From Astrophysics to Differential Privacy

Data science seminar DESY Zeuthen, April 2020 Matteo Giomi

About me

Science: time-domain astrophysics.



Now: privacy researcher.

Statice

Startup developing software for **privacy-preserving** data sharing.

Privacy

"Privacy is the ability of an individual to seclude themselves or information about themselves, and thereby express themselves selectively." (<u>wikipedia</u>)

Lack of privacy \rightarrow behavioural change.



Privacy in the digital era

Every interaction with technology creates data about the user.

- In the wrong hands data can be used for blackmail, social engineering, mass surveillance and the like.
- If used correctly, data can also lead to collective benefits.

Complete non-disclosure is not the best option. Also,

- We do share sensitive data with strangers (i.e. doctors)

How to share data in a privacy preserving way?

Example dataset

| medical condition | zip code | sex | birth year | race | phone |
|-------------------|----------|-----|------------|-------|-----------|
| chest_pain | 1203002 | f | 1964 | white | 015940192 |
| obesity | 1203505 | f | 1964 | white | 010405919 |
| short_breath | 1203106 | f | 1964 | white | 011500159 |
| heart_disease | 5403221 | m | 1965 | black | 010192042 |
| heart_disease | 5403221 | m | 1965 | black | 015909191 |
| heart_disease | 5403221 | m | 1965 | black | 015553436 |
| ovarian cancer | 3003202 | f | 1960 | white | 016901095 |
| ovarian cancer | 3003555 | f | 1960 | white | 017497297 |
| prostate cancer | 3003890 | m | 1960 | white | 018206810 |

Identifiers and quasi-identifiers

| | pho | ne | race | birth year | sex | zip | code | medica | I condition | |
|---|-------------------------------------|----|-------|-------------|------|-----|------|-----------|-------------|---|
| | 015940192 010405919 011500159 | | white | 1964 | f | 120 | 3002 | | chest_pain | |
| | | | white | 1964 | f | 120 | 3505 | | obesity | |
| | | | white | 1964 | f | 120 | 3106 | S | hort_breath | |
| Personally identifying ¹² information (PII) ⁹¹ | | 2 | black | "Ouasi" ide | 3221 | | he | Sensitive | | |
| | | 91 | black | | | 010 | 3221 | he | informatic | n |
| | 015553436 | | black | 1965 | m | 540 | 3221 | he | art_disease | |
| 01690109 01749729 | | 95 | white | 1960 | f | 300 | 3202 | ova | rian cancer | |
| | | 97 | white | 1960 | f | 300 | 3555 | ova | rian cancer | |
| | 0182068 | 10 | white | 1960 | m | 300 | 3890 | pros | tate cancer | |

"Sanitizing" a dataset

| | pho | ne | race | birth year | sex | zip | code | medica | al condition | |
|-------------------------------------|----------|-----|-------|------------|--------|---------|------|--------|--------------|---|
| | 15940192 | | white | 1964 | f | 120 | 3002 | | chest_pain | |
| | 0 04059 | 19 | white | 1964 | f | 1203505 | | | obesity | |
| | 01 500 | 59 | white | 1964 | f | 120 | 3106 | s | hort_breath | |
| Personally identifying ² | | 2 | black | "Quasi" id | entifi | ers | 3221 | he | Sensitive | |
| information (P | II) | 91 | black | Quubi iu | | 010 | 3221 | he | informatic | n |
| | 0155531 | 36 | black | 1965 | m | 540 | 3221 | he | art_disease | |
| | 0159010 | 95 | white | 1960 | f | 300 | 3202 | ova | arian cancer | |
| | 174972 | 297 | white | 1960 | f | 300 | 3555 | ova | arian cancer | |
| | 0182068 | 310 | white | 1960 | m | 300 | 3890 | pros | state cancer | |

Re-identification via linkage

Even without PII individuals can be re-identified by linking with external information.



"We are all special"

Given enough quasi-identifiers everyone is unique \rightarrow can be re-identified with certainty.



Rocher, L., et al. Estimating the success of re-identifications in incomplete datasets using generative models.

Removing PII is not enough

The notion of PII has no technical meaning: <u>everything is PII!</u>

However:

"We do not share information of data in any personally identifiable form.."

Arvind Narayanan and Vitaly Shmatikov, Myths and and fallacies of "Personally identifiable information"

Old solution: k-anonymity

race

Avoid unique joints: "any combination of quasi-identifiers must appear at least *k* times"

birth year sex zip code medical condition

| chest_pain | 1203* | * | 1964 | white |
|-----------------|-------|---|------|-------|
| obesity | 1203* | * | 1964 | white |
| short_breath | 1203* | * | 1964 | white |
| heart_disease | 5403* | * | 1965 | black |
| heart_disease | 5403* | * | 1965 | black |
| heart_disease | 5403* | * | 1965 | black |
| ovarian cancer | 3003* | * | 1960 | white |
| ovarian cancer | 3003* | * | 1960 | white |
| prostate cancer | 3003* | * | 1960 | white |

P. Samarati and L. Sweeney, Protecting Privacy when Disclosing Information: k-Anonymity and its Enforcement through Generalization and Suppression

No solution: k-anonymity

Problems with k-anonymity: lack of diversity, background knowledge

| phone | race | birth year | sex | zip code | | white | 1964 | * | 1203* | chest_pain |
|-----------|--------|------------|-----|----------|---|-------|------|---|-------|-----------------|
| 015940192 | white | 1964 | f | 1203002 | | white | 1964 | * | 1203* | obesity |
| | | | | | | white | 1964 | * | 1203* | short_breath |
| | | | | | | black | 1965 | * | 5403* | heart_disease |
| phone | race | birth year | sex | zip code | | black | 1965 | * | 5403* | heart_disease |
| 015909191 | black | 1965 | f | 5403014 | | black | 1965 | * | 5403* | heart_disease |
| 018206810 | white | 1960 | m | 3003890 | | white | 1960 | * | 3003* | ovarian cancer |
| 010200010 | millio | 1000 | | | | white | 1960 | * | 3003* | ovarian cancer |
| | | | | | × | white | 1960 | * | 3003* | prostate cancer |

race birth year sex zip code medical condition

Differential privacy

Suppose I have a secret and I'm considering whether to share my data. I wish that

Pr(someone guess my secret | data) ~ Pr(someone guess my secret)



Dwork C., et al., (2006) Calibrating Noise to Sensitivity in Private Data Analysis

Differential privacy

Suppose I have a secret and I'm considering whether to share my data. I wish that

Pr(someone guess my secret | data) ~ Pr(someone guess my secret)

The ε parameter is the 'privacy guarantee':

- The smaller the ε the stronger the privacy.

$$\frac{\Pr[f(D)=R]}{\Pr[f(D\prime)=R]} \leq e^{\varepsilon}$$

DP example: laplacian mechanism

Adds noise draws from laplacian distribution:

 $M(f, \varepsilon, D) = f(D) + Laplace(\Delta f/\varepsilon)$

- Δf is the <u>sensitivity</u> of the query f: maximum change in its output due to single records in the dataset.
- For a given epsilon, high-sensitivity queries requires more noise.



Toy example

Use laplacian mechanism to estimate the average income of a population.



Problems: outliers

To mask an outlier much more noise have to be added.



Differential privacy

DP is a property of the analysis and is based on the <u>addition of 'just enough' noise</u>.

Pros:

- Makes no assumption on the attacker.
- Robust against post-processing.
- Different DP analysis compose.

Cons:

- Outliers requires large amount of noise.
- Noise can be averaged out via multiple queries.

Privacy-preserving synthetic data

- Use deep generative models to learn the data-generating distribution.
- Sample from this distribution to obtain synthetic data.



Deep learning with differential privacy

Models learn by minimizing a loss function.

Minimization via stochastic gradient descent (SGD).

The model interacts with the data <u>exclusively via</u> <u>the gradients</u>

Fitting a model with DP:



- Each gradient is clipped to a maximum length (fix the sensitivity).
- Noise is added to each gradient component.

This makes sure that the model won't learn from individual training examples. M. Abadi et al, Deep Learning with Differential Privacy (2016)

What am I doing every day

Implement DP algorithms and develop privacy evaluations (attack is the best defence)



Transition to industry: applying for a job

Before you start:

- Time (3 to 6 months) and patience.
- Sleek 1-page CV, honest cover letter (sometimes not asked).
- Profile on linkedin, glassdoors and the like.

Typical job interview:

- Initial call.
- Code challenge / test task.
- Final interview.

Job interview is a skill to be learned: don't be afraid to apply generously.

Transition to industry: tips and tricks

This is what I gathered

- For position you really like, get in contact through company website.
- Sometimes it's hard to convince people that you can "deliver".

Good to have:

- Familiarity with usual data-science stack: numpy, pandas, seaborn, sklearn.
- Plenty of material on machine learning (Andrew Ng on yt, for example).
- Plenty of datasets and tasks on <u>kaggle</u>, start playing around.
- Luck.

Thanks for the attention