

Differentially private boxplots

Jairo DIAZ-RODRIGUEZ
York University, Canada

joint work with
Kelly RAMSAY

International Conference of Machine Learning.

June 14, 2025

Differential Privacy (DP)

\$10000



\$7500



\$7000



**Average public
salary**

Day 1: \$8666



\$8000



\$7500



\$10000

Differential Privacy (DP)

\$10000



\$7500



\$7000



**Average public
salary**

Day 1: \$8666

Day 2: \$8666



\$8000



\$7500



\$10000

Differential Privacy (DP)

\$10000



\$7500



\$7000



\$8000



\$7500



\$10000

**Average public
salary**

Day 1: \$8666

Day 2: \$8666

Day 3: \$8400

Differential Privacy (DP)

\$10000



\$7500



\$7000



Someone fired!



\$8000



\$7500



\$10000

Day 1: \$8666

Day 2: \$8666

Day 3: \$8400

Differential Privacy (DP)

\$10000



\$7500



\$7000



**DP Average public
salary**

Day 1: \$8666 + noise

Day 2: \$8666 + noise

Day 3: \$8400 + noise



\$8000



\$7500



\$10000

Differential Privacy (DP)

\$10000



\$7500



\$7000



\$8000



\$7500



\$10000

**DP Average public
salary**

Day 1: \$8510

Day 2: \$8680

Day 3: \$8450

Differential Privacy (DP)

\$10000



\$7500



\$7000



Someone fired?



\$8000



\$7500



\$10000

Day 1: \$8510

Day 2: \$8680

Day 3: \$8450

Differential Privacy (DP)

A randomized function \mathcal{A} , which operates on datasets, is ϵ -differentially private if for any two datasets D and D' that differ in a single observation (these are often called adjacent datasets), and for all sets S of possible outputs of \mathcal{A} :

$$P(\mathcal{A}(D) \in S) \leq e^\epsilon \times P(\mathcal{A}(D') \in S)$$

Differential Privacy (DP)

Low
privacy budget
(small ϵ)

Day 1: \$8590

Day 2: \$8680

Day 3: \$8550

High
privacy budget
(big ϵ)

Day 1: \$8650

Day 2: \$8675

Day 3: \$8420

Differential Privacy (DP)

Low
privacy budget
(small ϵ)

Day 1: \$8590

Day 2: \$8680

Day 3: \$8550

High
privacy budget
(big ϵ)

Day 1: \$8650

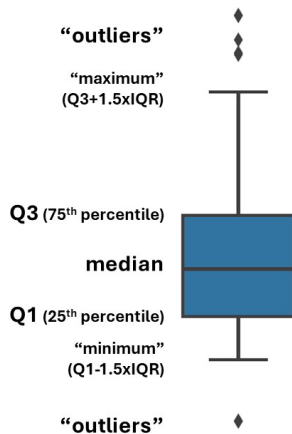
Day 2: \$8675

Day 3: \$8420

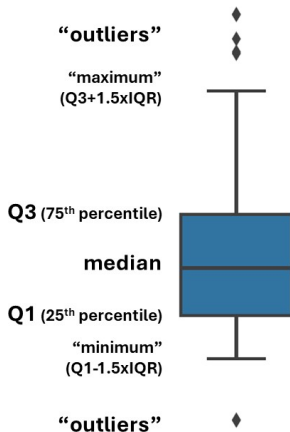
Challenge

Draw useful insights from datasets while *protecting individual privacy*

Boxplots



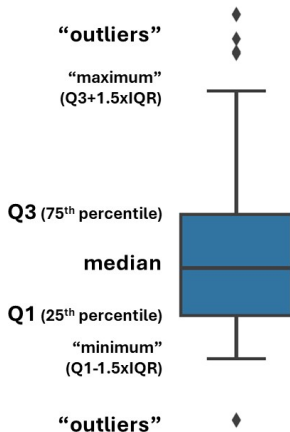
Boxplots



Data analysis with boxplots

- Location
- Scale
- Skewness
- Tails

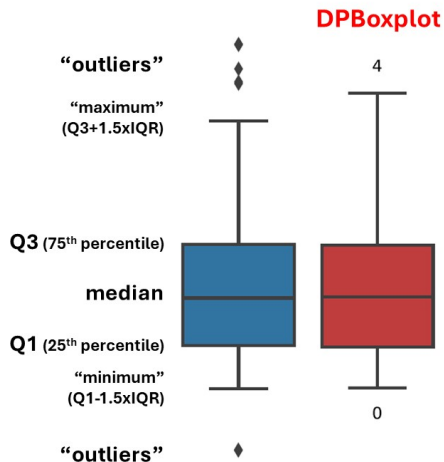
Boxplots



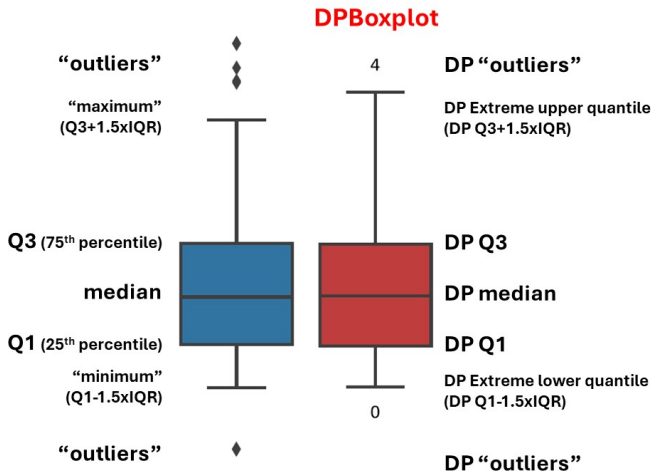
**Challenge for
differentially private
boxplots**

*Draw useful insights
from datasets while
protecting individual
privacy*

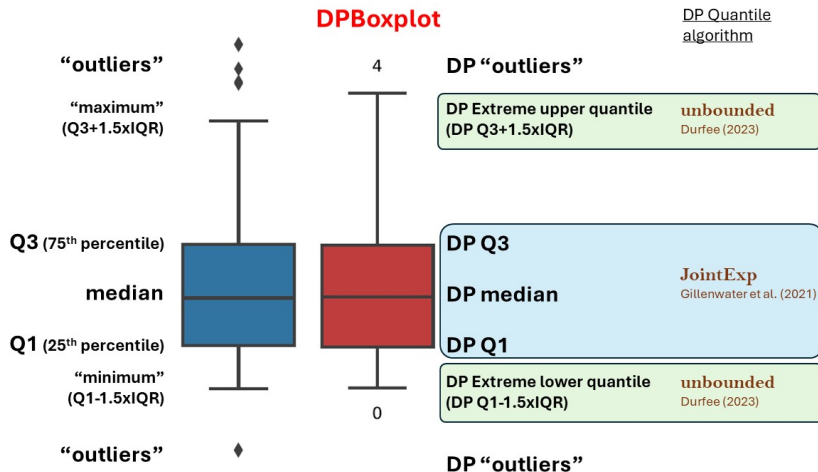
Differentially private boxplots



Differentially private boxplots



Differentially private boxplots



On DP Quantile estimation

Box

On DP Quantile estimation

Box

- We used JointExp quantile estimation (Gillenwater et al., 2021).

On DP Quantile estimation

Box

- We used JointExp quantile estimation (Gillenwater et al., 2021).
- **Main theoretical contributions:**
 - We proved a minimax lower bound for privately estimating a quantile which matches the upper bound.

On DP Quantile estimation

Box

- We used JointExp quantile estimation (Gillenwater et al., 2021).
- **Main theoretical contributions:**
 - We proved a minimax lower bound for privately estimating a quantile which matches the upper bound.
 - This implies that the **scale** and **location** of the proposed private boxplot are estimated optimally, up to logarithmic factors.

On DP Quantile estimation

Whiskers and outliers

On DP Quantile estimation

Whiskers and outliers

- We used unbounded quantile estimation (Durfee, 2023).

On DP Quantile estimation

Whiskers and outliers

- We used unbounded quantile estimation (Durfee, 2023).
- **Main theoretical contributions:**
 - JointExp is inconsistent for extreme quantiles.

On DP Quantile estimation

Whiskers and outliers

- We used unbounded quantile estimation (Durfee, 2023).
- **Main theoretical contributions:**
 - JointExp is inconsistent for extreme quantiles.
 - We proved that whiskers and number of outliers are weakly consistent for their population counterparts when using unbounded.

On DP Quantile estimation

Whiskers and outliers

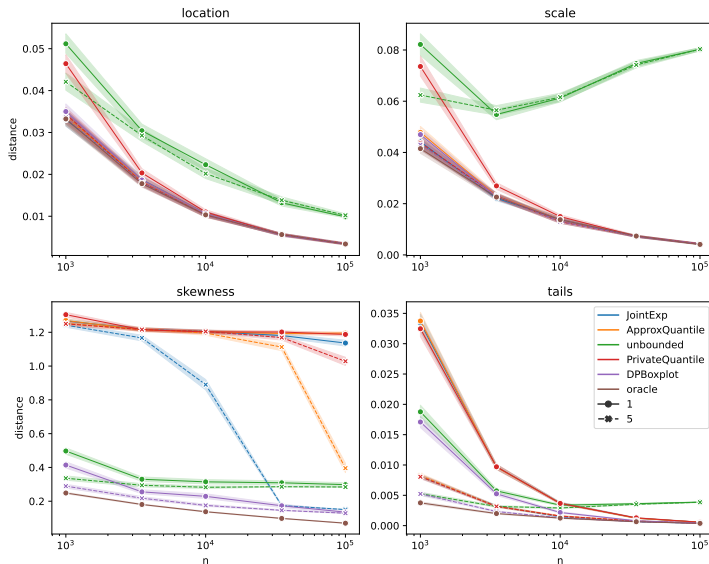
- We used unbounded quantile estimation (Durfee, 2023).
- **Main theoretical contributions:**
 - JointExp is inconsistent for extreme quantiles.
 - We proved that whiskers and number of outliers are weakly consistent for their population counterparts when using unbounded.
 - This implies that **skewness** and **tails** will be correctly portrayed by our proposed differentially private boxplot.

Simulation studies

We performed an exhaustive simulation study under multiple settings:

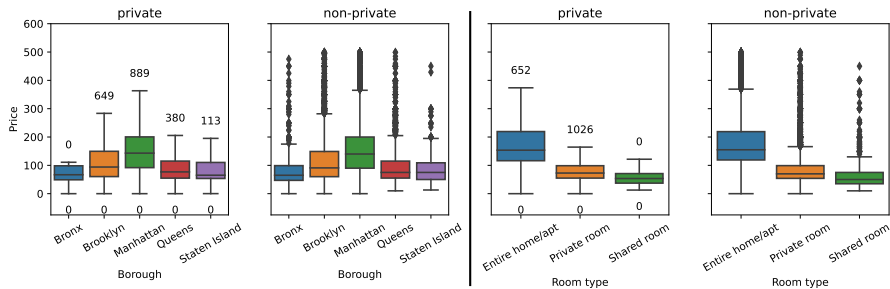
- Generating distributions: Normal, skew-normal, uniform, beta, real data.
- ϵ : 0.5, 1, 5, 10
- sample sizes: 1000, 2000, 10000, 20000, 100000
- We quantified similitude in **location**, **scale**, **skewness**, and **tails** between constructed private boxplot and the population counterpart.
- We compared with naive private boxplots constructed by using different private quantile estimation methods.

Simulation studies (summary)



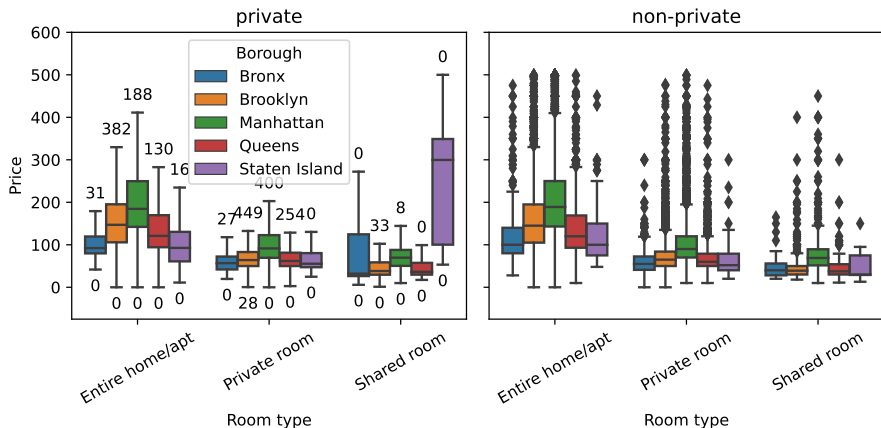
Real data

Do discernible patterns emerge in Airbnb listing prices across various boroughs in New York City and differing room types?



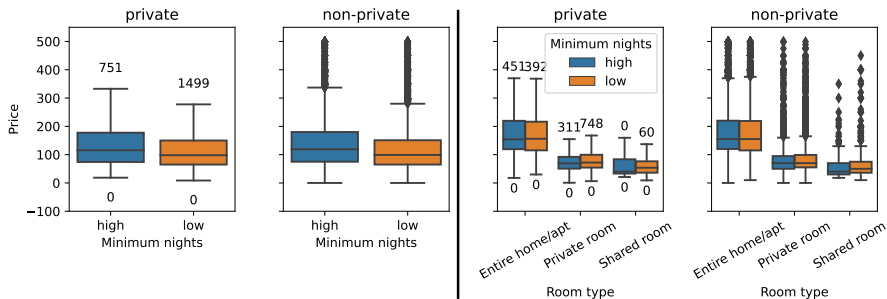
Real data

Do discernible patterns emerge in Airbnb listing prices across various boroughs in New York City and differing room types?



Real data

Are there observable trends in Airbnb listing prices concerning minimum nights required for reservation and the types of rooms offered?



Real data

Conclusions

- We observe relatively consistent patterns between the private and non-private boxplots.

Real data

Conclusions

- We observe relatively consistent patterns between the private and non-private boxplots.
- The primary visual disparities pertains to the positioning of whiskers on the boxplots and number of outliers. This underscores the recognized challenge associated with differentially private estimation of extreme quantiles and small sample sizes.

Real data

Conclusions

- We observe relatively consistent patterns between the private and non-private boxplots.
- The primary visual disparities pertains to the positioning of whiskers on the boxplots and number of outliers. This underscores the recognized challenge associated with differentially private estimation of extreme quantiles and small sample sizes.
- Discrepancies does not materially impede the analytical value of the visual findings.

Real data

Conclusions

- We observe relatively consistent patterns between the private and non-private boxplots.
- The primary visual disparities pertains to the positioning of whiskers on the boxplots and number of outliers. This underscores the recognized challenge associated with differentially private estimation of extreme quantiles and small sample sizes.
- Discrepancies does not materially impede the analytical value of the visual findings.
- This case study demonstrates the potential of differentially private, exploratory data analysis and confirms the efficacy of the differentially private boxplot, in accordance with our theoretical and simulated results.

Thank you.

Differentially private boxplots.

code: <https://github.com/jairoadiazr/DPBoxplot>