

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 its 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. 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 problems from areas such as chemistry [1, 2, 3], biology [4, 5], economics [6, 7], physics [8, 9] etc. During the past years many methods have been proposed to problems of this category such as neural networks [10], radial basis functions networks [11], support vector machines [12], etc. The proposed method, which is initially described in [13], constructs classification programs in a human readable format using the technique of Grammatical Evolution [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.

16 Also, this tool has a series of command line options to control the various
17 aspects of the software, such as mutation rate or number of chromosomes in
18 the genetic population.

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

24 **2. Software description**

25 *2.1. The proposed algorithm*

26 The main steps of the algorithm are:

27 **1. Initialization** step.

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

32 **2. Genetic** step

- 33 (a) **For** $i = 1, \dots, N_g$ **do**
 - 34 i. Create for every chromosome in the population a classification
35 program using Grammatical Evolution.
 - 36 ii. Calculate the fitness for every chromosome of the population
 - 37 iii. Execute the genetic operators of selection with rate P_S and
38 mutation with rate P
- 39 (b) **EndFor**

40 **3. Evaluation** step

- 41 (a) Create a classification program for the best chromosome in the
42 population.
- 43 (b) Apply the previous program to test set and report the induced
44 error.

45 The installation of the software is explained in detail in the relevant Wiki
46 page <https://github.com/itsoulos/GenClass/wiki>.

47 2.2. The executable *genclass*

48 The outcome of the software compilation and installation is the executable
49 **genclass** under the directory **bin**. The executable has the following com-
50 mand line parameters:

- 51 1. **-h**: The program prints a help screen.
- 52 2. **-c count**: The parameter **count** determines the number of chromo-
53 somes with default value 500.
- 54 3. **-f count**. Specify fold count for fold validation. Default value 0 (no
55 folding).
- 56 4. **-g gens**: The parameter **gens** determines the maximum number of
57 generations with default value is 200.
- 58 5. **-d count**. The parameter **d** determines the maximum number of threads
59 used by the OpenMp library. The default value is 16.
- 60 6. **-s srate**: The parameter **srate** specifies the selection rate with default
61 value 0.10 (10%).
- 62 7. **-m mrate**: The parameter **mrate** specifies the mutation rate with de-
63 fault value 0.05 (5%).
- 64 8. **-l size**: The parameter **size** determines the size of every chromosome
65 with default value 100.
- 66 9. **-p train_file**: The string parameter **train_file** specifies the file con-
67 taining the points that will be used as train data for the algorithm.
68 The file should conform to the format outlined in figure 1.
- 69 10. **-t test_file**: The string parameter **test_file** specifies the file containing
70 the test data for the particular problem. The file should be in the same
71 format as the **train_file**.
- 72 11. **-o method**: The string parameter **method** specifies the output method
73 for the executable. The available options are
 - 74 (a) **simple**. The program prints output only on termination.
 - 75 (b) **csv**. In every generation the program prints: number of genera-
76 tions, train error and test error.
 - 77 (c) **full**. The program prints in every generation detailed information
78 about the optimization procedure.
- 79 12. **-r seed**: The integer parameter **seed** specifies the seed for the random
80 number generator. It can assume any integer value.

81 2.3. The output files

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

86 3. Experiments

87 3.1. A typical example

88 Consider the Ionosphere dataset available from the Machine Learning
89 Repository. The ionosphere dataset contains data from the Johns Hop-
90 kins Ionosphere database. The dataset has been divided into two files,
91 `ionosphere.train` and `ionosphere.test` under directory `examples` of the
92 distribution. A typical run for the `GenClass` will be

```
93 ../bin/genclass -p ionosphere.train -t ionosphere.test -g 10 -o csv
```

94 The output of this command is shown in Figure 4.

95 3.2. Experiments

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

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

102 All the experiments were conducted 30 times using different seed for the
103 random generator each time and averages were taken. In all experiments we
104 have used the parameters shown in table 1. The following datasets were used

- 105 1. **Wine** dataset. The wine recognition dataset contains data from wine
106 chemical analysis.
- 107 2. **Glass** dataset. The dataset contains glass component analysis for glass
108 pieces that belong to 6 classes.
- 109 3. **Tae** dataset. The data consist of evaluations of teaching performance
110 for the University of Wisconsin-Madison.
- 111 4. **Spiral** dataset: The spiral artificial dataset contains 1000 two-dimensional
112 examples that belong to two classes (500 examples each). The num-
113 ber of the features is 2. The data in the first class are created using
114 the following formula: $x_1 = 0.5t \cos(0.08t)$, $x_2 = 0.5t \cos(0.08t + \frac{\pi}{2})$
115 and the second class data using: $x_1 = 0.5t \cos(0.08t + \pi)$, $x_2 =$
116 $0.5t \cos(0.08t + \frac{3\pi}{2})$
- 117 5. **Pima** dataset. The Pima Indians Diabetes dataset contains with two
118 categories: healthy and diabetic.
- 119 6. **Ionosphere** dataset. The ionosphere dataset (ION in the following
120 tables) contains data from the Johns Hopkins Ionosphere database.
- 121 7. **Appendictis** dataset, proposed in [20].

- 122 8. **Australian**, the dataset concerns credit card applications.
- 123 9. **Hayes roth** dataset. This dataset[21] contains 5 numeric-valued at-
124 tributes and 132 patterns.
- 125 10. **Alcohol**, a dataset about Alcohol consumption [22].
- 126 11. **Dermatology**. Dataset used for differential diagnosis of erythemato-
127 squamous diseases.
- 128 12. **Balance**. This data set was generated to model psychological experi-
129 mental results.
- 130 13. **Regions2** dataset. It is created from liver biopsy images of patients
131 with hepatitis C [25].
- 132 14. **Parkinsons**. This dataset is composed of a range of biomedical voice
133 measurements from 31 people, 23 with Parkinson’s disease (PD)[24].
- 134 15. **Wdbc**. The Wisconsin diagnostic breast cancer dataset (WDBC) con-
135 tains data for breast tumors.
- 136 16. **Popfailures**. This dataset contains records of simulation crashes en-
137 countered during climate model.
- 138 17. **Heart**. The task is to detect the absence or presence of heart disease.
- 139 18. **Ecoli**. The goal here is to predict the localization site of proteins.
- 140 19. **Haberman**. A dataset about breast cancer from a study at the Uni-
141 versity of Chicago’s Billings Hospital.
- 142 20. **HouseVotes**. This data set includes votes for each of the U.S. House
143 of Representatives Congressmen on the 16 key votes.
- 144 21. **Shuttle**. The task is to decide what type of control of a spacecraft
145 should be employed.
- 146 22. **Lymography**. The aim here is to detect the presence of a lymphoma
147 in patients.
- 148 23. **Mammographic**. This dataset be used to identify the severity (benign
149 or malignant) of a mammographic mass lesion.
- 150 24. **OptDigits**. Optical Recognition of Handwritten Digits data set.
- 151 25. **Page Blocks**. The dataset contains blocks of the page layout of a
152 document that has been detected by a segmentation process.
- 153 26. **Penbased**. This is a Pen-Based Recognition of Handwritten Digits
154 data set.
- 155 27. **Saheart**. The dataset is about to categorize persons if have a coronary
156 heart disease.
- 157 28. **Segment**. This database contains patterns from a database of 7 out-
158 door images (classes).
- 159 29. **Thyroid**. The goal in this dataset is to detect if a patient is normal
160 or suffers from hyperthyroidism or hypothyroidism.

161 30. **Eeg** datasets. As an real word example, consider an EEG dataset
 162 described in [26] is used here. The dataset consists of five sets (denoted
 163 as Z, O, N, F and S) each containing 100 single-channel EEG segments
 164 each having 23.6 sec duration. With different combinations of these sets
 165 the produced datasets are Z_F_S, ZO_NF_S, ZONF_S and Z_O_N_F_S.

166 The results from the experiments are displayed in table 2. It is evident, that
 167 the proposed method outperforms the other methods in terms of efficiency
 168 in the majority of the objective problems. Of course, the method could
 169 be slower than RBF or Neural network due to the usage of Grammatical
 170 Evolution, however this increase in time can be significantly reduced with
 171 the use of threads.

172 4. Impact

173 This software tool produces human readable classification rules in ANSI
 174 C++ as well as in Python format. This rules used to classify any classification
 175 problem without a priory knowledge about the nature of the problem. The
 176 implemented software utilizes the Grammatical Evolution technique to pro-
 177 duce the classification rules based on simple BNF grammar. The software is
 178 guided by a series of command line options that control many critical options
 179 of the underlying method such as the number of chromosomes, the mutation
 180 rate etc. Also, the software is designed to be use in multi core computational
 181 environments through the public available library of OpenMP. From a exten-
 182 sive series of experiments the proposed method proved to be very reliable on
 183 a series of simple and hard classification problems and hence the software can
 184 be used in different areas such as physics, chemistry, medicine etc. The user
 185 need to provide only the patterns of the classification problem in a simple
 186 text format and if it is required a series of command line options. Finally,
 187 the software can be easily extended to be more friendly for all users, through
 188 a graphical interface.

189 5. Conclusions

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

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

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.
- [7] Moshe Leshno, Yishay Spector, Neural network prediction analysis: The bankruptcy case, *Neurocomputing* **10**, pp. 125-147, 1996.
- [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
- [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 Mathematics, 2003.
- [12] I. Steinwart, A. Christmann, Support Vector Machines, Information Science and Statistics, Springer, 2008.

- 229 [13] Ioannis G. Tsoulos, Creating classification rules using grammatical evolution,
230 International Journal of Computational Intelligence Studies **9**,
231 pp. 161-171, 2020.
- 232 [14] M. O'Neill, C. Ryan, Grammatical evolution, IEEE Trans. Evol. Comput.
233 **5**, pp. 349–358, 2001.
- 234 [15] Alfonso Ortega, Rafael Sánchez, Manuel Alfonseca Moreno, Automatic
235 composition of music by means of grammatical evolution, APL '02 Proceedings
236 of the 2002 conference on APL: array processing languages: lore, problems, and
237 applications Pages 148 - 155.
- 238 [16] Michael O'Neill, Anthony Brabazon, Conor Ryan, J. J. Collins, Evolving
239 Market Index Trading Rules Using Grammatical Evolution, Applications of
240 Evolutionary Computing Volume 2037 of the series Lecture Notes in Computer
241 Science pp 343-352.
- 242 [17] M. O'Neill, C. Ryan, Grammatical Evolution: Evolutionary Automatic
243 Programming in a Arbitrary Language, Genetic Programming, vol. 4,
244 Kluwer Academic Publishers, Dordrecht, 2003
- 245 [18] J.J. Collins, C. Ryan, in: Proceedings of AROB 2000, Fifth International
246 Symposium on Artificial Life and Robotics, 2000.
- 247 [19] J. Alcalá-Fdez, A. Fernandez, J. Luengo, J. Derrac, S. García,
248 L. Sánchez, F. Herrera. KEEL Data-Mining Software Tool: Data
249 Set Repository, Integration of Algorithms and Experimental Analysis
250 Framework. Journal of Multiple-Valued Logic and Soft Computing **17**,
251 pp. 255-287, 2011.
- 252 [20] Weiss, Sholom M. and Kulikowski, Casimir A., Computer Systems That
253 Learn: Classification and Prediction Methods from Statistics, Neural
254 Nets, Machine Learning, and Expert Systems, Morgan Kaufmann Publishers
255 Inc, 1991.
- 256 [21] B. Hayes-Roth, B., F. Hayes-Roth. Concept learning and the recognition
257 and classification of exemplars. Journal of Verbal Learning and Verbal
258 Behavior **16**, pp. 321-338, 1977.
- 259 [22] Tzimourta, Katerina D. and Tsoulos, Ioannis and Bilerio, Thanasis and
260 Tzallas, Alexandros T. and Tsipouras, Markos G. and Giannakeas, Nikolaos,
261 Direct Assessment of Alcohol Consumption in Mental State Using
262 Brain Computer Interfaces and Grammatical Evolution, Inventions **3**,
263 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.

D
 M
 $x_{11} \quad x_{12} \quad \dots \quad x_{1D} \quad y_1$
 $x_{21} \quad x_{22} \quad \dots \quad x_{2D} \quad y_2$
 $\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots$
 $x_{M1} \quad x_{M2} \quad \dots \quad x_{MD} \quad y_M$

Table 1: Parameters for the experiments.

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

- 264 [23] M. O'Neill, C. Ryan, in: K. Miettinen, M.M. Mkel, P. Neittaanmki, J.
265 Periaux (Eds.), *Evolutionary Algorithms in Engineering and Computer*
266 *Science*, Jyväskylä, Finland, 1999, pp. 127–134.
- 267 [24] Little MA, McSharry PE, Hunter EJ, Spielman J, Ramig LO.
268 Suitability of dysphonia measurements for telemonitoring of
269 Parkinson's disease. *IEEE Trans Biomed Eng.* 2009;56(4):1015.
270 doi:10.1109/TBME.2008.2005954
- 271 [25] Giannakeas, N., Tsipouras, M.G., Tzallas, A.T., Kyriakidi, K., Tsianou,
272 Z.E., Manousou, P., Hall, A., Karvounis, E.C., Tsianos, V., Tsianos,
273 E. A clustering based method for collagen proportional area extraction
274 in liver biopsy images (2015) *Proceedings of the Annual International*
275 *Conference of the IEEE Engineering in Medicine and Biology Society,*
276 *EMBS, 2015-November*, art. no. 7319047, pp. 3097-3100.
- 277 [26] R.G. Andrzejak, K. Lehnertz, F. Mormann, C. Rieke, P. David, and C.
278 E. Elger, Indications of nonlinear deterministic and finite-dimensional
279 structures in time series of brain electrical activity: Dependence on
280 recording region and brain state, *Phys. Rev. E* **64**, pp. 1-8, 2001.

281 **Required Metadata**

282 **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 2: Results.

DATASET	NEURAL	RBF	GENCLASS
Alcohol	26.85%	46.63%	14.68%
Appendicitis	22.50%	12.23%	15.00%
Australian	30.48%	34.89%	14.57%%
Balance	8.19%	33.42%	0.00%
Dermatology	14.33%	62.34%	3.72%
Ecoli	53.60%	59.59%	24.54%
Glass	54.24%	50.16%	33.81%
Haberman	29.94%	25.10%	27.33%
Hayes Roth	35.16%	64.36%	25.39%
Heart	25.95%	31.20%	17.55%
HouseVotes	7.54%	6.13%	3.72%
Ionosphere	15.83%	16.22%	8.29%
Liverdisorder	33.82%	30.84%	32.06%
Lymography	25.57%	25.31%	19.29%
Mammographic	27.08%	21.38%	17.35%
OptDigits	50.37%	81.22%	24.33%
Page Blocks	7.10%	10.09%	3.82%
Parkinsons	19.44%	17.42%	9.47%
Penbased	42.91%	82.47%	25.32%
Pima	29.07%	25.78%	25.59%
Popfailures	6.57%	7.04%	7.04%
Regions2	33.32%	38.29%	19.52%
Saheart	33.22%	32.19%	31.30%
Segment	24.54%	59.68%	10.52%
Shuttle	31.29%	32.97%	0.15%
Spiral	42.60%	44.87%	36.90%
Tae	47.53%	60.07%	37.33%
Thyroid	4.28%	10.52%	1.89%
Wdbc	20.55%	7.27%	4.11%
Wine	46.63%	31.41%	4.71%
Z_F_S	14.17%	13.16%	7.00%
Z_O_N_F_S	78.73%	48.71%	32.20%
ZO_NF_S	9.97%	9.02%	2.60%
ZONF_S	3.48%	4.03%	1.60%

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
5,      12.57,      15.34
6,      12.57,      15.34
7,      12.57,      15.34
8,       12,       13.64
9,       12,       13.64
FINAL OUTPUT
EXPRESSION=
if (!(x7<log(cos(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 used for this code version	https://github.com/itsoulos/GenClass/
C3	Legal Code License	GNU General Public License (GPL)
C4	Code versioning system used	git
C5	Software code languages, tools, and services used	C++
C6	Compilation requirements, operating environments & dependencies	Linux
C7	If available Link to developer documentation/manual	https://github.com/itsoulos/GenClass/wiki
C8	Support email for questions	itsoulos@uoi.gr