

A Firm Foundation for Private Data Analysis

wbg231

January 2023

1 Introduction

- statistical data protection has a long history because it is important
- must focus on rigorous privacy guarantees
- accessing the database should reveal nothing about any individual in it. need to separate the utility of a database from the increased risk of harm due to joining the database
- this can be achieved with low distortion
- the key idea is to randomize responses so as to effectively hide the presence or absence of any individual over the course of the life time of the database
- what does it mean to preserve privacy and how and it be accomplished

how hard is hard

- let's think about common approaches and their shortcomings

large query set

- one idea is to prevent user from making queries about specific individuals
- should not be able to search on users at all.

query auditing

- each query to a database is evaluated in the context of the query history to determine if the response would disclose information if so then refuse the query
- this is bad because it is computationally infeasible and the fact a query won't be disclosed gives information

sub sampling

- can only query from a randomized subset of individuals
- but this is not secure for those who appear in a subsample

input perturbation

- either the data or the queries are modified before the response is given
- randomized response flip a coin if it is heads always respond negatively if it is true respond truthfully half the time
- this does not work well for complex data

adding random noise to the output

- will fail if done naively
- if the noise is added at the output and is mean zero can repeat a query and average to get an approximate true value

non private database

- a system is blatantly non private if an adversary can construct a candidate database that agrees with the real database D in some large proportion of entries by querying our system
- if noise is bounded by some upper bound ϵ (that is a response can only be ϵ away from the truth) then an adversary can reconstruct the database within $4/\epsilon$ positions of the true values
- bounded noise can be defeated by some number of queries
- so need to bound the number of queries

what is hard

- linkage attack released data linked with auxiliary data to capture information about the respondents other than what is released from the database
- need to take into account auxiliary information
- anything that can be learned about an individual from a statistical database should be learnable without access to it

Differential Privacy

- DP ensures that the ability of an adversary to inflict harm (or good) to any set of people should be the same independent of whether a person is in the dataset
- consider two data bases D, D' which differ by only one row
- a function K is ϵ -Differential private if $\forall(D, D')$ and all $S \subseteq \text{rank}(K)$

$$P(K(D) \in S) \leq e^\epsilon P(K(D') \in S)$$

- this holds trivially for anyone not in the dataset
- if any one participant were to be added to the database under this regime no output would become more or less likely

achieving DP

- for a function f the ℓ_1 sensitivity of f is given by

$$\delta f = \max_{D, D'} \|f(D) - f(D')\|_1$$

- that is the sensitivity is the max ℓ_1 difference between the function with or without any one person
- with Laplacian noise we have highest density at zero, and for any $z, z' : |z - z'| \leq 1$ the density is at most e^ϵ
- also symmetric about zero
- can be extended to more dimensions
- need to add Laplacian noise to each query even if the queries are chosen adaptively with each successive query dependent on the previous answer this will work
- our expected error is the same regardless of the size of the dataset
- i think the rest is kinda in the weeds
-