



Fingerprinting Relational Databases

Quality Evaluation and Impact on Learning Tasks

Tanja Šarčević

Outline



Introduction



Fingerprinting techniques



Robustness Evaluation



Data Utility Evaluation



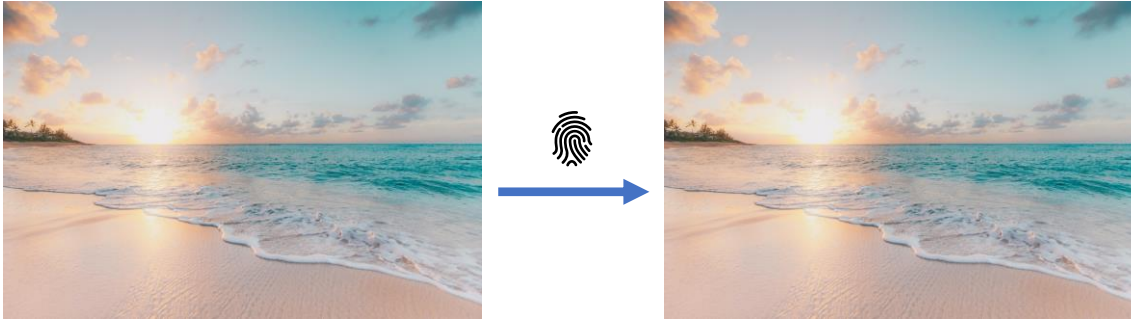
Conclusion and Future Work



Use case: medical data

- Medical data – sensitive and a valuable asset to the hospital
 - The owner wants to share data with researchers
 - The owner wants to have a claim of the ownership and trace the source of an unauthorised publishing
- **Idea:** hide a piece of information in the data before the distribution

Fingerprinting



Age	BloodPress.	Diabetes
32	64	1
31	66	0
50	72	1
48	70	0



Age	BloodPress.	Diabetes
33	64	1
31	68	0
50	72	1
47	70	0

Fingerprinting

FINGERPRINT - string of bits containing information about the owner and the recipient of the specific data copy

FINGERPRINTING - an information hiding technique that embeds the fingerprint into the data



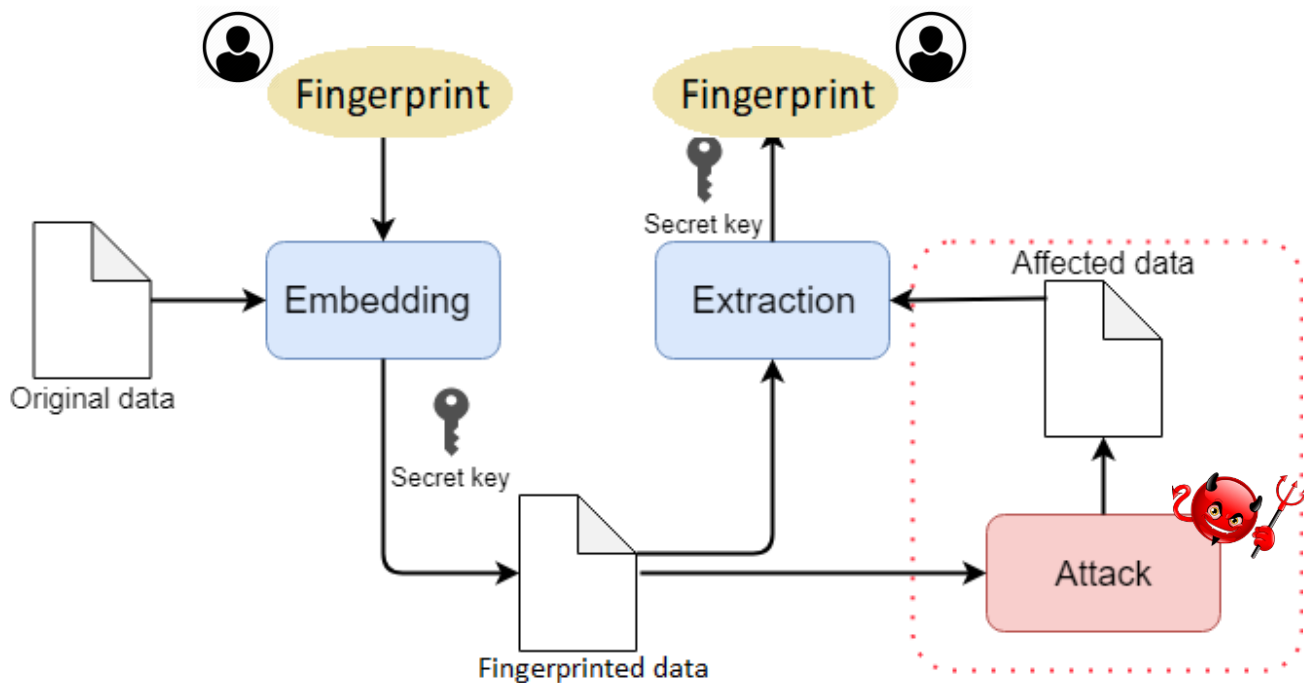
Fingerprinting

Requirements of the fingerprint:

- 1) recognizable by the owner
- 2) not detectable and consequently removable by the recipients
- 3) robust to attacks
- 4) does not change the utility of the data too much



Fingerprinting Schemes: workflow



Fingerprinting Schemes: Numerical data

AK (Agrawal and Kiernan) Scheme:

- Pseudorandom choice of a row, a column and a least significant bit of a value to be marked

Age	BloodPress.	Diabetes
33	64	1
31	66	0
50	71	1
48	70	0

Block Scheme:

- The data is first divided into blocks
- Pseudorandom choice of the value to be marked within every block

Age	BloodPress.	Diabetes
32	64	1
33	66	0
49	72	1
48	70	0

Two-level Scheme:

- 1st layer:** Pseudorandomly selects the values to be marked; this pattern identifies the owner
- 2nd layer:** Pseudorandomly selects the values to be marked; this pattern identifies the recipient

Age	BloodPress.	Diabetes
31	64	1
31	69	0
50	72	1
45	70	0

Numerical vs. Categorical Data

- Decimal:

32.3  → 32.7 😊

- Integer:

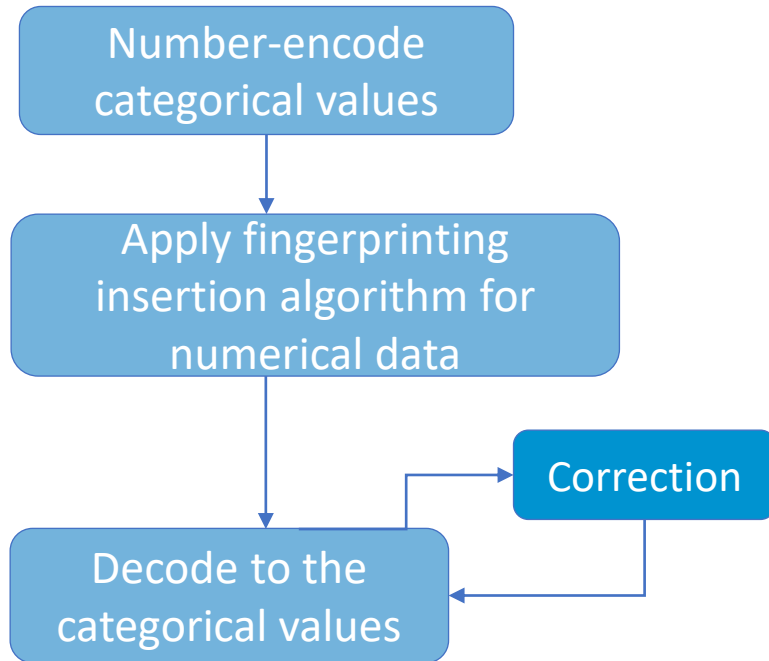
32  → 34 😊

- Non-numerical:

France  → Germany ?

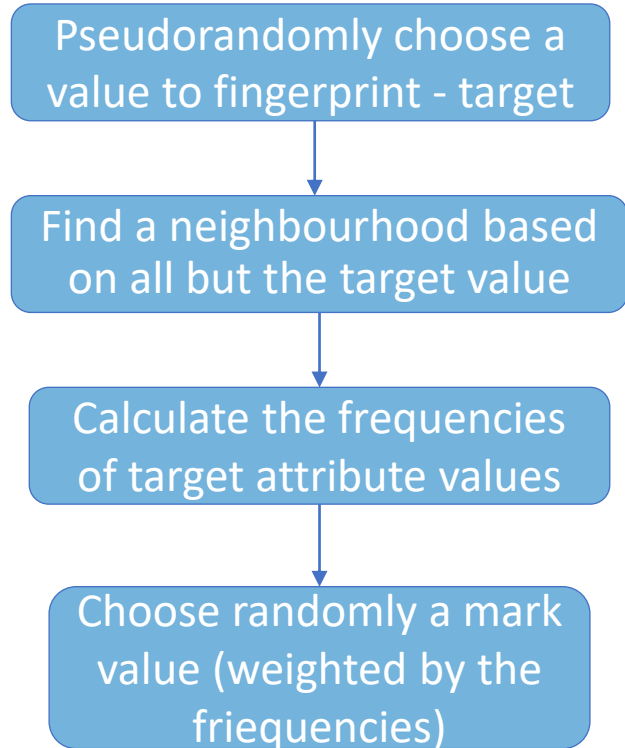
Fingerprinting categorical data: Approach #1

Insertion



- London, Paris, Vienna
- Encoding: 0,1,2 (00,01,**10**)
- Fingerprinting: 10 -> 11
- Decoding: 11 -> ??
- Correction: $11 \bmod 3 = 00$
- Decoding: 00 -> London

Fingerprinting categorical data: Approach #2



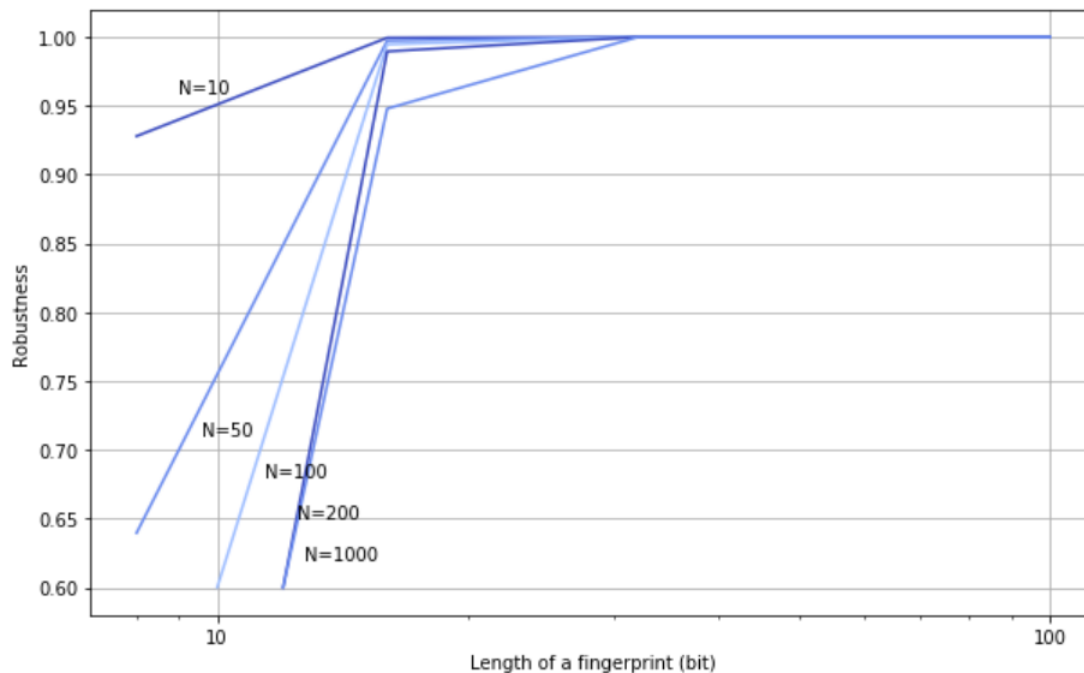
- Addressing the problem of semantic relations between categorical attributes that can be disturbed by fingerprinting
- *gender*: male, *pregnant*: yes
- The fingerprinted value will be a value that is likely to occur in a combination with other values from the row

Robustness Evaluation: Attacks

- Malicious operations on the fingerprinted dataset with the goal of
 - disabling the extraction of the correct fingerprint or
 - disabling association of a fingerprint with the correct recipient
- **Subset Attack**
- **Superset Attack**
- **Bit-flipping Attack**
- **Additive Attack**
- **Misdiagnosis False Hit Rate:** measures the likelihood of extracting a valid fingerprint from non-fingerprinted data

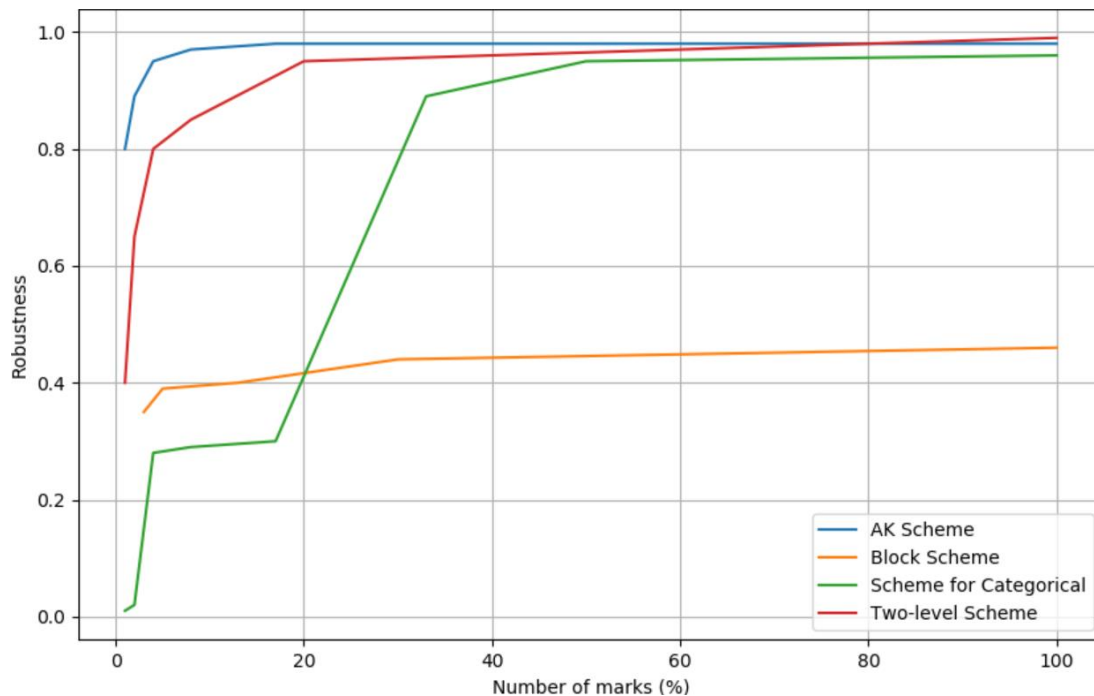
Robustness Evaluation

Misdiagnosis false hit



Robustness Evaluation

Subset Attack



Utility Evaluation

Mean and Variance

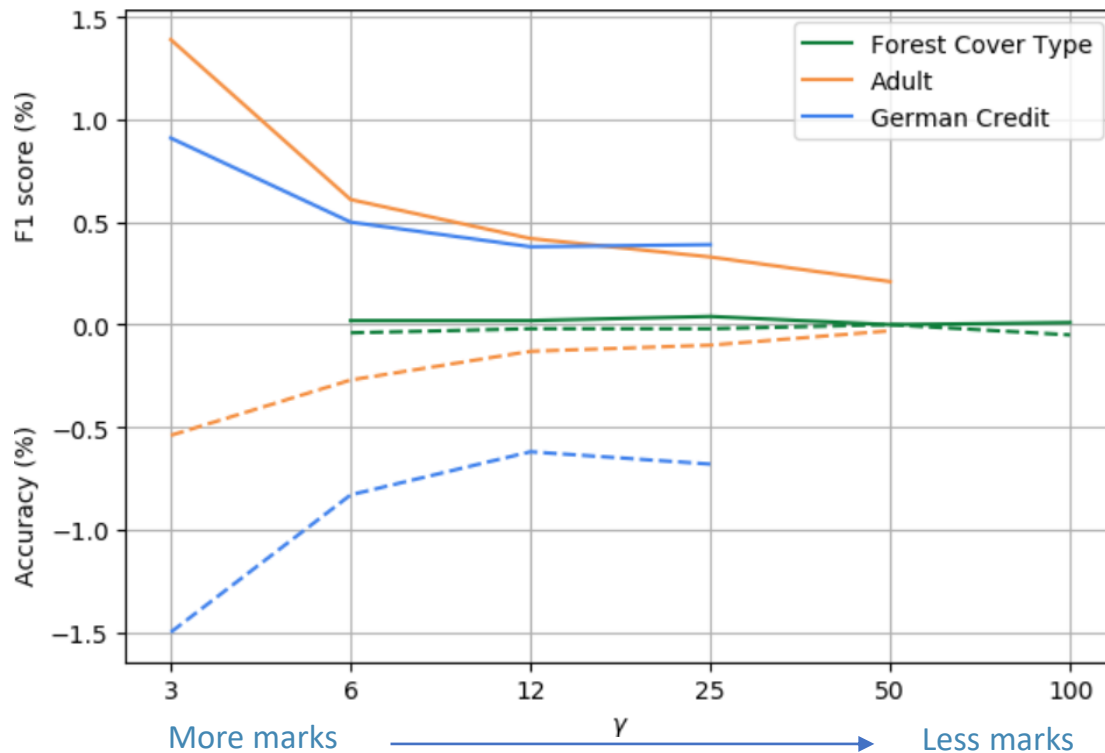
Less marks  More marks

Attribute	Variance	β ξ 30		25		15		10	
		4	8	4	8	4	8	4	8
Elevation	78391	0	+13	+1	+15	+1	+48	+1	+178
Aspect	12525	0	+7	0	+12	0	+35	0	+127
Slope	56	0	+12	0	+18	0	+48	0	0
HD-Hydrology	45177	0	+6	+1	+4	+1	+13	+2	0
VD-Hydrology	3398	0	+10	0	+15	0	+38	0	+87
HD-Roadways	2431276	0	+3	0	+3	0	+44	-2	0
Hillshade-9am	717	0	+11	0	+15	0	+41	0	+8
Hillshade-noon	391	0	+11	0	+16	0	+45	0	+200
Hillshade-3pm	1465	0	0	0	+13	0	+35	0	+160
HD-Fire-Points	1753493	0	0	0	-4	0	+54	0	+68

Number of
Least
Significant
Bits being
marked

Utility Evaluation

ML Classification



Conclusion and Future Work

Number of marks

\uparrow	ω	ξ	L
Misdiagnosis false hit	\uparrow		\uparrow
Subset Attack	\uparrow		\downarrow
Bit-flipping Attack	\uparrow	\uparrow	\downarrow
Additive Attack	\uparrow	\downarrow	
Utility	\downarrow	\downarrow	

- **Trade-off:** robustness of the scheme and the utility of fingerprinted data

Future work

- Fingerprinting scheme for categorical data
 - Further analysis on robustness of the neighbourhood-search approach
 - Blind scheme for fingerprinting relational data with categorical values



Thank you for
your attention