

Chapter XVIII

Repairing and Querying Inconsistent Databases

Gianluigi Greco
Università della Calabria, Italy

Sergio Greco
Università della Calabria, Italy

Ester Zumpano
Università della Calabria, Italy

ABSTRACT

The integration of knowledge from multiple sources is an important aspect in several areas such as data warehousing, database integration, automated reasoning systems, active reactive databases and others. Thus a central topic in databases is the construction of integration systems, designed for retrieving and querying uniform data stored in multiple information sources. This chapter illustrates recent techniques for computing repairs as well as consistent answers over inconsistent databases. Often databases may be inconsistent with respect to a set of integrity constraints, that is, one or more integrity constraints are not satisfied. Most of the techniques for computing repairs and queries over inconsistent databases work for restricted cases and only recently there have been proposals to consider more general constraints. In this chapter we give an informal description of the main techniques proposed in the literature.

INTRODUCTION

The problem of integrating heterogeneous sources has been deeply investigated in the fields of multidatabase systems (Breitbart, 1990), federated databases (Wiederhold, 1992) and, more recently, mediated systems (Ullman, 2000; Wiederhold, 1992). A large

variety of approaches has been proposed in the literature for performing data source integration. Many of them are embedded in more complex systems managing the interoperability and the cooperation of data sources characterized by heterogeneous representation formats.

The aim of data integration is to provide uniform integrated access to multiple heterogeneous information sources, which were designed independently for autonomous applications and whose contents are strictly related.

Integrating data from different sources consists of two main steps: first, the various relations are merged together, and second, some tuples are *removed* (or *inserted*) from the resulting database in order to satisfy integrity constraints.

In particular, there are several ways to integrate databases or possibly distributed information sources, but whatever integration architecture we choose, the heterogeneity of the sources to be integrated, often designed independently for autonomous applications, causes subtle problems. In particular, the database obtained from the integration process may be inconsistent with respect to integrity constraints, that is, one or more integrity constraints are not satisfied. Integrity constraints represent an important source of information about the real world. They are usually used to define constraints on data (functional dependencies, inclusion dependencies, etc.) and have, nowadays, a wide applicability in several contexts such as semantic query optimization, cooperative query answering, database integration and view update.

Since, the satisfaction of integrity constraints cannot generally be guaranteed, if the database is obtained from the integration of different information sources, in the evaluation of queries, we must compute answers which are consistent with the integrity constraints.

The following example shows a case of inconsistency.

Example 1. Consider the following database schema consisting of the single binary relation *Teaches* (*Course*, *Professor*) where the attribute *Course* is a key for the relation. Assume there are two different instances for the relations *Teaches*, $D1 = \{(c1, p1), (c2, p2)\}$ and $D2 = \{(c1, p1), (c2, p3)\}$.

The two instances satisfy the constraint that *Course* is a key, but from their union we derive a relation which does not satisfy the constraint since there are two distinct tuples with the same value for the attribute *Course*.

In the integration of two conflicting databases, simple solutions could be based on the definition of preference criteria such as a partial order on the source information or a majority criteria (Lin & Mendelzon, 1996). However, these solutions are not generally satisfactory and more useful solutions are those based on: 1) the computation of ‘repairs’ for the database; 2) the computation of consistent answers (Arenas et al., 1999).

The computation of repairs is based on the definition of minimal sets of insertion and deletion operations so that the resulting database satisfies all constraints. The computation of consistent answers is based on the identification of tuples satisfying integrity constraints and on the selection of tuples matching the goal.

For instance, for the integrated database of Example 1, we have two alternative repairs consisting of the deletion of one of the tuples $(c2, p2)$ and $(c2, p3)$. The consistent answer to a query over the relation *Teaches* contains the unique tuple $(c1, p1)$ so that we don’t know which professor teaches course $c2$.

40 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/repairing-querying-inconsistent-databases/9218

Related Content

Effects of Domain Familiarity on Conceptual Modeling Performance

Jihae Suhand Jinsoo Park (2017). *Journal of Database Management* (pp. 27-55).
www.irma-international.org/article/effects-of-domain-familiarity-on-conceptual-modeling-performance/182868

Metaschemas for ER, ORM and UML Data Models: A Comparison

Terry Halpin (2002). *Journal of Database Management* (pp. 20-30).
www.irma-international.org/article/metaschemas-orm-uml-data-models/3277

Integrity Checking and Maintenance in Relational and Deductive Database and Beyond

Davide Martinenghi, Henning Christiansen and Hendrik Decker (2007). *Intelligent Databases: Technologies and Applications* (pp. 238-285).
www.irma-international.org/chapter/integrity-checking-maintenance-relational-deductive/24236

Toward a Unified Model of Information Systems Development Success

Keng Siau, Yoanna Long and Min Ling (2010). *Journal of Database Management* (pp. 80-101).
www.irma-international.org/article/toward-unified-model-information-systems/39117

Big Data at Scale for Digital Humanities: An Architecture for the HathiTrust Research Center

Stacy T. Kowalczyk, Yiming Sun, Zong Peng, Beth Plale, Aaron Todd, Loretta Auvil, Craig Willis, Jiaan Zeng, Milinda Pathirage, Samitha Liyanage, Guangchen Ruan and J. Stephen Downie (2014). *Big Data Management, Technologies, and Applications* (pp. 270-294).
www.irma-international.org/chapter/big-data-at-scale-for-digital-humanities/85459