GenClass: A parallel tool for data classification based on Grammatical Evolution

Nikolaos Anastasopoulos⁽¹⁾, Ioannis G. Tsoulos⁽²⁾, Alexandros Tzallas⁽²⁾

(1) Department of Electrical and Computer Engineering, University of Patras, Greece
(2) Department of Informatics and Telecommunications, University of Ioannina, 47100
Arta, Greece

Abstract

A genetic programming tool is proposed here for data classification. The tool is based on Grammatical Evolution technique and it is designed to exploit multicore computing systems using the OpenMp library. The tool constructs classification programs in a C – like programming language in order to classify the input data, producing simple if – else rules and upon termination the tool produces the classification rules in C and Python files. Also, the user can use his own BNF grammar through a command line option. The tool is tested on a wide range of classification problems and the produced results are compared against traditional techniques for data classification, yielding very promising results.

Keywords: Genetic algorithm, Data classification, Grammatical evolution, Stochastic methods

1. Introduction

- Data classification finds many applications on a series of practical prob-
- lems from areas such as chemistry [1, 2, 3], biology [4, 5], economics [6, 7],
- 4 physics [8, 9] etc. During the past years, many methods have been pro-
- posed to problems of this category such as neural networks [10], radial basis
- 6 functions networks [11], support vector machines [12], etc. The proposed
- 7 method, which is initially described in [13], constructs classification programs
- s in a human readable format using the technique of Grammatical Evolution
- [14]. Grammatical evolution is an evolutionary process that has been applied
- [14]. Grammatical evolution is an evolutionary process that has been applied
- with success in many areas such as music composition [15], economics [16],
- symbolic regression [17], robot control [18] and caching algorithms [23].
- This article introduces a software tool (named GenClass) coded in ANSI C++, for data classification using classification rules in human readable

form. The rules are created with the Grammatical evolution method and can be applied in any classification problem, without a priory knowledge of the dataset. Also, this tool has a series of command line options to control the various aspects of the software, such as mutation rate or the number of chromosomes in the genetic population. Another similar software 18 is the Parallel Program Induction (PPI) available from the relevant URL 19 https://github.com/daaugusto/ppi, which constructs with Grammatical 20 Evolution programs for regression and classification, but the PPI software 21 has many dependencies on other software to be compiled. Also, GenClass 22 produces classification problems in human readable format and it can export 23 output in C++ and in Python. However, in the current version GenClass does not solve regression problems. Also, recently the software NNC has been 25 published [24], which incorporates the Grammatical Evolution to construct 26 neural networks for regression and classification problems. Even though the 27 GenClass software can not be used for regression problems in the current form 28 it has similar execution runs as the NNC software and also it can discover 29 the hidden relationships between the features of a dataset and the expected 30 output. Also, the software GenClass can use every core of a multicore environment and it has the ability to load a BNF grammar from a separate input 32 file. 33

The rest of this article is organized as follows: in section 2 the software is described in detail, in section 3 a series of experiments on some well - known classification datasets are demonstrated, in section 4 the impact of the software tool is discussed and finally in section 5 conclusions and guidelines for further expansion of the software are presented.

2. Software description

34

35

37

42

43

44

45

46

47

48

49

50

40 2.1. The proposed algorithm

The main steps of the algorithm are:

1. **Initialization** step.

- (a) Read the train data
- (b) Set N_G the maximum number of generations, N_C as the number of chromosomes, P_S the selection rate, P_M the mutation rate.
- (c) Initialize the chromosomes of the population.

2. Genetic step

- (a) For $i = 1, ..., N_g$ do
 - i. Create for every chromosome in the population a classification program using Grammatical Evolution.

- ii. Calculate the fitness for every chromosome of the population
- iii. Execute the genetic operators of selection with rate P_S and mutation with rate P

(b) EndFor

51

52

53

54

55

56

57

58 59

66

67

68

70

71

72

73

75

76

77

82

83

84

85

3. Evaluation step

- (a) Create a classification program for the best chromosome in the population.
- (b) Apply the previous program to test set and report the induced error
- The installation of the software is explained in detail in the main page https: //github.com/itsoulos/GenClass/.

2.2. The executable genclass

The outcome of the software compilation and installation is the executable genclass under the directory bin. The executable has the following command line parameters:

- 1. -h:The program prints a help screen.
 - 2. -c count: The parameter count determines the number of chromosomes with default value 500.
- 3. -f **count**. Specify fold count for fold validation. Default value 0 (no folding).
- 4. -i **grammarfile**. The string parameter **grammarfile** stands for the user defined grammar in BNF format. If this options is empty, then the default grammar will be used. An example of such a file is listed in Figure 5 and it is included under examples subdirectory of the distribution.
- 5. -g **gens**: The parameter **gens** determines the maximum number of generations with default value is 200.
- 6. -d **count**. The parameter d determines the maximum number of threads used by the OpenMp library. The default value is 16.
- 7. -s srate: The parameter srate specifies the selection rate with default value 0.10 (10%).
 - 8. -m mrate: The parameter mrate specifies the mutation rate with default value 0.05 (5%).
 - 9. -l **size**: The parameter **size** determines the size of every chromosome with default value 100.
- 10. -p train_file: The string parameter train_file specifies the file containing the points that will be used as train data for the algorithm.
 The file should conform to the format outlined in figure 1.

- 11. -t **test_file**: The string parameter **test_file** specifies the file containing the test data for the particular problem. The file should be in the same format as the **train_file**.
- 12. -o **method**: The string parameter method specifies the output method for the executable. The available options are
 - (a) simple. The program prints output only on termination.
 - (b) csv. In every generation the program prints: number of generations, train error and test error.
 - (c) full. The program prints in every generation detailed information about the optimization procedure.
- 13. -r seed: The integer parameter seed specifies the seed for the random number generator. It can assume any integer value.

101 2.3. The output files

95

96

97

98

The software produces upon termination two distinct files that contains the classification rules. The first file is written in ANSI C++ named classifier.c and an example is shown in Figure 2. The second file is written in Python named classifier.py and an example is displayed in Figure 3.

106 3. Experiments

3.1. A typical example

Consider the Ionosphere dataset available from the Machine Learning Repository. The ionosphere dataset contains data from the Johns Hopkins Ionosphere database. The dataset has been divided into two files, ionosphere.train and ionosphere.test under directory examples of the distribution. A typical run for the GenClass will be

113 ... / bin / genclass —p ionosphere.train —t ionosphere.test —g 10 —o csv
114 The output of this command is shown in Figure 4.

115 3.2. Experiments

121

In order to measure the efficiency of the proposed method a series of experiments were conducted on some common classification problems found in two major dataset databases:

- 1. UCI dataset repository, https://archive.ics.uci.edu/ml/index. php
 - 2. Keel repository, https://sci2s.ugr.es/keel/datasets.php[19].

All the experiments were conducted 30 times using different seed for the random generator each time and averages were taken. Also, to validate 123 the results 10 - fold cross validation was used in all experiments. In all experiments we have used the parameters shown in table 1. The following 125 datasets were used 126

- 1. Wine dataset. The wine recognition dataset contains data from wine chemical analysis.
 - 2. Glass dataset. The dataset contains glass component analysis for glass pieces that belong to 6 classes.
 - 3. Tae dataset. The data consist of evaluations of teaching performance for the University of Wisconsin-Madison.
- 4. Spiral dataset: The spiral artificial dataset contains 1000 two-dimensional examples that belong to two classes (500 examples each). The number of the features is 2. The data in the first class are created using the following formula: $x_1 = 0.5t \cos(0.08t), x_2 = 0.5t \cos(0.08t + \frac{\pi}{2})$ and the second class data using: $x_1 = 0.5t\cos(0.08t + \pi)$, $x_2 =$ $0.5t\cos\left(0.08t + \frac{3\pi}{2}\right)$
- 5. **Pima** dataset. The Pima Indians Diabetes dataset contains with two 139 categories: healthy and diabetic. 140
 - 6. **Ionosphere** dataset. The ionosphere dataset (ION in the following tables) contains data from the Johns Hopkins Ionosphere database.
 - 7. Appendictis dataset, proposed in [20].

127

128

129

130

131

132

133

134

135

136

137

138

141

142

143

144

147

148

149

156

157

158

160

- 8. Australian, the dataset concerns credit card applications.
- 9. Hayes roth dataset. This dataset [21] contains 5 numeric-valued at-145 tributes and 132 patterns. 146
 - 10. Alcohol, a dataset about Alcohol consumption [22].
 - 11. **Dermatology**. Dataset used for differential diagnosis of erythematosquamous diseases.
- 12. Balance. This data set was generated to model psychological experi-150 mental results. 151
- 13. **Regions2** dataset. It is created from liver biopsy images of patients 152 with hepatitis C [26]. 153
- 14. Parkinsons. This dataset is composed of a range of biomedical voice 154 measurements from 31 people, 23 with Parkinson's disease (PD)[25]. 155
 - 15. Wdbc. The Wisconsin diagnostic breast cancer dataset (WDBC) contains data for breast tumors.
- 16. **Popfailures**. This dataset contains records of simulation crashes encountered during climate model. 159
 - 17. **Heart**. The task is to detect the absence or presence of heart disease.

18. **Ecoli**. The goal here is to predict the localization site of proteins.

161

162

163

164

165

166

167

170

171

181

182

183

184

185

186

187

- 19. **Haberman**. A dataset about breast cancer from a study at the University of Chicago's Billings Hospital.
 - 20. **HouseVotes**. This data set includes votes for each of the U.S. House of Representatives Congressmen on the 16 key votes.
- 21. **Shuttle**. The task is to decide what type of control of a spacecraft should be employed.
- 168 22. **Lymography**. The aim here is to detect the presence of a lymphoma in patients.
 - 23. **Mammographic**. This dataset be used to identify the severity (benign or malignant) of a mammographic mass lesion.
- 24. OptDigits. Optical Recognition of Handwritten Digits data set.
- 173 25. **Page Blocks**. The dataset contains blocks of the page layout of a document that has been detected by a segmentation process.
- 175 26. **Penbased**. This is a Pen-Based Recognition of Handwritten Digits data set.
- 177 27. **Saheart**. The dataset is about to categorize persons if have a coronary heart disease.
- 28. **Segment**. This database contains patterns from a database of 7 outdoor images (classes).
 - 29. **Thyroid**. The goal in this dataset is to detect if a patient is normal or suffers from hyperthyroidism or hypothyroidism.
 - 30. **Eeg** datasets. As an real word example, consider an EEG dataset described in [27] is used here. The dataset consists of five sets (denoted as Z, O, N, F and S) each containing 100 single-channel EEG segments each having 23.6 sec duration. With different combinations of these sets the produced datasets are Z_F_S, ZO_NF_S, ZONF_S and Z_O_N_F_S.

The results from the experiments are displayed in Table 2. The column Neu-188 ral stands for the application of a Neural network with 10 processing nodes 189 and the usage of a BFGS variant due to Powell [28] as the training method. 190 The column Mlp-Gen stands for the application of a neural network with 10 191 processing nodes and the utilization of a genetic algorithm as the training 192 method. The genetic algorithm used has the same parameters as in the case 193 of GenClass method, in order to make the results comparable. The column 194 RBF stands for the application of an RBF network using 10 centers. Also, 195 the column RBF-Gen represents the application of an RBF network with 196 10 processing units that is trained using a genetic algorithm with the same 197 set of parameters as in GenClass. Also a statistical comparison is displayed 198 in Figure 6. It is evident, that the proposed method outperforms the other 199

methods in terms of efficiency in the majority of the objective problems. Of course, the method could be slower than RBF or Neural network due to the usage of Grammatical Evolution, however this increase in time can be significantly reduced with the use of threads. In addition, the software appears to be slightly better in some datasets than originally published [13] due to the use of far more chromosomes in the genetic population.

206 4. Impact

This software tool produces human readable classification rules in ANSI 207 C++ as well as in Python format. These rules are used to classify any classifi-208 cation problem without a priory knowledge about the nature of the problem. 209 The implemented software utilizes the Grammatical Evolution technique to 210 produce the classification rules based on simple BNF grammar. The soft-211 ware is guided by a series of command line options that control many critical 212 options of the underlying method such as the number of chromosomes, the 213 mutation rate etc. Also, the software is designed to be used in multi core com-214 putational environments through the publicly available library of OpenMP. 215 The proposed method proved to be very reliable on a series of simple and 216 hard classification problems and hence the software can be used in different 217 areas such as physics, chemistry, medicine etc. The user needs to provide 218 only the patterns of the classification problem in a simple text format and if 219 it is required a series of command line options. Also, the user can provide 220 his own grammar in BNF format in a separate text file. Finally, the software 221 can be easily extended to be more friendly for all users, through a graphical 222 interface. 223

5. Conclusions

228

230

A software which implements a method for data classifications was introduced. The method is based on grammatical evolution and the software is designed to be portable. Future versions of the software will include

- Input from various formats such as CSV, Json etc.
- Usage of improvement stopping rules for the genetic algorithm.
 - Additional output formats.

231 References

- [1] C. Güler, G. D. Thyne, J. E. McCray, K.A. Turner, Evaluation of graphical and multivariate statistical methods for classification of water chemistry data, Hydrogeology Journal 10, pp. 455-474, 2002
- [2] E. Byvatov ,U. Fechner ,J. Sadowski , G. Schneider, Comparison of Support Vector Machine and Artificial Neural Network Systems for Drug/Nondrug Classification, J. Chem. Inf. Comput. Sci. 43, pp 1882–1889, 2003.
- [3] Kunwar P. Singh, Ankita Basant, Amrita Malik, Gunja Jain, Artificial neural network modeling of the river water quality—A case study, Ecological Modelling **220**, pp. 888-895, 2009.
- [4] I. Guyon, J. Weston, S. Barnhill, V. Vapnik, Gene Selection for Cancer Classification using Support Vector Machines, Machine Learning 46, pp. 389-422, 2002.
- [5] R. Marabini, J.M. Carazo, Pattern recognition and classification of images of biological macromolecules using artificial neural networks, Biophysical Journal 66, pp. 1804-1814, 1994.
- [6] I. Kaastra, M. Boyd, Designing a neural network for forecasting financial and economic time series, Neurocomputing **10**, pp. 215-236, 1996.
- Moshe Leshno, Yishay Spector, Neural network prediction analysis: The bankruptcy case, Neurocomputing **10**, pp. 125-147, 1996.
- ²⁵² [8] S. R. Folkes, O. Lahav, S. J. Maddox, An artificial neural network approach to the classification of galaxy spectra, Montly Notices of the Royal Astronomical Society **283**, pp 651-665, 1996.
- [9] C.Z. Cai, W.L. Wang, Y.Z. Chen, Support vector machine classification of physical and biological datasets, Int. J. Mod. Phys. C 14, pp. 575, 2003
- ²⁵⁸ [10] K. Hornik, Multilayer feedforward networks are universal approximators, Neural Networks **2**, pp. 359-366, 1989
- [11] M.D. Buhmann, Radial basis functions, Cambridge monographs on Applied and Computational Mathemetics, 2003.
- ²⁶² [12] I. Steinwart, A. Christmann, Support Vector Machines, Information Science and Statistics, Springer, 2008.

- [13] Ioannis G. Tsoulos, Creating classification rules using grammatical evolution, International Journal of Computational Intelligence Studies 9, pp. 161-171, 2020.
- ²⁶⁷ [14] M. O'Neill, C. Ryan, Grammatical evolution, IEEE Trans. Evol. Comput. **5**, pp. 349–358, 2001.
- [15] Alfonso Ortega, Rafael Sánchez, Manuel Alfonseca Moreno, Automatic
 composition of music by means of grammatical evolution, APL '02 Proceedings of the 2002 conference on APL: array processing languages:
 lore, problems, and applications Pages 148 155.
- [16] Michael O'Neill, Anthony Brabazon, Conor Ryan, J. J. Collins, Evolving Market Index Trading Rules Using Grammatical Evolution, Applications of Evolutionary Computing Volume 2037 of the series Lecture
 Notes in Computer Science pp 343-352.
- [17] M. O'Neill, C. Ryan, Grammatical Evolution: Evolutionary Automatic
 Programming in a Arbitary Language, Genetic Programming, vol. 4,
 Kluwer Academic Publishers, Dordrecht, 2003
- [18] J.J. Collins, C. Ryan, in: Proceedings of AROB 2000, Fifth International
 Symposium on Artificial Life and Robotics, 2000.
- [19] J. Alcalá-Fdez, A. Fernandez, J. Luengo, J. Derrac, S. García,
 L. Sánchez, F. Herrera. KEEL Data-Mining Software Tool: Data
 Set Repository, Integration of Algorithms and Experimental Analysis
 Framework. Journal of Multiple-Valued Logic and Soft Computing 17,
 pp. 255-287, 2011.
- ²⁸⁷ [20] Weiss, Sholom M. and Kulikowski, Casimir A., Computer Systems That
 ²⁸⁸ Learn: Classification and Prediction Methods from Statistics, Neural
 ²⁸⁹ Nets, Machine Learning, and Expert Systems, Morgan Kaufmann Pub²⁹⁰ lishers Inc, 1991.
- [21] B. Hayes-Roth, B., F. Hayes-Roth. Concept learning and the recognition
 and classification of exemplars. Journal of Verbal Learning and Verbal
 Behavior 16, pp. 321-338, 1977.
- [22] Tzimourta, Katerina D. and Tsoulos, Ioannis and Bilero, Thanasis and
 Tzallas, Alexandros T. and Tsipouras, Markos G. and Giannakeas, Nikolaos, Direct Assessment of Alcohol Consumption in Mental State Using
 Brain Computer Interfaces and Grammatical Evolution, Inventions 3,
 pp. 1-12, 2018.

Figure 1: Data format. The integer value D determines the dimensionality of the problem and the value M determines the number of points in the file. Every subsequent line contains a pattern and the final column is the real output (category) for this pattern.

- [23] M. O'Neill, C. Ryan, in: K. Miettinen, M.M. Mkel, P. Neittaanmki, J.
 Periaux (Eds.), Evolutionary Algorithms in Engineering and Computer
 Science, Jyvskyl, Finland, 1999, pp. 127–134.
- Joannis G. Tsoulos, Alexandros Tzallas, Dimitris Tsalikakis, NNC: A tool based on Grammatical Evolution for data classification and differential equation solving, Software X 10, 2019.
- Little MA, McSharry PE, Hunter EJ, Spielman J, Ramig LO.
 Suitability of dysphonia measurements for telemonitoring of
 Parkinson's disease. IEEE Trans Biomed Eng. 2009;56(4):1015.
 doi:10.1109/TBME.2008.2005954
- [26] Giannakeas, N., Tsipouras, M.G., Tzallas, A.T., Kyriakidi, K., Tsianou, Z.E., Manousou, P., Hall, A., Karvounis, E.C., Tsianos, V., Tsianos, E. A clustering based method for collagen proportional area extraction in liver biopsy images (2015) Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBS, 2015-November, art. no. 7319047, pp. 3097-3100.
- R.G. Andrzejak, K. Lehnertz, F. Mormann, C. Rieke, P. David, and C. E. Elger, Indications of nonlinear deterministic and finite-dimensional structures in time series of brain electrical activity: Dependence on recording region and brain state, Phys. Rev. E **64**, pp. 1-8, 2001.
- [28] M.J.D Powell, A Tolerant Algorithm for Linearly Constrained Optimization Calculations, Mathematical Programming 45, pp. 547-566, 1989.

321 Required Metadata

322 Current code version

Figure 2: A typical produced C file of the tool. The double array x is the input and the value class the estimated class of the method.

```
#include <math.h>
int classifier (double *x)
{
    int CLASS = 0;
    if (!(x[6] < x[25] || x[4] < (((x[5] - ((-5.79)/x[0])) / (((-503.562)/x[22])/(x[20]/x[24])))/(x[0] + (x[23]/x[4])))
    ||!(x[23] < = x[2]))) CLASS=0;
    else CLASS=1;
    return CLASS;
}
```

Figure 3: A typical output file in Python language.

```
import ctypes
import numpy as np from typing
import List
def classifier(input: List):
    fun = ctypes.CDLL("./classifier.so")
    fun.classifier.argtypes = [ctypes.POINTER(ctypes.c_double)]
    fun.classifier.restype = ctypes.c_int
    a = np.array(input)
    input_ptr = a.ctypes.data_as(ctypes.POINTER(ctypes.c_double))
    return fun.classifier(input_ptr)
```

Table 1: Parameters for the experiments.

PARAMETER	VALUE
N_C	500
N_G	500
P_S	90%
P_M	5%

Table 2: Results DATASET NEURAL MLP-GEN **RBF** RBF GEN GENCLASS Alcohol 26.85%25.45%46.63%20.29% 14.68%Appendicitis 22.50% $\overline{21.30\%}$ 12.23%11.50%15.00% $\overline{34.89}\%$ Australian 30.48%29.83% 33.30% 14.57%Balance 8.19%8.02%33.42%12.96%0.00%Dermatology 14.33%8.75%62.34%34.94%3.72% $\overline{59.59}\%$ 53.60%52.86%56.64%24.54%Ecoli 54.24%49.30% 33.81% Glass 50.16%45.54%Haberman 29.94%26.98%25.10%24.06%27.33% $\overline{64.36}\%$ Hayes Roth 35.16%33.28%33.54%25.39% Heart 25.95%24.20%31.20%21.00%17.55%HouseVotes 7.54%7.28%6.13%5.29%3.72%15.83%11.30%16.22%Ionosphere 8.63%8.29%Liverdisorder 33.82%31.37%30.84%28.33%32.06%24.55%25.31%20.93%19.29%Lymography 25.57%Mammographic 27.08%16.74%21.38%17.16%17.35% $\overline{81.22}\%$ 50.37%45.96%70.93% 24.33% **OptDigits** Page Blocks 7.10%7.06%10.09% 9.28%3.82%Parkinsons 19.44%14.26%17.42%15.47%9.47%Penbased 41.79% 42.91%82.47%73.16%25.32%Pima 29.07%28.76% 25.78% 24.12% 25.59% **Popfailures** 6.57%6.38%7.04%4.95%7.04%Regions2 33.32%28.56%38.29%33.40%19.52%Saheart 33.22%31.40% 32.19% 29.86%31.30% 24.54% $\overline{21.01\%}$ 10.52%Segment 59.68%48.70%Shuttle 31.29%15.70%32.97%21.73%0.15%Spiral 42.60%41.57%44.87% 41.27%36.90% 55.07%Tae 47.53%47.22%60.07%37.33%Thyroid 4.28%4.16%10.52%7.42%1.89%Wdbc 6.76%7.27%20.55%6.22%4.11%Wine 46.63%12.43% 31.41%12.84%4.71% Z_F_S 14.17%10.68%13.16%9.65%7.00%Z_O_N_F_S 78.73%59.93%48.71%45.20%32.20%ZO_NF_S 2.60%9.97%7.33%9.02%6.81%

3.02%

4.03%

3.98%

1.60%

ZONF_S

3.48%

Figure 4: Typical output of the GenClass command. In every line the tool prints the generation number, the train error as well as the test error. At the end the tool prints the final classification rule and the estimated train and test error.

```
1,
        15.43,
                    19.32
  2,
          15.43,
                      19.32
  3,
          15.43,
                      19.32
  4,
          13.71,
                      17.05
          12.57,
                      15.34
  5,
          12.57,
  6,
                      15.34
  7,
          12.57,
                      15.34
  8,
             12,
                      13.64
  9,
             12,
                      13.64
  FINAL OUTPUT
  EXPRESSION=
  if (!(x7 < log(cos(((-788.787) + ((sin(x28)/sin(cos(((-7.17)/x34)))))
      +(-83.6))))))x6>x13&x7<log(x5))
  CLASS=0.00 else CLASS=1.00
TRAIN ERROR = 12.00\%
 CLASS ERROR = 13.64\%
```

Nr.	Code metadata description	
C1	Current code version	1.0
C2	Permanent link to code/repository	https://github.com/itsoulos/
	used for this code version	GenClass/
С3	Legal Code License	GNU General Public License (GPL)
C4	Code versioning system used	git
C5	Software code languages, tools, and	C++
	services used	
C6	Compilation requirements, operat-	Linux
	ing environments & dependencies	
C7	If available Link to developer docu-	https://github.com/itsoulos/
	mentation/manual	GenClass/wiki
C8	Support email for questions	itsoulos@uoi.gr

Figure 5: An example of user defined grammar.

```
<S>::=<BEXPR>
<BEXPR>::= <XLIST><BOOLOP><EXPR>
         | !(<BEXPR>)
          | <XLIST>BOOLOP>EXPR>&<BEXPR>
          | <XLIST><BOOLOP><EXPR>|<BEXPR>
<BOOLOP>::= >
           | >=
           | <
           | <=
<EXPR>::= (<EXPR><BINARYOP><EXPR>)
        | <FUNCTION>(<EXPR>)
        | <TERMINAL>
<BINARYOP>::=+
<FUNCTION>::= sin
              \cos
              exp
              log
<TERMINAL>::= <XLIST>
             | <DIGITLIST>.<DIGITLIST>
             (-<DIGITLIST>.<DIGITLIST>)
<DIGITLIST>::= <DIGIT>
              | <DIGIT><DIGIT>
              | <DIGIT><DIGIT>
<DIGIT>::= 0
           3
           4
           5
           6
           7
           9
\langle XLIST \rangle ::= x1
           x2
           x3
           x4
           x5
                            14
           x6
           x7
           x8
```

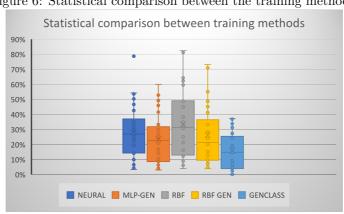


Figure 6: Statistical comparison between the training methods.