

Approximation Algorithms for Resource Allocation Optimization

by

Kewen Liao

A THESIS SUBMITTED FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY



THE UNIVERSITY
of ADELAIDE

School of Computer Science,
The University of Adelaide.

March, 2014

Copyright ©2014
Kewen Liao
All Rights Reserved

*To my parents and grandparents,
for their endless support.*

Contents

Contents	i
List of Figures	iii
List of Tables	iv
List of Acronyms	v
Abstract	viii
Declaration	ix
Preface	x
Acknowledgments	xii
1 Introduction	1
1.1 Research problems	2
1.2 Thesis aims and impact	5
1.3 Thesis results	8
1.4 Thesis structure	9
2 Preliminaries	10
2.1 Computational complexity	11
2.1.1 Computational problems	11
2.1.2 Languages, machine models, and algorithms	12
2.1.3 \mathcal{P} and \mathcal{NP}	14
2.2 LP and ILP	16
2.3 Approximation algorithms	20
2.3.1 Greedy and local search algorithms	22
2.3.2 LP rounding algorithms	24
2.3.3 Primal-dual algorithms	26

3	Discrete Facility Location	30
3.1	Uncapacitated Facility Location	31
3.1.1	LP formulations	31
3.1.2	Hardness results	33
3.1.3	Approximation algorithms	36
3.2	Other facility location problems	49
3.2.1	Fault-Tolerant Facility Location	51
3.2.2	Capacitated Facility Location	53
3.2.3	K Uncapacitated Facility Location	56
4	Unconstrained Fault-Tolerant Resource Allocation	59
4.1	Introduction	60
4.2	A greedy algorithm with ratio 1.861	63
4.3	A greedy algorithm with ratio 1.61	69
4.4	A hybrid greedy algorithm with ratio 1.52	75
4.5	A simple reduction to <i>UFL</i>	76
4.6	Capacitated <i>FTRA</i> _∞	78
4.7	Summary	79
5	Reliable Resource Allocation	80
5.1	Introduction	80
5.2	Primal-dual algorithms	83
5.3	The inverse dual fitting analysis	87
5.4	Minimum set cover: formalizing IDF	92
5.5	Reduction to <i>UFL</i>	94
5.6	Reduction to <i>FTRA</i> _∞	96
5.7	Summary	98
6	Constrained Fault-Tolerant Resource Allocation	100
6.1	Introduction	101
6.2	A unified LP-rounding algorithm	104
6.3	Reduction to <i>FTFL</i>	109
6.4	The uniform <i>FTRA</i>	113
6.5	The uniform <i>k-FTRA</i>	124
6.6	Summary	133
7	Conclusion	134
	Bibliography	137

List of Figures

1.1	A unified resource allocation model with input parameters \mathbf{p} , \mathbf{l} , \mathbf{r} , k , \mathbf{R} , and \mathbf{u} (to be explained in detail later) that capture various practical issues in resource allocation.	3
2.1	A star graph	24
3.1	Clusters constructed from G'	41
5.1	An example of the RRA model	82
6.1	An $FTRA$ instance with a feasible solution	101
6.2	Illustration of bounding the connection costs	109

List of Tables

3.1	<i>UFL</i> approximation results	36
-----	--	----

List of Acronyms

<i>Acronym</i>	<i>Meaning</i>
<i>UFL</i>	Uncapacitated Facility Location (problem)
<i>KM</i>	K-median (problem)
<i>ILP</i>	Integer linear program
<i>FTFL</i>	Fault-Tolerant Facility Location (problem)
<i>SCFL</i>	Soft Capacitated Facility Location (problem)
<i>k-UFL</i>	K-Uncapacitated Facility Location (problem)
<i>FTRA_∞</i>	Unconstrained Fault-Tolerant Resource Allocation (problem)
<i>CFTRA_∞</i>	Capacitated Unconstrained Fault-Tolerant Resource Allocation (problem)
<i>RRA</i>	Reliable Resource Allocation (problem)
<i>FTRA</i>	Constrained Fault-Tolerant Resource Allocation (problem)
<i>k-FTRA</i>	K-Constrained Fault-Tolerant Resource Allocation (problem)
<i>QoS</i>	Quality of service
<i>FPT</i>	Fixed parameter tractable
<i>CDN</i>	Content distribution/delivery network
<i>PM</i>	Physical machine
<i>VM</i>	Virtual machine
<i>OR</i>	Operations research
<i>TCS</i>	Theoretical computer science
<i>CC</i>	Computational complexity (theory)
<i>VCO</i>	Vertex cover optimization (problem)
<i>VCD</i>	Vertex cover decision (problem)
<i>TM</i>	Turing Machine
<i>FSC</i>	Finite state control
<i>RAM</i>	Random Access Machine
<i>CU</i>	Control unit
<i>PC</i>	Program counter
<i>SAT</i>	Boolean satisfiability (problem)

<i>Acronym</i>	<i>Meaning</i>
IS	Independent set (problem)
UVC	Unweighted vertex cover (problem)
LP	Linear program
CSC	Complementary slackness condition
<i>APX</i>	Approximable
<i>PTAS</i>	Polynomial time approximation scheme
<i>FPTAS</i>	Full polynomial time approximation scheme
TSP	Traveling salesman problem
AP-reduction	Approximation preserving reduction
PCP	Probabilistically checkable proof (theorem)
VC	Weighted vertex cover (problem)
SC	Set cover (problem)
JV	An approximation algorithm by Jain and Vazirani [75] for <i>UFL</i>
MP	An approximation algorithm by Mettu and Plaxton [113] for <i>UFL</i>
MMSV	An approximation algorithm by Mahdian <i>et al.</i> [104] for <i>UFL</i>
JMS	An approximation algorithm by Jain <i>et al.</i> [74] for <i>UFL</i>
CRR	Clustered randomized rounding (algorithm) in [41]
CSGA	Cost scaling and greedy augmentation (procedures)
MYZ	An approximation algorithm by Mahdian <i>et al.</i> [106] for <i>UFL</i>
<i>FLO</i>	Facility Location with Outliers (problems)
<i>MLFL</i>	Multi-level Facility Location (problems)
<i>OFL</i>	Online Facility Location (problems)
DR	Dependent rounding (technique)
LC	Laminar clustering (technique)
<i>CFL</i>	Capacitated Facility Location (problem)
<i>HCFL</i>	Hard Capacitated Facility Location (problem)
IG	Integrity gap (of an integer linear program)
<i>CFLS</i>	Capacitated Facility Location (problem) with splittable demands
<i>CFLU</i>	Capacitated Facility Location (problem) with unsplittable demands
<i>UniFL</i>	Universal Facility Location (problem)

<i>Acronym</i>	<i>Meaning</i>
<i>HCFLU</i>	Hard Capacitated Facility Location (problem) with unsplittable demands
<i>GAP</i>	Generalized assignment problem
LR	Lagrangian relaxation (technique)
LMP	Lagrangian multiplier preserving (property)
<i>FTFA</i>	Fault-Tolerant Facility Allocation (problem)
<i>FTFP</i>	Fault-Tolerant Facility Placement (problem)
SG-1	A star-greedy algorithm for $FTRA_{\infty}$
PD-1	A primal-dual algorithm for $FTRA_{\infty}$
SG-2	An improved star-greedy algorithm for $FTRA_{\infty}$
PD-2	An improved primal-dual algorithm for $FTRA_{\infty}$
MRR	Minimum reliability requirement
VLSI	Very-large-scale integration
PD-3	A primal-dual algorithm for RRA
APD-3	An accelerated primal-dual algorithm for RRA
IDF	Inverse dual fitting (technique)
ULPR	A unified LP-rounding algorithm for $FTRA$
PD-4	A primal-dual algorithm for $FTRA$
APD-4	An accelerated primal-dual algorithm for $FTRA$
SOC	Sum of contributions
AGA	Acceleration of greedy augmentation (procedure)
PK	Procedures for solving k - $FTRA$
BS	Binary search (procedure)
GP	Greedy pairing (procedure)
RR	Randomized rounding (procedure)
e.g.	For example
i.e.	That is
etc.	And so on
w.r.t.	With respect to
w.l.o.g.	Without loss of generality
s.t.	Such that

Abstract

Nowadays, data storage, server replicas/mirrors, virtual machines, and various kinds of services can all be regarded as different types of resources. These resources play an important role in today's computer world because of the continuing advances in information technology. It is usual that similar resources are grouped together at the same site, and can then be allocated to geographically distributed clients. This is the resource allocation paradigm considered in this thesis. Optimizing solutions to a variety of problems arising from this paradigm remains a key challenge, since these problems are \mathcal{NP} -hard.

For all the resource allocation problems studied in this thesis, we are given a set of sites containing facilities as resources, a set of clients to access these facilities, an opening cost for each facility, and a connection cost for each allocation of a facility to a client. The general goal is to decide the number of facilities to open at each site and allocate the open facilities to clients so that the total cost incurred is minimized. This class of the problems extends the classical \mathcal{NP} -hard facility location problems with additional abilities to capture various practical resource allocation scenarios.

To cope with the \mathcal{NP} -hardness of the resource allocation problems, the thesis focuses on the design and analysis of approximation algorithms. The main techniques we adopt are linear programming based, such as primal-dual schema, linear program rounding, and reductions via linear programs. Our developed solutions have great potential for optimizing the performances of many contemporary distributed systems such as cloud computing, content delivery networks, Web caching, and Web services provisioning.

Declaration

I, Kewen Liao, certify that this work 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. In addition, I certify that no part of this work will, in the future, be used in a submission for any other degree or diploma in any university or other tertiary institution without the prior approval of the University of Adelaide and where applicable, any partner institution responsible for the joint-award of this degree.

I give consent to this copy of my thesis, when deposited in the University Library, being made available for loan and photocopying, subject to the provisions of the Copyright Act 1968.

The author acknowledges that copyright of published works contained within this thesis (as listed on Page x) resides with the copyright holder(s) of those works.

I also give permission for the digital version of my thesis to be made available on the web, via the University's digital research repository, the Library catalogue, and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

Signature

Date

Preface

During my PhD study at the University of Adelaide from 2009 to 2013, I have produced four conference papers and two journal articles related to this thesis (see my homepage at <http://cs.adelaide.edu.au/~kewen> for my bio with a complete list of publications). The thesis topic is theoretical in nature and based on the content presented in the following papers.

Conference Publications:

- Kewen Liao and Hong Shen. Unconstrained and constrained fault-tolerant resource allocation. In *Proceedings of the 17th annual international conference on computing and combinatorics (COCOON)*, pages 555–566, Dallas, Texas, USA, 14-16 August 2011. Springer-Verlag, Berlin
- Kewen Liao and Hong Shen. Fast fault-tolerant resource allocation. In *Proceedings of 12th International Conference on Parallel and Distributed Computing, Applications and Technologies (PDCAT)*, pages 231–236, Gwangju, Korea, October 20-22 2011. IEEE
- Kewen Liao and Hong Shen. Approximating the reliable resource allocation problem using inverse dual fitting. In *Proceedings of the Eighteenth Computing: The Australasian Theory Symposium (CATS)*, page Vol. 128, Melbourne, Australia, January-February 2012. ACS, Sydney
- Kewen Liao, Hong Shen, and Longkun Guo. Improved approximation algorithms for constrained fault-tolerant resource allocation. In *Fundamentals of Computation Theory (FCT) - 19th International Symposium*, pages 236–247, Liverpool, UK, August 19-21 2013. Springer-Verlag, Berlin (extended version invited to the special issue of Theoretical Computer Science)

Journal Publications:

- Kewen Liao and Hong Shen. Lp-based approximation algorithms for reliable resource allocation. *The Computer Journal*, 57(1):154–164, 2014
- Kewen Liao, Hong Shen, and Longkun Guo. Constrained fault-tolerant resource allocation. *Theoretical Computer Science*, submitted in December 2013, at <http://arxiv.org/abs/1208.3835>, currently under review

Acknowledgments

This PhD study was a real challenge for me, and overcoming this challenge would not have been possible without the support of many people.

First and foremost, I would like to thank my principal supervisor Prof. Hong Shen, for his great supervision over the past few years. I am especially grateful for his trust in me working on tough theoretical problems. Without his constant advice and guidance, this thesis would not be even completed.

I am truly indebted and thankful to my supervisor A/Prof. Michael Sheng. He was also the supervisor of my honors thesis. He is the one who introduced me to research, inspired and encouraged me to pursue a PhD during my honors year. Without him, this thesis might not even exist.

I am sincerely grateful to my supervisor Emeritus Prof. Zbigniew Michalewicz, for his kind help, support, and encouragement in every aspect of my research.

I owe a huge debt of gratitude to my parents and grandparents, for their unconditional love, support, and encouragement. My most heartfelt thanks go to my parents for their spiritual and financial support during my eight years of study in Australia.

I would like to thank some of my colleagues and friends: Changjian, Donglai, Jeff, Lei, Li, Lina, Longkun, Mike, Scott, Shihong, Sim, Xiang, Xiaoqiang, Yidong, Yihong, Yong, Yongrui, Denny, Ke, Leo, Nick, Su, White, Xiaoming, Yibing, Yuguo, and Zilang, for their accompany, inspiration, and influence during my PhD journey.

Finally, I thank the University of Adelaide for providing me scholarship and the School of Computer Science for financially supporting my conference travels. I would also like to express my gratitude to the anonymous examiners of this thesis.