# Tutorials on Data Management

## Lesson 4: Data Collection Entry and Manipulation



CC image by Cobalt123 on Flickr

Data Entry and Manipulation

DataONE

# Lesson Topics

- Best Practices for Creating Data Files
- Data Entry Options
- Data Manipulation Options



CC image by JISC on Flickr

DataONE

# Learning Objectives

- Recognize inconsistencies that can make a dataset difficult to understand and/or manipulate
- Describe characteristics of stable data formats and list reasons for using these formats
- Identify data entry tools
- Identify validation measures that can be performed as data is entered
- Describe the basic components of a relational database

DataONE

# The Data Life Cycle

Plan

Analyze

Collect

Assure

Integrate

Discover

Describe

Preserve

DataONE

# Goals of Data Entry

- Create data sets that are:
  - Valid
  - Organized to support ease of use



CC image by Travis S on Flickr

DataONE

# Example: Poor Data Entry



- **Inconsistency between data collection events**
  - **Location of Date information**
  - **Inconsistent Date format**
  - **Column names**
  - **Order of columns**

# Example: Poor Data Entry



- Inconsistency between data collection events
  - Different site spellings, capitalization, spaces in site names—hard to filter
  - Codes used for site names for some data, but spelled out for others
  - Mean1 value is in Weight column
  - Text and numbers in same column – what is the mean of 12, "escaped < 15", and 91?

DataONE

# Best Practices



- Columns of data are consistent: only numbers, dates, or text
- Consistent Names, Codes, Formats (date) used in each column
- Data are all in one table, which is much easier for a statistical program to work with than multiple small tables which each require human intervention

DataONE

# Best Practices

- Create descriptive column names without spaces or special characters

    - Soil T30 ⭢ Soil_Temp_30cm

    - Species-Code ⭢ Species_Code (avoid using -,+,*,^ in column names. Some software may interpret these symbols as an operator)

- Use a descriptive file name.  For instance, a file named SEV_SmallMammalData_v.5.25.2010.csv indicates the project the data is associated with (SEV),  the theme of the data (SmallMammalData) and also when this version of the data was created (v.5.25.2010).   This name is much more helpful than a file named mydata.xls.

# Best Practices

- Missing data
  - Preferably leave field empty (NULL = no value)
  - In numeric fields, use a distinct value such as 9999 to indicate a missing value
  - In text fields, use NA ("Not Applicable" or "Not Available")
  - Use Data flags in a separate column to qualify missing value

| Date | Time | NO3_N_Conc | NO3_N_Conc_Flag |
|------|------|------------|-----------------|
| 20081011 | 1300 | 0.013 | |
| 20081011 | 1330 | 0.016 | |
| 20081011 | 1400 | | M1 |
| 20081011 | 1430 | 0.018 | |
| 20081011 | 1500 | 0.001 | E1 |

M1 = missing; no sample collected

E1 = estimated from grab sample

DataONE

# Best Practices

- Enter complete lines of data



Sorting an Excel file with empty cells is not a good idea!
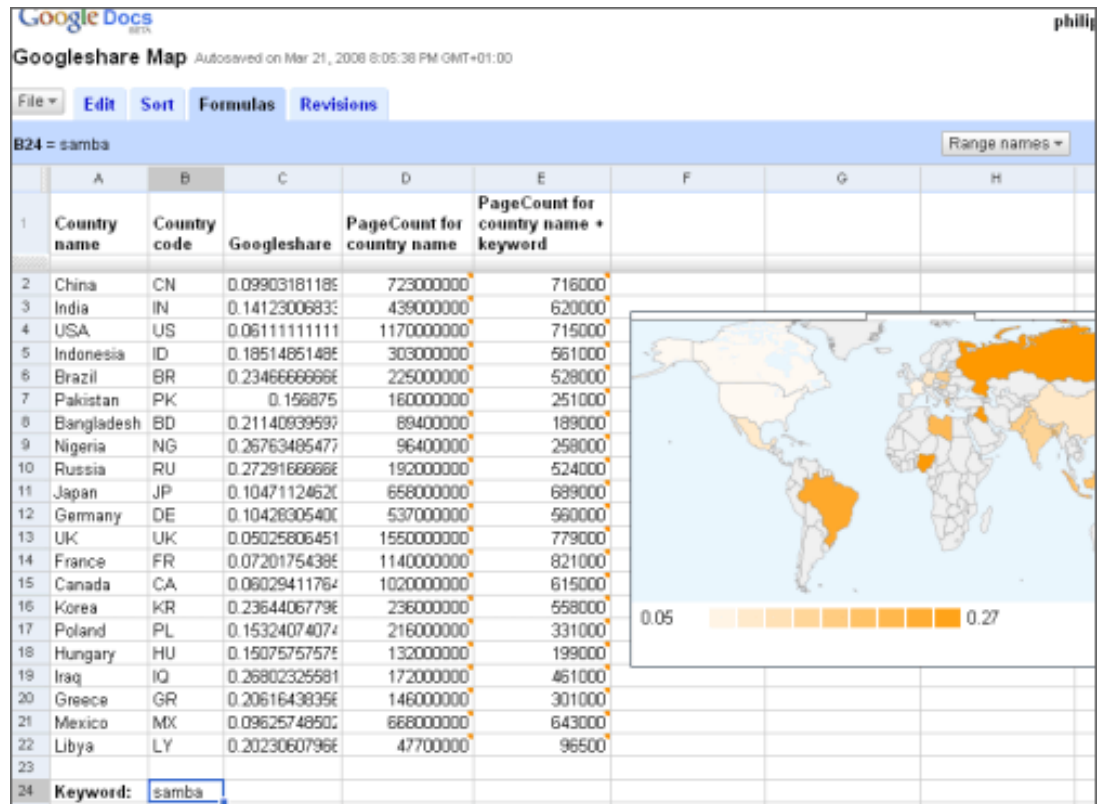
DataONE

# Best Practices

- For the long term, store data in a consistent format that can be read well in to the future and that can be used by any application now or in the future

- Appropriate file types include:
    - Non-proprietary: Open, documented standard
    - Common usage by research community: Standard representation (ASCII, Unicode)
    - Unencrypted
    - Uncompressed

- ASCII formatted files will be readable into the future
    - Use ASCII (comma-separated) for tabular data

DataONE

# References

1. Best Practices for Preparing Environmental Data Sets to Share and Archive. September 2010. Les A. Hook, Suresh K. Santhana Vannan, Tammy W. Beaty, Robert B. Cook, and Bruce E. Wilson.
   http://daac.ornl.gov/PI/BestPractices-2010.pdf

# Data Entry Tools

- Googledocs Forms
- Spreadsheets

# Googledocs Forms

Data Entry and Manipulation

# Data Entry Tools: Excel

# Excel: Data Validation

# Spreadsheet vs. Relational Database

- Great for charts, graphs, calculations
- Flexible about cell content type—cells in same column can contain numbers or text
- Lack record integrity--can sort a column independently of all others)
- Easy to use – but harder to maintain as complexity and size of data grows

- Easy to query to select portions of data
- Data fields are typed – For example, only integers are allowed in integer fields
- Columns cannot be sorted independently of each other
- Steeper learning curve than a spreadsheet

DataONE

# What is a relational database?

**Sample sites**

*siteID
site_name
latitude
longitude
description

**Samples**

*sampleID
siteID
sample_date
speciesID
height
flowering
flag
comments

**Species**

*speciesID
species_name
common_name
family
order

- A set of tables
- Relationships
- A command language

DataONE

# Database Features: Explicit control over data types

| Date | Site | Height | Flowering |
|---|---|---|---|
| <dates only> | <text only> | < real numbers only> | < 'y' and 'n' only> |
| | | | |
| | | | |
| | | | |

**Advantages**
- **quality control**
- **performance**

DataONE

# Relationships are defined between tables

| Date | Site | Species | Flowering? |
|---|---|---|---|
| 2/13/2010 | A | BOGR2 | y |
| 2/13/2010 | B | HODR | y |
| 4/15/2010 | B | BOER4 | y |
| 4/15/2010 | C | PLJA | n |

| Site | Latitude | Longitude |
|---|---|---|
| A | 34.1 | -109.3 |
| B | 35.2 | -108.6 |
| C | 32.6 | -107.5 |

Mix and Match data on the fly

| Date | Site | Species | Flowering? | Latitude | Longitude |
|---|---|---|---|---|---|
| 2/13/2010 | A | BOGR2 | y | 34.1 | -109.3 |
| 2/13/2010 | B | HODR | y | 35.2 | -108.6 |
| 4/15/2010 | B | BOER4 | y | 35.2 | -108.6 |
| 4/15/2010 | C | PLJA | n | 32.6 | -107.5 |

DataONE

# Powerful Command Language called Structured Query Language (SQL)

This table is called SoilTemp

| Date | Plot | Treatment | SensorDepth | Soil_Temperature |
|------|------|-----------|-------------|------------------|
| 2010-02-01 | C | R | 30 | 12.8 |
| 2010-02-01 | B | C | 10 | 13.2 |
| 2010-02-02 | C | R | 0 | 6.3 |
| 2010-02-02 | A | N | 0 | 15.1 |

SQL examples: Select Date, Plot, Treatment, SensorDepth, Soil_Temperature from SoilTemp where Date = '2010-02-01'

| Date | Plot | Treatment | SensorDepth | Soil_Temperature |
|------|------|-----------|-------------|------------------|
| 2010-02-01 | C | R | 30 | 12.8 |
| 2010-02-01 | B | C | 10 | 13.2 |

Select * from SoilTemp where Treatment='N' and SensorDepth='0'

| Date | Plot | Treatment | SensorDepth | Soil_Temperature |
|------|------|-----------|-------------|------------------|
| 2010-02-02 | A | N | 0 | 15.1 |

Data Entry and Manipulation

DataONE

# Data Entry with a Database

- Forms can be created that make entering data in to a relational database as easy as entering it in to Excel. The screenshot below shows embedded forms that were quickly generated in MS Access for adding data to three tables in a database of plant cover measurements

# Conclusion

- Be aware of Best Practices when designing data file structures
- Choose a data entry method that allows some validation of data as it is entered
- Consider investing time in learning how to use a database if datasets are large or complex



CC image by fo.ol on Flickr

# If you want to try a database:

- Consider trying one of these:
  - Personal, single-user databases can be developed in MS Access, which is stored as a file on the user's computer.  MS Access comes with easy GUI tools to create databases, run queries, and write reports.
  - A more robust database that is free, accommodates multiple users and will run on Windows or Linux is MySQL.   GUI interfaces for MySQL include phpMyadmin (free) and Navicat (inexpensive).

DataONE

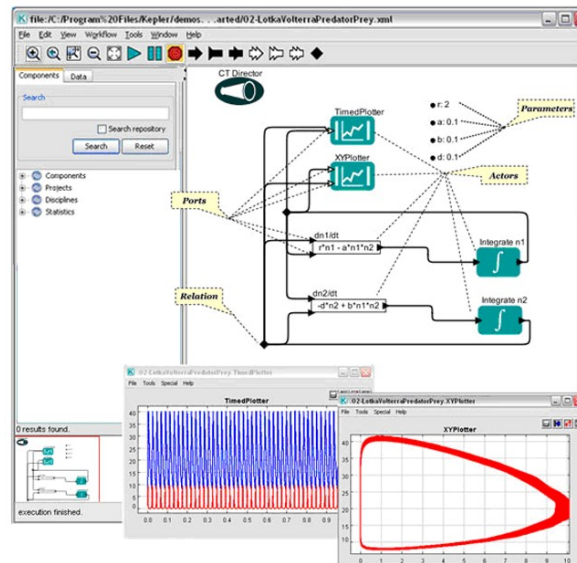# To learn more about designing a relational database:

- Database Design for Mere Mortals: A Hands-On Guide to Relational Database Design (2nd Edition)  by Michael J. Hernandez.  Addison-Wesley.  2003.

# Data Manipulation

- Useful for analyzing, subsetting and transforming data
- Can be used to quality assure data
- Options include SAS, SPSS, R, and Matlab
  - Not Free
    - SAS:  Has outstanding support
    - SPSS:  Has a user-friendly GUI
    - Matlab: Analysis and Visualization platform that has "toolboxes" available for different disciplines, such as modeling or genomic analyses

DataONE

# R

- Free (http://www.r-project.org/index.html)
- Produces publication quality graphics
- Lots of forums from which to get help
- Software (such as Kepler for developing workflows) will integrate analytical components written in R

DataONE

The full slide deck may be downloaded from:

http://www.dataone.org/education-modules

Suggested citation:

DataONE Education Module: Data Entry and Manipulation. DataONE. Retrieved Nov12, 2012. From http://www.dataone.org/sites/all/documents/L04_DataEntryManipulation.pptx

Data Entry and Manipulation

DataONE