# Fuzzy linear regression with global continuous optimization

By

## M. Hadi Mashinchi

I certify that the work in this thesis entitled **"Fuzzy Linear Regression with Global Continuous Optimization"** has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree to any other university or institution other than Macquarie University.

I also certify that the thesis is an original piece of research and it has been written by me. Any help and assistance that I have received in my research work and the preparation of the thesis itself have been appropriately acknowledged.

In addition, I certify that all information sources and literature used are indicated in the thesis.

M. Hadi Mashinchi

# Acknowledgements

I am very grateful to Prof. Mehmet Orgun for his guidance and support as my supervisor. His supervisory role, experience, and patience were inspirational to me. I also want to thank all the academics with whom I was in touch throughout my research career to date; Abhaya Nayak, Rajan Shankaran, Yan Wang, Mark Dras, Len Hamey, Scott McCallum, Xingquan (Hill) Zhu, Albert Zomaya, Frances Griffin, Mashaallah Mashinchi, Witold Pedrycz, Chi-Tsuen Yeh, and Nikola Kasabov. I am also grateful to David A. Barda who is both my friend and mentor and I have learnt a lot from him during my studies and working with him.

Special thanks go to the admin staff at Macquarie Univesity, graduate students and other people who have helped me throughout my research career, particularly Melina Chan, Jane Yang, Meredith McGregor, Donna Hua, Jackie Walsh, Sylvian Chow, Raina Kim, Aldrinine Creado, Peter Curran, Reza Sepahi, Armin Hezart, Lei Li and Kerstin Klemisch.

I would also like to thank Andy Connor, Sandhya Samarasinghe, Siti Mariyam HJ. Shamssudin, and Bernard DeBaets who hosted me at their institutions during my PhD studies. I am also very grateful to Macquarie University for funding my PhD research.

The last but not the least, my very special thanks go to my father, mother and brother whom without their never stopping friendship and support I would not be able to be what I am now.

# List of Publications

- M. H. Mashinchi, M. A. Orgun, M. Mashinchi, and W. Pedrycz. *A tabu-harmony search based approach to fuzzy linear regression*. IEEE Transactions on Fuzzy Systems 19(3), 432-448 (2011).

- M. H. Mashinchi, M. A. Orgun, and W. Pedrycz. *Hybrid optimization with improved tabu search*. Applied Soft Computing 11(2), 1993-2006 (2011).

- M. H. Mashinchi, L. Li, M. A. Orgun, and Y. Wang. *The prediction of trust rating based on the quality of services using fuzzy linear regression*. In IEEE International Conference on Fuzzy Systems, 1953-1959, IEEE Computer Society Press (2011).

- M. H. Mashinchi, Mehmet A. Orgun, and M. R. Mashinchi. *A least square approach for the detection and removal of outliers for fuzzy linear regression*. In World Congress on Nature & Biologically Inspired Computing, 134-139, IEEE Computer Society Press (2010).

- M. H. Mashinchi, M. A. Orgun, and M. Mashinchi. *Solving fuzzy linear regression with hybrid optimization*. In 16th International Conference on Neural Information Processing, Lecture Notes in Computer Science (LNCS), 336-343, Springer- Verlag Berlin Heidelberg (2009).

# Abstract

Finding the global optimum of an unknown system has attracted a great deal of interest in many engineering problems. In this setting, meta-heuristics are very common and efficient approaches for solving complex real-world problems in Global Continuous Optimization Problems (GCOPs) as they can approximate solutions without any need for mathematical assumptions such as differentiability. The application of global continuous optimization methods is essential in many engineering applications where an optimization problem has certain properties such as *unreliable derivatives* and/or *black-box nature*. Meta-heuristic based optimizations, as one of the promising approaches in global continuous optimization, have a slow rate of convergence. Hybridization frameworks are investigated as a potential way of enhancing the optimization speed, and the quality of solutions.

Fuzzy linear regression analysis is a powerful tool to model the input-output relationship for forcasting purposes or studying the behavior of the data. The existing challenges in fuzzy linear regression are, dealing with non-transparent fitness measures, outlier detection and spread increasing problem. The application of global continuous optimization is investigated to tackle these issues. We propose an Unconstrained Global Continuous Optimization (UGCO) method based on tabu search and harmony search to support the design of Fuzzy Linear Regression models (FLR). The proposed approach offers the flexibility of using any kind of an objective function based on the client's requirements or requests and the nature of the data set, and then attains its minimum error.

Fuzzy linear analysis may lead to an incorrect interpretation of data in case of being incapable of dealing with outliers. Both basic probabilistic and least squares approaches are sensitive to outliers. In order to detect outliers, we propose a two stage least squares approach based on global continuous optimization which outperforms some issues that exist in other methods. In both the first and second phases,

the minimization of the model fitting measurement is achieved by hybrid optimization which gives us the flexibility of using any type of model fitting measures regardless of being continuous, differentiable, or transparent.

Some of the fuzzy linear regression models suffer from constantly increasing spreads of the outputs with the increase in the magnitude of the inputs. Such models are known to have the so-called spread increasing problem. We introduce a model, obtained by the application of hybrid optimization, which is capable of having variable spreads for different input variables regardless of their magnitude. The proposed approach is also compared and contrasted with other models in terms of the number of parameters, the flexibility of spreads, and errors.

# Contents

# List of Figures

# List of Tables