

30 May 2015

General Assembly -- Boston

SQL Bootcamp

Let's run some queries

DAY I GOT COOKIE



Starting simple

```
1 SELECT
2     au_id, au_fname, au_lname
3 FROM
4     books.authors
5 ;
6
```

- Gives us selected fields for all records

```
1 +-----+-----+-----+
2 | au_id | au_fname | au_lname
3 +-----+-----+-----+
4 | A01  | Sarah    | Buchman
5 | A02  | Wendy    | Heydemark
6 | A03  | Hallie   | Hull
7 | A04  | Klee     | Hull
8 | A05  | Christian | Kells
9 | A06  |
10 | A07 | Paddy    | Kellsey
11 |          |          | O'Furniture
12 +-----+-----+-----+
```

Get different fields

```
1 SELECT
2     au_id, au_lname, city, state
3 FROM
4     books.authors
5 ;
6
```

Gives us a different selection of fields, again for all records

```
1 +-----+-----+-----+-----+
2 | au_id | au_lname | city   | state  |
3 +-----+-----+-----+-----+
4 | A01  | Buchman  | Bronx  | NY     |
5 | A02  | Heydemark | Boulder | CO     |
6 | A03  | Hull      | San Francisco | CA     |
7 | A04  | Hull      | San Francisco | CA     |
8 | A05  | Kells     | New York  | NY     |
9 | A06  | Kellsey   | Palo Alto | CA     |
10 | A07 | O'Furniture | Sarasota | FL     |
11 +-----+-----+-----+-----+
12
```

Simple Query

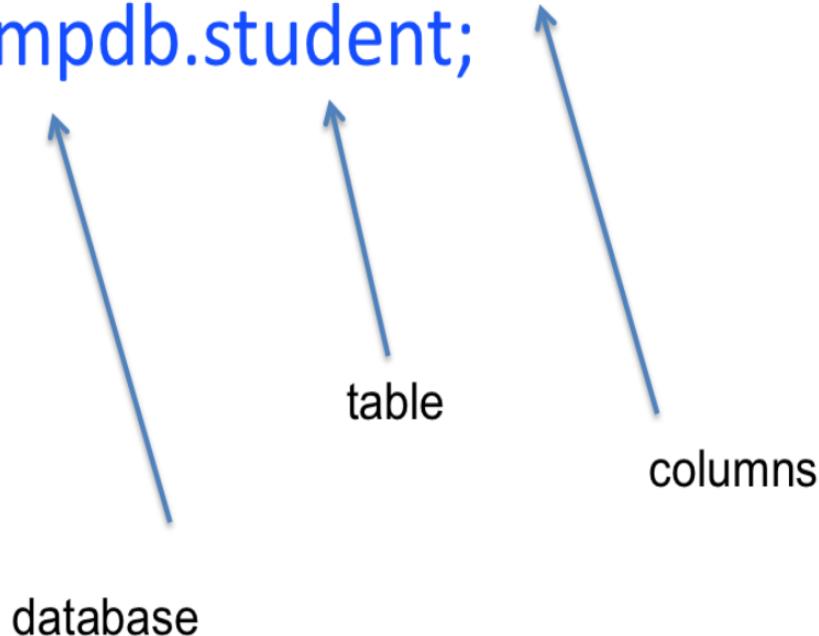
```
1 SELECT
2     name, sex, student_id
3 FROM
4     sampdb.student
5 LIMIT 5
6 ;
7
```

```
1 +-----+-----+-----+
2 | name  | sex   | student_id |
3 +-----+-----+-----+
4 | Megan | F     | 1          |
5 | Joseph| M     | 2          |
6 | Kyle   | M     | 3          |
7 | Katie  | F     | 4          |
8 | Abby   | F     | 5          |
9 +-----+-----+-----+
10
```

- Query gets all fields of all records in table
- LIMIT, however, truncates output -- useful when we want to peek, but don't want the whole table

Simple Query

```
SELECT name, sex, student_id  
FROM sampdb.student;
```



We could also use:

```
SELECT * FROM sampdb.student
```

- << * >> is a wildcard meaning “all columns”

Let's unpack that:

"database"

A program concept for holding "tables"

Comparable to a folder in your computer's file system

"table"

A program concept holding data about instances of a particular entity type - students, presidents, exams

These data are held in "records" or "rows"

"record"

Data about a particular instance eg, particular student or particular movie

Data is in a collection of fields or columns

Comparable to a row in a spreadsheet

"Column" / "field"

An attribute of an entity type about which data is collected - a student's name, or sex, or the state of origin of a president

Comparable to a cell in a spreadsheet row

Schema Exploration

How do we see what is available to us?

What databases are available?

```
1 SHOW DATABASES;  
2
```

| 1 | Database |
|----|--------------------|
| 2 | information_schema |
| 3 | books |
| 4 | csv_load |
| 5 | demo_app_db |
| 6 | entities |
| 7 | google_api |
| 8 | homework |
| 9 | join_sample |
| 10 | learning |
| 11 | mta |
| 12 | mysql |
| 13 | sampdb |
| 14 | |
| 15 | |
| 16 | |
| 17 | |

- Shows the databases available to us on the server

What tables are available?

```
1 SHOW TABLES IN books;
2
```

```
1 +-----+
2 | Tables_in_books |
3 +-----+
4 | au_orders
5 | authors
6 | dups
7 | employees
8 | empsales
9 | hier
10 | publishers
11 | roadtrip
12 | royalties
13 | telephones
14 | temps
15 | time_series
16 | title_authors
17 | titles
18 +-----+
19
```

- Shows the tables available to us within a particular database

What columns are available?

```
1 DESCRIBE books.authors;
2
3
```

| 1 | Field | Type | Null | Key | Default | Extra |
|----|----------|-------------|------|-----|---------|-------|
| 2 | au_id | char(3) | NO | PRI | NULL | |
| 3 | au_fname | varchar(15) | NO | | NULL | |
| 4 | au_lname | varchar(15) | NO | | NULL | |
| 5 | phone | varchar(12) | YES | | NULL | |
| 6 | address | varchar(20) | YES | | NULL | |
| 7 | city | varchar(15) | YES | | NULL | |
| 8 | state | char(2) | YES | | NULL | |
| 9 | zip | char(5) | YES | | NULL | |
| 10 | | | | | | |
| 11 | | | | | | |
| 12 | | | | | | |
| 13 | | | | | | |

- Shows columns in the table, their data types, whether they are primary keys, whether they can take NULL values

We can see how a table is created

```
1 SHOW CREATE TABLE books.authors\G
2
3
```

```
1 **** 1. row ****
2      Table: authors
3 Create Table: CREATE TABLE `authors` (
4       `au_id` char(3) NOT NULL,
5       `au_fname` varchar(15) NOT NULL,
6       `au_lname` varchar(15) NOT NULL,
7       `phone` varchar(12) DEFAULT NULL,
8       `address` varchar(20) DEFAULT NULL,
9       `city` varchar(15) DEFAULT NULL,
10      `state` char(2) DEFAULT NULL,
11      `zip` char(5) DEFAULT NULL,
12      PRIMARY KEY (`au_id`)
13 ) ENGINE=MyISAM DEFAULT CHARSET=latin1
14
```

- Shows columns in the table, their data types, whether they are primary keys, whether they can take NULL values

Practice

Find a table in books that has 'emp_id' as primary key

Find a table in books that has 'emp_name'

Write a query outputting emp_id and emp_name

What columns are available in sampdb.president?

Write a query outputting each President's first name, last name, date of birth and date of death

What is SQL?



Role of SQL in Data Analysis

1. Data analysis requires:

Getting, cleaning, loading, 'modeling'

Exploring, manipulating, extracting, analyzing

Reporting, visualization

2. SQL is best at analysis, but all help with all of this

Getting and Cleaning

1. Something else has to get the data to CSV or database
2. SQL string functions can help clean and reformat data
3. Tricks with LOAD and INSERT can remodel data

Reporting and Visualization

1. SQL has zero visualization capacity, but it can output data as needed for your graphics program
2. "Reporting" is here sort of a term of art

E.g. pivot tables, cross tabs, and so forth

SQL can do an awful lot, but isn't great at this

Implementations do provide base layers to reporting tools like Jasper, Pentaho, etc

Analysis and Transformation

1. SQL's stock in trade is filtering, aggregating and grouping on set-based criteria

Filter WHERE records meet some stated criteria

2. State equals Illinois, name in (Smith, Baker, Han)

3. Date before 2012-06-01

4. Score greater than 15

GROUP BY values in a field, or in a combination of fields

5. Uses that capacity to combine information in different sets / tables by matching records on some criteria

6. We can nest these abilities into "subqueries", generating criteria on the fly

As A Language - Application Usage

1. I.e., why we care

2. A syntax for:

Describing tables

Extracting and manipulating data from those tables

Plus dealing with the cruft of managing data: loading, indexing, etc

3. The action is in bullet #2, but we have to deal with the other stuff too

Extracting and Manipulating

1. Extracting:

We want to specify criteria for which records we display or analyze

Combine data extracted from several tables

2. Manipulating:

Transform values in fields

Calculate values for fields of groups of records

3. These fundamental operations can be extended to very complex transformations

What is SQL -- Architecture

A "Standard", Not a Program

- **There are many SQLs**

MySQL / Postgres / SQLite / Microsoft Access / Oracle / Sybase / etc etc etc

- **Each is an implementation of a standard**

"SQL" is basically a program's promise to behave in a defined way in response to a particular query for a particular set of data

"Client / Server"

- **Most implementations are actually two programs**

One is a front end / client, giving us access to a server and showing its outputs

- **MySQL Workbench**

The other is a back end / server, holding the data and processing the queries

- **MySQL**

MySQL Workbench is to MySQL as Chrome browser is to an HTTP server

Client / Server Implications

1. We can access existing data

We don't need data dumps to do analysis

2. We do need some sensible and safe access to the existing databases

If your production database has a million records, you can analyze them all without downloading them all

3. We can access remote data

If your servers are on Rackspace or Amazon you can still get at the data

4. Programs can access the data

Languages like Ruby, Python, PHP all have tools for interacting with databases

They act as clients, but getting data for (say) serving web pages rather than display to an analyst



Why Excel?

1. Good graphic interface orients us to the data and the tools

We can find tools just by clicking around the menus

2. Excel organizes models visually, so we can see the relationships among cells

We can see the cells flowing into a formula

We imagine extensions of our models by thinking through new visual relationships

3. Excel handles many issues we don't even know about

Until we use a tool with less hand-holding - like SQL

Why SQL? -- Analytic

1. Operated through text statements

Easy to identify and extend patterns

2. Replicable

If a new query resembles a prior one, just copy and revise the prior

3. Chainable

We can easily add conditions and "filters"

We can write a query and refer to its results in another query

4. Traceable

We can see how we got to a result by examining the query

5. Non-visual modeling allows us to fluidly work in more than two dimensions

Implications

1. Excel gets used like a Swiss Army Knife

We have problems, Excel gives us path to solve them

2. SQL can do the same things and more, BUT requires infrastructure, training and problem solving - i.e., brain damage

When we do need it, who has time for training?

Why SQL? -- Organizational

1. A lot of data already in your company's relational databases
2. Dev team familiar with SQL
3. SQL easier to automate
4. Easy to operate over a network

Manipulating Table Records

Behold, some data

Simple query

```
1 SELECT
2     first_name,
3     last_name,
4     state
5 FROM
6     sampdb.president
7 LIMIT
8     20
9 ;
10
```

Results

| | first_name | last_name | state |
|----|---------------|------------|-------|
| 4 | John | Adams | MA |
| 5 | John Quincy | Adams | MA |
| 6 | Chester A. | Arthur | VT |
| 7 | James | Buchanan | PA |
| 8 | George H.W. | Bush | MA |
| 9 | George W. | Bush | CT |
| 10 | James E. | Carter | GA |
| 11 | Grover | Cleveland | NJ |
| 12 | William J. | Clinton | AR |
| 13 | Calvin | Coolidge | VT |
| 14 | Dwight D. | Eisenhower | TX |
| 15 | Millard | Fillmore | NY |
| 16 | Gerald R. | Ford | NE |
| 17 | James A. | Garfield | OH |
| 18 | Ulysses S. | Grant | OH |
| 19 | Warren G. | Harding | OH |
| 20 | Benjamin | Harrison | OH |
| 21 | William H. | Harrison | VA |
| 22 | Rutherford B. | Hayes | OH |
| 23 | Herbert C. | Hoover | IA |

- No order imposed on the results

We can specify an order

Re-order by state

```
1 SELECT
2     first_name,
3     last_name,
4     state
5 FROM
6     sampdb.president
7 ORDER BY
8     state
9 LIMIT
10    20
11 ;
12
```

Results, Ordered

| | first_name | last_name | state |
|----|-------------|-----------|-------|
| 1 | William J. | Clinton | AR |
| 2 | Richard M. | Nixon | CA |
| 3 | George W. | Bush | CT |
| 4 | James E. | Carter | GA |
| 5 | Herbert C. | Hoover | IA |
| 6 | Ronald W. | Reagan | IL |
| 7 | Barack | Obama | IL |
| 8 | Abraham | Lincoln | KY |
| 9 | George H.W. | Bush | MA |
| 10 | John | Adams | MA |
| 11 | John F. | Kennedy | MA |
| 12 | John Quincy | Adams | MA |
| 13 | Harry S | Truman | MO |
| 14 | James K. | Polk | NC |
| 15 | Andrew | Johnson | NC |
| 16 | Gerald R. | Ford | NE |
| 17 | Franklin | Pierce | NH |
| 18 | Grover | Cleveland | NJ |
| 19 | Martin | Van Buren | NY |
| 20 | Franklin D. | Roosevelt | NY |

- Note that ordering applies only to `RESULT`

We can re-order by non-string, non-alphabetized fields

Re-order by birth

```
1 SELECT
2     first_name,
3     last_name,
4     birth
5 FROM
6     sampdb.president
7 ORDER BY
8     birth
9 LIMIT
10    10
11 ;
12
```

Results

| | first_name | last_name | birth |
|----|-------------|------------|------------|
| 4 | George | Washington | 1732-02-22 |
| 5 | John | Adams | 1735-10-30 |
| 6 | Thomas | Jefferson | 1743-04-13 |
| 7 | James | Madison | 1751-03-16 |
| 8 | James | Monroe | 1758-04-28 |
| 9 | Andrew | Jackson | 1767-03-15 |
| 10 | John Quincy | Adams | 1767-07-11 |
| 11 | William H. | Harrison | 1773-02-09 |
| 12 | Martin | Van Buren | 1782-12-05 |
| 13 | Zachary | Taylor | 1784-11-24 |
| 14 | | | |
| 15 | | | |

And we can reverse the order

Re-order by birth, reversed

```
1 SELECT
2     first_name,
3     last_name,
4     birth
5 FROM
6     sampdb.president
7 ORDER BY
8     birth DESC
9 LIMIT
10    10
11 ;
12
```

Results

| | first_name | last_name | birth |
|----|-------------|------------|------------|
| 4 | William J. | Clinton | 1946-08-19 |
| 5 | George W. | Bush | 1946-07-06 |
| 6 | James E. | Carter | 1924-10-01 |
| 7 | George H.W. | Bush | 1924-06-12 |
| 8 | John F. | Kennedy | 1917-05-29 |
| 9 | Gerald R. | Ford | 1913-07-14 |
| 10 | Richard M. | Nixon | 1913-01-09 |
| 11 | Ronald W. | Reagan | 1911-02-06 |
| 12 | Lyndon B. | Johnson | 1908-08-27 |
| 13 | Dwight D. | Eisenhower | 1890-10-14 |
| 14 | | | |
| 15 | | | |

We can re-order by any field

Re-order by first_name

```
1 SELECT
2     first_name,
3     last_name,
4     state
5 FROM
6     sampdb.president
7 ORDER BY
8     first_name
9 LIMIT
10    10
11 ;
12
```

Results

| | first_name | last_name | state |
|----|-------------|------------|-------|
| 4 | Abraham | Lincoln | KY |
| 5 | Andrew | Johnson | NC |
| 6 | Andrew | Jackson | SC |
| 7 | Benjamin | Harrison | OH |
| 8 | Calvin | Coolidge | VT |
| 9 | Chester A. | Arthur | VT |
| 10 | Dwight D. | Eisenhower | TX |
| 11 | Franklin | Pierce | NH |
| 12 | Franklin D. | Roosevelt | NY |
| 13 | George | Washington | VA |

We can re-order by several fields

Re-order by state and first_name

```
1 SELECT
2     state,
3     first_name,
4     last_name
5 FROM
6     sampdb.president
7 ORDER BY
8     state,
9     first_name
10 ;
11
```

Results

| 1 | state | first_name | last_name |
|----|-------|-------------|-----------|
| 2 | AR | William J. | Clinton |
| 3 | CA | Richard M. | Nixon |
| 4 | CT | George W. | Bush |
| 5 | GA | James E. | Carter |
| 6 | IA | Herbert C. | Hoover |
| 7 | IL | Barack | Obama |
| 8 | IL | Ronald W. | Reagan |
| 9 | KY | Abraham | Lincoln |
| 10 | MA | George H.W. | Bush |
| 11 | MA | John | Adams |
| 12 | MA | John F. | Kennedy |
| 13 | MA | John Quincy | Adams |
| 14 | MO | Harry S | Truman |
| 15 | NC | Andrew | Johnson |
| 16 | NC | James K. | Polk |
| 17 | NE | Gerald R. | Ford |
| 18 | . | | |
| 19 | . | | |
| 20 | . | | |
| 21 | | | |
| 22 | | | |

ORDER BY Several Fields Introduces "Tuples"

Tuples provide crucial insight to structure of our records

- A 'tuple' is a subset of a record's fields
- For a given record, a tuple is that record's values for the specified fields
- Tuples are ordered -- (first_name, last_name) is a different tuple than (last_name, first_name)
- Note that records are, strictly speaking, tuples themselves, as they are ordered collections of fields
- At a 'query' level, a tuple is a specification of the fields we will look at for each particular record

We use tuples to specify unique or similar records, which helps us filter data

Data Structure / Records Structure

Analysis assumes some structure in the data

Data structure?

1. No patterns == Noise
2. Our records' values fall into patterns
3. Finding those patterns and interpreting them is the point of analysis

Tuples give us a path into structure of our data

We can use tuples to indicate unique records

If a particular set of fields has no more than 1 record for each combination of values in the table, that tuple indicates unique records

State and first_name fail for president -- we have several records with VA and James

How can we find unique records without visual inspection?

Getting Tuple Combinations

SELECT DISTINCT: tuple combinations

```
1 SELECT DISTINCT
2   state,
3   first_name
4 FROM
5   sampdb.president
6 ORDER BY
7   state,
8   first_name
9
```

This returns the set of unique combinations of values for the tuple

No duplications

But can we tell if any have duplicates?

We need information about the records themselves -- a count of records by tuple field combination values

Just unique tuples

| | state | first_name |
|----|-------|-------------|
| 4 | AR | William J. |
| 5 | CA | Richard M. |
| 6 | CT | George W. |
| 7 | GA | James E. |
| 8 | IA | Herbert C. |
| 9 | IL | Barack |
| 10 | IL | Ronald W. |
| 11 | KY | Abraham |
| 12 | MA | George H.W. |
| 13 | MA | John |
| 14 | MA | John F. |
| 15 | MA | John Quincy |
| 16 | MO | Harry S |
| 17 | NC | Andrew |
| 18 | NC | James K. |
| 19 | NE | Gerald R. |
| 20 | NH | Franklin |
| 21 | NJ | Grover |
| 22 | NY | Franklin D. |
| 23 | NY | Martin |
| 24 | NY | Millard |
| 25 | NY | Theodore |
| 26 | ... | |
| 27 | | |
| 28 | | |

Aggregation

Aggregation Functions Describe Record Sets

COUNT(*) counts records

```
1 SELECT
2   COUNT(*)
3 FROM
4   sampdb.president
5 ;
6
```

Result

| | | | |
|---|---|------------|---|
| 1 | + | ----- | + |
| 2 | | COUNT(*) | |
| 3 | + | ----- | + |
| 4 | | 42 | |
| 5 | + | ----- | + |
| 6 | | | |

This returns the count of records returned by the query conditions

Here, the number of records in the table

COUNT()

COUNT()

- A "function" -- performs calculations based on arguments passed to it
- Previous query base, as an argument, the "all records" wild card
- An "aggregating" function -- it takes a set of records or values, and performs calculations on that set
- We will also see functions that act on values of only a single record
- COUNT() simply counts the (not NULL) values passed to it
- There is an exception for the * argument, COUNT() counts all records even if some have only NULL values

COUNT() Fields

COUNT([field name]) counts values in fields

```
1 SELECT
2     COUNT(first_name),
3     COUNT(death)
4 FROM
5     sampdb.president
6 ;
7
```

Result

| | COUNT(first_name) | COUNT(death) |
|---|-------------------|--------------|
| 4 | 42 | 38 |

COUNT() Distinct Tuples

COUNT() distinct tuples

```
1 SELECT
2   COUNT(DISTINCT state, first_name)
3   FROM
4     sampdb.president
5
6 ;
7
```

Result

| | | | |
|---|---|-----------------------------------|---|
| 1 | + | ----- | + |
| 2 | | COUNT(DISTINCT state, first_name) | |
| 3 | + | ----- | + |
| 4 | | 41 | |
| 5 | + | ----- | + |
| 6 | | | |

GROUP BY Subsets Data

GROUP BY divides records by unique tuples

Aggregation functions count within those subsets, not over the entire table

```
1 SELECT
2     state,
3     COUNT(*)
4 FROM
5     sampdb.president
6 GROUP BY
7     state
8 ORDER BY
9     state
10 ;
11
```

Result

| | state | COUNT(*) |
|----|-------|----------|
| 4 | AR | 1 |
| 5 | CA | 1 |
| 6 | CT | 1 |
| 7 | GA | 1 |
| 8 | IA | 1 |
| 9 | IL | 1 |
| 10 | KY | 1 |
| 11 | MA | 4 |
| 12 | MO | 1 |
| 13 | NC | 2 |
| 14 | NE | 1 |
| 15 | NH | 1 |
| 16 | NJ | 1 |
| 17 | NY | 4 |
| 18 | OH | 7 |
| 19 | PA | 1 |
| 20 | SC | 1 |
| 21 | TX | 2 |
| 22 | VA | 8 |
| 23 | VT | 2 |

GROUP BY Tuple

We can GROUP BY tuples

```
1 SELECT
2     state,
3     first_name,
4     COUNT(*)
5 FROM
6     sampdb.president
7 GROUP BY
8     state,
9     first_name
10 ORDER BY
11     COUNT(*) DESC,
12     state,
13     first_name
14 LIMIT
15     10
16 ;
17
```

Result

| | state | first_name | COUNT(*) |
|----|-------|-------------|----------|
| 4 | VA | James | 2 |
| 5 | AR | William J. | 1 |
| 6 | CA | Richard M. | 1 |
| 7 | CT | George W. | 1 |
| 8 | GA | James E. | 1 |
| 9 | IA | Herbert C. | 1 |
| 10 | IL | Ronald W. | 1 |
| 11 | KY | Abraham | 1 |
| 12 | MA | George H.W. | 1 |
| 13 | MA | John | 1 |

So we can count records within tuples

Exciting, right?

Doesn't seem like much

- But notice that the syntax depends heavily on context
- Aggregating functions -- eg, COUNT() -- behave very differently based on the presence and detail of a GROUP BY clause
- Our queries define sets of records, return particular records or information about subsets of those records
- GROUP BY drives the collection of records into subsets for summary
- Doesn't seem like much, but provides a foundation for our next steps through the syntax

Practice

Write a query counting presidents by first_name

Order that query in descending order

Limit the output to the first five names

Do the same -- all three steps -- for last_name

Filtering -- WHERE and HAVING

Filtering the results

```
1 SELECT
2   first_name,
3   last_name,
4   state
5 FROM
6   sampdb.president
7 WHERE
8   state = 'VA'
9 ;
10
```

| | first_name | last_name | state |
|----|------------|------------|-------|
| 4 | William H. | Harrison | VA |
| 5 | Thomas | Jefferson | VA |
| 6 | James | Madison | VA |
| 7 | James | Monroe | VA |
| 8 | Zachary | Taylor | VA |
| 9 | John | Tyler | VA |
| 10 | George | Washington | VA |
| 11 | Woodrow | Wilson | VA |

- Notice the reorganization of the text for clarity

Filtering by comparison

```
1 SELECT
2   first_name,
3   last_name,
4   birth,
5   state
6 FROM
7   sampdb.president
8 WHERE
9   birth <= '1900-01-01'
10 ORDER BY
11   birth
12 LIMIT 10
13 ;
14
```

| | first_name | last_name | birth | state |
|----|-------------|------------|------------|-------|
| 4 | George | Washington | 1732-02-22 | VA |
| 5 | John | Adams | 1735-10-30 | MA |
| 6 | Thomas | Jefferson | 1743-04-13 | VA |
| 7 | James | Madison | 1751-03-16 | VA |
| 8 | James | Monroe | 1758-04-28 | VA |
| 9 | Andrew | Jackson | 1767-03-15 | SC |
| 10 | John Quincy | Adams | 1767-07-11 | MA |
| 11 | William H. | Harrison | 1773-02-09 | VA |
| 12 | Martin | Van Buren | 1782-12-05 | NY |
| 13 | Zachary | Taylor | 1784-11-24 | VA |
| 14 | | | | |
| 15 | | | | |

- Here we take only those presidents born before 20th century

Filtering affects set GROUPed

```
1 SELECT
2     state,
3     COUNT(*)
4 FROM
5     sampdb.president
6 WHERE
7     birth <= '1900-01-01'
8 GROUP BY
9     state
10 ORDER BY
11     COUNT(*)
12 ;
13
```

| state | COUNT(*) |
|-------|----------|
| KY | 1 |
| NH | 1 |
| PA | 1 |
| MO | 1 |
| NJ | 1 |
| TX | 1 |
| IA | 1 |
| SC | 1 |
| MA | 2 |
| VT | 2 |
| NC | 2 |
| NY | 4 |
| OH | 7 |
| VA | 8 |

- We count here by state for presidents before 20th century

HAVING filters after GROUP

```
1 SELECT
2     state,
3     COUNT(*)
4 FROM
5     sampdb.president
6 WHERE
7     birth <= '1900-01-01'
8 GROUP BY
9     state
10 HAVING
11     COUNT(*) >= 2
12 ORDER BY
13     COUNT(*)
14 ;
15
```

| state | COUNT(*) |
|-------|----------|
| MA | 2 |
| VT | 2 |
| NC | 2 |
| NY | 4 |
| OH | 7 |
| VA | 8 |

• We count here by state for presidents before 20th century

Filter by list

```
1 SELECT
2   first_name,
3   last_name,
4   state
5 FROM
6   sampdb.president
7 WHERE
8   state
9     IN
10  (
11    'MA',
12    'VA'
13  )
14 ORDER BY
15   state
16 ;
17
```

| | first_name | last_name | state |
|----|-------------|------------|-------|
| 4 | John | Adams | MA |
| 5 | John Quincy | Adams | MA |
| 6 | George H.W. | Bush | MA |
| 7 | John F. | Kennedy | MA |
| 8 | George | Washington | VA |
| 9 | John | Tyler | VA |
| 10 | Zachary | Taylor | VA |
| 11 | James | Monroe | VA |
| 12 | James | Madison | VA |
| 13 | Thomas | Jefferson | VA |
| 14 | William H. | Harrison | VA |
| 15 | Woodrow | Wilson | VA |
| 16 | | | |
| 17 | | | |

- We use parentheses to define a list of state for which we want presidents

Filter by tuple

```
1 SELECT
2   first_name,
3   last_name,
4   state
5 FROM
6   sampdb.president
7 WHERE
8   (last_name, state)
9     IN
10    (
11      ('Adams', 'MA')
12    )
13 ;
14
```

| 1 | first_name | last_name | state |
|---|-------------|-----------|-------|
| 2 | John | Adams | MA |
| 3 | John Quincy | Adams | MA |
| 4 | | | |
| 5 | | | |
| 6 | | | |
| 7 | | | |

- And we can filter by tuple

Combining Conditions

```
1 SELECT
2   first_name,
3   last_name,
4   state
5 FROM
6   sampdb.president
7 WHERE
8   state
9     IN
10    (
11      'MA',
12      'CA',
13      'OH'
14    )
15
16   AND
17   birth > '1900-01-01'
18 ORDER BY
19   state
20 ;
21
```

| | first_name | last_name | state |
|---|-------------|-----------|-------|
| 4 | Richard M. | Nixon | CA |
| 5 | George H.W. | Bush | MA |
| 6 | John F. | Kennedy | MA |

- We can combine conditions

Practice

Female students

Students with scores for test event 3

Students with scores for test event 3, with scores greater than 75

Presidents alive during the Civil War

Presidents alive during the Second World War

Record Level Functions

CONCAT() combines strings

```
1 SELECT
2     first_name,
3     last_name,
4     CONCAT(first_name, ' ', last_name)
5 FROM
6     sampdb.president
7 WHERE
8     state = 'VA'
9 ;
10
```

- CONCAT() combines two strings into one
- Unlike aggregation function (e.g. COUNT() et al), record functions operate only on the values of each particular record

| 1 | 2 | 3 | 4 |
|----|------------|------------|------------------------------------|
| | first_name | last_name | CONCAT(first_name, ' ', last_name) |
| 4 | William H. | Harrison | William H. Harrison |
| 5 | Thomas | Jefferson | Thomas Jefferson |
| 6 | James | Madison | James Madison |
| 7 | James | Monroe | James Monroe |
| 8 | Zachary | Taylor | Zachary Taylor |
| 9 | John | Tyler | John Tyler |
| 10 | George | Washington | George Washington |
| 11 | Woodrow | Wilson | Woodrow Wilson |
| 12 | | | |
| 13 | | | |

CONCAT() doesn't need _field_ values

```
1 SELECT
2     last_name,
3     state,
4     CONCAT('Mickey', ' ', 'Mouse')
5 FROM
6     sampdb.president
7 WHERE
8     state = 'VA'
9 ;
10
```

- Strings passed to CONCAT can be hard-coded values
- Note that the call on CONCAT() produces the same output for each record
- And, that there is one record output for each occurrence record meeting the filter criterion

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|----|------------|-------|--------------------------------|---|---|---|---|---|----|----|----|----|
| | last_name | state | CONCAT('Mickey', ' ', 'Mouse') | | | | | | | | | |
| 4 | Harrison | VA | Mickey Mouse | | | | | | | | | |
| 5 | Jefferson | VA | Mickey Mouse | | | | | | | | | |
| 6 | Madison | VA | Mickey Mouse | | | | | | | | | |
| 7 | Monroe | VA | Mickey Mouse | | | | | | | | | |
| 8 | Taylor | VA | Mickey Mouse | | | | | | | | | |
| 9 | Tyler | VA | Mickey Mouse | | | | | | | | | |
| 10 | Washington | VA | Mickey Mouse | | | | | | | | | |
| 11 | Wilson | VA | Mickey Mouse | | | | | | | | | |

We use AS to create 'alias' labels

```
1 SELECT
2     first_name,
3     last_name,
4     CONCAT(first_name, ' ', last_name) AS full_name
5 FROM
6     sampdb.president
7 WHERE
8     state = 'VA'
9 ;
10
```

- This gives a cleaner name to the field created by the function output

| | first_name | last_name | full_name |
|----|------------|------------|---------------------|
| 4 | William H. | Harrison | William H. Harrison |
| 5 | Thomas | Jefferson | Thomas Jefferson |
| 6 | James | Madison | James Madison |
| 7 | James | Monroe | James Monroe |
| 8 | Zachary | Taylor | Zachary Taylor |
| 9 | John | Tyler | John Tyler |
| 10 | George | Washington | George Washington |
| 11 | Woodrow | Wilson | Woodrow Wilson |
| 12 | | | |

Field values can be simple terms in expressions

```
1 SELECT
2     event_id, student_id,
3     score,
4     score * 2 AS doubled
5 FROM
6     sampdb.score
7 LIMIT 10
8 ;
9
10
```

- Here, for each record output, the value of the score field is simply doubled

```
1 +-----+-----+-----+-----+
2 | event_id | student_id | score | doubled |
3 +-----+-----+-----+-----+
4 | 1 | 1 | 20 | 40 |
5 | 1 | 3 | 20 | 40 |
6 | 1 | 4 | 18 | 36 |
7 | 1 | 5 | 13 | 26 |
8 | 1 | 6 | 18 | 36 |
9 | 1 | 7 | 14 | 28 |
10 | 1 | 8 | 14 | 28 |
11 | 1 | 9 | 11 | 22 |
12 | 1 | 10 | 19 | 38 |
13 | 1 | 11 | 18 | 36 |
14 +-----+-----+-----+-----+
15
```

DATE types permit date arithmetic

```
1 SELECT
2     last_name,
3     birth,
4     death,
5     DATEDIFF(death, birth) AS days,
6     DATEDIFF(death, birth) / 365 AS years
7 FROM
8     sampdb.president
9 WHERE
10    death IS NOT NULL
11 LIMIT 10
12 ;
13
```

- DATEDIFF() gives the number of days between two dates
- We can use its value as an expression term to roughly calculate age in years

| | last_name | birth | death | days | years |
|----|------------|------------|------------|-------|---------|
| 4 | Adams | 1735-10-30 | 1826-07-04 | 33119 | 90.7370 |
| 5 | Adams | 1767-07-11 | 1848-02-23 | 29446 | 80.6740 |
| 6 | Arthur | 1829-10-05 | 1886-11-18 | 20863 | 57.1589 |
| 7 | Buchanan | 1791-04-23 | 1868-06-01 | 28163 | 77.1589 |
| 8 | Cleveland | 1837-03-18 | 1908-06-24 | 26030 | 71.3151 |
| 9 | Coolidge | 1872-07-04 | 1933-01-05 | 22099 | 60.5452 |
| 10 | Eisenhower | 1890-10-14 | 1969-03-28 | 28654 | 78.5041 |
| 11 | Fillmore | 1800-01-07 | 1874-03-08 | 27088 | 74.2137 |
| 12 | Ford | 1913-07-14 | 2006-12-26 | 34133 | 93.5151 |
| 13 | Garfield | 1831-11-19 | 1881-09-19 | 18202 | 49.8685 |
| 14 | | | | | |
| 15 | | | | | |

String functions: Return text and text positions

```
1 SELECT LOCATE(' ', 'cat on the mat') AS locate_example;
2
3 SELECT LEFT('cat on the mat', 3) AS left_example;
4
5 SELECT LEFT(
6   'cat on the mat',
7   LOCATE(' ', 'cat on the mat')
8 ) AS nested_example;
9
10 SELECT MID('cat on the mat', 5, 2) AS mid_example;
11
12
```

```
1 +-----+
2 | locate_example |
3 +-----+
4 | 4 |
5 +-----+
6 +-----+
7 | left_example |
8 +-----+
9 | cat |
10 +-----+
11 +-----+
12 | nested_example |
13 +-----+
14 | cat |
15 +-----+
16 +-----+
17 | mid_example |
18 +-----+
19 | on |
20 +-----+
21
```

Use string functions to parse text

```
1 SELECT
2     first_name,
3     LOCATE(' ', first_name) AS first_space,
4     LEFT(first_name, LOCATE(' ', first_name) -1 ) AS just_first,
5     LENGTH(first_name) AS string_length,
6     mid(
7         first_name,
8         LOCATE(' ', first_name),
9         LENGTH(first_name) - LOCATE(' ', first_name) + 1
10    ) AS middle_name
11 FROM
12     sampdb.president
13 ;
14
```

Parsed output

| 1 | first_name | first_space | just_first | string_length | middle_name |
|----|-------------|-------------|------------|---------------|-------------|
| 2 | George | 0 | | 6 | |
| 3 | John | 0 | | 4 | |
| 4 | Thomas | 0 | | 6 | |
| 5 | James | 0 | | 5 | |
| 6 | James | 0 | | 5 | |
| 7 | John Quincy | 5 | John | 11 | Quincy |
| 8 | Andrew | 0 | | 6 | |
| 9 | Martin | 0 | | 6 | |
| 10 | William H. | 8 | William | 10 | H. |
| 11 | John | 0 | | 4 | |
| 12 | James K. | 6 | James | 8 | K. |
| 13 | Zachary | 0 | | 7 | |
| 14 | ... | | | | |
| 15 | | | | | |
| 16 | | | | | |
| 17 | | | | | |
| 18 | | | | | |

Use Logic Functions to handle edge cases

```
1 SELECT
2     first_name,
3     LOCATE(' ', first_name) AS first_space,
4     IF(
5         LOCATE(' ', first_name) = 0,
6         first_name,
7         LEFT(first_name, LOCATE(' ', first_name) -1 )
8     ) as proper_first,
9     LEFT(first_name, LOCATE(' ', first_name) -1 ) AS just_first,
10    mid(
11        first_name,
12        LOCATE(' ', first_name),
13        LENGTH(first_name) - LOCATE(' ', first_name) + 1
14    ) AS middle_name
15 FROM
16     sampdb.president
17 ;
18 ;
```

Parsed output

| 1 | first_name | first_space | proper_first | just_first | middle_name |
|----|-------------|-------------|--------------|------------|-------------|
| 2 | George | 0 | George | | |
| 3 | John | 0 | John | | |
| 4 | Thomas | 0 | Thomas | | |
| 5 | James | 0 | James | | |
| 6 | James | 0 | James | | |
| 7 | John Quincy | 5 | John | John | Quincy |
| 8 | Andrew | 0 | Andrew | | |
| 9 | Martin | 0 | Martin | | |
| 10 | William H. | 8 | William | William | H. |
| 11 | John | 0 | John | | |
| 12 | James K. | 6 | James | James | K. |
| 13 | Zachary | 0 | Zachary | | |
| 14 | ... | | | | |
| 15 | | | | | |
| 16 | | | | | |
| 17 | | | | | |
| 18 | | | | | |

Practice Queries

Query returning just the first name of each president

Query counting presidents by month of birth. (The function expression MONTH(birth) will return the month of that field.)

Query returning just the 'proper' first name of each president

Query counting presidents by proper first name

The initials of each president

Query counting presidents by last initial

Functions and expressions can filter

```
1 SELECT
2     last_name,
3     birth,
4     death,
5     DATEDIFF(death, birth) AS days,
6     DATEDIFF(death, birth) / 365 AS years
7 FROM
8     sampdb.president
9 WHERE
10    death IS NOT NULL
11    -- fails:
12    -- AND years < 70
13    -- works:
14    AND DATEDIFF(death, birth) / 365 < 70
15 LIMIT 10
16 ;
17
```

- Notice that we can't use the alias 'years'
- We can use that same expression, though

| | last_name | birth | death | days | years |
|---|-----------|------------|------------|-------|---------|
| 4 | Arthur | 1829-10-05 | 1886-11-18 | 20863 | 57.1589 |
| 5 | Coolidge | 1872-07-04 | 1933-01-05 | 22099 | 60.5452 |

Logical Functions

CASE is the SQL standard

This works in both MySQL and Postgres

```
1 SELECT
2     CASE WHEN 2 > 1
3             THEN 'it is true'
4             ELSE 'it is false'
5     END AS case_example
6 ;
7
```

```
1 +-----+
2 | case_example |
3 +-----+
4 | it is true   |
5 +-----+
6
```

MySQL provides a more familiar IF()

This won't work in Postgres

You could define a Postgres function for it

```
1 SELECT
2     IF( 2 > 1,
3         'It is true',
4         'It is false'
5     ) AS if_example
6 ;
7 -- IF() function not available in postgres?
8
```

```
1 +-----+
2 | if_example |
3 +-----+
4 | It is true |
5 +-----+
6
```

CASE statements allow multiple conditions

We needn't nest CASE, we can add conditions

```
1 SELECT
2     CASE
3         WHEN 10 <= 5
4             THEN 'less than five'
5         WHEN 10 > 5 AND 10 < 8
6             THEN 'between five and eight'
7         WHEN 10 < 10
8             THEN 'between eight and ten'
9         ELSE 'it is ten or greater'
10    END AS case_example
11 ;
12
```

```
1 +-----+
2 | case_example
3 +-----+
4 | it is ten or greater |
5 +-----+
6
```

Combining Tables -- JOIN

Finding Megan's scores

Find Megan's student_id

```
1 SELECT
2     name
3     , student_id
4 FROM
5     sampdb.student
6 WHERE
7     name = 'Megan'
8 ;
9
```

```
1 +-----+-----+
2 | name | student_id |
3 +-----+-----+
4 | Megan |           1 |
5 +-----+-----+
6
```

Get the scores for that student_id

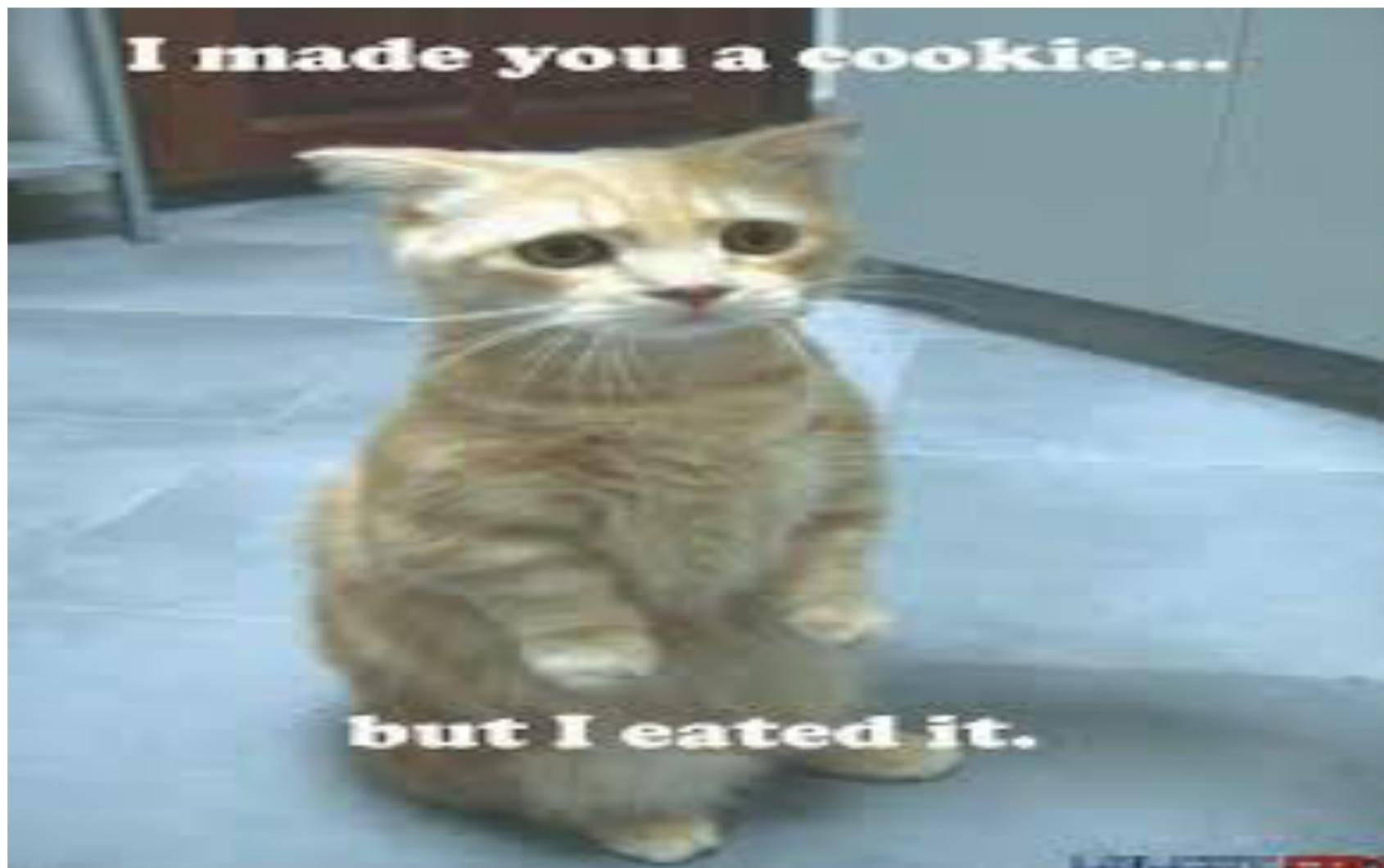
```
1 SELECT
2     student_id
3     , event_id
4     , score
5
6 FROM
7     sampdb.score
8
9 WHERE
10    student_id = 1
11 ;
12
```

| | student_id | event_id | score |
|----|------------|----------|-------|
| 4 | 1 | 1 | 20 |
| 5 | 1 | 2 | 17 |
| 6 | 1 | 3 | 88 |
| 7 | 1 | 5 | 15 |
| 8 | 1 | 6 | 100 |
| 9 | | | |
| 10 | | | |

But we would rather have the name on that table

```
1 SELECT
2     sampdb.student.name
3     , sampdb.student.student_id
4     , sampdb.score.event_id
5     , sampdb.score.score
6 FROM
7     sampdb.student
8 INNER JOIN
9     sampdb.score
10 ON
11     sampdb.student.student_id =
12         sampdb.score.student_id
13 WHERE
14     sampdb.student.name = 'Megan'
15 ;
16
```

| 1 | name | student_id | event_id | score |
|----|-------|------------|----------|-------|
| 2 | Megan | 1 | 1 | 20 |
| 3 | Megan | 1 | 2 | 17 |
| 4 | Megan | 1 | 3 | 88 |
| 5 | Megan | 1 | 5 | 15 |
| 6 | Megan | 1 | 6 | 100 |
| 7 | | | | |
| 8 | | | | |
| 9 | | | | |
| 10 | | | | |



Two Really Simple Tables

| | | | |
|---|----|---|----|
| 1 | i1 | | c1 |
| 2 | | | |
| 3 | | | |
| 4 | 1 | a | |
| 5 | 2 | b | |
| 6 | 3 | c | |
| 7 | | | |
| 8 | | | |

| | | | |
|---|----|---|----|
| 1 | i2 | | c2 |
| 2 | | | |
| 3 | | | |
| 4 | | | |
| 5 | 2 | c | |
| 6 | 3 | b | |
| 7 | 4 | a | |
| 8 | | | |

Inner Join

```
1 SELECT * FROM join_sample.t1 INNER JOIN join_sample.t2;  
2
```

```
1 +-----+-----+-----+-----+  
2 | i1 | c1 | i2 | c2 |  
3 +-----+-----+-----+-----+  
4 | 1 | a | 2 | c |  
5 | 2 | b | 2 | c |  
6 | 3 | c | 2 | c |  
7 | 1 | a | 3 | b |  
8 | 2 | b | 3 | b |  
9 | 3 | c | 3 | b |  
10 | 1 | a | 4 | a |  
11 | 2 | b | 4 | a |  
12 | 3 | c | 4 | a |  
13 +-----+-----+-----+-----+  
14
```

Limiting Inner Join Results

```
1 SELECT
2     *
3 FROM
4     join_sample.t1
5 INNER JOIN
6     join_sample.t2
7 WHERE
8     join_sample.t1.i1 = join_sample.t2.i2
9 ;
10
```

Left Join

Get all rows from the "left" table, each with that row from the "right" table that matches on the specified fields

```
1 SELECT
2   *
3 FROM
4   join_sample.t1
5 LEFT JOIN
6   join_sample.t2
7   ON
8     join_sample.t1.i1 = join_sample.t2.i2
9 ;
10
```

| | i1 | c1 | i2 | c2 |
|---|----|----|------|------|
| 4 | 1 | a | NULL | NULL |
| 5 | 2 | b | 2 | c |
| 6 | 3 | c | 3 | b |

Breaking down the query

A join combines rows from different tables

The "ON" clause specifies conditions for combining records

Here, records with the same student_id value are combined into a single record

The SELECT clause still specifies which fields are displayed from that combined record

INNER JOIN specifies which table to join data from

The addition of "scr" and "st" in the specification of tables provides alias values for reference by the other clauses

Without these the query wouldn't know which table it should find a field in

```
1 SELECT
2     st.name AS           Name,
3     scr.student_id AS   Id,
4     AVG(scr.score) AS   Average,
5     COUNT(scr.score) AS "# Tests"
6 FROM
7     sampdb.score scr
8 INNER JOIN
9     sampdb.student st
10    ON
11        scr.student_id = st.student_id
12 GROUP BY
13        scr.student_id
14 ORDER BY
15        Average DESC
16 LIMIT 5
17 ;
18
```

Practice

Join all rows in sampdb.student with all rows in sampdb.score

As above, but filter to match each score to the proper student

Output employee name and their sales

More on JOINS

Data is spread over several tables

Title Data

| 1 | Field | Type | Null | Key | Default | Extra |
|----|------------|--------------|------|-----|---------|-------|
| 2 | title_id | char(3) | NO | PRI | NULL | |
| 3 | title_name | varchar(40) | NO | | NULL | |
| 4 | type | varchar(10) | YES | | NULL | |
| 5 | pub_id | char(3) | NO | | NULL | |
| 6 | pages | int(11) | YES | | NULL | |
| 7 | price | decimal(5,2) | YES | | NULL | |
| 8 | sales | int(11) | YES | | NULL | |
| 9 | pubdate | date | YES | | NULL | |
| 10 | contract | smallint(6) | NO | | NULL | |
| 11 | | | | | | |
| 12 | | | | | | |
| 13 | | | | | | |
| 14 | | | | | | |

Title and Author do NOT join directly

Authors

| 1 | Field | Type | Null | Key | Default | Extra |
|----|----------|-------------|------|-----|---------|-------|
| 2 | au_id | char(3) | NO | PRI | NULL | |
| 3 | au_fname | varchar(15) | NO | | NULL | |
| 4 | au_lname | varchar(15) | NO | | NULL | |
| 5 | phone | varchar(12) | YES | | NULL | |
| 6 | address | varchar(20) | YES | | NULL | |
| 7 | city | varchar(15) | YES | | NULL | |
| 8 | state | char(2) | YES | | NULL | |
| 9 | zip | char(5) | YES | | NULL | |
| 10 | | | | | | |
| 11 | | | | | | |
| 12 | | | | | | |
| 13 | | | | | | |

Title and Authors

| 1 | Field | Type | Null | Key | Default | Extra |
|---|---------------|--------------|------|-----|---------|-------|
| 2 | title_id | char(3) | NO | PRI | NULL | |
| 3 | au_id | char(3) | NO | PRI | NULL | |
| 4 | au_order | smallint(6) | NO | | NULL | |
| 5 | royalty_share | decimal(5,2) | NO | | NULL | |
| 6 | | | | | | |
| 7 | | | | | | |
| 8 | | | | | | |
| 9 | | | | | | |

- Notice that this table serves to connect authors to their titles

Getting books written

```
1 -- List books each author wrote
2 -- Fehily 7.7
3
4 SELECT
5   a.au_id,
6   a.au_fname,
7   a.au_lname,
8   ta.title_id,
9   ta.au_id AS ta_author_id
10 FROM
11   books.authors a
12 INNER JOIN
13   books.title_authors ta
14   ON a.au_id = ta.au_id
15 ORDER BY
16   a.au_id ASC,
17   ta.title_id ASC
18 ;
19
```

| | au_id | au_fname | au_lname | title_id | ta_author_id |
|----|-------|-----------|-----------|----------|--------------|
| 4 | A01 | Sarah | Buchman | T01 | A01 |
| 5 | A01 | Sarah | Buchman | T02 | A01 |
| 6 | A01 | Sarah | Buchman | T13 | A01 |
| 7 | A02 | Wendy | Heydemark | T06 | A02 |
| 8 | A02 | Wendy | Heydemark | T07 | A02 |
| 9 | A02 | Wendy | Heydemark | T10 | A02 |
| 10 | A02 | Wendy | Heydemark | T12 | A02 |
| 11 | A03 | Hallie | Hull | T04 | A03 |
| 12 | A03 | Hallie | Hull | T11 | A03 |
| 13 | A04 | Klee | Hull | T04 | A04 |
| 14 | A04 | Klee | Hull | T05 | A04 |
| 15 | A04 | Klee | Hull | T07 | A04 |
| 16 | A04 | Klee | Hull | T11 | A04 |
| 17 | A05 | Christian | Kells | T03 | A05 |
| 18 | A06 | | Kellsey | T08 | A06 |
| 19 | A06 | | Kellsey | T09 | A06 |
| 20 | A06 | | Kellsey | T11 | A06 |
| 21 | | | | | |
| 22 | | | | | |

Getting books titles

```
1 SELECT
2     t.title_id,
3     t.title_name,
4     ta.title_id
5 FROM
6     books.titles t
7 INNER JOIN
8     books.title_authors ta
9 ON
10    t.title_id = ta.title_id
11 ;
12
```

| | title_id | title_name | title_id |
|----|----------|-------------------------------------|----------|
| 4 | T01 | 1977 | T01 |
| 5 | T02 | 200 Years of German Humor | T02 |
| 6 | T03 | Ask Your System Administrator | T03 |
| 7 | T04 | But I Did It Unconsciously | T04 |
| 8 | T04 | But I Did It Unconsciously | T04 |
| 9 | T05 | Exchange of Platitudes | T05 |
| 10 | T06 | How About Never? | T06 |
| 11 | T07 | I Blame My Mother | T07 |
| 12 | T07 | I Blame My Mother | T07 |
| 13 | T08 | Just Wait Until After School | T08 |
| 14 | T09 | Kiss My Boo-Boo | T09 |
| 15 | T10 | Not Without My Faberge Egg | T10 |
| 16 | T11 | Perhaps It's a Glandular Problem | T11 |
| 17 | T11 | Perhaps It's a Glandular Problem | T11 |
| 18 | T11 | Perhaps It's a Glandular Problem | T11 |
| 19 | T12 | Spontaneous, Not Annoying | T12 |
| 20 | T13 | What Are The Civilian Applications? | T13 |
| 21 | | | |
| 22 | | | |

A multi-table JOIN connects these in one query

Two JOINs in one query

```
1 SELECT
2   ttl.title_name,
3   -- ttl.title_id,
4   -- ta.au_id,
5   auth.au_fname,
6   auth.au_lname
7
8 FROM
9   books.titles AS ttl
10 INNER JOIN
11   books.title_authors AS ta
12 INNER JOIN
13   books.authors AS auth
14 WHERE
15   ttl.title_id = ta.title_id
16   AND ta.au_id = auth.au_id
17 ;
18
19
```

- Remember, INNER JOIN makes a 'super table'
- Here, this table combines three base tables
- Of that table, rows are returns on matching title and author ids
- Note that no fields from title_authors are displayed
- title_authors here serves to 'bridge' the one table to another
- This structure allows the data model to handle authors with multiple titles, and titles with multiple authors

Title and Author in one table

```
1 SELECT
2     ttl.title_name,
3     auth.au_fname,
4     auth.au_lname
5
6 FROM
7     books.titles AS ttl
8
9 INNER JOIN
10    books.title_authors AS ta
11 ON
12    ttl.title_id = ta.title_id
13
14 INNER JOIN
15    books.authors AS auth
16 ON
17    ta.au_id = auth.au_id
18
19 ;
20
21
```

Same query, using ON rather than WHERE

```
1 SELECT
2   ttl.title_name,
3   auth.au_fname,
4   auth.au_lname
5
6 FROM
7   books.titles AS ttl
8
9 INNER JOIN
10  books.title_authors AS ta
11 ON
12  ttl.title_id = ta.title_id
13
14 INNER JOIN
15  books.authors AS auth
16 ON
17  ta.au_id = auth.au_id
18
19 ;
20
21
```

We can join a table to itself

Accesses a table's data in reference to itself

```
1 SELECT
2     au_1.au_id,
3     au_1.au_fname,
4     au_1.au_lname,
5     au_1.state
6 FROM
7     books.authors au_1
8 INNER JOIN
9     books.authors au_2
10    ON
11        au_1.state = au_2.state
12 WHERE
13     au_2.au_id = 'A04'
14
15 ;
16
```

| au_id | au_fname | au_lname | state |
|-------|----------|----------|-------|
| A03 | Hallie | Hull | CA |
| A04 | Klee | Hull | CA |
| A06 | | Kellsey | CA |

Self joins can enable hierarchy extraction

Consider a table of employee / boss data

| 1 | emp_id | emp_name | boss_id |
|----|--------|-------------------|---------|
| 2 | E01 | Lord Copper | NULL |
| 3 | E02 | Jocelyn Hitchcock | E01 |
| 4 | E03 | Mr. Salter | E01 |
| 5 | E04 | William Boot | E03 |
| 6 | E05 | Mr. Corker | E03 |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |

- The table is self-referential, in that each boss_id value is found in employee_id
- We can use a self-join to get a table of who reports to whom

Hierarchy extraction (cont)

```
1 SELECT
2   e1.emp_name AS "Employee Name",
3   e2.emp_name AS "Boss Name",
4   e1.boss_id,
5   e2.emp_id
6 FROM
7   books.employees e1
8 INNER JOIN
9   books.employees e2
10 ON
11   e1.boss_id = e2.emp_id
12 ;
13
```

| | Employee Name | Boss Name |
|---|-------------------|-------------|
| 4 | Jocelyn Hitchcock | Lord Copper |
| 5 | Mr. Salter | Lord Copper |
| 6 | William Boot | Mr. Salter |
| 7 | Mr. Corker | Mr. Salter |

Here, the ON condition selects those rows that match an employee with their boss

Practice Problems

From books, employee name with employee sales

Subqueries -- Referring to SQL results

Subqueries embed one query within another

These embedded queries are 'subqueries'

We have to watch the table returned by the subquery carefully

The returned table has to have the fields and row expected of it by the calling query

Two varieties, 'correlated' and 'uncorrelated'

'correlated' is cooler but harder, we'll do 'uncorrelated' first

Get names of publishers of biography

Step one: Get pub_id of the publishers

```
1 SELECT
2     pub_id
3 FROM
4     books.titles
5 WHERE
6     type = 'biography'
7 ;
8
9
```

| pub_id |
|--------|
| P01 |
| P03 |
| P01 |
| P01 |

Get names of publishers of biography -- cont.

Step two: Get names of those publishers

```
1 SELECT
2     pub_name
3 FROM
4     books.publishers
5 WHERE
6     pub_id IN ('P01', 'P03')
7 ;
8
9
```

```
1 +-----+
2 | pub_name |
3 +-----+
4 | Abatis Publishers |
5 | Schadenfreude Press |
6 +-----+
7
```

Get names of publishers of biography -- cont.

But we can do this in one query

```
1 SELECT
2     pub_name
3 FROM
4     books.publishers
5 WHERE
6     pub_id IN (
7         SELECT
8             pub_id
9             FROM
10            books.titles
11            WHERE
12                type = 'biography'
13        )
14 ;
15
16
```

```
1 +-----+
2 | pub_name |
3 +-----+
4 | Abatis Publishers |
5 | Schadenfreude Press |
6 +-----+
7
```

Another example

Get authors who have not written a book

```
1 SELECT
2     au_id, au_fname, au_lname
3 FROM
4     books.authors
5 WHERE
6     au_id NOT IN
7     (SELECT au_id FROM books.title_authors)
8 ;
9
```

| | au_id | au_fname | au_lname |
|---|-------|----------|-------------|
| 1 | A07 | Paddy | O'Furniture |
| 2 | | | |
| 3 | | | |
| 4 | | | |
| 5 | | | |
| 6 | | | |

More complex example

Say we want a list of students scoring above on grade event 3

Getting the score is easy:

```
1 SELECT
2     event_id,
3     AVG(score)
4 FROM
5     sampdb.score
6 WHERE
7     event_id = 3
8 GROUP BY
9     event_id
10 LIMIT 5
11 ;
12
```

Filtering is easy:

```
1 SELECT
2     student_id AS student,
3     event_id AS event,
4     score
5 FROM
6     sampdb.score
7 WHERE
8     event_id = 3
9     AND score > 78.2258
10 ;
11
```

Results

Results

| 1 | student | event | score |
|----|---------|-------|-------|
| 2 | 1 | 3 | 88 |
| 3 | 2 | 3 | 84 |
| 4 | 5 | 3 | 97 |
| 5 | 6 | 3 | 83 |
| 6 | 7 | 3 | 88 |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |

But:

What if the data change?

We want results for another event?

Better: refer directly to the average

```
1 SELECT
2     student_id,
3     event_id,
4     score
5 FROM
6     sampdb.score
7 WHERE
8     event_id = 3
9     AND score > (
10         SELECT
11             AVG(score)
12         FROM
13             sampdb.score
14         WHERE
15             event_id = 3
16         GROUP BY
17             event_id
18     )
19 LIMIT 5;
20
```

| | student_id | score |
|----|------------|-------|
| 4 | 1 | 88 |
| 5 | 2 | 84 |
| 6 | 5 | 97 |
| 7 | 6 | 83 |
| 8 | 7 | 88 |
| 9 | | |
| 10 | | |

Now:

- query automatically changes output for any data change
- we can get a similar list simply by changing the event_id we are seeking
- BUT -- what if we change event_id in one place, and not the other . . . ?

Better still: Remove duplication

```
1 SELECT
2   student_id,
3   score
4 FROM
5   sampdb.score AS scr
6 WHERE
7   event_id = 3
8   AND score > (
9     SELECT
10    AVG(score)
11   FROM
12     sampdb.score
13   WHERE
14     event_id=scr.event_id
15   GROUP BY
16     event_id
17   )
18 LIMIT 5;
19
```

| | student_id | score |
|----|------------|-------|
| 4 | 1 | 88 |
| 5 | 2 | 84 |
| 6 | 5 | 97 |
| 7 | 6 | 83 |
| 8 | 7 | 88 |
| 9 | | |
| 10 | | |

Now:

- The subquery's event_id filter looks to each particular record for the value of event_id used in its filter
- Thus our query only refers to the desired event_id once
- e.g., we cannot make the 'event_id values out of step' mistake
- This is a 'correlated subquery', and we will spend more time on them later

Correlated queries send values from outer query to inner

```
1 -- List the books that have sales greater than or
2 -- equal to the average sales for books of its type
3 -- feinly 8.13
4
5 SELECT
6   candidate.title_id,
7   candidate.type,
8   candidate.sales,
9   (SELECT ROUND(AVG(sales), 1)
10  FROM books.titles average
11  WHERE average.type = candidate.type) AS type_average
12 FROM
13   books.titles candidate
14 WHERE
15   sales >=
16   (SELECT AVG(sales)
17  FROM books.titles average
18  WHERE average.type = candidate.type)
19 ;
20
```

| 1 | title_id | type | sales | type_average |
|----|----------|------------|---------|--------------|
| 4 | T02 | history | 9566 | 6866.3 |
| 5 | T03 | computer | 25667 | 25667.0 |
| 6 | T05 | psychology | 201440 | 102854.7 |
| 7 | T07 | biography | 1500200 | 537173.7 |
| 8 | T09 | children | 5000 | 4547.5 |
| 9 | T13 | history | 10467 | 6866.3 |
| 10 | | | | |
| 11 | | | | |

Each score event has its own average and std. dev

```
1 SELECT
2     event_id,
3     AVG(score) AS average,
4     STD(score) AS standard_deviation
5 FROM
6     sampdb.score
7 GROUP BY
8     event_id
9 ;
10
```

| | event_id | average | standard_deviation |
|----|----------|---------|--------------------|
| 4 | 1 | 15.1379 | 3.7849 |
| 5 | 2 | 14.1667 | 3.7424 |
| 6 | 3 | 78.2258 | 9.2063 |
| 7 | 4 | 14.0370 | 3.8344 |
| 8 | 5 | 14.1852 | 3.3448 |
| 9 | 6 | 80.1724 | 12.4404 |
| 10 | | | |
| 11 | | | |

Each student_id in score has multiple events

```
1 SELECT
2     student_id,
3     COUNT(event_id) as events
4 FROM
5     sampdb.score
6 GROUP BY
7     student_id
8 LIMIT 10;
9
10
```

| | student_id | events |
|----|------------|--------|
| 4 | 1 | 5 |
| 5 | 2 | 5 |
| 6 | 3 | 6 |
| 7 | 4 | 5 |
| 8 | 5 | 6 |
| 9 | 6 | 6 |
| 10 | 7 | 6 |
| 11 | 8 | 6 |
| 12 | 9 | 6 |
| 13 | 10 | 6 |
| 14 | | |
| 15 | | |

We can find statistics for each of those events

```
1 SELECT
2   student_id,
3   event_id,
4   score,
5   (  SELECT ROUND(AVG(score), 1)
6     FROM sampdb.score AS s2
7     WHERE s1.event_id = s2.event_id
8     GROUP BY event_id
9   ) AS e_avg,
10  (  SELECT ROUND(STD(score), 1)
11    FROM sampdb.score AS s2
12    WHERE s1.event_id = s2.event_id
13    GROUP BY event_id
14  ) AS e_dev
15 FROM
16   sampdb.score AS s1
17 WHERE
18   student_id = 1
19 ;
20
21
22
```

| | student_id | event_id | score | e_avg | e_dev |
|----|------------|----------|-------|-------|-------|
| 4 | 1 | 1 | 20 | 15.1 | 3.8 |
| 5 | 1 | 2 | 17 | 14.2 | 3.7 |
| 6 | 1 | 3 | 88 | 78.2 | 9.2 |
| 7 | 1 | 5 | 15 | 14.2 | 3.3 |
| 8 | 1 | 6 | 100 | 80.2 | 12.4 |
| 9 | | | | | |
| 10 | | | | | |

Get those statistics, for every event, for every record

```
1  SELECT
2    score, student_id AS stdnt, event_id AS vnt,
3    event_average AS vnt_avg, std_dev AS std,
4    ( score - event_average ) / std_dev AS z_score
5  FROM
6  (
7    SELECT
8      s1.score, s1.student_id, s1.event_id,
9      ( SELECT AVG(score)
10        FROM sampdb.score AS s2
11        WHERE s2.event_id = s1.event_id
12        GROUP BY event_id
13      ) AS event_average,
14      ( SELECT STD(score)
15        FROM sampdb.score AS s2
16        WHERE s2.event_id = s1.event_id
17        GROUP BY event_id
18      ) AS std_dev
19      FROM sampdb.score AS s1
20  ) AS tmp
21
22 LIMIT 5
23 ;
24
```

Result

| 1 | score | stdnt | vnt | vnt_avg | std | z_score |
|----|-------|-------|-----|---------|--------|-------------|
| 2 | 20 | 1 | 1 | 15.1379 | 3.7849 | 1.28458961 |
| 3 | 20 | 3 | 1 | 15.1379 | 3.7849 | 1.28458961 |
| 4 | 18 | 4 | 1 | 15.1379 | 3.7849 | 0.75618024 |
| 5 | 13 | 5 | 1 | 15.1379 | 3.7849 | -0.56484320 |
| 6 | 18 | 6 | 1 | 15.1379 | 3.7849 | 0.75618024 |
| 7 | | | | | | |
| 8 | | | | | | |
| 9 | | | | | | |
| 10 | | | | | | |

Correlated Subqueries work record by record

The 'inner' queries are run for each record in the 'outer' query

Definitions

- 'Outer' query is the query calling on result
- 'Inner' query supplies result
- Records in separate are referenced by aliases

Outer query checks inner query for each of its own records

1. Outer gets record
2. Outer consults inner, 'correlated' query
3. Inner query takes values from record in process for outer query
4. Inner passes result out to outer query

Are duplicative inner queries cached?

- It depends
- Performance is complex, requires first mastering syntax

A convoluted problem:

Say we want a list of those presidents from the three states that have sent the most presidents

We can get the list by a query:

```
1 SELECT
2     state,
3     COUNT(*) AS state_count
4 FROM
5     sampdb.president
6 GROUP BY
7     state
8 ORDER BY
9     state_count DESC
10 LIMIT 3
11 ;
12
```

| state | state_count |
|-------|-------------|
| VA | 8 |
| OH | 7 |
| MA | 4 |

First, get criteria right

Revise our condition for state inclusion

```
1 SELECT
2   MIN(st_count)
3 FROM
4 (
5   SELECT
6     state,
7     COUNT(*) AS st_count
8   FROM
9     sampdb.president
10  GROUP BY
11    state
12  ORDER BY
13    state_count DESC
14  LIMIT 3
15 ) AS tmp
16 ;
17
```

| | | | |
|---|---|---------------|---|
| 1 | + | ----- | + |
| 2 | | MIN(st_count) | |
| 3 | + | ----- | + |
| 4 | | 4 | |
| 5 | + | ----- | + |
| 6 | | | |

Revise listing query

Uses presidents count criteria for state inclusion

```
1 SELECT
2     state, COUNT(*) as state_count
3 FROM
4     sampdb.president
5 GROUP BY
6     state
7 HAVING
8     COUNT(*) >=
9     (
10    SELECT
11        MIN(state_count)
12    FROM
13        (
14            SELECT
15                state,
16                COUNT(*) AS state_count
17            FROM
18                sampdb.president
19            GROUP BY
20                state
21            ORDER BY
22                state_count DESC
23            LIMIT 3
24        ) AS t3
25    )
26 ORDER BY state_count DESC
27 ;
28
```

| 1 | state | state_count |
|---|-------|-------------|
| 2 | VA | 8 |
| 5 | OH | 7 |
| 6 | MA | 4 |
| 7 | NY | 4 |

- This gives us a 'dynamic' derivation of the state list we need
- HAVING
- - similar to WHERE
- - appears after GROUP BY
- - refers to, and filters on, the grouped records and aggregated fields
- So WHERE filters the records that are aggregated into groups, and HAVING filters the groups themselves

Putting all this together provides our answer

```
1 SELECT last_name, state
2 FROM sampdb.president
3 WHERE
4   state IN
5   (
6     SELECT state
7     FROM sampdb.president
8     GROUP BY state
9     HAVING
10    COUNT(*) >=
11    (
12      SELECT
13        MIN(state_count)
14      FROM
15      (
16        SELECT COUNT(*) AS state_count
17        FROM sampdb.president
18        GROUP BY state
19        ORDER BY state_count DESC
20        LIMIT 3
21      ) AS t3
22    )
23  )
24 ORDER BY
25  state, last_name, first_name
26 ;
27
```

| | last_name | state |
|----|------------|-------|
| 4 | Adams | MA |
| 5 | Adams | MA |
| 6 | Bush | MA |
| 7 | Kennedy | MA |
| 8 | Fillmore | NY |
| 9 | Roosevelt | NY |
| 10 | Roosevelt | NY |
| 11 | Van Buren | NY |
| 12 | Garfield | OH |
| 13 | Grant | OH |
| 14 | Harding | OH |
| 15 | Harrison | OH |
| 16 | Hayes | OH |
| 17 | McKinley | OH |
| 18 | Taft | OH |
| 19 | Harrison | VA |
| 20 | Jefferson | VA |
| 21 | Madison | VA |
| 22 | Monroe | VA |
| 23 | Taylor | VA |
| 24 | Tyler | VA |
| 25 | Washington | VA |
| 26 | Wilson | VA |
| 27 | | |

- Inner query lists the top three states and their counts
- Next extracts the minimum value of those counts
- Next out lists the states with that count or higher
- Last query takes presidents whose states fall in that list

Subqueries in SELECT field expressions

```
1 SELECT
2 ( SELECT
3     ROUND(AVG(score),1)
4     FROM sampdb.score
5     WHERE event_id = 1
6     GROUP BY event_id
7 ) AS 1_scr,
8 ( SELECT
9     ROUND(AVG(score),1)
10    FROM sampdb.score
11    WHERE event_id = 2
12    GROUP BY event_id
13 ) AS 2_scr,
14 ( SELECT
15     ROUND(AVG(score),1)
16     FROM sampdb.score
17     WHERE event_id = 3
18     GROUP BY event_id
19 ) AS 3_scr
20
21 -- Notice the SELECT
22 -- expression is
23 -- just a field list
24
25 ;
26
27
28
```

| | 1_scr | 2_scr | 3_scr |
|---|-------|-------|-------|
| 4 | 15.1 | 14.2 | 78.2 |
| 5 | | | |
| 6 | | | |

Subqueries can be used in field expressions

- Here we use them to replicate an Excel pivot table
- Notice that this doesn't handle "arbitrary" columns
- We have to designate a field in our output for each column
- The fundamental problem is we can specify sets of records, but not fields
- There are hacks around this, but they're complicated

HAVING

Seen earlier in subquery

HAVING applies criteria to output of aggregations

- Calculated after WHERE and GROUP BY
- WHERE filters input records to initial query
- HAVING filters resultant aggregations

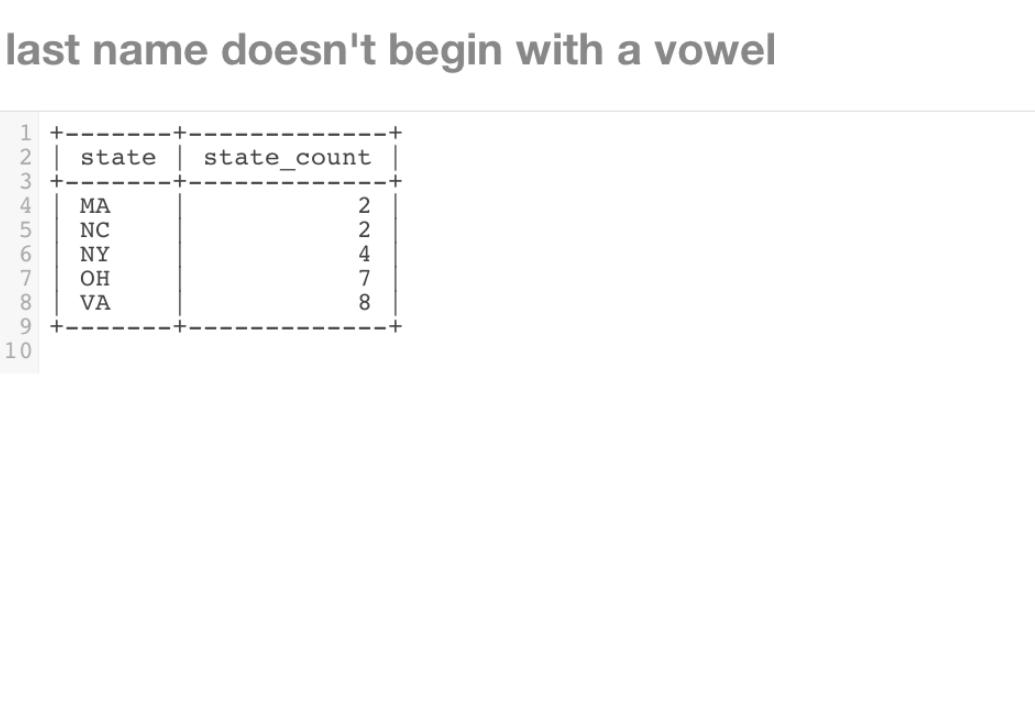
HAVING Example

States with more than one President whose last name doesn't begin with a vowel

```
1 SELECT
2   state,
3   COUNT(*) AS state_count
4 FROM
5   sampdb.president
6 WHERE
7   LEFT(last_name, 1)
8   NOT IN
9   ('A', 'E', 'I',
10  'O', 'U')
11 GROUP BY
12 state
13 HAVING
14   state_count >= 2
15 ;
16
17
```

| state | state_count |
|-------|-------------|
| MA | 2 |
| NC | 2 |
| NY | 4 |
| OH | 7 |
| VA | 8 |

10



- WHERE clause strips out the vowel last names
- HAVING gets the states with more than one (now not-vowel) president
- Splitting up the WHERE and HAVING criteria gives a finer control over the output set

Practice

Revise query to get below average scores for event 2

Employees with above average sales

Presidents who lived longer than average

The age of the third longest lived President

The name of the third longest lived President

Queries for Analysis



Using queries for analysis

1. Notice that all of our queries return (e.g., output) tables

The output is a set of records, with defined fields

2. We can perform further operations on those returned tables, and join them to one another

3. By cascading queries, we can filter and massage our data to get where we want

"views" give us a tool for holding queries "in place"

Create a "base table"

- "Stores" the result of the SELECT query in a "table" available to other queries
- Uses functions to create an "age" column from birth and death

```
1 DROP VIEW IF EXISTS sampdb.president_age;
2 CREATE VIEW
3     sampdb.president_age
4 AS
5     SELECT
6         last_name,
7         state,
8         ROUND(DATEDIFF(death, birth) / 365, 1) as AGE
9     FROM
10        sampdb.president ;
11
12    SELECT * FROM sampdb.president_age LIMIT 5;
13
14
```

| | last_name | state | AGE |
|---|------------|-------|------|
| 4 | Washington | VA | 67.9 |
| 5 | Adams | MA | 90.7 |
| 6 | Jefferson | VA | 83.3 |
| 7 | Madison | VA | 85.3 |
| 8 | Monroe | VA | 73.2 |

Derive a second table from our base

- Uses our prior query to create a table of aggregated values

```
1 DROP VIEW IF EXISTS sampdb.president_aggregates;
2 CREATE VIEW
3     sampdb.president_aggregates
4 AS
5     SELECT
6         state,
7         COUNT(state) as StateCount,
8         AVG(age) as AverageAge,
9         MAX(age) as MaxAge,
10        MIN(age) as MinAge
11    FROM
12        sampdb.president_age
13    GROUP BY
14        state
15    ORDER BY
16        StateCount DESC
17;
18    SELECT * FROM sampdb.president_aggregates LIMIT 5;
19
20
```

| state | StateCount | AverageAge | MaxAge | MinAge |
|-------|------------|------------|--------|--------|
| VA | 8 | 72.83750 | 85.3 | 65.7 |
| OH | 7 | 62.87143 | 72.5 | 49.9 |
| MA | 4 | 72.63333 | 90.7 | 46.5 |
| NY | 4 | 69.32500 | 79.7 | 60.2 |
| TX | 2 | 71.45000 | 78.5 | 64.4 |

Combine the two to get our result

- Now we can join our Age table with our Aggregates table to get oldest presidents by State

```
1 SELECT
2     sub.state,
3     sub.last_name,
4     sub.age
5 FROM
6     sampdb.president_age sub
7 INNER JOIN
8     sampdb.president_aggregates sub2
9 ON
10    sub.state = sub2.state
11    AND sub.age = sub2.MaxAge
12 LIMIT 5
13 ;
14
15
```

| state | last_name | AGE |
|-------|-----------|------|
| MA | Adams | 90.7 |
| VA | Madison | 85.3 |
| SC | Jackson | 78.3 |
| NY | Van Buren | 79.7 |
| NH | Pierce | 64.9 |

Cascading through queries to result

| last_name | first_name | suffix | city | state | birth | death |
|------------|-------------|--------|---------------------|-------|------------|------------|
| Washington | George | NULL | Wakefield | VA | 1732-02-22 | 1799-12-14 |
| Adams | John | NULL | Braintree | MA | 1735-10-30 | 1826-07-04 |
| Jefferson | Thomas | NULL | Albemarle County | VA | 1743-04-13 | 1826-07-04 |
| Madison | James | NULL | Port Conway | VA | 1751-03-16 | 1836-06-28 |
| Monroe | James | NULL | Westmoreland County | VA | 1758-04-28 | 1831-07-04 |
| Adams | John Quincy | III | Braintree | MA | 1767-07-11 | 1848-02-23 |



| last_name | state | age |
|------------|-------|------|
| Washington | VA | 67.9 |
| Adams | MA | 90.7 |
| Jefferson | VA | 83.3 |
| Madison | VA | 85.3 |
| Monroe | VA | 73.2 |
| Adams | MA | 80.7 |
| Jackson | SC | 78.3 |
| Van Buren | NY | 79.7 |
| Harrison | VA | 68.2 |
| Tyler | VA | 71.9 |
| Polk | NC | 53.7 |
| Taylor | VA | 65.7 |
| Fillmore | NY | 74.2 |



| state | StateCount | AverageAge | MaxAge | MinAge |
|-------|------------|------------|--------|--------|
| VA | 8 | 72.83750 | 85.3 | 65.7 |
| OH | 7 | 62.87143 | 72.5 | 49.9 |
| MA | 4 | 72.63333 | 90.7 | 46.5 |
| NY | 4 | 69.32500 | 79.7 | 60.2 |
| TX | 2 | 71.45000 | 78.5 | 64.4 |
| NC | 2 | 59.15000 | 66.6 | 52.7 |



| state | last_name | age |
|-------|-----------|------|
| MA | Adams | 90.7 |
| VA | Madison | 85.3 |
| SC | Jackson | 78.3 |
| NY | Van Buren | 79.7 |
| NH | Pierce | 64.9 |
| PA | Buchanan | 77.2 |



Can we do without the views?

- Here we use a subquery
- We're still combining tables with one another, but here we generate one of them on the fly

```
1 # Gets oldest president, by state, without intervening views
2 # - better than only matching age to max age values b/c it avoids the
3 # possibility of duplicate max age values
4
5 SELECT
6     state,
7     last_name AS oldest,
8     ROUND(DATEDIFF(death, birth) / 365, 1) AS age
9
10 FROM
11     sampdb.president
12
13 WHERE
14     (
15         state,
16         ROUND(DATEDIFF(death, birth) / 365, 1)
17     ) IN
18     (
19         (
20             SELECT
21                 state,
22                 MAX(ROUND(DATEDIFF(death, birth) / 365, 1)) AS age
23             FROM
24                 sampdb.president
25             GROUP BY
26                 state
27     )
28
29 ;
```



Annoying but necessary

1. No fundamental concepts involved here - just file management
2. Understanding this process does orient us to the fundamentals of what the computer is doing
3. We will go over this fast and revisit in exercises

Next week we'll go over problems in those

First, get data to database

1. Get data into CSV format

Excel can do this for us

On Mac:

2. "Comma Separated Values" will have CR line terminators

3. "Windows CSV" will have CRLF terminators

4. Line terminators? Hang on

Create a database

Simple query creating a home for our table

```
1 DROP DATABASE IF EXISTS learning;  
2 CREATE DATABASE learning;  
3
```

Create a table

- Very simple table

Notice homogeneity of data type (all varchar)

```
1 DROP TABLE IF EXISTS learning.from_excel;
2 CREATE TABLE learning.from_excel (
3
4     date varchar(20),
5     category varchar(20),
6     event_id varchar(20)
7 );
8
```

Load data

Let's break this down

```
1 #    Loads a CSV file, saved from Excel with "Comma Separated Values" option
2 #    -      not "Windows" csv, not "MS-DOS" csv
3 #    -      Excel 2008 for Mac, 12.3.5
4
5 #    file command output: ASCII text, with CR line terminators
6
7 LOAD DATA LOCAL
8
9 #    EDIT THIS PATH:
10 INFILe "/Users/gimli/Desktop/from_excel.csv"
11
12 INTO TABLE learning.from_excel
13
14 FIELDS
15     TERMINATED BY ','
16     ENCLOSED BY '\"'
17
18 LINES
19     TERMINATED BY '\n'
20     #    If that gets strange results, comment it out and try this:
21     #-TERMINATED BY '\r\n'
22
23 IGNORE 1 LINES;
24
```

A moment on file data structure

1. Loading CSV to SQL forces us to consider how our data is represented in a file
2. Excel handles a variety of file storage formats for us - it will guess which is the proper one
3. **SQL will not**

It is a power user tool

Stuff like file format has defaults but if these don't work we must figure things out for ourselves

File data structure: Problem

1. The file is a string of bytes - in our simple case, ASCII data

Unicode? Wrong class

2. Some of those bytes are data, some are metadata

How does MySQL know where one field begins and another ends?

Where one record ends another begins?

Solved by convention and definition

1. Fields "delimited", or "terminated", by commas

When MySQL sees a comma, it starts a new field

2. Fields are also "enclosed" by double quotes

When MySQL sees a double quote, it considers all that follows a field until it sees another double quote

3. Not redundant

This structure allows us to have commas in our data

If an enclosure is missing MySQL will just go with the field terminator

Double quotes in data can be handled by escaping mechanism

(e.g. \" is not an enclosure) but we haven't time for that

Solved by convention and definition - 2

1. Lines (e.g. records, rows) terminated by a line termination character(s)

In Windows, generally carriage return + line feed, or CR LF

\r\n

In Mac and Linux, line feed

\n

2. These characters are **VERY ANNOYING** because we almost never see them

Our programs all have facilities for interpreting them for us

3. We can figure this out or just guess

NEVER MIND -- JUST DO THIS

1. Create a database
2. Create a table
3. Save as CSV
4. Run this
5. Check results, if bad fiddle with this and try again

```
1 #      Loads a CSV file, saved from Excel with "Comma Separated Values" option
2 #      - not "Windows" csv, not "MS-DOS" csv
3 #      - Excel 2008 for Mac, 12.3.5
4
5 #      file command output: ASCII text, with CR line terminators
6
7 LOAD DATA LOCAL
8
9 #      EDIT THIS PATH:
10 INFILE "/Users/gimli/Desktop/from_excel.csv"
11
12 INTO TABLE learning.from_excel
13
14 FIELDS
15     TERMINATED BY ','
16     ENCLOSED BY '\"'
17
18 LINES
19     TERMINATED BY '\n'
20     # If that gets strange results, comment it out and try this:
21     #-TERMINATED BY '\r\n'
22
23 IGNORE 1 LINES;
24
```

Data Editing

Properly typing data?

```
1 DROP TABLE IF EXISTS csv_load.temperature;
2
3 CREATE TABLE csv_load.temperature (
4     date      DATE,
5     city      VARCHAR(20),
6     temperature INT(2)
7 );
8
```

```
1 +-----+-----+-----+
2 | date | city | temperature |
3 +-----+-----+-----+
4 | 2011-10-01 | San Francisco | 62 |
5 | 2011-10-01 | New York | 63 |
6 | 2011-10-01 | Austin | 72 |
7 | 2011-10-02 | San Francisco | 59 |
8 | 2011-10-02 | New York | 58 |
9 +-----+-----+-----+
10
```

- VERY similar to the simple.csv load example, we've just added data typing
- Note that the values for DATE are in YYYY-MM-DD format; if they weren't, we'd have to load as CHAR() and do some data cleaning
- Also, we have to designated a size for our CHAR and INT types
- Any potential problems with our INT type definition?

These operations CHANGE BASE DATA

Prior operations were all outputs

Temporary results, generated to our specification, without affecting the stored data

Changes to stored data are serious

Be very careful with these operations

Be extremely careful in production

INSERT Table Data

Adding President Obama to presidents table

```
1 INSERT INTO
2     sampdb.president
3
4     SET
5         last_name    =    'Obama',
6         first_name   =    'Barack',
7         suffix       =    NULL,
8         city         =    'Houston',
9         state        =    'TX',
10        birth        =    '1961-08-04',
11        death        =    NULL
12 ;
13
14 ;
```

Alternate INSERT Syntax

1. That entry is wrong, so let's add the right one
2. **VALUES** provides another syntax for adding data

```
1 INSERT INTO
2     sampdb.president
3 VALUES  ('Obama', 'Barack', NULL, 'Honolulu', 'HI', '1961-08-04', NULL)
4 ;
5
```

Editing database data

1. We could also fix the problem with UPDATE

```
1 UPDATE
2     sampdb.president
3 SET
4     state = 'HI',
5     city = 'Honolulu'
6 WHERE
7     last_name = 'Obama'
8 ;
9
10
11
```

DELETE -- Removing Records

- Now we have two entries, let's remove one

```
1 DELETE FROM sampdb.president WHERE last_name = 'Obama' LIMIT 1;
```

2

DELETE -- Removing Many Records

1. Without the LIMIT, DELETE takes all records

```
1 DELETE FROM sampdb.president WHERE last_name = 'Obama';
2
```

Removing LOTS of Records

Be careful with DELETE!!

```
1 DELETE FROM sampdb.president;
2
```

```
1 +-----+
2 | COUNT(*) |
3 +-----+
4 |          0 |
5 +-----+
6
```

Here, not a big deal:

```
1 INSERT INTO sampdb.president
2 VALUES
3     ('Washington', 'George', NULL, 'Wakefield', 'VA', '1732-02-22', '1799-12-14'),
4     ('Adams', 'John', NULL, 'Braintree', 'MA', '1735-10-30', '1826-07-04'),
5     ('Jefferson', 'Thomas', NULL, 'Albemarle County', 'VA', '1743-04-13', '1826-07-04'),
6     ('Madison', 'James', NULL, 'Port Conway', 'VA', '1751-03-16', '1836-06-28'),
7     ('Monroe', 'James', NULL, 'Westmoreland County', 'VA', '1758-04-28', '1831-07-04'),
8     . . .
9 ;
```

Databases can protect against these mistakes

```
1 DROP TABLE IF EXISTS sampdb.president;
2 CREATE TABLE sampdb.president
3 (
4     last_name  VARCHAR(15) NOT NULL,
5     first_name VARCHAR(15) NOT NULL,
6     suffix      VARCHAR(5)  NULL,
7     city        VARCHAR(20) NOT NULL,
8     state       VARCHAR(2)  NOT NULL,
9     birth       DATE NOT NULL,
10    death       DATE NULL,
11    UNIQUE (last_name, first_name)
12 );
13
```

UNIQUE provides a data consistency rule

Prevents insertion of a record with a first_name, last_name tuple matching an extant record

Thus prevents addition of data inconsistent with data already entered

"Consistent" != "Right"

If we enter the Hawaii values first, the Illinois ones are prevented

But, we are forced to solve the problem in some way that doesn't risk mistakes like two duplicating records

Summary of Query Components

SELECT

- Specifies the columns output in the query result
- Can be a simple field retrieval from a base table, OR
- An expression of some kind -- function return, mathematical, subquery
- We use SELECT expressions to specify and manipulate output

FROM

- Specifies the source table from which the query draws
- Simplest case is a 'base table' holding data
- BUT because queries themselves always return tables, FROM can refer to another query -- e.g. a subquery
- FROM can also refer to a VIEW -- a shorthand stored on the server referring to a specified query

WHERE

- Filters records in source of query
- Multiple conditions can be specified using AND and OR
- Filters can compare tuples of source record with tuples, and lists of tuples. This allows us to use queries nested in WHERE to specify criteria for inclusion of records from an outer query.
-

GROUP BY

- Specifies subsetting of records for aggregation functions in SELECT
- Records are put into subsets according to their values for the tuple specified in GROUP BY
- The query output will have a row for each unique combination of values for that tuple in the base data

HAVING

- Provides filtration of records output by a query using an aggregating function
- Output records are filtered by their values -- generally those for some aggregating function

ORDER BY

- Provides control over the order in which query results are displayed
- Can specify order with tuples, each subsequent term ordering records within the order provided by the prior term(s)

JOIN

- Provides means for linking records in one table with those in other tables
- Crucial, as proper data modeling distributes data over a number of tables
- JOIN allows us to recombine those abstracted data into the outputs we want
- Works by creating joined records, among those of specified tables, of those records in those tables whose values in specified fields match
- (Note that JOIN can use comparisons other than equality.)

Primary and Foreign Keys

Primary Keys specify unique records in a table

- Primary key field or fields are specified when table is created
- Database server will not enter data that duplicates primary key values
- Other tables can reference entities in the table by these primary key fields

Foreign keys link a table to a 'parent'

The 'child' table specifies a foreign key at creation

- Foreign keys can be one or many fields
- All entries to the child table must have entries, matching on the foreign key values, in the parent table
- These values need not be unique in the child table
- The foreign key 'links' the child table to the parent
- The presence of parent entries matching child entries is called 'referential integrity'
- Database server protects referential integrity by refusing child entries without matching parents
- Server can also be setup to automatically delete all child records of a parent record when that parent record is deleted

Normal Forms

Modeling "norms" protect data clarity

- Poorly organized data models are hard to understand
- Worse, data gets lost, or duplicated
- Creating risk that data become inconsistent with the model
- The problem gets worse as data becomes increasingly involved
- 'Norms' provide rules that prevent many inconsistencies and losses

With all this syntax, we don't care about base table organization

- We are accustomed to Excel-like data organization, where the 'base' tables are some commonly used views
- and we derive other views of the data from those
- But in SQL, we always look at some view of the base data anyway
- We have to split our expectations of how we view the data from the way in which it is organized in base tables
- Base tables should be organized for clarity of the data model, NOT for some particular view of the data

First Normal Form

Data should be atomic, eg values shouldn't be combined

"George Washington" shouldn't be a value, because it combines first and last name

Unless we never will search on or retrieve first and last name

No repeating columns

No columns like author1, author2, author3

Why? We don't dependency on order, nor further processing

Values should be retrievable by name, not by further parsing or processing

We don't want values established by order -- again, we want values retrievable by field name

Second Normal Form

First Normal, plus:

No 'partial dependencies'

Cannot determine a nonkey column value with only a portion of the primary key

Such values should be remodeled to a table where the nonkey values are related to the (new) table's full primary key

Why? We want tables focused on entities

A partial dependency says that there are subentities to the entities in the table

Those should be removed to a table where they are the entities

This avoids duplication of information, confusion over entity definitions, and location of information

Third Normal Form

Second Normal, plus:

No nonkey column depends on the value of another nonkey column

Values are driven by the primary key fields

Again, such values should be removed to a table of their own

Why? Again, we want to avoid subentities

Basic idea: Primary Keys to rule them all

Data should be broken down so that any datum is found by searching on some table on some primary key, and only by that table and primary key

Because primary keys are table-unique, this organizes data so there are no duplications

Note that this can slow performance, because we have to join tables to provide denormalized views

But the model avoids repetition and inconsistency of data

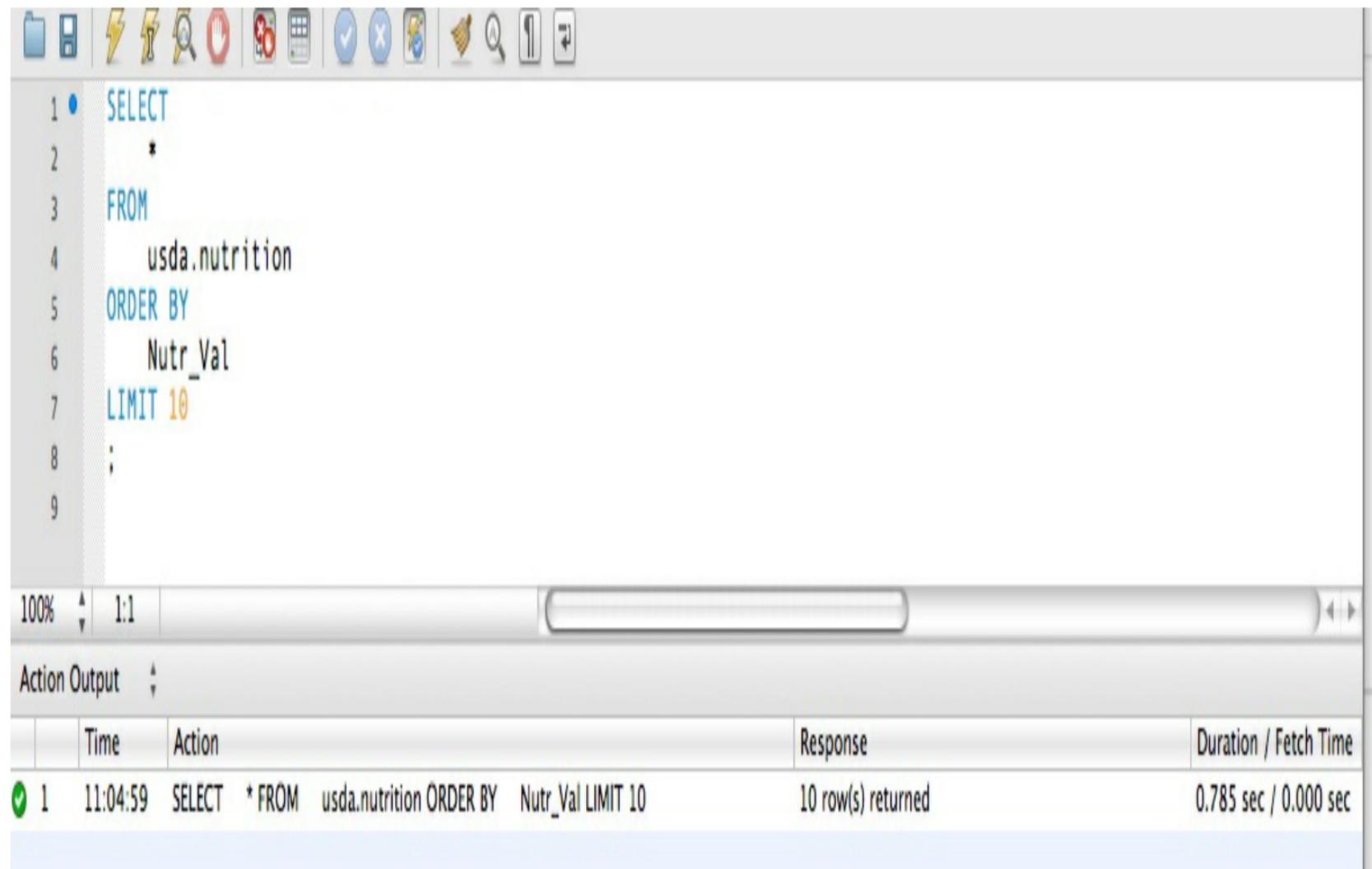
As models become more complex, they outgrow our ability to easily hold them entirely in our heads and notice problems

The normal forms protect us from inconsistent and out of date data

Performance and Indexing

Consider a large table

ORDER BY takes a long time to run on 590,000 records



The screenshot shows a database query interface with the following details:

- Toolbar:** Includes icons for file operations, search, and various database functions.
- Code Editor:** Displays a SQL query with line numbers:

```
1 • SELECT
2 *
3 FROM
4     usda.nutrition
5 ORDER BY
6     Nutr_Val
7 LIMIT 10
8 ;
9
```
- Execution Results:** Shows the execution progress at 100% completion and a 1:1 ratio.
- Action Output:** A dropdown menu currently set to "Action Output".
- Table:** A results table with the following columns:

| | Time | Action | Response | Duration / Fetch Time |
|---|----------|---|--------------------|-----------------------|
| 1 | 11:04:59 | SELECT * FROM usda.nutrition ORDER BY Nutr_Val LIMIT 10 | 10 row(s) returned | 0.785 sec / 0.000 sec |

Indexing makes it faster

An index presorts a specific field on a table

That speeds up search operations enormously

```
1 ALTER TABLE
2   usda.nutrition
3 ADD INDEX
4   nutr_val (Nutr_Val)
5 ;
6
```

1

Faster!

Takes no time at all!

| Action Output | | | |
|---------------|----------|---|-----------------------|
| | Time | Action | Response |
| | | | Duration / Fetch Time |
| 1 | 11:04:59 | SELECT * FROM usda.nutrition ORDER BY Nutr_Val LIMIT 10 | 10 row(s) returned |
| 2 | 11:23:08 | SELECT * FROM usda.nutrition ORDER BY Nutr_Val LIMIT 10 | 10 row(s) returned |

Why indexing works

An index is a pre-ordered sort of a particular field

- Done and stored by the database server
- In an unordered list, the search must check every entry until it hits a match
- In an ordered list, dictionary, the search skips large blocks of entries as before or after the target
- E.g., a dictionary is faster than an unsorted pile of words
- We don't index everything because indexing has storage and performance costs
- For every index on a table, the database must reorder the index for a new entry
- Each index takes up space, in storage and memory

Workflow Integration

How would we chart the Titanic data?

People want graphs, not tables!

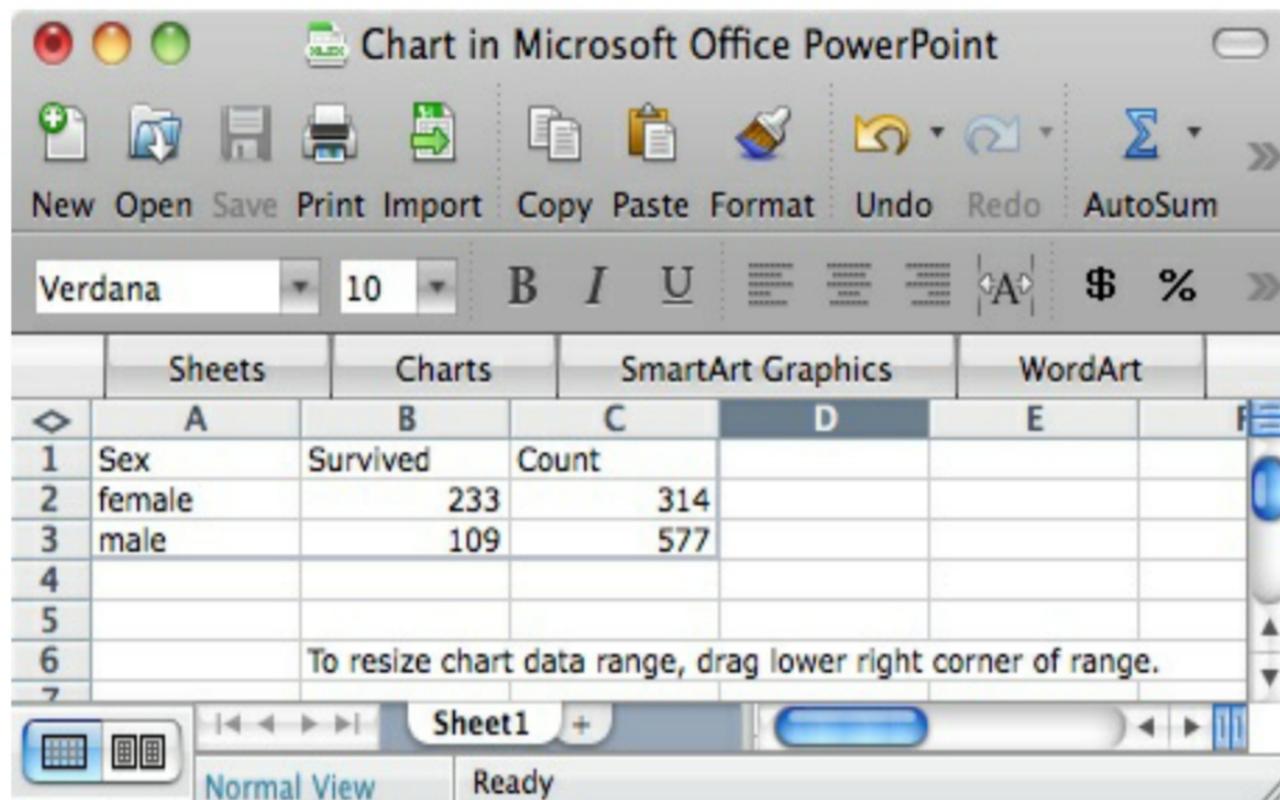
But, SQL doesn't provide visualization tools

We have to move the data to a program that does

Fortunately, moving data is nothing new

Powerpoint Data Entry

We already paste data from Excel to Powerpoint for graphing



The screenshot shows a Microsoft Office PowerPoint window titled "Chart in Microsoft Office PowerPoint". The ribbon menu is visible with tabs for "New", "Open", "Save", "Print", "Import", "Copy", "Paste", "Format", "Undo", "Redo", and "AutoSum". The font and size are set to "Verdana" and "10". The "Format" tab is selected. The "Sheets" tab is active in the ribbon, and the "Charts" tab is visible. A data table is displayed in the center of the slide:

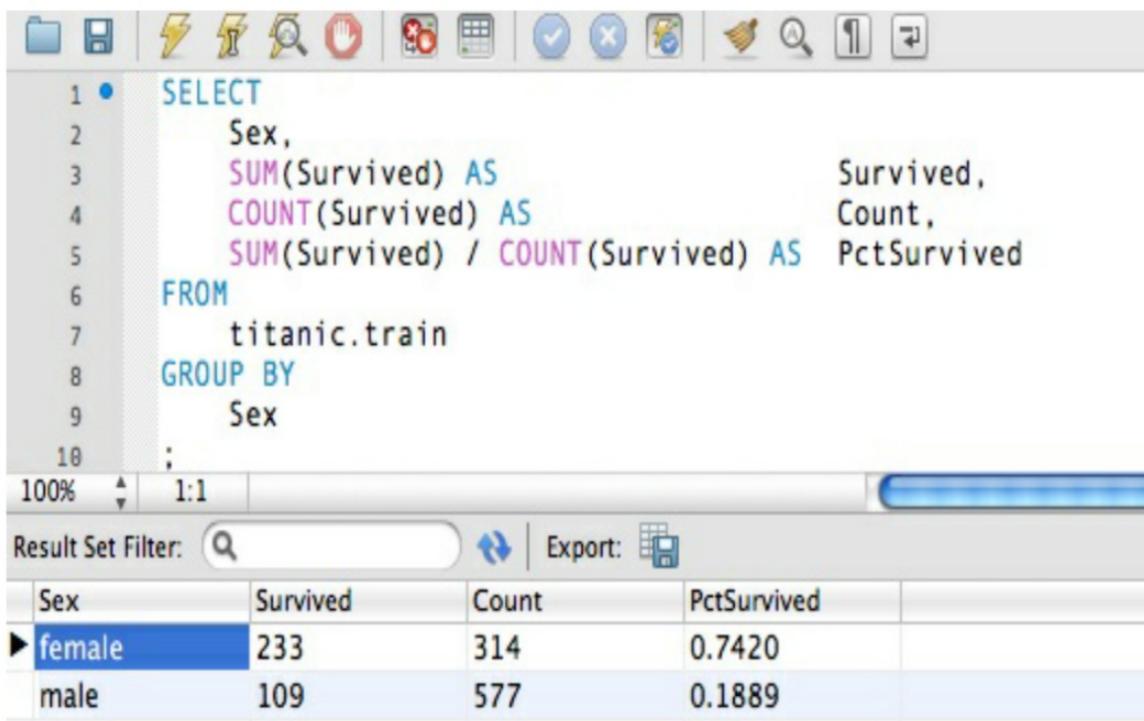
| | A | B | C | D | E | F |
|---|--------|----------|-------|-----|---|---|
| 1 | Sex | Survived | Count | | | |
| 2 | female | | 233 | 314 | | |
| 3 | male | | 109 | 577 | | |
| 4 | | | | | | |
| 5 | | | | | | |
| 6 | | | | | | To resize chart data range, drag lower right corner of range. |
| 7 | | | | | | |

The status bar at the bottom shows "Sheet1" and "Normal View".

- We can do the same thing from Workbench
- And, many other GUI clients
- Or, from Excel

From Workbench:

Use CSV export



The screenshot shows a SQL query in the SQL tab and its results in the Results tab of Oracle SQL Workbench.

SQL Tab:

```
1 • SELECT
2     Sex,
3         SUM(Survived) AS Survived,
4         COUNT(Survived) AS Count,
5         SUM(Survived) / COUNT(Survived) AS PctSurvived
6     FROM
7         titanic.train
8     GROUP BY
9         Sex
10    ;
```

Results Tab:

| Sex | Survived | Count | PctSurvived |
|--------|----------|-------|-------------|
| female | 233 | 314 | 0.7420 |
| male | 109 | 577 | 0.1889 |

- Note that here the query doesn't specify column headers -- the exported CSV includes a header row
- Clicking 'Export' opens a dialog for saving results to CSV
- These can be opened for copy and paste in Excel -- or directly into the Powerpoint graph

Queries can export direct to CSV

```
1 -- LACKS COLUMN HEADERS
2 SELECT
3   "Sex",
4   "Survived",
5   "Count",
6   "PctSurvived"
7
8 UNION
9
10 SELECT
11   Sex,
12   SUM(Survived) AS Survived,
13   COUNT(Survived) AS Count,
14   SUM(Survived) / COUNT(Survived) AS PctSurvived
15 FROM
16   titanic.train
17 GROUP BY
18   Sex
19 -- BEWARE permissions issues
20 INTO OUTFILE
21   '/tmp/survival_chart.csv'
22 FIELDS
23   TERMINATED BY ','
24   OPTIONALLY ENCLOSED BY "'"
25 LINES
26   TERMINATED BY '\n';
27 ;
28
```

What is UNION?

SQL statement concatenating two query results into one

- Must have the same number of columns
- Combines the rows of one query with another
- We use this here to put a header row into our exported data

Here, paste serves as a crude 'interface'

An interface is a point where one component provides input to another

- SQL serves as an analytic 'component'
- Powerpoint / D3 / Tableau / et cetera serve as graphing components
- We manually pass the output of our analysis to visualization through copy and paste
- Programming languages (Python, Ruby, PHP) provide direct / automated interface

This manual process is a hassle

- BUT it does mean that we can use ANY visualization tool
- We just have to figure out its interface, and we can write our queries to meet that requirement

Data Cleaning

Google Finance Data

Google NASDAQ:AAPL gordon@practicalhorseshoeing.com

Finance Apple Inc. (NASDAQ:AAPL)

Company

| | Daily prices | | | | | |
|--|--------------|--------|--------|--------|--------|------------|
| | Date | Open | High | Low | Close | Volume |
| Summary | Jan 29, 2014 | 503.95 | 507.37 | 498.62 | 500.75 | 17,991,828 |
| News | Jan 28, 2014 | 508.76 | 515.00 | 502.07 | 506.50 | 38,119,083 |
| Option chain | Jan 27, 2014 | 550.07 | 554.80 | 545.75 | 550.50 | 20,602,736 |
| Related companies | Jan 24, 2014 | 554.00 | 555.62 | 544.75 | 546.07 | 15,483,491 |
| Historical prices | Jan 23, 2014 | 549.94 | 556.50 | 544.81 | 556.18 | 14,425,478 |
| Financials | Jan 22, 2014 | 550.91 | 557.29 | 547.81 | 551.51 | 13,602,762 |
| Markets | Jan 21, 2014 | 540.99 | 550.07 | 540.42 | 549.07 | 11,750,792 |
| News | Jan 17, 2014 | 551.48 | 552.07 | 539.90 | 540.67 | 15,489,527 |
| Portfolios | Jan 16, 2014 | 554.90 | 556.85 | 551.68 | 554.25 | 8,210,190 |
| Stock screener | Jan 15, 2014 | 553.52 | 560.20 | 551.66 | 557.36 | 14,067,517 |
| Google Domestic Trends | Jan 14, 2014 | 538.22 | 546.73 | 537.66 | 546.39 | 11,962,053 |
| | Jan 13, 2014 | 529.91 | 542.50 | 529.88 | 535.73 | 13,551,549 |
| | Jan 10, 2014 | 539.83 | 540.80 | 531.11 | 532.94 | 10,902,952 |
| Recent Quotes <small>(Turn on)</small> | Jan 9, 2014 | 546.80 | 546.86 | 535.35 | 536.52 | 9,986,457 |
| | Jan 8, 2014 | 538.81 | 545.56 | 538.69 | 543.46 | 9,240,955 |
| You have no recent quotes | Jan 7, 2014 | 544.32 | 545.96 | 537.92 | 540.04 | 11,347,538 |

Daily prices Jan 31, 2013 - Jan 30, 2014

Historical chart



Feb 1, 2013 Jan 29, 2014

Export

We can't load this directly

```
1 DROP TABLE IF EXISTS csv_load.aapl_initial;
2
3 CREATE TABLE csv_load.aapl_initial (
4     date      VARCHAR(10),
5     -- date      DATE,
6     open      DECIMAL(5, 2),
7     high      DECIMAL(5, 2),
8     low       DECIMAL(5, 2),
9     close     DECIMAL(5, 2),
10    volume    INTEGER(12)
11
12 );
13
```

| | date | open | high | low | close | volume |
|----|-----------|--------|--------|--------|--------|----------|
| 4 | 29-Jan-14 | 503.95 | 507.37 | 498.62 | 500.75 | 17991828 |
| 5 | 28-Jan-14 | 508.76 | 515.00 | 502.07 | 506.50 | 38119083 |
| 6 | 27-Jan-14 | 550.07 | 554.80 | 545.75 | 550.50 | 20602736 |
| 7 | 24-Jan-14 | 554.00 | 555.62 | 544.75 | 546.07 | 15483491 |
| 8 | 23-Jan-14 | 549.94 | 556.50 | 544.81 | 556.18 | 14425478 |
| 9 | 22-Jan-14 | 550.91 | 557.29 | 547.81 | 551.51 | 13602762 |
| 10 | 21-Jan-14 | 540.99 | 550.07 | 540.42 | 549.07 | 11750792 |
| 11 | 17-Jan-14 | 551.48 | 552.07 | 539.90 | 540.67 | 15489527 |
| 12 | 16-Jan-14 | 554.90 | 556.85 | 551.68 | 554.25 | 8210190 |
| 13 | 15-Jan-14 | 553.52 | 560.20 | 551.66 | 557.36 | 14067517 |
| 14 | 14-Jan-14 | 538.22 | 546.73 | 537.66 | 546.39 | 11962053 |
| 15 | 13-Jan-14 | 529.91 | 542.50 | 529.88 | 535.73 | 13551549 |
| 16 | 10-Jan-14 | 539.83 | 540.80 | 531.11 | 532.94 | 10902952 |
| 17 | 9-Jan-14 | 546.80 | 546.86 | 535.35 | 536.52 | 9986457 |
| 18 | 8-Jan-14 | 538.81 | 545.56 | 538.69 | 543.46 | 9240955 |
| 19 | 7-Jan-14 | 544.32 | 545.96 | 537.92 | 540.04 | 11347538 |
| 20 | 6-Jan-14 | 537.45 | 546.80 | 533.60 | 543.93 | 14765593 |
| 21 | 3-Jan-14 | 552.86 | 553.70 | 540.43 | 540.98 | 14043410 |
| 22 | 2-Jan-14 | 555.68 | 557.03 | 552.02 | 553.13 | 8398851 |
| 23 | 31-Dec-13 | 554.17 | 561.28 | 554.00 | 561.02 | 7974196 |
| 24 | | | | | | |
| 25 | | | | | | |

- We can't load date data in this format

We can parse the data

```
1 DROP VIEW IF EXISTS csv_load.date_parsed;
2
3 CREATE VIEW csv_load.date_parsed
4
5 AS
6
7 SELECT
8   date,
9   LEFT(date, LOCATE('-', date) - 1 ) as day,
10  RIGHT(date, LENGTH(date) - LOCATE('-', date, 4) ) as yr,
11  MID(
12    date,
13    LOCATE('-', date) + 1,
14    ( LOCATE('-', date, 4) - LOCATE('-', date) - 1 )
15  ) as mnt,
16  open,
17  high,
18  low,
19  close,
20  volume
21
22 FROM
23  csv_load.aapl_initial
24 ;
25
```

| | date | day | yr | mnt | open | high | low | close | volume |
|----|-----------|-----|----|-----|--------|--------|--------|--------|----------|
| 1 | 29-Jan-14 | 29 | 14 | Jan | 503.95 | 507.37 | 498.62 | 500.75 | 17991828 |
| 2 | 28-Jan-14 | 28 | 14 | Jan | 508.00 | 510.00 | 505.07 | 506.50 | 18119083 |
| 3 | 27-Jan-14 | 27 | 14 | Jan | 551.07 | 554.80 | 545.75 | 550.50 | 20602736 |
| 4 | 24-Jan-14 | 24 | 14 | Jan | 554.00 | 555.62 | 544.75 | 545.07 | 15483491 |
| 5 | 23-Jan-14 | 23 | 14 | Jan | 549.94 | 556.50 | 544.81 | 556.18 | 14425478 |
| 6 | 22-Jan-14 | 22 | 14 | Jan | 550.91 | 557.29 | 547.81 | 551.51 | 13602762 |
| 7 | 21-Jan-14 | 21 | 14 | Jan | 540.99 | 550.07 | 540.42 | 549.07 | 11750792 |
| 8 | 17-Jan-14 | 17 | 14 | Jan | 551.48 | 552.07 | 539.90 | 540.67 | 15489527 |
| 9 | 16-Jan-14 | 16 | 14 | Jan | 554.90 | 556.85 | 551.68 | 554.25 | 8210190 |
| 10 | 15-Jan-14 | 15 | 14 | Jan | 553.52 | 560.20 | 551.66 | 557.36 | 14067517 |
| 11 | 14-Jan-14 | 14 | 14 | Jan | 538.22 | 546.73 | 537.66 | 546.39 | 11962053 |
| 12 | 13-Jan-14 | 13 | 14 | Jan | 529.91 | 542.50 | 529.88 | 535.73 | 13551549 |
| 13 | 10-Jan-14 | 10 | 14 | Jan | 539.83 | 540.80 | 531.11 | 532.94 | 10902952 |
| 14 | 9-Jan-14 | 9 | 14 | Jan | 546.80 | 546.86 | 535.35 | 536.52 | 9986457 |
| 15 | 8-Jan-14 | 8 | 14 | Jan | 538.81 | 545.56 | 538.69 | 543.46 | 9240955 |
| 16 | 7-Jan-14 | 7 | 14 | Jan | 544.32 | 545.96 | 537.92 | 540.04 | 11347538 |
| 17 | 6-Jan-14 | 6 | 14 | Jan | 537.45 | 546.46 | 533.98 | 543.18 | 14765193 |
| 18 | 3-Jan-14 | 3 | 14 | Jan | 552.86 | 553.70 | 540.43 | 540.98 | 14043410 |
| 19 | 2-Jan-14 | 2 | 14 | Jan | 555.68 | 557.03 | 552.02 | 553.13 | 8398851 |
| 20 | 31-Dec-13 | 31 | 13 | Dec | 554.17 | 561.28 | 554.00 | 561.02 | 7974196 |
| 21 | | | | | | | | | |
| 22 | | | | | | | | | |
| 23 | | | | | | | | | |
| 24 | | | | | | | | | |
| 25 | | | | | | | | | |

- We use string functions to break up -- parse -- the date string

We can reformat the data

```
1 DROP VIEW IF EXISTS csv_load.date_reformat;
2
3 CREATE VIEW
4     csv_load.date_reformat
5 AS
6 SELECT
7     CONCAT(
8         -- year term
9         CONCAT( '20', yr),
10        '-',
11        -- month term
12        CASE
13            WHEN mnt = 'Jan' THEN '01'
14            WHEN mnt = 'Feb' THEN '02'
15            WHEN mnt = 'Mar' THEN '03'
16            WHEN mnt = 'Apr' THEN '04'
17            WHEN mnt = 'May' THEN '05'
18            WHEN mnt = 'Jun' THEN '06'
19            WHEN mnt = 'Jul' THEN '07'
20            WHEN mnt = 'Aug' THEN '08'
21            WHEN mnt = 'Sep' THEN '09'
22            WHEN mnt = 'Oct' THEN '10'
23            WHEN mnt = 'Nov' THEN '11'
24            WHEN mnt = 'Dec' THEN '12'
25        END,
26        '-',
27        -- day term
28        CONCAT( IF(LENGTH(day) = 1, '0', ''), day
29    )
30    ) AS date,
31    open,
32    high,
33    low,
34    close,
35    volume
36
37 FROM
38     csv_load.date_parsed
39;
40
41
42
```

| 1 | date | open | high | low | close | volume |
|----|------------|--------|--------|--------|--------|----------|
| 2 | 2014-01-29 | 504.95 | 507.37 | 498.62 | 500.75 | 17991828 |
| 3 | 2014-01-28 | 508.76 | 515.80 | 509.07 | 506.50 | 38119083 |
| 4 | 2014-01-27 | 550.07 | 558.80 | 545.75 | 550.50 | 20602736 |
| 5 | 2014-01-24 | 554.00 | 555.62 | 544.75 | 546.07 | 15483491 |
| 6 | 2014-01-23 | 549.94 | 556.50 | 544.81 | 556.18 | 14425478 |
| 7 | 2014-01-22 | 550.91 | 557.29 | 547.81 | 551.51 | 13602762 |
| 8 | 2014-01-21 | 540.99 | 550.07 | 540.42 | 549.07 | 11750792 |
| 9 | 2014-01-17 | 551.48 | 552.07 | 539.90 | 540.67 | 15489527 |
| 10 | 2014-01-16 | 554.90 | 556.85 | 551.68 | 554.25 | 8210190 |
| 11 | 2014-01-15 | 553.52 | 560.20 | 551.66 | 557.36 | 14067517 |
| 12 | 2014-01-14 | 538.22 | 546.73 | 537.66 | 546.39 | 11962053 |
| 13 | 2014-01-13 | 529.91 | 542.50 | 529.88 | 535.73 | 13551549 |
| 14 | 2014-01-10 | 539.83 | 540.80 | 531.11 | 532.94 | 10902952 |
| 15 | 2014-01-09 | 546.80 | 546.86 | 535.35 | 536.52 | 9986457 |
| 16 | 2014-01-08 | 538.81 | 545.56 | 538.69 | 543.46 | 9240955 |
| 17 | 2014-01-07 | 544.32 | 545.96 | 537.92 | 540.04 | 11347538 |
| 18 | 2014-01-06 | 537.45 | 546.80 | 533.60 | 533.93 | 1414939 |
| 19 | 2014-01-03 | 552.86 | 553.70 | 540.43 | 540.98 | 14043410 |
| 20 | 2014-01-02 | 555.68 | 557.03 | 552.02 | 553.13 | 8398851 |
| 21 | 2013-12-31 | 554.17 | 561.28 | 554.00 | 561.02 | 7974196 |
| 22 | | | | | | |
| 23 | | | | | | |
| 24 | | | | | | |
| 25 | | | | | | |

- We use CASE to transpose string months to digits
- and CONCAT to form a proper date string for load as DATE

This, we can load

```
1 INSERT INTO
2     csv_load.aapl
3 SELECT
4     date,
5     open,
6     high,
7     low,
8     close,
9     volume
10    FROM
11        csv_load.date_reformat
12    ;
13
```

| 1 | date | open | high | low | close | volume |
|----|------------|--------|--------|--------|--------|----------|
| 2 | 2014-01-29 | 504.95 | 507.37 | 498.62 | 500.75 | 17991828 |
| 3 | 2014-01-28 | 508.76 | 515.80 | 507.07 | 506.50 | 38119083 |
| 4 | 2014-01-27 | 550.07 | 558.80 | 545.75 | 550.50 | 20602736 |
| 5 | 2014-01-24 | 554.00 | 555.62 | 544.75 | 546.07 | 15483491 |
| 6 | 2014-01-23 | 549.94 | 556.50 | 544.81 | 556.18 | 14425478 |
| 7 | 2014-01-22 | 550.91 | 557.29 | 547.81 | 551.51 | 13602762 |
| 8 | 2014-01-21 | 540.99 | 550.07 | 540.42 | 549.07 | 11750792 |
| 9 | 2014-01-17 | 551.48 | 552.07 | 539.90 | 540.67 | 15489527 |
| 10 | 2014-01-16 | 554.90 | 556.85 | 551.68 | 554.25 | 8210190 |
| 11 | 2014-01-15 | 553.52 | 560.20 | 551.66 | 557.36 | 14067517 |
| 12 | 2014-01-14 | 538.22 | 546.73 | 537.66 | 546.39 | 11962053 |
| 13 | 2014-01-13 | 529.91 | 542.50 | 529.88 | 535.73 | 13551549 |
| 14 | 2014-01-10 | 539.83 | 540.80 | 531.11 | 532.94 | 10902952 |
| 15 | 2014-01-09 | 546.80 | 546.86 | 535.35 | 536.52 | 9986457 |
| 16 | 2014-01-08 | 538.81 | 545.56 | 538.69 | 543.46 | 9240955 |
| 17 | 2014-01-07 | 544.32 | 545.96 | 537.92 | 540.04 | 11347538 |
| 18 | 2014-01-06 | 537.45 | 546.80 | 533.60 | 543.93 | 1414939 |
| 19 | 2014-01-03 | 552.86 | 553.70 | 540.43 | 540.98 | 14043410 |
| 20 | 2014-01-02 | 555.68 | 557.03 | 552.02 | 553.13 | 8398851 |
| 21 | 2013-12-31 | 554.17 | 561.28 | 554.00 | 561.02 | 7974196 |
| 22 | | | | | | |
| 23 | | | | | | |
| 24 | | | | | | |
| 25 | | | | | | |

- This string loads properly as DATE

Now we can do date based analysis

```
1 SELECT
2     MIN(date) as week_start,
3     MAX(high) AS weekly_high,
4     MIN(low) AS weekly_low
5 FROM
6     csv_load.aapl
7 GROUP BY
8     week(date)
9 ORDER BY
10    week_start
11 LIMIT
12    20
13 ;
14
```

| | week_start | weekly_high | weekly_low |
|----|------------|-------------|------------|
| 4 | 2013-02-01 | 554.80 | 448.35 |
| 5 | 2013-02-04 | 478.81 | 442.00 |
| 6 | 2013-02-11 | 484.94 | 459.92 |
| 7 | 2013-02-19 | 462.73 | 442.82 |
| 8 | 2013-02-25 | 455.12 | 429.98 |
| 9 | 2013-03-04 | 435.43 | 419.00 |
| 10 | 2013-03-11 | 444.23 | 425.14 |
| 11 | 2013-03-18 | 462.10 | 441.20 |
| 12 | 2013-03-25 | 469.95 | 441.62 |
| 13 | 2013-04-01 | 443.70 | 419.68 |
| 14 | 2013-04-08 | 437.99 | 422.49 |
| 15 | 2013-04-15 | 427.89 | 385.10 |
| 16 | 2013-04-22 | 418.77 | 391.28 |
| 17 | 2013-04-29 | 453.23 | 420.00 |
| 18 | 2013-05-06 | 465.75 | 450.48 |
| 19 | 2013-05-13 | 457.90 | 418.90 |
| 20 | 2013-05-20 | 448.35 | 430.10 |
| 21 | 2013-05-28 | 457.10 | 439.40 |
| 22 | 2013-06-03 | 454.43 | 432.77 |
| 23 | 2013-06-10 | 449.08 | 428.50 |

- WEEK() returns which week of the year a date falls within
- We can use this value to group our data, and extract weekly highs and lows

Data Analysis -- More

Let's reprise data loading

We have this CSV data on Titanic survivors

```
1 PassengerId,Survived,Pclass,Name,Sex,Age,SibSp,Parch,Ticket,Fare,Cabin,Embarked
2 1,0,3,"Braund, Mr. Owen Harris",male,22,1,0,A/5 21171,7.25,,S
3 2,1,1,"Cumings, Mrs. John Bradley (Florence Briggs Thayer)",female,38,1,0,PC 17599,71.2833,C85,C
4 3,1,3,"Heikkinen, Miss. Laina",female,26,0,0,STON/O2. 3101282,7.925,,S
5 4,1,1,"Futrelle, Mrs. Jacques Heath (Lily May Peel)",female,35,1,0,113803,53.1,C123,S
6 5,0,3,"Allen, Mr. William Henry",male,35,0,0,373450,8.05,,S
7 6,0,3,"Moran, Mr. James",male,,0,0,330877,8.4583,,Q
8 7,0,1,"McCarthy, Mr. Timothy J",male,54,0,0,17463,51.8625,E46,S
9 8,0,3,"Palsson, Master. Gosta Leonard",male,2,3,1,349909,21.075,,S
10 9,1,3,"Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)",female,27,0,2,347742,11.1333,,S
11
```

Create database, table

Familiar?

- These tell the database server to create containers for the data
- We name the table 'train' as part of an analytic technique -- 'training' a model on a subset of the available data

```
1 DROP DATABASE IF EXISTS titanic;
2 CREATE DATABASE titanic;
3
```

```
1 DROP TABLE IF EXISTS titanic.train;
2 CREATE TABLE titanic.train (
3     PassengerId INT(4),
4     Survived     INT(1),
5     Pclass       INT(3),
6     Name         VARCHAR(50),
7     Sex          CHAR(6),
8     Age          INT(3),
9     SibSp        INT(2),
10    Parch        VARCHAR(50),
11    Ticket       INT(5),
12    Fare         DECIMAL(10,2),
13    Cabin        VARCHAR(10),
14    Embarked     CHAR(2)
15 );
16
```

Load the data

Familiar?

```
1 LOAD DATA LOCAL
2
3 #   EDIT THIS PATH:
4 INFILE "/Users/gimli/Work/teaching/practical_sql/class_materials/deck_templates/december2013/class_04/titanic_setup/data/train.csv"
5
6 INTO TABLE titanic.train
7
8 FIELDS
9   TERMINATED BY ','
10  ENCLOSED BY '\"'
11
12 LINES
13   #-TERMINATED BY '\n'
14   # If that gets strange results, comment it out and try this:
15   TERMINATED BY '\r\n'
16
17 IGNORE 1 LINES;
18
```

This tells the server where to find the file, where to put the data, and how to identify particular fields and records

Now we have a table

```
1 DESCRIBE titanic.train;  
2
```

| Field | Type | Null | Key | Default | Extra |
|-------------|---------------|------|-----|---------|-------|
| PassengerId | int(4) | YES | | NULL | |
| Survived | int(1) | YES | | NULL | |
| Pclass | int(3) | YES | | NULL | |
| Name | varchar(50) | YES | | NULL | |
| Sex | char(6) | YES | | NULL | |
| Age | int(3) | YES | | NULL | |
| SibSp | int(2) | YES | | NULL | |
| Parch | varchar(50) | YES | | NULL | |
| Ticket | int(5) | YES | | NULL | |
| Fare | decimal(10,2) | YES | | NULL | |
| Cabin | varchar(10) | YES | | NULL | |
| Embarked | char(2) | YES | | NULL | |

And the table holds data

```
1 SELECT * FROM titanic.train LIMIT 10;
```

| PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare |
|-------------|----------|--------|---|--------|-----|-------|-------|--------|------|
| 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22 | 1 | 0 | 0 | 7. |
| 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Thayer) | female | 38 | 1 | 0 | 0 | 71. |
| 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26 | 0 | 0 | 0 | 7. |
| 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35 | 1 | 0 | 113803 | 53. |
| 5 | 0 | 3 | Allen, Mr. William Henry | male | 35 | 0 | 0 | 373450 | 8. |
| 6 | 0 | 3 | Moran, Mr. James | male | 0 | 0 | 0 | 330877 | 8. |
| 7 | 0 | 1 | McCarthy, Mr. Timothy J | male | 54 | 0 | 0 | 17463 | 51. |
| 8 | 0 | 3 | Palsson, Master. Gosta Leonard | male | 2 | 3 | 1 | 349909 | 21. |
| 9 | 1 | 3 | Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg) | female | 27 | 0 | 2 | 347742 | 11. |
| 10 | 1 | 2 | Nasser, Mrs. Nicholas (Adele Achem) | female | 14 | 1 | 0 | 237736 | 30. |

Analysing Data

Let's explore that

Use DISTINCT and GROUP BY to see what's there

```
1 SELECT DISTINCT
2   Survived,
3   Pclass,
4   Sex,
5   Parch,
6   Embarked
7 FROM
8   titanic.train
9 LIMIT 30
10 ;
11
```

| | Survived | Pclass | Sex | Parch | Embarked |
|----|----------|--------|--------|-------|----------|
| 1 | 0 | 3 | male | 0 | S |
| 2 | 1 | 1 | female | 0 | C |
| 3 | 1 | 3 | female | 0 | S |
| 4 | 1 | 1 | female | 0 | S |
| 5 | 0 | 3 | male | 0 | Q |
| 6 | 0 | 1 | male | 0 | S |
| 7 | 0 | 3 | male | 1 | S |
| 8 | 1 | 3 | female | 2 | S |
| 9 | 1 | 2 | female | 0 | C |
| 10 | 1 | 3 | female | 1 | S |
| 11 | 0 | 3 | male | 5 | S |
| 12 | 0 | 3 | female | 0 | S |
| 13 | 1 | 2 | female | 0 | S |
| 14 | 0 | 3 | male | 1 | S |
| 15 | 0 | 3 | female | 0 | S |
| 16 | 1 | 2 | female | 0 | S |
| 17 | 0 | 3 | male | 1 | Q |
| 18 | 1 | 2 | male | 0 | S |
| 19 | 1 | 3 | female | 0 | C |
| 20 | 0 | 2 | male | 0 | S |
| 21 | 1 | 3 | female | 0 | Q |
| 22 | 1 | 1 | male | 0 | S |
| 23 | 0 | 3 | female | 1 | S |
| 24 | 1 | 3 | female | 5 | S |
| 25 | 0 | 3 | male | 0 | C |
| 26 | 0 | 1 | male | 2 | S |
| 27 | 0 | 1 | male | 0 | C |
| 28 | 1 | 3 | male | 0 | C |
| 29 | 0 | 2 | female | 0 | S |
| 30 | 1 | 2 | female | 2 | C |
| 31 | 0 | 1 | male | 1 | C |
| 32 | 1 | 2 | female | 2 | S |
| 33 | 0 | 3 | male | 2 | S |
| 34 | | | | | |
| 35 | | | | | |

Who survived?

We use GROUP BY and output transformation

```
1 SELECT
2     Sex,
3     SUM(Survived) AS Survived,
4     COUNT(Survived) AS Count,
5     SUM(Survived) / COUNT(Survived) AS PctSurvived
6 FROM
7     titanic.train
8 GROUP BY
9     Sex
10 ;
11
```

```
1 +-----+-----+-----+
2 | Sex   | Survived | Count | PctSurvived |
3 +-----+-----+-----+
4 | female |      233 |    314 |      0.7420 |
5 | male   |      109 |    577 |      0.1889 |
6 +-----+-----+-----+
7
```

We can sort the fares into bins

Use CASE to assign fare values into "buckets"

```
1 DROP VIEW IF EXISTS titanic.binned;
2 CREATE VIEW titanic.binned
3 AS
4     SELECT
5         *,
6         CASE
7             WHEN Fare <= 50
8                 THEN 1
9             WHEN Fare > 50 AND Fare <= 100
10                THEN 2
11             WHEN Fare > 100 AND Fare <= 150
12                 THEN 3
13             WHEN Fare > 150 AND Fare <= 200
14                 THEN 4
15             ELSE
16                 5
17         END AS FareBin
18     FROM
19         titanic.train
20 ;
21 ;
```

- CASE allows us to assign values to tiers

Now we've a VIEW with bins

```
1 SELECT
2     PassengerId,
3     Fare,
4     FareBin
5 FROM
6     titanic.binned
7 LIMIT 10
8 ;
9
```

| | PassengerId | Fare | FareBin |
|----|-------------|-------|---------|
| 4 | 1 | 7.25 | 1 |
| 5 | 2 | 71.28 | 2 |
| 6 | 3 | 7.93 | 1 |
| 7 | 4 | 53.10 | 2 |
| 8 | 5 | 8.05 | 1 |
| 9 | 6 | 8.46 | 1 |
| 10 | 7 | 51.86 | 2 |
| 11 | 8 | 21.08 | 1 |
| 12 | 9 | 11.13 | 1 |
| 13 | 10 | 30.07 | 1 |
| 14 | | | |
| 15 | | | |

Now we can add bins to our survival analysis

```
1 SELECT
2             Sex,
3             FareBin,
4             SUM(Survived) AS SumSurvived,
5             COUNT(Survived) AS CountSurvived,
6             SUM(Survived) / COUNT(Survived) AS PctSurvived
7 FROM
8     titanic.binned
9 GROUP BY
10            Sex,
11            FareBin
12 ORDER BY
13            Sex,
14            FareBin
15 ;
16
17
```

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|----|--------|---------|-------------|---------------|-------------|---|---|---|----|----|----|----|----|----|
| | Sex | FareBin | SumSurvived | CountSurvived | PctSurvived | | | | | | | | | |
| 4 | female | 1 | 151 | 227 | 0.6652 | | | | | | | | | |
| 5 | female | 2 | 50 | 53 | 0.9434 | | | | | | | | | |
| 6 | female | 3 | 15 | 15 | 1.0000 | | | | | | | | | |
| 7 | female | 4 | 5 | 7 | 0.7143 | | | | | | | | | |
| 8 | female | 5 | 12 | 12 | 1.0000 | | | | | | | | | |
| 9 | male | 1 | 82 | 504 | 0.1627 | | | | | | | | | |
| 10 | male | 2 | 20 | 54 | 0.3704 | | | | | | | | | |
| 11 | male | 3 | 4 | 9 | 0.4444 | | | | | | | | | |
| 12 | male | 4 | 1 | 2 | 0.5000 | | | | | | | | | |
| 13 | male | 5 | 2 | 8 | 0.2500 | | | | | | | | | |

We can add further "dimensions"

```
1 SELECT
2             Sex,
3             FareBin,
4             Embarked,
5             SUM(Survived) AS SumSurvived,
6             COUNT(Survived) AS CountSurvived,
7             SUM(Survived) / COUNT(Survived) AS PctSurvived
8 FROM
9     titanic.binned
10 GROUP BY
11     Sex,
12     FareBin,
13     Embarked
14 ORDER BY
15     Sex,
16     FareBin,
17     Embarked
18 ;
19
20 ;
```

Digging deeper into the data

| 1 | Sex | FareBin | Embarked | SumSurvived | CountSurvived | PctSurvived |
|----|--------|---------|----------|-------------|---------------|-------------|
| 2 | female | 1 | C | 27 | 36 | 0.7500 |
| 3 | female | 1 | Q | 26 | 35 | 0.7429 |
| 4 | female | 1 | S | 98 | 156 | 0.6282 |
| 5 | female | 2 | | 2 | 2 | 1.0000 |
| 6 | female | 2 | C | 19 | 19 | 1.0000 |
| 7 | female | 2 | Q | 1 | 1 | 1.0000 |
| 8 | female | 2 | S | 28 | 31 | 0.9032 |
| 9 | female | 3 | C | 11 | 11 | 1.0000 |
| 10 | female | 3 | S | 4 | 4 | 1.0000 |
| 11 | female | 4 | S | 5 | 7 | 0.7143 |
| 12 | female | 5 | C | 7 | 7 | 1.0000 |
| 13 | female | 5 | S | 5 | 5 | 1.0000 |
| 14 | male | 1 | C | 18 | 71 | 0.2535 |
| 15 | male | 1 | Q | 3 | 40 | 0.0750 |
| 16 | male | 1 | S | 61 | 393 | 0.1552 |
| 17 | male | 2 | C | 8 | 13 | 0.6154 |
| 18 | male | 2 | Q | 0 | 1 | 0.0000 |
| 19 | male | 2 | S | 12 | 40 | 0.3000 |
| 20 | male | 3 | C | 1 | 6 | 0.1667 |
| 21 | male | 3 | S | 3 | 3 | 1.0000 |
| 22 | male | 4 | S | 1 | 2 | 0.5000 |
| 23 | male | 5 | C | 2 | 5 | 0.4000 |
| 24 | male | 5 | S | 0 | 3 | 0.0000 |

SQL allows more "dimensions" than Excel

- We can imagine a 2 x 2 pivot table, and even a third dimension
- But SQL allows us to subdivide the data on as many fields as we can imagine
- We can use this to simply filter out records (e.g., let's look only at females in fare bin 3)
- Or, to look outcomes based on subcategories
- We could also write queries calculating simple linear regression values

Summing Up

Takeaways

- Model is "verbal", not visual -- which requires adjustment from Excel based habits
- Syntax is 'Germanic', expressions later in queries can change the effect of preceding statements
- Practice builds intuition and conceptual orientation

Takeaways

- Base data is immutable and not changed -- we use queries to describe desired outputs
- Notice that this frees us to arrange the base data in ways most amenable to query, rather than according to some particular view
- What matters about base table arrangement is our ability to describe outputs from it, rather than its immediate / initial appearance
- We conduct more complex analysis by performing queries on output of queries, eg subqueries and joins
- Mastering the basic syntax allows us to devise more complicated expressions, passing one analysis output on to another

Get the Deck

gordon@practicalhorseshoeing.com