



**Information
Systems
Group**

Hasso Plattner Institut | Universität Potsdam

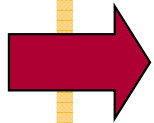
Duplicate Detection

6.6.2013

Felix Naumann

Overview

2



- Duplicate detection
- Similarity measures
- Algorithms
- Data sets and evaluation
- Data fusion



Duplicate Detection

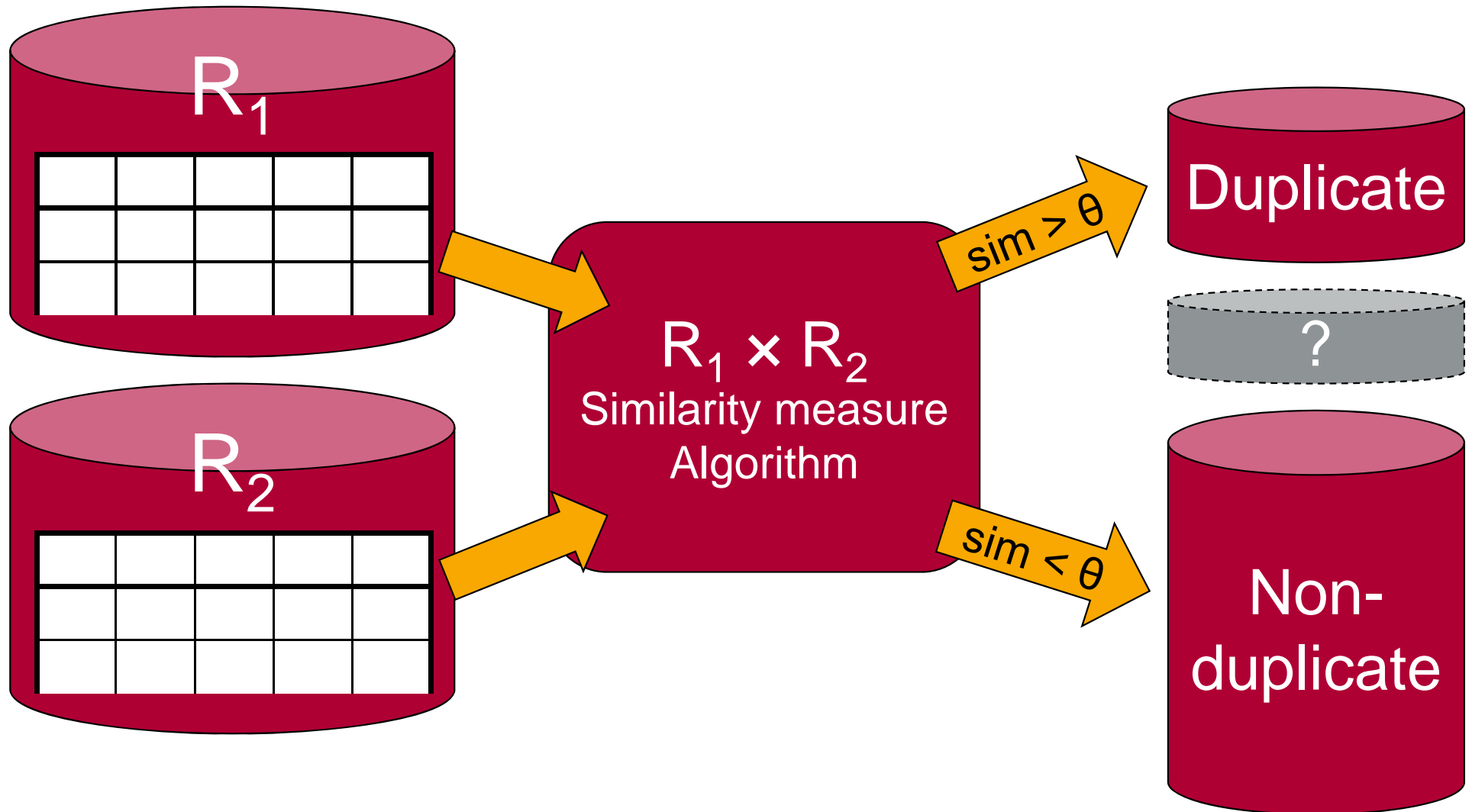
3

Duplicate detection is the discovery of multiple representations of the same real-world object.

- Problem 1: Representations are not identical.
 - *Fuzzy duplicates*
- Solution: Similarity measures
 - Value- and record-comparisons
 - Domain-dependent or domain-independent
- Problem 2: Data sets are large.
 - Quadratic complexity: Comparison of every pair of records.
- Solution: Algorithms
 - E.g., avoid comparisons by partitioning.

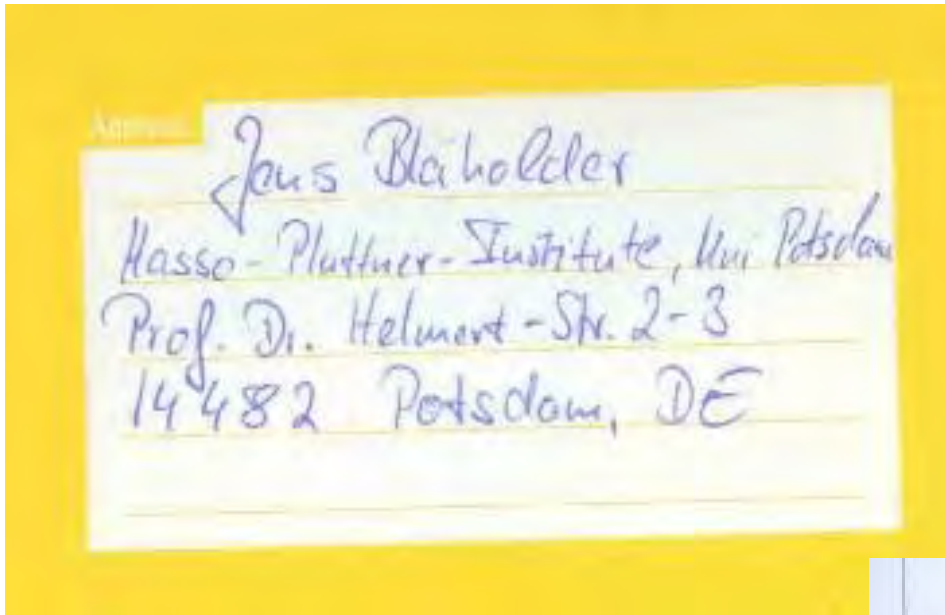
Duplicate Detection

4

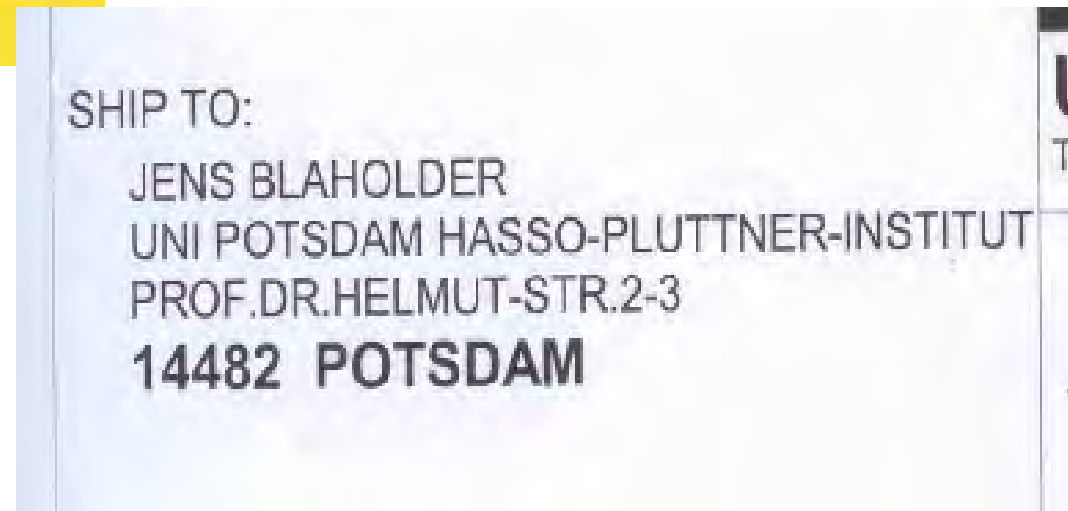


Origins of duplicates

5



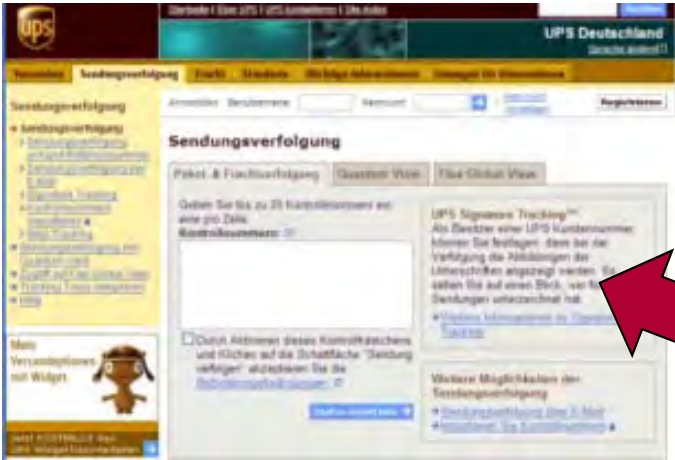
Original



Scanned

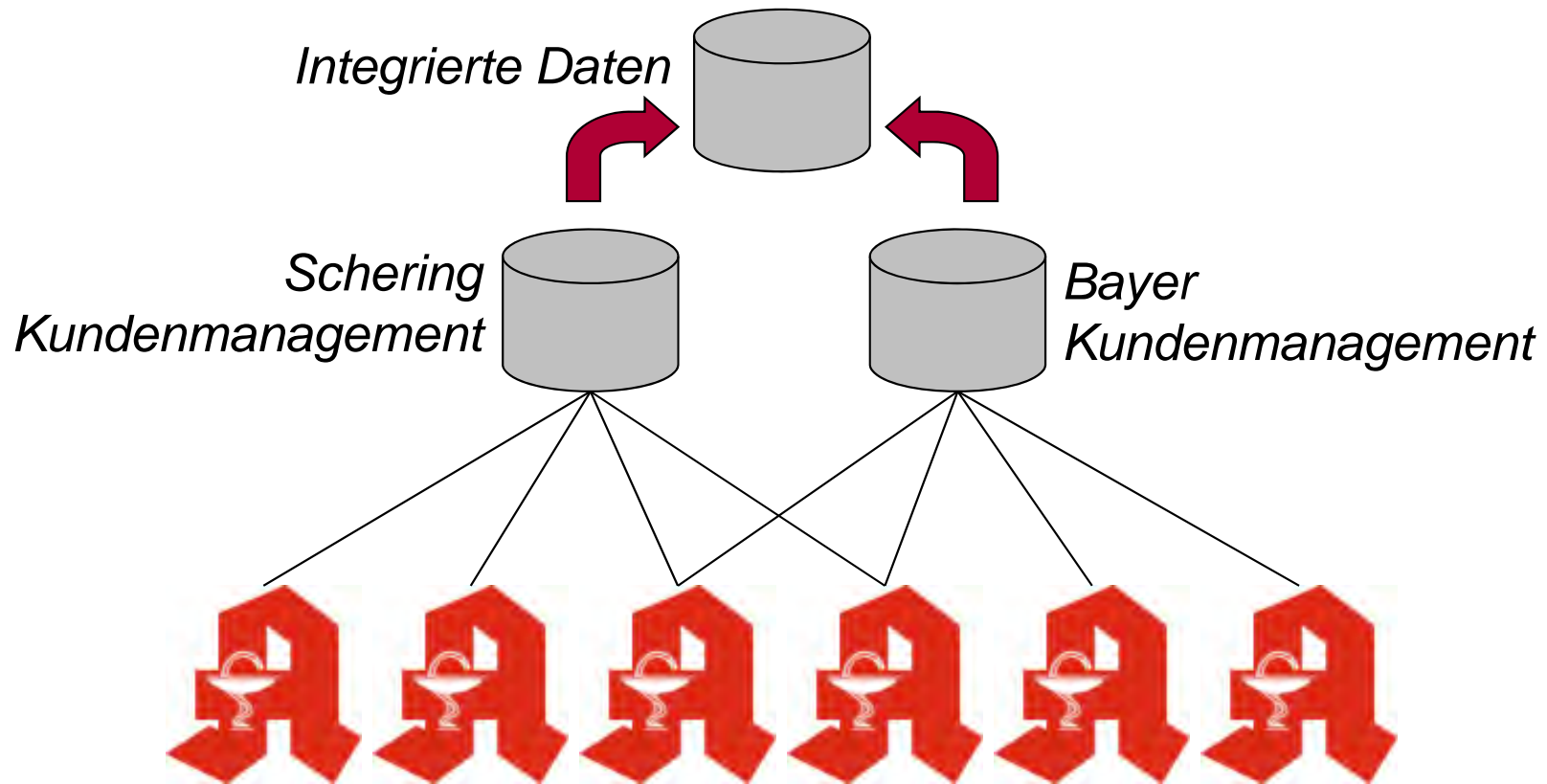
Origins of duplicates

6



Origins of duplicates

7



Difficult names

8

488941 britney spears	29 britent spears	9 brinttany spears	5 brney spears	3 britiy sp
40134 brittany spears	29 brittnany spears	9 britanay spears	5 broitney spears	3 britmeny
36315 brittney spears	29 britttany spears	9 britinany spears	5 brotny spears	3 britneeyy
24342 britany spears	29 btiney spears	9 britn spears	5 bruteny spears	3 britnehy
7331 britny spears	26 birttney spears	9 britnew spears	5 btiyney spears	3 britnely
6633 briteny spears	26 breitney spears	9 britneyn spears	5 btrittney spears	3 britnesy
2696 britteny spears	26 brinity spears	9 britrney spears	5 gritney spears	3 britnetty
1807 briney spears	26 britenay spears	9 brtiny spears	5 spritney spears	3 britnex e
1635 brittny spears	26 britneyt spears	9 brtittney spears	4 bittny spears	3 britneyxx
1479 brintey spears	26 brittan spears	9 brtny spears	4 bnritney spears	3 britnity
1479 britanny spears	26 brittne spears	9 brytny spears	4 brandy spears	3 britntey
1338 britiny spears	26 btittany spears	9 rbitney spears	4 brbritney spears	3 britnyey
1211 britnet spears	24 beitney spears	8 birtiny spears	4 breating spears	3 britterny
1096 britiney spears	24 birteny spears	8 bithney spears	4 breetney spears	3 brittneey
991 britaney spears	24 brightney spears	8 brattany spears	4 bretiney spears	3 britttney
991 britnay spears	24 brintiny spears	8 breitny spears	4 brfitney spears	3 brittynyey
811 brithney spears	24 britanty spears	8 breteny spears	4 briattany spears	3 brityen s
811 brtiney spears	24 britenny spears	8 brightny spears	4 brieteny spears	3 briytney
664 birtney spears	24 britini spears	8 brintay spears	4 briety spears	3 brltney s
664 brintney spears	24 britnwy spears	8 brinttey spears	4 briitny spears	3 broteny s
664 briteney spears	24 brittni spears	8 briotney spears	4 briittany spears	3 brtaney s
601 bitney spears	24 brittnie spears	8 britanys spears	4 brinie spears	3 brtiiany
601 brinty spears	21 biritney spears	8 britley spears	4 brinteney spears	3 brtinay s
544 brittaney spears	21 birtany spears	8 britneyb spears	4 brintne spears	3 brtinney
544 brittnay spears	21 biteny spears	8 britnrey spears	4 britaby spears	3 brtitany
364 britey spears	21 bratney spears	8 britnty spears	4 britaey spears	3 brtiteny
364 brittyny spears	21 britani spears	8 brittner spears	4 britainey spears	3 brtnet sp
329 brtney spears	21 britanie spears	8 brottany spears	4 britinie spears	3 brytiny s
269 bretney spears	21 briteany spears	7 baritney spears	4 britinney spears	3 btney spe
269 britneys spears	21 brittay spears	7 birntey spears	4 britmney spears	3 drittney
244 britne spears	21 brittinay spears	7 biteney spears	4 britnear spears	3 pretney s
244 brytney spears	21 brtany spears	7 bitiny spears	4 britnel spears	3 rbritney
220 breatney spears	21 brtiany spears	7 breateny spears	4 britneuy spears	2 barittany
220 britiany spears	19 birney spears	7 brianty spears	4 britnewy spears	2 bbbritney
199 britnney spears	19 birtney spears	7 brintye spears	4 britnmey spears	2 bbitney s
163 britny spears	19 britnaey spears	7 britianny spears	4 brittaby spears	2 bbritny s
147 britny spears	19 britny spears	7 britny spears	4 britny spears	2 britny spe

Motivation

9

- Possible effects

- Example: Portfolio Management Offers
- Credit maximum not detected
- Too low inventory levels
- No quantity discount for multiple orders
- Total revenue of preferred customers unknown
- Multiple mailings of same catalog to same household

Customer	Revenue
BMW	20.000
BaMoWe	5.000.000
Bayerische Motorenwerke	300.000
...	...

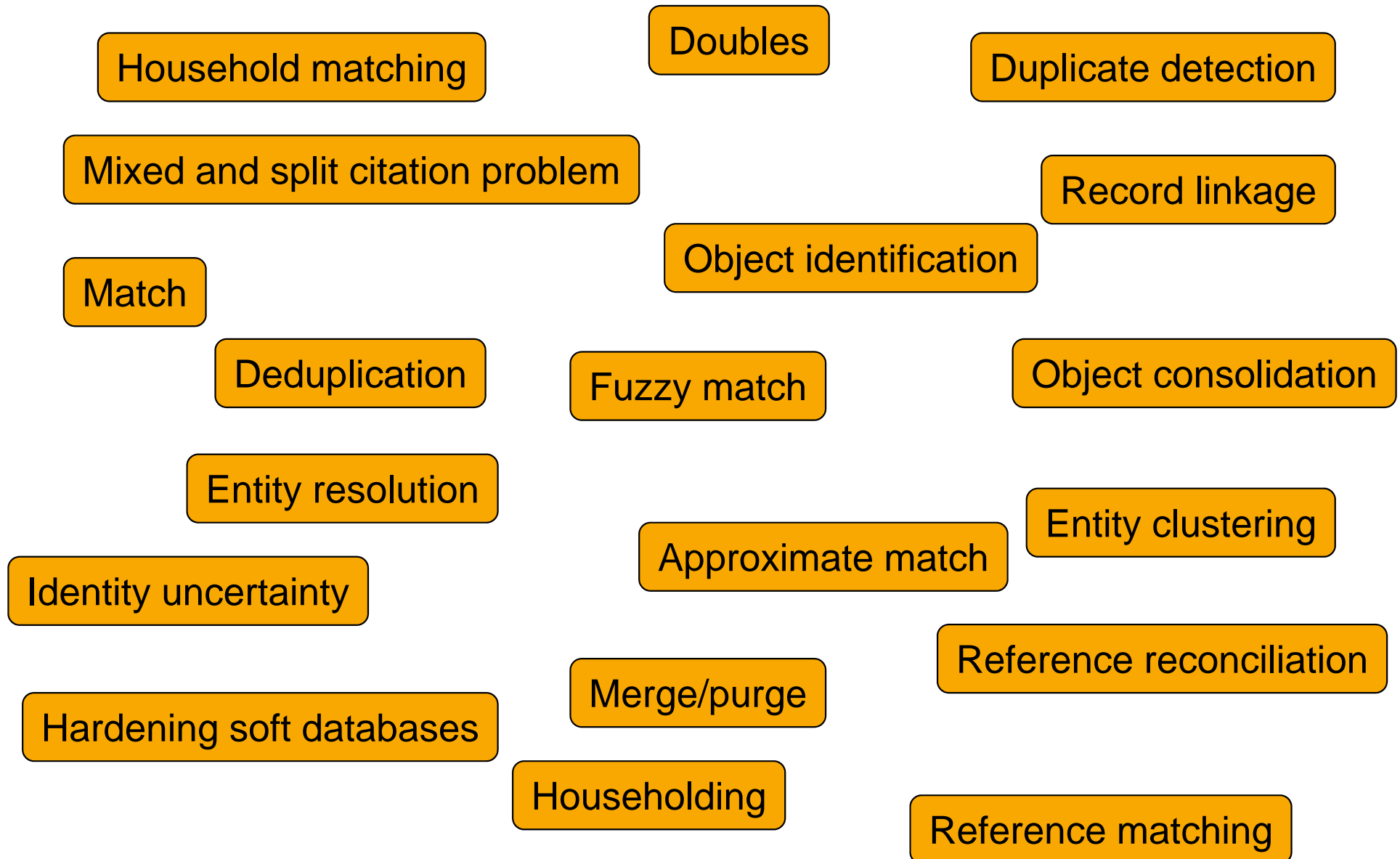
- General problems

- Additional, unnecessary IT expenses
- Low customer satisfaction
- Potentials and dangers not detected
- Poor quality financial data



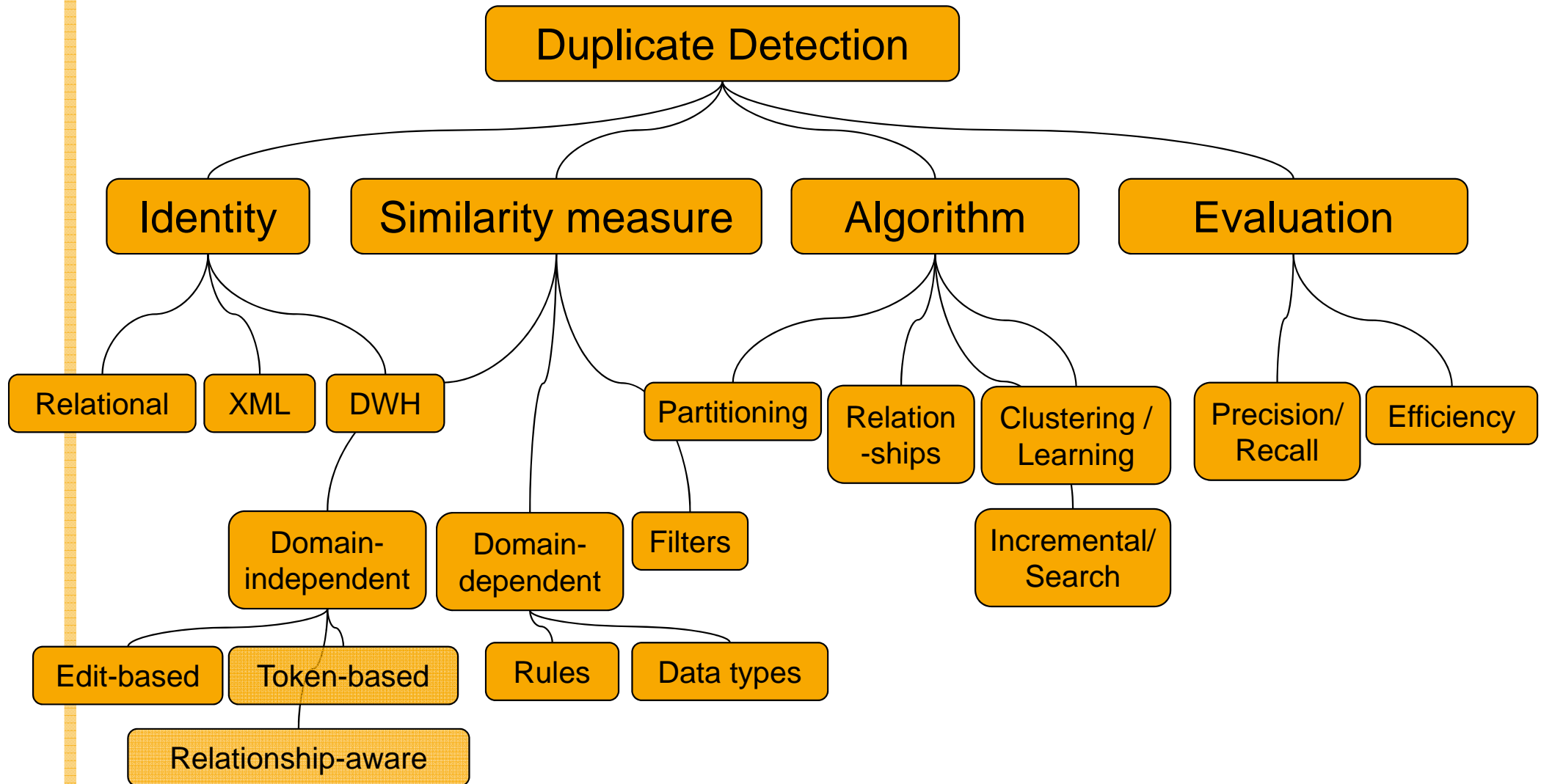
Ironically, "Duplicate Detection" has many Duplicates

10



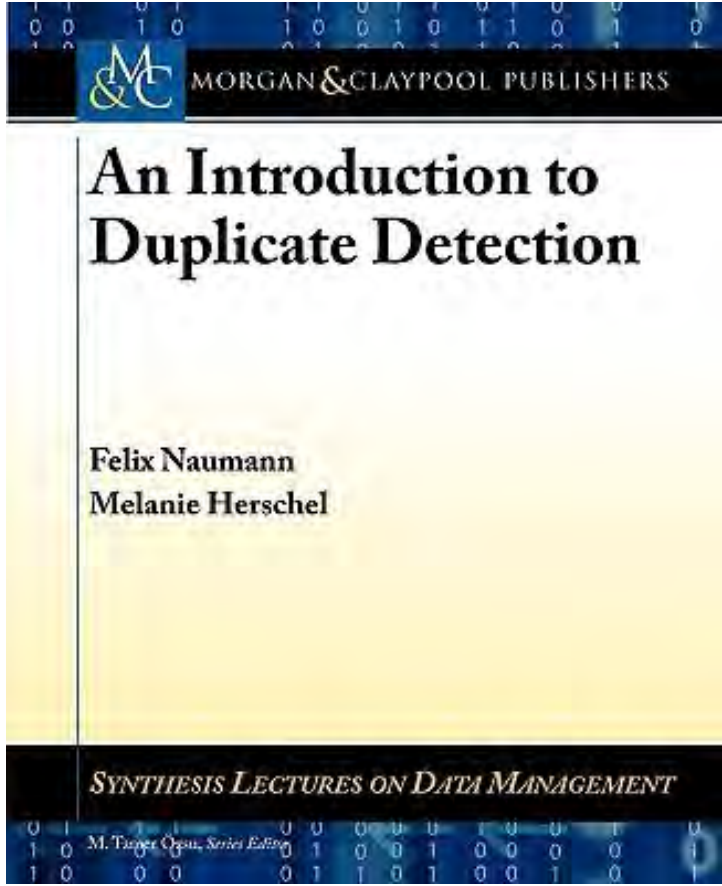
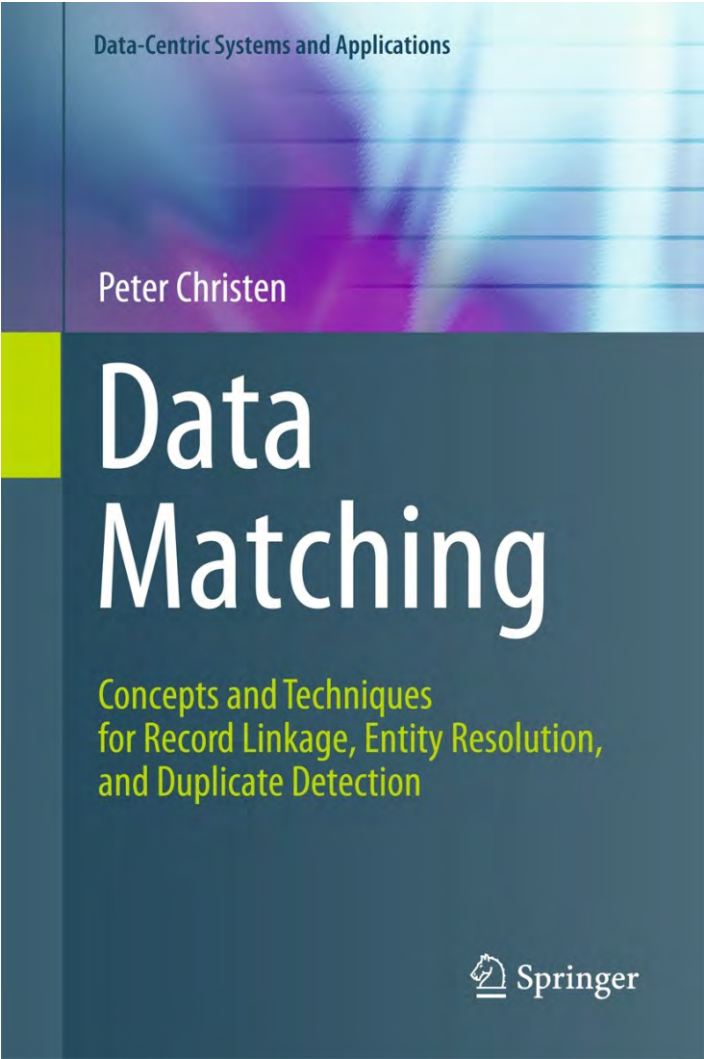
Duplicate Detection – Research

11



Literature

12



Overview

13

- Duplicate detection
- Similarity measures
- Algorithms
- Data sets and evaluation
- Data fusion



Token-based Similarity Measures


14

- Tokens
 - Words / Terms
 - n-grams
- Jaccard
 - $|\{\text{common tokens}\}| / |\{\text{all tokens}\}|$
- TFIDF [Cohen et al. 2003]
 - Term frequency: tf
 - Inverse document frequency: idf
 - TFIDF: $\log(\text{tf} + 1) \times \log(\text{idf})$
 - Common words have low weight
 - Cosine similarity of term vectors weighted by tfidf
- And many more
[Koudas Srivastava 2005]

Edit-based Similarity Measures

15

- Jaro [Jaro 1989] / Jaro-Winkler [Winkler 1999]
 - Common letters within $\frac{1}{2}$ string length
 - Transposed letters
- Edit-distance / Levenshtein-distance [Levenshtein 1965]
 - Minimum number of edits from one word to the other
 - Domain-specific costing
 - Dynamic Programming
- Soundex
 - 4-letter code for each word
 - `SOUNDEX('Farwick ')` = F620
- ...



Frass, Fricke,
Fahruschi,
Feuerhake

Domain-dependent Similarity Measures

16

■ Data Types

- Special similarity for dates
- Special similarity for numerical attributes
- ...

■ Rules

- [Hernandez Stolfo 1998], [Lee et al. 2000]
- **Given two records, r1 and r2.**
IF last name of r1 = last name of r2,
AND first names differ slightly,
AND address of r1 = address of r2
THEN r1 is equivalent to r2.

Relationship-aware Similarity Measures

17

- Idea: Not only values of the records, but values of related records are relevant for similarity.

- Persons: spouse, children, employer
- Movies: actors
- CDs: songs
- Customers: orders, addresses
- Dimensions in a DWH [Ananthakrishna et al. 2002]

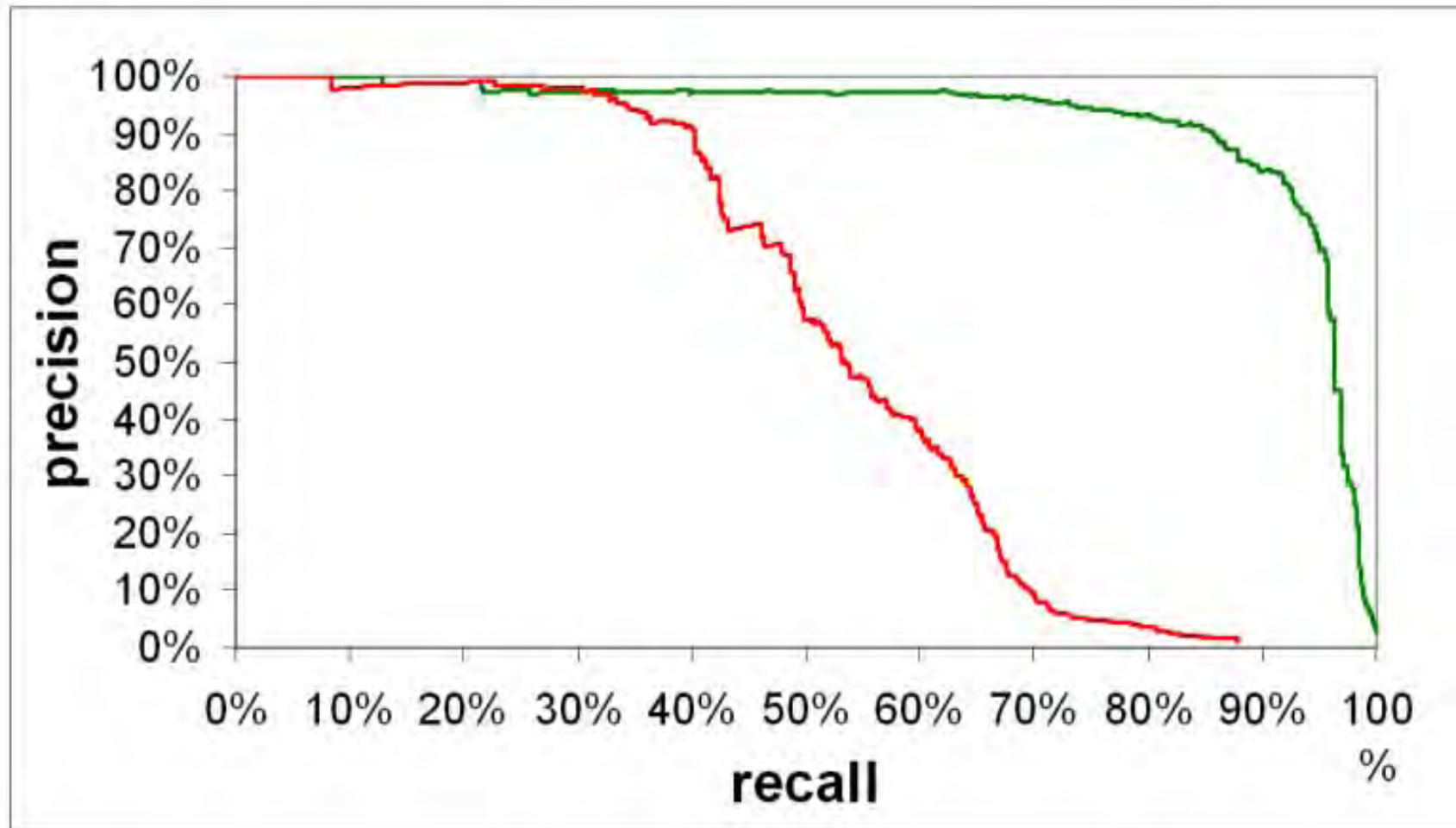
ID	Country
1	USA
2	United States
3	Unitd States

ID	City	Country
1	New York	1
2	Los Angeles	1
3	Now York	2
4	Los Angeles	2
5	New York	3
6	Los Angels	3

ID	Street
1	First Ave
2	High St.
3	Broadway
4	Embarca
5	Broadway
6	Second S
7	P St.
8	Pennsylv
9	Sunset B
10	Santa Mc
11	Ocean Av

Relationship-aware Similarity Measures – Evaluation

18



— without actors — with actors

Overview

19

- Duplicate detection
- Similarity measures
- Algorithms
- Data sets and evaluation
- Data fusion



Record Pairs as Matrix

20

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

Number of comparisons: All pairs

21

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

400
comparisons

Reflexivity of Similarity

22

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1	■																			
2		■																		
3			■																	
4				■																
5					■															
6						■														
7							■													
8								■												
9									■											
10										■										
11											■									
12												■								
13													■							
14														■						
15															■					
16																■				
17																	■			
18																		■		
19																			■	
20																				■

380
comparisons

Symmetry of Similarity

23

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	
1		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2			1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
3				1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
4					1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
5						1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
6							1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
7								1	1	1	1	1	1	1	1	1	1	1	1	1	1
8									1	1	1	1	1	1	1	1	1	1	1	1	1
9										1	1	1	1	1	1	1	1	1	1	1	1
10											1	1	1	1	1	1	1	1	1	1	1
11												1	1	1	1	1	1	1	1	1	1
12													1	1	1	1	1	1	1	1	1
13														1	1	1	1	1	1	1	1
14															1	1	1	1	1	1	1
15																1	1	1	1	1	1
16																	1	1	1	1	1
17																		1	1	1	1
18																			1	1	1
19																				1	1
20																					1

190 comparisons

Complexity

24

- Problem: Too many comparisons!
 - 10.000 customers
 - => 49.995.000 comparisons
 - ◇ $(n^2 - n) / 2$
 - ◇ Each comparison is already expensive.

- Idea: Avoid comparisons...
 - ... by filtering out individual records.
 - ... by partitioning the records and comparing only within a partition.

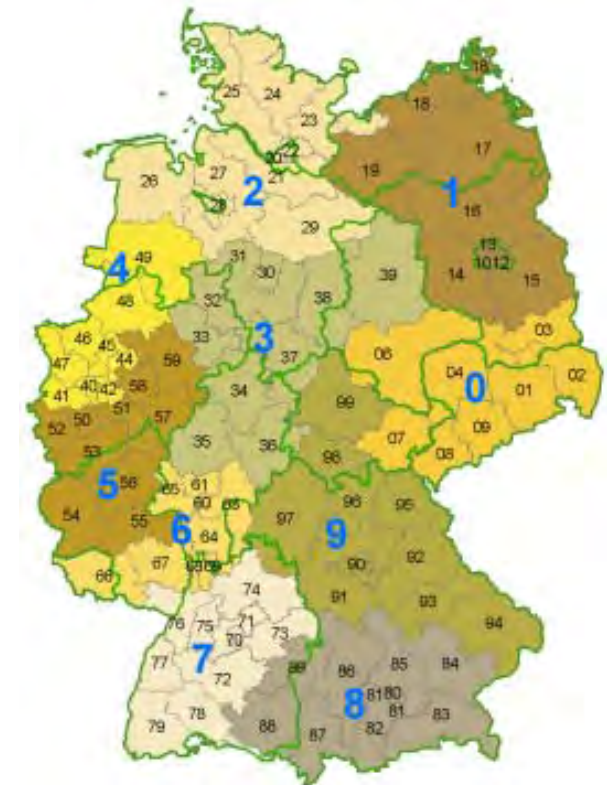


Partitioning / Blocking

25

- Partition the records (horizontally) and compare pairs of records only within a partition.
 - Partitioning by first two zip-digits
 - ◇ Ca. 100 partitions in Germany
 - ◇ Ca. 100 customers per partition
 - ◇ => 495.000 comparisons
 - Partition by first letter of surname
 - ...

- Idea: Partition multiple times by different criteria.
 - Then apply transitive closure on discovered duplicates.



Source: wikipedia.de

Records sorted by ZIP

26

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

190
comparisons

Blocking by ZIP

27

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1		■	■	■	■															
2			■	■	■															
3				■	■															
4					■															
5						■	■	■												
6							■	■												
7								■												
8									■	■	■									
9										■	■	■								
10											■	■								
11												■								
12													■	■	■	■				
13														■	■	■	■			
14															■	■	■			
15																■	■			
16																	■			
17																		■		
18																			■	■
19																				■
20																				

32 comparisons

Sorted Neighborhood

[Hernandez Stolfo 1998]

28

■ Idea

- Sort tuples so that similar tuples are close to each other.
- Only compare tuples within a small neighborhood (window).

1. Generate key

- E.g.: SSN+“first 3 letters of name” + ...

2. Sort by key

- Similar tuples end up close to each other.

3. Slide window over sorted tuples

- Compare all pairs of tuples within window.

■ Problems

- Choice of key
- Choice of window size

■ Complexity: At least 3 passes over data

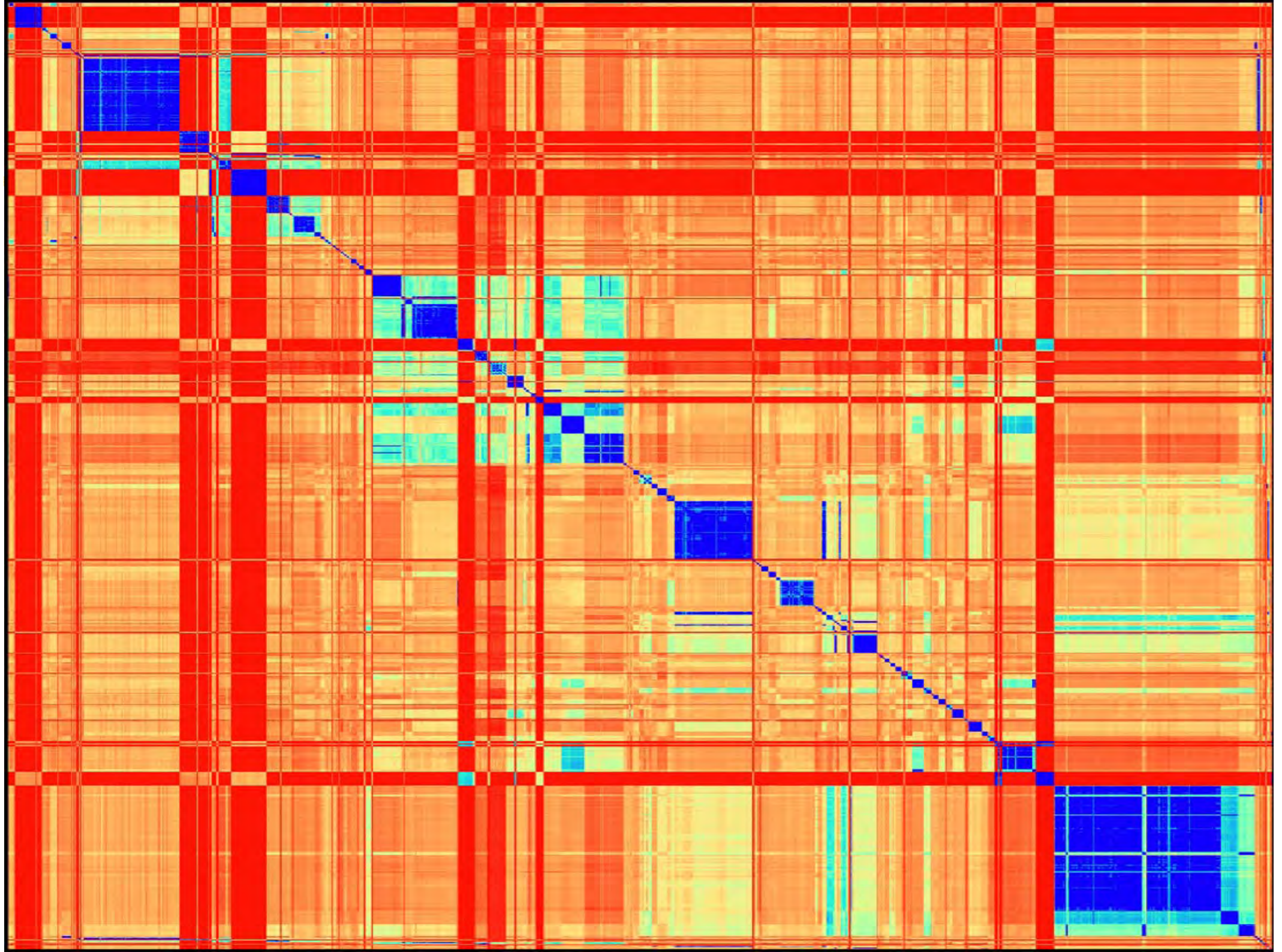
- Sorting!

SNM by ZIP (window size 4)

29

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1		■	■	■	■															
2			■	■	■															
3				■	■	■														
4					■	■	■													
5						■	■	■												
6							■	■	■											
7								■	■	■										
8									■	■	■									
9										■	■	■								
10											■	■	■							
11												■	■	■						
12													■	■	■					
13														■	■	■				
14															■	■	■			
15																■	■	■		
16																	■	■	■	
17																		■	■	■
18																			■	■
19																				■
20																				

54 comparisons



Overview

31

- Duplicate detection
- Similarity measures
- Algorithms
- Data sets and evaluation
- Data fusion



Precision & Recall

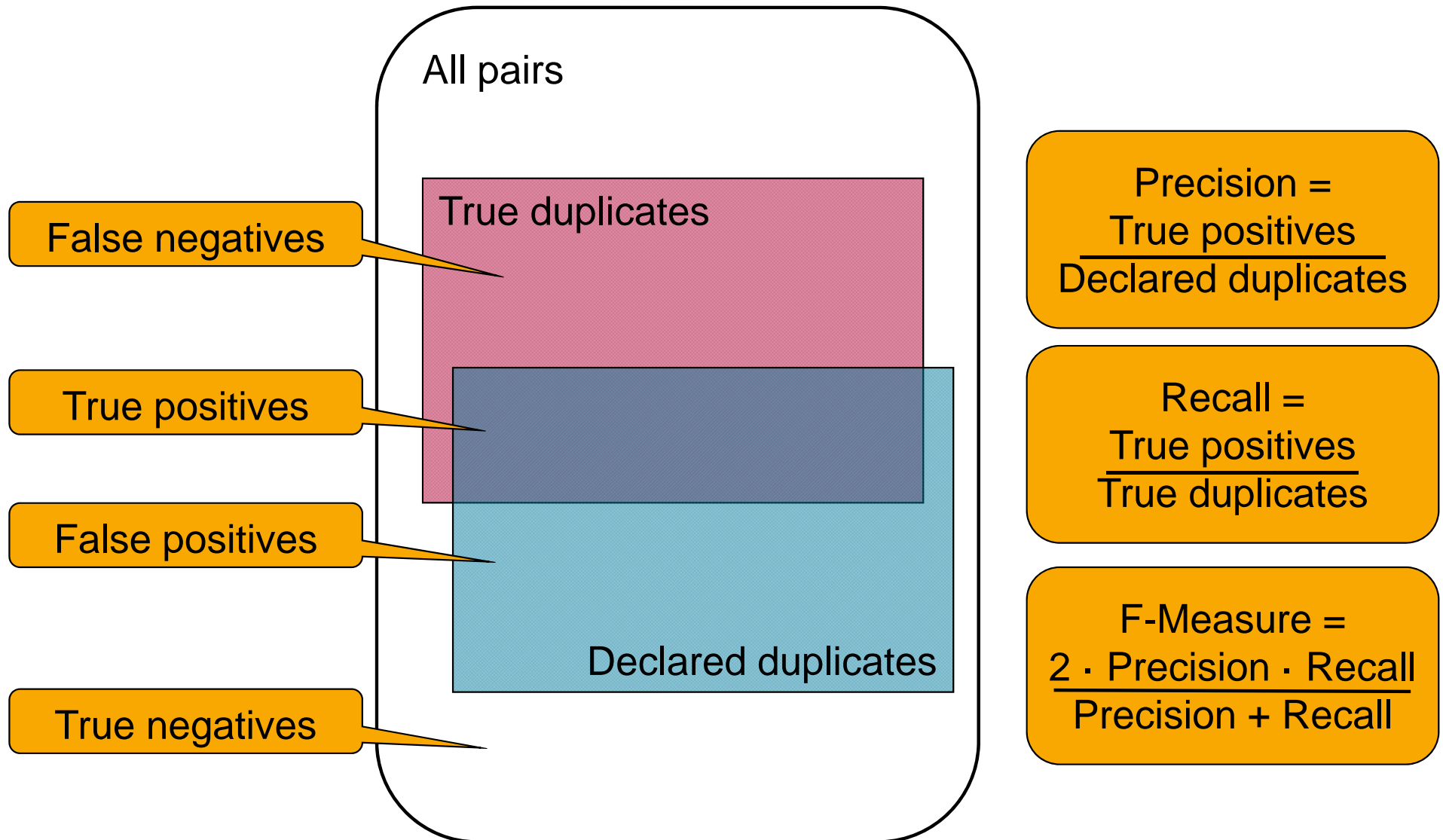
32

- True positives (TP): Correctly declared duplicates
- False positives (FP): Incorrectly declared duplicates
- True negatives (TN): Correctly avoided pairs
- False negatives (FN): Missed duplicates

- Precision = $TP / (TP + FP)$
 - = TP / declared dups
 - Proportion of found matches that are correct
- Recall = $TP / (TP + FN)$
 - = TP / all dups
 - Proportion of correct matches that are found

Precision & Recall

33



F Measure

34

- *Harmonic mean* of recall and precision

$$F = \frac{1}{\frac{1}{2} \left(\frac{1}{R} + \frac{1}{P} \right)} = \frac{2RP}{R + P}$$

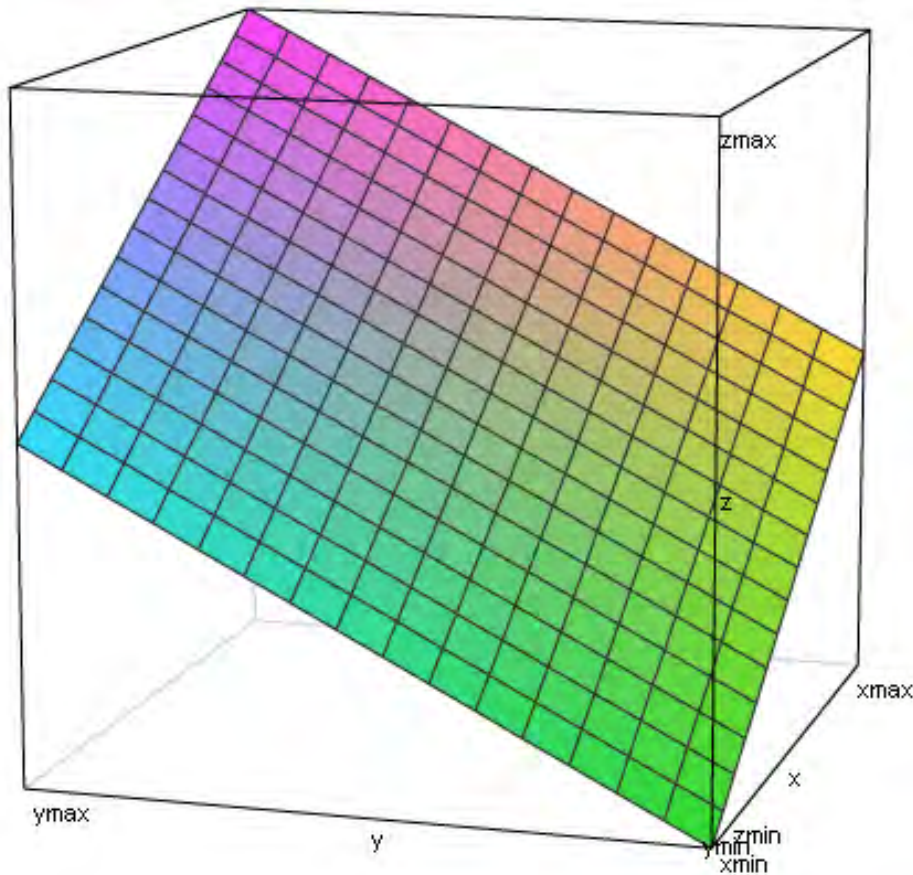
- Harmonic mean emphasizes the importance of small values, whereas arithmetic mean is affected more by outliers that are unusually large.

- More general form: Weighted harmonic mean

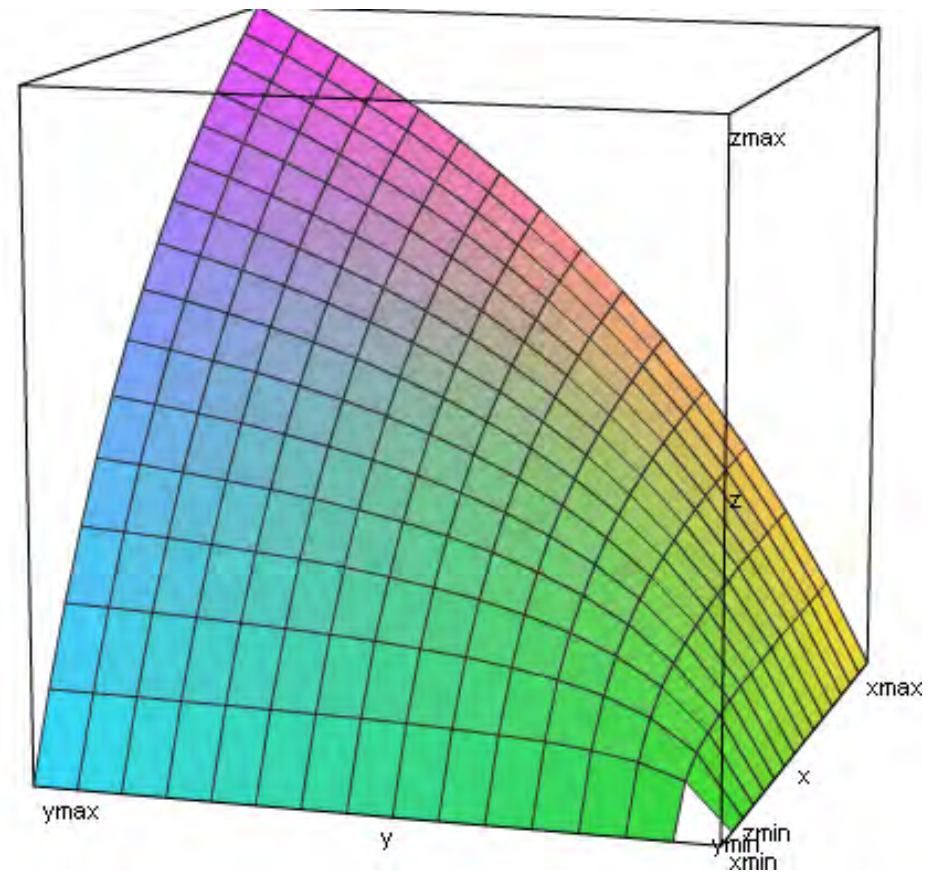
$$F_{\alpha} = \frac{RP}{\alpha R + (1 - \alpha)P}$$

- Thus, harmonic mean is $F_{1/2}$

Arithmetic mean („Average“) vs. Harmonic mean („F-Measure“)



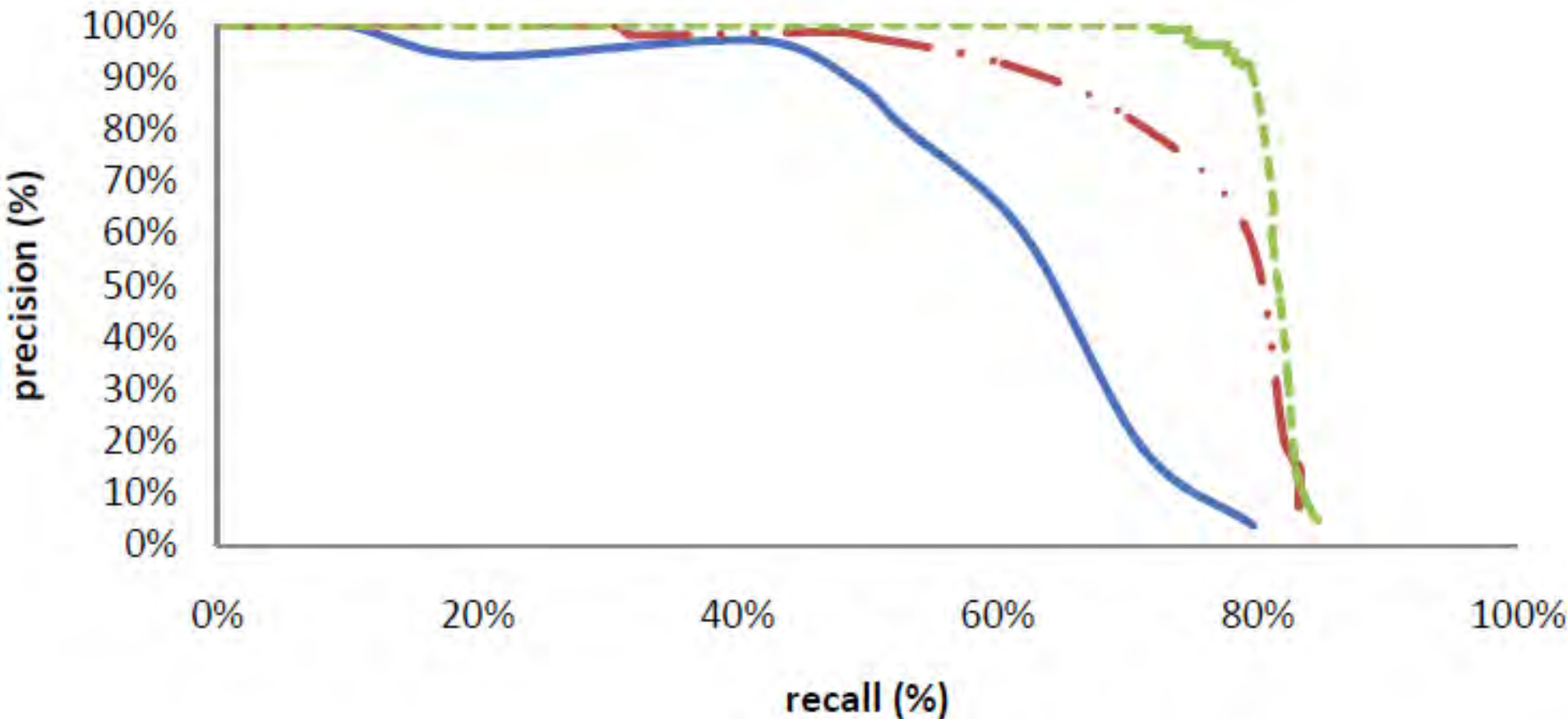
$$z = \frac{1}{2} (x + y)$$



$$z = \frac{2 (x \cdot y)}{(x + y)}$$

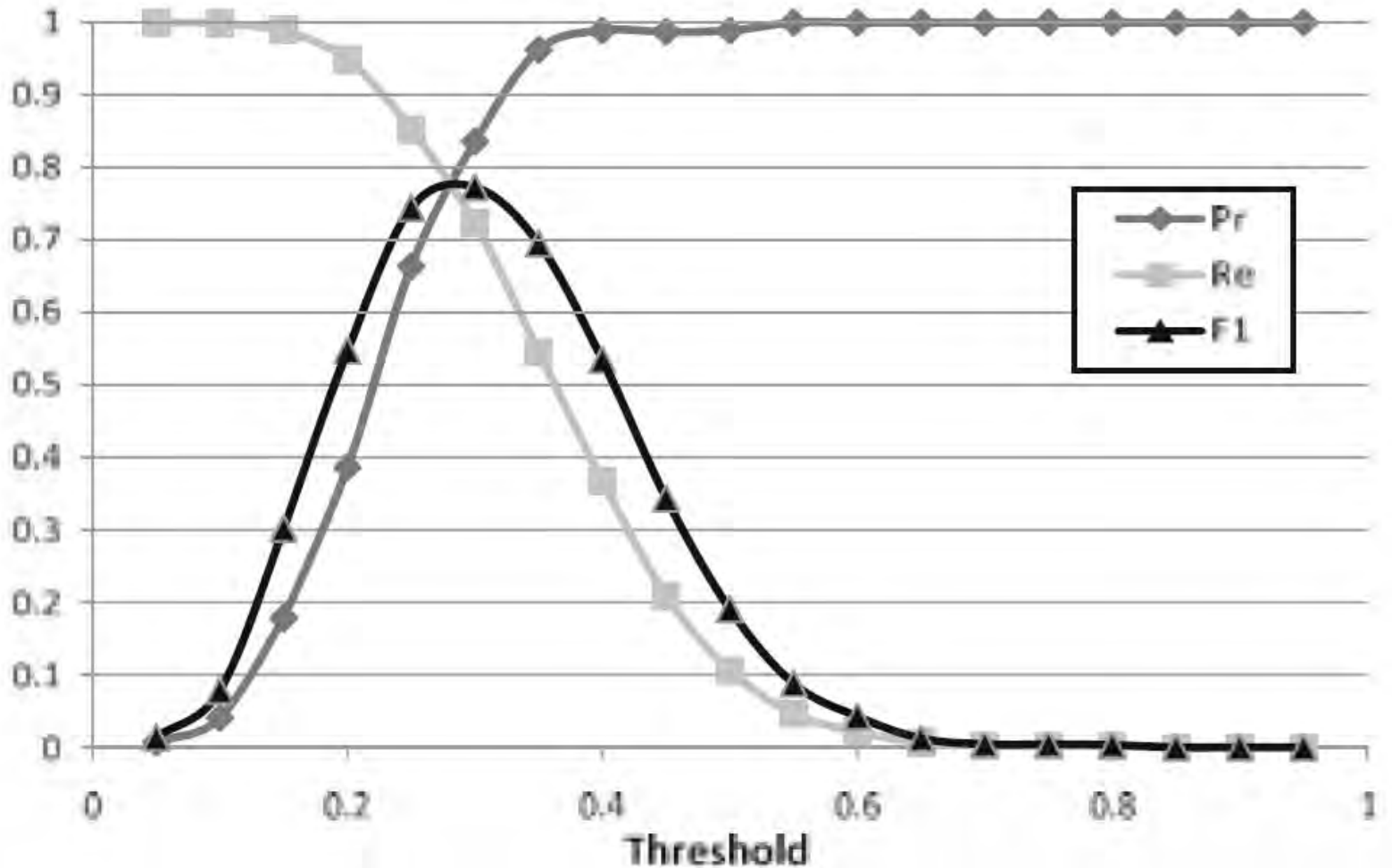
Recall-precision-diagram

36



F1 Graph

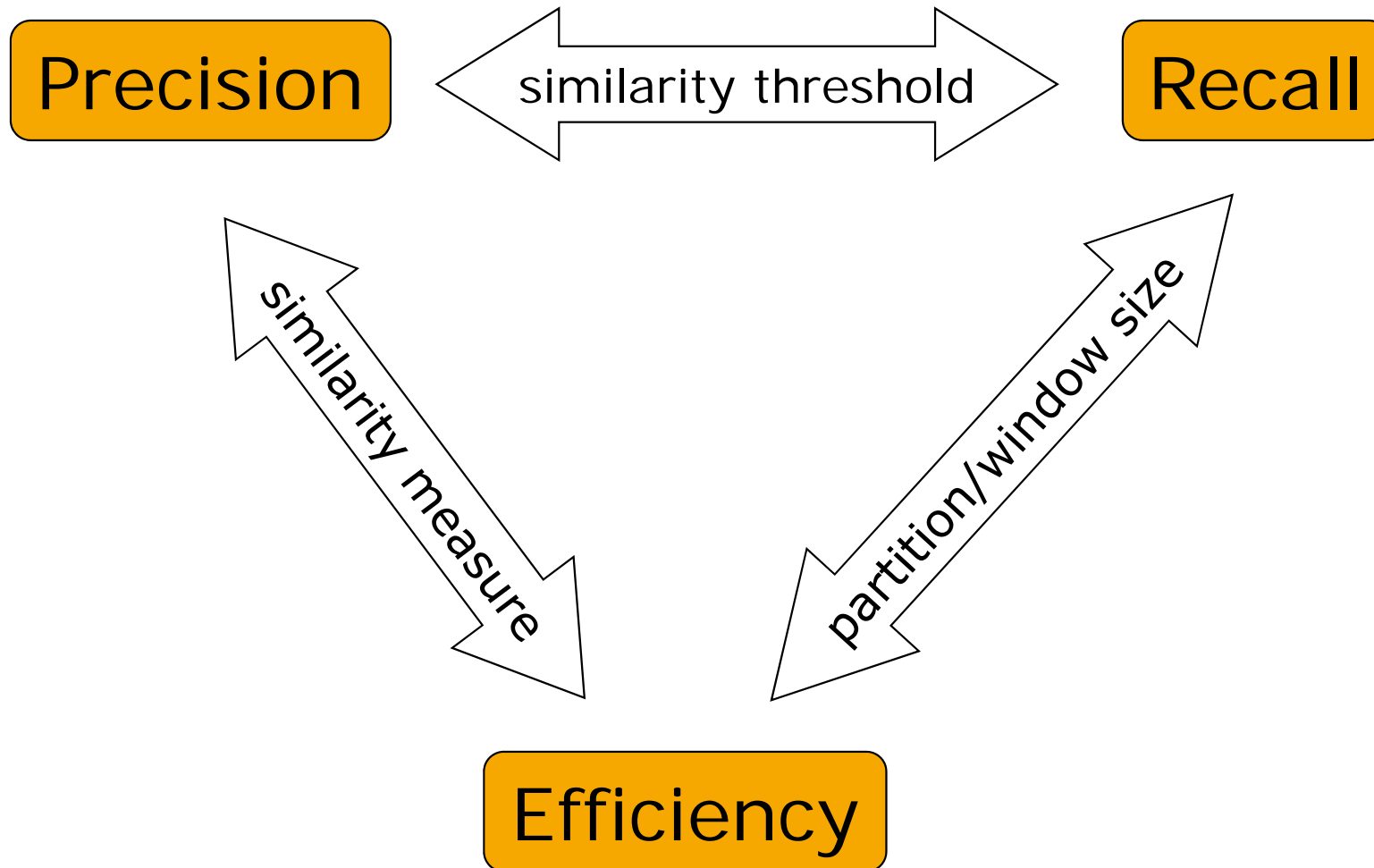
37



From Creating probabilistic databases from duplicated data
Oktie Hassanzadeh · Renée J. Miller (VLDBJ)

Evaluating Duplicate Detection

38



Other effectiveness measures

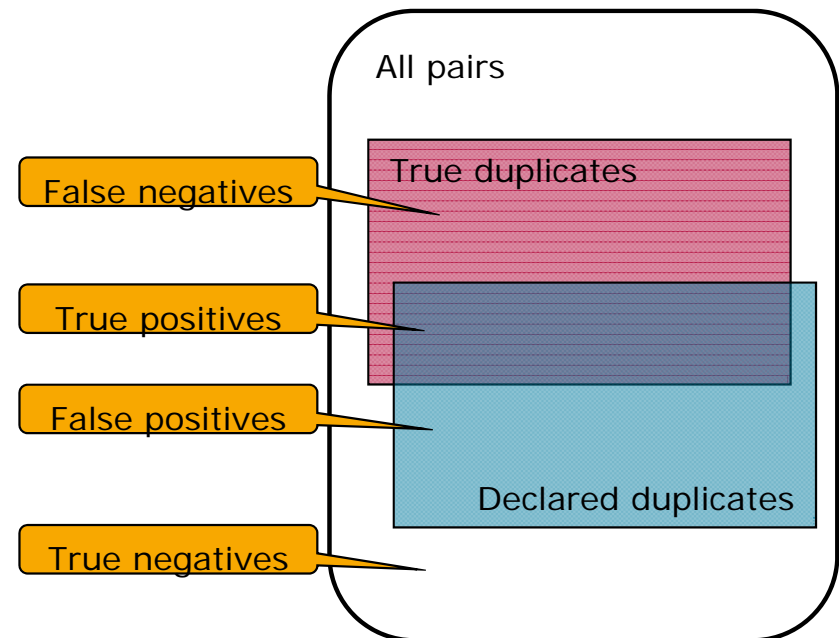
39

■ Accuracy

- $(TP + TN) / (TP + FP + TN + FN)$
- Used for balanced classes
- For duplicate detection, TN usually dominates overall result

■ Specificity

- $TN / (TN + FP)$
- = $TN / \text{true non-matches}$
- Again: TN dominates overall result



Complexity measures

40

- Problem: Measure not only quality of similarity measure, but also that of algorithm
 - Solution 1: Runtime measurements
 - ◇ But: Different hardware, difficult repeatability
 - Solution 2: Measure how well/poor algorithms filter candidates

- Reduction ratio

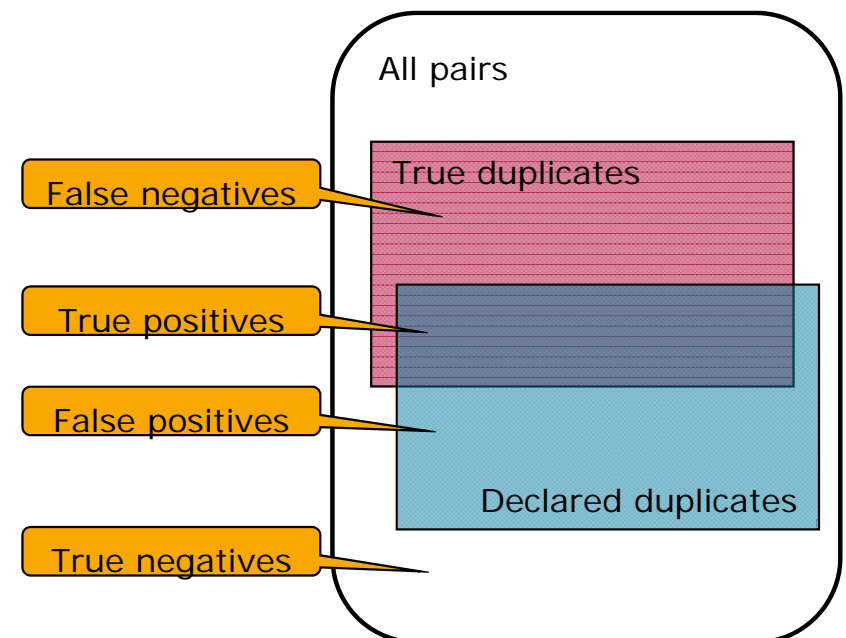
- $1 - ((TP + TN) / (FN + TP + FP + TN))$
 $= 1 - \text{accuracy}$

- Pairs completeness

- $TP / (FN + TP) = \text{recall}$

- Pairs quality

- $TP / (TP + TN)$



Data sets to evaluate deduplication

41

■ Requirements

- Real-world application
- Interestingly large
- Interestingly dirty
- Gold standard
- Publicly available
- (Relational)
- (Understandable)

■ Hardly available

- Privacy issues
- Security issues
- Embarrassment
- Data fiefdoms



Small datasets with gold standard

CORA

- 1878 bibliographic references in XML format
- http://www.hpi.uni-potsdam.de/naumann/projekte/repeatability/datasets/cora_dataset.html

DBLP

- 50,000 bibliographic references in XML format
- http://www.hpi.uni-potsdam.de/naumann/projekte/repeatability/datasets/dblp_dataset.html

Restaurants

- 864 restaurants with 112 duplicates
- <http://www.cs.utexas.edu/users/ml/riddle/data.html>

Whirl datasets

- 11 smaller datasets with a single string attribute
- <http://www.cs.purdue.edu/commugrate/data/whirl/match/>

Large datasets without gold standard

Places

- 1.4 million POIs from Facebook, Gowalla, Foursquare
- <https://www.hpi.uni-potsdam.de/naumann/trac/LuSim/wiki/quellen>

WheelMap

- 120,000 places/things in Germany

FreeDB

- 1.9 million CDs, dirty, some duplicate clusters quite large
- original: http://www.freedb.org/en/download_database.10.html
- derived: http://www.hpi.uni-potsdam.de/naumann/projekte/repeatability/datasets/cd_datasets.html

CITeseerX

- 1.3 million publications in CSV
- <http://asterix.ics.uci.edu/data/csx.raw.txt.gz>

Data generation

44

- For lack of gold standard: create one
- Data base
 - Real-world data sets (without or without enough duplicates)
 - Real-world values (from dictionaries)
 - Synthetic strings
- Data corruption: Duplicate and modify some percentage of tuples
 - Duplication: Cluster sizes?
 - Data values
 - ◇ Insert/remove/transpose/change certain letters
 - ◇ Delete values
 - ◇ Swap values (within tuple, from dictionary, across tuples)
- General suspicion: Similarity measure and candidate selection is geared towards known types of errors.

Data generators

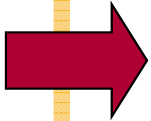
45

- UIS Database Generator
 - Generates a list of randomly perturbed names and US mailing addresses.
 - Written by Mauricio Hernández.
 - <http://www.cs.utexas.edu/users/ml/riddle/data/dbgen.tar.gz>
- FEBRL-Generator
 - Part of a cleansing suite
 - Dictionaries with frequencies
 - <http://sourceforge.net/projects/febrl/>
- Dirty XML Generator
 - http://www.hpi.uni-potsdam.de/naumann/projekte/completed_projects/dirtyxml.html

Overview

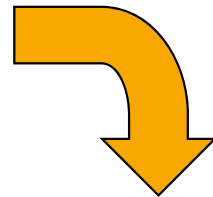
46

- Duplicate detection
- Similarity measures
- Algorithms
- Data sets and evaluation
- Data fusion

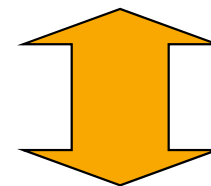
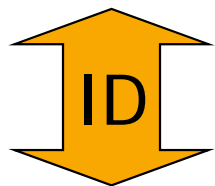




Data Fusion

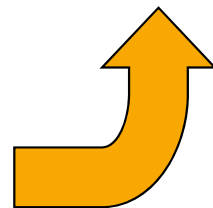
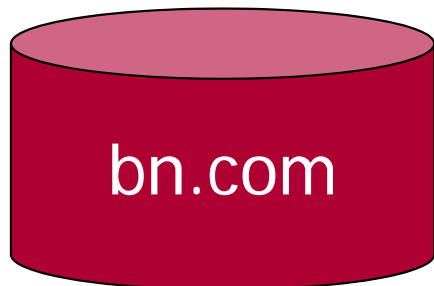
47



0766607194	H. Melville		\$3.98	
------------	-------------	--	--------	---

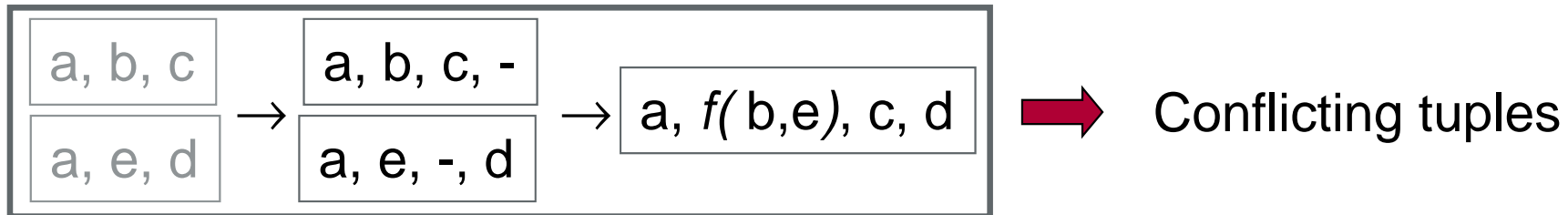
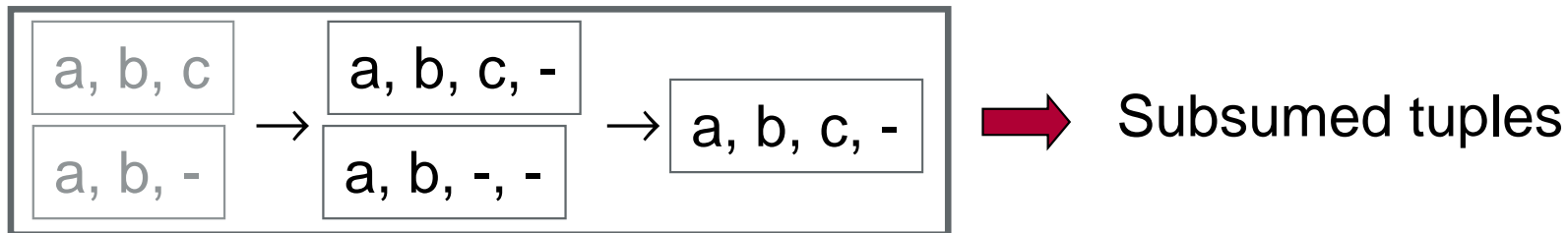
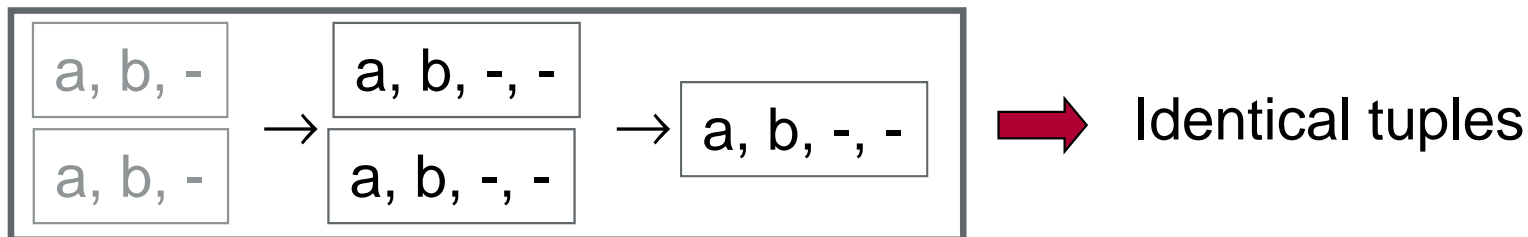
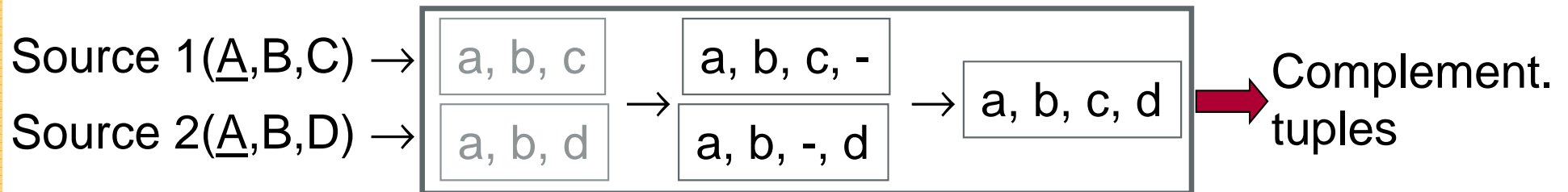


0766607194	Herman Melville	Moby Dick	\$5.99	 
------------	-----------------	-----------	--------	---



“Proper” Data Fusion

48



Conflict Resolution Functions

49

Min, Max, Sum, Count, Avg, StdDev	Standard aggregation
Random	Random choice
First, Last	Choose first/last value; depends on order
Longest, Shortest	Choose longest/shortest value
Choose(source)	Choose value from a particular source
ChooseDepending(col, val)	Choose depending on val in other column col
Vote	Majority decision
Coalesce	Choose first non-null value
Group, Concat	Group or concatenate all values
MostRecent	Choose most recent (up-to-date) value
MostAbstract, MostSpecific	Use a taxonomy / ontology
....

Visualization of Integrated Data

50

HumMer-Demo [Window Title]

File Extra Help

0. Sources
 1. Matching
 2. Duplicate Definition
 3. Duplicate Detection
 4. Conflict
 5. Result

Result

Choose the fusion implementation to use **default** [Dropdown] **Execute** [Button]

back next [Buttons]

#	CLU...	TITLE	VERSI...	COUN...	YEAR	ORIGI...	GENRE	DIREC.
		VOTE	COALESCE	COALES...	MAX	COALES...	LAST	COALES.
13	87	HOPE FLOATS	engl...	USA	1998	Hop...	Unterhaltu...	Fore...
14	84	GOOD WILL H...	engl...	USA	1998	Goo...	Drama	Gus .
15	83	GODZILLA	engl...	USA	1998	God...	Fantasy, S...	Rola.
16	80	Gadjo Dilo Gadjo Dilo GADJO DILO	franz... franz.&r... ↓	F/Rum ↓	1998 1998 1997	Gadj... Gadjo ... ↓	Unterhaltu... Unterhaltung Drama	Ton.
17	77	Deconstructin...	engl...	USA	1998	Dec...	Komödie/...	Woo.
18	74	City Of Angels	engl...	USA	1998	City ...	Drama	Brad.
19	69	BOOGIE NIGH...	engl...	USA	1998	Boo...	biografisc...	Paul .
20	65	Antz	engl...	USA	1998	Antz	Animation,...	Darn.
21	57	SPIDER	↓	↓	2002	↓	Drama	↓
22	51	SECRETARY	↓	↓	2002	↓	Komödie	↓
23	49	S.F.W.	↓	↓	1994	↓	Komödie	↓
24	31	Intolerable Cr...	↓	↓	2003	↓	Komödie	↓
25	25	GANGSTER N...	↓	↓	2000	↓	Gangsterfi...	↓
26	24	From Hell	↓	↓	2001	↓	↓	↓
27	17	DEATHWATCH	↓	↓	2002	↓	Kriegsfilm	↓
28	15	CHARLOTTE ...	↓	↓	2001	↓	Melodram	↓
29	11	Big Fish	↓	↓	2003	↓	Drama	↓

Rows: 0:99

Duplicate Contradiction Uncertainty Unique [Legend]

Start Over [Button] Done [Button] Back [Button] Next [Button]

Tool-based Data Fusion

51

Fuzzy Fuzlon Additional Information Test/Debug

Gruppen 0 bis 50 von 39449 Filtermodus Wert:

fdb.gr...	TITLE	SALUT...	FIRST...	LASTN...	COMP...	COUN...	STREET	STREE...	ZIP	CITY	ADR1	ADR2	ADR3	ADR4	ADR5
1253		Frau u...		Koste...	Daimle...	D	Alt-Mo...	96 A	10559	Berlin	Frau u...	Kosten...	Daimle...	Alt-Mo...	
1333		Herr	Markus	Bauer		D	HPC V...		10878	Berlin	Herr	Marku...	HPC V...		D - 10...
1782		Herr	Frank	Leusc...		D	Arenh...		12103	Berlin	Herr	Frank ...	Arenh...		D - 12...
1874		Monsieur	Frank	Eichler		D	Falken...	78A	13589	Berlin	Monsieur	Frank ...	Falken...		D - 13...
2159		Herr	Horst	Fucks		D	Nassa...		10717	Berlin	Herr	Horst ...	Nassa...		D - 10...
2196	Dr.	Frau	Christa	Schün...		D	Uhlan...	121	10717	Berlin	Frau	Dr. Ch...	Uhlan...		D - 10...
2217	Dr.	Familie		Hofma...		D	Hede...	13	10969	Berlin	Familie	Dr. H...	Hede...		D - 10...
2498		Frau	Julia	Görsc...		D	Regin...	20	13409	Berlin	Frau	Julia G...	Regin...	13409 ...	
2552		Herr		Ehlers		D	Über F...		10587	Berlin	Herr	Ehlers	Über F...	10587 ...	

10. Gruppe :

fdb.group	TITLE	SALUTATION	FIRSTNAME	LASTNAME	COMPAN...	COUNTRY_CODE	STREET	STREET_NUMBER	ZIP	
1874		Monsieur	Frank	Eichler		D	Falkenseer Chaussee	78A	13589	Berlin
1874		Firma		Eichler		D	Flkenseer Chaussee 78A		13589	Berlin
1874		Herr	Frank	Eichler	Zres	D	Falkenseer Chaussee 78 A		13589	Berlin
Erster	..	GLOBAL...	Erster	Erster	Erst...	Erster	gemischte Schreibweise	Erster	Erster	Erster
1874		Monsieur	Frank	Eichler	Zres			78A	13589	Berlin

References

52

- [Ananthakrishna et al. 2002] R. Ananthakrishna, S. Chaudhuri, V. Ganti: Eliminating Fuzzy Duplicates in Data Warehouses, *Proc. of the 28th VLDB Conference*, Hong Kong, China, pages 586-597, 2002.
- [Dey et al. 2002] D. Dey, S. Sarkar, P. De: A Distance-based Approach to Entity Reconciliation in Heterogeneous Databases, *IEEE Transactions on Knowledge and Data Engineering (TKDE)* 14(3): 567-582, 2002.
- [Gravano et al. 2001] L. Gravano, P. Ipeirotis, H. Jagadish, N. Koudas, S. Muthukrishnan, D. Srivastava: Approximate String Joins in a Database (Almost) for Free, *Proc. of the 27th VLDB Conference*, Roma, Italy, pages 491-500, 2001.
- [Hernandez Stolfo 1995] M. Hernandez, S. Stolfo: The Merge/Purge Problem for Large Databases, *Proc. ACM SIGMOD Conference 1995*, San Jose, USA, pages 127-138, 1995.
- [Hernandez Stolfo 1998] M. Hernandez, S. Stolfo: Real-world Data is Dirty: Data Cleansing and the Merge/Purge, *Journal of Data Mining and Knowledge Discovery*, 2(1):9-37, 1998.
- [Jaro 1989] M. Jaro: Advances in Record Linkage Methodology as Applied to Matching the 1985 Census of Tampa, Florida, *Journal of the American Statistical Association* 84(406):414-420, 1989.
- [Low et al. 2001] W. Low, M. Lee, T. Ling: A Knowledge-based Approach for Duplicate Elimination in Data Cleaning, *Information Systems* 26(8):585-606, 2001.
- [Monge 2000] A. Monge: Matching Algorithms within a Duplicate Detection System, *IEEE Data Engineering Bulletin*, 23(4):14-20, 2000.
- [Navarro 2001] G. Navarro: A Guided Tour of Approximate String Matching, *ACM Computing Surveys* 31(1):31-88, 2001.
- [Weis Naumann 2005] Melanie Weis, Felix Naumann: DogmatiX Tracks down Duplicates in XML. In Proc. of the ACM Conference on Management of Data (SIGMOD) 2005.
- [GL94] Outerjoins as Disjunctions, Cesar A. Galindo-Legaria, SIGMOD 1994 conference
- [Cod79] E. F. Codd: Extending the Database Relational Model to Capture More Meaning. *TODS* 4(4): 397-434 (1979)
- [Ull89] Jeffrey D. Ullman: Principles of Database and Knowledge-Base Systems, Volume II. Computer Science Press 1989.