# Tidy-ing Your Data: Simple Steps for Reproducible Research

Jeanette Clark

National Center for Ecological Analysis and Synthesis

jclark@nceas.ucsb.edu

Synthesis Research

Training

Data Science Infrastructure

Synthesis Research

Training

Data Science Infrastructure

# NCEAS Learning Hub

**National Center for Ecological Analysis and Synthesis**

A knowledge-sharing community where researchers can learn the latest data skills and technologies to increase efficiency, productivity, transparency, and collaborative capacity.

*Courses:* Fee-based and grant-supported intensive data science workshops

*Mentored Programs:* Experiential residential and remote learning programs to build skills in data and open science

*Resources:* Extensive online curricula, webinars, training materials and best practices

*Partnerships:* Customized workshops and collaborative initiatives in data science training

- Some simple guidelines for effective data management

- How to recognize and tidy untidy data

- Using tidy data in analysis

NCEAS
National Center for Ecological Analysis and Synthesis

# Data management is for everyone!



Audrey McCombs - https://notebooks.dataone.org/networked-lod/week-5-the-really-cool-thing/

Your data don't need to be of a particular type, size, or complexity before you start implementing data management practices
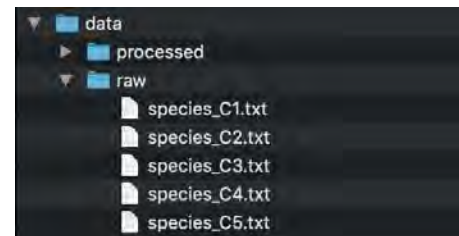
**NCEAS**
National Center for Ecological Analysis and Synthesis

# Data management is for everyone!



You don't have to be using a relational database system to benefit from the concepts of relational data models (aka tidy data)

# Simple Guidelines for Data Management

- Use a scripted program

- Nonproprietary formats

- Keep a raw version of data

- Descriptive names

- Header line

- Plain ASCII text

.csv, .txt



Borer, E. T. et al, (2009), Some Simple Guidelines for Effective Data Management. The Bulletin of the Ecological Society of America. doi:10.1890/0012-9623-90.2.205

NCEAS
National Center for Ecological Analysis and Synthesis

# Simple Guidelines for Data Management

- Design to add rows, not columns

- Each column should contain only one type of information

- Record a single piece of data only once; separate information collected at different scales into different tables. In other words, use a relational model

**NCEAS**
National Center for Ecological Analysis and Synthesis

# Data model diversity

- There are lots of data models besides tabular data
  - multiband raster
  - matrices
  - spatial vector

NCEAS
National Center for Ecological Analysis and Synthesis

# Recognizing untidy data

# Characteristics of tidy data

## Observations

- Separate tables for each entity measured

NCEAS
National Center for Ecological Analysis and Synthesis

# Recognizing untidy data

# Characteristics of tidy data

## Observations

- Separate tables for each entity measured
- Each row represents a single observed entity

# Recognizing untidy data

# Characteristics of tidy data

## Observations

- Separate tables for each entity measured
- Each row represents a single observed entity

## Variables

- All values in a column are of the same type

# Recognizing untidy data

# Characteristics of tidy data

## Observations

- Separate tables for each entity measured
- Each row represents a single observed entity
- Observations (rows) are all unique

## Variables

- All values in a column are of the same type
- All columns pertain to the same observation (row)
- Each column represents either an identifying or measured variable

NCEAS
National Center for Ecological Analysis and Synthesis

# Characteristics of tidy data



variables        observations        values

R for data science: import, tidy, transform, visualize, and model data. H Wickham, G Grolemund – 2016. https://r4ds.had.co.nz/

# Recognizing untidy data

| id | date | site | elev | sp1code | sp1height | sp2code | sp2height |
|----|------|------|------|---------|-----------|---------|-----------|
| 1 | 2017-10-10 | 1 | 3.7 | DAPU | 4.6 | DAMA | 4.5 |
| 2 | 2017-09-05 | 2 | 3.2 | DAMA | 3.5 | DAPU | 3.9 |

- Each row contains observations about multiple entities (site characteristics and species observations)

- A new species observation would add a column (wide format)

# Tidying our data

| id | date | site | elev | sp1code | sp1height | sp2code | sp2height |
|----|------|------|------|---------|-----------|---------|-----------|
| 1 | 2017-10-10 | 1 | 3.7 | DAPU | 4.6 | DAMA | 4.5 |
| 2 | 2017-09-05 | 2 | 3.2 | DAMA | 3.5 | DAPU | 3.9 |

- What are the observed entities?
  - plant species
  - site characteristics
- What are the variables associated with those observations?
  - height
  - elevation

**NCEAS**
National Center for Ecological Analysis and Synthesis

# Tidying our data

| id | date | site | spcode | height |
|---|---|---|---|---|
| 1 | 2017-10-10 | 1 | DAPU | 4.6 |
| 2 | 2017-09-05 | 2 | DAMA | 3.5 |
| 3 | 2017-10-10 | 1 | DAMA | 4.5 |
| 4 | 2017-09-05 | 2 | DAPU | 3.9 |

| site | name | elev |
|---|---|---|
| 1 | Taku | 3.7 |
| 2 | Lituya | 3.2 |

- Individual species observations
  - identifying variables: id, date, site, spcode
  - measured variables: height

- Site observations where species occurred
  - identifying variables: site, name
  - measured variables: elev

# Tidying our data

| id | date | site | elev | sp1code | sp1height | sp2code | sp2height |
|---|---|---|---|---|---|---|---|
| 1 | 2017-10-10 | 1 | 3.7 | DAPU | 4.6 | DAMA | 4.5 |
| 2 | 2017-09-05 | 2 | 3.2 | DAMA | 3.5 | DAPU | 3.9 |

| id | date | site | spcode | height |
|---|---|---|---|---|
| 1 | 2017-10-10 | 1 | DAPU | 4.6 |
| 2 | 2017-09-05 | 2 | DAMA | 3.5 |
| 3 | 2017-10-10 | 1 | DAMA | 4.5 |
| 4 | 2017-09-05 | 2 | DAPU | 3.9 |

| site | name | elev |
|---|---|---|
| 1 | Taku | 3.7 |
| 2 | Lituya | 3.2 |

- Add rows not columns

- Separate information collected at different scales into different tables

- Record a single piece of data only once

# Benefits of normalized data



| id | date | site | elev | sp1code | sp1height | sp2code | sp2height |
|----|------------|------|------|---------|-----------|---------|-----------|
| 1  | 2017-10-10 | 1    | 3.7  | DAPU    | 4.6       | DAMA    | 4.5       |
| 2  | 2017-09-05 | 2    | 3.2  | DAMA    | 3.5       | DAPU    | 3.9       |

| id | date | site | spcode | height |
|----|------------|------|--------|--------|
| 1  | 2017-10-10 | 1    | DAPU   | 4.6    |
| 2  | 2017-09-05 | 2    | DAMA   | 3.5    |
| 3  | 2017-10-10 | 1    | DAMA   | 4.5    |
| 4  | 2017-09-05 | 2    | DAPU   | 3.9    |

| site | name | elev |
|------|--------|------|
| 1    | Taku   | 3.7  |
| 2    | Lituya | 3.2  |

- Search and filter rows

- Describe columns more precisely

- Optimize storage

- Enforce data integrity
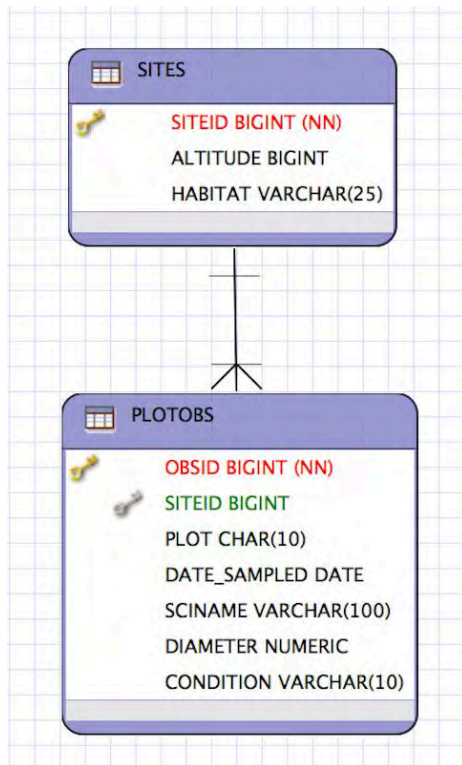


NCEAS
National Center for Ecological Analysis and Synthesis

# Using normalized data

| id | date | site | spcode | height |
|---|---|---|---|---|
| 1 | 2017-10-10 | 1 | DAPU | 4.6 |
| 2 | 2017-09-05 | 2 | DAMA | 3.5 |
| 3 | 2017-10-10 | 1 | DAMA | 4.5 |
| 4 | 2017-09-05 | 2 | DAPU | 3.9 |

| site | name | elev |
|---|---|---|
| 1 | Taku | 3.7 |
| 2 | Lituya | 3.2 |

- **Primary key**
  - unique identifier for each observation within an entity

- **Foreign Key**
  - reference to a primary key in another table

NCEAS
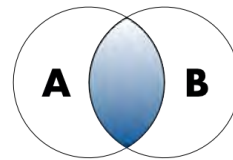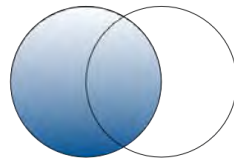National Center for Ecological Analysis and Synthesis

# Entity-relationship diagrams



- Draw relationships between tables concisely

- Used in database management systems

# Merging normalized data
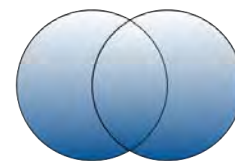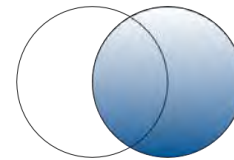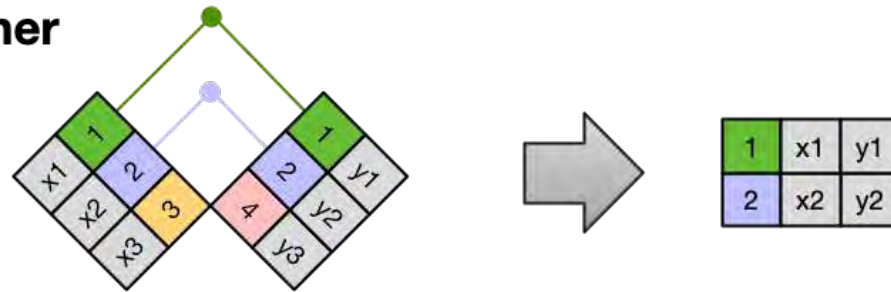
# Merging normalized data



R for data science: import, tidy, transform, visualize, and model data. H Wickham, G Grolemund – 2016. https://r4ds.had.co.nz/

# Merging normalized data

# Merging normalized data

# Merging normalized data

# Merging normalized data

Left join

| id | date | site | spcode | height |
|---|---|---|---|---|
| 1 | 2017-10-10 | 1 | DAPU | 4.6 |
| 2 | 2017-09-05 | 2 | DAMA | 3.5 |
| 3 | 2017-10-10 | 1 | DAMA | 4.5 |
| 4 | 2017-09-05 | 2 | DAPU | 3.9 |

| site | name | elev |
|---|---|---|
| 1 | Taku | 3.7 |
| 2 | Lituya | 3.2 |

| id | date | site | spcode | height | name | elev |
|---|---|---|---|---|---|---|
| 1 | 2017-10-10 | 1 | DAPU | 4.6 | Taku | 3.7 |
| 2 | 2017-09-05 | 2 | DAMA | 3.5 | Lituya | 3.2 |
| 3 | 2017-10-10 | 1 | DAMA | 4.5 | Taku | 3.7 |
| 4 | 2017-09-05 | 2 | DAPU | 3.9 | Lituya | 3.2 |

NCEAS
National Center for Ecological Analysis and Synthesis

# A not-so-reproducible workflow

# Building a reproducible workflow



Field data sheets

Tidy, raw data

| id | date | site | spcode | height |
|----|------------|------|--------|--------|
| 1 | 2017-10-10 | 1 | DAPU | 4.6 |
| 2 | 2017-09-05 | 2 | DAMA | 3.5 |
| 3 | 2017-10-10 | 1 | DAMA | 4.5 |
| 4 | 2017-09-05 | 2 | DAPU | 3.9 |

| site | name | elev |
|------|--------|------|
| 1 | Taku | 3.7 |
| 2 | Lituya | 3.2 |

| id | date | site | spcode | height | name | elev |
|----|------------|------|--------|--------|--------|------|
| 1 | 2017-10-10 | 1 | DAPU | 4.6 | Taku | 3.7 |
| 2 | 2017-09-05 | 2 | DAMA | 3.5 | Lituya | 3.2 |
| 3 | 2017-10-10 | 1 | DAMA | 4.5 | Taku | 3.7 |
| 4 | 2017-09-05 | 2 | DAPU | 3.9 | Lituya | 3.2 |

Quality controlled, derived data, often merged or summarized

Figures, tables, maps

NCEAS National Center for Ecological Analysis and Synthesis

# When to start?



- Thinking about your data model **early** helps you be more efficient at every stage of the data lifecycle

- Its never too late to tidy things up!


NCEAS
National Center for Ecological Analysis and Synthesis

**Reproducible Research Techniques for Snythesis**

A five day immersion into widely adopted R-based tools for open science

DataONE    NCEAS

**Details**

**Dates:**
February 3-7, 2020
May 11-15, 2020

**Location:**
NCEAS
Santa Barbara, CA

**Cost:**
$2100
Includes: 5 days of instruction, refreshments and lunch.

**www.nceas.ucsb.edu/learning-hub/short-course**

NCEAS
National Center for Ecological Analysis and Synthesis

# Questions?

## www.nceas.ucsb.edu/learning-hub/short-course

Jeanette Clark
National Center for Ecological Analysis and Synthesis
jclark@nceas.ucsb.edu