# Data Cleaning for Analysis and Publication Using the OpenRefine Software Package

**Arctic Data Center**, CUAHSI, DataONE, Environmental Data Initiative, GBIF, **iDigBio**, NEON, Neotoma

Jeanette Clark, Deborah Paul

**#datahelpdesk**
Ecological Society of America 2020
Career Central
August
**http://bit.ly/datahelpesa2020**

# Advancing the Digitization of Biological Collections
## iDigBio Hub and Thematic (Museum) Collection Networks



(Meta)data Aggregator
**Community Building**
Collections Data

PEOPLE IN THE LOOP
*people graphic by Dorothy Allard*

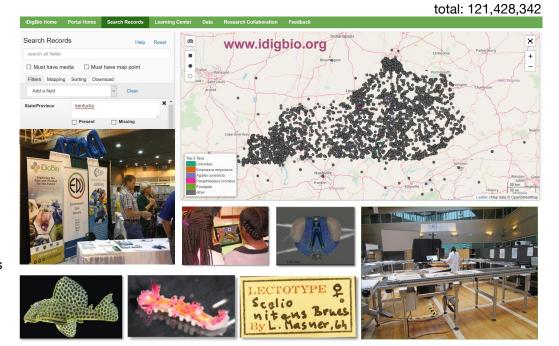**Digitization**
Workflows & Protocols
Dissemination

**Research Use**
Cyberinfrastructure
Tool collaboration
Portal development
ENM workshop
Research focus
Data quality
APIs

**Training**
Biodiversity informatics
Data skills and literacy
Collections software
Imaging
Project Management

total: 121,428,342

www.idigbio.org

**Education Outreach**
Citizen Science
K-12 materials
Undergraduate
Fossil Clubs
Mentor teachers

**Methods**
Workshops
Webinars
Symposia
Conferences
Working Groups
Short Courses
Adobe Connect
Listservs
Publications
Social Media **@idigbio**

3

# What do we mean by "Clean" Data?

More Data from More Sources =

- Structural Issues
- Inconsistent/unclear missing values
- Mixed data in single columns
- Mixed data types in single columns
- Ambiguous data values
- Data you can't use

Journal of eScience Librarianship
putting the pieces together: theory and practice

**Common Errors in Ecological Data Sharing**

Karina E. Kervin,[1] William K. Michener,[2] Robert B. Cook[3]

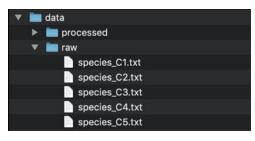[1] University of Michigan, Ann Arbor, MI, USA
[2] University of New Mexico, Albuquerque, NM, USA
[3] Oak Ridge National Laboratory, Oak Ridge, TN, USA

# 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

# 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, create a relational database.

Borer, E. T. et al, (2009), Some Simple Guidelines for Effective Data Management. The Bulletin of the Ecological Society of America

# Recognizing untidy data

| | | main trunks | reiterated trunks | limbs | branches | leaves | | | | | dry masses (kg) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| species | tree | kg | kg | kg | kg | kg | | type | species | main trunk | reiteration | limb | branch | leaf | TOTAL | % total |
| SESE | Atlas | 255144.9 | 46020.6 | 5477.7 | 13433.2 | 1101.2 | | tree | SESE | 3569312 | 213247 | 53714 | 230945 | 17192 | 4084409 | 95.3491 |
| SESE | Ballantine | 221966.4 | 7651.6 | 5922.9 | 11210.0 | 1084.8 | | tree | PSME | 135815 | 0 | 0 | 8338 | 961 | 145114 | 3.3876 |
| SESE | Bell | 253246.4 | 5454.3 | 5792.6 | 48500.7 | 1043.4 | | tree | THSE | 31799 | 0 | 0 | 6343 | 864 | 39006 | 0.9105 |
| SESE | Broken Top | 130928.9 | 4805.2 | 1608.1 | 5137.4 | 729.9 | | tree | ACMA | 4444 | 0 | 0 | 925 | 264 | 5634 | 0.1315 |
| SESE | Buena Vista | 128833.0 | 3486.5 | 0.0 | 8552.1 | 518.4 | | tree | UMCA | 2921 | 0 | 0 | 937 | 273 | 4131 | 0.0964 |
| SESE | Demeter | 155896.0 | 11085.6 | 3204.2 | 10054.1 | 768.7 | | shrub | RUSP | 0 | 0 | 0 | 1974 | 686 | 2660 | 0.0620 |
| SESE | Epimetheus | 226987.0 | 12915.7 | 1797.2 | 13585.2 | 1029.4 | | fern | POMU | 0 | 0 | 0 | 0 | 1271 | 1271 | 0.0296 |
| SESE | Iluvatar | 349586.6 | 65003.9 | 12315.6 | 13987.0 | 1461.8 | | shrub | VAOV | 0 | 0 | 0 | 526 | 26 | 552 | 0.0129 |
| SESE | Kronos | 134154.1 | 12204.4 | 7232.7 | 5036.1 | 597.3 | | shrub | COCO | 0 | 0 | 0 | 284 | 6 | 289 | 0.0067 |
| SESE | Pleiades I | 182385.2 | 3735.0 | 1935.2 | 10846.6 | 762.2 | | fern | POSC | 0 | 0 | 0 | 107 | 89 | 196 | 0.0045 |
| SESE | Pleiades II | 235838.8 | 11183.4 | 4306.0 | 11306.5 | 877.7 | | tree | RHPU | 100 | 0 | 0 | 44 | 18 | 162 | 0.0037 |
| SESE | Prometheus | 239414.0 | 25228.9 | 1612.6 | 12458.2 | 1086.0 | | herb | OXOR | 0 | 0 | 0 | 0 | 112 | 112 | 0.0026 |
| SESE | Rhea | 143710.4 | 487.8 | 730.1 | 5524.2 | 691.2 | | shrub | VAPA | 0 | 0 | 0 | 94 | 4 | 99 | 0.0023 |
| SESE | Zeus | 243365.7 | 2885.5 | 1620.4 | 19104.7 | 954.3 | | tree | PISI | 0 | 0 | 0 | 1 | 0 | 1 | 0.0000 |
| SESE | 3 | 1761.3 | 0.0 | 0.0 | 87.6 | 41.4 | | tree | CHLA | 0 | 0 | 0 | 1 | 0 | 1 | 0.0000 |
| SESE | 4 | 6312.0 | 356.0 | 73.5 | 214.1 | 43.8 | | shrub | GASH | 0 | 0 | 0 | 0 | 0 | 0 | 0.0000 |
| SESE | 5 | 206.0 | 0.0 | 0.0 | 8.7 | 2.5 | | shrub | SACA | 0 | 0 | 0 | 0 | 0 | 0 | 0.0000 |
| SESE | 6E | 18697.4 | 0.0 | 0.0 | 1055.2 | 66.3 | | | | 3744390 | 213247 | 53714 | 250519 | 21767 | 4283636 | |
| SESE | 6W | 14651.5 | 7.7 | 0.0 | 626.3 | 49.6 | | | | | | | | | | proportion |
| SESE | 11 | 614.4 | 0.0 | 0.0 | 28.1 | 17.0 | | | | main trunk | reiteration | limb | branch | leaf | total | geophytic |
| SESE | 12 | 232.1 | 0.0 | 0.0 | 11.2 | 10.3 | | | SESE geo | 3569312 | 213247 | 53714 | 230945 | 17192 | 4084409 | 1.00 |
| SESE | 18 | 15632.0 | 0.0 | 0.0 | 946.3 | 106.8 | | | SESE epi | 0 | 0 | 0 | 0 | 0 | 0 | |
| SESE | 19 | 11805.5 | 0.0 | 0.0 | 770.1 | 80.3 | | | PSME geo | 135815 | 0 | 0 | 8338 | 961 | 145114 | 1.00 |
| SESE | 20 | 309.5 | 0.0 | 0.0 | 12.5 | 5.9 | | | PSME epi | 0 | 0 | 0 | 0 | 0 | 0 | |
| SESE | 22 | 25618.3 | 0.0 | 0.0 | 1504.0 | 120.2 | | | TSHE geo | 31740 | 0 | 0 | 6332 | 860 | 38932 | 0.99 |
| SESE | 23 | 463.7 | 0.0 | 0.0 | 18.9 | 4.5 | | | TSHE epi | 59 | 0 | 0 | 12 | 4 | 74 | |

# Characteristics of Tidy Data

**Observations**

- **Separate tables for each entity measured**

# 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

| species | tree | main trunks kg | reiterated trunks kg | limbs kg | branches kg | leaves kg | | type | species | main trunk | reiteration | limb | branch | leaf | TOTAL | % total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | dry masses (kg) | | | |
| SESE | Atlas | 255144.9 | 46020.6 | 5477.7 | 13433.2 | 1101.2 | | tree | SESE | 3569312 | 213247 | 53714 | 230945 | 17192 | 4084409 | 95.3491 |
| SESE | Ballantine | 221966.4 | 7651.6 | 5922.9 | 11210.0 | 1084.8 | | tree | PSME | 135815 | 0 | 0 | 8338 | 961 | 145114 | 3.3876 |
| SESE | Bell | 253246.4 | 5454.3 | 5792.6 | 48500.7 | 1043.4 | | tree | THSE | 31799 | 0 | 0 | 6343 | 864 | 39006 | 0.9105 |
| SESE | Broken Top | 130928.9 | 4805.2 | 1608.1 | 5137.4 | 729.9 | | tree | ACMA | 4444 | 0 | 0 | 925 | 264 | 5634 | 0.1315 |
| SESE | Buena Vista | 128833.0 | 3486.5 | 0.0 | 8552.1 | 518.4 | | tree | UMCA | 2921 | 0 | 0 | 937 | 273 | 4131 | 0.0964 |
| SESE | Demeter | 155896.0 | 11085.6 | 3204.3 | 10054.1 | 768.7 | | shrub | RUSP | 0 | 0 | 0 | 1974 | 686 | 2660 | 0.0620 |
| SESE | Epimetheus | 226987.0 | 12915.7 | 1797.2 | 13585.2 | 1029.4 | | fern | POMU | 0 | 0 | 0 | 0 | 1271 | 1271 | 0.0296 |
| SESE | Iluvatar | 349586.6 | 65003.9 | 12315.6 | 13987.0 | 1461.8 | | shrub | VAOV | 0 | 0 | 0 | 526 | 26 | 552 | 0.0129 |
| SESE | Kronos | 134154.1 | 12204.4 | 7232.7 | 5036 | | | | | | 0 | 0 | 284 | 6 | 289 | 0.0067 |
| SESE | Pleiades I | 182385.2 | 3735.0 | 1935.2 | 10846 | | | | | | 0 | 0 | 107 | 89 | 196 | 0.0045 |
| SESE | Pleiades II | 235838.8 | 11183.4 | 4306.0 | 11306 | | | | | | 0 | 0 | 44 | 18 | 162 | 0.0037 |
| SESE | Prometheus | 239414.0 | 25228.9 | 1612.6 | 12458 | | | | | | 0 | 0 | 0 | 112 | 112 | 0.0026 |
| SESE | Rhea | 143710.4 | 487.8 | 730.1 | 5524 | | | | | | 0 | 0 | 94 | 4 | 99 | 0.0023 |
| SESE | Zeus | 243365.7 | 2885.5 | 1620.4 | 1910 | | | | | | 0 | 0 | 1 | 0 | 1 | 0.0000 |
| SESE | 3 | 1761.3 | 0.0 | 0.0 | 87 | | | | | | | 0 | 1 | 0 | 1 | 0.0000 |
| SESE | 4 | 6312.0 | 356.0 | 73.5 | 214 | | | | | | | 0 | 0 | 0 | 0 | 0.0000 |
| SESE | 5 | 206.0 | 0.0 | 0.0 | | | | | | | | 0 | 0 | 0 | 0 | 0.0000 |
| SESE | 6E | 18697.4 | 0.0 | 0.0 | 1055 | | | | | | 213247 | 53714 | 250519 | 21767 | 4283636 | |
| SESE | 6W | 14651.5 | 7.7 | 0.0 | 62 | | | | | | | | | | | proportion |
| SESE | 11 | 614.4 | 0.0 | 0.0 | 2 | | | | | | eration | limb | branch | leaf | total | geophytic |
| SESE | 12 | 232.1 | 0.0 | 0.0 | 11.2 | 10.3 | | | SESE geo | 3569312 | 213247 | 53714 | 230945 | 17192 | 4084409 | 1.00 |
| SESE | 18 | 15632.0 | 0.0 | 0.0 | 946.3 | 106.8 | | | SESE epi | 0 | 0 | 0 | 0 | 0 | 0 | |
| SESE | 19 | 11805.5 | 0.0 | 0.0 | 770.1 | 80.3 | | | PSME geo | 135815 | 0 | 0 | 8338 | 961 | 145114 | 1.00 |
| SESE | 20 | 309.5 | 0.0 | 0.0 | 12.5 | 5.9 | | | PSME epi | 0 | 0 | 0 | 0 | 0 | 0 | |
| SESE | 22 | 25618.3 | 0.0 | 0.0 | 1504.0 | 120.2 | | | TSHE geo | 31740 | 0 | 0 | 6332 | 860 | 38932 | 0.99 |
| SESE | 23 | 463.7 | 0.0 | 0.0 | 18.9 | 4.5 | | | TSHE epi | 50 | 0 | 0 | 13 | 4 | 74 | |

All the same variable? No.

# Characteristics of Tidy Data

## Observations

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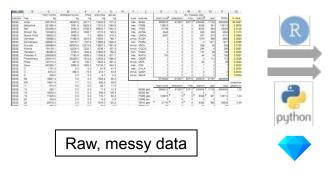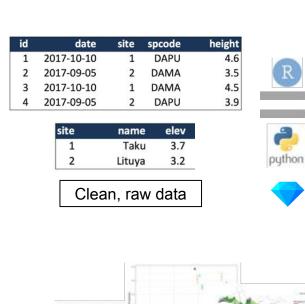- **Observations (rows) are all unique**

## Variables

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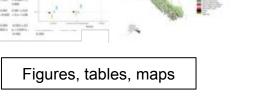- **Each column represents either an identifying or measured variable**

14

# A not-so-reproducible workflow

# Building a reproducible workflow



Raw, messy data

Clean, raw data

Merged/summarized derived data

Figures, tables, maps

# Where data come from matters! (a sample)

- Excel
  Issues that may arise from these data sources are our focus today
  - Automatic conversion of gene names to dates or floating point numbers*
  - Date values can be converted when transferring data between operating systems and applications
- Text (e.g. CSV) & Excel
  - Free-form structure - lack of enforcement of column-row structure, type consistency
- Text (e.g. CSV)
  - Inconsistent structure - quotes, commas, missing values, spaces
- Database
  - Enforced structure - tables, column typing
  - Specialized methods for interaction (pros and cons to this)

* Ziemann, M., Eren, Y. & El-Osta, A. Gene name errors are widespread in the scientific literature. *Genome Biology* **17**, (2016).

# The ESAUSSEE Data Help Desk

## who we are and how to find us

Amber Budden, @aebudden, @DataONE_org, aebudden@epscor.unm.edu

Deborah Paul, @idbdeb, @idigbio, dpaul@fsu.edu

Dmitry Schigel, @dschigel, @GBIF, dschigel@gbif.org

Karl Benedict, @kbene, kbene@unm.edu, president@esipfed.org

Kristen Vanderbilt, @vanderbik, @EDIgotdata, krvander@fiu.edu

Kyle Copas, @kylecopas, @GBIF, kcopas@gbif.org

Laura Brenskelle, @lbrensk, @idigbio, lbrensk@ufl.edu

Margaret O'Brien @ , @EDIgotdata, margaret.obrien@ucsb.edu

Megan Jones, @MeganAHJones, @NEON_sci, mjones01@battelleecology.org

Rebekah Wallace, www.eddmaps.org, bekahwal@uga.edu

# Data lessons compiled - inspired by workshop
*Georeferencing for Research Use of Museum Collections Data*

- **Data mapped to standards**
  - supports use and re-use (e.g. Darwin Core DwC, Ecological Metadata Language EML)
  - standards help with data validation and cleaning
- **Data have issues**
  - what are some you have experienced
  - need to be addressed before applying research methods
  - keep raw data raw
  - track your changes
- **Data visualization is key**
  - QGIS lessons
  - Open Refine
  - R, etc.

Carabidae (beetles) of California

20

# (Fun!) features and functions in Open Refine

- runs on your computer (not in the cloud)
- data formats supported
- raw data
- column manipulation
- text facet
- routine cleaning (white space)
- clustering
- step-wise editable task script
- APIs
- regular expressions
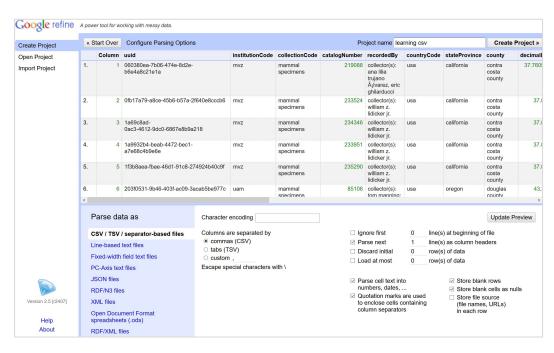- export
- share project files

# Open Refine - getting started is quick and easy

- download and install
- launch
- import your data
- your raw data is NOT touched
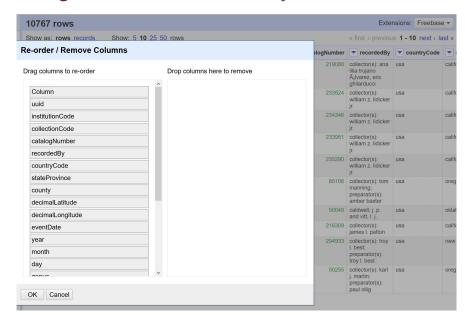- supported data formats
- subset data
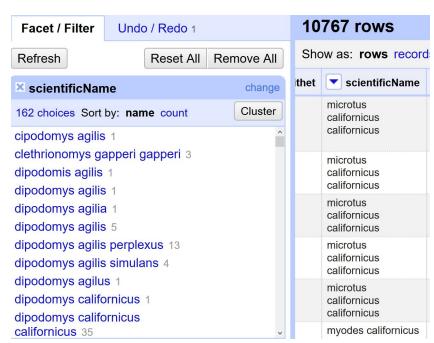
# Open Refine - managing columns

- **reorganize columns** easily

# Open refine - text facet
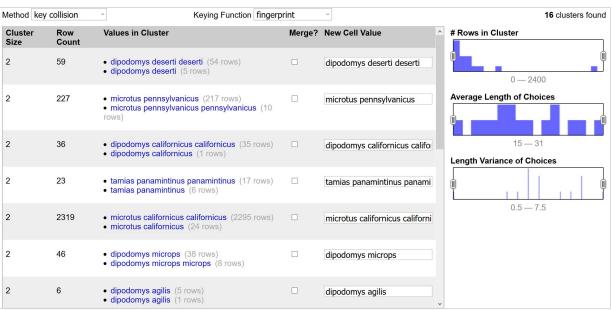*lists* and *counts* the distinct values in a column

# Open Refine - the magic of clustering algorithms

*or how to find issues that abc sort won't
and fix them all at once - no hunting*

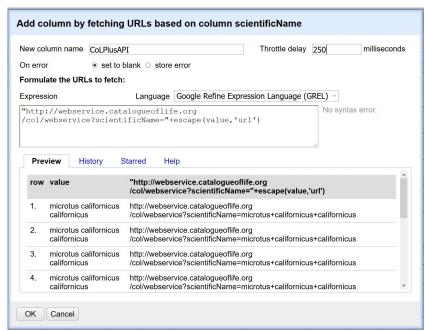# Open Refine - manages pesky white spaces

# Open Refine - add data to your data using APIs
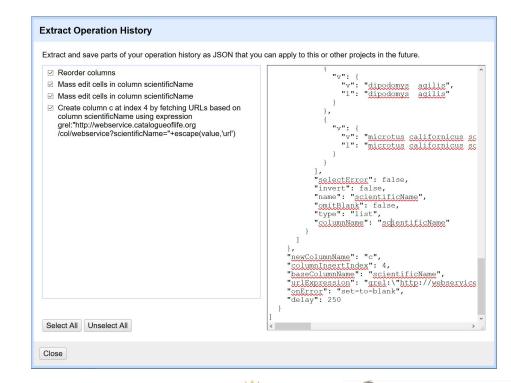*that's application programming interface*

# Open Refine - saves your steps

- ❏ *supports reproducibility*

- ❏ *tracks your work for you*

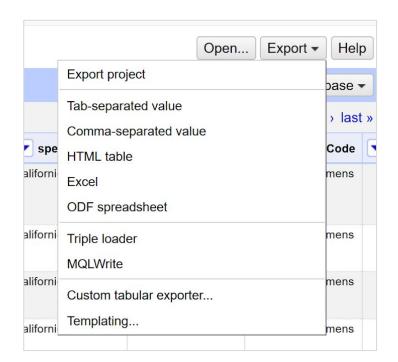- ❏ *easy to go back to earlier steps with confidence*

# Open Refine - export your data, share project files

*select the format
export subsets too
and project files*

# Open Refine - make some friends

- share this tool with students, friends, families, colleagues
- **imagine future tools, think beyond spreadsheets**

## Increase Reproducibility and Productivity using tools like Open Refine

Magic is here.
Ask for it, plan for it.

# Looking for next steps now?
## *R, Open Refine, and Data Management resources*

- The #datahelpdesk is ready to offer data assistance!
- #CareerCentral Q and A: Wednesday, August 5th, 9:30-10:30 PDT (12:30-1:00 EDT)
- Data Help Desk Wiki **https://bit.ly/datahelpesa2020**
- Data Carpentry lessons
- on Twitter #ESA2020 #datahelpdesk