0x27 Privacy Exercise Session





Data de-anonymization

Netflix Prize

- An open competition for the best algorithm to predict user ratings for films
- Netflix Prize Dataset: Netflix released anonymous ratings of 500,000 Netflix users

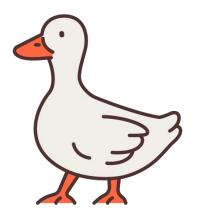
Netflix data de-anonymization

- "How to Break Anonymity of the Netflix Prize Dataset" Narayanan and Shmatikov
- correlated rankings with IMDB records to de-anonymize users





Ex1. What are Donald's favorite movies



Setting

- Goal: de-anonymize "anonymized" datasets
- Download zip file from Moodle
- "anonymized" COM-402 databases similar to Netflix DB

```
com402-1.csv
com402-2.csv
com402-3.csv
sha256(salt | email), sha256(salt | movie), date, rating
```

- Public IMDB databases
 - imdb-1.csvimdb-2.csvemail,movie,date, rating
 - o imdb-3.csv
- Task: find out what movies a user with email `donald.trump@whitehouse.gov`
 has rated.
- The answer for true movies are given in files:
 - user-1.csv
 - user-2.csv
 - user-3.csv

3 parts

Ex1-1: Dates are giving it away

Each user rated the movie at the same date in both datasets

Ex1-2: More realistic

 Dates are randomized, reflecting the fact that you might not rate a movie on Netflix and on IMDb on the same day.

Ex1-3: Even more realistic (optional)

 Dates are within 14 days, and are following a triangular distribution using Python's `random.choices` and weights: `[1, 2, ..., 14, 13, ..., 1]`.

Ex. 1.1)

$\mathrm{IMDB}_1 \subset \mathrm{COM}\text{-}402_1$

```
email, sha256(salt | email),
movie, sha256(salt | movie),
date, rating date, rating
```

- You can test with the user-*.csv file if the guesses are true
 - o assert movie_guesses == true_movies
 - o print(movie_guesses)

Ex. 1.1) - solution idea

- If the date and stars match in public and anonymized dataset, then record the candidate plaintext email and movie
 - o hash2movie[anon_entry.movie].append(pub_entry.movie)
 - o user2hash[pub_entry.user].append(anon_entry.user)
- Guess the victim's email hash as the most common candidate.
 - Get_most_common helper method
- Filter the ratings made by this victim's email hash

Ex. 1.2)

$IMDB_2 \subset COM-402_2$

```
email, sha256(salt | email),
movie, sha256(salt | movie),
date, randomize(date),
rating
```

Where $\operatorname{randomize}:\mathcal{D}\to\mathcal{D}$ maps each date (uniquely) to some random date

Ex. 1.2) - solution idea

- Match movie hashes to plaintexts by frequency.
 - Sort_by_freq method in helpers.py
 - o movie2hash[movie_name] = movie_hash
 - o hash2movie[movie_hash] = movie_name
- Get the victim's movie hashes from the public data.
 - filter_ratings_by_user(public, victim_email)
- Identify the victim's email hash:
 - Map the email hashes to the corresponding movie hashes from the anonymized data.
 - find the hash that has all the movie hashes from above.
 - If multiple possibilities, then the generated datasets are not good.

Ex. 1.3) [OPTIONAL]

$\text{IMDB}_3 \subset \text{COM-402}_3$

```
email, sha256(salt | email),
movie, sha256(salt | movie),
date, real_randomize(date),
rating
```

Where $real_randomize : \mathcal{D} \to \mathcal{D}$ maps randomly maps each date according to a triangular distribution within 14 days of the initial date

Ex2. Differentially Private Queries

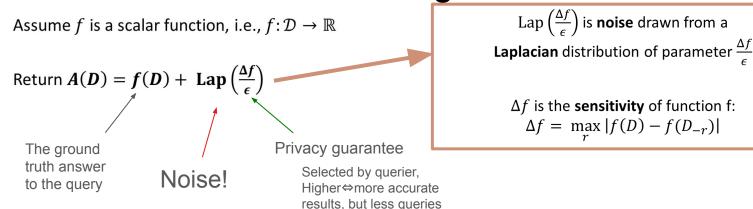
[optional]

Ex2



- Playing the role of the company: IMDB
- Researchers send you queries, you respond in a differentially private way
- Query: "For a given movie, how many reviews are above a threshold?"
 - o get count(movie name, rating threshold, epsilon)
- Each researcher gets their own **privacy budget** ε_{total} : the total epsilon a researcher can use across queries.
- for each query, an epsilon value can be specified
 - the higher epsilon is, the higher the accuracy of the response will be, but will allow for less queries. and vice versa, a lower epsilon will allow for more queries.

Laplace mechanism for adding noise



- ullet In your database, a user can only rate one movie once. This means sensitivity Δf =1
- np.random.laplace(loc=0, scale=1. / epsilon)

Sequential composition

- Sequential composition property: If algorithms A₁, A₂, ..., A_k use independent randomness and each A_i satisfies ε_i-differential privacy, respectively. Then (A₁, A₂, ..., A_k) is (ε₁+,ε₂+...+ ε_k)-differentially private
 - can keep account of spent privacy budget and ensure that the queries do not exceed it.

Testing with verify.py

b (base) → solution git:(main) x /usr/local/anaconda3/bin/python Running tests... > Basic privacy budget accounting PASSED 🌖 > Privacy budget depletion control PASSED ঙ > DP noise distribution PASSED 🌖 > Multiple query behavior (I) PASSED 🌖 > Multiple query behavior (II) PASSED 3