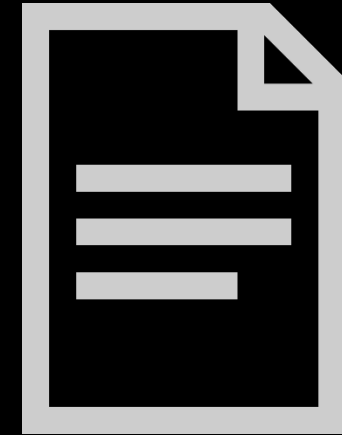# Intro to SQL

Fernando Rios, PhD.
Research Data Management Specialist, UA Libraries

May 2, 2025

# About me

- Academic
  - Geographic information systems & science
  - Computational hydrogeology
  - Computational physics
- IT
  - Software development
  - Database administration

- Research data management
  - Data organization, metadata
  - Data management plans
  - Data policy compliance
- Manage U of A's research data repository, ReDATA
  - Funder and journal data sharing requirements
  - University data retention policies
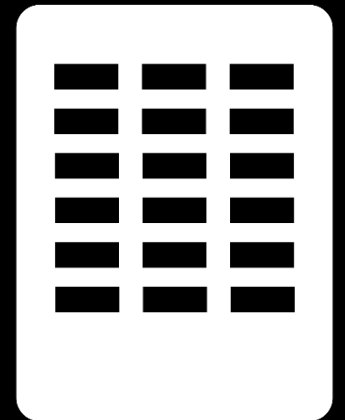  - Data curation for more reusable data
- Teaching

# Data Storage – Un- and semi-structured information

- Measurements about stuff, information about a groups of things, time series
  - Physical lab notebooks
  - Images, videos
  - Individual tables in an unstructured spreadsheet
  - Information in proprietary formats

- Unorganized – not easily machine-readable
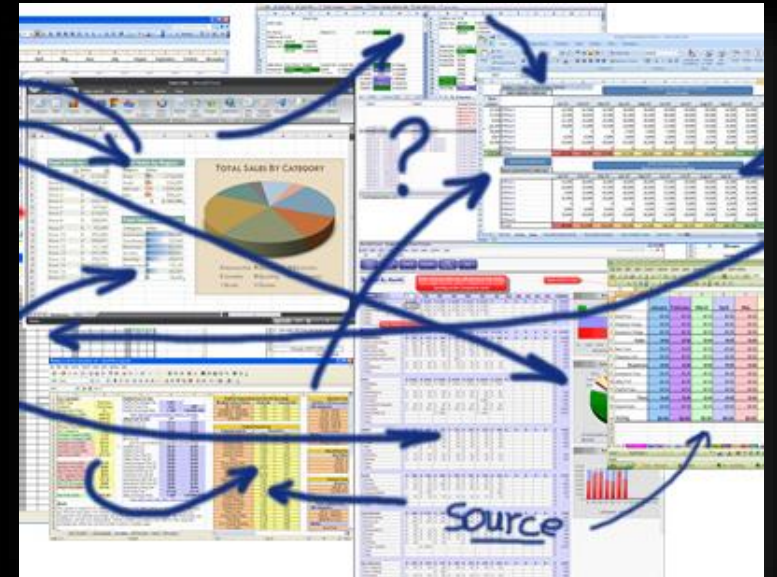- Limited analysis capabilities
- Doesn't scale well

# Structured Spreadsheets

- Standardizes item characteristics, records them in a tabular structure

- Few constraints = Need to follow best-practices  https://osf.io/vew32

- Calculations in Excel error-prone as complexity increases

- When to use
  - Simple column/row summaries
  - Doing a quick filter or plot
  - Simple column calculations / one-off pivot tables
  - Exploring data
  - Low consequences

# Databases

- Adds more structure, constraints
- Compact "formulas" and repeatable processes = always desirable
- When to use
  - Data can be reasonably put in a table format
  - Analyzing tens of thousands of rows
  - When your formulas start getting too complex
  - Need to ensure data integrity
  - Management capabilities

# Databases – key terms



| id | ISSN-L | ISSNs | PublisherId | Journal_Title |
|---|---|---|---|---|
| 0 | 2056-9890 | 2056-9890 | 1 | Acta Crystallographica Section E Crystallographic Communications |
| 1 | 2077-0472 | 2077-0472 | 2 | Agriculture |
| 2 | 2073-4395 | 2073-4395 | 2 | Agronomy |
| 3 | 2076-2615 | 2076-2615 | 2 | Animals |
| 4 | 2076-3417 | 2076-3417 | 2 | Applied Sciences |
| 5 | 2306-5354 | 2306-5354 | 2 | Bioengineering |
| 6 | 2079-7737 | 2079-7737 | 2 | |
| 7 | 2079-6374 | 2079-6374 | 2 | |

Field • Table • Record • Value

- Values: store a single piece of information
- Fields: single kind and type of information
  - temperature, age, address, etc.
  - integer, text, date, etc.
- Record: set of related fields containing specific values
- Usually have more than one table that are related

# SQL, RDBMS?

- Structured Query Language (SQL)
- Relational database systems (RDBMS)
  - Structured relationships between tables + management layer
  - Use SQL to manipulate data in the DB
  - More common in business than science
    - Metadata could be in a DB – even if actual data isn't

# What can we do with RDBMS + SQL?

- Aggregating, summarizing, combining, filtering, adding
- Robust, reproducible
- Ensure better data quality
- Scalable and fast
- Management capabilities
  - Access controls
  - Concurrency
  - Auditing
  - Backups
  - Distributed data

# Other kinds of databases

- Key-value
  - Simple, fast
  - Python dictionaries, Redis
- Document-based
  - Less structured, more heterogeneous
  - MongoDB
- Graph databases
  - Entities (nodes), predicates (edge), objects (another node). Complex relationships, knowledge graphs
  - Node4j
- Wide-column
  - Rows can have different columns within the same table
  - Google BigTable

# Let's get started

https://tinyurl.com/s4u42xb3

# Relationships – key concepts



Database schema

**articles**

| | |
|---|---|
| id | INTEGER |
| Title | TEXT |
| Authors | TEXT |
| DOI | TEXT |
| URL | TEXT |
| Subjects | TEXT |
| ISSNs | TEXT |
| Citation | TEXT |
| LanguageId | INTEGER |
| LicenceId | INTEGER |
| Author_Count | INTEGER |
| First_Author | TEXT |
| Citation_Count | INTEGER |
| Day | INTEGER |
| Month | INTEGER |
| Year | INTEGER |

**journals**

| | |
|---|---|
| id | INTEGER |
| ISSN-L | TEXT |
| ISSNs | TEXT |
| PublisherId | INTEGER |
| Journal_Title | TEXT |

**languages**

| | |
|---|---|
| id | INTEGER |
| Language | TEXT |

**licences**

| | |
|---|---|
| id | INTEGER |
| Licence | TEXT |

**publishers**

| | |
|---|---|
| id | INTEGER |
| Publisher | TEXT |

**Cardinality**
* = many
0..1 = zero or one

**Primary key**
uniquely identify a row

**Foreign key**
Links rows in one table to those in another

# Joins



**Journals**

| Journal | PubID |
|---------|-------|
| ISPRS | 2 |
| JDDT | 4 |
| JDMDS | 5 |
| JLPEA | 2 |
| OPCJ | 6 |

**Publishers**

| Id | Name |
|----|------|
| 1 | Consejo Superior de Investigaciones Científicas |
| 2 | MDPI AG |
| 3 | Soc. of Pharmaceutical Technocrats |
| 4 | Int. Union of Crystallography |

**(Inner) Join**

| Journal | PubID | Id | Name |
|---------|-------|----|------|
| ISPRS | 2 | 2 | MDPI AG |
| JDDT | 4 | 4 | Int. Union of Crystallography |
| JLPEA | 2 | 2 | MDPI AG |

# Key points of RDBMS

- Data integrity
  - Keeps data separate from analysis
  - Reduces accidents
  - Data types help with quality control
- Data is stored in related tables
  - Reduce redundancy, increase data quality
- Many RDBMS with different strengths
  - SQLite – small compact. Embedded devices
  - MySQL, PostgreSQL – general purpose
  - MS SQL Server, Oracle – Strong management capabilities, enterprise solution

# Do I have to use SQL?