

## Motivation

Cumulative distribution functions (CDFs) and empirical CDFs (eCDFs) are widely used to summarize distributions. The Dvoretzky-Kiefer-Wolfowitz (DKW) inequality shows the eCDF is a good approximation of the true.

If the underlying data is sensitive, we can use privacy-preserving approximations of the CDF. We want to:

- Balance privacy and utility/accuracy.
- Be compatible with federated and online updating/learning.

**Goal:** Design a differentially private CDF approximation that allows for continuous updating.



Privacy comes from a randomized map Q(y|x). We say Q guarantees  $(\epsilon, \delta)$ -differential privacy [1-3] if

 $Q(\mathcal{S}|x) \le e^{\epsilon}Q(\mathcal{S}|x') + \delta$ 

for all measurable subsets  $S \subseteq Y$  and all  $x, x' \in X$  with  $x \sim x'$  differ in a single  $x_i$ .



### **Prior work on private CDF estimation**

Given scalar data  $\mathcal{D} = \{x_i \in [A, B]: i \in A\}$ [*n*]}, estimate the eCDF

$$\widehat{F}_X(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(x_i \le t).$$

Two standard algorithms:

- Histogram Queries (HQ) bins the data, computes a differentially private histogram, and interpolates a CDF.
- **Adaptive Quantiles (AQ)** queries splits of data to sequentially estimate values of CDF (cf. binary search).

HQ and AQ methods have limitations in various scenarios.

# **Differentially Private Distribution Estimation Using Functional Approximation**

## ECE Department, Rutgers University, Piscataway, New Jersey 08854

## **CDF** estimation as signal approximation

### **Approach: function approximation to** estimate the CDF

The eCDF is a very structured signal: piecewise constant, bounded, monotonic....

Signal processing has lots of techniques for structured signal approximation (Fourier Series, wavelets, polynomials).

The first three partial sums of the Fourier series for a square wave. Image sourced from Wikipedia.

### **Polynomial Projection (PP) with Postprocessing**

**Step 1**: Estimate the empirical CDF using a series of polynomial functions

$$\widehat{F}_X(t) \approx \sum c_k P_k(t).$$

- Choose the polynomial space  $\mathcal{P}$  spanned by the first K + 1 Legendre polynomials  $\{P_0, P_1, \dots, P_K\}$ , where each  $P_k$  is a polynomial of degree of k.
- Obtain the optimal CDF estimate in the polynomial space by projecting the empirical CDF onto this space.

**<u>Step 2</u>**: Add noise to the coefficients for privacy protection.

• These coefficients are associated with the moments of the samples  $\frac{1}{n}\sum x_i^m$ , providing a straightforward approach for sensitivity computation.

**Step 3**: Isotonic regression is used as a post-processing method to ensure the estimated CDF is non-decreasing.

## **Prior work on private CDF estimation**

Theorem (Upper Bound for  $|| F^* - \widetilde{F} ||_2$ ).

Let  $F^*$  be the true CDF for a random variable with  $x \in [-1,1]$ . If  $\check{F}$  is the optimal approximation of  $F^*$  in the polynomial space  $\mathcal{P}$  and  $|| F^* - \check{F} ||_2 \leq$  $\alpha$ , then with probability at least

 $1 - 2 \exp\left(-\frac{N(\eta - \alpha)^2}{16}\right) - 2(K + 1)$ 

we have  $|| F^* - \tilde{F} ||_2 \le \eta$  for  $\eta > \alpha > 0$ .

If the space is well-chosen, implying that  $\check{F}$  represents  $F^*$  well, then the DP  $\tilde{F}$ can also approximate  $F^*$  well.

Work supported by the USA NIH under Award 2R01DA040487: COINSTAC 2.0 (PI: V. Calhoun)





## Ye Tao, Anand D. Sarwate

How should we introduce privacy?



1) exp
$$(-\frac{(\eta-\alpha)^2}{4(K+1)^4\sigma^2}),$$

**<u>Result 1</u>**: PP outperforms HQ and is comparable to AQ, even better than it for certain distributions, such as Beta distribution in centralized settings.



Symposium (CSF). IEEE, 2017, pp. 263–275.





### Results

**<u>Result 2</u>**: PP outperforms both HQ and AQ in various scenarios, such as decentralized settings and those involving newly collected data.

- [1] Cynthia Dwork and Aaron Roth, "The algorithmic foundations of differential privacy," Foundations and Trends<sup>®</sup> in Theoretical Computer Science, vol. 9, no. 3–4, pp. 211–407, 2014. [2] Peter Kairouz, Sewoong Oh, and Pramod Viswanath, "The composition theorem for differential
- privacy," in International Conference on Machine Learning. PMLR, 2015, pp. 1376–1385.
- [3] Ilya Mironov, "Renyi differential privacy," in 2017 IEEE 30th Computer Security Foundations
- [4] D. G. Luenberger, Optimization by vector space methods. John Wiley & Sons, 1997.