

Data-Centric Systems and Applications

Antonio Badia

# SQL for Data Science

Data Cleaning, Wrangling and  
Analytics with Relational Databases



Springer

# **Data-Centric Systems and Applications**

---

## **Series Editors**

Michael J. Carey, University of California, Irvine, CA, USA

Stefano Ceri, Politecnico di Milano, Milano, Italy

## **Editorial Board Members**

Anastasia Ailamaki, École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

Shivnath Babu, Duke University, Durham, NC, USA

Philip A. Bernstein, Microsoft Corporation, Redmond, WA, USA

Johann-Christoph Freytag, Humboldt Universität zu Berlin, Berlin, Germany

Alon Halevy, Facebook, Menlo Park, CA, USA

Jiawei Han, University of Illinois, Urbana, IL, USA

Donald Kossmann, Microsoft Research Laboratory, Redmond, WA, USA

Gerhard Weikum, Max-Planck-Institut für Informatik, Saarbrücken, Germany

Kyu-Young Whang, Korea Advanced Institute of Science & Technology, Daejeon, Korea (Republic of)

Jeffrey Xu Yu, Chinese University of Hong Kong, Shatin, Hong Kong

Intelligent data management is the backbone of all information processing and has hence been one of the core topics in computer science from its very start. This series is intended to offer an international platform for the timely publication of all topics relevant to the development of data-centric systems and applications. All books show a strong practical or application relevance as well as a thorough scientific basis. They are therefore of particular interest to both researchers and professionals wishing to acquire detailed knowledge about concepts of which they need to make intelligent use when designing advanced solutions for their own problems.

Special emphasis is laid upon:

- Scientifically solid and detailed explanations of practically relevant concepts and techniques  
(what does it do)
- Detailed explanations of the practical relevance and importance of concepts and techniques  
(why do we need it)
- Detailed explanation of gaps between theory and practice  
(why it does not work)

According to this focus of the series, submissions of advanced textbooks or books for advanced professional use are encouraged; these should preferably be authored books or monographs, but coherently edited, multi-author books are also envisaged (e.g. for emerging topics). On the other hand, overly technical topics (like physical data access, data compression etc.), latest research results that still need validation through the research community, or mostly product-related information for practitioners (“how to use Oracle 9i efficiently”) are not encouraged.

More information about this series at <http://www.springer.com/series/5258>

Antonio Badia

# SQL for Data Science

Data Cleaning, Wrangling and Analytics  
with Relational Databases



Springer

Antonio Badia  
Computer Engineering & Computer Science  
University of Louisville  
Louisville, KY, USA

ISSN 2197-9723 ISSN 2197-974X (electronic)  
Data-Centric Systems and Applications  
ISBN 978-3-030-57591-5 ISBN 978-3-030-57592-2 (eBook)  
<https://doi.org/10.1007/978-3-030-57592-2>

© Springer Nature Switzerland AG 2020

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG.  
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

# Preface

Data Science (or Data Analytics, or whatever one prefers to call it) is a ‘hot’ topic right now. There is an explosion of courses on the subject, especially online: many universities and several for-profit and non-profit organizations (Coursera, edX, Udacity, Udemy, DataCamp, and many others) offer on-campus and online courses, certification, and degrees. The coverage of these offerings is quite diverse, reflecting the fact that Data Science is still a young and evolving field. However, many courses seem to coalesce around a few topics (Machine Learning, mostly) and tools (R, Python, and SQL, mostly). What few of these courses offer is a textbook.

There are already many books on databases and SQL, but almost all of them focus on the traditional curriculum for Computer Science majors or Information Systems majors (there are a few exceptions, like [11] and [17]). In contrast, the present book explains SQL *within the context* of Data Science and is more in line with what is being taught in these new courses. This book introduces the different parts of SQL as they are needed for the tasks usually carried out during data analysis. Using the framework of the *data life cycle*, it focuses on the steps that are given the short shift in traditional textbooks, like data loading, cleaning, and pre-processing.

This book is for anyone interested in Data Science and/or databases. It should prove useful to anyone taking any of the abovementioned courses, online or on-campus, as well as to students working on their own. It assumes very little from the reader; it just demands a bit of ‘computer fluency,’ but no background on databases or data analysis. In general, all concepts are introduced intuitively and with a minimum of specialized jargon. It contains an appendix (Appendix A) meant to help students without prior experience with databases, with instructions on how to download and install the two open-source database systems (MySQL and Postgres) that we use for examples throughout the book. All readers of the book are encouraged to install both systems and follow the book along with a computer in order to practice, do the exercises, and play around—simply reading the book alone is going to be much less useful than *using* it.

The book is organized as follows: Chapter 1 describes the *Data Life Cycle*, the sequence of stages, from data acquisition and ingestion until archiving, that data goes through as it is prepped for analysis and then actually analyzed, together with

the different activities that take place at each stage. It also explains the different ways that datasets can be organized, and the different types of data one may have to deal with. Many students have an intuitive understanding of the concepts in this chapter, but it is useful to have it all together in one place and to give a name to each concept for later reference. Chapter 2 gets into databases proper, explaining how relational databases organize data. The chapter also explains how data in tables *should* look like (what Hadley Wickham has called *tidy* data [19]), a point which is not traditionally emphasized and can lead to severe problems down the road. Non-traditional data, like XML and text, are also covered. Chapter 3 introduces SQL *queries*, the SQL commands that allow us to ask questions about the data. Unlike traditional textbooks, queries and their parts are described around typical data analysis tasks (data exploration, cleaning, and transformation). These tasks are vital for a proper examination of the data but are frequently overlooked in Data Mining and Machine Learning textbooks. Chapter 4 introduces some basic techniques for Data Analysis. Even though this is not the focus of the book, the chapter shows that SQL can be used for some simple analyses without too much complication.

After this part, which constitutes the core of the book, Chap. 5 introduces additional SQL constructs that come in handy in a variety of situations. This chapter completes the coverage of SQL queries so that readers get an overview of all the main aspects of this important topic. Chapter 6 briefly explains how to use SQL from within R and from within Python programs. This chapter is not an introduction to R (or to Python) and, unlike other chapters in the book, does assume that the reader is already familiar with at least the basics of R and Python. It focuses on how these languages can interact with a database, and how what has been learned about SQL can be leveraged to make life easier when using R or Python.

The book also contains another appendix (besides the one already mentioned), which introduces some basic approaches for handling very large datasets. The purpose of this appendix is to demystify the ideas behind the vague label *Big Data* and give the readers basic guidance on how to use their newly acquired skills in this world.

As in many textbooks, none of what this one contains is new. This book covers the same (or very similar) content to what can be found in many sources, especially online. What this book does is to put it all together under one roof and to give it some order and structure. In many blogs and sites, the material is presented as an answer to a particular question (how do you...?), which may be useful to someone with a specific need but gives the impression that learning SQL is about a bag of tricks. Here, the material is logically organized using the idea of the data life cycle so that all the concepts introduced can be understood as parts of a coherent whole.

Data Science itself is a relatively new and still changing field, but it has deep roots, as it uses approaches and techniques from well-established fields, mostly math (statistics, linear algebra, and others) and computer science (databases, machine learning, and others). As a result, the same concept is sometimes given different names by different authors in different textbooks. Whenever I am aware of this, I

have given a list of known names so that readers with different backgrounds can relate what is in here with what they already know.

The goal of the book is to introduce some basic concepts to a wide variety of readers and provide them a good foundation on which they can build. After going through this book, readers should be able to profitably learn more about Data Mining, Machine Learning, and database management from more advanced textbooks and courses. It is my hope that most of them feel that they have been given a springboard from which they are in a good position to dive deeper into the fascinating world of data analysis.

Louisville, KY, USA  
July 2020

Antonio Badia



# Contents

- 1 The Data Life Cycle**..... 1
  - 1.1 Stages and Operations in the Data Life Cycle ..... 2
  - 1.2 Types of Datasets ..... 6
    - 1.2.1 Structured Data ..... 7
    - 1.2.2 Semistructured Data..... 9
    - 1.2.3 Unstructured Data ..... 16
  - 1.3 Types of Domains ..... 19
    - 1.3.1 Nominal/Categorical Data ..... 20
    - 1.3.2 Ordinal Data ..... 21
    - 1.3.3 Numerical Data..... 21
  - 1.4 Metadata ..... 24
  - 1.5 The Role of Databases in the Cycle ..... 29
- 2 Relational Data**..... 31
  - 2.1 Database Tables ..... 32
    - 2.1.1 Data Types ..... 32
    - 2.1.2 Inserting Data ..... 36
    - 2.1.3 Keys..... 38
    - 2.1.4 Organizing Data into Tables ..... 43
  - 2.2 Database Schemas ..... 48
    - 2.2.1 Heterogeneous Data..... 49
    - 2.2.2 Multi-valued Attributes ..... 50
    - 2.2.3 Complex Data ..... 53
  - 2.3 Other Types of Data ..... 60
    - 2.3.1 XML and JSON Data ..... 61
    - 2.3.2 Graph Data ..... 64
    - 2.3.3 Text ..... 67
  - 2.4 Getting Data In and Out of the Database ..... 69
    - 2.4.1 Importing and Loading Data..... 69
    - 2.4.2 Updating Data ..... 72
    - 2.4.3 Exporting Data ..... 74

<b>3</b>	<b>Data Cleaning and Pre-processing</b>	77
3.1	The Basic SQL Query	77
3.1.1	Joins	83
3.1.2	Functions	89
3.1.3	Grouping	95
3.1.4	Order	101
3.1.5	Complex Queries	103
3.2	Exploratory Data Analysis (EDA)	105
3.2.1	Univariate Analysis	107
3.2.2	Multivariate Analysis	120
3.2.3	Distribution Fitting	129
3.3	Data Cleaning	132
3.3.1	Attribute Transformation	134
3.3.2	Missing Data	144
3.3.3	Outlier Detection	150
3.3.4	Duplicate Detection and Removal	152
3.4	Data Pre-processing	156
3.4.1	Restructuring Data	159
3.5	Metadata and Implementing Workflows	165
3.5.1	Metadata	167
<b>4</b>	<b>Introduction to Data Analysis</b>	171
4.1	What Is Data Analysis?	171
4.2	Supervised Approaches	172
4.2.1	Classification: Naive Bayes	173
4.2.2	Linear Regression	179
4.2.3	Logistic Regression	184
4.3	Unsupervised Approaches	185
4.3.1	Distances and Clustering	185
4.3.2	The kNN Algorithm	191
4.3.3	Association Rules	193
4.4	Dealing with JSON/XML	198
4.5	Text Analysis	202
4.6	Graph Analytics: Recursive Queries	212
4.7	Collaborative Filtering	218
<b>5</b>	<b>More SQL</b>	221
5.1	More on Joins	221
5.2	Complex Subqueries	225
5.3	Windows and Window Aggregates	229
5.4	Set Operations	238
5.5	Expressing Domain Knowledge	241

<b>6 Databases and Other Tools</b> .....	243
6.1 SQL and R .....	243
6.1.1 DBI .....	244
6.1.2 dbplyr .....	247
6.1.3 sqldf .....	250
6.1.4 Packages: Advanced Data Analysis .....	254
6.2 SQL and Python .....	254
6.2.1 Python and Databases: DB-API .....	255
6.2.2 Libraries and Further Analysis .....	259
<b>A Getting Started</b> .....	261
A.1 Downloading and Installing Postgres and MySQL .....	261
A.2 Getting the Server Started .....	262
A.3 User Management .....	264
<b>B Big Data</b> .....	269
B.1 What Is Big Data? .....	269
B.2 Data Warehouses .....	271
B.3 Cluster Databases .....	276
B.4 The Cloud .....	278
<b>References</b> .....	281
<b>Index</b> .....	283