

Data integration

Why database integration?

- Combine database when two companies merge
- enhance information using data from different sources
- access data in legacy databases

Related Issues

- Data linkage
 - Identifying records referring to the same real-world entity
 - Computing data similarity
 - Applicability of different similarities
- Info fusion
 - Extract info from different data source
- Data cleaning
 - Remove noise in original data
 - Remove noise in integrated data

inconsistency and redundancy
- Data quality
 - Data augmentaion
 - Data constraints
 - Data provenance
- Data privacy
 - Share data with the assurance that private info can not be derived

Difference between Data integration and Info fusion

- Common goal: Integrate and organize data from multiple sources in order to present a unified view of data to derive actionable insights
- Data integration: focuses on combining data to create a bigger and consistent data set
- Info fusion: focuses on deriving insight from real-time streaming data with semantic context from other big data source

Global Info systems

- Federated database(FDB)
 - Semi-autonomous database systems
 - global view is provided
- Multi-databases(MDB)
 - Autonomous database systems, no or limited global view is provided
- Interoperable info systems
 - Definition
 - Only API provided to communicate with different database
 - No global virtual view
 - Interoperability
 - Ability for an application to access multiple distinct systems
 - Not necessarily for database only
 - Interoperable systems
 - Exchange messages and requests
 - Receive services and operate as a unit towards a common goal
 - Different types of interoperability
 - Syntactic
 - languages and data formats
 - Semantic
 - meaning
 - System
 - machines, networks, database models
 - Structural
 - data structures and data models

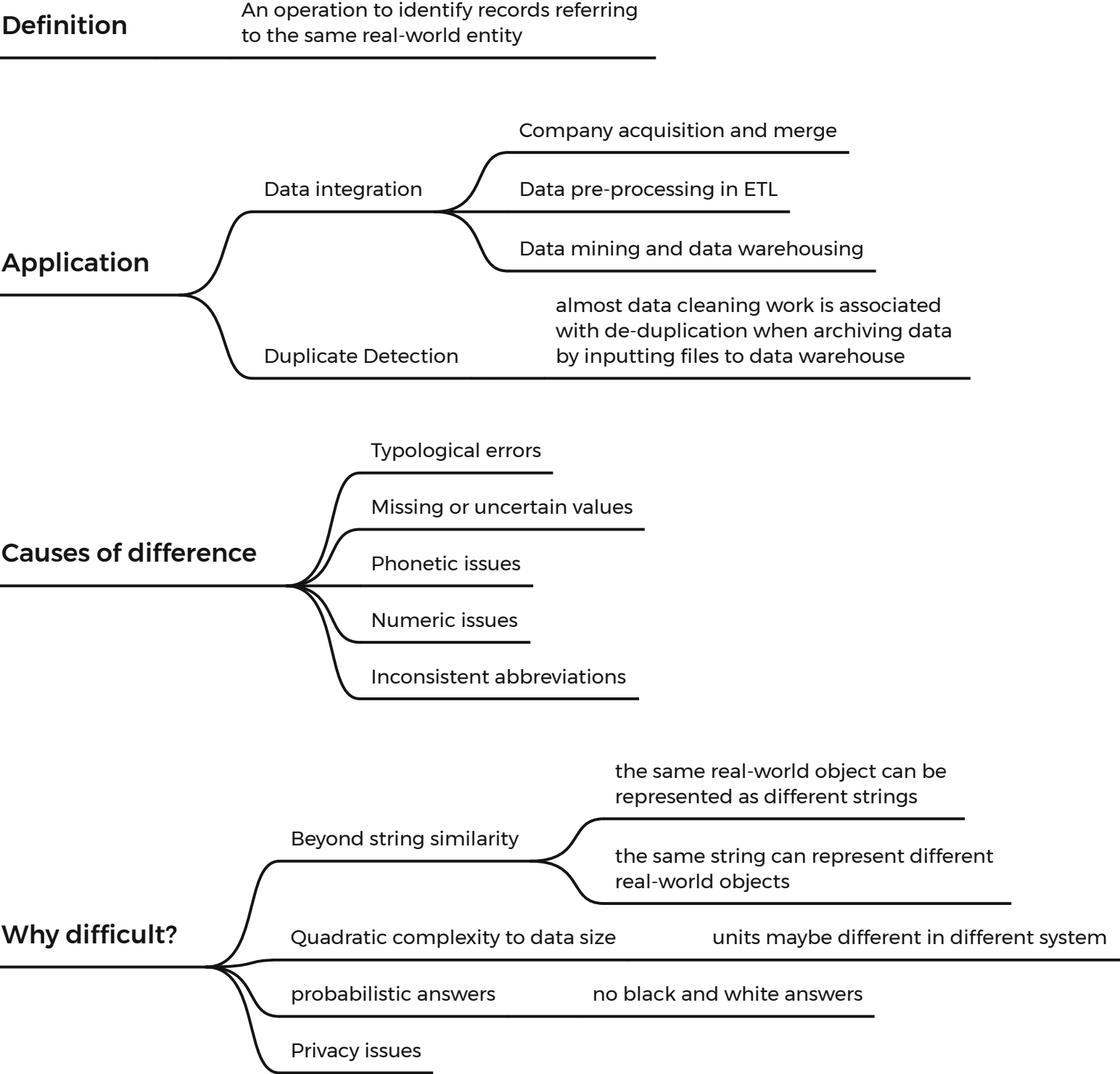
Challenges in DB integration

- Schema heterogeneity
- Type heterogeneity
- Value heterogeneity
- Semantic heterogeneity

Process of DB Integration

- Schema mapping
 - mapping of structures
- Data mapping
 - matching based on content
- Data fusion
 - reconciliation of mismatching content

Data linkage



Data quality

Governance

What the organisation does and what it should become in the future

Management

How the organisation will reach those goals and aspirations

solutions

organizational

architectural

computational

defining and enforcing data constraints

Doing entity resolution effectively and efficiently

managing missing values and uncertain data

Data quality dimensions

Integrity

Meaningless

Accuracy

Erroneous

the closeness between a value v and value v' considered as the correct representation

Completeness

Missing

the sufficiency of data for the task at hand

Currency

Obsolete, out of date

How promptly data is updated

Volatility: frequency with which data vary in time

Timeliness: how current the data are for the task at hand

Representational consistency

inconsistent

Accessibility

unavailable

Uncertainty

reliability, trust

Recognize the problem

Data quality problem

Data acquisition

Data integration

Data utilization

Measure its cost

How many errors?

Data quality cost

Costs caused by low Data quality

Direct costs

verification costs

Re-entry costs

Compensation costs

Indirect costs

Costs based on lower reputation

Costs based on wrong Decisions or actions

Sunk investment costs

Costs of improving or assuring data quality

Prevention cost

training costs

monitoring costs

standard development and deployment cost

Detection costs

Analysis costs

reporting costs

Repair costs

Repair planning costs

Repair implementation costs

Basic steps for governance

Devise strategy for improvement

Short term

Define metrics and its relationship to business impact to identify which data to improve

Produce baseline

Use same metric to measure change from baseline

Sustain improvements through ongoing monitoring

long term

Establish process owner and management team

describe process qualitatively and understand requirements

Establish measurement system

Establish process control and conformance to requirements

Identify improvement opportunities

Select opportunities and set objectives

make and sustain improvements

Aim towards the governance maturity

Data warehouse

What is data warehouse

Definition of DW

A storage architecture designed to hold data extracted from transaction systems, operational data stores and external sources

data are combined in an aggregate, summary form suitable for enterprise-wide data analysis and reporting for predefined business needs.

data are periodically updated and stored for read-only applications

A data warehouse is a set of facts perceived by a number of dimensions

Characteristic of DW

Subject oriented

Nonvolatile

not usually subject tot changes

integrated

data is consistent and integrated from multiple courses

Time variant

historical data is recorded for analytical applications

Why we need data warehouse

Traditional database issues

Traditional database applications consist of both updates and queries, some queries are large scale aggregation reports which can take long time to generate on-the-fly

Database updates and queries must lock data resources, large scale aggregation reports lock many resources for a long time

Benefits of data warehouse

Organizations are analyzing current and historical data to identify useful patterns and support business strategies

Emphasis is on complex,interactive, exploratory analysis of very large datasets created by integrating data from across all parts of an enterprise

Difference between DW and DB

Integrated data spanning long time periods, often augmented with summary information

very large volume

Interactive response times expected for complex queries

Ad-hoc updated uncommon(write-once and Read forever)

DW design

Star schema

advantage

faster query processing speed

Snowflake schema

fact constellation

A set of fact tables that share some dimension tables

advantage

less data integrity problem

less space consumption

DW implementation

Data cube

OLAP

Pivot

Rotate data cube to show a different orientation of axes

Roll up

Move up concept hierarchy, grouping into larger units along a dimension with generalization

Drill-down

Disaggregate to a finer-grained view to show more details

Slice and dice

Perform projection operations on the dimensions

Cuboid

A cuboid is a denotation of one of the 2^d summarized views, which can be used for materialized view

Advantages

OLAP queries are typically aggregate queries

Pre-aggregation is essential for interactive response time

Pre-calculate expensive joins

Speed up online OLAP queries

Disadvantages

It increases storage cost

The content of the materialized views must be maintained when the underlying details tables are modified

Trade off between query performance and accessibility to up-to-date data

Database application

Big data application

- Connecting Dots — From small to big
- Discovering specifics — From big to small — How do we find outliers, predict trends, provide summaries, and give explanations from a dataset.
- Inferencing — Knowing unknown — Everything is related to everything else
- Key values — The smallest unit of big data
 - Every key-value is unique with timestamp
 - Key values are associated to each other
 - All key values can be drawn as a

Network science

- The curse of dimensionality — Adding extra dimension to a data space will exponentially increase the volume of data space
- Scale-free Network — the characteristic of network are independent of the size of network
- Different type of networks
- Computing issues of complex networks

Explaining outliers in aggregate queries

Effective storage of big data

- Column-based VS Row-based data database stroages
- Applicability of row and column storages
- Compare and contrast of row and column storages
 - Row storage
 - Efficient when many columns of a single row are required at the same time
 - Well-suited for OLTP-like workloads which are more heavily loaded with interactive transactions
 - Column storage
 - Efficient on aggregation operation over many row but only for a smaller subset of all columns of data
 - Efficient when inserting new values of a column for all rows at once
 - Well-suited for OLAP-like workloads
- NoSQL(not only SQL)
 - Technology
 - Cloud Platform — A viable alternative to relational databases operating on cluster servers
 - No schema — Different types of data is collected, stored, accessed without a schema
 - Data fusion — A data integration technique for multi-source data
 - Flexible access — A query model accessing data without using traditional SQL
 - Property
 - To add/delete/query massive arrays and still allow for persistence and fault tolerance
 - To store objects using key-values
 - To implement large data query on MapReduce framework
 - System
 - CouchDB — A document-oriented database that can be queried and indexed in a MapReduce fashion using JavaScript.
 - MongoDB — A scalable, high-performance, open source, document-oriented database system
 - Hadoop — Hadoop develops open-source software for reliable, scalable, distributed computing
 - HBase — Hadoop database system supports random, real-time, read/write access to big data.

Distributed database

