# INSPECTRE: Privately Estimating the Unseen

Jayadev Acharya, ECE, Cornell University
Gautam Kamath, CSAIL, MIT
Ziteng Sun, ECE, Cornell University
**Huanyu Zhang, ECE, Cornell University**

## Property Estimation

- $p$: unknown discrete distribution
- $f(p)$: some property of distribution, e.g. entropy.
- $\alpha$: accuracy
- **Input:** i.i.d. samples $X_1^n$ from $p$
- **Output:** $\hat{f} : X_1^n \to \mathbb{R}$ such that w.p. at least $2/3$,

$$\left| \hat{f}(X_1^n) - f(p) \right| < \alpha$$

.

## Privacy should be concerned.

Data may contain sensitive information.

- In medical studies, data may contain health records or disease history.
- In map application, position information indicates users' residence.

**Differential Privacy:** $\hat{f}$ is $\epsilon$-differentially private (DP) if for any $X_1^n$ and $Y_1^n$, with $d_{ham}(X_1^n, Y_1^n) \leq 1$, for all measurable $S$,

$$\frac{\Pr\left(f(X_1^n) \in S\right)}{\Pr\left(f(Y_1^n) \in S\right)} \leq e^{\epsilon}.$$

## Private Property Estimation

Given i.i.d. samples from an unknown distribution $p$, the goals are:

- *Accuracy*: estimate $f(p)$ up to $\pm\alpha$ with probability $> \frac{2}{3}$.
- *Privacy*: estimator must satisfy $\epsilon$-differential privacy.

We are interested in the following properties:

- **Entropy**, $H(p)$: the Shannon entropy.
- **Support Coverage**, $S_m(p)$: expected number of distinct symbols in $m$ draws from $p$.
- **Support Size**, $S(p)$: # symbols with non-zero probability.

## Main Theorem

Informally, our upper bounds show that the cost of privacy in these settings is often **negligible** compared to the non-private statistical task. Furthermore, our upper bounds are **almost tight** in all parameters.
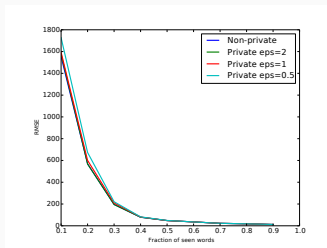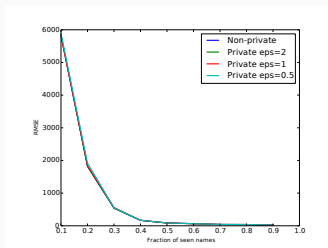
## Laplace Mechanism

Our algorithms use *Laplace Mechanism*.

- Compute a non-private estimate of the property;
- Privatize this estimator by adding Laplace noise
  $X \sim Lap(\Delta_{n,\hat{f}}/\epsilon)$.

We find estimators with **low sensitivity** for all these problems.

# Evaluation on real data

- Support coverage estimation
- Comparison on performance of private and non-private estimator
- The dataset: 2000 US Census data, and Hamlet

# The End

Details in paper online!

https://arxiv.org/abs/1803.00008