

# Spreadsheet vs Database

---

Avoiding Data Disasters

Anne Pajon, CRUK-CI

# Spreadsheet

## The good...

- It's easy to **browse** data.
- It's easy to **manually enter and edit** data.
- It's easy to **share** copies of files.
- You have fine control over **visual presentation**.
- It has a very **flexible structure**.
- Formulas make it a **living document**.
- It has a **built-in suite of helpers** for charts, comments, spell checking, etc.
- It's relatively **easy to learn**.

## The not so good...

- It **lacks data integrity**. Because every cell is unique, things can get very **inconsistent**. What you see doesn't necessarily represent the underlying data. A number is not necessarily a number. Data is not necessarily data.
- It's not very good for **working with multiple datasets** in combination.
- It's not very good for answering **detailed questions** with your data.
- It **doesn't scale**. As the amount of data increases, performance suffers, and the visual interface becomes a liability instead of a benefit. It also has fixed limits on how big a spreadsheet and its cells can be.
- **Collaborating is hard**. It's hard to control versions and have a “master” set of data, especially when many people are working on the same project.

# Enter relational databases

## What is relational database?

It consists of a “**server**” that stores all your data (think of a huge library) and a **mechanism for querying** it (think of a reference librarian).

The querying is where **SQL** comes in, SQL stands for **Structured Query Language**, and it is a syntax for requesting things from the database. It's the language the reference librarian speaks.

The “relational” part is a hint that these databases care about **relationships between data**.

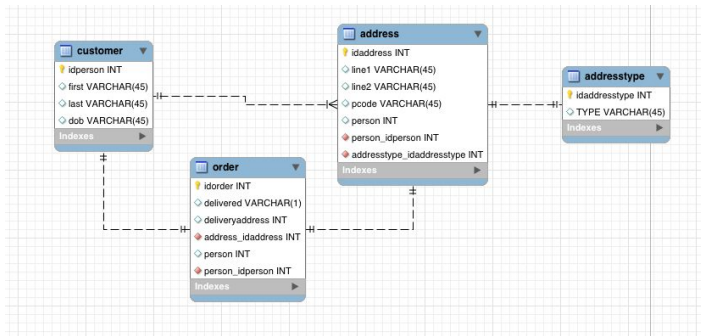


# The database mantra

Everything in its proper place.

A database encourages you to **store things logically**. Sometimes it forces you to.

Every database consists of **tables and relationships** between them. Think of a table like a single worksheet in an Excel file, except with more ground rules. A database table consists of **columns** and **rows**.



ID	First Name	Column	Name	Street Address	City	State
1	Tracey		am	7 East Walker Dr.	Raleigh	NC
2	Lucinda		George	789 Brewer St.	Cary	NC
3	Jerrold		Smith	211 St. George Ave.	Raleigh	NC
4	Brett		Newkirk	47 Hillsborough St.	Raleigh	NC
5	Chloe		Jones	23 Solo Ln.	Raleigh	NC
6	Quinton		Boyd	4 Cypress Cr.	Durham	NC
7	Alex		Hinton	1011 Hodge Ln.	Cary	NC
8	Nisha		Hall	123 Huntington St.	Raleigh	NC
9	Hillary		Clayton	2516 Newman	Raleigh	NC
10	Kiara		Williams	9014 Miller Ln.	Durham	NC
11	Katy		Jones	456 Denver Rd.	Cary	NC
12	Beatrix		Joslin	85 North West St.	Raleigh	NC
13	Mariah		Allen	12 Jupe	Raleigh	NC
14	Jennifer		Hill	2100 Field Ave.	Raleigh	NC
15	Jaleel		Smith	123 Hill Top Drive	Garner	NC

## Columns

Every column is given a **name** (like 'Address') and a defined **column type** (like 'Integer,' 'Date,' 'Date+Time', or 'Text').

*Columns define the structure of your data.*

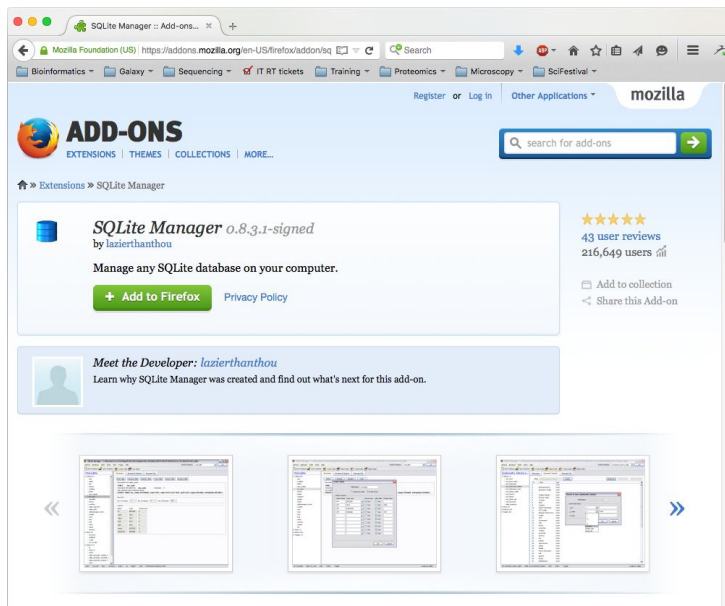
## Rows

Rows are the actual data in the table. Once you establish the column structure, you can add in as many rows as you like.

*Every row has a value for every column.*

# Where to start?

**SQLite** is a good way to get started. You can install the “**SQLite Manager**” addon for Firefox and do everything within the browser.



# File Management

---

Avoiding Data Disasters

Anne Pajon, CRUK-CI

# Use descriptive and informative **file names**



# File Names ... Best Practices

- Do not name all your data files '**data.xls**'
- Include any information that will allow you to distinguish your files from one another
  - Project or experiment name or acronym
  - Location/spatial coordinates
  - Researcher name/initials
  - Date or date range of experiment
  - Type of data
  - Conditions
  - Version number of file
  - Three-letter file extension for application-specific files
- Choose a consistent naming scheme across all your files
- Include in the directory a '**README.txt**' file that explains your naming format along with any abbreviations or codes you have used



# File Names ... Other Tips (1)

- **Avoid special characters** such as ~ ! @ # \$ % ^ & \* ( ) ` ; < > ? , [ ] { } ' " and |
- **Use short file names**, long ones do not work well with all types of software
- A **good format for date** designations is **YYYYMMDD** or **YYMMDD**
  - All of your files stay in chronological order, even over the span of many years
- Use **leading zeros** for clarity and to make sure files sort in sequential order
  - For example, use "001, 002, ...010, 011 ... 100, 101, etc." instead of "1, 2, ...10, 11 ... 100, 101, etc."

## File Names ... Other Tips (2)

- **Do not use spaces.** Some software will not recognize file names with spaces, and file names with spaces must be enclosed in quotes when using the command line. Other options include:
  - Underscores, e.g. `file_name.xxx`
  - Dashes, e.g. `file-name.xxx`
  - No separation, e.g. `filename.xxx`
  - Camel case, where the first letter of each section of text is capitalized, e.g. `FileName.xxx`

Choose **file formats** that will ensure long-term access



# File Formats ... Best Practices

- Save data in a **non-proprietary** (open) file format when possible
  - Usable on diverse platforms and by multiple applications
  - Export your data as tab separated file
- Unencrypted
- Uncompressed
- In common usage by the research community
- Preferred formats
  - Tabular data: CSV, TXT

Track different  
**versions** of your  
documents



# Data versioning

Versioning refers to saving new copies of your files when you make changes allowing you to reverse or roll back those changes or retrieve specific versions of your files later

- Simple file versioning
- Simple software options
- Advanced software options

# Simple File versioning

- Manually save new versions when you make significant changes
  - Include a version number, e.g. "v1," "v2," or "v2.1" into file names
- This works well, only if...
  - No need to keep lots of different versions
  - Only one person working on these files
  - Always access these files from one location

# Simple Software Options

- Use Google Drive's word processing, spreadsheet and presentation
  - Any time you edit files, new versions are saved as you go
  - Version information includes who was editing the file and when the new version was created

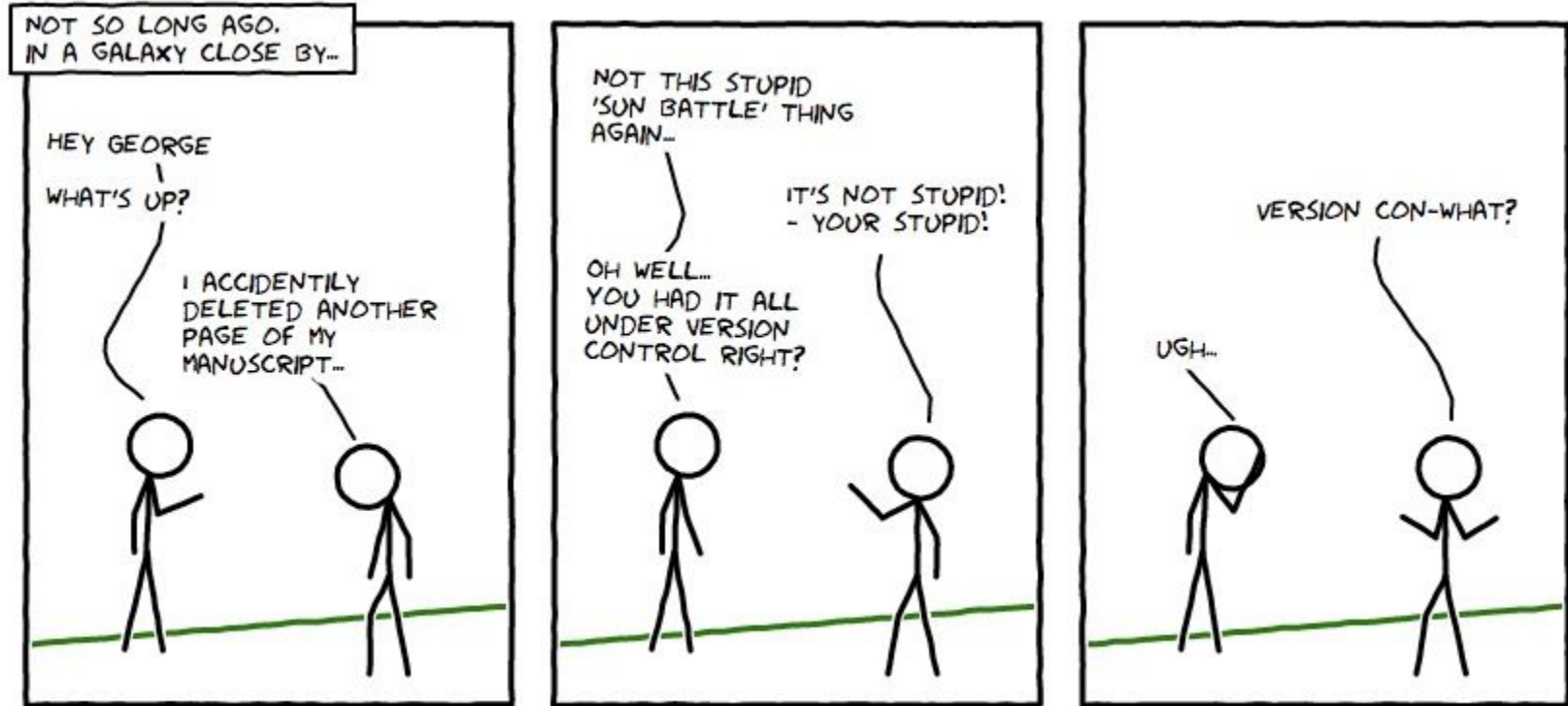


- Dropbox free version keeps track of all of the changes you make for 30 days
  - The paid Packrat version keeps track of every change you ever make to the files in your Dropbox





# Advanced Software Options



# Advanced Software Options

**Version control** systems like subversion and git are frequently used for groups writing software and code, but can be used for any kind of files or projects. Many people share their git repositories on GitHub.

**Version control** is the management of changes to documents, computer programs, and other collections of information. Changes are usually identified by a number named the "**revision number**".

Each revision is associated with a timestamp and the person making the change. Revisions can be compared, restored, and with some types of files, merged.



# Simple collaboration from your desktop

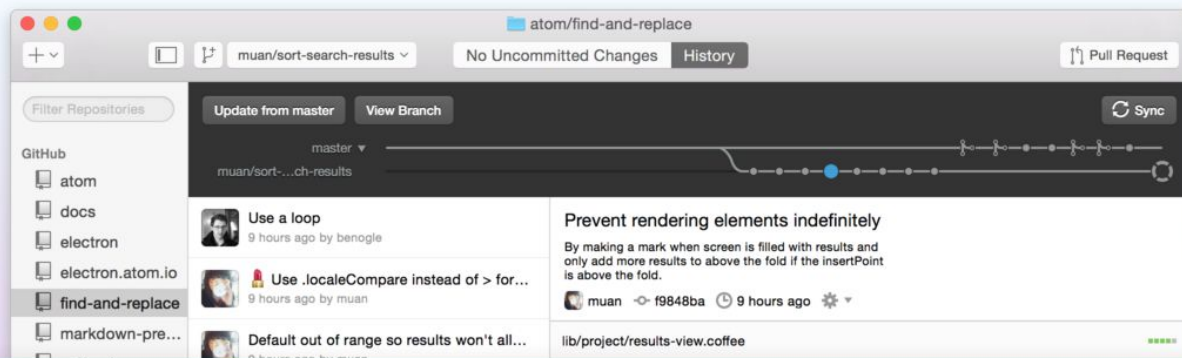
GitHub Desktop is a seamless way to contribute to projects on **GitHub** and **GitHub Enterprise**.

Available for [Mac](#) and [Windows](#)

**Download GitHub Desktop**

OS X 10.9 or later

By clicking the Download button you agree to the [End-User License Agreement](#)



Your GitHub workflow in one native app



Clone repositories



Create branches



Commit changes



Share code

# Training material on github

<https://github.com/>

**Search for repository:**

avoid-data-disaster  
or bioinformatics-core-shared-  
training



bioinformatics-core-shared-training / avoid-data-disaster

Unwatch 4 ★ Star 0 🍴 Fork 0

Code Issues 0 Pull requests 0 Wiki Pulse Graphs Settings

how to organise and keep your data tidy <http://bioinformatics-core-shared-training.github.io/avoid-data-disaster> — Edit

12 commits 2 branches 0 releases 1 contributor

Branch: master New pull request New file Find file HTTPS https://github.com/bioinf Download ZIP

markdunning add open refine project Latest commit 0d58eeb 9 minutes ago

images	add cruk banner	13 days ago
AvoidingDataDisastersFileMgmt.pdf	add file management talk	3 days ago
Presentation2.pptx	Add draft presentations	13 days ago
README.md	couple of tweaks	21 days ago
example1.google-refine.tar.gz	add open refine project	9 minutes ago
example1.xlsx	update example1 and add modified version	4 days ago
example1_modified.xlsx	update example1 and add modified version	4 days ago
example2.xlsx	Add draft presentations	13 days ago
principles.pdf	Add draft presentations	13 days ago
principles_nonotes.pdf	update Andy's slides	2 hours ago

README.md

## Avoiding data disasters

# Reference

- Data best practices  **STANFORD UNIVERSITY LIBRARIES**
  - <http://library.stanford.edu/research/data-management-services/data-best-practices>
- Excel vs Databases 
  - <http://schoolofdata.org/2013/11/07/sql-databases-vs-excel/>