

1 Selection of Polynomial Features by Genetic Algorithms

2 Francisco Coelho^{a,c,*}, João Pedro Neto^{b,c}

3 ^aDept. Informática, Universidade de Évora, Rua Romão Ramalho 58, 7000-671 Évora

4 ^bDept. Informática, Faculdade de Ciências da Universidade de Lisboa, Campo Grande
5 1749-016 Lisboa

6 ^cLaboratory of Agent Modelling (LabMAg)

7 Abstract

8 Many applications require models that have no acceptable linear approximation
9 and many nonlinear models are defined by polynomials. The use of genetic
10 algorithms to find polynomial models is decades old but still poses challenges
11 due to the complexity of the search and different definitions of optimal solution.

This paper describes a two-step method that uses genetic algorithms and linear regression to find empirical polynomial regressions. Experiments on common datasets show that, discounting the training computational effort, this method is quite competitive.

12 *Keywords:* polynomial regression, genetic algorithm, feature extraction,
13 dimensionality reduction

14 1. Introduction

15 With notable exceptions (*e.g.* neural networks) machine learning regression
16 techniques are based on linear models. The linearity assumption has many
17 advantages including reduced computational complexity and strong theoretical
18 framework. However nonlinearity is unavoidable in many application scenarios,
19 specially those with phase transitions or feedback loops, so common in ecology,
20 cybernetics, robotics and other areas.

21 Polynomials, one of the most studied subjects in mathematics, generalize
22 linear functions and define, perhaps, the simplest and most used nonlinear mod-
23 els. Applications include colorimetric calibration (Mendes and Carvalho, 2005),

*Corresponding author

Email addresses: fc@di.uevora.pt (Francisco Coelho), jpn@di.fc.ul.pt (João Pedro Neto)
Preprint submitted to Elsevier July 26, 2013

24 explicit formulæ for turbulent pipe flows (Davidson et al., 1999), computational
 25 linguistics (Sánchez et al., 2009) and more recently, analytical techniques for
 26 cultural heritage materials (Cséfalvayová et al., 2010), liquid epoxy molding pro-
 27 cess (Chan et al., 2011), B-spline surface reconstruction (Gálvez et al., 2012),
 28 product design (Chan et al., 2012) or forecasting cyanotoxins presence in water
 29 reservoirs (García Nieto et al., 2013). This not only illustrates the wide spec-
 30 trum of applications but, additionally, work in each one of these polynomial
 31 models uses, at some point, a genetic algorithm.

32 Genetic algorithms (GA) where, arguably, one of the hottest topics of re-
 33 search in the recent decades and with good reason since they outline an opti-
 34 mization scheme easy to conceptualize and with very broad application. If a
 35 nonlinear (or otherwise) model requires parameterization GAs provide a simple
 36 and often effective approach to search for locally optimal parameters. Research
 37 related to genetic algorithms abound and spans from the 1950s seminal work
 38 of Nils Aall Barricelli (Barricelli, 1962) in the Institute for Advanced Study of
 39 Princeton to today’s principal area of study for thousands of researchers, cov-
 40 ered in hundreds of conferences, workshops and other meetings. Perhaps the
 41 key impulse to GAs come from John Holland’s work and his book “Adaptation
 42 in Natural and Artificial Systems” (Holland, 1975).

43 One interesting take on genetic algorithms, named *genetic programming* by
 44 John Koza (Koza, 1992), proposed the use of GAs to search the syntactic struc-
 45 ture of complex functions. This syntactic structure search is keen to the central
 46 ideas of deep learning (Bengio et al., 2013; Bengio, 2009), a subarea of ma-
 47 chine learning actually producing quite promising results (*e.g.* in Tarlow et al.
 48 (2013)). It is also related to the work presented in this paper in the sense that,
 49 unlike linear models that have a simple structure, $y = \sum_i \beta_i x_i$, nonlinear (in
 50 particular polynomial) models pose an additional “structure” search problem.

51 The idea of using GAs to find a polynomial regression is not new (Maertens
 52 et al., 2006; Yu and Lin, 2008; Wu et al., 2009) but still generates original
 53 research (Hofwing et al., 2011; Cetisli and Kalkan, 2011). In line with that
 54 research this work describes a general method to find a polynomial regression of

55 a given dataset. The optimal regression minimizes a cost function that accounts
 56 for both the root-mean-square error (error) and a regularization factor to avoid
 57 over-fitting.

58 It turns out that, discarding the computational cost of training, the poly-
 59 nomial regression method presented here, Genetic Algorithms for Polynomi-
 60 als (GAPOLY), provides a quite competitive regression method. Indeed, it is
 61 only systematically out-performed by random forests, an *ensemble* method.

62 The remainder of this paper is organized as usual: the next section describes
 63 the details of our method and is followed by a presentation of some performance
 64 results. The last section draws some conclusions and points future research
 65 tasks.

66 2. Genetic Algorithms for Polynomials

67 This section is dedicated to the description of an algorithm to find a poly-
 68 nomial regression from a given dataset. It starts with a brief introduction and
 69 outline of the algorithm and proceeds into core details as the encoding used
 70 to represent individual polynomial instances in the GA populations and the
 71 regularization of the cost function.

A usual representation of polynomials is

$$p(x_1, \dots, x_m) = \sum_i \theta_i q_i$$

where each q_i is a monomial, $q_i = \prod_j x_j^{\alpha_{ij}}$, the exponents are non-negative integers, $\alpha_{ij} \in \mathbb{N}_0$, and the coefficients are real valued, $\theta_i \in \mathbb{R}$. For example $p(x_1, x_2, x_3) = 2x_1 + x_2x_3 + \frac{1}{2}x_1^2x_3$ has monomials $q_1 = x_1$, $q_2 = x_2x_3$ and $q_3 = x_1^2x_3$, coefficients $\theta_1 = 2$, $\theta_2 = 1$ and $\theta_3 = 1/2$ and exponents $\alpha_{1,1} = 1$, $\alpha_{2,2} = 1$, $\alpha_{2,3} = 1$, $\alpha_{3,1} = 2$, $\alpha_{3,3} = 1$ and all other $\alpha_{ij} = 0$. The exponents alone are a matrix that defines the monomial structure of the polynomial, $A = [\alpha_{ij}]$. For the example above

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 2 & 0 & 1 \end{bmatrix} \sim \begin{bmatrix} x_1 \\ x_2x_3 \\ x_1^2x_3 \end{bmatrix}$$

72 where each row defines a monomial and each column represents a variable.
 73 Changing the order of the rows doesn't change the polynomial whereas changing
 74 the order of the columns corresponds to changing the respective variables.

75 This representation of polynomials makes the problem of structure search
 76 very clear: except for the trivial cases, the number of possible monomials given
 77 n variables and a maximum joint degree d grows exponentially with either n or
 78 d . But more importantly, the polynomial regression problem can be naturally
 79 split into two subproblems:

- 80 1. For a given set of monomials $\mathcal{P} = \{q_1, \dots, q_k\}$, find the regression coeffi-
 81 cients $\theta_1, \dots, \theta_k$ that minimize the error on a given dataset;
- 82 2. Find the fittest set of monomials, *i.e.* the polynomial that minimizes the
 83 error on the same dataset;

84 More precisely, concerning the first problem, let \mathcal{D} be a dataset with n obser-
 85 vations of variables Y, X_1, \dots, X_m and $\mathcal{P} = \{q_1, \dots, q_k\}$ a set of k monomial
 86 expressions over X_1, \dots, X_m . Define the hypothesis¹

$$h_{\Theta, \mathcal{P}}(x_1, \dots, x_m) = \sum_{j=1}^k \theta_j q_j|_{X_i=x_i, \forall 1 \leq i \leq m}$$

87 and let the cost

$$J_{\text{fit}}(\Theta; \mathcal{P}, \mathcal{D}) = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - h_{\Theta, \mathcal{P}}(x_1^{(i)}, \dots, x_m^{(i)}) \right)^2} \quad (1)$$

88 be the usual root-mean-square error (error) function. Now the first problem can
 89 be stated as: *Given a dataset \mathcal{D} and a set of monomials \mathcal{P} , find parameters Θ*
 90 *that minimize $J_{\text{fit}}(\Theta; \mathcal{P}, \mathcal{D})$.*

91 It turns out that this problem can be solved as a usual linear regression
 92 problem by expanding \mathcal{D} with columns that replicate the monomials in \mathcal{M} .

¹The expression $q|_{X=x}$ reads “replace all instances of X by x in q ”.

Algorithm 1 GAPOLY uses linear regression to find monomial coefficients that minimize the error over a dataset and GAs to explore the space of polynomials. The role of the linear regression step is to produce a fitting error that sorts the population. At exit the error of the fittest instance is bounded by ϵ .

```

function GAPOLY( $D, pop_0, \epsilon$ )
     $pop \leftarrow pop_0; err \leftarrow 1.0 + \epsilon$ 
    while  $err > \epsilon$  do
         $pop \leftarrow \text{ITERATEGA}(pop)$ 
         $pop \leftarrow \text{SORT}(pop, key = J)$   $\triangleright$  Sort by regression cost,  $J$ 
         $err \leftarrow J(\text{FIRST}(pop))$ 
    end while
    return  $\text{FIRST}(pop)$ 
end function

```

93 The second problem is treated in the GA setting: Let \mathcal{D} be a dataset as
 94 above and Q a set of polynomials. For each polynomial $p \in Q$ let \mathcal{P}_p be the set
 95 of monomials in p (without the coefficients) and compute the fitness

$$\phi_p = \min_{\Theta} J_{\text{fit}}(\Theta; \mathcal{P}_p, \mathcal{D})$$

96 by solving the first problem. With a fitness of every instance, a GA will apply
 97 genetic operators (usually mutation and crossover) to evolve the population Q
 98 until a reasonable approximation of a local minimum is found. Notice that the
 99 properties of GAs and linear regression entail that the composition of GAs with
 100 linear regression, as defined in Algorithm 1, converges to a polynomial that is a
 101 local minimum of the fitness function, encapsulated in the error function J_{fit} .

102 Subsection 2.1 describes the encoding of individual polynomial instances as
 103 chromosomes and other parameters of the utilized GA implementation. The
 104 regularization of the cost function is discussed in subsection 2.2.

105 2.1. Polynomial Encoding

106 The encoding for any polynomial will be as follows:

1. an initial segment detailing which monomials are active (the first monomial is always active), this is represented in unary description, i.e., each monomial is identified by a single bit
2. the remaining bits are split into k sets of equal size, each one representing a monomial
3. each monomial is split into m sets of d size each, i.e., a variable
4. for each variable, the remaining bits are the binary description of the variable degree, i.e., the maximum exponent is given by $2^d - 1$

Let's see an example: consider polynomial $x_1^3x_3+x_3^7+x_1x_2$ with $k = 4, m = 3$ and $d = 3$. One possible encoding would be:

110 - 011,000,001 ; 000,000,111 ; 001,001,000 ; 110,010,101

(for reading purposes the semicolons separate monomials, the commas separate variables)

The first three bits inform that the second and third monomials are active while the fourth is not (as said, the first monomial is always active). This last monomial does not enter neither in the polynomial regression nor in the error (fitness) evaluation. However, it acts as a kind of junk DNA, becoming active when, in a future mutation or crossover, the third bit of the entire sequence flips from 0 to 1.

Let's interpret the first monomial description, 011,000,001. It is divided by three since $m = 3$. The first triple 011 is the binary description of the exponent of variable x_1 which is 3, so the first monomial includes x_1^3 . The second triple, 000, means that x_2 is not part of the monomial. The third triple 001 says that variable x_3 has exponent 1. The first monomial is $x_1^3x_3$.

Notice that all binary descriptions give rise to valid polynomials. However this is not a bijective mapping. For each polynomial there are multiple representations. For example, $x_1 + x_2$ and $x_2 + x_1$ have different representations. The authors considered that more complex mappings in order to force a one to one mapping would impact negatively in the algorithm's performance.

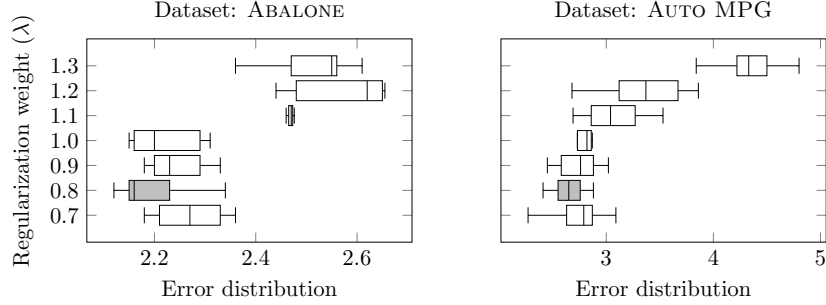


Figure 1: Error distribution by regularization exponent for the Abalone and Auto MPG datasets. The box plots summarize the error values of 10 runs for each value of λ . The smallest overall error is shown in gray and, in both cases, was achieved with $\lambda = 0.8$.

2.2. Cost Function

The polynomial regression error considered so far is based on the ability to predict the test set after the polynomial regression has found the appropriate coefficients θ_i for each one of the monomials q_i .

This error function tends to prefer more complex polynomials, namely in the number of monomials which provides the regression algorithm for more fitting possibilities. One way to balance this is to provide a regularization term into the error function. Our proposal is to include a multiplicative factor proportional to the number of monomials. Thus, the error from equation 1 defines

$$J_{\text{reg}}(\Theta; \mathcal{P}, \mathcal{D}) = \lambda^k J_{\text{fit}}(\Theta; \mathcal{P}, \mathcal{D}) \quad (2)$$

where k is the number of monomials in the polynomial. When $\lambda > 1$ polynomials with more monomials are penalized.

Somewhat unexpectedly after some experiences it was found that lower values for λ provide better, even if marginal, results. Figure 1 shows regularization results for the Abalone and Auto MPG datasets with ten runs for each λ . The following section includes information about these datasets.

The typical inflection point lies around $\lambda = 0.8$. The dataset results for applying the proposed regression method use the regularization parameter with this value.

154 2.3. Genetic Operators

155 To perform the genetic algorithm it was used the R package `genalg` (Wil-
156 lighagen, 2012). The operators were the standard ones: (a) crossover, *i.e.*, a
157 pair of solutions from the previous generations are combined by splitting and
158 mixing their respective representations, and (b) mutation, changing the values
159 of single bits; the mutation chance applied in the datasets was 5%. There was
160 also elitism between generations, *i.e.*, 20% of the best solutions survive to the
161 next generation.

162 3. Experimental Results

163 The results were found using R programming language (R Core Team, 2013)².
164 To compare this paper’s proposed algorithm we applied the exact same train and
165 test samples using several well-known learning algorithms for regression, namely:
166 Linear Regression, Support Vector Machines (Meyer et al., 2012), Regression
167 Trees (Therneau et al., 2013), Conditional Inference Trees (Hothorn et al., 2006;
168 Strobl et al., 2007, 2008) and Random Forests (Liaw and Wiener, 2002)

169 In order to train and test the performance of GAPOLY several mainstream
170 datasets were used. For each dataset, we selected 70% for training purposes
171 and the remaining observations to make a test set in order to compute the
172 estimated error. To achieve more robust results, each dataset were processed
173 25 times, each one with different samples for the train and test sets. For the
174 datasets with attribute values of different magnitudes, a preliminary scaling was
175 executed. The results below are box plots for the test set error predictions over
176 these different runs.

177 ARTIFICIAL this is an artificial dataset with four numeric features, x_1, \dots, x_4 ,
178 where x_1, x_3 are outcomes from Poisson random variables, and x_2, x_4 from

²The datasets and the R code used to produce the results and plots in this paper are available online at <https://github.com/jpneto/GenAlgPoly>.

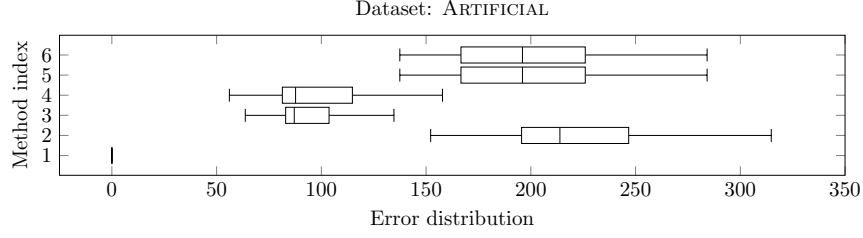


Figure 2: Results for Artificial dataset. As shown, GAPOLY was able to find the exact polynomial structure that generates the dataset, reducing the test set prediction error to zero. The regression methods depicted are: 1. GAPOLY, 2. Random Forests, 3. Linear Regression, 4. SVM, 5. Regression Trees and 6. Conditional Inference Trees

179 Normal random variables. The dependent variable y is given by expression
180 $x_2x_4^2 + x_1^2x_3 + 5$. The dataset includes $n = 50$ observations.

181 This dataset was used in order to verify if GAPOLY was able to find the
182 polynomial relation, which the algorithm did (cf. figure 2). The genetic
183 algorithm run with a population of $n = 100$ solutions with 50 iterations
184 for each run.

185 In the next datasets, the population of the genetic algorithm had size $n =$
186 250 with 100 iterations for each run.

187 HOUSING: This data set concerns the task of predicting housing values in areas
188 of Boston. There are $m = 13$ continuous attributes and the dependent
189 variable is the median value of owner-occupied homes in \$1000's. There
190 are $n = 506$ observations.

191 Just as an example of the model the GAPOLY algorithm outputs: in this
192 dataset the best polynomial,

$$y = -0.12x_6x_9^4 + 0.78x_6 - 0.4x_9^2x_{13} + 0.17x_{13}^2 - 0.044$$

193 where the attributes mean: x_6 , RM average number of rooms per dwelling;
194 x_9 , RAD index of accessibility to radial highways; x_{13} , Lower status of the
195 population.

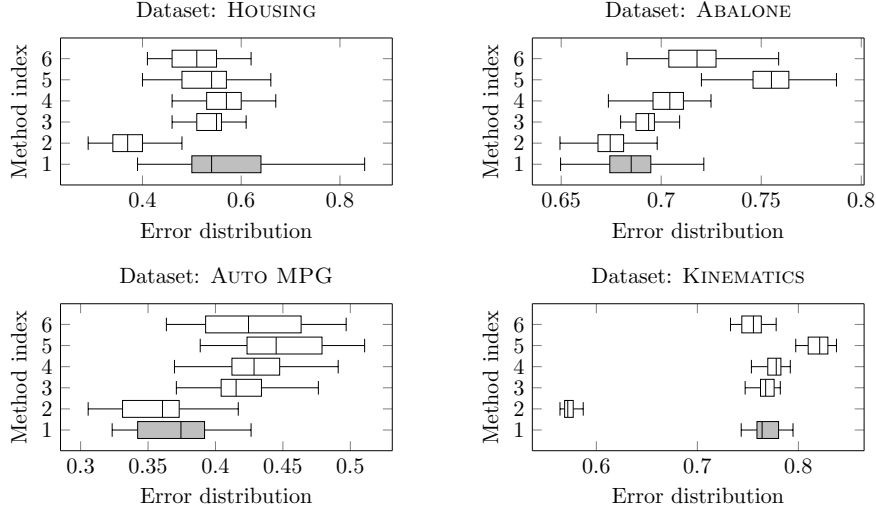


Figure 3: Summary results for different regression methods on the Housing, Abalone, Auto MPG and Kinematics datasets. Random Forests, an ensemble method (“2” in the plots), has the better results. However GAPOLY achieves similar errors in two of the datasets. For the remaining two GAPOLY produces similar errors to the other classic regression methods. The regression methods depicted in these figures are: 1. GAPOLY, 2. Random Forests, 3. Linear Regression, 4. SVM, 5. Regression Trees and 6. Conditional Inference Trees

For comparison, if we access the mean decrease in accuracy found by Random Forests, the most important attributes are — in decreasing order — x_{13}, x_6, x_5 (but x_5 is already considered 4 times less important than x_6). Both algorithms agree in two of their three most important attributes.

ABALONE This dataset can be used to predict the age of a abalone shell using the given $m = 8$ numeric attributes concerning several physical measurements. There are $n = 4177$ observations.

AUTO MPG This dataset is used to predict fuel consumption in miles per gallon, based on two discrete and five continuous attributes ($m = 7$). There are $n = 398$ observations.

KINEMATICS This dataset is concerned with the realistic simulation of the forward kinematics of an 8 link robot arm. The task is to predict the distance

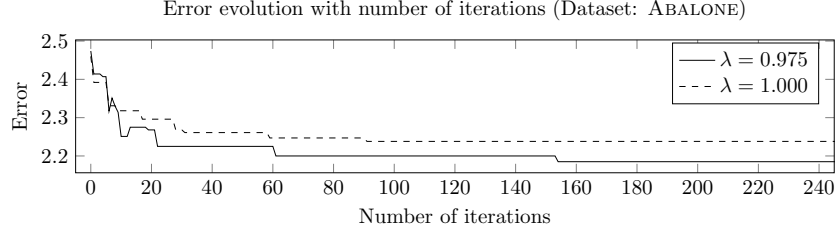


Figure 4: Error progress for Abalone dataset during a single execution of the genetic algorithm. The figure shows the fitness evolution for two different regularization values. The population for both consisted of 200 polynomials. The error values seem to stabilize around iteration 250.

of the end-effector from a target using $m = 8$ continuous attributes. There are $n = 8192$ observations.

3.1. Convergence speed

The GA quickly proceeds in the first 50 to 100 generations to reasonable error rates. Then, it proceeds slower achieving best solutions with marginal error reduction. Since the entire learning process takes some time, in the current R implementation, placing a limit between 50 to 100 generations already achieves good results, relative to higher iteration values. Figure 4 shows a typical error evolution for the dataset Abalone given two different values for λ .

4. Conclusion and Future Work

The proposed method has competitive results comparing with some well-known regression methods. Only Random Forest outperforms GAPOLY systematically (it also outperforms all the other regression algorithms). One exception — besides the artificial dataset that uses a straightforward polynomial relation — is the Auto MPG dataset where GAPOLY has comparable results. This is evidence that applying standard genetic operators for polynomial model searching is a viable tool for regression purposes.

For complexity considerations GAPOLY demands some processing time. On a current quad-core computer, processing the Kinematics dataset (with 8k ob-

227 servations) takes approximately 5 minutes. The processing time can probably
228 be speeded by one to two orders in magnitude if the process is implemented in
229 a low level programming language like C++. However, speed optimization was
230 not the focus of this article.

231 A cross-validation procedure can be implemented to refine the appropriate
232 parameter values to achieve better errors. Namely, the regularization parameter,
233 λ , can be tested with several values, instead of being fixed at 0.8. Other
234 parameters like mutation chance or the amount of elitism could also be tested.
235 However, these type of tests need a low-level, fast implementation of GAPOLY.

236 Acknowledgements

237 The authors are grateful to the Fundação para a Ciência e Tecnologia (FCT)
238 and the R&D laboratory LabMAG for the financial support given to this work,
239 under the strategic project PEST-OE/EEI/UI0434/2011.

240 The datasets used herein were selected from Luís Torgo's data repository,
241 <http://www.dcc.fc.up.pt/~ltorgo/Regression/DataSets.html>. Most come
242 originally from UCI ML repository, <http://archive.ics.uci.edu/ml/>.

243 The authors wish to thank professor André Falcão for motivation and useful
244 discussions around the article.

245 Barricelli, N.A., 1962. Numerical testing of evolution theories. part i: Theoret-
246 ical introduction and basic tests. *Acta Biotheoretica* 16, 69–98.

247 Bengio, Y., 2009. Learning deep architectures for AI. *Foundations and trends*
248 *in Machine Learning* 2, 1–127.

249 Bengio, Y., Courville, A., Vincent, P., 2013. Representation learning: A review
250 and new perspectives .

251 Cetisli, B., Kalkan, H., 2011. Polynomial curve fitting with varying real powers.
252 *Electronics and Electrical Engineering* 112, 117–122.

253 Chan, K.Y., Dillon, T.S., Kwong, C.K., 2011. Modeling of a liquid epoxy
 254 molding process using a particle swarm optimization-based fuzzy regression
 255 approach. *Industrial Informatics, IEEE Transactions on* 7, 148–158.

256 Chan, K.Y., Kwong, C., Dillon, T.S., 2012. Development of product design
 257 models using fuzzy regression based genetic programming, in: *Computational*
 258 *Intelligence Techniques for New Product Design*. Springer, pp. 111–128.

259 Cséfalvayová, L., Pelikan, M., Kralj Cigić, I., Kolar, J., Strlič, M., 2010. Use of
 260 genetic algorithms with multivariate regression for determination of gelatine
 261 in historic papers based on FT-IR and NIR spectral data. *Talanta* 82, 1784–
 262 1790.

263 Davidson, J., Savic, D., Walters, G., 1999. Method for the identification of
 264 explicit polynomial formulae for the friction in turbulent pipe flow. *Journal*
 265 *of Hydroinformatics* 1, 115–126.

266 Gálvez, A., Iglesias, A., Puig-Pey, J., 2012. Iterative two-step genetic-algorithm-
 267 based method for efficient polynomial b-spline surface reconstruction. *Infor-*
 268 *mation Sciences* 182, 56–76.

269 García Nieto, P.J., Alonso Fernández, J., de Cos Juez, F., Sánchez Lasheras, F.,
 270 Díaz Muñoz, C., 2013. Hybrid modelling based on support vector regression
 271 with genetic algorithms in forecasting the cyanotoxins presence in the trasona
 272 reservoir (northern spain). *Environmental research* .

273 Hofwing, M., Strömberg, N., Tapankov, M., 2011. Optimal polynomial re-
 274 gression models by using a genetic algorithm, in: *Proceedings of the Second*
 275 *International Conference on Soft ComputingTechnology in Civil, Structural*
 276 *and Environmental Engineering Conference,(Crete, Greece), 2011*009.

277 Holland, J.H., 1975. *Adaptation in natural and artificial systems: An introduc-*
 278 *tory analysis with applications to biology, control, and artificial intelligence.*
 279 *U Michigan Press.*

280 Hothorn, T., Hornik, K., Zeileis, A., 2006. Unbiased recursive partitioning: A
281 conditional inference framework. *Journal of Computational and Graphical*
282 *Statistics* 15, 651–674.

283 Koza, J.R., 1992. *Genetic Programming: vol. 1, On the programming of com-*
284 *puters by means of natural selection. volume 1.* MIT press.

285 Liaw, A., Wiener, M., 2002. Classification and regression by randomforest. *R*
286 *News* 2, 18–22. URL: <http://CRAN.R-project.org/doc/Rnews/>.

287 Maertens, K., De Baerdemaeker, J., Babuška, R., 2006. Genetic polynomial re-
288 gression as input selection algorithm for non-linear identification. *Soft Com-*
289 *puting* 10, 785–795.

290 Mendes, L., Carvalho, P.d., 2005. Adaptive polynomial regression for colorimet-
291 ric scanner calibration using genetic algorithms, in: *Intelligent Signal Process-*
292 *ing, 2005 IEEE International Workshop on, IEEE.* pp. 22–27.

293 Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., Leisch, F., 2012. e1071:
294 Misc functions of the department of statistics (e1071), tu wien URL: [http:](http://CRAN.R-project.org/package=e1071)
295 [//CRAN.R-project.org/package=e1071](http://CRAN.R-project.org/package=e1071). r package version 1.6-1.

296 R Core Team, 2013. *R: A language and environment for statistical computing*
297 URL: <http://www.R-project.org/>.

298 Sánchez, L., Otero, J., Couso, I., 2009. Obtaining linguistic fuzzy rule-based
299 regression models from imprecise data with multiobjective genetic algorithms.
300 *Soft Computing* 13, 467–479.

301 Strobl, C., Boulesteix, A.L., Kneib, T., Augustin, T., Zeileis, A., 2008. Condi-
302 tional variable importance for random forests. *BMC Bioinformatics* 9. URL:
303 <http://www.biomedcentral.com/1471-2105/9/307>.

304 Strobl, C., Boulesteix, A.L., Zeileis, A., Hothorn, T., 2007. Bias in random forest
305 variable importance measures: Illustrations, sources and a solution. *BMC*
306 *Bioinformatics* 8. URL: <http://www.biomedcentral.com/1471-2105/8/25>.

- 307 Tarlow, D., Sutskever, I., Zemel, R.S., 2013. Stochastic k-neighborhood selec-
 308 tion for supervised and unsupervised learning. *Journal of Machine Learning*
 309 *Research* .
- 310 Therneau, T., Atkinson, B., Ripley, B., 2013. rpart: Recursive partitioning R
 311 package version 4.1-1.
- 312 Willighagen, E., 2012. genalg: R based genetic algorithm.
- 313 Wu, C.H., Tzeng, G.H., Lin, R.H., 2009. A novel hybrid genetic algorithm
 314 for kernel function and parameter optimization in support vector regression.
 315 *Expert Systems with Applications* 36, 4725–4735.
- 316 Yu, T.L., Lin, W.K., 2008. Optimal sampling of genetic algorithms on poly-
 317 nomial regression, in: *Proceedings of the 10th annual conference on Genetic*
 318 *and evolutionary computation*, ACM. pp. 1089–1096.