

Building complex DP algorithms using composition

Privacy & Fairness in Data Science

CompSci 590.01 Fall 2018



DUKE
COMPUTER SCIENCE

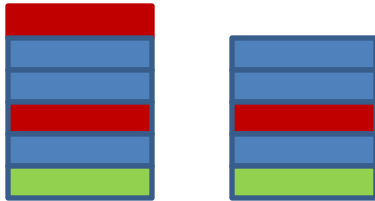
Outline

- Recap
 - Laplace Mechanism
- Composition Theorems
- Optimizing accuracy of DP algorithms
 - Utilizing Parallel Composition
 - Postprocessing & Inference
 - Strategy Selection
 - Data dependent noise

Differential Privacy

[Dwork ICALP 2006]

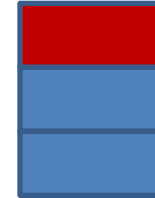
For every pair of inputs
that differ in one row



D_1

D_2

For every output ...

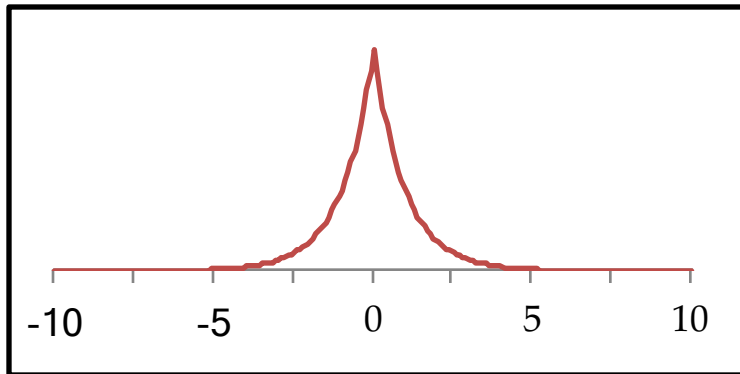
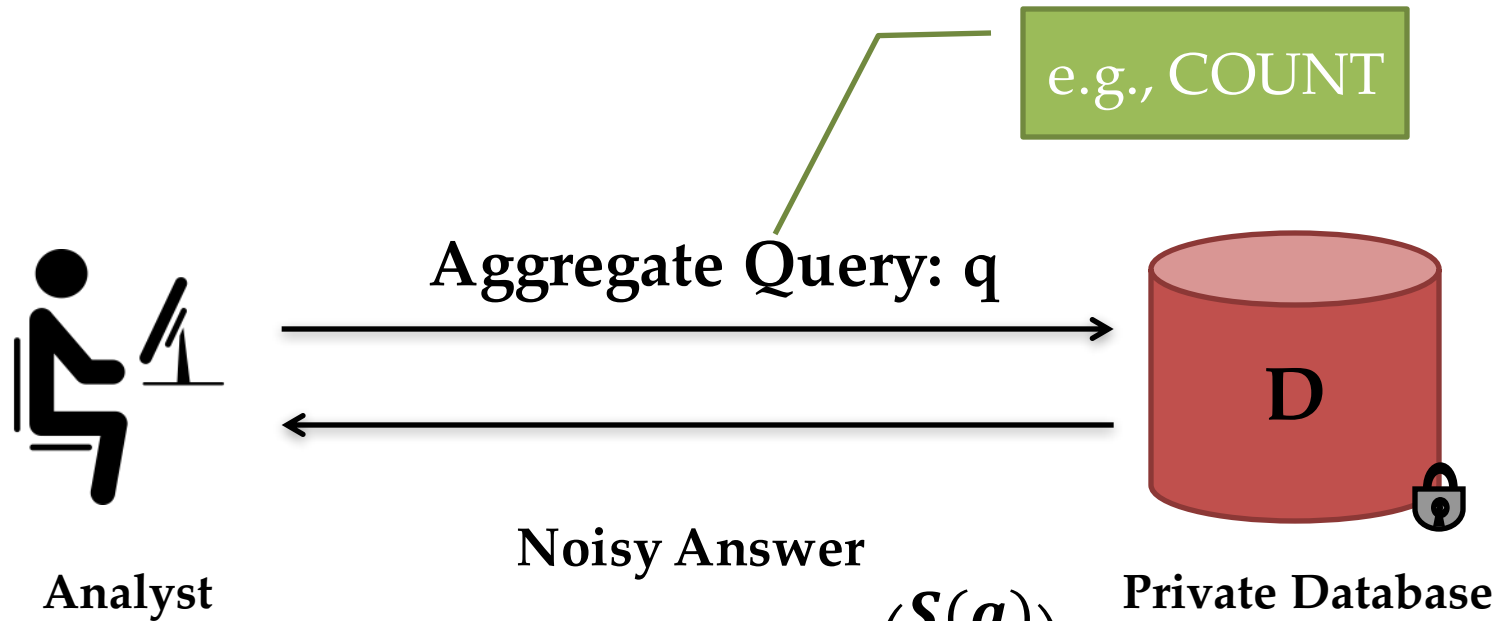


O

Adversary should not be able to distinguish
between any D_1 and D_2 based on any O

$$\forall \Omega \in \text{range}(A), \ln \left(\frac{\Pr[A(D_1) \in \Omega]}{\Pr[A(D_2) \in \Omega]} \right) \leq \varepsilon, \quad \varepsilon > 0$$

Laplace mechanism



Laplace Mechanism

Theorems:

$$E \left((\tilde{q}(D) - q(D))^2 \right) = 2 \left(\frac{S(q)}{\varepsilon} \right)^2$$

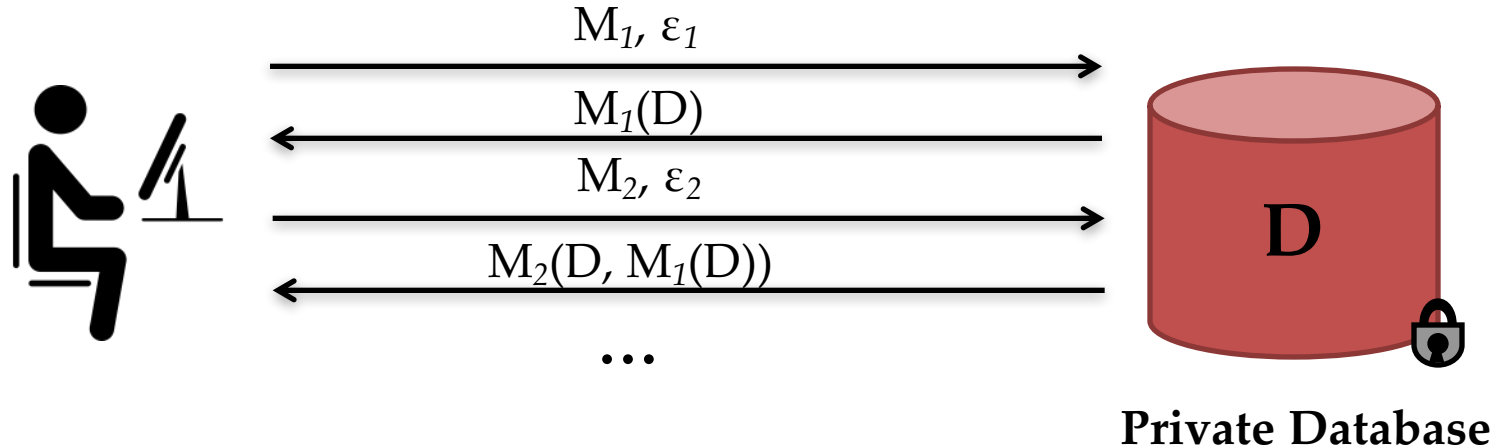
Error is *data independent*
Depends on q and ε , but not on D

$$Pr \left[|\tilde{q}(D) - q(D)| \geq \frac{S(q)}{\varepsilon} \ln \left(\frac{1}{\delta} \right) \right] \leq \delta$$

Outline

- Recap
 - Laplace Mechanism
- Composition Theorems
- Optimizing accuracy of DP algorithms
 - Utilizing Parallel Composition
 - Postprocessing & Inference
 - Strategy Selection
 - Data dependent noise

Sequential Composition

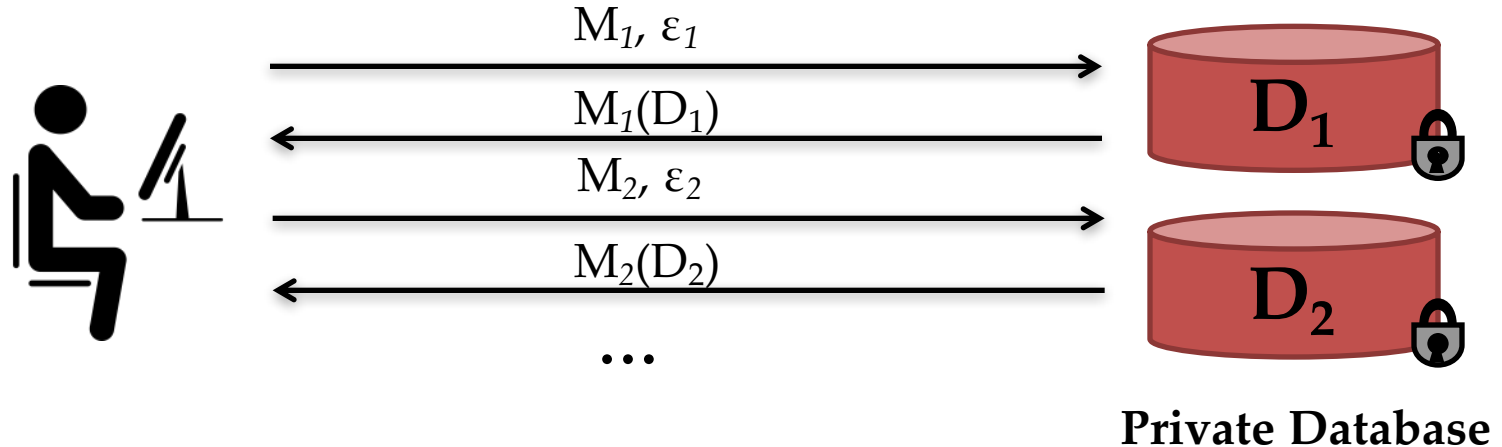


- If M_1, M_2, \dots, M_k are algorithms that access a private database D such that each M_i satisfies ϵ_i -differential privacy,

then the combination of their outputs satisfies ϵ -differential privacy with

$$\epsilon = \epsilon_1 + \dots + \epsilon_k$$

Parallel Composition

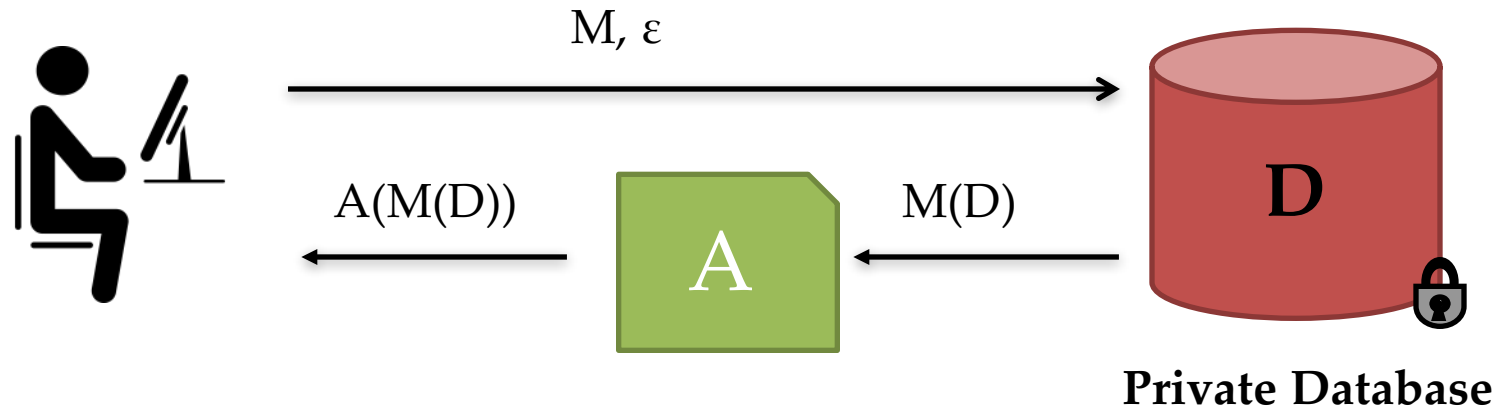


- If M_1, M_2, \dots, M_k are algorithms that access disjoint databases D_1, D_2, \dots, D_k such that each M_i satisfies ϵ_i -differential privacy,

then the combination of their outputs satisfies ϵ -differential privacy with

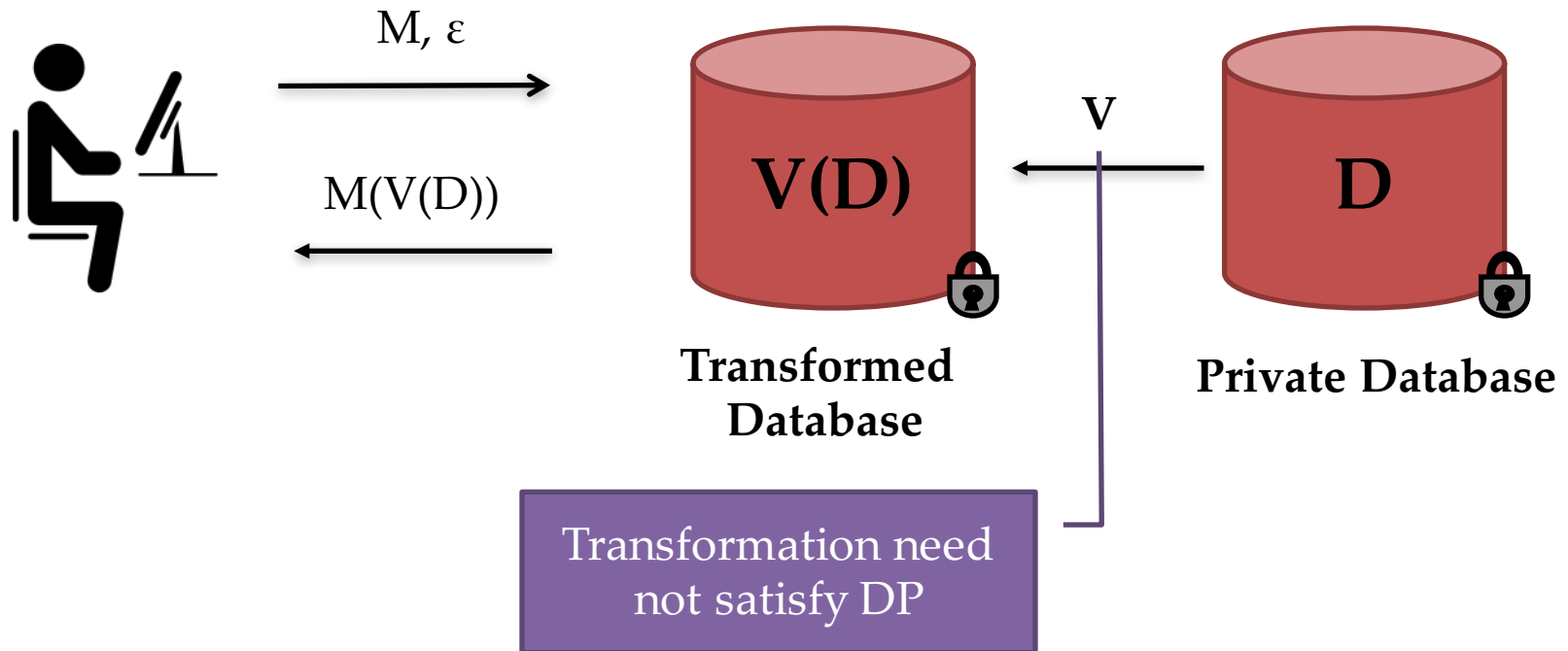
$$\epsilon = \max(\epsilon_1, \dots, \epsilon_k)$$

Postprocessing



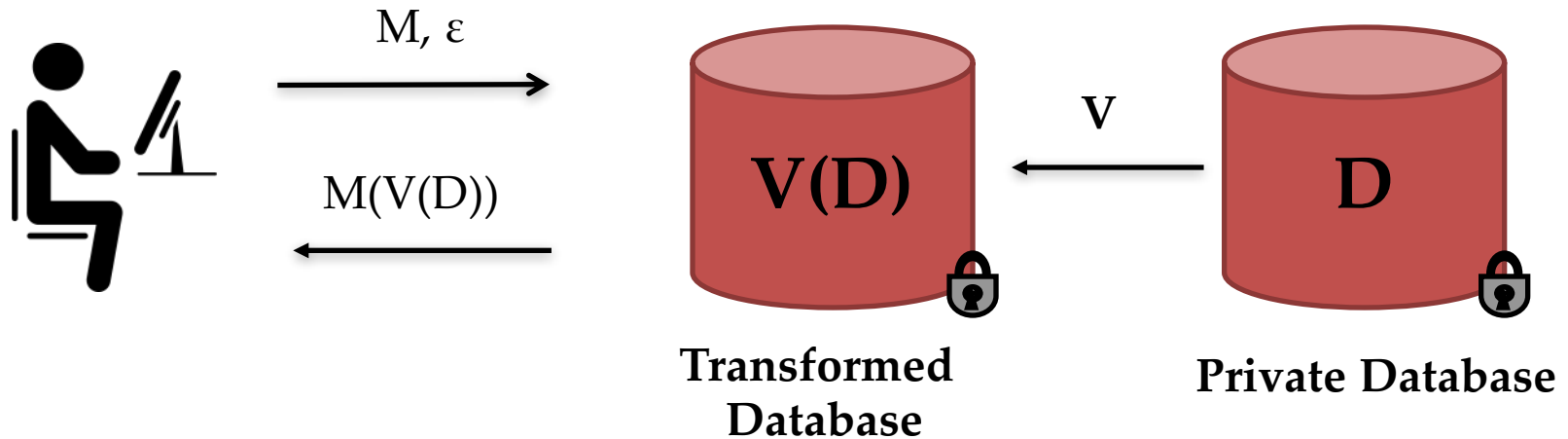
- If M is an ϵ -differentially private algorithm, any additional post-processing $A \circ M$ also satisfies ϵ -differential privacy.

Transformations & Stability



- σ_V : Stability of the transformation
 - Maximum number of rows in V that can change due to changing a single row in D

Transformations & Stability



- Executing an ϵ -differentially private algorithm M on a transformation of a database $V(D)$ satisfies $\epsilon \cdot \sigma_V$ -differential privacy.
- σ_V : Stability of the transformation
 - Maximum number of rows in V that can change due to changing a single row in D

Transformations & Stability

- V_1 : For each row $(x_1, x_2, x_3) \rightarrow (x_1, x_2+x_3)$

Stability = 1

- V_2 : Each row in D is a tweet $(id, \{words\})$. For each row in D , generate k rows with first k words $\{(id, word_1), \dots, (id, word_k)\}$

Stability = k

- V_3 : Sample each row with probability p .

Stability = 1 ... but can prove $2p\epsilon$ -differential privacy*

*Adam Smith, [Differential Privacy and Secrecy of the Sample](#)

Outline

- Recap
 - Laplace Mechanism
- Composition Theorems
- Optimizing accuracy of DP algorithms
 - Utilizing Parallel Composition
 - Postprocessing & Inference
 - Strategy Selection
 - Data dependent noise

Problem

Sex	Height	Weight
M	6'2"	210
F	5'3"	190
F	5'9"	160
M	5'3"	180
M	6'7"	250

Queries:

- # Males with BMI < 25
- # Males
- # Females with BMI < 25
- # Females

- Design an ϵ -differentially private algorithm that can answer all these questions.
- What is the total error?

Algorithm 1

Return:

- # Males with BMI < 25 + Lap($4/\epsilon$)
- # Males + Lap($4/\epsilon$)
- # Females with BMI < 25 + Lap($4/\epsilon$)
- # Females + Lap($4/\epsilon$)

Privacy

- BMI can be computed by transforming each row $(s, h, w) \rightarrow (s, \text{bmi})$. This is stability 1.
- Sensitivity of count = 1. So each query is answered using a $\epsilon/4$ -DP algorithm.
- By sequential composition, we get ϵ -DP.

Utility

Error:

$$\sum E \left((\tilde{q}(D) - q(D))^2 \right)$$

Total Error:

$$2 \left(\frac{4}{\varepsilon} \right)^2 \times 4 = \frac{128}{\varepsilon^2}$$

Algorithm 2

Compute:

- $\widetilde{q}_1 = \# \text{ Males with BMI} < 25 + \text{Lap}(1/\varepsilon)$
- $\widetilde{q}_2 = \# \text{ Males with BMI} > 25 + \text{Lap}(1/\varepsilon)$
- $\widetilde{q}_3 = \# \text{ Females with BMI} < 25 + \text{Lap}(1/\varepsilon)$
- $\widetilde{q}_4 = \# \text{ Females with BMI} > 25 + \text{Lap}(1/\varepsilon)$

Return

- $\widetilde{q}_1, \widetilde{q}_1 + \widetilde{q}_2, \widetilde{q}_3, \widetilde{q}_3 + \widetilde{q}_4$

Privacy

- Sensitivity of count = 1. So each query is answered using a ϵ -DP algorithm.
- q_1, q_2, q_3, q_4 are counts on disjoint portions of the database. Thus by *parallel composition* releasing $\widetilde{q}_1, \widetilde{q}_2, \widetilde{q}_3, \widetilde{q}_4$ satisfies ϵ -DP.
- By the *postprocessing theorem*, releasing $\widetilde{q}_1, \widetilde{q}_1 + \widetilde{q}_2, \widetilde{q}_3, \widetilde{q}_3 + \widetilde{q}_4$ also satisfies ϵ -DP.

Utility

Error:

$$\sum E \left((\tilde{q}(D) - q(D))^2 \right)$$

Total Error:

$$2 \left(\frac{1}{\varepsilon} \right)^2 + 2 \cdot 2 \left(\frac{1}{\varepsilon} \right)^2 + 2 \left(\frac{1}{\varepsilon} \right)^2 + 2 \cdot 2 \left(\frac{1}{\varepsilon} \right)^2 = \frac{12}{\varepsilon^2}$$

\widetilde{q}_1

$\widetilde{q}_1 + \widetilde{q}_2$

\widetilde{q}_3

$\widetilde{q}_3 + \widetilde{q}_4$

Utility

Tighter privacy analysis gives better accuracy for the same level of privacy

Total Error:

$$\underbrace{2 \left(\frac{1}{\varepsilon} \right)^2}_{\widetilde{q}_1} + \underbrace{2 \cdot 2 \left(\frac{1}{\varepsilon} \right)^2}_{\widetilde{q}_1 + \widetilde{q}_2} + \underbrace{2 \left(\frac{1}{\varepsilon} \right)^2}_{\widetilde{q}_3} + \underbrace{2 \cdot 2 \left(\frac{1}{\varepsilon} \right)^2}_{\widetilde{q}_3 + \widetilde{q}_4} = \frac{12}{\varepsilon^2}$$

Generalized Sensitivity

- Let $f: \mathcal{D} \rightarrow \mathbb{R}^d$ be a function that outputs a vector of d real numbers. The sensitivity of f is given by:

$$S(f) = \max_{D, D': |D \Delta D'|=1} \|f(D) - f(D')\|_1$$

where $\|\mathbf{x} - \mathbf{y}\|_1 = \sum_i |x_i - y_i|$

Generalized Sensitivity

- $q_1 = \# \text{ Males with BMI} < 25$
- $q_2 = \# \text{ Males with BMI} > 25$
- $q = \# \text{ Males with BMI}$
- Let f_1 be a function that answers both q_1, q_2
- Let f_2 be a function that answers both q_1, q
- Sensitivity of $f_1 = 1$
- Sensitivity of $f_2 = 2$
- An alternate privacy proof for Alg 2 is to show that the generalized sensitivity of $\widetilde{q}_1, \widetilde{q}_2, \widetilde{q}_3, \widetilde{q}_4$ is 1.

Outline

- Recap
 - Laplace Mechanism
- Composition Theorems
- Optimizing accuracy of DP algorithms
 - Utilizing Parallel Composition
 - Postprocessing & Inference
 - Strategy Selection
 - Data dependent noise

Improving utility of Alg 2

Compute:

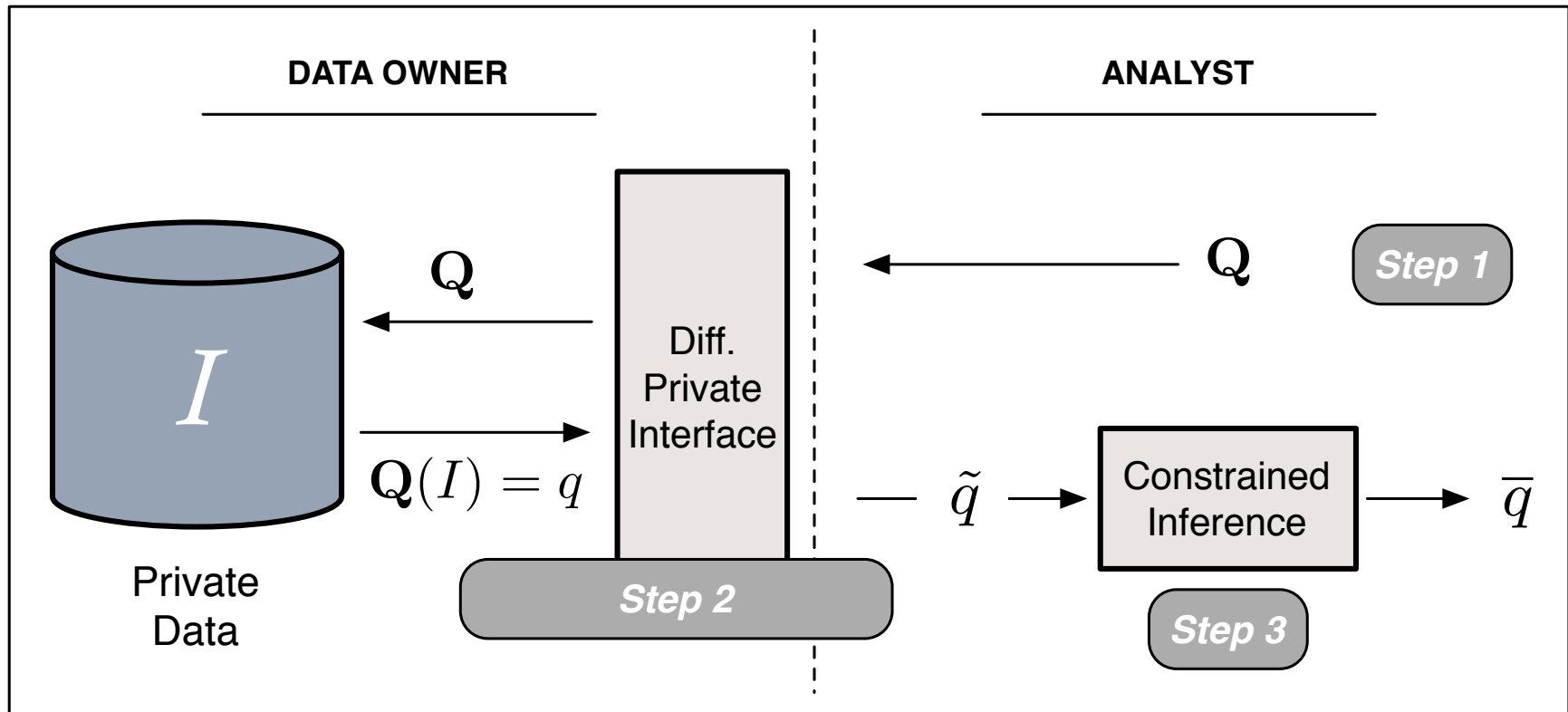
- $\widetilde{q}_1 = \# \text{ Males with BMI} < 25 + \text{Lap}(1/\varepsilon)$
- $\widetilde{q}_2 = \# \text{ Males with BMI} > 25 + \text{Lap}(1/\varepsilon)$

Return

- $\widetilde{q}_1, \widetilde{q}_1 + \widetilde{q}_2$

We know $q_1 \leq q_1 + q_2$,
but $P[\widetilde{q}_1 > \widetilde{q}_1 + \widetilde{q}_2] > 0$

Constrained Inference



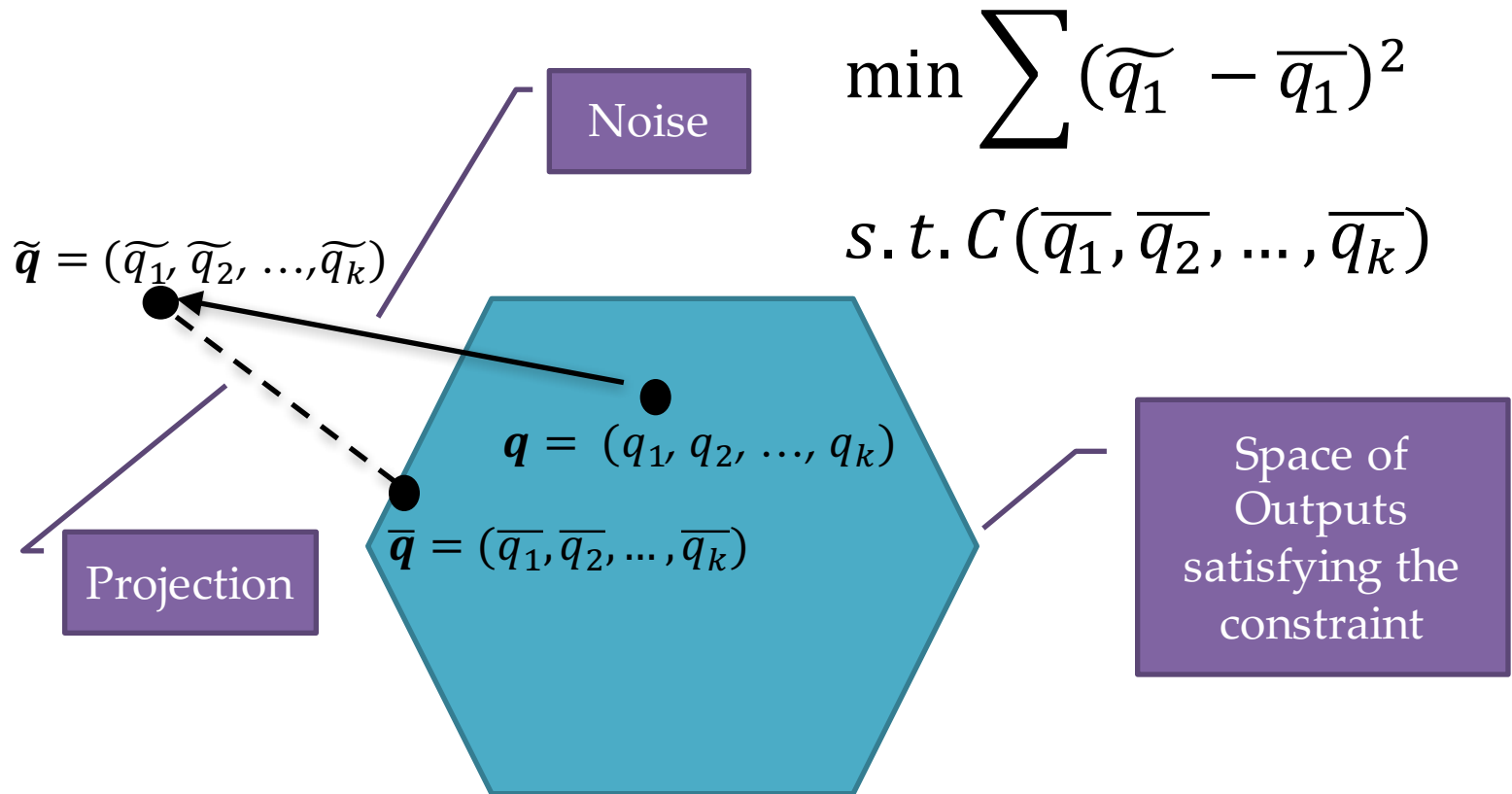
Constrained Inference

- q_1, q_2, \dots, q_k be a set of queries
- $\widetilde{q}_1, \widetilde{q}_2, \dots, \widetilde{q}_k$ be the noisy answers
- Constraint $C(q_1, q_2, \dots, q_k) = 1$ holds on true answers (for all typical databases), but does not hold on noisy answers.
- Goal: Find $\overline{q}_1, \overline{q}_2, \dots, \overline{q}_k$ that are:
 - Close to $\widetilde{q}_1, \widetilde{q}_2, \dots, \widetilde{q}_k$
 - Satisfy the constraint $C(\overline{q}_1, \overline{q}_2, \dots, \overline{q}_k)$

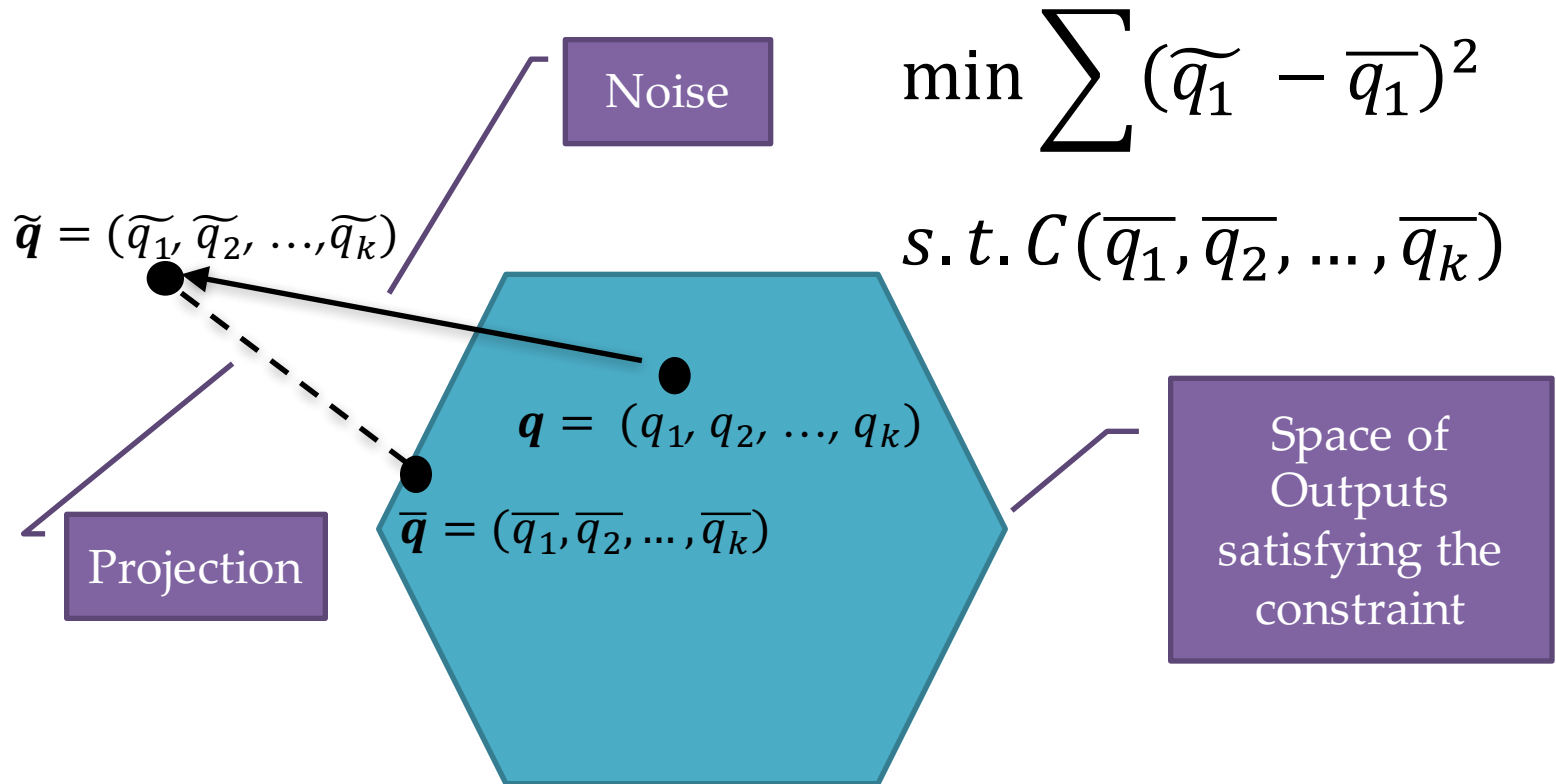
Least Squares Optimization

$$\begin{aligned} \min \quad & \sum (\widetilde{q_1} - \overline{q_1})^2 \\ \text{s. t. } & \mathcal{C}(\overline{q_1}, \overline{q_2}, \dots, \overline{q_k}) \end{aligned}$$

Geometric Interpretation



Geometric Interpretation



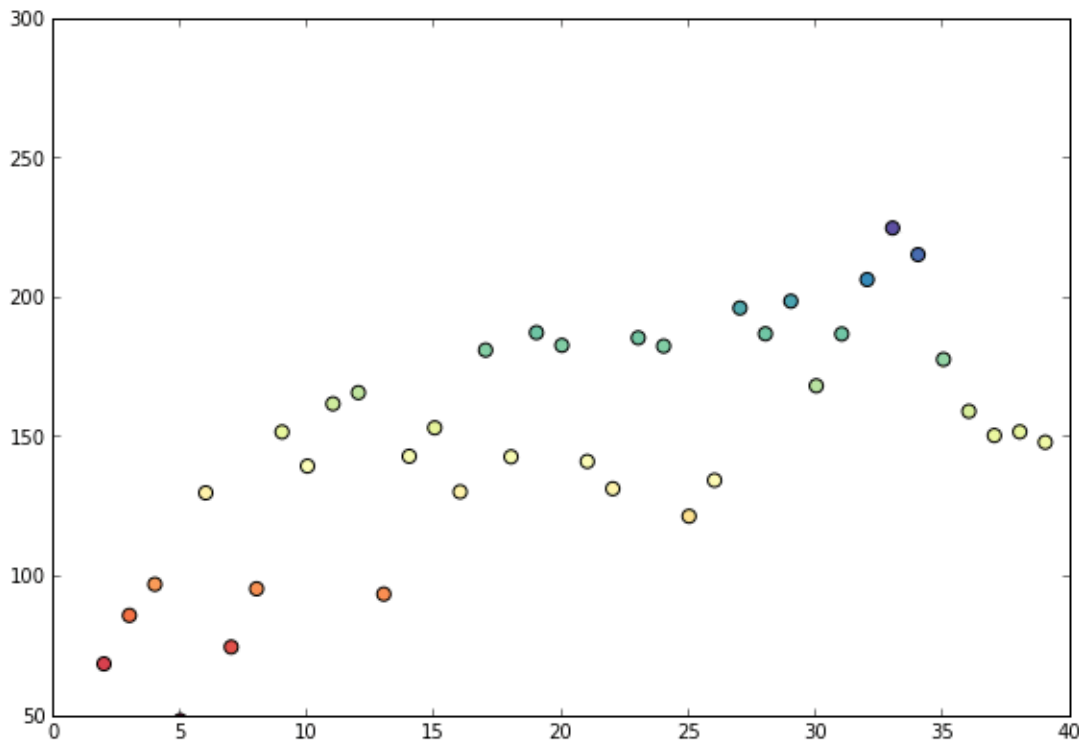
Theorem: $\|\mathbf{q} - \bar{\mathbf{q}}\|_2 \leq \|\mathbf{q} - \tilde{\mathbf{q}}\|_2$ when the constraints form a convex space

Ordering Constraint

Isotonic Regression:

$$\min \sum (\widetilde{q}_1 - \overline{q}_1)^2$$

$$s.t. \overline{q}_1 \leq \overline{q}_1 \leq \dots \leq \overline{q}_k$$



Outline

- Recap
 - Laplace Mechanism
- Composition Theorems
- Optimizing accuracy of DP algorithms
 - Utilizing Parallel Composition
 - Postprocessing & Inference
 - Strategy Selection
 - Data dependent noise

Problem

Sex	Height	Weight
M	6'2"	210
F	5'3"	190
F	5'9"	160
M	5'3"	180
M	6'7"	250

Queries:

- # people with height in [5'1", 6'2"]
- # people with height in [2'0", 4'0"]
- # people with height in [3'3", 7'0"]
- ...

- Design an ϵ -differentially private algorithm that can answer all range queries.
- What is the total error?

Problem

- Let $\{v_1, \dots, v_k\}$ be the domain of an attribute
- Let $\{x_1, \dots, x_k\}$ be the number of rows with values v_1, \dots, v_k
- Range Query: $q_{ij} = x_i + x_{i+1} + \dots + x_j$
- Goal: Answer all range queries

Strategy 1:

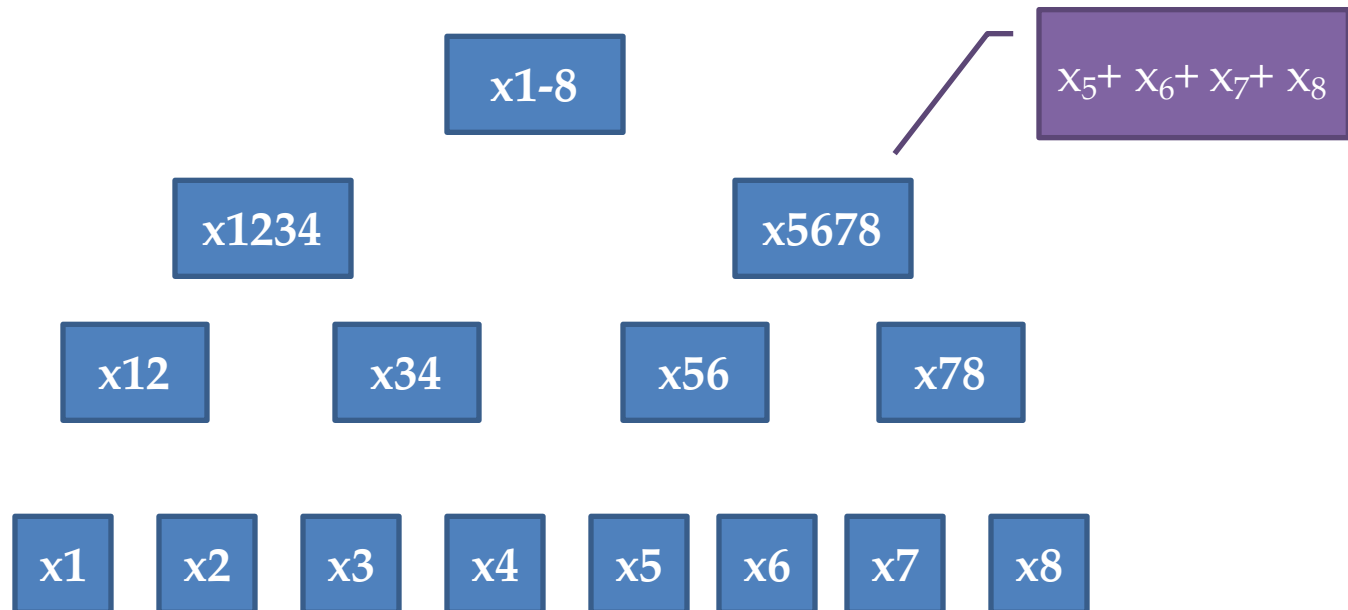
- Answer all range queries using Laplace mechanism
- Sensitivity: $O(k^2)$
- Total Error: $O(k^4/\epsilon^2)$

Strategy 2:

- Estimate each individual x_i using Laplace mechanism
- Answer: $q_{ij} = \tilde{x}_i + \tilde{x}_{i+1} + \dots + \tilde{x}_j$
- Error in each \tilde{x}_i : $O(1/\varepsilon^2)$
- Error in q_{1k} : $O(k/\varepsilon^2)$
- Total Error: $O(k^3/\varepsilon^2)$

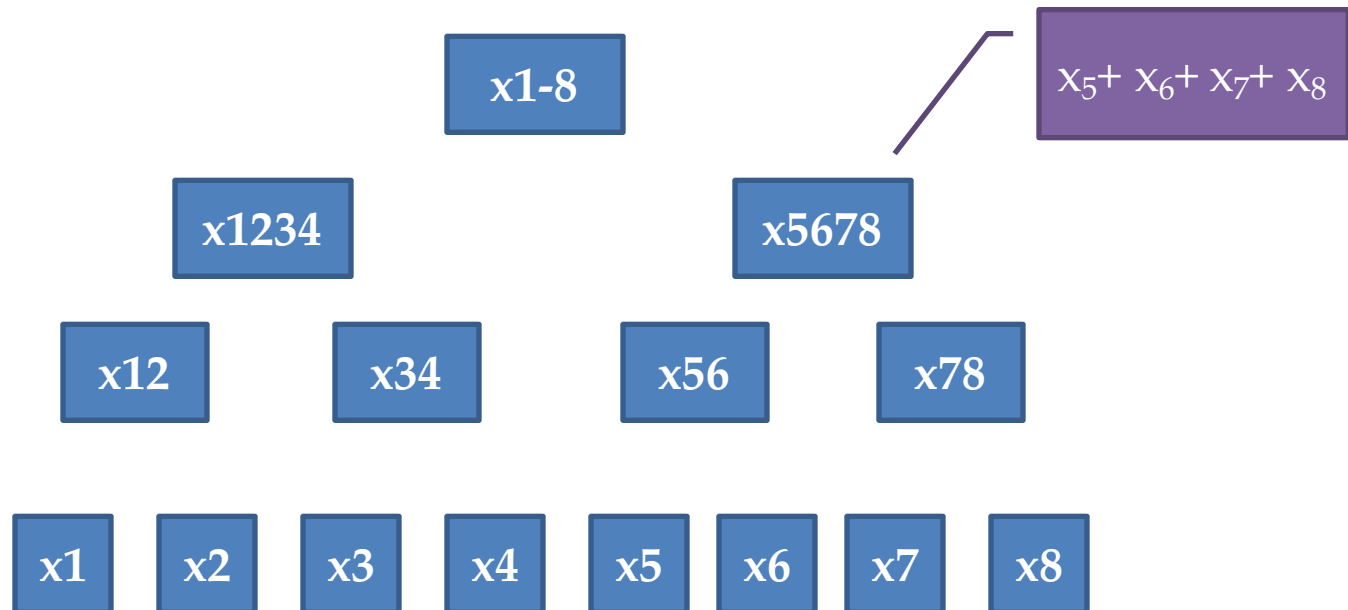
Strategy 3: Hierarchy

- Estimate all the counts in the tree below using Laplace mechanism



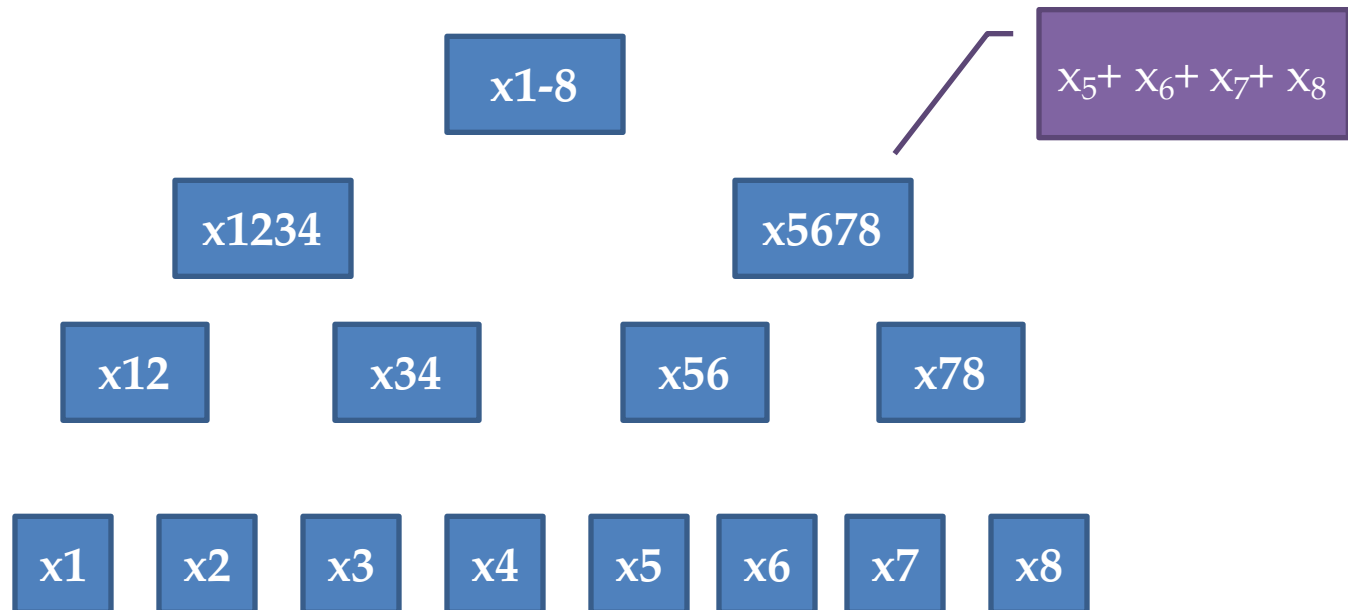
Strategy 3: Hierarchy

- Sensitivity: $\log k$
- Every range query can be answered by summing up at most $2 \log k$ nodes in the tree.



Strategy 3: Hierarchy

- Error in each node: $O((\log k)^2 / \varepsilon^2)$
- Max error on a range query: $O((\log k)^3 / \varepsilon^2)$
- Total Error: $O(k^2 (\log k)^3 / \varepsilon^2)$



Strategy 3: Hierarchy

- Error in each node: $O((\log k)^2 / \varepsilon^2)$
- Max error on a range query: $O((\log k)^3 / \varepsilon^2)$
- Total Error: $O(k^2 (\log k)^3 / \varepsilon^2)$
- Error can be further reduced using constrained inference
 - Here the constraint is that parent counts should not be smaller than child counts.

Strategy based mechanisms

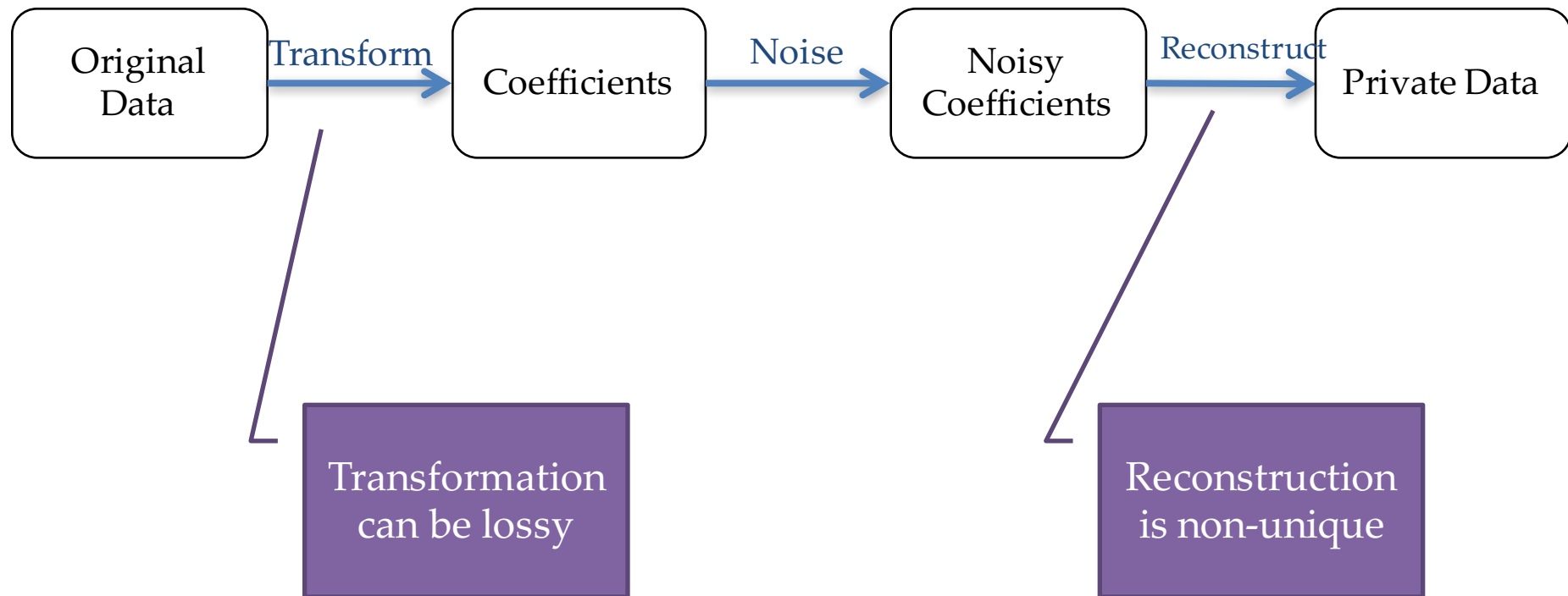


- Can think of nodes in the tree as coefficients.
- Other algorithms use other transformations
 - Wavelets, Fourier coefficients
- Should be able to *losslessly* reconstruct the original data/query answers.
- **General Idea:**
 - Apply transform
 - Add noise to the transformed space (based on sensitivity)
 - Reconstruct original data/query answers from noisy coefficients

Outline

- Recap
 - Laplace Mechanism
- Composition Theorems
- Optimizing accuracy of DP algorithms
 - Utilizing Parallel Composition
 - Postprocessing & Inference
 - Strategy Selection
 - Data dependent noise

Data dependent noise mechanisms



[LHMY14] Li et al. A data- and workload-aware algorithm for range queries under differential privacy. In PVLDB, 2014.

Data dependent noise mechanisms

- Use a data dependent sensitivity measure called Smooth sensitivity.

K. Nissim, S. Raskhodnikova, A. Smith, “Smooth Sensitivity and sampling in private data analysis”, STOC 2007

Summary

- Composition theorems help build complex algorithms using simple building blocks
 - Sequential composition
 - Parallel composition
 - Postprocessing
 - *There are more advanced forms of composition.*

Summary

- For the same privacy budget, a better designed algorithm can extract more utility
 - When possible use parallel composition
 - Inference on constraints between queries can reduce error
 - Answering a different *strategy* of queries can help reduce error