

CSE660

Differential Privacy

October 4, 2017

Marco Gaboardi

Room: 338-B

gaboardi@buffalo.edu

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

(ϵ, δ) -Differential Privacy

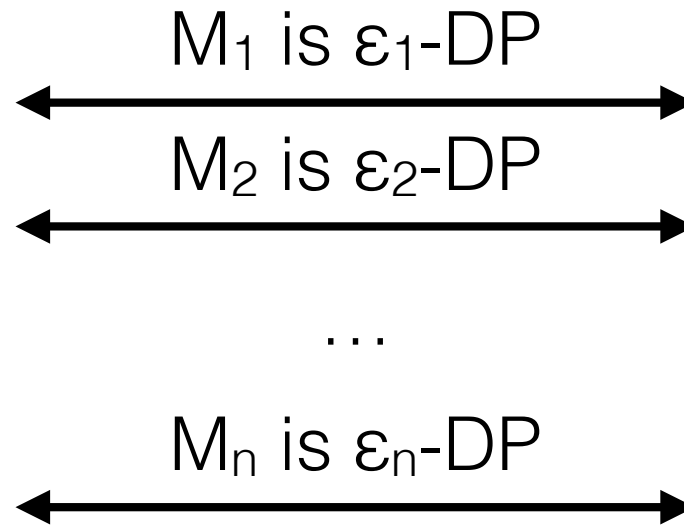
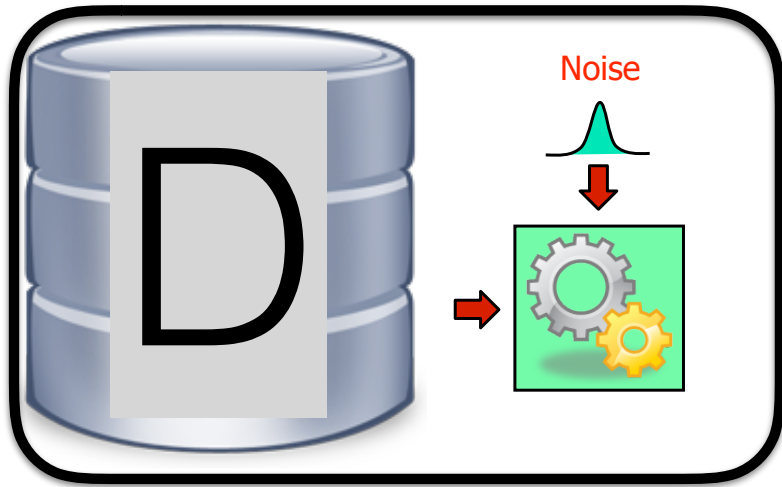
Definition

Given $\epsilon, \delta \geq 0$, a probabilistic query $Q: X^n \rightarrow R$ is (ϵ, δ) -differentially private iff

for all adjacent database b_1, b_2 and for every $S \subseteq R$:

$$\Pr[Q(b_1) \in S] \leq \exp(\epsilon) \Pr[Q(b_2) \in S] + \delta$$

Composition



The overall process is $(\epsilon_1 + \epsilon_2 + \dots + \epsilon_n)$ -DP

Multiple queries

Question: how much perturbation do we have if we want to answer n counting queries with Laplace under ϵ -DP?

Multiple queries

Question: how much perturbation do we have if we want to answer n counting queries with Laplace under ϵ -DP?

We can split the privacy budget uniformly:

$$\epsilon = \frac{\epsilon_{\text{global}}}{n}$$

Laplace accuracy: with high probability we have:

$$\left| q(D) - r \right| \leq O\left(\frac{1}{\epsilon n}\right)$$

Multiple queries

Question: how much perturbation do we have if we want to answer n counting queries with Laplace under ϵ -DP?

By putting them together (hiding some details) we have as a max error

$$O\left(\frac{n}{\epsilon_{\text{global}}n}\right) = O\left(\frac{1}{\epsilon_{\text{global}}}\right)$$

Notice that if we don't renormalize this is of the order of

$$O\left(\frac{n}{\epsilon_{\text{global}}}\right)$$

bigger than the sample error.

Advanced Composition

Question: how much perturbation do we have if we want to answer n queries under (ϵ, δ) -DP?

We have (by hiding many details) as a max error

$$O\left(\frac{1}{\epsilon_{\text{global}} \sqrt{n}}\right)$$

If we don't renormalize this is of the order of

$$O\left(\frac{\sqrt{n}}{\epsilon_{\text{global}}}\right)$$

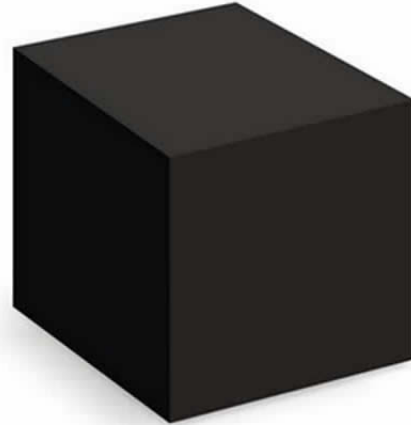
comparable to the sample error.

[DworkRothblumVadhan 10, SteinkeUllman 16]

Multiple queries

Question: Can we do better?

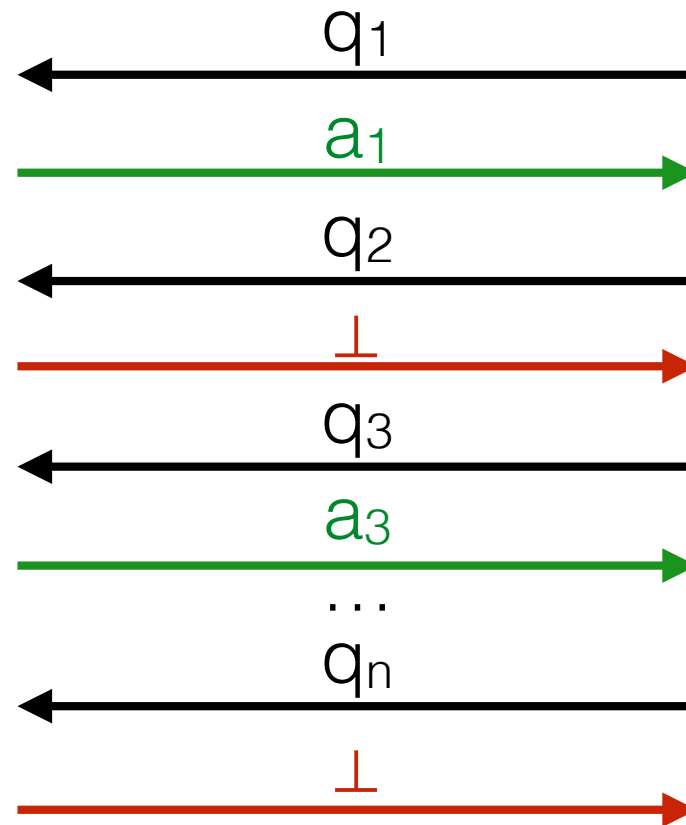
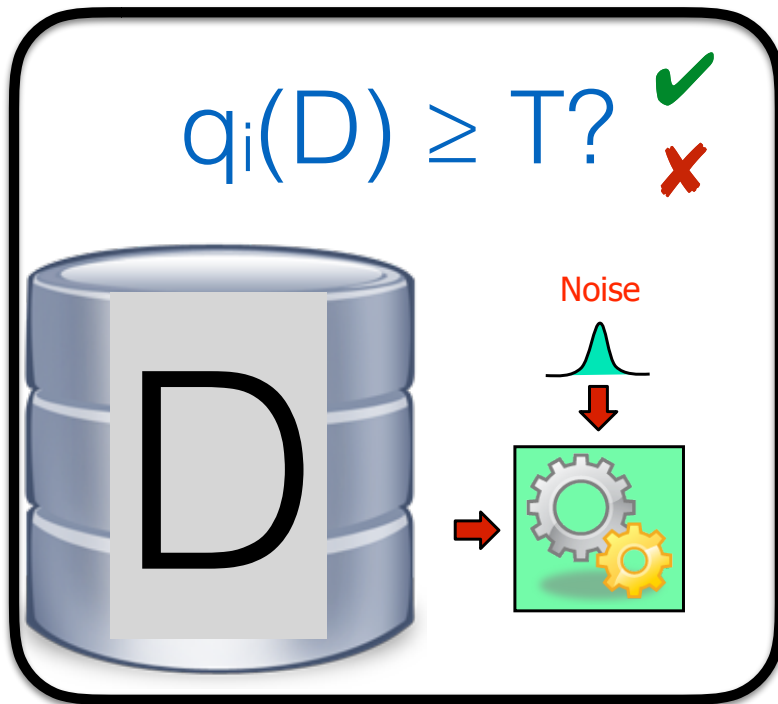
Composition



We always need to think before applying composition to whether we have other options!

Sparse vector

$\text{SparseVector}(D, q_1, \dots, q_n, T, \epsilon)$



How can we achieve epsilon-DP by paying only for the queries above T ?

An example: above threshold

Algorithm 1 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , and a threshold T . Output is a stream of responses a_1, \dots

AboveThreshold($D, \{f_i\}, T, \epsilon$)

Let $\hat{T} = T + \text{Lap}\left(\frac{2}{\epsilon}\right)$.

for Each query i do

Let $\nu_i = \text{Lap}\left(\frac{4}{\epsilon}\right)$

if $f_i(D) + \nu_i \geq \hat{T}$ then

Output $a_i = \top$.

Halt.

else

Output $a_i = \perp$.

end if

end for

An example: above threshold

Algorithm 1 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , and a threshold T . Output is a stream of responses a_1, \dots

AboveThreshold($D, \{f_i\}, T, \epsilon$)

Let $\hat{T} = T + \text{Lap}\left(\frac{2}{\epsilon}\right)$.

for Each query i do

Let $\nu_i = \text{Lap}\left(\frac{4}{\epsilon}\right)$

if $f_i(D) + \nu_i \geq \hat{T}$ then

Output $a_i = \top$.

Halt.

else

Output $a_i = \perp$.

end if

end for

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

Sparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\sigma = \frac{2c}{\epsilon}$. **Else Let** $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Output $a_i = \top$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

Sparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ Let $\sigma = \frac{2c}{\epsilon}$. Else Let $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i do

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ then

Output $a_i = \top$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ then

Halt.

end if

end for

Different setup for
 $\delta=0$ or $\delta \neq 0$

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

Sparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\sigma = \frac{2c}{\epsilon}$. **Else Let** $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Output $a_i = \top$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

$\text{Sparse}(D, \{f_i\}, T, c, \epsilon, \delta)$

If $\delta = 0$ Let $\sigma = \frac{2c}{\epsilon}$. Else Let $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let **count = 0**

for Each query i do

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ then

Output $a_i = \top$.

Let **count = count + 1.**

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if **count $\geq c$ then**

Halt.

end if

end for

A counter for how many queries above the threshold.

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

Sparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\sigma = \frac{2c}{\epsilon}$. **Else Let** $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Output $a_i = \top$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

Sparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\sigma = \frac{2c}{\epsilon}$. **Else Let** $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Output $a_i = \top$.

Let count = count + 1

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

We reset the threshold

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

Sparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\sigma = \frac{2c}{\epsilon}$. **Else Let** $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Output $a_i = \top$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

Multiple Above Thresholds

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

Sparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\sigma = \frac{2c}{\epsilon}$. **Else Let** $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Output $a_i = \top$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

The privacy analysis just uses composition.

Numeric Multiple Above Thresholds

Algorithm 3 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

NumericSparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\epsilon_1 \leftarrow \frac{8}{9}\epsilon, \epsilon_2 \leftarrow \frac{2}{9}\epsilon$. **Else Let** $\epsilon_1 = \frac{\sqrt{512}}{\sqrt{512+1}}\epsilon, \epsilon_2 = \frac{2}{\sqrt{512+1}}$

If $\delta = 0$ **Let** $\sigma(\epsilon) = \frac{2c}{\epsilon}$. **Else Let** $\sigma(\epsilon) = \frac{\sqrt{32c \ln \frac{2}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma(\epsilon_1))$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma(\epsilon_1))$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Let $v_i \leftarrow \text{Lap}(\sigma(\epsilon_2))$

Output $a_i = f_i(D) + v_i$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma(\epsilon_1))$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

Numeric Multiple Above Thresholds

Algorithm 3 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots

NumericSparse($D, \{f_i\}, T, c, \epsilon, \delta$)

If $\delta = 0$ **Let** $\epsilon_1 \leftarrow \frac{8}{9}\epsilon, \epsilon_2 \leftarrow \frac{2}{9}\epsilon$. **Else Let** $\epsilon_1 = \frac{\sqrt{512}}{\sqrt{512+1}}\epsilon, \epsilon_2 = \frac{2}{\sqrt{512+1}}$

If $\delta = 0$ **Let** $\sigma(\epsilon) = \frac{2c}{\epsilon}$. **Else Let** $\sigma(\epsilon) = \frac{\sqrt{32c \ln \frac{2}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma(\epsilon_1))$

Let count = 0

for Each query i **do**

Let $\nu_i = \text{Lap}(2\sigma(\epsilon_1))$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ **then**

Let $v_i \leftarrow \text{Lap}(\sigma(\epsilon_2))$

Output $a_i = f_i(D) + v_i$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma(\epsilon_1))$

else

Output $a_i = \perp$.

end if

if count $\geq c$ **then**

Halt.

end if

end for

We add fresh noise
before returning
the result of the query

Accuracy for AboveThreshold

Definition 3.9 (Accuracy). We will say that an algorithm which outputs a stream of answers $a_1, \dots, \in \{\top, \perp\}^*$ in response to a stream of k queries f_1, \dots, f_k is (α, β) -accurate with respect to a threshold T if except with probability at most β , the algorithm does not halt before f_k , and for all $a_i = \top$:

$$f_i(D) \geq T - \alpha$$

and for all $a_i = \perp$:

$$f_i(D) \leq T + \alpha.$$

Accuracy for AboveThreshold

Theorem 3.24. For any sequence of k queries f_1, \dots, f_k such that $|\{i < k : f_i(D) \geq T - \alpha\}| = 0$ (i.e. the only query close to being above threshold is possibly the last one), $\text{AboveThreshold}(D, \{f_i\}, T, \epsilon)$ is (α, β) accurate for:

$$\alpha = \frac{8(\log k + \log(2/\beta))}{\epsilon}.$$

Accuracy for AboveThreshold

Theorem 3.24. For any sequence of k queries f_1, \dots, f_k such that $|\{i < k : f_i(D) \geq T - \alpha\}| = 0$ (i.e. the only query close to being above threshold is possibly the last one), $\text{AboveThreshold}(D, \{f_i\}, T, \epsilon)$ is (α, β) accurate for:

$$\alpha = \frac{8(\log k + \log(2/\beta))}{\epsilon}.$$

Proof. Observe that the theorem will be proved if we can show that except with probability at most β :

$$\max_{i \in [k]} |\nu_i| + |T - \hat{T}| \leq \alpha$$

If this is the case, then for any $a_i = \top$, we have:

$$f_i(D) + \nu_i \geq \hat{T} \geq T - |T - \hat{T}|$$

Accuracy for AboveThreshold

Theorem 3.24. For any sequence of k queries f_1, \dots, f_k such that $|\{i < k : f_i(D) \geq T - \alpha\}| = 0$ (i.e. the only query close to being above threshold is possibly the last one), $\text{AboveThreshold}(D, \{f_i\}, T, \epsilon)$ is (α, β) accurate for:

$$\alpha = \frac{8(\log k + \log(2/\beta))}{\epsilon}.$$

Proof. Observe that the theorem will be proved if we can show that except with probability at most β :

$$\max_{i \in [k]} |\nu_i| + |T - \hat{T}| \leq \alpha$$

If this is the case, then for any $a_i = \top$, we have:

$$f_i(D) + \nu_i \geq \hat{T} \geq T - |T - \hat{T}|$$

or in other words:

$$f_i(D) \geq T - |T - \hat{T}| - |\nu_i| \geq T - \alpha$$

Accuracy for AboveThreshold

Theorem 3.24. For any sequence of k queries f_1, \dots, f_k such that $|\{i < k : f_i(D) \geq T - \alpha\}| = 0$ (i.e. the only query close to being above threshold is possibly the last one), $\text{AboveThreshold}(D, \{f_i\}, T, \epsilon)$ is (α, β) accurate for:

$$\alpha = \frac{8(\log k + \log(2/\beta))}{\epsilon}.$$

Similarly, for any $a_i = \perp$ we have:

$$f_i(D) < \hat{T} \leq T + |T - \hat{T}| + |\nu_i| \leq T + \alpha$$

We will also have that for any $i < k$: $f_i(D) < T - \alpha < T - |\nu_i| - |T - \hat{T}|$, and so: $f_i(D) + \nu_i \leq \hat{T}$, meaning $a_i = \perp$. Therefore the algorithm does not halt before k queries are answered.

Accuracy for AboveThreshold

Theorem 3.24. For any sequence of k queries f_1, \dots, f_k such that $|\{i < k : f_i(D) \geq T - \alpha\}| = 0$ (i.e. the only query close to being above threshold is possibly the last one), $\text{AboveThreshold}(D, \{f_i\}, T, \epsilon)$ is (α, β) accurate for:

$$\alpha = \frac{8(\log k + \log(2/\beta))}{\epsilon}.$$

Recall that if $Y \sim \text{Lap}(b)$, then: $\Pr[|Y| \geq t \cdot b] = \exp(-t)$. Therefore we have:

$$\Pr[|T - \hat{T}| \geq \frac{\alpha}{2}] = \exp\left(-\frac{\epsilon\alpha}{4}\right)$$

Setting this quantity to be at most $\beta/2$, we find that we require $\alpha \geq \frac{4 \log(2/\beta)}{\epsilon}$

Similarly, by a union bound, we have:

$$\Pr[\max_{i \in [k]} |\nu_i| \geq \alpha/2] \leq k \cdot \exp\left(-\frac{\epsilon\alpha}{8}\right)$$

Setting this quantity to be at most $\beta/2$, we find that we require $\alpha \geq \frac{8(\log(2/\beta) + \log k)}{\epsilon}$. These two claims combine to prove the theorem. \square