



CS2080: Applied Privacy for Data Science

Intro to Membership Inference Attacks

James Honaker, Priyanka Nanayakkara, Salil Vadhan
School of Engineering & Applied Sciences
Harvard University

February 3, 2025



Takeaway Message on Reconstruction Attacks

- Every statistic released yields a (hard or soft) constraint on the dataset.
 - Sometimes have nonlinear or logical constraints \Rightarrow use fancier solvers (e.g. SAT or SMT solvers)
- Releasing too many statistics with too much accuracy necessarily determines almost the entire dataset.
- This works in theory and in practice (see readings, ps2).

How to Defend Against Reconstruction

- **Q:** what is a way that we can release many pretty-accurate estimates of proportions (counts divided by n) on a dataset while ensuring that an adversary can only reconstruct a small fraction of our dataset?
- **A:** subsample $k \ll n$ rows at random, and answer all queries using just the k rows.
 - If k is large enough (e.g. $k = 1000$), each individual proportion should be approximately preserved whp
 - We're only giving the adversary information about k rows, no info to reconstruct the others.

Subsampling vs. Reconstruction

- **Q:** If the adversary is just trying to reconstruct a single sensitive bit per individual, what fraction of the dataset should we expect the adversary to reconstruct if we subsample k rows and answer arbitrarily many counts?

- **Guess 1:** $\approx \frac{k}{n}$

- **Guess 2:** $\approx \frac{k}{n} + \frac{n-k}{n} \cdot \frac{1}{2} = \frac{1}{2} + \frac{1}{2} \cdot \frac{k}{n}$.

- **A:** $\approx \frac{k}{n} + \frac{n-k}{n} \cdot \max\{p, 1-p\} = \max\{p, 1-p\} + \underbrace{\min\{p, 1-p\} \cdot k/n}_{\text{gain over baseline}}$

if the sensitive attribute is 1 in a p fraction of the population and the adversary has no other prior information about whether individuals in the dataset are 1 or 0.

- **Q:** is subsampling a satisfactory privacy defense?