

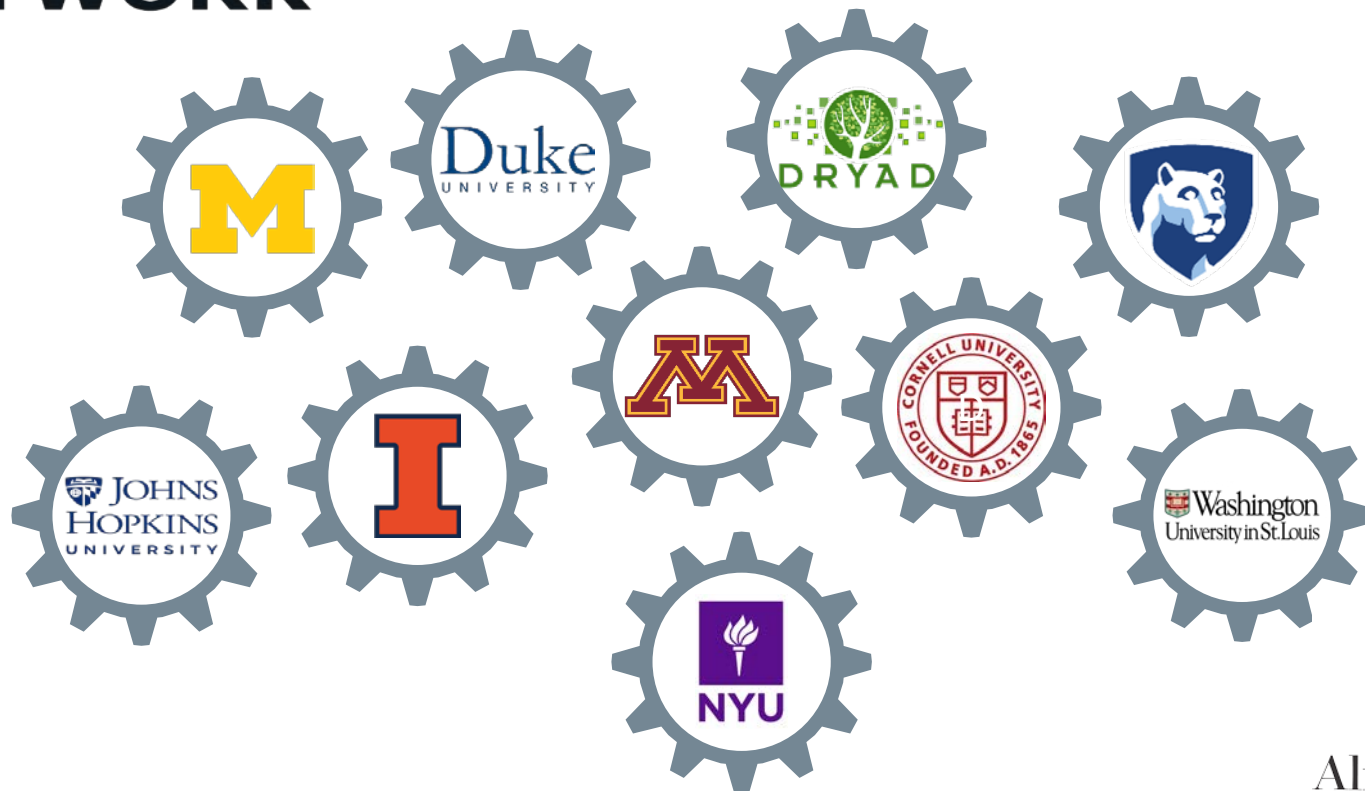
The Data Curation Network

Cynthia Hudson Vitale
Pennsylvania State University



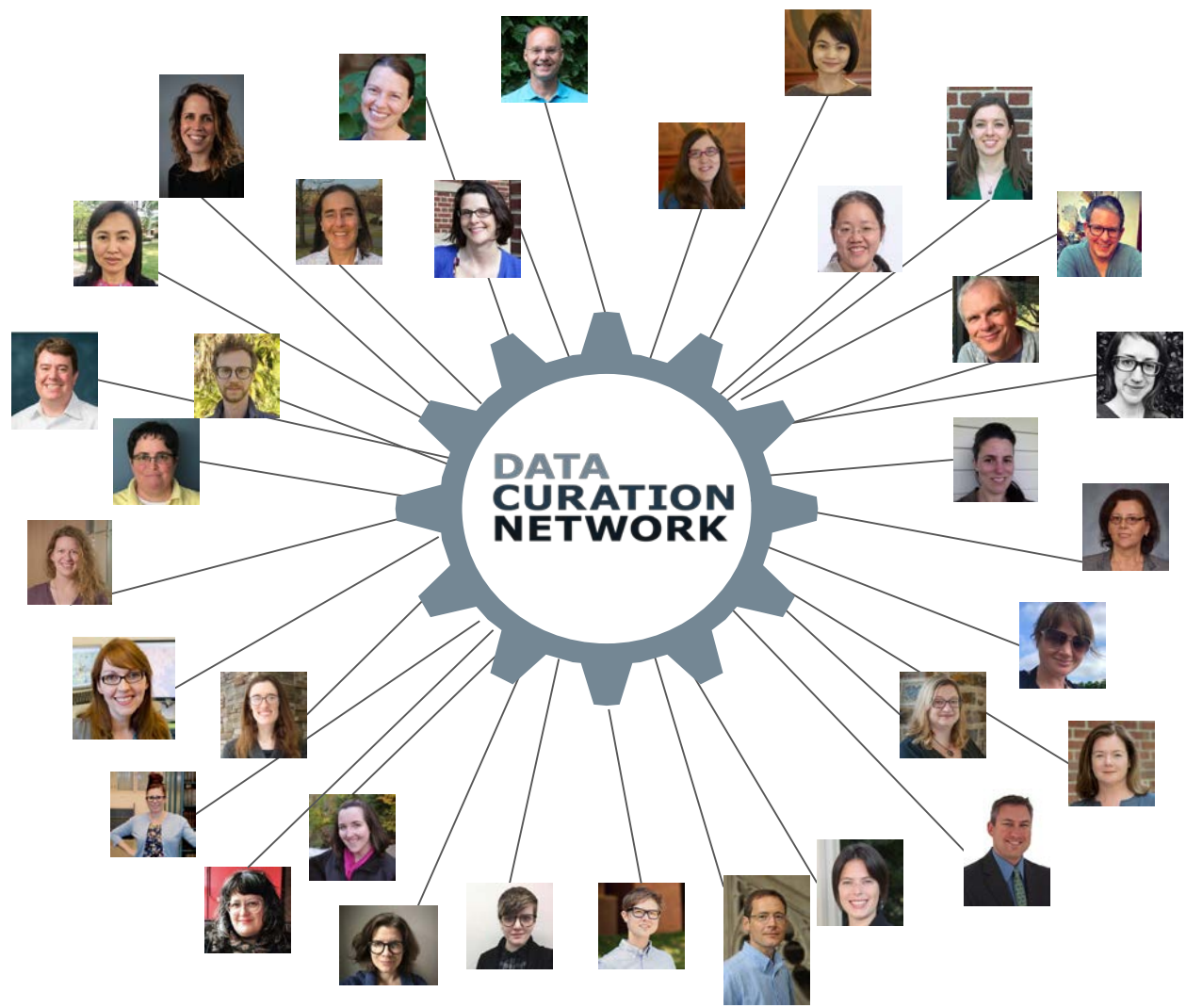
Data Curation Network All Hands Meeting 2019

DATA CURATION NETWORK



Alfred P. Sloan
FOUNDATION

**DATA
CURATION
NETWORK**



Mission

“The Data Curation Network will enable researchers that are faced with a growing number of requirements to ethically share their research data in ways that make it findable, accessible, interoperable and reusable (FAIR).”

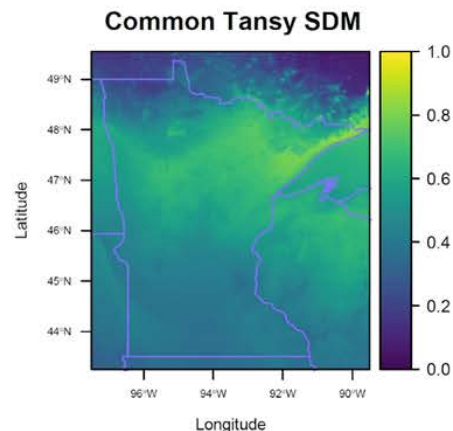
What is data curation?

Data curators enrich research data publications and ensure the data are FAIR, by:

- Finding and adding missing files and documentation
 - Screening for privacy disclosure risk
 - Detecting and fixing code and other quality assurance issues
 - Transforming file formats for long term access
 - Arranging and describing files
 - Reviewing and augmenting metadata
-

Species Distribution Models and Joint Species Distribution Models of Nine Invasive Species in North America

Lake, Thomas, A; Briscoe Runquist, Ryan, D; Moeller, David, A (2019)



Persistent link to this item

<https://doi.org/10.13020/z71w-jx69>

<http://hdl.handle.net/11299/206482>

Services

[Full Metadata \(xml\)](#)

[View Usage Statistics](#)

Title

Species Distribution Models and Joint Species Distribution Models of Nine Invasive Species in North America

Published Date

2019-08-27

Authors

Lake, Thomas, A

Briscoe Runquist, Ryan, D

Moeller, David, A

Group

Moeller Lab

Author Contact

Briscoe Runquist, Ryan, D (rbriscoe@umn.edu)

Type

Dataset

Abstract

Species Distribution Models (SDM) and Joint Species Distribution Models (JSDM) for nine invasive species in North America. Species SDMs include Brown Knapweed (*Centaurea jacea*), Black Swallowwort (*Cynanchum louiseae*), Common Tansy (*Tanacetum vulgare*), Common Teasel (*Dipsacus fullonum*), Wild Parsnip (*Pastinaca sativa*), and Dalmatian Toadflax (*Linaria dalmatica*). Species JSDMs include Japanese Hops (*Humulus japonicus*), Oriental Bittersweet (*Celastrus orbiculatus*), and Narrowleaf Bittercress (*Cardamine impatiens*). These models aim to predict the current and future habitat suitability of nine invasive species in North America. Models were constructed with the

Species Distribution Models and Joint Species Distribution Models of Nine Invasive Species in North America

Lake, Thomas, A; Briscoe Runquist, Ryan, D; Moeller, David, A (2019)

Data Curator:Melinda Kernik

Spatial Data Analyst and Curator, University of Minnesota Libraries and U-spatial

- Identified 4 missing files and 2 unnecessary files
- Created folder organization and hierarchy
- Clarified inconsistent file naming and variable naming conventions
- Added citations for data reuse from EddMaps and GBIF
- Developed a data dictionary for variable names

Dear Thomas,

Thank you for your submission to the Data Repository for the U of M (DRUM). I have been looking over a copy of your submission and have some tests on the copies of the files. I have some questions and suggestions.

1) "Hum_Jap_FullState_19Feb2019.gri" is a zipped folder. Is it supposed to be a .gri file? The JSDM documentation says the zipped folder. Would you please provide documentation?

2) There are two files named "Hum_Jap_FullState_19Feb2019.gri". If so, should I delete one of them?

3) The xxx_coordinates.csv files are not added to the documentation. There is also one xml in the folder. Should it be added to the documentation?

4) The "SDM DRUM Data" folder is included. Are there any additional documentation or process for creating them?

5) Variable definitions. A little more information about the input files. For example, what data are you using? The version of MaxEnt.

Hi,

Thank you for your help.

1. The file Hum_Jap_FullState_19Feb2019.gri should be included in the dataset. Every .gri file should have an associated .GRI file. I have attached the Hum_Jap_FullState_19Feb2019.GRI file.

2. The two copies of "Hum_Jap_FullState_19Feb2019.gri" should be removed.

3. The species_coordinate_definitions.xml file should be removed.

4. I will fill out and return the documentation for files, source version used.

5. My input on the folder organization is that most users interested in understanding the distribution of the species would want to see the folder organization.

Thanks for your help,

-Thomas

Hello Thomas,

Thanks for the additional documentation and .gri file! I am still working on checking/fixing some of the larger subfolder uploads and will return to it early next week. A few more questions:

1) Regarding occurrences files:
a. There is a mixture of xx_occ, xx_bittersweet, coordinates.csv, and japanese_hops_occurrences.csv. For most of the species, they are named "xx_occ". Three files contain a mixture of "xx_bittersweet" and "xx_hops".
b. What is the source of this information?

Good Morning,

Here are my comments regarding your questions.

1a. The three .csv files that contain additional fields (bittersweet_occurrences.csv, japanese_hops_occurrences.csv, and japanese_hops_coordinates.csv) are used to construct the JSDM models. The other .csv files should only contain two variables (lon, lat coordinate pairs), and are used to construct the SDM models. These files are correct.

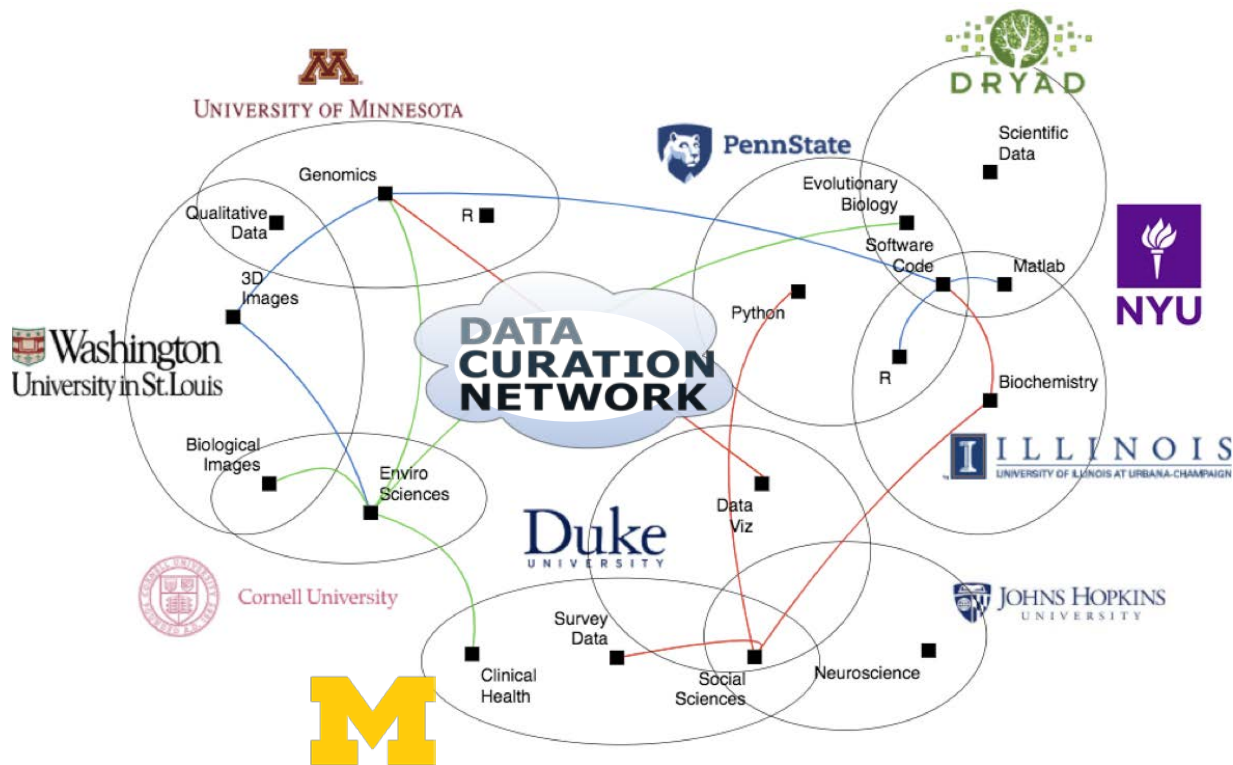
1b. The sources used to obtain species (lon, lat) coordinate pairs are from the online databases EddMaps and GBIF.

2. The .rda files are only the JSDM models (not predictions), and should be in the same folder as the JSDM contents. They can go in the same folder as the JSDM raster projection files, as part of the larger JSDM project.

3. I believe the "19" at the end of the BioClim19.grd/gri refers to the total number of bioclimatic variables in the dataset. Every bioclimatic dataset (current and future) contains the 19 bioclimatic variables as referenced on worldclim.

4a. The ENM Eval Model Results.csv file represents the results of each individual model run. These are the results of how changing model parameters influences the output model statistics. This is sort of an autogenerated output from MaxEnt. I modified the original code from the R ENMEval package (<https://github.com/bobmuscarella/ENMEval>) to output additional columns to the autogenerated output file (ie TSS bin, Boyce bin, tp, fn, tn, fn, sens bin ...).

DCN Expert Network



DCN CURATE Steps

DCN Curators will take **CURATE** steps for each data set, that in

C Check data files and read documentation

U Understand the data (try to), if not...

R Request missing information or changes

A Augment the submission with metadata for findability

T Transform file formats for reuse and long-term preservation

E Evaluate and rate the overall submission for FAIRness.

Table A1. Draft checklist of DCN CURATE steps and FAIRness scorecard

CURATE Actions	Curation Checklist
Check data files and read documentation <ul style="list-style-type: none">Review the content of the data files (e.g., open and run the files or code).Verify all metadata provided by the author and review the available documentation.	<input type="checkbox"/> Files open as expected <ul style="list-style-type: none"><input type="checkbox"/> Issues _____ <input type="checkbox"/> Code runs as expected <ul style="list-style-type: none"><input type="checkbox"/> Produces minor errors<input type="checkbox"/> Does not run and/or produces many errors <input type="checkbox"/> Metadata quality is rich, accurate, and complete <ul style="list-style-type: none"><input type="checkbox"/> Metadata has issues _____ <input type="checkbox"/> Documentation Type (<i>circle</i>) Readme / Codebook / Data Dictionary / Other: _____ <ul style="list-style-type: none"><input type="checkbox"/> Missing/None<input type="checkbox"/> Needs work
Understand the data (or try to) <ul style="list-style-type: none">Check for quality assurance and usability issues such as missing	<i>Varies based on file formats and subject domain. For example....</i>

DCN Growth and Sustainability

- Curated 50 data sets since Jan 1, 2019!
- 2 new members in Year 2
- Aim to add two more in Year 3
- Canada and Dutch groups planning stages to launch their own network
- Exploring fiscal and administrative models to support beyond grant

