# Using Statistical Analysis to Calculate the Size of Biodiversity Literature

**Alicia Esquivel**

NDSR Resident, Chicago Botanic Garden

June 6, 2017 | Inaugural Digital Data in Biodiversity Research Conference

BHL
Biodiversity Heritage Library

*Share your thoughts on social media using* #BHLib #BHLNDSR

# Inspiring Discovery through Free Access to Biodiversity Knowledge

**10 years** of inspiring discovery

*through* free & open access

to biodiversity literature & archives

*from the* 15th-21st centuries

## Mission

The Biodiversity Heritage Library improves research methodology by collaboratively making biodiversity literature openly available to the world as part of a global biodiversity community.

BHL
Biodiversity Heritage Library

# NDSR

NDSR's mission is "to build a dedicated community of professionals who will advance our nation's capabilities in managing, preserving, and making accessible the digital record of human achievement."

# Content Analysis

## What's in BHL?

119, 682 titles
201,703 volumes
51,917,032 pages
3,732,986 species


*as of 5/12/17

## What's not in BHL?

Everything else



**BHL**
Biodiversity Heritage Library
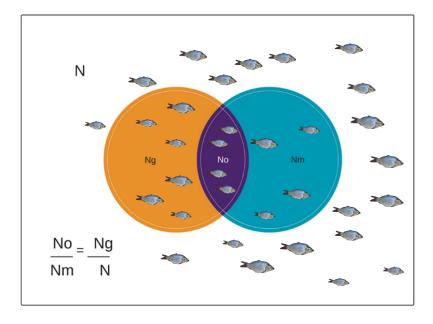
# Capture-Mark-Recapture (CMR)

N= unknown population
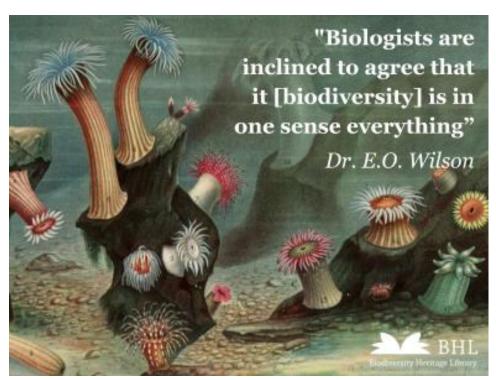Ng= first capture
Nm= second capture
No= recaptured

# CMR in Lit Reviews

- Kastner, M., S. E. Straus, K. A. McKibbon, and C. H. Goldsmith (2009). "The Capture-mark-recapture Technique Can Be Used as a Stopping Rule When Searching in Systematic Reviews." Journal of Clinical Epidemiology.

- Lane, D., Dykeman, J., Ferri, M., Goldsmith, C., Stelfox, H. (2013). "Capture-mark-recapture as a tool for estimating the number of articles available for systematic reviews in critical care medicine." Journal of Critical Care.

- Khabsa M, Giles CL (2014) "The Number of Scholarly Documents on the Public Web." PLOS ONE 9(5): e93949. https://doi.org/10.1371/journal.pone.0093949

- Ariño, AH.(2010) "Approaches to estimating the universe of natural history collections data." Biodiversity Informatics, 7: 81–92.

Biodiversity Heritage Library

# CMR for BHL

1. Define area of study
2. Determine how to search for organism
3. Determine which statistical model to use



"Biologists are inclined to agree that it [biodiversity] is in one sense everything"
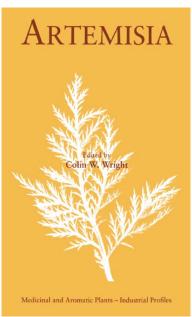
Dr. E.O. Wilson

# CMR for BHL

1. Define area of study
2. Determine how to search for organism
3. Determine which statistical model to use



*Artemisia frigida* © 1998 Gary A. Monroe
from Encyclopedia of Life

*Artemisia vulgaris* © 2008 Zoya Akulova
from Encyclopedia of Life

# CMR for BHL

1. Define area of study
2. Determine how to search for organism
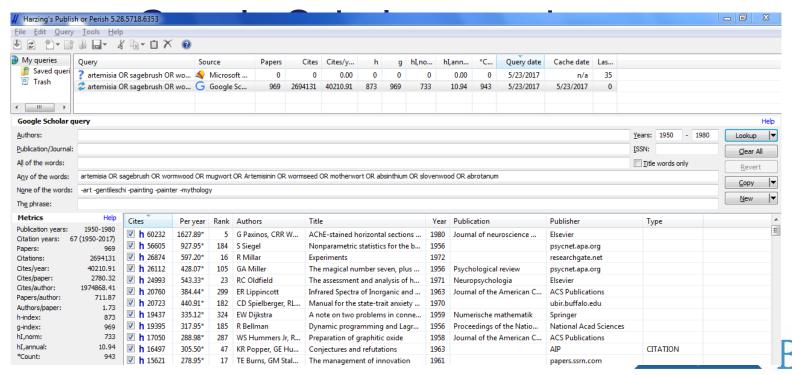3. Determine which statistical model to use

WorldCat results for: **kw: artemisia or kw: artemisinin or kw: sagebrush or kw: mugwort or kw: wormwood or kw: southernwood or kw: wormseed or kw: motherwort or kw: absinthium or (kw: qinghao and kw: su) or kw: slovenwood or kw: abrotanum not kw: art not kw: gentileschi not kw: painting and yr: 1950-1980 and (dt= "bks" or dt= "ser" or dt= "com" or dt= "art" or dt= "url").** (Save Search)

Records found: **1,942** (English: **1,565**) Rank by: **Number of Libraries**

WorldCat

# CMR for BHL

1. Define area of study
2. Determine how to search for organism
3. Determine which statistical model to use

# CMR for BHL

1. Define area of study
2. Determine how to search for organism
3. Determine which statistical model to use

Biodiversity Heritage Library uses taxonomic intelligence tools, including Global Names Recognition and Discovery (GNRD) developed by Global Names Architecture, to locate, verify, and record scientific names located within the text of each digitized page. Note: The text used for this identification is uncorrected OCR, so may not include all results expected or visible in the page.

Artemisia (38189)

Artemisia × decipiens (Vaccari) Giacom. & Pignatti (4)

Artemisia abaensis (2)

Artemisia abaensis Y.R. Ling & S.Y. Zhao (2)

Artemisia abaensis Y.R.Ling & S.Y.Zhao (20)

Artemisia abrotanifolia Salisb. (1)

Artemisia abrotanifolium (1)

Artemisia abrotanoides (22)

Artemisia abrotanoides Nutt. (2)

Artemisia abrotanum (601)

Artemisia abrotanum L. (29)

BHL
Biodiversity Heritage Library

# CMR for BHL

1.Define area of study
2.Determine how to search for organism
3.Determine which statistical model to use

Compare bibliographies and apply:
- Lincoln Petersen
- Seber probabilistic
- Poisson regression

# Thank You!

*Questions?*

**Alicia Esquivel**

6/6/17 | Inaugural Digital Data in
Biodiversity Research Conference

## Stay Connected with BHL!

Follow @BioDivLibrary on social media



Join our Mailing List: library.si.edu/bhl-newsletter-signup



BHL
Biodiversity Heritage Library