

Paper Title

Aditya Garg
New York University
NY, USA
aditya.garg@nyu.edu

Aashutosh Trivedi
New York University
NY, USA
aashutosh.trivedi@nyu.edu

Ansa Mary Ephraim
New York University
NY, USA
ame378@nyu.edu

Abstract—

Considering the huge number of research papers published and the even bigger number of citations in the same, we thought about analyzing this data. There are multiple dimensions in which research is being carried out in the field of computer science and this analysis could help come up with some interesting predictions like research peers and critical conferences. While our experimental data is from field of computer science, we believe the same analysis can be used over any other field also.

Keywords— *citation, co-authorship, analytics*

I. INTRODUCTION

We access citation database of conferences over a period of time. Some analytics we hope to derive are

- 1) Potential Hot Topic : Paper with the least citation distance from all the published papers.
- 2) Find citation distance and collaboration distance between authors
- 3) Help suggest should two authors work together in future.
- 4) Suggest/predict conferences to submit entries to for authors

II. MOTIVATION

(Write a couple of paragraphs describing why you think this analytic is important. Why should people care about this analytic?)

III. RELATED WORK

Analysis of SIGMOD's Co-Authorship Graph :

The first input gathered from this paper was a valuable data source : a DBLP XML data file that contains authorship info and panel members. A few points that were analyzed were : number of authors per paper, number of SIGMOD papers per author, the longest path distance of a collaboration graph. They found that in no year was the co-authorship graph a single connected component. For the single largest connected component every year, they also analyzed the size, the clustering

coefficient, and tried to see if it exhibited the small world phenomenon. The paper also explored the existence of hubs in the graphs, and measured their centrality.

Some Analyses of Erdos~ Collaboration Graph :

This was not exactly a well published paper, but it was hard to find popular publications that relate to our project as closely. It sets out some features of analysis that can be performed on large network graphs, and specifically carries them out on Erdos graphs.

It analyzed features regarding the co-authorship , ie the mean, median, average degree, maximum and maximizer of vertex degrees in Erdos graphs. It analyzed the top authors based on number of co-authors. It also described the concept of *cores* and went on to analyze authors and number of co-authors in the main core, total number of co-authors, average core and average degree of all their co-authors, and their collaborativeness. Another interesting idea they proposed was that of 'Lords' - vertices that have strong influence on their neighbourhoods. A recursive process, much like page ranking, finally assigns the 'power' to each vertex. Another analytic used was block modelling.

The Chasms of CSCW: A Citation Graph Analysis of the CSCW Conference :

The authors look at papers from the CSCW conference and tries to prove the following hypothesis.

- H1: There is a strongly connected core of the CSCW conference : Proven to some extent
- H2a: The CSCW conference is divided into several thematic clusters : Clearly evident from results of clustering
- H2b: Social science and computer science papers will reside in different clusters : Again evident from clustering
- H3: There are chasm-papers in the CSCW conference that are cited outside the conference significantly more than within it : Proven to some extent.

The paper divides the corpus into clusters using the betweenness centrality algorithm implemented in JUNG framework, which iteratively removes edges from citation

graph. It also defines a success function and chasm potential to find chasm papers.

Analysis of Papers from Twenty-Five Years of SIGIR Conferences: What Have We Been Doing for the Last Quarter of a Century?:

The authors try to look at all the papers published in last 25 years in SIGIR conferences and perform a content analysis on the same. They try to determine how the trending topics have changed over time, which topics have come and gone. They have categories like databases, evaluation, Probabilistic and language models, Conceptual IR, users and search and general among others. They also created co-authorship graph (including cleaning up author names) and perform some analysis on that. These include authors with maximum papers, authors with greatest number of collaborators, and analyze the path between the authors, Erdos-type analysis. Authors then try to predict the hottest topic for next year and the co-authorship combination.

The structure of scientific collaboration networks: M. E. J. Newman:

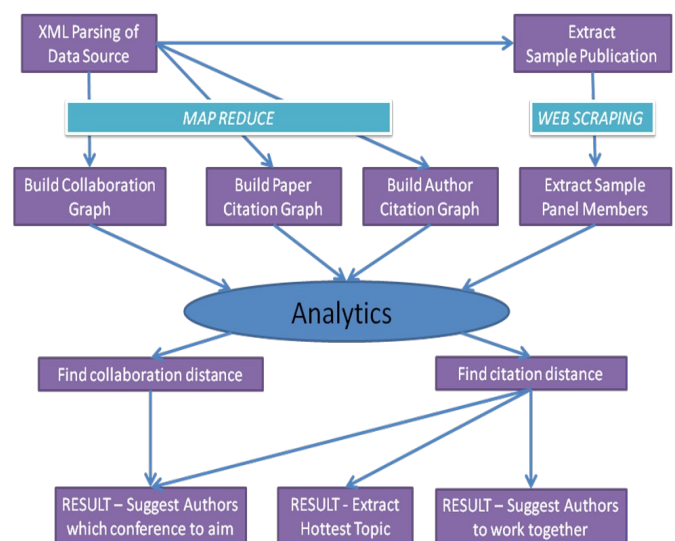
The author considers 2 scientists to be connected to each other if they have co-authored a paper together. The paper argues that most people who have written a paper together will know one another quite well. Thus it is a moderately stringent definition. The author has constructed collaboration graphs for scientists in a variety of fields. The data come from four databases: MEDLINE (which covers published papers on biomedical research), the Los Alamos e-Print Archive (preprints primarily in theoretical physics), SPIRES (published papers and preprints in high-energy physics), and NCSTRL (preprints in computer science). In each case, he has examined papers that appeared in a 5-year window, from 1995 to 1999 inclusive. The sizes of the databases range from 2 million papers for MEDLINE to 13,000 for NCSTRL. The distance is calculated based on the co-authoring of the paper. 2 scientists are said to be at distance 1 if they have co-authored a paper together. Any authors who have worked with these co-authors come at a distance 2 from the original author and so on. The paper takes a number of different scenarios such as the number of authors (and the errors such as same author giving different initials or different authors having same initials), Mean Papers per Author and Authors per Paper and number of collaborators. The paper also suggests that with a very large database and huge number of connections, the first degree connection takes up around 3/4 part of the database. Also the second connection group is much smaller than the first degree connection group. The author concludes that the maximum distance between 2 scientists is an average around 6 irrespective of the size of the database.

Understanding Importance of Collaborations in Co-authorship Networks: A Supportiveness Analysis Approach

This paper supports the theory that 2 authors co-authoring a paper means that one author supports the research of the other author. Thus it provides various such supportiveness measures. The definition of supportiveness here is given as "For an author a, the supportiveness from author b to a is used to measure how close the collaborations from b to a." The authors develop efficient methods to extract top n most supportive authors in co-authorship networks. They model the supportiveness ranking problem as a reverse k nearest neighbor (k-RNN for short) searching problem on graphs. To better model the co-authorship relation, we use hyper graphs in this paper. The authors use harmonic distance measure to calculate the support. It shows that support can differ in different directions based on the number of papers written in total and the number of papers co-authored. Thus the closeness and distance is calculated to be the harmonic mean of the contribution from one author to the other. The k nearest neighbor algorithm is used to find the author with maximum support. The algorithm finds the neighbor having the least harmonic distance and thus concludes that it can be a neighbor. Here one author can have multiple or one nearest neighbor, depending on different factors. The authors also have tried to expand the algorithm to work for not just a single vertex but also a group of vertices. They have discussed the efficiency of both problems on large co-authoring networks and come up with very interesting solutions.

Ic. DESIGN

RTBD Project – Citation and Collaboration Network Analysis



§.

§I. RESULTS

(Future... In this section, you can describe: Your experimental setup/issues with data/performance/etc. Describe your experiments, describe what you learned. Did you prove or disprove your hypothesis? Were some results unexpected? Why?)

§II. FUTURE WORK

(Future... Given time, how would you expand your analytic? Could it be applied to other areas? Etc...)

§III. CONCLUSION

(Future... One or two paragraphs about the value/accuracy/goodness of your analytic.)

ACKNOWLEDGMENT

(This section is optional. It can be used to thank the people/companies/organizations who have made data available to you, for example. You can list any HPC

people who were particularly helpful, if you used the NYU HPC.)

REFERENCES

- [1] T. White. Hadoop: The Definitive Guide. O'Reilly Media Inc., Sebastopol, CA, May 2012.
- [2] J. Dean and S. Ghemawat. MapReduce: Simplified data processing on large clusters. In proceedings of 6th Symposium on Operating Systems Design and Implementation, 2004.
- [3] S. Ghemawat, H. Gobioff, S. T. Leung. The Google File System. In Proceedings of the nineteenth ACM Symposium on Operating Systems Principles – SOSP '03, 2003.
- [4] Mario A. Nascimento, Jorg Sander and Jeffrey Pound. Analysis of SIGMOD's Co-Authorship Graph
- [5] Vladimir Batagelj and Andrew Mrvar. Some Analyses of Erdos' Collaboration Graph.
- [6] Michal Jacovi, Vladimir Soroka, Gail Gilboa-Freedman, Sigalit Ur, Elad Shahar, Natalia Marmasse. The Chasms of CSCW : A Citation Graph Analysis of the CSCW Conference.
- [7] Yi Han, Bin Zhou, Jian Pei, Yan Jia. Understanding Importance of Collaborations in Co-authorship Networks: A Supportiveness Analysis Approach.
- [8] M.E.J. Newman. The structure of scientific collaboration networks.
- [9] Alan F. Smeaton, Gary Keogh, Cathal Gurrin, Kieran McDonald and Tom Sødring. Analysis of Papers from Twenty-Five Years of SIGIR Conferences: What Have We Been Doing for the Last Quarter of a Century ?