

Introduction to Data Cleaning

Helena Galhardas
DEI/IST

1

References

- No single reference!
- “Data Quality: Concepts, Methodologies and Techniques”, C. Batini and M. Scannapieco, Springer-Verlag, 2006 (Chapts. 1, 2, and 4)
- Slides “Data Quality and Data Cleansing” course, Felix Naumann, Winter 2014/15
- “Foundations of Data Quality Management”, W. Fan and F. Geerts, 2012
- Oliveira, P. (2009). “Detecção e correcção de problemas de qualidade de dados: Modelo, Sintaxe e Semântica”. PhD thesis, U. do Minho.

2

So far...

- We've studied how to perform:
 - String matchingefficiently and effectively.
- We've seen how string matching is important in data integration
- Now, we'll see how string matching is important in **data cleaning**

3

Example (1)

Table R

Name	SSN	Addr
Jack Lemmon	430-871-8294	Maple St
Harrison Ford	292-918-2913	Culver Blvd
Tom Hanks	234-762-1234	Main St
...

Table S

Name	SSN	Addr
Tom Hanks	234-162-1234	Main Street
Kevin Spacey	-	Frost Blvd
Jack Lemon	430-817-8294	Maple Street
...

- Find records from different datasets that could be the same entity

Example (2)

```
<country>
  <name> United States of America </name>
  <cities> New York, Los Angeles, Chicago </
    cities>
  <lakes>
    <name> Lake Michigan </name>
  </lakes>
</country>
```

and

```
<country>
  United States
  <city> New York </city>
  <city> Los Angeles </city>
  <lakes>
    <lake> Lake Michigan </lake>
  </lakes>
</country>
```

are the same
object?

Example (3)

P. Bernstein, D. Chiu: Using Semi-Joins to Solve Relational Queries. JACM 28(1): 25-40(1981)

Philip A. Bernstein, Dah-Ming W. Chiu, Using Semi-Joins to Solve Relational Queries, Journal of the ACM (JACM), v.28 n.1, p.25-40, Jan. 1981

- These two bibliographic references concern the same publication!

The three examples refer to the same problem that is known under different names:

- ❑ approximate duplicate detection
- ❑ record linkage
- ❑ entity resolution
- ❑ merge-purge
- ❑ data matching ...

It is one of the data quality problems addressed by **data cleaning**

Outline

- Introduction to data cleaning
- Application contexts of data cleaning
- Data quality dimensions
- Taxonomy of data quality problems
- Data quality process
- Main data quality tools
- Real-world examples

Why Data Cleaning?

Data in the real world is **dirty**

incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data

- e.g., occupation=""

noisy: containing errors (spelling, phonetic and typing errors, word transpositions, multiple values in a single free-form field) or outliers

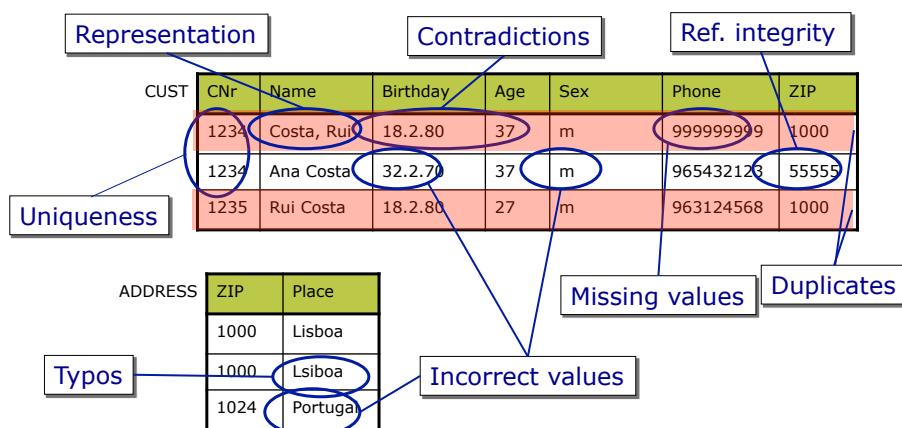
- e.g., Salary="-10"

inconsistent: containing discrepancies in codes or names (synonyms and nicknames, prefix and suffix variations, abbreviations, truncation and initials)

- e.g., Age="42" Birthday="03/07/1997"
- e.g., was rating "1,2,3", now rating "A, B, C"
- e.g., discrepancy between approximate duplicate records

9

Data Quality Problems (Dirty Data)



10

Impact of Data Quality Problems

- **Incorrect prices** in inventory retail databases [English 1999]
 - Costs for consumers 2.5 billion \$
 - 80% of barcode-scan-errors to the disadvantage of consumer
- **IRS 1992**: almost 100,000 tax refunds not deliverable [English 1999]
- 50% to 80% of computerized **criminal records in the U.S.** were found to be inaccurate, incomplete, or ambiguous. [Strong et al. 1997a]
- **US-Postal Service**: of 100,000 mass-mailings up to 7,000 undeliverable due to incorrect addresses [Pierce 2004]

**IRS might
be after you
— to mail
you a check**

Incorrect addresses
stall nearly 1,500
Tennessee refunds

By **BONNA de la CRUZ**
Staff Writer

Now that Tilcia L. Menifee knows that she'll be getting \$500 in a tax refund from Uncle Sam, she can do some Christmas shopping, she said.

Why Is Data Dirty?

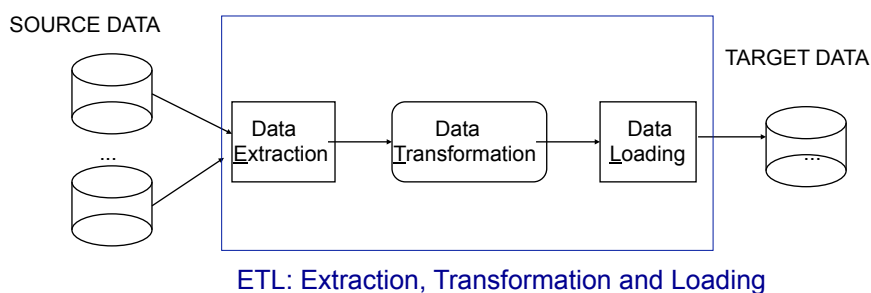
- **Incomplete data** comes from:
 - non available data value when collected
 - different criteria between the time when the data was collected and when it is analyzed
 - human/hardware/software problems
- **Noisy data** comes from:
 - data collection: faulty instruments
 - data entry: human or computer errors
 - data transmission
- **Inconsistent (and duplicate) data** comes from:
 - Different data sources, so non-uniform naming conventions/data codes
 - Functional dependency and/or referential integrity violation

Application contexts

- **Integrate data** from different sources
 - E.g., populating a DW from different operational data stores or a mediator-based architecture
- **Eliminate errors and duplicates** within a single source
 - E.g., duplicates in a file of customers
- **Migrate data** from a source schema into a different fixed target schema
 - E.g., discontinued application packages
- **Convert poorly structured data** into structured data
 - E.g., processing data collected from the Web

13

When materializing the integrated data (data warehousing)...



70% of the time in a data warehousing project is spent with the ETL process

14

Why is Data Cleaning Important?

Activity of converting source data into target data without errors, duplicates, and inconsistencies, i.e.,

Cleaning and Transforming to get...

High-quality data!

- No quality data, **no quality decisions!**
 - Quality decisions must be based on good quality data (e.g., duplicate or missing data may cause incorrect or even misleading statistics)

15

Quality

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

Robert Pirsig



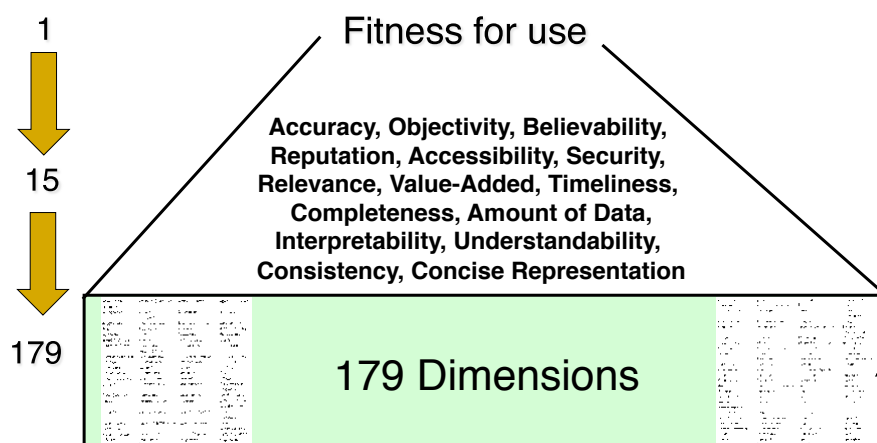
16

Outline

- Introduction to data cleaning
- Application contexts of data cleaning
- **Data quality dimensions**
- Taxonomy of data quality problems
- Data quality process
- Main data quality tools
- Real-world examples

17

What is Data of Good Quality?



Category	IQ Criteria	TDQM	MBIS	Weikum	DWQ	SCOUG	Chen
Content-related Criteria	Accuracy	Yes	Yes	Yes	Yes	Yes	Yes
	Documentation					Yes	
	Relevancy	Yes	Yes		Yes		Yes
	Value-Added	Yes				Yes	
	Completeness	Yes	Yes	Yes	Yes	Yes	Yes
Technical Criteria	Interpretability	Yes			Yes		
	Timeliness	Yes	Yes	Yes	Yes	Yes	Yes
	Reliability			Yes			
	Latency			Yes			Yes
	Performability			Yes		Yes	
	Response time		Yes	Yes			Yes
	Security	Yes		Yes	Yes		
	Accessibility	Yes	Yes	Yes	Yes	Yes	
Intellectual Criteria	Price		Yes	Yes		Yes	
	Customer Support					Yes	
	Believability	Yes	Yes	Yes	Yes	Yes	
Instantiation related Criteria	Reputation	Yes	Yes		Yes		
	Objectivity	Yes					
	Verifiability			Yes			
	Amount of data	Yes	Yes				Yes
	Understandability	Yes	Yes				
	Concise represent.	Yes					
	Consistent represent.	Yes	Yes	Yes	Yes	Yes	

Data Quality Dimensions (classical)

Accuracy

- Refers to the closeness of values in a database to the true values of the entities that the data in the database represent; if it is not 100% that means that there are errors in data

Example: "Jhn" vs. "John"

Completeness

- Concerns whether the database has complete information to answer queries
- Partial knowledge of the records in a table or of the attributes in a record

Currency

- Aims at identifying the current values of entities represented by tuples in a database and to answer queries using those values

Example: Residence (Permanent) Address: out-dated vs. up-to-dated

Consistency

- Refers to the validity and integrity of data representing real-world entities; if it is violated, leads to discrepancies and conflicts in the data

Example: ZIP Code and City inconsistent

Accuracy

- Closeness between a value v and a value v' , considered as the correct representation of the real-world phenomenon that v aims to represent.

- Ex: for a person name "John", $v' = \text{John}$ is correct, $v = \text{Jhn}$ is incorrect

Syntactic accuracy: closeness of a value v to the elements of the corresponding definition domain D

- Ex: if $v = \text{Jack}$, even if $v' = \text{John}$, v is considered syntactically correct, because it is an admissible value in the domain of people names.
- Measured by means of **comparison functions** (e.g., edit distance) that evaluate the distance between v and the values of the domain

Semantic accuracy: closeness of the value v to the true value v'

- Measured with a <yes, no> or <correct, not correct> domain
- Coincides with **correctness**
- The corresponding true value has to be known

21

Ganularity of accuracy definition

- Accuracy may refer to:
 - a single value of a relation attribute
 - an attribute or column
 - a relation
 - the whole database

22

Metrics for quantifying accuracy

- **Weak accuracy error**
 - Characterizes accuracy errors that do not affect identification of tuples
- **Strong accuracy error**
 - Characterizes accuracy errors that affect identification of tuples
- **Percentage of accurate tuples**
 - Characterizes the fraction of accurate tuples matched with a reference table

23

Completeness

- “The extent to which data are of sufficient breadth, depth, and scope for the task in hand.”
- Three types:
 - **Schema completeness**: degree to which concepts and their properties are not missing from the schema
 - **Column completeness**: evaluates the missing values for a specific property or column in a table.
 - **Population completeness**: evaluates missing values with respect to a reference population

24

Completeness of relational data

- The **completeness of a table** characterizes the extent to which the table represents the real world.
- Can be characterized with respect to:
 - **The presence/absence and meaning of null values**
Example: In Person(name, surname, birthdate, email), if email is null may indicate the person has no mail (no incompleteness), email exists but is not known (incompleteness), it is not known whether Person has an email (incompleteness may not be the case)
 - **Validity of open world assumption (OWA) or closed world assumption (CWA)**
 - **OWA**: assumes that in addition to missing values, some tuples representing real-world entities may also be missing
 - **CWA**: assumes the database has collected all the tuples representing real-world entities, but the values of some attributes in those tuples are possible missing

25

Metrics for quantifying completeness (1)

- **Model without null values with OWA**
 - Needs a **reference relation** $ref(r)$ for a relation r , that contains all the tuples that satisfy the schema of r

$$C(r) = |r| / |ref(r)|$$

Example: according to a registry of Lisbon municipality, the number of citizens is 2 million. If a company stores data about Lisbon citizens for the purpose of its business and that number is 1,400,000 then $C(r) = 0,7$

26

Metrics for quantifying completeness (2)

- **Model with null values with CWA:** specific definitions for different granularities:
 - **Values:** to capture the presence of null values for some fields of a tuple
 - **Tuple:** to characterize the completeness of a tuple wrt the values of all its fields:
 - Evaluates the % of specified values in the tuple wrt the total number of attributes of the tuple itself
- Example: `Student(stID, name, surname, vote, examdate)`
- Equal to 1 for (6754, Mike, Collins, 29, 7/17/2004)
- Equal to 0.8 for (6578, Julliane, Merrals, NULL, 7/17/2004)

27

Metrics for quantifying completeness (3)

- **Attribute:** to measure the number of null values of a specific attribute in a relation
 - Evaluates % of specified values in the column corresponding to the attribute wrt the total number of values that should have been specified.
- Example: For calculating the average of votes in `Student`, a notion of the completeness of `vote` should be useful
- **Relations:** to capture the presence of null values in the whole relation
 - Measures how much info is represented in the relation by evaluating the content of the info actually available wrt the maximum possible content, i.e., without null values.

28

Time-related dimensions

Currency: concerns how promptly data are updated

- Example: if the residential address of a person is updated (it corresponds to the address where the person lives) then the currency is high

Volatility: characterizes the frequency with which data vary in time

- Example: Birth dates (volatility zero) vs stock quotes (high degree of volatility)

Timeliness: expresses how current data are for the task in hand

- Example: The timetable for university courses can be current by containing the most recent data, but it cannot be timely if it is available only after the start of the classes.

29

Metrics of time-related dimensions

■ Last update metadata for currency

- Straightforward for data types that change with a fixed frequency

■ Length of time that data remain valid for volatility

■ Currency + check that data are available before the planned usage time for timeliness

30

Consistency

- Captures the **violation of semantic rules** defined over a set of data items, where data items can be tuples of relational tables or records in a file
 - **Integrity constraints** in relational data
 - Domain constraints, key definitions, inclusion and functional dependencies

31

Other dimensions

- **Interpretability**: concerns the documentation and metadata that are available to correctly interpret the meaning and properties of data sources
- **Synchronization** between different time series: concerns proper integration of data having different time stamps.
- **Accessibility**: measures the ability of the user to access the data from his/her own culture, physical status/functions, and technologies available.

32

Outline

- Introduction to data cleaning
- Application contexts of data cleaning
- Data quality dimensions
- **Taxonomy of data quality problems**
- Data quality process
- Main data quality tools
- Real-world examples

33

Taxonomy of data quality problems [Oliveira 2009]

- Value-level
- Value-set (attribute/column) level
- Record level
- Relation level
- Multiple relations level

34

Value level

Missing value: value not filled in a not null attribute

- Ex: birth date = ''

Syntax violation: value does not satisfy the syntax rule defined for the attribute

- Ex: zip code = 27655-175; syntactical rule: xxxx-xxx

Spelling error

- Ex: city = 'Lsboa', instead of 'Lisbon'

Domain violation: value does not belong to the valid domain set

- Ex: age = 240; age: {0, 120}

35

Value-set and Record levels

Value-set level

- **Existence of synonyms:** attribute takes different values, but with the same meaning
 - Ex: emprego = 'futebolista'; emprego = 'jogador futebol'
- **Existence of homonyms:** same word used with diff meanings
 - Ex: same name refers to different authors of a publication
- **Uniqueness violation:** unique attribute takes the same value more than once
 - Ex: two clients have the same ID number
- **Integrity constraint violation**
 - Ex: sum of the values of percent attribute is more than 100

Record level

- **Integrity constraint violation**
 - Ex: total price of a product is different from price plus taxes

36

Relation level

Heterogeneous data representations: different ways of representing the same real world entity

- Ex: name = 'John Smith'; name = 'Smith, John'

Functional dependency violation

- Ex: (2765-175, 'Estoril') and (2765-175, 'Oeiras')

Existence of approximate duplicates

- Ex: (1, André Fialho, 12634268) and (2, André Pereira Fialho, 12634268)

Integrity constraint violation

- Ex: sum of salaries is superior to the max established

37

Multiple tables level

Heterogeneous data representations

- Ex: one table stores meters, another stores inches

Existence of synonyms

Existence of homonyms

Different granularities: same real world entity represented with diff. granularity levels

- Ex: age: {0-30, 31-60, > 60}; age: {0-25, 26-40, 40-65, >65}

Referential integrity violation

Existence of approximate duplicates

Integrity constraint violation

38

Outline

- Introduction to data cleaning
- Application contexts of data cleaning
- Data quality dimensions
- Taxonomy of data quality problems
- **Data quality process**
 - Main data quality tools
 - Real-world examples

39

Data Quality Process

1. **Data Quality Auditing (Assessment)**
 - Data Profiling
 - Data Analysis
2. **Data Quality Improvement**
 - Data Cleaning
 - Data Enrichment

40

Data quality auditing

- Constituted by:
 - **Data profiling** – analysing data sources to identify data quality problems
 - **Data analysis** – statistical evaluation, logical study and application of data mining algorithms to define data patterns and rules
- **Main goals:**
 - To obtain a definition of the data: **metadata** collection
 - To check violations to metadata definition
 - To detect other data quality problems that belong to a given taxonomy
 - To supply recommendations in what concerns the data cleaning task

41

Data Profiling

- **Data source discovery**
 - Metadata
- **Schema discovery**
 - Schema matching and mapping
 - Profiling for metadata (keys, foreign keys, data types, ...)
- **Data discovery**
 - Column-level: Null-values, domains, patterns, value distributions / histograms
 - Table-level: Data mining, rules

Typical techniques used in data quality auditing

- **Dictionaries of words:** so that attribute values are compared with one or more dictionaries of the domain
 - Ex: wordnet
- **Algorithms to detect functional dependencies and their violations**
- **Algorithms to detect duplicates**
 - String matching for string fields
 - Character-based
 - Token-based
 - Phonetic algorithms
 - Record matching
 - Rule-based
 - Probabilistic
 - ...

Nome	Cod.Postal	Localidade
Maria	2765	Estoril
António	2765	S.João Estoril
José	2780	Oeiras
Andreia	1000	Lisboa
Manuela	2865	Setúbal


Localidade=>Cod.Postal

43

Data quality improvement

- Includes often:
 - **Data transformation** – set of operations that source data must undergo to fit target schema
 - **Data cleaning**– detecting, removing and correcting dirty data (including **approximate duplicate elimination**)
 - **Data enrichment**– use of additional information to improve data quality
- **Main goal:**
 - To **correct** the data quality problems detected during the data quality auditing process

44

Typical techniques used in data cleaning and transformation

- Dictionaries of words
- Libraries of pre-defined cleaning functions
- Machine learning techniques
- Techniques for consolidating approximate duplicates

45

Methodology for data cleaning

1. Extraction of the individual fields that are relevant
2. Standardization of record fields
3. Correction of data quality problems at value level
 - Missing values, syntax violation, etc
4. Correction of data quality problems at value-set level and record level
 - Synonyms, homonyms, uniqueness violation, integrity constraint violation, etc
5. Correction of data quality problems at relation level
 - Violation of functional dependencies, duplicate elimination, etc
6. Correction of data quality problems problems at multiple relations level
 - Referential integrity violation, duplicate elimination, etc
- User feedback
 - To solve instances of data quality problems not addressed by automatic methods
- Effectiveness of the data cleaning and transformation process must be always measured for a sample of the data set

46

Data Cleaning Tasks

1. Extraction from sources
 - Technical and syntactic obstacles
2. Transformation
 - Schematic obstacles
3. Standardization
 - Syntactic and semantic obstacles
4. Duplicate detection
 - Similarity functions
 - Algorithms
5. Data fusion / consolidation
 - Semantic obstacles
6. Loading into warehouse / presenting to user

47

Human Interaction is Needed

- Components to implement
 - Wrappers for technical heterogeneity
 - Schema integration based on correspondences
 - Similarity measure for schema elements
 - Similarity measure for records
- Knobs to turn
 - Thresholds for similarity measures
 - Partition size / window size
- Expert guidance
 - Rule selection / rule specification
 - Schema matching
 - Duplicate detection
 - Data fusion

48

Outline

- Introduction to data cleaning
- Application contexts of data cleaning
- Data quality dimensions
- Taxonomy of data quality problems
- Data quality process
- **Main data quality tools**
- Real-world examples

49

Existing technology for ensuring data quality

Ad-hoc programs written in a programming language like C or Java or using an RDBMS proprietary language

- Programs difficult to optimize and maintain

RDBMS mechanisms for guaranteeing integrity constraints

- Do not address important data instance problems

Data transformation workflow scripts using a **data cleaning/profiling tool**

50

Existing technology for ensuring data quality

Ad-hoc programs written in a programming language like C or Java or using an RDBMS proprietary language

- Programs difficult to optimize and maintain

RDBMS mechanisms for guaranteeing integrity constraints

- Do not address important data instance problems

➤ **Data transformation workflow scripts using an data cleaning/profiling tool**

51

Criteria for comparing commercial data quality tools (1)

Debugger:

Data lineage: data lineage or provenance identifies the set of source data items that produced a given data item

Breakpoints: breakpoints is an intentional stopping or pausing place in a cleaning program put in place for debugging purposes

Edit values: the user can edit values during debugging

52

Criteria for comparing commercial data quality tools (2)

Profiling:

Rules: A rule is a business logic that defines conditions applied to data. They are used to validate the data and to measure data quality

Filters: A filter is used to split the data tuples in different groups. Each group should be validated by a different set of rules.

53

Criteria for comparing commercial data quality tools (3)

Execution:

User involvement: Support for user interaction in a data cleaning process

Incremental updates: The ability to incrementally update data targets, instead of rebuilding them from scratch every time

54

Commercial Data Cleaning Tools(2014) (1/3)

Tools	Debugger			Profiling		Execution	
	Data lineage	Breakpoints	Edit values	Rules	Filters	User involvement	Incremental updates
Informatica PowerCenter	Y	Y	Y	Y	Y	N	Y
IBM Information Server	Y	Y	N	Y	Y	N	Y
Talend Open Studio	N	Y	N	Y	Y	N	Y
Oracle Data Integrator	Y	Y	N	Y	Y	N	Y
SQL Server Integration Services	Y	Y	N	Y	N	N	Y
SAS Data Integration Studio	Y	N	N	Y	Y	N	Y
Pentaho Data Integration	N	N	N	Y	N	N	Y
Clover ETL	N	N	Y	Y	Y	N	Y

55

Criteria for comparing commercial data quality tools (4)

Extensibility:

Create operators: the user can define new operators

Modify operators: the user can modify standard operators

User Interface:

Drag and drop: the user can define data quality processes using a drag and drop interface

Editor: the user can define and edit data quality processes modeled as workflows using a graphical interface

56

Commercial Data cleaning tools (2014) (2/3)

Tools	Extensibility		User Interface	
	Create Operators	Modify Operators	Drag and Drop	Grahical Editor
Informatica PowerCenter	Y (Java)	N	Y	Y
IBM Information Server	Y (Java)	N	Y	Y
Talend Open Studio	Y (Java, Groovy)	Y (Java)	Y	Y
Oracle Data Integrator	Y	Y	Y	Y
SQL Server Integration Services	Y (C#, VB)	N	Y	Y
SAS Data Integration Studio	Y (SAS)	Y (SAS)	Y	Y
Pentaho Data Integration	Y (Javascript)	N	Y	Y
Clover ETL	Y (CTL)	N	Y	Y

57

Criteria for comparing commercial data quality tools (5)

Scalability:

Grid: the tool can run a cleaning process on a collection of computer resources from multiple locations

Partitioning: the user can partition the data and run each partition independently (on different CPUs or cores)

Pushdown optimization: the tool translates the transformation logic into SQL queries and sends the SQL queries to the database. The database engine executes the SQL queries to process the transformations

Others:

Free version: the tool has a free version

58

Commercial Data Cleaning Tools (2014) (3/3)

Tools	Scalability			Others
	Grid	Partitioning	Pushdown Optimization	Free version
Informatica PowerCenter	Y	Y	Y	Y
IBM Information Server	Y	Y	Y	N
Talend Open Studio	Y	N	Optional ELT	Y
Oracle Data Integrator	Y	N	ELT	Y
SQL Server Integration Services	N	Y	-	Y (IST)
SAS Data Integration Studio	Y	Y	Y	Y (IST)
Pentaho Data Integration	Y	Y	N	Y
Clover ETL	Y	Y	N	Y

59

Research Data cleaning tools (2014) (1/2)

Tools	Detection DQ problems		Repair DQ problems		
	Constraints	Satistical	Search	ML/St	Data Transformations
Cleenex	QCs	N	N	N	Y
Llunatic	Egds	N	Y	N	N
Nadeef	CFDs, MDs	N	Y	N	N
Guided data repair	CFDs	N	Y	Y	N
Scare	N	Y	N	Y	N
Eracer	N	Y	N	Y	N
Continuous data cleaning	FDs	N	Y	Y	N

60

Criteria for comparing research data cleaning tools (1)

Detection:

Constraints – use of rules or/and conditions

- EGDs - equality generating dependencies
- QCs - quality constraints
- CFDs - Conditional functional dependencies
- MDs - Matching dependencies

Statistical – dirty tuples are detected based on simple statistics or in complex data analysis

61

Criteria for comparing research data cleaning tools (1)

Repair:

Search: The system explores the space of possible clean tables and heuristically selects the best table

ML/St: The system uses machine learning and/or statistical models to infer data values or to prune the search

Data transformations: The system models the data cleaning process as a data transformation graph

62

Criteria for comparing research data cleaning tools (3)

User Interface:

Graphical interface: the system provides a visualizing tool and menus to interact

User edition: the system allows the user to edit data values

Others:

Scalability: the system execution time grows linearly with the number of input tuples

Streaming: the system receives tuples and processes each of them treat them individually (opposed to batch processing)

Extensible: the system allows the user to modify and/or insert new algorithms

63

Research Data cleaning tools (2014) (1/2)

Tools	User Interface		Others		
	Graphical Interface	User edition	Extensible	Streaming	Scalability
Cleenex	Y	Y	Matching algorithms	N	N
Llunatic	Y	Y	Cost Managers	N	Y
Nadeef	Y	N	Repair algorithms	N	N
Guided data repair	N	Y	N	N	N
Scare	N	N	N	N	Y
Eracer	N	N	N	N	N
Continuous data cleaning	N	N	N	Y	Y

64

Outline

- Introduction to data cleaning
- Application contexts of data cleaning
- Data quality dimensions
- Taxonomy of data quality problems
- Data quality process
- Main data quality tools
- **Real-world examples**

65

Death by Typo

‘Resurrected,’ but still wallowing in red tape

Government records incorrectly kill off thousands, and there’s no easy fix

By Alex Johnson and Nancy Amons

Reporters

MSNBC and NBC News

updated 6:21 p.m. ET Feb. 29, 2008

For a dead woman, Laura Todd is awfully articulate.

“I don’t think people realize how difficult it is to be dead when you’re not,” said Todd, who is very much alive and kicking in Nashville, Tenn., even though the federal government has said otherwise for many years.

Todd’s struggle started eight years ago with a typo in government records. The government has reassured her numerous times that it has cleared up the confusion, but the problems keep coming.

[Story continues below ↓](#)

Video



[Launch](#)

Does this woman look dead to you?

The government says Toni Anderson is dead, but she insists she is very much alive. David MacAnally of NBC affiliate WTHR reports from Muncie, Ind.

NBC News Channel

66

Google searches for Britney Spears

488941 britney spears	29 brlntey spears	9 britnatty spears	5 bmeiy spears	3 britny spears	2 brirenyy spears
40134 britanny spears	29 brittany spears	9 britatny spears	5 brotnay spears	3 britmny spears	2 britrny spears
36315 brittney spears	29 brittany spears	9 britatny spears	5 brotny spears	3 britneey spears	2 britrtyny spears
24342 britany spears	29 britny spears	9 brltin spears	5 brulteny spears	3 britrny spears	2 britrtney spears
7331 britny spears	26 brittney spears	9 britnew spears	5 btinye spears	3 britrly spears	2 britain spears
6633 britney spears	26 brittney spears	9 britney spears	5 brittney spears	3 britny spears	2 britane spears
2696 brittney spears	26 brlnly spears	9 britny spears	5 gritny spears	3 britrnty spears	2 britanay spears
1807 briney spears	26 brittany spears	9 britny spears	5 spritney spears	3 britnx spears	2 britania spears
1635 brittney spears	26 britney spears	9 brittany spears	4 britlay spears	3 britnyxxx spears	2 britann spears
1479 britney spears	26 britny spears	9 brittany spears	4 britny spears	3 britnnay spears	2 britana spears
1479 brittany spears	26 brittney spears	9 byrtn spears	4 brandy spears	3 britrntey spears	2 britannie spears
1338 britny spears	26 brittany spears	9 rbttney spears	4 brrbtney spears	3 britrntys spears	2 britannt spears
1211 britnet spears	24 beltney spears	8 britny spears	4 breatiny spears	3 britrntty spears	2 britannu spears
1096 brittney spears	24 birtley spears	8 bittney spears	4 breeetny spears	3 britrnttey spears	2 britantj spears
991 britaney spears	24 brightney spears	8 brattany spears	4 bretiny spears	3 britrntty spears	2 britantj spears
991 britany spears	24 brittany spears	8 breitny spears	4 britfny spears	3 britrntty spears	2 britetny spears
811 brittney spears	24 britanty spears	8 breteny spears	4 briattany spears	3 britryen spears	2 britetany spears
811 britney spears	24 britenny spears	8 brighity spears	4 brieteny spears	3 britrntye spears	2 britenet spears
664 britney spears	24 brtlil spears	8 brintay spears	4 briety spears	3 britrny spears	2 briteny spears
664 brittney spears	24 brtliny spears	8 brittney spears	4 britlay spears	3 brotney spears	2 britenys spears
604 britney spears	24 brittany spears	8 britny spears	4 brittlay spears	3 brittlay spears	2 brittany spears
601 britny spears	24 brittily spears	8 britvny spears	4 brinie spears	3 brittlay spears	2 britin spears
601 britny spears	21 brittney spears	8 brittley spears	4 brittney spears	3 brittlay spears	2 brittany spears
544 brittany spears	21 brittlay spears	8 britneb spears	4 brittne spears	3 brittney spears	2 britmy spears
544 brittany spears	21 bitely spears	8 brittney spears	4 britaby spears	3 brittany spears	2 britmat spears
364 britey spears	21 bratney spears	8 brtity spears	4 britay spears	3 brittany spears	2 britnat spears
364 brittney spears	21 britani spears	8 brittner spears	4 britayne spears	3 brtnet spears	2 britnbey spears
329 britny spears	21 britanie spears	8 brottany spears	4 britlie spears	3 birtly spears	2 britndy spears
289 britney spears	21 britany spears	7 baritney spears	4 brittney spears	3 bity spears	2 britnh spears
289 britneys spears	21 brittay spears	7 birtney spears	4 brittmy spears	3 brittney spears	2 britnety spears
244 britne spears	21 brittany spears	7 bitney spears	4 brittmy spears	3 pretney spears	2 britney6 spears
244 brytney spears	21 brittany spears	7 bity spears	4 britnel spears	3 brittney spears	2 britneye spears
220 brittney spears	21 brittany spears	7 brittany spears	4 britney spears	3 brittney spears	2 brittneh spears
120 britney spears	19 britny spears	7 briantly spears	4 brittney spears	3 britrbtney spears	2 brittymn spears
120 brittney spears	19 brittany spears	7 birttney spears	4 brittmy spears	2 blbitney spears	2 brittvyvy spears
163 britny spears	19 brittany spears	7 brittany spears	4 brittby spears	2 brittyny spears	2 brittney spears
147 britney spears	19 britney spears	7 britty spears	4 brittby spears	2 brittyny spears	2 brittney spears
147 brittney spears	19 brittany spears	7 britny spears	4 brittney spears	2 brittyny spears	2 brittney spears
147 brittany spears	19 brittany spears	7 britneu spears	4 brittney spears	2 brittyny spears	2 brittney spears

Directmarketing by The Economist

nomist
Tel. 030 713 91 91
Dr. Felix Naumann
72 A R-Breitscheid-Str
Potsdam
14482
GERMANY

If undelivered please return to:
BTB Mailflight Wolsley Road Kempton Beds MK2 7UA

2 ROYAL MAIL
POSTAGE PAID GB
Q 550 000

POSTAGE PAID GB

2 ROYAL MAIL
POSTAGE PAID GB
HQ 5878

POSTAGE PAID GB
HQ 5878

QWMX0071362
Felix Naumann
Rudolf-Breitscheid-Str 72A
Potsdam
14482
GERMANY

If undelivered please return to:
RTB MailPost Wireless Road Kingston, Rhode 02881-2114

111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052 1053 1054 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076 1077 1078 1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096 1097 1098 1099 1100 1101 1102 1103 1104 1105 1106 110

QWMQ0071368
Dr Felix Naumann
72 A R.-Breitscheid-Str
Potsdam
14482
GERMANY

|||||

QWMX0071362
Felix Naumann
Rudolf-Breitscheid-Str 72A
Potsdam
14482
GERMANY

FIFA registration form (2010)

The image displays several screenshots of the FIFA registration form (2010) interface. The main form is on the left, and three smaller dropdown menus are shown on the right. The main form includes fields for Nationality, Country of Residence, Mother Tongue, Preferred FIFA Language, Secondary FIFA Language, Organisation Name, Organisation Role (Prof), and Notes (Max 2000 chars). The dropdown menus show lists of countries and regions, including German Democratic Republic, Germany, Germany Federal Republic, Ghana, Gibraltar, Great Britain, All Ireland (all-Ireland pre 1921), American Samoa, Andorra, Angola, Wales, Yemen, Yemen PDR, Yugoslavia, Zaire, Zambia, and Zimbabwe.

69

German Umlaute

dblp.uni-trier.de

Search Results for 'dessloch'

- ♦ [Stefan Deßloch](#)
- ♦ [Stefan Dessloch](#)

DBLP: [[Home](#)] Search: [Author](#), [Title](#) | [Conferences](#) | [Journals](#)

Michael Lev (lev@uni-trier.de) Thu Jan 31 10:44:06 2008

70

Next lecture

- Data Matching

Follow me on [LinkedIn](#) for more :
Steve Nouri
<https://www.linkedin.com/in/stevenouri/>
