

# Responsible Data Science

## Differential privacy

*March 23, 2022*

---

**Prof. Julia Stoyanovich**

Center for Data Science &  
Computer Science and Engineering  
New York University

**query sensitivity  
& composition**

# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$



**lower  $\epsilon$  = stronger privacy**



- Example 1: counting queries
  - “How many elements in  $D$  satisfy property  $P$ ?” **What’s  $\Delta q$  ?**
  - “What fraction of the elements in  $D$  satisfy property  $P$ ?”
- Example 2: max / min
  - “What is the maximum employee salary in  $D$  ?” **What’s  $\Delta q$  ?**

**Intuition: for a given  $\epsilon$ , the higher the sensitivity, the more noise we need to add to meet the privacy guarantee**

# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$

query $q$	query sensitivity $\Delta q$
select count(*) from D	1
select count(*) from D where sex = Male and age > 30	?

# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$

query $q$	query sensitivity $\Delta q$
select count(*) from D	1
select count(*) from D where sex = Male and age > 30	1
select MAX(salary) from D	?

# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$

query $q$	query sensitivity $\Delta q$
select count(*) from D	1
select count(*) from D where sex = Male and age > 30	1
select MAX(salary) from D	$MAX(salary) - MIN(salary)$
select gender, count(*) from D group by gender	?

# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$

query $q$	query sensitivity $\Delta q$
select count(*) from D	1
select count(*) from D where sex = Male and age > 30	1
select MAX(salary) from D	$MAX(salary) - MIN(salary)$
select gender, count(*) from D group by gender	1 (disjoint groups, presence or absence of one tuple impacts only one of the counts)

# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$

query  $q$

query sensitivity  $\Delta q$

select gender, count(\*)  
from D group by gender

**1 (disjoint groups)**, presence or absence of one tuple impacts only one of the counts)

an arbitrary list of  $m$  counting queries

?



# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$

query  $q$

query sensitivity  $\Delta q$

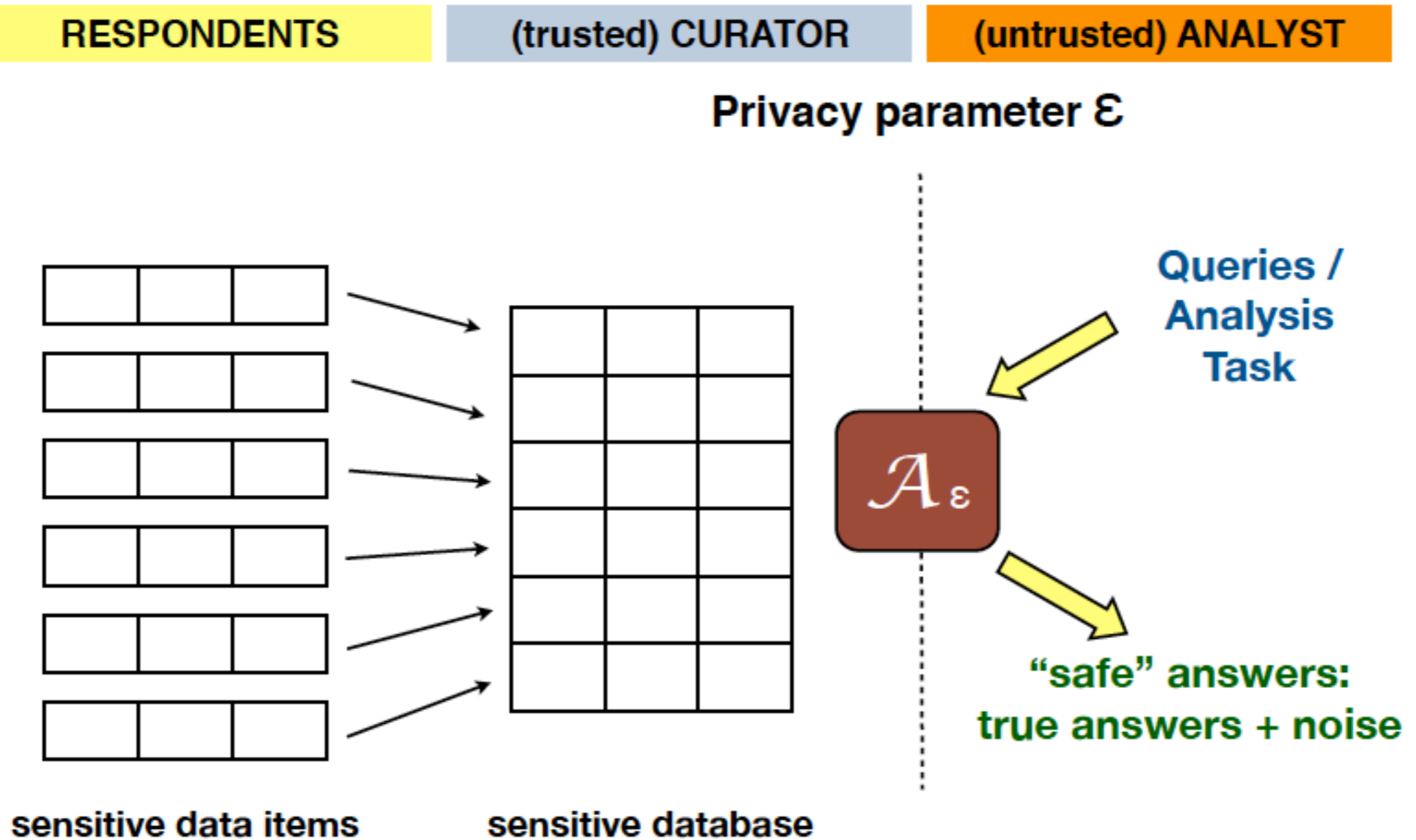
select gender, count(\*)  
from D group by gender

**1** (**disjoint groups**, presence or absence of one tuple impacts only one of the counts)

an arbitrary list of  $m$  counting queries

**$m$**  (no assumptions about the queries, and so a single individual may change the answer of **every query** by 1)

# Adding noise

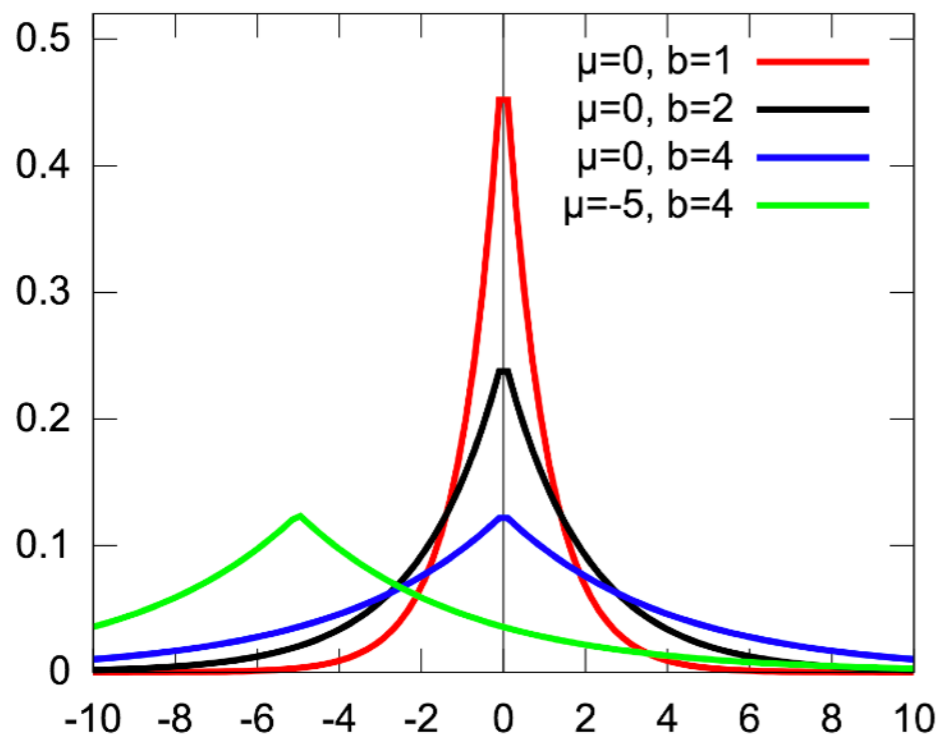


# Adding noise

Use the **Laplace mechanism** to answer  $q$  in a way that's  $\epsilon$ -differentially private

$$M(\epsilon) : q(D) + \text{Lap}\left(\frac{\Delta q}{\epsilon}\right)$$

The Laplace distribution, centered at 0 with scale  $b$ , denoted **Lap(b)**, is the distribution with probability density function:

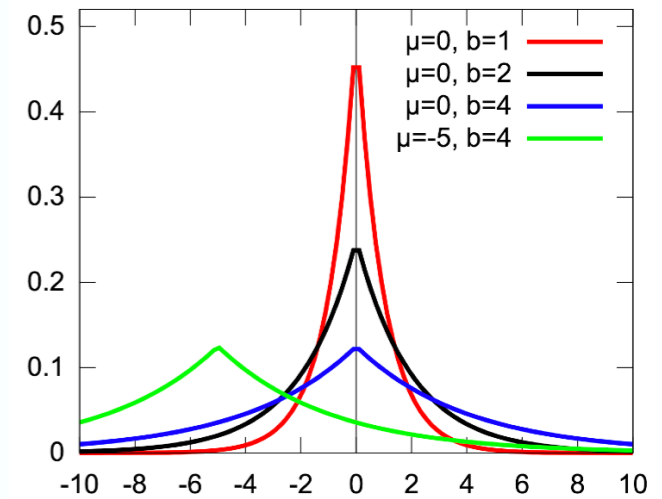
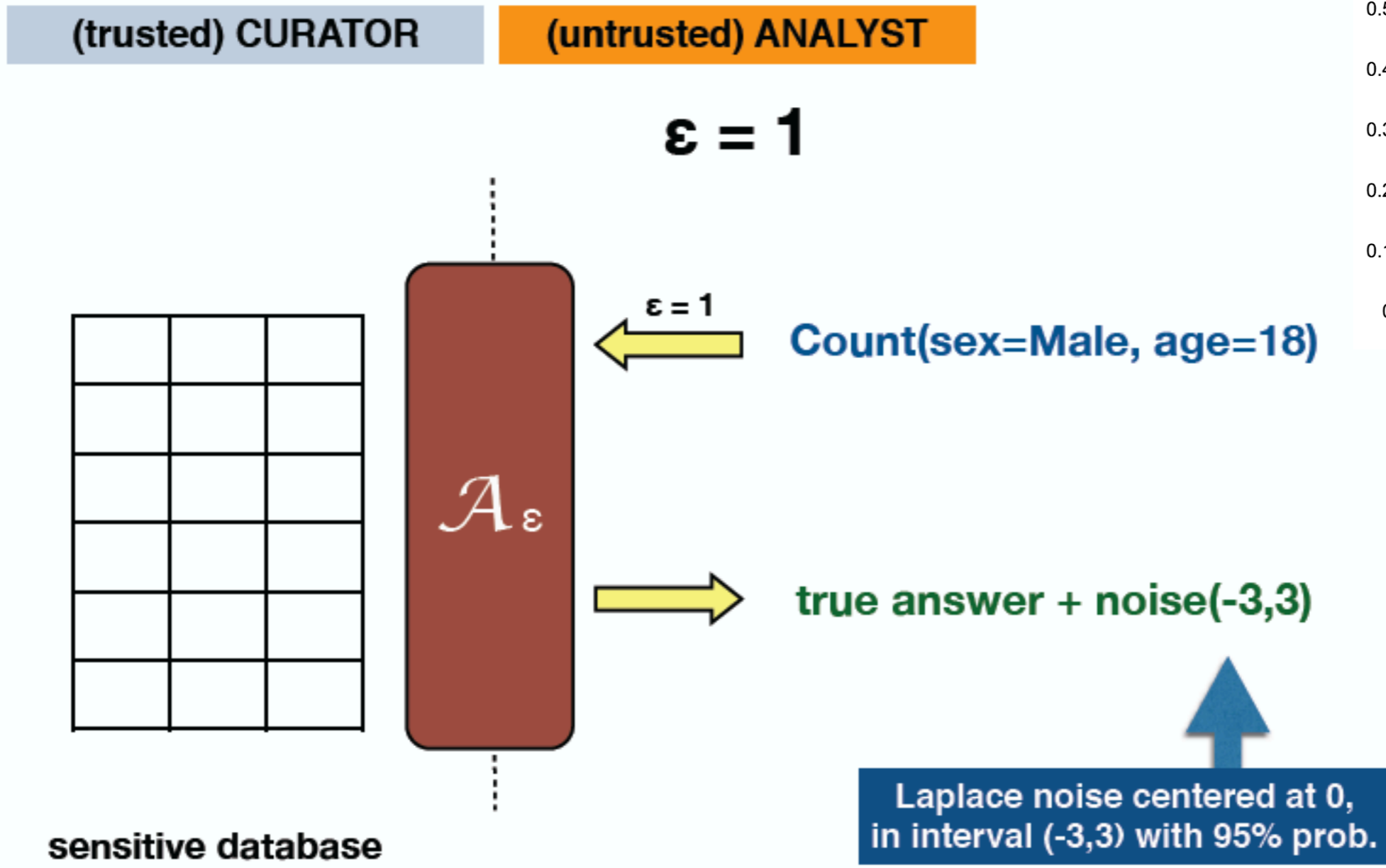


[https://en.wikipedia.org/wiki/Laplace\\_distribution](https://en.wikipedia.org/wiki/Laplace_distribution)

fix sensitivity  $\Delta q$ , verify that more noise is added for lower  $\epsilon$

↓ **lower  $\epsilon$  = stronger privacy** ↑

# Adding noise



# Query sensitivity

The  $\ell_1$  sensitivity of a query  $q$ , denoted  $\Delta q$ , is the maximum difference in the result of that query on a pair of neighboring databases

$$\Delta q = \max_{D, D'} |q(D) - q(D')|$$

query  $q$

query sensitivity  $\Delta q$

## parallel composition

select gender, count(\*)  
from D group by gender

**1** (**disjoint groups**, presence or absence of one tuple impacts only one of the counts)

## sequential composition

an arbitrary list of  $m$  counting queries

**$m$**  (no assumptions about the queries, and so a single individual may change the answer of **every query** by 1)

# Sequential composition

- Consider 4 queries executed in sequence
  - Q1: select count(\*) from D under  $\epsilon_1 = 0.5$
  - Q2: select count(\*) from D where sex = Male under  $\epsilon_2 = 0.2$
  - Q3: select count(\*) from D where sex = Female under  $\epsilon_3 = 0.25$
  - Q4: select count(\*) from D where age > 20 under  $\epsilon_4 = 0.25$
- $\epsilon = \epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4 = 1.2$  That is: all queries together are  $\epsilon$ -differentially private for  $\epsilon = 1.2$ . **Can we make a stronger guarantee?**
- This works because **Laplace noise is additive**

More generally: set a **cumulative privacy budget**, and split it between all queries, pre-processing, other data manipulation steps of the pipeline

# Parallel composition

- If the inputs are disjoint, then the result is  $\epsilon$ -differentially private for  $\epsilon = \max(\epsilon_1, \dots, \epsilon_k)$ 
  - Q1: select count(\*) from D under  $\epsilon_1 = 0.5$
  - Q2: select count(\*) from D where sex = Male under  $\epsilon_2 = 0.2$
  - Q3: select count(\*) from D where sex = Female under  $\epsilon_3 = 0.25$
  - Q4: select count(\*) from D where age > 20 under  $\epsilon_4 = 0.25$
- $\epsilon = \epsilon_1 + \max(\epsilon_2, \epsilon_3) + \epsilon_4 = 1$  That is: all queries together are  $\epsilon$ -differentially private for  $\epsilon = 1$ .

# Composition and consistency

- Consider again 4 queries executed in sequence
  - Q1: select count(\*) from D under  $\varepsilon_1 = 0.5$  returns **2005**
  - Q2: select count(\*) from D where sex = Male under  $\varepsilon_2 = 0.2$  returns **1001**
  - Q3: select count(\*) from D where sex = Female under  $\varepsilon_3 = 0.25$  returns **995**
  - Q4: select count(\*) from D where age > 20 under  $\varepsilon_4 = 0.25$  returns **1789**

Assuming that there are 2 genders in D, Male and Female, there is **no database consistent with these statistics!**

Also don't want any negative counts + may want to impose datatype checks, e.g., no working adults with age = 5 etc.

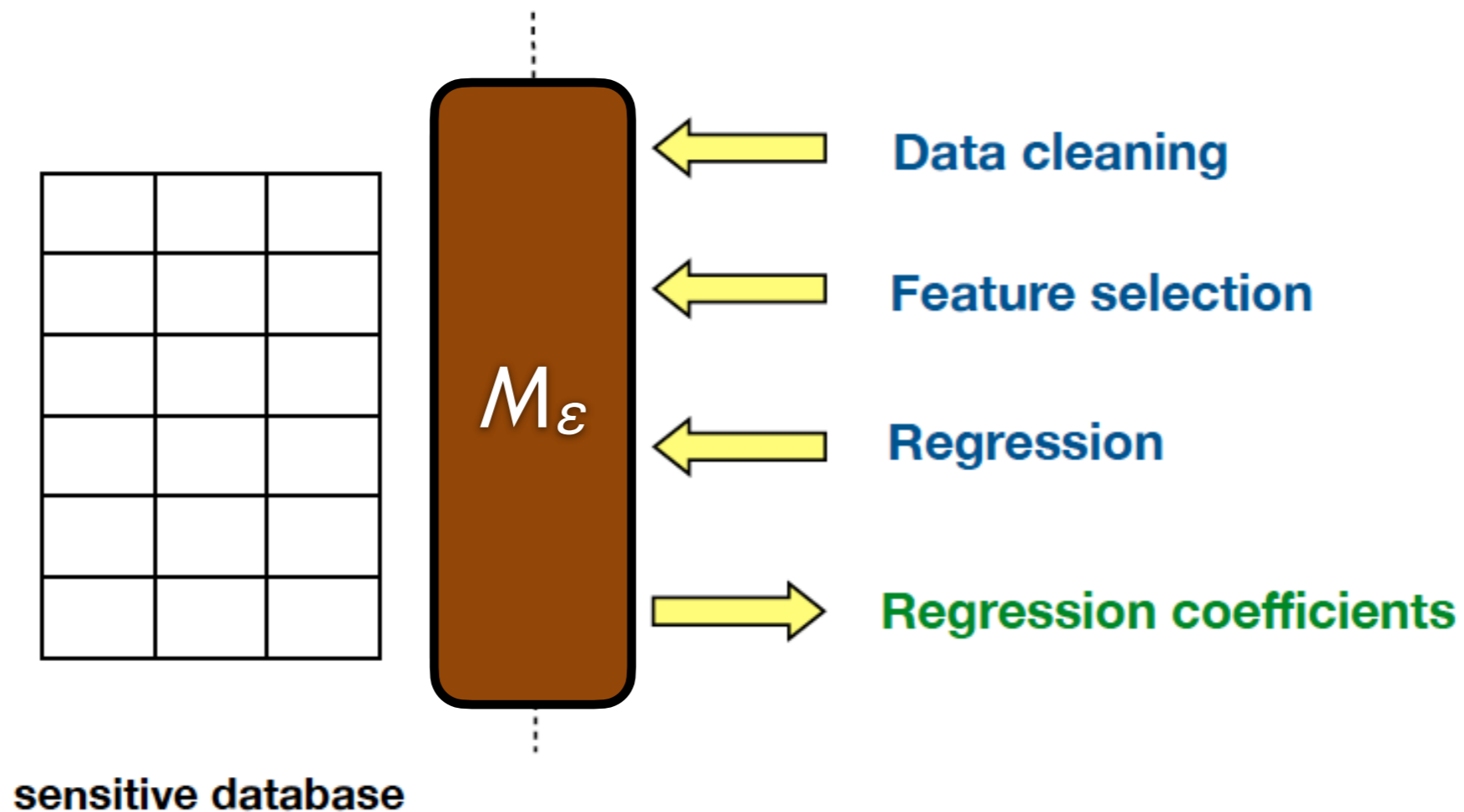


# Entire workflow must be DP

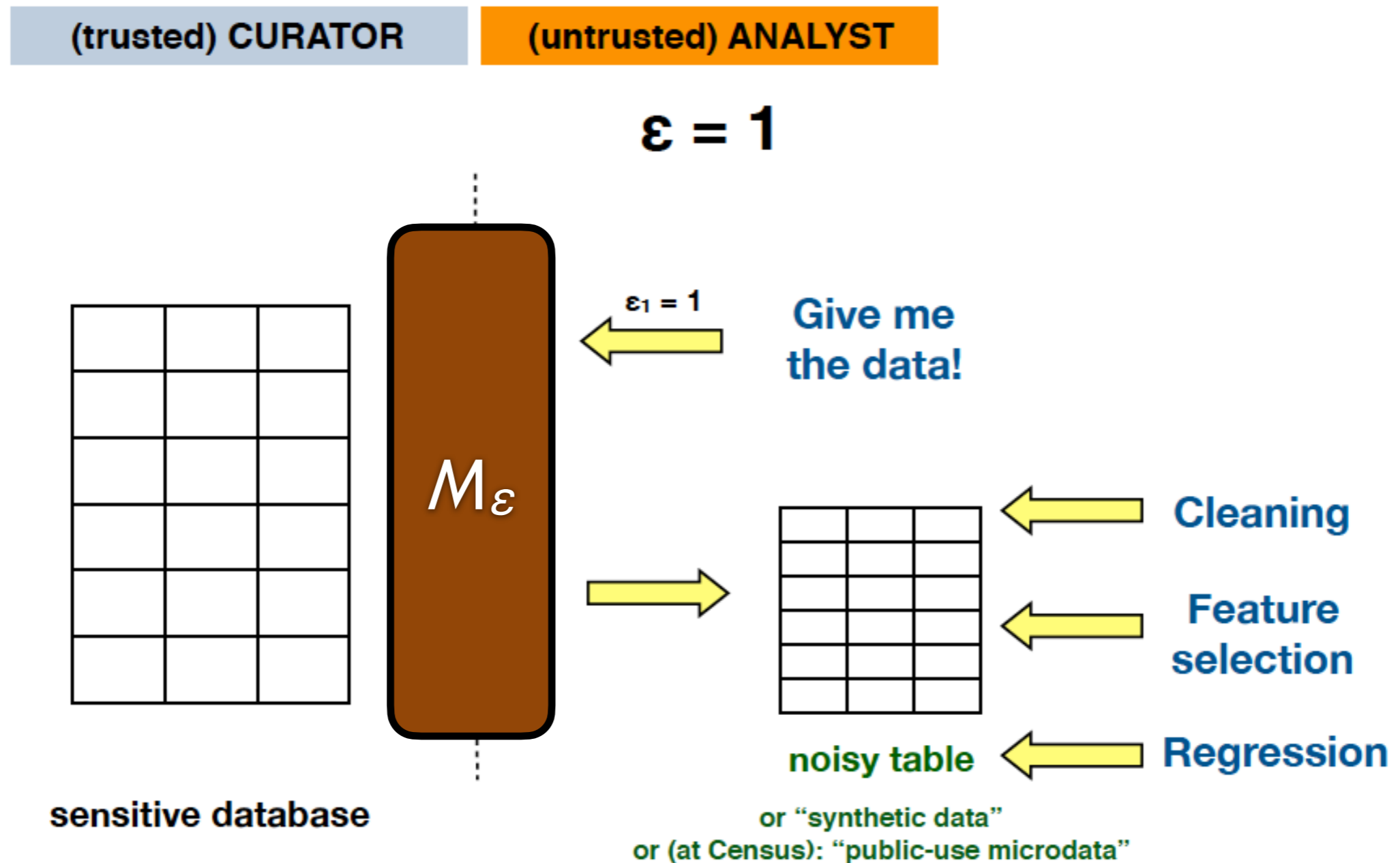
(trusted) CURATOR

(untrusted) ANALYST

$\epsilon = 1$



# Privacy-preserving synthetic data



*DP synthetic  
data generation*

# DP synthetic data

## Lots of advantages

- Consistency is not an issue
- Analysts can treat synthetic data as a regular dataset, run existing tools
- No need to worry about the privacy budget
- Can answer as many queries as they want, and any kind of a query they want, including record-level queries

## What's the catch?

**Recall the Fundamental Law of Information Recovery. It tells us that we cannot answer all these queries accurately and still preserve privacy!**

Therefore, when releasing synthetic data, we need to document it with which queries it supports well

# Data Synthesizer



input

UID	sex	race	MarriageSta	DateOfBirth	age	juv_fel	cour_decile	score
1	1	0	1	4/18/47	69	0	1	1
2	2	0	2	1/22/82	34	0	3	3
3	3	0	2	5/14/91	24	0	4	4
4	4	0	2	1/21/93	23	0	8	8
5	5	0	1	1/22/73	43	0	1	1
6	6	0	1	8/22/71	44	0	1	1
7	7	0	3	7/23/74	41	0	6	6
8	8	0	1	2/25/73	43	0	4	4
9	9	0	3	6/10/94	21	0	3	3
10	10	0	3	6/1/88	27	0	4	4
11	11	1	3	8/22/78	37	0	1	1
12	12	0	2	12/2/74	41	0	4	4
13	13	1	3	6/14/68	47	0	1	1
14	14	0	2	1/25/85	31	0	3	3
15	15	0	4	1/25/79	37	0	1	1
16	16	0	2	6/22/90	25	0	10	10
17	17	0	3	12/24/84	31	0	5	5
18	18	0	3	1/8/85	31	0	3	3
19	19	0	2	6/28/51	64	0	6	6
20	20	0	2	11/29/94	21	0	9	9
21	21	0	3	8/6/88	27	0	2	2
22	22	1	3	3/22/95	21	0	4	4
23	23	0	4	1/23/92	24	0	4	4
24	24	0	3	1/10/73	43	0	1	1
25	25	0	1	8/24/83	32	0	3	3
26	26	0	1	2/8/89	27	0	3	3
27	27	1	3	9/3/79	36	0	3	3
28	28	1	1	9/3/79	36	0	3	3

Data  
Describer



summary

age	int	min=23	32%	40
		max=60	mis	20
				0
name	str	length	no	
		10 to 98	mis	
sex	str	cat	10%	60
			mis	30
				0

Data  
Generator



output

UID	sex	race	MarriageSta	DateOfBirth	age	juv_fel	cour_decile	score
1	1	0	1	4/18/47	69	0	1	1
2	2	0	2	1/22/82	34	0	3	3
3	3	0	2	5/14/91	24	0	4	4
4	4	0	2	1/21/93	23	0	8	8
5	5	0	1	1/22/73	43	0	1	1
6	6	0	1	8/22/71	44	0	1	1
7	7	0	3	7/23/74	41	0	6	6
8	8	0	1	2/25/73	43	0	4	4
9	9	0	3	6/10/94	21	0	3	3
10	10	0	3	6/1/88	27	0	4	4
11	11	1	3	8/22/78	37	0	1	1
12	12	0	2	12/2/74	41	0	4	4
13	13	1	3	6/14/68	47	0	1	1
14	14	0	2	1/25/85	31	0	3	3
15	15	0	4	1/25/79	37	0	1	1
16	16	0	2	6/22/90	25	0	10	10
17	17	0	3	12/24/84	31	0	5	5
18	18	0	3	1/8/85	31	0	3	3
19	19	0	2	6/28/51	64	0	6	6
20	20	0	2	11/29/94	21	0	9	9
21	21	0	3	8/6/88	27	0	2	2
22	22	1	3	3/22/95	21	0	4	4
23	23	0	4	1/23/92	24	0	4	4
24	24	0	3	1/10/73	43	0	1	1
25	25	0	1	8/24/83	32	0	3	3
26	26	0	1	2/8/89	27	0	3	3
27	27	1	3	9/3/79	36	0	3	3
28	28	1	1	9/3/79	36	0	3	3

Model  
Inspector



comparison

age	int	min=23	32%	40
		max=60	mis	20
				0
name	str	length	no	
		10 to 98	mis	
sex	str	cat	10%	60
			mis	30
				0

# Data Synthesizer



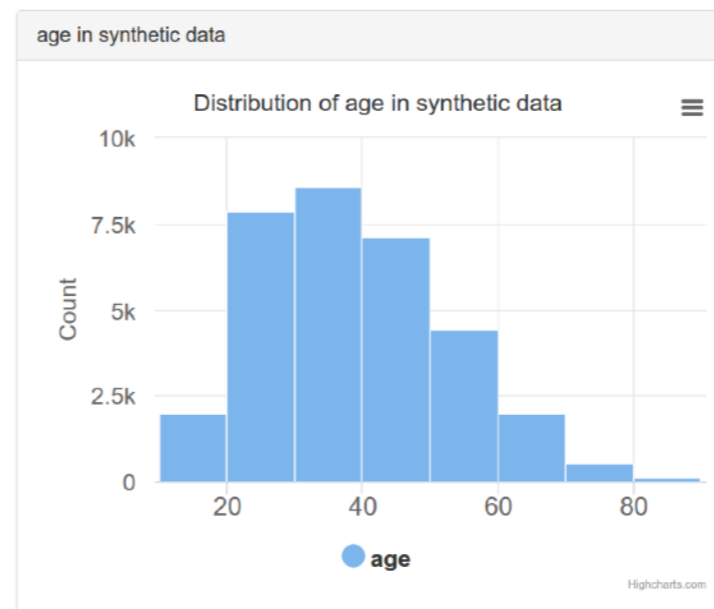
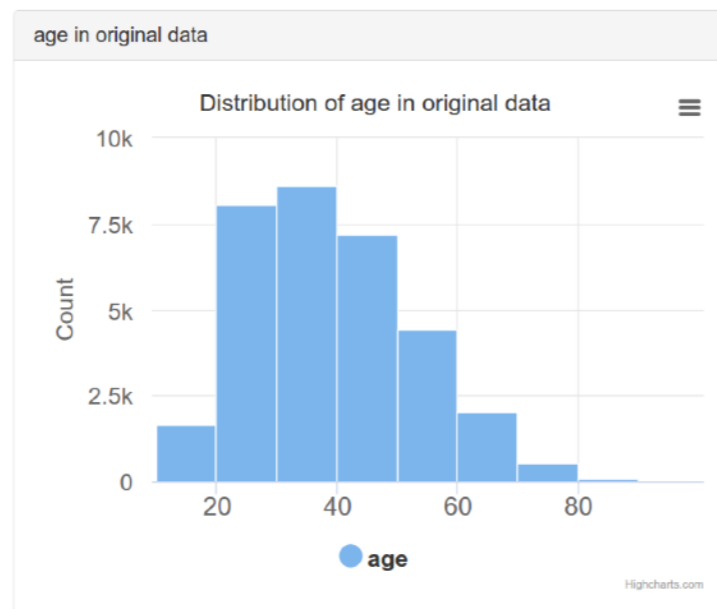
- Main goal: **usability first**
  - user is the data owner
  - the tool picks up data types from the input file: categorical / string / numerical (integer, float) / date-time
  - the tool computes the frequency of missing values per attribute
  - user can then inspect the result, over-ride what was learned about an attribute, e.g., whether it's categorical, or what its datatype is
- The tool generates an output dataset of a specified size, in one of three modes
  - **random** - type-consistent random output
  - **independent attribute** - learn a noisy histogram for each attribute
  - **correlated attribute** - learn a noisy Bayesian network (BN)

# Data Synthesizer: Independent attributes



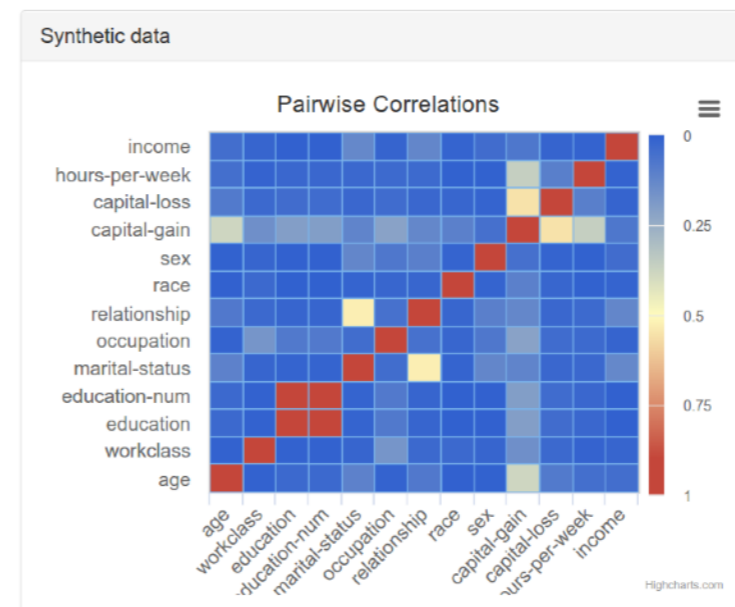
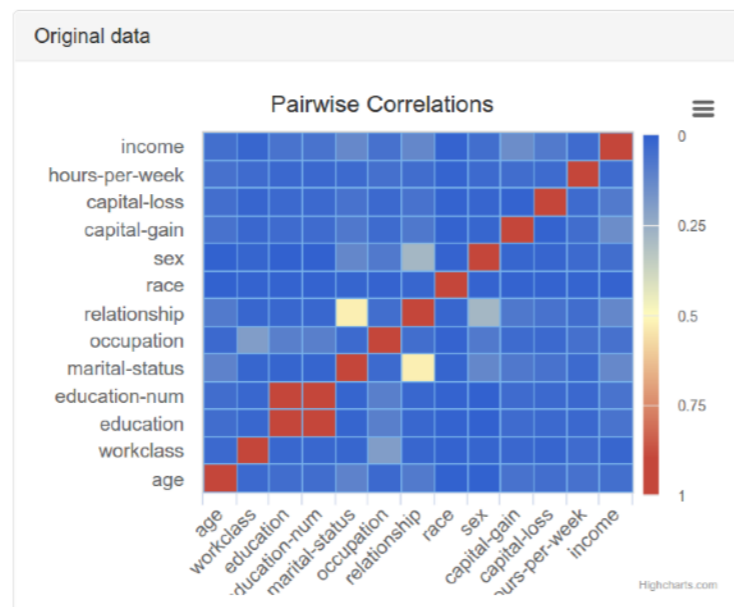
Given the over-all privacy budget  $\epsilon$ , and an input dataset of size  $n$ . Allocate  $\epsilon/d$  of the budget to each attribute  $\mathbf{A}_i$  in  $\{\mathbf{A}_1, \dots, \mathbf{A}_d\}$ . Then for each attribute:

- Compute the  $i$ th histogram with  $t$  bins ( $t=20$  by default), with query  $q_i$
- The sensitivity  $\Delta q_i$  of this (or any other) histogram query is  $2/n$  **Why?**
- So, each bin's noisy probability is computed by adding  $Lap\left(\frac{2d}{\epsilon n}\right)$



# Data Synthesizer: Correlated attributes

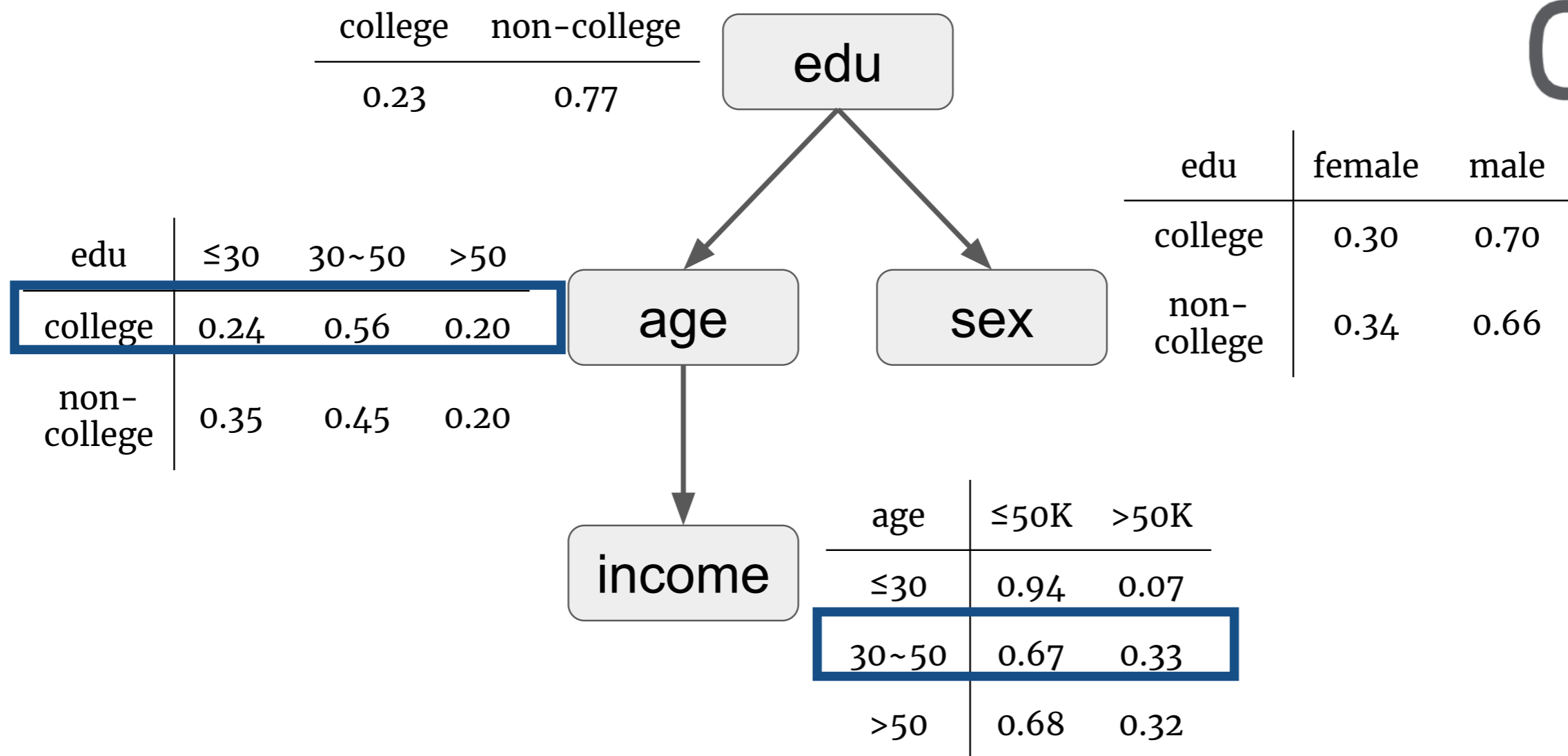
- Learn a differentially private Bayesian network (BN)
- Use the method called **PrivBayes** [Zhang, Cormode, Procopiuc, Srivastava, Xiao, 2016]
- Privacy budget is split equally between (a) network structure computation and (b) populating the conditional probability tables of each BN node
- User inputs privacy budget  $\epsilon$  and the maximum number of parents for a BN node  $k$  - you'll play with these settings as part of HW2
- The tool treats a missing attribute value as one of the values in the attribute's domain (not shown in the examples in the next two slides)





# Data Synthesizer: Correlated attributes

$K=1$  not a causal DAG, a regular Bayesian network!

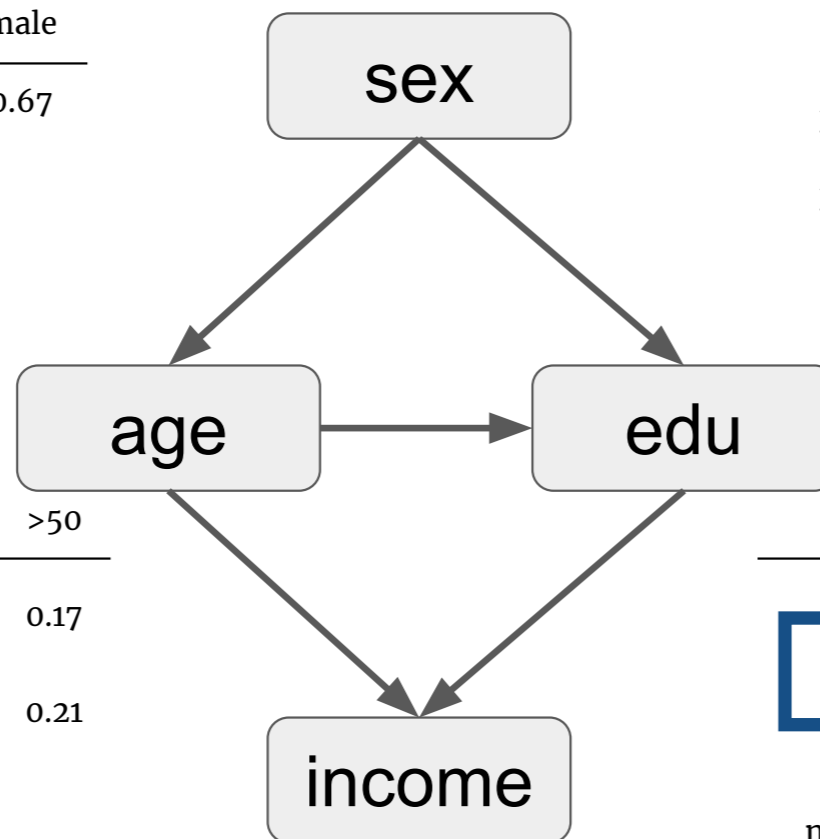


# Data Synthesizer: Correlated attributes

not a causal DAG, a regular Bayesian network!

K=2

female	male
0.33	0.67



sex	≤30	30~50	>50
female	0.40	0.43	0.17
male	0.29	0.59	0.21

age	sex	college	non-college
≤30	female	0.18	0.82
≤30	male	0.16	0.84
30~50	female	0.25	0.75
30~50	male	0.28	0.72
>50	female	0.17	0.83
>50	male	0.25	0.75

edu	age	≤50K	>50K
college	≤30	0.83	0.17
college	30~50	0.45	0.55
college	>50	0.41	0.59
non-college	≤30	0.96	0.04
non-college	30~50	0.76	0.24
non-college	>50	0.75	0.25

