

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/281775100>

# Evaluating the potential of Genetic Programming as an exploratory data analysis in soil science

Research · September 2015

---

CITATIONS

0

---

READS

51

2 authors:



**Lorenzo Menichetti**

Swedish University of Agricultural Sciences

28 PUBLICATIONS 286 CITATIONS

SEE PROFILE



**Alberto Paolo Tonda**

French National Institute for Agricultural Resea...

87 PUBLICATIONS 259 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Evolutionary Algorithms for Network Stress Tests [View project](#)



Modelling C accumulation in non-equilibrium forest ecosystems characterized by land use changes [View project](#)

All content following this page was uploaded by [Lorenzo Menichetti](#) on 15 September 2015.

The user has requested enhancement of the downloaded file.

# Evaluating the potential of Genetic Programming as an exploratory data analysis in soil science

L. Menichetti<sup>a\*</sup> and A. Tonda<sup>b</sup>

<sup>a</sup>Swedish University of Agricultural Sciences, Department of Soil and Environment, P.O.Box 7014, 75007 Uppsala, Sweden

<sup>b</sup>UMR 782 GMPA, INRA, 1 Av. Lucien Brétignères, 78850, Thiverval-Grignon, France

\*corresponding author, tel. +46768549268, e-mail: Lorenzo.Menichetti@slu.se

## Abstract

Genetic Programming is a powerful optimization technique, able to deliver high-quality results in several real-world problems. One of its most successful applications is symbolic regression, where the objective is to find a suitable expression to model the underlying relationship between data points, with no aprioristic assumptions. In this paper, we propose the application of a Genetic Programming technique to a dataset on soil respiration and soil properties, in order to investigate possible influences of soil properties on soil respiration through symbolic regression. The best candidate models obtained by the technique are then studied to determine possible differences in the relationships related to environmental factors. Recurring patterns in the best solutions proposed by the search algorithm are identified, and the suitability of symbolic regression in soil science is evaluated and discussed. Genetic Programming proves to be an extremely promising data mining technique for soil scientists, as it is able to uncover relationships that could otherwise remain hidden, while remaining completely neutral and bias-free. We suggest its application for routine data analysis, as the technique presents particular interest for environmental modeling and development of pedotransfer functions.

## 1. Introduction

As new field methods are developed, making field measurements cheaper and denser, and new studies are published, the amount of data available to the scientific community grows more than linearly over time. An unprecedented amount of data is now at disposal of ecosystem scientists, and there is a need for methods able to treat it in a comprehensive and objective way. This objective involves the use of new algorithms and data mining procedures, as the field slowly adopts more and more automatic processes.

34 Semi-empirical relationships are widely exploited in soil science. For example, it is common to  
35 predict soil properties, which would be too costly or difficult to estimate otherwise, through the  
36 use of pedotransfer functions (Bouma, 1989) that exploit easily measurable variables. Because  
37 of the high global concern for climate change and the emission of greenhouse gases in the  
38 atmosphere, another variable that received a lot of attention in the last decades is soil  
39 respiration. Its relation with different soil edaphic properties and soil processes is nevertheless  
40 still not completely clear. While it is widely accepted that soil respiration is linked with  
41 temperature and moisture conditions (see Lloyd & Taylor, 1994 and Moyano *et al.*, 2011 for  
42 some examples), there is a lack of understanding on how site-specific properties can modify  
43 these relationships in the field. Part of the observed error in these relationships is probably  
44 accountable to yet unknown links between soil respiration and environment.

45 One possible approach to the issue is to explore the space of possible dependencies between  
46 elements, while being as unbiased as possible towards the shape of the solution. The rise in  
47 complexity of available data has led the machine learning community to develop refined  
48 methods able to uncover relationships between variables in huge datasets. For natural laws,  
49 evolutionary-based computation has been successfully used to detect hidden dependencies,  
50 especially in field of physics (Schmidt & Lipson, 2009). While the expertise of human scientists  
51 is irreplaceable, machine learning can be exploited to obtain a large number of *candidate*  
52 *solutions*, that is, equations proposing a connection between variables.

53 Evolutionary algorithms have been sometimes applied to soil science problems for the  
54 development of pedotransfer functions (Crowe et al, 2006, Padarian et al, 2012) and more often  
55 to hydrology problems (Johari et al., 2011, Pedroso et al, 2011), but the potential of the  
56 technique in soil science is still largely unexplored. In this paper, we propose to apply a state-of-  
57 the-art evolutionary algorithm to a real-world dataset obtained by crossing two freely available  
58 datasets on soil respiration and soil properties. The most promising solutions obtained through  
59 the automatic approach are then examined, and recurring patterns are found, hinting at possible

60 strong, uncovered relationship between variables. While further experiments are required to  
61 draw more definite conclusions, preliminary results show great promises for the coupling of  
62 automatic approaches and human expertise.

63

## 64 2. Material and Methods

### 65 2.1 Evolutionary Algorithms and Symbolic Regression

66 The term *Evolutionary Algorithms* (EAs) groups a great variety of bio-inspired stochastic meta-  
67 heuristics for optimization, loosely inspired by the paradigm of Neo-Darwinian natural evolution.  
68 In EAs, an *individual* is defined as a candidate solution for a given problem. A *population* of  
69 solutions is randomly created, and then evaluated with a *fitness function*, that examines their  
70 efficacy with regards to a target problem. The fittest individuals are then selected for  
71 *reproduction*, usually performed by slightly altering some elements of the solution (*mutation*) or  
72 by mixing the information contained in two individuals (*crossover*). The result of the reproduction  
73 step is a new generation of candidate solutions, which are subsequently evaluated with the  
74 fitness function. The worst individuals are removed from the population, and the loop resumes  
75 from reproduction, until a user-defined *stop condition* is reached.

76 After the seminal work on *Genetic Algorithms* (Holland, 1975) carried on by Holland during the  
77 60s, where solutions are modeled as bit strings, other independent research lines led by Fogel  
78 and Schwefel gave birth to *Evolutionary Programming* (Fogel, 1962) and *Evolution Strategies*  
79 (Schwefel, 1965), powerful algorithms focused on real-value optimization. At the beginning of  
80 the 90s, John Koza presented *Genetic Programming* (GP) (Koza, 1992), an EA whose  
81 individuals are modeled as trees: the expressive power of this idea made it possible to approach  
82 extremely complex problems, where the shape of a solution could range from a network layout  
83 to a complete Assembly-language program.

84 Thanks to the development of GP, the EA community tried to answer to the pressing practical  
85 need for improved forms of scientific data mining (Clery & Voss, 2005 and Valdés-Pérez, 1999)  
86 with the *symbolic regression* technique. In symbolic regression, the objective is to find a  
87 mathematical expression linking variables' values in a dataset, without making assumptions on  
88 the structure of the expression itself. Candidate equations to solve the problem are modeled as  
89 trees, while the fitness function usually aims at minimizing the absolute or squared difference  
90 from experimental data. From the first promising results (Koza, 1992), a research line led by  
91 Schmidt and Lipson produced an extremely efficient GP-based algorithm (Schmidt & Lipson,  
92 2009), able to deliver high-quality solutions in small amounts of time. The derived software,  
93 *Eureqa Formulize* (<http://formulize.nutonian.com/>, accessed on 25 September 2013), is now  
94 considered the state of the art in the field.

95

## 96 **2.2 The dataset**

97 In this study, we use data from the updated soil respiration database (SRDB) (Bond-Lamberty &  
98 Thomson, 2012): in particular, we consider variables describing soil respiration, mean annual  
99 temperature and mean annual precipitations. Latitude and longitude specified in the  
100 corresponding study are used for combining this dataset with the harmonized world soil  
101 database (HWSD) (FAO/IIASA/ISRIC/ISSCAS/JRC, 2012). Topsoil and subsoil gravel, sand, silt  
102 and clay content, topsoil and subsoil pH and cation exchange capacity (CEC), topsoil and  
103 subsoil soil organic carbon (SOC) content are taken from the HWSD, while soil respiration data  
104 and the environment ecologic identification are obtained from the SRDB.

105 The target variable for the study is soil respiration, normalized by the SOC content in the topsoil,  
106 as the latter is already known to explain most of the observed variation.

107 In order to improve the effectiveness of the search algorithm, outliers outside two times the  
108 interquartile range are removed. Data are then normalized, so that each variable has mean 0

109 and variance 1, and then multiplied by 100 in order to obtain a medium magnitude. The dataset  
110 is then randomly divided between a training set (80% of the samples) and a validation set (20%  
111 of the samples).

112

### 113 2.3 Data treatment

114 The target expression is:

$$R_{norm} = f(lat, long, T, P, Gravel_{tops}, Sand_{tops}, Silt_{tops}, Clay_{tops}, BD_{tops}, pH_{tops}, CEC_{tops}, \\ Gravel_{subs}, Sand_{subs}, Silt_{subs}, Clay_{subs}, BD_{subs}, pH_{subs}, SOC_{subs}, CEC_{subs})$$

115 where *tops* denotes topsoil and *subs* denotes subsoil. The term  $R_{norm}$  denotes soil respiration,  
116 normalized by topsoil SOC content ( in g C m<sup>-2</sup> per unit % of SOC content), the term  $T$  the mean  
117 annual air temperature (in ° C). The two terms *lat* and *long* denote latitude and longitude,  
118 respectively. The terms *Gravel*, *Sand*, *Silt*, *Clay* and *BD* denote gravel, sand, silt and clay  
119 percentage and bulk density, respectively. The term *pH* denotes the soil pH measured in H<sub>2</sub>O,  
120 the term *SOC* denotes the soil organic carbon content in percentage and the term *CEC* denotes  
121 the cation exchange capacity in cmol kg<sup>-1</sup>. The following basic functions are used as building  
122 blocks during the GP search: *constant*, *integer constant*, *input variable*, *addition*, *subtraction*,  
123 *multiplication*, *division*, *negation*, *sine*, *cosine*, *tangent*, *exponential*, *natural logarithm*, *factorial*,  
124 *power*, *square root*, *minimum*, *maximum*, *modulo*, *floor* and *ceiling*. After the search, the ten  
125 best solutions proposed by the software are tested against the validation dataset, and residuals  
126 for each point are computed. Residuals are then plotted, divided by ecosystem group.  
127 As a measurement of the fit of the possible models, we consider the following indicators: mean  
128 error (ME), mean absolute error (MAE), root mean squared error (RMSE), normalized root  
129 mean squared error (NRMSE), percent bias (PBIAS), Nash-Sutcliffe Efficiency (NSE), index of  
130 agreement (d), Pearson's correlation coefficient (r) and coefficient of determination (R<sup>2</sup>).

131 Candidate solutions are also visually compared through a principal component analysis (PCA)  
132 (Venables & Ripley, 2002) on the residuals.  
133 The machine used to run the search is a 64-bit workstation with 64 GB of RAM, mounting 2 Intel  
134 Xeon 2-Ghz E5-2650 processors, using a total of 16 cores and 32 threads. The software used  
135 for the experiments shows several statistics to detect convergence: in this case, we observe  
136 *maturity*, a metric that describes diversity inside the population. When an EA is close to  
137 convergence, most of the candidate solutions inside the population closely resemble each other,  
138 with minimal differences between them: in such a condition, the EA is focusing on exploitation of  
139 a small part of a search space, and it is unlikely to produce dramatically different solutions. We  
140 stop the experiment when the maturity score of the population reaches 90%, after about 25  
141 hours of computation. It is important to notice that the same results could have been achieved  
142 on a standard desktop computer in a reasonable amount of time (around one week).

143

## 144 3. Results

### 145 3.1 The selected candidate solutions

146 Our search evaluated  $2.4 \times 10^{12}$  solutions over approximately 8 million of generations. We  
147 selected the 10 best solutions presented by the search algorithm according to the best  
148 compromise between complexity (the size of the function) and squared error minimization. The  
149 selected solutions are the following:

150

$$151 \quad R_{norm} = 1.79 \cdot pH_{tops} + mod(Clay_{tops}, 0.63) - 0.31 \cdot Silt_{subs} - 1.7 \cdot pH_{subs} - 0.21 \cdot T^2 \quad (1)$$

$$152 \quad R_{norm} = 0.24 + 2.17 \cdot pH_{tops} + 1.60 \cdot \min(\min(pH_{tops}^2, 1.56 - pH_{tops}), T) - T - 2.09 \cdot pH_{subs} \quad (2)$$

$$153 \quad R_{norm} = 0.27 + 2.14 \cdot pH_{tops} + 1.62 \cdot \min(\min(pH_{tops}^2, 1.12 + Clay_{subs} - Silt_{tops}), T) - T - 2.04 \cdot$$

154

$$Ph_{subs} \quad (3)$$

155  $R_{norm} = 1.65 \cdot pH_{tops} + 0.25 \cdot Clay_{tops} + 0.16 \cdot BD_{tops} - 1.68 \cdot pH_{subs} - 0.22 \cdot T^2(4)$

156  $R_{norm} = 1.70 \cdot pH_{tops} + 0.19 \cdot BD_{tops} + 2.05 \cdot Clay_{tops} \cdot \max(0.12, Gravel_{subs}) - 1.73 \cdot pH_{subs} -$   
 157  $0.22 \cdot T^2(5)$

158  $R_{norm} = 0.41 + 1.73 \cdot pH_{tops} + 0.26 \cdot Clay_{tops} + 0.19 \cdot BD_{tops} - 1.76 \cdot pH_{subs} - 0.24 \cdot T^2(6)$

159  $R_{norm} = pH_{tops} + 0.32 \cdot Clay_{subs} + 0.21 \cdot BD_{tops} - pH_{subs} - 0.22 \cdot T^2(7)$

160  $R_{norm} = \text{mod}(BD_{tops} + 1.61 \cdot Clay_{tops}, 0.83) - 0.15 \cdot T^2(8)$

161  $R_{norm} =$

162  $0.27 + 2.14 \cdot pH_{tops} + 1.61 \cdot \min(\min(pH_{tops}^2, 1.48 + 1.29 \cdot Gravel_{subs} - pH_{subs} - pH_{tops} \cdot$   
 163  $Gravel_{subs}), T) - T - 2.03 \cdot pH_{subs}(9)$

164  $R_{norm} = 1.64 \cdot pH_{tops} + \text{mod}(Clay_{subs}, 0.53) - 1.61 \cdot pH_{subs} - 0.22 \cdot T^2(10)$

165 Considering the following parameters:  $R_{norm}$ ,  $T$ ,  $BD_{tops}$ ,  $Silt_{tops}$ ,  $Clay_{subs}$  and  $Silt_{subs}$ .

166

### 167 3.2 The fit of solutions by ecosystem groups

168 The residuals of the selected solutions look similar when the selected solutions are tested  
 169 against the validation dataset (Fig. 1). In savanna ecosystems, the variation of residuals is quite  
 170 high, and the obtained functions do not seem to perform well.

171 The only functions that seem to reproduce data in the validation dataset are Eq. 2 and Eq. 3.

172 The PCA analysis of the residuals (Figure 2) does not find relevant differences by ecosystem  
 173 group, but helps to highlight the differences between Eq. 2, Eq. 3 and all the others.

174 The variables selected by the search do not include latitude, cation exchange capacity or mean  
 175 annual precipitation, and all the variability is explained according to mean annual temperature,  
 176 pH and soil texture. The two most performing functions, Eq. 2 and Eq. 3, do not include  
 177 exponential terms for the mean annual temperature, and both are almost linear, differently from  
 178 the others.



## 179 4. Discussion

### 180 4.1 The candidate solutions

181 All the selected equations present an  $R^2$  value on the training dataset between 0.46 and 0.50,  
182 but do not perform accordingly on the validation dataset. Eq. 2 and Eq. 3 are the only two  
183 solutions that can be considered to explain some of the variability in the validation dataset. Both  
184 functions suggest a linear relationship between soil respiration, topsoil pH and temperature,  
185 while introducing also a small nonlinear factor for topsoil pH. The better fit of Eq. 3 seems to be  
186 related to the inclusion of soil texture in the function.

187 The bad fit for most of the functions on the validation dataset, together with the relatively good fit  
188 on the training dataset, can be explained considering the specificity of the constant terms  
189 proposed by the algorithm. Furthermore, in the machine learning community, there is evidence  
190 that GP models with a high degree of complexity might overfit the training set, introducing terms  
191 that increase the fitting by a minimal amount, exploiting specific characteristics of the dataset  
192 that do not generalize well (Rosca, 1996). The information on the possible relationships  
193 between the data that all the selected functions carry is nevertheless potentially valuable, as  
194 many of the relationships that have been found might contain relevant information on the shape  
195 of potential dependencies between variables.

196 In general, the algorithm discards most of the chemical information contained in the CEC  
197 values, and retains pH as the only chemical variable. Temperature is present in all the selected  
198 functions, sometimes in a linear form and more often in an exponential form. Soil texture  
199 appears quite often, but never using coarse fractions of the topsoil as a predictor, and just  
200 seldom considering the gravel content of subsoil (that could be a proxy of other variables as  
201 water infiltration or aeration). Sand is never used, while finer fractions seem to play a role in  
202 predicting soil respiration, probably because of their interaction with soil organic matter.

203

### 204 **4.3 Suitability of the method in the context of soil science**

205 The symbolic regression algorithm finds several potential correlations in the dataset. The first  
206 benefit of this technique is to find hidden relationships between data in a way that is totally  
207 neutral toward the solution and carries absolutely no human bias.

208 Although only two of the selected solutions could be used for predictions, the main asset of the  
209 technique in our case concerns the exploration of possible relationships rather than predictions,  
210 and in this respect the technique presents a good potential. The identification of potential  
211 relationships between variables in a mathematical form and in a way that it is not biased by the  
212 beliefs of the experimenter is an invaluable asset for any model study, and might be significantly  
213 superior to traditional correlation analyses. The suggestion for possible numerical  
214 transformations contained in the best equations found by the EA can represent an important aid  
215 for modelers, although at the moment the technique should be followed by a second phase of  
216 “traditional” modeling with a human expert. We must anyway consider that the accuracy of the  
217 technique is extremely dependent on the number of generations, and therefore any increase in  
218 computing power (foreseeable in a near future on common desktop machines, or already  
219 achievable with relatively cheap infrastructures such as rented cloud grids or clusters) could  
220 increase such accuracy.

221

## 222 **5. Conclusions**

223 The EA-based search identifies a set of solutions performing relatively well in predicting soil  
224 respiration over the training dataset, although performances with the validation dataset are  
225 comparable only in a few cases. The selected solutions contain, nevertheless, relevant  
226 information on possible relationships between the predicted variables and all potential  
227 predictors.

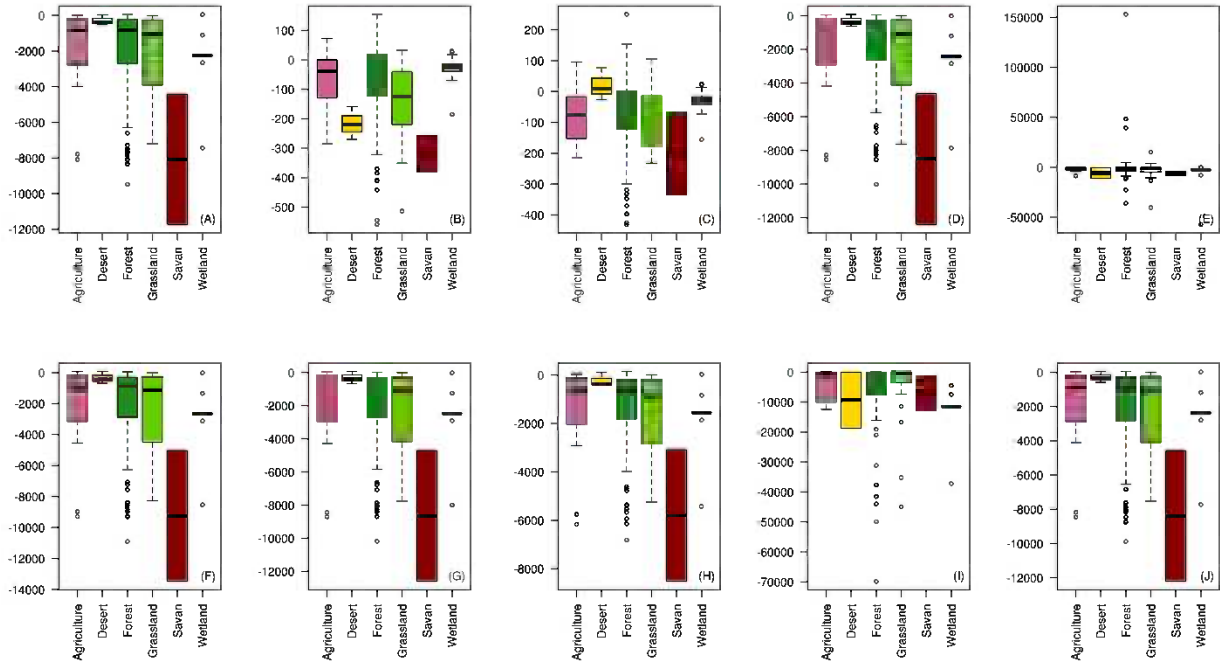
228 The main benefit of this technique is the totally unbiased estimation of possible links between  
229 the variables. The technique explores the most promising part of all possible combinations of  
230 numerical transformations to apply on the data, inside a subset of transformation functions  
231 defined by the user. This allows for a much deeper assessment of correlations between the  
232 variables than traditional techniques of correlation analysis. Still, as an asset over other  
233 machine learning techniques, the EA-based search retains complete transparency to the user.  
234 Solutions found by symbolic regression, although not directly usable for mechanistic modeling,  
235 are a useful tool for data interpretation and could be used for the development of a more  
236 mechanistic model. We therefore advocate for the adoption of symbolic regression techniques  
237 in the early part of the routine analysis workflow of soil related datasets, as an explorative data  
238 mining technique, and particularly as an explorative method for modeling purposes and for the  
239 development of pedotransfer functions.

240

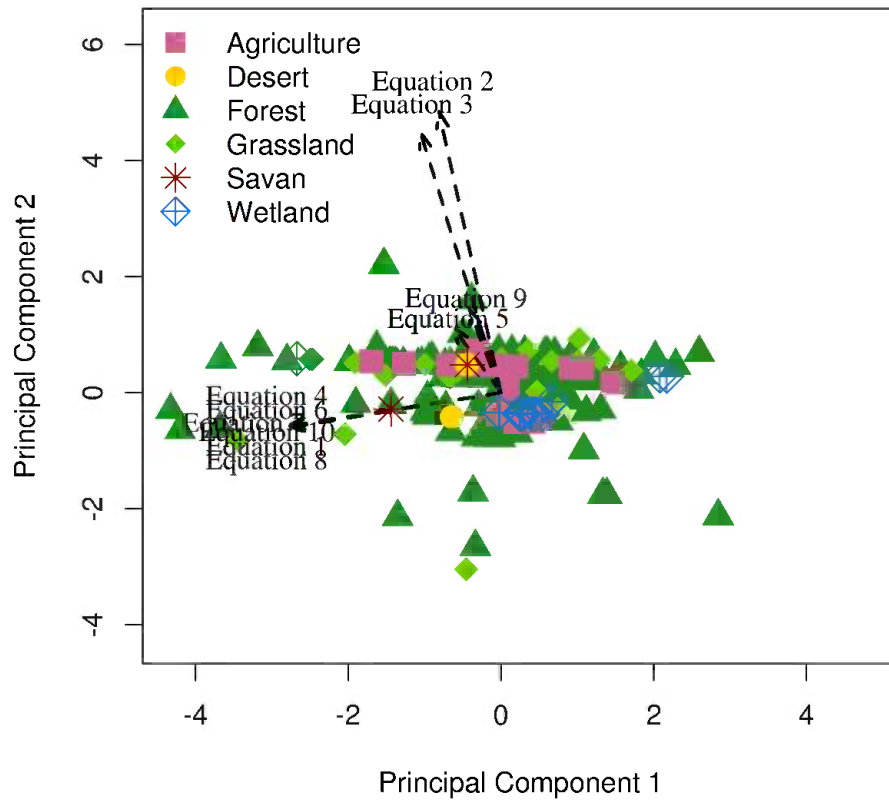
## 241 **References**

- 242 Bond-Lamberty, B.P. & Thomson, A.M., 2012. A Global Database of Soil Respiration Data, Version 2.0.  
243 Data set. Available on-line [<http://daac.ornl.gov>] from Oak Ridge National Laboratory Distributed Active  
244 Archive Center, Oak Ridge, Tennessee, U.S.A. <http://dx.doi.org/10.3334/ORNLDAAC/1070>  
245 Crowe, A., Mcclean, C., & Cresser, M., 2006. An application of genetic algorithms to the robust estimation  
246 of soil organic and mineral fraction densities. *Environmental Modelling & Software*, 21, 1503–1507.  
247 Koza, J. R., 1992. *Genetic programming: on the programming of computers by means of natural*  
248 *selection*. MIT Press, Cambridge, MA, USA  
249 Clery, D. & Voss, D., 2005. All for one and one for all. *Science* 308, 809  
250 FAO/IIASA/ISRIC/ISSCAS/JRC, 2012. *Harmonized World Soil Database (version 1.2)*. FAO, Rome, Italy  
251 and IIASA, Laxenburg, Austria.  
252 Fogel, L. J., 1962. Autonomous automata. *Industrial Research* 4, 14-19.  
253 Holland, J. H., 1992. *Adaptation in natural and artificial systems: An introductory analysis with*  
254 *applications to biology, control, and artificial intelligence*. MIT Press Cambridge, MA, USA  
255 Johari, A., Javadi, A. A., & Habibagahi, G., 2011. Modelling the mechanical behaviour of unsaturated  
256 soils using a genetic algorithm-based neural network. *Computers and Geotechnics*, 38, 2–13.  
257 Padarian, J., Minasny, B., & McBratney, A., 2012. Using genetic programming to transform from  
258 Australian to USDA/FAO soil particle-size classification system. *Soil Research*, 50, 443–446.  
259 Pedroso, D. M., & Williams, D. J., 2011. Automatic calibration of soil–water characteristic curves using  
260 genetic algorithms. *Computers and Geotechnics*, 38, 330–340.  
261 Lloyd, J., & Taylor, J., 1994. On the Temperature Dependence of Soil Respiration. *Functional Ecology*, 8,  
262 315–323.  
263 Moyano, F., Vasilyeva, N., Bouckaert, L., Cook, F., Craine, J., Curiel Yuste, J., Don, A., Epron, D.,  
264 Formanek, P., Franzluebbers, A., Ilsted, U., Kätterer, T., Orchard, V., Reichstein, M., Rey, A., Ruamps,  
265 L., Subke, J.-A., Thomsen, I. K. & Chenu, C., 2011. The moisture response of soil heterotrophic  
266 respiration: interaction with soil properties. *Biogeosciences Discuss*, 8, 11577–11599.

267 Rosca, J. P., 1996. Generality versus size in genetic programming. In Proceedings of the First Annual  
 268 Conference on Genetic Programming, pp. 381-387. MIT Press.  
 269 Schwefel, H.-P., 1956. Cybernetic Evolution as Strategy for Experimental Research in Fluid Mechanics  
 270 (Diploma Thesis in German). Hermann Föttinger-Institute for Fluid Mechanics, Technical University of  
 271 Berlin  
 272 Valdés-Pérez, R. E., 1999. Discovery tools for science apps. Communications of the ACM 42.11 (1999):  
 273 37-41.  
 274 Venables, W., & Ripley, B., 2002. Modern Applied Statistics with S. Springer-Verlag.  
 275



276  
 277 Figure 1: Residuals of the selected functions. A) Equation 1, B) Equation 2, C) Equation 3, D)  
 278 Equation 4, E) Equation 5, F) Equation 6, G) Equation 7, H) Equation 8



279  
 280 Figure 2: PCA analysis of the residuals

281  
 282

283 **Tables**

	Function 1	Function 2	Function 3	Function 4	Function 5	Function 6	Function 7	Function 8
ME	-2101.6	-73.6	-63.2	-2138.0	-1244.5	-2324.8	-2175.7	-1470.1
MAE	2103.4	101.7	90.4	2141.1	4906.5	2327.7	2178.3	1474.2
RMSE	3433.3	148.4	126.9	3505.4	13412.4	3811.7	3566.0	2417.6
NRMSE %	3684.7	159.3	136.2	3762.1	14394.5	4090.8	3827.1	2594.7
PBIAS %	28932.6	1012.7	869.6	49411.5	30274.3	53729.6	50283.1	33974.7
NSE	-1362.3	-1.6	-0.9	-1409.0	-20166.4	-1666.1	-1458.1	-669.7
d	0.0	0.4	0.6	0.0	0.0	0.0	0.0	0.0
r	0.0	0.0	0.2	0.0	0.0	0.0	0.0	0.0
R2	-2101.6	-73.6	-63.2	-2138.0	-1244.5	-2324.8	-2175.7	-1470.1

284  
 285 Table 1: the goodness of fit indicators considered for each function. ME = mean error, MAE =  
 286 mean absolute error, RMSE = root mean squared error, NRMSE = normalized root mean  
 287 squared error (-100% <= nrms <= 100% ), PBIAS = percent bias, NSE = Nash-Sutcliffe  
 288 Efficiency, d = index of agreement (0 <= d <= 1), r = Pearson's correlation coefficient, R<sup>2</sup> =

289 coefficient of determination ( $0 \leq R^2 \leq 1$ ). These indexes have been calculated on the whole  
290 dataset, without removing the outliers.  
291