Combinatorial optimization methods for the (α, β) -k Feature Set Problem

Amir Salehipour



Faculty of Engineering and Built Environment School of Electrical Engineering and Computing University of Newcastle, NSW, Australia

Combinatorial optimization methods for the (α, β) -k Feature Set Problem

Amir Salehipour

March 2019

This research was supported by an Australian Government Research Training Program (RTP) Scholarship. © Copyright by

Amir Salehipour March 2019

Statement of Originality

The thesis contains no material which has been accepted 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 hereby certify that the work embodied in the thesis is my own work, conducted under normal supervision. 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.

Amir Salehipour

28 March 2019

Dedication

To my parents, my wife and my son, and my sisters.

Acknowledgment

My first acknowledgment must go to my lovely wife, Dr Leila Moslemi Naeni. I want to thank her for all the support that she gave me during these long four years of my research.

Secondly, I would like to thank my supervisors, Prof. Pablo Moscato and Prof. Regina Berretta for providing me with this opportunity to pursue my PhD studies at the University of Newcastle, Australia. I would like to thank all the past and present members of the CIBM team (Center for Bioinformatics, Biomarker Discovery and Information-Based Medicine, the University of Newcastle) for their friendly support, chats, and helpful ideas. Starting with ex-member Dr Carlos Riveros, an amazing supportive person, and a true friend, who helped me many times with programming, and modeling. Then, Dr Renato Vimieiro, Dr Mateus Rocha de Paula and Dr Ahmed Shamsul Arefin for their friendship. Finally, thanks to all my colleagues and friends that supported me thorough this journey: Ademir, Amer, Claudio, Francia, Heloisa, Inna, Lukas, Luke, Marta, Mohammad, Nader, Nisha and Shannon.

I am very thankful to the University of Newcastle for my scholarships and for the opportunity to be part of an incredible academic environment. My Doctorate was fully funded by the University of Newcastle's Postgraduate Research Scholarships; those financial supports are very much appreciated.

I would like to thank my mother and my father, for the love and support that they have given me from the other side of the world. You gave me the strength and motivation to continue and finish my PhD.

> Amir Salehipour March 2019

Contents

1	Intr	roduction	1
	1.1	Introduction	1
	1.2	Optimization in Bioinformatics	2
	1.3	Feature selection in Bioinformatics	3
	1.4	Research question and research goals	5
	1.5	Thesis structure	6
	1.6	Conclusion	7
2	\mathbf{Res}	earch Problem and Literature Review	9
	2.1	Introduction	9
	2.2	Problem statement	10
	2.3	Feature selection methods	15
	2.4	Research motivation	18
	2.5	Conclusion	19
3	Mat	thematical Models and Properties	21
	3.1	Introduction	21
	3.2	Definitions and notations	23
	3.3	A bipartite graph representation $\ldots \ldots \ldots$	26
	3.4	Illustrative examples	27
	3.5	Mathematical models	29
		3.5.1 An integer program for the Min k (α, β)-k Feature Set Problem	29
		3.5.2 An integer program for the Max β $(\alpha,\beta)\text{-}k$ Feature Set Problem $\ .$	36
		3.5.3 An integer program for the Max Cover $(\alpha,\beta)\text{-}k$ Feature Set Problem	37
	3.6	Bounds	38
	3.7	Mathematical properties	40
	3.8	Conclusion	43
4	Solı	ition Methods for the Min k (α, β)-k Feature Set Problem	45
	4.1	Introduction	46
	4.2	The Set k-Cover Problem	47

	4.3	The Min k (α, β) -k Feature Set Problem
	4.4	Lower bounds
	4.5	A greedy construction algorithm 50
	4.6	A removal local search
	4.7	An exact+heuristic algorithm
		4.7.1 Obtaining a lower bound 54
		4.7.2 Obtaining a feasible solution
		4.7.3 Improving the feasible solution
	4.8	Computational results
		4.8.1 Computational results of real-world instances
		4.8.2 Computational results of standard instances of the Set Cover Problem . 61
		4.8.3 Computational results of random instances
	4.9	Conclusion
5	Solu	ution Methods for the Max β and Max Cover (α, β)-k Feature Set Prob-
	lem	s 95
	5.1	Introduction
	5.2	Literature review
	5.3	Pre-processing methods
	5.4	An exact+heuristic algorithm
		5.4.1 Initial solutions
		5.4.2 Improved solutions
	5.5	Computational results
		5.5.1 Computational results of real-world instances
		5.5.2 Computational results of random instances
	5.6	The Max Cover (α, β) -k Feature Set Problem
		5.6.1 Proposed solution method
		5.6.2 Computational results of real-world instances
		5.6.3 Computational results of random instances
	5.7	Conclusion
6	Con	cluding Remarks and Future Research 117
	6.1	Theoretical and computational contributions and outcomes 117
	6.2	Future research directions

List of Figures

2.1	A bipartite graph to represent the (α, β) -k Feature Set Problem	14
3.1	An undirected bipartite graph for the (α, β) -k Feature Set Problem	26
4.1	Solution representation for heuristic algorithms	50
4.2	Performance of solution methods for solving instances of SCP $(\alpha = \alpha_{min})$	65
4.3	Performance of solution methods for solving instances of SCP $(\alpha = \alpha_{med})$	66
4.4	Performance of solution methods for solving instances of SCP $(\alpha = \alpha_{max})$	67
4.5	Performance of solution methods for solving randomly generated instances	86
5.1	Optimality gap versus upper bound gap for real-world instances of Max Cover	
	(α, β) -k Feature Set Problem	112
5.2	Computation time of obtaining feasible, optimal, and upper bound solutions for	
	real-world instances of Max Cover $(\alpha,\beta)\text{-}{\bf k}$ Feature Set Problem	113
5.3	Gap between integer upper bound and objective function value of Max Cover	
	$(\alpha,\beta)\text{-}k$ Feature Set Problem	114

List of Tables

2.1	An example of a dataset with two classes of data	12
2.2	Building an instance of the (α, β) -k Feature Set Problem	13
3.1	Mathematical notations of the (α, β) -k Feature Set Problem	25
3.2	An example of (α, β) -k Feature Set Problem	25
3.3	The selected list of features for dataset DS	28
3.4	The selected list of features for dataset ADMF	30
3.5	An example of obtaining a lower bound for the Max β ($\alpha,\beta)\text{-}k$ Feature Set Problem	40
4.1	Real-world instances for the Min k (α, β)-k Feature Set Problem	58
4.2	Summary of computational results for solving real-world instances of Min k	
	(α, β) -k Feature Set Problem	59
4.3	Detailed computational results for solving real-world instances of Min k (α, β)-k	
	Feature Set Problem	60
4.4	Standard instances of the Set Cover Problem	62
4.5	Summary of computational results for solving weighted instances of Min k (α, β)-	
	k Feature Set Problem	62
4.6	Maximum computation time of LAGRASP algorithm	64
4.7	Paired-sample <i>t</i> -tests for comparing gap of solution methods	69
4.8	Pair-wise Wilcoxon Signed Rank tests for comparing times of solution methods	71
4.9	Detailed computational results for solving weighted instances of Min k $(\alpha,\beta)\text{-k}$	
	Feature Set Problem $(\alpha = \alpha_{min})$	72
4.10	Detailed computational results for solving weighted instances of Min k $(\alpha,\beta)\text{-k}$	
	Feature Set Problem $(\alpha = \alpha_{med})$	76
4.11	Detailed computational results for solving weighted instances of Min k $(\alpha,\beta)\text{-k}$	
	Feature Set Problem $(\alpha = \alpha_{max})$	80
4.12	Summary of computational results for solving random instances of Min k (α, β)-k	
	Feature Set Problem	84
4.13	Percent of heuristic solutions matching with CPLEX	85
4.14	Detailed computational results for solving random instances of Min k $(\alpha,\beta)\text{-k}$	
	Feature Set Problem	87

5.1	Real-world instances for the Max β (α , β)-k Feature Set Problem	104
5.2	Summary of computational results for solving real-world instances of Max β	
	$(\alpha,\beta)\text{-}{\bf k}$ Feature Set Problem $\hfill \ldots \hfill \hfill \ldots \hfill \hfill \ldots \hfill \hfill \ldots \hfill \hfill \hfill \hfill \ldots \hfill \hfil$	105
5.3	Detailed computational results for solving real-world instances of Max β $(\alpha,\beta)\text{-k}$	
	Feature Set Problem	106
5.4	Summary of computational results for solving random instances of Max β (α , β)-	
	k Feature Set Problem	107
5.5	Summary of computational results for solving real-world instances of Max Cover	
	$(\alpha,\beta)\text{-}{\bf k}$ Feature Set Problem $\hfill \ldots \hfill \hfill \hfill \ldots \hfill \hf$	110
5.6	Detailed computational results for solving real-world instance of Max Cover	
	$(\alpha,\beta)\text{-}{\bf k}$ Feature Set Problem $\hfill \ldots \hfill \hfill \ldots \hfill \hfill \ldots \hfill \hfill \ldots \hfill \hfill \hfill \hfill \ldots \hfill \hfil$	111
5.7	Summary of computational results for solving random instances of Max β (α , β)-	
	k Feature Set Problem	113
6.1	Contributions and outcomes of the research thesis	120

List of Algorithms

4.1	The multi Column Row Cover Construction heuristic (mCRCC) for Min k (α, β)-	
	k Feature Set Problem	52
4.2	The Removal Local Search (RLS) for Min k $(\alpha,\beta)\text{-}k$ Feature Set Problem	53
4.3	The exact+heuristic (EH) for Min k (α, β)-k Feature Set Problem	54
4.4	A linear programming relaxation-based algorithm	55
4.5	An algorithm for building feasible solutions for Min k $(\alpha,\beta)\text{-}k$ Feature Set Problem	56
5.1	Procedure of generating feasible lower bound solutions for the Max β (α , β)-k	
	Feature Set Problem	99
5.2	The exact+heuristic (EH) for Max β (α , β)-k Feature Set Problem	101
5.3	Procedure of generating a partially built solution for Max β (α , β)-k Feature Set	
	Problem	102

Abstract

This PhD research thesis proposes novel and efficient combinatorial optimization-based solution methods for the (α, β) -k Feature Set Problem. The (α, β) -k Feature Set Problem is a combinatorial optimization-based feature selection approach proposed in 2004, and has several applications in computational biology and Bioinformatics. The (α, β) -k Feature Set Problem aims to select a minimum cost set of features such that similarities between entities of the same class and differences between entities of different classes are maximized.

The developed solution methods of this research include heuristic and exact methods. While this research focuses on utilizing exact methods, we also developed mathematical properties, and heuristics and problem-driven local searches and applied them in certain stages of the exact methods in order to guide exact solvers and deliver high quality solutions. The motivation behind this stems from computational difficulty of exact solvers in providing good quality solutions for the (α, β) -k Feature Set Problem. Our proposed heuristics deliver very good quality solutions including optimal, and that in a reasonable amount of time.

The major contributions of the presented research include: 1) investigating and exploring mathematical properties and characteristics of the (α, β) -k Feature Set Problem for the first time, and utilizing those in order to design and develop algorithms and methods for solving large instances of the (α, β) -k Feature Set Problem; 2) extending the basic modeling, algorithms and solution methods to the weighted variant of the (α, β) -k Feature Set Problem (where features have a cost); and, 3) developing algorithms and solution methods that are capable of solving large instances of the (α, β) -k Feature Set Problem in a reasonable amount of time (prior to this research, many of those instances pose a computational challenge for the exact solvers).

To this end, we showed the usefulness of the developed algorithms and methods by applying them on three sets of 346 instances, including real-world, weighted, and randomly generated instances, and obtaining high quality solutions in a short time. To the best of our knowledge, the developed algorithms of this research have obtained the best results for the (α, β) -k Feature Set Problem. In particular, they outperform state-of-the-art algorithms and exact solvers, and have a very competitive performance over large instances because they always deliver feasible solutions, and obtain new best solutions for a majority of large instances in a reasonable amount of time.

Awards, Publications, and Outcomes

Part of the material presented in this research thesis has been already presented, and published in peer-reviewed conferences. The list of publications and presentations is provided below. It is worth mentioning that during my PhD studies I won one major award of the University of Newcastle for conducting an outstanding original research on the (α, β) -k Feature Set Problem.

Awards

- Winner of the "2015 Research Poster Prize Competition", awarded by Faculty of Engineering and Built Environment, University of Newcastle, 2015.
- 2. Two PGRSS travel grants (approximately, \$4,500), awarded by the Faculty of Engineering and Built Environment, University of Newcastle, 2016.

Conference proceedings

- "Tight lower bounds and a hybrid heuristic for a problem of selecting features", orally presented at EURO 2016 International Conference (peer-reviewed). Poznan, 3 – 6 July 2016.
- "An optimization approach towards selecting features in biological datasets", orally presented at 24th National Conference of the Australian Society for Operations Research (peer-reviewed). Canberra, 16 – 18 November 2016.

Working papers

- 1. "Efficient solution methods for the Min k (α, β)-k Feature Set Problem" (75% complete; will be submitted to an international journal in next few months).
- 2. "A heuristic algorithm for the (α, β) -k Feature Set Problem" (80% complete; will be submitted to a top tier journal in the next few weeks).