

Study from Data Anonymization Competition of Online Retail Data

Hiroaki Kikuchi
Meiji University, Japan

Issues in Anonymization

- 1. No real **dataset**
 - Data owner won't publish confidential dataset. Inconsistent Quasi-identifiers
- 2. No standard **metrics** for quantifying risk
 - Complicated models. Risk depends on many factors, e.g. dataset, technical skill, availability of background data. Utility depends on use case (but which is unknown when collecting data)
- 3. No standard model of **adversary**
 - “mildly motivated adversary” vs. “highly motivated adversary”

Competition PWSCUP 2015, 2016

- Privacy Workshop
- Organized by IPSJ, CSEC SIG

	2015	2016
Venue	Nagasaki (Brick Hall)	Akita (Castel Hotel)
When	Oct. 21, 22	Oct. 11, 12
Participants	13 Teams (20 in total)	15 Teams (42 in Total)
Dataset	NSTAC synthesized data	UCI Dataset, Online Retail

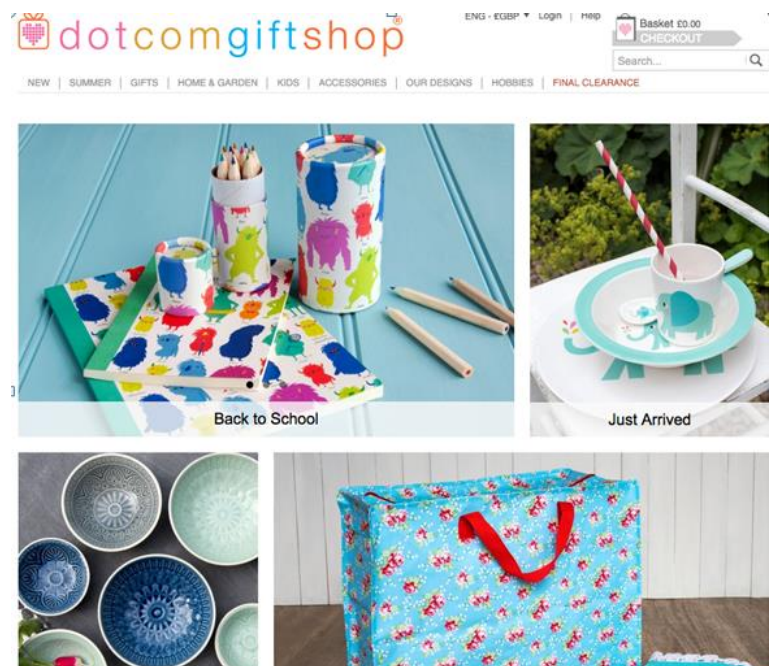


Our Approach

- 1. Common Dataset
 - We have used “pseudo microdata” synthesized by governmental agency, NSTAC, in 2015, and UCI Online Retail in 2016.
- 2. Quantifying risk
 - We focus on “records re-identification” risk and defines baseline **utility functions** and **some re-identification algorithms**. With arbitrary techniques, the best anonymization dataset is determined.
- 3. Adversary Model
 - We adopt Josef Domingo’s “*maximum-knowledge attacker*” model.

Dataset 'Online Retail'

- Available from UCI Machine Learning Repository
<https://archive.ics.uci.edu/ml/datasets/Online+Retail>
- Real payment transaction of UK Online Shop
 - One year transactions from 2010 Dec.
 - Gift shop
 - 540,000 records



Dataset 'Online Retail'

■ Master M

- $n = 400$ customers
- From 36 countries

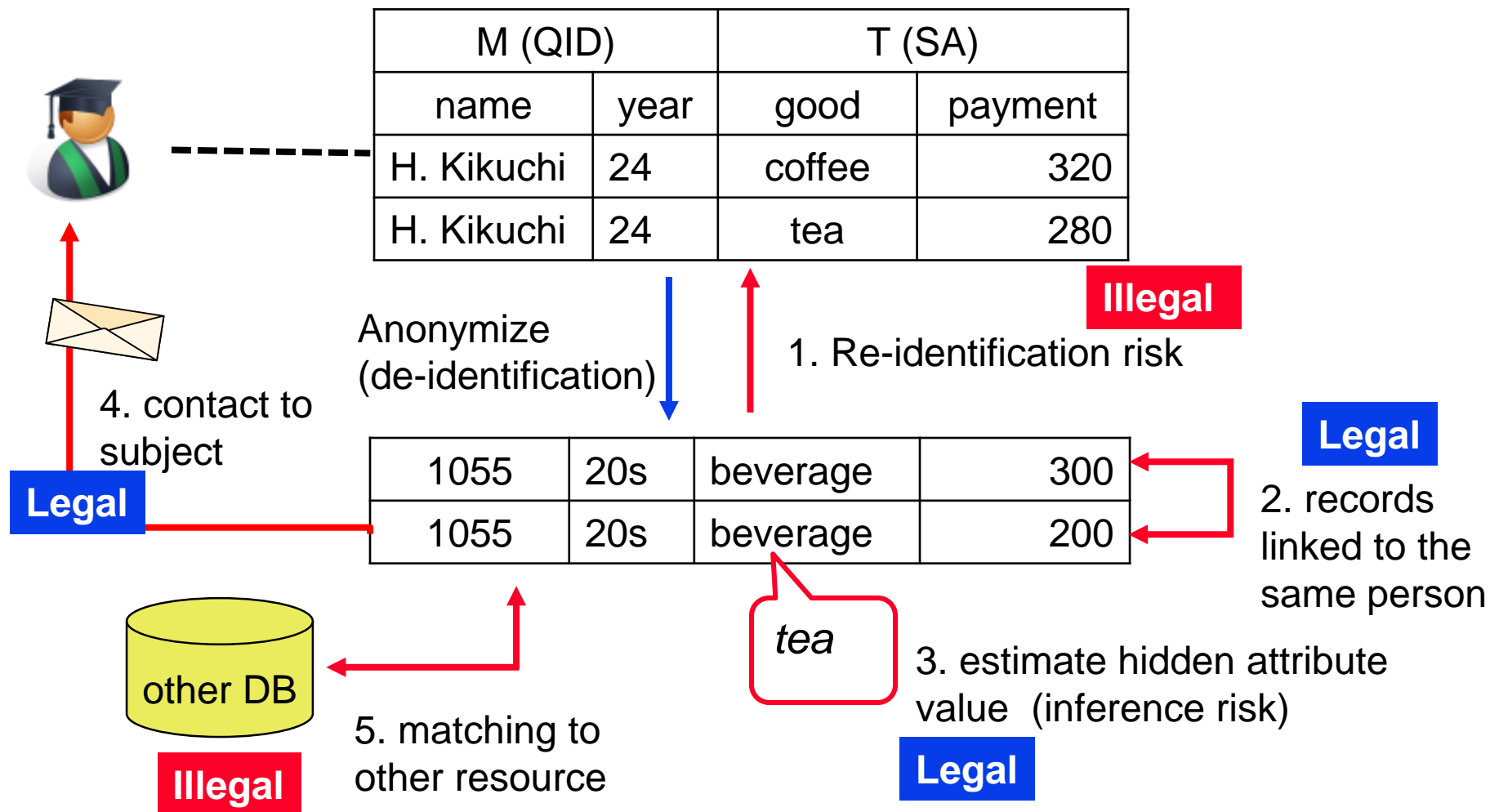
Customer ID	Sex	Birthday	Nationality
Online retail	synthesized		Online retail
12360	M	1876/2/24	Australia
12361	F	1954/2/14	Belgium
12362	F	1963/12/2	Belgium
12364	F	1960/9/16	Belgium

■ Transaction T

- $m = 38,087$ records
- 2,781 goods (stock code)

Customer ID	Invoice ID	Date	Time	Stock Code	Unit Price	Quantity
12362	544203	2011/2/17	10:30	21913	3.75	4
12362	544203	2011/2/17	10:30	22431	1.95	6
12361	545017	2011/2/25	13:51	22630	1.95	12
12361	545017	2011/2/25	13:51	22326	2.95	6

Privacy Risks (in Japan)



The Game

Master M

C. ID	Sex	Birthday	Country
12346	f	1960/12/25	UK
12347	f	1957/5/15	Iceland
12348	m	1947/2/19	Finland

Transaction T

C. ID	Date	Stock
12347	2010/12/7	85116
12347	2010/12/7	22375
12346	2011/1/18	23166

Anonymization (pseudonym, perturbation, suppression)

Anonymized M'

C. ID	Sex	Birthday	Country
10	<i>m</i>	1947/01/01	UK
20	<i>f</i>	1960/01/01	UK
30	<i>f</i>	1960/01/01	UK

Anonymized T'

C. ID	Date	Stock
10	2010/12/1	85123A
30	2010/12/1	85123A
30	2010/12/7	20000
20	2011/1/18	20000

Record index

Q	P
3	3
2	1
2	2

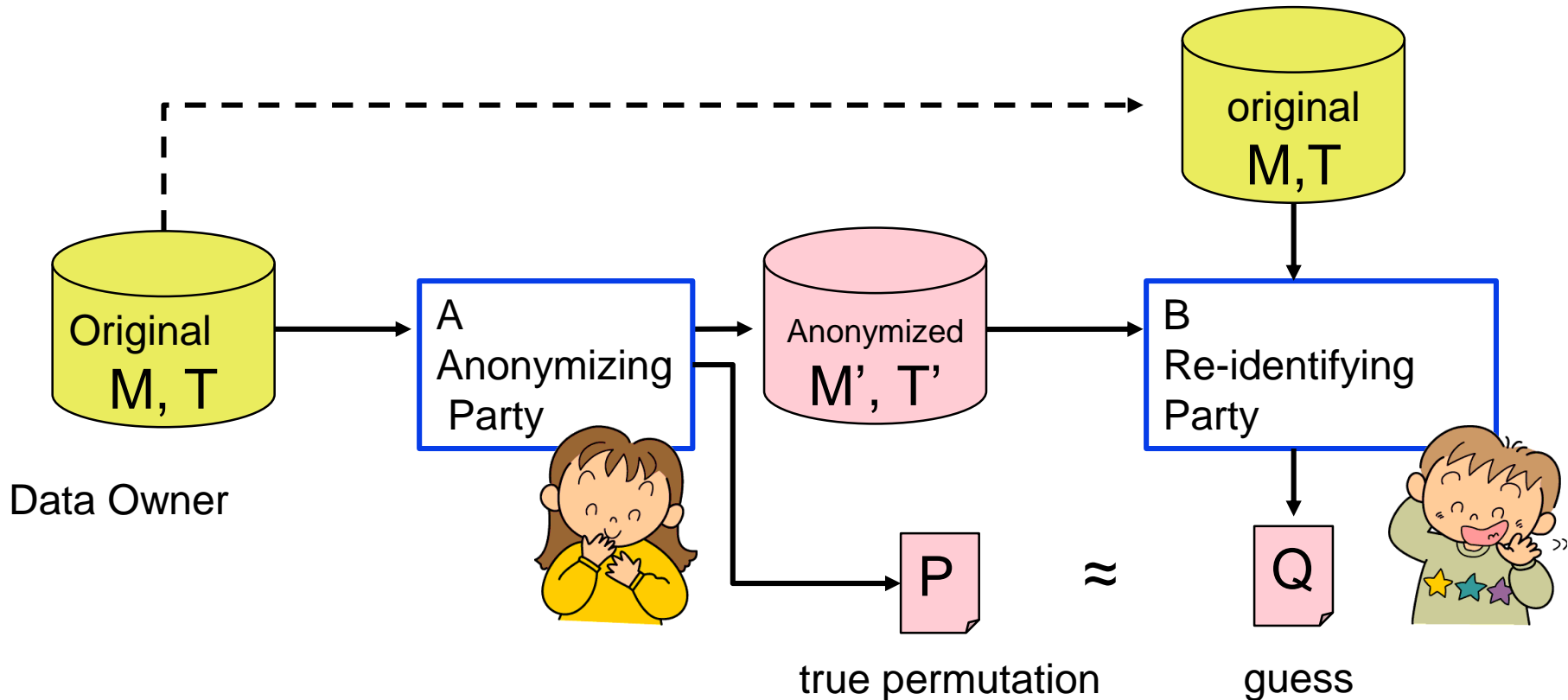


Estimated index Q

$$\text{Re-identification rate } \text{Re-ID}(P, Q) = \frac{\# \text{ Correct records}}{n'} = 2/3$$

Adversary Model

- Maximum Knowledge Adversary Model



Use cases and Utility

- 1. RFM Analysis

- Classification of customers based on **R**ecency (last purchase), **F**requency (of purchase), **M**onetary (Amount of payment)

U3: ut-rfm

- 2. Association Rule mining

- Association rule of stock code

U4: ut-top_item

- 3. Cross tabulation

- Accumulation of payment for several categories, sex, age, countries.

U1: ut-cmae

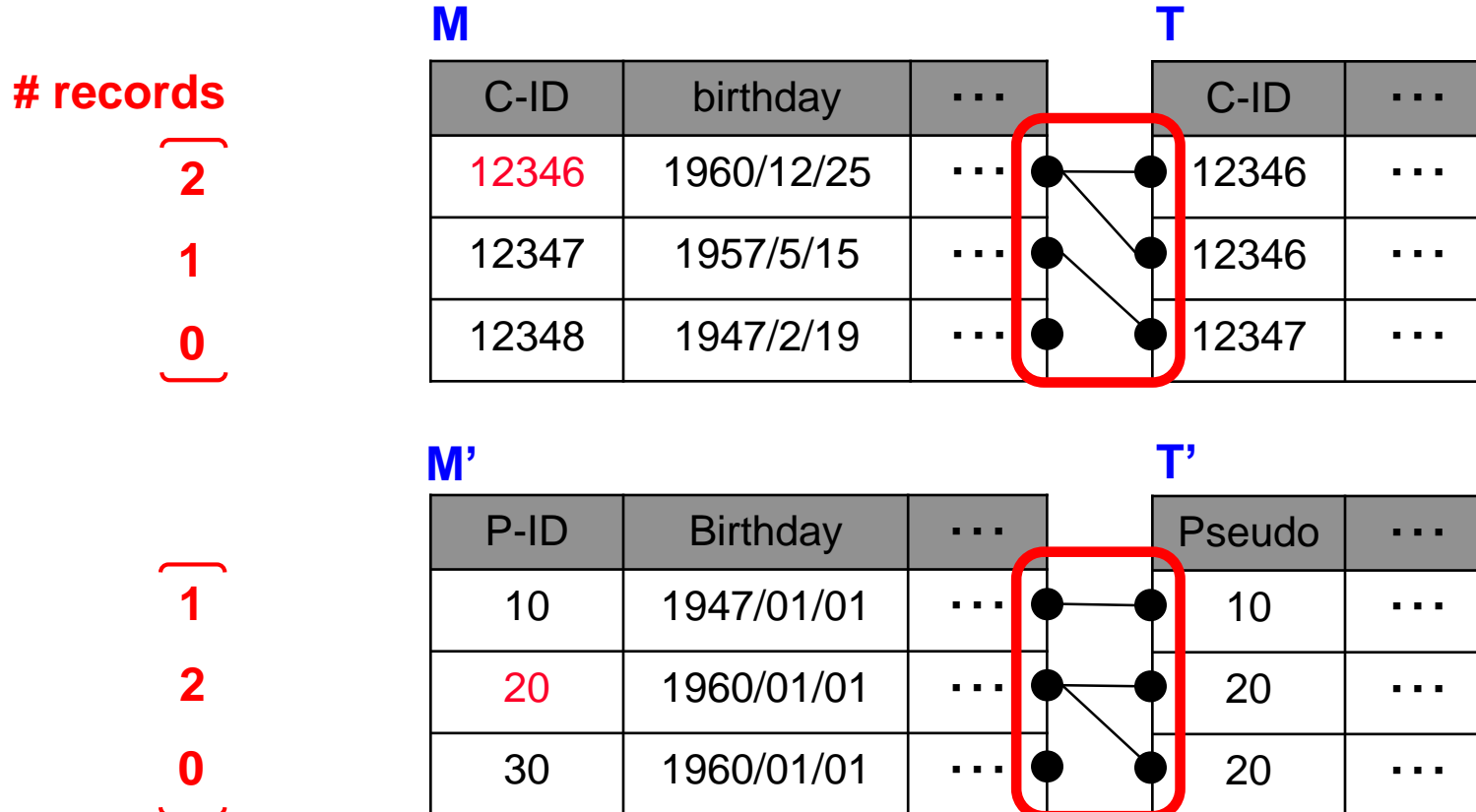
U2: ut-cmae2

Sample Re-identification

No	Algorithm	Description	M	T
E1	Re-birthday.py	Find the shortest birthday	✓	
E2	Re-eqi.rb	Find exact match	✓	✓
E3	Re-sort.rb	Sort and match	✓	
E4	Re-sort.rb	Sort by M and match	✓	
E5	Re-recnum.py	Find the shortest # recipients		✓
E6	Re-eqtr.rb	Find the same T		✓
E7	Re-tnum.rb	Sort by # records		✓
E8	Re-voting.py	Voting by birth, mean time, payment		✓
E9	Re-meantime.py	Find the shortest mean time		✓
E10	Re-ret.jar	Find similar set of goods		✓
E11	Re-sort2.tb	Sort by time and match		✓
E12	Re-search.rb	Find the shortest total payments		✓
E13	Re-totprice.py	Find the nearest set of goods		✓

E7 re-tnum-bi (best re-id score)

- ❑ Step 1: count # records in T for each customer
- ❑ Step 2: sort C-ID and P-ID by # records and birthday
- ❑ Step 3: match two sorted sequence and output Q



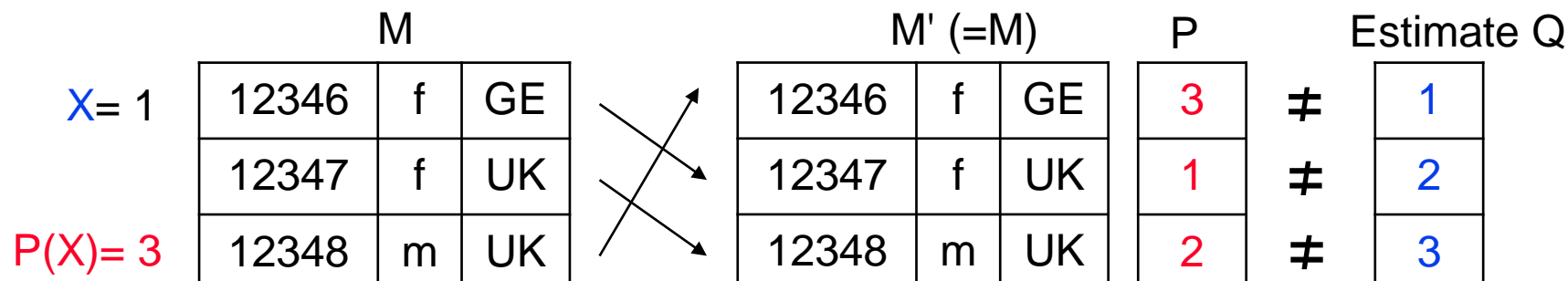
Competition rule

■ Rule Ver. 1.3

- (1) Each team submits one anonymized data.
- (2) Reject **cheating** anonymization
- (3) Each team is allowed to re-identify the anonymized data submitted by others in hour.
- (4) Winner is determined by grade defined by $U + E$, the sum of minimum utilities and the minimum security (max re-identification rate).
- (5) Best Re-identification is award to team who succeeds to re-identetify the winner's data.

The “Cheating”

■ Cheating anonymization



■ Cheating detection

□ Y1 (subset) > 50,000

□ Y2 (Jaccard) > 0.7

Y1:

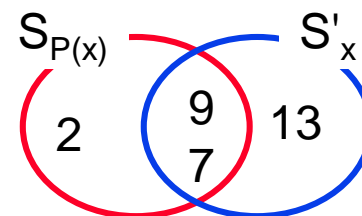
$\mu_{P(X)}$ = Total monthly payment of $P(X)$ = 305

μ_X = Total monthly payment of X = 405

Y2:

S'_x = set of goods paid by X

$S_{P(X)}$ = set of goods paid by $P(X)$



Two Phases of games

■ Phase 1

□ Online (web based)

7/27-8/16 (2.5 weeks)	Team entry
8/25-9/20 (3 weeks)	To submit anonymized data (update any times)
9/26-10/3 (1 week)	To submit estimated permutation (10 times per data)

□ Weight 1

■ Phase 2

□ Onsite

10/11 (10 min)	To submit one anonymized data (one time)
(1 hour)	To submit estimated permutation (10 times per data)

□ Weight 9

Browser address bar: <https://pwscup.personal-data.biz/login/WorldCup/index>

Google Translate: Translated to: English | Show original

By NIFTY Powered. Test team Logout

PWS CUP 匿名加工・再識別コンテスト

Contest ▾ Rankings ▾ Ranking (2016) ▾ Configuration ▾ manual ▾

0.006 sec

Contest
-> Anonymize

An
re-i

In the "anonymous and re-identification contest", to the anonymous data that participants submitted, other participants will attempt to re-identify. By re-identified by the researchers to each other, to verify the safety of anonymous data.



Anonymization and re-identification contest

In the "anonymous and re-identification contest", to the anonymous data that participants submitted, other participants will attempt to re-identify. By re-identified by the researchers to each other, to verify the safety of anonymous data.

Drag & Drop

M_***.CSV

T_***.CSV

P_***.CSV

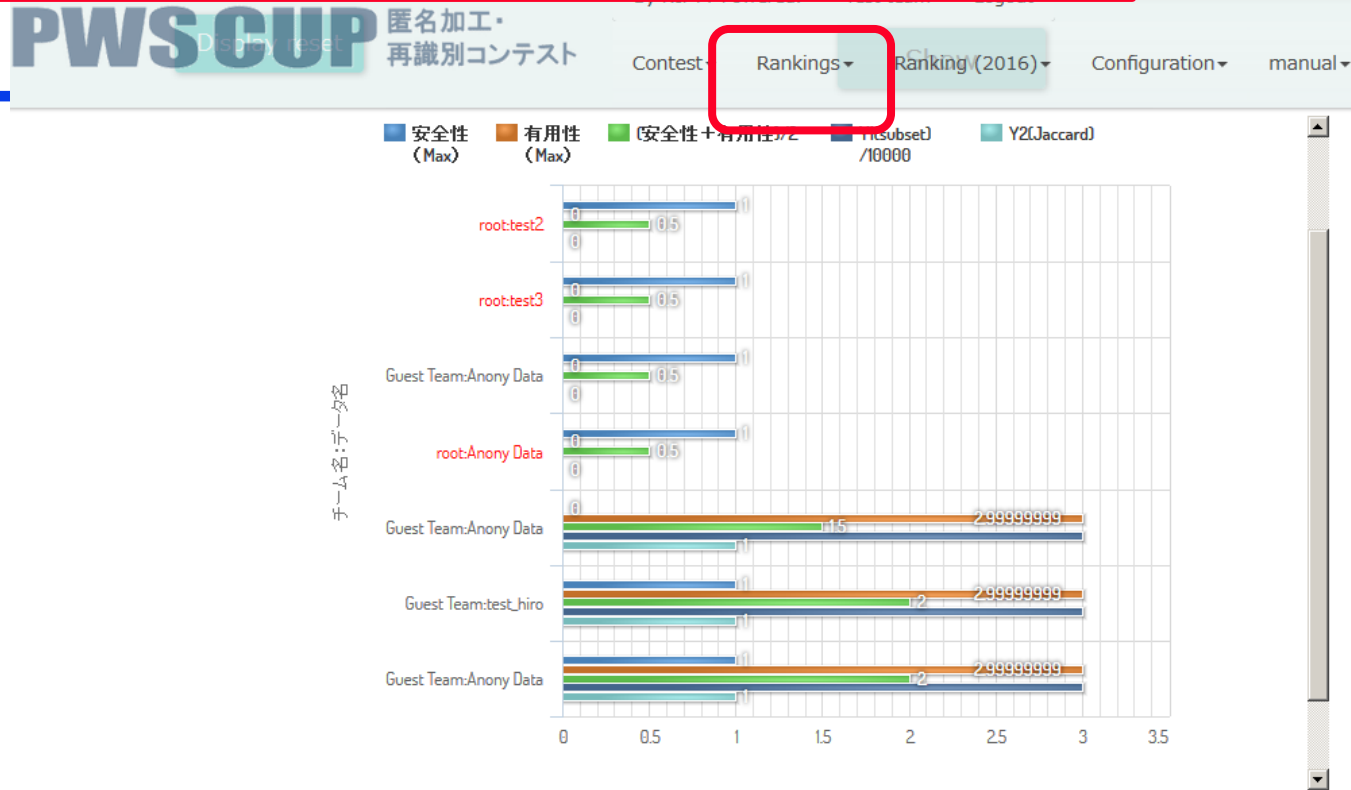
The screenshot shows a web browser window with a search bar and navigation menu. Below the menu, there is a green box containing the text "Please enter the data registered name" and a text input field with the placeholder "under 10 words". To the right of the input field, the text "0.001 sec" is displayed. Below the input field, there is a large green dashed box labeled "Upload area" with the text "(M, T, P) to drag and drop into this area. select the file by clicking." and a folder icon with an upward arrow. At the bottom of the page, there is a "hand in" button.

The screenshot shows a file explorer window with a search bar and a list of files. The files are M_400.csv, P_400.csv, and T_400.csv. A red arrow points from the files to the "Upload area" in the main screenshot.

Automated Risk Evaluation

Use Program	Re-identification program description	rate]	Result detail
E1-birthday	Re-identify the customer ID distance of the date of birth to each other becomes the minimum ※ preliminary round of the defect corrected	396/400 [0.99000000]	[inspection result]
E2-eqi	I guess the record the transaction is fully consistent with the master of the attributes (except for the temporary ID). If not random	400/400 [1.00000000]	[inspection result]
E3-sort	(Sex, date of birth, country) in the sort	400/400 [1.00000000]	[inspection result]
E4-sort2	Sorted by date of birth	400/400 [1.00000000]	[inspection result]
E5-recnum	Record number of matching (re-identification to the nearest customer distance of the record number of each other transaction)	168/400 [0.42000000]	[inspection result]
E6-eqtr	I guess the record that the transaction is complete match. If not random	400/400 [1.00000000]	[inspection result]
E7-tnum	Sorted by the number of transactions	400/400 [1.00000000]	[inspection result]

the Ranking Available



Sort reset

↓ clicked a column to sort in ascending or descending order.

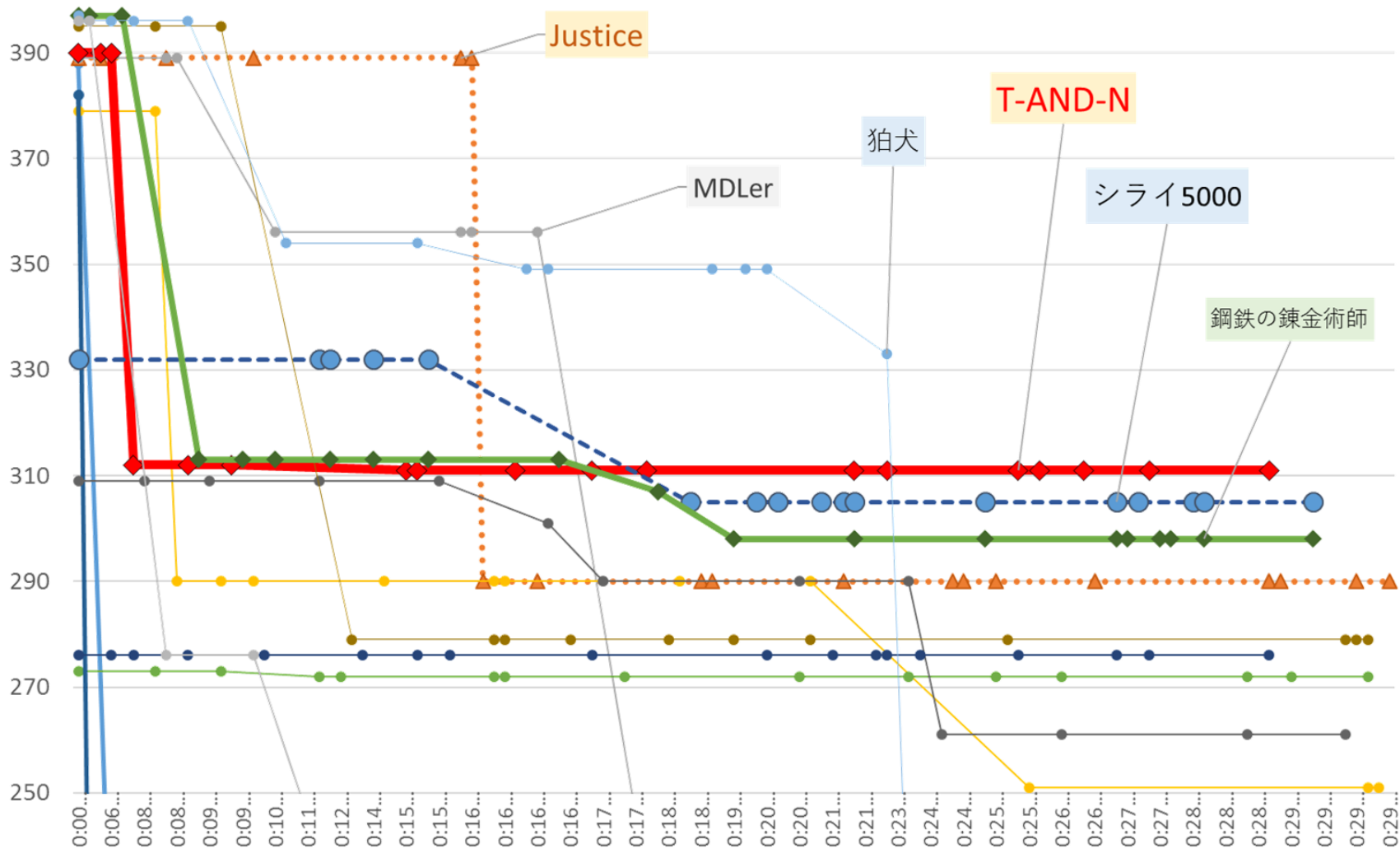
useorg	Safety (Max)	Usefulness (Max)	(Safety + usefulness) / 2	Y1 (Subset) / 10000	Y2 (Jaccard)
検索語...	検索語...	検索語...	検索語...	検索語...	検索語...
root: test3	1.00000000	0.00000000	0.50000000	0.00000000	0.00000000
root: test2	1.00000000	0.00000000	0.50000000	0.00000000	0.00000000

Top of page
Back

Oct. 11, Pwscup Final

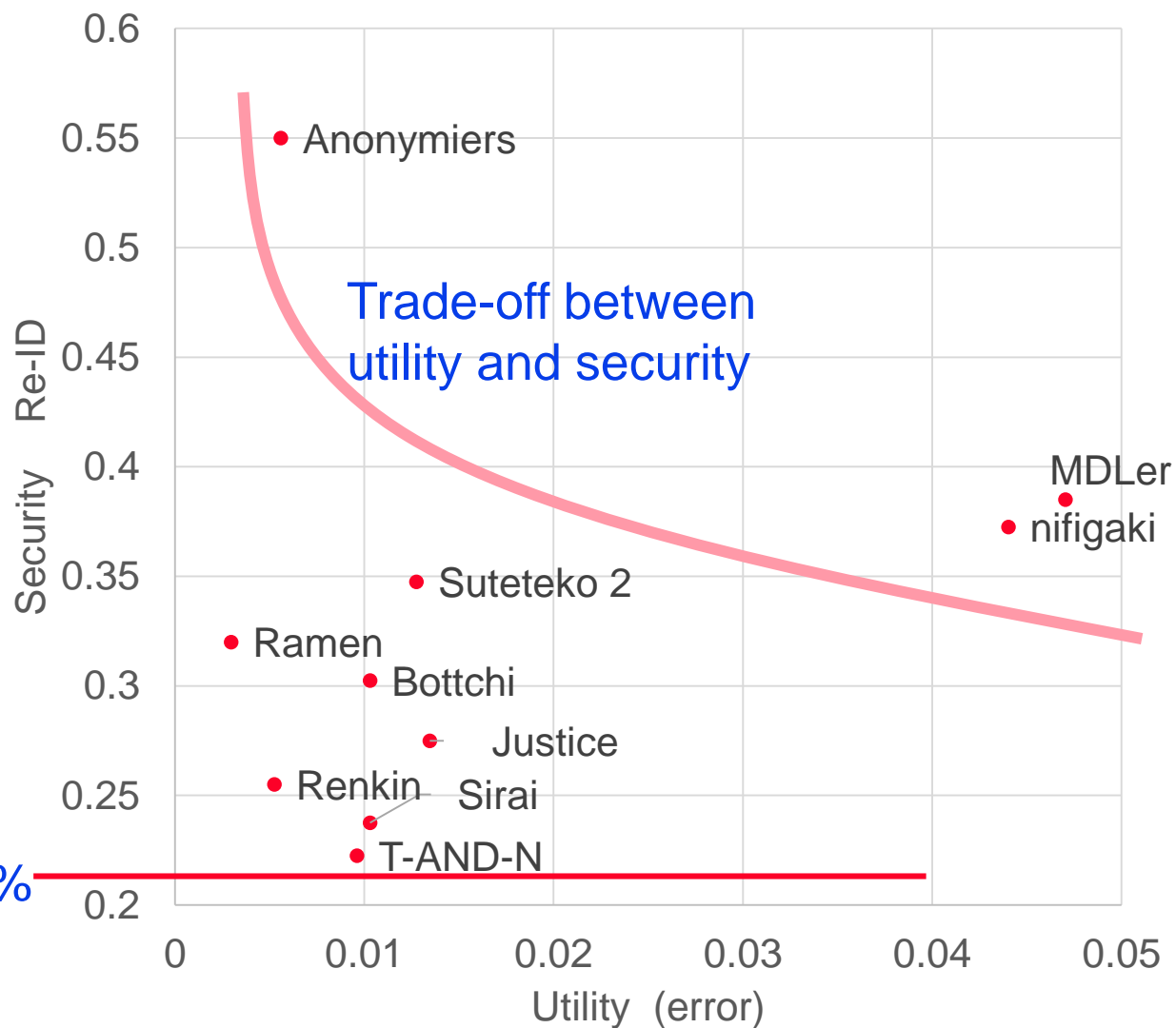


(Onsite) Rank Transition

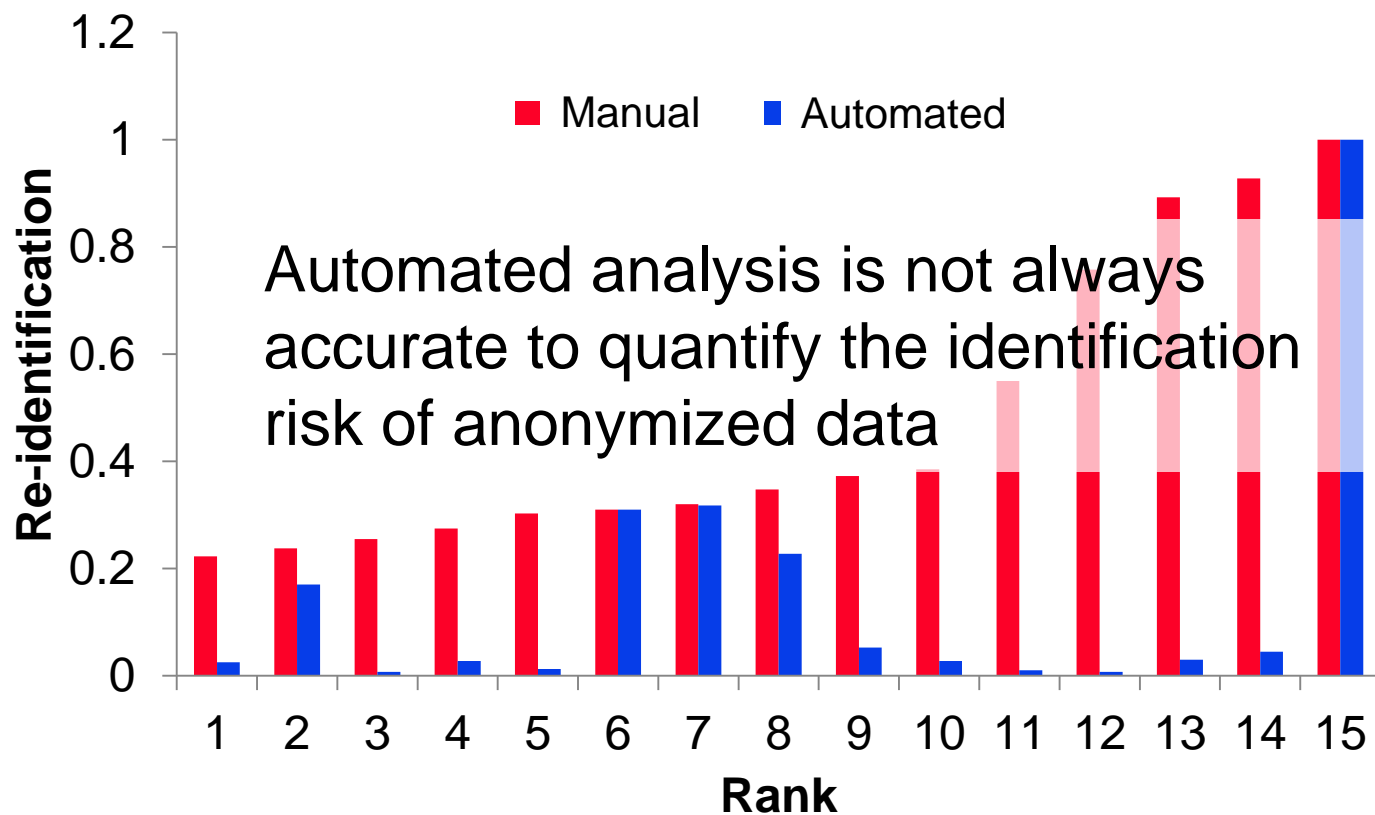


Competition Result (Top 10 teams)

No	team
1	T-AND-N
2	Shirai 5000
3	Renkin
4	Justice
5	Bottchi
6	Ramen
7	Suteteko 2
8	nifigaki
9	MDLer
10	Anonymiers



Automated and Manual re-id.



Conclusions

- Data anonymization competition 2016 with real online retail data was done successfully.
- Average re-identification is 188 (47%) out of 400 customers. The best (minimum) re-identification ratio is 22%.
- Mean Automated re-identification was 18%, manual re-identification was 47%.
 - Kikuchi, et.al, “A Study from the Data Anonymization Competition Pwscup 2015”, DPM 2016, LNCS 9963.
 - Kikuchi, et. Al, “Ice and Fire: Quantifying the Risk of Re-identification and Utility in Data Anonymization”, IEEE AINA 2016.