# Introduction to Data Science

Daniel Gutierrez, Data Scientist

Los Angeles, Calif.

# Course Outcomes

- Accessing data sets from R

- Review of commonly accessed data sources

- Provide reusable code snippets for accessing data

- Learn how to write data files

# Lesson Objectives

- Accessing data sources
- Downloading files from the web
- Comma separated value (CSV)
- Excel
- JSON
- Web page scraping
- SQL databases
- SQL equivalents in R
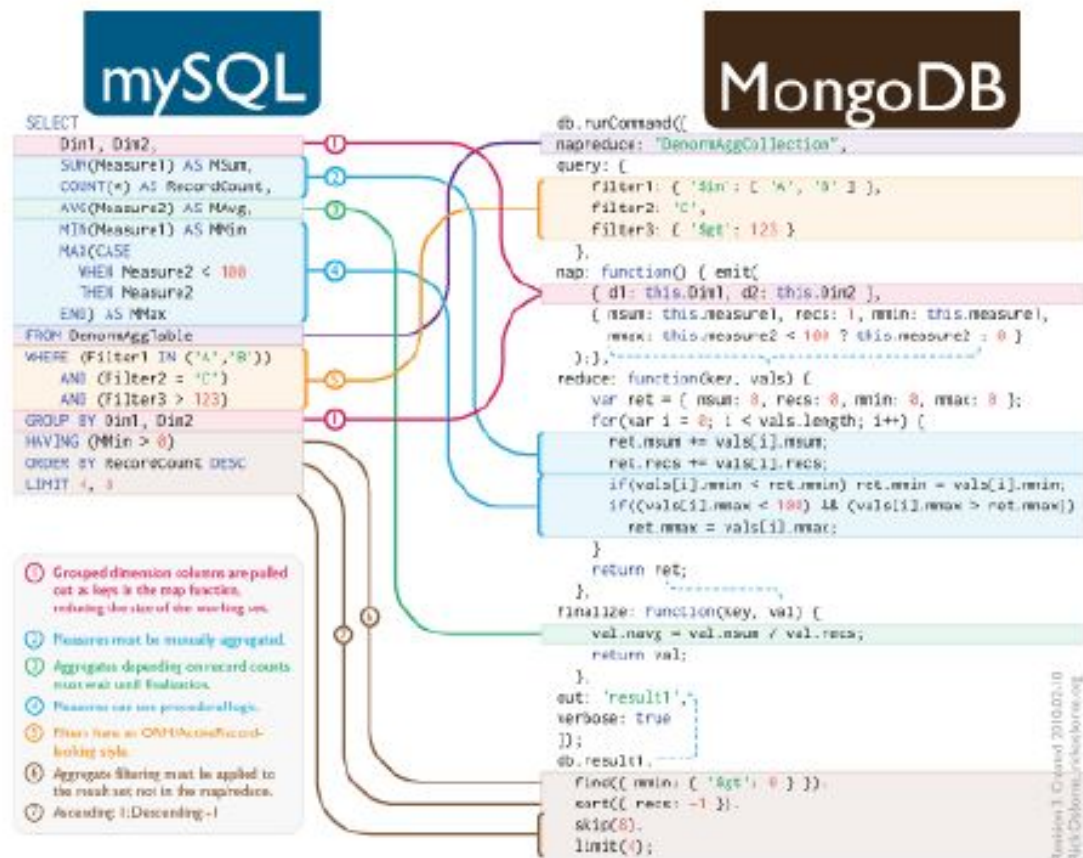- Writing data

# Data Access

## What you wish data looked like

# What does data really look like?

```
------------------------------ ALLERGIES ------------------------------          ------------------------------ MEDICATION HISTORY ----------------

ast Updated: 01 Dec 2011 @ 0851                                                  Last Updated: 11 Apr 2011 @ 1737

                                                                                 Medication: AMLODIPINE BESYLATE 10MG TAB
llergy Name:        TRIMETHOPRIM                                                 Instructions: TAKE ONE TABLET BY MOUTH TAKE ONE-HALF TABLET FOR :
ocation:            DAYT29                                                       GRAPEFRUIT JUICE--
ate Entered:        09 Mar 2011                                                  Status: Active
eaction:                                                                         Refills Remaining: 3
llergy Type:        DRUG                                                         Last Filled On: 20 Aug 2010
A Drug Class:       ANTI-INFECTIVES,OTHER                                        Initially Ordered On: 13 Aug 2010
bserved/Historical: HISTORICAL                                                   Quantity: 45
omments:            The reaction to this allergy was MILD (NO SQUELAE)           Days Supply: 90
                                                                                 Pharmacy: DAYTON
llergy Name:        TRAMADOL                                                     Prescription Number: 2718953
ocation:            DAYT29
ate Entered:        09 Mar 2011                                                  Medication: IBUPROFEN 600MG TAB
eaction:            URINARY RETENTION                                            Instructions: TAKE ONE TABLET BY MOUTH FOUR TIMES A DAY WITH FOOD
llergy Type:        DRUG                                                         Status: Active
A Drug Class:       NON-OPIOID ANALGESICS                                        Refills Remaining: 3
bserved/Historical: HISTORICAL                                                   Last Filled On: 20 Aug 2010
omments:            gradually worsening difficulty emptying bladder              Initially Ordered On: 01 Jul 2010
```

## Definition of data

" Data are values of qualitative or quantitative variables, belonging to a **set of items**. "

http://en.wikipedia.org/wiki/Data

- Set of items – sometimes called the population; the set of objects you are interested in
- Variables – a measurement or characteristic of an item
- Qualitative item – country of origin, gender, department, etc.
- Quantitative item – Q1 sales, salary, square feet, etc.

# Raw versus processed data

**Raw data**

- The original source of the data

- Often hard to use for data analyses

- Data analysis *includes* processing

- Raw data may only need to be processed once

http://en.wikipedia.org/wiki/Raw_data

**Processed data**

- Data that is ready for analysis

- Processing can include merging, subsetting, transforming, etc.

- There may be standards for processing

- All steps should be recorded

http://en.wikipedia.org/wiki/Computer_data_processing

# The tidy data

1. Each variable you measure should be in one column

2. Each different observation of that variable should be in a different row

3. There should be one table for each "kind" of variable

4. If you have multiple tables, they should include a column in the table that allows them to be linked

*Some other important tips*

- Include a row at the top of each file with variable names.

- Make variable names human readable AgeAtDiagnosis instead of AgeDx

- In general data should be saved in one file per table.

# Code module

- WEEK 5-1 Code module – downloading files
- WEEK 5-2 Code module – reading CSV and Excel
- WEEK 5-3 Code module – reading JSON files
- WEEK 5-4 Code module – SQL databases
- WEEK 5-5 Code module – SQL equivalents in R
- WEEK 5-6 Code module – writing data files

# Summary

- In WEEK 5 of Introduction to Data Science we continue to add useful items to our data science toolbox. This time, we added tools to access different data sources
- We saw how to download files from the web
- We read in files in CSV and Excel format
- We read in files in JSON format
- We saw how to scrape data off web pages
- We read in data from a SQL database
- We saw how to do SQL equivalents using R
- We wrote a new data file