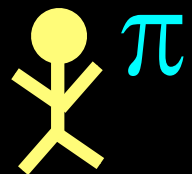


# Differential Privacy Tutorial

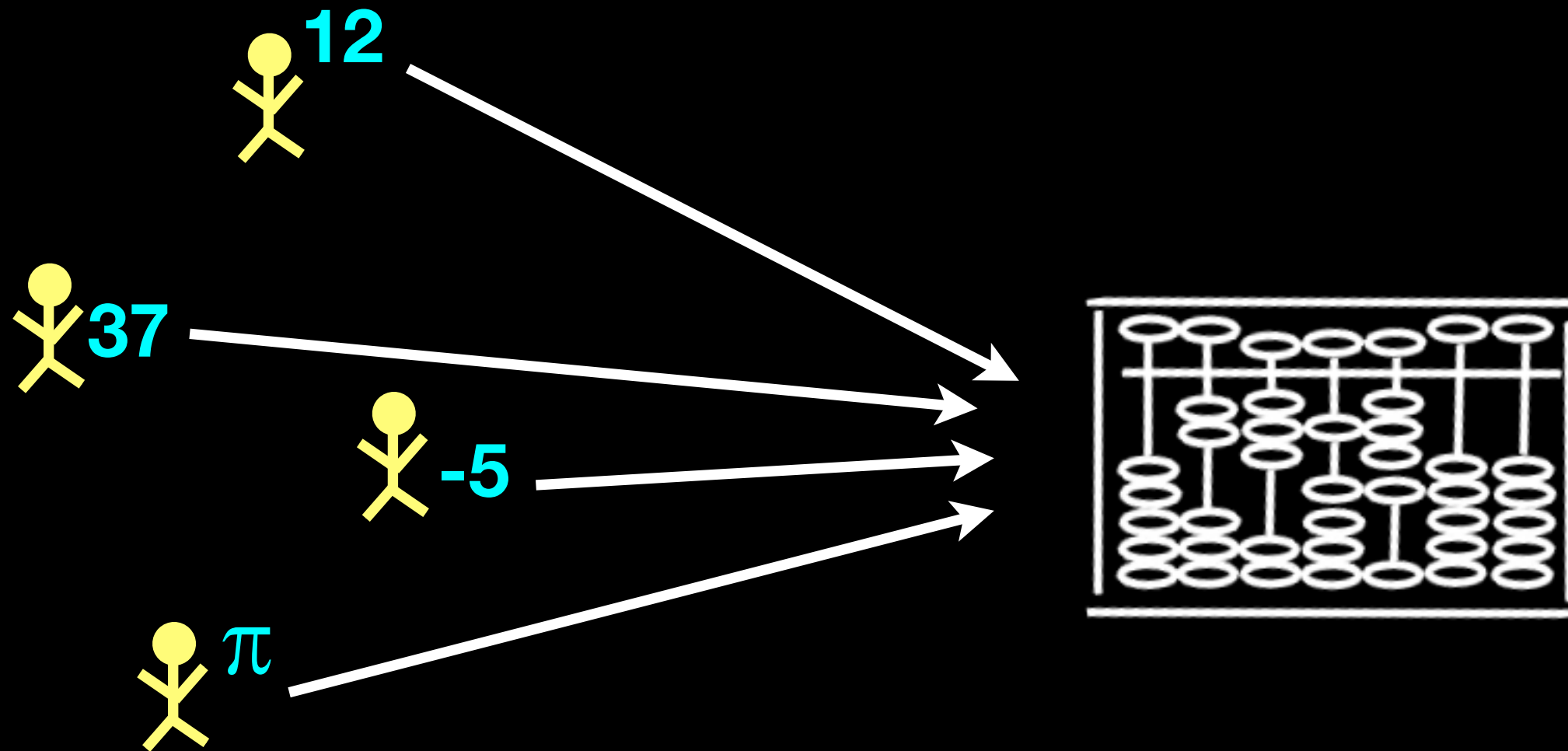
Simons Institute Workshop on Privacy and Big Data

Katrina Ligett  
Caltech

individuals have lots of interesting data...



individuals have lots of interesting data...



...we want to compute on it













Form

**1040**

Department of the Treasury—Internal Revenue Service

**U.S. Individual Income Tax Return**

**2008**







Form

**1040**

Department of the Treasury—Internal Revenue Service

**U.S. Individual Income Tax Return**

**2008**



**PatientCare Portal**





Form

**1040**

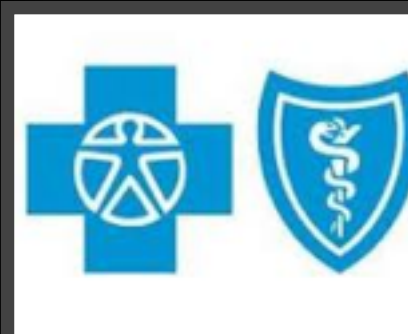
Department of the Treasury—Internal Revenue Service

**U.S. Individual Income Tax Return**

**2008**



**PatientCare Portal**







Form

**1040**

Department of the Treasury—Internal Revenue Service

**U.S. Individual Income Tax Return**

**2008**



**PatientCare Portal**





# Google™

data privacy|

data privacy day	75,400,000 results
data privacy laws	17,600,000 results
data privacy act	11,100,000 results
data privacy policy	60,400,000 results
data privacy safe harbor	332,000 results
data privacy breaches	1,320,000 results
data privacy legislation	980,000 results
data privacy audit	684,000 results
data privacy through optimal k-anonymization	4,200 results
data privacy laws us	71,900,000 results

[close](#)

Form **1040** Department of the Treasury—Internal Revenue Service **2008**  
**U.S. Individual Income Tax Return**







data privacy|

data privacy day	75,400,000 results
data privacy laws	17,600,000 results
data privacy act	11,100,000 results
data privacy policy	60,400,000 results
data privacy safe harbor	332,000 results
data privacy breaches	1,320,000 results
data privacy legislation	980,000 results
data privacy audit	684,000 results
data privacy through optimal k-anonymization	4,200 results
data privacy laws us	71,900,000 results

[close](#)

Form **1040** Department of the Treasury—Internal Revenue Service **2008**  
**U.S. Individual Income Tax Return**



- finding statistical correlations
  - genotype/phenotype associations
  - correlating medical outcomes with risk factors or events
- publishing aggregate statistics
- noticing events/outliers
  - intrusion detection
  - disease outbreaks
- datamining/learning tasks
  - use customer data to update strategies



**Web**

**Images**

**Video**

**News**

**Local**

**Shopping**

**more »**

aol search debacle|

**Search**



**Web**

Images

Video

News

Local

Shopping

more »

aol search debacle|

**Search**

### [Stats: Who's to blame for AOL's search debacle?](#)

Let the fingerpointing begin A friend of ousted **AOL** advertising executive Mike Kelly takes issue with our assignment of blame.

[gawker.com/302054/whos-to-blame-for-aols-search-debacle](http://gawker.com/302054/whos-to-blame-for-aols-search-debacle) - [Similar pages](#)

### [AOL Proudly Releases Massive Amounts of Private Data](#)

Yet Another Update: **AOL**: This was a screw up Further Update: Sometime after 7 pm the download link went down as well, but ...

[www.techcrunch.com/2006/08/06/aol-proudly-releas...](http://www.techcrunch.com/2006/08/06/aol-proudly-releas...) - 123k - [Similar pages](#)

### [AOL search data scandal - Wikipedia, the free encyclopedia](#)

The **AOL search** data scandal was the result of a research project by **AOL**. .... **AOL** apologizes for release of user **search** data | CNET News.com; ^ **AOL search** ...

[en.wikipedia.org/wiki/AOL\\_search\\_data\\_scandal](http://en.wikipedia.org/wiki/AOL_search_data_scandal) - 45k - [Similar pages](#)

Web Images Video News Local Shopping more »

aol search debacle

Search

No “personally identifiable information” was released

Katrina Ligett	data privacy
Katrina Ligett	aol search debacle
Katrina Ligett	Ligett DBLP
Katrina Ligett	computer science news
Katrina Ligett	Caltech rankings
Katrina Ligett	weather Pasadena
Jane Smith	youtube
Jane Smith	free tv download
Jane Smith	streaming tv
Chris Jones	childrens books
Chris Jones	dr seuss
Chris Jones	“the cat and the hat”
Chris Jones	gifts for children

aol search debacle

Search

No “personally identifiable information” was released

John Doe	tax forms
user195023	data privacy
user195023	aol search debacle
user195023	Ligett DBLP
user195023	computer science news
user195023	Caltech rankings
user195023	weather Pasadena
Jane Smith	youtube
Jane Smith	free tv download
Jane Smith	streaming tv
Chris Jones	childrens books
Chris Jones	dr seuss
Chris Jones	“the cat and the hat”
Chris Jones	gifts for children
Chris Jones	

Web Images Video News Local Shopping more »

aol search debacle

Search

No “personally identifiable information” was released

John Doe	tax forms
user195023	data privacy
user195023	aol search debacle
user195023	Ligett DBLP
user195023	computer science news
user195023	Caltech rankings
user195023	weather Pasadena
Jane Smith	youtube
Jane Smith	free tv download
Jane Smith	streaming tv
Chris Jones	childrens books
Chris Jones	dr seuss
Chris Jones	“the cat and the hat”
Chris Jones	gifts for children
Chris Jones	

Web Images Video News Local Shopping more »

aol search debacle

Search

No “personally identifiable information” was released

user195023	data privacy
user195023	aol search debacle
user195023	Ligett DBLP
user195023	computer science news
user195023	Caltech rankings
user195023	weather Pasadena

Ad hoc solutions are risky!

youtube  
free tv download  
streaming tv  
childrens books  
dr seuss  
“the cat and the hat”  
gifts for children



Web Images Video News Local Shopping more »

aol search debacle

Search

No “personally identifiable information” was released

user195023	data privacy
user195023	aol search debacle
user195023	Ligett DBLP
user195023	computer science news
user195023	Caltech rankings
user195023	weather Pasadena

Ad hoc solutions are risky!

Huge opportunity for formalism.

This doesn't apply to me! I don't want to publish the whole dataset!

individuals hold data...

...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

individuals hold data...

...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

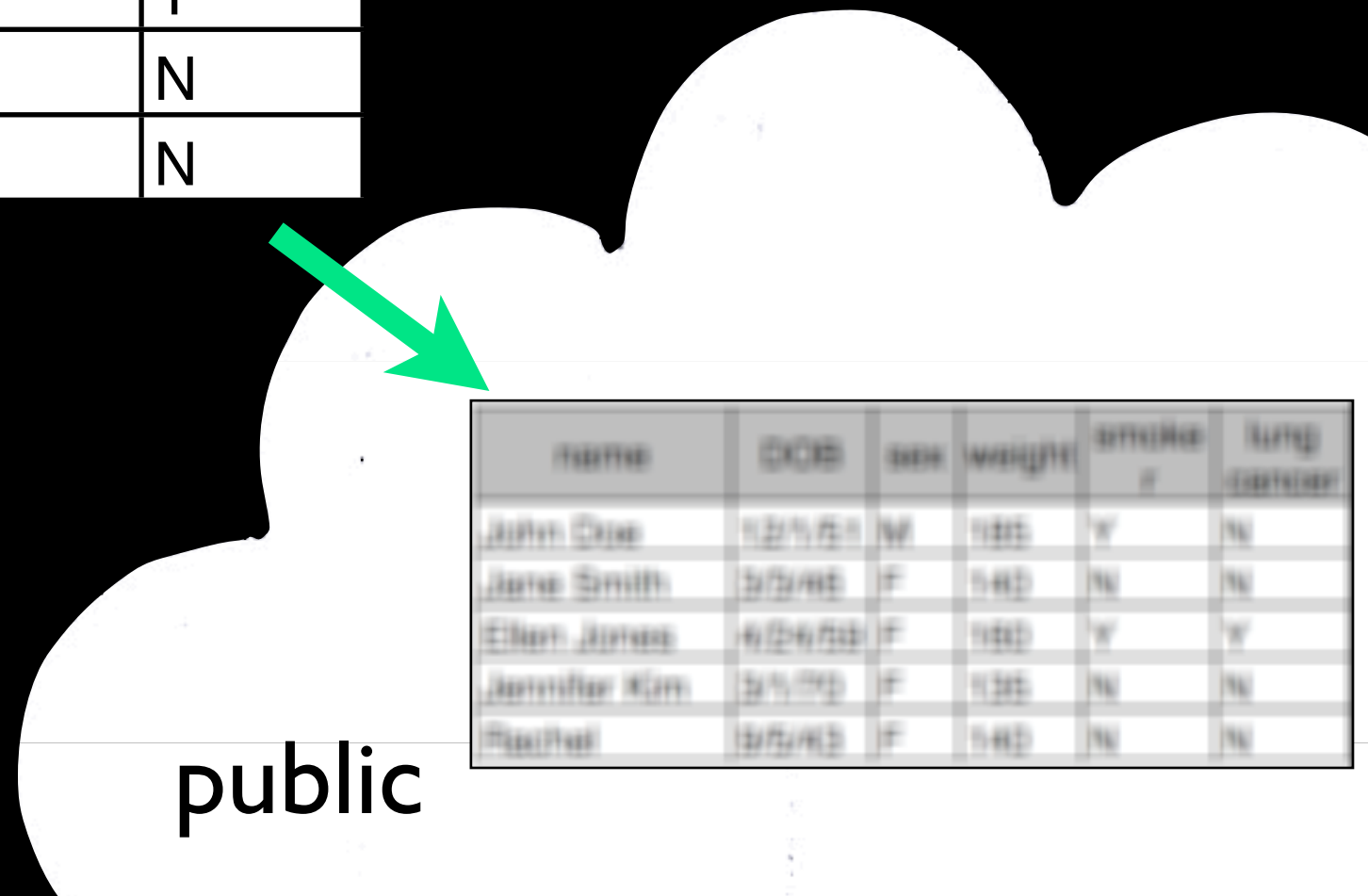


public

# individuals hold data...

## ...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel	9/5/43	F	140	N	N

individuals hold data...

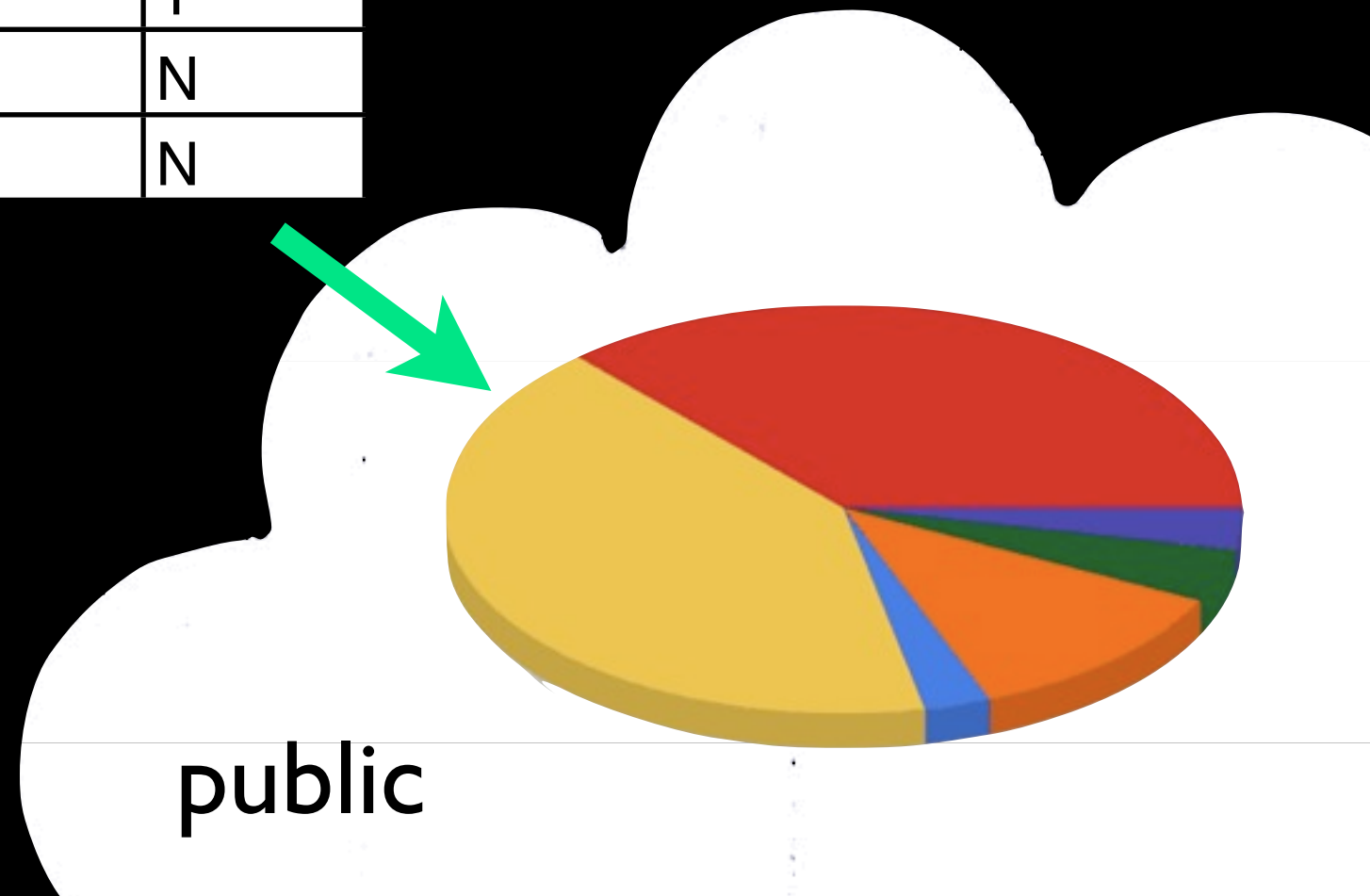
...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

individuals hold data...

...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



individuals hold data...

...what if it's sensitive?

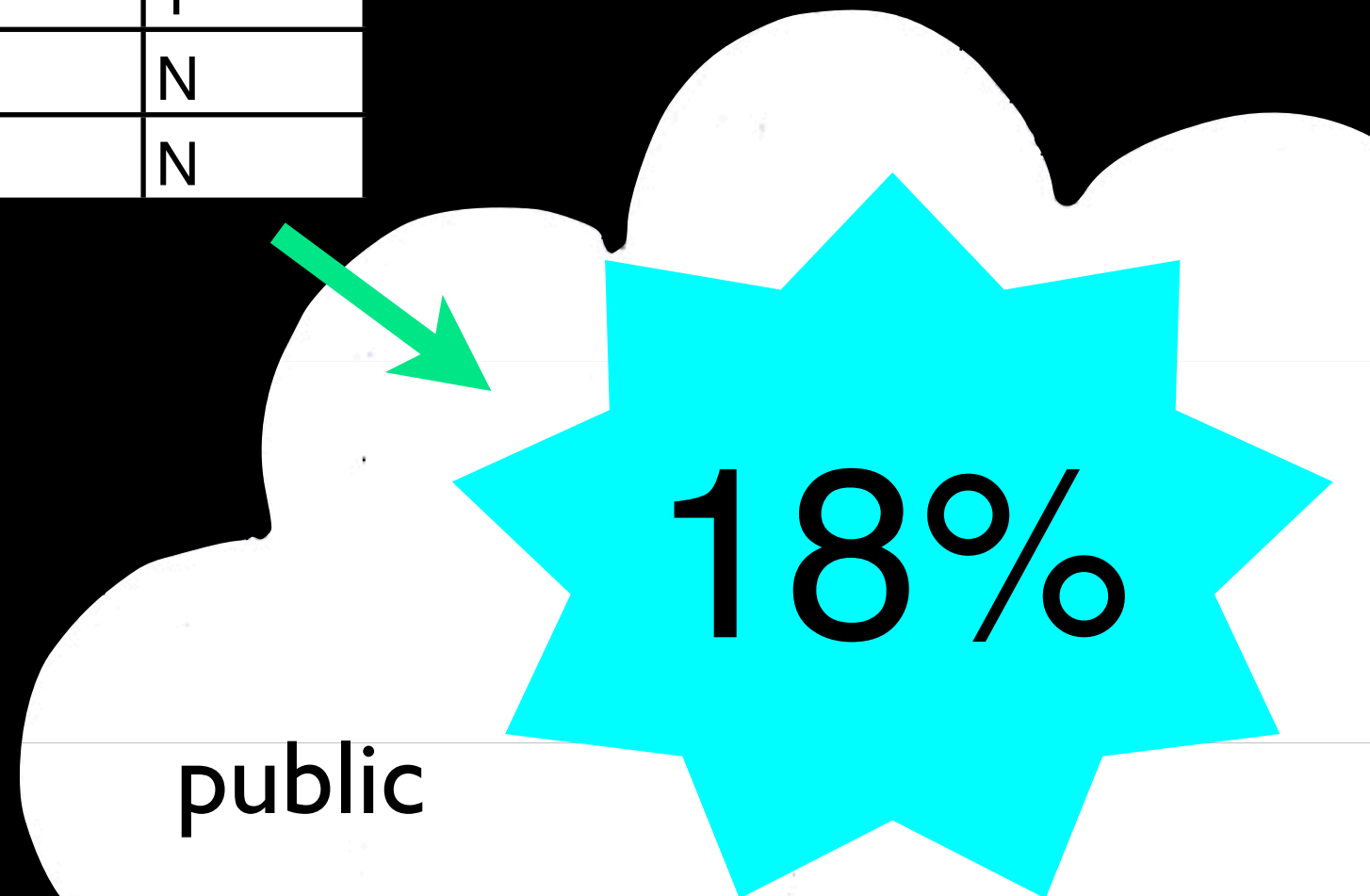
name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



individuals hold data...

...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



individuals hold data...

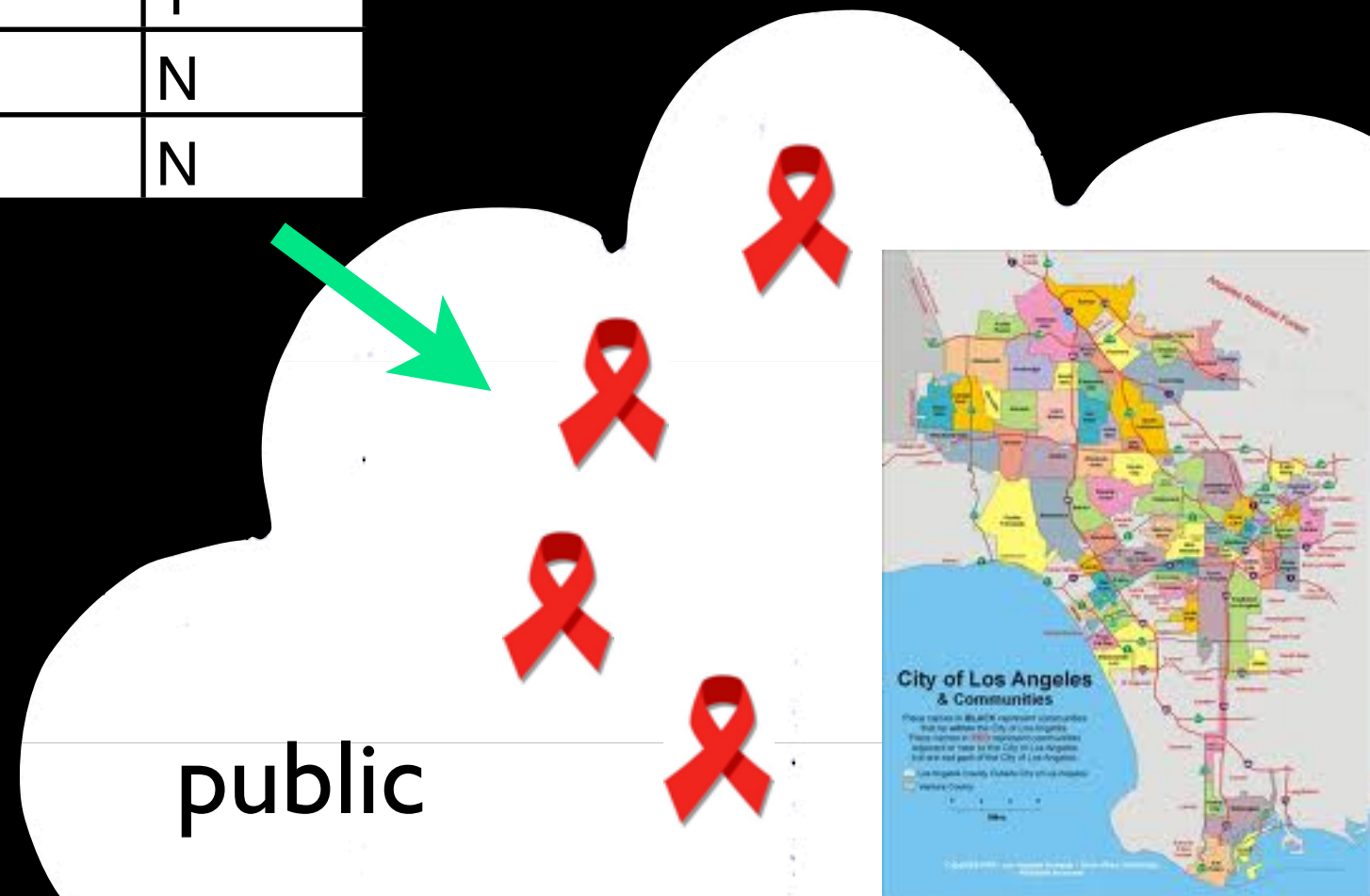
...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

# individuals hold data...

# ...what if it's sensitive?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



This doesn't apply to me! I don't want to publish the whole dataset!

This doesn't apply to me! I don't want to  
publish the whole dataset!

not so fast...

This doesn't apply to me! I don't want to publish the whole dataset!

not so fast...

see, e.g., Korolova 2011's Facebook microtargeting attack

The Facebook logo, consisting of the word "facebook" in white lowercase letters on a blue rectangular background.

facebook

This doesn't apply to me! I don't want to publish the whole dataset!

not so fast...

see, e.g., Korolova 2011's Facebook microtargeting attack

... must pay attention to *all* uses of sensitive data

The Facebook logo, consisting of the word "facebook" in white lowercase letters on a blue rectangular background.

facebook

what to promise about output?



what to promise about output?

access to the output should  
not enable one to learn  
anything about an individual  
that could not be learned  
without access

what to promise about output?

access to the output should  
not enable one to learn  
anything about an individual  
that could not be learned  
without access

is this  
possible?

what to promise about output?

access to the output should  
not enable one to learn  
anything about an individual  
that could not be learned  
without access

is this  
possible?

hint:  
*either* privacy or utility  
separately is easy

# what if wanted to do a study about smoking and cancer?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

# what if wanted to do a study about smoking and cancer?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



public

# what if wanted to do a study about smoking and cancer?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

public

there is a correlation of xxx

# what if wanted to do a study about smoking and cancer?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith			140	N	N
Ellen				Y	Y
					N
					N

but what if someone knew Alice is a smoker?

public

there is a correlation of xxx

what to promise about output?

access to the output should  
not enable one to learn  
anything about an individual  
that could not be learned  
without access



what to promise about output?

access to the output should  
not enable one to learn  
anything about an individual  
that could not be learned  
without access

not possible!

# what to promise about output?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

# what to promise about output?

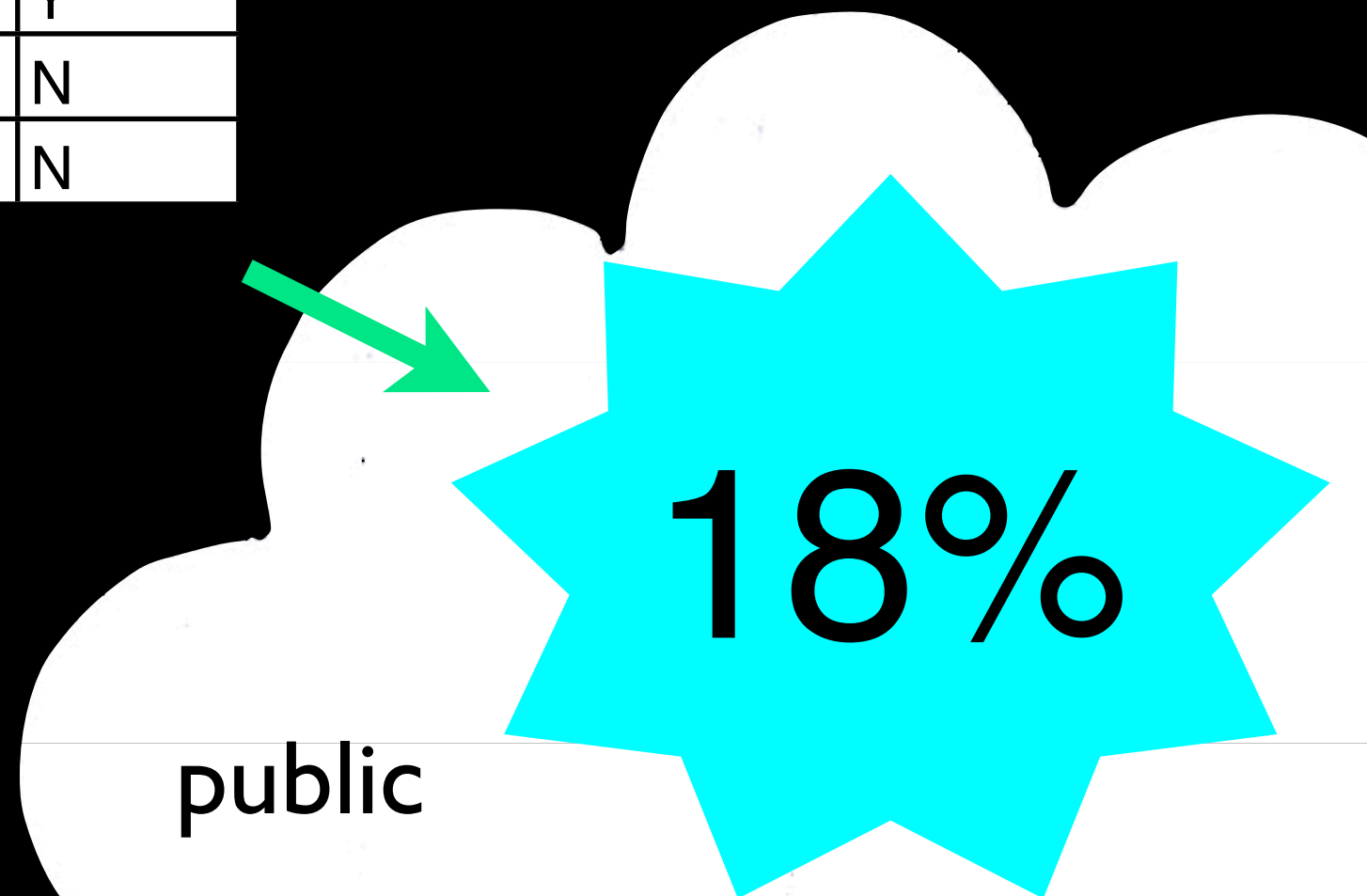
## think of output as randomized

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

# what to promise about output?

## think of output as randomized

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



# what to promise about output?

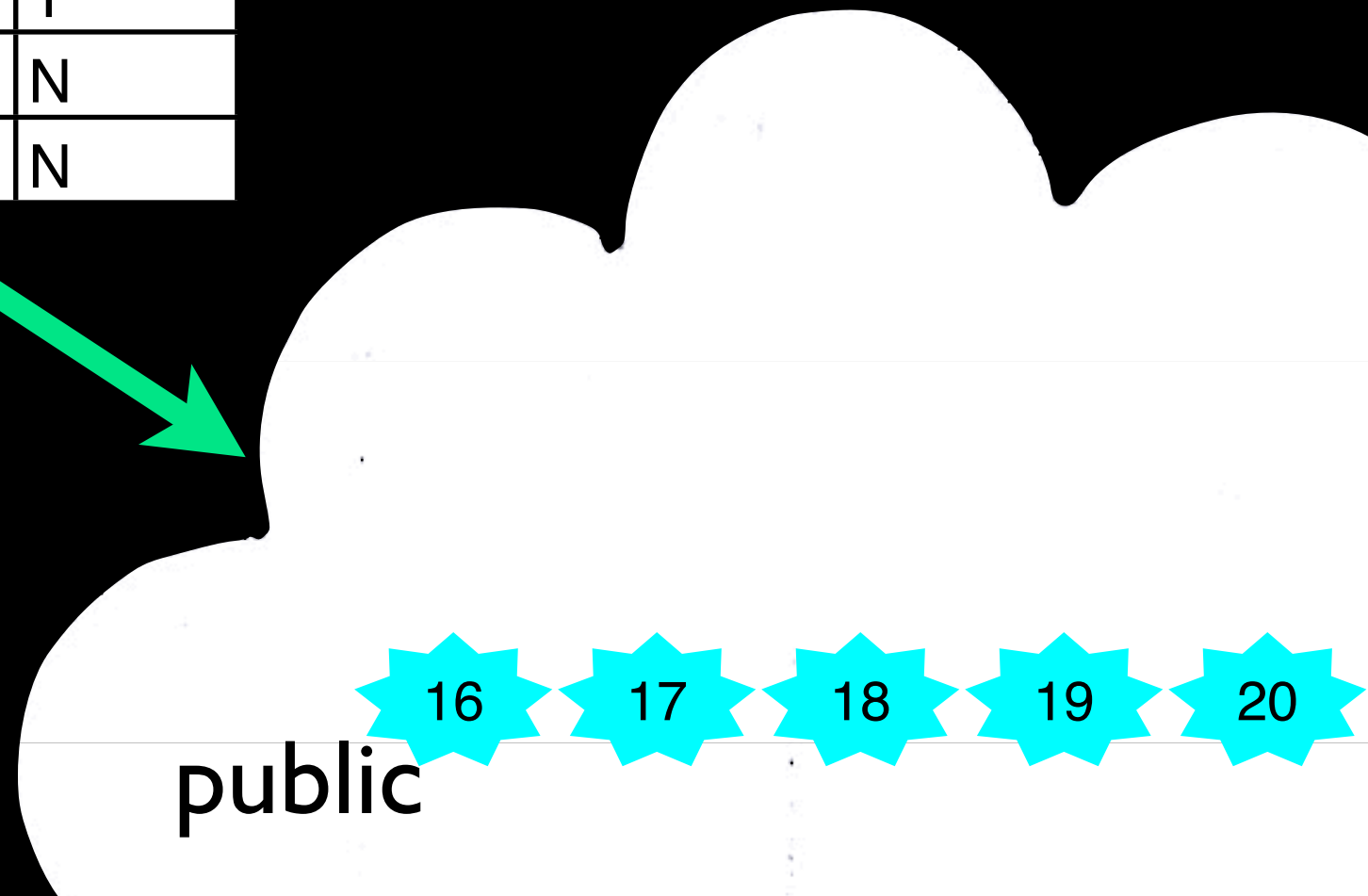
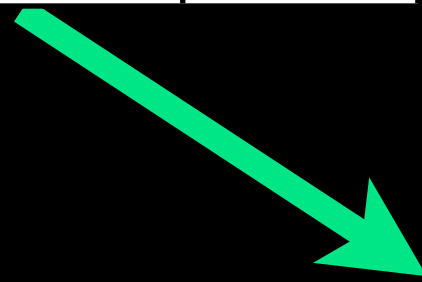
## think of output as randomized

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

# what to promise about output?

## think of output as randomized

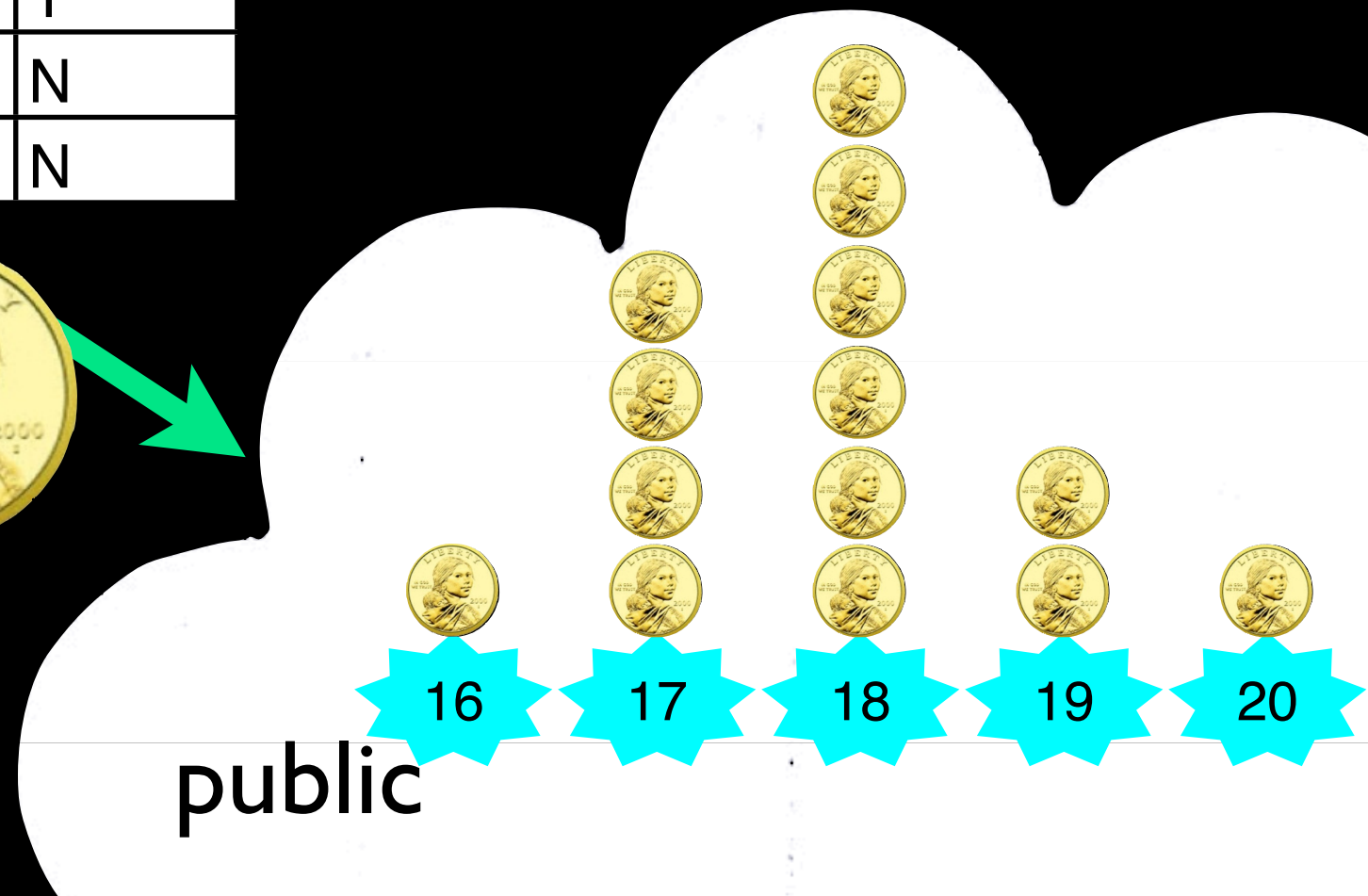
name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



# what to promise about output?

## think of output as randomized

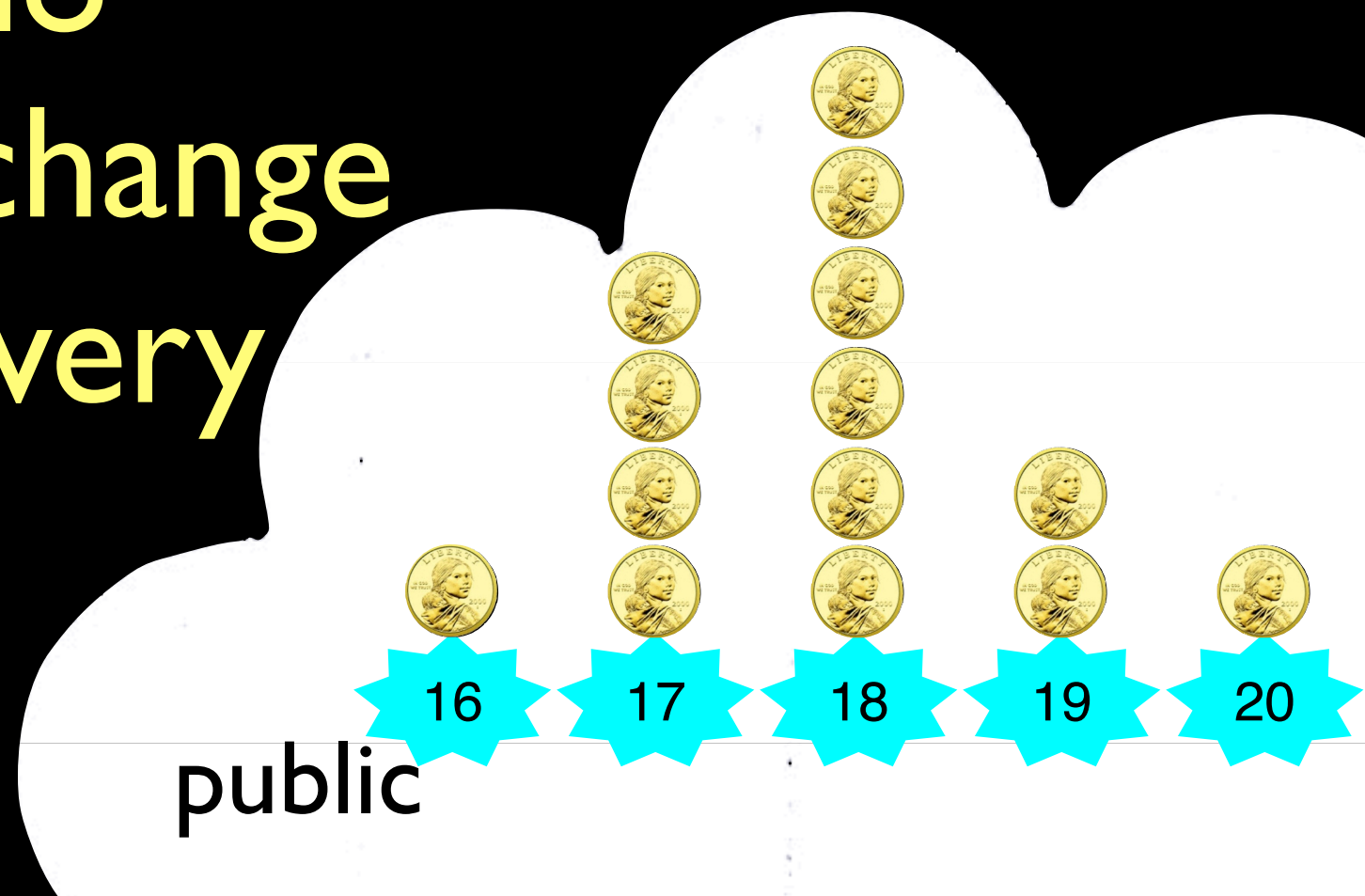
name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



what to promise about output?

think of output as randomized

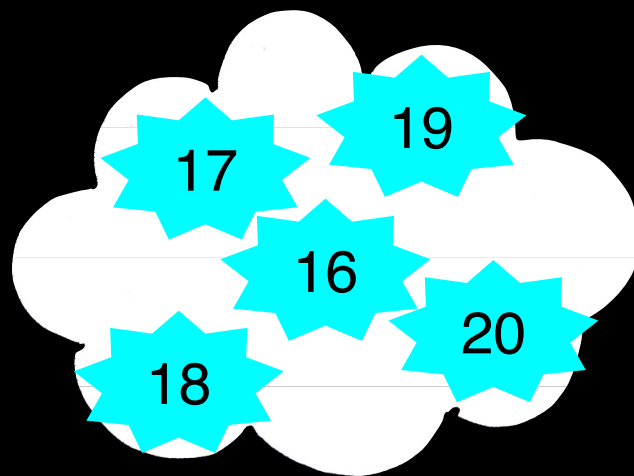
promise: if you leave  
the database, no  
outcome will change  
probability by very  
much





# more formally...

- database **D** a set of rows, one per person
- sanitizing algorithm **M** probabilistically maps **D** to event or object in outcome space



A white cloud containing three overlapping tables representing a database **D**.

name	DOB	sex	weight	smoke	lung
John Doe	12/15/1980	M	185	N	N
Jane Smith	3/24/85	F	140	N	N
John Doe	12/15/1980	M	185	Y	Y
Jane Smith	3/24/85	F	140	N	N
John Doe	12/15/1980	M	185	Y	N
Jane Smith	3/24/85	F	140	N	N

name	DOB	sex	weight	smoke	lung
John Doe	12/15/1980	M	185	Y	N
Jane Smith	3/24/85	F	140	N	N
Ellen Jones	4/24/55	F	160	Y	Y
Jennifer Kim	3/11/72	F	130	N	N
Rachel	9/5/92	F	140	N	N

name	DOB	sex	weight	smoke	lung
John Doe	12/15/1980	M	185	Y	N
Jane Smith	3/24/85	F	140	N	N
Ellen Jones	4/24/55	F	160	Y	Y
Jennifer Kim	3/11/72	F	130	N	N
Rachel	9/5/92	F	140	N	N



# differential privacy

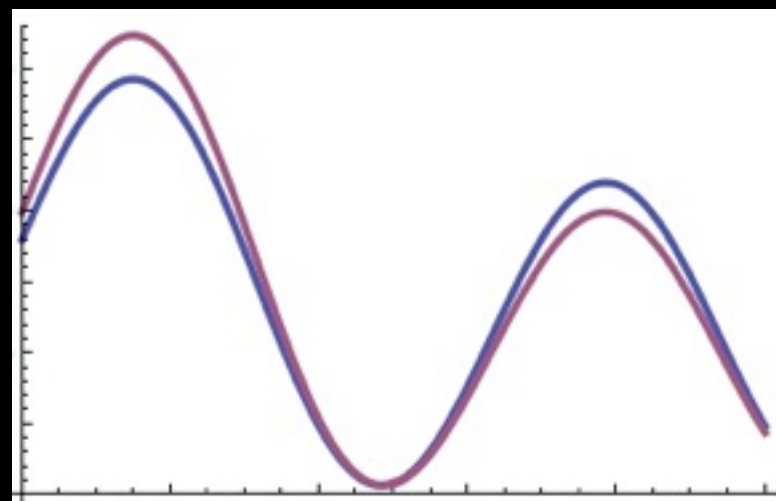
[DinurNissim03, DworkNissimMcSherrySmith06]

$\epsilon$ -Differential Privacy for mechanism  $M$ :

for any two neighboring data sets  $D_1, D_2$ ,

any  $C \in \text{range}(M)$ ,

$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$



# differential privacy

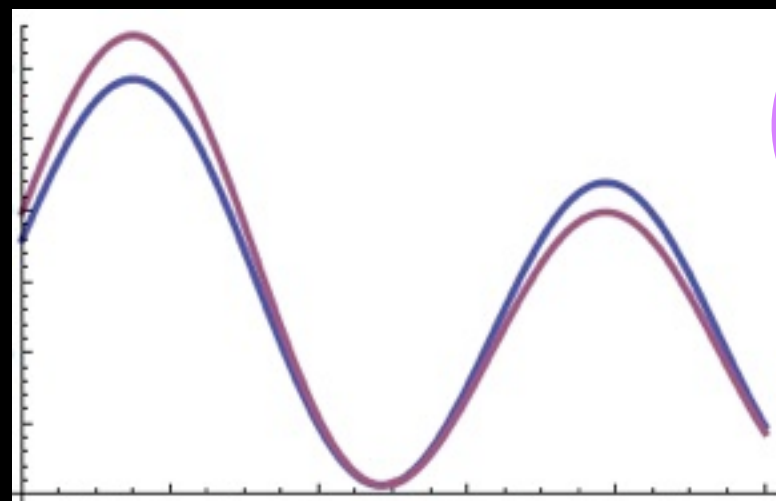
[DinurNissim03, DworkNissimMcSherrySmith06]

$\epsilon$ -Differential Privacy for mechanism  $M$ :

for any two neighboring data sets  $D_1, D_2$ ,

any  $C \in \text{range}(M)$ ,

$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$

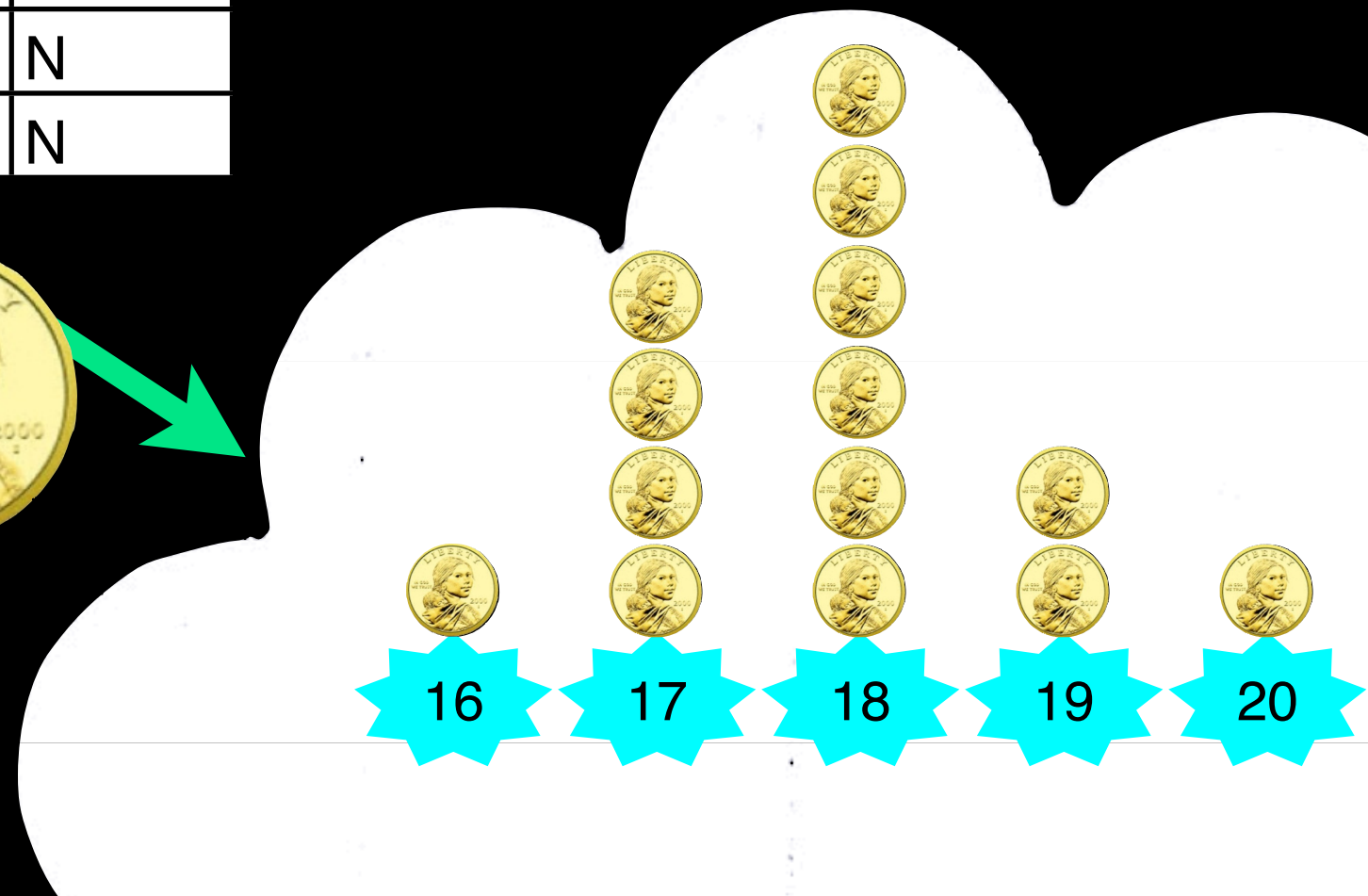


$$e^\epsilon \sim (1 + \epsilon)$$

# differential privacy

$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$

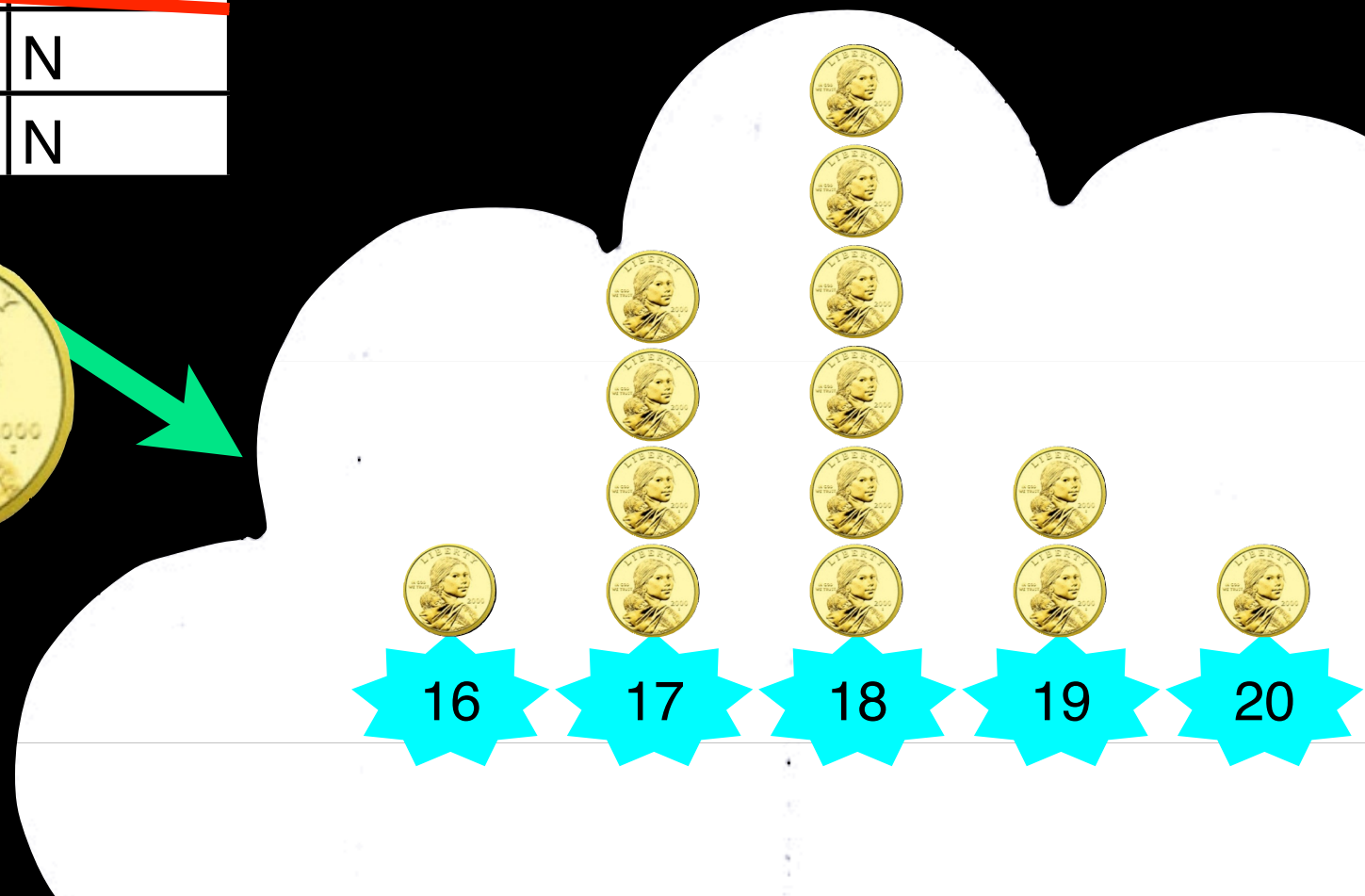
name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



# differential privacy

$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$

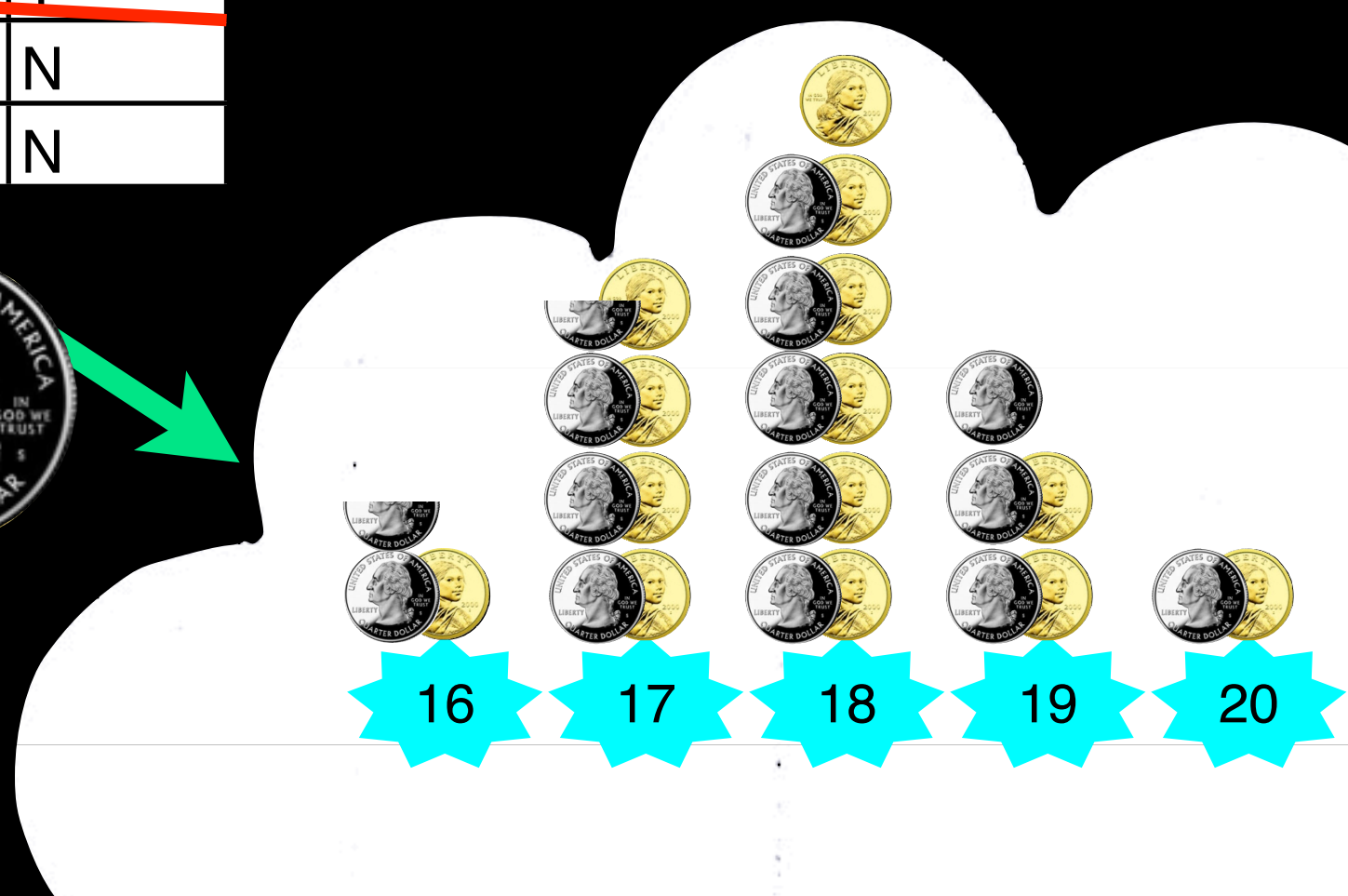
name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
<del>Ellen Jones</del>	<del>4/24/59</del>	<del>F</del>	<del>160</del>	<del>Y</del>	<del>Y</del>
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



# differential privacy

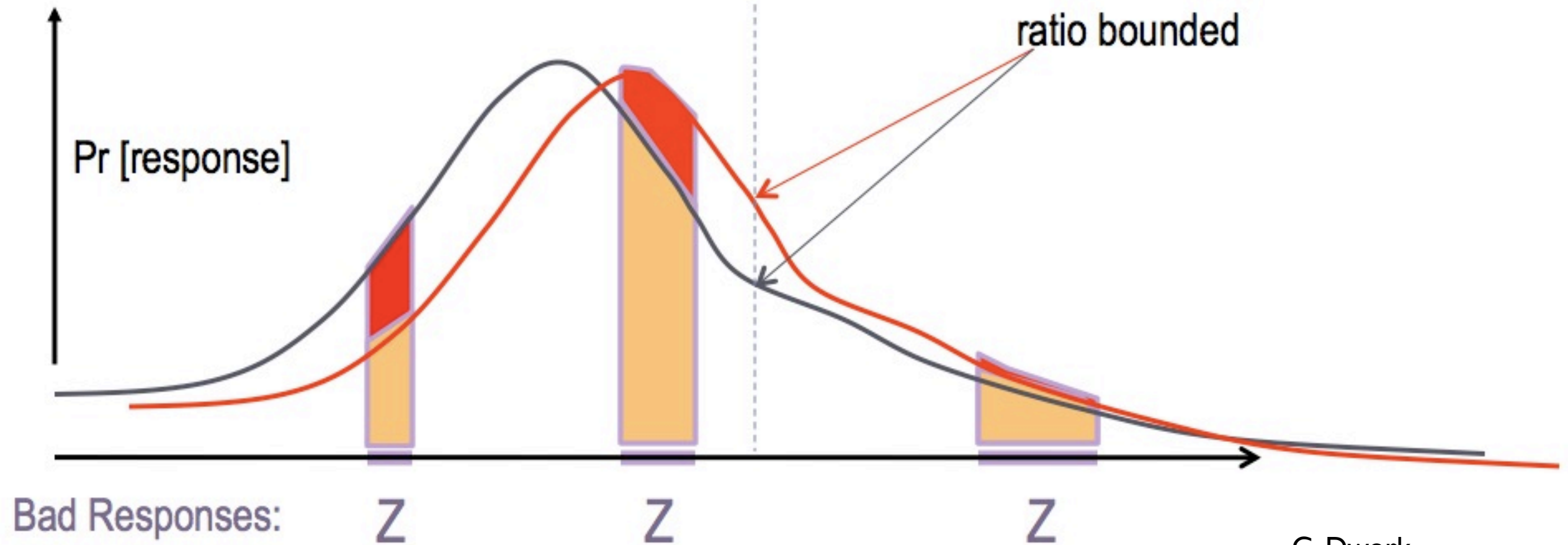
$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
<del>Ellen Jones</del>	<del>4/24/59</del>	<del>F</del>	<del>160</del>	<del>Y</del>	<del>Y</del>
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



# differential privacy

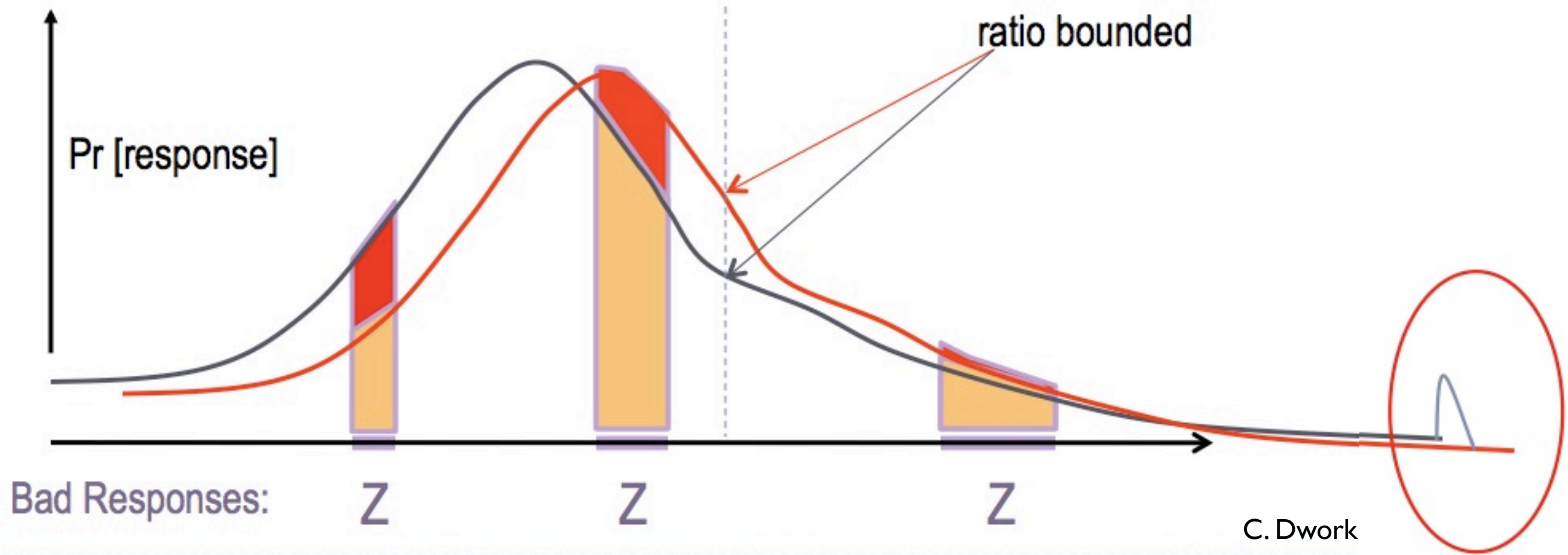
$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$





# $(\epsilon, \delta)$ -differential privacy

$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C] + \delta$$





# differential privacy

$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$

Is a statistical property of **mechanism** behavior

- unaffected by auxiliary information
- independent of adversary's computational power

# differential privacy

$$\Pr[M(D_1) \in C] \leq e^\epsilon \Pr[M(D_2) \in C]$$

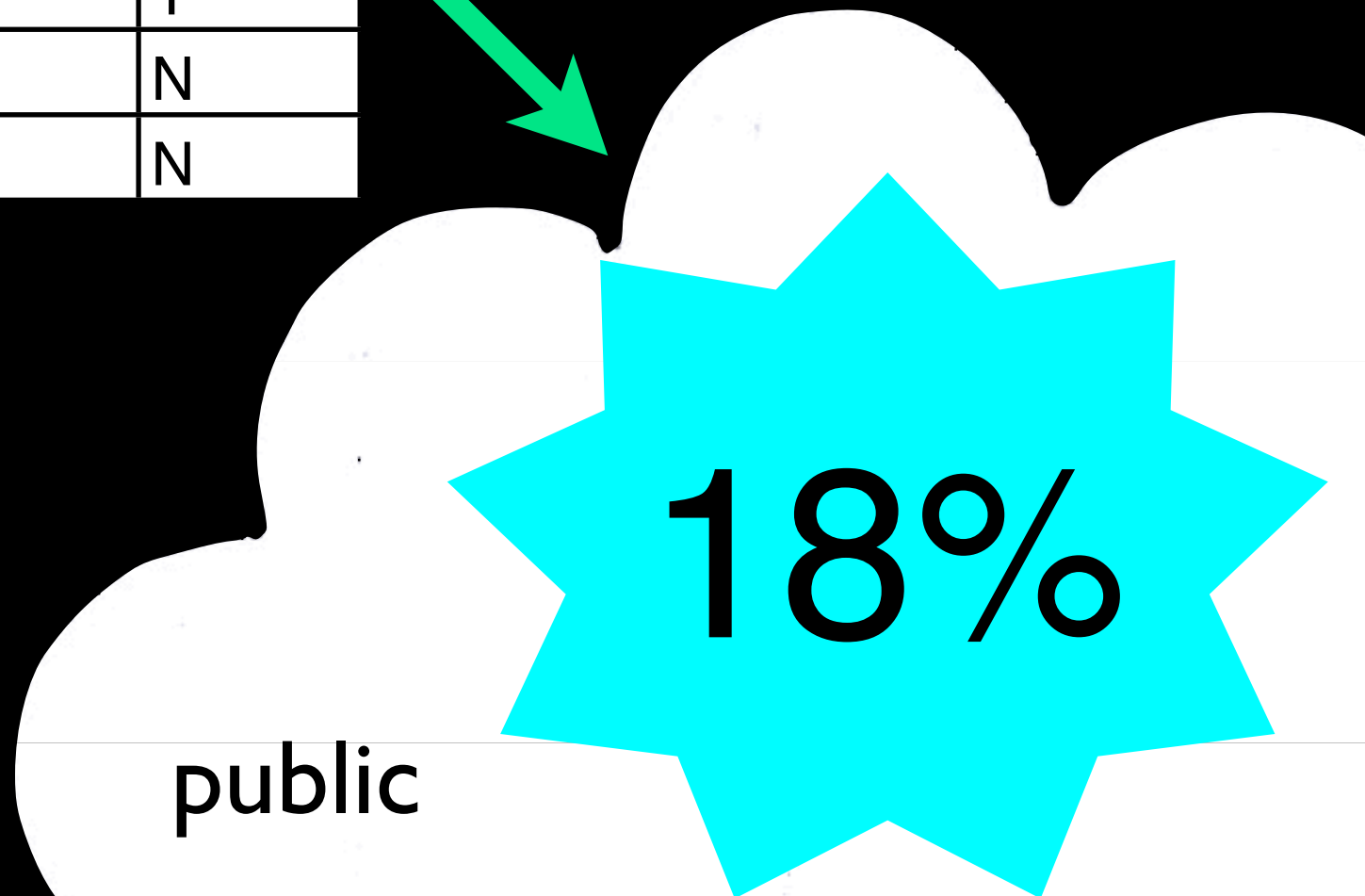
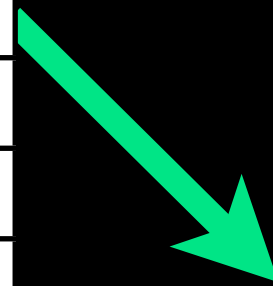
promise: if you leave  
the database, no  
outcome will change  
probability by very  
much

is this achievable?

yes!

if your output is a number...

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



public

# if your output is a number...

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



18%

add noise with particular shape

public

# sensitivity of a function $f$

$$\Delta f = \max_{D_1, D_2} |f(D_1) - f(D_2)|$$

for neighboring data sets  $D_1, D_2$

# sensitivity of a function $f$

$$\Delta f = \max_{D_1, D_2} |f(D_1) - f(D_2)|$$

for neighboring data sets  $D_1, D_2$

- measures how much one person can affect output



# sensitivity of a function $f$

$$\Delta f = \max_{D_1, D_2} |f(D_1) - f(D_2)|$$

for neighboring data sets  $D_1, D_2$

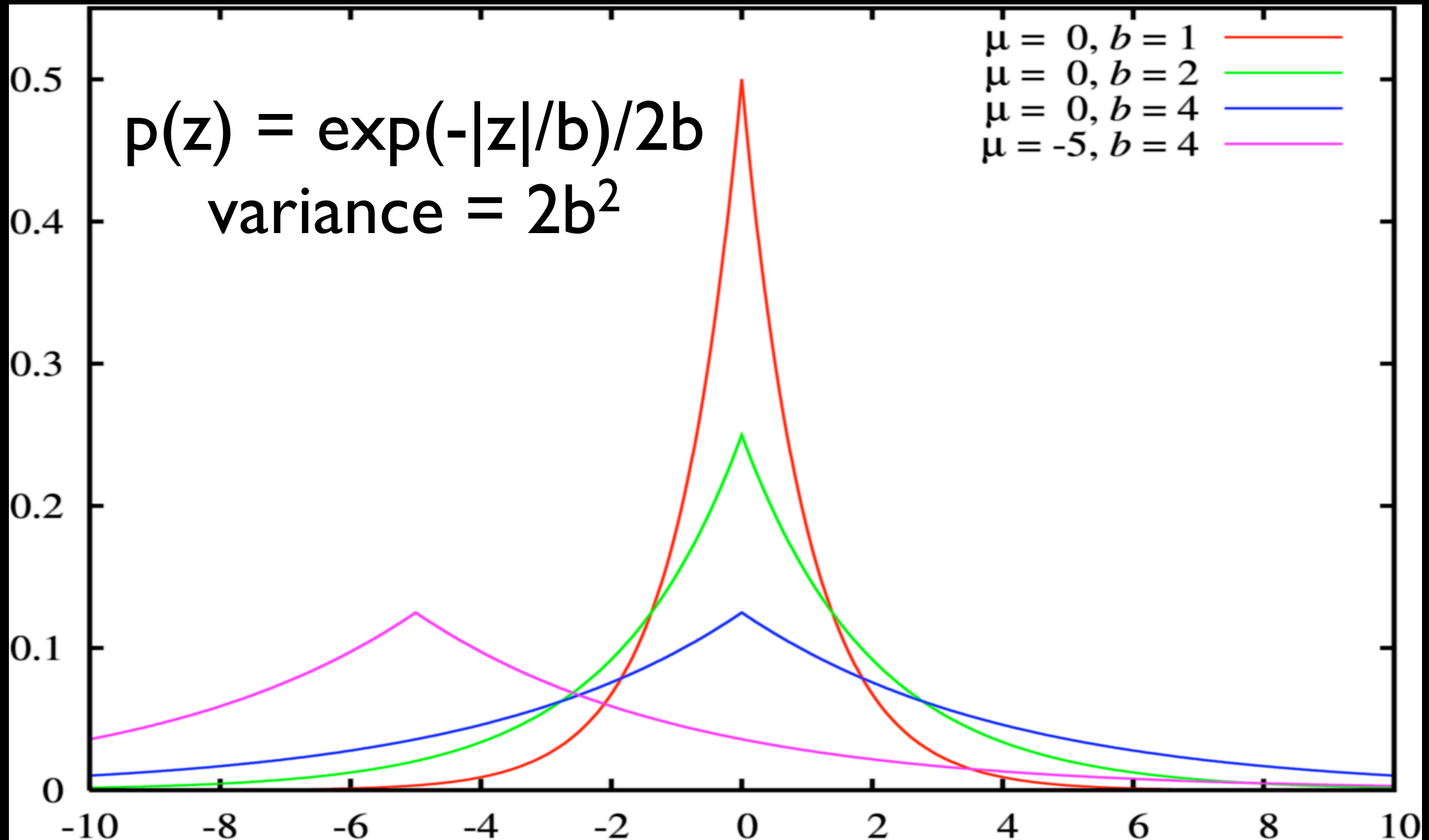
- measures how much one person can affect output
- sensitivity is 1 for **counting queries** that count number of rows satisfying a predicate

# scale noise with sensitivity

$$\Delta f = \max_{D_1, D_2} |f(D_1) - f(D_2)|$$

[DMNS06]: on query  $f$ , can add scaled symmetric noise  $\text{Lap}(b)$  with  $b = \Delta f/\epsilon$ , to achieve  $\epsilon$ -differential privacy.

# Laplace distribution Lap(b)



# applying the Laplace mechanism

# applying the Laplace mechanism

single **counting query**: how many people in the database satisfy predicate **P**?

- sensitivity = 1
- can add noise  $\text{Lap}(1/\epsilon)$

what if want more than one  
query? ...composition

- an  $\varepsilon_1$ -DP mechanism, followed by an  $\varepsilon_2$ -DP mechanism, is  $(\varepsilon_1 + \varepsilon_2)$ -DP

# what if want more than one query? ...composition

- an  $\varepsilon_1$ -DP mechanism, followed by an  $\varepsilon_2$ -DP mechanism, is  $(\varepsilon_1 + \varepsilon_2)$ -DP
- can also sum both the epsilons and the deltas for  $(\varepsilon, \delta)$ -DP



# what if want more than one query? ...composition

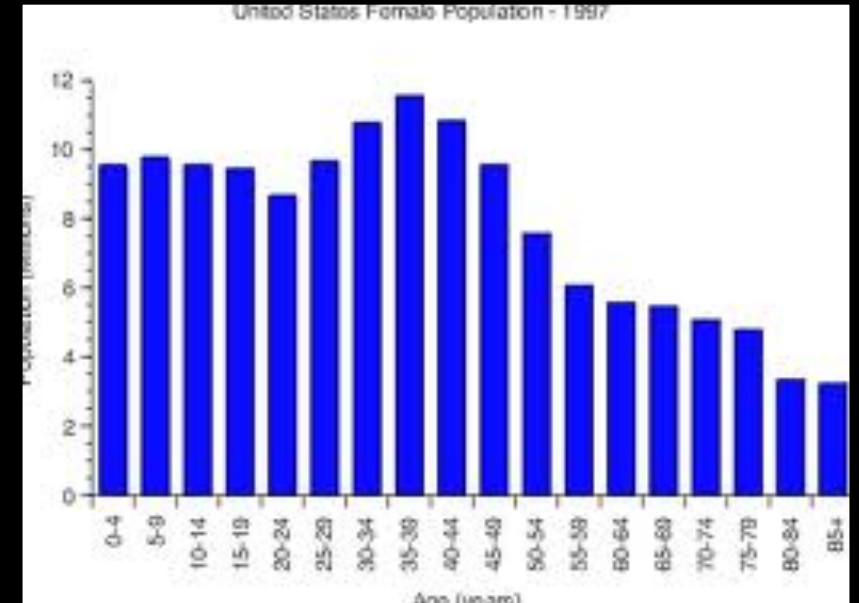
- an  $\epsilon_1$ -DP mechanism, followed by an  $\epsilon_2$ -DP mechanism, is  $(\epsilon_1 + \epsilon_2)$ -DP
- can also sum both the epsilons and the deltas for  $(\epsilon, \delta)$ -DP
- more sophisticated argument:  $k$  runs of  $(\epsilon, \delta)$ -DP mechanisms gives  $(\epsilon', k\delta + \delta')$ -DP for  $\epsilon' = (2k \ln(1/\delta'))^{1/2}\epsilon + k\epsilon(e^\epsilon - 1)$

# applying the Laplace mechanism

vector-valued queries of dimension  $d$

- can apply composition and add noise  $\text{Lap}(d\Delta f/\epsilon)$  in each component of output, where  $\Delta f$  is the sensitivity of each component

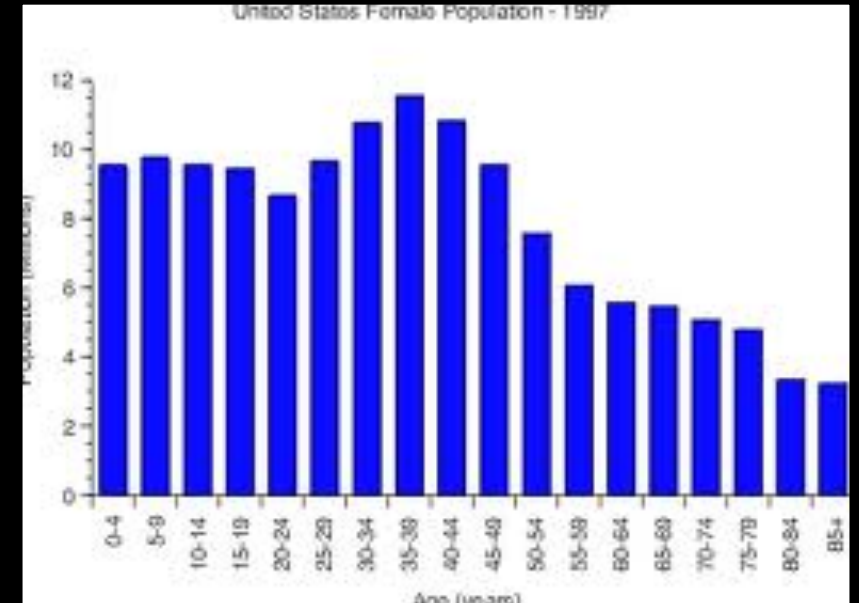
# applying the Laplace mechanism



## histogram queries

- could again use noise  $\text{Lap}(d/\epsilon)$

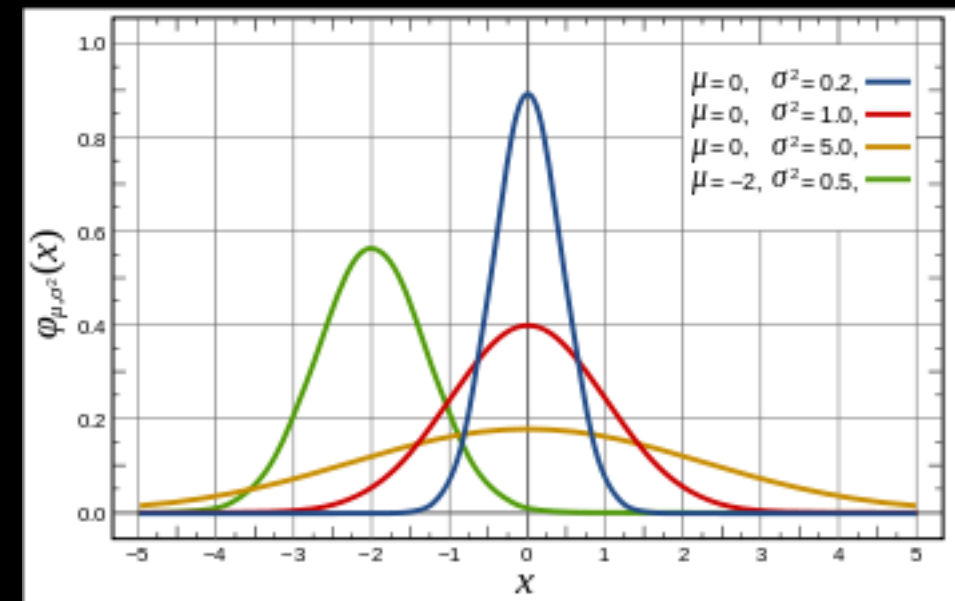
# applying the Laplace mechanism



## histogram queries

- could again use noise  $\text{Lap}(d/\epsilon)$
- but actually only need  $\sim \text{Lap}(1/\epsilon)$ , since sensitivity generalizes as  $\max L_1$  distance

# Gaussian mechanism



[DKMMN06]: Gaussian noise gives  $(\epsilon, \delta)$ -DP  
with

$$\sigma \geq (2 \ln(2/\delta))^{1/2} / \epsilon * (\max L_2 \text{ distance})$$

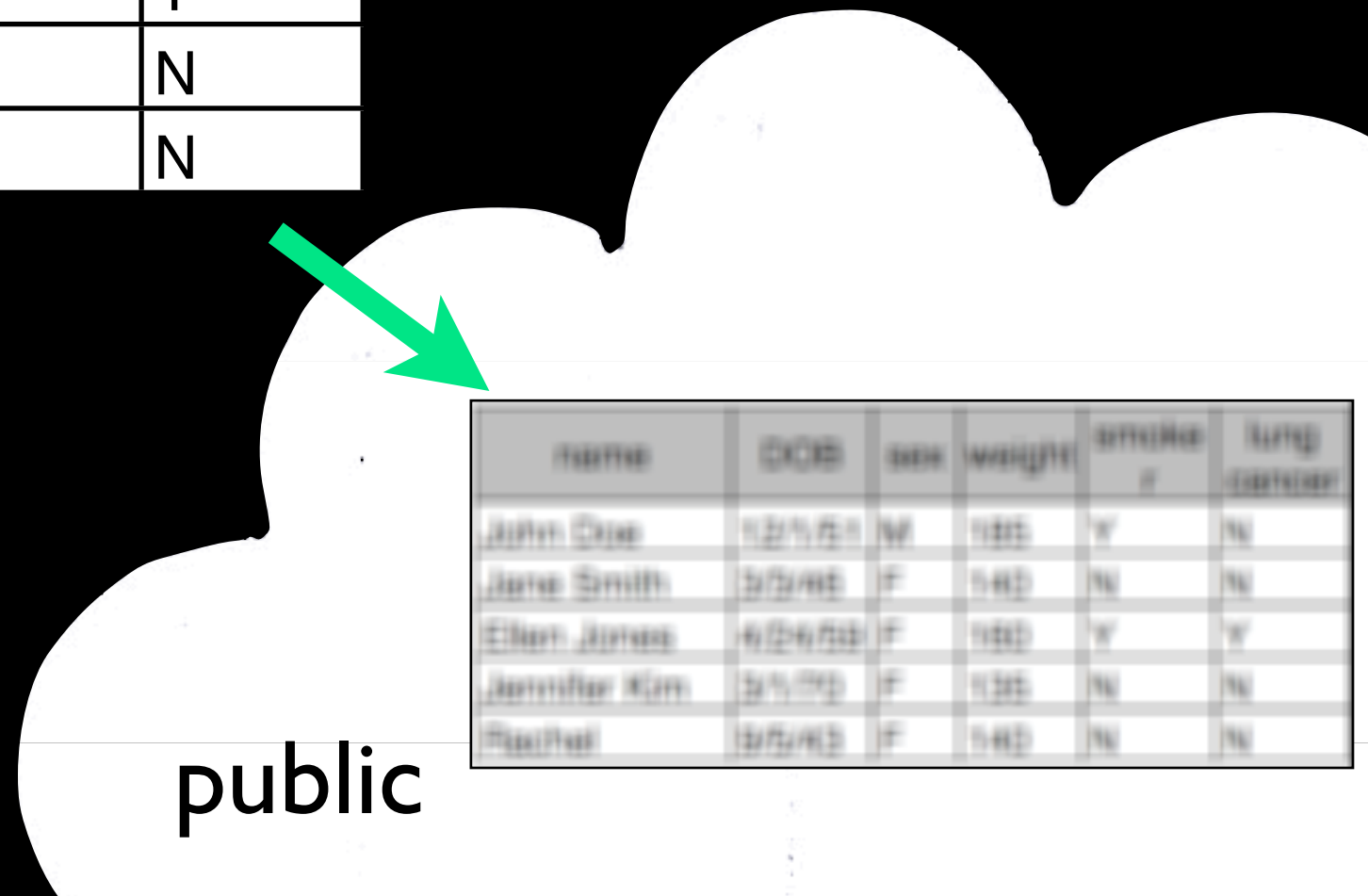
Ok, but I wanted to use my database for more than a handful of statistics...

Ok, but I wanted to use my database for more than a handful of statistics...

Data can be “big” in two dimensions: more rows makes privacy easier (lower sensitivity); more columns makes it harder (more queries to preserve)

# handling an exponential number of queries

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N



handling an expected  
number of queries

name	DOB	sex	weight	smoke	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

what fraction over age 50? what fraction smoke and have lung cancer? what fraction of males over 150 lbs?  
...

public

name	DOB	sex	weight	smoke	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel	9/5/43	F	140	N	N

# a brief history of synthetic data (theory)

BLR08:  $\epsilon$ -DP, error  $\log^{1/3} |Q| n^{2/3}$

# a brief history of synthetic data (theory)

BLR08:  $\epsilon$ -DP, error  $\log^{1/3} |Q| n^{2/3}$

DNRRV09:  $(\epsilon, \delta)$ -DP, error  $|Q|^{\circ(1)} n^{1/2}$

# a brief history of synthetic data (theory)

BLR08:  $\epsilon$ -DP, error  $\log^{1/3} |Q| n^{2/3}$

DNRRV09:  $(\epsilon, \delta)$ -DP, error  $|Q|^{\circ(1)} n^{1/2}$

DRV10:  $(\epsilon, \delta)$ -DP, error  $\text{polylog} |Q| n^{1/2}$

# a brief history of synthetic data (theory)

BLR08:  $\epsilon$ -DP, error  $\log^{1/3} |Q| n^{2/3}$

DNRRV09:  $(\epsilon, \delta)$ -DP, error  $|Q|^{\circ(1)} n^{1/2}$

DRV10:  $(\epsilon, \delta)$ -DP, error  $\text{polylog} |Q| n^{1/2}$

HR10:  $(\epsilon, \delta)$ -DP, error  $\log |Q| n^{1/2}$

# a brief history of synthetic data (theory)

BLR08:  $\epsilon$ -DP, error  $\log^{1/3} |Q| n^{2/3}$

DNRRV09:  $(\epsilon, \delta)$ -DP, error  $|Q|^{\circ(1)} n^{1/2}$

DRV10:  $(\epsilon, \delta)$ -DP, error  $\text{polylog} |Q| n^{1/2}$

HR10:  $(\epsilon, \delta)$ -DP, error  $\log |Q| n^{1/2}$

HLM12: simple & matches best bounds

# a brief history of synthetic data (theory)

BLR08:  $\epsilon$ -DP, error  $\log^{1/3} |Q| n^{2/3}$

DNRRV09:  $(\epsilon, \delta)$ -DP, error  $|Q|^{\circ(1)} n^{1/2}$

DRV10:  $(\epsilon, \delta)$ -DP, error  $\text{polylog} |Q| n^{1/2}$

HR10:  $(\epsilon, \delta)$ -DP, error  $\log |Q| n^{1/2}$

HLM12: simple & matches best bounds

# a brief history of synthetic data (theory)

BLR08:  $\epsilon$ -DP, error  $\log^{1/3} |Q| n^{2/3}$

DNRRV09:  $(\epsilon, \delta)$ -DP, error  $|Q|^{\circ(1)} n^{1/2}$

DRV10:  $(\epsilon, \delta)$ -DP, error  $\text{polylog} |Q| n^{1/2}$

HR10:  $(\epsilon, \delta)$ -DP, error  $\log |Q| n^{1/2}$

HLM12: simple & matches best bounds

Can (sometimes) do much better than naive noise addition, with much more sophisticated techniques



# exponential mechanism [MT07]

select an element  $C \in \text{range}(M)$  with

probability  $\sim \exp(\varepsilon u(D, C)/(2 \Delta u))$

where  $u$  is a “scoring function”

# exponential mechanism [MT07]

select an element  $C \in \text{range}(M)$  with

probability  $\sim \exp(\epsilon u(D, C) / (2 \Delta u))$

where  $u$  is a “scoring function”

privacy obvious

# exponential mechanism [MT07]

select an element  $C \in \text{range}(M)$  with

probability  $\sim \exp(\epsilon u(D, C) / (2 \Delta u))$

where  $u$  is a “scoring function”

privacy obvious

utility... depends

[BLR08]

combines

- **exponential mechanism** [MT07] for sampling complex output space
- sample complexity bounds from learning theory to guarantee existence of good output

[BLR08]

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Michael Ray	3/2/81	M	200	Y	N
Fran Michaels	9/9/54	F	155	N	N
Rachel Kim	1/21/77	F	130	Y	Y
Michelle Lo	2/2/83	F	135	N	N
Nira Waters	9/5/43	F	140	N	N
Jennifer Kim	3/1/70	F	135	N	N
Lisa Smith	9/5/43	F	140	N	N
Tony Miller	12/1/51	M	210	Y	N
Eve Casey	3/3/46	F	140	N	N
Paul Larson	4/24/59	F	160	Y	Y
Noelle Mason	3/1/70	F	130	N	N
Rachel Waters	9/5/43	F	140	Y	N
Shirley Wu	3/1/70	F	150	N	N
Rachel Waters	9/5/43	F	140	N	Y
Lawrence Vay	12/1/51	M	185	Y	N
Laura Rich	3/3/46	F	140	N	N

# [BLR08]

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Michael Ray	3/2/81	M	200	Y	N
Fran Michaels	9/9/54	F	155	N	N
Rachel Kim	1/21/77	F	130	Y	Y
Michelle Lo	2/2/83	F	135	N	N
Nira Waters	9/5/43	F	140	N	N
Jennifer Kim	3/1/70	F	135	N	N
Lisa Smith	9/5/43	F	140	N	N
Tony Miller	12/1/51	M	210	Y	N
Eve Casey	3/3/46	F	140	N	N
Paul Larson	4/24/59	F	160	Y	Y
Noelle Mason	3/1/70	F	130	N	N
Rachel Waters	9/5/43	F	140	Y	N
Shirley Wu	3/1/70	F	150	N	N
Rachel Waters	9/5/43	F	140	N	Y
Lawrence Vay	12/1/51	M	185	Y	N
Laura Rich	3/3/46	F	140	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

# [BLR08]

Size  $\tilde{O}(VCDIM(Q)/\epsilon^2)$

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Michael Ray	3/2/81	M	200	Y	N
Fran Michaels	9/9/54	F	155	N	N
Rachel Kim	1/21/77	F	130	Y	Y
Michelle Lo	2/2/83	F	135	N	N
Nira Waters	9/5/43	F	140	N	N
Jennifer Kim	3/1/70	F	135	N	N
Lisa Smith	9/5/43	F	140	N	N
Tony Miller	12/1/51	M	210	Y	N
Eve Casey	3/3/46	F	140	N	N
Paul Larson	4/24/59	F	160	Y	Y
Noelle Mason	3/1/70	F	130	N	N
Rachel Waters	9/5/43	F	140	Y	N
Shirley Wu	3/1/70	F	150	N	N
Rachel Waters	9/5/43	F	140	N	Y
Lawrence Vay	12/1/51	M	185	Y	N
Laura Rich	3/3/46	F	140	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

# [BLR08]

Size  $\tilde{O}(VCDIM(Q)/\epsilon^2)$

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Michael Ray	3/2/81	M	200	Y	N
Fran Michaels	9/9/54	F	155	N	N
Rachel Kim	1/21/77	F	130	Y	Y
Michelle Lo	2/2/83	F	135	N	N
Nira Waters	9/5/43	F	140	N	N
Jennifer Kim	3/1/70	F	135	N	N
Lisa Smith	9/5/43	F	140	N	N
Tony Miller	12/1/51	M	210	Y	N
Eve Casey	3/3/46	F	140	N	N
Paul Larson	4/24/59	F	160	Y	Y
Noelle Mason	3/1/70	F	130	N	N
Rachel Waters	9/5/43	F	140	Y	N
Shirley Wu	3/1/70	F	150	N	N
Rachel Waters	9/5/43	F	140	N	Y
Lawrence Vay	12/1/51	M	185	Y	N
Laura Rich	3/3/46	F	140	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	sex	weight	smoker	lung cancer
WWW	4/13/48	F	125	N	N
XXX	4/22/61	F	135	Y	Y
YYY	1/11/74	F	150	N	Y
ZZZ	10/5/44	F	150	N	N





# [HLM12]

- simple to describe and to implement
- actually implemented and tested it
- state of the art in theory, performs well in practice (and quickly, despite bad worst-case news)

# [HLM12]

- simple to describe and to implement
- actually implemented and tested it
- state of the art in theory, performs well in practice (and quickly, despite bad worst-case news)

...will hear more about multiplicative weights-based techniques in Jon's talk

I have to know all my queries in advance?!

# interactive mechanisms

- so far, have discussed creating synthetic data where must know query set in advance
- tools exist to answer similar number of queries on the fly (correlating randomness across queries)

It seems like DP would add too much noise for my application.

It seems like DP would add too much noise for my application.

... stop and think about what this means

DP connected to *robustness* of computation to presence or absence of individuals



DP connected to *robustness* of computation to presence or absence of individuals

- computation not robust? (should worry!)

DP connected to *robustness* of computation to presence or absence of individuals

- computation not robust? (should worry!)
- need more data (individuals) to get desired privacy-utility tradeoff (should think)

DP connected to *robustness* of computation to presence or absence of individuals

- computation not robust? (should worry!)
- need more data (individuals) to get desired privacy-utility tradeoff (should think)
- expect on “real” data will be robust (we can do something about this!)

# robustness (an aside)

- robustness in DP sense not identical to statistical robustness---DP is worst-case rather than wrt to distribution
- there are connections (will mention shortly)

# expect study robust on actual data

- idea 1: bootstrapping

# bootstrap aggregation

- given training dataset, create many new training sets of smaller size by sampling uniformly with replacement
- fit your model (estimate your statistic) on each
- combine (e.g., voting, averaging)

# bootstrap intuition

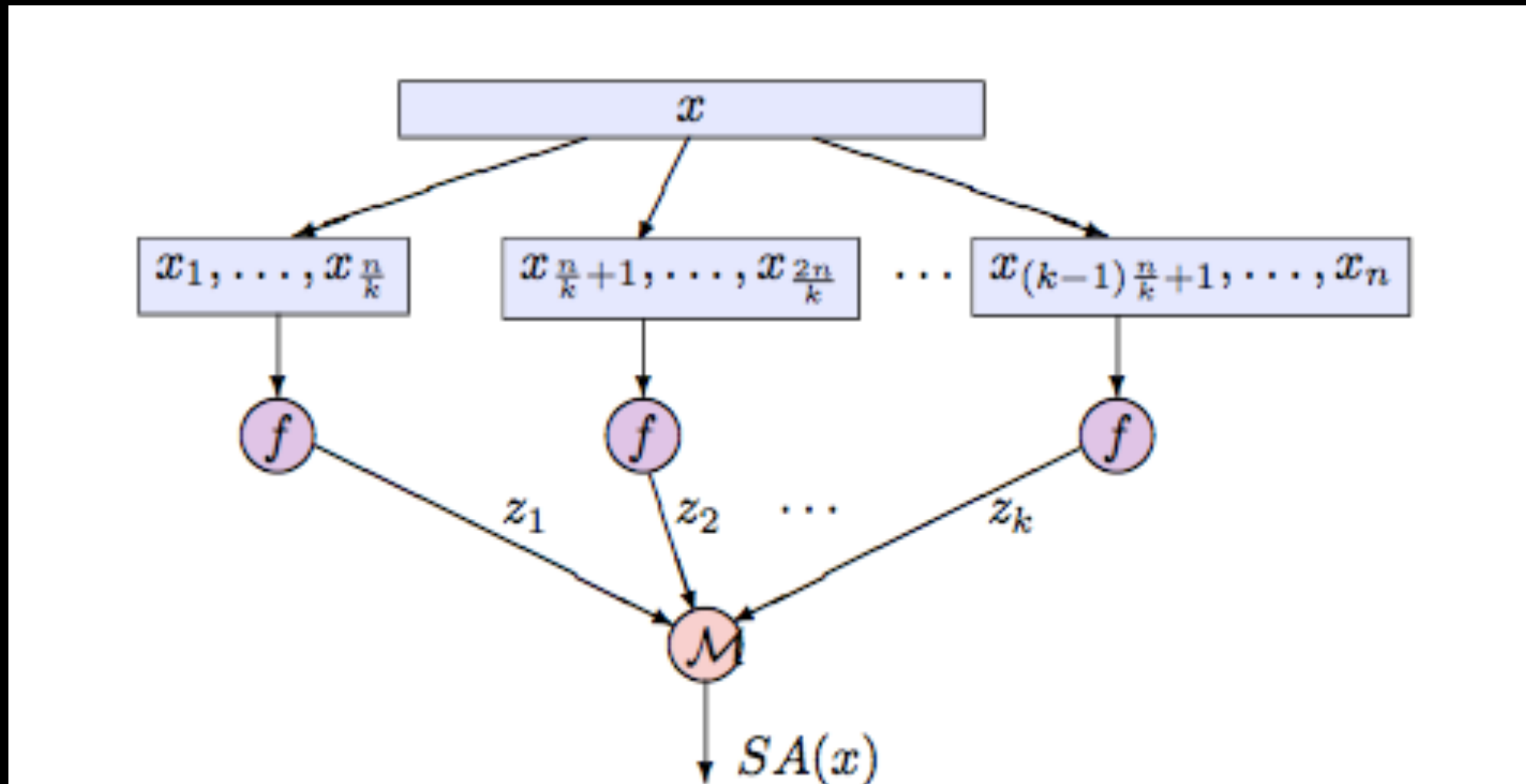
- a “robust” statistic should be stable on most reasonably sized subsets of the data
- if statistic is somewhat unstable, aggregation “smooths” result
- if statistic was very stable, loss in precision should be small

# bootstrap for privacy

- if function is not low-sensitivity but suspect it's usually stable, not clear how to guarantee DP
- if *aggregation* preserves privacy, get DP guarantee even when aggregating non-DP estimates



# bootstrap for privacy = subsample and aggregate [NRS07]



# subsample and aggregate: good news

- can use any DP aggregation function (as long as choice doesn't depend on data)
- private aggregation just requires adding noise scaled to sensitivity of the aggregation function
- privacy is immediate!

# subsample and aggregate: bad news

- may be difficult to bound worst-case sensitivity of aggregation function
- default bound is max of its range (quite bad)
- may be difficult to get good utility guarantees

# subsample and aggregate: applications

- underlying function might be selecting best model from among set of  $m$  options; could aggregate with a noisy max

# subsample and aggregate: applications

- underlying function might be selecting best model from among set of  $m$  options; could aggregate with a noisy max
- similarly, could output a set of significant features (as for LASSO)

# expect study robust on actual data

- idea 1: bootstrapping
- idea 2: check robustness before proceeding

# check robustness

- would like to be able to test in DP manner whether computation “should” proceed
- “should”: e.g., whether desired function is robust (low-sensitivity) on *actual data*
- if not, halt

# local sensitivity of function $f$ on database $D$

$$\max_{D'} |f(D) - f(D')|_1$$

for  $D'$  neighboring data set of  $D$



# local sensitivity of function $f$ on database $D$

$$\max_{D'} |f(D) - f(D')|_1$$

for  $D'$  neighboring data set of  $D$

- measures how much one person can affect output *on this data*

# propose-test-release [DL09]

- *propose* a bound on local sensitivity

# propose-test-release [DL09]

- *propose* a bound on local sensitivity
- *test* in DP manner whether actual data satisfies bound

# propose-test-release [DL09]

- *propose* a bound on local sensitivity
- *test* in DP manner whether actual data satisfies bound
- if fails, halt

# propose-test-release [DL09]

- *propose* a bound on local sensitivity
- *test* in DP manner whether actual data satisfies bound
- if fails, halt
- if passes, *release* function with noise tailored to *local sensitivity*

# notes on propose-test-release

- *test* could be “what is  $L_1$  distance to closest database that fails local sensitivity bound?”

# notes on propose-test-release

- *test* could be “what is  $L_1$  distance to closest database that fails local sensitivity bound?”
- *test only* has sensitivity 1

# notes on propose-test-release

- *test* could be “what is  $L_1$  distance to closest database that fails local sensitivity bound?”
- *test* only has sensitivity 1
- can use conservative local sensitivity threshold



# notes on propose-test-release

- *test* could be “what is  $L_1$  distance to closest database that fails local sensitivity bound?”
- *test* only has sensitivity 1
- can use conservative local sensitivity threshold
- could still sometimes fail; can't get  $\epsilon$ -DP

DP output need not be noisy!

# DP output need not be noisy!

PTR can be used to privately check whether distance to nearest unstable data set is far, and if so release the *true*  $f(x)$

# robustness, revisited

# robustness, revisited

- robustness wrt adding/removing a few points from dataset

# robustness, revisited

- robustness wrt adding/removing a few points from dataset
- robustness wrt subsampling

# subsample & aggregate + propose-test-release

[ST13]: can modify subsample & aggregate so that outputs true  $f(x)$  with high probability when  $f$  is subsampling stable on  $x$

# subsample & aggregate + propose-test-release

[ST13]: can modify subsample & aggregate so that outputs true  $f(x)$  with high probability when  $f$  is subsampling stable on  $x$

shows that DP model selection only increases sample complexity of model selection by  $O(\log(1/\delta)/\epsilon)$



# DP statistics

# DP statistics

- connections to robustness; interquartile distance, median, linear regression [DworkLei09]

# DP statistics

- connections to robustness; interquartile distance, median, linear regression [DworkLei09]
- M-estimators [Lei11, NekipelovYakovlev11]

# DP statistics

- connections to robustness; interquartile distance, median, linear regression [DworkLei09]
- M-estimators [Lei11, NekipelovYakovlev11]
- for almost any estimator that is asymptotically normal on i.i.d. samples, DP adds asymptotically no additional perturbation [Smith11]

# DP statistics

- connections to robustness; interquartile distance, median, linear regression [DworkLei09]
- M-estimators [Lei11, NekipelovYakovlev11]
- for almost any estimator that is asymptotically normal on i.i.d. samples, DP adds asymptotically no additional perturbation [Smith11]
- convergence rate of DP estimators tied to gross error sensitivity [ChaudhuriHsu12]

# DP statistics

- connections to robustness; interquartile distance, median, linear regression [DworkLei09]
- M-estimators [Lei11, NekipelovYakovlev11]
- for almost any estimator that is asymptotically normal on i.i.d. samples, DP adds asymptotically no additional perturbation [Smith11]
- convergence rate of DP estimators tied to gross error sensitivity [ChaudhuriHsu12]
- minimizing convex loss functions [ChaudhuriMonteleoniSarwate11, Rubinstein et al. 2012, KiferSmithThakurta12]

# DP statistics

- connections to robustness; interquartile distance, median, linear regression [DworkLei09]
- M-estimators [Lei11, NekipelovYakovlev11]
- for almost any estimator that is asymptotically normal on i.i.d. samples, DP adds asymptotically no additional perturbation [Smith11]
- convergence rate of DP estimators tied to gross error sensitivity [ChaudhuriHsu12]
- minimizing convex loss functions [ChaudhuriMonteleoniSarwate11, Rubinstein et al. 2012, KiferSmithThakurta12]
- model selection [SmithThakurta13]

# DP statistics

- connections to robustness; interquartile distance, median, linear regression [DworkLei09]
- M-estimators [Lei11, NekipelovYakovlev11]
- for almost any estimator that is asymptotically normal on i.i.d. samples, DP adds asymptotically no additional perturbation [Smith11]
- convergence rate of DP estimators tied to gross error sensitivity [ChaudhuriHsu12]
- minimizing convex loss functions [ChaudhuriMonteleoniSarwate11, Rubinstein et al. 2012, KiferSmithThakurta12]
- model selection [SmithThakurta13]
- empirical investigations [VuSlavkovic09, ChaudhuriMonteleoniSarwate11, AbowdSchneiderVilhuber13]



I need to look at the data before I know what statistics to run.

I need to look at the data before I know what statistics to run.

not DP

I need to look at the data before I know what statistics to run.

not DP

interactive (or hybrid interactive/  
noninteractive ) mechanisms?

I need to look at the data before I know what statistics to run.

not DP

interactive (or hybrid interactive/  
noninteractive ) mechanisms?

big data force us to formalize “looking at the  
data”

# summary

# summary

- privacy easy to get wrong; DP provides compelling definition and useful dose of paranoia

# summary

- privacy easy to get wrong; DP provides compelling definition and useful dose of paranoia
- powerful tools exist (some with no cost of privacy, and some with no noise!)

# summary

- privacy easy to get wrong; DP provides compelling definition and useful dose of paranoia
- powerful tools exist (some with no cost of privacy, and some with no noise!)
- powerful intuition from notions of robustness



# summary

- privacy easy to get wrong; DP provides compelling definition and useful dose of paranoia
- powerful tools exist (some with no cost of privacy, and some with no noise!)
- powerful intuition from notions of robustness
- many nearly ready (and quite relevant) to common big data applications

# summary

- privacy easy to get wrong; DP provides compelling definition and useful dose of paranoia
- powerful tools exist (some with no cost of privacy, and some with no noise!)
- powerful intuition from notions of robustness
- many nearly ready (and quite relevant) to common big data applications
- no ready-to-use, commercial- grade applications: need demand!

# Differential Privacy Tutorial

Katrina Ligett  
[katrina@caltech.edu](mailto:katrina@caltech.edu)

“Privacy and Data-Based Research,” with Ori Heffetz.  
Available on SSRN.