



UNIVERSITI PUTRA MALAYSIA

**APPLICATION OF OPTIMIZATION METHODS FOR SOLVING
CLUSTERING AND CLASSIFICATION PROBLEMS**

PARVANEH SHABANZADEH

IPM 2011 3

**APPLICATION OF OPTIMIZATION METHODS FOR SOLVING
CLUSTERING AND CLASSIFICATION PROBLEMS**

By

PARVANEH SHABANZADEH

**Thesis Submitted to the School of Graduate Studies, Universiti Putra
Malaysia, in Fulfilment of the Requirements for the Degree of
Doctor of Philosophy**

March 2011



Abstract of thesis presented to the Senate of University Putra Malaysia in fulfilment of the requirement for the degree of Doctor of Philosophy.

**APPLICATION OF OPTIMIZATION METHODS FOR SOLVING
CLUSTERING AND CLASSIFICATION PROBLEMS**

By

PARVANEH SHABANZADEH

March 2011

Chairman: Professor. Malik Hj. Abu Hassan, PhD

Faculty: Institute for Mathematical Research

Cluster and classification analysis are very interesting data mining topics that can be applied in many fields. Clustering includes the identification of subsets of the data that are similar. Intuitively, samples within a valid cluster are more similar to each other than they are to a sample belonging to a different cluster. Samples in the same cluster have the same label. The aim of data classification is to set up rules for the classification of some observations that the classes of data are supposed to be known. Here, there is a collection of classes with labels and the problem is to label a new observation or data point belonging to one or more classes of data. The focus of this thesis is on solving clustering and classification problems. Specifically, we will focus on new optimization methods for solving clustering and classification problems. First we briefly give some data analysis background. Then a review of different methods currently available that can be used to solve clustering and classification problems is also given.

Clustering problem is discussed as a problem of non-smooth, non-convex optimization and a new method for solving this optimization problem is developed. This optimization problem has a number of characteristics that make it challenging: it has many local minimum, the optimization variables can be either continuous or categorical, and there are no exact analytical derivatives. In this study we show how to apply a particular class of optimization methods known as pattern search methods to address these challenges. This method does not explicitly use derivatives, and is particularly appropriate when functions are non-smooth. Also a new algorithm for finding the initial point is proposed. We have established that our proposed method can produce excellent results compared to those previously known methods. Results of computational experiments on real data sets present the robustness and advantage of the new method. Next the problem of data classification is studied as a problem of global, non-smooth and non-convex optimization; this approach consists of describing clusters for the given training sets. The data vectors are assigned to the closest cluster and correspondingly to the set, which contains this cluster and an algorithm based on a derivative-free method is applied to the solution of this problem. The proposed method has been tested on real-world datasets. Results of numerical experiments have been presented which demonstrate the effectiveness of the proposed algorithm.

Abstrak tesis untuk dibentangkan kepada Senat Universiti Putra Malaysia bagi memenuhi syarat ijazah Doktor Falsafah.

PELAKSANAAN KAEDAH PENGOPTIMUMAN UNTUK MENYELESAIKAN
MASALAH CLUSTERING DAN KLASIFIKASI

Oleh

PARVANEH SHABANZADEH

Mac 2011

Pengerusi: Profesor. Malik Hj. Abu Hassan, PhD

Institusi: Instituti Penyelidikan Matematik

Analisis berkelompok dan pengelasan adalah topik data lombong yang menarik yang boleh digunakan dalam banyak bidang. Berkelompok termasuk pencaman bagi subset untuk data yang serupa. Sampel dalam suatu kelompok yang sah adalah lebih serupa antara satu dengan lain daripada sample kepunyaan suatu kelompok berbeza. Sampel dalam kelompok yang sama mempunyai sama label. Tujuan pengelasan data adalah menyediakan petua untuk pengelasan bagi beberapa pengamatan yang kelas data diandaikan dah tahu. Di sini, ada suatu pungutan kelas dengan label dan masalah adalah melabelkan suatu pengamatan baharu atau titik data dipunyai oleh satu atau lebih kelas data. Tumpuan tesis ini adalah menyelesaikan masalah berkelompok dan penelasan. Khususnya, kita akan tumpukan ke atas kaedah pengoptimuman baharu untuk menyelesaikan masalah berkelompok dan pengelasan. Mula mula kita berikan beberapa latarbelakang data analisis. Kemudian kita berikan suatu sorotan kaedah

berbeza yang sediaada masa kini yang boleh digunakan untuk menyelesaikan masalah berkelompok dan pengelasan. Masalah berkelompok dibincangkan sebagai suatu masalah penngoptimuman tak licin, tak cembung dan satu kaedah baharu bagi menyelesaikan masalah pengoptimuman ini di bangunkan. Masalah pengoptimuman ini mempunyai suatu bilangan ciri yang membuatnya mencabar: ia ada banyak minimum setempat, pembolehubah pengoptimuman boleh jadi samada selanjar atau berkategori, dan tiada terbitan beranalitik. Dalam pengajian ini kita tunjukkan bagaimana untuk menggunakan suatu kelas khusus bagi kaedah pengoptimuman yang dikenali sebagai kaedah gelintaran corak untuk mengusulkan cabaran ini. Kaedah ini tidak menggunakan terbitan secara tak tersirat dan adalah sesuai apabila fungsi adalah tak licin. Juga suatu algoritma baharu untuk mencari titik awal dicadangkan. Kita telah buktikan bahawa kaedah baharu yang dicadangkan boleh mengeluarkan keputusan cemerlang berbanding dengan kaedah terdahulu. Keputusan ujikaji pengkomputeran di atas set data nyata mengesahkan keteguhan dan kebaikan bagi kaedah baharu. Kemudian masalah pengelasan data dikaji sebagai masalah pengoptimuman sejagat, tak licin dan tak cembung; pendekatan ini mengandungi kelompok perihalan untuk set latihan yang diberi. Vektor data diumpukkan kepada kelompok terhampir dan berpadanan kepada set yang mengandungi kelompok ini dan algoritma berasaskan ke atas kaedah bebas terbitan digunakan kepada penyelesaian masalah ini. Kaedah yang dicadang telah diuji ke atas set data dunia nyata. Keputusan ujikaji berangka telah dikemukakan dan menunjukkan keberkesanan algoritma yang dicadangkan.

ACKNOWLEDGEMENTS

At first I want to thank Allah for all of things that he has given in my life and then I offer my sincerest gratitude to my chairman, Professor Dr. Malik Hj. Abu Hassan who has supported me throughout my thesis from the initial to the final level with his patience and knowledge whilst allowing me the room to work in my own way. I would like to express my deep and sincere gratitude to my co-supervisor, Dr. Leong Wah June. His wide knowledge and his logical way of thinking have been of great value for me. His encouraging, detailed and constructive comments have enabled me to develop an understanding of the subject. I am also grateful to Dr. Mansor Monsi for serving in the supervisory committee.

I offer my regards to the Head of Institute, academic and general staffs of the Institute for Mathematical Research, University Putra Malaysia, who supported me in any respect during the completion of the project.

I owe great thanks to my precious parents for all things that they gave me or taught me. Without their encouragements, understandings and supports I would never have made any success and also it would have been impossible for me to finish this work. I owe my loving thanks to my dear sister and darling brother, for their loving supports and personal guidance.

I am greatly indebted to my darling law parents for financial and spiritual supporting them without their supports I would never have finished my PhD.

I owe special thanks to my dear's husband Kamyar Shameli and my sweetheart daughter Armita for all their helping in during my study without them I would never able to finish my PhD.

APPROVAL SHEET 1

I certify that Examination Committee has met on date of viva voce to conduct the final examination of Parvaneh Shabanzadeh on her Degree of Doctor of Philosophy thesis entitled “Investigation of Clustering and Classification Problems via Optimization Method” in accordance with Universiti Pertanian Malaysia (Higher Degree) Act 1980 and Universiti Pertanian Malaysia (Higher Degree) Regulations 1981. The Committee recommends that the student be awarded the Degree of Doctor of Philosophy.

Member of the Examination committee were as follows

FUDZIAH BINTI ISMAIL, PhD

Assoc. Prof
Faculty of Science
Universiti Putra Malaysia
(Chairman)

DATO' MOHAMED SULEIMAN, PhD

Professor
Faculty of Science
Universiti Putra Malaysia
(Internal Examiner)

NORIHAN MD. ARIFFIN, PhD

Assoc. Prof
Faculty of Science
Universiti Putra Malaysia
(Internal Examiner)

DIPTESH GHOSH, PhD

Professor
Indian Institute of Management
Vastrapur, Ahmedabad, Gujarat, India
(External Examiner)

HASANAH MOHD GHAZALI, PhD

Professor and Dean
School of Graduate Studies
Universiti Putra Malaysia

Date:



This thesis was submitted to the Senate of Universiti Putra Malaysia and has been accepted as fulfilment of the requirement for the degree of Doctor of Philosophy. The members of Supervisory committee were as follows:

Malik B Hj Abu Hassan, PhD

Professor
Faculty of Science
Universiti Putra Malaysia
(Chairman)

Leong Wah June, PhD

Assoc. Prof
Faculty of Science
Universiti Putra Malaysia
(Member)

Mansor B Monsi, PhD

Science Lecturer
Faculty of Science
Universiti Putra Malaysia
(Member)

HASANAH MOHD GHAZALI , PhD

Professor and Dean
School of Graduate Studies
Universiti Putra Malaysia

Date:



DECLARATION

I hereby declare that the thesis is my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously, and it is not concurrently, submitted for any other degree at University Putra Malaysia or at any institutions.

PARVANEH SHABANZADEH

Date: 29 March 2011

TABLE OF CONTENTS

	Page
ABSTRACT	ii
ABSTRAK	iv
ACNOWLEDGEMENTS	vi
APPROVAL	vii
DECLARATION	ix
LIST OF TABLES	xiii
LIST OF FIGUERS	xv
LIST OF NOTATIONS	xvi

CHAPTER

1	INTRODUCTION	1
1.1	Background of the problem	1
1.2	Classification	3
1.2.1	Classification Applications	4
1.3	Clustering	5
1.3.1	Definition of Clusters	8
1.3.2	Clustering Applications	9
1.4	Objective of the Research	10
1.5	Outline of Thesis	11
2	LITERATURE REVIEW	13
2.1	Introduction	13
2.2	Records and Attributes	13
2.3	Distance and Similarity Measures	14
2.4	Preprocessing	19
2.4.1	Scaling	19
2.5	Background of Clustering	22
2.5.1	Connectionist techniques	26
2.5.2	k-Methods	27
2.5.3	Search-based Clustering Algorithms	30



2.6	Background of Classification	37
2.6.1	Support Vector Machine	39
2.7	Feature Selection	41
2.8	Conclusion	42
3	OPTIMIZATION ALGORITHM FOR SOLVING CLUSTERING PROBLEMS	43
3.1	Introduction	43
3.2	Optimization approach to clustering problem	46
3.2.1	An Optimization Algorithm for Solving Clustering Problems	47
3.3	Pattern Search Method	49
3.3.1	Generalized Pattern Search Algorithm	50
3.4	Results and Discussion	54
3.5	Conclusion	59
4	A NEW OPTIMIZATION METHOD FOR SOLVING CLUSTERING PROBLEMS	60
4.1	Introduction	60
4.2	The Optimization Method to Clustering	62
4.3	An Optimization Algorithm for Solving Clustering Problem	64
4.3.1	A New Algorithm for Finding the Starting Point	66
4.4	Solving Optimization Problems	67
4.5	Results of Numerical Experiments and Discussions	72
4.6	Conclusions	82
5	OPTIMIZATION METHOD FOR SOLVING CLASSIFICATION PROBLEM	84
5.1	Introduction	84
5.2	The optimization algorithm to classification problem	86
5.3	Feature selection algorithm	90
5.4	Method of local optimization	93

5.5	Results and Discussions	96
5.6	Conclusion	105
6	A NEW METHOD FOR SOLVING SUPERVISED DATA CLASSIFICATION PROBLEMS	107
6.1	Introduction	107
6.2	New Optimization Algorithm for Solving Classification Problem	109
6.3	Solving Optimization Problems	113
6.3.1	The MADS algorithm	114
6.4	Results of numerical experiments	116
6.5	Conclusion	126
7	CONCLUSION AND RECOMMENDATION FOR FUTURE RESEARCH	127
7.1	Conclusion	127
7.2	Future Studies	128
	REFERENCES	130
	APPENDIX	140
	BIODATA OF STUDENT	143
	LIST OF PUBLICATIONS	144