

Privacy-Preserving Data Analysis & Security by Design

Daniel Kraschewski

Big Techday 11

2018-05-18



personal data is omnipresent

- internet browsing history
- cell phone movements
- smart metering, smart homes, IoT
- social media, cloud
- ...

Photo: "Base." by Instant Vantage (CC BY-SA 2.0), clipped to fit page layout

customer analytics



city planning



medical surveys



social surveys



Photos by @nordwood, @hellocolor, @rawpixel, @jaseess on Unsplash

GDPR (from 2018-05-25)

- strong notion of consent
 - informed
 - freely given
 - specific
 - unambiguous
 - clear affirmative act
- high fines (up to 4% of annual turnover) for data privacy violations



Privacy-Preserving Data Analysis

- derive large-scale statistical insights
- still preserve individual's privacy
- security by design
- provable security/privacy



privacy-preserving results

- aggregated statistics must not reveal personal information
- proper formal notion of anonymity
- Differential Privacy, Laplace Mechanism

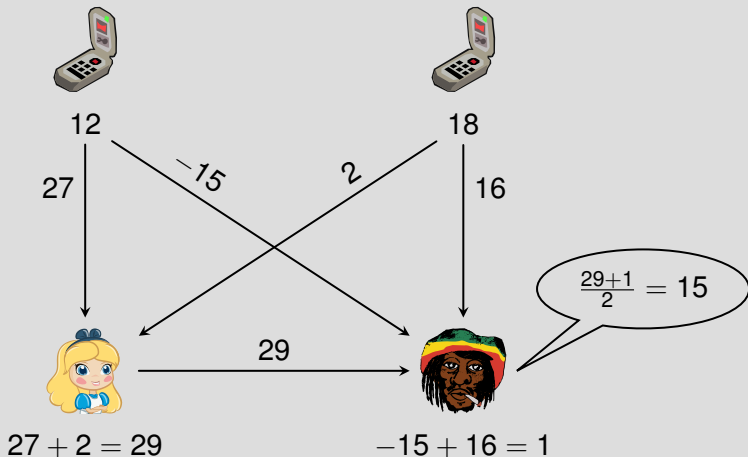
privacy-preserving computation

- personal data must be unaccessible even for system owner
- Cryptography
- Secret Sharing, Secure Multi-Party Computation

privacy-preserving environment

- no unauthorized third party may access any data
- IT-Security
- Access Control, encrypted & authenticated channels

Example: Privacy-Preserving Averaging



general setting

- mutually mistrusting parties P_1, P_2, \dots
- secret inputs x_1, x_2, \dots
- want to compute some agreed on function value $f(x_1, x_2, \dots)$
- nothing but $f(x_1, x_2, \dots)$ should be revealed about x_1, x_2, \dots

a universal solution

- write f as arithmetic circuit
- transform each x_i into unintelligible *Secret Sharing*
- evaluate f gate-by-gate, preserving *Secret Sharing*
- recombine result

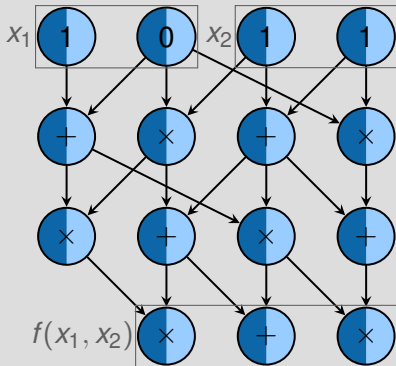


Photo by Sahmeditor on Wikimedia Commons
(CC BY-SA 2.0), clipped to fit page layout

Shared Addition



x_A, y_A

x_B, y_B

z_A

z_B

s.t. $z_A + z_B = x + y$

easy: $z_A = x_A + y_A$ and $z_B = x_B + y_B$

Shared Multiplication



x_A, y_A

x_B, y_B

z_A

z_B

s.t. $z_A + z_B = x \times y$

problematic: $z = x_A y_A + x_A y_B + x_B y_A + x_B y_B$

missing building block



v_A

v_B

w_A

w_B

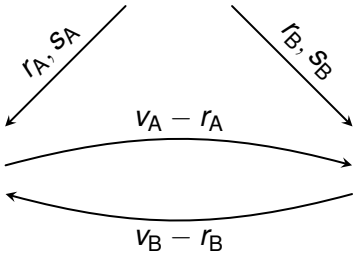
s.t. $w_A + w_B = v_A v_B$

random r_A, r_B, S_A
 $S_B := r_A r_B - S_A$



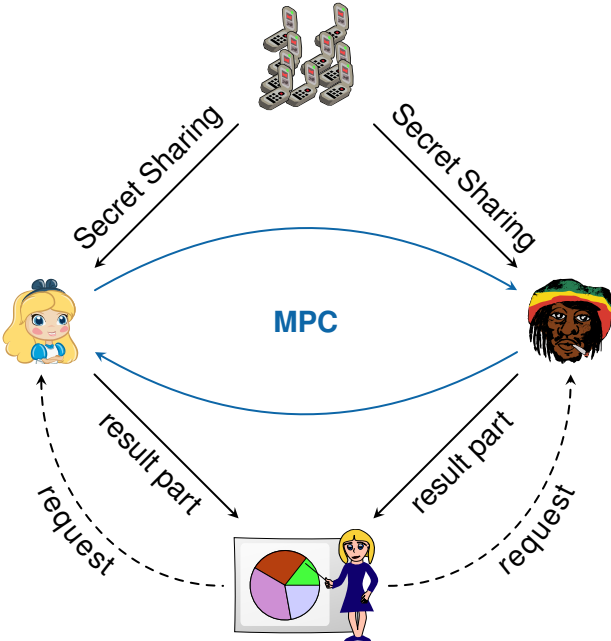
Invariants

- $r_A r_B = S_A + S_B$
- $V_A r_B = S_A + W_B$
- $V_A V_B = W_A + W_B$



$W_A := S_A + V_A(V_B - r_B)$

$W_B := S_B + r_B(V_A - r_A)$

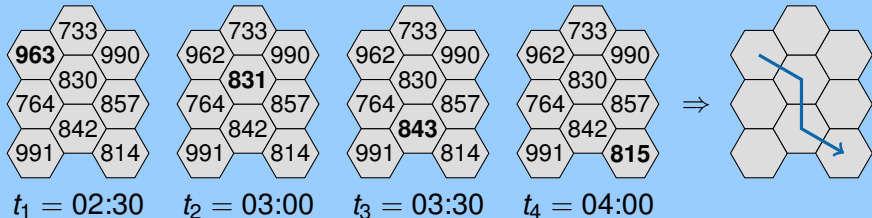


k -Anonymity

- published data must coincide with at least k individuals

De-anonymization attack on correlated data¹

- published data: number of people in mobile cell at time t_i



- trajectory recovery $\hat{=}$ optimization problem
 - higher “costs” for sudden/far movements
 - higher “costs” for irregular movements and/or movements at night
- 50% – 91% accuracy, depending on space-time resolution

¹Fengli Xu, Zhen Tu, Yong Li, Pengyu Zhang, Xiaoming Fu, Depeng Jin: Trajectory Recovery From Ash: User Privacy Is NOT Preserved in Aggregated Mobility Data, 26th International Conference on World Wide Web (WWW 2017)

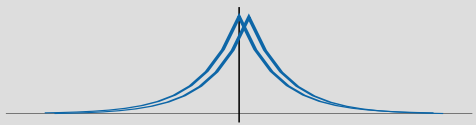
ϵ -Differential Privacy

- statistical similarity: $\kappa(\text{real data}) \approx \kappa(\text{real data} \setminus \text{me})$ up to factor e^ϵ

Laplace Mechanism

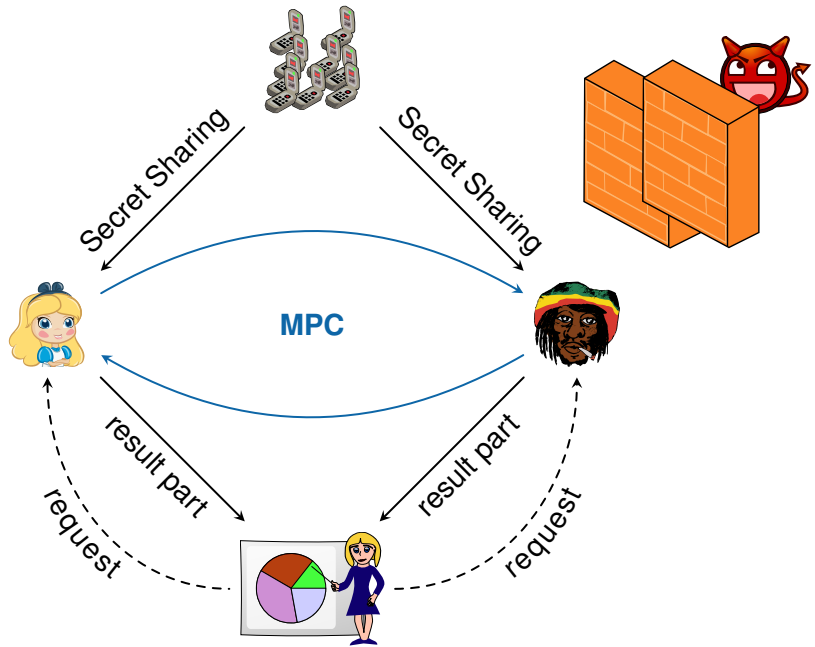
- calculate histogram
- add Laplace noise
- output noisy group sizes

Laplace Distributions



Example histogram ($\frac{1}{10}$ -Differential Privacy)

Age	Sex	Diagnosis	count	noise	result	
< 35	f	infection	48	9	57	
< 35	f	NCD	61	-1	60	
< 35	m	infection	75	-5	70	
< 35	m	NCD	44	-7	37	
<hr/>						
≥ 35	f	infection	165	6	171	
≥ 35	f	NCD	127	-4	123	
≥ 35	m	infection	228	2	230	
≥ 35	m	NCD	168	-2	166	



integration

simplicity

Strictness

- fail-safe defaults
- need-to-know principle
- principle of least privilege

Robustness

- separation of duties
- multi-factor/layered security
- forward secrecy

Consistency

- complete security-model
- defense in depth
- homogeneity/uniformity



Photo by BalticServers.com on Wikimedia Commons (CC BY-SA 3.0), clipped to fit page layout

Summary

- large-scale statistics can be calculated in a privacy-preserving way
- security/privacy by design, not just by contract
- security/privacy mathematically defined and provable
- though, inefficient universal solutions

Improvements

- less generic, optimized MPC constructions
- less MPC, more IT-Security (e.g., “self-sealing” hardware)
- tailored Differential Privacy mechanisms
- ...

Thank you for your attention!



Dr. Daniel Kraschewski
Senior Consultant

TNG Technology Consulting GmbH
Betastraße 13a
85774 Unterföhring

tel: +49 89 2158 9960

fax: +49 89 2158 9969

daniel.kraschewski@tngtech.com