

BPD Implementations

Swapnil
Mona

Github Link

Github link - <https://github.com/MonaMaNotAvailable/BostonPoliceDataAnalysis-CS151DataPrivacySecurity>

😊 Exact count = 4139 😈 Noisy Count = 4133 with epsilon 4.0

| HighestNeighAllegations | |
|-------------------------|-----------|
| neighborhood | frequency |
| Dorchester | 1051 |
| Hyde Park | 848 |
| West Roxbury | 703 |
| Roxbury | 306 |
| Roslindale | 263 |
| Mattapan | 225 |
| South Boston | 212 |
| East Boston | 133 |
| Jamaica Plain | 109 |
| Allston/Brighton | 105 |
| Charlestown | 63 |
| South End | 41 |
| Central Boston | 32 |
| Fenway/Kenmore | 32 |
| Back Bay/Beacon Hill | 16 |

| neighborhood | frequency |
|----------------------|-----------|
| Dorchester | 1051 |
| Hyde Park | 847 |
| West Roxbury | 703 |
| Roxbury | 306 |
| Roslindale | 262 |
| Mattapan | 225 |
| South Boston | 211 |
| East Boston | 132 |
| Jamaica Plain | 109 |
| Allston/Brighton | 105 |
| Charlestown | 63 |
| South End | 41 |
| Central Boston | 31 |
| Fenway/Kenmore | 31 |
| Back Bay/Beacon Hill | 16 |

BostonPoliceDataAnalysis-CS151DataPrivacySecurity

This Repo comprises three different datasets:

- [Boston Police Department Crime Hub Data](#)
- [The Woke Windows Project](#)
- [Boston Cop Track](#)

Tools:

- [SmartNoise SDK: Tools for Differential Privacy on Tabular Data](#)
- [Matplotlib](#)

Design Docs:

1. Checkpoint 3 - Analysis done with different queries:
https://docs.google.com/document/d/1-G3i7vvoT-xE0hwLoJcmTejouE_otZFtw2ELXEazik/edit?usp=sharing
2. Checkpoint 4 - Implementation with OpenDp SmartNoise:
https://docs.google.com/document/d/1fQjWv4tuH1c_cX26mAo6Qt5zfODX_8M0SJMOovoJ_YQ/edit?usp=sharing

5 Different SQL Queries

| Query | Query History |
|-------|---|
| 1 | <code>SELECT allegation, COUNT(*) AS frequency</code> |
| 2 | <code>FROM bpd_allegations</code> |
| 3 | <code>GROUP BY allegation</code> |
| 4 | <code>ORDER BY frequency DESC;</code> |

Figure 1. Count of Allegations by Type

| Query | Query History |
|-------|---|
| 1 | <code>SELECT neighborhood, COUNT(*) AS frequency</code> |
| 2 | <code>FROM bpd_allegations</code> |
| 3 | <code>WHERE neighborhood IS NOT NULL</code> |
| 4 | <code>GROUP BY neighborhood</code> |
| 5 | <code>ORDER BY frequency DESC;</code> |

Figure 2. Allegation Counts by Neighborhood

| Query | Query History |
|-------|--|
| 1 | <code>SELECT</code> |
| 2 | <code>first_name,</code> |
| 3 | <code>last_name,</code> |
| 4 | <code>COUNT(allegation)/2 AS allegation_count</code> |
| 5 | <code>FROM</code> |
| 6 | <code>AllegationCountOnOfficers.AllegationCountOnOfficers</code> |
| 7 | <code>WHERE</code> |
| 8 | <code>active = TRUE</code> |
| 9 | <code>GROUP BY</code> |
| 10 | <code>first_name,</code> |
| 11 | <code>last_name</code> |
| 12 | <code>ORDER BY allegation_count DESC</code> |

Figure 3. Adjusted Allegation Counts per Officer

| Query | Query History |
|-------|--|
| 1 | <code>SELECT title, AVG(total) AS avgSalaries</code> |
| 2 | <code>FROM officers</code> |
| 3 | <code>WHERE title IS NOT NULL</code> |
| 4 | <code>GROUP BY title</code> |
| 5 | <code>HAVING AVG(total) IS NOT NULL</code> |
| 6 | <code>ORDER BY avgSalaries DESC;</code> |

Figure 4. Average Salaries of Officers by Title or Rank

| Query | Query History |
|-------|--|
| 1 | <code>SELECT NEIGHBORHOOD, COUNT(*) AS totalVictims</code> |
| 2 | <code>FROM PersonShot.PersonShot</code> |
| 3 | <code>WHERE NEIGHBORHOOD IS NOT NULL</code> |
| 4 | <code>GROUP BY NEIGHBORHOOD</code> |
| 5 | <code>ORDER BY totalVictims DESC;</code> |

Figure 5. Shootings Counts by Neighborhood

Usage of Privacy Budget

- First, total epsilon is divided in a sequential manner between all these queries, i.e. Total epsilon 50 \Rightarrow 50/5 \Rightarrow 10 per query
- High epsilon is chosen to provide more accuracy.
- Low epsilon is chosen to provide more noise.

| | Q1 | Q2 | Q3 | Q4 | Q5 | Total |
|-----------------|------|------|------|-------|------|-------|
| Lower Bound | 0.01 | 0.1 | 0.05 | 9.0 | 1.0 | 10.16 |
| Upper Bound | 5.0 | 4.0 | 3.0 | 60.0 | 2.2 | 74.2 |
| Assigned Budget | 3.37 | 2.70 | 2.02 | 40.43 | 1.48 | 50 |

Table 1: Range of epsilon for different queries

Epsilon vs Error & Epsilon vs Runtime

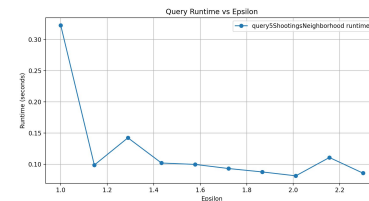
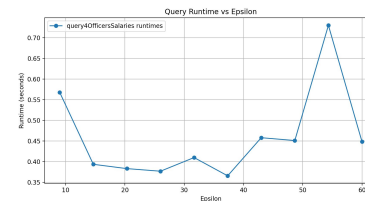
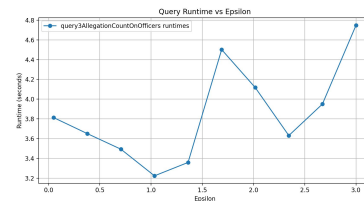
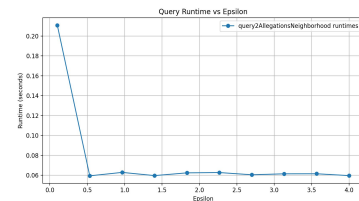
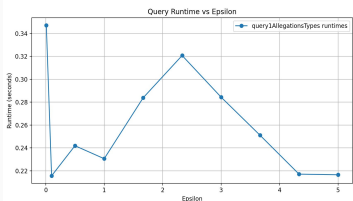
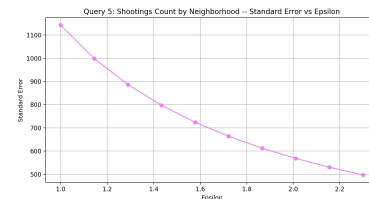
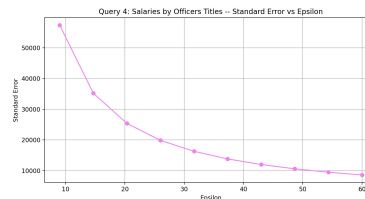
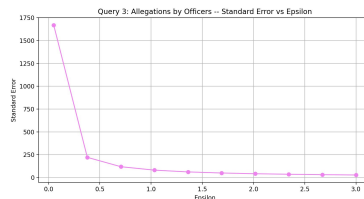
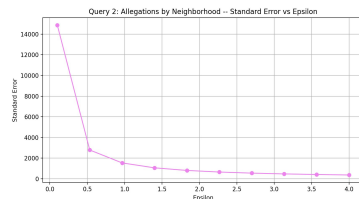
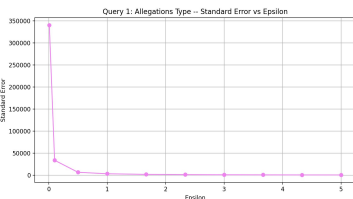
Q1

Q2

Q3

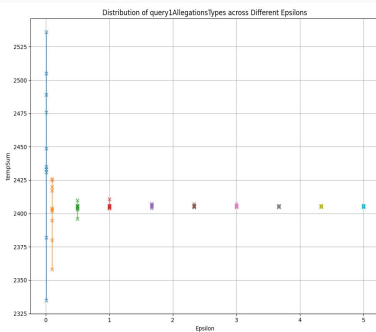
Q4

Q5

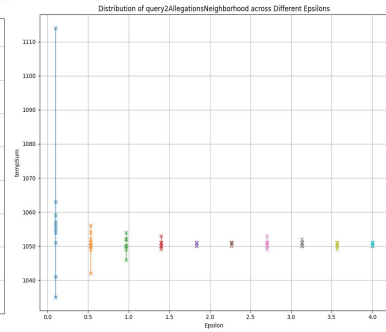


Distribution of Results after 10 Executions

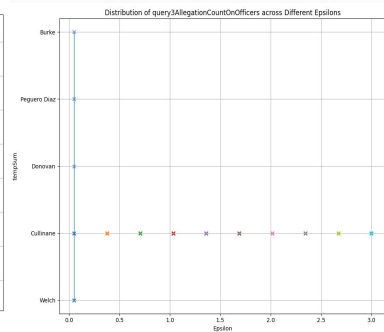
Q1



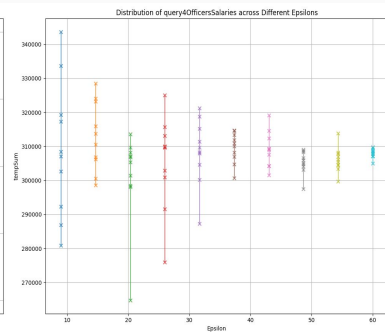
Q2



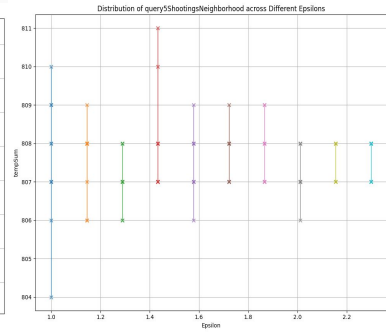
Q3



Q4

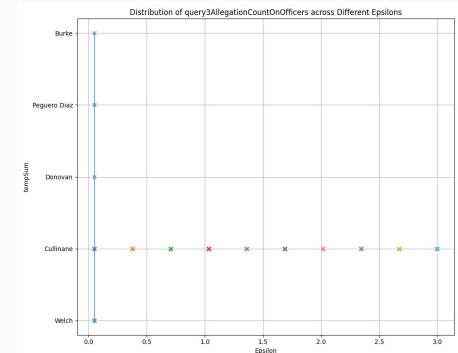
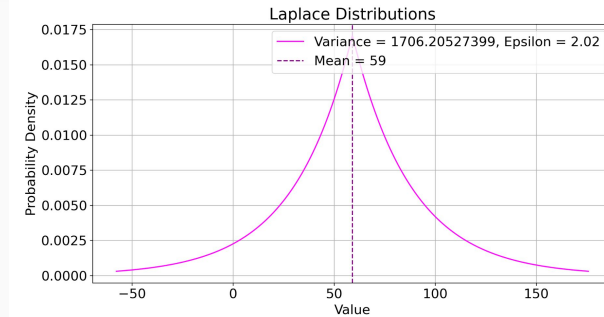


Q5



Comparing the analytical result from checkpoint 3

- Query 3: Number of allegation count per police officer
- Empirical Result:
 - Error: **~2.09 (e=40)**, **~ 8.34 (e=10)** -> final count in the range of (57,61) and (51,67)
- Actual Result:
 - Error: **27.81 (e=3)**, **1668.77 (e=0.05)** -> as precise as possible but with noise added
 - Actual privacy budget allocated \Rightarrow **Error: 41.31 (e = 2.02)**



Thank you for listening!

Any question?

