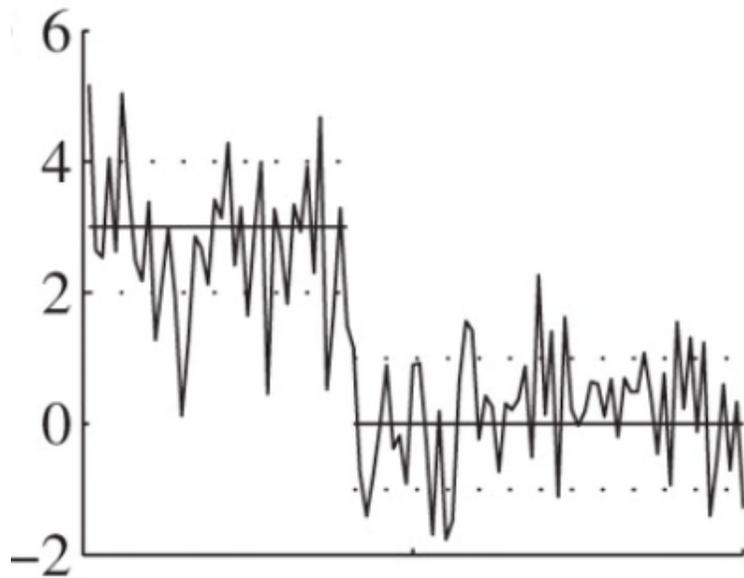


Privately Detecting Changes in Unknown Distributions

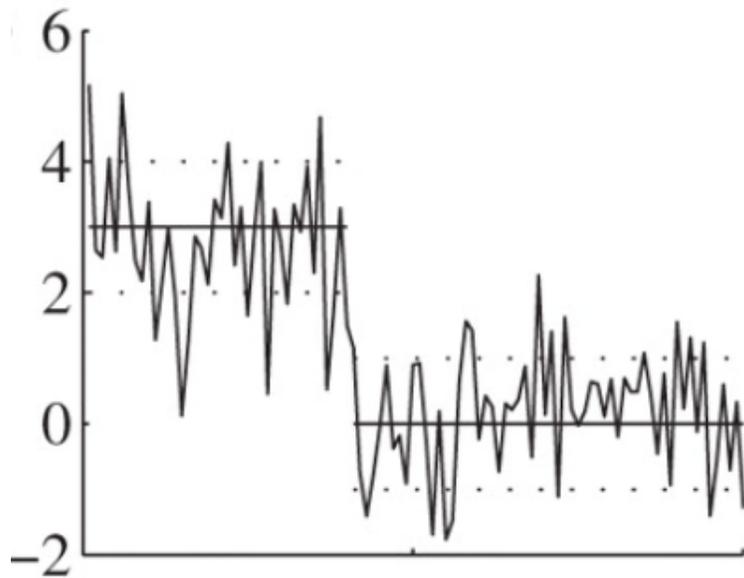
Wanrong Zhang, Georgia Tech

joint work with Rachel Cummings, Sara Krehbiel, Yuliia Lut

Motivation I: Smart-home IoT devices



Motivation II: Disease outbreaks



Change-point problem:

Identify distributional changes in stream of **highly sensitive** data

Model:

Data points $x_1, \dots, x_{k^*} \sim P_0$ (pre-change)

$x_{k^*}, \dots, x_n \sim P_1$ (post-change)

Need formal **privacy** guarantees for change-point detection algorithms

Question:

Estimate the unknown change time k^*

Previous work: parametric model [CKM+18] (P_0 and P_1 **known**)

Our work: nonparametric model (P_0 and P_1 **unknown**)

Differential privacy [DMNS '06]

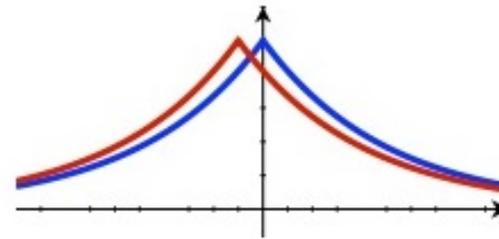
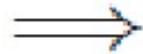
Bound the maximum amount that one person's data can change the distribution of an algorithm's output

An algorithm $M: T^n \rightarrow R$ is **ϵ -differentially private** if \forall neighboring $x, x' \in T^n$ and $\forall S \subseteq R$,

$$P[M(x) \in S] \leq e^\epsilon P[M(x') \in S]$$

$(t_1, \dots, t_i, \dots, t_n)$

$(t_1, \dots, t'_i, \dots, t_n)$



- S as set of “bad outcomes”
- Worst-case guarantee

Privately Detecting Changes in Unknown Distributions

1. Offline setting: dataset known in advance
2. Online setting: data points arrive one at a time
3. Drift change detection (in paper)
4. Empirical results (in paper)

Privately Detecting Changes in Unknown Distributions

1. ***Offline setting: dataset known in advance***
2. Online setting: data points arrive one at a time
3. Drift change detection (in paper)
4. Empirical results (in paper)

Mann-Whitney test [MW '47]

Datasets: $x_1, x_2, \dots, x_k \sim P_0$ and $x_{k+1}, x_{k+2}, \dots, x_n \sim P_1$

$$H_0: P_0 = P_1, H_1: P_0 \neq P_1$$

$$\text{Test statistic: } V(k) = \frac{1}{k(n-k)} \sum_{j=k+1}^n \sum_{i=1}^k I(x_i > x_j)$$

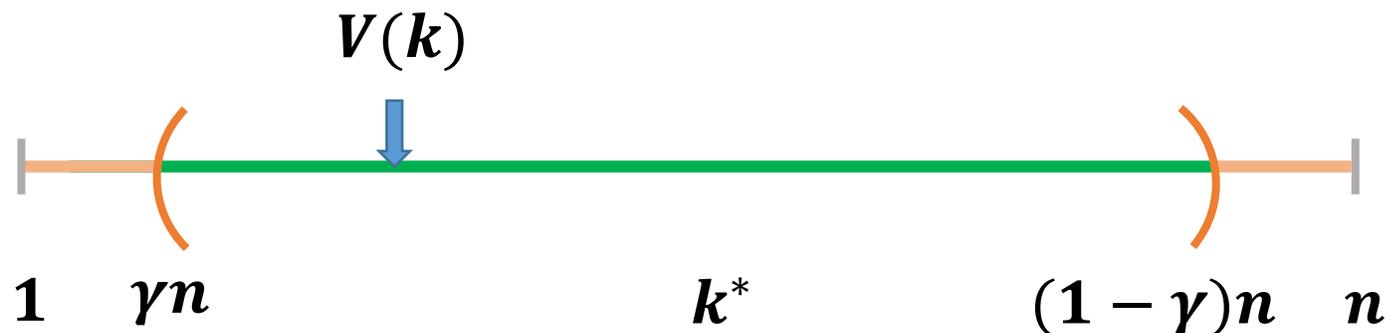
Under H_1 , require $a := Pr_{x \sim P_0, y \sim P_1} [x > y] \neq \frac{1}{2}$

Number of such pairs (x_i, x_j) such that $x_i > x_j$

Non-private nonparametric change-point detection [Darkhovsky '79]

1. For every $k \in [\gamma n], \dots, [(1 - \gamma)n]$
2. Compute $V(k)$
3. Output $\hat{k} = \operatorname{argmax}_k V(k)$

Can we compute $V(k)$ or $\operatorname{argmax}_k V(k)$ privately?



Adding differential privacy

Differentially private algorithms add noise that scale with the *sensitivity* of a query.

Query sensitivity: The sensitivity of real-valued query f is:

$$\Delta f = \max_{X, X' \text{ neighbors}} |f(X) - f(X')|.$$

Laplace Mechanism: The mechanism $M(f, X, \epsilon) = f(X) + \text{Lap}(\frac{\Delta f}{\epsilon})$ is ϵ -differentially private.

Offline PNCPD

= Mann-Whitney + ReportNoisyMax

Private Nonparametric Change-Point Detector: $PNCPD(X, \epsilon, \gamma)$

1. Input: database, privacy parameter ϵ , constraint parameter γ
2. for $k \in [\gamma n], \dots, \lfloor (1 - \gamma)n \rfloor$
3. Compute statistic $V(k)$
4. Sample $Z_k \sim \text{Lap}\left(\frac{2}{\epsilon \gamma n}\right)$
5. Output $\tilde{k} = \text{argmax}_k (V(k) + Z_k)$

Main results: OfflinePNCPD

Theorem: $\text{OfflinePNCPD}(X, \epsilon, \gamma)$ is ϵ -differentially private and with probability $1 - \beta$, it outputs private change-point estimator \tilde{k} with error at most

$$|\tilde{k} - k^*| < O\left(\frac{1}{\epsilon\gamma^4(a - 1/2)^2}\right)^{1.01} \cdot \log\frac{1}{\beta}$$

- Previous non-private analysis [Darkhovsky '76]

$$|\hat{k} - k^*| < O(n^{2/3})$$

- Our improved non-private analysis:

$$|\hat{k} - k^*| < O\left(\frac{1}{\gamma^4(a - 1/2)^2} \log\frac{1}{\beta}\right) = O(1)$$

Privately Detecting Changes in Unknown Distributions

1. Offline setting: dataset known in advance
2. ***Online setting: data points arrive one at a time***
3. Drift change detection (in paper)
4. Empirical results (in paper)

Online setting

More challenging: must detect change quickly without much post-change data

High Level Approach:

1. Privately detect online when $V(k) > T$ in the center of a sliding window of last n data points.
2. Run OfflinePNCPD on the identified window.

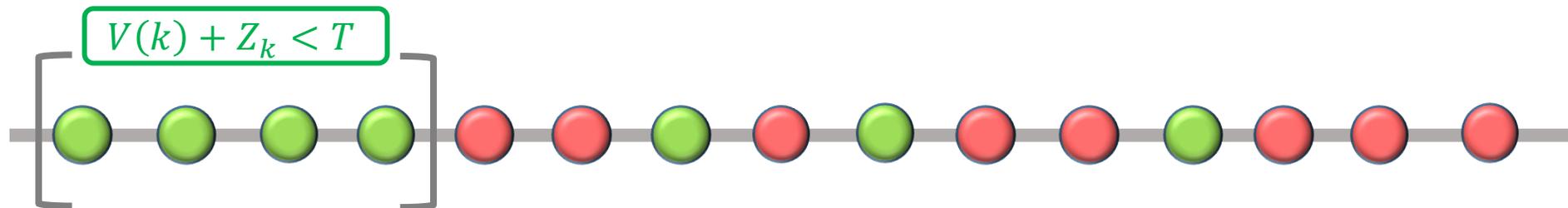
Have DP algorithm
(AboveNoisyThreshold)
for this

Online setting

More challenging: must detect change quickly without much post-change data

Our Approach:

1. Run AboveNoisyThreshold on Mann-Whitney queries in the center of a sliding window of last n data points.
2. Run OfflinePNCPD on the identified window.

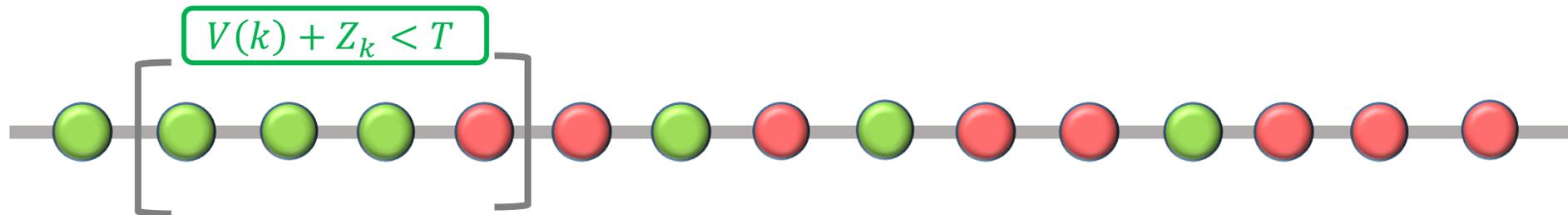


Online setting

More challenging: must detect change quickly without much post-change data

Our Approach:

1. Run AboveNoisyThreshold on Mann-Whitney queries in the center of a sliding window of last n data points.
2. Run OfflinePNCPD on the identified window.

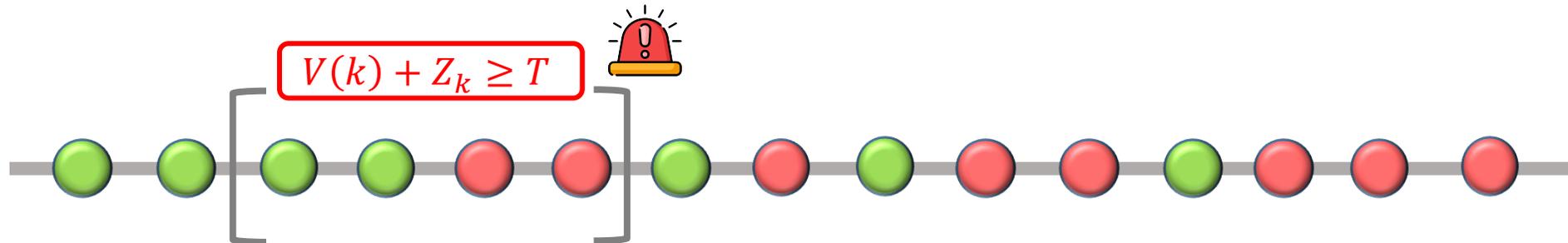


Online setting

More challenging: must detect change quickly without much post-change data

Our Approach:

1. Run AboveNoisyThreshold on Mann-Whitney queries in the center of a sliding window of last n data points.
2. Run OfflinePNCPD on the identified window.



OnlinePNCPD

1. Input: database $X = \{x_1, \dots\}$, privacy parameter ϵ , threshold T
2. Let $\hat{T} = T + \text{Lap}\left(\frac{8}{\epsilon n}\right)$
3. For each new data point x_k :
4. Compute Mann-Whitney statistic $V(k)$ in center of last n data points
5. Sample $Z_k \sim \text{Lap}\left(\frac{16}{\epsilon n}\right)$
6. If $V(k) + Z_j > \hat{T}$, then
7. Run OfflinePNCPD on last n data points with $\epsilon/2$
8. Else, output \perp

Main result: OnlinePNCPD

Theorem: $\text{OnlinePNCPD}(X, T, \epsilon, \gamma)$ is ϵ -differentially private. For appropriate threshold T , with probability $1 - \beta$, it outputs private change-point estimator \tilde{k} with error at most

$$|\tilde{k} - k^*| < O\left(\frac{1}{\epsilon} \log \frac{n}{\beta}\right)$$

where n is the window size.

Choice of T

- Can't raise alarm too early (False positive: $T > T_L$)
- Can't fail to raise alarm at true change (False negative: $T < T_H$)

Privately Detecting Changes in Unknown Distributions

1. Offline setting: dataset known in advance
2. Online setting: data points arrive one at a time
3. Drift change detection (in paper)
4. Empirical results (in paper)

References

- Cummings, R., Krehbiel, S., Mei, Y., Tuo, R., & Zhang, W. Differentially private change-point detection. In *Advances in Neural Information Processing Systems*, NeurIPS'18 pp. 10848-10857, 2018
- Dwork, C., McSherry, F., Nissim, K., & Smith, A. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pp. 265-284, 2006.
- Dwork, C., Roth, A. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4), 211-407, 2014.
- Darkhovsky, B. A nonparametric method for the a posteriori detection of the "disorder" time of a sequence of independent random variables. *Theory of Probability & Its Applications*, 21(1):178-183, 1976.
- Mann, H.B. and Whitney, D.R. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics*, pp 50-60, 1947.

Privately Detecting Changes in Unknown Distributions

Wanrong Zhang, Georgia Tech

joint work with Rachel Cummings, Sara Krehbiel, Yuliia Lut