

Post-processing outputs for better utility

CompSci 590.03

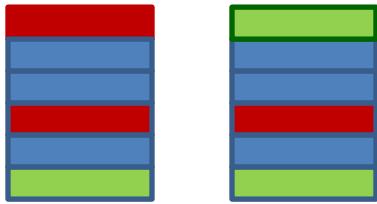
Instructor: Ashwin Machanavajjhala

Announcement

- Project proposal submission deadline is **Fri, Oct 12 noon.**

Recap: Differential Privacy

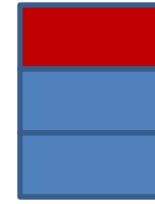
For every pair of inputs that differ in one value



D_1

D_2

For every output ...

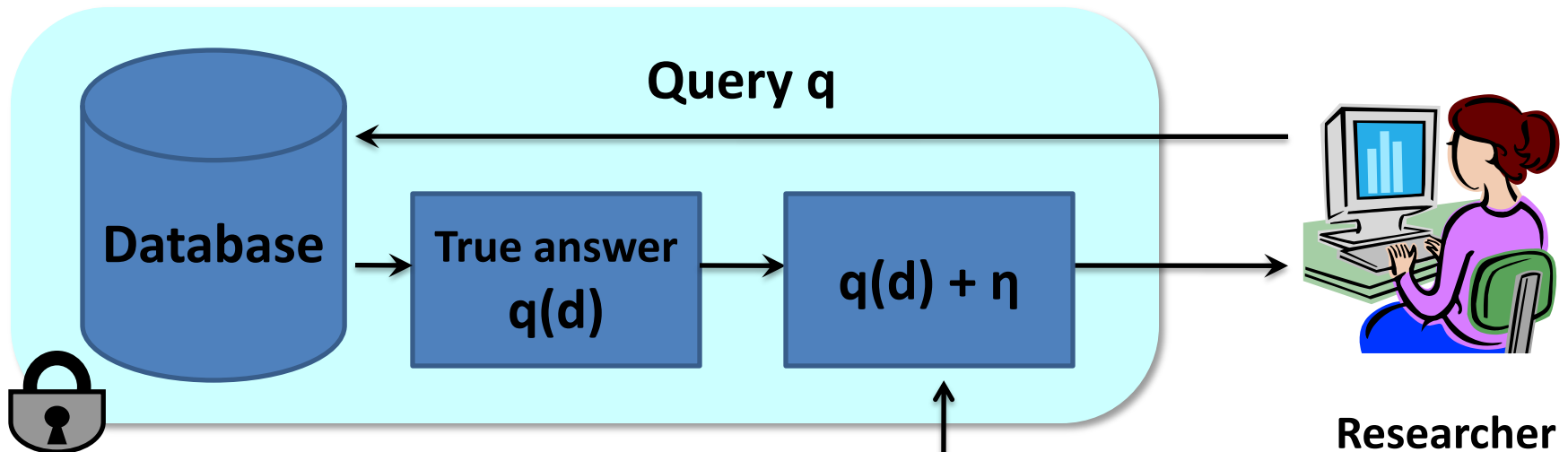


O

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

$$\log \left(\frac{\Pr[A(D_1) = O]}{\Pr[A(D_2) = O]} \right) < \epsilon \quad (\epsilon > 0)$$

Recap: Laplacian Distribution



Privacy depends on the λ parameter

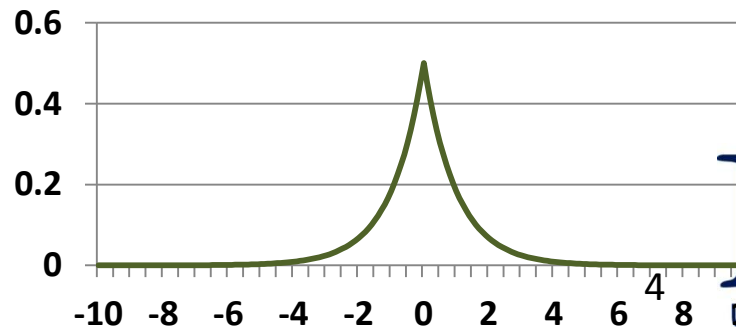
$$h(\eta) \propto \exp(-\eta / \lambda)$$

Mean: 0,

Variance: $2 \lambda^2$

Lecture 10 : 590.03 Fall 12

Laplace Distribution – Lap(λ)



Recap: Laplace Mechanism

Thm: If **sensitivity** of the query is **S**, then the following guarantees ϵ -differential privacy.

$$\lambda = S/\epsilon$$

Sensitivity: Smallest number s.t. for any d, d' differing in one entry,

$$|| q(d) - q(d') || \leq S(q)$$

Histogram query: Sensitivity = 2

- Variance / error on each entry = $2 \times 4 / \epsilon^2 = \mathbf{O(1/\epsilon^2)}$

This class

- What is the optimal method to answer a batch of queries?

How to answer a batch of queries?

- Database of values $\{x_1, x_2, \dots, x_k\}$
- Query Set:
 - Value of x_1 $\eta_1 = x_1 + \delta_1$
 - Value of x_2 $\eta_2 = x_2 + \delta_2$
 - Value of $x_1 + x_2$ $\eta_3 = x_1 + x_2 + \delta_3$
- But we know that η_1 and η_2 should sum up to η_3 !

Two Approaches

- **Constrained inference**
 - Ensure that the returned answers are consistent with each other.

- **Query Strategy**
 - Answer a different set of *strategy* queries A
 - Answer original queries using A

 - **Universal Histograms**
 - **Wavelet Mechanism**
 - **Matrix Mechanism**

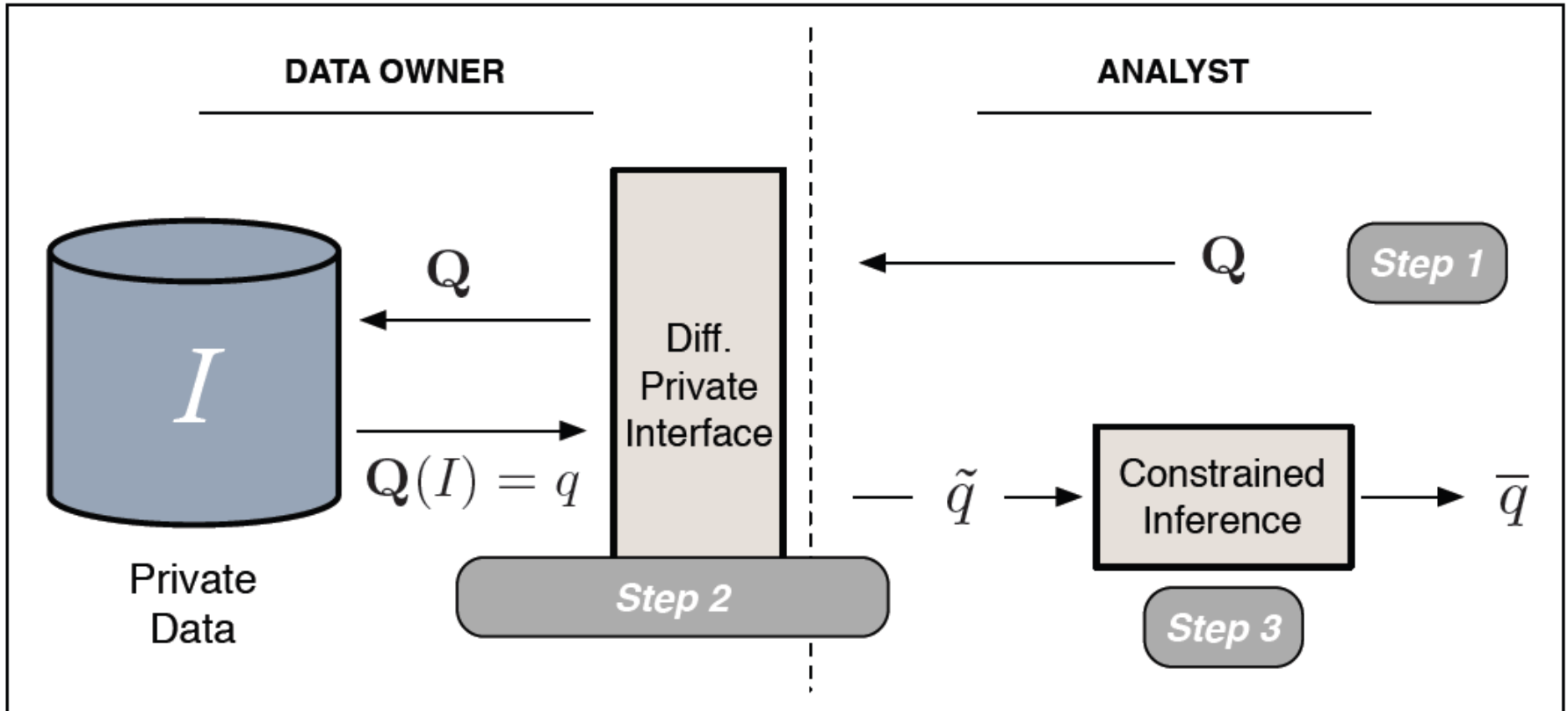
Two Approaches

- **Constrained inference**
 - Ensure that the returned answers are consistent with each other.

- **Query Strategy**
 - Answer a different set of *strategy* queries A
 - Answer original queries using A

 - **Universal Histograms**
 - **Wavelet Mechanism**
 - **Matrix Mechanism**

Constrained Inference



Constrained Inference

- Let x_1 and x_2 be the original values. We observe noisy values η_1 , η_2 and η_3
- We would like to reconstruct the best estimators y_1 (for x_1) and y_2 (for x_2) from the noisy values.
- That is, we want to find the values of y_1 , y_2 such that:

$$\begin{aligned} & \min (y_1 - \eta_1)^2 + (y_2 - \eta_2)^2 + (y_3 - \eta_3)^2 \\ & \text{s.t., } y_1 + y_2 = y_3 \end{aligned}$$

Constrained Inference [Hay et al VLDB 10]

DEFINITION 2.4 (MINIMUM L_2 SOLUTION). *Let \mathbf{Q} be a query sequence with constraints $\gamma_{\mathbf{Q}}$. Given a noisy query sequence $\tilde{q} = \tilde{\mathbf{Q}}(I)$, a minimum L_2 solution, denoted \bar{q} , is a vector \bar{q} that satisfies the constraints $\gamma_{\mathbf{Q}}$ and at the same time minimizes $\|\tilde{q} - \bar{q}\|_2$.*

Sorted Unattributed Histograms

- Counts of diseases
 - (without associating a particular count to the corresponding disease)
- Degree sequence: List of node degrees
 - (without associating a degree to a particular node)
- Constraint: The values are sorted

Sorted Unattributed Histograms

True Values 20, 10, 8, 8, 8, 5, 3, 2

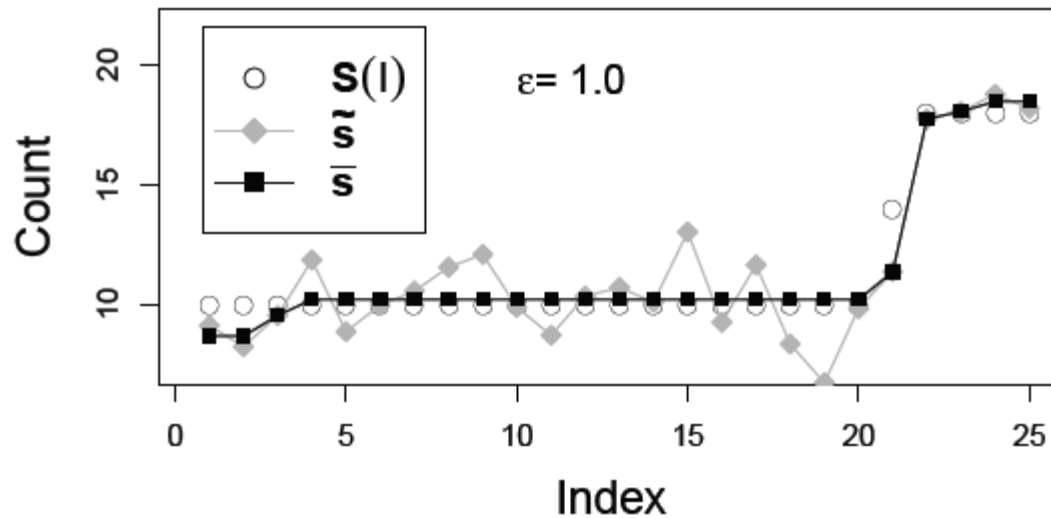
Noisy Values 25, 9, 13, 7, 10, 6, 3, 1 (noise from Lap(1/ε))

$$\text{minimize } \|\tilde{s} - \bar{s}\|_2 \quad \text{s.t. } \forall i, \bar{s}_i \leq \bar{s}_{i-1}$$

$$\bar{s}_k = \min_{j \in [k, n]} \max_{i \in [1, j]} \frac{(\tilde{s}_i + \tilde{s}_{i+1} + \dots + \tilde{s}_j)}{j - i + 1}$$

Proof:?

Sorted Unattributed Histograms



Sorted Unattributed Histograms

- n : number of values in the histogram
- d : number of distinct values in the histogram
- n_i : number of times i^{th} distinct value appears in the histogram.

THEOREM 2. *There exist constants c_1 and c_2 independent of n and d such that*

$$\text{error}(\bar{\mathbf{S}}) \leq \sum_{i=1}^d \frac{c_1 \log^3 n_i + c_2}{\epsilon^2}$$

Thus $\text{error}(\bar{\mathbf{S}}) = O(d \log^3 n / \epsilon^2)$ whereas $\text{error}(\tilde{\mathbf{S}}) = \Theta(n / \epsilon^2)$.

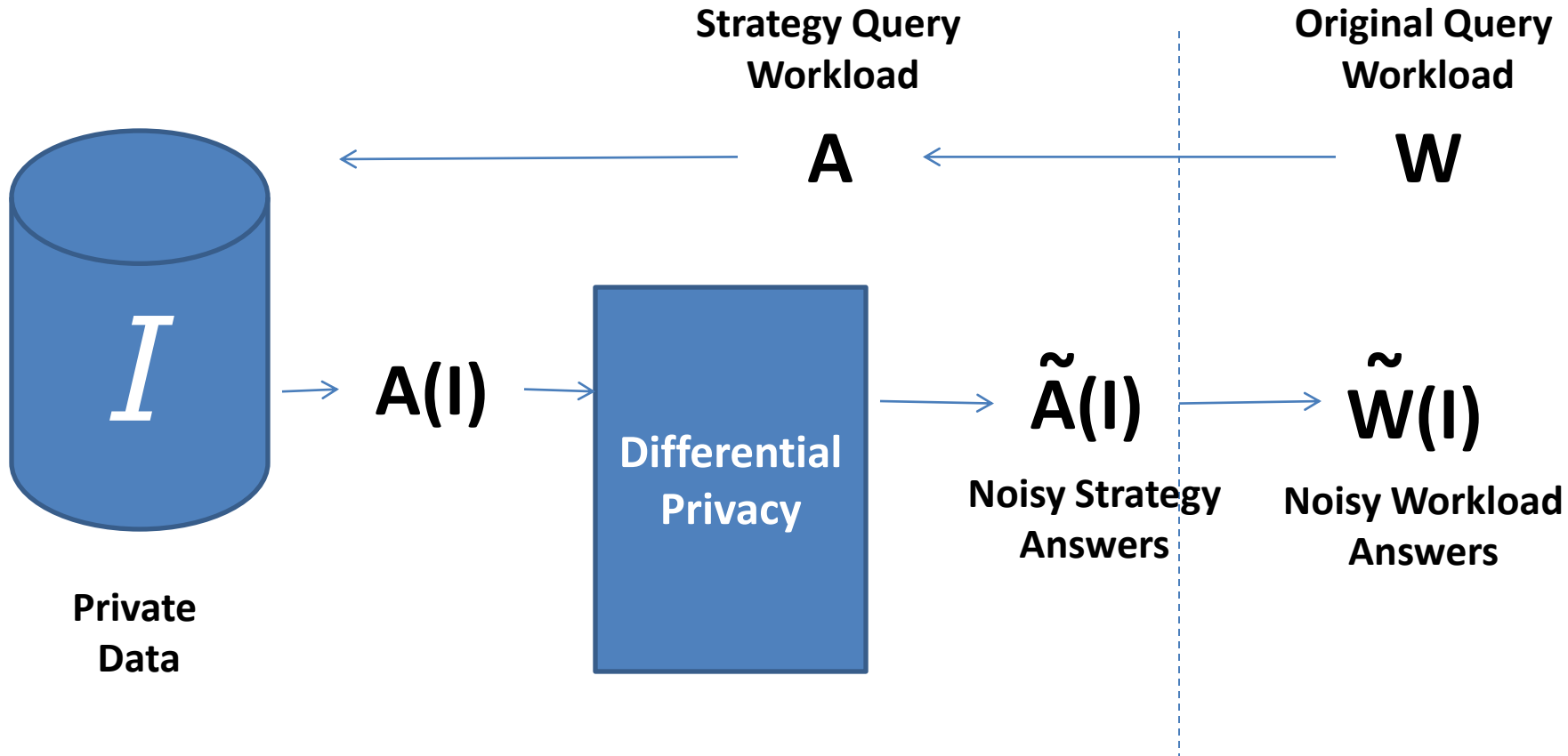
Two Approaches

- **Constrained inference**
 - Ensure that the returned answers are consistent with each other.

- **Query Strategy**
 - Answer a different set of *strategy* queries A
 - Answer original queries using A

 - **Universal Histograms**
 - **Wavelet Mechanism**
 - **Matrix Mechanism**

Query Strategy



Range Queries

- Given a set of values $\{x_1, x_2, \dots, x_n\}$
- Range query: $q(j,k) = x_j + \dots + x_k$

Q: Suppose we want to answer all range queries?

Strategy 1: Answer all range queries using Laplace mechanism

- $O(n^2/\epsilon^2)$ total error.
- May reduce using constrained optimization ...

Range Queries

- Given a set of values $\{x_1, x_2, \dots, x_n\}$
- Range query: $q(j,k) = x_j + \dots + x_k$

Q: Suppose we want to answer all range queries?

Strategy 1: Answer all range queries using Laplace mechanism

- Sensitivity = $O(n^2)$
- $O(n^4/\epsilon^2)$ total error across all range queries.
- May reduce using constrained optimization ...

Range Queries

- Given a set of values $\{x_1, x_2, \dots, x_n\}$
- Range query: $q(j,k) = x_j + \dots + x_k$

Q: Suppose we want to answer all range queries?

Strategy 2: Answer all x_i queries using Laplace mechanism
Answer range queries using noisy x_i values.

- $O(1/\epsilon^2)$ error for each x_i .
- $\text{Error}(q(1,n)) = O(n/\epsilon^2)$
- Total error on all range queries : $O(n^3/\epsilon^2)$

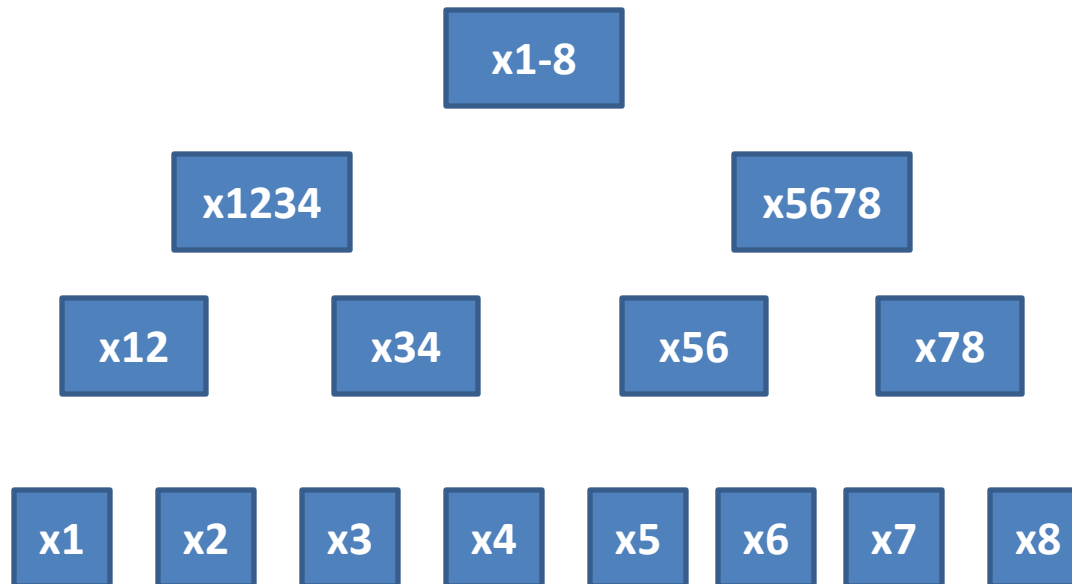
Universal Histograms for Range Queries

[Hay et al VLDB 2010]

Strategy 3:

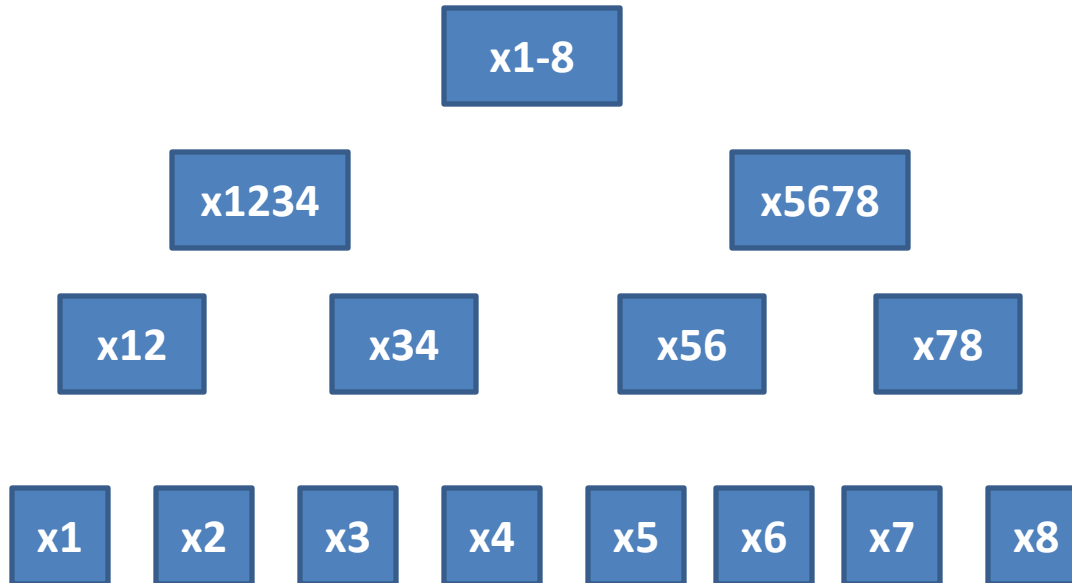
Answer *sufficient statistics* using Laplace mechanism

Answer range queries using noisy sufficient statistics.



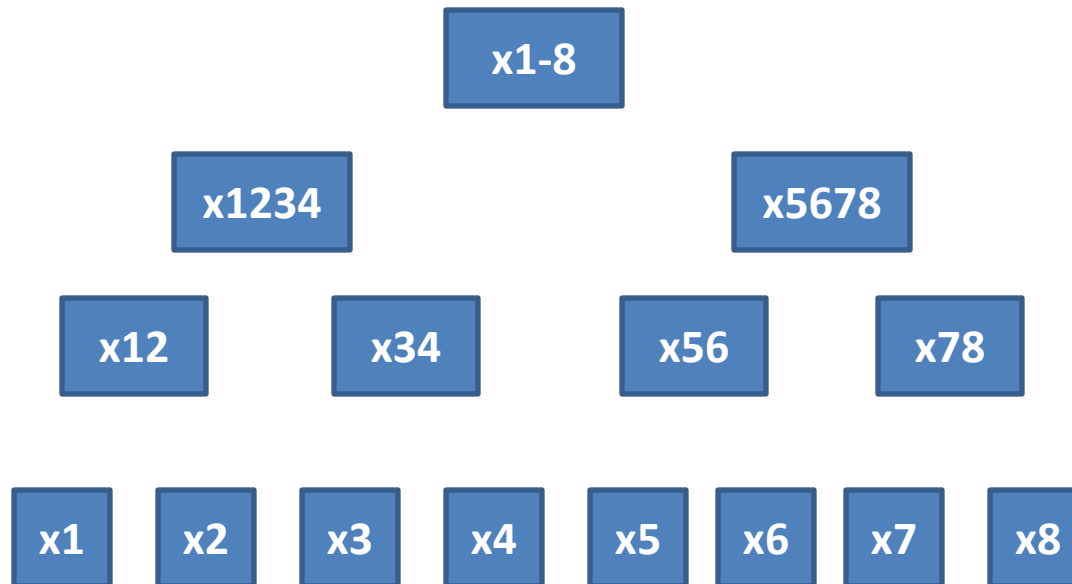
Universal Histograms for Range Queries

- Sensitivity: $\log n$
 - $q(2,6) = x_2 + x_3 + x_4 + x_5 + x_6$
 $= x_2 + x_{34} + x_{56}$
- Error = $2 \times 5 \log^2 n / \epsilon^2$
Error = $2 \times 3 \log^2 n / \epsilon^2$



Universal Histograms for Range Queries

- Every range query can be answered by summing at most $\log n$ different noisy answers
- Maximum error on any range query = $O(\log^3 n / \epsilon^2)$
- Total error on all range queries = $O(n^2 \log^3 n / \epsilon^2)$



Universal Histograms & Constrained Inference

[Hay et al VLDB 2010]

- Can further reduce the error by enforcing constraints
 $x_{1234} = x_{12} + x_{34} = x_1 + x_2 + x_3 + x_4$

$$\begin{aligned} & \textit{minimize} \sum_v (\tilde{c}(v) - \bar{c}(v))^2 \\ & \textit{s.t.} \bar{c}(v) = \sum_{u \in \textit{child}(v)} \bar{c}(u) \end{aligned}$$

- 2-pass algorithm to compute a consistent version of the counts

Universal Histograms & Constrained Inference

[Hay et al VLDB 2010]

- Pass 1: (Bottom Up)

$$z(v) = \begin{cases} \tilde{c}(v), & \text{if leaf node} \\ \alpha \cdot c(v) + (1 - \alpha) \sum_{u=\text{child}(v)} z(u) \end{cases}$$

- Pass 2: (Top down)

$$\bar{c}(v) = \begin{cases} z(v), & \text{if root node} \\ z(v) + \frac{1}{2} \left(\bar{c}(v) - \sum_{u=\text{child}(v)} z(u) \right) \end{cases}$$

Universal Histograms & Constrained Inference

- Resulting consistent counts
 - Have lower error than noisy counts (upto 10 times smaller in some cases)
 - Unbiased estimators
 - Have the least error amongst all unbiased estimators

Next Class

- **Constrained inference**
 - Ensure that the returned answers are consistent with each other.

- **Query Strategy**
 - Answer a different set of ***strategy*** queries A
 - Answer original queries using A

 - **Universal Histograms**
 - **Wavelet Mechanism**
 - **Matrix Mechanism**