

UNIVERSITY OF WISCONSIN MILWAUKEE

Network Analysis of Scientific  
Collaboration and Co-authorship of the  
Trifecta of Malaria, Tuberculosis and  
HIV/AIDS in Benin.

by

Gbedegnon Roseric Azondekon

A thesis submitted in partial fulfillment for the degree of  
Doctor of Philosophy (PhD) in Biomedical and Health Informatics

in the  
College of  
Engineering and Applied Sciences

March 2018

## ABSTRACT

Despite the international mobilization and increase in research funding, Malaria, Tuberculosis and HIV/AIDS are three infectious diseases that have claimed more lives in sub Saharan Africa than any other place in the World. Meanwhile, research collaborations have peaked in an ultimate effort to dramatically decrease the mortality and morbidity of those diseases on the continent. Consortia, research network and research centers both in Africa and around the world team up in a multidisciplinary and transdisciplinary approach to boost efforts to curb these diseases. Other studies have already reported a universal rise in terms of scientific collaborations. Despite the progress in research, very little is known on the dynamics of research collaboration in the fight of these Infectious Diseases in Africa resulting in a lack of information on the relationship between African research collaborators. Understanding the structure of these complex networks is capital since it can help improve research prioritization, identification of prolific researchers, better design, strategic planning and implementation of research program], and promote cooperation and translational research initiatives. In this doctoral thesis proposal, we propose to document, describe and analyze the scientific collaboration, and co-authorship of the research conducted in the Republic of Benin on Malaria, Tuberculosis and HIV/AIDS.

## *Acknowledgements*

First, I would like to express my sincere gratitude to my advisor Prof. Susan McRoy for the continuous support of my Ph.D study, for her patience, motivation, and immense knowledge. Her guidance helped me in all the time of research and writing of this dissertation.

I would also like to offer my special thanks to Dr Charles Welzig and the Welzig Neuroscience and Neurotechnology lab at the Medical College of Wisconsin. I benefited from the computational resource available in the lab to run the analyses and the computational intensive simulations. My thanks also go to Zachary James Harper for his availability. Without his precious support it would not be possible to conduct this research.

My special thanks are extended to Dr Spencer (Chiang Ching) Huang of the Joseph Zilber School of Public Health. He has been crucial to the successful continuation of my PhD journey at the University of Wisconsin Milwaukee. He provided me an opportunity to join, and who gave access to the laboratory and research facilities.

Finally, I would like to thank the rest of my thesis committee: Prof Christine Cheng, and Dr. Rohit Kate, and Dr. Zhang Qing, for their insightful comments, encouragement, and crucial recommendations and suggestions which widen my research from various perspectives.

# Contents

<b>Abstract</b>	i
<b>Acknowledgements</b>	ii
<b>List of Figures</b>	vi
<b>List of Tables</b>	xi
<b>Abbreviations</b>	xiii
<b>General Introduction</b>	1
<b>1 Literature Review</b>	5
1.1 Brief Overview of Malaria, Tuberculosis and HIV/AIDS . . . . .	5
1.2 Network Analysis of Scientific Research collaboration . . . . .	9
1.3 General and specific Objectives . . . . .	17
1.4 Gap in the literature . . . . .	19
<b>2 Methodology</b>	20
2.1 Overview . . . . .	20
2.2 Data Collection . . . . .	21
2.2.1 Text Mining and Network generation . . . . .	21
2.2.2 Author Name Disambiguation . . . . .	22
2.3 Descriptive Data Analysis . . . . .	23
2.3.1 Characterizing Network cohesion . . . . .	25
2.4 Modeling of Network Data . . . . .	27
2.4.1 Mathematical Modeling . . . . .	27
2.4.2 Statistical Modeling . . . . .	29

2.4.2.1	Stochastic Block Model . . . . .	29
2.4.2.2	Exponential Random Graph Model . . . . .	30
2.4.2.3	Temporal Exponential Random Graph Model . . . . .	32
2.4.2.4	Latent Network Model . . . . .	33
<b>3</b>	<b>Results: The Malaria Co-authorship Network</b>	<b>36</b>
3.1	Data . . . . .	36
3.2	Descriptive Data Analysis . . . . .	38
3.2.1	Network Cohesion . . . . .	40
3.3	Modeling . . . . .	43
3.3.1	Mathematical Modeling . . . . .	43
3.3.2	Statistical Modeling . . . . .	46
3.3.2.1	Stochastic Block Model . . . . .	46
3.3.2.2	Exponential Random Graph Model . . . . .	48
3.3.2.3	Temporal Exponential Random Graph Model . . . . .	52
3.3.2.4	Latent Network Model . . . . .	58
3.4	Discussion and Conclusion . . . . .	61
<b>4</b>	<b>Results: The HIV/AIDS Co-authorship Network</b>	<b>68</b>
4.1	Data . . . . .	68
4.2	Descriptive Data Analysis . . . . .	70
4.2.1	Network Cohesion . . . . .	71
4.3	Modeling . . . . .	74
4.3.1	Mathematical Modeling . . . . .	74
4.3.2	Statistical Modeling . . . . .	78
4.3.2.1	Stochastic Block Model . . . . .	78
4.3.2.2	Exponential Random Graph Model . . . . .	81
4.3.2.3	Temporal Exponential Random Graph Model . . . . .	83
4.3.2.4	Latent Network Model . . . . .	87
4.4	Discussion and Conclusion . . . . .	90
<b>5</b>	<b>Results: The Tuberculosis Co-authorship Network</b>	<b>93</b>
5.1	Data . . . . .	93
5.2	Descriptive Data Analysis . . . . .	95
5.2.1	Network Cohesion . . . . .	97
5.3	Modeling . . . . .	100
5.3.1	Mathematical Modeling . . . . .	100
5.3.2	Statistical Modeling . . . . .	103
5.3.2.1	Stochastic Block Model . . . . .	103
5.3.2.2	Exponential Random Graph Model . . . . .	105

5.3.2.3	Temporal Exponential Random Graph Model . . . . .	107
5.3.2.4	Latent Network Model . . . . .	111
5.4	Discussion and Conclusion . . . . .	113
<b>6</b>	<b>AuthorVis: A Co-authorship Visualization and Scientific Collaboration Prediction tool</b>	<b>117</b>
6.1	Background . . . . .	117
6.2	Related work . . . . .	118
6.3	Design and Architecture . . . . .	120
6.3.1	Data . . . . .	121
6.3.2	Network Visualization . . . . .	121
6.3.3	Web Framework . . . . .	122
6.4	Deployment . . . . .	123
	<b>General Conclusion</b>	<b>124</b>
	<b>Bibliography</b>	<b>127</b>

# List of Figures

3.1	Degree distribution of the Malaria Co-authorship network . . . . .	39
3.2	Evolution of the published Malaria related documents, authors and collaborations from January 1996 to December 2016 . . . . .	39
3.3	Malaria Co-authorship network – Main component. Authors (vertices) of the same color belong to the same research community or cluster . . . . .	42
3.4	Monte-Carlo simulations: Number of detected communities by the random graph models . . . . .	44
3.5	Monte-Carlo simulations: Number of detected communities by the Watts-Strogatz and the Barabási-Albert models . . . . .	45
3.6	Summary of the goodness-of-fit of the SBM analysis on the Malaria co-authorship network. . . . .	47
3.7	Distribution of national, international and regional authors by communities detected by the SBM in the Malaria network. . . . .	48
3.8	Summary of the goodness-of-fit of the SBM analysis highlighting interactions between the top 5 larger classes of the Malaria co-authorship network.	49

3.9	ERGM goodness-of-fit of final model 4 assessment.	53
3.10	Topological structure of the different snapshots of the malaria co-authorship network.	54
3.11	Goodness-of-fit assessment for the final Malaria TERGM Model 4 with temporal dependencies.	57
3.12	Visualizations of the Malaria co-authorship network with layouts determined according to the inferred latent eigenvectors in the LNM models.	59
3.13	ROC curves comparing the goodness-of fit of the Malaria co-authorship network for three different eigenmodels, specifying (i) no pair specific covariates (blue), (ii) nodal covariates (red), and (iii) nodal and dyadic covariates (green), respectively.	60
4.1	Evolution of the published HIV related documents, authors and collaborations from January 1996 to December 2016	70
4.2	Degree distribution of the HIV/AIDS Co-authorship network	71
4.3	Log-Average Neighbor degree Distribution of the HIV/AIDS Co-authorship network	72
4.4	Topological Structure of the HIV/AIDS Co-authorship Network. Authors (vertices) of the same color belong to the same research community or cluster	75
4.5	Monte-Carlo simulations of the HIV/AIDS network: Number of detected communities by the random graph models	76

4.6	Monte-Carlo simulations of the HIV/AIDS network: Number of detected communities by the Watts-Strogatz and the Barabási-Albert models . . . . .	77
4.7	Summary of the goodness-of-fit of the SBM analysis on the HIV/AIDS co-authorship network. . . . .	78
4.8	Distribution of national, international and regional authors by communities detected by the SBM in the HIV/AIDS network. . . . .	79
4.9	Summary of the goodness-of-fit of the SBM analysis highlighting interactions between the largest classes of the HIV/AIDS co-authorship network. .	80
4.10	ERGM goodness-of-fit of final model 3 assessment on the HIV/AIDS co-authorship network. . . . .	83
4.11	Topological structure of the different snapshots of the HIV/AIDS co-authorship network. . . . .	84
4.12	Goodness-of-fit assessment for the final HIV/AIDS TERGM Model 4 with temporal dependencies of the HIV/AIDS co-authorship network. . . . .	87
4.13	Visualizations of the HIV/AIDS co-authorship network with layouts determined according to the inferred latent eigenvectors in the LNM models. . . . .	89
4.14	ROC curves comparing the goodness-of fit of the HIV/AIDS co-authorship network for the model specifying (i) no pair specific covariates (blue) and the model specifying (ii) nodal covariates (red). . . . .	90
5.1	Evolution of the published TB related documents, authors and collaborations from January 1996 to December 2016 . . . . .	95

5.2	Degree distribution of the TB Co-authorship network . . . . .	96
5.3	Log-Average Neighbor degree Distribution of the TB Co-authorship network	96
5.4	Topological Structure of the Tuberculosis Co-authorship Network. Authors (vertices) of the same color belong to the same research community or cluster	99
5.5	Monte-Carlo simulations of the TB network: Number of detected commu- nities by the random graph models . . . . .	101
5.6	Monte-Carlo simulations of the TB network: Number of detected commu- nities by the Watts-Strogatz and the Barabási-Albert models . . . . .	101
5.7	Summary of the goodness-of-fit of the SBM analysis on the Tuberculosis co-authorship network. . . . .	103
5.8	Distribution of national, international and regional authors by communities detected by the SBM in the TB network. . . . .	104
5.9	ERGM goodness-of-fit of final model 3 assessment on the TB co-authorship network. . . . .	107
5.10	Topological structure of the different snapshots of the TB co-authorship network. . . . .	108
5.11	Goodness-of-fit assessment for the final TB TERGM Model 4 with tempo- ral dependencies of the TB co-authorship network. . . . .	110
5.12	Visualizations of the TB co-authorship network with layouts determined according to the inferred latent eigenvectors in the LNM models. . . . .	112

5.13 ROC curves comparing the goodness-of fit of the TB co-authorship network for the model specifying (i) no pair specific covariates (blue) and the model specifying (ii) nodal covariates (red). . . . .	113
6.1 Screenshots of the Shiny application interface. . . . .	120
6.2 Screenshot of the co-authorship network visualization interface. . . . .	122

# List of Tables

3.1	Malaria Bibliographic Search Queries. . . . .	37
3.2	List of the most important authors and collaborations in the Malaria Co-authorship network . . . . .	40
3.3	ERGM of the co-authorship Malaria network. . . . .	52
3.4	Temporal ERGM of Malaria Co-authorship Network. . . . .	56
4.1	HIV/AIDS Bibliographic Search Queries. . . . .	69
4.2	List of the most important authors and collaborations in the HIV/AIDS Co-authorship network . . . . .	73
4.3	ERGM of the HIV/AIDS Co-authorship Network. . . . .	81
4.4	Temporal ERGM of the HIV/AIDS Co-authorship Network. . . . .	86
5.1	TB Bibliographic Search Queries. . . . .	94
5.2	List of the most important authors and collaborations in the Tuberculosis Co-authorship network . . . . .	98
5.3	ERGM of the TB Co-authorship Network. . . . .	105

5.4 Temporal ERGM of the TB Co-authorship Network. . . . .	110
------------------------------------------------------------	-----

# Abbreviations

<b>AIC</b>	Akaike's Information Criterion
<b>AIDS</b>	Acquired Immune Deficiency Syndrome
<b>AND</b>	Author Name Disambiguation
<b>ARV</b>	Antiretroviral
<b>AUC</b>	Area Under the Curve
<b>AWS</b>	Amazon Web Services
<b>BIC</b>	Bayesian Information Criterion
<b>CD4</b>	Cluster of Differentiation 4
<b>CGI</b>	Common Gateway Interface
<b>CI</b>	Confidence Interval
<b>CPU</b>	Central Processing Unit
<i>df</i>	degree of freedom
<b>DOI</b>	Digital Object Identifier
<b>ELISA</b>	Enzyme-Linked Immunosorbent Assay
<b>ERGM</b>	Exponential Random Graph Model
<b>Fig.</b>	Figure
<b>GLM</b>	Generalized Linear Model
<b>HIV</b>	Human Immunodeficiency Virus
<b>ICL</b>	Integration Classification Likelihood
<b>JSON</b>	JavaScript Object Notation
<b>LNM</b>	Latent Network Model

<b>MeSH</b>	Medical Subject Headings
<b>MCMC</b>	Markov Chain Monte Carlo
<b>MCMLE</b>	Monte Carlo Maximum Likelihood Estimation
<b>MDG6</b>	Millenium Development Goal 6
<b>MLE</b>	Maximum Likelihood Estimation
<b>MPLE</b>	Maximum PseudoLikelihood Estimation
<b>PA</b>	Preferential Attachment
<b>PR</b>	Precision Recall
<b>REF</b>	Reference
<b>ROC</b>	Receiver Operating Characteristics
<b>SAOM</b>	Stochastic Actor-Oriented Model
<b>SBM</b>	Stochastic Block Model
<b>SCI</b>	Science Citation Index
<b>SE</b>	Standard Error
<b>SNA</b>	Social Network Analysis
<b>SW</b>	Small World
<b>TB</b>	Tuberculosis
<b>TERGM</b>	Temporal Exponential Random Graph Model
<b>US</b>	United States
<b>WHO</b>	World Health Organization
<b>WOS</b>	World Of Science

*For my family...*

# General Introduction

Infectious diseases have long claimed the lives of millions of people worldwide. They disproportionately affect the developing nations where 90% of the deaths are caused by very few diseases among which Malaria, Tuberculosis (TB) and HIV/AIDS [1]. Malaria, TB and HIV/AIDS remain the three major public health concerns in Sub Saharan Africa where they are responsible for high mortality, morbidity rates and impact negatively on the socioeconomic way of life of the populations [2, 3]. These three diseases have been given special attention at the Millenium Declaration in its 6<sup>th</sup> Goal of Millenium Development [4]. Initiatives such as the US President's Malaria Initiative, the Global Fund for Malaria, TB and HIV/AIDS and the President's Emergency Plan for AIDS have led to the investment of more than 70 million of US dollars to encourage Research and Development, Private-Public partnership as well as to reinforce the activities of non-governmental organizations within the healthcare systems of the affected countries [5–7].

The Global Fund disbursement in 2010 peaked at over 1.45 billion dollars for HIV/AIDS, 416 million dollars for TB and 714 million dollars for Malaria [8, 9]. With these financial supports at hand, efforts have led to a sharp increase of public health interventions and

## *General Introduction*

---

many positive public health outcomes in terms of the reduction of mortality and morbidity related to those diseases [10]. For example, in Benin, such increase in public health interventions translated in the financing, successful implementation and sustainability of the entomological surveillance of malaria for more than six years since 2008 [11]. Encouraged and motivated by the success stories in controlling these diseases, some authors formulated the ambitious zero incidence goal of TB and HIV and the zero death goal of the three diseases by 2015 [12].

After the declaration of the Millenium Challenge Goal 6 in 2000, significant progress has been made in the treatment and prevention of Malaria, TB and HIV/AIDS, leading to the reverse of the mortality and morbidity due to these three diseases. Nevertheless, sub Saharan Africa still carries the burden of these diseases. For example, in 2009, 2.6 million new cases and 1.8 million of death related to HIV were estimated out of which 68% and 72% of respectively new cases and deaths were in Africa [13]. TB cases were estimated at 9.4 million and 1.3 million deaths out of which HIV-positive cases make up 12% of all cases and 23% of all TB deaths [14]. Although the rapid expansion of vector control strategies worldwide, malaria was responsible of 225 million cases and 781,000 death in 2009 out of which over 90% were in Africa [9].

In the Republic of Benin, between 2000 and 2013, the impact of the increase in funding has led to an annual decrease in the incidence of 7.6%, 0.6% and 5.2% respectively in HIV/AIDS, TB and Malaria. Similar results were obtained in terms of prevalence with a decrease of 1.3% in HIV/AIDS and 0.8% in TB. Annual death rates decreased also at about 3.1%, 1.2% and 5.3% respectively in HIV/AIDS, TB and Malaria [9, 13, 14].

## *General Introduction*

---

Successful scientific collaborations have led to the eradication of chickenpox and the near eradication of poliomyelitis through the development of vaccines [15]. For Malaria and HIV/AIDS, the development of a vaccine has proven significantly difficult to develop despite the decades of active research that has not been successful so far [16–18]. This is why researchers need to form continuous and sustainable collaborations through intensive network practices that go beyond the regional boundaries [19]. Scientific collaborations give researchers the opportunity to work and learn from each other. Such collaborations are further needed to overcome the overgrowing challenge of co-infections of HIV/AIDS and Tuberculosis [20, 21]. In the republic of Benin, Malaria, TB and HIV/AIDS have become a common aspect of the public health system. The three are the main impediments of economic and social progress that are characteristics of poverty. According to a 2000 World Health Organization (WHO) press report, malaria slows economic growth on the African continent by 1.3% each year [22]. And it is known that Tuberculosis and HIV/AIDS patients experienced severe economic burden in terms of access to health care, treatment and diagnosis [23]. The situation is further compounded by the poorly developed immunity among the children and the elderlies, and the predominant malnutrition problem experienced by a majority of the population [15]. The disappointing aspect is that the extensive research conducted has not prevented these three diseases from outpacing the proposed solutions and the progress made [24].

Therefore, in this thesis to document, we document, describe and analyze the different aspects of scientific research collaboration of the three leading infectious diseases in the

## *General Introduction*

---

Republic of Benin. The social network analysis of research collaboration approach is chosen to reveal undiscovered knowledge on effort of researchers in working together towards the reduction of the burden of Malaria, TB and HIV/AIDS. Modern times have rendered research and scientific collaborations irreplaceable policy formulations processes. This is because research collaboration forms a stable basis for the provision of evidence based information in the formulation of fundamental principles and guidelines for the elaboration of public health strategies, particularly in developing countries like Benin. For this reason, this thesis focuses on the Network analysis of the scientific collaborations through co-authorship network analysis.

# Chapter 1

## Literature Review

### 1.1 Brief Overview of Malaria, Tuberculosis and HIV/AIDS

AIDS is a health condition caused by the Human Immunodeficiency Virus (HIV) [25, 26].

HIV infects and attacks the cells that are responsible for the immune system in the body (CD4 cells) that provide protection against infections and illness. The virus infects the human host by making him vulnerable and unable to fight future infections [27]. The virus eventually weakens and kills the CD4 cells resulting in a weak immune system and vulnerability to diseases. HIV is transmitted through body fluids exchange, and the infection exists in four stages. The first stage is the primary infection stage and lasts within 2 to 4 weeks. It is characterized by flu-like symptoms, and the infected person is highly contagious. The second stage is the asymptomatic stage that may last for about ten years, and the infected person does not display significant symptoms of the infections.

## *Literature Review*

---

The third stage is the symptomatic stage. At this stage, the virus weakens the immune system, and the infected person suffers from both mild and chronic symptoms as the infected person suffers opportunistic diseases. Illnesses like malaria and TB in HIV infected subjects, are experienced in a severe manner. The fourth stage is AIDS; it causes death within two years if left untreated [27, 28].

According to the World Health Organization (WHO) the signs for HIV/AIDS change through the stages of infections as the disease progresses. To determine whether a person is infected, an HIV test needs to be conducted. ELISA method based HIV testing is one of the most common antibody-based testing method characterized by 99% accuracy rate [25]. It is recommended that a HIV negative test result should be confirmed after three months because the immune system can sometimes take up to 12 weeks to develop the tested antibodies [29]. It is however possible to get false negative results during the 12 weeks window period. The antiretroviral (ARV) drug therapy is initiated when the infected person reaches the third or fourth stage of infection to suppress the virus and boost the immune system. Such measures are taken because there is currently, no cure for HIV and the early initiation of the therapy may result in drug resistance [30–32].

TB is a highly infectious disease that is caused by *Mycobacterium tuberculosis*. The disease exists in active and inactive forms. The active form, also known as the open disease causes the infected person to suffer and to be highly infectious. The inactive/latent TB infection is not infectious, and the infected individual does not suffer from the signs and symptoms associated with the active disease. Healthy individuals with latent infection have approximately 10% probability of getting active TB disease over their life. Chances

## *Literature Review*

---

of infection are high in the first two years after the exposure to the bacteria, and in the case where the host develops any form of lung or immune system damage [33, 34]. On the other hand, in HIV infected individuals co-infected with TB, there exists a 10% annual chance of developing active TB [35–37]. Active TB in adults may result from re-infection with a new strain of TB or perhaps a reaction to the latent infection. Consequently, researchers insist that silica inhalation, HIV infection, and silicosis are responsible for the high risk of TB infection in the working adults' population [36, 38]. TB symptoms are characterized by a chronic cough, night-time fevers, profuse sweating, and significant weight loss within a short time. However, studies show that people with TB can be infectious prior to showing the symptoms or complaining of any form of pulmonary discomfort. In the worst case scenario, TB goes beyond the pulmonary and infects other parts of the body, especially for people infected with HIV. HIV complicates the manifestation of TB in terms of its symptoms and signs in 70% of the HIV/AIDS infected population suffering from TB [36]. Studies indicate that people with undetected open TB disease are the leading cause of TB infections. Even though TB is a treatable disease, the treatment procedure is extremely aggressive. The treatment procedure for first-time patients entails administration of a six-months dose under close medical supervision termed as directly observed therapy. The other challenge in the treatment is that there are approximately 25% TB-drugs resistance cases worldwide every year [39, 40]. Approximately 80% of people with TB can be cured of their active TB infection, however, HIV and Silicosis increases the risk of reinfection by 20%. The infection among individuals with silicosis, may cumulatively contribute to lung damage and work inability. Additionally, the HIV/AIDS

## *Literature Review*

---

increases the risk of opportunistic infections, which may result in a poor outcome for the TB treatment [41–43].

Malaria is an infectious disease caused by the Plasmodium parasites. Even though, malaria is predominantly found in the tropical regions, 48% of the instances of infections have been experienced in the Northern and Southern parts of America, Asia, and Africa, putting approximately 50% of the world's population at risk. The malaria pathogens are *Plasmodium ovale*, *Plasmodium malariae*, *Plasmodium vivax*, and *Plasmodium falciparum* which is the deadliest. The distribution of the disease matches that of its vectors, the female mosquitoes of the genus Anopheles [44, 45]. In the sub-Saharan African countries, the vector of the disease is *Anopheles gambiae s.l.* Malaria has a range of symptoms and signs that manifest differently from one person to another. The most common symptoms are fevers, gastrointestinal symptoms, and fatigue, headaches, and muscle aches. The malaria pathogen infects two hosts, the Anopheles mosquito, and the infected human. When the infected mosquito feeds from an individual, it injects sporozoites into the circulatory system of the bitten person. The sporozoites reside in the liver cells until they become mature schizonts. The schizonts rupture upon maturity and release merozoites, which infect the red blood cells [46]. The two most used malaria test are rapid tests using an instant result kit akin to the home pregnancy test device, and the blood smear test that is examined under the microscope for the presence of red blood cells that are infected by the parasite. Treatment entails administration of drugs that range in types. While some malaria drug prescriptions may have a three days dosage, others may have up to one week dosage [47, 48].

## 1.2 Network Analysis of Scientific Research collaboration

Collaboration in science is essential to research and development, knowledge discovery, technology and innovation. It occupies a predominant place in scientometrics. According to Leydesdorff and Milojevic [49], scientometrics uses quantitative and computational methods to analyzing and measuring science, communication in science and science policy. According to the same authors, the field of scientometrics emerged from Eugene Garfield's idea to improve Information Retrieval [50], followed by the creation of the Science Citation Index (SCI) in the 1960s, and the availability of scientific databases references publications. The discipline of Scientometrics is aimed at providing guidance to several research issues involving the measurement of science impact, the measurement of impact journals and institutional units, theories of citation, and the mapping of science. We aim at focusing on the mapping of science since it is essential to understanding the dynamic of science, informing policy decisions, and identifying important fields, research groups, specialties based on evidence from the literature [49]. Such goals can be achieved by mapping publications authors and analyzing patterns of collaborations between them. Since the publication of the first co-authored paper in 1665, scientific co-authorship has spread significantly throughout the scientific realm and the number of co-authored scientific publications have tremendously increased [51]. According to Wagner [52], the increase in international scientific co-authorship has been of a fast growth. International co-authorship originates from international collaborations between scientists. In general,

## *Literature Review*

---

international collaborations have more visibility than national collaborations and often result in publications in high impact journals [53].

The paradigm of co-authorship network is rooted in network theory. In a co-authorship network, the set of nodes is represented by the researchers and the set of edges describes the relationship between them. An edge between two researchers in such a network means that they both coauthor a publication. Unlike citation networks, the scientific community has dedicated less attention to co-authorship networks because of the long tradition of citation network analysis in bibliometric [19, 54]. Nevertheless, the analyses of how complex co-authorship networks form and evolve in time is crucial for identifying leading researchers in a particular scientific domain, and describe their extant to collaborate with their peers as well as the impact of their research [55]. An example of such investigation is illustrated in Newman scientific collaboration paper series on Biomedical research, physics and computer science co-authorship network [19, 54, 56, 57].

Taking publication as units, the analyses of scientific collaboration facilitate the study of trans and inter-disciplinary research by focusing on the dynamics of the collaboration networks [58]. In addition, these networks can provide important information regarding cooperation patterns among authors and their status and location in the structures of the scientific community [59]. Furthermore, Mali et al. [60] argue that studies such as co-authorship social network studies are highly relevant for funding organizations for promising and emerging topics support in science.

Although many authors have proposed different features for classifying co-authorship networks [61–63], the categorization features of Andrade et al. [61] identifies three levels

## *Literature Review*

---

of classification of scientific collaboration: the cross-disciplinary level with the intradisciplinarity and interdisciplinarity subdimensions, the cross-sectoral level with the intramural and extramural research collaboration subdimensions and the cross-national level including the national and international scientific collaboration subdimensions. For a full description of each level of scientific collaboration, we refer the reader to Mali et al. [60]. The methods of co-authorship network studies have emerged from social network analysis and graph theory. Such studies heavily relied upon access to scientific collaboration data sources such as SCOPUS, the Web Of Science, PubMed, Medline or even Google Scholar. In general, Mali et al. [60] identify three methodological approaches to studying scientific co-authorship networks:

- (i) basic analysis of network properties using temporal data (usually in the form of a time-series of snapshots), (ii) deterministic approaches to the analysis of scientific co-authorship networks, and (iii) statistical modeling of network dynamics

Although the identification of the three methodological approaches, an important body of the literature involving co-authorship network studies use the basic analysis of network properties such as network degree, density, path, path length, shortest path and the global clustering coefficient. Many scientific collaboration network studies have adopted this basic analysis methodological approach to scientific co-authorship investigation. In the next paragraphs, we present and discuss the purpose, methods and the results of some of those studies.

## *Literature Review*

---

Newman [19] investigated scientific network collaboration in biomedical research, physics and computer science. In this study, the author has collected data from four databases and presented distribution of collaboration networks, demonstrate the presence of clustering and highlights differences between the scientific fields under investigation. According to his findings, Newman [19] concludes on the "small words" nature of such networks in which scientists are only separated by shorter paths. In a second paper published the same year, Newman [56] provided a deeper analysis of the networks using the same data. He presented a variety of statistical properties of the networks, identified giant collaborative components and study centrality and connectedness measures. In Newman [57], the author introduced typical distances between scientists in his analyses providing therefore insights in the strength of collaboration in each network. In a last paper in the same series, the author summarizes the results of the three previous studies and showed how patterns of collaboration varied between scientists within a scientific field over time. In another study, Hou et al. [64] focuses his scientific collaboration analysis specifically on the field of scientometrics, analyzing data retrieved from the Science Citation Index (SCI) over a period expanding from 1978 to 2014. In addition to methods of Social Network Analysis (SNA), the authors have used co-occurrence analysis, cluster analysis and frequency analysis of words to describe the microstructure of the scientometrics network, reveal the major collaborative clusters and identify the center of the scientometrics collaborative network. Similarly, to Newman's publications, this paper uses basic network analysis based on network properties such as degree, centrality and betweenness metrics. Unlike Newman's studies, it also accounted for citation data. Yet another paper reported the

## *Literature Review*

---

collaborative patterns in co-authorship network in the scientific discipline of reproductive biology [65]. This study conducted a bibliometric analysis on 4,702 papers published in the field from 2003 to 2005. Although their analysis was basic, it does not make use of any network property metrics but was rather, mainly descriptive. Nevertheless, they did identify important components by applying a clustering algorithm. A similar bibliometric analysis is also reported by Toivanen and Ponomariov [66] who investigated the research collaboration patterns in the African regional systems. Their data consist in publications from African institutions from 2005 to 2009. The authors adopted an empirical clustering method based on the geographic regions within the African research context. Their research uncovers the dynamic nature of African collaborative efforts despite the lack of research capabilities, the structural weaknesses, and the uneven integration of resources. South Africa proved to be the emerging hub as it holds critical network function for collaborative research in the African context.

Some researchers have studied scientific network co-authorship across a scientific discipline in specific institutions or organizations. For example, Bellanca [67] used basic network analysis to measure interdisciplinary research by describing three co-authorship networks of researchers in Biology and chemistry departments at the University of York. They discovered fewer interdisciplinary research between biologists and chemists within the University but more interdisciplinary links between biology and mathematics, bioinformatics, biophysics and biochemistry. Their findings are potentially important for the development of strategies to promote interdisciplinary research within the University. Another study conducted in a Spanish institution analyzed collaboration between Spanish

## *Literature Review*

---

authors [68]. After retrieving 448 published papers between 1998 and 2007, the authors used basic network analysis to their network and identify group of authors as well as their relationship with others. In their future directions, the authors recommend that a dynamic time series analysis method as the next step to better understand their co-authorship network.

In some other studies, the authors focus their research on a single country, across a specific scientific discipline. Ghafouri et al. [69] propose a sociogram analysis to social co-authorship network of Iranian researchers, in an attempt to help improve research prioritization, research centers establishment, teams and new curricula in the field of emergency medicine. According to their results, they concluded on a poorly connected, loose and sparse co-authorship network in the field of emergency medicine in Iran. While their study was keyword based and might have not included all papers, they recommended the rethink of research prioritization, the establishment of new research centers more emergency medicine specialists to Iranian policy makers. Yet another Iranian study by Salamati & Soheili [70] focuses on the field of violence, assessing scientific outputs from Iranian researchers from 1972 to 2014. Using the network properties listed above in addition to other properties such as closeness, betweenness, eigenvector metrics, they identify structural holes, active authors, analyzed the structural indices of their network and evaluate the trend of published articles. One important limitation of their study was the attempt to manually standardize Iranian authors' names and the keyword based search leading to the lack of comprehensiveness of the search results. A similar study of Iranian researchers on Medical Parasitology is also reported by Sadoughi et al. [71].

## *Literature Review*

---

The study also uses basic network analysis to identify prolific researchers in the field of Medical parasitology by collecting 1048 published documents of all types in the field from 1972 to 2013. The study aims at identifying aspects of scientific collaboration to help policy makers in the medical parasitology research area. A Brazilian study reported in the literature uses the same methodological approach to generate new tools to help the Brazilian research fund to better select and prioritize research proposals [72]. The authors search scientific databases on seven neglected tropical diseases, generate and analyze co-authorship networks of each disease. Their results generate new information leading to better design and strategic planning and implementation of a research funding program. This study supports the claim that traditional criteria to fund research such as research productivity or impact factor of scientific journals are not valuable indicators for grant selection in low productivity neglected tropical diseases research areas. This Brazilian study is one of the few that focused on co-authorship network in the fields of neglected tropical diseases and the vast field of tropical infectious disease. Another study investigated the state of scientific collaboration on Chagas disease research [55]. Their goal was to promote cooperative and translational research initiatives by analyzing the scientific literature on Chagas disease in the Medline database between 1940 and 2009. On a total of 13,989 documents retrieved, the authors applied bibliometrics, social network analysis, and clustering methods to analyze the evaluation of collaboration patterns and to identify influential research groups. The results revealed a dramatic increase in research collaborations. As in Newman [19], this study concluded that the co-authorship network of Chagas disease constitutes a "small world" network characterized by a high

## *Literature Review*

---

degree of clustering. Another important remark is the scarcity of African co-authorship network studies. Our review only identified the study by Toivanen and Ponomariov [66] who focuses on research collaboration patterns in the African regional systems with less insights into specific research areas.

In their entirety, the studies reviewed above used descriptive, basic social analysis methods and bibliometrics as analysis methods. Recently, Zhang [73] proposes a complex approach to social network analysis, focusing only on link prediction, one of the network topology inference questions. Her approach focuses on the development of a computationally efficient solution based on machine learning techniques. She tested her approach on different datasets including a citation network, a co-authorship network and a protein-protein network. Quite often, these methods are not perfect since they failed to correctly tease out unreliable nodes from reliable ones, compromising the reliability of the network. However, new methodological approaches to scientific co-authorship network analysis are emerging to address those limitations. For example, Oliveira et al. [74] proposed a Bayesian approach to the analysis of such networks. Yet another limitation worth noting is that none of the studies reviewed above applied dynamic network analyses such as dynamic time series analysis or longitudinal network analysis [60].

There has been a steady increase in published sources relating to Malaria TB, and HIV/AIDS. The increasing data sources have successfully contributed to accurate estimates and in-depth understanding of the trends of the diseases that have provided grounds to validate the global funding to fight TB, HIV/AIDS, and malaria. In the year 2000, a Global Fund was set up to fight the three diseases. The fund was administered by

a non-governmental organization established by the World Health Organization (WHO). Ever since the year 2002 to 2016, these organizations have invested approximately \$19 billion in the control of the three infectious diseases. The investment saved 4.9 million lives. In the year 2008, the USA department of health under President Bush introduced the president's malaria initiative [7] that provided \$1.2 million funding in a bid to reduce malaria-related deaths in sub-Saharan Africa. The effort included the provision of top-notch malaria treatment and prevention services in the highly affected African nations.

### **1.3 General and specific Objectives**

The purpose of this research is to analyze the structure and dynamics of scientific collaborations and co-authorship in the fields of Malaria, Tuberculosis and HIV/AIDS research areas over the last 20 years in the Republic of Benin. Our results will help improve grant and research resource allocation to funding and help research organizations and national control programs to promote and encourage transdisciplinary and interdisciplinary research in the country. In addition, our results recommend new approaches and important tools to support the Beninese national control programs via better strategic planning and implementation of public health policies, research and development. We also aim at the prototyping and evaluation of an online research collaboration tool to help funding organizations promote multidisciplinary, research collaboration and co-authorship in the republic of Benin. More specifically, we address the following research questions:

## *Literature Review*

---

- What is the structure of scientific research collaboration networks in Benin over the last 20 years in Malaria, TB and HIV/AIDS research?
- Who are the most prolific authors, scientific research groups within each field?
- How have transdisciplinary and interdisciplinary research evolved over the last two decades in the Republic of Benin?
- What are the characteristics and the dynamics of the current co-authorship research collaborations in Benin in the three research areas?

This thesis fills the gap in the current literature, and reveal the role of the collaborative research in the prevailing research networks. Our research is designed to meet the following specific objectives:

1. To identify the most productive and prolific scientific research groups and authors within each research area.
2. To document and describe the structure of Malaria, TB, HIV/AIDS co-authorship networks and their characteristics, how they evolve over time in Benin over the last two decades.
3. To unravel the mechanistic phenomenon explaining the formation and trends of these networks over time.
4. To predict and recommend future research collaboration ties in Benin in the three research areas.

5. To develop and evaluate a scientific collaboration and co-authorship tool to disseminate the findings of this research and help design future policies

## 1.4 Gap in the literature

Despite the increasing financing effort and increasing number of published reports, the literature does not provide sufficient data regarding co-authorship networks of scientific research collaborations and their dynamics in the fields of malaria and TB and HIV/AIDS research in Africa, and particularly in Benin. Knowing such information is crucial to consolidate the progress made at controlling those diseases, support cooperative and translational research initiatives [55]. Lack of such information makes it difficult for policy makers to sustain the important progress made to reduce morbidity and mortality rates.

# Chapter 2

## Methodology

### 2.1 Overview

To attain objective 1, our methodological approach consists in performing descriptive analysis of the network data and a bibliometric analysis of each co-authorship network, following the methodology used by Newman et al. [19] and Ghafouri et al. [69]. For objective 2, we use clustering methods, and shortest paths algorithms as explained by Newman [56, 57]. Next, we apply mathematical modeling to attain objective 3. Regarding objective 4, we apply advanced statistical modeling including dynamic or longitudinal network analysis methods as recommended by Mali et al. [60]. Finally, we predict future research collaboration ties using the best performing statistical model.

## 2.2 Data Collection

Our research utilized secondary data collection techniques using the systematic literature search. The search was conducted on papers published by Beninese authors, or papers published about Malaria, TB and HIV/AIDS involving authors affiliated to Beninese research institutions. The documents of interest were the ones that provide evidence based information regarding scientific networks and collaborative research in evaluating the challenge presented by the three infectious diseases in Benin. All published documents under consideration included at least one Author with affiliation to a Beninese research institution. No restriction was placed upon the document types. Peer-reviewed articles were selected from systematic bibliographic search on papers indexed in Thompson's Institute for Scientific Information Web of Science (WOS) (formerly known as the Web of Knowledge). Full citations information containing the authors' names, their institutional affiliations, the year of publication, as well as the number of times the document was cited were recorded as a bibliographic corpus in text format. The searches were restricted to only research published between January 1, 1996 and December 31, 2016.

### 2.2.1 Text Mining and Network generation

From each bibliographic text files, we constructed a corpus of the published documents using Tethne v0.8 [75], a python library for parsing bibliographic data. Using NetworkX [76], another python library, we generated undirected multigraph co-authorship networks containing parallel edges. Vertices were defined by several attributes including name,

affiliation, city, country, number of publication and total number of times cited. Edges too, had attributes associated with them such as a unique identifier, the number of times a pair of authors was cited and the number of publications of a pair of authors. We normalized and disambiguated the information collected such as researchers' names, research center denominations, and any other information that appeared ambiguous.

### **2.2.2 Author Name Disambiguation**

One common challenge in collecting bibliometric data is the matching problem. Multiple names can refer to the same author. A well-known approach to solving this issue is termed as Author Name Disambiguation (AND). While many AND methods have been reported in the literature [77, 78], we performed a fuzzy matching machine learning technique of AND. We used Dedupe, a python library to disambiguate authors' names and assign a unique identification number to each author. We manually annotated 10% of the names and then trained the algorithm to automatically disambiguate the remaining of the entries. Dedupe is interactive and adjusts further annotations as the disambiguation process evolves. Dedupe is based on the work of Bilenko [79] and has been developed by Gregg Forest and Derek Eder. For more information on Dedupe, we refer the reader to the authors' Github repository available at <https://github.com/datamade/dedupe>. We evaluated our AND fuzzy matching machine learning method by computing Precision and recall metrics.

## 2.3 Descriptive Data Analysis

Using **igraph**, a network analysis package developed in R, we computed the following vertex and dyadic centrality measures:

- Degree of the vertices in the network defined as the number of ties to a given author.

After converting the multigraph network in a weighted graph where weights are the number of authorship between two authors, the strength of the vertices was also computed.

- Betweenness: it is the number of shortest paths between alters that go through a particular author. It relates to the perspective that importance relates to where a vertex is located with respect to the paths in the network graph. According to Freeman [80], it is defined as:

$$c_B(v) = \frac{\sigma(s, t|v)}{\sum_{s \neq t \neq v \in V} \sigma(s, t)} \quad (2.1)$$

where  $\sigma(s, t|v)$  is the total number of shortest paths between vertices  $s$  and  $t$  that pass through vertex  $v$ , and  $\sigma(s, t)$  is the total number of shortest paths between  $s$  and  $t$  regardless of whether or not they pass through  $v$ .

- Closeness: the number of steps required for a particular author to access every other authors in the network. It captures the notion that a vertex is central if it is close to many other vertices. Considering a network  $G = (V, E)$  where  $V$  is the set of vertices and  $E$ , the set of edges, the closeness centrality  $c_{Cl}(v)$  of a vertex  $v$  is

defined as:

$$c_{Cl}(v) = \frac{1}{\sum_{u \in V} dist(v, u)} \quad (2.2)$$

where  $dist(v, u)$  is defined as the geodesic distance between the vertices  $u, v \in V$ .

- Eigenvectors: degree to which an author is connected to other well connected authors in the network. It seeks to capture the idea that the more central the neighbors of a vertex are, the more central that vertex itself is. According to Bonacich [81] and Katz [82], the Eigenvector centrality measure is defined as:

$$c_{E_i}(v) = \alpha \sum_{\{u, v\} \in E} c_{E_i}(u) \quad (2.3)$$

Where the vector  $\mathbf{c}_{E_i} = (c_{E_i}(1), \dots, c_{E_i}(N_v))^T$  is the solution to the eigenvalue problem  $\mathbf{A}\mathbf{c}_{E_i} = \alpha^{-1}\mathbf{c}_{E_i}$ , where  $\mathbf{A}$  is the adjacency matrix for the network  $G$ . According to Bonacich [81], an optimal choice of  $\alpha^{-1}$  is the largest eigenvalue of  $\mathbf{A}$

- Brokerage: degree to which an actor occupies a brokerage position across all pairs of alters.
- Edge betweenness centrality extends from the notion of vertex centrality. It reflects the number of shortest paths traversing that edge. This centrality measure was computed to assess which co-authorship collaboration ties are important for the flow of information.

### 2.3.1 Characterizing Network cohesion

The extent to which subsets of authors are cohesive with respect to their relation in the co-authorship network was assessed through network cohesion. Specifically, we determined if collaborators (co-authors) of a given author tend to collaborate as well, and what subset of collaborating authors tend to be more productive in the network. While there are many techniques to determine network cohesion, we chose local triads and global giant components. In addition, we conducted cliques detection and clustering or communities detection on each network:

- Cliques: According to Kolaczyk and Csárdi [83], cliques are defined as complete subgraphs such that all vertices within the subset are connected by edges. We computed the number of maximal cliques and assessed their size.
- Density: Defined as the frequency of realized edges relative to potential edges, the density of a subgraph  $H$  in  $G$  provides a measure of how close  $H$  is to be a clique in  $G$ . Density values varie between 0 and 1:

$$den(H) = \frac{|E_H|}{|V_H|(V_H - 1)/2} \quad (2.4)$$

- Relative frequency: we assess the relative frequency of  $G$  by computing its transitivity defined as:

$$cl_T = \frac{3\tau_\Delta(G)}{\tau_3(G)} \quad (2.5)$$

where  $\tau_\Delta(G)$  is the number of triangles in  $G$ , and  $\tau_3(G)$  is the number of connected triples (sometimes referred to as 2-star). This measure is also referred to as the fraction of transitive triples. It represents a measure of global clustering of  $G$  summarizing the relative frequency with which connected triples close to form triangles [83].

- **Connectivity, Cuts, and Flows:** We investigated the concepts of vertex and edge cuts derived from the concept of vertex (edge) connectivity. The vertex (edge) connectivity of a graph  $G$  is the largest integer such that  $G$  is k-vertex- (edge-) connected [83]. These measures helped assess the most important authors for information flow and the long-term sustainability of each network. Since co-authorship networks are undirected graphs, the concept of weak and strong connectivity was irrelevant. A graph  $G$  is said to be connected if every vertex in  $G$  is reachable from every other vertex. Usually, one of the connected components can dominate the others, hence the concept of giant component.
- **Graph Partitioning:** Regularly framed as community detection problem, we applied graph partitioning to find subsets of vertices that demonstrate a 'cohesiveness' with respect to their underlying relational patterns. Cohesive subsets of vertices generally are well connected among themselves and are well separated from the other vertices in the graph. Two established methods of graph partitioning are Hierarchical clustering (agglomerative vs divisive) and Spectral clustering [83]. In our research, we applied agglomerative Hierarchical Clustering to the co-authorship networks.

## 2.4 Modeling of Network Data

The purposes of network graph modeling are to test significance of the characteristics of observed network graphs, and to study proposed mechanisms of real-world networks such as degree distributions and small-world effects [83]. A model for a network graph is a collection of possible graphs  $\mathcal{G}$  with a probability distribution  $\mathbb{P}_\theta$  defined as:

$$\{\mathbb{P}_\theta(G), G \in \mathcal{G} : \theta \in \Theta\} \quad (2.6)$$

where  $\theta$  is a vector of parameters ranging over values in  $\Theta$ .

Given an observed co-authorship network graph  $G^{obs}$  and some structural characteristics  $\eta(\cdot)$ , our goal is to assess if  $\eta(G^{obs})$  is unusual. We then compare  $\eta(G^{obs})$  to collection of values  $\{\eta(G) : G \in \mathcal{G}\}$ . If  $\eta(G^{obs})$  is too extreme with respect to this collection, then we have enough evidence to assert that  $\eta(G^{obs})$  is not a uniform draw from  $\mathcal{G}$ .

Given the computationally expensive calculations involved in modeling in general, and the expected large size of our network, we parallelized all processing.

### 2.4.1 Mathematical Modeling

We applied different mathematical models for network graphs including:

- Classical Random Graph Models: First established by Erdős and Rényi [84–86], it specifies a collection of graphs  $\mathcal{G}$  with a uniform probability  $\mathbb{P}(\cdot)$  over  $\mathcal{G}$ . A variant

of this model called the Bernoulli Random Graph Model was also defined by Gilbert [87].

- Generalized Random Graph Models: These models emanated from the generalization of Erdős and Rényi's formulation, defining a collection of graphs  $\mathcal{G}$  with prespecified degree sequence.
- Mechanistic Network Graph Models: These models mimic real-world phenomena and include Small-World Models commonly referred to as "six-degree separation". It was introduced by Watts and Strogatz [88] and have since received a lot of interests in the existing literature especially in Neuroscience. Small-world networks usually exhibit high levels of clustering and small distances between vertices. Classical models are not fit to better represent such behaviors since they usually display low levels of clustering and small distance between vertices. Examples of known small-world networks include the network of connected proteins or the transcriptional networks of genes [89]. A variant of Small-World models is the Preferential Attachment Models defined based on the popular principle of "the rich get richer". Preferential attachment models gained fascination after the work of Barabási and Albert who studied the growth of the World Wide Web [90]. Examples of Preferential Attachment networks include that of World Wide Web and the scientific citation network [91, 92]. An important characteristic of these models is that as time tend to infinity, there degree distribution tends to follow a power law.

## 2.4.2 Statistical Modeling

Although mathematical models tend to be simpler than statistical models, the latter allow model fitting and assessment. In order to unravel the mechanistic phenomenon explaining the structure of our co-authorship networks, we fit a variety of statistical models to the network data. For each model, we accounted for an important social network principle referred to as homophily which is defined in our network as the tendency of similar authors to collaborate. Another very important social network principle we also accounted for, is the one of structural equivalence which is the similarity of network positions on the formation of collaboration ties in a given network. We hypothesized that tie formation in each co-authorship network (i) is dependent on certain authors' characteristics, (ii) is dependent on the concept of distance in latent space, and (iii) collaboration type and/or membership to a certain research community or class determines collaboration tie formation. We applied various static and temporal statistical network models to each co-authorship network to verify our hypotheses. The purpose of this approach to network modeling is to unveil structural patterns driving collaboration tie formation in each co-authorship network in Benin.

### 2.4.2.1 Stochastic Block Model

Blockmodeling is a statistical method to identify, in a given network, clusters or classes of authors that share structural characteristics [93, 94]. Each such cluster forms a position. The units within a cluster have the same or similar connection patterns. Given a graph

$G = (V, E)$  and its adjacency matrix  $\mathbf{Y}$ , for two distinct nodes  $i, j \in V$ , the block model defined by Kolaczyk and Csárdi [83], specifies that each element  $Y_{ij}$  of  $\mathbf{Y}$  is conditional on the class label  $q$  and  $r$  of the vertices  $i$  and  $j$ . The model has the form:

$$Pr(\mathbf{Y} = \mathbf{y}) = \left( \frac{1}{\kappa} \right) \exp \left\{ \sum_{q,r} \theta_{qr} L_{qr}(\mathbf{y}) \right\} \quad (2.7)$$

where  $L_{qr}$  is the number of edges in the observed graph  $\mathbf{y}$  connecting vertices of classes  $q$  and  $r$ . Stochastic block model (SBM) originated from the ideas that equivalent units can be grouped together. There are three definitions of equivalences which are structural, automorphic and regular [60]. In practice, the differences in types of equivalence tend to blur when stochastic block modeling is applied to real networks. We used SBM to both model our observed network but also as a model based clustering technique. After fitting the SBM, we extract the posterior probability of class membership and determined the class membership of each vertex class assignment based on the maximum a posteriori criterion. Class membership was added to the network as an additional nodal attribute. The R package **mixer** was used to fit the SBM. **Mixer** used the Integration Classification Likelihood (ICL) criterion to select the number of classes fit to the observed network.

#### 2.4.2.2 Exponential Random Graph Model

Also referred to as p\* models, Exponential Random Graph Models (ERGMs) are probability models for network designed in analogy to Generalized Linear Models (GLMs) [83]. ERGM have gain increasing interests especially in modeling social networks. Robins et

## *Methodology*

---

al. [95] provides a nice introduction to ERGM as well as a general framework for ERGM creation which we closely followed here. We used ERGM to investigate how local processes affect collaboration tie formation between authors in our network. We modeled the network ties, the dependent variable as a function of nodal and dyadic attributes (covariates) such as the number of times an author was cited, the number of publications, the number of collaborators, the collaboration type as well as its community membership as determined by the SBM.

Given a random graph  $G = (V, E)$ , for two distinct nodes  $i, j \in V$ , we define a random binary variable  $Y_{ij}$  such that  $Y_{ij} = 1$  if there is an edge  $e \in E$  between  $i$  and  $j$ , and  $Y_{ij} = 0$  otherwise. Since co-authorship networks are by definition undirected networks,  $Y_{ij} = Y_{ji}$  and the matrix  $\mathbf{Y} = [Y_{ij}]$  represents the random adjacency matrix for  $G$ . The general formulation of ERGM is therefore:

$$Pr(\mathbf{Y} = \mathbf{y}) = \left( \frac{1}{\kappa} \right) \exp \left\{ \sum_H \theta_H g_H(\mathbf{y}) \right\} \quad (2.8)$$

where each  $H$  is a configuration, a set of possible edges among a subset of the vertices in  $G$  and  $g_H(\mathbf{y}) = \prod_{y_{ij} \in H} y_{ij}$  is the network statistic corresponding to the configuration  $H$ ;  $g_H(\mathbf{y}) = 1$  if the configuration is observed in the network  $\mathbf{y}$ , and is 0 otherwise.  $\theta_H$  is the parameter corresponding to the configuration  $H$  (and is non-zero only if all pairs of variables in  $H$  are assumed to be conditionally dependent);  $\kappa$  is a normalization constant. In order to obtain the best model, several models containing nodal, dyadic and structural terms were fit to the observed network data. The first model we fit is a naive model containing only the ERGM "edge" term. This model is nothing but the Bernoulli random

graph model [84]. We then fit another model containing only nodal and/or dyadic terms. Third, we fit a structural model containing only high-order terms representing network statistics such as triangles, k-stars, geometrically weighted edge-wise shared partner distribution and many more [83, 95]. Ideally, we expect the best model to contain nodal and dyadic covariates as well as high order ERGM terms. Model log-likelihood, the Akaike's Information (AIC) and the Bayesian Information (BIC) criteria were used to select the best model. After checking for model diagnostics whenever necessary, we finally evaluated the best model (lowest AIC or BIC and highest likelihood) by assessing its goodness-of-fit to the observed network. We expect each model to converge within a maximum of 1,000 iterations. The R package **ergm** was used to fit the models.

#### 2.4.2.3 Temporal Exponential Random Graph Model

The Temporal Exponential Random Graph Model (TERGM) is an extension of the ERGM described in section 2.4.2.2 proposed by Hanneke, Fu, and Xing [96] from the work of Robins and Pattison [97]. The TERGM was designed with the idea of accounting for inter-temporal dependence in longitudinally collected network data. TERGM was applied to each co-authorship network following from the work of Leifeld, Cranmer, and Desmarais [98]. For a full description of the TERGM, we refer the reader to Leifeld, Cranmer, and Desmarais [98].

Each network is subset in different temporal snapshots. In general, when the temporal network is overly dense or sparse early on or in later time periods, the TERGM tends to fit different time spans differently [98]. To avoid such an issue, the cumulative network was

## *Methodology*

---

subset in a certain way that balanced the number of edges across the years. This strategy improved the robustness and convergence of our models. We modeled the network ties, the dependent variable as a function of nodal, dyadic variables, and dyadic stability and delay reciprocity memory terms. To check whether there is a linear trend in collaboration tie formation, we also included a linear time covariate in the model. We accounted for network structural predictors and homophily on the type of collaboration. Model log-likelihood, the Akaike's Information (AIC) and the Bayesian Information (BIC) criteria were used to select the best model (final model) corresponding to the lowest AIC or BIC, and highest log-likelihood.

To evaluate the extent to which the final model captures the endogenous properties and processes of the observed network, we checked for model diagnostics, assessing the within-sample and out-of-sample goodness-of-fit. For the out-of-sample goodness-of-fit, we estimated the model on the first network snapshots leaving out the last network snapshot in the series. We simulated 1000 networks from the model and assessed how the simulated network predicted the left out network. As described by Desmarais and Cranmer [99], we also provided a micro-interpretation of the final TERGM.

All models were fit using the Markov Chain Monte Carlo Maximum Likelihood Estimation (MCMC-MLE) for TERGMs implemented in the **btergm** R package.

### **2.4.2.4 Latent Network Model**

Designed in analogy to Mixed Models, Latent Network Models (LNM) allow the incorporation of latent or unobserved variables in network modeling. These models specifically

account for structural equivalence, to model hidden factors or information not available in the network. Kolaczyk and Csárdi [83] provide a formulation of LNM. Given the adjacency matrix  $\mathbf{Y}$  of a graph  $G = (V, E)$ , for each element  $Y_{ij}$  of  $\mathbf{Y}$ , the latent variable model is of the form:

$$Y_{ij} = h(\theta, z_i, z_j, \epsilon_{ij}) \quad (2.9)$$

where  $\theta$  is a constant, the  $\epsilon_{ij}$  are independent and identically distributed pair-specific effects, and  $h$  is a symmetric function. The model assumes that each vertex  $i \in V$  has a latent variable  $z_i$ . Considering observed covariates  $\mathbf{Z}$ , the probability of forming an edge between two nodes  $i$  and  $j$  ( $i, j \in V$ ) is independent of all other vertex pairs given values of latent variables, and is defined as:

$$Pr(\mathbf{Y}|\mathbf{Z}, \theta) = \prod_{i \neq j} Pr(Y_{ij}|z_i, z_j, \theta) \quad (2.10)$$

We specified latent effects according to an approach suggested by Hoff [100] and based upon the principles of eigen-analysis. The R package **eigenmodel** developped by Hoff [101] was used to fit the LNM to the observed network. We fit LNM with both no pair-specific and pair-specific covariates such as the type of collaboration and community assignment from the SBM. The rationale of fitting the pair-specific models with those two variables is supported by our third hypothesis which states that collaboration ties in each co-authorship network are driven by homophily in terms of community membership and/or collaboration type. We also fit other pair-specific covariates model using nodal and dyadic covariates. We visualized and compared the co-authorship network using

### *Methodology*

---

3 dimensional layouts determined according to the inferred latent eigenvectors in each model. Finally, we used a 5-fold cross-validation method to assess the goodness-of-fit of each model which we compared using ROC curves via the R package **ROCR**.

# **Chapter 3**

## **Results: The Malaria Co-authorship Network**

### **3.1 Data**

The data collection was carried on papers indexed in Thompson's Institute for Scientific Information Web Of Science (formerly known as the Web of Knowledge). The search was conducted using combinations of Malaria related MeSH terms including "malaria", "Anopheles", "Plasmodium" and "vector". We restricted the search to the period from 1996 to 2016 and to "Benin" for country. We further screened the papers in order to only select those published by Beninese authors, or papers published on Malaria involving at least one author affiliated to a Beninese research institution. No restriction was placed upon the document types. We first started querying with each term independently, we

## *Results: The Malaria Co-authorship Network*

---

then combined the other terms so the query return the maximum number of results. The Full citations information containing the authors' names, their institutional affiliations, the year of publication, as well as the number of times the document was cited were recorded as a bibliographic corpus in text format. After a second screening only research that have met the above listed inclusion criteria and that were published between January 1, 1996 and December 31, 2016 were selected in this study.

The final query set (Table 3.1) returned 685 records. After screening, 424 documents met the selection criteria. On average, there was 10.67 authors per published document.

After the Author Name Disambiguation, we identified 1792 unique authors with a precision of 99.87% and a recall of 95.46%. The generated multigraph co-authorship network therefore contained 1792 vertices (authors) and 116,388 parallel edges (collaborations).

TABLE 3.1. Malaria Bibliographic Search Queries.

Set	Queries	Results
#1	TOPIC: (malaria) OR TOPIC: (mosquito), Refined by: COUNTRIES/TERRITORIES: (BENIN)	513
#2	TOPIC: (malaria) OR TOPIC: (mosquito) OR TOPIC: (anopheles), Refined by: COUNTRIES/TERRITORIES: (BENIN)	529
#3	TOPIC: (malaria) OR TOPIC: (mosquito) OR TOPIC: (anopheles) OR TOPIC: (plasmodium) OR TOPIC: (bednet), Refined by: COUNTRIES/TERRITORIES: (BENIN)	544
#4	TOPIC: (malaria) OR TOPIC: (mosquito) OR TOPIC: (anopheles) OR TOPIC: (plasmodium) OR TOPIC: (net) OR TOPIC: (vector), Refined by: COUNTRIES/TERRITORIES: (BENIN)	685
Final Set	#1 OR #2 OR #3 OR #4	685

## 3.2 Descriptive Data Analysis

The degrees of the multigraph network range between 1 and 1338 with an average degree distribution of 106.46. We noted in addition, a substantial number of vertices with low degrees (Fig. 3.1). There was also a non-trivial number of vertices with higher order of degree magnitudes. A log scale distribution of the degrees demonstrate that the vertex degrees tend to follow a heavy-tail distribution.

After we convert the multigraph network in a weighted graph, it results in a simple graph of 1792 vertices and 95,787 weighted edges. Mean Closeness centrality ranges between  $3.118 \times 10^{-7}$  and  $5.152 \times 10^{-6}$  with a median of  $5.112 \times 10^{-6}$ . This measure suggests a highly right-skewed distribution. Betweenness measures range between 0 and 245600 with a median of 1985. A network visualization with the vertices' size proportional to betweenness centrality measures clearly reveals the presence of broker authors (Table 3.2). The median Eigenvectors median is 0.005, its mean is estimated at 0.09. Eigenvectors measures reveal the presence of multiple cluttered authors suggesting the presence of closed collaboration groups. Table 3.2 presents a list of the 10 authors with the highest Eigenvectors values.

The computation of edge betweenness identifies co-authorship collaborations that are important for the flow of information. In Table 3.2, We present the top 10 most important collaborations for the flow of information in the Malaria Co-authorship network in Benin.

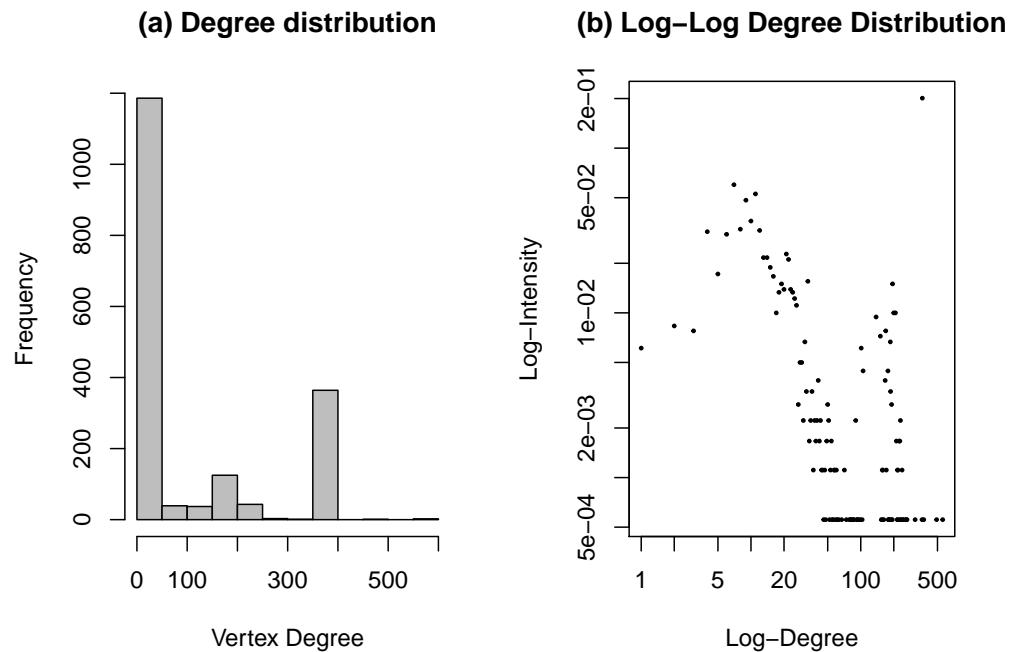


FIGURE 3.1: Degree distribution of the Malaria Co-authorship network

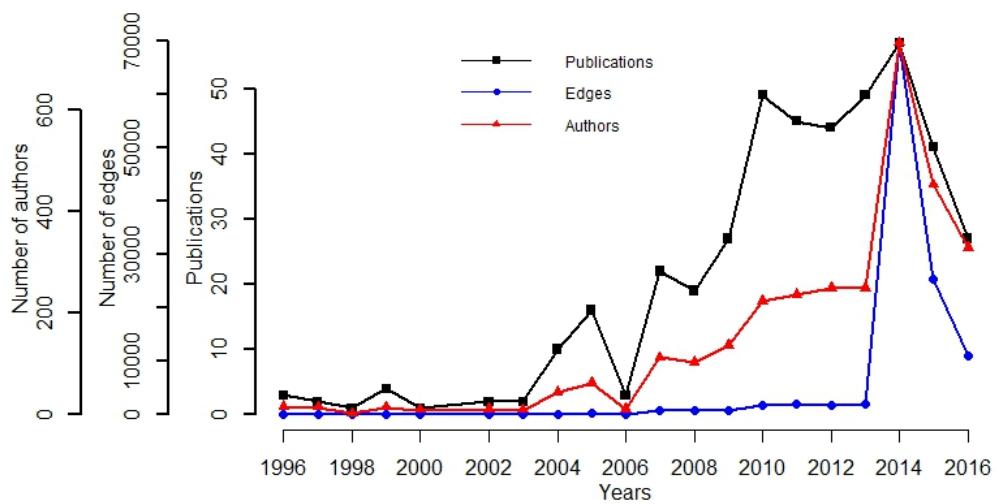


FIGURE 3.2: Evolution of the published Malaria related documents, authors and collaborations from January 1996 to December 2016

*Results: The Malaria Co-authorship Network*

---

TABLE 3.2. List of the most important authors and collaborations in the Malaria Co-authorship network

<b>Top 10 Brokers</b>
MASSOUGBODJI ACHILLE
HAY SIMON I
KAREMA CORINE
SANNI AMBALIOU
KENGNE ANDRE PASCAL
AKOGBETO MARTIN
NDAM NICASE TUIKUE
MALIK ELFATIH M
DABIRE K ROCH
DELORON PHILIPPE
<b>Top 10 most connected authors (Top 10 network hubs)</b>
MASSOUGBODJI ACHILLE
KAREMA CORINE
GONZALEZ RAQUEL
MENENDEZ CLARA
DALESSANDRO UMBERTO
OGUTU BERNHARDS R
FAUCHER JEANFRANCOIS
BASSAT QUIQUE
MARTENSSON ANDREAS
HAY SIMON I
<b>Top 10 most important edges for information flow</b>
DABIRE K ROCH – KENGNE ANDRE PASCAL
BALDET THIERRY – KENGNE ANDRE PASCAL
AKOGBETO MARTIN – MALIK ELFATIH M
AVLESSI FELICIEN – MOUDACHIROU MANSOUROU
AKOGBETO MARTIN – AVLESSI FELICIEN
MASSOUGBODJI ACHILLE – RAHIMY MOHAMED CHERIF
DIABATE ADOULAYE – KENGNE ANDRE PASCAL
GARCIA ANDRE – SANNI AMBALIOU
KAREMA CORINE – MALIK ELFATIH M
HAY SIMON I – MALIK ELFATIH M
<b>Weak articulation points</b>
NOEL VALERIE
DJOGBENOU LUC
ZOHOUN I
SANNI AMBALIOU
EDORH ALEODJRODO PATRICK
ALLABI AUREL
HOUNKONNOU MAHOUTON NORBERT
FAYOMI BENJAMIN
KINDEGAZARD DOROTHEE A
DJOUAKA ROUSSEAU
RAHIMY MOHAMED CHERIF
BALDET THIERRY
DOSSOUGBETE L
GARCIA ANDRE
MASSOUGBODJI ACHILLE
AKOGBETO MARTIN

### 3.2.1 Network Cohesion

A total of 365 maximal cliques are identified in the network among which 9 cliques of size 2, 14 cliques of size 3, 155 cliques of size 8, and 142 cliques of size 7. Larger maximal

## *Results: The Malaria Co-authorship Network*

---

cliques sizes range from 102 authors to 365 authors and are all found once across the network.

The malaria co-authorship network has a density of 0.0596 and a transitivity of 0.965 indicating that 96.5% of the connected triples in the network are close to form triangles.

The transitivity metrics is a measure of the global clustering of the network.

The network is not connected and a census of all the connected components within the network reveals the existence of a giant component that dominates all the other connected components. This giant component includes 94% (1686 vertices) of all the vertices in the network with none of the other components alone carrying less than 1% of the vertices in the network (Fig. 3.3).

The assessment of information flow in the network via cut vertices reveal the existence of 16 authors as the most vulnerable vertices in the network. Table 3.2 lists the authors that constitute the weak articulation points in the malaria co-authorship network. Cut vertices are crucial to the sustainability of networks [83].

The agglomerative hierarchical clustering method identifies 23 research communities (or clusters) in the network. Sizes of the clusters range between 2 and 570 with large research communities containing between 202 and 569 authors. Medium size research communities contain between 10 and 62 authors. Only 7 out of the 23 research communities identified are part of the giant component. Figure 3.3 displays the giant component of the network with each different colors representing each of the 7 research communities.



FIGURE 3.3: Malaria Co-authorship network – Main component.  
Authors (vertices) of the same color belong to the same research community or cluster

### 3.3 Modeling

#### 3.3.1 Mathematical Modeling

The hierarchical clustering method of community detection algorithm has identified 23 different clusters/communities in the co-authorship network out of which 7 form a giant component. One of the question of interest in this section is whether the number of communities detected is expected or not. We performed 1,000 Monte Carlo based simulations to test the significance of this observed characteristics on the malaria co-authorship network. Figure 3.4 clearly demonstrates that the number of communities detected is unusual from the perspective of both Classical random graphs and generalized random graphs ( $p\text{-value} < 0.0001$ ). From the Classical random graph model, the expected number of communities is 3.934 (95%CI: 3.90 – 3.97). Similarly, the expected number of communities from the generalized random graph model is 7.501 (95%CI: 7.39 – 7.61).

Figure 3.5 displays the number of detected research communities using the Barabási-Albert’s preferential attachment and the Watts-Strogatz models. Surprisingly enough, the observed number of communities is also extreme per both models ( $p\text{-value} < 0.0001$ ). The expected number from the Watts-Strogatz model simulations is 3.056 (95%CI: 3.04 – 3.07) and 45.569 (95%CI: 45.42 – 45.72) from the Barabási-Albert model simulations.

We also compared the clustering coefficient and the average shortest-path length. The observed clustering coefficient is 0.9645. Surprisingly, there is substantially more clustering in our malaria co-authorship network than expected from all 4 mathematical models

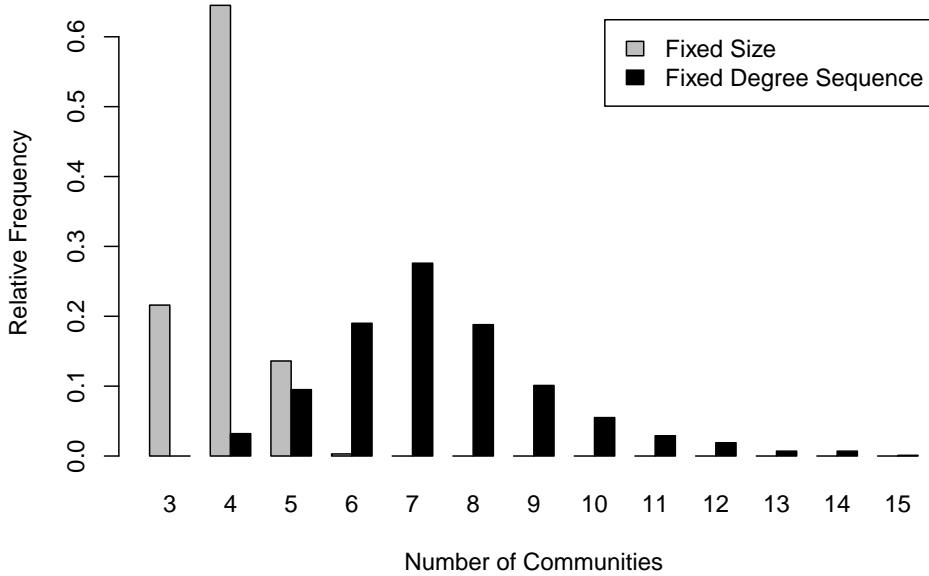


FIGURE 3.4: Monte-Carlo simulations: Number of detected communities by the random graph models

( $p\text{-value} < 0.0001$ ). The expected clustering coefficient is 0.0596 (95%CI: 0.05963068 – 0.05964648) and 0.4334 (95%CI: 0.4333912 – 0.4334522) respectively for the classic random graph and the generalized random graph models.

Similarly, The Watts-Strogatz Small World model expected clustering is 0.7464 (95%CI: 0.7464326 – 0.7464356).

We observed an average shortest-path length of 2.99 in the malaria co-authorship network. This observed shortest-path length is significantly larger than what is expected from the random graph models ( $p\text{-value} < 0.0001$ ) and significantly lower than what is expected from Watts-Strogatz small world model and the Barabási-Albert preferential attachment model ( $p\text{-value} < 0.0001$ ).

The average shortest-path length is 1.94 (95%CI: 1.941955 – 1.941960) and 2.26 (95%CI:

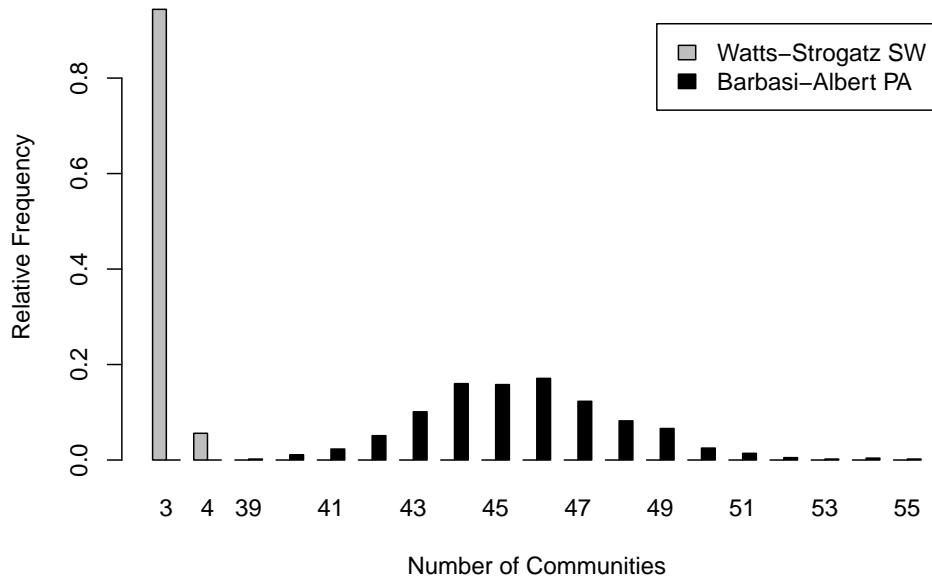


FIGURE 3.5: Monte-Carlo simulations: Number of detected communities by the Watts-Strogatz and the Barabási-Albert models

2.259468 2.259586) respectively for the classic random graph and the generalized random graph models.

For the Watts-Strogatz small world and the Barabási-Albert models, the average shortest-path length is respectively 3.83 (95%CI: 3.81 – 3.86) and 9.17 (95%CI: 9.14 – 9.21).

All simulations were also performed on the giant component of the network and led to similar outcomes.

### 3.3.2 Statistical Modeling

#### 3.3.2.1 Stochastic Block Model

The ICL plot on figure 3.6 shows that the malaria co-authorship network has been fit with 39 classes by the SBM with a degree of latitude of 30 to 39 classes being reasonable. The degree distribution of the fitted SBM (blue curve) provides a decent description of the observed distribution (yellow histogram). In the inter/intra class probabilities network, the vertices correspond to the 39 classes detected by the SBM. The vertex sizes are proportional to the number of authors assigned to each class. Each vertex is further broken down in a pie chart with each portion reflecting the relative proportion of the types of collaboration. Yellow represents the proportion of authors of international affiliations, orange represents regional authors who are affiliated with African institutions other than Beninese institutions, and green for authors affiliated to Beninese research institutions. In general, we observe a dominance of international and regional researchers over national researchers across all detected clusters.

A close look at the reorganized adjacency matrix, reveals the presence of 4 larger classes (classes number 2, 4, 10 and 27) and 35 other classes of smaller sizes. One of the larger class (class 27) displays a tendency of its members to only establish collaboration ties between themselves. This class seems to have the characteristics of a clique. Examination of the distribution of each class by their type of collaboration (Figure 3.7) indicates that this class of authors (class 27) is primarily made of international contributors to the malaria research effort in Benin. Although members of this class seem to have rare

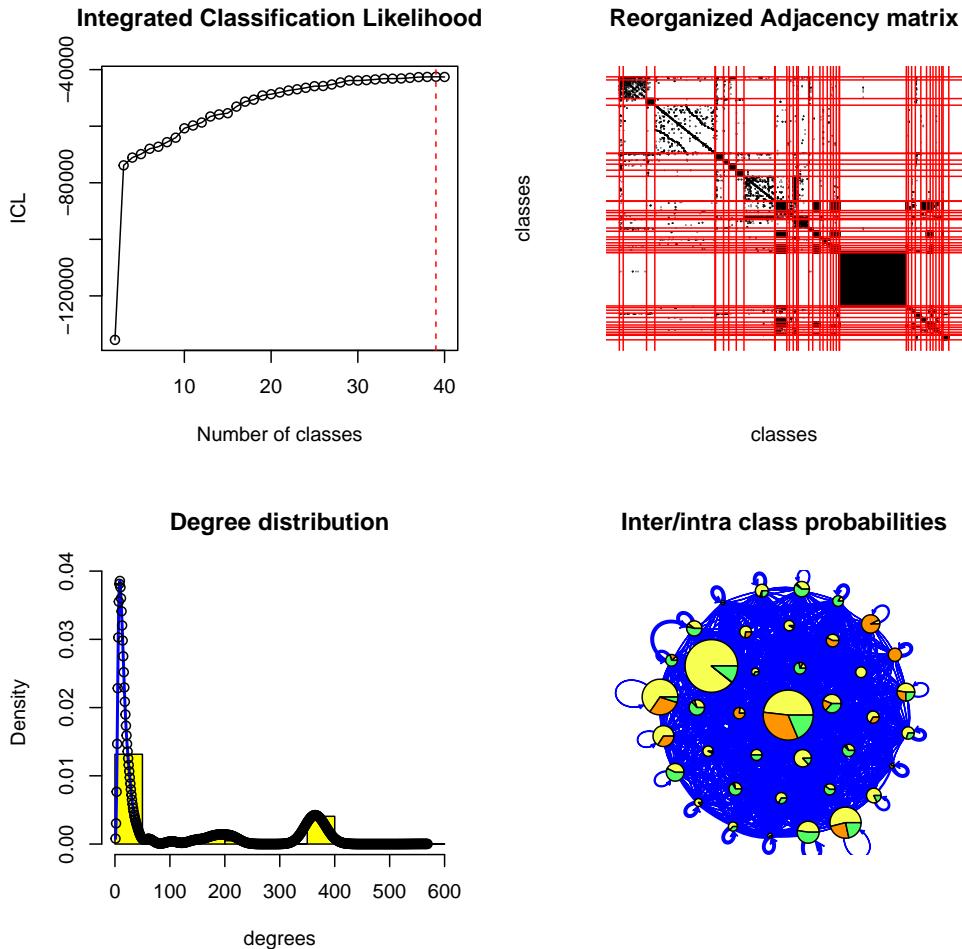


FIGURE 3.6: Summary of the goodness-of-fit of the SBM analysis on the Malaria co-authorship network.

collaboration ties with members of other classes, we also notice the presence of very few broker authors as national liaisons between this class 27 and another larger class (class 2). Though, it also appears in the other three larger classes that the authors tend to primarily collaborate within their respective classes, they also tend to collaborate with authors of other classes.

Figure 3.7 also shows that the co-authorship malaria network in Benin is dominated by international researchers with national contributors be unevenly distributed across the

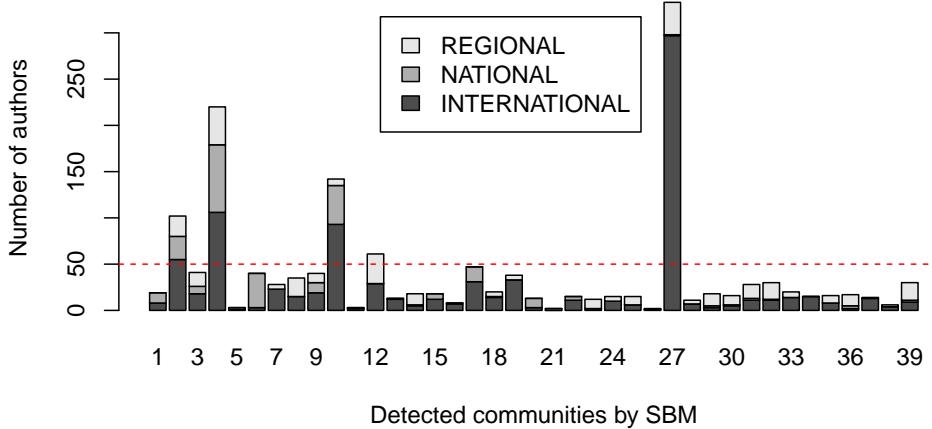


FIGURE 3.7: Distribution of national, international and regional authors by communities detected by the SBM in the Malaria network.

detected research communities. In order to better explain the inter/intra class interactions, we highlight in figure 3.8, the main classes driving the structure of the network. We present the results from the SBM on the classes with 50 authors or more. This reorganization clearly confirmed the presence of a clique of mainly international contributors who tend to collaborate rarely outside their class. The larger size here (Figure 3.8) is very diverse and contains all regional contributors to the malaria research effort. The presence of 3 smaller cliques which collaborate intensively between themselves is worth noting as well (See inter/intra class probabilities network on figure 3.8).

### 3.3.2.2 Exponential Random Graph Model

Table 3.3 summarizes the results of the different models we fit to the observed network. Model 1 is analogous to the null model in a typical General Linear Model (GLM). The

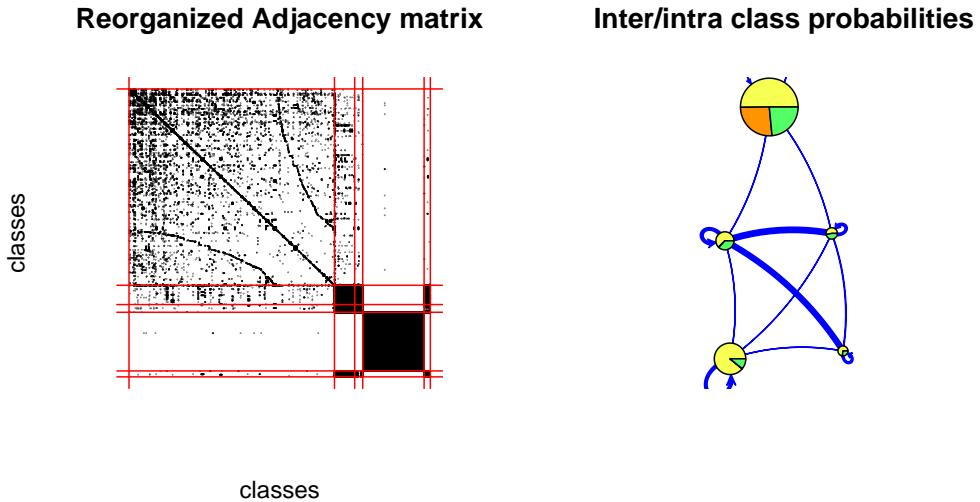


FIGURE 3.8: Summary of the goodness-of-fit of the SBM analysis highlighting interactions between the top 5 larger classes of the Malaria co-authorship network.

probability of any two authors establishing a collaboration tie is therefore expressed as the inverse logit of the edge coefficient. The inverse logit of a coefficient  $x$  is defined as  $logit^{-1}(x) = 1/(1 + exp(-x))$ . The conditional log-odds for a collaboration between authors in the network is  $-2.76$ . The associated probability of any two authors establishing a collaboration tie is therefore 5.96%. Let's recall that this probability is the same as the density of the malaria co-authorship network. Since, our network is characterized by a high transitivity, we modeled the triangle ERGM term along with the edge term in model 2. We see some improvements in the model performance with a significantly positive but small triangle effect on the collaboration tie formation (Coefficient = 0.08,  $p < 0.001$ ).

In model 3, we describe the co-authorship network as a function of the number of collaborations, the number of publications, and the number of citations of authors inside the network. We also include confounding homophily on cluster assignment from the SBM

## *Results: The Malaria Co-authorship Network*

---

and on the collaboration type. Compared to models 1 and 2, model 3 has tremendously improved (See AIC and BIC in table 3.3). The edge effect has decreased (Coefficient =  $-7.98$ ,  $p < 0.001$ ) with the associated conditional probability (given all other terms in the model) equal to 0.03%. We observed a small, though positively significant effect of the number of collaborators and the number of publications on the odds of collaboration tie formation between any two authors. One unit increase in the number of collaborators increases the odds of collaboration tie by 2% while one unit increase in the number of publications increases the odds of establishing a collaboration tie by 12.75%. On the other hand, model 3 has found a very small but significant negative effect of the number of times an author was cited on the odds of collaboration tie formation. One unit increase in the number of citation of a given author was associated with 1% decrease in the odds of collaboration between two authors conditional on all the other terms in the model. It clearly appears that the process underlying the malaria co-authorship network is driven by homophily on cluster assignment or membership to a specific research community and the type of collaboration. The conditional probability of two authors collaborating adjusted by the homophily on their membership to a research community is estimated at 8.32% compared to the baseline probability of 0.03% given all other terms in model 3. Adjusted by the collaboration type, the same probability is estimated at 0.05% conditional on all other terms in the model. The overall conditional probability adjusting for all terms in model 3 is estimated at 14.06% which is a lot greater than the 5.95% estimated from model 1.

### *Results: The Malaria Co-authorship Network*

---

In model 4, we introduced factor attributes on the collaboration type in order to investigate the likelihood of researchers affiliated to Beninese institutions to establish international and regional or African collaboration ties. While model 4 slightly improved upon model 3, it displays minor changes in the coefficient of the terms it has in common with model 3. Overall, compared to researchers with international research affiliations, researchers affiliated to Beninese research institutions have 37.7% average decrease in the odds of establishing collaboration ties. On the other hand, researchers affiliated to other African research institutions have 78.6% increase in the odds of establishing a collaboration tie than researchers affiliated to international research institutions. In other words, in model 4, the probability for researchers affiliated to international institutions to establish a collaboration tie is estimated at 14.19%, that of researchers affiliated to Beninese institutions is 10.72%, and that of researchers affiliated to African institutions other than Beninese institutions is 22.79%.

None of the structural models containing high order ERGM terms, nor the models containing the dyadic attribute terms converged after the maximum of 1,000 iterations making estimates from these models unreliable. This observation justifies the reason why we do not present the results from these models in table 3.3. The inability of model containing structural terms to converge also makes it impossible for us to assess model degeneracy as recommended by Handcock et al. [103].

Figure 3.9 presents the goodness-of fit of model 4. The observed properties are depicted by the black lines. Gray lines with circles represent the 95% confidence intervals for the simulated network properties. Goodness-of-fit is asserted when the black lines lie

## *Results: The Malaria Co-authorship Network*

---

TABLE 3.3. ERGM of the co-authorship Malaria network.

	Model 1	Model 2	Model 3	Model 4
	Estimate (SE)	Estimate (SE)	Estimate (SE)	Estimate (SE)
Network structural predictor				
Intercept(edge)	-2.76 (0.00)***	-5.00 (0.01)***	-7.98 (0.02)***	-8.22 (0.02)***
Triangle	-	0.08 (0.00)***	-	-
Number of collaborations	-	-	0.02 (0.00)***	0.01 (0.00)***
Number of publications	-	-	0.12 (0.00)***	0.13 (0.00)***
Number of times cited	-	-	-0.01 (0.00)***	-0.01 (0.00)***
Homophily on cluster assignment	-	-	5.58 (0.02)***	5.68 (0.02)***
Homophily on collaboration type	-	-	0.46 (0.01)***	0.61 (0.00)***
Factor attribute effect (collaboration type)				
International	-	-	-	<i>REF</i>
National	-	-	-	-0.32 (0.02)***
Regional	-	-	-	0.58 (0.01)***
Number of iterations	6	18	8	9
Akaike's Information Criterion (AIC)	725268	660444	220964	217026
Bayesian Information Criterion (BIC)	725280	660469	221038	217125
Model Log Likelihood	-362633 (df = 1)	-330220 (df = 2)	-110475.9 (df = 6)	-108505.2 (df = 8)

*REF* = reference, *SE* = Standard Error, *df* = degree of freedom

\*\*\**p* < .001

\*\**p* < .01

\**p* < .05

in-between the confidence intervals lines. The wide range of degree distribution of our co-authorship network makes it difficult to assess model fit in terms of degree distribution. But it is clear that in general, model 4 fits poorly to the observed network despite the highly significant estimates obtained. We therefore have strong evidence confirming that there is likely something other than the terms included in this model that are driving the structure of the network, possibly additional attributes our study did not control for. The following section attempts to address this shortcoming.

### 3.3.2.3 Temporal Exponential Random Graph Model

The observed cumulative network was subset in seven snapshots representing respectively the following time spans: 1996 – 2006, 2007 – 2009, 2010 – 2011, 2012 – 2013, 2014, 2015

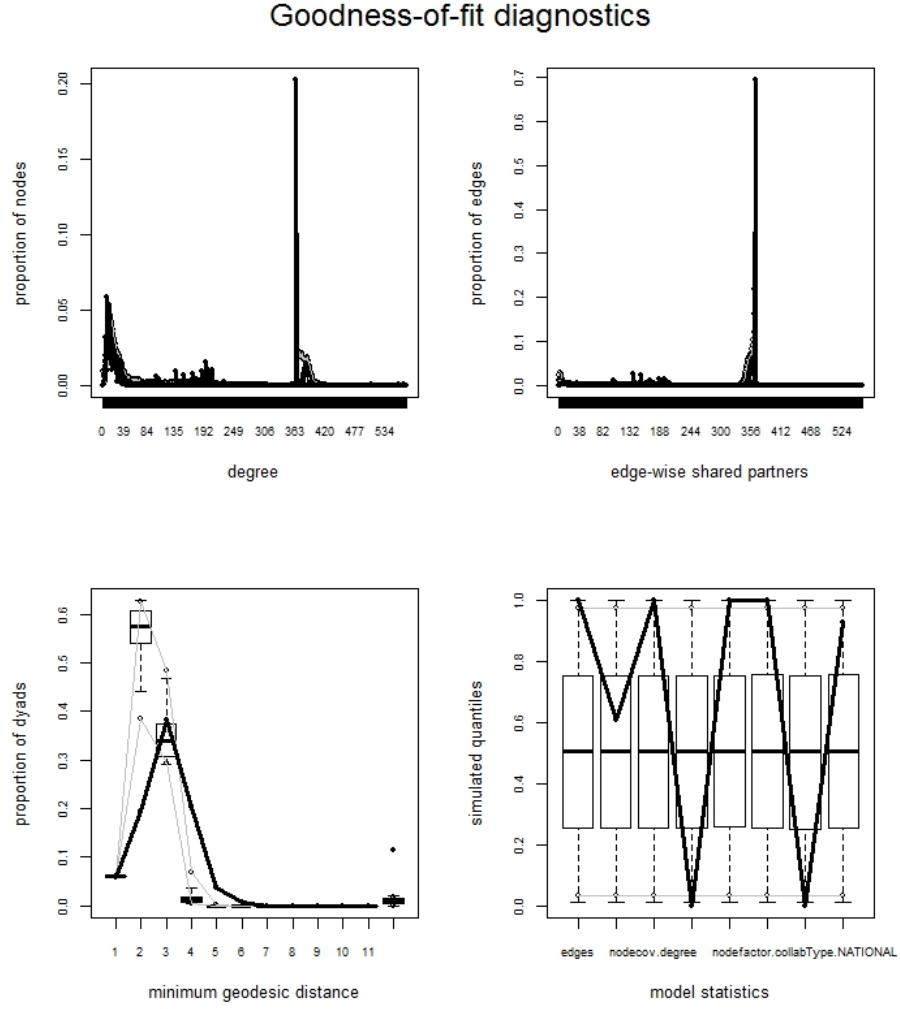


FIGURE 3.9: ERGM goodness-of-fit of final model 4 assessment.

and 2016. Figure 3.10 displays the topological structure of the snapshots of the different time steps.

Table 3.4 summarizes the results of the different temporal models we fit to the observed network. Models 1, 2, and 3 are equivalent to a pooled ERGM across the 7 different time points (Fig. 3.10). The null model of the TERGM (model 1) suggests that the baseline log-odds for collaboration tie formation between authors in the network is  $-4.66$ . This coefficient is equivalent to a baseline probability of 0.9% for any two authors in the

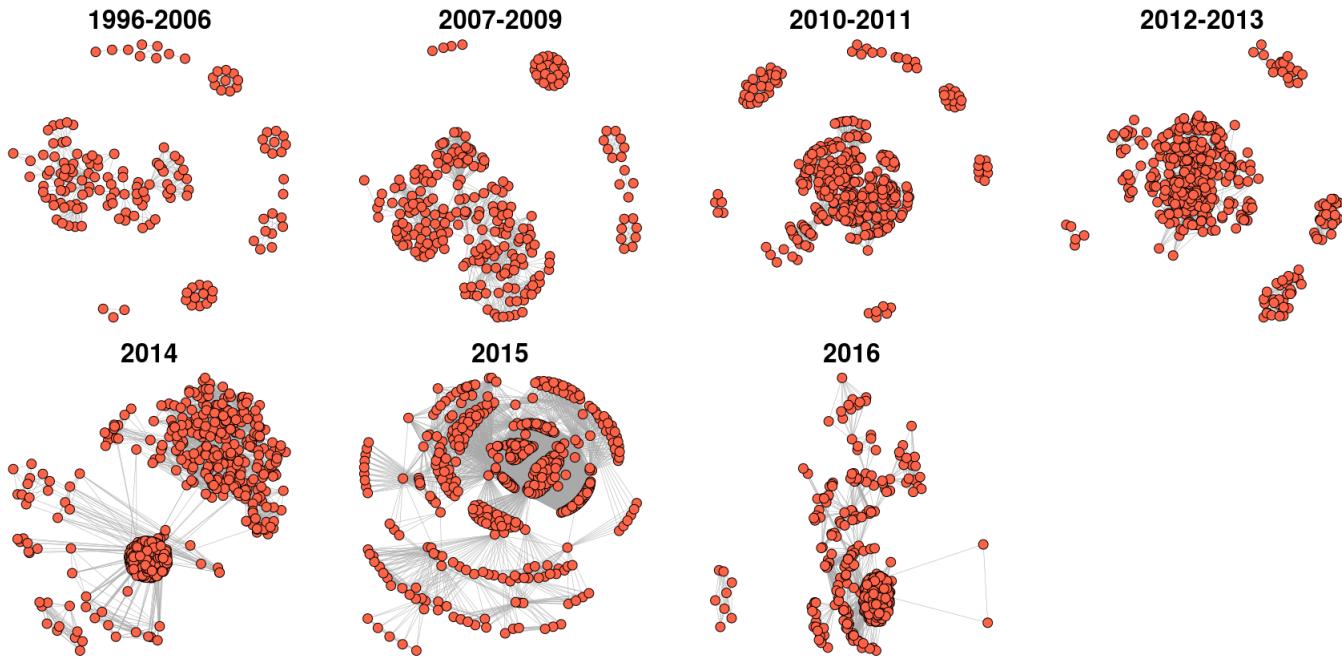


FIGURE 3.10: Topological structure of the different snapshots of the malaria co-authorship network.

network to establish a stable collaboration tie. This probability is significantly lower than the 5.96% baseline probability of collaboration tie establishment reported by the ERGM (section 3.3.2.2).

Model 2 of the TERGM describes the co-authorship network as a function of the number of collaborations, the number of publications, and the number of citations of authors inside the network. It is also adjusted by homophily on cluster assignment from the SBM and on the collaboration type. Compared to model 1, model 2 has slightly improved (See AIC and BIC in table 3.4). The edge effect has decreased (Coefficient =  $-10.14$ ,  $p < 0.001$ ) with the associated conditional probability (given all other terms in the model) equal to 0.004%. We observed a relatively high positively significant effect of the homophily on cluster assignment on the odds of collaboration tie formation between

## *Results: The Malaria Co-authorship Network*

---

any two authors. Adjusting for the other variables in model 2, authors of the same research groups/communities are 4.96 times as likely to collaborate than authors that belong to different research groups. The effect of the other attributes in model 2 are minor. When we adjust for attribute effect on the collaboration type, we obtained model 3 which is slightly better than model 2. Relatively to model 2, the edge effect decreases more followed by an even stronger effect of the homophily on cluster assignment of the authors in the network (Coefficient =  $-5.06$ ,  $p < 0.001$ ).

After introducing temporal dependencies terms, we obtained model 4 which tremendously improved compared to models 1,2 and 3. Model 4 confirms the observation made in section 3.3.2.2 that the the process underlying the malaria co-authorship network is driven by homophily on cluster assignment or membership to a specific research community and the type of collaboration. It further confirms that the linear trend suspected observed in figure 3.2 is significantly associated with the odds of collaboration tie formation in the Malaria co-authorship network. Model 4 suggests that the baseline conditional probability of any two authors to collaborate is estimated at 0.02% given all other terms in the model. The coefficient associated to the dyadic stability term is 1.07 meaning that the odds of existent and non existent collaboration ties at one time point to remain the same at the next time point increased on average by 65.7%. In other words, the odds of new collaboration ties and non-ties to occur from one time point to another is 34.3%. In addition, the TERGM showed that the probability of sustainable collaboration tie formation among international researchers is 12.13% versus 12.24% for researchers affiliated with national institutions ( $p > 0.05$ ). However, this probability significantly increases to 20.26% for

## *Results: The Malaria Co-authorship Network*

---

researchers affiliated to African research institutions other than those in Benin. These probabilities confirm the results from the ERGM final model with respect to the higher probability of tie formation between researchers affiliated to African institutions other than Beninese institutions. None of the structural temporal models containing high order TERGM terms, nor the models containing the dyadic attribute terms converged after the maximum of 1,000 iterations making estimates from these models untrustful.

TABLE 3.4. Temporal ERGM of Malaria Co-authorship Network.

	Model 1	Model 2	Model 3	Model 4
	Estimate (SE)	Estimate (SE)	Estimate (SE)	Estimate (SE)
Network structural predictor				
Intercept(edge)	-4.66 (0.00)***	-10.14 (0.02)***	-10.45 (0.02)***	-8.65 (0.05)***
Number of collaborations	-	0.03 (0.00)***	0.03 (0.00)***	0.03 (0.00)***
Number of times cited	-	-0.03 (0.00)***	-0.02 (0.00)***	-0.03 (0.00)***
Number of publications	-	0.45 (0.00)***	0.46 (0.00)***	0.45 (0.00)***
Homophily on cluster assignment	-	4.96 (0.02)***	5.06 (0.02)***	4.79 (0.02)***
Homophily on collaboration type	-	0.44 (0.01)***	0.56 (0.01)***	0.54 (0.01)***
Factor attribute effect (collaboration type)				
International	-	-	REF	REF
National	-	-	-0.10 (0.02)***	0.01 (0.02)
Regional	-	-	0.55 (0.01)***	0.60 (0.01)***
Temporal dependencies				
Dyadic stability	-	-	-	1.07 (0.01)***
Linear trends	-	-	-	-0.18 (0.01)***
Akaike's Information Criterion (AIC)	94681198	93740511	93737596	67005816
Bayesian Information Criterion (BIC)	94681230	93740624	93737742	67005991
Model Log Likelihood	-47340597	-46870248	-46868789	-33502897

REF = reference, SE = Standard Error

\*\*\* $p < .001$

\*\* $p < .01$

\* $p < .05$

Figure 3.11 presents the goodness-of-fit assessment for the TERGM model 4. We can see that this model containing temporal dependencies fits better to the observed Malaria co-authorship network than the final ERGM model 4. While the first five subfigures compare the distribution of endogenous network statistics between the observed network and the simulated ones, the last subfigure presents the Receiver Operating Characteristics

## Results: The Malaria Co-authorship Network

---

(ROC) and precision-recall (PR) curves. The ROC for model 4 is depicted by the dark red curve compared to the ROC of a random graph depicted by the light red curve. Similarly, the dark blue curve represents the PR of model 4 versus the light blue curve representing the PR of a random graph [98]. It clearly appears that the final TERGM model 4 outperformed the random null model with an Area Under the Curve (AUC) value estimated at 79.98%.

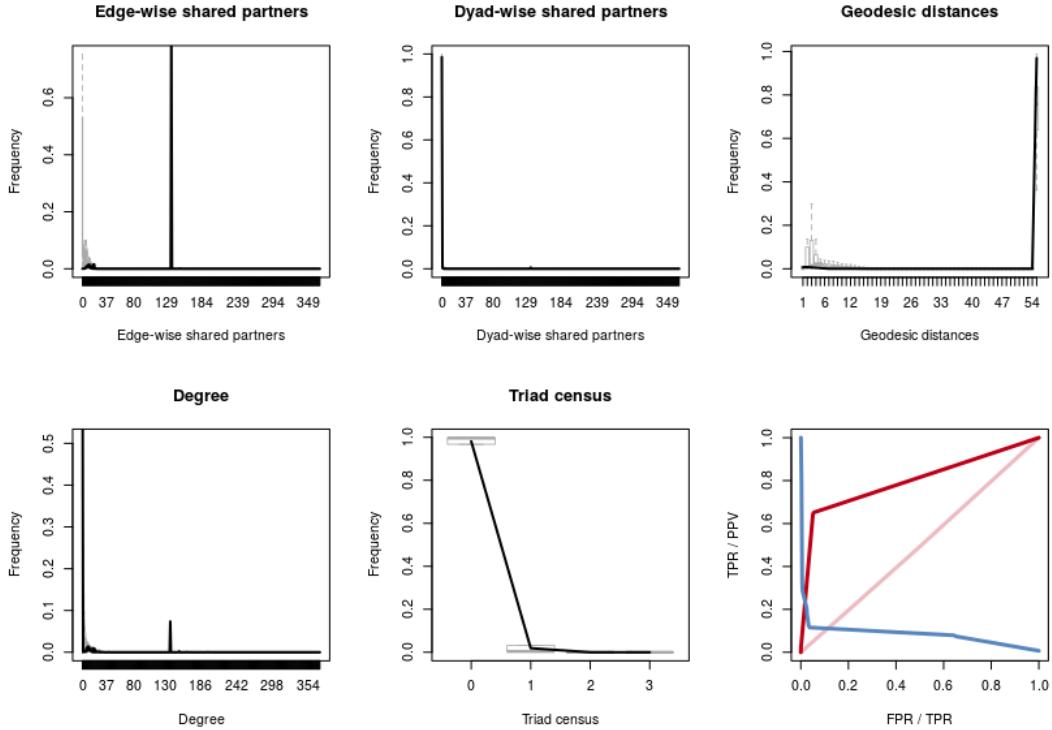


FIGURE 3.11: Goodness-of-fit assessment for the final Malaria TERGM Model 4 with temporal dependencies.

### 3.3.2.4 Latent Network Model

Figure 3.12 presents a 3-dimensional visualization of the Malaria co-authorship network, with layouts determined according to the inferred latent eigenvectors from the no pair-specific model (on top), the model containing nodal covariates (middle), and the model containing nodal and dyadic covariates (bottom). Blue vertices represent authors affiliated to Beninese research institutions, Red vertices are authors affiliated to international institutions, Gold vertices represent authors affiliated to African research institutions other than Benin, and White vertices represent authors with no determined affiliations. Node sizes are proportional to the betweenness value of each vertex. Looking at the three visualizations, it clearly appears that the first two visualizations are somewhat similar while the third is different. In fact, in the first two visualizations, the authors are clustered in mainly three clusters. We can see that all the authors affiliated to Beninese research institutions (in blue) are clustered in one cluster while authors with international affiliations (in red) and regional authors (in gold) are distributed across all three main clusters. These observations suggest a significant geography effect on the odds of collaboration tie establishment in the malaria co-authorship network.

The first two visualizations also highlight key brokers that liaison between clusters. In the third visualization, on the other hand, there appears to be only one main cluster. This last observation suggests that the nodal covariates and mainly homophily on research community membership and type of affiliation explain much less coarse-scale network compared to dyadic covariates. Indeed, when the dyadic covariates are added to the model, there is less structure left to be captured by the latent variables. These results

## *Results: The Malaria Co-authorship Network*

---

compensate the lack of-fit of the ERGM model and confirmed our findings in the previous section.

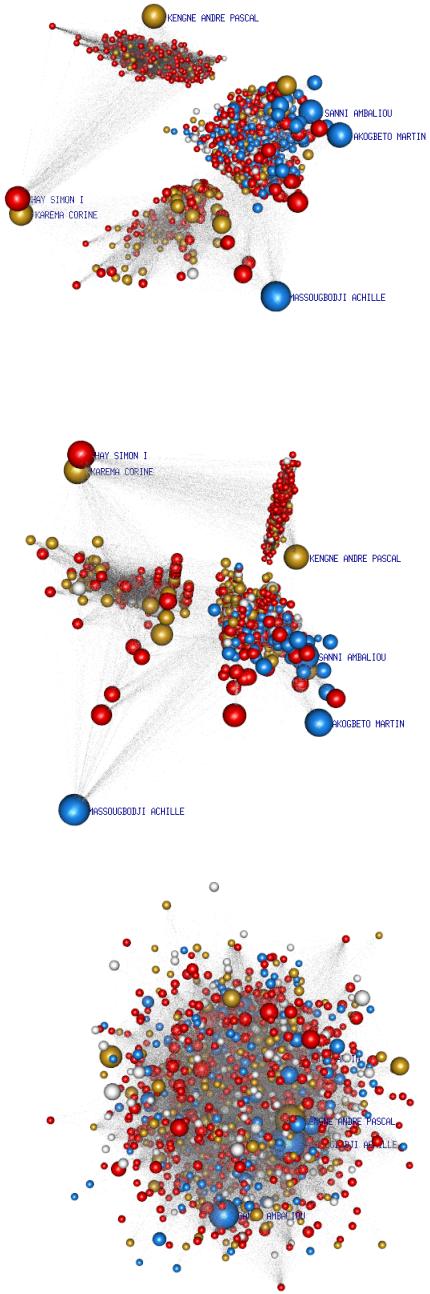


FIGURE 3.12: Visualizations of the Malaria co-authorship network with layouts determined according to the inferred latent eigenvectors in the LNM models.

The ROC curves on figure 3.13 show that the first two models appear to be comparable

in their performance from the perspective of edge status prediction with an Area Under the Curve (AUC) being roughly 98.8%.

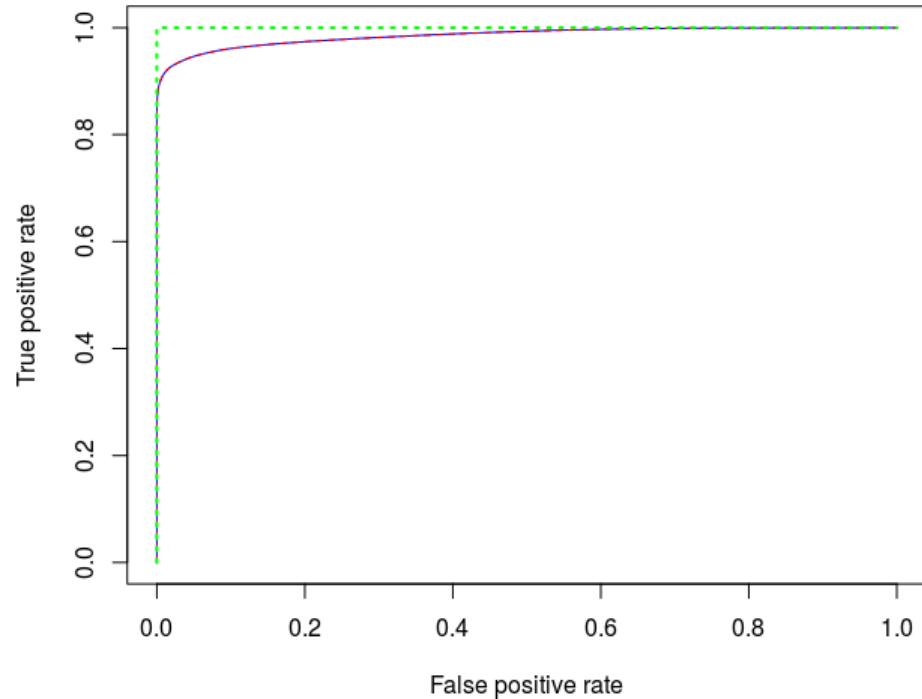


FIGURE 3.13: ROC curves comparing the goodness-of fit of the Malaria co-authorship network for three different eigenmodels, specifying (i) no pair specific covariates (blue), (ii) nodal covariates (red), and (iii) nodal and dyadic covariates (green), respectively.

### 3.4 Discussion and Conclusion

In this chapter, we provide insights in the structural characteristics of the malaria co-authorship network in the Republic of Benin over a relatively long period. The 20 years of data collected coincides with the onset of active malaria research from 1996 until today. The significant increase in malaria research and collaborations (figure 3.7) between the authors over the years is an expected finding given the regain and renewed interest in malaria control and elimination goals set forth [102, 104]. Our results show that the mechanism underlying the formation of the malaria co-authorship network in Benin is not random. It further demonstrates that the malaria research collaboration network in Benin is a complex network that seems to display small-world properties (often referred to as "six degrees of separation").

The non-trivial number of authors with higher order of magnitudes confirms the presence of closed research groups where collaborative research likely happens only among members. In other words, interdisciplinary collaboration tends to occur at higher levels between prolific researchers with the majority of the collaborations happening between researchers from the same scientific communities. Prominent authors with important collaborations tend to collaborate with similar authors, young or less prolific authors tend to collaborate with both prolific authors and authors with very few collaborations. Similar findings were reported by Janet Okamoto [105] who studied scientific collaboration on a much smaller scale. Key brokers facilitate scientific collaborations within and outside their scientific community [67]. Betweenness centrality measures identify such

## *Results: The Malaria Co-authorship Network*

---

brokers who are important hubs for inter and transdisciplinary research. Many of the main brokers proved to also be the most connected and the most central authors confirming the presence of long publishing tenure authors in our network [106]. The flow of information in this network in Benin is slow as it only relies on 16 authors representing less than 1% of all the authors in the network. Such a low information flow was also reported by Salamatia and Soheili [70] in a 2016 study on a co-authorship analysis of Iranian researchers in the field of violence. Generally, the most important authors in a co-authorship network are the ones with the highest degree of collaborations [107, 108]. However, to the long-term sustainability of the malaria research network in Benin, the 16 authors identified as cut vertices are the most important authors. In other words, the removal of less than 1% of the authors from the network would lead to its collapse. Such a collapse would undoubtedly be detrimental to the future of malaria research in Benin. This finding clearly confirms the conclusion of Toivanen and Ponomariov [66] that the African research collaboration network is vulnerable to structural weaknesses and uneven integration.

Small-world networks are known to have small shortest path distance and a high clustering coefficient. Although this co-authorship network seems to display such properties, the Monte-Carlo simulations revealed that the observed network has unexpected properties compared to classic small-world networks. A study of co-authorship network conducted on Chagas disease has found similar findings [55]. Unlike our study, the authors of the Chagas disease co-authorship study did not deepen their analysis to confirm the small-world nature of their observed network. Other mechanisms such as preferential attachment

## *Results: The Malaria Co-authorship Network*

---

have been found to explain the structure of international scientific collaboration network [109]. Unlike those studies, our network displayed unexpected properties that are more extreme than the 4 mathematical models we simulated. Our network has significantly larger shortest path distance and significantly higher clustering than expected from the 4 mathematical models presented here. One observation we are sure of is that none of the random graph models used here tend to explain the growth and the structure of the malaria co-authorship network in Benin. We therefore claim without any doubt that the structure and growth of our network is not random confirming the presence of hidden factors explaining the current structure of the network. Assessing such factors and the extent to which they influence scientific collaborations is important for the future of malaria research and its long-term sustainability. Unfortunately, none of the proposed mathematical models seem to accurately describe the observed structure of the network. To address these limitations, advanced statistical modeling was used to further explain the structure of the network.

Our first approach to modeling our network relied on the use of SBM. In addition of being a model based clustering method, the SBM identified important organizational and interactional patterns in the network. It identified a large clique of mainly international researchers with little or no collaborations with other research groups. The overwhelming dominance of regional and international players in the network is consistent with previous observations by Onyancha and Maluleka [110] who concluded on a much higher likelihood of Sub-Saharan African countries to collaborate with non-African states.

## *Results: The Malaria Co-authorship Network*

---

Overall, the ERGM and TERGM show that the mechanistic phenomenon driving collaboration ties in the malaria research in Benin is influenced by homophily on the type of affiliation (national, international or regional) and on membership to a research group or cluster, verifying therefore our third hypothesis. The models clearly show that the dominance of the Beninese malaria research arena by international and regional players, and further demonstrates the lower likelihood of local Beninese researchers to establish international collaboration ties compared to regional researchers. This latter finding has been confirmed by the LNM which also confirms our second hypothesis. The ERGM and the TERGM revealed that factors such as number of publications, number of citations and number of collaborations are associated to higher likelihood to establishing collaboration ties, confirming therefore our first hypothesis.

It is worth noting that many of the studies on co-authorship network analysis are descriptive in nature. This study is one of the rare co-authorship network analysis to model a co-authorship network using advanced statistical models. ERGM is the leading approach to modeling network [111]. The literature has reported application of this model in studying various social network such as the analysis of friendship and obesity [112, 113], the exploration of the association between hormone and social network structure [114]. Similarly to friendship networks, the use of ERGM to model co-authorship networks is easily justified. However, the size of our network prevented the fitting of complex models including dyadic and structural terms. In addition, our best ERGM model failed to adequately fit the observed network data. This lack of goodness-of fit, according to Hunter, Goudreau and Handcock [115], could be improved by including the geometrically weighted edgewise

### *Results: The Malaria Co-authorship Network*

---

shared partner, geometrically weighted dyadic shared partner, and geometrically weighted degree network statistics to our model. Although, we follow such recommendations by including these structural network statistics to our final model, the ERGM model failed to converge after a maximum of 1,000 iterations. At about 750 iterations, we noticed that the processing became both computationally intensive and expensive in terms of CPU time and memory usage. In a recently published paper, Schmid and Desmarais [111] acknowledged the difficulty of fitting network which size is of the order 1,000 vertices using ERGM. They recommended that using the maximum pseudolikelihood estimation (MPLE) instead of the Monte Carlo maximum likelihood (MCMLE) could tremendously reduce computation time. Having followed these recommendations too, the ERGM model containing dyadic and structural terms still failed to converge. By finally including temporal dependencies and fitting a temporal ERGM, we have tremendously improved the fitness with a predictive performance of roughly 80%. Nevertheless, we suspect that the number of edges, the large size of the network added to the possibility of hidden/latent variables might justify the failure of the models containing the dyadic and structural endogenous terms to converge. We remedy this situation by applying LNM to the observed network data.

All three latent network models (LNM) proved to have tremendous fit to the observed network. The fact that the third latent network model include all nodal and dyadic variables could explain its perfect fit. In addition, there was less structure left to be captured by the latent variables in this model. On the other hand, the first two latent network models gave us confident in validating the results of the TERGM.

### *Results: The Malaria Co-authorship Network*

---

A study by Kronegger et al. [116] conducted an investigation aiming at describing the collaboration in Slovenian scientific communities using data from four different disciplines. Their methodological approach is consistent with ours. The main difference is their application of Stochastic Actor-Oriented Model (SAOM) on the dynamics of their co-authorship networks. Since the SAOM is an actor-oriented modeling method and we are interesting in tie prediction here, we relied rather on a tie-oriented approach by applying the TERGM to our network data.

Our results suggest that the regain in Malaria research funding has appealed to research groups all around the world, hence the explosion in publications number and research collaborations. As the disease continues to be main public health concern in the Republic of Benin, it is essential to consolidate the knowledge generated from the numerous studies on the disease and reinforce the different communities involved in the research effort. In addition, there is an urgent need to reinforce the malaria research network in Benin by continuously supporting, stabilizing the identified key brokers and most productive authors, and promoting the junior scientists in the field. However, we observed a tendency of the international researchers to only collaborate among themselves. Although the rise in scientific collaboration between advanced and developing nations [117], the latter observation may limit effective and sustainable technology transfer in Benin. It is possible that some of the isolated cliques within the network have top-notch research capabilities and skills researchers affiliated to Beninese institutions can acquire, should the research groups be more inclusive. Unfortunately, our visualizations showed that broker authors

---

*Results: The Malaria Co-authorship Network*

---

that liaison those closed groups to national researchers tend to be regional or international researchers as well. We therefore recommend, that policies should be designed, at international, regional and country level, to diversify research groups operating in any Sub-Saharan African countries. Such policies will ultimately enable effective technology transfer, multidisciplinarity, and promote junior African researchers to advance the search of a solution to the Malaria problem in Africa and particularly, in Benin.

# Chapter 4

## Results: The HIV/AIDS Co-authorship Network

### 4.1 Data

A literature search was conducted in the Web Of Science (WOS) using combinations of HIV/AIDS related MeSH terms including "HIV", "AIDS", "VIH" and "HIV Infections". We restricted the search to the period from 1996 to 2016 and to "Benin" for country. We further screened the papers in order to only select those published by Beninese authors, or papers published on HIV/AIDS from Benin. No restriction was placed upon the document types. We first started querying with each term independently, we then combined the other terms so the query return the maximum number of results. The Full citations information containing the authors' names, their institutional affiliations, the

### *Results: The HIV/AIDS Co-authorship Network*

---

year of publication, as well as the number of times the document was cited were recorded as a bibliographic corpus in text format. After a second screening only research that have met the above listed inclusion criteria and that were published between January 1, 1996 and December 31, 2016 were selected.

The final query set (Table 4.1) returned 237 records. After a rigorous screening process, 102 documents met the selection criteria. On average, there was 9.47 authors per published document.

TABLE 4.1. HIV/AIDS Bibliographic Search Queries.

Set	Queries	Results
#1	TOPIC: (HIV AIDS) Refined by: COUNTRIES/TERRITORIES: (BENIN)	52
#2	TOPIC: (HIV AIDS) AND ADDRESS: (BENIN)	107
#3	TOPIC: (HIV) OR TOPIC: (AIDS) AND ADDRESS: (BENIN), Refined by: COUNTRIES/TERRITORIES: (BENIN)	182
#4	TOPIC: (HIV) OR TOPIC: (VIH) OR TOPIC: (AIDS) AND ADDRESS: (BENIN), Refined by: COUNTRIES/TERRITORIES: (BENIN)	182
Final Set	#1 OR #2 OR #3 OR #4	237

The Author Name Disambiguation process led to the identification of 516 unique authors with a precision of 99.88% and a recall of 82.54%. The generated multigraph co-authorship network therefore contained 516 vertices (authors) and 5,114 parallel edges (collaborations). As displayed in figure 4.1, we can see the significant increase in publications, scientific collaborations and the number of authors involved in HIV/AIDS research from 2008 until 2016. This general upward trend seems to be linear from the year 2008 to 2016.

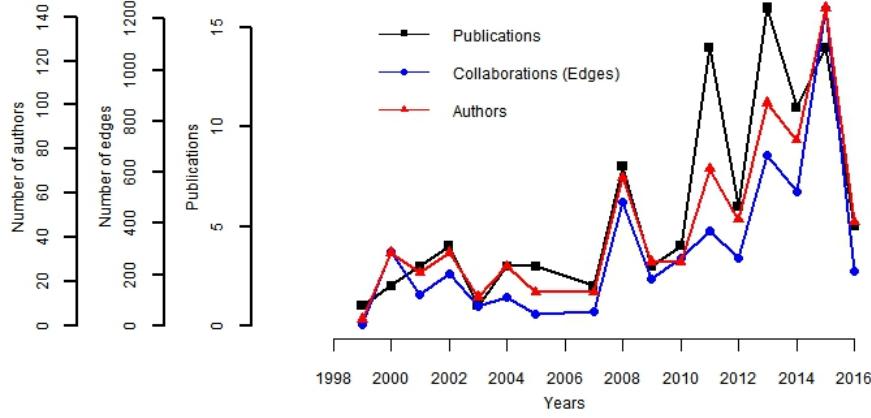


FIGURE 4.1: Evolution of the published HIV related documents, authors and collaborations from January 1996 to December 2016

## 4.2 Descriptive Data Analysis

For the multigraph network, the degree distribution ranged between 1 and 403 with an average degree distribution of 19.82 and a median of 12. In addition, there was a substantial number of vertices with low degrees (Fig. 4.2). The log scale distribution of the degrees on figure 4.3 reveals that there was also a non-trivial number of vertices with higher order of degree magnitudes. This observation confirms the tendency of vertex degrees to follow a heavy-tail distribution suspected on figure 4.2.

After we convert the multigraph network in a weighted graph, it results in a simple graph of 516 vertices and 3,966 weighted edges. Closeness centrality ranges between  $3.76 \times 10^{-6}$  and  $3.19 \times 10^{-5}$  with a median of  $3.13 \times 10^{-5}$ . Betweenness measures range between 0 and 49,280 with a median of 426.2. A network visualization with the vertices' size proportional to betweenness centrality measures clearly reveals the presence of broker

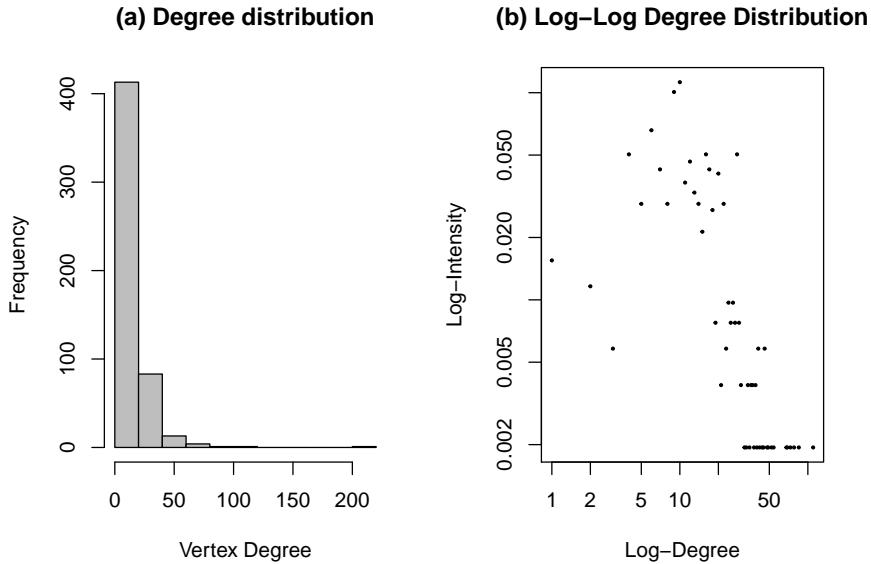


FIGURE 4.2: Degree distribution of the HIV/AIDS Co-authorship network

authors (Figure 4.4 and Table 4.2). The median Eigenvectors is 0.202 with a mean of 0.045. The eigenvectors measures confirm the presence of author hubs in the network suggesting the presence of closed collaboration groups. Table 4.2 presents a list of the 10 author hubs with the highest Eigenvectors values.

Edge betweenness centrality measures identify co-authorship collaboration ties that are important for the flow of information. Table 4.2 presents the top 10 most important collaboration ties for the flow of information in the HIV/AIDS Co-authorship network in Benin.

#### 4.2.1 Network Cohesion

In total, 29 maximal cliques were detected in the network among which 2 cliques of size 24, 1 clique of size 23 and 4 cliques of size 3. Larger maximal cliques sizes range from 14

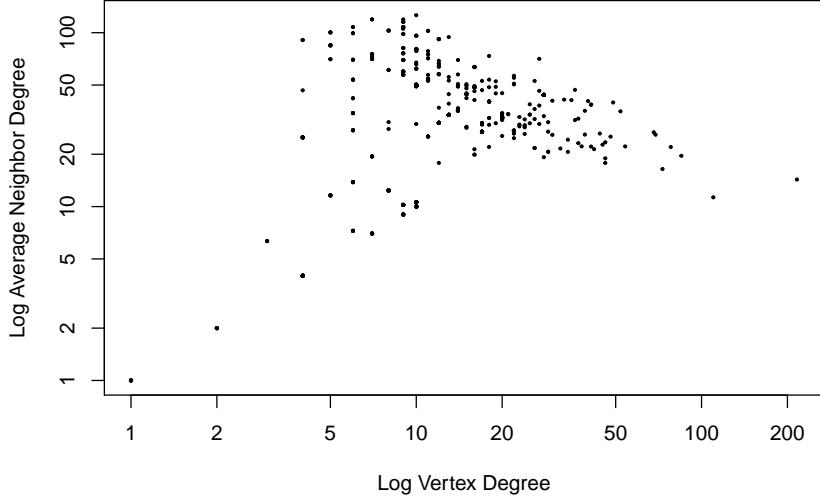


FIGURE 4.3: Log-Average Neighbor degree Distribution of the HIV/AIDS Co-authorship network

authors to 25 authors.

The HIV/AIDS co-authorship network has a density of 0.0298 indicating that the baseline probability of collaboration tie formation is 2.98%. The network also has a transitivity of 0.482 meaning that 48.2% of the connected triples in the network are close to form triangles. The transitivity metrics is a measure of the global clustering of the network. The network is not connected and a census of all the connected components within the network reveals the existence of a giant component that dominates all the other connected components. The giant component of the HIV/AIDS co-authorship network includes 88.6% (457 vertices) of all the vertices in the network with the other components alone carrying less than 1% of the vertices (Fig. 4.4).

Information flow assessment of this network via cut vertices confirms the existence of 8 authors as the most vulnerable vertices in the network. Table 4.2 lists the authors

*Results: The HIV/AIDS Co-authorship Network*

---

TABLE 4.2. List of the most important authors and collaborations in the HIV/AIDS Co-authorship network

<b>Top 10 Brokers</b>
ZANNOU DJIMON MARCEL
ALARY MICHEL
LEROY VALERIANE
AZONDEKON ALAIN
ANAGOUNOU SEVERIN
ADE GABRIEL
AZONKOUANOU ANGELE
NDOYE IBRA
NDOUR MARGUERITE
AFFOLABI D
<b>Top 10 most connected authors (Top 10 network hubs)</b>
ZANNOU DJIMON MARCEL
ALARY MICHEL
ANAGOUNOU SEVERIN
LOWNDES CATHERINE M
LABBE ANNIECLAUDE
DABIS FRANCOIS
MINANI ISAAC
BEHANZIN LUC
DIABATE SOULEYMANE
EKOUEVI DIDIER K
<b>Top 10 most important edges for information flow</b>
ZANNOU DJIMON MARCEL – LEROY VALERIANE
ZANNOU DJIMON MARCEL – NDOUR MARGUERITE
ALARY MICHEL – AZONKOUANOU ANGELE
ZANNOU DJIMON MARCEL – NDOYE IBRA
ANAGOUNOU SEVERIN – ADE GABRIEL
ZANNOU DJIMON MARCEL – WACHINOU ABLO PRUDENCE
ZANNOU DJIMON MARCEL – DALMEIDA MARCELLINE
AZONDEKON ALAIN – ADE GABRIEL
AZONKOUANOU ANGELE – AZONDEKON ALAIN
ZANNOU DJIMON MARCEL – COFFIE PATRICK A
<b>Weak articulation points</b>
ATADOKPEDE FELIX
NDOUR MARGUERITE
DALMEIDA MARCELLINE
AZONDEKON ALAIN
GANDAHO PROSPER
AFFOLABI D
ADE GABRIEL
ZANNOU DJIMON MARCEL

that constitute the weak articulation points in the HIV/AIDS co-authorship network.

The identification of cut vertices is a measure of the vulnerability of the HIV/AIDS co-authorship network [83].

Via the agglomerative hierarchical clustering method, we identify 24 different research communities (or clusters) which sizes range between 1 and 108 authors. Large research communities contain between 71 and 108 authors. Medium size research communities contain between 10 and 55 authors. Out of the 24 research clusters or communities detected, 12 are part of the giant component. Figure 4.4 displays the structure of the network with each different colors representing each of the 24 research communities.

## 4.3 Modeling

### 4.3.1 Mathematical Modeling

We performed 1,000 Monte Carlo based simulations to test the significance of the observed characteristics of the HIV/AIDS co-authorship network. Figure 4.5 clearly demonstrates that the number of communities detected is unusual from the perspective of both Classical random graphs and generalized random graphs ( $p\text{-value} < 0.0001$ ). From the Classical random graph model, the expected number of communities was 5.574 (95%CI: 5.53 – 5.62). Similarly, the expected number of communities from the generalized random graph model is 6.65 (95%CI: 6.60 – 6.70).

Figure 4.6 displays the number of detected research communities using the Barabási-Albert’s preferential attachment and the Watts-Strogatz models. The observed number of communities was extreme per both models ( $p\text{-value} < 0.0001$ ). The expected number from the Watts-Strogatz model simulations is 3.181 (95%CI: 3.16 – 3.21) and 22.8

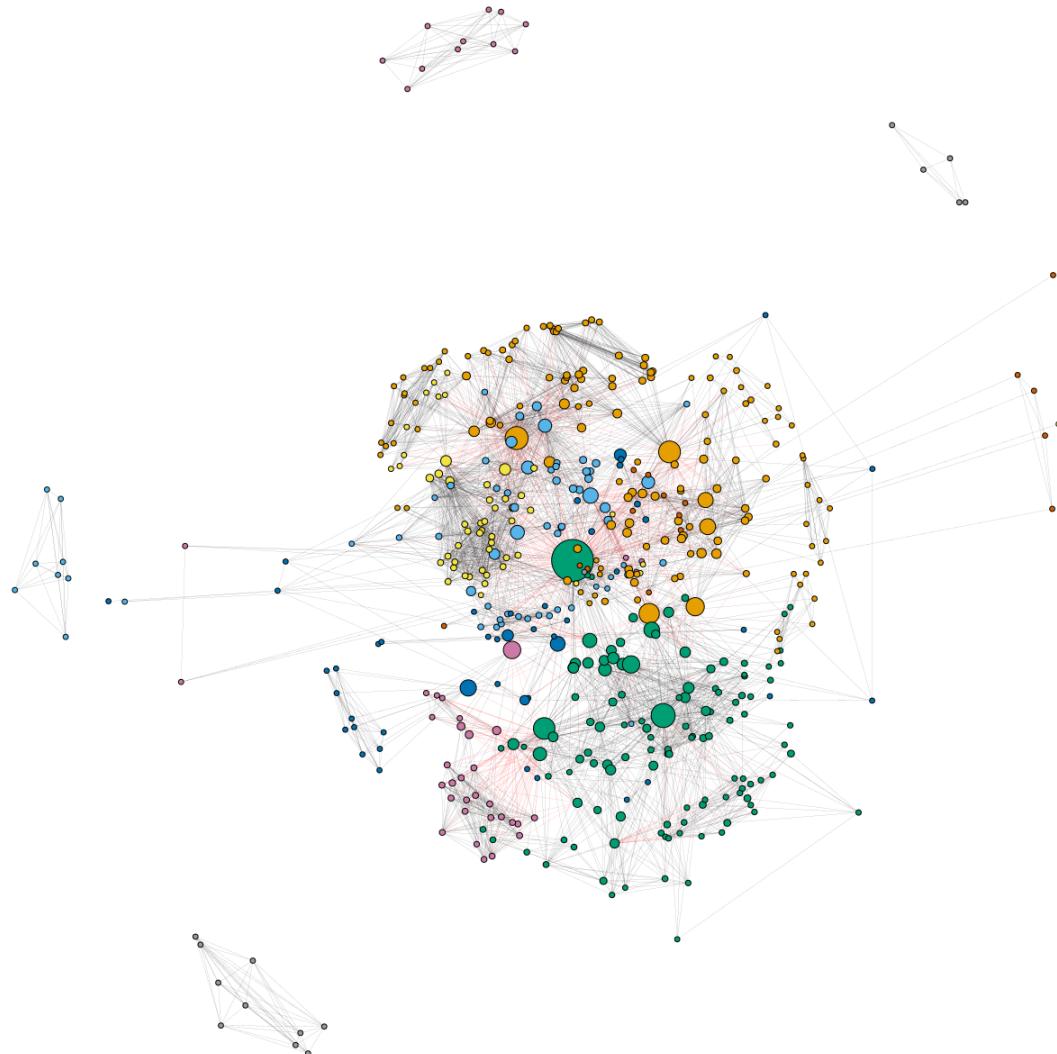


FIGURE 4.4: Topological Structure of the HIV/AIDS Co-authorship Network. Authors (vertices) of the same color belong to the same research community or cluster

(95%CI: 22.7 – 23.0) from the Barabási-Albert model simulations. We also compared the clustering coefficient and the average shortest-path length. Let's recall that the observed clustering coefficient is 0.482. On one hand, there was substantially more clustering in

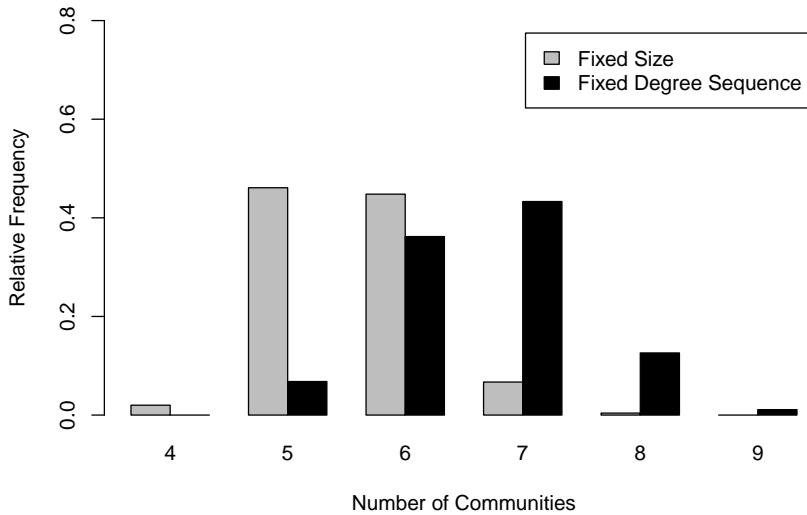


FIGURE 4.5: Monte-Carlo simulations of the HIV/AIDS network: Number of detected communities by the random graph models

our HIV/AIDS co-authorship network than expected from both random graph models ( $p\text{-value} < 0.0001$ ). The expected clustering coefficients was 0.0365 (95%CI: 0.0363 – 0.0365) and 0.0842 (95%CI: 0.0841 – 0.0843) respectively for the classic random graph and the generalized random graph models.

On the other hand, there was substantially less clustering in our HIV/AIDS co-authorship network than expected by the Watts-Strogatz Small World model which expected clustering was 0.72615 (95%CI: 0.72611 – 0.72618).

We observed an average shortest-path length of 2.75 in the HIV/AIDS co-authorship network. This observed shortest-path length is significantly larger than what was expected from the random graph models ( $p\text{-value} < 0.0001$ ) and significantly lower than what was expected from Watts-Strogatz small world model and the Barabási-Albert preferential attachment model ( $p\text{-value} < 0.0001$ ).

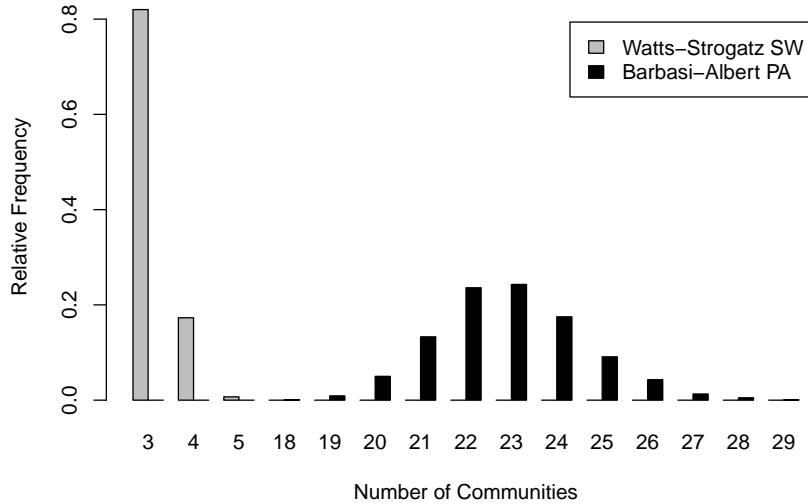


FIGURE 4.6: Monte-Carlo simulations of the HIV/AIDS network: Number of detected communities by the Watts-Strogatz and the Barabási-Albert models

The average shortest-path length was 2.49069 (95%CI: 2.49062 – 2.49077) and 2.381 (95%CI: 2.380 2.381) respectively for the classic random graph and the generalized random graph models.

For the Watts-Strogatz small world and the Barabási-Albert preferential attachment models, the average shortest-path length is respectively 5.31 (95%CI: 5.28 – 5.36) and 7.35 (95%CI: 7.31 – 7.38).

We performed the same simulations on the giant component of the network with similar results leading to similar outcomes.

### 4.3.2 Statistical Modeling

#### 4.3.2.1 Stochastic Block Model

The SBM identifies 26 classes with a degree of latitude of 17 to 26 classes being reasonable (See ICL plot on figure 4.7).

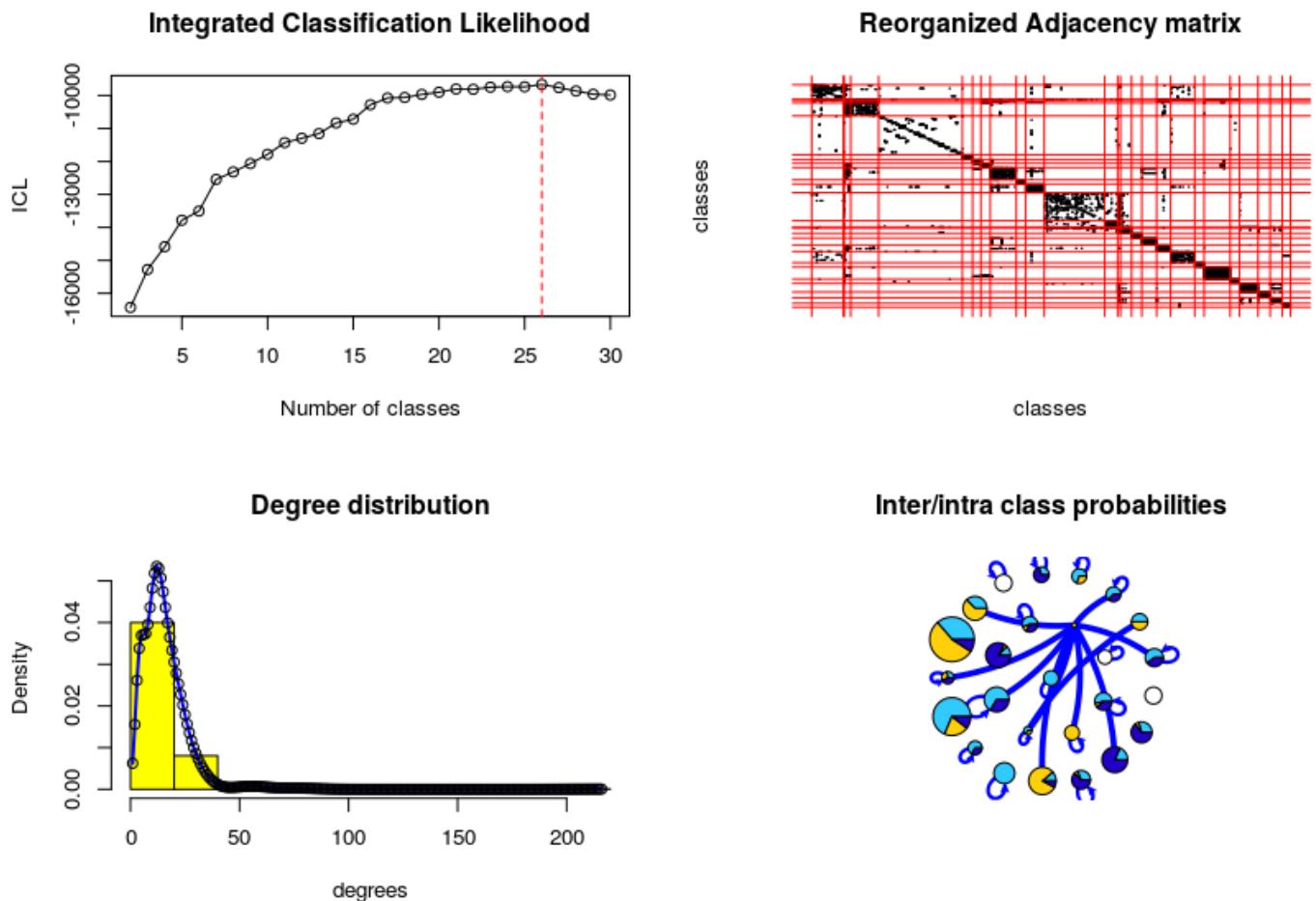


FIGURE 4.7: Summary of the goodness-of-fit of the SBM analysis on the HIV/AIDS co-authorship network.

Regarding the degree distribution, the fitted SBM describes well the observed degree distribution. On the network depicting the inter/intra probabilities between the classes,

## *Results: The HIV/AIDS Co-authorship Network*

---

the vertices represent the 26 identified classes, with each one of them divided into a pie chart displaying the proportion of authors of international affiliations (lightblue), authors of regional or other African affiliations (darkblue), and authors affiliated to Beninese research institutions (yellow). Generally, the dominance across the classes of international and regional players is observed. In addition, we observe denser ties between medium size and smaller size classes.

A close examination of the pie charts reveals that almost all the classes are heterogeneous. We note the presence of 2 large classes which are classes 5 and 12 (See reorganized adjacency matrix on figure 4.7). Class 5 is dominated by researchers with Beninese affiliations but appears sparser than class 5 which is dominated by international authors (Figure 4.7).

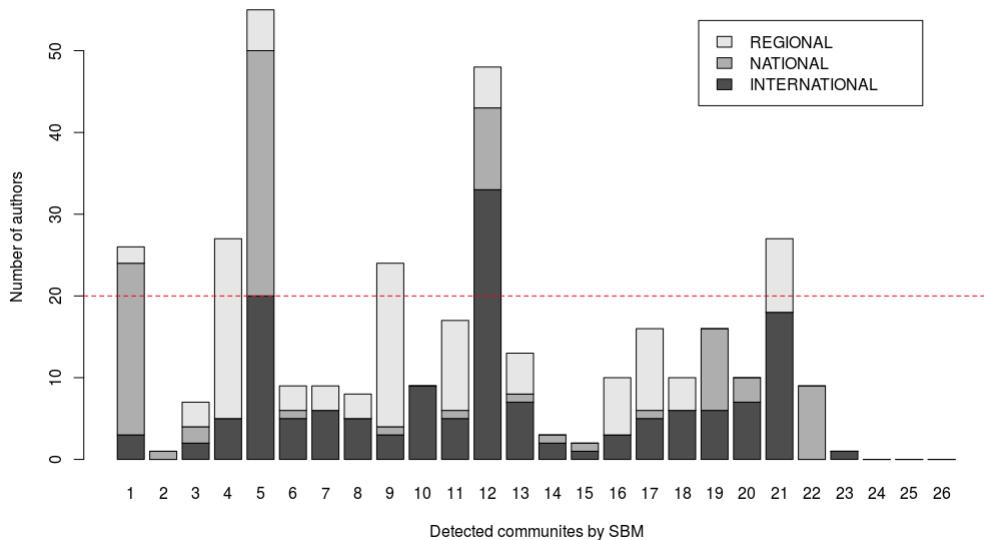


FIGURE 4.8: Distribution of national, international and regional authors by communities detected by the SBM in the HIV/AIDS network.

On figure 4.9, we present the SBM results emphasizing the largest classes (with more

*Results: The HIV/AIDS Co-authorship Network*

---

than 20 members). Here, we can confirm that smaller classes tend to collaborate more among themselves and intra-class collaborations tend to occur more.

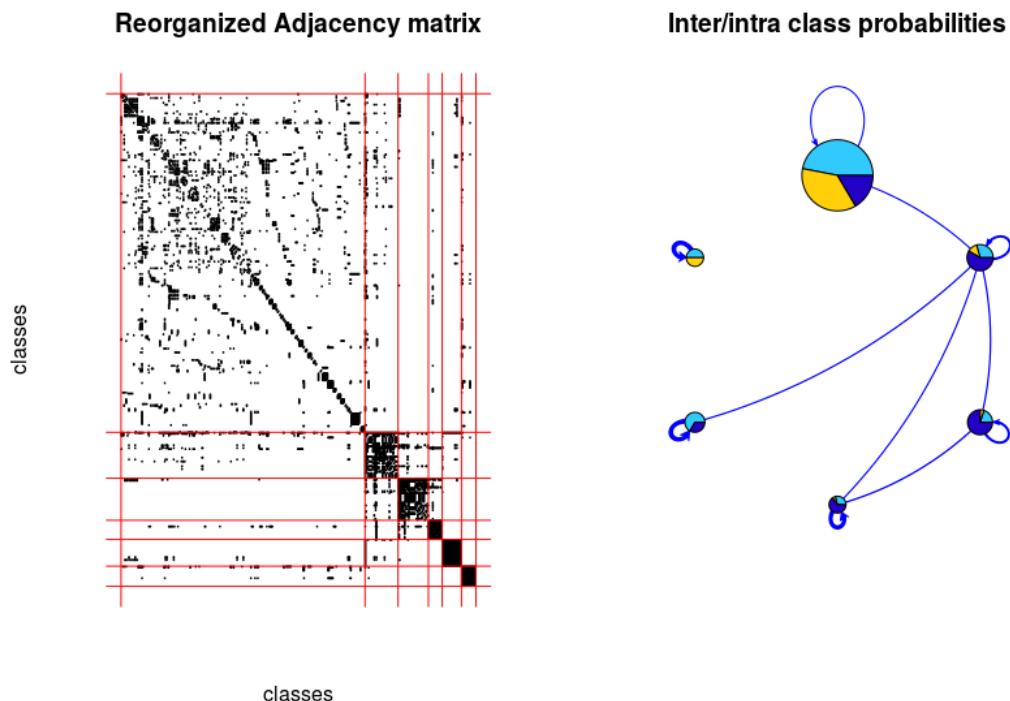


FIGURE 4.9: Summary of the goodness-of-fit of the SBM analysis highlighting interactions between the largest classes of the HIV/AIDS co-authorship network.

*Results: The HIV/AIDS Co-authorship Network*

---

	Model 1	Model 2	Model 3
	Estimate (SE)	Estimate (SE)	Estimate (SE)
Network structural predictor			
Intercept(edge)	-3.48 (0.02)***	-7.51 (0.06)***	-7.55 (0.07)***
Number of times cited	-	0.00 (0.00)***	0.00 (0.00)***
Number of collaborations	-	0.08 (0.00)***	0.08 (0.00)***
Number of publications	-	-0.29 (0.01)***	-0.28 (0.01)***
Homophily on cluster assignment	-	5.01 (0.05)***	5.02 (0.05)***
Homophily on collaboration type	-	0.77 (0.05)***	0.72 (0.05)***
Factor attribute effect (collaboration type)			
International	-	-	REF
National	-	-	-0.05 (0.04)
Regional	-	-	0.21 (0.03)***
AIC	35668.54	18956.20	18912.74
BIC	35678.34	19014.98	18991.12
Log Likelihood	-17833.27	-9472.10	-9448.37

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

TABLE 4.3. ERGM of the HIV/AIDS Co-authorship Network.

#### 4.3.2.2 Exponential Random Graph Model

Different models were fit with the ERGM method (Table 4.3). Model 1, the null model, contains only the "edge" term. The inverse logit of the coefficient associated with this term is 0.0298 which is the baseline probability of collaboration ties establishment and also the density of the HIV/AIDS co-authorship network.

In model 2, we included all nodal variables, a homophily term on collaboration type and on cluster assignment determined from the SBM. Model 2 improved tremendously compared to model 1 (See AIC, BIC and model likelihood in table 4.3). We note a decrease

---

*Results: The HIV/AIDS Co-authorship Network*

---

in the edge effect (Coefficient =  $-7.51$ ,  $p < 0.001$ ) with the associated conditional probability (given all the other terms in the model) estimated at 0.05%. For the remaining terms in model 2, we observed a positive and significant effect except for the number of publications. Model 3 differs from model 2 in that it includes a factor term on the collaboration type with a substantial improvement compared to model 2. Model 3 is therefore chosen as our last model. Regarding the number of publication, a one unit increase in the number of publication is associated with 32.3% average decrease in the odds of collaboration ties establishment. Model 3 further proves that the process underlying the structure of the HIV/AIDS co-authorship network in Benin is mainly driven by homophily on cluster assignment or membership to a research community or group (Coefficient =  $5.02$ ,  $p < 0.001$ ). The conditional probability of any two authors belonging to the same research group to collaborate is estimated at 7.38% compared to the baseline probability of 2.98%. The same probability changes to 14.06% after adjustment by the collaboration type, and 11.82% after adjusting for the number of citations, collaborations and publications. Compared to research affiliated to international institutions, researchers affiliated to Beninese institutions have 5.1% average decrease in the odds of collaboration tie establishment. This average decrease is not statistically significant ( $p > 0.05$ ). For researchers affiliated to institutions other than Beninese institutions, the odds of collaboration tie establishment increases on average by 18.9% compared to internationally affiliated researchers. Overall, model 3 estimated the probability of collaboration ties formation at 11.8% for international researchers, 11.3% for national researchers and 14.2% for regional players.

## Results: The HIV/AIDS Co-authorship Network

---

Since none of the models containing endogenous ERGM terms and/or the dyadic variables, attained convergence, we do not present those results in table 4.3.

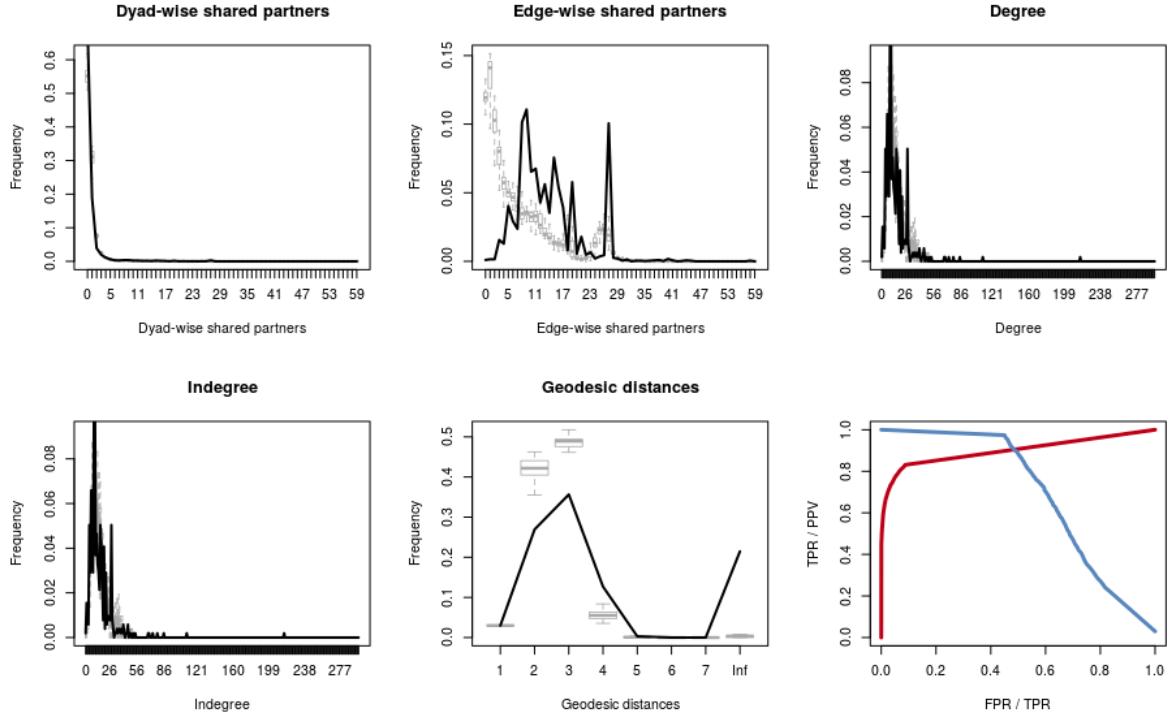


FIGURE 4.10: ERGM goodness-of-fit of final model 3 assessment on the HIV/AIDS co-authorship network.

Figure 4.10 presents the goodness-of-fit of the final model 3. It appears that the ERGM fits well the observed HIV/AIDS co-authorship network in terms of edge-wise, dyad-wise shared partners, degree, geodesic distances, triad census. In addition, 89.9% of the time, model 3 accurately predicted new collaboration ties among the authors ( $AUC = 89.9\%$ ).

### 4.3.2.3 Temporal Exponential Random Graph Model

We subset the cumulative observed network in six snapshots according to the following time spans: 1996 – 2001, 2002 – 2008, 2009 – 2010, 2011 – 2012, 2013 – 2014 and 2015 –

## *Results: The HIV/AIDS Co-authorship Network*

---

2016. In figure 4.11, we show the topological structure of the network snapshots for the different time steps.

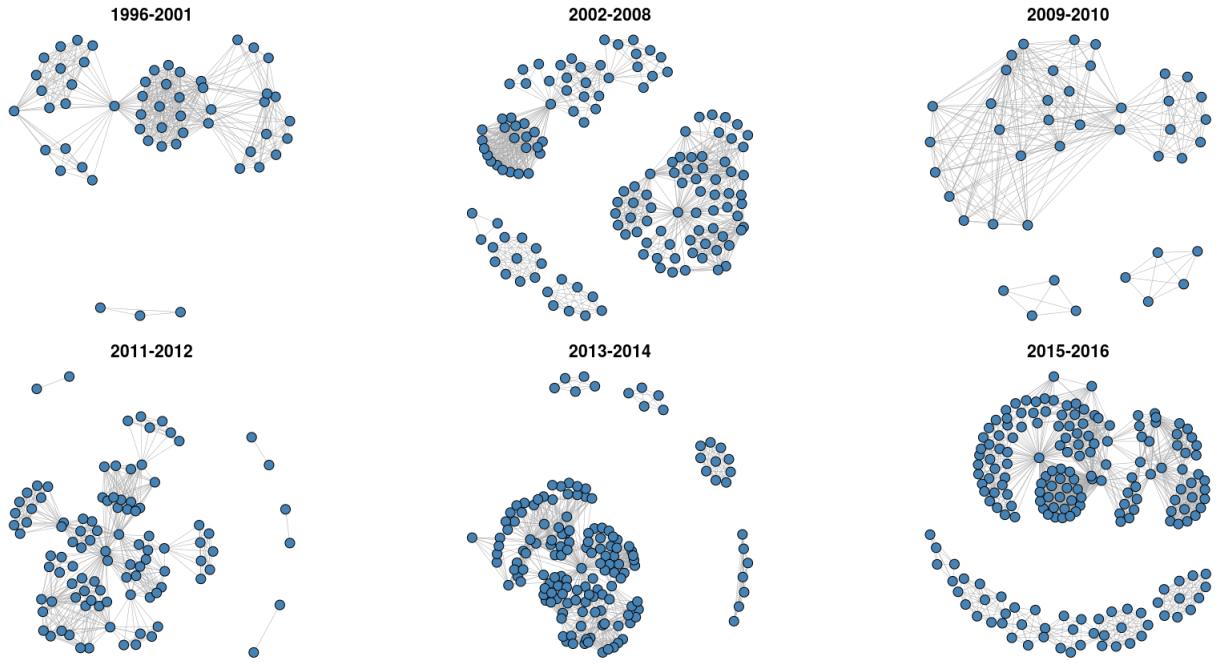


FIGURE 4.11: Topological structure of the different snapshots of the HIV/AIDS co-authorship network.

Table 4.4 summarizes the results of the different temporal models fit to the observed

network. The coefficient for the edge term in the null pooled ERGM model 1 is estimated

at  $-5.18$  with an associated baseline pooled probability of collaboration tie formation of

0.56%. This probability is lower than the density of the cumulative network.

After adjusting for the nodal variables and the homophily terms, model 2 improved

slightly over the null model 1. Model 3 adjusted model 2 by including a factor attribute

effect on the collaboration type with a slight improvement over model 2. Model 3 con-

tains the same terms as the final model of the ERGM in the previous section. Unlike the

final model of the ERGM, we observed here a significant decrease of 33.6% in the odds of

researchers affiliated with Beninese institutions to collaborate compared to international researchers. This effect is maintained after adjusting for the temporal dependencies in model 4.

Model 4 displays a tremendous improvement over model 3, and is hence our final temporal model. The results of model 4 confirm the observation made in section 4.3.2.2 that the process of collaboration tie establishment in the HIV/AIDS network is mainly driven by homophily on collaboration type and on membership to research groups or communities. Both temporal dependencies effects are significant in the final model. We observed a significantly positive dyadic stability effect accompanied with a significantly negative linear trends effect. For dyadic stability, the coefficient is 0.37 meaning that the odds of existent and non existent collaboration ties at one time point to remain the same at the next time point increased on average by 30.9%. In other words, the odds of new collaboration ties and non-ties to occur from one time point to another is 69.1%. Overall, the probability of international authors to establish a stable collaboration tie is 7.94% versus 6.30% and 9.62% respectively for national and regional researchers.

The goodness-of-fit assessment of the final TERGM model 4 is presented in figure 4.12. The first five subfigures comparing the distribution of endogenous network statistics between the observed network and the simulated ones show a good fit of the final model to the observed network data. The AUC of the ROC curve in the six subfigures is 79.9% for model 4 in predicting ties in the last snapshot. While this performance is lower than the performance of the final ERGM model 3 from the previous section, the walktrap and edge betweenness modularity distributions from model 4 predicted well the observed ones.

*Results: The HIV/AIDS Co-authorship Network*

---

TABLE 4.4. Temporal ERGM of the HIV/AIDS Co-authorship Network.

	Model 1	Model 2	Model 3	Model 4
	Estimate (SE)	Estimate (SE)	Estimate (SE)	Estimate (SE)
<b>Network structural predictor</b>				
Intercept(edge)	-5.18 (0.02)***	-8.73 (0.05)***	-8.68 (0.06)***	-7.86 (0.09)***
Number of times cited	-	0.00 (0.00)***	0.00 (0.00)***	0.00 (0.00)***
Number of collaborations	-	0.12 (0.00)***	0.11 (0.00)***	0.10 (0.00)***
Number of publications	-	-0.10 (0.01)***	-0.06 (0.01)***	-0.03 (0.01)
Homophily on cluster assignment	-	4.60 (0.05)***	4.61 (0.05)***	4.46 (0.05)***
Homophily on collaboration type	-	0.52 (0.04)***	0.50 (0.04)***	0.59 (0.04)***
<b>Factor attribute effect (collaboration type)</b>				
International	-	-	REF	REF
National	-	-	-0.29 (0.03)***	-0.25 (0.04)***
Regional	-	-	0.14 (0.03)***	0.21 (0.03)***
<b>Temporal dependencies</b>				
Dyadic stability	-	-	-	0.37 (0.04)***
Linear trends	-	-	-	-0.08 (0.02)***
AIC	5591754.39	5563258.81	5563125.93	3715452.45
BIC	5591781.15	5563352.48	5563246.37	3715595.64
Log Likelihood	-2795875.19	-2781622.40	-2781553.96	-1857715.22

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

Finally, the walktrap community comembership prediction displays an AUC of 80%.

## Results: The HIV/AIDS Co-authorship Network

---

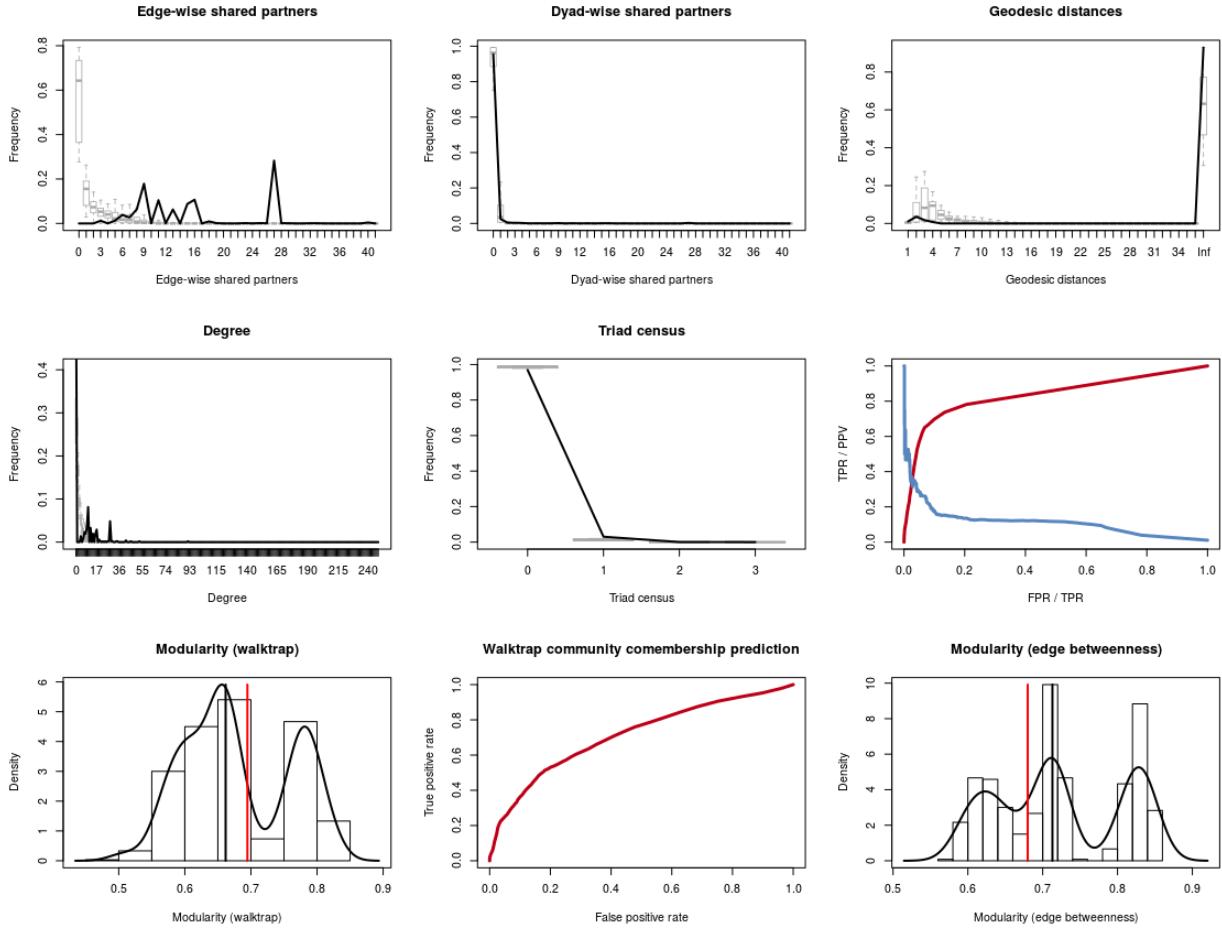


FIGURE 4.12: Goodness-of-fit assessment for the final HIV/AIDS TERGM Model 4 with temporal dependencies of the HIV/AIDS co-authorship network.

### 4.3.2.4 Latent Network Model

In figure 4.13, we present the 3-dimensional visualization of the HIV/AIDS co-authorship with layouts determined according to the inferred latent eigenvectors from the no pair-specific model (on top), the model containing nodal covariates (middle), and the model containing nodal and dyadic covariates (bottom). Blue vertices represent authors affiliated to Beninese research institutions, Red vertices are authors affiliated to international institutions, Gold vertices represent authors affiliated to African research institutions

other than Benin, and White vertices represent authors with no determined affiliations.

Node sizes are set to be proportional to the betweenness value of each vertex, with bigger nodes emphasizing key broker authors in the network.

The first visualization represents the LNM model with no pair-specific covariates. It shows mainly two clusters with little demarcation. We can see that there is a heterogeneity with regard of the spatial distribution of the nodes. After adjusting for the nodal covariates (second visualization), the clustering of the nodes appears less apparent. This results seem to suggest the non-significant role of geography in the establishment of collaboration ties in the HIV/AIDS co-authorship network.

After adding dyadic variables to the model, the resulting visualization shows that there is less structure left to be captured by the latent variables. This observation can explain the failure of our ERGM and TERGM containing dyadic covariates to converge. It also confirms our ERGM and TERGM findings.

We assess the goodness-of-fit of the LNM models. The ROC curves on figure 4.14 shows that the LNM model containing the nodal covariates ( $AUC = 0.966$ ) outperforms the null model ( $AUC = 0.898$ ).

*Results: The HIV/AIDS Co-authorship Network*

---

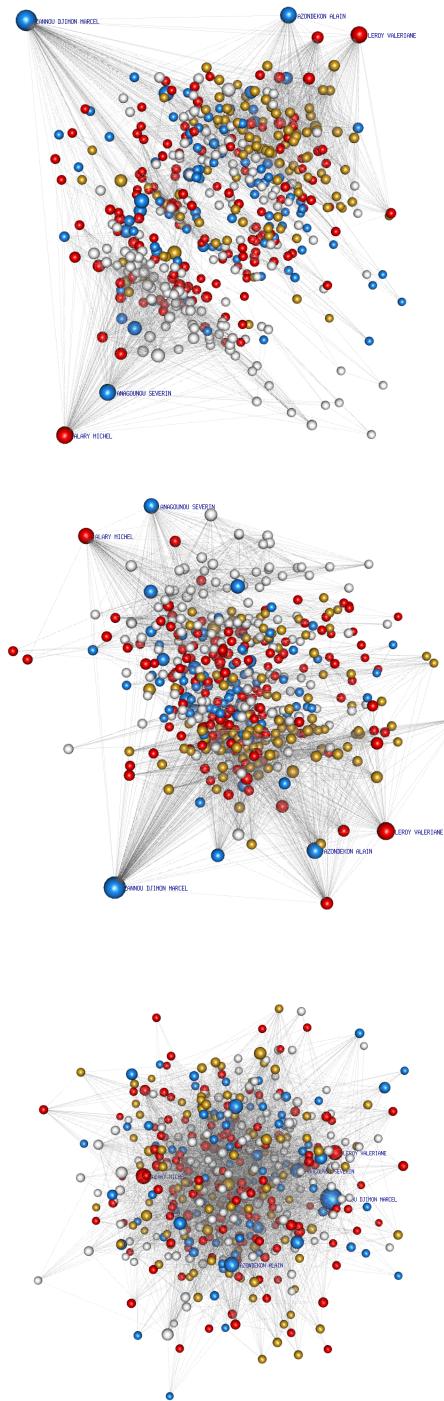


FIGURE 4.13: Visualizations of the HIV/AIDS co-authorship network with layouts determined according to the inferred latent eigenvectors in the LNM models.

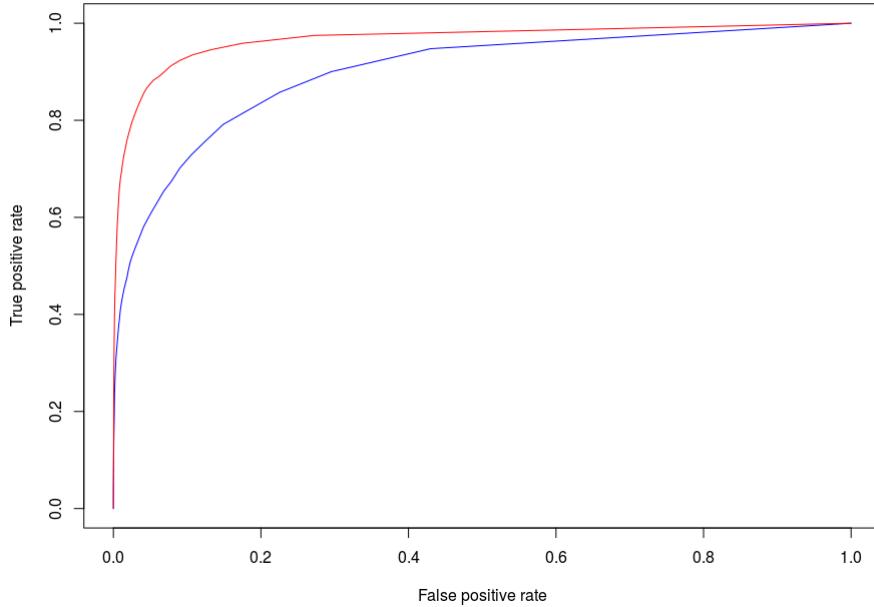


FIGURE 4.14: ROC curves comparing the goodness-of fit of the HIV/AIDS co-authorship network for the model specifying (i) no pair specific covariates (blue) and the model specifying (ii) nodal covariates (red).

## 4.4 Discussion and Conclusion

This chapter deciphers the HIV/AIDS co-authorship network over the last 20 years. The results from the descriptive analyses in this chapter are similar to the descriptive analyses results from chapter 3. As the malaria co-authorship network, the HIV/AIDS co-authorship network in Benin is a complex network, as it exhibits unexpected properties that are more extreme than the 4 mathematical models used for the Monte-Carlo based simulations. The observed characteristics disproved previous studies supporting the idea that co-authorship have small world properties [55] or are preferential attachment networks [109]. In fact, unlike our methodology, those studies mainly used descriptive methods and did not apply advance statistical methods to test their network properties.

---

*Results: The HIV/AIDS Co-authorship Network*

---

The HIV/AIDS co-authorship network in Benin has a low density with a highly right-skewed node degree distribution. Compared to the malaria co-authorship network, the relatively low transitivity provides evidence of less hierarchy - well connected authors in this network tend to connect with poorly connected ones. This also indicates that this network is less assortative than the malaria co-authorship network, with prolific and non-tenure authors connected similar authors. As in Salamati and Soheili [70], The flow of information in the HIV/AIDS network in Benin is slow as it only relies on 8 authors representing less than 1% of all the authors in the network. The removal of these authors from the network would lead to its collapse. Such a structural vulnerability is not just inherent to the HIV/AIDS co-authorship network, as it is a global observation already reported elsewhere [66]. Since the mathematical models we applied fell short to thoroughly explain the mechanistic phenomenon explaining the growth and the structure of the network, we suspect hidden factors which we attempted to model using advanced statistical models. As our first modeling approach, the SBM identified heterogeneous classes with no dominance of regional, national or international players, despite a much higher likelihood of Sub-Saharan African countries to collaborate with non-African states [110].

Based on the results from our ERGM and TERGM models, in the HIV/AIDS co-authorship network authors are more likely to establish collaboration ties within their research groups or communities. Unfortunately, we were not able to control for transitivity as all the models adjusting for this term failed to converge. We suspect the size and complexity of this network to have prevented the convergence of such models, even after 1,000 iterations [111].

---

*Results: The HIV/AIDS Co-authorship Network*

---

Although marginal, factors such as number of publications, number of citations and number of collaborations are associated to higher likelihood to establishing collaboration ties, confirming therefore our first hypothesis. Adding temporal dependencies to our ERGM models tremendously improved the fitness of the model to the observed network data, but at a cost of decreased performance compared to the model without temporal dependencies.

The LNM complements the ERGM and TERGM by adding another layer of analysis. With the LNM, we are able to visualize the effect of geography on the structure of the network. The lack of clear cluster demarcation suggests that distance does not play a significant role in collaboration tie formation in the HIV/AIDS network.

Our results confirm that the regain in HIV/AIDS research funding has led to a significant increase in publications number and research collaborations in Benin. In order to consolidate the knowledge generated, there is an urgent need to reinforce the HIV/AIDS research network in Benin given its vulnerability and uneven integration. Identified key brokers and most productive authors need to continuously be supported, and identified junior scientists in the field be promoted.

# Chapter 5

## Results: The Tuberculosis Co-authorship Network

### 5.1 Data

The literature search was conducted in the Web Of Science using combinations of TB related MeSH terms including "Tuberculosis", "Mycobacterium", "Infection". We restricted the search to the period from 1996 to 2016 and to "Benin" for country. We further screened the papers in order to only select those published by Beninese authors, or papers published on TB from Benin. All published documents under considerations included at least one Author from Benin. No restriction was placed upon the document types. We first started querying with each term independently, we then combined the other terms so the query return the maximum number of results. The Full citations

## *Results: The Tuberculosis Co-authorship Network*

---

information containing the authors' names, their institutional affiliations, the year of publication, as well as the number of times the document was cited were recorded as a bibliographic corpus in text format. After a second screening only research that have met the above listed inclusion criteria and that were published between January 1, 1996 and December 31, 2016 were selected.

The final query set (Table 5.1) returned 109 records. After a rigorous screening process carried out by all the authors, 37 documents met the selection criteria. On average, there was 9.38 authors per published document.

TABLE 5.1. TB Bibliographic Search Queries.

Set	Queries	Results
#1	TOPIC: (Tuberculosis) AND ADDRESS: (BENIN)	109
#2	TOPIC: (Tuberculosis), Refined by: COUNTRIES/TERRITORIES: (BENIN )	77
#3	TOPIC: (Mycobacterium Tuberculosis), AND ADDRESS: (Benin)	77
#4	TOPIC: (Tuberculosis) OR TOPIC: (Infection) AND ADDRESS: (BENIN), Refined by: COUNTRIES/TERRITORIES: (BENIN)	89
Final Set	#1 OR #2 OR #3 OR #4	109

After the Author Name Disambiguation, we identified 173 unique authors with a precision of 99.99% and a recall of 99.99%. The generated multigraph co-authorship network therefore contained 173 vertices (authors) and 1,937 parallel edges (collaborations). As displayed in figure 5.1, we can see the significant increase in publications, scientific collaborations and the number of authors involved in TB research from 2008 until 2016. This

general upward trend seems to be linear from the year 2008 to 2016.

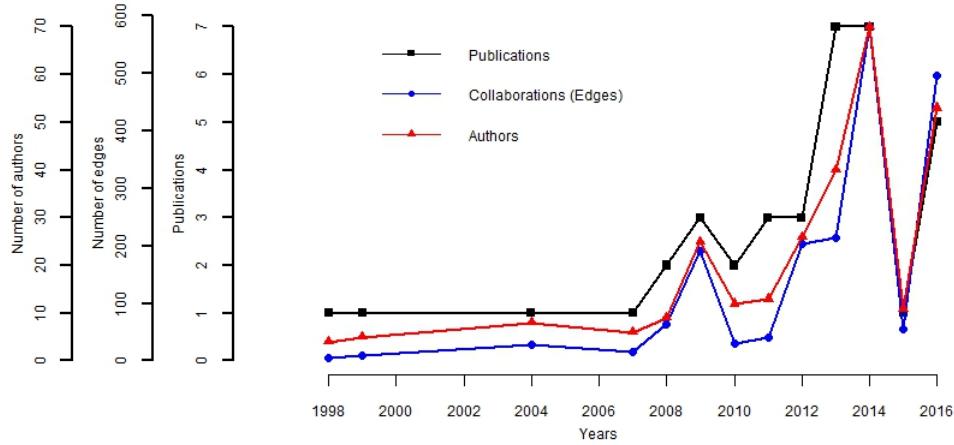


FIGURE 5.1: Evolution of the published TB related documents, authors and collaborations from January 1996 to December 2016

## 5.2 Descriptive Data Analysis

For the multigraph network, the degree distribution ranged between 2 and 165 with an average degree distribution of 17.36 and a median of 15. In addition, there was a substantial number of vertices with low degrees (Fig. 5.2). The log scale distribution of the degrees on figure 5.3 reveals that there was a tendency of prolific authors to collaborate with less prolific authors.

After converting the multigraph network in a weighted graph, the network results in a simple graph of 173 vertices and 1,502 weighted edges. Closeness centrality ranges between  $3.68 \times 10^{-5}$  and  $3.28 \times 10^{-4}$  with a median of  $3.18 \times 10^{-4}$ . Betweenness measures

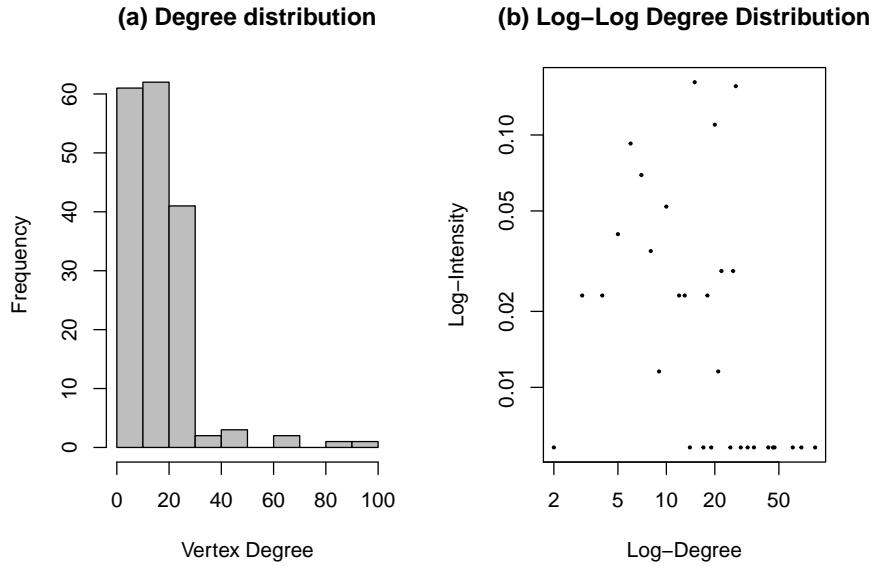


FIGURE 5.2: Degree distribution of the TB Co-authorship network

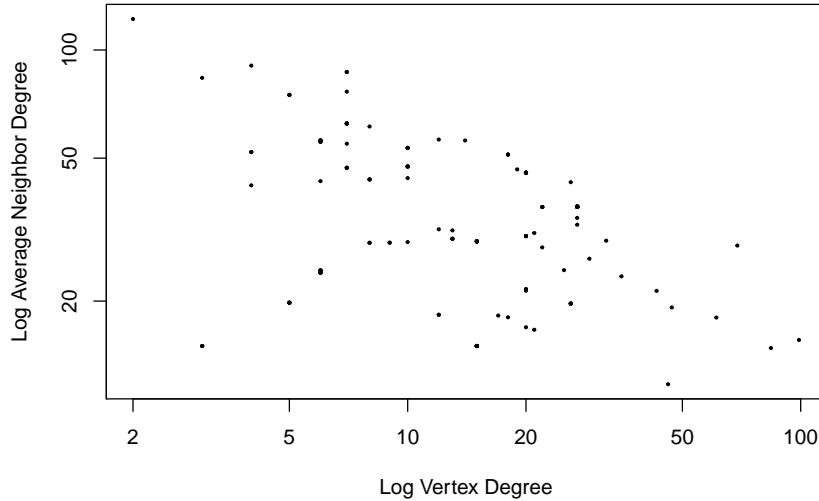


FIGURE 5.3: Log-Average Neighbor degree Distribution of the TB Co-authorship network

range between 0 and 3,077 with a median of 12.49. A network visualization with the vertices' size proportional to betweenness centrality measures clearly reveals the presence of broker authors (Figure 5.4 and Table 5.2). The median Eigenvectors is 0.087 and a

mean of 0.138. Eigenvectors measures reveal the presence of author hubs in the network suggesting the presence of closed collaboration groups. Table 5.2 presents a list of the 10 author hubs with the highest Eigenvectors values.

Edge betweenness centrality measures identify co-authorship collaboration ties that are important for the flow of information. Table 5.2 presents the top 10 most important collaboration ties for the flow of information in the TB Co-authorship network in Benin.

### 5.2.1 Network Cohesion

Overall, 28 maximal cliques were detected in the network among which 1 clique of size 10, 2 cliques of size 5, and 2 cliques of size 4. The largest clique has size 10.

The TB co-authorship network has a density of 0.10095 indicating that the baseline probability of collaboration tie formation is 10.095%. The network also has a transitivity of 0.6305 meaning that 63.05% of the connected triples in the network are close to form triangles. The transitivity metrics is a measure of the global clustering of the network.

The network is not connected and a census of all the connected components within the network reveals the existence of a giant component that dominates all the other connected components. The giant component of the TB co-authorship network includes 90.8% (157 vertices) of all the vertices in the network with the other components alone carrying less than 0.1% of the vertices in the network (Fig. 5.4).

*Results: The Tuberculosis Co-authorship Network*

---

TABLE 5.2. List of the most important authors and collaborations in the Tuberculosis Co-authorship network

<b>Top 10 Brokers</b>
AFFOLABI DISSOU
GNINAFON MARTIN
DE JONG BOUKE C
TREBUCQ ARNAUD
ODOUN MATHIEU
ANAGONOU SEVERIN
WACHINOU PRUDENCE
FAIHUN FRANK
KASSA FERDINAND
ADE SERGE
<b>Top 10 most connected authors (Top 10 network hubs)</b>
GNINAFON MARTIN
AFFOLABI DISSOU
ANAGONOU SEVERIN
MERLE CORINNE S C
TREBUCQ ARNAUD
OLLIARO PIERO L
RUSTOMJEE ROXANA
LO MAME BOCAR
LIENHARDT CHRISTIAN
HORTON JOHN
<b>Top 10 most important edges for information flow</b>
ODOUN MATHIEU – GNINAFON MARTIN
FAIHUN FRANK – DE JONG BOUKE C
ODOUN MATHIEU – TREBUCQ ARNAUD
ZELLWEGER J P – GNINAFON MARTIN
TREBUCQ ARNAUD – ADJONOU CHRISTINE
ODOUN MATHIEU – WACHINOU PRUDENCE
AFFOLABI DISSOU – BAHSOW OUMOU
AFFOLABI DISSOU – TOUNDOH N
AFFOLABI DISSOU – BEKOU W
AFFOLABI DISSOU – MAKPENON A
<b>Weak articulation point</b>
WACHINOU PRUDENCE

---

Information flow assessment of the network via cut vertices reveals the existence of a single author as the most vulnerable vertex in the network (Table 5.2). The cut vertex constitute the weak articulation point of the TB co-authorship network. Cut vertices represent a measure of the vulnerability of the network [83].

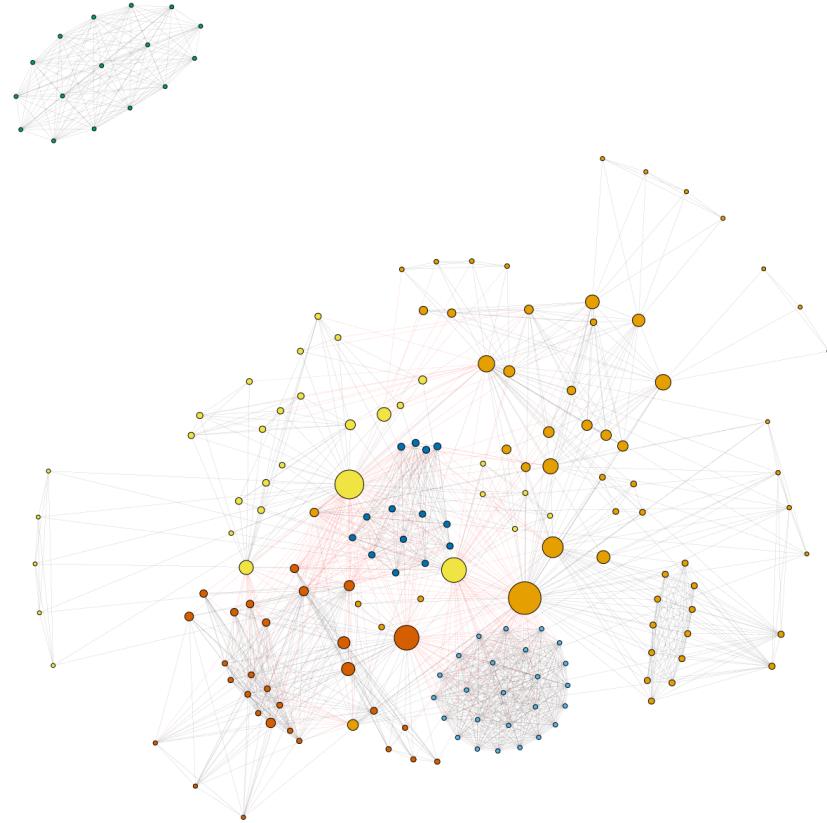


FIGURE 5.4: Topological Structure of the Tuberculosis Co-authorship Network. Authors (vertices) of the same color belong to the same research community or cluster

The agglomerative hierarchical clustering method identifies 6 different research communities (or clusters) in the network. Sizes of the clusters range between 14 and 58 authors. Out of the 6 research clusters or communities detected, 5 are in the giant component. Figure 5.4 displays the giant component of the network with each different colors representing each of the 5 research communities.

## 5.3 Modeling

### 5.3.1 Mathematical Modeling

From the hierarchical clustering method of community detection, 6 different clusters/-communities were detected in the co-authorship network out of which 5 form a giant component. One of the question of interest in this section is whether the number of communities detected is expected or not. To provide an answer to this question, we performed 1,000 Monte Carlo based simulations to test the significance of this observed characteristics of the TB co-authorship network. Figure 5.5 clearly demonstrates that the number of communities detected is unusual from the perspective of both Classical random graphs and generalized random graphs ( $p\text{-value} < 0.0001$ ). From the Classical random graph model, the expected number of communities was 4.734 (95%CI: 4.70 – 4.77). Similarly, the expected number of communities from the generalized random graph model is 5.34 (95%CI: 5.29 – 5.38).

Figure 5.6 displays the number of detected research communities using the Barabási-Albert's preferential attachment and the Watts-Strogatz models. Here too, the observed number of communities was extreme per both models ( $p\text{-value} < 0.0001$ ). The expected number from the Watts-Strogatz model simulations is 3.017 (95%CI: 3.01 – 3.03) and 13.77 (95%CI: 13.70 – 13.85) from the Barabási-Albert model simulations.

We also compared the clustering coefficient and the average shortest-path length. Let's

*Results: The Tuberculosis Co-authorship Network*

---

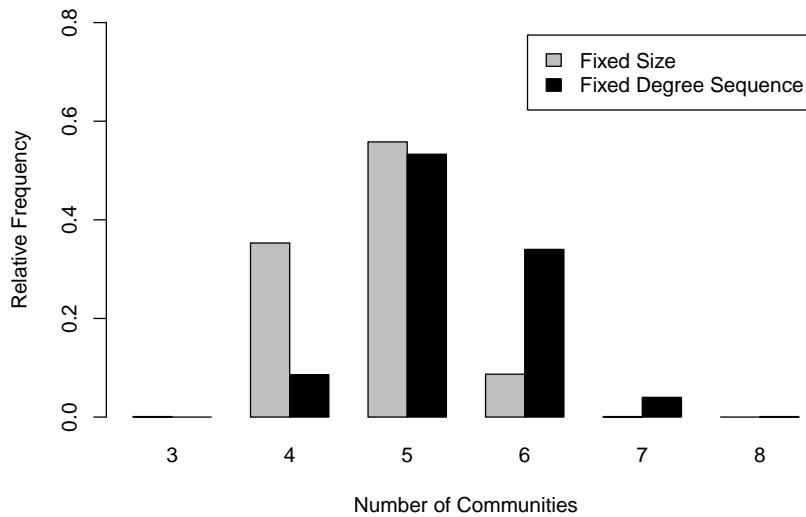


FIGURE 5.5: Monte-Carlo simulations of the TB network: Number of detected communities by the random graph models

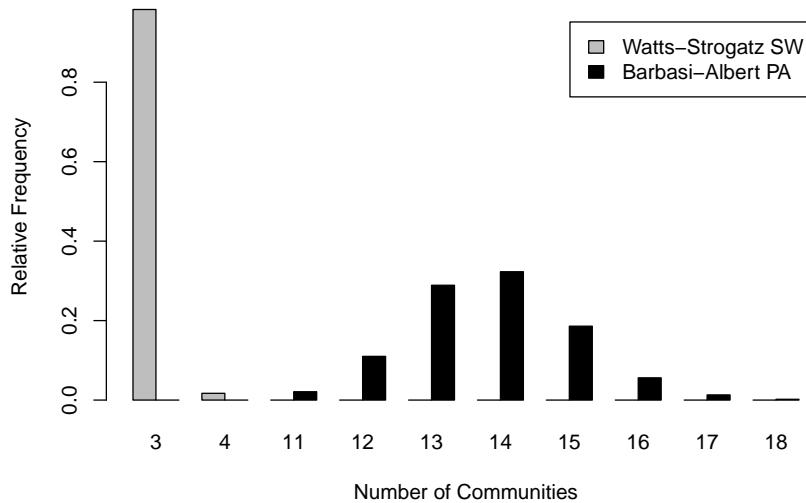


FIGURE 5.6: Monte-Carlo simulations of the TB network: Number of detected communities by the Watts-Strogatz and the Barabási-Albert models

recall that the observed clustering coefficient is 0.614. On one hand, there was substantially more clustering in our TB co-authorship network than expected from both random graph models ( $p\text{-value} < 0.0001$ ). The expected clustering coefficients was 0.10087

*Results: The Tuberculosis Co-authorship Network*

---

(95%CI: 0.10068 – 0.10107) and 0.1937 (95%CI: 0.1934 – 0.1939) respectively for the classic random graph and the generalized random graph models.

On the other hand, there was substantially less clustering in our TB co-authorship network than expected the Watts-Strogatz Small World model which expected clustering was 0.7259 (95%CI: 0.7258 – 0.7260).

We observed an average shortest-path length of 2.126 in the TB co-authorship network. This observed shortest-path length is significantly larger than what was expected from the random graph models ( $p\text{-value} < 0.0001$ ) and significantly lower than what was expected from Watts-Strogatz small world model and the Barabási-Albert preferential attachment model ( $p\text{-value} < 0.0001$ ).

The average shortest-path length was 2.0548 (95%CI: 2.0546 – 2.0550) and 2.072 (95%CI: 2.0715 – 2.0726) respectively for the classic random graph and the generalized random graph models.

For the Watts-Strogatz small world and the Barabási-Albert models, the average shortest-path length is respectively 2.623 (95%CI: 2.616 – 2.631) and 6.06 (95%CI: 6.03 – 6.09).

We performed the same simulations on the giant component of the network with similar results leading to the same outcomes.

### 5.3.2 Statistical Modeling

#### 5.3.2.1 Stochastic Block Model

The SBM identifies 14 classes with a degree of latitude of 9 to 14 classes being reasonable (See ICL plot on figure 5.7).

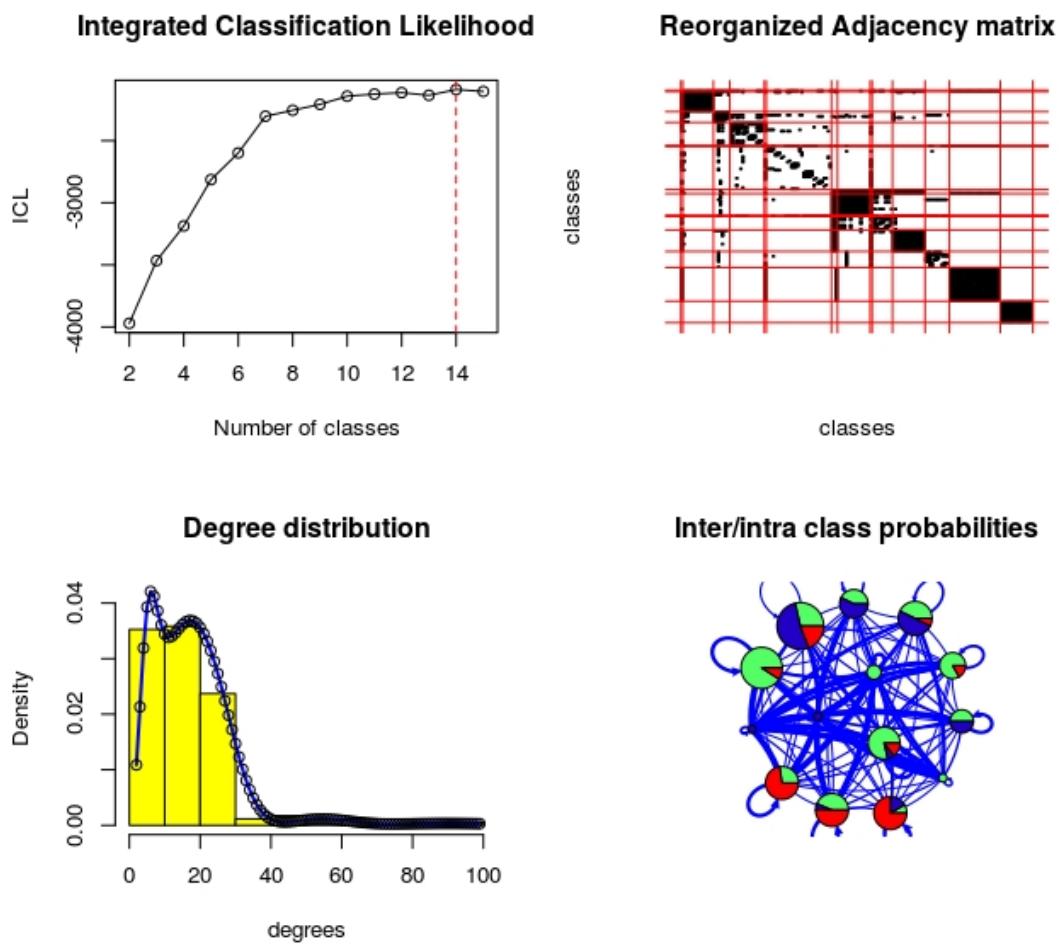


FIGURE 5.7: Summary of the goodness-of-fit of the SBM analysis on the Tuberculosis co-authorship network.

The fitted SBM describes well the observed degree distribution. The vertices in the network depicting the inter/extraclass probabilities represent the 14 identified classes, with each

## *Results: The Tuberculosis Co-authorship Network*

---

one of them divided into a pie chart displaying the proportion of authors of international affiliations (lightgreen), authors of regional or other African affiliations (red), and authors affiliated to Beninese research institutions (blue). Generally, the dominance across the classes of international and regional players is observed. From the inter/intra probability network shows denser ties inter class ties. Looking at the pie charts, we can see that the classes are heterogeneous of almost all the classes with most of the classes having the same sizes (5.7). Figure 5.8 presents the distribution of the classes by affiliation types.

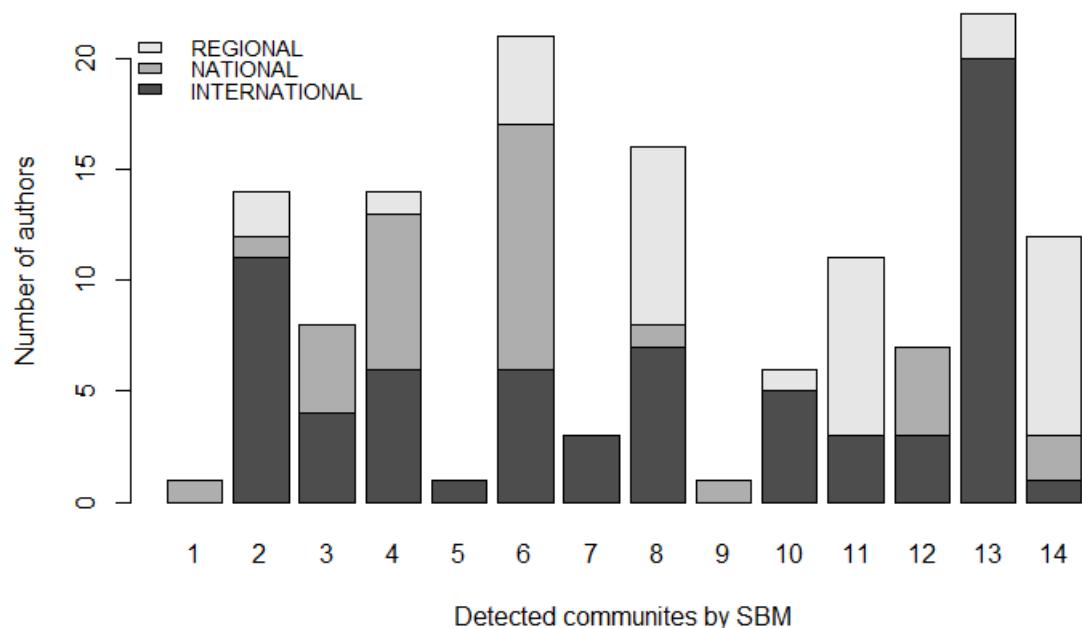


FIGURE 5.8: Distribution of national, international and regional authors by communities detected by the SBM in the TB network.

*Results: The Tuberculosis Co-authorship Network*

---

	Model 1	Model 2	Model 3
	Estimate (SE)	Estimate (SE)	Estimate (SE)
Network structural predictor			
Intercept(edge)	-2.19 (0.03)***	-7.84 (0.16)***	-7.86 (0.17)***
Number of times cited	-	0.01 (0.00)***	0.01 (0.00)***
Number of collaborations	-	0.08 (0.00)***	0.07 (0.00)***
Number of publications	-	-0.05 (0.01)**	0.01 (0.02)
Homophily on cluster assignment	-	6.02 (0.13)***	6.12 (0.14)***
Homophily on collaboration type	-	0.83 (0.10)***	0.90 (0.10)***
Factor attribute effect (collaboration type)			
International	-	-	<i>REF</i>
National	-	-	-0.40 (0.09)***
Regional	-	-	0.22 (0.08)**
AIC	9737.42	3776.48	3747.34
BIC	9745.03	3822.12	3808.20
Log Likelihood	-4867.71	-1882.24	-1865.67

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

TABLE 5.3. ERGM of the TB Co-authorship Network.

### 5.3.2.2 Exponential Random Graph Model

We fit multiple ERGM method (Table 5.3). In the null model (model 1), the inverse logit of the coefficient associated with the intercept (edge term) is 0.10 which is the baseline probability of collaboration ties establishment and also the density of the TB co-authorship network.

Model 2 including all nodal variables, a homophily term on collaboration type and on cluster assignment improved tremendously compared to model 1 (See AIC, BIC and model likelihood in table 5.3). We note a decrease in the edge effect (Coefficient = -7.84,  $p < 0.001$ ) with the associated conditional probability (given all the other terms in the model) estimated at 0.039%. For the remaining terms in model 2, we observed a positive

## *Results: The Tuberculosis Co-authorship Network*

---

and significant effect except for the number of publications. Model 3 including the collaboration type as factor term, improved substantially compared to model 2. We therefore chose model 3 as our final model. One unit increases the number of citation, increases the odds of collaboration ties establishment by 1%. A one unit increase in the number of collaborations is associated with a 7.25% increase in the odds of collaboration ties establishment. The coefficient associated with the number of publications is insignificant. Model 3 further proves that the process underlying the structure of the TB co-authorship network in Benin is mainly driven by homophily on cluster assignment or membership to a research community or group (Coefficient = 6.12,  $p < 0.001$ ). The conditional probability of any two authors belonging to the same research group is estimated at 14.93% compared to the baseline probability of 10%. The same probability changes to 30.15% after adjustment by the collaboration type, and 32.08% after adjusting for the number of citations, collaborations and publications. Compared to research affiliated to international institutions, researchers affiliated to Beninese institutions have 49.2% average decrease in the odds of collaboration tie establishment. This average decrease is not statistically significant ( $p > 0.05$ ). For researchers affiliated to institutions other than Beninese institutions, the odds of collaboration tie establishment increases on average by 24.05% compared to internationally affiliated researchers. Overall, model 3 estimated the probability of collaboration ties formation at 32.08% for international researchers, 24.05% for national researchers and 37.05% for regional players.

Unfortunately, none of the models containing endogenous ERGM terms and/or the dyadic variables, attained convergence, we do not present those results in table 5.3.

## Results: The Tuberculosis Co-authorship Network

---

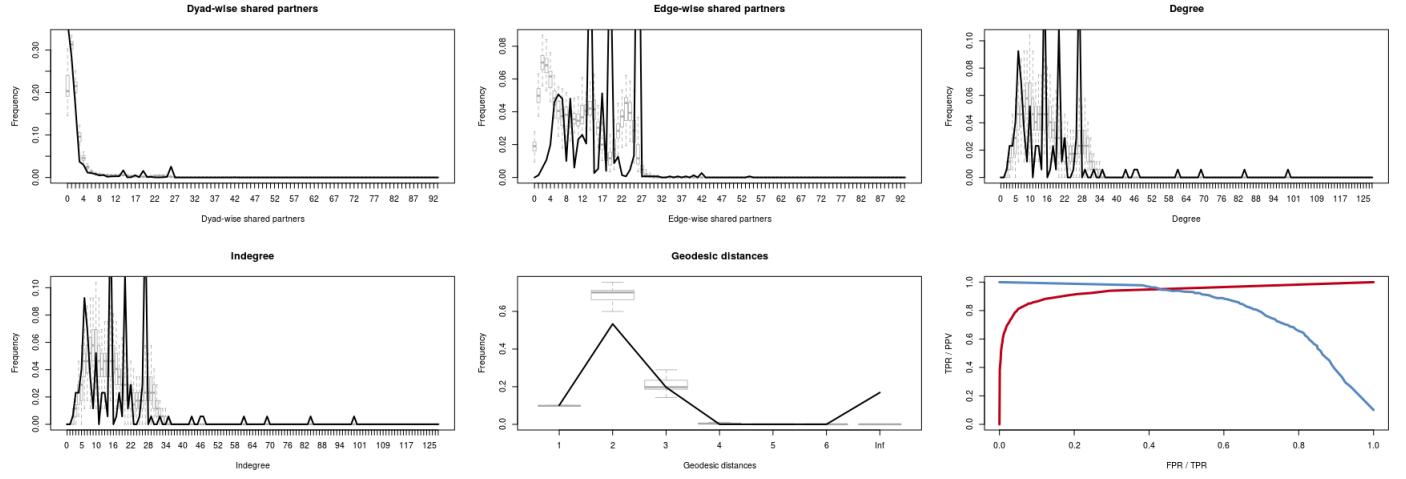


FIGURE 5.9: ERGM goodness-of-fit of final model 3 assessment on the TB co-authorship network.

Figure 5.9 presents the goodness-of-fit of the final model 3. It appears that the ERGM fits somewhat poorly the observed TB co-authorship network in terms of edge-wise, dyad-wise shared partners, degree, geodesic distances, triad census. Meanwhile, it displays a 93.7% for the ROC model (in red) and 80.9% for the Precision Recall (PR) model.

### 5.3.2.3 Temporal Exponential Random Graph Model

We subset the cumulative observed network in cinq snapshots according to the following time spans: 1996 – 2008, 2009 – 2011, 2012 – 2013, 2014 – 2015 and 2016. In figure 5.10, we show the topological structure of the network snapshots for the different time steps.

Table 5.4 summarizes the results of the different temporal models fit to the observed network. The coefficient for the edge term in the null pooled ERGM model 1 is estimated at  $-3.75$  with an associated baseline pooled probability of collaboration tie formation of 2.30%, which is lower than the density of the observed cumulative TB network.

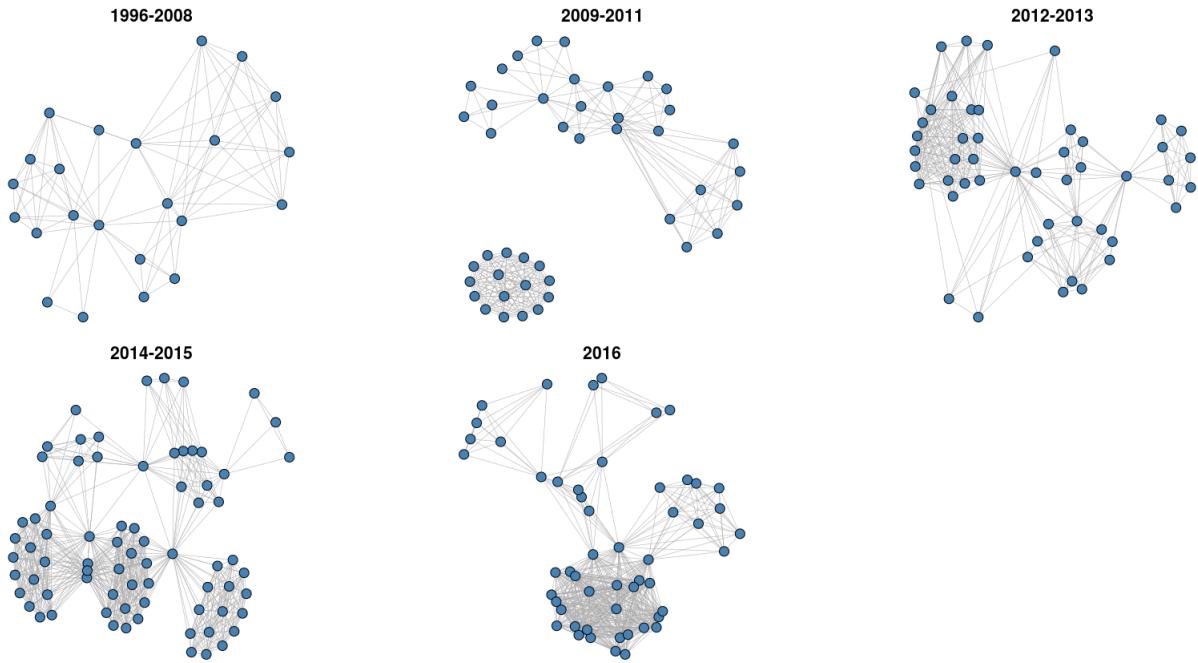


FIGURE 5.10: Topological structure of the different snapshots of the TB co-authorship network.

Model 2 adjusting for the nodal variables and the homophily terms improved slightly over the null model 1. Model 3 adjusted model 2 by including a factor attribute effect on the collaboration type with a slight improvement over model 2. Unlike the final model of the ERGM, we observed in model 3, a significant decrease of 23.4% in the odds of researchers affiliated with Beninese institutions to collaborate compared to international researchers. This percentage decrease changes to 40.5% after adjusting for the temporal dependencies in model 4.

We chose Model 4 as our final model because it significantly improved over model 3. The results of model 4 confirm our observation from the ERGM results that the process of collaboration tie establishment in the TB network is mainly driven by homophily on collaboration type and on membership to research groups or communities.

### *Results: The Tuberculosis Co-authorship Network*

---

Temporal dependencies effects proved significant in the final model. A significantly positive dyadic stability effect accompanied with a significantly negative linear trends effect is observed. For dyadic stability, the coefficient is 0.44 meaning that the odds of existent and non existent collaboration ties at one time point to remain the same at the next time point increased on average by 35.6%. In other words, the odds of new collaboration ties and non-ties to occur from one time point to another is 64.4%. Overall, the probability of international authors to establish a stable collaboration tie is 15.71% versus 11.71% and 16.11% respectively for national and regional researchers.

The goodness-of-fit assessment of the final TERGM model 4 is presented in figure 5.11. Regarding the endogenous network statistics, we observe a better fit of the final TERGM model 4 compared to the final ERGM model 3. In other words, the simulated network by model 4 show a good fit to the observed TB network data. The AUC of the ROC curve of model 4 is estimated at 83.2% meaning that 83.2% of the times, model 4 accurately predicting ties in the last snapshot. While this performance is lower than the performance of the final ERGM model 3 from the previous section, the walktrap and edge betweenness modularity distributions from model 4 predicted well the observed ones. Finally, the walktrap community comembership prediction displays an AUC of 71.4%.

## Results: The Tuberculosis Co-authorship Network

TABLE 5.4. Temporal ERGM of the TB Co-authorship Network.

	Model 1	Model 2	Model 3	Model 4
	Estimate (SE)	Estimate (SE)	Estimate (SE)	Estimate (SE)
Network structural predictor				
Intercept(edge)	-3.75 (0.02)***	-10.07 (0.15)***	-10.01 (0.16)***	-8.62 (0.28)***
Number of times cited	-	0.00 (0.00)*	0.00 (0.00)	-0.00 (0.00)**
Number of collaborations	-	0.14 (0.00)***	0.14 (0.00)***	0.16 (0.00)***
Number of publications	-	0.68 (0.03)***	0.72 (0.03)***	0.57 (0.03)***
Homophily on cluster assignment	-	5.24 (0.11)***	5.23 (0.11)***	5.40 (0.13)***
Homophily on collaboration type	-	0.69 (0.08)***	0.69 (0.08)***	0.73 (0.09)***
Factor attribute effect (collaboration type)				
International	-	-	REF	REF
National	-	-	-0.21 (0.07)**	-0.34 (0.08)***
Regional	-	-	0.03 (0.07)	0.03 (0.08)
Temporal dependencies				
Dyadic stability	-	-	-	0.44 (0.07)***
Linear trends	-	-	-	-0.36 (0.06)***
AIC	431184.00	419860.54	419853.82	253170.25
BIC	431205.66	419936.36	419951.30	253284.48
Log Likelihood	-215590.00	-209923.27	-209917.91	-126574.12

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

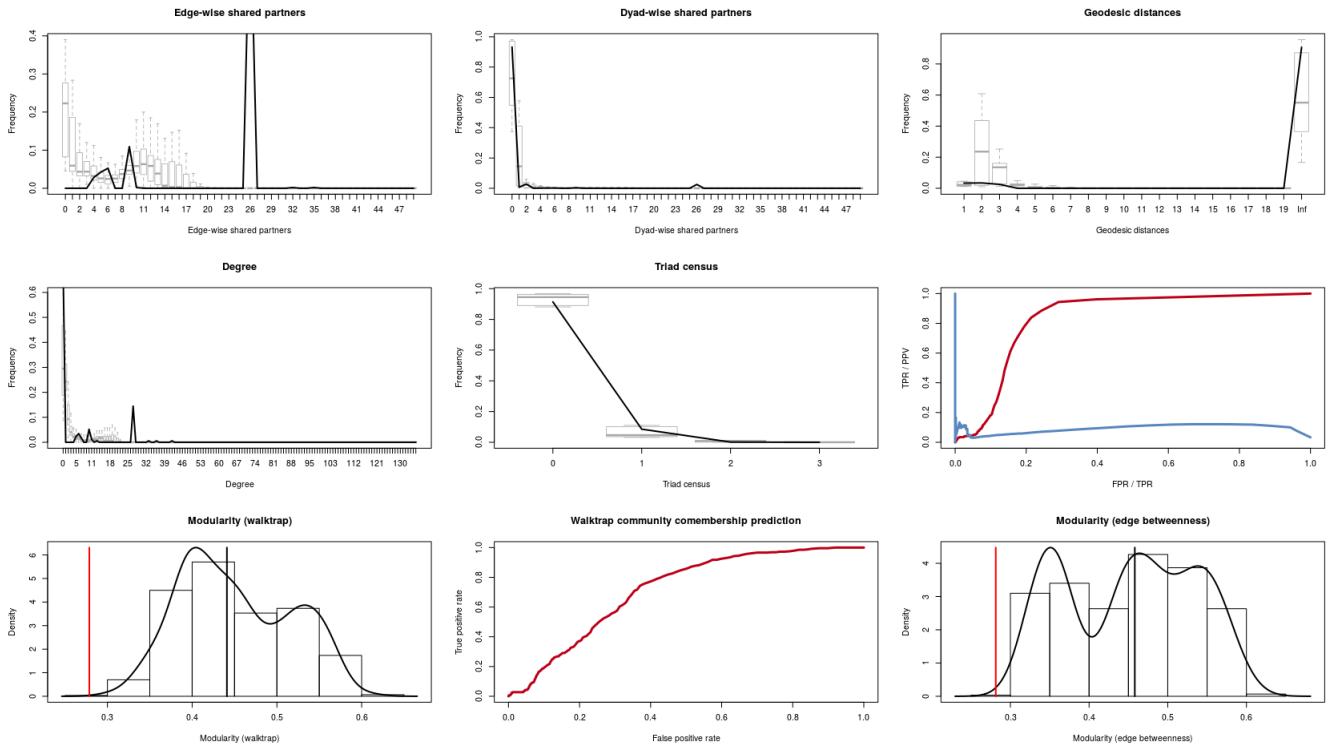


FIGURE 5.11: Goodness-of-fit assessment for the final TB TERGM Model 4 with temporal dependencies of the TB co-authorship network.

#### 5.3.2.4 Latent Network Model

On the 3-dimensional visualization of the TB co-authorship network presented on figure 5.12, the layouts are determined according to the inferred latent eigenvectors from the no pair-specific model (on top), the model containing nodal covariates (middle), and the model containing nodal and dyadic covariates (bottom). Blue vertices represent authors affiliated to Beninese research institutions, Red vertices are authors affiliated to international institutions, Gold vertices represent authors affiliated to African research institutions other than Benin, and White vertices represent authors with no determined affiliations. Node sizes are set to be proportional to the betweenness value of each vertex, with bigger nodes emphasizing key broker authors in the network.

The first visualization represents the null LNM model with no pair-specific covariates. It shows mainly three clusters. The largest cluster appears more spatially heterogeneous than the other two. It is also the largest cluster that contains the majority of the authors affiliated with Beninese research institutions. The other two clusters seem to be dominated respectively by international and regional researchers. This model displays fits reasonably well to the observed TB network ( $AUC = 0.912$ ). This observation suggest a significant effect of geography in the odds of collaboration tie establishment. After adjusting for the nodal covariates (second visualization), there is less structure left to be captured by the latent variables and the clustering is no more apparent. Adding dyadic attributes to the model leads to similar outcome despite an increase in terms of performance ( $AUC = 0.974$ ).

*Results: The Tuberculosis Co-authorship Network*

---

On figure 5.13, we present the ROC curves of each of the LNM models containing the nodal covariates and the null model.

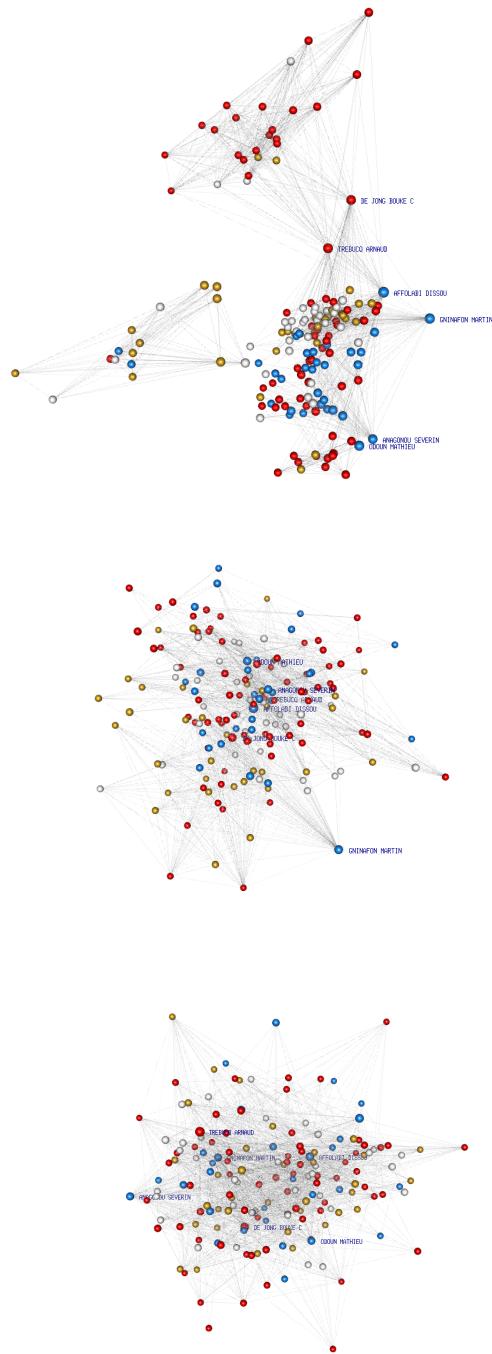


FIGURE 5.12: Visualizations of the TB co-authorship network with layouts determined according to the inferred latent eigenvectors in the LNM models.

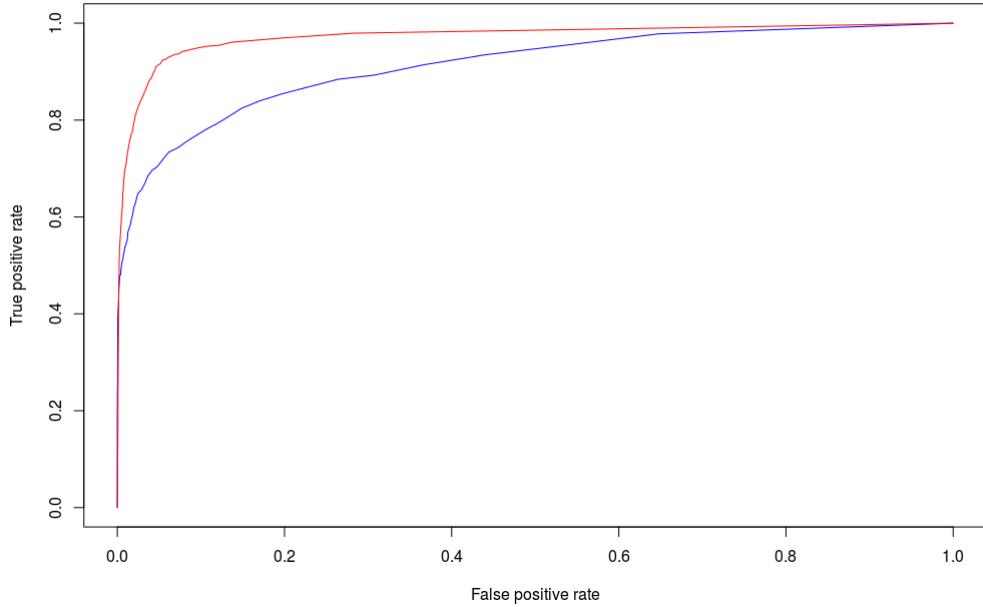


FIGURE 5.13: ROC curves comparing the goodness-of fit of the TB co-authorship network for the model specifying (i) no pair specific covariates (blue) and the model specifying (ii) nodal covariates (red).

## 5.4 Discussion and Conclusion

This chapter provides insights in the structural characteristics of the TB co-authorship network in Benin over the last 20 years. The evolution of the number of publications, authors and collaboration ties suggests a linear growth over the investigation period. We expected such findings given the place of TB in the public health concerns of Benin and the intensive effort towards the reduction of the incidence and the numerous campaigns of sensitization [124]. The findings from the descriptive analysis suggest that the mechanism underlying the formation of the TB co-authorship network in Benin is not random. However, we found inconclusive evidence of small world properties that further Monte-Carlo simulations disproved. The presence of closed research groups is suspected given

## *Results: The Tuberculosis Co-authorship Network*

---

the non-trivial number of authors with higher order of magnitudes. The observed trend of prolific authors in the TB network to collaborate with less prolific ones is another indication suggesting that TB research is a low productivity research field in Benin. Only 37 published documents were found relevant to the present study. In fact, none of the top 10 key brokers in our TB co-authorship network, was in the list of the top most connected authors and therefore would suggest the relative absence of long publishing tenure authors in the network [106].

The flow of information in the TB co-authorship network in Benin is slow as it only relies on a single author. A study by Salamatia and Soheili [70] on a co-authorship analysis of Iranian researchers in the field of violence reported similar but less extreme findings. For Bales et al. [107, 108], the most important authors in co-authorship networks generally tend to be the ones with the highest degree of collaborations. For information flow, cut vertices provide a better approach to identifying vertices that are important to the long-term sustainability of co-authorship networks [83]. The only author identified as a cut vertex is therefore the most important author for information flow.

Our observed network has unexpected properties compared to classic small-world networks. Our TB co-authorship network displays properties that are more extreme than those of small-world and preferential attachment networks contradicting previous studies reporting co-authorship network as having small-world or preferential attachment properties [55, 109].

As the first advanced statistical model we applied to this network, the SBM identified heterogeneous classes with higher probabilities towards inter class ties establishment. This

### *Results: The Tuberculosis Co-authorship Network*

---

observation is different from what we observed for the malaria and the HIV/AIDS co-authorship network which both display low inter class probabilities and higher intra class probabilities of tie formation.

As in the malaria and the co-authorship network, the ERGM and TERGM results suggest that authors within the TB co-authorship network are more likely to establish collaboration ties within their research groups or communities. Although marginal, factors such as number of publications, number of citations and number of collaborations are associated to higher likelihood to establishing collaboration ties, confirming therefore our first hypothesis. Adding temporal dependencies to our ERGM models tremendously improved the fitness of the model to the observed network data, but at a cost of decreased performance compared to the model without temporal dependencies.

We expected the ERGM and TERGM models containing ERGM structural to converge for the TB co-authorship network given its relatively small size. Unfortunately, as for the malaria and the HIV/AIDS co-authorship networks, adding such terms to the the models proved computationally expensive. None of the models converged after 1,000 iterations. We therefore, suspect the complexity of the network to have prevented the convergence of the models [111].

With the LNM, we complement the ERGM and TERGM by adding an extra layer of analysis. Visualizing the effect of geography on the structure of the network, we notice that none of the nodal or dyadic covariates played a significant role in the spatial distribution of the network. Such an observation contradicts that of the malaria and the HIV/AIDS co-authorship network. The cluster demarcation observed with the null LNM

*Results: The Tuberculosis Co-authorship Network*

---

model suggests that distance does play a significant role in collaboration tie formation in the TB co-authorship network.

As the co-infection TB-HIV/AIDS continues to be an important aspect of the public health strategies in the Republic of Benin, consolidating the knowledge generated from the many TB-related research is crucial. Furthermore, public health policies must empower and reinforce the different research groups or communities involved in the research effort. Our results suggest a need for a continuous support to the TB research network, considering its low productivity status in Benin. Such actions will help stabilize the research groups already involved in TB research and promote the junior scientists in the field. We finally believe that such measures will ultimately insure the long-term sustainability of the TB co-authorship and collaborative research network in Benin.

# Chapter 6

## AuthorVis: A Co-authorship

## Visualization and Scientific

## Collaboration Prediction tool

### 6.1 Background

In this chapter, we propose a co-authorship network exploration, and link prediction tool specific to the three networks investigated in this dissertation. While many network visualization solutions have already been published, most of them are not specifically adapted to co-authorship network [125–128]. Even those designed for visualizing co-authorship network have several limitations among others, their inability to satisfactorily display large networks, the lack of interactivity in the display, and the inability for the

end user to control the display [125].

Here, we present a tool that not only provides a visualization of each of the networks but allow the end user to query the network. Our integrates bibliometrics information to the visualization. With our proposed model, all the authorship information are embedded within the network, and at the fingertip of the end user. In the visualization interface, users can select a particular node or author to emphasize its subnetwork, hover over a node to display author's information or select an edge between two nodes/authors to display information related to materials co-authored by the two nodes defining that particular edge.

## 6.2 Related work

Various authors have proposed diverse tools for the visualizations of co-authorship networks. One of such tools has been reported by Liu and colleagues [127] who proposed an author navigator application for visual examination of co-authorhip networks. In their conception of the toolkits, the authors combined a web based application tool for the interactive navigation of the network and a Java based backend swing application for the management of CGI requests. To support Brazilian researchers, Barbosa and colleagues proposed **VRRC**, a web based tool for the visualization and recommendation of co-authorship network [129]. According to its developers, **VRRC** provides an interactive visualization, an overview of the collaborations over time, and recommendations to initiate new collaborations and reinforce existing ones. **VICI**, another co-authorship

visualization tool was proposed by Odoni and colleagues [126]. **VICI** combined a Python based backend system for the extraction and management of the network data and a web based frontend using Flask [130] to display the network. The visualization of the network was finally rendered using the Javascript D3.js [131] library. **NeL<sup>2</sup>**, a general purpose tool for the visualization of networks as a layered network diagram was proposed by Nakazono, Misue, and Tanaka [125]. They applied their tool to the visualization of co-authorship networks to visualize transitions in the network over a period of time, as well as various co-authorship data.

Another framework, the WebRelievo system was proposed for the visualization of the evolutionary processes of Web pages [132]. Other techniques were also proposed for the visualization of co-citation networks [133], and for the visualization of the relationship of scientific literature [134].

Here, we propose **AuthorVis**, a co-authorship visualization and scientific collaboration tool for Malaria, TB and HIV/AIDS research in Benin. In addition to providing the same features as the aforementioned tools, with **AuthorVis**, we propose a different approach to co-authorship network visualization. Our approach integrates network structure and network data, hence requires no data management using traditional database framework. In addition, our visualization allows the end-user to navigate the network with an interactive navigation panel, but also integrate published materials within the visualization interface.

## 6.3 Design and Architecture

**AuthorVis** is built to a Shiny dashboard with an R based backend system that managed each co-authorship network data as an igraph object [135]. The backend system allows the end user to query the system through the Shinyboard. It is combined to a Javascript web based frontend that displays the network graph, and handle user interactions with the network.

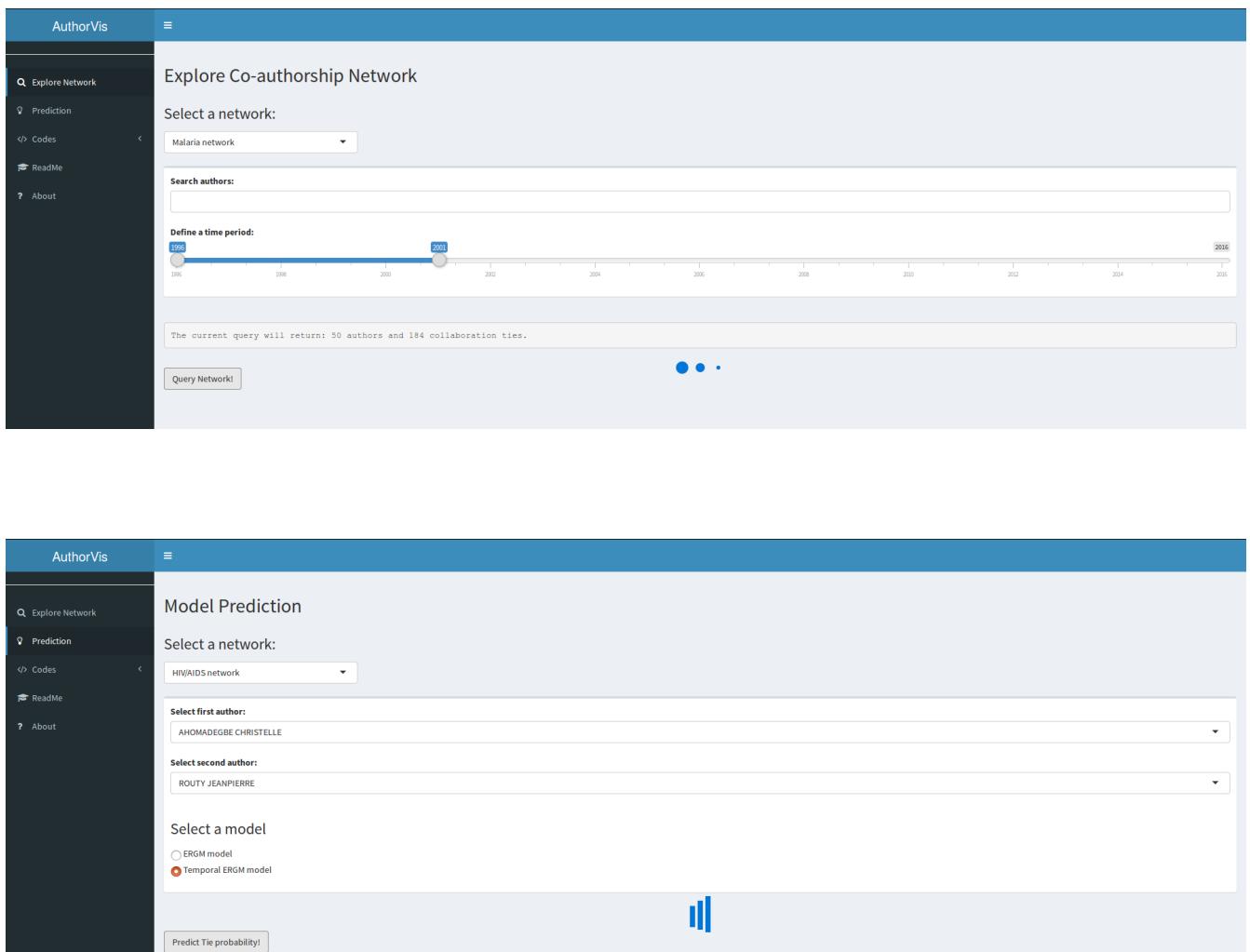


FIGURE 6.1: Screenshots of the Shiny application interface.

### 6.3.1 Data

Currently, **AuthorVis** is designed specifically for the visualization of the Malaria, Tuberculosis and HIV/AIDS collaborative network in Benin. We refer the reader to section 2.2 for details on the collection and treatment of the co-authorship data. On the server end, each network data is maintained as an igraph object. Each submitted user query is interpreted and incorporated in an igraph function to extract the network data. Another igraph object is generated as a result and converted into a JSON data using a specific Python script.

### 6.3.2 Network Visualization

The frontend network visualisation is built using the Javascript D3.js [131]. We built in a navigation panel allowing the user to interact with the network and control physics of network [136]. We incorporated several Javascript functions to design an intuitive and user friendly visualization interface. A mouse hover over a vertex displays a tooltip of details on the author represented by the vertex while a double-click on a vertex highlights the subnetwork of the identified network. We made edges clickable. Once an edge is clicked, the list of published materials co-authored by the two vertices defining the clicked edge is displayed in a panel on the right hand side. All published materials listed can be traced back to their publication page on the web via their DOI or the WOS accession number with a single click.

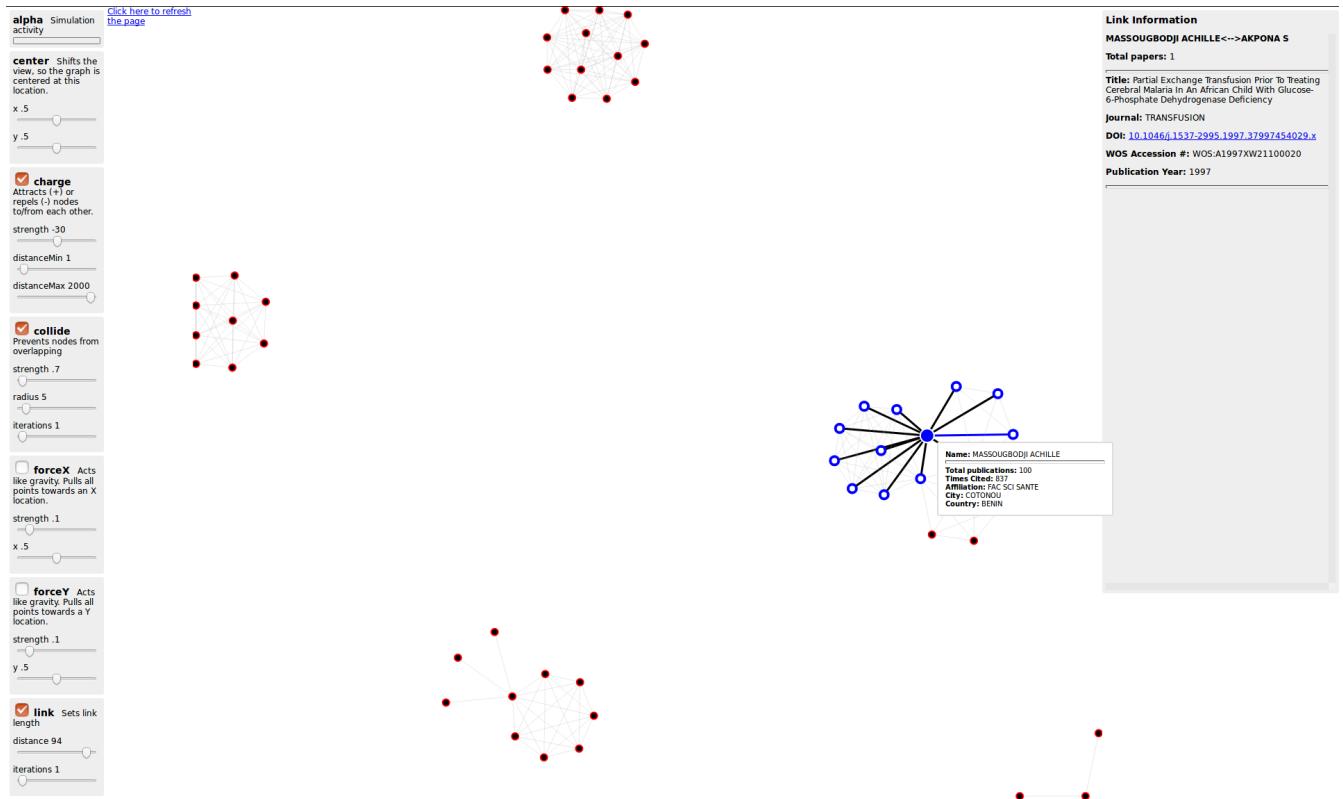


FIGURE 6.2: Screenshot of the co-authorship network visualization interface.

### 6.3.3 Web Framework

The whole system is built into a Shiny dashboard thanks to the R package **Shinyboard** [137, 138]. Using the dashboard, the user can choose to use the prediction tool menu, or query and visualize the network data. We also provide within the dashboard all our Shiny codes and a link to our Git directory containing all our source codes.

The prediction tool is model based and used the ERGM models to calculate a probability of collaboration between two authors. A micro-interpretation of the model is provided based on the user query [99].

When the user chooses the visualization menu, an interface allows him to submit his

query to the system. The query is interpreted and processed on the server end and the user is automatically prompted to a new visualization page.

## 6.4 Deployment

The visualization front-end is maintained by an Python http-server. The system is packed in a Docker container to facilitate its use and installation. The project source files can be forked or cloned from Github at <https://github.com/rosericazondekon/authorvis>. We also made it accessible online via an AWS server.

# General Conclusion

In this dissertation, we have documented and described the collaborative pattern in Malaria, HIV/AIDS and TB research in Benin. Our findings suggest that each one of the collaborative research network of Malaria, HIV/AIDS and TB has a complex structure. We modeled these complex structures to predict the establishment of future collaboration ties. We implemented the models in a shiny-based application for co-authorship visualization and scientific collaboration prediction tool which we named **AuthorVis**. The application of temporal or dynamic modeling techniques is the major strength of our research along with its application of not only descriptive methods but also robust network analysis methods such as inferential methods like Monte-Carlo simulations, unlike most studies on co-authorship analysis. Our data mining strategy involved a robust machine learning algorithm that helped address the crucial issue of the disambiguation of authors names and assign a unique identification to each of them. This technique maintained a good quality of the data collected throughout the pre-processing and analysis steps. To the best of our knowledge, our study is the first to describe the malaria research collaborations network via co-authorship network analysis in Benin. It is also

## *General Conclusion*

---

the first to apply statistical network models to investigate co-authorship network in a specific research area in an African country.

The fact that we collected data only from the Web Of Science can be considered as an important limitation of this study. However, according to Falagas and colleagues [139], who compared PubMed, Scopus, Web Of Science and Google Scholar in their paper, the Web Of Science appears as a reasonable scientific database source for our analysis. In addition, it proved to cover a wide range of both old and recently published papers. Falagas and colleagues [139] found PubMed to be the optimal choice in terms of scientific database. For that reason we did run the same bibliographic search in PubMed. Unfortunately, the Web Of Science returns more relevant data than PubMed. Yet another limitation is inherent to the nature of all co-authorship studies. Collaborators, in a co-authorship network, do not often come from the same scientific discipline, or do not play the same roles on a particular research project. The data we collected did not allow us to accurately assess or even infer the disciplines each author came from or their specific contribution in the published document.

## **Future Directions**

There are several future directions. Our work be extended to the entire African collaboration network in Malaria, HIV/AIDS and TB. Since collaborations usually are often initiated between individuals, labs or even countries, the analysis of bipartite co-authorship networks is an interesting direction to our study. Currently, **AuthorVis** is specifically

## *General Conclusion*

---

built for Malaria, TB and HIV/AIDS in Benin. Future development may extend the tool to other research domain. Adding a general purpose module to **AuthorVis** for the visualization of any user-input co-authorship network is an interesting venture since it will also require the integration of a data pre-processing module to facilitate the disambiguation and deduplication of co-authorship information. Furthermore, incorporating a layered structured network visualization [125] functionality to the visualization in order to display temporal changes in the evolution of the co-authorship network is another interesting direction. It can, in addition be designed into a real-time, cross-domain, and cross-collection co-authorship visualization interface capable of automatically searching the literature. Outside of the realm of co-authorship analyses, the same idea of network visualization can be extended to other important disciplines such as Neuroscience, making possible the visualization of real time connectivity dynamics between well identified regions of interest in the brain during certain resting, memory or motor tasks.

# Bibliography

- [1] Jonathan R Davis and Joshua Lederberg. *Emerging Infectious Diseases from the Global to the Local Perspective: Workshop Summary*. National Academies Press, 2001. ISBN 0-309-07184-4. 00021.
- [2] John Luke Gallup and Jeffrey D Sachs. The economic burden of malaria. *The American journal of tropical medicine and hygiene*, 64(1 suppl):85–96, 2001. ISSN 0002-9637. 01510.
- [3] M. Vitoria, R. Granich, C. F. Gilks, C. Gunneberg, M. Hosseini, W. Were, M. Ravaglione, and K. M. De Cock. The Global Fight Against HIV/AIDS, Tuberculosis, and Malaria: Current Status and Future Perspectives. *American Journal of Clinical Pathology*, 131(6):844–848, June 2009. ISSN 0002-9173, 1943-7722. doi: 10.1309/AJCP5XHDB1PNAEYT. 00003.
- [4] UN General Assembly. United Nations millennium declaration. *United Nations General Assembly*, 2000. 00156.
- [5] Margaret Arthur. Institute for Health Metrics and Evaluation. *Nursing Standard*, 28(42):32–32, 2014. ISSN 0029-6570. 00000.

## *Bibliography*

---

[6] Christopher J L Murray, Katrina F Ortblad, Caterina Guinovart, Stephen S Lim, Timothy M Wolock, D Allen Roberts, Emily A Dansereau, Nicholas Graetz, Ryan M Barber, Jonathan C Brown, Haidong Wang, Herbert C Duber, Mohsen Naghavi, Daniel Dicker, Lalit Dandona, Joshua A Salomon, Kyle R Heuton, Kyle Foreman, David E Phillips, Thomas D Fleming, Abraham D Flaxman, Bryan K Phillips, Elizabeth K Johnson, Megan S Coggeshall, Foad Abd-Allah, Semaw Ferede Abera, Jerry P Abraham, Ibrahim Abubakar, Laith J Abu-Raddad, Niveen Me Abu-Rmeileh, Tom Achoki, Austine Olufemi Adeyemo, Arsène Kouablan Adou, José C Adsuar, Emilie Elisabet Agardh, Dickens Akena, Mazin J Al Kahbouri, Deena Alasfoor, Mohammed I Albittar, Gabriel Alcalá-Cerra, Miguel Angel Alegretti, Zewdie Aderaw Alemu, Rafael Alfonso-Cristancho, Samia Alhabib, Raghib Ali, Francois Alla, Peter J Allen, Ubai Alsharif, Elena Alvarez, Nelson Alvis-Guzman, Adansi A Amankwaa, Azmeraw T Amare, Hassan Amini, Walid Ammar, Benjamin O Anderson, Carl Abelardo T Antonio, Palwasha Anwari, Johan Ärnlöv, Valentina S Arsic Arsenijevic, Ali Artaman, Rana J Asghar, Reza Assadi, Lydia S Atkins, Alaa Badawi, Kalpana Balakrishnan, Amitava Banerjee, Sanjay Basu, Justin Beardsley, Tolesa Bekele, Michelle L Bell, Eduardo Bernabe, Tariku Jibat Beyene, Neeraj Bhala, Ashish Bhalla, Zulfiqar A Bhutta, Aref Bin Abdulhak, Agnes Binagwaho, Jed D Blore, Berrak Bora Basara, Dipan Bose, Michael Brainin, Nicholas Breitborde, Carlos A Castañeda-Orjuela, Ferrán Catalá-López, Vineet K Chadha, Jung-Chen Chang, Peggy Pei-Chia Chiang, Ting-Wu Chuang, Mercedes Colomar, Leslie Trumbull

## *Bibliography*

---

Cooper, Cyrus Cooper, Karen J Courville, Benjamin C Cowie, Michael H Criqui, Rakhi Dandona, Anand Dayama, Diego De Leo, Louisa Degenhardt, Borja Del Pozo-Cruz, Kebede Deribe, Don C Des Jarlais, Muluken Dessalegn, Samath D Dharmaratne, Uğur Dilmen, Eric L Ding, Tim R Driscoll, Adnan M Durrani, Richard G Ellenbogen, Sergey Petrovich Ermakov, Alireza Esteghamati, Emerito Jose A Faraon, Farshad Farzadfar, Seyed-Mohammad Fereshtehnejad, Daniel Obadare Fijabi, Mohammad H Forouzanfar, Urbano Fra.Paleo, Lynne Gaffikin, Amiran Gamkrelidze, Fortuné Gbètoho Gankpé, Johanna M Geleijnse, Bradford D Gessner, Katherine B Gibney, Ibrahim Abdelmageem Mohamed Ginawi, Elizabeth L Glaser, Philimon Gona, Atsushi Goto, Hebe N Gouda, Harish Chander Gugnani, Rajeev Gupta, Rahul Gupta, Nima Hafezi-Nejad, Randah Ribhi Hamadeh, Mouhanad Hammami, Graeme J Hankey, Hilda L Harb, Josep Maria Haro, Rasmus Havmoeller, Simon I Hay, Mohammad T Hedayati, Ileana B Heredia Pi, Hans W Hoek, John C Hornberger, H Dean Hosgood, Peter J Hotez, Damian G Hoy, John J Huang, Kim M Ibburg, Bulat T Idrisov, Kaire Innos, Kathryn H Jacobsen, Panniyammakal Jeemon, Paul N Jensen, Vivekanand Jha, Guohong Jiang, Jost B Jonas, Knud Juel, Haidong Kan, Ida Kankindi, Nadim E Karam, André Karch, Corine Kakizi Karema, Anil Kaul, Norito Kawakami, Dhruv S Kazi, Andrew H Kemp, Andre Pascal Kengne, Andre Keren, Maia Kereselidze, Yousef Saleh Khader, Shams Eldin Ali Hassan Khalifa, Ejaz Ahmed Khan, Young-Ho Khang, Irma Khonelidze, Yohannes Kinfu, Jonas M Kinge, Luke Knibbs, Yoshihiro Kokubo, S Kosen, Barthelemy Kuate Defo, Veena S Kulkarni,

## *Bibliography*

---

Chanda Kulkarni, Kaushalendra Kumar, Ravi B Kumar, G Anil Kumar, Gene F Kwan, Taavi Lai, Arjun Lakshmana Balaji, Hilton Lam, Qing Lan, Van C Lansingh, Heidi J Larson, Anders Larsson, Jong-Tae Lee, James Leigh, Mall Leinsalu, Ricky Leung, Yichong Li, Yongmei Li, Graça Maria Ferreira De Lima, Hsien-Ho Lin, Steven E Lipshultz, Shiwei Liu, Yang Liu, Belinda K Lloyd, Paulo A Lotufo, Vasco Manuel Pedro Machado, Jennifer H MacLachlan, Carlos Magis-Rodriguez, Marek Majdan, Christopher Chabila Mapoma, Wagner Marcenes, Melvin Barrientos Marzan, Joseph R Masci, Mohammad Taufiq Mashal, Amanda J Mason-Jones, Bongani M Mayosi, Tasara T Mazorodze, Abigail Cecilia Mckay, Peter A Meaney, Man Mohan Mehndiratta, Fabiola Mejia-Rodriguez, Yohannes Adama Melaku, Ziad A Memish, Walter Mendoza, Ted R Miller, Edward J Mills, Karzan Abdulmuhsin Mohammad, Ali H Mokdad, Glen Liddell Mola, Lorenzo Monasta, Marcella Montico, Ami R Moore, Rintaro Mori, Wilkister Nyaora Moturi, Mitsuru Mukaigawara, Kinnari S Murthy, Aliya Naheed, Kovin S Naidoo, Luigi Naldi, Vinay Nangia, K M Venkat Narayan, Denis Nash, Chakib Nejjari, Robert G Nelson, Sudan Prasad Neupane, Charles R Newton, Marie Ng, Muhammad Imran Nisar, Sandra Nolte, Ole F Norheim, Vincent Nowaseb, Luke Nyakarahuka, In-Hwan Oh, Takayoshi Ohkubo, Bolajoko O Olusanya, Saad B Omer, John Nelson Opio, Orish Ebere Orisakwe, Jeyaraj D Pandian, Christina Papachristou, Angel J Paternina Caicedo, Scott B Patten, Vinod K Paul, Boris Igor Pavlin, Neil Pearce, David M Pereira, Aslam Pervaiz, Konrad Pesudovs, Max Petzold, Farshad Pourmalek, Dima Qato,

## *Bibliography*

---

Amado D Quezada, D Alex Quistberg, Anwar Rafay, Kazem Rahimi, Vafa Rahimi-Movaghar, Sajjad Ur Rahman, Murugesan Raju, Saleem M Rana, Homie. Global, regional, and national incidence and mortality for HIV, tuberculosis, and malaria during 1990–2013: A systematic analysis for the Global Burden of Disease Study 2013. *The Lancet*, 384(9947):1005–1070, September 2014. ISSN 01406736.

doi: 10.1016/S0140-6736(14)60844-8. 00405.

[7] Craig Stoops. President's Malaria Initiative. Technical report, DTIC Document, 2008. 00000.

[8] Global Fund. Making a difference: Global fund results report 2011. *The Global Fund, Geneva*, 2011. 00005.

[9] World Health Organization. World malaria report 2010. *Geneva: World Health Organization View Article Google Scholar*, 2012. 00161.

[10] Lawrence M Barat. Four malaria success stories: How malaria burden was successfully reduced in Brazil, Eritrea, India, and Vietnam. *The American journal of tropical medicine and hygiene*, 74(1):12–16, 2006. ISSN 0002-9637. 00117.

[11] Martin C. Akogbéto, Rock Y. Aïkpon, Roseric Azondékon, Gil G. Padonou, Razaki A. Ossè, Fiacre R. Agossa, Raymond Beach, and Michel Sèzonlin. Six years of experience in entomological surveillance of indoor residual spraying against malaria transmission in Benin: Lessons learned, challenges and outlooks.

*Malaria Journal*, 14(1), December 2015. ISSN 1475-2875. doi: 10.1186/s12936-015-0757-5. 00002.

## *Bibliography*

---

- [12] Joint United Nations Programme on HIV/AIDS. *Getting to Zero: 2011–2015 strategy*. 2010. 00035.
- [13] Joint United Nations Programme on HIV/AIDS. *Global Report: UNAIDS Report on the Global AIDS Epidemic 2010*. UNAIDS, 2010. ISBN 92-9173-871-9. 01036.
- [14] World Health Organization. *Global Tuberculosis Control: WHO Report 2010*. World Health Organization, 2010. ISBN 92-4-156406-7. 00010.
- [15] Dean T Jamison. *Disease and Mortality in Sub-Saharan Africa*. World Bank Publications, 2006. ISBN 0-8213-6398-0. 00259.
- [16] Carole A Long and Fidel Zavala. Malaria vaccines and human immune responses. *Current Opinion in Microbiology*, 32:96–102, August 2016. ISSN 13695274. doi: 10.1016/j.mib.2016.04.006. 00003.
- [17] Fausto Titti, Aurelio Cafaro, Flavia Ferrantelli, Antonella Tripiciano, Sonia Moretti, Antonella Caputo, Riccardo Gavioli, Fabrizio Ensoli, Marjorie Robert-Guroff, Susan Barnett, and Barbara Ensoli. Problems and emerging approaches in HIV/AIDS vaccine development. *Expert Opinion on Emerging Drugs*, 12(1):23–48, March 2007. ISSN 1472-8214, 1744-7623. doi: 10.1517/14728214.12.1.23. 00036.
- [18] B. D. Walker and D. R. Burton. Toward an AIDS Vaccine. *Science*, 320(5877): 760–764, May 2008. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.1152622. 00407.

## Bibliography

---

- [19] M. E. J. Newman. The structure of scientific collaboration networks. *Proceedings of the National Academy of Sciences*, 98(2):404–409, January 2001. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.98.2.404. 04061.
- [20] Elizabeth L. Corbett, Catherine J. Watt, Neff Walker, Dermot Maher, Brian G. Williams, Mario C. Raviglione, and Christopher Dye. The Growing Burden of Tuberculosis: Global Trends and Interactions With the HIV Epidemic. *Archives of Internal Medicine*, 163(9):1009, May 2003. ISSN 0003-9926. doi: 10.1001/archinte.163.9.1009. 02947.
- [21] Neel R. Gandhi, N. Sarita Shah, Jason R. Andrews, Venanzio Vella, Anthony P. Moll, Michelle Scott, Darren Weissman, Claudio Marra, Umesh G. Laloo, and Gerald H. Friedland. HIV Coinfection in Multidrug- and Extensively Drug-Resistant Tuberculosis Results in High Early Mortality. *American Journal of Respiratory and Critical Care Medicine*, 181(1):80–86, January 2010. ISSN 1073-449X, 1535-4970. doi: 10.1164/rccm.200907-0989OC. 00248.
- [22] World Health Organization. Economic costs of malaria are many times higher than previously estimated. In *Economic Costs of Malaria Are Many Times Higher than Previously Estimated*. 2000. 00012.
- [23] Linda M. Richter, Knut Lönnroth, Chris Desmond, Robin Jackson, Ernesto Jaramillo, and Diana Weil. Economic Support to Patients in HIV and TB Grants in Rounds 7 and 10 from the Global Fund to Fight AIDS, Tuberculosis and

## *Bibliography*

---

- Malaria. *PLoS ONE*, 9(1):e86225, January 2014. ISSN 1932-6203. doi: 10.1371/journal.pone.0086225. 00015.
- [24] Chinua Akukwe. *Don't Let Them Die: HIV/AIDS, Malaria, Tuberculosis and the Healthcare Crisis in Africa*. Adonis & Abbey Pub Limited, 2006. ISBN 1-905068-24-7. 00005.
- [25] Sam M Mbulaiteye, Kishor Bhatia, Clement Adebamowo, and Annie J Sasco. HIV and cancer in Africa: Mutual collaboration between HIV and cancer programs may provide timely research and public health data. *Infectious Agents and Cancer*, 6(1):16, 2011. ISSN 1750-9378. doi: 10.1186/1750-9378-6-16. 00038.
- [26] U. D'Alessandro, B.O. Olaleye, W. McGuire, M.C. Thomson, P. Langerock, S. Bennett, and B.M. Greenwood. A comparison of the efficacy of insecticide-treated and untreated bed nets in preventing malaria in Gambian children. *Transactions of the Royal Society of Tropical Medicine and Hygiene*, 89 (6):596–598, November 1995. ISSN 00359203. doi: 10.1016/0035-9203(95)90401-8.
- [27] Joint United Nations Programme on HIV/AIDS and World Health Organization. *AIDS Epidemic Update, December 2006*. World Health Organization, 2007. ISBN 92-9173-542-6. 00519.
- [28] Alan Whiteside. *HIV/AIDS: A Very Short Introduction*, volume 174. Oxford University Press, 2008. ISBN 0-19-280692-0. 00119.

## *Bibliography*

---

- [29] Centers for Disease Control and Prevention. Revised recommendations for HIV testing of adults, adolescents, and pregnant women in health-care settings. *Annals of Emergency Medicine*, 49(5):575–577, 2007. ISSN 0196-0644. 08464.
- [30] Bluma Brenner and Mark A. Wainberg. We need to use the best antiretroviral drugs worldwide to prevent HIV drug resistance:. *AIDS*, 30(17):2725–2727, November 2016. ISSN 0269-9370. doi: 10.1097/QAD.0000000000001234. 00000.
- [31] Alexandra Calmy, Fernando Pascual, and Nathan Ford. HIV drug resistance. *New England Journal of Medicine*, 350(26):2720–2721, 2004. ISSN 0028-4793. 00038.
- [32] François Clavel and Allan J Hance. HIV drug resistance. *New England Journal of Medicine*, 350(10):1023–1035, 2004. ISSN 0028-4793. 00817.
- [33] Stefan HE Kaufmann and Paul van Helden. *Handbook of Tuberculosis: Clinics, Diagnostics, Therapy and Epidemiology*. Wiley-VCH, 2008. ISBN 3-527-31888-7. 00010.
- [34] Alimuddin Zumla. Handbook of tuberculosis. *The Lancet Infectious Diseases*, 9 (12):736, 2009. ISSN 1473-3099. 00000.
- [35] MC Raviglione, AD Harries, R Msiska, David Wilkinson, and P Nunn. Tuberculosis and HIV: Current status in Africa. *AIDS (London, England)*, 11: S115, 1997. ISSN 0269-9370. 00252.

## Bibliography

---

- [36] SK Sharma, Alladi Mohan, and Tamilarasu Kadiravan. HIV-TB co-infection: Epidemiology, diagnosis & management. *Indian Journal of Medical Research*, 121(4):550–567, 2005. ISSN 0971-5916.
- [37] Z Toossi, Hirsch Mayanja-Kizza, CS Hirsch, KL Edmonds, T Spahlinger, DL Hom, H Aung, P Mugyenyi, JJ Ellner, and CW Whalen. Impact of tuberculosis (TB) on HIV-1 activity in dually infected patients. *Clinical & Experimental Immunology*, 123(2):233–238, 2001. ISSN 1365-2249. 00175.
- [38] Lia D'Anibrosio, Antonio Spanevello, and Rosella Centis. Epidemiology of TB. *Tuberculosis*, 58:14, 2014. ISSN 1849840288. 00000.
- [39] Centers for Disease Control and Prevention (CDC). Emergence of Mycobacterium tuberculosis with extensive resistance to second-line drugs—worldwide, 2000-2004. *MMWR. Morbidity and mortality weekly report*, 55(11):301, 2006. ISSN 1545-861X. 00651.
- [40] World Health Organization. Multidrug and extensively drug-resistant TB. 2010. 00915.
- [41] Robert L Cowie. The epidemiology of tuberculosis in gold miners with silicosis. *American journal of respiratory and critical care medicine*, 150(5):1460–1462, 1994. ISSN 1073-449X. 00176.
- [42] Emmanuel M Mulenga, Hugh B Miller, Thomson Sinkala, Tracy A Hysong, and Jefferey L Burgess. Silicosis and tuberculosis in Zambian miners. *International journal of occupational and environmental health*, 2013. 00014.

## Bibliography

---

- [43] D Rees and J Murray. Silica, silicosis and tuberculosis [State of the Art Series. Occupational lung disease in high-and low-income countries, Edited by M. Chan-Yeung. Number 4 in the series]. *The International Journal of Tuberculosis and Lung Disease*, 11(5):474–484, 2007. ISSN 1027-3719. 00132.
- [44] Marianne E Sinka, Michael J Bangs, Sylvie Manguin, Yasmin Rubio-Palis, Theeraphap Chareonviriyaphap, Maureen Coetzee, Charles M Mbogo, Janet Hemingway, Anand P Patil, and William H Temperley. A global map of dominant malaria vectors. *Parasites & vectors*, 5(1):1, 2012. ISSN 1756-3305. 00189.
- [45] Robert W. Snow, Carlos A. Guerra, Abdisalan M. Noor, Hla Y. Myint, and Simon I. Hay. The global distribution of clinical episodes of Plasmodium falciparum malaria. *Nature*, 434(7030):214–217, March 2005. ISSN 0028-0836, 1476-4679. doi: 10.1038/nature03342. 02814.
- [46] S. P. James and P. Tate. New Knowledge of the Life-Cycle of Malaria Parasites. *Nature*, 139(3517):545–545, March 1937. ISSN 0028-0836. doi: 10.1038/139545a0. 00080.
- [47] P.L. Alonso, S.W. Lindsay, J.R.M. Armstrong Schellenberg, K. Keita, P. Gomez, F.C. Shenton, A.G. Hill, P.H. David, G. Fegan, K. Cham, and B.M. Greenwood. A malaria control trial using insecticide-treated bed nets and targeted chemoprophylaxis in a rural area of The Gambia, West Africa. *Transactions of the Royal Society of Tropical Medicine and Hygiene*, 87:37–44, June 1993. ISSN 00359203. doi: 10.1016/0035-9203(93)90174-O.

## Bibliography

---

- [48] Katherine E. Battle, Donal Bisanzio, Harry S. Gibson, Samir Bhatt, Ewan Cameron, Daniel J. Weiss, Bonnie Mappin, Ursula Dalrymple, Rosalind E. Howes, Simon I. Hay, and Peter W. Gething. Treatment-seeking rates in malaria endemic countries. *Malaria Journal*, 15(1), December 2016. ISSN 1475-2875. doi: 10.1186/s12936-015-1048-x.
- [49] Loet Leydesdorff and Staša Milojević. Scientometrics. *arXiv preprint arXiv:1208.4566*, 2012. 00467.
- [50] Garfield Eugene. Citation Indexing, Its Theory and Application in Science, Technology, and Humanities. 1979.
- [51] Terttu Luukkonen, Olle Persson, and Gunnar Sivertsen. Understanding patterns of international scientific collaboration. *Science, Technology & Human Values*, 17 (1):101–126, 1992. ISSN 0162-2439. 00457.
- [52] Caroline S. Wagner. Six case studies of international collaboration in science. *Scientometrics*, 62(1):3–26, January 2005. ISSN 0138-9130, 1588-2861. doi: 10.1007/s11192-005-0001-0. 00193.
- [53] Wolfgang Glänzel and András Schubert. Analysing scientific networks through co-authorship. In *Handbook of Quantitative Science and Technology Research*, pages 257–276. Springer, 2004. 00493.
- [54] M. E. J. Newman. Coauthorship networks and patterns of scientific collaboration. *Proceedings of the National Academy of Sciences*, 101(Supplement 1):5200–5205, April 2004. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.0307545100. 01352.

## Bibliography

---

- [55] Gregorio González-Alcaide, Jinseo Park, Charles Huamaní, Joaquín Gascón, and José Manuel Ramos. Scientific authorships and collaboration network analysis on Chagas disease: Papers indexed in PubMed (1940-2009). *Revista do Instituto de Medicina Tropical de São Paulo*, 54(4):219–228, 2012. ISSN 0036-4665. 00028.
- [56] M. E. J. Newman. Scientific collaboration networks. I. Network construction and fundamental results. *Physical Review E*, 64(1), June 2001. ISSN 1063-651X, 1095-3787. doi: 10.1103/PhysRevE.64.016131.
- [57] M. E. J. Newman. Scientific collaboration networks. II. Shortest paths, weighted networks, and centrality. *Physical Review E*, 64(1), June 2001. ISSN 1063-651X, 1095-3787. doi: 10.1103/PhysRevE.64.016132. 02213.
- [58] Katy Börner, Chaomei Chen, and Kevin W. Boyack. Visualizing knowledge domains. *Annual Review of Information Science and Technology*, 37(1):179–255, January 2003. ISSN 1550-8382. doi: 10.1002/aris.1440370106. 00981.
- [59] Andrea Scharnhorst, Katy Börner, and Peter van den Besselaar, editors. *Models of Science Dynamics: Encounters between Complexity Theory and Information Sciences*. Understanding complex systems. Springer, Heidelberg ; New York, 2012. ISBN 978-3-642-23067-7. 00050.
- [60] F. Mali, L. Kronegger, P. Doreian, and A. Ferligoj. *Dynamic Scientific Co-Authorship Networks*. Understanding Complex Systems. 2012. 00050.
- [61] Helga Bermeo Andrade, Ernesto de los Reyes López, and Tomas Bonavia Martín. Dimensions of scientific collaboration and its contribution to the academic

## *Bibliography*

---

- research groups' scientific quality. *Research Evaluation*, 18(4):301–311, October 2009. ISSN 09582029, 14715449. doi: 10.3152/095820209X451041. 00038.
- [62] Juan D Rogers, Barry Bozeman, and Ivan Chompalov. Obstacles and opportunities in the application of network analysis to the evaluation of R&D. *Research Evaluation*, 10(3):161–172, December 2001. ISSN 09582029, 14715449. doi: 10.3152/147154401781777033. 00000.
- [63] Diane H. Sonnenwald. Scientific collaboration. *Annual Review of Information Science and Technology*, 41(1):643–681, 2007. ISSN 00664200. doi: 10.1002/aris.2007.1440410121. 00418.
- [64] Haiyan Hou, Hildrun Kretschmer, and Zeyuan Liu. The structure of scientific collaboration networks in Scientometrics. *Scientometrics*, 75(2):189–202, May 2008. ISSN 0138-9130, 1588-2861. doi: 10.1007/s11192-007-1771-3. 00186.
- [65] Gregorio González-Alcaide, Rafael Aleixandre-Benavent, Carolina Navarro-Molina, and Juan Carlos Valderrama-Zurián. Coauthorship networks and institutional collaboration patterns in reproductive biology. *Fertility and Sterility*, 90(4):941–956, October 2008. ISSN 00150282. doi: 10.1016/j.fertnstert.2007.07.1378. 00035.
- [66] Hannes Toivanen and Branco Ponomariov. African regional innovation systems: Bibliometric analysis of research collaboration patterns 2005–2009. *Scientometrics*, 88(2):471–493, August 2011. ISSN 0138-9130, 1588-2861. doi: 10.1007/s11192-011-0390-1. 00035.

## *Bibliography*

---

- [67] L. Bellanca. Measuring interdisciplinary research: Analysis of co-authorship for research staff at the University of York. *Bioscience Horizons*, 2(2):99–112, June 2009. ISSN 1754-7431. doi: 10.1093/biohorizons/hzp012.
- [68] R. Aleixandre-Benavent, G. González-Alcaide, A. Alonso-Arroyo, M. Bolaños-Pizarro, L. Castelló-Cogollos, and J.C. Valderrama-Zurián. Coauthorship Networks and Institutional Collaboration in Farmacia Hospitalaria. *Farmacia Hospitalaria (English Edition)*, 32(4):226–233, January 2008. ISSN 21735085. doi: 10.1016/S2173-5085(08)70044-3. 00003.
- [69] H.B. Ghafouri, H. Mohammadhassanzadeh, F. Shokraneh, M. Vakilian, and S. Farahmand. Social network analysis of Iranian researchers on emergency medicine: A sociogram analysis. *Emergency Medicine Journal*, 31(8):619–624, 2014. doi: 10.1136/emermed-2012-201781. 00008.
- [70] P. Salamati and F. Soheili. Social network analysis of Iranian researchers in the field of violence. *Chinese Journal of Traumatology - English Edition*, 19(5):264–270, 2016. doi: 10.1016/j.cjtee.2016.06.008. 00000.
- [71] F. Sadoughi, A. Valinejadi, M. Serati Shirazi, and R. Khademi. Social network analysis of Iranian researchers on medical parasitology: A 41 year co-authorship survey. *Iranian Journal of Parasitology*, 11(2):204–212, 2016. 00001.
- [72] Carlos Medicis Morel, Suzanne Jacob Serruya, Gerson Oliveira Penna, and Reinaldo Guimarães. Co-authorship Network Analysis: A Powerful Tool for Strategic Planning of Research, Development and Capacity Building Programs on

## *Bibliography*

---

Neglected Diseases. *PLoS Neglected Tropical Diseases*, 3(8):e501, August 2009.

ISSN 1935-2735. doi: 10.1371/journal.pntd.0000501. 00096.

- [73] Qing Zhang. *Complex Network Analysis for Scientific Collaboration Prediction and Biological Hypothesis Generation*. PhD thesis, University of Wisconsin-Milwaukee, Milwaukee, WI, USA, 2014. 00000.

- [74] Sandra Cristina Oliveira, Juliana Cobre, and Taiane de Paula Ferreira. A Bayesian approach for the reliability of scientific co-authorship networks with emphasis on nodes. *Social Networks*, 48:110–115, January 2017. ISSN 0378-8733. doi: 10.1016/j.socnet.2016.06.005. 00000.

- [75] Erick Peirson, Aaron Baker, Ramki Subramanian, Abhishek Singh, and Yogananda Yalugoti. Tethne v0.8, 2016.

- [76] Daniel A Schult and P Swart. Exploring network structure, dynamics, and function using NetworkX. volume 2008, pages 11–16, 2008. 01142.

- [77] Anderson A Ferreira, Marcos André Gonçalves, and Alberto HF Laender. A brief survey of automatic methods for author name disambiguation. *Acm Sigmod Record*, 41(2):15–26, 2012. ISSN 0163-5808. 00124.

- [78] C Lee Giles, Hongyuan Zha, and Hui Han. Name disambiguation in author citations using a k-way spectral clustering method. pages 334–343. IEEE, 2005. ISBN 1-58113-876-8. 00277.

## Bibliography

---

- [79] Mikhail Yuryevich Bilenko. *Learnable Similarity Functions and Their Application to Record Linkage and Clustering*. PhD thesis, University of Texas at Austin, Austin, TX, USA, 2006. 00019.
- [80] Linton C Freeman. A set of measures of centrality based on betweenness. *Sociometry*, pages 35–41, 1977. ISSN 0038-0431. 05656.
- [81] Phillip Bonacich. Factoring and weighting approaches to status scores and clique identification. *Journal of Mathematical Sociology*, 2(1):113–120, 1972. ISSN 0022-250X. 01823.
- [82] Leo Katz. A new status index derived from sociometric analysis. *Psychometrika*, 18(1):39–43, 1953. ISSN 0033-3123. 01961.
- [83] Eric D Kolaczyk and Gábor Csárdi. *Statistical Analysis of Network Data with R*, volume 65. Springer, 2014. 00792.
- [84] Paul Erdős and Alfréd Rényi. On random graphs, I. *Publicationes Mathematicae (Debrecen)*, 6:290–297, 1959. 03506.
- [85] Paul Erdos and Alfréd Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.*, 5(1):17–60, 1960. 03290.
- [86] Paul Erdős and Alfréd Rényi. On the strength of connectedness of a random graph. *Acta Mathematica Academiae Scientiarum Hungarica*, 12(1-2):261–267, 1964. ISSN 0001-5954. 00797.

## Bibliography

---

- [87] Edgar N Gilbert. Random graphs. *The Annals of Mathematical Statistics*, 30(4):1141–1144, 1959. ISSN 0003-4851. 00932.
- [88] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*, 393(6684):440–442, 1998. ISSN 0028-0836. 32864.
- [89] Vera Van Noort, Berend Snel, and Martijn A Huynen. The yeast coexpression network has a small-world, scale-free architecture and can be explained by a simple model. *EMBO reports*, 5(3):280–284, 2004. ISSN 1469-221X. 00213.
- [90] Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999. ISSN 0036-8075. 27102.
- [91] Réka Albert, Hawoong Jeong, and Albert-László Barabási. Internet: Diameter of the world-wide web. *nature*, 401(6749):130–131, 1999. ISSN 0028-0836. 05099.
- [92] Hawoong Jeong, Zoltan Néda, and Albert-László Barabási. Measuring preferential attachment in evolving networks. *EPL (Europhysics Letters)*, 61(4):567, 2003. ISSN 0295-5075. 00487.
- [93] François Lorrain and Harrison C. White. Structural equivalence of individuals in social networks. *The Journal of Mathematical Sociology*, 1(1):49–80, January 1971. ISSN 0022-250X, 1545-5874. doi: 10.1080/0022250X.1971.9989788.
- [94] Patrick Doreian, Vladimir Batagelj, and Anuska Ferligoj. *Generalized Blockmodeling*. Cambridge University Press, Cambridge, 2004. ISBN 978-0-511-58417-6. doi: 10.1017/CBO9780511584176.

## Bibliography

---

- [95] Garry Robins, Pip Pattison, Yuval Kalish, and Dean Lusher. An introduction to exponential random graph ( $p^*$ ) models for social networks. *Social networks*, 29(2):173–191, 2007. ISSN 0378-8733. 00000.
- [96] Steve Hanneke, Wenjie Fu, and Eric P Xing. Discrete temporal models of social networks. *Electronic Journal of Statistics*, 4:585–605, 2010. ISSN 1935-7524.
- [97] Garry Robins and Philippa Pattison. Random graph models for temporal processes in social networks. *Journal of Mathematical Sociology*, 25(1):5–41, 2001. ISSN 0022-250X.
- [98] Philip Leifeld, Skyler J Cranmer, and Bruce A Desmarais. Temporal Exponential Random Graph Models with xergm: Estimation and Bootstrap Confidence Intervals. *Journal of Statistical Software*, 2015. 00000.
- [99] Bruce A. Desmarais and Skyler J. Cranmer. Micro-Level Interpretation of Exponential Random Graph Models with Application to Estuary Networks: Desmarais/Cranmer: Micro-Level Interpretation of ERGM. *Policy Studies Journal*, 40(3):402–434, August 2012. ISSN 0190292X. doi: 10.1111/j.1541-0072.2012.00459.x.
- [100] Peter Hoff. Modeling homophily and stochastic equivalence in symmetric relational data. pages 657–664, 2008.
- [101] Peter Hoff. Eigenmodel: Semiparametric factor and regression models for symmetric relational data. *R package version*, 1, 2012.

## *Bibliography*

---

- [102] Pedro L Alonso, Graham Brown, Myriam Arevalo-Herrera, Fred Binka, Chetan Chitnis, Frank Collins, Ogobara K Doumbo, Brian Greenwood, B Fenton Hall, and Myron M Levine. A research agenda to underpin malaria eradication. *PLoS Med*, 8(1):e1000406, 2011. ISSN 1549-1676. 00409.
- [103] Mark S Handcock, Garry Robins, Tom AB Snijders, Jim Moody, and Julian Besag. Assessing degeneracy in statistical models of social networks. Technical report, Citeseer, 2003.
- [104] Joel G Breman. Eradicating malaria. *Science progress*, 92(1):1–38, 2009. ISSN 0036-8504. 00050.
- [105] The Centers for Population Health and Health Disparities Evaluation Working Group and Janet Okamoto. Scientific collaboration and team science: A social network analysis of the centers for population health and health disparities. *Translational Behavioral Medicine*, 5(1):12–23, March 2015. ISSN 1869-6716, 1613-9860. doi: 10.1007/s13142-014-0280-1. 00000.
- [106] Eldon Y. Li, Chien Hsiang Liao, and Hsiuju Rebecca Yen. Co-authorship networks and research impact: A social capital perspective. *Research Policy*, 42(9):1515–1530, November 2013. ISSN 00487333. doi: 10.1016/j.respol.2013.06.012. 00082.
- [107] Michael E Bales, Stephen B Johnson, and Chunhua Weng. Social network analysis of interdisciplinarity in obesity research. volume 870, 2008. 00015.

## *Bibliography*

---

- [108] Michael E Bales, Stephen B Johnson, Jonathan W Keeling, Kathleen M Carley, Frank Kunkel, and Jacqueline A Merrill. Evolution of coauthorship in public health services and systems research. *American journal of preventive medicine*, 41(1):112–117, 2011. ISSN 0749-3797. 00012.
- [109] Caroline S. Wagner and Loet Leydesdorff. Network structure, self-organization, and the growth of international collaboration in science. *Research Policy*, 34(10):1608–1618, December 2005. ISSN 00487333. doi: 10.1016/j.respol.2005.08.002. 00657.
- [110] Omwoyo Bosire Onyancha and Jan Resenga Maluleka. Knowledge production through collaborative research in sub-Saharan Africa: How much do countries contribute to each other’s knowledge output and citation impact? *Scientometrics*, 87(2):315–336, May 2011. ISSN 1588-2861. doi: 10.1007/s11192-010-0330-5.
- [111] Christian S Schmid and Bruce A Desmarais. Exponential Random Graph Models with Big Networks: Maximum Pseudolikelihood Estimation and the Parametric Bootstrap. *arXiv preprint arXiv:1708.02598*, 2017.
- [112] Thomas W. Valente, Kayo Fujimoto, Chih-Ping Chou, and Donna Spruijt-Metz. Adolescent Affiliations and Adiposity: A Social Network Analysis of Friendships and Obesity. *Journal of Adolescent Health*, 45(2):202–204, August 2009. ISSN 1054139X. doi: 10.1016/j.jadohealth.2009.01.007.

## *Bibliography*

---

- [113] Kayla de la Haye, Garry Robins, Philip Mohr, and Carlene Wilson. Obesity-related behaviors in adolescent friendship networks. *Social Networks*, 32(3):161–167, July 2010. ISSN 03788733. doi: 10.1016/j.socnet.2009.09.001.
- [114] Olga Kornienko, Katherine H Clemans, Dorothée Out, and Douglas A Granger. Hormones, behavior, and social network analysis: Exploring associations between cortisol, testosterone, and network structure. *Hormones and behavior*, 66(3):534–544, 2014. ISSN 0018-506X.
- [115] David R Hunter, Steven M Goodreau, and Mark S Handcock. Goodness of Fit of Social Network Models. *Journal of the American Statistical Association*, 103(481):248–258, March 2008. ISSN 0162-1459, 1537-274X. doi: 10.1198/016214507000000446.
- [116] Luka Kronegger, Franc Mali, Anuška Ferligoj, and Patrick Doreian. Collaboration structures in Slovenian scientific communities. *Scientometrics*, 90(2):631–647, February 2012. ISSN 1588-2861. doi: 10.1007/s11192-011-0493-8.
- [117] Caroline S Wagner, Irene Brahmakulam, Brian Jackson, Anny Wong, and Tatsuro Yoda. Science and technology collaboration: Building capability in developing countries. Technical report, RAND CORP SANTA MONICA CA, 2001.
- [118] World Health Organization. Global tuberculosis report 2016. 2016. ISSN 924156539X.
- [119] World Health Organization. MDG 6: Combat HIV/AIDS, malaria and other diseases. *Updated February*, 2013.

## *Bibliography*

---

- [120] Dennis Falzon, Holger J Schünemann, Elizabeth Harausz, Licé González-Angulo, Christian Lienhardt, Ernesto Jaramillo, and Karin Weyer. World Health Organization treatment guidelines for drug-resistant tuberculosis, 2016 update. *European Respiratory Journal*, 49(3):1602308, 2017. ISSN 0903-1936.
- [121] John B Lynch. Multidrug-resistant tuberculosis. *Medical Clinics*, 97(4):553–579, 2013. ISSN 0025-7125.
- [122] John Scott. *Social Network Analysis*. Sage, 2017. ISBN 1-5264-1225-X.
- [123] Shahadat Uddin, Liaquat Hossain, Alireza Abbasi, and Kim Rasmussen. Trend and efficiency analysis of co-authorship network. *Scientometrics*, 90(2):687–699, 2012. ISSN 0138-9130.
- [124] World Health Organization. *Atlas of African Health Statistics 2016: Health situation analysis of the African Region*. 2016. ISSN 9290232919.
- [125] Nagayoshi Nakazono, Kazuo Misue, and Jiro Tanaka. NeL 2: Network drawing tool for handling layered structured network diagram. pages 109–115. Australian Computer Society, Inc., 2006. ISBN 1-920682-41-4.
- [126] Fabian Odoni, Wolfgang Semar, and Elena Mastrandrea. Visualisation of Collaboration in Social Collaborative Knowledge Management Systems. In *Understanding Information Spaces. Proceedings of the 15th International Symposium of Information Science (ISI 2017)*, pages 386–388, Berlin, 13-15 March 2017.

## *Bibliography*

---

- [127] Xiaoming Liu, Johan Bollen, Michael L. Nelson, Herbert Van de Sompel, Jeremy Hussell, Rick Luce, and Linn Marks. Toolkits for visualizing co-authorship graph. page 404. ACM Press, 2004. ISBN 978-1-58113-832-0. doi: 10.1145/996350.996470.
- [128] Zdenek Horak, Milos Kudelka, Vaclav Snasel, Ajith Abraham, and Hana Rezankova. Forcoa.NET: An interactive tool for exploring the significance of authorship networks in DBLP data. pages 261–266. IEEE, October 2011. ISBN 978-1-4577-1133-6 978-1-4577-1132-9 978-1-4577-1131-2. doi: 10.1109/CASON.2011.6085955.
- [129] Eduardo M. Barbosa, Mirella M. Moro, Giseli Rabello Lopes, and J. Palazzo M. de Oliveira. VRRC: Web based tool for visualization and recommendation on co-authorship network (abstract only). page 865. ACM Press, 2012. ISBN 978-1-4503-1247-9. doi: 10.1145/2213836.2213975.
- [130] Miguel Grinberg. *Flask Web Development: Developing Web Applications with Python*. " O'Reilly Media, Inc.", 2014. ISBN 1-4919-4761-6.
- [131] Michael Bostock. D3. js. *Data Driven Documents*, 492:701, 2012.
- [132] Masashi Toyoda and Masaru Kitsuregawa. A system for visualizing and analyzing the evolution of the web with a time series of graphs. pages 151–160. ACM, 2005. ISBN 1-59593-168-6.

## *Bibliography*

---

- [133] Chaomei Chen and Leslie Carr. Visualizing the evolution of a subject domain: A case study. pages 449–452. IEEE Computer Society Press, 1999. ISBN 0-7803-5897-X.
- [134] Cesim Erten, Stephen G Kobourov, Vu Le, and Armand Navabi. Simultaneous Graph Drawing: Layout Algorithms and Visualization Schemes. *J. Graph Algorithms Appl.*, 9(1):165–182, 2005.
- [135] Gabor Csardi and Tamas Nepusz. The igraph software package for complex network research. *InterJournal, Complex Systems*, 1695(5):1–9, 2006.
- [136] Mark Newman. The physics of networks. *Physics today*, 61(11):33–38, 2008. ISSN 0031-9228.
- [137] Winston Chang, Joe Cheng, JJ Allaire, Yihui Xie, and Jonathan McPherson. Shiny: Web Application Framework for R. R package version 1.0. 3. 2017. URL <https://CRAN.R-project.org/package=shiny>.
- [138] W Chang and Barbara Borges Ribeiro. Shinydashboard: Create Dashboards with ‘Shiny’. *R package version 0.5*, 1, 2015.
- [139] M. E. Falagas, E. I. Pitsouni, G. A. Malietzis, and G. Pappas. Comparison of PubMed, Scopus, Web of Science, and Google Scholar: Strengths and weaknesses. *The FASEB Journal*, 22(2):338–342, September 2007. ISSN 0892-6638, 1530-6860. doi: 10.1096/fj.07-9492LSF. 00997.

## *Bibliography*

---

- [140] Catherine Matias and Vincent Miele. Statistical clustering of temporal networks through a dynamic stochastic block model. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, August 2016. ISSN 13697412. doi: 10.1111/rssb.12200. 00018.