

# The Role of Differential Privacy in GDPR Compliance

Rachel Cummings and Deven Desai, Georgia Tech

(talk by Yatharth Dubey)

FATREC – Oct. 6, 2018

# EU General Data Protection Regulation (GDPR)



GDPR asserts that individuals “shall have the right to obtain [...] the erasure of personal data concerning him or her without undue delay”

# Compliance?



Disconnect between legal language and machine learning:  
What algorithms can I run on my data?

# Personal data vs aggregates

- Personal data must be deleted upon user request
- Aggregates may be stored longer for “collection and the processing of personal data necessary for statistical surveys or for the production of statistical results”

# Aggregate data must be anonymous

Former Chief Privacy Officer of Microsoft asserted that GDPR-compliant aggregate data:

1. must not be “directly linked to identifying data”
2. must not be a “known, systemic way to (re)identify data”
3. must not “relate to a specific person”

# Alternative: pseudonymization

- Pseudonymization is “processing of personal data in such a manner that the personal data can no longer be attributed to a specific data subject without the use of additional information”
- GDPR allows for pseudonymization of aggregate data
- Allows for linkage attacks, hopes they don't happen



# Memoization in ML

- Many learning algorithms **memoize** individual data entries during training inadvertently by imbedding personal data in the learned model
- Deep learning algorithms for word prediction leaked SSNs from the training corpus [CLKES '18]
  - Complete: “My Social Security Number is...”

Unsurprising for ML, bad for privacy

# Memoization and GDPR

- Machine learned model that memoizes personal data cannot be an aggregate
- Individual data has not been de-identified and/or can be re-identified
  - Can extract Personally Identifying Information (PII) from model

Need formal guarantee to prevent memoization



# Differential privacy [DMNS '06]

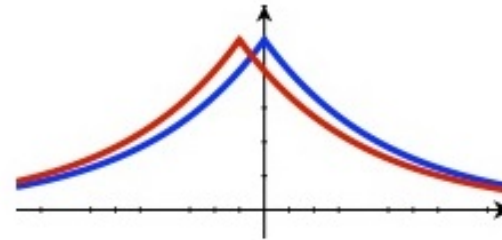
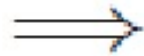
Bound the “maximum amount” that one person’s data can change the output of a computation

An algorithm  $M: T^n \rightarrow R$  is  $(\epsilon, \delta)$ -**differentially private** if  $\forall$  neighboring  $x, x' \in T^n$  and  $\forall S \subseteq R$ ,

$$P[M(x) \in S] \leq e^\epsilon P[M(x') \in S] + \delta$$

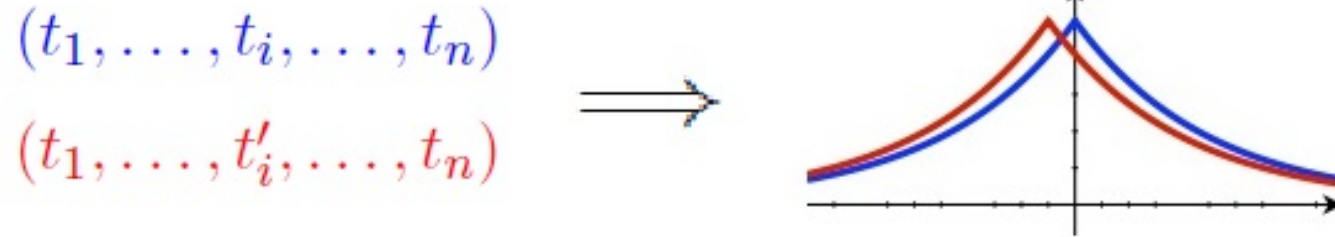
$(t_1, \dots, t_i, \dots, t_n)$

$(t_1, \dots, t'_i, \dots, t_n)$



- $S$  as set of “bad outcomes”
- Worst-case guarantee

# Differential privacy [DMNS '06]



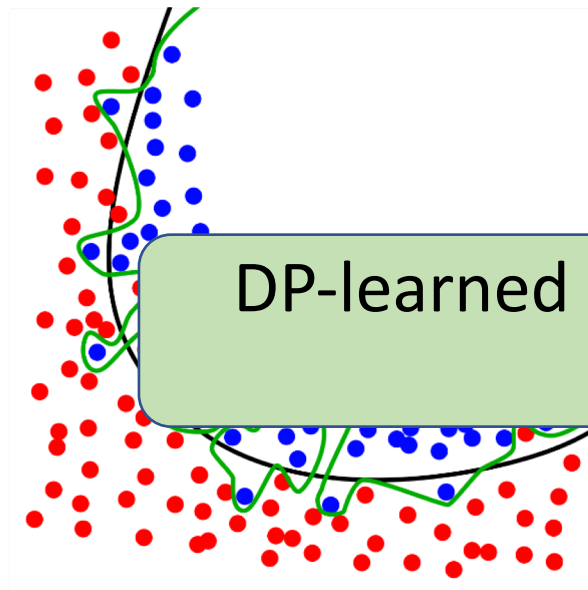
“You will not be affected... by allowing your data to be used... **no matter what other information sources are available.**”

DP addresses the paradox of learning nothing about an **individual** while learning useful information about a **population.**”

The Algorithmic Foundations of Differential Privacy, Dwork and Roth.

# DP formally prevents memoization

- Constrained to learn the same thing without your data
  - e.g., won't output your SSN if it wasn't in corpus
- Theorem [DFHPRR '15][CLNRW '16]\*: An  $\epsilon$ -differentially private algorithm cannot overfit its training set by more than  $\epsilon$ .



DP-learned models are GDPR-compliant aggregate data

\*Lots of interesting details missing!

# Future policy challenges

## 1. How to set $\epsilon$ ?

- Theory:  $\epsilon$  is small constant  $\ll 1$  (e.g., 0.01) or diminishing in  $n$  (e.g.,  $O(1/\sqrt{n})$ )
- Practice:  $\epsilon$  is large (e.g., Apple uses  $\epsilon=42$ , Census uses  $\epsilon \approx 9$ )
- Trade-off between privacy and accuracy

## 2. How to set $\delta$ ?

- $\delta$  is failure probability of privacy guarantee
- Allowing  $\delta = o(\exp(-n))$  can significantly reduce  $\epsilon$  for same accuracy level
- Practical GDPR-compliant data erasures = encrypt data and delete key
- Cryptographically small failure probability acceptable?

# The Role of Differential Privacy in GDPR Compliance

Rachel Cummings and Deven Desai, Georgia Tech

(talk by Yatharth Dubey)

FATREC – Oct. 6, 2018