

MANIPULATING DATA WITH STYLE IN SQL

AN INTRODUCTION TO SQL, THE INTERFACE LANGUAGE TO
MOST OF THE WORLD'S STRUCTURED DATA, AND
PRACTICES FOR READABLE AND REUSABLE SQL CODE



RYAN B. HARVEY

OCTOBER 14, 2014

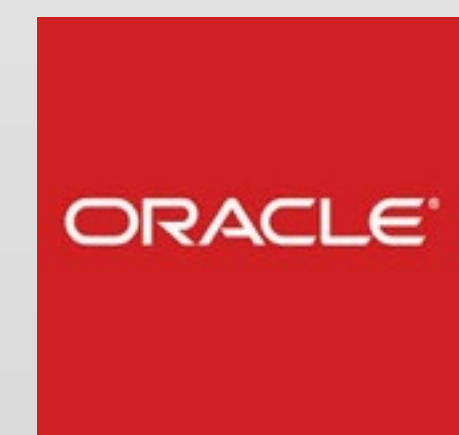


RELATIONAL DATA

- RELATIONAL DATA IS ORGANIZED IN TABLES CONSISTING OF COLUMNS AND ROWS
- FIELDS (COLUMNS) CONSIST OF A COLUMN NAME AND DATA TYPE CONSTRAINT
- RECORDS (ROWS) IN A TABLE HAVE A COMMON FIELD (COLUMN) STRUCTURE AND ORDER
- RECORDS (ROWS) ARE LINKED ACROSS TABLES BY KEY FIELDS

INTRO TO SQL

- SQL (“STRUCTURED QUERY LANGUAGE”) IS A DECLARATIVE DATA DEFINITION AND QUERY LANGUAGE FOR RELATIONAL DATA
- SQL IS AN ISO/IEC STANDARD WITH MANY IMPLEMENTATIONS IN COMMON DATABASE MANAGEMENT SYSTEMS (A FEW BELOW)



WHICH DATABASE SYSTEM SHOULD I USE?

1. USE THE ONE YOUR DATA IS IN
2. UNLESS YOU NEED SPECIFIC THINGS
(PERFORMANCE, FUNCTIONS, ETC.),
USE THE ONE YOU KNOW BEST
3. IF YOU NEED OTHER STUFF OR YOU'VE NEVER
USED A DATABASE BEFORE:
 - A. SQLITE: FOSS, ONE FILE DB, EASY/LIMITED
 - B. POSTGRESQL: FOSS, ENTERPRISE-READY

THE ABOVE ARE MY OPINIONS BASED ON EXPERIENCE. OTHERS MAY DISAGREE, AND THAT'S OK.

SQL: WORKING WITH OBJECTS

- DATA DEFINITION LANGUAGE (DB OBJECTS)
 - **CREATE** (TABLE, INDEX, VIEW, FUNCTION, ...)
 - **ALTER** (TABLE, INDEX, VIEW, FUNCTION, ...)
 - **DROP** (TABLE, INDEX, VIEW, FUNCTION, ...)

SQL: WORKING WITH ROWS

- DATA MANIPULATION LANGUAGE (RECORDS)
AKA QUERY LANGUAGE
- `SELECT ... FROM ...`
- `INSERT INTO ...`
- `UPDATE ... SET ...`
- `DELETE FROM ...`



PostgreSQL



FEATURE COMPARISON

SQL: SELECT STATEMENT

- `SELECT <COL_LIST> FROM <TABLE> ...`
- MERGING/COLUMN BINDING: `JOIN` CLAUSE
- ROW BINDING: `UNION` CLAUSE
- FILTERING: `WHERE` CLAUSE
- AGGREGATION: `GROUP BY` CLAUSE
- AGGREGATED FILTERING: `HAVING` CLAUSE
- SORTING: `ORDER BY` CLAUSE

INTRO TO RELATIONAL ALGEBRA

- BASIC OPERATORS

SELECT	σ	WHERE, HAVING
PROJECT	Π	<COL_LIST>
RENAME	ρ	AS

- JOIN OPERATORS: INNER/OUTER, CARTESIAN

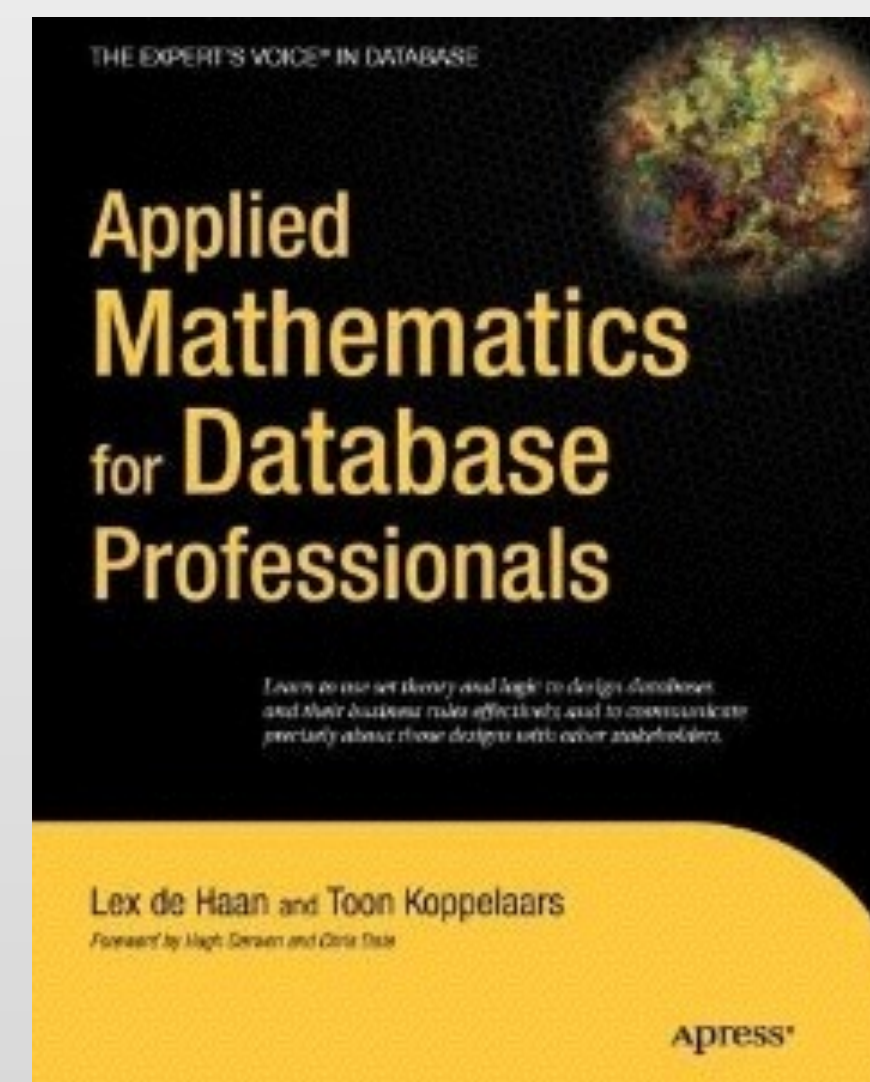
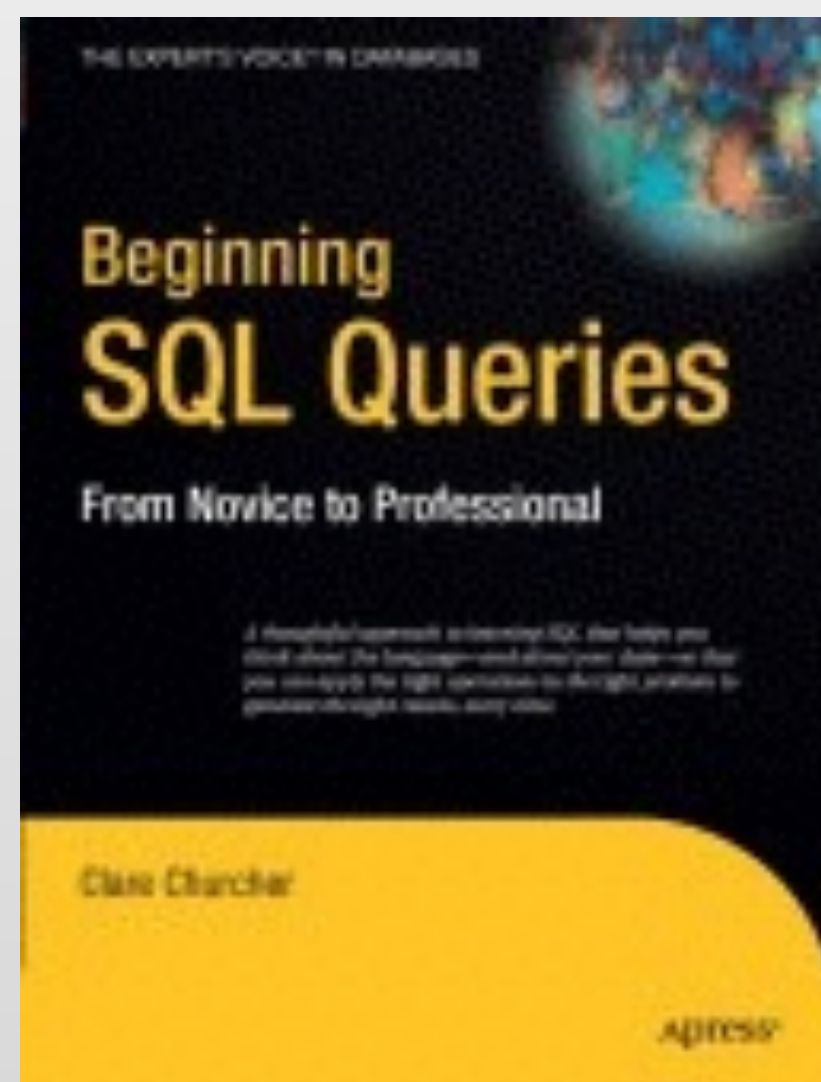
- SET OPERATORS: UNION, INTERSECT, SET MINUS, AND, OR, ETC.

- SELECT NAME, ID FROM T1 WHERE ID < 3 AND DOB < DATE '2004-01-01'

$$\Pi_{NAME, ID} \sigma_{ID < 3 \wedge DOB < (1/1/2004)} (T1)$$

SQL BEGINNER RESOURCES

- BASIC SQL COMMANDS REFERENCE:
[HTTP://WWW.CS.UTEXAS.EDU/~MITRA/
CSFALL2013/CS329/LECTURES/SQL.HTML](http://www.cs.utexas.edu/~mitra/csfall2013/cs329/lectures/sql.html)



SQL: COMMON TABLE EXPRESSIONS (CTEs)

- `WITH <NAME> [(<COL_LIST>)] AS (SELECT ...)`
- `SELECT <COL_LIST> FROM <TABLE OR CTE> ...`
- MERGING/COLUMN BINDING: `JOIN` CLAUSE
- ROW BINDING: `UNION` CLAUSE
- FILTERING: `WHERE` CLAUSE
- AGGREGATION: `GROUP BY` CLAUSE
- AGGREGATED FILTERING: `HAVING` CLAUSE
- SORTING: `ORDER BY` CLAUSE

Same
as
before!



FEATURE COMPARISON

SQL: VIEWS FROM SELECTs

P
P
D
C

- `CREATE VIEW <NAME> AS ...`
- `SELECT <COL_LIST> FROM <TABLE> ...`
 - MERGING/COLUMN BINDING: `JOIN` CLAUSE
 - ROW BINDING: `UNION` CLAUSE
 - FILTERING: `WHERE` CLAUSE
 - AGGREGATION: `GROUP BY` CLAUSE
 - AGGREGATED FILTERING: `HAVING` CLAUSE
 - SORTING: `ORDER BY` CLAUSE



PostgreSQL



FEATURE COMPARISON

SQL: FUNCTIONS FROM VIEWS

- `CREATE FUNCTION` `<NAME>` (`<PARAMS>`) `AS ...`
- `SELECT ... <PARAMS> ...`
 - MERGING/COLUMN BINDING: `JOIN` CLAUSE
 - ROW BINDING: `UNION` CLAUSE
 - FILTERING: `WHERE` CLAUSE
 - AGGREGATION: `GROUP BY` CLAUSE
 - AGGREGATED FILTERING: `HAVING` CLAUSE
 - SORTING: `ORDER BY` CLAUSE



SQL: TUNING WITH EXPLAIN

- `EXPLAIN` <OPTIONS> `SELECT ...` ← Same as before!
- ROWS SCANNED: `COST` OPTION
- WORDY RESPONSE: `VERBOSE` OPTION
- OUTPUT FORMATTING: `FORMAT` OPTION
- ACTUALLY RUN IT: `ANALYZE` OPTION
- RUNTIME (ONLY WITH `ANALYZE`): `TIMING` OPTION
- (`EXPLAIN` IS NOT PART OF THE SQL STANDARD, BUT MOST IMPLEMENTATIONS SUPPORT IT)

SQL: TUNING USING INDEXES

- `CREATE INDEX <NAME> ON <TABLE>`
`(<COL_LIST|EXPRESSION>) ...`
- `UNIQUE` INDICES FOR KEY FIELDS
- USE FUNCTIONS IN EXPRESSIONS:
`LOWER(<TEXT_COL>), INT(<NUM_COL>)`
- SPECIFY ORDERING (`ASC`, `DESC`, `NULLS FIRST`,
ETC.) AND METHOD (`BTREE`, `HASH`, `GIST`, ETC.)
- PARTIAL INDEXES VIA `WHERE` CLAUSE

What's in
your
WHERE
clause?

SQL IN OTHER LANGUAGES

(OR, ACCESSING DATA IN DATABASES VIA SQL IN OTHER LANGUAGES)

- R WITH LIBRARIES
 - RPOSTGRESQL, DPLYR
- PYTHON WITH MODULES
 - PSYCOPG2, SQLALCHEMY



SQL IN OTHER LANGUAGES

(OR, OPERATING ON OTHER LANGUAGES' DATA STRUCTURES VIA SQL)

- R WITH LIBRARIES
 - RSQLITE, SQLDF
- PYTHON WITH MODULES
 - PANDAS, PANDASQL



Mostly,
Data
Frames.

SLIDES AND CODE ARE AVAILABLE ON GITHUB AT [NIHONJINRXS/POLYGLOT-OCTOBER2014!](https://github.com/nihonjinrxs/polyglot-october2014)

The screenshot shows the GitHub interface for the repository `nihonjinrxs / polyglot-october2014`. The repository has 7 commits, 1 branch, 0 releases, and 1 contributor. The main content area displays a list of files and folders with their commit history:

File/Folder	Commit Message	Time Ago
<code>R</code>	initial commit: license, readme, slide deck, data, R and python code.	17 hours ago
<code>data</code>	initial commit: license, readme, slide deck, data, R and python code.	17 hours ago
<code>python</code>	initial commit: license, readme, slide deck, data, R and python code.	17 hours ago
<code>slides</code>	Added PDF and HTML versions of slide deck.	17 hours ago
<code>sql</code>	Completed SQL examples, including CTEs, views, functions, explain & i...	17 minutes ago
<code>.gitignore</code>	initial commit: license, readme, slide deck, data, R and python code.	17 hours ago
<code>LICENSE</code>	initial commit: license, readme, slide deck, data, R and python code.	17 hours ago
<code>README.md</code>	Minor typo edits to README.md.	3 hours ago

The `README.md` file is expanded, showing the title **Manipulating Data in Style with SQL** and the following text:

An introduction to SQL, the interface language to most of the world's structured data, and practices for readable and reusable SQL code

This repository contains materials for my talk at the Polyglot Programming DC meetup on October 14, 2014, which is loosely based on material from my talks at the Data Wranglers DC meetups on June 4, 2014 (materials at [nihonjinrxs/dwdc-june2014](#)) and August 6, 2014 (materials at [nihonjinrxs/dwdc-](#)

On the right side of the repository page, there are links for `Code`, `Issues`, `Pull Requests`, `Wiki`, `Pulse`, `Graphs`, and `Settings`. At the bottom, there are buttons for `Clone in Desktop` and `Download ZIP`.



RYAN B. HARVEY

[HTTP://DATASCIENTIST.GURU](http://datascientist.guru)

RYAN.B.HARVEY@GMAIL.COM

[@NIHONJINRXS](#)

[+RYAN.B.HARVEY](#)

EMPLOYMENT & AFFILIATIONS*

IT PROJECT MANAGER

OFFICE OF MANAGEMENT AND BUDGET

EXECUTIVE OFFICE OF THE PRESIDENT

DATA SCIENTIST & SOFTWARE ARCHITECT

KITCHOLOGY INC.

RESEARCH AFFILIATE

NORBERT WIENER CENTER FOR HARMONIC ANALYSIS & APPLICATIONS

COLLEGE OF COMPUTER, MATHEMATICAL & NATURAL SCIENCES

UNIVERSITY OF MARYLAND AT COLLEGE PARK

Thank you!
Questions?

* MY REMARKS, PRESENTATION AND PREPARED MATERIALS ARE MY OWN, AND DO NOT REPRESENT THE VIEWS OF MY EMPLOYERS.