

ANALYSIS OF RESEARCH TOPICS AND COLLABORATIVE NETWORK OF THAI GOVERNMENT SUPPORTED SCHOLARS

SIRIKORN SANTIROJANAKUL

College of Arts, Media and Technology, Chiang Mai University,
239 Huaykaew Rd., Suthep, Muang, Chiang Mai, Thailand
E-mail: sirikorn.s@cmu.ac.th

Abstract

The aims of this study are to investigate the existing scholars' interest and the collaborative network who work in the government institution in north Thailand. Data sets were gathered from the government's system and Scopus. The benefits of this paper are to create an understanding of the research topic and the structure of the scholar community. This research used text mining techniques, word cloud and social network analysis to discover the published research in the works of Science, Technology, Engineering and Mathematics (STEM). Word cloud was used for listing the latent research keywords. The results show that the structure of network relationship of single researchers is 1-N node. STEM researcher has the average degree of 1.954.

Keywords: Research, Social network analysis, STEM, Text mining, Word cloud.

1. Introduction

1.1. Scholars research in STEM

Government aims to create the systems and mechanisms for developing the national human capital in various national projects. Government supports scholarship for educational institutions and individual students by funding them to study in Thailand and abroad. The agenda is to restore the country back at the level of effective performance that can cooperate in making high contribution with other foreign entities. The scholarly graduates consist of individual with high knowledge, skills, ability, and capability in development and being innovate. The human capital management process has intentions to fulfil the performance management role in creating human capital to achieve the country's competitive advantages. For example: associate with their commitment, talent management, develop capability, transformation, leadership, engagement, motivate, training, reward and retirement [1]. It is essential for government organizations to estimate the skills, talents, and expertise of scholars for human capital management across the industries and other sectors [2].

Many previous studies have focused on the effectiveness and relationship of research, innovation, human capital and social capital. The results identify the knowledge management, entrepreneur and culture as being the driver of innovation outcome. The managing innovation is a complex process and requires a deep understanding of input, process and outcome [3]. The Government Human Capital Knowledge Management of Republic Indonesia (NUSANTARA) use the knowledge management model to solve the bureaucratic reform of the human capital management issues and the results have a readiness index that is eligible [4]. In the case of the Thai government, it requires the existing human capital as national scholars to enhance the Science, Technology, Engineering, and Mathematics (STEM) research and help industries to create a more competitive advantage in a knowledge-based economy. The great development of human capital will reach the S-curve model and industries cluster: super cluster and other targeted clusters that are using advanced technology and future industry model.

The government scholarship student in this study is a person who received a scholarship from the Ministry of Science and Technology and the Royal Golden Jubilee Ph.D. Program (RGJ Ph.D.); until his graduation and work as a lecturer and researcher in a government organization. To raise the level in performance of the country, the RGJ Ph.D. identified the number of researchers as the following:

- Phase I (1998-2012): To produce 5,000 Ph.D. researchers and 5,000 publications
- Phase II (2008-2022): To produce 20,000 Ph.D. researchers and 20,000 publications

The objectives of RGJ Ph.D. for funding are 1) produce numerous research-oriented PhD graduates and research outputs with international standards, 2) promote collaborations between Thai and foreign researchers, 3) strengthen graduate education in the Thai university system, 4) save the budget for overseas study of PhD students [5]. Each institution has its own data base system for maintaining the scholars' existing information before and during their studying. On the other hand, they cannot update the scholars' data after he/she graduated and is officially working. The main reasons are that they work in different locations throughout Thailand and each department of the Ministry of Science and Technology and the RGJ Ph.D. have a great amount of responsibilities with the new large number of scholars in each fiscal

year. The objective of this paper is to investigate the pattern of relationship among the researchers in STEM. This study uses the text mining technique to extract and visualize the collaboration graph. In this paper the text mining techniques are used to study the research topic and the structure of the collaborative network. Topic Modelling are effectively used to interpret the information from the large documents and serves as an efficient method for revealing a dataset consisting of 3,627 research papers from the journals and transactions of Biomedical Engineering and presented in graphics [6]. In addition, text mining has been used in various ways to identify research trends through research paper analyses in the field of technology and exploring patent data [7]. Currently, text mining is a research method that had been widely used to discover valuable information latent in the texts. The next section reviews literature relating to text mining, word cloud and social network analysis on the analysis published papers and collaborative of STEM scholars. Section 3 explains the steps of text mining from data collection through data analysis. Section 4 presents the results and discussion that were discovered from the systematic analysis. The conclusion is given in Section 5.

2. Literature Review

2.1. Text mining

Text has always been the default way of storing information that have been around for hundreds of years in various systems and format. Today, the benefit of text mining come with the large amount of valuable information latent in texts that are not available in classical structured data formats [8]. The enormous text message in various system is difficult to understand and deal with because of its high volume and overlapping data. Some scholars' defined text mining techniques to extract useful information by processing a large number of unstructured text to reveal useful fragments, model, direction, future trend and rule [9].

Analyzing texts is a great process to create more important sentiment than extracting structured data because of the sheer volume of valuable information of almost any imaginable types contained in them. The visualization tools and analytics models are suitable for applying in information search, knowledge management, statistical analysis, cognitive science, case event correlation, information filtering, research, and decision science that can develop potential information of enormous database, data warehouse, and data lake [9, 10].

Techniques of text mining are based on the presence of word or phrase via statistical analysis. Here, the new mining technique proposed heuristic and evolutionary algorithms including a hierarchical clustering of different data that is categorical, numerical, and multidimensional [11]. Many researches used text mining techniques for analyzing the consumer sentiment on their products and services. On the other hand, some researches apply the tool on the academic domain to explore about the facts and opinions that can provide valuable information. The visualize tools and techniques can be a significant tool to create more understanding of the features and developments of the large dataset various system.

2.2. Text visualization as word cloud

Word clouds is an easy and quick visual description of words from any written material built on its frequency that show meaningful descriptions from the text

source [10]. Word clouds are increasingly used in the public and private sector as a tool for representation that is derived from the collection of written texts. There are many Word cloud approaches to creating informative or interesting. Word cloud was analyzed by applying qualitative content analysis that can create outstanding keywords emerging from data mining and text mining techniques [12]. Some study proposed Molecular clouds (MCs) evolving into galaxies whereas the names of several text visualization present a simple algorithm MCs like word clouds [13]. Moreover, some paper focus on tag clouds that are visualization displaying the content of document or database as a group of tag key words; indexed by frequency in the research results that are displayed in different typeface color and size [14]. Also, the multimedia database of broadcast company can combine image and textual information retrieved from the query results- FaceCloud that is linked to an individual video [15]. Currently, the English online word cloud generators are available for users with simple requests [12]. Word cloud application was created from various tools and techniques including, Python, R language, Standard Query Language (SQL) and word2vec. Table 1 summarizes the application of word cloud in the last year with different domain and data's source.

Table 1. The application of word cloud technique in various domain.

Domain	Methodology	Source of data	Result
Education: Papers Reviews, 2017 [8]	- Select the highest score article - Using R language and visualized word clouds	Comments from Altmertics.com	Present the research hotspot and associate methods
Education: eLearning overload manual evaluation, 2017 [10]	- Used superlative model and construct word clouds	Student's open-ended answers	Students and teachers have favorable acceptance
Online product review, 2017 [16].	- Investigate buyers' post-purchase behavior on feedback ratings	eBay	Statistical results of buyer and seller
Hotel: Customer satisfaction, 2017 [17]	- Used online comments	TripAdvisor.com	Some clearer managerial implication to understand
Electronic health records, 2017 [18]	- Extracted SQL with LIKE and wildcards. - Using R package	Residents' data in clinical narratives system	Create physiotherapy corpus and generate a new body of knowledge (care and lives)

2.3. Social network analysis

Social network and database system are crucial sources for generating data which directly or indirectly gives a lot of information. The extracting of latent data and mining information is an intense challenge for researchers that require the appropriate technology to discover metadata and get valuable information by analyzing and visualizing the results. From the previous survey, researchers have a lot of great techniques that present various perspective analysis including web mining, link analysis, network analysis, content analysis and others [19]. Social Network Analysis (SNA) was proposed as the best significant technique to study

and analyze the large social structure [20]. The stakeholders are called “nodes” or “actor” and the relationship between the actors are line or ties [21].

SNA investigate the interactions among actors in a network formed throughout their relationship as a binary point of view [22, 23]. The overall social network can be summarized and analyzed from a mathematical analysis to reveal underlying properties such as patterns of coordination and hierarchical structure from the large or complex dataset [24]. If a relationship between two actors was setting in two different strengths in weighted and unweighted networks the obtained results that are extracted from weighted are more accurate and detailed [22]. Moreover, SNA can model the structure of a network consisting of ties between nodes.

The traditional goal of SNA is to identify and characterize the structure of a network display in visualization as graph of relationship [25]. Researchers can retrieve content from database and accesses through sanctioned API (Application Programming Interface). The API response is the meaningful queries to collect relevant data [19]. The Scopus APIs expose curated abstracts and citation data from all scholarly journals indexed by Scopus, Elsevier's abstract and citation database [26]. Some studies have analyzed the journal impact factor, average number of cited references, self-citing, the co-authoring, author ordering, etc. and some researchers choose to study the structures of social network among scientific researchers.

The cooperation research networks have been studied in different research fields, including computer science, social sciences, engineering, and natural sciences [27, 28]. The network characteristics in analyzing research collaboration network consist of degree/strength, characteristic path length, clustering coefficient and diversity. An understanding of the generative mechanisms of highly dynamic research collaborations can promote more benefit collaboration among researchers in their community [29].

Prior research has shown some evidence in contribution of SNA in many field. SNA can help researchers to comprehend that the adopted research has grown sustainable and the collaboration of interdisciplinary teams are across the geography and disciplinary areas. Inspired from these insights, this paper models the solution by inferring published papers of STEM scholars as a large dataset to cluster the latent research information. Furthermore, the work represents the researching keywords and authors by using word cloud and social network analysis in graph visualization techniques.

3. Research Methodology

In carrying out the research design the application of a text mining is essential for analytic such as word cloud that can identify the knowledge structures of scholars' research and to pinpoint the scholars' research opportunities. In this study, two analytical techniques were employed: word cloud and social network analysis. The data was gathered from a combination of different data base system and data set of their education information and existing papers. This research was related to scholars working in different parts of Thailand. The suitable sampling method of this research was the cluster sampling that can be related with the industrial cluster. For this case, a sample group of scholars working at universities located in the northern region of Thailand were selected from all of the contacts in the

government's database system and Scopus from December, 2017. This framework was divided into three steps:

Step 1. Data acquisition

The purpose of this step is retrieving the two main data set: scholars' education information and published papers. The educational data was retrieved from the existing government's database system (http://dpst.ipst.ac.th/index.php/2017-11-29-04-36-35/student_detail). After that, first name and last name were retrieved from the relevant published papers and their co-author from Scopus with Python programming. The Scopus API allows requesting information with different levels of detail that programmer can request full data to select the most valuable information.

Step 2. Data preparation

This step aims to combine information obtained from two database systems for step 3. The data set items are first name, last name, education information, department, group type, e-mail, journal ID, published year, co-author, number of author, title, keywords, abstract, publication type, citation, etc. To characterize the data set for word cloud, SNA, and importing them into the new database system the researcher had to clean and transform the obtained data to be converted into the CSV format or an Excel file.

Step 3. Data analysis

Bibliographic information on article from step two are the data set that is suitable in preparing for analysis with different techniques. Therefore, word clouds generated the keyword visualization for five data sets: one for overall keywords and four for analyzing the keywords of each categories-STEM. This study was done to perform the relationship between the actors of SNA by using the authors' last name. Final step is extracting other results with the SQL statement in phpMyAdmin.

4. Results and Discussion

The result section is composed of three subsections: the first presents the characterization of the data sets, the second reveals published research keywords from word cloud, and the third subsection displays the core collaboration network of researchers by using SNA.

4.1. Characterization of the data sets

375 scholars' data were retrieved from the government's database system working in the north of Thailand as a researcher or lecturer in the government institution. In the next step, a questionnaire was sent via e-mail and the feedback was crosschecked. Figure 1 presents the number of scholars separated into four categories S-T-E-M. 19% of scholars quit their job from the government institution, moved out of the northern region of Thailand, retired, and for some other reasons. 81% of the scholars are working at the government institutions. The total number of active scholars is 304 persons that was divided into four groups: the science group is the biggest group with 65% of all active scholars, technology group is 5% of all active scholars, engineering group is 17% of all active scholars, and mathematics is 13% of all active scholars. The number of research is likely to increase in the same trend.

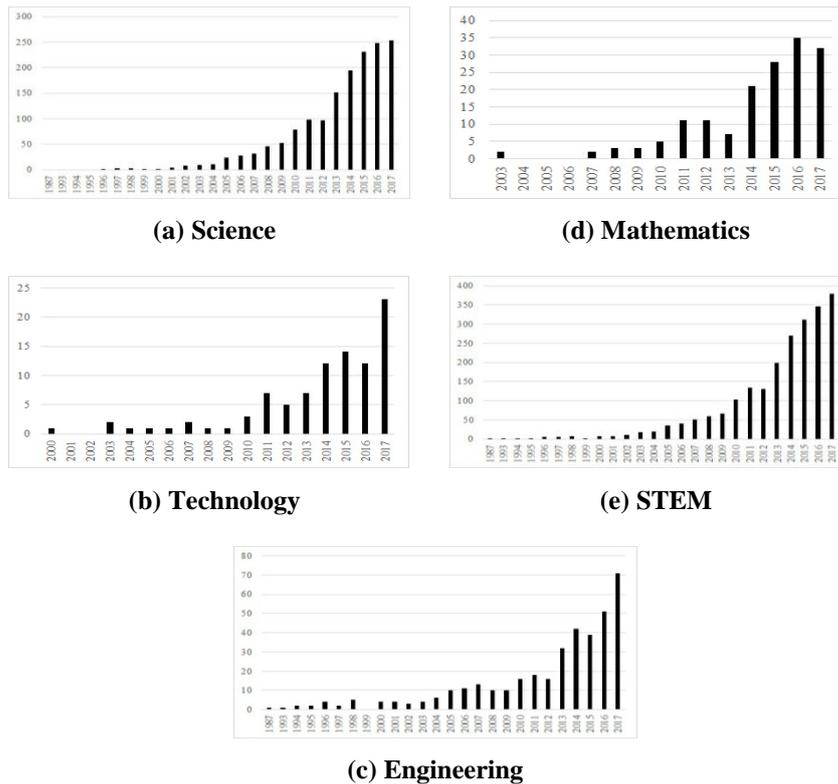


Fig. 1. Statistics of data set of research from 1987 – 2018.

4.2. Reveal published research keywords

In this case we used the specific open source software for visualization which is known as wordclouds.com (<https://www.wordclouds.com>). This technique can use cluster analysis to identify groups of word or document-clustering or factoring. The final data set that is ready for word analysis was comprised of 2,256 records set. 1,330 persons are co-authored papers and 225 persons are scholars. This study found that 74.01% had published papers and 25.99% did not have any publication information in Scopus. This analysis included visualization and comparison of keywords, which will lead to the identification of latent keywords in an overview perspective and focus on each category as shown in Table 2. In Fig. 2, the word cloud of each research areas is displayed. Research in the STEM domain: Science group, Technology group, Engineering group and Mathematics group.

As mentioned previously the section number on scientific scholars are more numerous than all other groups. Therefore, four keywords of scientific scholars are the top five of STEM keywords. Latent information from published article has been discovered for their pass work. The Science and Engineering group work mainly on the study about properties and analysis of various issues. Observations on the Technology group have shown that they emphasize on the functions and systems. Also, the most interesting thing is that the research of the Mathematical group focusses on energy,

Engineering, and Mathematics, respectively. In Fig. 3, the SNA of each research areas are displayed. Some statistics will help to understand that the information has never been revealed in the systematic analysis. The schematic of the collaboration network used the open-source visualization platform, Gephi 0.9.2. The analytic was separated into two views: node overview [30] and edge overview [31]. The information is presented in graphics and statistics. Each node presents a researcher, and the edge between the two studies indicate that they have co-authored at least on paper. The node thickness (degree) depends on the number of different co-authors of each published article. The Science domain has a large concentration. Moreover, each category has a value that is very similar in the modularity [32, 33]. The type of structure for network relationship of groups of researchers in this case is a Swarm Cluster and the structure for network relationship of single researchers is 1-N node.

Table 3. Metrics and statistics obtained of collaborative network.

Categories	Average Degree	Max Degree	Average Path Length	Modularity
STEM	1.954	22	10.841	0.963
Science	0.587	17	10.251	0.974
Technology	0.84	8	1.582	0.907
Engineering	0.753	10	3.737	0.954
Mathematics	0.821	9	2.371	0.948

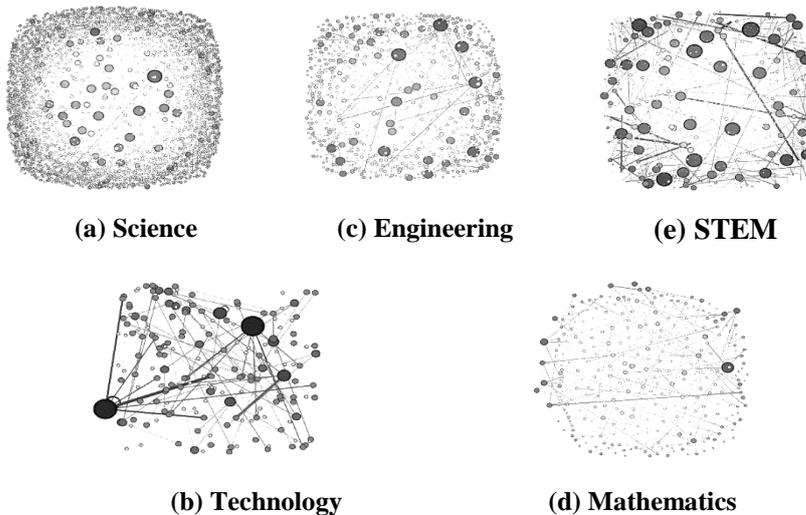


Fig. 3. Collaborative network.

The STEM has the highest average degree of 1.954 meaning that the researcher in this field is likely to work as a two-people team, while other researchers from other disciplines tend to work alone. For STEM and Science, the average path length is significantly higher than others implying that the researcher is likely to publish the research with the same college.

5. Conclusions

This study used text mining analytics techniques to identify the knowledge structures of scholar' published articles that emphasized on the STEM domain. A wide array of methodologies that support in collecting data and identifying the flow of research with text mining techniques have been applied in diverse areas such as medical information, transportation, and criminal investigations [7]. The research summary analysis discovered that Thai scholars have a considerable number of researchers in the scientific field. However, number of Thai scholars in the technology sector are very small; only 5% of who are based in the northern region of Thailand.

The policy that encourages scholar to study at the doctoral degree merged with the concrete research policy that forces the researchers to develop innovation and research which can lead to more publications for Thailand as well. Through word cloud and SNA, an empirical approach can be improved in a systematic way. Social media tools can improve the research process by recording the online series of scholarly behaviors and present the implicit information about research papers [8]. This study aimed to investigate the use of reliable sources such as government's database system and Scopus that can create a strong support for evaluating the past direction of national scholars' interested research field. Both the word cloud and social network analysis play an important role in this paper.

The concept provides the keywords topic modelling and reveals the relationship among researchers and their co-author in conferences or journals from a prospective database system. Word could be used for generating keywords and give a list of research keywords of scholars. After knowing the sheer valuable result, based on the scholar current situation and STEM research direction, there should be a further study on describing how to approach and establish the various policy of government focusing on STEM research to generate a great deal of contribution to super cluster and other clusters in a specific industry that is relying on the country's direction and making an announcement for clustering.

Abbreviations

CSV	Comma Separated Values
NUSANTARA	Government Human Capital Knowledge Management of Republic Indonesia
SNA	Social Network Analysis
STEM	Science, Technology, Engineering and Mathematics

References

1. Cahyaningsih, E.; Senses, D.I.; and Sari, W.P. (2015). Defining knowledge of government human capital management: A qualitative study. *2015 International Conference on Information Technology Systems and Innovation (ICITSI)*, 1-6.
2. Horesh, R.; Varshney, K.R.; and Yi, J. (2016). Information retrieval, fusion, completion, and clustering for employee expertise estimation. *2016 IEEE International Conference on Big Data (Big Data)*, 1385-1393.
3. Wardhani, A.R.; Acur, N.; and Mendibil, K. (2016). Human capital, social capital and innovation outcome: A systematic review and research agenda.

- 2016 *IEEE International Conference on Industrial Engineering and Engineering Management (IEEM)*, 355-359.
4. Sensesuse, D.I.; Cahyaningsih, E.; Hidayat, E.; Sukmasetya, P.; and Wibowo, W.C. (2017). Implementation of government human capital knowledge management of republic Indonesia case study government ministries. *International Conference on Information Technology Systems and Innovation (ICITSI)*, 139-144.
 5. Fund, T.R. (2015). *The royal golden jubilee (RGJ-Ph.D. programme)*. (Thailand Research Fund) Retrieved January 05, 2018, from <http://rgj.trf.or.th>
 6. Sendhikumar, S.; Srivani, M.; and Mahalakshmi, G. (2017). Generation of word clouds using document topic models. *2017 Second International Conference on Recent Trends and Challenges in Computational Models (ICRTCCM)*, 306-308.
 7. Kim, S.K.; and Oh, J. (2017). Information science techniques for investigating research area: A case study in telecommunications policy. *The Journal of Supercomputing*, 1-28.
 8. Li, J.; Shin, S.Y.; and Lee, H.C. (2017). Text mining and visualization of papers reviews using R language. *Journal of information and communication convergence engineering*, 15(3), 170-174.
 9. Chang, S.T.; Huang, H.W.; Wu, L.Y.; Syu, C.C.; Lin, C.J.; Liu, Y.; and Li, H.H. (2016). Using text cloud technology to build predictive models of disease-taking osteoporosis case for example. *2016 International Conference on Machine Learning and Cybernetics (ICMLC)*, 644-648.
 10. Jayashankar, S.; and Sridaran, R. (2017). Superlative model using word cloud for short answer evaluation in eLearning. *Education and Information Technologies*, 22(5), 2383-2402.
 11. Thilagavathy, R.; and Sabitha, R. (2017). Using cloud effectively in concept based text mining using grey wolf self organizing feature map. *Cluster Computing*, 1-11.
 12. Jin, Y. (2017). Development of word cloud generator software based on Python. *Procedia Engineering*, 17, 788-792.
 13. He, S.; Cao, Y.; and Xiong, H. (2017). Generate Galaxy-like word cloud using molecular cloud evolution. *2017 9th International Conference on Intelligent Human-Machine Systems and Cybernetics (IHMSC)*.
 14. Dhou, K.; Hadzikadic, M.; and Faust, M. (2017). Typeface size and weight and word location influence on relative size judgments in tag clouds. *Journal of Visual Languages and Computing*, 44, 97-105.
 15. Renoust, B.; Ren, H.; Melancon, G.; Viaud, M.L.; and Satoh, S. (2017). FaceCloud: Heterogeneous cloud visualization of multiplex networks for multimedia archive exploration. *MM '17 Proceedings of the 2017 ACM on Multimedia Conference*, 1235-1236.
 16. Yap, C.S.; Ong, M.Y.; and Ahmad, R. (2017). Online product review, product knowledge, attitude, and online purchase behavior. *International Journal of E-Business Research*, 13(3), 20.
 17. Cherapanukorn, V.; and Charoenkwan, P. (2017). Word cloud of online hotel reviews in Chiang Mai for customer satisfaction analysis. *2017 International Conference on Digital Arts, Media and Technology*, 146-151.

18. Delespierre, T.; Denormandie, P.; Bar-Hen, A.; and Josseran, L. (2017). Empirical advances with text mining of electronic health records. *BMC Medical Informatics and Decision Making*, 17,127.
19. Goswami, A.; and Kumar, A. (2017). Challenges in the analysis of online social networks: A data collection tool perspective. *Wireless Personal Communications*, 97(3), 4015-4061.
20. Lee, C.Y.; Chong, H.Y.; Liao, P.C.; and Wang, X. (2018). Critical review of social network analysis applications in complex project management. *Journal of Management in Engineering*, 34(2).
21. Zhang, S.; and Fang, Y. (2018). Research on construction project organization based on social network analysis. *Wireless Personal Communications*, 1-11.
22. Andrade, R.L.; and Rego, L.C. (2018). The use of nodes attributes in social network analysis with an application to an international trade network. *Physica A: Statistical Mechanics and its Applications*, 491, 249-270.
23. Lovric, M.; Re, R.D.; Vidale, E.; Pettenella, D.; and Mavsar, R. (2018). Social network analysis as a tool for the analysis of international trade of wood and non-wood forest products. *Forest Policy and Economics*, 86, 45-66.
24. McCurdie, T.; Sanderson, P.; and Aitken, L.M. (2018). Applying social network analysis to the examination of interruptions in healthcare. *Applied Ergonomics*, 67, 50-60.
25. Jorgensen, T.; Forney, K.; Hall, J.; and Giles, S. (2018). Using modern methods for missing data analysis with the social relations model: A bridge to social network analysis. *Social Networks*, 54, 26-40.
26. Elsevier B.V. (2018). *Elsevier Developers*. (Elsevier B.V.) Retrieved January 7, 2018, from https://dev.elsevier.com/sc_apis.html.
27. Montoya, F.G.; Alcaide, A.; Banos, R.; and Manzano-Agugliaro, F. (2018). A fast method for identifying worldwide scientific collaborations using the Scopus database. *Telematics and Informatics*, 35(1), 168-185.
28. Kim, Y.; and Cho, N.W. (2018). Research trends in social network analysis using topic modeling and network analysis. *An International Journal of Research and Surveys, ICIC Express Letters*, 12(1), 71-78.
29. Bian, J.; Xie, M.; Topaloglu, U.; Hudson, T.; Eswaran, H.; and Hogan, W. (2014). Social network analysis of biomedical research collaboration networks in a CTSA institution. *Journal of Biomedical Informatics*, 52, 130-140.
30. Latapy, M. (2008). Main-memory triangle computations for very large (Sparse (Power-Law)) graphs. *Theoretical Computer Science*, 407(1-3), 458-473.
31. Brandes, U. (2001). A Faster Algorithm for Betweenness Centrality. *Journal of Mathematical Sociology*, 25(2), 163-177.
32. Blondel, V.; Guillaume, J.L.; Lambiotte, R.; and Lefebvre, E. (2008). Fast unfolding of communities in large networks . *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10):P10008.
33. Lambiotte, R.; Delvenne, J.C.; and Barahona, M. (2015). Laplacian dynamics and multiscale modular structure in networks. *IEEE Transactions on Network Science and Engineering*, 1(2), 76-90.