

# Group and attack: Auditing differential privacy

**Johan Lokna**

**Anouk Paradis**

Dimitar I. Dimitrov

Martin Vechev

ETH Zurich, Switzerland

lokjoh@gmail.com, {anouk.paradis, dimitar.iliev.dimitrov, martin.vechev}@inf.ethz.ch

# Why Audit Differential Privacy?

## Definition

$M$  is  $(\epsilon, \delta)$ -DP iff **for all**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] \leq \exp(\epsilon) \Pr[M(a') \in S] + \delta$$

# Why Audit Differential Privacy?

## Definition

$M$  is  $(\epsilon, \delta)$ -DP iff **for all**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] \leq \exp(\epsilon) \Pr[M(a') \in S] + \delta$$

Privacy decreases with increasing  $(\epsilon, \delta)$

# Why Audit Differential Privacy?

## Definition

M is  $(\epsilon, \delta)$ -DP iff **for all**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] \leq \exp(\epsilon) \Pr[M(a') \in S] + \delta$$


## Theory - Gauss Mechanism

$M_{\epsilon, \delta}(a) :$

- $\sigma = f(\epsilon, \delta)$
- $\eta \sim \mathcal{N}(0, \sigma)$
- return  $\text{count}(a) + \eta$

# Why Audit Differential Privacy?

## Definition

M is  $(\epsilon, \delta)$ -DP iff **for all**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] \leq \exp(\epsilon) \Pr[M(a') \in S] + \delta$$


## Theory - Gauss Mechanism

$M_{\epsilon, \delta}(a) :$

$$\left| \begin{array}{l} \sigma = f(\epsilon, \delta) \\ \eta \sim \mathcal{N}(0, \sigma) \\ \text{return count}(a) + \eta \end{array} \right.$$


## Implementation

```
1 def dp_count(count):  
2     std = compute_std(epsilon, delta)  
3     noise = np.random.normal(scale=std)  
4     return count + noise  
5
```

# Why Audit Differential Privacy?

## Definition

$M$  is  $(\epsilon, \delta)$ -DP iff **for all**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] \leq \exp(\epsilon) \Pr[M(a') \in S] + \delta$$


## Theory - Gauss Mechanism

$M_{\epsilon, \delta}(a) :$

$$\begin{cases} \sigma = f(\epsilon, \delta) \\ \eta \sim \mathcal{N}(0, \sigma) \\ \text{return count}(a) + \eta \end{cases}$$


## Implementation



```
1 def dp_count(count):  
2     std = compute_std(epsilon, delta)  
3     noise = np.random.normal(scale=std)  
4     return count + noise  
5
```

# Why Audit Differential Privacy?

## Definition

$M$  is  $(\epsilon, \delta)$ -DP iff **for all**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] \leq \exp(\epsilon) \Pr[M(a') \in S] + \delta$$


## Theory - Gauss Mechanism

$M_{\epsilon, \delta}(a) :$

$$\left| \begin{array}{l} \sigma = f(\epsilon, \delta) \\ \eta \sim \mathcal{N}(0, \sigma) \\ \text{return count}(a) + \eta \end{array} \right.$$


## Auditing

**There exists**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] > \exp(\epsilon) \Pr[M(a') \in S] + \delta$$


## Implementation



```
1 def dp_count(count):  
2     std = compute_std(epsilon, delta)  
3     noise = np.random.normal(scale=std)  
4     return count + noise  
5
```

# Why Audit Differential Privacy?

## Definition

$M$  is  $(\epsilon, \delta)$ -DP iff **for all**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] \leq \exp(\epsilon) \Pr[M(a') \in S] + \delta$$


## Theory - Gauss Mechanism

$M_{\epsilon, \delta}(a)$  :

- $\sigma = f(\epsilon, \delta)$
- $\eta \sim \mathcal{N}(0, \sigma)$
- return  $\text{count}(a) + \eta$



## Auditing with Delta-Siege

**There exists**  $(a, a') \in \mathcal{N}$  and  $S$ :

$$\Pr[M(a) \in S] > \exp(\epsilon) \Pr[M(a') \in S] + \delta$$


## Implementation



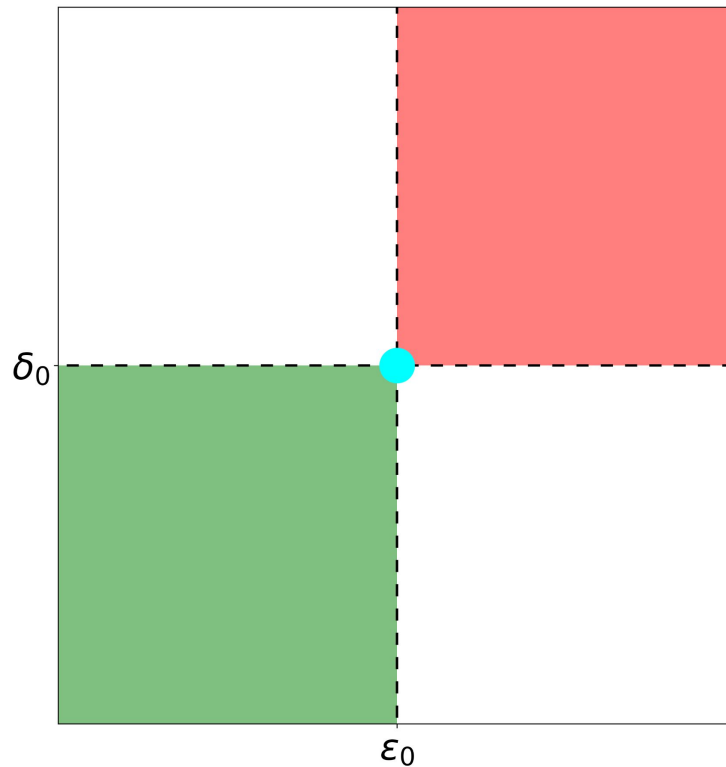
```
1 def dp_count(count):
2     std = compute_std(epsilon, delta)
3     noise = np.random.normal(scale=std)
4     return count + noise
5
```



# Characterizing a Violation

## Mechanism (●)

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```



$(\epsilon_0, \delta_0)$

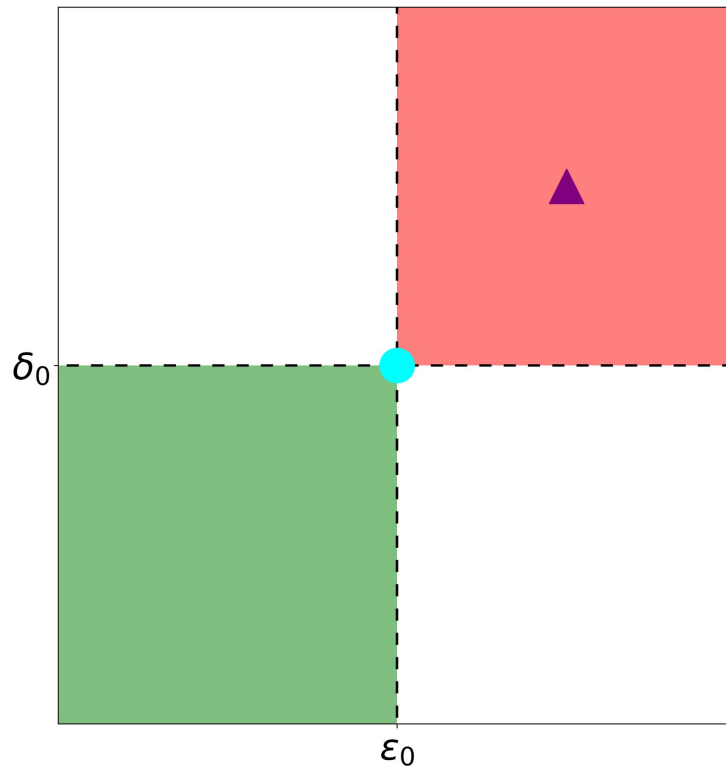
# Characterizing a Violation

## Mechanism (●)

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

## Empirical Estimate

▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$



●  $(\epsilon_0, \delta_0)$     ▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$

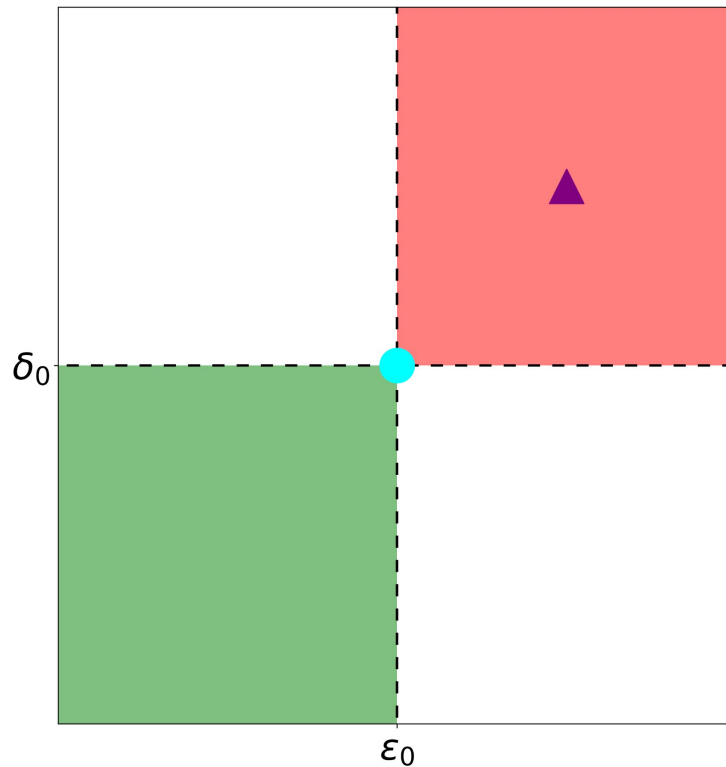
# Characterizing a Violation

## Mechanism ( ● )

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

## Empirical Estimate

▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$



●  $(\epsilon_0, \delta_0)$  ▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$

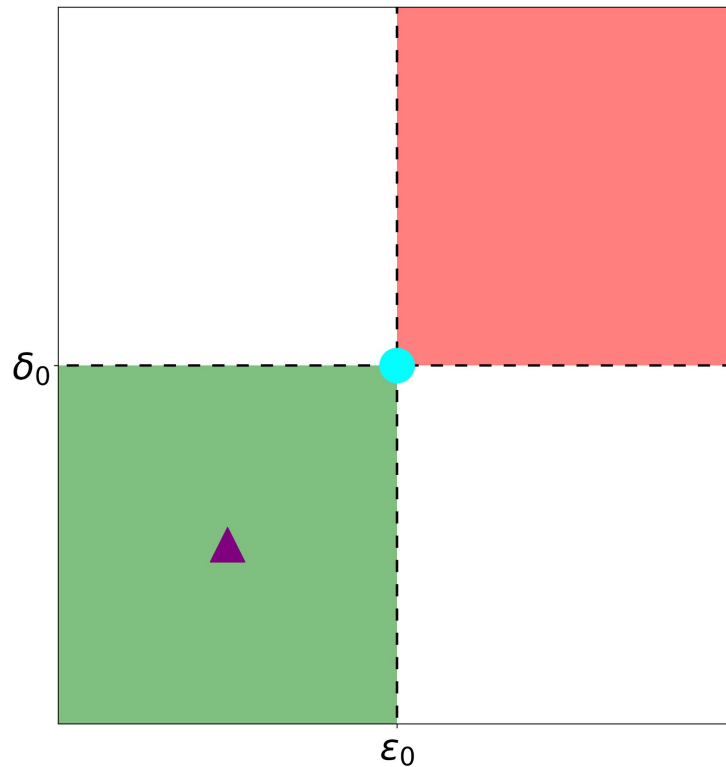
# Characterizing a Violation

## Mechanism ( ● )

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

## Empirical Estimate

▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$



●  $(\epsilon_0, \delta_0)$  ▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$

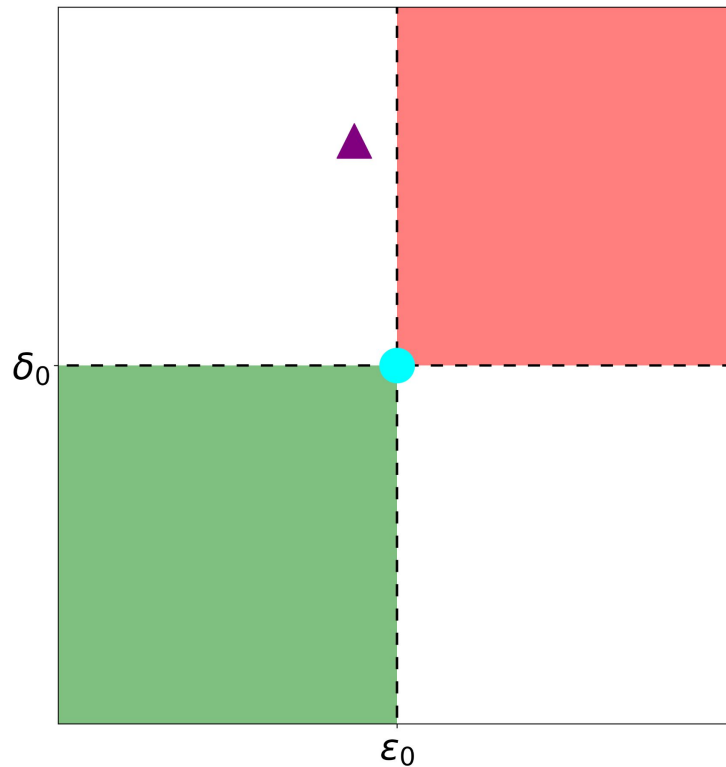
# Characterizing a Violation

## Mechanism ( ● )

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

## Empirical Estimate

▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$



●  $(\epsilon_0, \delta_0)$     ▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$

# Mechanism

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

# Mechanism

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

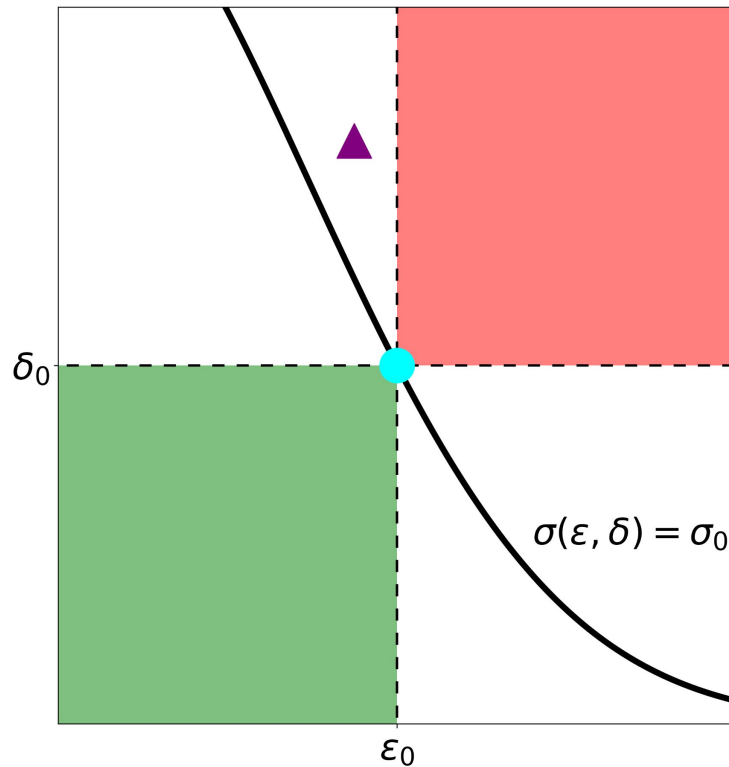
# Characterizing a Violation

## Mechanism (●)

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

## Empirical Estimate

▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$



●  $(\epsilon_0, \delta_0)$  ▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$



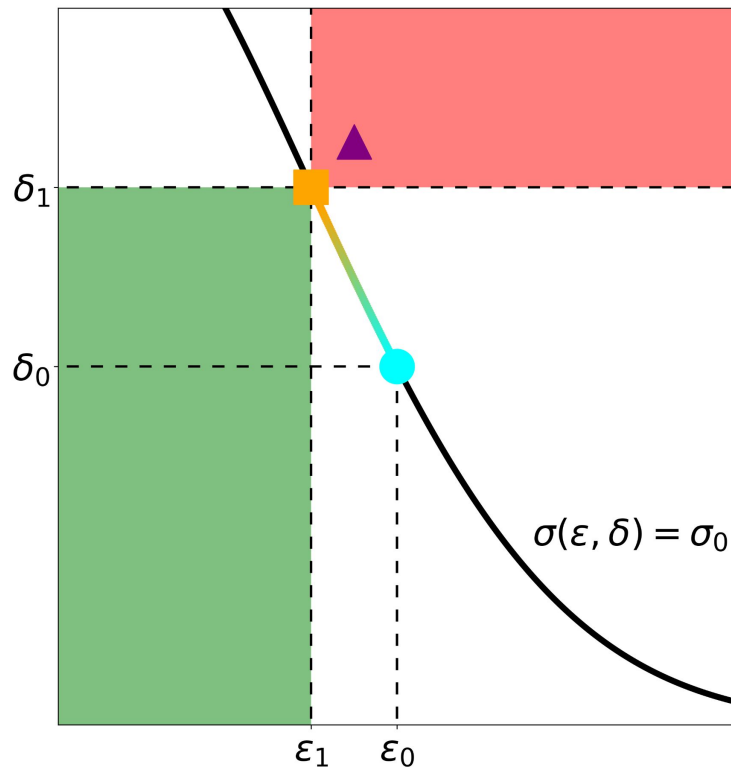
# Characterizing a Violation

## Mechanism (■)

```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

## Empirical Estimate

▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$



●  $(\epsilon_0, \delta_0)$     ▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$     ■  $(\epsilon_1, \delta_1)$

# Characterizing a Violation

## Mechanism (●)

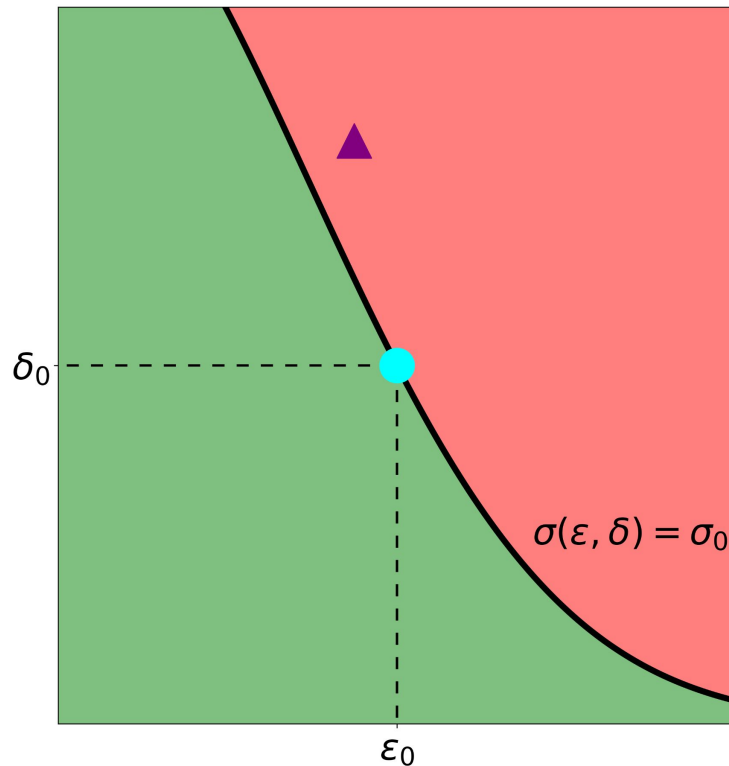
```
1 import numpy as np
2
3 epsilon = 1e-1
4 delta = 1e-6
5
6 def dp_count(count):
7     std = compute_std(epsilon, delta)
8     noise = np.random.normal(scale=std)
9     return count + noise
10
```

## Empirical Estimate

▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$

## Violation Found If

$\sigma(\text{▲}) < \sigma(\text{●})$



●  $(\epsilon_0, \delta_0)$     ▲  $(\hat{\epsilon}_0, \hat{\delta}_0)$

## Finding ▲

**Goal:** ▲ such that  $\Pr[M(a) \in S] \gg \Pr[M(a') \in S]$

## Finding ▲

**Goal:** ▲ such that  $\Pr[M(a) \in S] \gg \Pr[M(a') \in S]$

Bichsel, Benjamin, et al. "Dp-sniper: Black-box discovery of differential privacy violations using classifiers."(2021)

## Finding ▲

**Goal:** ▲ such that  $\Pr[M(a) \in S] \gg \Pr[M(a') \in S]$

Bichsel, Benjamin, et al. "Dp-sniper: Black-box discovery of differential privacy violations using classifiers."(2021)

0. Pick  $a, a'$

## Finding ▲

**Goal:** ▲ such that  $\Pr[M(a) \in S] \gg \Pr[M(a') \in S]$

Bichsel, Benjamin, et al. "Dp-sniper: Black-box discovery of differential privacy violations using classifiers."(2021)

0. Pick  $a, a'$

1. Train a classifier  $p(b) \approx \Pr[A=a \mid M(A)=b]$

## Finding ▲

**Goal:** ▲ such that  $\Pr[M(a) \in S] \gg \Pr[M(a') \in S]$

Bichsel, Benjamin, et al. "Dp-sniper: Black-box discovery of differential privacy violations using classifiers."(2021)

0. Pick  $a, a'$
1. Train a classifier  $p(b) \approx \Pr[A=a \mid M(A)=b]$
2. Define  $S = \{b \mid p(b) > 0.99\}$

## Finding the Root Cause

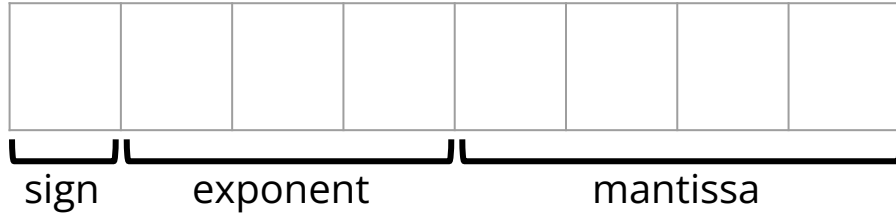
```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$



## Finding the Root Cause

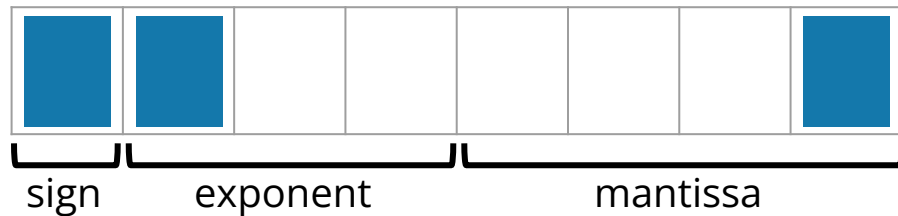


```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

## Finding the Root Cause

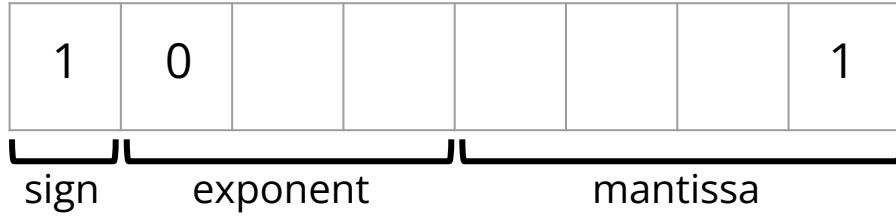


```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

## Finding the Root Cause



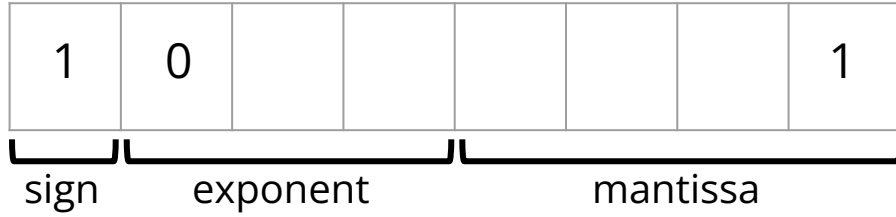
$$b = -1.xx1e^{\leq 0}$$

```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

# Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

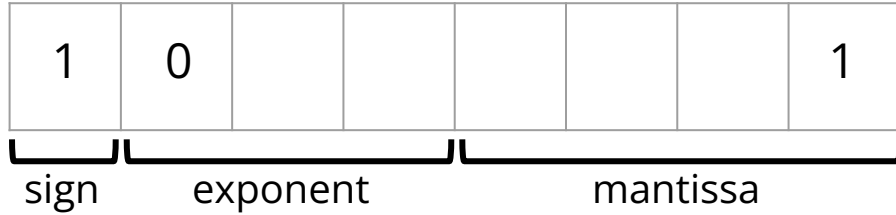
**Case a':**

```
noise = np.random.normal(scale=std)
return count + noise
```

count(a) = 0, count(a') = 1

p ← linear regression

## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

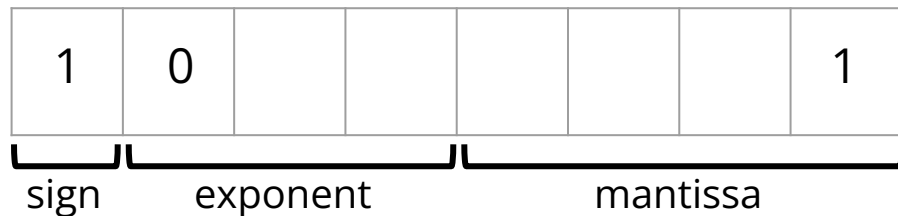
Example noise =  $-1.1111e^0$

```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$

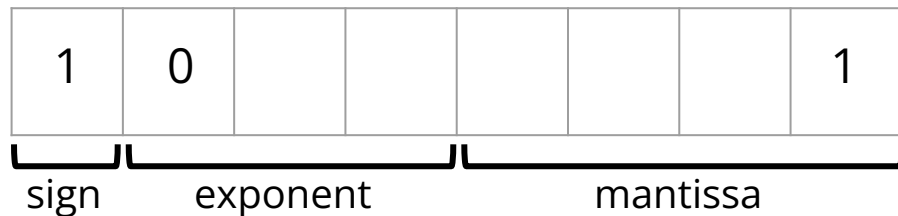
$$\begin{array}{rcl} \text{count}(a') & 1.0000e^0 \\ + \text{ noise} & \underline{-1.1111e^0} \end{array}$$

```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$

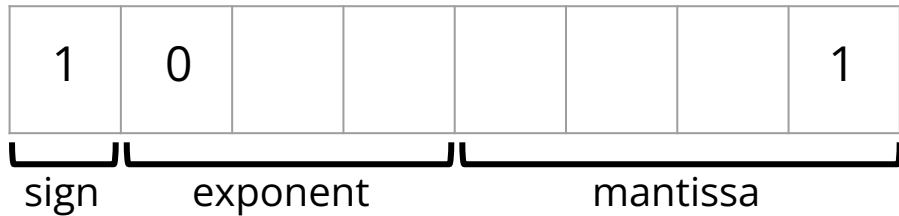
$$\begin{array}{rcl} \text{count}(a') & 1.0000e^0 & \\ + \text{ noise} & \underline{-1.1111e^0} & \\ & -0.1111e^0 & \end{array}$$

```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$

$$\begin{array}{r} \text{count}(a') \quad 1.0000e^0 \\ + \text{noise} \quad \underline{-1.1111e^0} \\ \hline \quad \quad -0.1111e^0 \\ \quad \quad -1.1110e^{-1} \end{array}$$

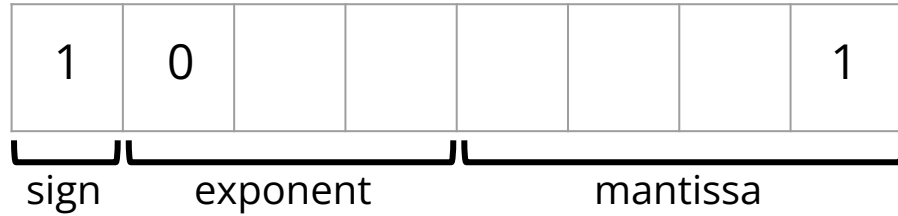
```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$




## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$

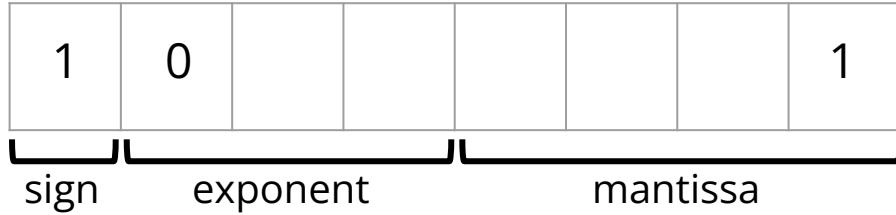
$$\begin{array}{r} \text{count}(a') \quad 1.0000e^0 \\ + \text{noise} \quad \underline{-1.1111e^0} \\ \hline \quad \quad -0.1111e^0 \\ \quad \quad -1.1110e^1 \end{array}$$


```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$


## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$

$$\begin{array}{r} \text{count}(a') \quad 1.0000e^0 \\ + \text{noise} \quad \underline{-1.1111e^0} \\ \hline \quad \quad -0.1111e^0 \\ \quad \quad -1.1110e^1 \end{array}$$


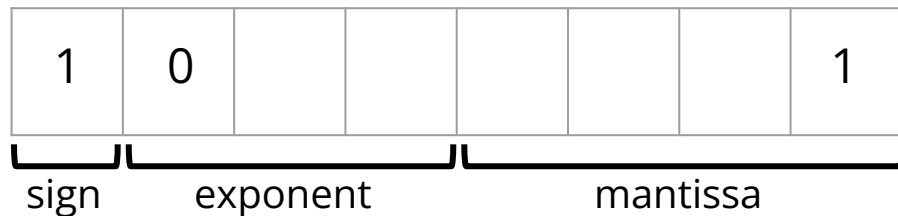
```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

### Case a:


## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$

$$\begin{array}{rcl} \text{count}(a') & 1.0000e^0 & \\ + \text{ noise} & \underline{-1.1111e^0} & \\ & -0.1111e^0 & \\ & -1.1110e^1 & \end{array}$$


```
noise = np.random.normal(scale=std)
return count + noise
```

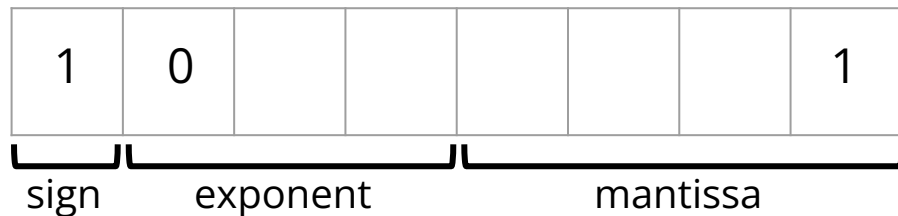
$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

### Case a:

$$\begin{array}{rcl} \text{count}(a') & 0.0000e^0 & \\ + \text{ noise} & \underline{-1.1111e^0} & \\ & -1.1111e^0 & \end{array}$$


## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$


$$\begin{array}{rcl} \text{count}(a') & 1.0000e^0 & \\ + \text{ noise} & \underline{-1.1111e^0} & \\ & -0.1111e^0 & \\ & -1.1110e^1 & \end{array}$$


```
noise = np.random.normal(scale=std)
return count + noise
```

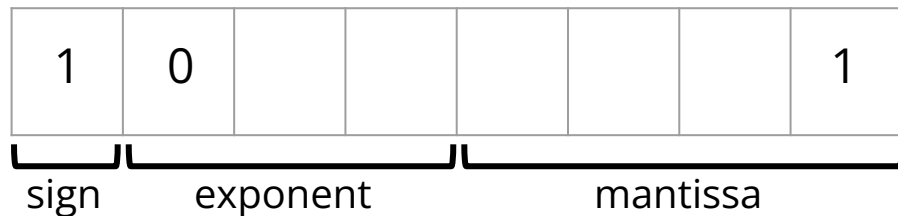
$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

### Case a:

$$\begin{array}{rcl} \text{count}(a') & 0.0000e^0 & \\ + \text{ noise} & \underline{-1.1111e^0} & \\ & -1.1111e^0 & \end{array}$$



## Finding the Root Cause



$$b = -1.xx1e^{\leq 0}$$

### Case a':

Example noise =  $-1.1111e^0$


$$\begin{array}{r} \text{count}(a') \quad 1.0000e^0 \\ + \text{noise} \quad \underline{-1.1111e^0} \\ \hline -0.1111e^0 \\ -1.1110e^1 \end{array}$$


```
noise = np.random.normal(scale=std)
return count + noise
```

$\text{count}(a) = 0, \text{count}(a') = 1$

$p \leftarrow \text{linear regression}$

### Case a:

$$\begin{array}{r} \text{count}(a') \quad 0.0000e^0 \\ + \text{noise} \quad \underline{-1.1111e^0} \\ \hline -1.1111e^0 \end{array}$$


There are fixes!

Desfontaines, Damien, and Samuel Haney. "How to Break, Then Fix, Differential Privacy on Finite Computers." (2023).

# Auditing Results

Mechanism	
<b>OpenDP</b>	Laplace Gauss
<b>DiffPrivLib (IBM)</b>	Laplace Float Gauss Analytic Float Gauss Discrete Gauss
<b>PyDP</b>	Laplace Gauss
<b>Opacus (Fb)</b>	Gauss
<b>MST</b>	
<b>AIM</b>	

# Auditing Results

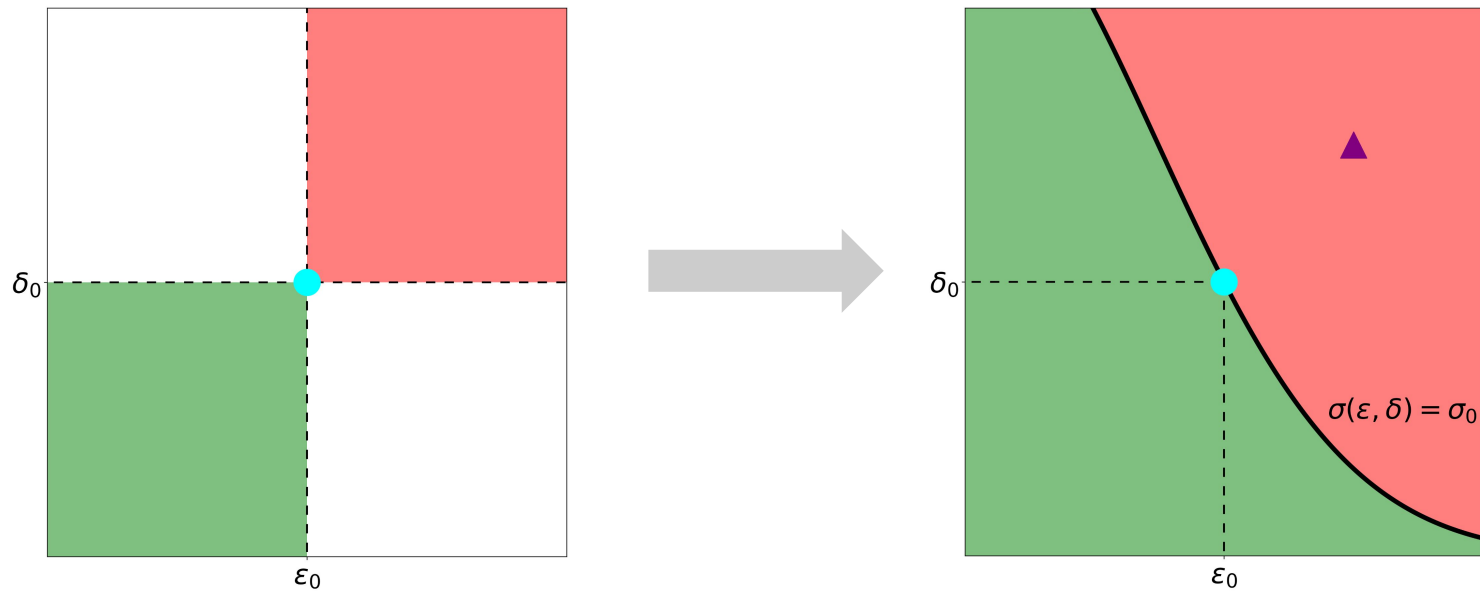
Mechanism		Violation?
<b>OpenDP</b>	Laplace Gauss	
<b>DiffPrivLib (IBM)</b>	Laplace Float Gauss Analytic Float Gauss Discrete Gauss	✓ ✓ ✓
<b>PyDP</b>	Laplace Gauss	✓
<b>Opacus (Fb)</b>	Gauss	✓
<b>MST</b>		✓
<b>AIM</b>		✓

# Auditing Results

Mechanism		Violation?
<b>OpenDP</b>	Laplace Gauss	
<b>DiffPrivLib (IBM)</b>	Laplace Float Gauss Analytic Float Gauss Discrete Gauss	✓ ✓ ✓
<b>PyDP</b>	Laplace Gauss	✓
<b>Opacus (Fb)</b>	Gauss	✓
<b>MST</b>		✓
<b>AIM</b>		✓



# Summary



New vulnerabilities found

Root cause analysis



eth-sri/Delta-Siege