# Citation Networks and Impact Analysis
## *National Library of Medicine*
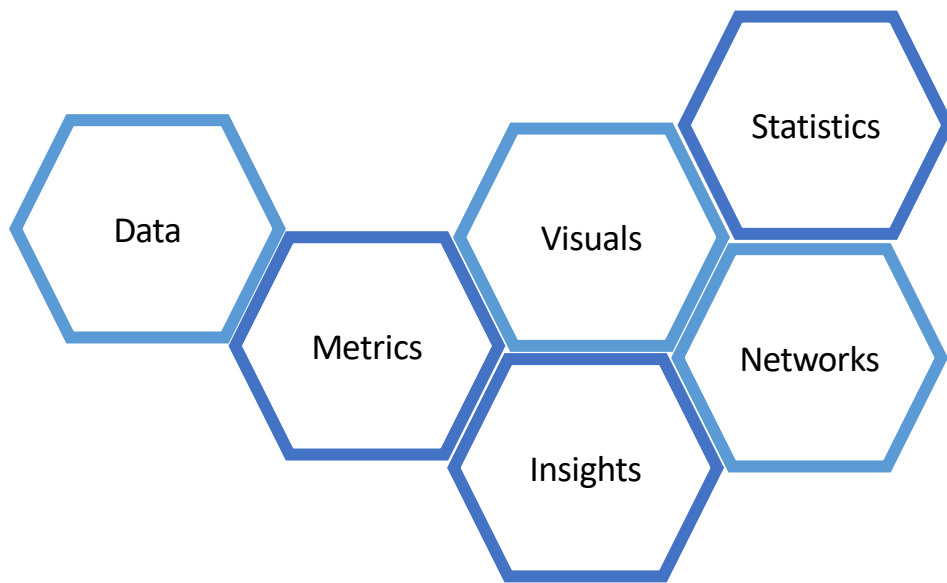
Isaac Robinson (CDF)
*Harvard University, Mathematics, 2022*

Special Thanks:
Anna Calcagno (NLM)
Kenneth Wilkins (NIDDK)

NIH National Institutes of Health

# PROJECT OBJECTIVES

1. Aggregate public and private data sources to provide wholistic and selectively accessible information
2. Create multiplex citation graph pipeline for future research
3. Suggest new impact metrics for funding institutions with non-traditional risk profiles
4. Package the above into easily digestible formats for technical and non-technical users

## Data Science
### What are our priorities?

### Understanding

- Statistical analysis
- Network analysis
- Metric
- Machine learning
- Important for drawing insights and actionable results
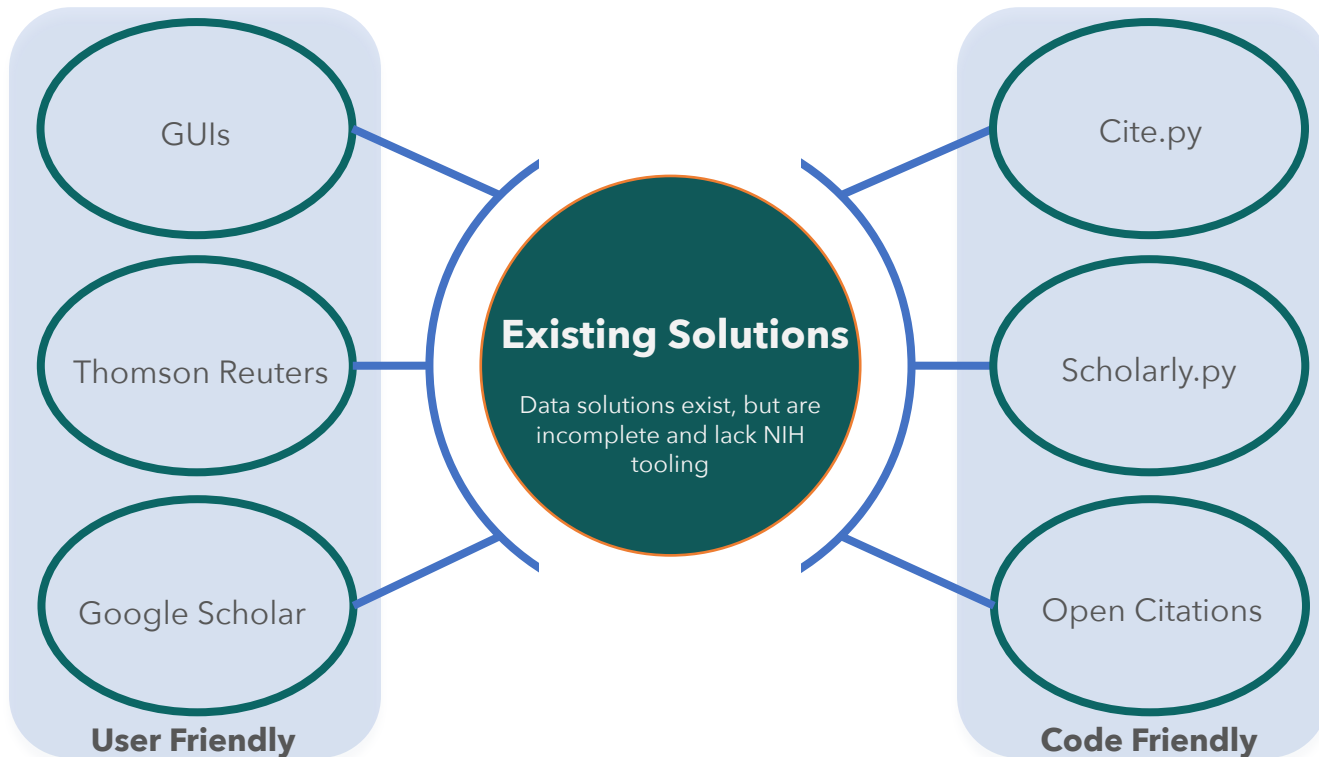
### Demonstrating

- Visualizations
- Tooling
- Interactivity and usability
- Important for making actionable results actually happen

### Key Takeaway

Data analysis can be a bit obtuse to those outside of a project, but those people are essential for generating insights.

NIH National Institutes of Health

# Step 1: Getting the Data

# EXISTING OPTIONS



**User Friendly**
- GUIs
- Thomson Reuters
- Google Scholar

**Existing Solutions**

Data solutions exist, but are incomplete and lack NIH tooling

**Code Friendly**
- Cite.py
- Scholarly.py
- Open Citations

What to do :/

# Introducing iCiteNLM

1. A library for curating all the data you need to build citation graphs

2. Accesses multiple pipelines including Google Scholar, Open Citations, and more.

3. Provides general tooling as well as NIH-specific functions, such as PMID to DOI conversions

4. Lots of room to expand and include NIH-specific resources like InCite

# Introducing iCiteNLM

**Curate**

A library for curating all the data you need to build citation graphs

**Private-Public**

Accesses multiple open access pipelines including Google Scholar and Open Citations, as well as support for private options
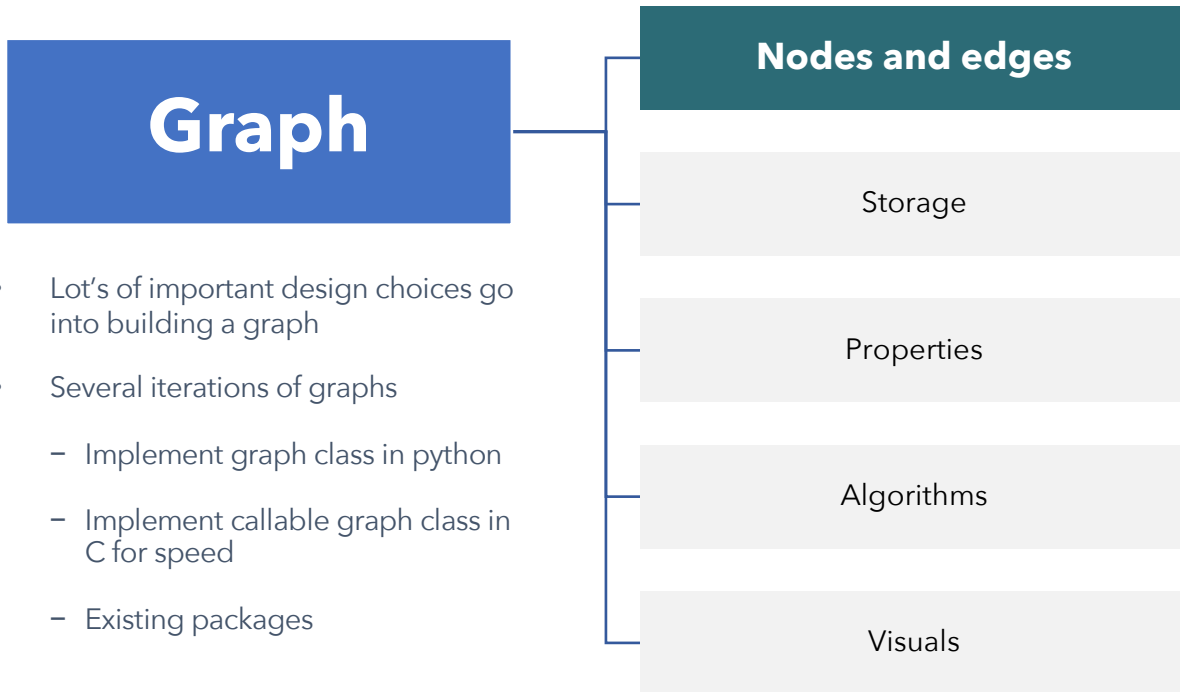
**By NIH, for NIH**

Provides general tooling as well as NIH-specific functions, such as PMID to DOI conversions

Now add it to the graph!

**Step 2: Create the Graph**

# Graph Basics

**Graph**

- Lot's of important design choices go into building a graph

- Several iterations of graphs

  - Implement graph class in python

  - Implement callable graph class in C for speed

  - Existing packages
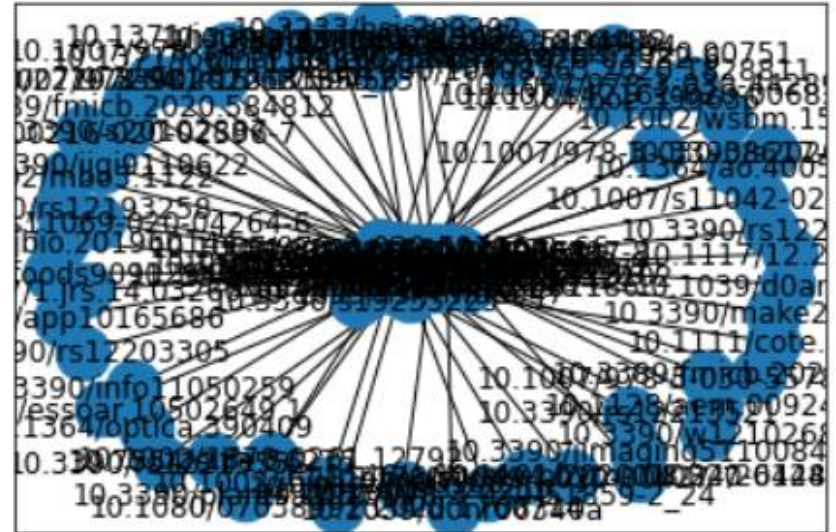
**Nodes and edges**

Storage

Properties

Algorithms

Visuals
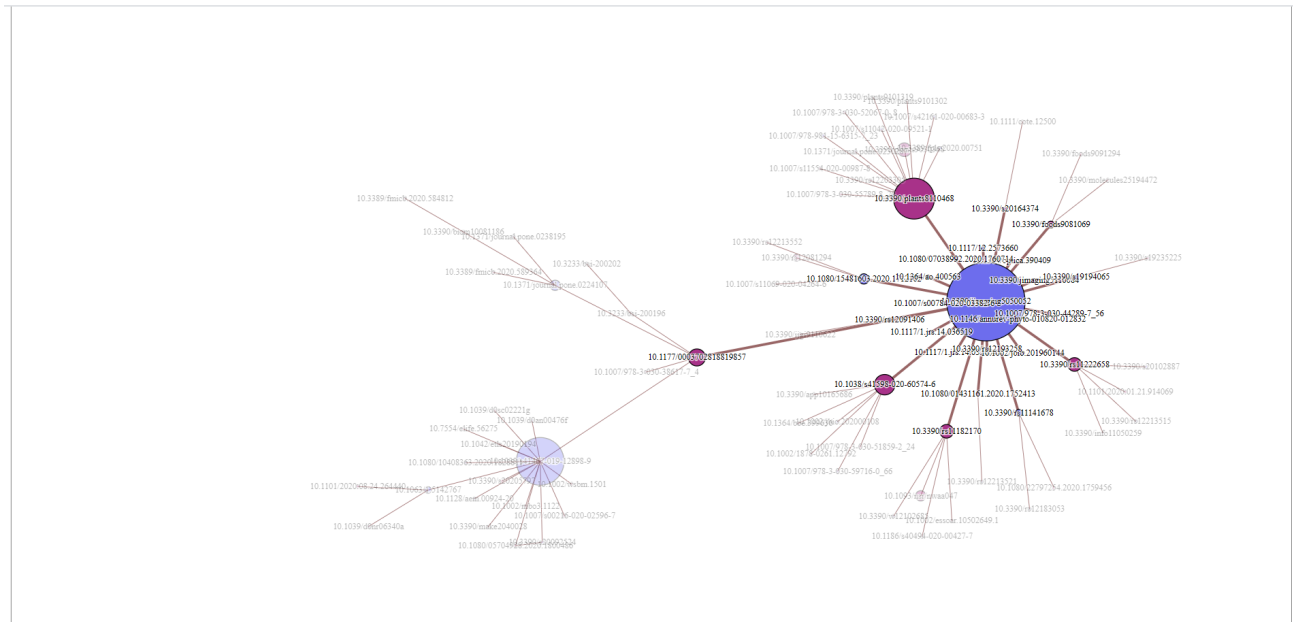
# Visualizing the Graph

# Pre-Built Solutions

- Most graphs used for **data science** are only displayed for exploratory purposes

- The accepted solution is matplotlib

- Hard to understand, only manipulatable using code, limited controllability

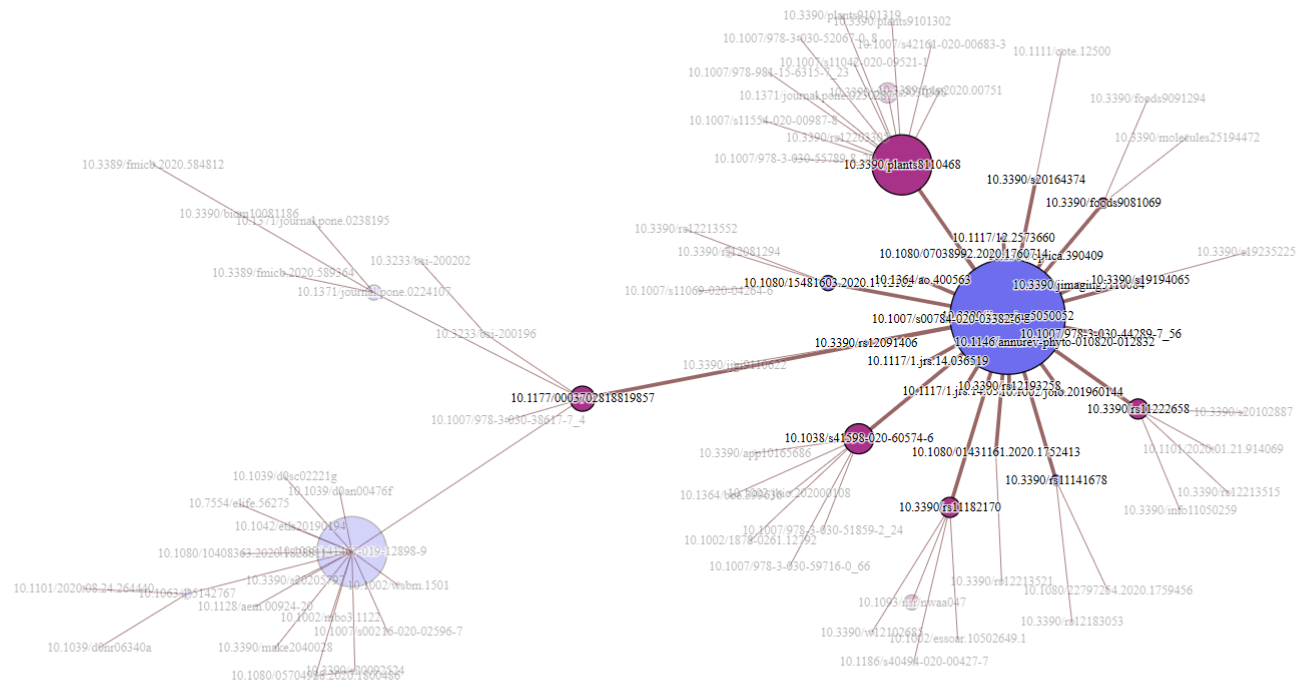- Poor node positioning

# My Solution: Web-Server, HMTL, D3

- Looked for examples of the prettiest visuals and then found what they used
  - NYT + D3
  - Fell in love with the library

- A DOM-manipulator, making it easy to launch (as easy as launching a website)

- Any desired features can be added, or embed as part of a dash



Citations: 23, Journal: J. Imaging

10.3390/plants9101319
10.3390/plants9101302
10.1007/978-3-030-52067-0_4
10.1007/s42161-020-00683-3
10.1111/cote.12500
10.1007/s11104-020-09521-1
10.1007/978-981-15-6315-7_23
10.1371/journal.pone.0233233
10.3389/fpls.2020.00751
10.3390/foods9091294
10.1007/s11554-020-00987-3
10.3390/rs12263887
10.3390/molecules25194472
10.1007/978-3-030-55789...
10.3390/plants8110468
10.3390/s20164374
10.3389/fmicb.2020.584812
10.3390/s20081069
10.1117/12.2573660
10.3390/rs12213552
10.1080/07038992.2020.1760714
10.3390/plants8101319
10.3390/biem10081186
10.1371/journal.pone.0238195
10.3390/rs12981294
10.1080/15481603.2020.1...
10.1364/ao.400563
10.3390/jimaging6100084
10.3390/s19194065
10.3233/tsi-200202
10.3390/s19235225
10.3389/fmicb.2020.589564
10.1371/journal.pone.0224107
10.1007/s11069-020-04284-6
10.1007/s00784-020-03582-6
10.3390/s20051052
10.1007/978-3-030-44289-7_56
10.3233/tsi-200196
10.3390/rs12091406
10.1146/annurev-phyto-010820-012832
10.3390/rs9110622
10.1117/1.jrs.14.036519
10.1177/0003702818819857
10.1007/978-3-030-38617-7_4
10.1117/1.jrs.13.032...
10.1002/joc.201960144
10.3390/rs12193258
10.3390/rs11222658
10.3390/s20102887
10.1038/s41598-020-60574-6
10.3390/app10165686
10.1080/01431161.2020.1752413
10.1101/2020.01.21.914069
10.3390/rs11141678
10.3390/rs12213515
10.1039/d0sc02221g
10.1039/d0an00476f
10.1364/boe.399768
10.1002/2020.00108
10.3390/info11050259
10.7554/elife.56275
10.3390/rs11182170
10.1042/etls20190194
10.3390/s20092898-9
10.1007/978-3-030-51859-2_24
10.1080/10408363.2020.1...
10.1002/wsbm.1501
10.1007/978-3-030-59716-0_66
10.3390/rs12213521
10.1101/2020.08.24.264440
10.1040635142767
10.1128/aem.00924-20
10.3390/s20153734
10.1080/22797254.2020.1759456
10.1039/d0nr06340a
10.1002/mbo3.1125
10.3390/rs12183053
10.3390/maize2040028
10.1007/s00286-020-02596-7
10.1093/rb/rbwaa047
10.1002/essoar.10502649.1
10.1080/05704928.2020.1800486
10.1186/s40494-020-00427-7

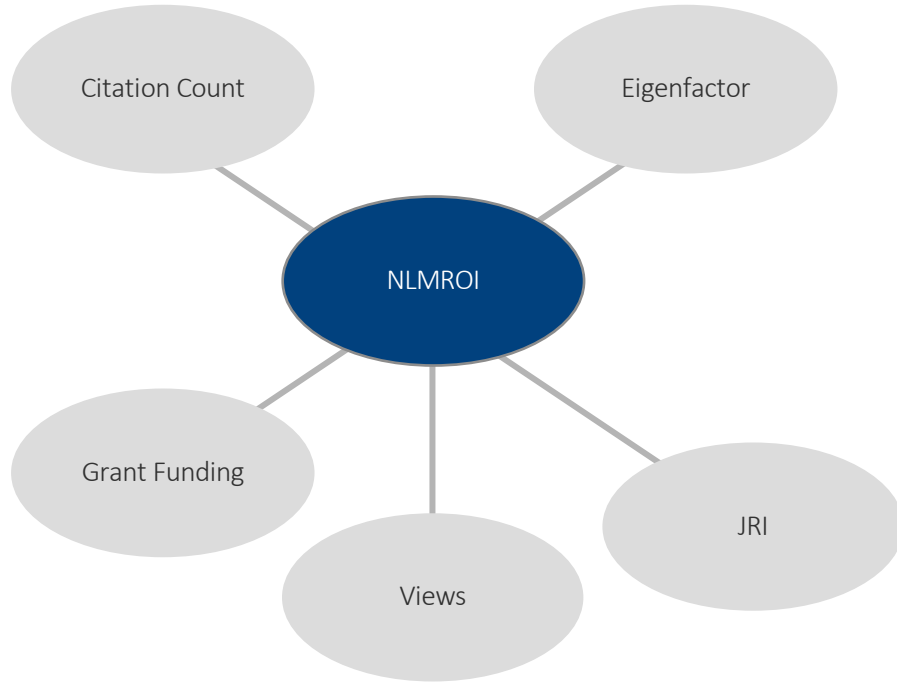Citations: 23, Journal: J. Imaging

**Step 3: Metrics**

# Introducing NLMROI

1. Citation count

2. Eigenfactor

3. JRI

4. Views and Web2.0 Metrics

5. Grant funding

The risk profile for governmental organizations necessitates a different metric

# Existing Metrics

1. Citation count
2. Eigenfactor
3. JRI
4. Views and Web2.0 Metrics
5. Grant funding

# Introducing NLMROI

$$NLMROI = \frac{\sum C_i E_{c_i} + \sum \sqrt{C_j} E_{cj} + \sum \sqrt[3]{C_k} E_{c_k}}{\sum Funding}$$

- Government-specific metric
- Consistent with existing literature and metrics including Eigenfactor, JRI, Impact Factor etc.
- Helps account for differences in research areas (still best to compare within one area)

Now add it to the graph!

**Step 4: Usability**

# Outputs

1. Graph written in D3 and hosted locally means it is easily web-publishable
   1. Open-source/publicly available options of iCiteNLM makes data privacy less of an issue and means anyone can access and see this data
2. For those interested, a host of algorithms were implemented on the graph including clustering, k-cliques, shortest-path, centrality, and more
   1. These can all be accessed through the python module
3. Excel file has been formatted for easy analysis

# Next Steps

# Only the beginning

1.  Any metric you can imagine can be added to the graph

    1.  Parallel graphs for co-authorship

    2.  Weighing for small vs large teams

2.  Include funding information for 2nd, 3rd degree nodes

    1.  Validate metric using studies and experts

3.  Host graph on data servers and expand to higher degrees

4.  Any other graph analysis technique you would like!

    1.  Causal inference

    2.  Dominating sets

    3.  Network flow

# Other Tail Ends

1. Topic modeling with Latent Dirichlet Allocation
2.  NER with Bert-based machine learning models
3. Kleinberg burst-detection
4. Excel files

CIVIC
DIGITAL
FELLOWSHIP

Thank you to everyone for this opportunity!

Please reach out if you have questions/comments, or if you think any of these techniques could be useful in other areas!

Email: isaac_robinson@college.harvard.edu

# Sources

**Leveraging Citation Networks to Visualize Scholarly Influence Over Time**
J. Portenoy, J. Hullman, J.D. West (2016)
arXiv: 1611.07135

**Static ranking of scholarly papers using article-level Eigenfactor (ALEF)**
I. Wesley-Smith, C. T. Bergstrom, and J. D. West (2016)
The 9th ACM International Conference on Web Search and Data Mining (WSDM)

**CiteSpace II: Detecting and visualizing emerging trends and transient patterns in scientific literature**
C Chen
Journal of the American Society for information Science …, 2006

**Author-Level Eigenfactor Metrics: Evaluating the Influence of Authors Institutions and Countries Within the SSRN Community.**
J.D. West, M.C. Jensen, R.J. Dandrea, G.J. Gordon, and C.T. Bergstrom, (2013)
Journal of the American Society of Information Science and Technology 64: 787-801

**Journal impact evaluation: a webometric perspective**
Thelwall, M.
Scientometrics 92, 429–441 (2012)

# Sources and Papers

**Eigenfactor: ranking and mapping scientific knowledge**
J.D. West, (2010)
PhD Dissertation. University of Washington
**Big Macs and Eigenfactor Scores: Don't Let Correlation Coefficients Fool You**
J.D. West, T.C. Bergstrom, C.T. Bergstrom, (2010)
Journal of the American Society for Information Science & Technology. 61(9): 1800-1807

**How to improve the use of metrics: Learn from Game Theory**
J.D. West, (2010)
Nature 465:871-872 [pdf]
**Network Analysis and Indicators**
Ding Y., Rousseau R., Wolfram D.

**Causal Inference for Social Network Data**
Elizabeth L. Ogburn Department of Biostatistics

**Searching for intellectual turning points: Progressive knowledge domain visualization**
C Chen
Proceedings of the National Academy of Sciences, 2004

# Sources and Papers

**Impact factor: a valid measure of journal quality?**
Saha, Somnath et al.
Journal of the Medical Library Association : JMLA vol. 91,1 (2003): 42-6.

**Response to "Big Macs and Eigenfactor Scores: The Correlation Conundrum"**
J.D. West, T.C. Bergstrom, C.T. Bergstrom, (2010)
Journal of the American Society for Information Science & Technology 61:2592

**Dominating Sets in Social Network Graphs**
Laura L. Kelleher, Margaret B. Cozzens,
Mathematical Social Sciences,
Volume 16, Issue 3, 1988