

Publication Mining using Social Network Analysis and Probabilistic Topic Modeling

-: DATA: -CHALLENGE



Supriya Chinthavali¹, Dakotah Maguire², Ethan Hicks³, Philip Hicks³ Oak Ridge National Laboratory¹, ORAU², HERE³

ABSTRACT

The amount of scientific information available on-line has increased at an exceptional rate, with recent estimates reporting that a new paper is getting published every 20 seconds.

The 4 main challenges in the area of publication mining are

- 1.Identify the individual or group of individuals who appear to be the expert in a particular field or subfield.
- 2.Identify topics that have been researched across all publications
- 3. Visualize the geographic distribution of the topics in the publications
- 4. Identify how topics have shifted over time

The data mimers team (ORNL+UTK) provides solutions to the four challenges mentioned above using machine learning and graph modelling techniques called and social network analysis probabilistic topic models to mine the scientific publications and identify key characteristics and patterns that can be used by human researchers to develop useful knowledge and further enhance scientific discovery.

Our approach involves: 1. Constructing various collaboration/citation networks and enable analysis of graph-based collaboration networks that can identify experts within a given area. using computing community detection algorithms to answer analytic questions posed by this data challenge.

2.Developing a visual-analytic tool that enables a user to understand the key topics that have been researched over the previous 50 years using an unsupervised bayesian probabilistic model known as Dynamic Topic Model based on LDA(Linear Dirichlet Allocation).

POINTS OF CONTACTS

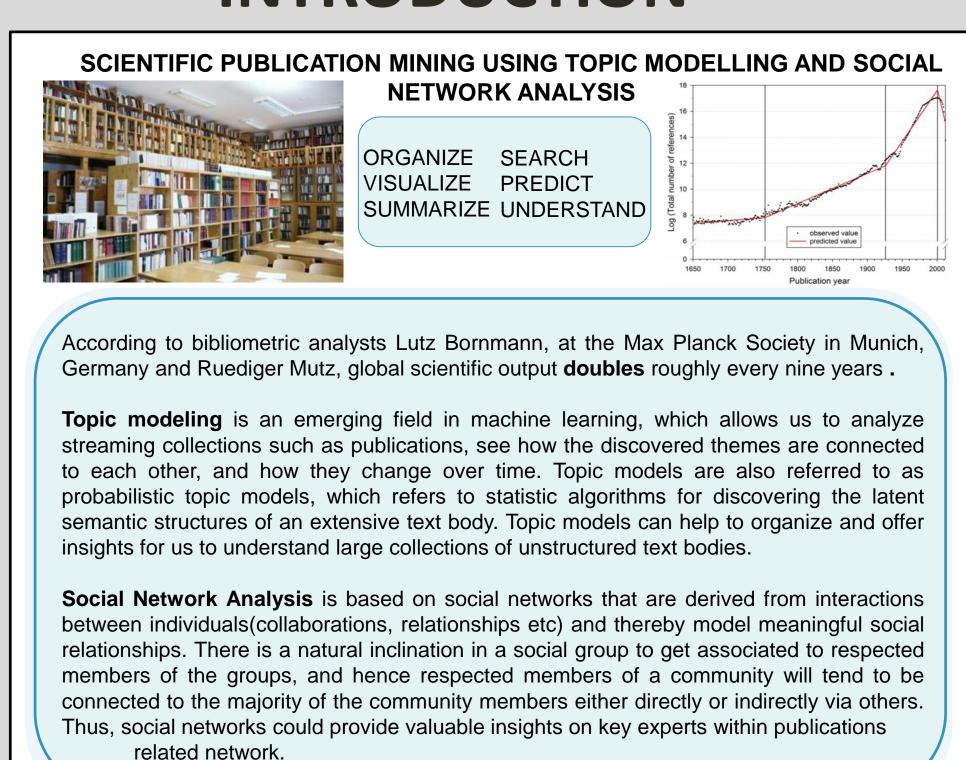
Supriya Chinthavali

Data Visualization and Machine Learning Team, Computational Data Analytics Group, Computational Sciences and Mathematics Division, Oak Ridge National Laboratory

(865) 574-7592

chinthavalis@ornl.gov

INTRODUCTION



TOPIC MODELLING

CHALLENGES AND OUR APPROACH



Our approach involves:

1.Construction of several collaboration

data and applying community detection

algorithms to answer analytic questions

2. Developing a **visual-analytic** tool that

50 years based on the underlying theme

structure and also provide graph-based

collaboration networks that can identify

enables a user to organize the documents

that have been researched over the previous

posed by this data challenge.

experts within a given area.

based networks using citation and authors

Collaboration Networks (direct citation, co-citation, co-author networks)

Dynamic Topic Modelling (LDA based

Social Network Analysis

model)

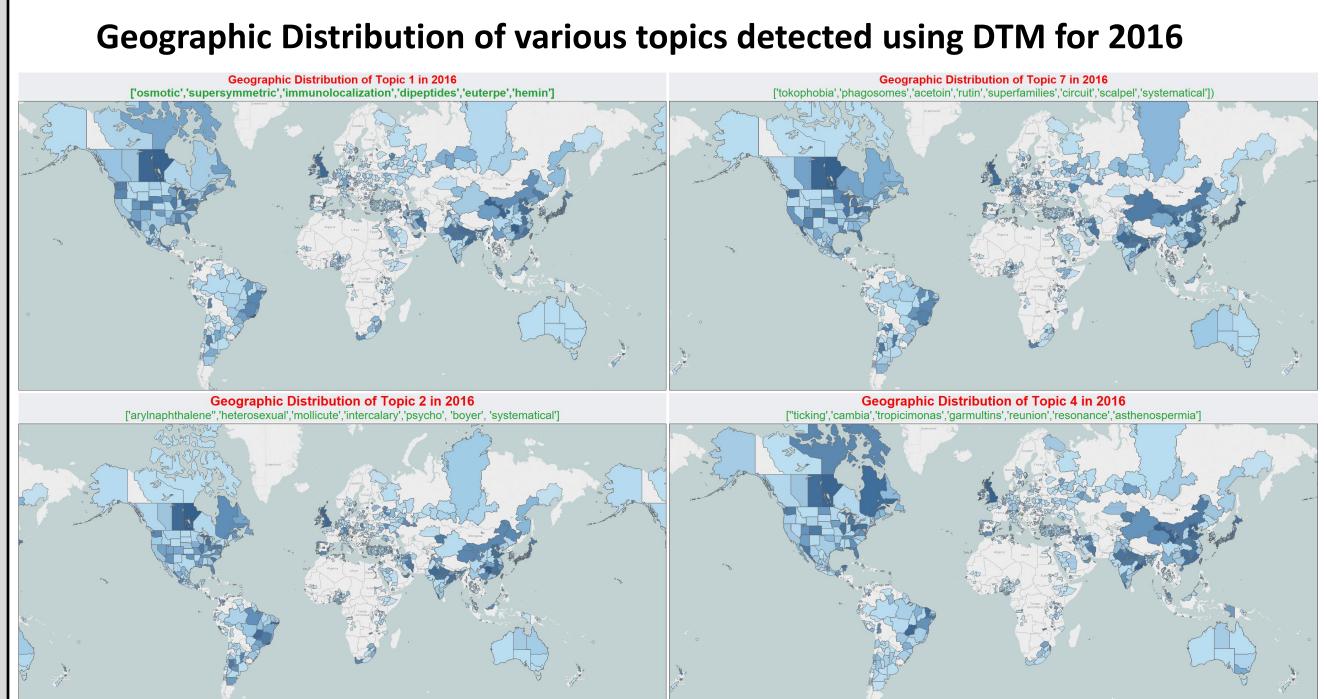
Network Visualization using Cystoscope and Information vis using Tableau

Visualization of number of publications for 2016

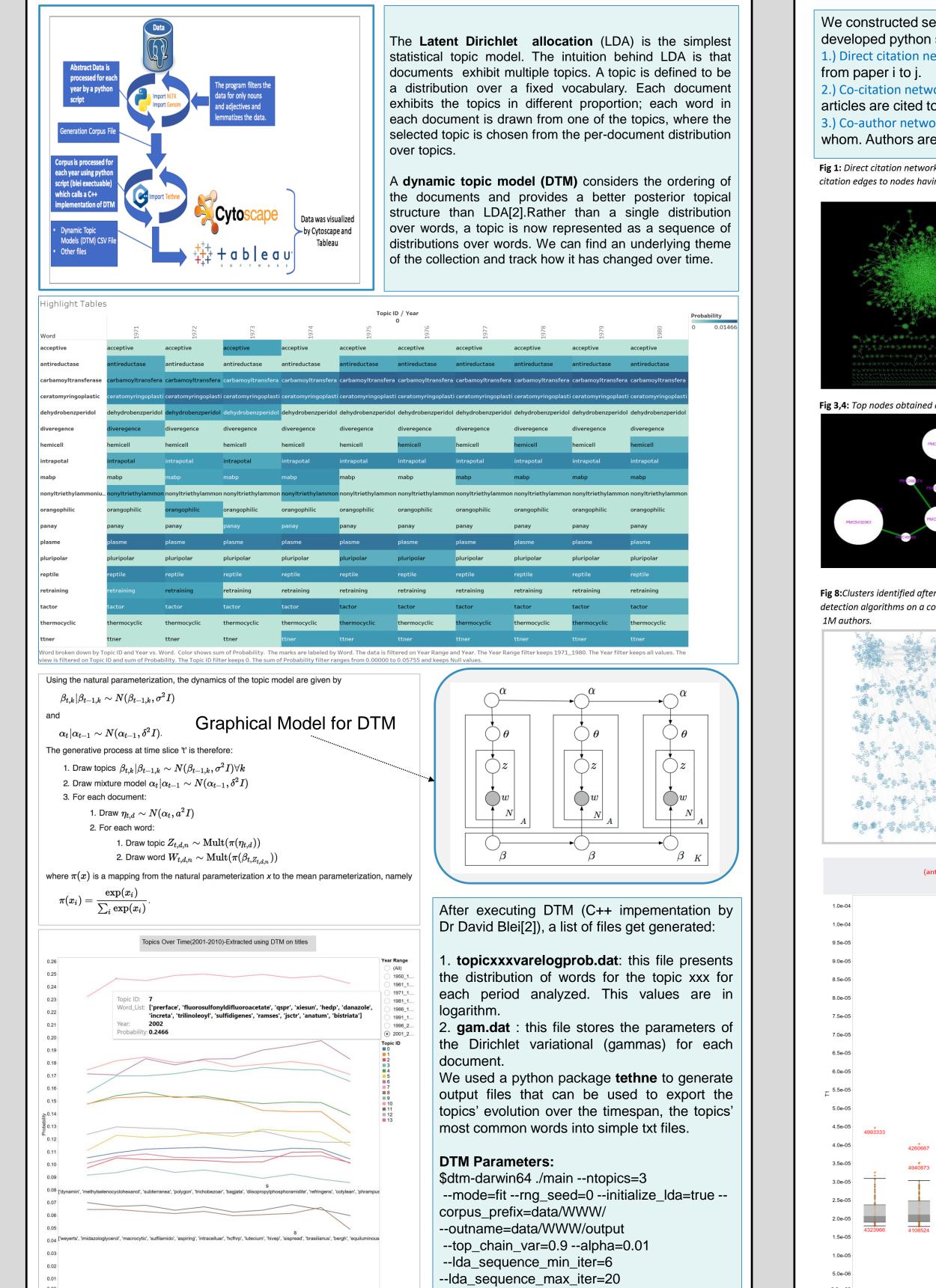


Using R, we take the output files of DTM (gam.dat) and combine the probabilities of words in each topic to get topic distribution of the documents. We then join this with the affiliation data to plot the geographic distribution of topics on a world map.

Note: Due to missing or uninterpretable data (e.g. various admin levels), 7.5% of the publications for 2016 were not aeocoded.

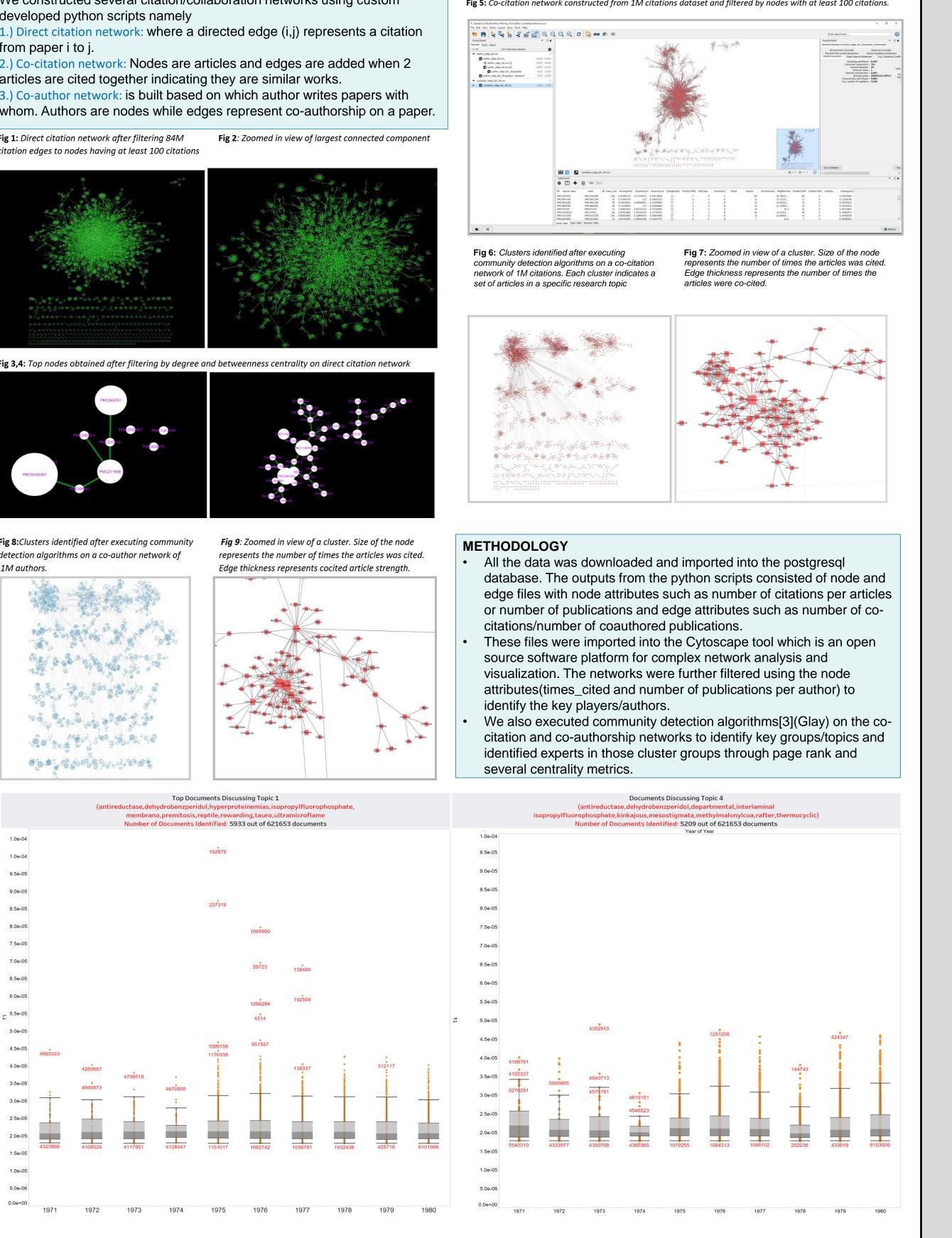


VISUALISATION: GEOGRAPHIC DISTRIBUTION



--lda_max_em_iter=20

EXPERT IDENTIFICATION USING SNA

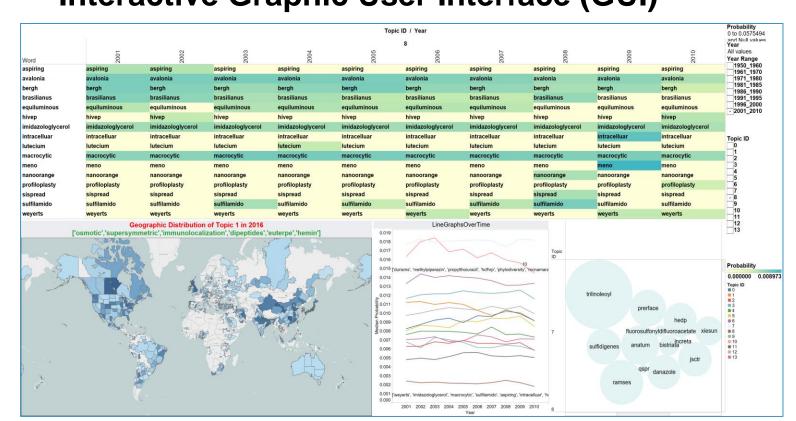


REFERENCES & FUTURE WORK

CONCLUSION

We have developed an interactive **visual-analysis** tool that enables a general user/ subject matter expert to understand the key topics that have been researched over the previous 50 years and organize the documents accordingly using an unsupervised bayesian probabilistic model known as Dynamic Topic Model based on LDA (Linear Dirichlet Allocation).

Interactive Graphic User Interface (GUI)



FUTURE WORK

- Experimentation with parallel versions of DTM for faster training of the model using abstracts as opposed to titles(for years > 2010)
- Using a memory fat machine to compute the collaboration networks for the entire datasets instead of subsets
- Implement algorithms to help clean-up affiliation address data that will allow better visualization of geographic distributions
- Include context information within the topic models such authorship information or geographical coordinates associated with documents.

REFERENCES

1.Blei, David M. "Probabilistic topic models." Communications of the ACM 55.4 (2012): 77-84. 2.Blei, David M., and John D. Lafferty. "Dynamic topic models." Proceedings of the 23rd international conference on Machine learning. ACM, 2006.

3.Gang Su, Allan Kuchinsky, John H. Morris, David J. States, Fan Meng; GLay: community structure analysis of biological networks. Bioinformatics 2010; 26 (24): 3135-3137.

doi:10.1093/bioinformatics/btg596

4.http://socialcomputing.ing.puc.cl/uploads/DynamicTopicModellingTutorial.pdf 5.https://www.cognizant.com/whitepapers/identifying-key-opinion-leaders-using-social-network-an alysis-codex1234.pdf 6.http://nealcaren.web.unc.edu/a-sociology-citation-network/

7.https://github.com/magsilva/dtm

8.http://web.mit.edu/seyda/www/Papers/ECIR_extended.pdf

9.Newman, MarkEJ. "Coauthorship networks and patterns of scientific collaboration." Proceedings of the national academy of

sciences101.suppl 1 (2004): 5200-5205