Overlapping community detection in complex networks with memetic algorithms

by

Ademir Cristiano Gabardo

Ph.D. Thesis

submitted in partial fulfilment of the requirements for the Degree of

Doctor of Philosophy

Supervisor: Professor Regina Berretta Co-supervisor: Professor Pablo Moscato



The University of Newcastle

Faculty of Engineering and Built Environment School of Electrical Engineering and Computing

July, 2018

ii

© Copyright

by

Ademir Cristiano Gabardo

2018

iv

Overlapping community detection in complex networks with memetic algorithms

Statement of Originality

I hereby certify that the work embodied in the thesis is my own work, conducted under normal supervision.

The thesis contains no material which has been accepted, or is being examined, for the award of any other degree or diploma in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text. I give consent to the final version of my thesis being made available worldwide when deposited in the University's Digital Repository, subject to the provisions of the Copyright Act 1968 and any approved embargo.

Ademir Cristiano Gabardo 31st July 2018

vi

Acknowledgments

First, I would like to express my sincere gratitude to my advisers, Professor Regina Berretta and Professor Pablo Moscato, for their valuable knowledge and continuous support of my PhD study. Their guidance helped me throughout my research and while writing this thesis. I could not have imagined having better advisers for my PhD.

This research was supported by the Brazilian National Council for Scientific and Technological Development (CNPQ).¹

A research thesis may have a single author. However, this thesis is unquestionably not a journey of a single person there is so much to consider and there are so many to thank. These brief pages are not enough to express the enormous gratitude I have for all those somehow affected by my decision to pursue a PhD. With love and passion, I acknowledge my wife, ♡Juliana Gabardo♡, who supported me during this long journey. Without her love, I certainly would not have survived the challenging paths of academic research.

I also acknowledge my mother, Maria Francisca Gabardo, for her enormous courage in crossing the oceans and coming to Australia to bring me her heartwarming mother's love, and my father, Jair Gabardo, for supporting my decision and encouraging me to move forward. With deep respect, I acknowledge my father-in-law and my mother-in-law, who promptly supported our decision to move abroad.

I acknowledge my son, Adrian Gabardo, for embarking on this great adventure by my side and endure the challenges of moving overseas.

I am certain that the most valuable treasure in life are the good friends we meet along the way, and acknowledge my friends from the Priority Research Centre for Bioinformatics, Biomarker Discovery and Information-Based Medicine and Life and Economy Applications of Data Science research groups: Amer Abu Zaher, Amir

¹CNPq-http://cnpq.br

Salehipour, Claudio Sanhueza Lobos, Francia Jimenez, Heloisa Milioli, Jake Fitzsimmons, Inna Tischenko, Leila Moslemi Naeni, Luke Mathieson, Lukasz Olech, Marta A. Jimenez, Nader Mahmoudi, Natalie de Vries, Nasimul Noman and Shannon Fenn. I thank them for the friendly environment, the countless coffee breaks and the support.

Although distant, but always present, I acknowledge my friend and mentor, Professor Heitor Silvério Lopes, who always had a word of support and encouragement.

I also extend my deep respect and esteem to the other friends have also been part of this journey, especially Luiz Fernando Nunes, Lauro Cesar Galvão, Chidambaram Chidambaram and Osvalcir Peters.

> Ademir Cristiano Gabardo The University of Newcastle July 2018

Contents

A	cknov	vledgm	ents	/ii
Li	st of ⁻	Tables		iii
Li	st of I	Figures	×	′ii
Ał	ostrac	:t		vi
1	Intro	oductio	n	1
	1.1	Netwo	orks and applications	1
		1.1.1	Community structures in networks	4
		1.1.2	Overlapping community structures in networks	5
	1.2	Resea	rch problem	6
	1.3	Resea	rch objectives	7
	1.4	Organ	isation of the thesis	8
	1.5	List of	f publications and presentations	9
2	Con	cepts a	nd definitions	1
	2.1	Netwo	orks and graphs	11
	2.2	Graph	properties	4
		2.2.1	Node degree	4
		2.2.2	Degree distribution	5
		2.2.3	Centrality measures	5
		2.2.4	Node neighbourhood	9
		2.2.5	Clique	9
		2.2.6	Connected components	0
		2.2.7	Clustering coefficient	0
		2.2.8	Graph density	<u>21</u>
	2.3	Line g	raph	.2
	2.4	Netwo	ork models	3
		2.4.1	Erdös and Rényi model and random graphs	3

		2.4.2	Watts and Strogatz small-world model and the 'six degrees	0.4
		0.4.0		. 24
	0 5	2.4.3	Barabasi and Albert model and scale-free networks	. 26
	2.5			. 28
		2.5.1		. 29
	0.6	2.5.2		. 31
	2.6	Metric		. 32
		2.6.1		. 32
		2.6.2		. 34
	0 7	2.6.3	Other metrics for community detection	. 35
	2.7	Sumn	nary	. 37
3	A re	view of	f methods and benchmarks for community detection	. 39
	3.1	Metho	ods for overlapping community detection	. 39
		3.1.1	Clique-based methods	. 40
		3.1.2	Local-expansion methods	. 42
		3.1.3	Label-propagation methods	. 46
		3.1.4	Link-clustering methods	. 48
		3.1.5	Other strategies for overlapping community detection	. 52
	3.2	Metah	neuristics for community detection	. 52
		3.2.1	Genetic algorithms for community detection	. 53
		3.2.2	Memetic algorithms for community detection	. 59
		3.2.3	Other metaheuristic approaches	. 63
	3.3	Bench	marks for community detection	. 66
		3.3.1	Real-world benchmark networks for community detection .	. 66
		3.3.2	The Girvan and Newman benchmark	. 68
		3.3.3	LFR synthetically generated benchmark networks	. 68
	3.4	Summ	nary	. 71
4	MAI	DOC: A	memetic algorithm to detect overlapping communities	. 73
	4.1	The g	eneral framework of MADOC	. 74
	4.2	Line g	raphs for link-clustering	. 75
	4.3	Individ	dual representation	. 75
	4.4	A mul	ti-agent ternary tree population	. 76
	4.5	Popul	ation Initialisation	. 78
		4.5.1	Creating individuals with local-expansion	. 78
		4.5.2	Creating individuals with label-propagation	. 80
	4.6	Genet	ic operators	. 82
		4.6.1	Modularity recombination crossover	. 82
		4.6.2	Neighbouring communities mutation	. 83

	4.7	Select	ion and replacement	. 85
	4.8	Local-	search algorithms	. 86
		4.8.1	A local-seach algorithm for modularity maximisation	. 86
		4.8.2	A deterministic-update annealing for modularity	
			maximisation	. 87
	4.9	Comp	utational experiments and results	. 89
		4.9.1	Computational environment and parameters	. 89
		4.9.2	A comparison between the hill-climbing local-search and the	
			deterministic-update annealing	. 90
		4.9.3	Modularity maximisation results	. 91
		4.9.4	Experiments in synthetically generated overlapping	
			benchmark networks	. 93
	4.10	Concl	usion and remarks	. 95
5	Co-p	ourchas	ses for luxury brands at Amazon.com: A case study	. 97
	5.1	Introd		. 97
	5.2	Const	ructing the network	. 98
	5.3	Modu	larity maximization over the line graph	. 100
	5.4	Result	ts and discussion	. 101
	5.5	Concl	usion	. 106
6	Co-p	ourchas	ses of photographic products: A case study	. 107
	6.1	Introd		. 107
	6.2	Co-pu	rchasing networks	. 108
	6.3	Const	ructing the network	. 109
	6.4	Modu	larity maximization over the line graph	. 111
	6.5	Result	ts and discussion	. 111
	6.6	Concl	usion	. 116
7	MEN	1ELink:	: A memetic algorithm for link partition density maximisatior	ו 117
	7.1	The g	eneral framework of MEMELink	. 117
	7.2	Individ	Jual representation	. 118
	7.3	A mul	ti-agent ternary tree population	. 119
	7.4	Popula	ation initialisation	. 120
		7.4.1	Creating individuals with local-expansion	. 121
		7.4.2	Creating individuals with label-propagation	. 121
		7.4.3	Eigenvector centrality	. 122
	7.5	Genet	ic operators	. 123
		751	Uniform crossover	123
		7.0.1		. 120

	7.6	Local-	search	. 124
	7.7	Post-p	processing	. 125
	7.8	Comp	utational experiments and results	. 127
		7.8.1	Computational environment and parameters	. 127
		7.8.2	Instances of the LFR synthetically generated benchmark net-	
			works	. 128
		7.8.3	A comparison between MADOC and MEMELink	. 129
		7.8.4	The effect of network size (n) on NMI and time consumption	131
		7.8.5	The effect of mixing parameter μ and number of community	
			assignments (om) on NMI	. 133
		7.8.6	The effect of overlapping extension (on) and community	
			assignments (om) on NMI	. 136
		7.8.7	Comparing MEMELink with other algorithms	. 137
	7.9	Concl	usion and remarks	. 139
8	A st	orm of	swords: A case study	141
Ŭ	81	Introd		. 141
	8.2	The ne	etwork of characters in the novel 'A storm of swords'	142
	8.3	Result	ts of MADOC and discussion	143
	8.4	Result	ts of MEMELink and discussion	147
	8.5	Concl	usion	. 150
•	-			
9	Con	clusion	is and future work	. 151
	9.1	Concl		. 151
		9.1.1	MADOC	. 151
		9.1.2		. 152
		9.1.3	Overlapping communities of luxury brands	. 153
		9.1.4	Overlapping communities from a brand-centric point of view	. 153
		9.1.5	Overlapping communities of characters in a novel	. 154
	9.2	Future	e work	. 154
Bi	bliogı	raphy		. 157
	-			

List of Tables

2.1	A symmetric adjacency matrix representing the graph in Figure 2.3	13
3.1	Algorithms that employ cliques for overlapping community detec- tion. Where h_{max} is the number of number of pairs of maximal cliques that are neighbours, s_{max} is the number of maximal cliques, and $ C $ is the number of cliques.	41
3.2	Seed-expansion and local-expansion algorithms for overlapping com- munity detection where m represents the number of edges, n repres-	
3.3	ents the number of nodes and s is the number of communities Label-propagation algorithms for overlapping community detection where t is the number of iterations	45 47
3.4	Link-clustering algorithms for overlapping community detection where t is the number of iterations, s is the number of communities, k_{max} is the maximum degree of a node in the graph, n is the number of nodes, m is the number of edges in a graph and p_s is the population size	50
3.5	Genetic algorithms for disjoint and overlapping community detec-	56
3.6	Memetic algorithms for disjoint and overlapping community detec- tion ordered by year of publication or submission.	60
3.7	Bio-inspired and other approaches metaheuristics for disjoint and overlapping community detection.	64
3.8	Real-world networks frequently used to assess the efficiency of com- munity detection algorithms.	67
4.1	Comparison results for the minimum, average and maximum modu- larity between the local-search (LS) algorithm and the deterministic- update annealing (SA) for 50 individual runs for eight distinct bench- mark networks. Columns show the benchmark, the minimum, aver- age, maximum and the best-known modularity for both algorithms and the standard deviation.	90

4.2	Minimum, average, maximum and the standard deviation of the time in seconds spent by local-search and deterministic-update anneal-
4.3	ing for a single run in one individual
4.4	The average time spent in seconds and the standard deviation for 50 executions of MADOC in real-world benchmark networks 93
4.5	Synthetically generated networks used to assess MADOC's efficiency. Columns show the network, number of nodes and edges, mixing parameter μ , number of overlapping nodes (on) , number of com- munity assignments for the overlapping nodes (om) and the num- ber of communities $\#C$. The table also shows the number of nodes and edges in the corresponding line graph $L(G)$ and the maximum modularity produced by MADOC in $L(G)$
5.1	Network statistics for the graph G and its respective line graph $L(G)$. 100
5.2	Product count and percent coverage for the overlapping communit- ies C1 to C7 recovered from $L(G)$ to $G. \ldots \ldots$
5.3	Distribution of categories and gender according to the number of communities in which the products are assigned. The numbers at the left of the columns represent the product count, the targeted gender for the products are identified as <i>M</i> for male, <i>F</i> for female and <i>U</i> for unisex. The categories are ordered according to the number of products
5.4	The top 10 brands with the highest frequency among communities C1 to C7. (For short, in this table, Victoria's Secret = V.S., Dolce & Gabbana = D&G, Paris Hilton = P.H., Carolina Herrera = C.H., Calvin Klein = C.K, Paco Rabanne = P.R. and Emporio Armani = E.A.) 103
5.5	Columns indicate: Community A, its size, the percentage of nodes which overlaps with Community B, Community B size, the percent- age of nodes which overlaps with Community A and, the number of overlapping nodes between Communities A and B
6.1	Network statistics for the graph G and its respective line graph $L(G)$. 111
6.2	The number of products, the minimum, the average, the maximum price and the standard deviation for products in communities C1 to C11 for the Kodak co-purchasing network

6.3	The brands with largest number of products in Kodak's co-purchasing network.	. 113
7.1	The gain in NMI achieved with the post-processing for the bench- mark networks LFR1 to LFR6	. 126
7.2	LFR benchmark networks with 128 to 5,000 nodes, mixing parameter μ = 0.1 or 0.3, number of overlapping nodes corresponding to approximately 10% or 20% of the number of nodes in the network and community assignments (<i>om</i>) varying from 2 to 8. # <i>C</i> indicates the number of communities.	. 129
7.3	Comparative results between MADOC and MEMELink in a collection of 21 benchmark networks detailed in Table 7.2. Columns show the network, the average and maximum NMI and the standard deviation for both algorithms	120
7.4	Comparative results for the algorithms MADOC, SLPA and MEMELink for the benchmark networks LFR1 to LFR6, shown in Table 7.2. Column show the network and the average, maximum and the standard de- viation of NMI for both algorithms. We also report the average and maximum link partition density <i>D</i> and the standard deviation achieved	. 130
7.5	by MEMELink	. 132
	munities in the benchmark networks LFR1 to LFR6 for the algorithms SLPA, MEMELink and MADOC.	. 133
7.6	Comparative results of MEMELink and SLPA for the networks LFR7 to LFR34 detailed in Table 7.2. Columns show the network and the average, maximum and the standard deviation of link partition dens- ity <i>D</i> achieved by MEMELink and, the average, maximum and the	
7.7	standard deviation of NMI achieved by MEMELink and SLPA Comparative results of MEMELink and SLPA for the networks LFR14 to LFR20 and LFR35 to LFR41 shown in Table 7.2, varying the number of overlapping nodes (<i>on</i>) and the number of community assignments (<i>om</i>) for the overlapping nodes. Columns show the network and the average, the maximum and the standard deviation for the link partition density <i>D</i> achieved by MEMELink, and the average, the maximum and the standard deviation for the NMI achieved by SLPA and MEMELink.	. 135
8.1	The Pearson's, Spearman's and Kendall's correlation between the number of communities to which a character is associated with and	
	metrics of strength and centrality	. 145

8.2 The most central characters of the network "Storm of Swords" ordered according to the number of communities in which they appear. . . . 146

List of Figures

1.1	(Left), a 3D model of the worm <i>Caenorhabditis Elegans</i> . (Right) The network of protein-encoding genes of the worm <i>Caenorhabditis elegans</i> [257]. The colours represent the communities detected by our memetic algorithm presented in Chapter 4, showing groups of genes that interact with one another.	2
1.2	The network of 500 airports with the largest amount of traffic in the United States, the connections represent the traffic of passengers between airports. The size of the nodes represents the amount of traffic in the airport, and the colours represent the communities de- tected by our memetic algorithm presented in Chapter 4, showing groups of airports with intense traffic between each other	2
1.3	A network of 128 nodes divided into four communities represented by distinct colours and node shapes showing a higher density of con- nections inside communities and few connections between com- munities. The network is a modified instance of the Girvan and New- man benchmark presented in Section 3.3.2, and is edited and col-	5
	oured with the software yEd graph editor.	4
1.4	Network with 13 nodes organised in two overlapping communities. The network shows groups of students affiliated to different scholar clubs. Nodes six and seven are assigned to both communities, and are therefore overlapping. Network and figure are generated with the	
	software yEd graph editor.	6
2.1	The seven bridges of Königsberg are represented by a to g ; the island of <i>Kneiphof</i> is represented by A , and B , C and D , showing distinct portions of land divided by the Pregel River. The lines connecting A, B , C and D are the graph representation of the problem Euler proposed	10
22	Adapted from: The elements of network models - Brandes I IIrik et	١Z
<i>∠.</i> ∠	al. 'What is network science?' Network Science 1.01 (2013): 1-15.	12

2.3	(Left) A graph of six nodes joined by seven edges. (Right) The set of nodes and edges for the graph on the left.	12
2.4	 a) An unweighted undirected graph. b) A weighted undirected graph. c) An unweighted directed graph. d) A weighted directed graph. e) A mixed graph with undirected and directed edges, parallel edges and self-loop. 	13
2.5	The adjacency list corresponding to the graph of six nodes and seven edges from Figure 2.3.	14
2.6	(Left) The degree distribution of a random graph; the trend line shows the binomial distribution. (Right) The degree distribution of a real- world network representing the 500 airports with the largest amount of traffic in the US; the trend line shows a power-law distribution, where a few nodes have a high number of connections and most	
	nodes have few connections.	16
2.7	A graph with nine nodes and 15 edges where node 5 is the most cent- ral node connecting the two groups of nodes highlighted by distinct	16
<u> </u>	COIOURS.	10
2.0	of node 3 comprising nodes $\{2, 3, 4, 5\}$, highlighted in yellow.	19
2.9	A complete graph with five nodes in which every node connects to the other nodes in the network, forming a clique.	19
2.10	A graph with 11 nodes and 12 edges showing three connected components.	20
2.11	(Left) A sparse graph with 12 nodes and 21 edges with graph density $D(G) = 0.318$. (Centre) A dense graph with 12 nodes and 43 edges with graph density $D(G) = 0.652$. (Right) A completely connected graph with 12 nodes and 66 edges with graph density $D(G) = 1$	21
2.12	(Left) A graph G with eight nodes and 16 edges. (Right) The corresponding line graph $L(G)$ of G with 16 nodes and 54 edges. The edges between the nodes 1, 5 and 8 are highlighted in colour in the graph, and the corresponding nodes 1,5 and 5,8, and the edge between these pades are highlighted in colour in the line graph.	00
2.13	(Left) A graph with 79 nodes and 217 edges with average degree $\langle k \rangle$ = 5.494 generated with the Erdös and Rényi model. (Right) The bar graph of its degree distribution showing a <i>normal distribution</i> where	<u>Z</u> Z
	the majority of nodes have a degree similar to the average degree of the graph.	24

2.14	(Left) A graph with 75 nodes and 75 edges and average degree $\langle k \rangle =$ 2 generated with the Watts and Strogatz model. (Right) The respect- ive degree distribution for this graph. The majority of the nodes have the degree similar to the average degree in the network, and few	
2.15	nodes have a higher degree connecting to many nodes (Left) An example of a small-world graph with complex interconnections. (Right) The group of <i>X</i> 's acquaintances feedback into his own circle, normally eliminating new contacts. Figure adapted from; Travers, Jeffrey, and Stanley Milgram. 'The small-world problem', <i>Psy</i> -	25
2.16	chology Today 1 (1967): 61–67	26
2.17	$P(k) \sim k^{-3}$	27
	soft overlapping communities.	29
2.18	A network with 19 nodes divided into four disjoint communities	30
2.19	A network with 15 nodes divided into three overlapping communities.	31
3.1	A timeline highlighting the appearance of the major contributions in the overlapping community detection over the years: clique-based, local-expansion, label-propagation and link-clustering methods for overlapping community detection.	39
3.2	Two overlapping communities comprising two five-node cliques shar- ing four nodes.	41
3.3	The closed neighbourhood of the node '1' comprising by the node,	1 1
3.4	A small network after the first interaction of the algorithm. Each node shows the chance that it has to receive the label of its neigh- bour. Figure adapted from 'Finding overlapping communities in net- works by label propagation' by Steve Gregory, 2010, <i>New Journal of</i>	44
3.5	<i>Physics</i> , 103018, v.12, n.10	46
	L(G). Therefore, node five is assigned to both communities in G .	49

3.6	A network of 10 nodes and 25 edges divided into two edge com- munities. Nodes are assigned to communities of their adjacent edges. Node six is adjacent to edges in both communities. Therefore, node six is an overlapping node.	50
3.7	(Left) An example of the Girvan and Newman benchmark network divided into four communities highlighted by colour. (Right) At the top, a line graph showing the degree distribution; all nodes have the same degree $k = 16$. At the bottom is a bar chart showing the community size distribution—the four communities are of equal size with 32 nodes.	69
3.8	(Left) A LFR benchmark network with the following parameters: $n = 500$, $\tau_1 = 1$, $\tau_2 = 2$, $\langle k \rangle = 10$, max $\langle k \rangle = 25$ and $\mu = 0.1$. (Right) On top, the degree distribution for this network. Right, at the bottom, the community size distribution for this network. The communities are highlighted by distinct colours.	70
4.1	(Left) A graph $G = (V, E)$. (Centre) The line graph $L(G)$ of G divided into two communities that obtained the maximum modularity high- lighted by colours. (Right) The overlapping communities of nodes in G recovered from the edge communities.	75
4.2	(Left) An example of a graph G with seven nodes and 12 edges. (Centre) The corresponding $L(G)$ of G , with 12 nodes and 33 edges. In paren- thesis, we kept the labels of the nodes corresponding to the edges in G. (Right) The string-coding representation of the individual for $L(G)$;	
	the colours of nodes represent the disjoint communities of $L(G)$.	76
4.3 4.4	The multi-agent tertiary tree population employed by MADOC Communities of nodes generated by the node seed expansion method shown in Algorithm 4.2 for the Zachary karate club network. Square nodes represent seeds; the distinct colours highlight the community	//
	structure.	79
4.5	Communities of nodes generated by the label propagation method shown in Algorithm 4.3 for the Zachary karate club network. Square nodes represent seeds; the distinct colours highlight the community	
4.6	structure	81
	structures for a network recombined as a new individual (<i>offspring</i>)	00
47	A network divided into four communities labelled as C1 C2 C3 and	83
1.7	C4 and highlighted by distinct colours.	84

MADOC
A record of the database showing the product 'Pasha De Cartier By Cartier For Men' stored within the MongoDB document structure and
The distribution gender and categories on the Versace co-purchasing network.
Community size distribution for communities C1 to C7
Communities of products C1 to C7 in the co-purchasing network highlighted by distinct colors.
Distribution of nodes according to the number of communities they participate.
The predominant category of products and the targeted gender of items for communities C1 to C7.
The network of brands of co-purchases of luxury items on Amazon.com. The weight of the edges, shown by its thickness, represents the num- ber of co-purchases between two brands. The network is divided into four communities (a , b , c and d) highlighted by colours 105
A sample record of the database showing the product 'Kodak Ink Jet Photo Stickers' stored within the MongoDB document structure and its respective fields
(Left) Community size distribution for Kodak's co-purchasing net- work. (Right) The average, the minimum and the maximum price for products in USD \$ for communities C1 to C11 in the Kodak co-
purchasing network
Communities C1 to C11 highlighted by distinct colours in the Kodak co-purchasing network of photographic material
(Left) The communities C2 and C4 highlighted by colours with the overlapping nodes highlighted by a darker colour. (Centre) The word- cloud of products found in the overlapping section. (Right) The word- cloud of brands in the overlapping intersection of communities C2 and C4
(Left) The communities C11 and C4 highlighted by colours with the

6.6	(Left) The communities C1 and C11 highlighted by colours with the overlapping nodes highlighted by a darker colour. (Centre) The word-cloud of products found in the overlapping section. (Right) The word-cloud of brands in the overlapping intersection of communities C1 and C11.
6.7	(Left) The community C7 highlighted by colour. (Centre) The word- cloud of products in this community. (Right) The word-cloud of brands in the community C7
6.8	A detailed look at some nodes that are in a specific overlapping loc- ation. (Left) The central node clearly belongs to both communities, the algorithm was able to depict this structure while preserving the community structure. (Centre) The central node is overlapping differ- ent communities, where it shares several edges. (Right) the nodes sharing connections with other communities are overlapping, while the small community on the right is also depicted
7.1	(Left) A network comprised of 14 edges with two link communities a) and b) where node 4 is the overlapping node. (Right) The corresponding genotype using the LAR 119
7.2	A tertiary tree of 13 agents, comprising a pocket solution and six
7.3	Current solutions each
7.4	A small network with two well-defined, non-overlapping communit- ies, <i>a</i> and <i>b</i> , connected by a bridging edge
7.5	The effect of the number of overlapping nodes (on) , number of community assignments (om) and the mixing parameter μ over the NMI. 131
7.6	(Left) The maximum and average NMI achieved for networks LFR1 to LRF6 with 128 up to 4,096 nodes. (Right) The average time in seconds for the same networks
7.7	A comparison between the results of MEMELink and the SLPA al- gorithm for the four groups of benchmark networks, varying the mix- ing parameter μ and the number of community assignments for the overlapping nodes (<i>om</i>)

7.8	The effect of overlapping extension (<i>on</i>) and the number of com- munity assignments (<i>om</i>) over the NMI for SLPA and MEMELink 137
7.9	A comparison between the results achieved by MEMELink and the results presented by Xie et al. for nine distinct methods for overlapping community detection. Figure adapted from Xie, J., Kelley, S., & Szymanski, B. K. (2013). Overlapping community detection in networks: The state-of-the-art and comparative study. <i>ACM computing surveys (csur)</i> , 45(4), 43
8.1	The weighted network formed by the 107 characters and their inter- actions in the book 'storm of swords'
8.2	(Left) The solid line shows the degree distribution for the network 'A storm of swords', and the dotted blue line shows the trending line following a power law distribution. (Right) The solid line shows the centrality for each node in the network; additionally, the dotted blue line shows the trending line following a power law distribution
8.3	(Left) The count and percentage of characters according to the num- ber of communities. (Centre) The number of nodes for each com- munity. (Right) The chart with the distribution of the number of nodes and the number of community assignments
8.4	The eight overlapping communities detected by MADOC in the 'Storm of Swords' network
8.5	Communities C1 to C8 detected by MADOC highlighted by distinct colours in the 'Storm of Swords' network.
8.6	(Left) The count and percentage of characters according to the num- ber of communities. (Centre) The number of nodes for each com- munity. (Right) The chart with the distribution of the number of nodes and the number of community assignments
8.7	The nine overlapping communities detected by MEMELink in the 'Storm of Swords' network
8.8	Communities C1 to C9 detected by MEMELink highlighted by distinct colours in the 'Storm of Swords' network.
8.9	(Left) A comparison of the community sizes detected by MEMELink (in green) and MADOC (in blue). (Right) The number of communit- ies per node in the overlapping community structure detected by MEMELink and MADOC

List of Algorithms

3.1	The general framework of a genetic algorithm	53
3.2	The general framework of a simple local-search based memetic al-	
	gorithm.	59
4.1	The general framework of MADOC.	74
4.2	A node seed-expansion method to initialise the individuals in MADOC's	
	population.	79
4.3	The label-propagation method used to initialise individuals in MADOC's	
	population	80
4.4	Modularity recombination crossover used by MADOC	82
4.5	The neighbouring communities mutation used by MADOC	83
4.6	The basic hill-climbing local-search employed by MADOC	86
4.7	Deterministic-update annealing for modularity maximisation used by	
	MADOC as a local-search strategy	88
7.1	The general framework of MEMELink	118
7.2	Uniform crossover employed in MEMELink	123
7.3	The neighbouring mutation used by MEMELink.	124
7.4	The local-search strategy adopted by MEMELink.	125

ABSTRACT

Intricate relationships exist among the billions of individuals who form our society. The interactions that occur between thousands of genes within our cells comprise our biology. Millions of financial institutions perform billions of transactions daily, creating complex networks of entities. These are just a few examples of the many complex systems surrounding us.

Network science comprehends the collection, management, analysis, interpretation and presentation of complex systems that employ networks. Collections of interconnected items that are usually represented by a graph showing a set of nodes joined by edges.

Complex networks often present community structures where nodes preferentially link to one another. Examples of community structures include groups of friends in society, groups of co-functioning genes in gene networks and groups of similar products in co-purchasing networks, among many others. Detecting the community structure in networks offers important information about the organisation and functioning of such groups. For many phenomena represented by networks, communities can be overlapping, with nodes participating in multiple communities. For example, a person participates in several social organisations; a gene is related to different biological functions; a product can be sold in different markets.

Revealing the community structure in complex networks is no trivial task and can lead to a non-deterministic polynomial-time hardness (NP-hard) computational problem. In this thesis, we approach the overlapping community detection problem using memetic algorithms, metaheuristics that employ a population-based search and local-search inspired by Darwinian principles of natural evolution and Dawkins's notion of a meme defined as a unit of cultural evolution. We detail the construction of two different memetic algorithms, present computational results, compare our methods with other state-of-the-art metaheuristics and present applications of our methods as case studies.