

PDS, 1.12.21

7) Pseudonymization

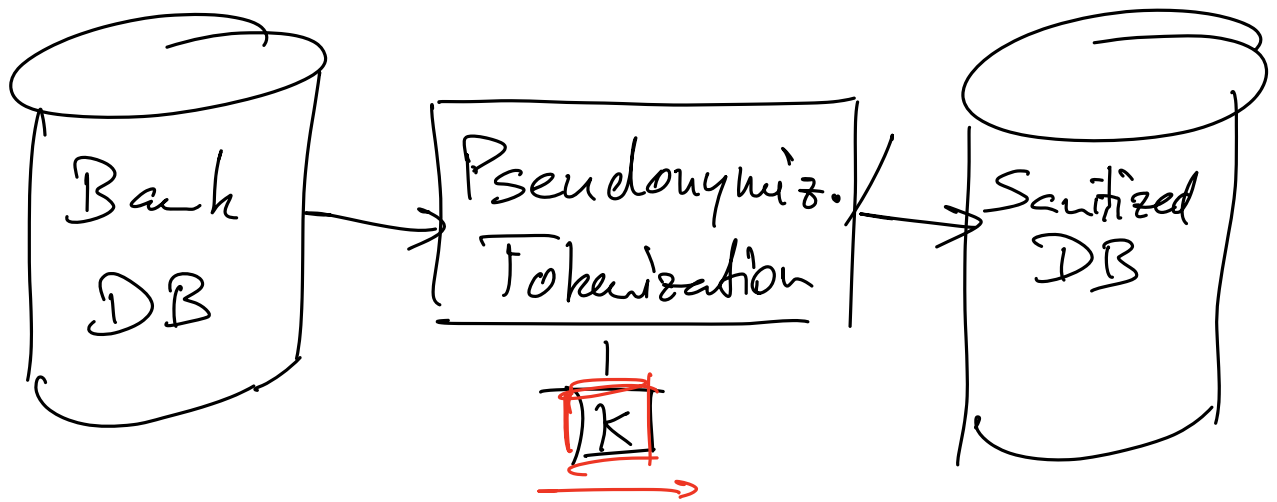
Recall GDPR:

- Personal data: personal, identifiable
- **Pseudonymous data**: can no longer be linked to a person, info w/ add. info.
- Anonymous data: can no longer be linked to data subject

Focus on pseudonymisation methods

- preserve (some) structure and utility of dataset
- not change format
- cryptographic guarantees

Ex.



- testing
- analytics
- legal attacks

foreign law

Methods overview

- Encryption whole dataset
- Redaction or masking
- Inversible tokenization
(collision-free, one-way func.)

• Simple

↓ complex

- no utility
- backup

- loss of utility

- preserves struct. & referential integrity
- may preserve format

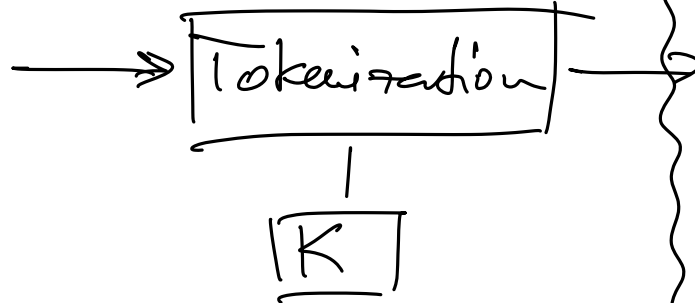
- Reversible tokenization
(using symmetric-key encryption)

- preserves struct. & referential integrity
- may preserve format

Security model

Dataset

- table
- fields



Pseudonymized dataset

trusted

untrusted

Basic method: Table lookup

- random values as tokens
- stores all mapping
- reversible
- costly

7.2) Irreversible tokeization

Hash function

$$H: \{0, 1\}^* \longrightarrow \{0, 1\}^k$$

...
maps arbitrarily long strings to k -bit strings,
e.g. $k=256$ in SHA-256

Security

- Collision-resistance: it is infeasible to find $x \neq x'$ s.t.

$$\Downarrow \quad H(x) = H(x').$$

- Second pre-image resistance: Given x , it is infeasible to find some $x' \neq x$

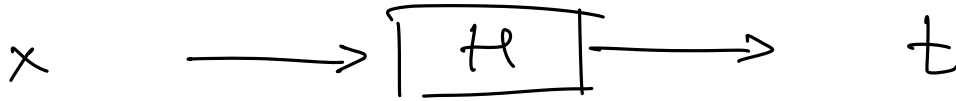
$$\Downarrow \quad \text{s.t. } H(x) = H(x')$$

- One-way: given some $(\dots) h \in \{0, 1\}^k$ it is infeasible to find x s.t. $H(x) = h$.

Hash-based tokenization

Dataset

Pseudonyms



- Irreversible
 - Format : k -bit string
 - Reconstruction attack via enumeration of possible inputs
-
- Pseudorandom functions

PRF

$$F : \{0,1\}^x \times \{0,1\}^m \rightarrow \{0,1\}^k$$

\uparrow
key k

Security:

An adversary A may repeatedly obtain $F_k(x)$ or R , where $R \leftarrow \{0,1\}^k$;

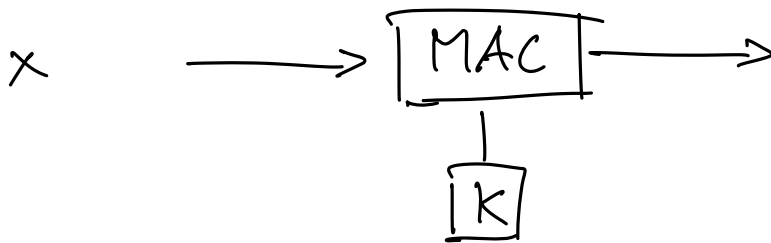
A cannot guess with prob. significantly above $\frac{1}{2}$ which is the case.

Implementations of PRFs

- MAC (message authentication code)
 - ↳ HMAC
- Blockcipher (AES)

PRF-based tokenization

Dataset



Pseudonyms

t

- Irreversible, even with key
- Outputs a **k-bit string**
- Additional security over the hash-based tokenization due key k

Preserving the format

- Database fields have structure
credit-card no., IBAN, ISBN ...
- Crypto primitive outputs k -bit string..
a number in $\mathcal{T} = \{0, \dots, 2^k - 1\}$
- Structure determined by a set
 $\mathcal{M} \subseteq \mathcal{T}$
with efficient test $x \in \mathcal{M}$

Cycle-walking algorithm

// cryptogr. primitive $F: \mathcal{T} \rightarrow \mathcal{T}$
(or $\{0,1\}^* \rightarrow \mathcal{T}$)

// is collision-free

// $|\mathcal{T}| = O(|\mathcal{M}|)$

map (F, s)

$t \leftarrow s$ // $s \in \mathcal{M}$

repeat

$t \leftarrow F(t)$

until $t \in \mathcal{U}$

return t

unmap (F^{-1}, t)

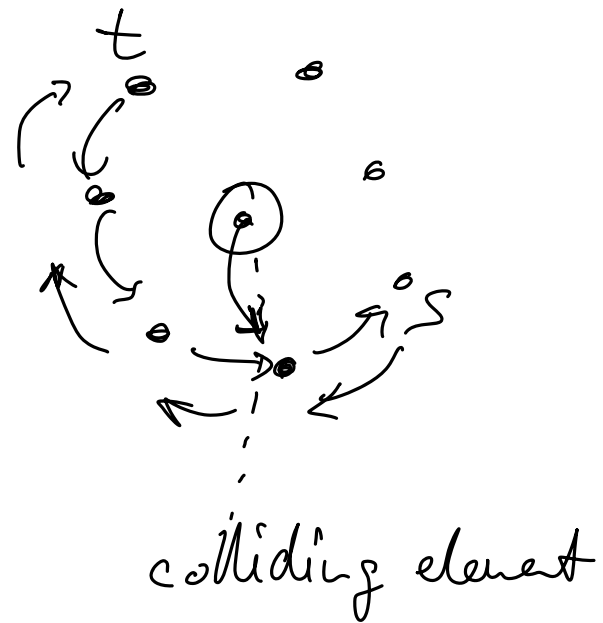
$s \leftarrow t$

repeat

$s \leftarrow F^{-1}(s)$

until $s \in \mathcal{U}$

return s



Intuition: F is a perm. on \mathcal{T} ,
we walk on a cycle of F

7.3) Reversible tokenization

How to encrypt on a small domain?

- Typically, block cipher uses 128-bit strings

- Need a PRP (pseudo-rand. perm.)
on the much smaller \mathcal{T}

$$\mathcal{T} = \{0, 1, \dots, N-1\}$$

Alg. FE : Feistel-Encrypt

- Use $a, b \in \mathbb{N}$ with $a \cdot b \geq N$
and $T = \{0, \dots, N-1\}$
- PRF $B : \mathcal{K} \times \{0, 1\}^* \rightarrow \{0, 1\}^k$

Alg. FE-Enc. (z, a, b, m) $z \in \mathcal{K}$ key
 $m \in T$ msg.

$$L_0 \leftarrow m \text{ div } b$$

$$R_0 \leftarrow m \text{ mod } b$$

for $i = 1, \dots$, rounds do $// \geq 3$

if i odd then $s \leftarrow a$ else $s \leftarrow b$ fi

$$L_i \leftarrow R_{i-1} \dots$$

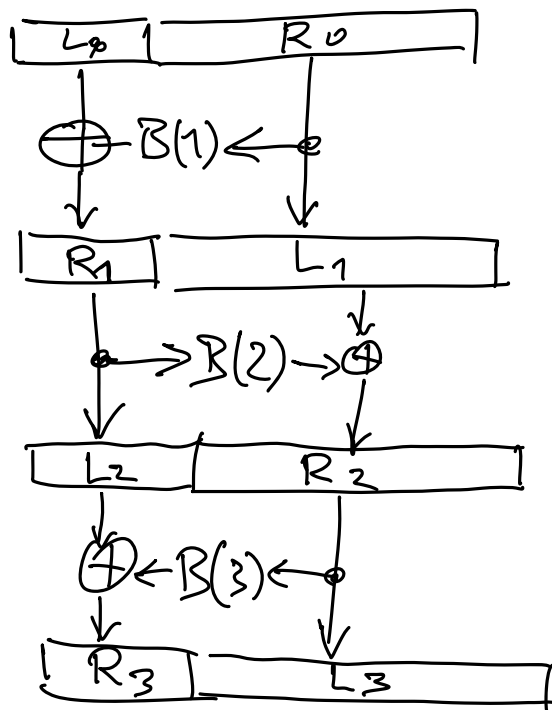
$$R_i \leftarrow (L_{i-1} \oplus B(z, i \| \dots \| R_{i-1})) \text{ mod } s$$

return $s \cdot L_{\text{rounds}} + R_{\text{rounds}}$
 $\in \{0, \dots, a \cdot b - 1\}$

Alg. FE-Dec (z, a, b, c) works accordingly

Ex.

$a \ll b$
.....



$i = 1$

$i = 2$

$i = 3$

- Format-preserving encryption uses Alg. FE plus cycle-walking algorithm