

Discovery of Data Models using Genetic Programming

W.J.L Nuwan Wijayaweera

109161N



University of Moratuwa, Sri Lanka
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

Faculty of Information Technology

University of Moratuwa

October 2012

Discovery of Data Models using Genetic Programming

W.J.L Nuwan Wijayaweera

109161N



University of Moratuwa, Sri Lanka
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

Dissertation submitted to the Faculty of Information Technology, University of Moratuwa, Sri Lanka for the partial fulfillment of the requirements of the Degree of MSc in Artificial Intelligence

October 2012

Declaration

I declare that this dissertation does not incorporate, without acknowledgment, any material previously submitted for a Degree or a Diploma in any University and to the best of my knowledge and belief, it does not contain any material previously published or written by another person or myself except where due reference is made in the text. I also hereby give consent for my dissertation, if accepted, to be made available for photocopying and for interlibrary loans, and for the title and summary to be made available to outside organization.

W.J.L.N Wijayaweera

Name of Student

Signature of Student

Date:



University of Moratuwa, Sri Lanka
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

Supervised by

Prof. Asoka S. Karunananda

Name of Supervisor(s)

Signature of Supervisor(s)

Date:

Dedication

To My Grandmother



University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

Acknowledgements

First and foremost I would like to express my sincere gratitude to my supervisor Prof. Asoka S. Karunananda for the continuous guidance and support. I would have been lost without his guidance, encouragement, sound advice, and good ideas. I am grateful to all of those who supported me in any respect during this project.

Lastly, I would like to thank my grandmother, my family and all my friends and loved ones for all their love and encouragement.

W.J.L.N Wijayaweera

October 2012



University of Moratuwa, Sri Lanka
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

Abstract

The field of Genetic Programming in Artificial Intelligence strives to get a computer to solve a problem without explicitly coding a solution by a programmer. Genetic Programming is a relatively new technology, which comes under automatic programming. After the initial work by John R. Koza in genetic programming, many researches have been done to discover data models in various datasets. These works have been rather domain specific and little attention have been given to develop generic framework for modeling and experimenting with genetic programming solutions for real world problems.

A project has been launched to develop a visual environment to design and experiment with genetic programming solutions for real world problems. It is named as GPVLab to mean its ability to facilitate the discovery of data models for real world problems through a wide range of experiments. GPVLab takes any numerical dataset as the input. Also the user should select the reference column before the main process. If user has not specified a column, the system will automatically take the last column as the reference column. This system has two possible ways to feed datasets into the system. There is an inbuilt facility to manually enter all the data. Furthermore this system facilitates data loading from comma separated value (*.csv) files. Output of this system is an evaluable expression. Users can initiate the main process by feeding data and selecting the reference column. Then the system runs the genetic programming process by generating populations of expressions and evaluating them to find their fitness. Finally the system determines the best fit model for the dataset. If the system has found a perfect solution or the maximum number of generations has exceeded, then the system stops the process and output the best so far model. This system outputs an evaluable expression as the data model in reverse polish notation (RPN). If the resultant expression only contains basic functions, GPVLab automatically converts the resultant evaluable expression in RPN into more human readable infix notation. System is developed for any person who needs to discover a model out of a collected numerical dataset. Advanced users with the knowledge of Genetic Algorithms or Genetic Programming can use advanced settings for better results. Nevertheless the default settings will work for most of the problems. Knowledge about Genetic Programming is not a necessity.

The system has been developed using C# language with .NET framework 4.0. GPVLab also extends the AForge.NET framework to accommodate data with arbitrary number of attributes and to remove noise in data. The system has an option to use two function sets for discovering process, namely basic and extended. Basic function set contains operators such as addition, subtraction, multiplication and division. Extended function set has the ability to generate models consisting of square root (sqrt), sine (sin), cosine (cos), logarithms (ln) and exponential (exp) in addition to the basic operators. Upon completion of the discovering process, the system immediately allows the user to evaluate the model by providing required parameters. The system has the facility to save resultant models and access and evaluate them via library as required. Model library is developed using SQL compact Edition database, which does not require SQL Server instance to operate. Hence, this software is highly portable and can be installed or run in any computer with .NET Framework 4.0 installed.

GPVLab has been compared with WEKA as the main evaluation. A real world noisy dataset with eight columns has been used as the main input dataset. This main experiment has proved that the error rate of the solution generated by WEKA falls between -93.74% and 52.00% but the error rate of the solution generated by GPVLab falls between -24.15% and 24.51%. Further GPVLab has successfully discovered data models in simple datasets including square root of a number, addition of three numbers and a dataset with ten columns which has a known data model. All these solutions were achieved in less than 150 generations. The experiment of finding the square root function has been done using the extended function set and it directly provided the answer using 'sqrt' function within the first generation. Experiments performed by tweaking advanced settings showed that all the required facilities are there in GPVLab to experiment with genetic programming problems. Furthermore the results obtained through users with no knowledge about genetic algorithms or genetic programming, proved that this can be a really good tool for the researches in non technical fields as well. GPVLab has achieved all the objectives of this project. According to the main evaluation it is evident that GPVLab can generate solutions which provide better results in 56% of the time. It is concluded that GPVLab can be used to model genetic programming application very conveniently.

Contents

	Page
Chapter 1 - Introduction	01
1.1 Introduction	01
1.2 Background and Motivation	01
1.3 Aim	02
1.4 Objectives	02
1.5 Resource Requirements	03
1.6 Summary	03
Chapter 2 - Related Work in Discovering Data Models	04
2.1 Introduction	04
2.2 Current Approaches	04
2.3 Genetic Algorithms and Automatic Programming	05
2.4 Genetic Programming	07
2.5 Genetic Programming in Symbolic Regression	09
2.6 Similar Products	10
2.7 Summary	10
Chapter 3 – Genetic Programming	11
3.1 Introduction	11
3.2 Genetic Programming	11
3.3 Summary	12
Chapter 4 - Discovery of Data Models using Genetic Programming	13
4.1 Introduction	13
4.2 Proposed Solution	13
4.2.1 Inputs	13
4.2.2 Output	13
4.2.3 Process	14
4.2.4 Users	14
4.2.5 Features	14
4.3 Summary	15

Chapter 5 - Design of Data Model Discovery System	16
5.1 Introduction	16
5.2 Analysis and Design	16
5.2.1 Initialization	16
5.2.1.1 Preparing the Dataset	17
5.2.1.2 Selecting the Reference Column	17
5.2.1.3 Adjusting Advanced Settings	17
5.2.2 Genetic Programming Process	18
5.2.3 Evaluation of Resultant Expression	18
5.3 Summary	19
Chapter 6 - Implementation of GPVLab	20
6.1 Introduction	20
6.2 Initialization	20
6.2.1 Utility Classes	21
6.2.2 Preparing the Dataset	23
6.2.3 Selecting the Reference Column	24
6.2.4 Adjusting Advanced Settings	24
6.3 Genetic Programming Process	25
6.4 Evaluation of Resultant Expression	32
6.5 Summary	34
Chapter 7 – Evaluation	35
7.1 Introduction	35
7.2 Main Evaluation	35
7.3 Evaluate Discovery of Data Models	42
7.4 Evaluate Noise Reduction	47
7.5 Using GPVLab for Genetic Programming Experiments	55
7.6 Applicability in Non Technical Domains	56
7.7 Summary	57

Chapter 8 – Conclusion and Further work	58
8.1 Introduction	58
8.2 Conclusion	58
8.2.1 Discover Data Models from Any Numerical Dataset	58
8.2.2 Discover Data Models from Noisy Datasets	58
8.2.3 A Visual Environment to Experiment with Genetic Programming Problems	59
8.2.4 Usability without Genetic Programming Knowledge	59
8.3 Problems Encountered	60
8.4 Limitations	60
8.5 Further Work	60
8.6 Summary	61
Reference	62
Appendix A: Detailed Design Diagram	64
A.1 Introduction	64
A.2 A Detailed Design Diagram of GPVLab	64
Appendix B: Genetic Programming Process	65
B.1 Introduction	65
B.2 Genetic Programming Process	65
Appendix C: How GPVLab Works	67
C.1 Introduction	67
C.2 How GPVLab Works	67
Appendix D: Main Dataset for Evaluation	75
D.1 Introduction	75
D.2 Main Dataset for Evaluation	75
Appendix E: Dataset for Noise Reduction Experiment	78
E.1 Introduction	78
E.2 Dataset for Noise Reduction Experiment	78

List of Figures

	Page
Figure 5.1 – Top level design diagram of the proposed system	17
Figure 7.1 – Explorer window after completion – Main Evaluation	37
Figure 7.2 – WEKA Explorer window after completion – Main Evaluation	39
Figure 7.3 – Explorer window after completion – Addition of numbers	43
Figure 7.4 – Explorer window after completion – Square root of a number	45
Figure 7.5 – Explorer window after completion – Ten column dataset	46
Figure 7.6 – Evaluate Data Model dialog – Ten column dataset	47
Figure 7.7 – Explorer window after noise reduction experiment.	50
Figure 7.8 – Explorer window – Ten column dataset with noise	54
Figure A.1 – A detailed design diagram of GPVLab	64
Figure C.1 – GPVLab Desktop Icon	67
Figure C.2 – A screenshot of GPVLab main window	68
Figure C.3 – A screenshot of GPVLab explorer window	69
Figure C.4 – A screenshot of “Evaluate Data Model” dialog	71
Figure C.5 – A screenshot of “Add to Library” dialog	72
Figure C.6 – A screenshot of “Model Library” window	73
Figure C.7 – GPVLab integrated help window	73
Figure C.8 – Generated executable files and related screens	74

List of Tables

	Page
Table 7.1: First ten rows of input dataset – Main Evaluation	36
Table 7.2: Results of Main Evaluation	40
Table 7.3: First ten rows of input dataset – Addition of numbers	42
Table 7.4: First ten rows of input dataset – Square root of numbers	43
Table 7.5: Results with default settings – Square root of numbers	44
Table 7.6: First ten rows of input dataset – Ten columns dataset	46
Table 7.7: First ten rows of input dataset – Noise reduction	48
Table 7.8: Experiment results without noise reduction	49
Table 7.9: Evaluation results of the resultant expression	51
Table 7.10: First ten rows for noise reduction – Ten columns dataset	53
Table 7.11: Experiment results without noise reduction- Ten column dataset	53
Table 7.12: Results obtained through GPVLab users	56
Table B.1: Example input dataset	65
Table B.2: Example initial population	66
Table B.3: Example new population	66
Table C.1: Sample dataset	67
Table D.1: Main Dataset for Evaluation (1960 – 2010)	75
Table E.1: Life expectancy vs. exchange rate (1960 – 2010)	78