

CSE660

Differential Privacy

August 28, 2017

Marco Gaboardi

Room: 338-B

gaboardi@buffalo.edu

<http://www.buffalo.edu/~gaboardi>

Differential Privacy

Data



Aol.

Statistics over Data



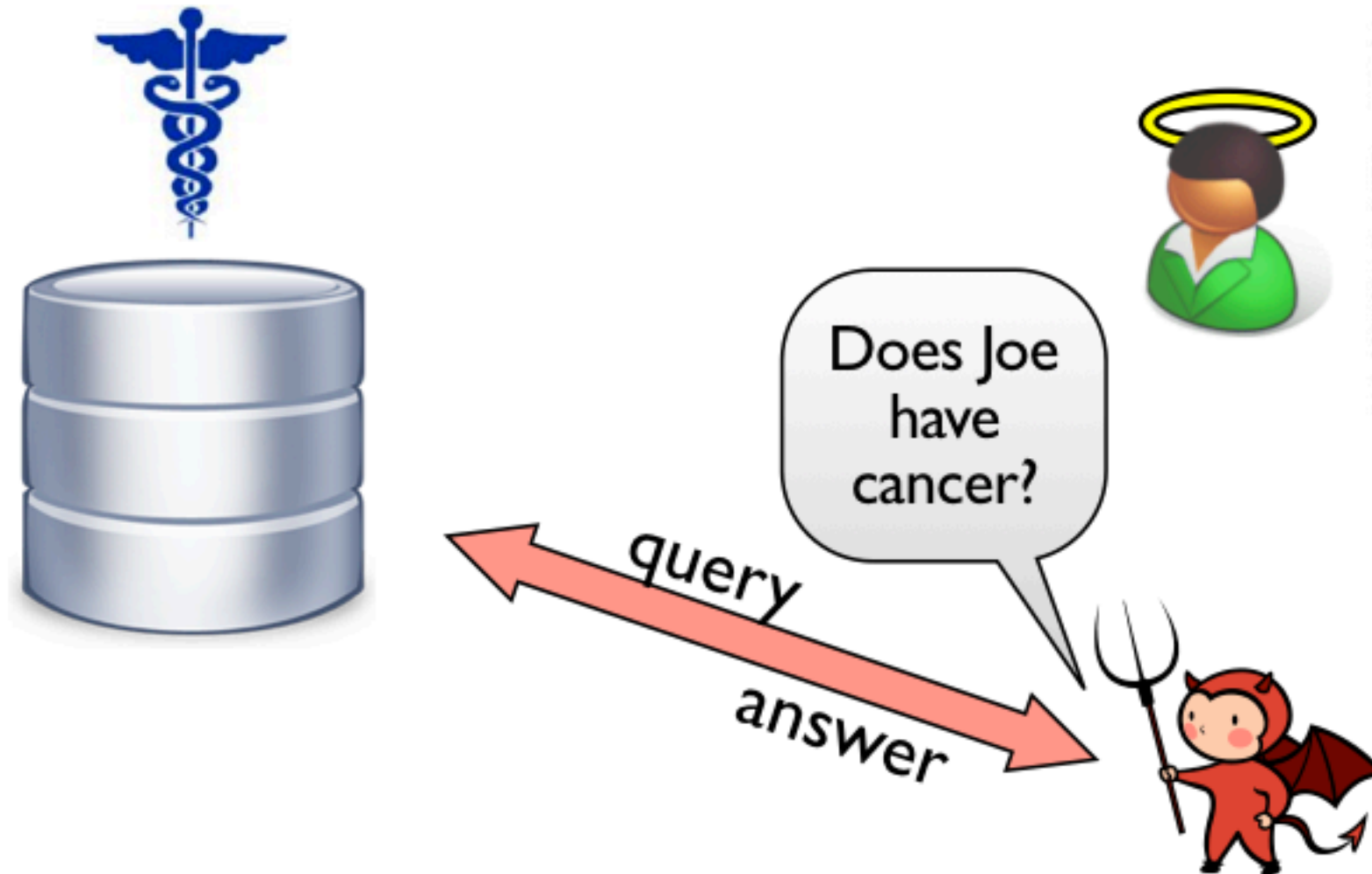
NETFLIX

Google

Private Queries?



Private Queries?



Methodological weakness in using correlation coefficients for assessing the interchangeability of analytic data between samples collected under different sampling conditions - the example of matrix metalloproteinase 2 determined in serum and plasma samples

Letter to the Editor

Dear Sirs,

We read with interest the article by [Author Name] et al. (1) in which the authors reported the results of a study on the interchangeability of analytic data between samples collected under different sampling conditions. The authors used correlation coefficients to assess the interchangeability of analytic data between samples collected under different sampling conditions. However, the use of correlation coefficients is not the most appropriate method for this purpose. The most appropriate method for this purpose is the use of the coefficient of variation (CV). The CV is a measure of the relative variability of a data set and is calculated as the standard deviation divided by the mean. The CV is a dimensionless quantity and is therefore independent of the units of the data. The CV is a more appropriate measure of variability than the standard deviation because it takes into account the magnitude of the mean. The CV is also a more appropriate measure of variability than the standard deviation because it is not affected by the units of the data. The CV is a more appropriate measure of variability than the standard deviation because it is not affected by the units of the data. The CV is a more appropriate measure of variability than the standard deviation because it is not affected by the units of the data.

Private Queries?



Does Joe
have
cancer?



Letter to the Editor

Methodological weakness in using correlation coefficients for assessing the interchangeability of analytic data between samples collected under different sampling conditions - the example of matrix metalloproteinase 2 determined in serum and plasma samples

Dear Sirs,

I have read with interest the article by [Author Name] et al. (1) in your journal, which deals with the methodological weakness in using correlation coefficients for assessing the interchangeability of analytic data between samples collected under different sampling conditions - the example of matrix metalloproteinase 2 determined in serum and plasma samples.

The authors state that the correlation coefficient (r) is a measure of the strength of the linear relationship between two variables. In this case, the variables are the concentration of matrix metalloproteinase 2 in serum and plasma samples. The authors claim that the correlation coefficient is a reliable measure of the interchangeability of analytic data between samples collected under different sampling conditions.

However, I believe that the use of correlation coefficients in this context is methodologically weak. Correlation coefficients are only a measure of the strength of the linear relationship between two variables. They do not take into account the underlying biological and analytical factors that may influence the results. For example, the concentration of matrix metalloproteinase 2 in serum and plasma samples may be influenced by a variety of factors, including the patient's clinical status, the sampling method, and the analytical method used.

In my opinion, a more robust method for assessing the interchangeability of analytic data between samples collected under different sampling conditions would be to use a method that takes into account the underlying biological and analytical factors. For example, a method that uses a combination of correlation coefficients and other statistical measures, such as the coefficient of variation, would be more appropriate.

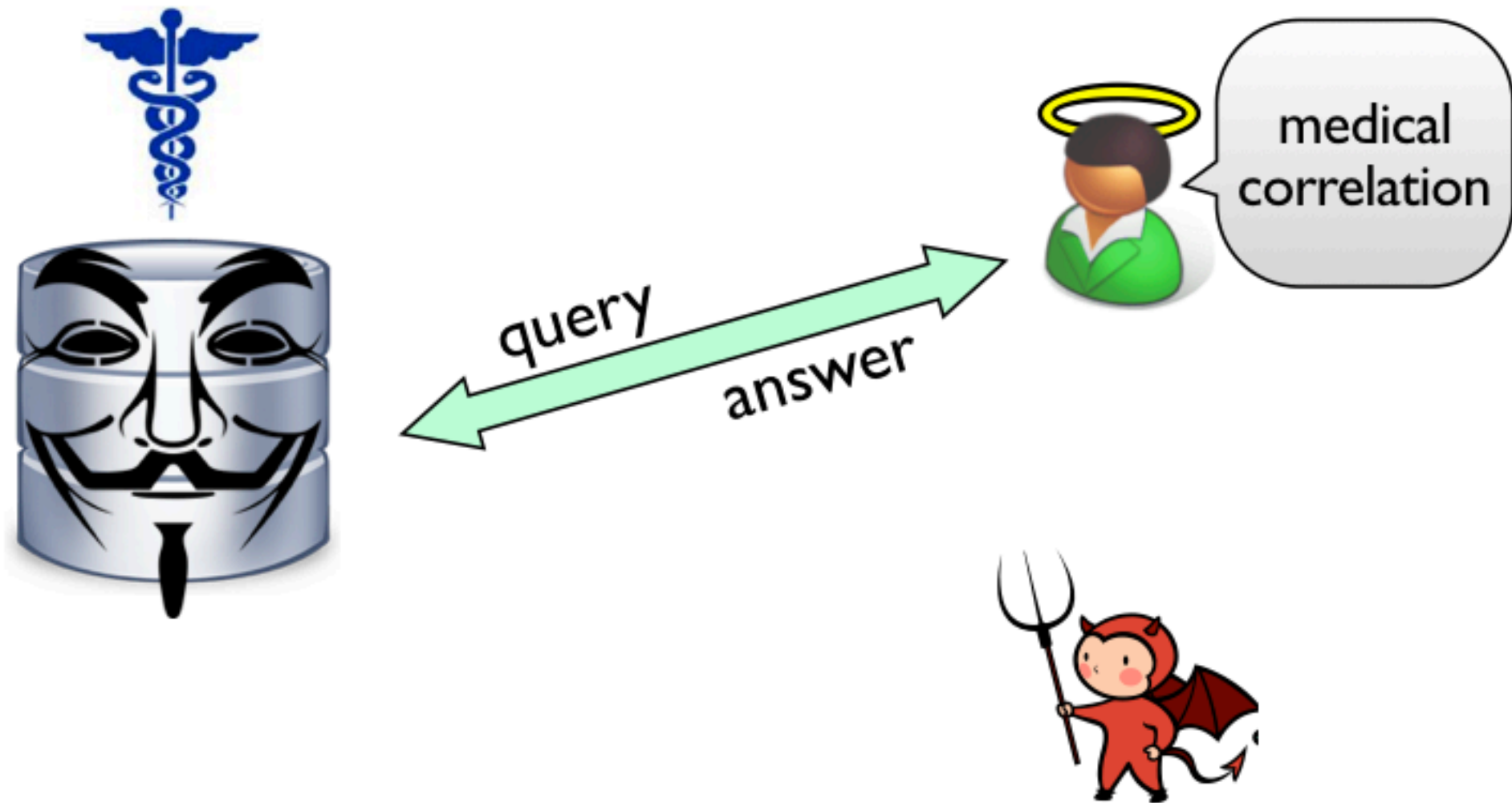
I would like to thank you for your attention to this matter. I am sure that your journal will continue to provide high-quality research in the field of clinical chemistry.

Sincerely,
[Author Name]

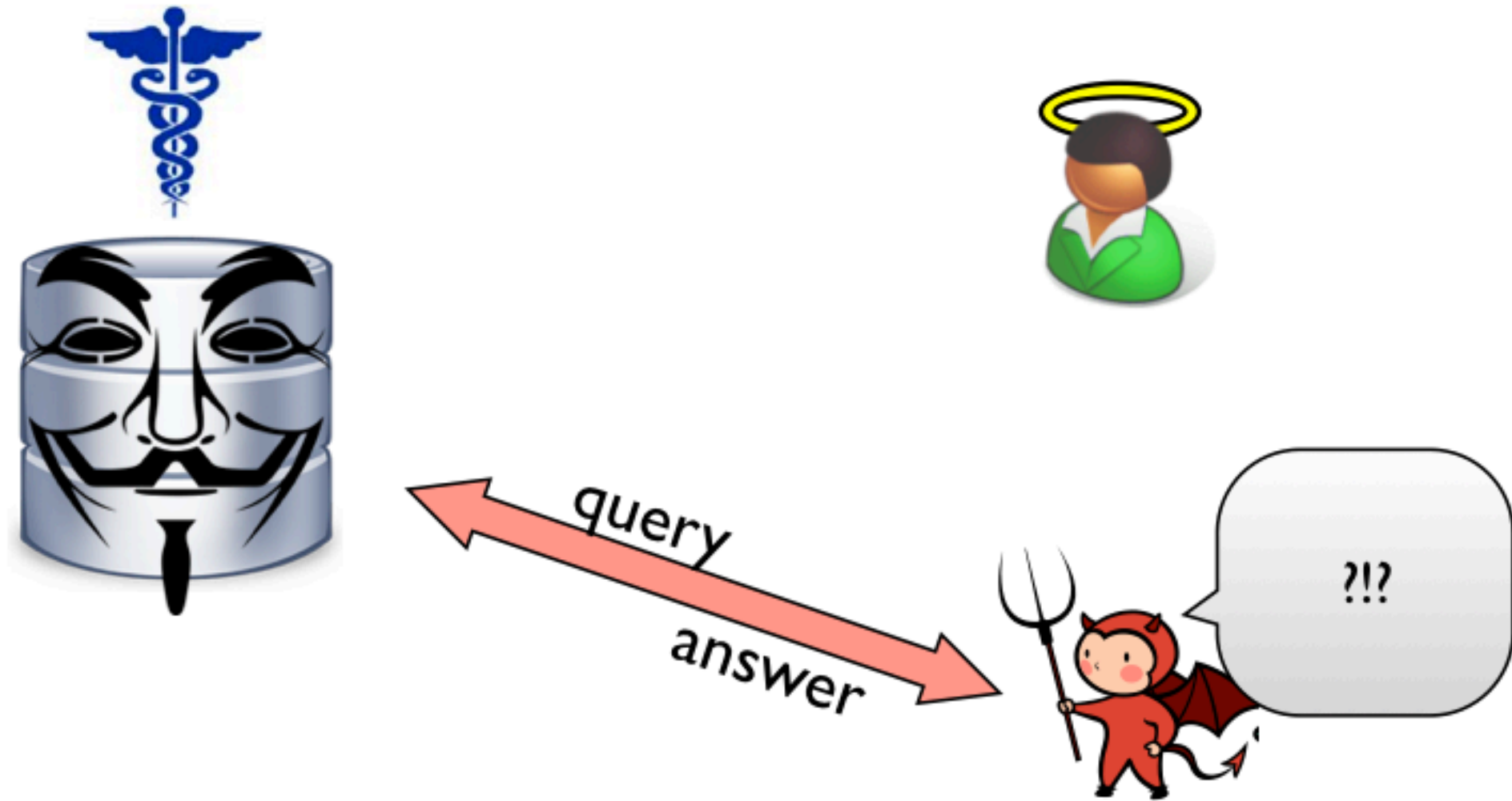
Anonymization?



Anonymization?



Anonymization?

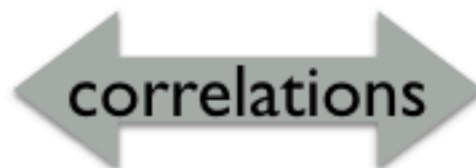


Some stories...went wrong!

Aol.

Attacks on Anonymization

(Narayanan, Shmatikov: Robust De-anonymization of Large Sparse Datasets. IEEE Symposium on Security and Privacy 2008)



Additional Data

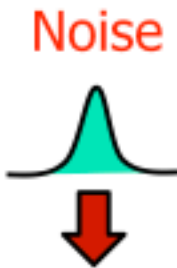


Anonymous Data

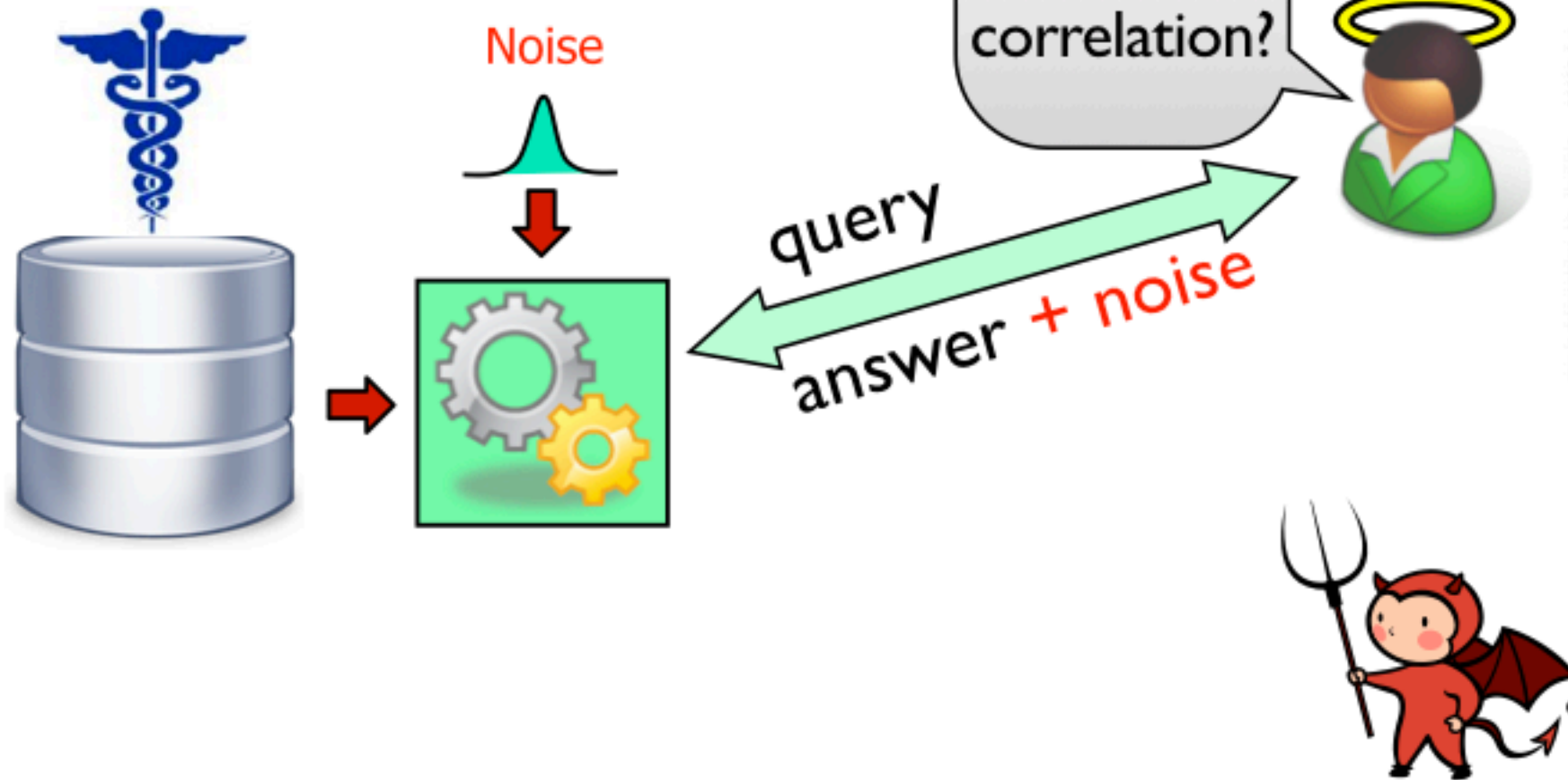


**A Possible Solution:
randomization**

Adding noise

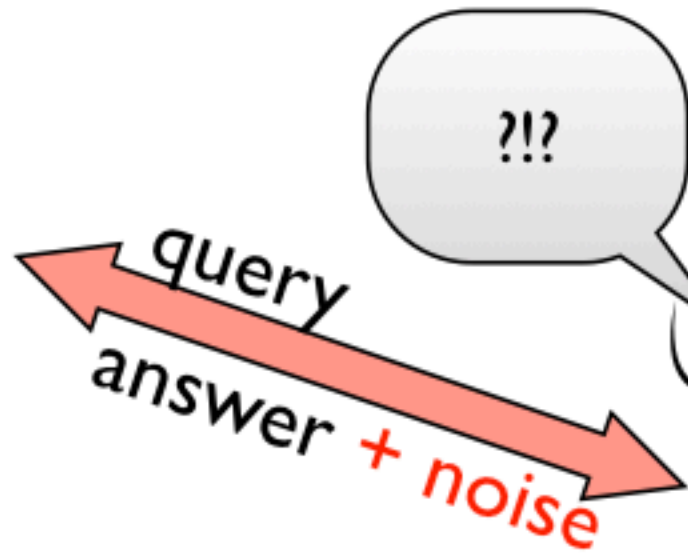
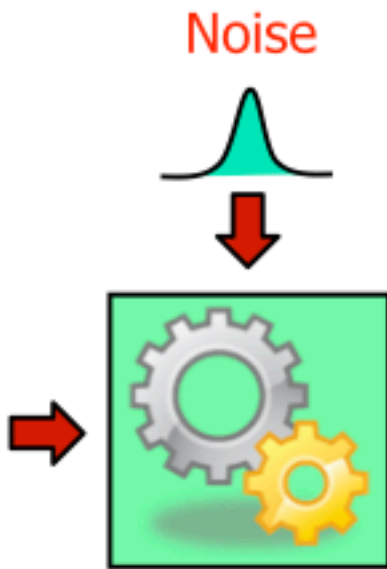


Adding noise



Methodological weakness in using correlation coefficients for assessing the interchangeability of analysis data between samples collected under different sampling conditions - the example of matrix metalloproteinase 2 determined in serum and plasma samples

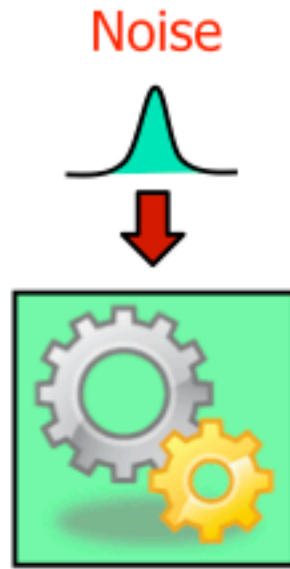
Adding noise



Methodological weakness in using correlation coefficients for assessing the interchangeability of analysis data between samples collected under different sampling conditions - the example of matrix metalloproteinase 2 determined in serum and plasma samples



Adding noise



Methodological weakness in using correlation coefficients for assessing the interchangeability of analysis data between samples collected under different sampling conditions – the example of matrix metalloproteinase 2 determined in serum and plasma samples



Data analyst



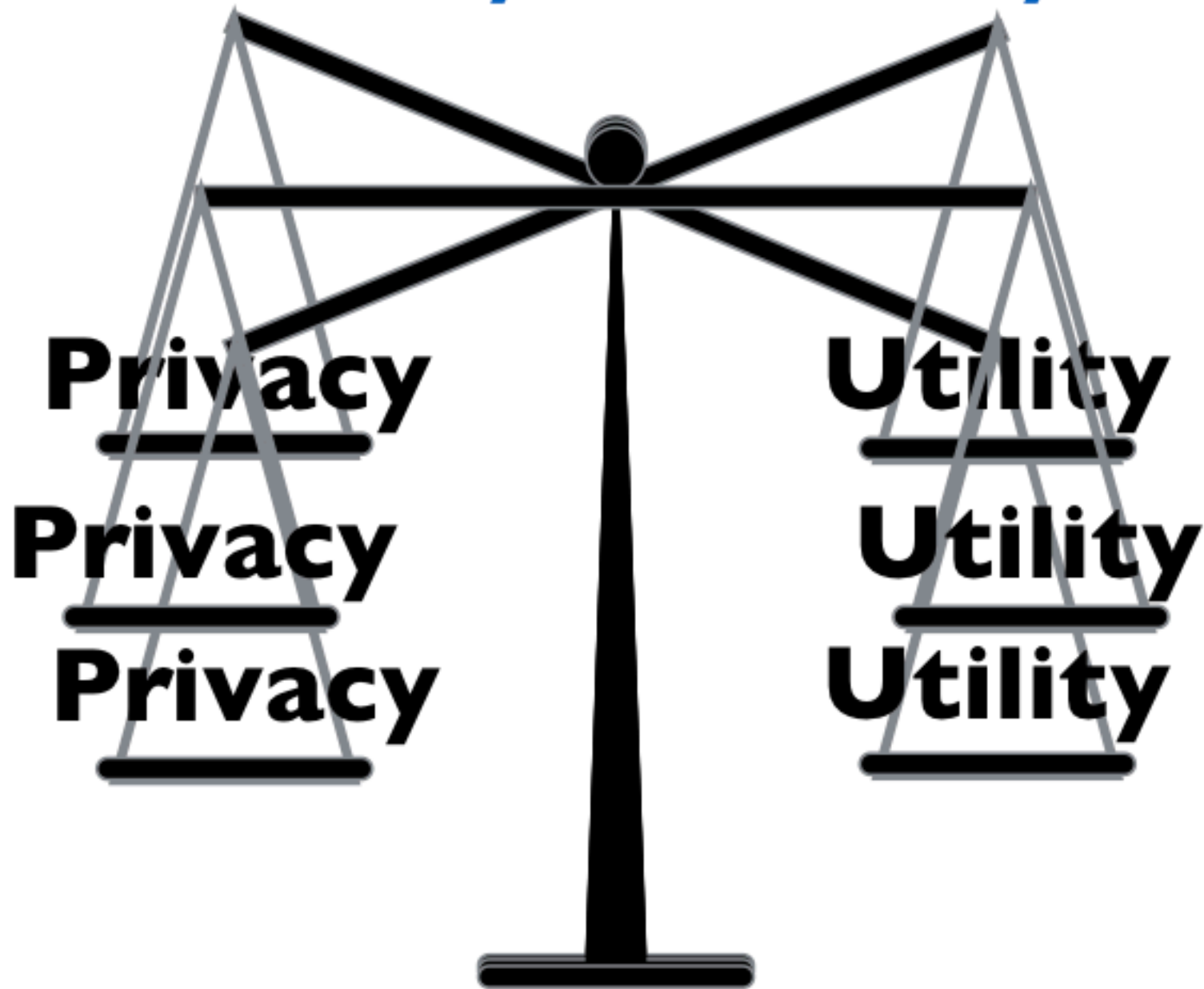
Fundamental Law of Information Reconstruction

The release of **too many** overly **accurate** statistics gives privacy violations.



[DinurNissim02]

Privacy vs Utility



Differential privacy:
understanding the mathematical and
computational meaning of this trade-off.

[DworkMcSherryNissimSmith06]

Some Official Users

- US Census Bureau - onTheMap, new releases in 2020
- Google - RAPPOR tool for Chrome
- Apple - typing statistics reports
- LeapYear - startup
- Uber / Amazon starting
- ...

CSE 660: this class

Syllabus for the course

Location: Davis 338A

Time: MW 10:00 - 11:20

Office Hours: Monday 1:00 - 3:00 or by appointment

Discussion forums: Piazza

class website:

<http://www.acsu.buffalo.edu/~gaboardi/teaching/CSE660-fall17.html>

Course load:

- participation in class and on Piazza,
- completing the assignments,
- working on a project and presenting the results.

Weekly Assignments

- They will consist mostly in programming exercise and theoretical problems,
- The recommended programming language is Python but also R is accepted.
- The assignments will be posted after the Wednesday class and they need to be submitted before the next Wednesday class by email to the instructor with subject “CSE660 Assignment #xx - first names”.
- Each assignment will be graded with a score:
 - 0 - clearly incorrect solution,
 - 1 - solution mostly correct but with some problem,
 - 2 - solution correct.
- Assignments can be submitted by teams of two people.

Academic integrity

Academic integrity is a fundamental university value. Through the honest completion of academic work, students sustain the integrity of the university while facilitating the university's imperative for the transmission of knowledge and culture based upon the generation of new and innovative ideas.

- Reference to the university Graduate Academic Integrity policy and any additional instructor requirements and comments regarding academic dishonesty.
<http://grad.buffalo.edu/Academics/Policies-Procedures/Academic-Integrity.html>
- Any violation of Academic integrity will have consequences according to the Department and University policies.

Final Projects

Projects can take different forms depending on the interest of each student but all the project must have a research component. Some examples of what would constitute a good project are:

- design of a new differentially private algorithm for a specific task
- design or implementation of (part of) a new programming language, system, or tool for differential privacy
- implementation of a previously published work and an experimental comparison with other works,
- investigate new applications of differential privacy

Final Projects

I will provide some specific ideas for possible projects after the first class but other ideas may be accepted if well motivated and discussed with me.

You may work on your project alone or with others. Groups can be composed by at most two students. Each group is invited to meet with me regularly (3-4 times during the term) to check on the advancements and directions of the project.

The deadline for choosing a project is **October 4**.

Final Grading

50% - submitted assignments,
40% - project advancement and
presentation,
10% - engagement and participation in
class and on Piazza

Reference Material

Cynthia Dwork and Aaron Roth,
“The Algorithmic Foundations of Differential Privacy,” 2014
Linked from the class website.

Salil Vadhan,
“The Complexity of Differential Privacy” 2016.
Linked from the class website.

Slides used during class.

Other resources that will be communicated when needed.

Outline of the class

Week 1

Introduction, motivation and privacy limitations. Definition of Differential Privacy and the curator model.

Week 2

Basic mechanisms: Randomized Response, Laplace Mechanism,

Week 3

Basic properties following from the definition, Exponential Mechanism and comparison with the other basic mechanisms.

Week 4

The Report Noisy max algorithm. The Sparse Vector technique.

Week 5

Formalizing privacy proofs using an approximate probabilistic coupling argument.

Week 6

Releasing Many Counting Queries with Correlated Noise. The smallDB algorithm.

Outline of the class

Week 7

The MWEM algorithm. The DualQuery algorithm.

Week 8

Studying the experimental accuracy. Adaptivity and adaptive MWEM.

Week 9

PAC learning and private PAC learning

Week 10

The local model for differential privacy.

Week 11

More algorithms for the local model.

Week 12

Differentially Private Hypothesis Testing

Week 13

Differential Privacy and Generalization in Adaptive Data Analysis

Week 14

Project presentations

<http://www.acsu.buffalo.edu/~gaboardi/teaching/CSE660-fall17.html>

Questions?

Is this data private?

	D1	D2	D3	D4	D5	D6	D6	D7	D8	D9	D10	D11	D12	D13	D14	D15
I1	0	0	0	1	0	0	0	0	1	0	0	1	1	0	0	1
I2	1	0	1	1	1	0	1	0	1	0	1	0	0	1	0	0
I3	0	1	0	1	1	1	0	1	0	0	0	1	0	0	1	0
I4	1	0	1	0	0	1	1	0	1	1	0	0	0	0	1	1
I5	0	0	0	1	1	0	1	1	0	1	0	1	0	1	0	0
I6	0	0	1	1	0	1	1	0	1	1	0	0	1	0	1	0
I7	1	1	0	0	1	0	1	1	1	0	1	0	1	0	0	1
I8	0	0	0	0	1	1	0	0	0	0	0	0	1	0	0	0
I9	0	1	0	0	1	0	1	1	0	1	1	1	0	1	1	0
I10	1	0	1	0	0	1	1	0	0	0	0	0	0	1	0	1
I11	0	1	0	1	1	0	0	1	0	1	0	1	0	1	1	0
I12	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
I13	1	1	1	0	1	1	1	1	0	0	1	0	1	0	1	0
I14	0	1	1	0	0	0	0	1	0	0	0	1	0	0	1	0
I15	0	1	0	1	0	1	1	0	1	0	1	0	1	0	0	1

How about if we also have this data?

	ID	Name
1	I1	Alice
2	I2	Bob
3	I3	Cynthia
4	I4	Dan
5	I5	Eve
6	I6	Frank
7	I7	Guy
8	I8	Hannah
9	I9	Ivan
10	I10	Jon
11	I11	Ken
12	I12	Lou
13	I13	Mike
14	I14	Noa
15	I15	Omer

	ID	Disease
1	D1	AMAN
2	D2	Behcet
3	D3	Celiac
4	D4	Dermatitis
5	D5	Evans synd.
6	D6	Fibrosis
7	D7	Graves' dis.
8	D8	Henoch-Schonlein
9	D9	IGA Neph.
10	D10	Juv. Diabetes
11	D11	Kawasaki dis.
12	D12	Lichen planus
13	D13	Myositis
14	D14	Narcolepsy
15	D15	Optic Neuritis

How about if we also have this data?

	D2	D3	D5	D6	D8	D10	D12	D14	D15
Alice	0	1	1	0	1	1	0	0	0
Cynthia	1	0	1	1	0	0	0	1	0
Eve	0	0	1	0	0	0	0	0	0
Frank	0	1	0	1	1	0	1	1	0
Guy	1	0	1	0	1	1	1	0	1
Ivan	1	0	1	0	0	1	0	1	0
Jon	0	1	0	1	0	0	0	0	1
Lou	0	0	0	0	0	1	0	0	0
Omer	1	0	0	1	1	1	1	0	1

Attacks on Anonymization

(Narayanan, Shmatikov: Robust De-anonymization of Large Sparse Datasets. IEEE Symposium on Security and Privacy 2008)



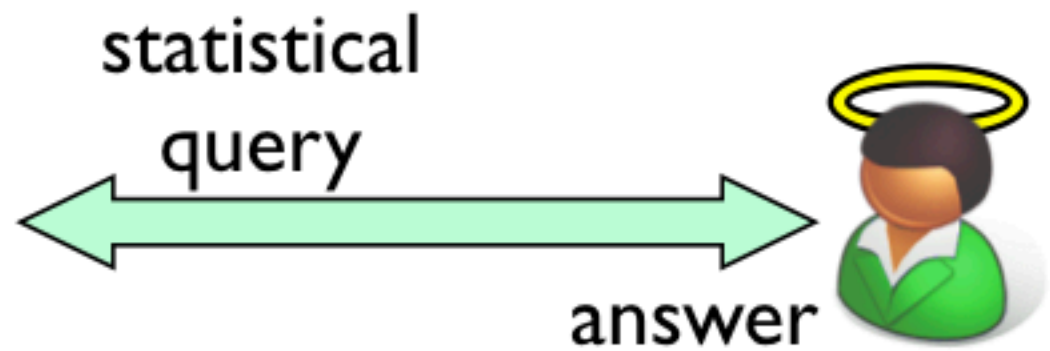
Additional Data



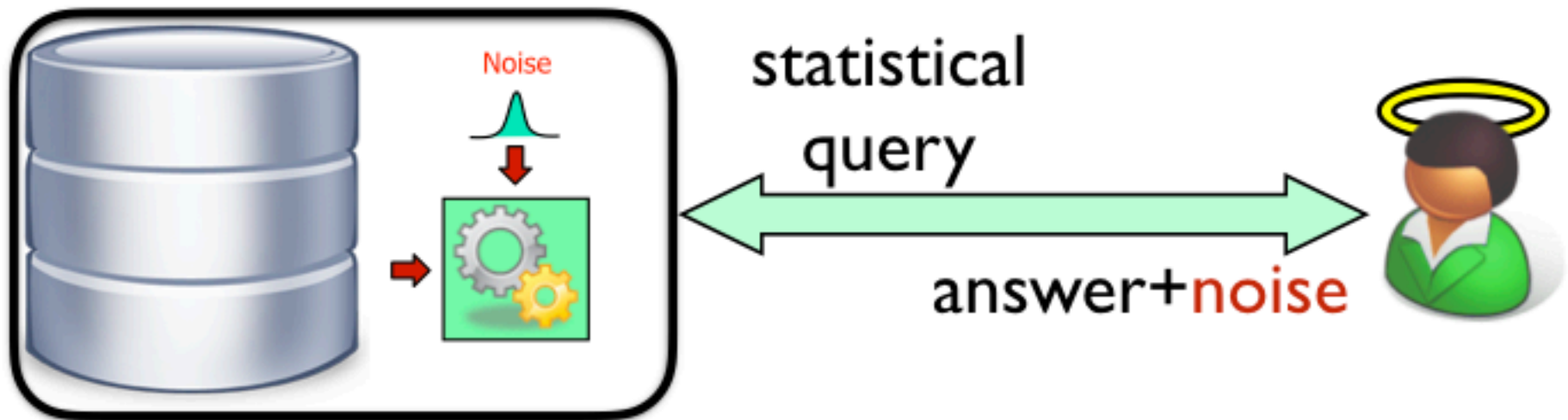
Anonymous Data



Statistical database



Private Statistical database



Database

- We will not be interested in implementation details of a database (I will sometimes use the word dataset),
- For us a database will just be a single table



	D1	D2	D3	D4
I1	0	0	0	1
I2	1	0	1	1
I3	0	1	0	1
I4	1	0	1	0
I5	0	0	0	1
I6	0	0	1	1
I7	1	1	0	0
I8	0	0	0	0
I9	0	1	0	0
I10	1	0	1	0
I11	0	1	0	1

Database

- We can think about a database as a list of records from some universe set:

$$D \in \mathcal{X}^n$$

- Sometimes we will think to them as functions

$$D(k) \in \mathcal{X}$$

- and sometimes we will write elements explicitly

$$(d_1, \dots, d_n) \in \mathcal{X}^n$$

Counting Queries

- A **counting query** $q : \mathcal{X}^n \rightarrow [0, 1]$ is a function counting the fraction of people in a dataset satisfying the **predicate** $q : \mathcal{X} \rightarrow \{0, 1\}$
- In symbols:

$$q(D) = \frac{1}{n} \sum_{i=1}^n q(d_i)$$

- Notice that we take a normalized count, which also corresponds to the average.

Example 1

Let's consider an arbitrary universe domain \mathcal{X} and let's consider the following predicate for $y \in \mathcal{X}$

$$q_y(x) = \begin{cases} 1 & \text{if } y = x \\ 0 & \text{otherwise} \end{cases}$$

we call a **point function** the associated counting query

$$q_y : \mathcal{X}^n \rightarrow [0, 1]$$

Question: Suppose that we answer all the point function queries for $y \in \mathcal{X}$. What well know statistics do we obtain?

Example 1

$$X = \{0, 1\}^3$$

$$D \in X^{10} =$$

	D1	D2	D3
I1	0	0	0
I2	1	0	1
I3	0	1	0
I4	1	0	1
I5	0	0	0
I6	0	0	1
I7	1	1	0
I8	0	0	0
I9	0	1	0
I10	1	0	1

$$q_{000}(D) = .3$$

$$q_{001}(D) = .1$$

$$q_{010}(D) = .2$$

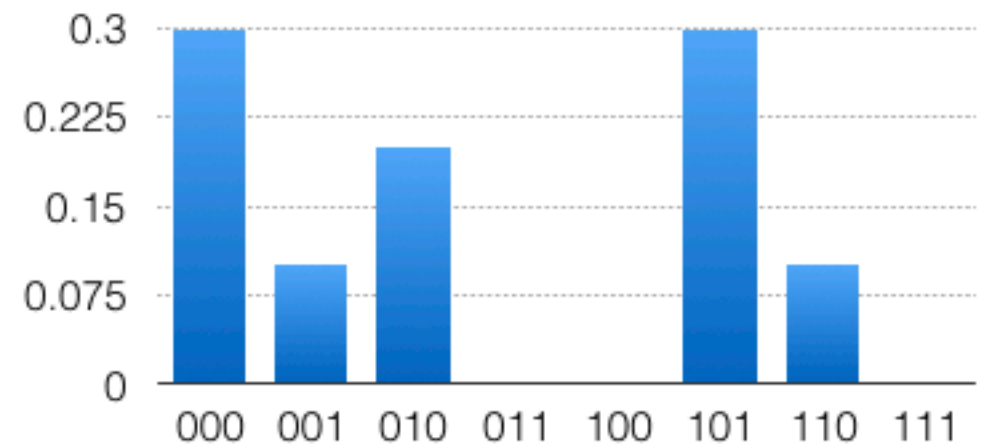
$$q_{011}(D) = 0$$

$$q_{100}(D) = 0$$

$$q_{101}(D) = .3$$

$$q_{110}(D) = .1$$

$$q_{111}(D) = 0$$



Example 1

Question: Suppose that we answer all the point function queries for $y \in \mathcal{X}$. What well known statistics do we obtain?

Answer: Histogram of the universe and of the database.

Example II

Let's consider an arbitrary **ordered** universe domain \mathcal{X} and let's consider the following predicate for $y \in \mathcal{X}$

$$q_y(x) = \begin{cases} 1 & \text{if } x \leq y \\ 0 & \text{otherwise} \end{cases}$$

we call a **threshold function** the associated counting query

$$q_y : \mathcal{X}^n \rightarrow [0, 1]$$

Question: Suppose that we answer all the threshold function queries for $y \in \mathcal{X}$. What well know statistics do we obtain?

Example II

$X = \{0, 1\}^3$
with order
given by the
corresponding
binary encoding.

$D \in X^{10} =$

	D1	D2	D3
I1	0	0	0
I2	1	0	1
I3	0	1	0
I4	1	0	1
I5	0	0	0
I6	0	0	1
I7	1	1	0
I8	0	0	0
I9	0	1	0
I10	1	0	1

$$q_{000}(D) = .3$$

$$q_{001}(D) = .4$$

$$q_{010}(D) = .6$$

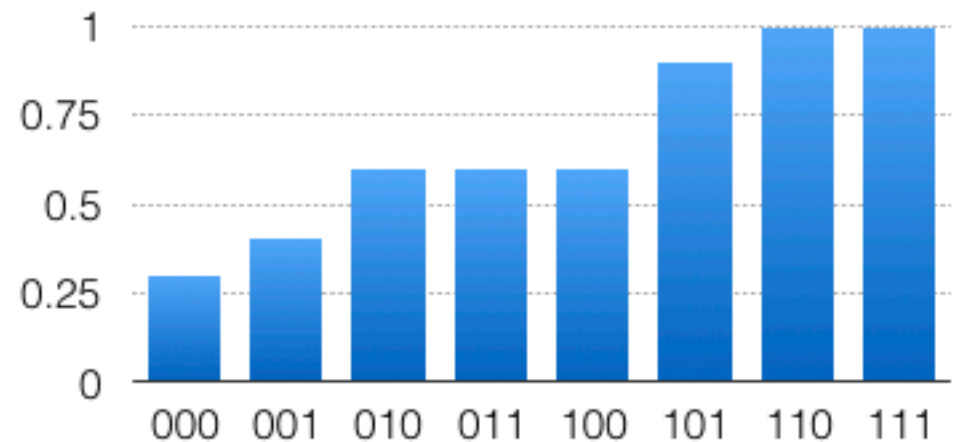
$$q_{011}(D) = .6$$

$$q_{100}(D) = .6$$

$$q_{101}(D) = .9$$

$$q_{110}(D) = 1$$

$$q_{111}(D) = 1$$



Example II

Question: Suppose that we answer all the threshold function queries for $y \in \mathcal{X}$. What well known statistics do we obtain?

Answer: CDF of the universe and of the database.

Example III

Let's consider the universe domain $\mathcal{X} = \{0, 1\}^d$ and let's consider the following predicate for an index $1 \leq j \leq d$

$$q_j(x) = x_j$$

we call an **attribute mean function** the associated counting query

$$q_j : \mathcal{X}^n \rightarrow [0, 1]$$

Question: Which statistics does correspond to releasing all the attribute mean functions?

Example III

$$X = \{0, 1\}^3$$

$$D \in X^{10} =$$

	D1	D2	D3
I1	0	0	0
I2	1	0	1
I3	0	1	0
I4	1	0	1
I5	0	0	0
I6	0	0	1
I7	1	1	0
I8	0	0	0
I9	0	1	0
I10	1	0	1
margin	4	3	4

$$q_1(D) = .4$$

$$q_2(D) = .3$$

$$q_3(D) = .4$$

Example III

Question: Which statistics does correspond to releasing all the attribute mean functions?

Answer: (1-way) Marginals of the distribution

Example IV

Let's consider the universe domain $\mathcal{X} = \{0, 1\}^d$ and let's consider $\vec{v} \in \{1, \bar{1}, \dots, d, \bar{d}\}^k$ with $1 \leq k \leq d$ and

$$q_{\vec{v}}(x) = q_{v_1}(x) \wedge q_{v_2}(x) \wedge \dots \wedge q_{v_k}(x)$$

where $q_j(x) = x_j$ and $q_{\bar{j}}(x) = \neg x_j$

We call a **conjunction** or k-way marginal the associated counting query

$$q_{\vec{v}} : \mathcal{X}^n \rightarrow [0, 1]$$

Question: Which statistics does correspond to releasing conjunctions?

Example IV

$$X = \{0, 1\}^3$$

$$D \in X^{10} =$$

	D1	D2	D3
I1	0	0	0
I2	1	0	1
I3	0	1	0
I4	1	0	1
I5	0	0	0
I6	0	0	1
I7	1	1	0
I8	0	0	0
I9	0	1	0
I10	1	0	1

$$k=2$$

$$q_{12}(D) = .1$$

$$q_{1/2}(D) = .3$$

$$q_{13}(D) = .3$$

$$q_{1/3}(D) = .1$$

$$q_{/12}(D) = .2$$

$$q_{/13}(D) = .1$$

$$q_{/1/2}(D) = .4$$

$$q_{/1/3}(D) = .5$$

	D1	/D1
D2	0.1	0.2
/D2	0.3	0.4

Example IV

Question: Which statistics does correspond to releasing conjunctions?

Answer: contingency tables

Linear Queries

- A **linear query** $q : \mathcal{X}^n \rightarrow [0, 1]$ is a function averaging the value of a function $q : \mathcal{X} \rightarrow [0, 1]$ over the elements of the dataset.
- In symbols:

$$q(D) = \frac{1}{n} \sum_{i=1}^n q(d_i)$$

Example 1

Let's consider the domain $\mathcal{X} = \{0, 1\}^d$ and let's consider the following predicate for $y \in \mathcal{X}$

$$q_y(x) = \begin{cases} .5 * y_1 & \text{if } y = x \\ 0 & \text{otherwise} \end{cases}$$

Example 1

$$X = \{0, 1\}^3$$

$$D \in X^{10} =$$

	D1	D2	D3
I1	0	0	0
I2	1	0	1
I3	0	1	0
I4	1	0	1
I5	0	0	0
I6	0	0	1
I7	1	1	0
I8	0	0	0
I9	0	1	0
I10	1	0	1

$$q_{000}(D) = 0$$

$$q_{100}(D) = 0$$

$$q_{001}(D) = 0$$

$$q_{101}(D) = .15$$

$$q_{010}(D) = 0$$

$$q_{110}(D) = .05$$

$$q_{011}(D) = 0$$

$$q_{111}(D) = 0$$

Sum queries

- Let's denote by $I \subseteq [n]$ a subset I of $\{0, \dots, n\}$

- A **sum query** $q_I : 0, 1^k \rightarrow \mathbb{N}^k$ is defined as

$$q_I(D) = \sum_{i \in I} d_i$$

Example

$$X = \{0, 1\}^3 \quad D \in X^{10} =$$

	D1	D2	D3
I1	0	0	0
I2	1	0	1
I3	0	1	0
I4	1	0	1
I5	0	0	0
I6	0	0	1
I7	1	1	0
I8	0	0	0
I9	0	1	0
I10	1	0	1

$$q_{\{1,2,3\}}(D) = (1, 1, 1)$$

$$q_{\{1,2,4\}}(D) = (2, 0, 2)$$

$$q_{\{5,8\}}(D) = (0, 0, 0)$$

$$q_{\{2,4,7,10\}}(D) = (4, 1, 3)$$

Question: Is releasing the result of counting queries private?