Research Issues CIn/UFPE

Ana Carolina Salgado

Summary

Schema Quality
Name Resolution
SPEED Project

Schema Quality Analysis in a Data Integration System

Maria da Conceição Moraes Batista Ana Carolina Salgado

Contents

Introduction

Schema quality

- Minimality
- Type consistency
- Schema completeness

Main contributions

This work

Our goal:

Quality of query execution in data integration systems

Our hypothesis:

The construction of good schemas with high IQ scores will improve query execution

Our proposal:

- The specification of schema IQ criteria minimality, type consistency and schema completeness
- IQ analysis of the integrated schema in a data integration system

Schema issues in data integration

Schemas in Data Integration

The user submits queries to the integrated schema

- A set of views over a number of data sources
- The data integration system reformulates a user query into queries that refers directly to schemas on the sources.
 - Schema mappings: correspondences between data sources and integrated schema elements

Schema Representation

- The X-Entity Model
 - E-R extension for XML data
- X-Entity:
 - Entity = XML Elements
 - Relationships = XML relationships
 - Contains: a XML element contains other XML element
 - Refers: a XML element refers other XML element
 - Attributes

The same real world concept: semantical equivalence (≡) already defined!!!

Minimality

Minimality

- The extent in which the schema is compactly modeled without redundancies
- The more minimal the integrated schema is, the least redundancies it contains, and, consequently, the more efficient the query execution
- A schema is <u>minimal</u> if all relevant domain concepts are described only once
 - No redundant elements Redundancy
 evaluation

Redundancy

Attribute Redundancy:

An attribute A₁ in schema S_m is considered <u>redundant</u>, i.e. <u>Red(A₁, S_m) = 1</u>

if $\exists A_2$ in schema S_m and $A_1 \equiv A_2$

Otherwise, the attribute A₁ in schema S_m is considered <u>non redundant</u>, i.e. Red(A₁, S_m) = 0

Redundancy

Entity Redundancy (ER)

The sum of attributes redundancy defines the entity redundancy:

$$\operatorname{Red}(E_{k},S_{m}) = \frac{\sum_{i=1}^{a_{k}} \operatorname{Red}(A_{ki},E_{k})}{a_{k}}$$

• where $\sum_{k=1}^{a_k} \text{Red}(A_{ki}, E_k)$ is the total number of redundant attributes in entity E_K and a_k is the total number of attributes in E_k .

Redundancy

Relationship Redundancy (RR)

- A relationship between two entities is redundant if there are other semantically equivalent relationships which paths are connecting the same two entities
- The minimality concept is based in the redundancy evaluation

Minimality

- The schema redundancy is measured by the sum of all redundancy values: entity redundancy (ER) and relationships redundancy (RR)
- The <u>schema minimality</u> is measured by the formula:

$$Mi_{S_m} = 1 - [ER(S_m) + RR(S_m)]$$

Schema Redundancy

- The extent in which the attributes corresponding to the same real world concept are represented with the same data type across all schemas
- A schema is <u>type consistent</u> if all its equivalent attributes are represented with the same data type
- The query becomes more efficient: less type conversions executed by the wrappers in composing results

How to evaluate:

- To determine which alternative data type is preferable (standard)
- A schema element is consistent if it adheres to the standard data type

- The type consistency metric is based in:
 - The number of semantically equivalent attributes in schema that adhere to the standard data type defined for the attribute
- Attribute Type Consistency
 - A given attribute A_{pj} is <u>consistent</u> in schema S_p i.e., Cst(A_{pj},S_p)= 1

if every semantically equivalent attribute to $A_{\rm pj}$ appears in the schema $S_{\rm p}$ with the standard data type of attribute $A_{\rm pj}$

Otherwise A_{pj} is <u>inconsistent</u>, i.e., Cst(A_{pj},S_p)= 0

The overall schema type consistency score in a given data integration system (Cst(S_m, D)) is obtained by:

$$\mathsf{Cst}(\mathsf{S}_{\mathsf{m}}, \mathsf{D}) = \frac{\sum_{k=1}^{\mathsf{n}_{\mathsf{m}}} \sum_{j=1}^{\mathsf{a}_{\mathsf{k}}} \mathsf{Cst}(\mathsf{A}_{\mathsf{k}j}, \mathsf{S}_{\mathsf{m}})}{\sum_{k=1}^{\mathsf{n}_{\mathsf{m}}} \mathsf{a}_{\mathsf{k}}} , \text{ where }$$

- $\sum_{k=1}^{n_m} \sum_{j=1}^{a_k} Cst(A_{kj}, S_m)$ is the sum of attribute consistency in S_m
- n_m is the total number of entities in Đ
 a_k is the number of attributes of the entity E_{k18}

Schema Completeness

Schema Completeness

- The <u>schema completeness</u> is the percentage of domain concepts modeled in the integrated schema when related to the concepts represented in *data source schemas*
- The overall schema completeness degree in a schema S_x in a data integration system D is:

SC(
$$S_x$$
) = $\frac{S_x}{S_x}$, where

- $\sigma_{\mathfrak{F}}$ • $\sigma_{\mathfrak{s}_{x}}$ is the number of distinct concepts in the schema S_{x} ;

Contributions

- Specification of three relevant schema IQ criteria - minimality, type consistency and schema completeness
- Analysis of system schema elements according to the specified minimality and type consistency criteria of an integrated schema.
- Specification, implementation and tests of the IQ Manager in a data integration system with a real health care

application

Thanks !!!