

Assignment Project Exam Help

Data Integration — 1

<https://powcoder.com>

Faculty of Information Technology, Monash University, Australia

Add WeChat FIT5196 week 10 powcoder

Assignment Project Exam Help

1 Recap

2 Data Integration

- Definition & Application
- Schema Integration

<https://powcoder.com>

Add WeChat powcoder

3 Summary

Outliers: the definition

- Types of outliers
 - ▶ Univariate outlier
 - ▶ Multivariate outlier

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

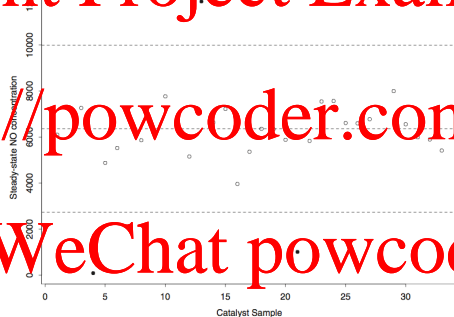
Outliers: Recap

- Types of outliers
 - ▶ Univariate outlier

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



- ▶ Multivariate outlier

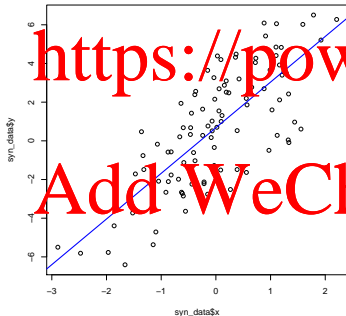
Outliers: Recap

- Types of outliers

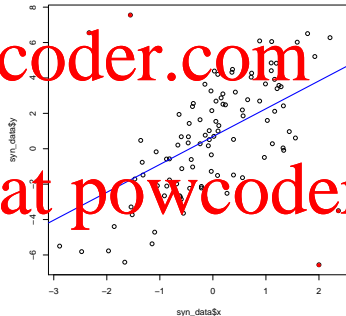
- ▶ Univariate outlier

- ▶ Multivariate outlier: for example bivariate outlier

Without Outliers



With Outliers



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Outliers: Recap

- Univariate outlier detection

- ▶ The 3σ edit rule
- ▶ The Hampel identifier
- ▶ The standard boxplot outlier rule

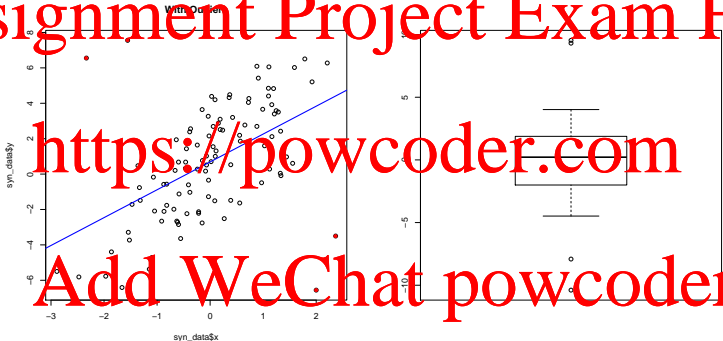
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

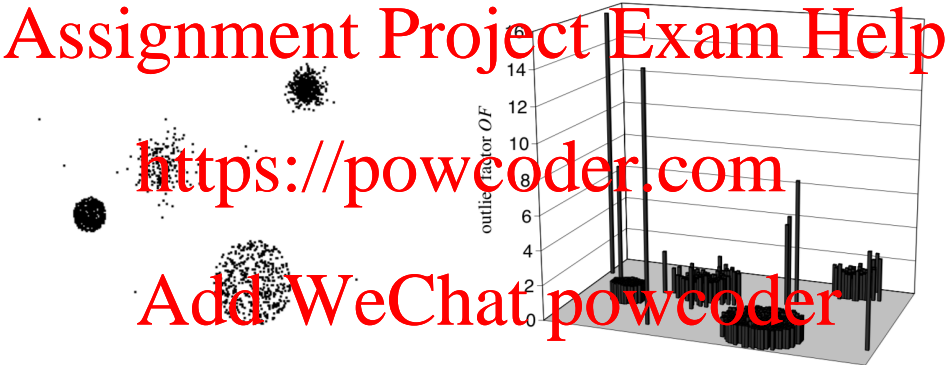
Outliers: Recap

- Multivariate outlier detection based on linear regression



Outliers: Recap

- Multivariate outlier detection based on LOF

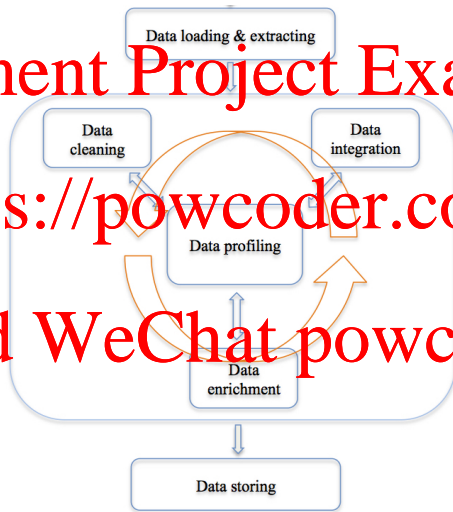


Data Wrangling Process

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Outline

Assignment Project Exam Help

1 Recap

2 Data Integration

- Definition & Application
- Schema Integration

3 Summary

<https://powcoder.com>
Add WeChat powcoder

Data Integration: Definition

- A process in which heterogeneous data is retrieved and combined as an incorporated form and structure.

- ▶ FullServe: an American-based ISP provide Internet access to homes,
- ▶ EuroCard: an ISP acquired by FullServe

Employee Database

FullTimeEmps(ssn, emplID, firstName,
middleName, lastName)
Hire(mplID, hireDate, recruiter)
TempEmployees(ssn, hireStart,
hireEnd, name, hourlyRate)

Resume Database

Interviews(interviewDate, plID, recruiter,
hireDecision, hireDate)
CVs(plID, resume)

Training Database

Courses(courseID, name, instructor)
Enrollments(courseID, emplID, date)

Services Database

Services(packName, textDescription)
Customers(name, EZipCode, streetAddr,
phone)
Contracts(custID, packName, startDate)

Sales Database

Products(prodName, prodID)
Sales(prodID, customerID,
custName, address)

HelpLine Database

Calls(date, agent, custID, text, action)

FIGURE 1.1 Some of the databases a company like FullServe may have. For each database, we show some of the tables and for each table, some of its attributes. For example, the Employee database has a table FullTimeEmps with attributes ssn, emplID, firstName, middleName, and lastName.

Data Integration: Definition

- A process in which heterogeneous data is retrieved and combined as an incorporated form and structure.

Assignment Project Exam Help

Employee Database

EmpID, firstNum, MiddleInitial,
last Name, salary)

Hire(ID, hireDate, recruiter)

Resume Database

interviews, ID, date, location,
recruiter, salary)

CVs(candID, resume)

Credit Card Database

Cards(CustID, cardNum,
expiration, currentBalance)

Customers(CustID, name,
address)

HelpLine Database

Calls(date, agent, custID,
description, followup)

FIGURE 1.2 Some of the databases of EuroCard. Note that EuroCard organizes its data quite differently from FullServe. For example, EuroCard does not distinguish between full-time and part-time employees. FullServe records the hire data of employees in the Resume database and the Employee database, while EuroCard only records the hire date in the Employee database.

Data Integration: Definition

- Data fusion: the integration of data and knowledge from several sources

- ▶ The Human Resources Department needs to be able to query for all of its employees, whether in the United States or in Europe.
- ▶ Combining data from the HelpLine database and the Sales database will help FullServe identify issues in their products and services early on,.
- ▶ Customer support hotline: a customer representative might need to know customers' Internet service, products purchased, call notes.

Assignment Project Exam Help
<https://powcoder.com>

Add WeChat powcoder

Data Integration: Definition

- Goal: create a single representation that provides a more accurate description than any of the individual data sources

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Data Integration: Application

Google's knowledge graph

Assignment Project Exam Help



Data Integration: Application

Map mashup: HousingMaps

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Map mashup: TrendMaps shows the latest trend in twitter.

<https://p0werc0der.com>

Add WeChat powcoder

Data Integration: Application

Product/service comparison portals: TheTracktor, PriceGrabber.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Why Data Integration is Challenging

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Why Data Integration is Challenging

- Heterogeneous data

- ▶ Data coming from different sources is often developed independently (e.g., different schema, different objectives)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Why Data Integration is Challenging

- Heterogeneous data
 - ▶ Data coming from different sources is often developed independently (e.g., different schema, different objectives)
- Various formats
 - ▶ Text, web logs, social networks, sensors, astronomy, genomics, medical records, surveillance, etc.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Why Data Integration is Challenging

- Heterogeneous data
 - ▶ Data coming from different sources is often developed independently (e.g., different schema, different objectives)
- Various formats
 - ▶ Text, web logs, social networks, sensors, astronomy, genomics, medical records, surveillance, etc.
- Incompatible Taxonomies
 - ▶ Different object identity and separate schema
 - Different definitions of a customer, an account, etc.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Why Data Integration is Challenging

- Heterogeneous data
 - ▶ Data coming from different sources is often developed independently (e.g., different schema, different objectives)
- Various formats
 - ▶ Text, web logs, social networks, sensors, astronomy, genomics, medical records, surveillance, etc.
- Incompatible Taxonomies
 - ▶ Different object identity and separate schema
 - Different definitions of a customer, an account, etc.
- Time synchronisation
 - ▶ Each source might have a time window that is different from each other.
 - ▶ Synchronisation of data collected in different time windows

Why Data Integration is Challenging

- Dealing with legacy data

- ▶ Historical data stored in legacy form, such as IMS, spreadsheets, and other ad hoc structure
- ▶ Combine historical data with modern data

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Why Data Integration is Challenging

- Dealing with legacy data

- ▶ Historical data stored in legacy form, such as IMS, spreadsheets, and other ad hoc structure
- ▶ Combine historical data with modern data

- Abstraction levels

- ▶ Different data sources might provide data at different level of abstraction, e.g.,
 - suburb level v.s. state level
 - annual v.s. weekly

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Why Data Integration is Challenging

- Dealing with legacy data

- ▶ Historical data stored in legacy form, such as IMS, spreadsheets, and other ad hoc structure
- ▶ Combine historical data with modern data

- Abstraction levels

- ▶ Different data sources might provide data at different level of abstraction, e.g.,
 - suburb level v.s. state level
 - annual v.s. weekly

- Data Quality

- ▶ Data is often erroneous, and combining data often aggravates the problems. Erroneous data has potentially devastating impact on the overall quality of the integrated data.

Why Data Integration is Challenging

- Dealing with legacy data

- ▶ Historical data stored in legacy form, such as IMS, spreadsheets, and other ad hoc structure
- ▶ Combine historical data with modern data

- Abstraction levels

- ▶ Different data sources might provide data at different level of abstraction, e.g.,
 - suburb level v.s. state level
 - annual v.s. weekly

- Data Quality

- ▶ Data is often erroneous, and combining data often aggravates the problems. Erroneous data has potentially devastating impact on the overall quality of the integrated data.

- The number of sources

- ▶ e.g., web-scale integration.

Where can we get the Data?

- Government and political data
- Social data
- Web resources
- Weather Data
- News data
- Preprocessed data for data analysis tasks
- Economics and Fiances

Method: Web scraping

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Integration Process

Assignment Project Exam Help

Integration Process

<https://powcoder.com>

Add WeChat powcoder

Schema
Integration

Data-level
Integration

Outline

Assignment Project Exam Help

1 Recap

2 Data Integration

- Definition & Application
- Schema Integration

<https://powcoder.com>

3 Summary

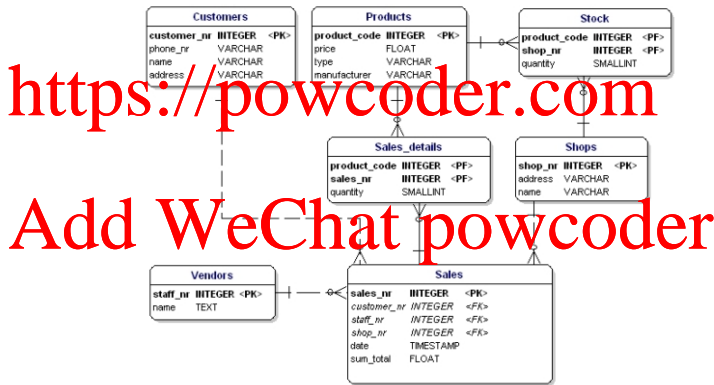
Add WeChat powcoder

What does the schema look like?

- Relational databases:

- A schema specifies a set of tables

- A table contains a set of attributes associated with their data types



This figure is from <http://www.datanamic.com/support/lt-dez005-introduction-db-modeling.html>

What does the schema look like?

- Data models like XML and JSON

- A schema is defined as a set of tags, classes and properties

Assignment Project Exam Help

<https://powcoder.com>

Add What powcoder

```

<us-patent-grant data-bbox="362 259 739 808">
  <us-patent-grant>
    <us-bibliographic-data-grant>
      <publication-reference>
        <document-id>
          <country>US</country>
          <doc-number>8709336</doc-number>
          <kind>B2</kind>
          <date>20110222</date>
        </document-id>
        </publication-reference>
        <application-reference appl-type="utility">
          <document-id>
            <country>US</country>
            <doc-number>1254086</doc-number>
            <date>20090812</date>
          </document-id>
          </application-reference>
          <priority-claims>
            <us-term-of-grant>
              <us-term-extension>35</us-term-extension>
            </us-term-of-grant>
            <classifications-ipc>
              <classification-ipc>
                <ipc-version-indicator>
                  <date>20060101</date>
                </ipc-version-indicator>
                <classification-level>
                  <section>
                    <class>
                      <subclass>
                        <main-group>19</main-group>
                        <subgroup>
                          <classification-value>1</classification-value>
                        </subgroup>
                      </class>
                    </subclass>
                  </section>
                </classification-level>
                <action-date>
                  <date>20110222</date>
                </action-date>
                <generating-office>
                  <country>US</country>
                </generating-office>
                <classification-status>
                  <classification-status>B</classification-status>
                </classification-status>
                <classification-data-source>
                  <classification-data-source>H</classification-data-source>
                </classification-data-source>
                </classifications-ipc>
                <classification-national>
                  <country>US</country>
                  <main-classification>200 118</main-classification>
                  </classification-national>
                </classification-national>
              </classification-ipc>
            </classifications-ipc>
          </priority-claims>
        </document-id>
      </publication-reference>
    </us-bibliographic-data-grant>
  </us-patent-grant>
</us-patent-grant>

```

- Data science

- A data schema is defined as the representation of the data arrangement, relationships and contents.

Why do we need schema integration?

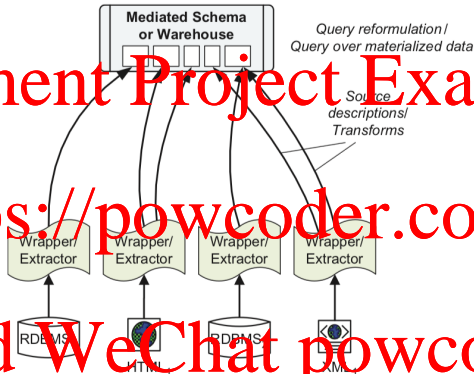


FIGURE 1.4 The basic architecture of a general-purpose data integration system. Data sources can be relational, XML, or any store that contains structured data. The *wrappers* or *loaders* request and parse data from the sources. The *mediated schema* or *central data warehouse* abstracts all source data, and the user poses queries over this. Between the sources and the mediated schema, *source descriptions* and their associated *schema mappings*, or a set of *transformations*, are used to convert the data from the source schemas and values into the global representation.

Schema Mapping

- The linkage between each data source and the mediate schema is done through semantic mapping

Assignment Project Exam Help

- ▶ Specifies how attributes in the sources correspond to attributes in the mediated schema (when such correspondences exist)
- ▶ Specifies how the different groupings of attributes into tables are resolved.
- ▶ Specifies how to resolve schema conflict from different sources

<https://powcoder.com>

Add WeChat powcoder

Problems with Schema Integration

Structure conflicts

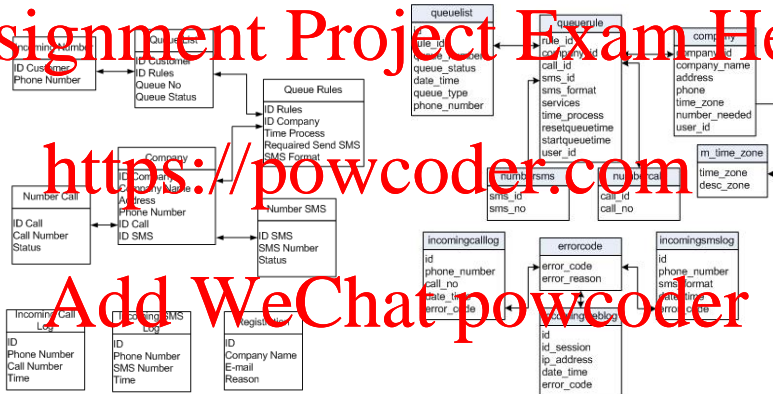
- Inconsistencies in the data structure among schemas, which include
 - Different data source origins: Data can be represented in a structure form (e.g., XML, HTML, JSON, semistructured, or completely unstructured data).
- Inconsistencies among the set of elements inside the different schemas

<https://powcoder.com>

Add WeChat powcoder

Problems with Schema Integration

Structure conflicts



Figures are from <http://www.urremote.com/un/ethering-the-queue-2>

Problems with Schema Integration

Naming conflicts

- homonyms vs synonyms

- ▶ The same name is used for different objects.
- ▶ Different names are used for the same object.

- Examples

- ▶ Homonyms: ID can refer to customer ID, product ID, store ID, etc.
- ▶ Synonyms: Customer ID and Client ID can refer to the same real world object, i.e., customer/client.

Add WeChat powcoder

Problems with Schema Integration

Naming conflicts

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

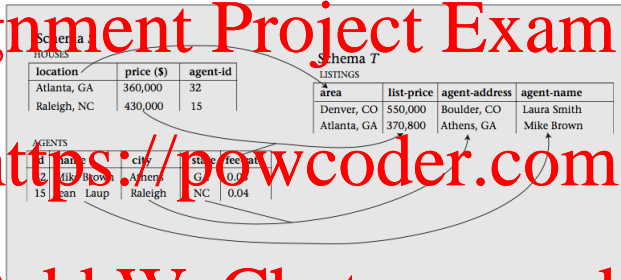


Figure 2: The Schemas of Two Relational Databases S and T on House Listing, and the semantic correspondence between them

Figure is from "Semantic-Integration Research in the Database community" by AnHai Doan and Alon Y. Halevy

- "area" can refer to different real-world entities, e.g., location or square-feet area
- "area" and "location" can refer to the same real-world entity, e.g., the location of the house

Problems with Schema Integration

Entity resolution/conflict resolution

- Different units:

- ▶ Temperature units: Celsius and Fahrenheit
- ▶ Currencies

- Data type heterogeneity

- ▶ Same kind of attributes with different data types
 - ▶ phone number can be stored as string in one database and integer in another database

- Value heterogeneity

- ▶ The use of Abbreviations: Professor v.s. Prof, Street v.s. St, Road v.s. Rd

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Problems with Schema Integration

Entity resolution/conflict resolution

- Semantic heterogeneity: differences in meaning and interpretation of data values¹

- ▶ Naming
 - Case sensitivity
 - Synonyms/Homonyms
 - Acronyms
- ▶ Generalisation/Specialisation: one schema may refer to "phone" but the other schema has multiple elements such as "home phone", "work phone" and "cell phone"
- Level of abstraction: different aggregation levels for an attributes
 - ▶ Address can be split into multiple fields, street number, street name, suburb, city, post-code, etc.
- Different points of time
 - ▶ Fortnight and monthly payment

¹https://en.wikipedia.org/wiki/Semantic_heterogeneity

Schema Integration: semantic matching

Semantic matching: relates a set of elements in schema S to a set of elements in schema T.

Assignment Project Exam Help

DVD-RENTAL

Movies(id, title, year)

Products(mid, releaseDate, releaseCompany, basePrice, rating, saleLocID)

Locations(lid, name, taxRate)

AGGREGATOR

Items(name, releaseDate, classification, price)

FIGURE 5.1 Example of two database schemas. Schema DVD-RENTAL belongs to DVD vendor, while

AGGREGATOR belongs to a shopping site that aggregates products from multiple vendors.

<https://powcoder.com>

Figure is from chapter 5 of "Principles of data integration"

Add WeChat powcoder

- One-to-One match
 - ▶ $\text{Movies.title} \approx \text{Items.name}$
 - ▶ $\text{Movies.year} \approx \text{Items.year}$
 - ▶ $\text{Product.rating} \approx \text{Items.classification}$
- One-to-Many match
 - ▶ $\text{Items.price} \approx \text{Products.basePrices} \times (1 + \text{Locations.taxRate})$

Schema Integration: Name-Based Matcher

Name-Based Matcher: compares the names of attributes (or column headers) in the hope that the names convey the true semantics of the elements.

- Split names according to certain delimiters, such as capitalization, numbers, or special symbols.
 - ▶ ClientName \Rightarrow Client Name
 - ▶ saleLocID \Rightarrow Sale Loc ID
- Expand known abbreviations or acronyms
 - ▶ loc \Rightarrow location
 - ▶ cust \Rightarrow customer
 - ▶ St \Rightarrow Street
 - ▶ DOB \Rightarrow Date of Birth
- Expand a string with its synonyms
 - ▶ Location \Rightarrow Address
 - ▶ Cost \Rightarrow Price
- Expand a string with its hypernyms
 - ▶ product \Rightarrow book, DVD, etc.
- Remove articles, propositions, and conjunctions
 - ▶ Exclude words like “in”, “at”

Schema Integration: Name-Based Matcher

Name-Based Matcher: compares the names of attributes (or column headers) in the hope that the names convey the true semantics of the elements.

DVD VENDOR

Movies(id, title, year)

Products(mid, releaseDate, releaseCompany, basePrice, rating, saleLocID)

Locations(lid, name, taxRate)

AGGREGATOR

Items(name, releaseInfo, classification, price)

(a)

name-based matcher: name \approx (name: 1, title: 0.2)

releaseInfo \approx (releaseDate: 0.5, releaseCompany: 0.5)

price \approx (basePrice: 0.8)

(b)

data-based matcher: name \approx (name: 0.2, title: 0.8)

releaseInfo \approx (releaseDate: 0.7)

classification \approx (rating: 0.6)

price \approx (basePrice: 0.2)

(c)

average combiner: name \approx (name: 0.6, title: 0.5)

releaseInfo \approx (releaseDate: 0.6, releaseCompany: 0.25)

classification \approx (rating: 0.3)

price \approx (basePrice: 0.5)

(d)

FIGURE 5.3 (a) Two schemas (reproduced from Figure 5.1); (b)-(c) the similarity matrices produced by two matchers for the above two schemas; and (d) the combined similarity matrix.

Schema Integration: Instance-Based Matcher

Instance-Based Matcher makes use of the data values.

- Rule-based matching method

- ▶ Hand-crafted rules exploit schema information such as element names, data types, structures, number of subelements, and integrity constraints.

- ▶ For DVD-vendor database:

- All possible classification: G, PG, PG-13, R, etc

- Given a new attribute, if most of its values appear in the list above.

- ▶ Advantages:

- Relatively inexpensive, do not require training

- ▶ Disadvantages:

- Cannot exploit data instances effectively (e.g., value format, frequently occurring values, etc.)

Schema Integration: Instance-Based Matcher

Instance-Based Matcher makes use of the data values.

- Learning-based matching method: learning techniques that can exploit both schema and data information.
 - ▶ Classification-based methods
 - ▶ (semi-)automated but Needs training

Example 1.6

If s_i is address, then positive examples may include “Madison WI” and “Mountain View CA,” and negative examples may include “(608) 695 9813” and “Lord of the Rings.” Now suppose that element t_j is location and that we have access to three data instances of this element: “Milwaukee WI,” “Palo Alto CA,” and “Philadelphia PA.” Then the classifier C may predict confidence scores 0.9, 0.7, and 0.5, respectively. In this case we may return the average confidence score of 0.7 as the similarity score between $s_i = \text{address}$ and $t_j = \text{location}$.

<https://powcoder.com>

Add WeChat powcoder

Summary

- Recap of outlier detection

- Data integration

- Schema integration

- Naming conflict
- Structure conflict
- Entity resolution

- Semantic mapping

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder