

Informationsqualität
Antrittsvorlesung am 26.4.2007

Felix Naumann
Hasso-Plattner-Institut
Fachgebiet Informationssysteme

Überblick

2

- ➔ ■ Informationsqualität
- Informationsintegration
- Duplikaterkennung
 - Ähnlichkeit
 - Algorithmen
- Ausblick



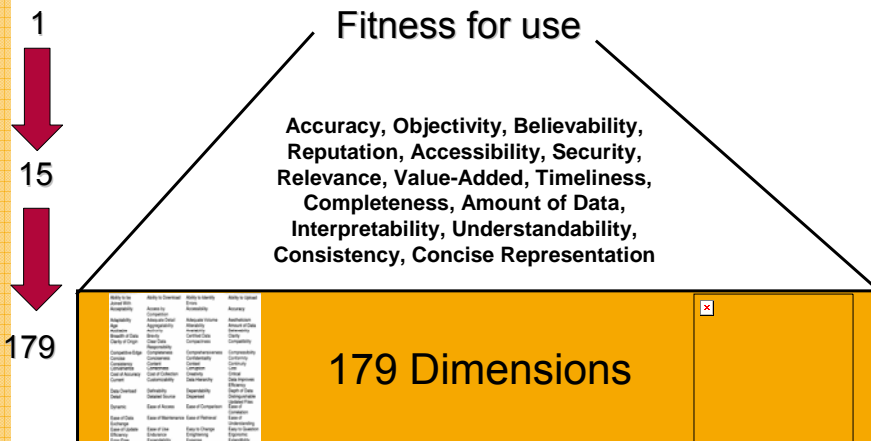
"Even though quality cannot be defined, you know what it is."

Robert Pirsig



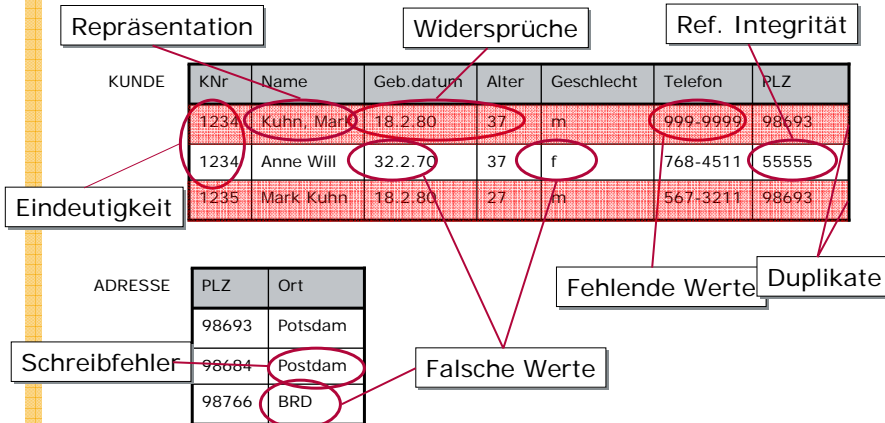
Robert M. Pirsig
**Zen und die Kunst
ein Motorrad
zu warten**

Fischer



Datenqualität: Probleme

5



Felix Naumann | Antrittsvorlesung | April 2007

DQ-Probleme: Auswirkungen

6

- Fehlerhafte Warenpreise in Artikel-DB des US-Einzelhandels [English 1999]
 - Kosten für Konsumenten 2.5 Mrd \$
 - 80% der Barcode-Scan-Fehler zulasten der Konsumenten
- US-Finanzbehörde 1992: knapp 100.000 Steuererstattungsbescheide nicht zustellbar [English 1999]
- 50-80% der Einträge im US-Vorstrafenregister ungenau, unvollständig oder fehlerhaft [Strong et al. 1997a]
- US-Post: von 100.000 Massen-Postsendungen bis zu 7.000 aufgrund von Adressfehlern nicht zustellbar [Pierce 2004]

Goodyear reveals \$100 million error ^{10/23/03 USA}
 Goodyear said late Wednesday that it will restate earnings for the past five years, decreasing income by as much as \$100 million because an accounting system caused billing errors. The tiremaker is delaying the release of its third-quarter earnings, expected this morning, until mid-November. Shares closed up 2 cents to \$6.83 before the announcement; in after-hours trading, shares plummeted 27%, or \$1.83, to \$5.

Felix Naumann | Antrittsvorlesung | Apr...

DQ-Probleme: Auswirkungen

7

- 2006 werden die Fortune 1000 Unternehmen mehr Geld wegen DQ-Problemen ausgeben als für ERP, CRM und BI zusammen [Gartner]
- Mehr als 35% aller IT-Projekte schlagen fehl aufgrund von DQ-Problemen; die jährlichen Kosten betragen allein in der USA 2-4 Mrd. \$ [Meta Group]
- DQ ist einer der wichtigsten Erfolgsfaktoren in DWH- und CRM-Projekten [PriceWaterhouseCoopers]
- Datensammlungen in der Folge des Tsunami 2004
 - Todesfälle und Verletzungen
 - Zerstörung von Häusern und Eigentum
 - Zuweisung von Hilfsmitteln und Unterstützung

Quelle: <http://www.informationquality.org/publiclyexposediqproblems.cfm>

Felix Naumann | Antrittsvorlesung | April 2007

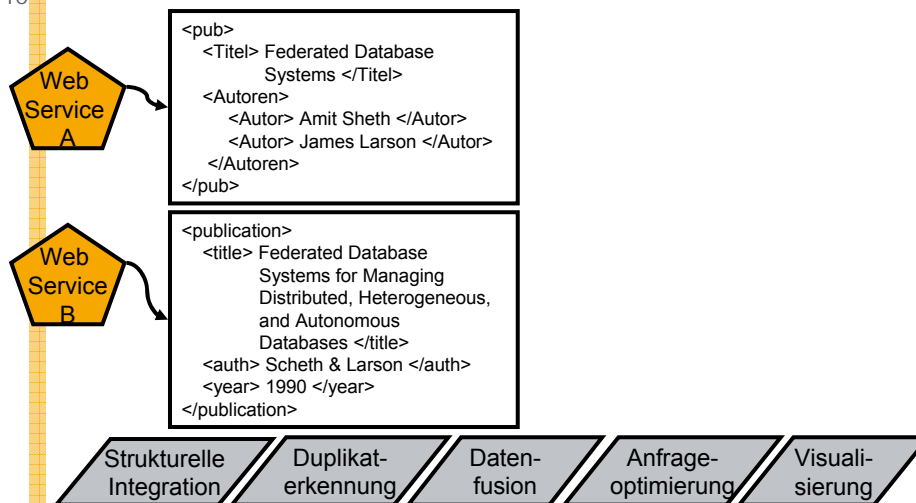
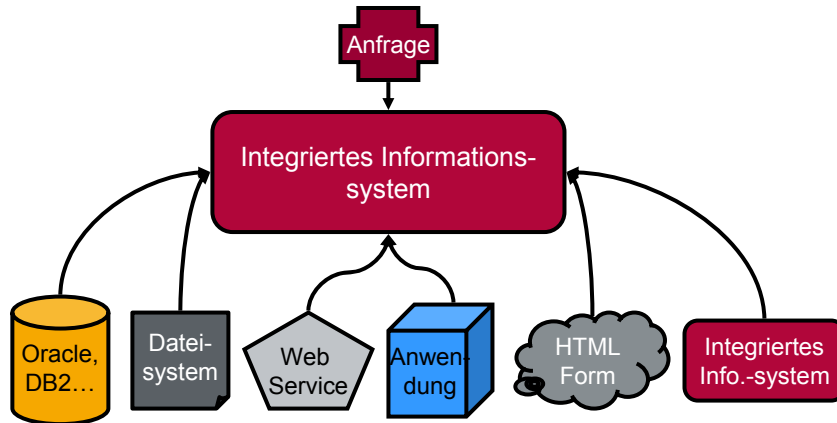
Überblick

8

- Informationsqualität
- Informationsintegration
- Duplikaterkennung
 - Ähnlichkeit
 - Algorithmen
- Ausblick

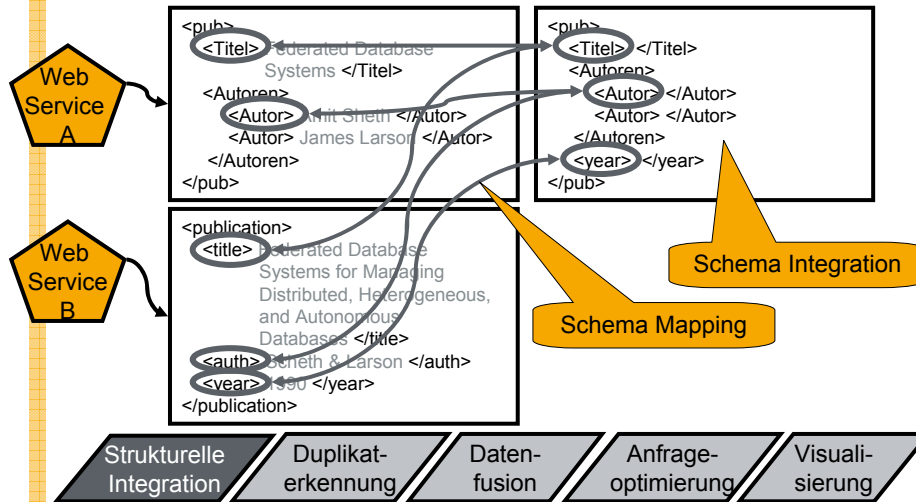


Felix Naumann | Antrittsvorlesung | April 2007



Informationsintegration

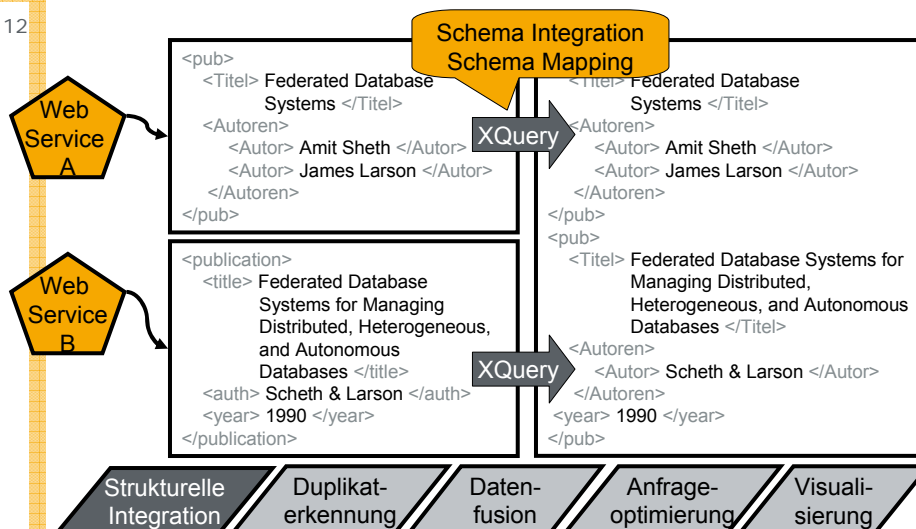
11



Felix Naumann | Antrittsvorlesung | April 2007

Informationsintegration

12



Felix Naumann | Antrittsvorlesung | April 2007

Informationsintegration

13

Web Service A

```
<pub>
  <Titel> Federated Database
  Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
```

Web Service B

```
<publication>
  <title> Federated Database
  Systems for Managing
  Distributed, Heterogeneous,
  and Autonomous
  Databases </title>
  <auth> Scheth & Larson </auth>
  <year> 1990 </year>
</publication>
```

```
<pub>
  <Titel> Federated Database
  Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
<pub>
  <Titel> Federated Database Systems for
  Managing Distributed,
  Heterogeneous, and Autonomous
  Databases </Titel>
  <Autoren>
    <Autor> Scheth & Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```



Felix Naumann | Antrittsvorlesung | April 2007

Informationsintegration

14

Web Service A

```
<pub>
  <Titel> Federated Database
  Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
```

Web Service B

```
<publication>
  <title> Federated Database
  Systems for Managing
  Distributed, Heterogeneous,
  and Autonomous
  Databases </title>
  <auth> Scheth & Larson </auth>
  <year> 1990 </year>
</publication>
```

```
<pub>
  <Titel> Federated Database
  Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
<pub>
  <Titel> Federated Database Systems for
  Managing Distributed,
  Heterogeneous, and Autonomous
  Databases </Titel>
  <Autoren>
    <Autor> Scheth & Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```



Felix Naumann | Antrittsvorlesung | April 2007

Informationsintegration

15

Web Service A

Web Service B

```
<pub>
  <Titel> Federated Database Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
<pub>
  <Titel> Federated Database Systems for
  Managing Distributed,
  Heterogeneous, and Autonomous
  Databases </Titel>
  <Autoren>
    <Autor> Scheth & Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```

```
<pub>
  <Titel> Federated Database Systems for
  Managing Distributed,
  Heterogeneous, and
  Autonomous Databases </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```



Felix Naumann | Antrittsvorlesung | April 2007

Informationsintegration

16

Web Service A
1sec.

Web Service B
5sec

```
<pub>
  <Titel> Federated Database Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
<pub>
  <Titel> Federated Database Systems for
  Managing Distributed,
  Heterogeneous, and Autonomous
  Databases </Titel>
  <Autoren>
    <Autor> Scheth & Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```

```
<pub>
  <Titel> Federated Database Systems for
  Managing Distributed,
  Heterogeneous, and
  Autonomous Databases </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```



Felix Naumann | Antrittsvorlesung | April 2007

Informationsintegration

17

Web Service A

1sec.

```
<pub>
  <Titel> Federated Database Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
<pub>
  <Titel> Federated Database Systems for Managing Distributed, Heterogeneous, and Autonomous Databases </Titel>
  <Autoren>
    <Autor> Scheth & Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```

Web Service B

5sec.

```
<pub>
  <Titel> Federated Database Systems for Managing Distributed, Heterogeneous, and Autonomous Databases </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
  <year> 1990 </year>
</pub>
```

WS B

WS A

WS B



Überblick

18

- Informationsqualität
- Informationsintegration
- ➔ ■ Duplikaterkennung
 - Ähnlichkeit
 - Algorithmen
- Ausblick



Duplikaterkennung

19

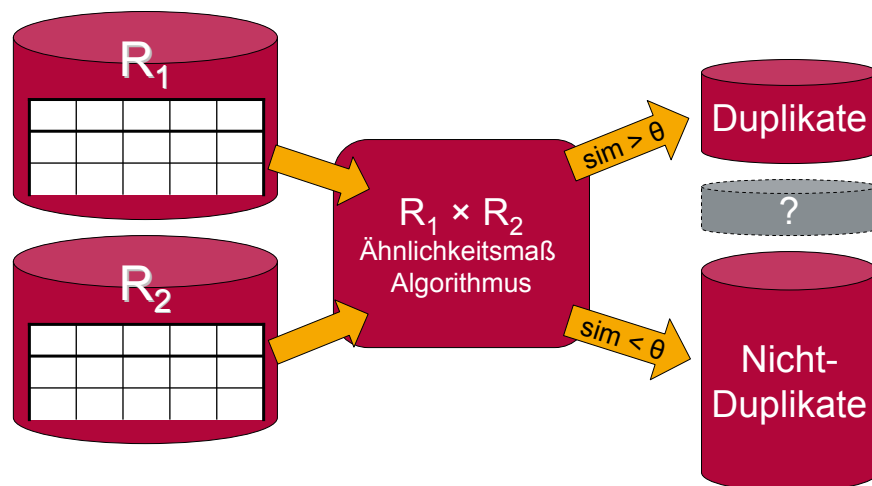
Duplikaterkennung ist das Finden mehrerer Repräsentationen desselben Realweltobjekts.

- Problem 1: Repräsentationen sind nicht identisch.
 - *Fuzzy duplicates*
- Lösung: Ähnlichkeitsmaße
 - Wert- und Datensatzvergleiche
 - Domänenunabhängig oder -abhängig
- Problem 2: Die Datenmenge ist groß.
 - Quadratischer Aufwand: Jedes Paar muss verglichen werden.
- Lösung: Algorithmen
 - Z.B. Vergleiche durch Partitionierung vermeiden

Felix Naumann | Antrittsvorlesung | April 2007

Duplikaterkennung

20



Felix Naumann | Antrittsvorlesung | April 2007

Wirkungen von Duplikaten

21

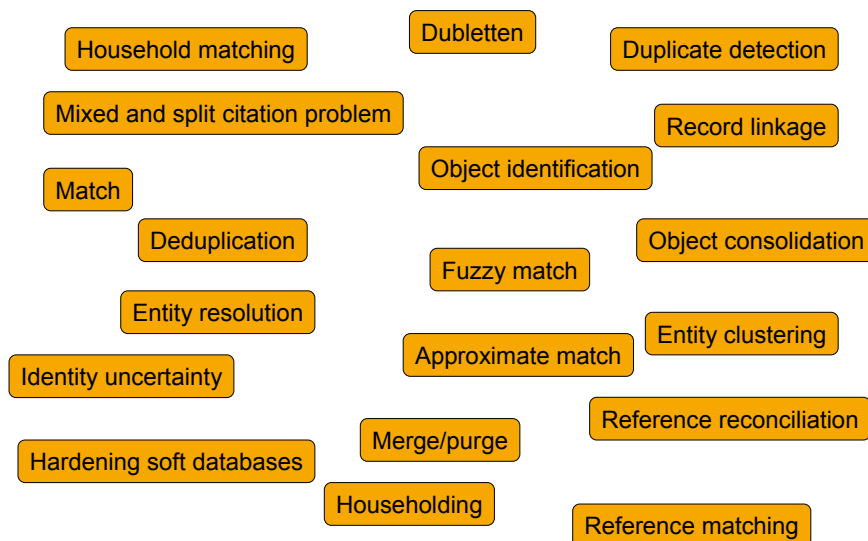
- Mehrfache Zusendung von Katalogen
- Rechnungen werden doppelt bezahlt
- Banken
 - Überschreiten des Kreditlimits wird nicht erkannt
- Lagerhaltung / Einkauf
 - Zu niedriger Lagerbestand einzelner Waren wird ausgewiesen.
 - Kein Ausnutzen von Mengenrabatten bei Bestellungen
- Gesamtumsatz eines Kunden bleibt unbekannt.
- Mehraufwand in der IT
- Sinkende Kundenzufriedenheit
- Potenziale und Gefahren nicht erkannt
- Inkorrekte Kennzahlen

Kunde	Umsatz
BMW	20.000
BaMoWe	5.000.000
Bayerische Motorenwerke	300.000
...	...

Felix Naumann | Antrittsvorlesung | April 2007

“Duplikaterkennung” hat viele Duplikate

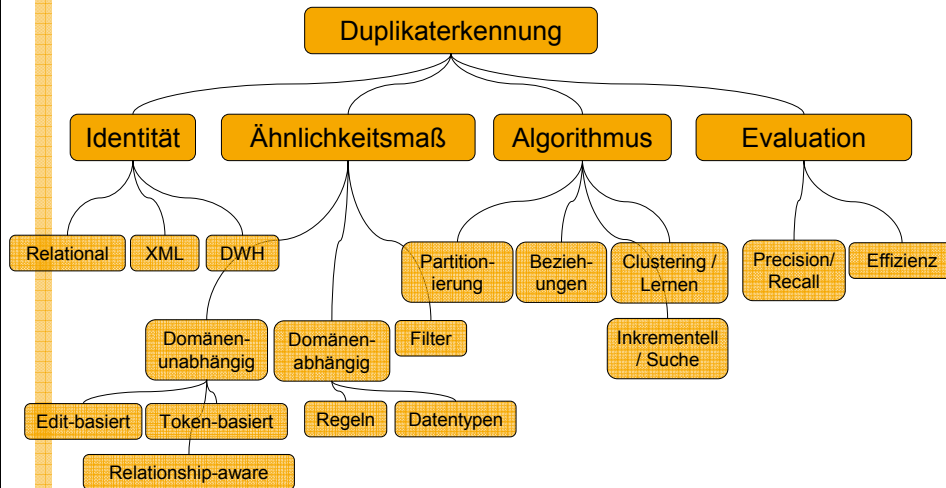
22



Felix Naumann | Antrittsvorlesung | April 2007

Duplikaterkennung

23



Felix Naumann | Antrittsvorlesung | April 2007

Überblick

24

- Informationsqualität
- Informationsintegration
- Duplikaterkennung
 - Ähnlichkeit
 - Algorithmen
- Ausblick



Felix Naumann | Antrittsvorlesung | April 2007

Token-basierte Ähnlichkeitsmaße

25

- Tokens
 - Words / Terms
 - n-grams
- Jaccard
 - $|\{\text{gemeinsame token}\}| / |\{\text{alle token}\}|$
- TFIDF
 - *Term frequency* (tf)
 - *Inverse document frequency* (idf)
 - TFIDF: $\log(\text{tf} + 1) \times \log \text{idf}$
 - Häufige Wörter haben niedriges Gewicht.
 - Ähnlichkeit ist Kosinus der Termvektoren, gewichtet durch TFIDF.
- ...

Felix Naumann | Antrittsvorlesung | April 2007

Edit-basierte Ähnlichkeitsmaße

26

- Edit-Distanz / Levenshtein-Distanz [Levenshtein 1965]
 - Minimale Anzahl von Edit-Operationen um den einen String in den anderen umzuwandeln.
 - Domänenspezifische Kosten
- Jaro [Jaro 1989] / Jaro-Winkler [Winkler 1999]
 - Common letters within $\frac{1}{2}$ string length
 - Transposed letters
- Soundex
 - 4-stelliger Code für jedes Wort
 - SOUNDEX('Farwick ') = F620
- ...

Frass, Fricke,
Fahruschi,
Feuerhake

Felix Naumann | Antrittsvorlesung | April 2007

Domänenabhängige Ähnlichkeitsmaße

27

Datentypen

- Spezielle Ähnlichkeitsmaße für Datumsangaben
- Spezielle Ähnlichkeitsmaße für numerische Attribute
- ...

Regeln

- [Hernandez Stolfo 1998], [Lee et al. 2000]
- Gegeben zwei Datensätze r1 und 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
    
```

Felix Naumann | Antrittsvorlesung | April 2007

Beziehungs-basierte Ähnlichkeitsmaße

28

Idee: Nicht nur die Werte der Datensätze selbst, sondern die Werte verwandter Datensätze sind relevant.

- Personen: Ehepartner, Kinder, Arbeitgeber
- Filme: Schauspieler
- CDs: Songs
- Kunden: Bestellungen, Adressen
- Dimensionen im DWH

ID	Land
1	USA
2	United States
3	Unitd States

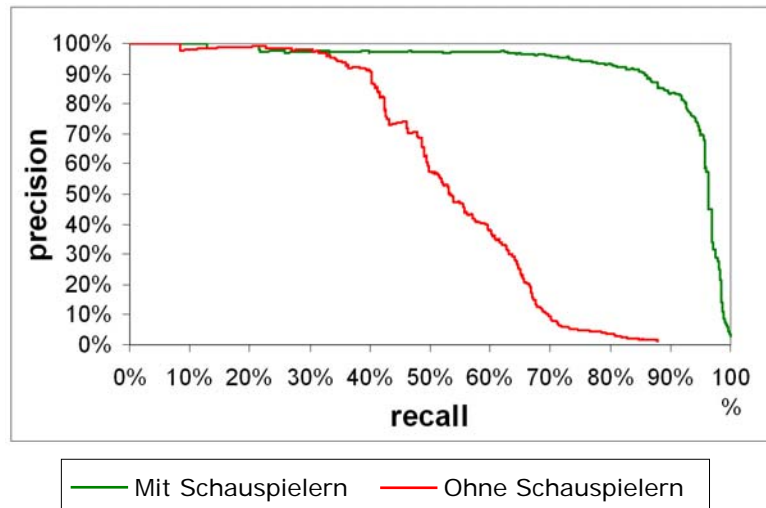
ID	Stadt	Land_ID
1	New York	1
2	Los Angeles	1
3	Now York	2
4	Los Angeles	2
5	New York	3
6	Los Angels	3

ID	Strasse
1	First Ave
2	High St.
3	Broadwa
4	Embarca
5	Broadwa
6	Second S
7	P St.
8	Pennsylv
9	Sunset B
10	Santa M
11	Ocean A

Felix Naumann | Antrittsvorlesung | April 2007

Beziehungs-basierte Ähnlichkeitsmaße – Evaluation

29



Felix Naumann | Antrittsvorlesung | April 2007

Überblick

30

- Informationsqualität
- Informationsintegration
- Duplikaterkennung
 - Ähnlichkeit
 - Algorithmen
- Ausblick



Felix Naumann | Antrittsvorlesung | April 2007

Komplexität

31

Problem: Zu viele Vergleiche

- 10.000 Kunden => 49.995.000 Vergleiche
 - $(n^2 - n) / 2$
 - Jeder Vergleich ist teuer.

Idee: Vergleiche vermeiden

- Durch Herausfiltern einzelner Datensätze
- Durch Partitionierung (Heuristiken)

Felix Naumann | Antrittsvorlesung | April 2007

Partitionierung

32

Idee: Partitioniere Datensätze und vergleiche Paare nur innerhalb der Partition

- Partitionierung nach ersten zwei Ziffern der PLZ
 - Ca. 100 Partitionen
 - Ca. 100 Kunden / Partition
 - Insgesamt 495.000 Vergleiche
- Partitionierung nach erstem Buchstaben des Nachnamens
- ...



Quelle: wikipedia.de

Idee: Partitioniere mehrfach nach unterschiedlichen Kriterien

- Transitive Hülle

Felix Naumann | Antrittsvorlesung | April 2007

Sorted Neighborhood [Hernandez Stolfo 1998]

33

Sortierte Nachbarschaft

- Sortiere Datensätze so, dass ähnliche Datensätze nah beieinander liegen.
- Vergleiche nur Datensätze, die innerhalb einer Nachbarschaft liegen.

Algorithmus

1. Sortierschlüssel erzeugen
 - Wahl des Schlüssels
2. Sortieren
3. Fenster über sortierte Liste schieben und nur innerhalb des Fensters vergleichen
 - Wahl der Fenstergröße
4. Transitive Hülle bilden

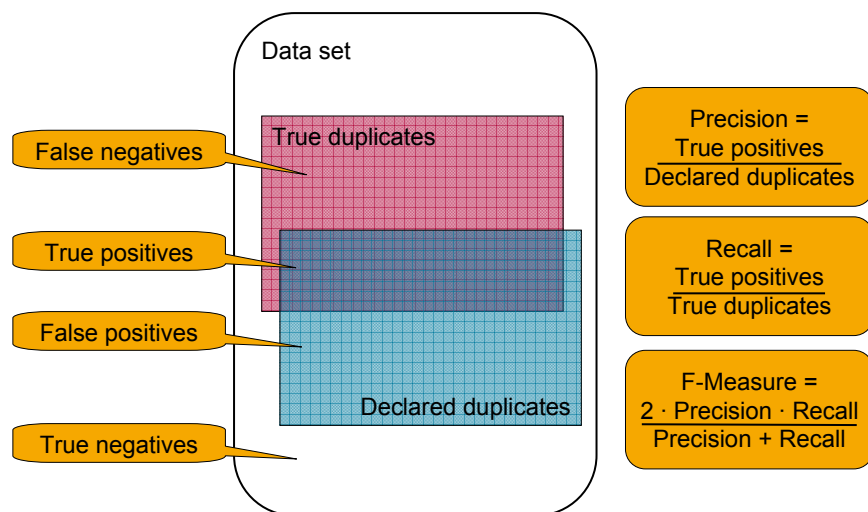
Komplexität ist I/O-bound

- 3 Durchläufe durch die Daten

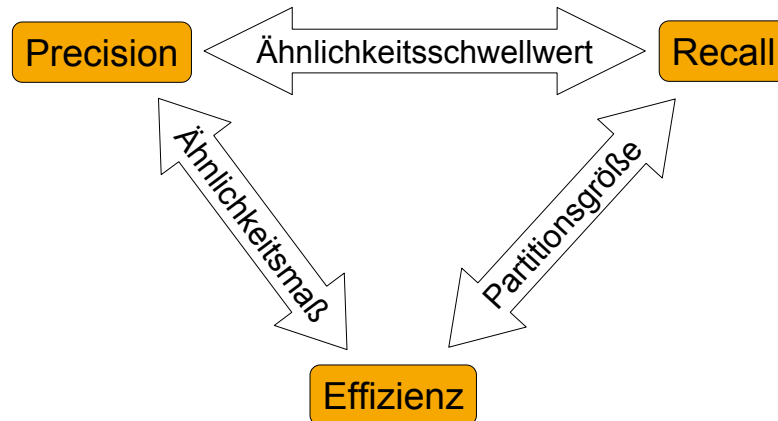
Felix Naumann | Antrittsvorlesung | April 2007

Precision & Recall

38



Felix Naumann | Antrittsvorlesung | April 2007



Duplikatfreiheit ist nur eine Dimension der Informationsqualität.

- Weitere „harte“ Qualitätsdimensionen
 - Vollständigkeit (*completeness*)
 - Genauigkeit (*accuracy*)
 - Aktualität (*timeliness*)
- Herausforderungen
 - Messung
 - Aggregation
 - Verwendung

Überblick

41

- Informationsqualität
- Informationsintegration
- Duplikaterkennung
 - Ähnlichkeit
 - Algorithmen
- ➔ ■ Ausblick



Felix Naumann | Antrittsvorlesung | April 2007

Datenfusion

42



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



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



Felix Naumann | Antrittsvorlesung | April 2007

Die Arbeitsgruppe Informationssysteme

45

- Wissenschaftliche Mitarbeiter / Doktoranden
 - Alexander Albrecht: Personal Information Management
 - Jana Bauckmann: Data Profiling, Aladin
 - Jens Bleiholder: Data Fusion, HumMer & FuSem
 - Frank Käufer: Matching, Forschungskolleg
 - Armin Roth: Peer-Daten-Management, System P
 - Melanie Weis: Duplicate Detection
- Studentische Hilfskräfte
 - Karsten Draba: HumMer
 - Christoph Böhm: Ranking, SPRINT
 - NN: Aladin Projekt
 - Matthias Weidlich: System P
- <http://www.hpi.uni-potsdam.de/~naumann/>

