacy

## Definition (Differential Privacy [Dwo06, DMNS06])

$\mathfrak{M}$ satisfies $\epsilon$-differential privacy if $P(\mathfrak{M}(i_1) \in S) \leq e^\epsilon P(\mathfrak{M}(i_2) \in S)$ for all measurable $S \subset \mathfrak{O}$ and all neighboring input databases $i_1, i_2 \in \mathfrak{I}$.

- There has been interest in relaxing differential privacy. For example: For example:

$$P(\mathfrak{M}(i_1) \in S) \leq e^\epsilon P(\mathfrak{M}(i_2) \in S) + \delta$$

## Definition (A Generic Version)

$\mathfrak{M}$ is a privacy mechanism if $G\left[P(\mathfrak{M}(i_1) \in S), P(\mathfrak{M}(i_2) \in S)\right] = T$ for all measurable $S \subset \mathfrak{O}$ and all neighboring input databases $i_1, i_2 \in \mathfrak{I}$.

- What other predicates can be used?

# Relaxations of Differential Privacy

## Definition (A Generic Version)

$\mathfrak{M}$ is a privacy mechanism if $G\left[P(\mathfrak{M}(i_1) \in S), P(\mathfrak{M}(i_2) \in S)\right] = T$ for all measurable $S \subset \mathfrak{O}$ and all neighboring input databases $i_1, i_2 \in \mathfrak{I}$.

- In principle, $G$ could be any predicate:
  - $G(a, b) = T$ if $a - b$ is rational.
  - $G(a, b) = T$ if $a < b^2$.
  - $G(a, b) = T$ if $b = (1 + \cos(2\pi a))/2$
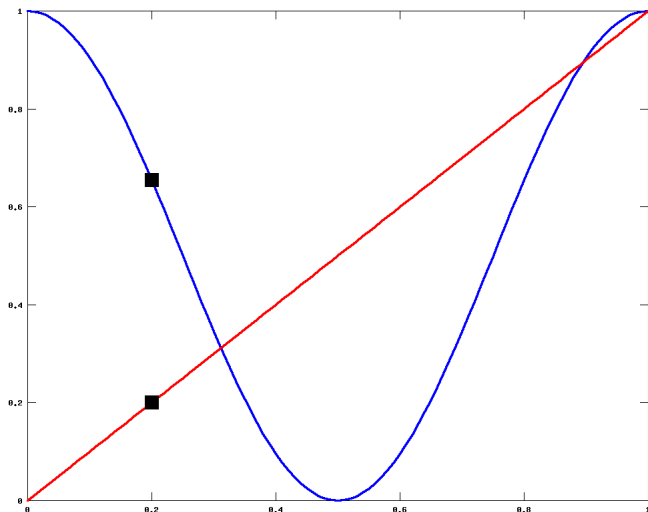- Choice and Transformation Invariance Axioms limit the possibilities.

# Example

$a = P(\mathfrak{M}(i_1) \in S)$ $\qquad$ $b = P(\mathfrak{M}(i_2) \in S)$ $\qquad$ $b = (1 + \cos(2\pi a))/2$

# Relaxations of Differential Privacy

## Definition (A Generic Version)

$\mathfrak{M}$ is a privacy mechanism if $G\left[P(\mathfrak{M}(i_1) \in S), P(\mathfrak{M}(i_2) \in S)\right] = T$ for all measurable $S \subset \mathfrak{O}$ and all neighboring input databases $i_1, i_2 \in \mathfrak{I}$.

- Replacing $G[a, b]$ with $G^*[a, b] \equiv G[a, b] \land G[1 - a, 1 - b]$ does not change privacy definition.
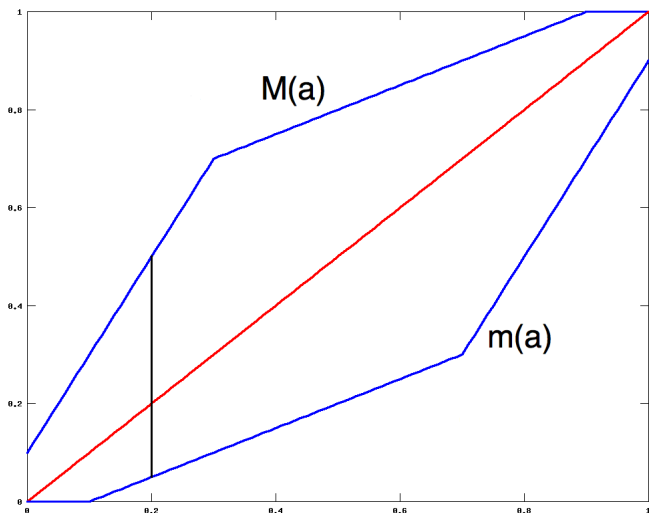
## Theorem

*Axioms of Transformation Invariance and Choice provide necessary and sufficient conditions on $G^*[a, b]$. There exists a well-behaved upper envelope $M(a)$ and lower envelope $m(a)$ that determine $G^*$.*

# See paper for details

$$a = P(\mathfrak{M}(i_1) \in S) \qquad b = P(\mathfrak{M}(i_2) \in S)$$



- $M(a)$ is
  - continuous*
  - concave
  - strictly increasing*
- $m(a)$ is determined by $M(a)$

# Summary

### Definition (A Generic Version)

$\mathfrak{M}$ is a privacy mechanism if $G\left[P(\mathfrak{M}(i_1) \in S), P(\mathfrak{M}(i_2) \in S)\right] = T$ for all measurable $S \subset \mathfrak{O}$ and all neighboring input databases $i_1, i_2 \in \mathfrak{I}$.

- Axioms imply a nice intuitive form for predicate $G$.
- For every $a$, there is interval of allowable $b$ values
- Interval endpoints vary nicely with $a$.
- Makes sense intuitively
    - But no need for intuition after axioms are selected
    - Avoids faulty/incomplete intuition

# Axioms for Utility?

- Privacy axioms limit the privacy mechanisms we can consider.
- How to choose among allowable mechanisms?
  - $\mathfrak{M}$ as a column stochastic matrix:
  - Column $i$ of $\mathfrak{M}$ is $P_{\mathfrak{M}}(\cdot \mid i)$.
- $\mu(\mathfrak{M})$ – how good is a privacy mechanism $\mathfrak{M}$?
  - How much information does it contain?
  - How useful are the outputs?
- Do we understand utility well enough?

# Example: Expected Utility

- Conducting a survey: Is this your favorite conference venue?
- Sensitive question, people may not respond truthfully.
- Idea: allow respondent to lie with certain probability (randomized response [War65]).
- Utility: expected loss (?)
    - I get a loss of 1 every time they lie (0 loss for truth)
    - I believe 75% of population could not imagine a better conference venue
    - Expected loss what do I believe my average (expected) loss is?

# Example: Expected Utility

- Is this your favorite conference venue?
- Subjective prior belief: 75% yes

Privacy Mechanism $\mathfrak{M}_2$

|     | True Answer | |
| --- | --- | --- |
|     | Yes | No |
| Yes | 1 | 1 |
| No | 0 | 0 |

$$E[\text{Loss}] = 1 \times 1/4$$
$$= 1/4$$

Privacy Mechanism $\mathfrak{M}_1$

|     | True Answer | |
| --- | --- | --- |
|     | Yes | No |
| Yes | 2/3 | 1/3 |
| No | 1/3 | 2/3 |

$$E[\text{Loss}] = 1 \times 3/4 \times 1/3$$
$$+1 \times 1/4 \times 1/3$$
$$= 1/3$$

- Mechanism $\mathfrak{M}_2$ has lower expected loss
- Yet contains no information
- $\mathfrak{M}_2(\text{true answer}) = \mathcal{A}(\mathfrak{M}_1(\text{true answer}))$

# Example: Expected Utility

- User has a prior distribution over the input space $\mathfrak{I}$.
- Output space $\mathfrak{O} = \mathfrak{I}$.
- User has a loss function $L(i, j)$.
- Create mechanism with smallest expected loss.

### Theorem ([GRS09])

*Under suitable conditions on $\mathfrak{I}$ and $L$, the geometric mechanism is universal – for any prior, the optimal mechanism is achieved by applying a many-to-one deterministic function to the output of geometric mechanism.*

- In general, cannot recover geometric mechanism from "optimal" mechanism.
- $\therefore$ "Optimal" mechanism contains less information than geometric mechanism.
    - "Optimal" mechanism should not be considered optimal.
    - Expected utility may not be an appropriate measure of utility.

# How to measure utility

- We should take a step back and think about what properties our utility measures should have.

### Definition (Sufficiency partial order)

Privacy mechanism $\mathfrak{M}_1$ is sufficient for $\mathfrak{M}_2$ ($\mathfrak{M}_2 \prec \mathfrak{M}_1$) if there exists a randomized algorithm $\mathcal{A}$ such that $\mathfrak{M}_2 = A \circ \mathfrak{M}_1$.

### Axiom (Sufficiency)

If $\mathfrak{M}_2 \prec \mathfrak{M}_1$ then $\mu(\mathfrak{M}_2) \leq \mu(\mathfrak{M}_1)$

### Definition (Sufficient Covering Set)

A set $S$ of privacy mechanisms is a covering set if every mechanism in $S$ is maximally sufficient and: $\forall \mathfrak{M}, \exists \mathfrak{M}^* \in S$ such that $\mathfrak{M} \prec \mathfrak{M}^*$

- Utility metric $\mu$ should choose some $\mathfrak{M}^* \in S$.

## Examles - finite input/output spaces

$$\mathfrak{M} = \begin{pmatrix} P(O_1 \mid *) \\ P(O_2 \mid *) \\ P(O_3 \mid *) \\ P(O_4 \mid *) \end{pmatrix} = \begin{pmatrix} P(O_1 \mid i_1) & P(O_1 \mid i_2) & P(O_1 \mid i_3) \\ P(O_2 \mid i_1) & P(O_2 \mid i_2) & P(O_2 \mid i_3) \\ P(O_3 \mid i_1) & P(O_3 \mid i_2) & P(O_3 \mid i_3) \\ P(O_4 \mid i_1) & P(O_4 \mid i_2) & P(O_4 \mid i_3) \end{pmatrix}$$

# Examples

- $|det\ \mathfrak{M}|$
  - For finite input space and output space of the same size.
  - Measures how much $\mathfrak{M}$ shrinks the unit hypercube (identity matrix).
  - Piecewise multilinear.
- Negative Dobrushin's coefficient of ergodicity.
  - $-\min_{j,k} \sum \min(m_{i,j}, m_{i,k})$
  - Finds the two columns that are hardest to distinguish.
  - Finds the two inputs hardest to distinguish.
  - Another measure of how the matrix contracts the input space [CDZ93].
- Branching Measures.
  - $\sum_i F(r_i)$
  - $r_i$ are the rows
  - $F$ is convex and $F(cx) = cF(x)$.
  - Example:

$$F(x_1, \ldots, x_n) = \sum_{i=1}^{n} x_i \log \frac{x_i}{x_1 + \cdots + x_n}$$

# Maximally Sufficient Mechanisms
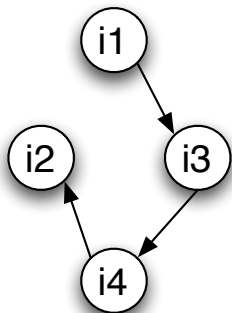
## Definition (Sufficient Covering Set)

A set $S$ of privacy mechanisms is a covering set if every mechanism in $S$ is maximally sufficient and: $\forall \mathfrak{M}, \exists \mathfrak{M}^* \in S$ such that $\mathfrak{M} \prec \mathfrak{M}^*$

- What do they look like?
- For finite input spaces, output space is finite but larger.
- Neighboring databases form a connected graph of input space.
- For each output $o_1$, its row subgraph must be a spanning tree*.
- Output space can be identified with a set of graphs.
    - Output space is a set of spanning trees* of input space.
    - Edges correspond to equality constraints in differential privacy.
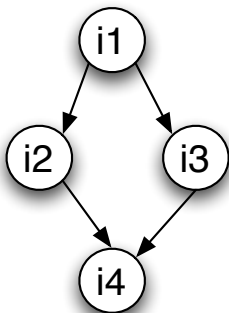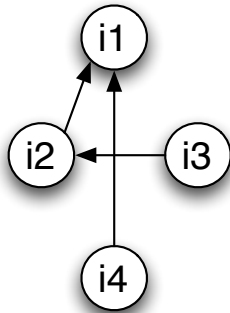    - Can also be interpreted as a restricted set of likelihood functions.

# Output Space

O1                    O2                    O3



P(O1 | * )          P(O2 | * )          P(O3 | * )

# Insights

- Output of a privacy mechanism many not correspond to a query answer.
    - Input: heads or tails
    - Output: red or blue or green
- Output of a privacy mechanism many not correspond to synthetic data.
    - May not have "attributes"
    - May not have "rows"
- You will need to postprocess the output for what you want to do.
- Use the likelihood principle.
- Goal: find a mechanism that allows greatest flexibility for postprocessing.

# Take home message

- Axioms are our building blocks.
    - Easier to understand and argue about than privacy definitions and utility measures.
    - Abstraction allows for generality.
    - Allows for comparison of privacy definitions.
- Shouldn't specify privacy definition directly, let axioms disqualify sets of randomized algorithms.
- Use axioms to choose the best mechanisms via utility.
- Output space may not correspond to query answers or synthetic data.
    - Because of potentially many different uses for the data.
- Need statistical postprocessing tools to work with resulting data.

📄 Joel E. Cohen, Yves Derriennic, and Gh. Zbaganu.

Majorization, monotonicity of relative entropy and stochastic matrices.

Contemporary Mathematics, 149, 1993.

📄 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith.

Calibrating noise to sensitivity in private data analysis.

In Theory of Cryptography Conference, pages 265–284, 2006.

📄 Cynthia Dwork.

Differential privacy.

In ICALP, 2006.

📄 Arpita Ghosh, Tim Roughgarden, and Mukund Sundararajan.

Universally utility-maximizing privacy mechanisms.

In STOC, 2009.

📄 S. L. Warner.

Randomized response: A survey technique for eliminating evasive answer bias.

Journal of the American Statistical Association, 1965.