

# A method for regularization of evolutionary polynomial regressions

Francisco Coelho<sup>a,\*</sup>, João Pedro Neto<sup>b</sup>

<sup>a</sup>*Departamento de Informática, Universidade de Évora, Rua Romão Ramalho 58, 7000-671 Évora, Portugal*

<sup>b</sup>*Departamento de Informática, Faculdade de Ciências da Universidade de Lisboa, Campo Grande, 1749-016 Lisboa, Portugal*

---

## Abstract

While many applications require models that have no acceptable linear approximation, the simpler nonlinear models are defined by polynomials. The use of genetic algorithms to find polynomial models from data is known as Evolutionary Polynomial Regression. This paper introduces Evolutionary Polynomial Regression with Regularization, an algorithm that extends the EPR method and describes a set of experiences on common datasets that compare both flavors of EPR and other methods including Linear Regression, Regression Trees and Support Vector Regression.

The empiric conclusion of those experiments is that EPR with regularization is able to achieve better fitting than other non-ensemble methods and it has shorter computation time than plain EPR.

*Keywords:* evolutionary polynomial regression, regularization, feature extraction

---

## 1. Introduction

With notable exceptions (*e.g.* neural networks) machine learning regression techniques produce linear models. The linearity assumption has many advantages including reduced computational complexity and strong theoretical framework. However nonlinearity is unavoidable in many application scenarios, specially those with phase transitions or feedback loops, so common in engineering, ecology, cybernetics and

---

\*Corresponding author: Tel.: +351-919-006-379  
Email address: fc@di.uevora.pt (Francisco Coelho)

7 other areas. The kernel trick in Support Vector Machines (SVM) ([24, 18, 1]) al-  
8 levates this problem by allowing special non-linear transformations of the feature-  
9 space. The condition such transformations must meet is known as the *kernel trick*,  
10  $k(x, x') = \langle \varphi(x), \varphi(x') \rangle$ , where  $\varphi$  is the feature-space transformation and  $\langle \cdot, \cdot \rangle$  denotes  
11 inner product. The “trick” consists on computing the kernel  $k(x, x')$  while avoiding the  
12 computation of the inner product and the transformations  $\varphi(x), \varphi(x')$ . A special case  
13 of polynomial transformation, the *polynomial kernel*,  $k(x, x') = \langle x, x' \rangle^d$  is commonly  
14 used in regression and classification tasks with SVMs. However general polynomial  
15 transformations do not verify the kernel trick.

16 Polynomials, one of the most studied subjects in mathematics, generalize linear  
17 functions and define, perhaps, the simplest and most used nonlinear models. Appli-  
18 cations include colorimetric calibration [20], explicit formulæ for turbulent pipe flows  
19 [9], computational linguistics [23] and more recently analytical techniques for cultural  
20 heritage materials [8], liquid epoxy moulding process [6], B-spline surface reconstruc-  
21 tion [11], product design [7] or forecasting cyanotoxins presence in water reservoirs  
22 [12]. These examples not only illustrate the wide spectrum of applications but, addi-  
23 tionally, each one uses, at some point, Genetic algorithms (GA).

24 Evolutionary algorithms, including GA, were, arguably, one of the hottest topics  
25 of research in the recent decades and with good reason since they outline an optimiza-  
26 tion scheme easy to conceptualize and with very broad application. If a nonlinear (or  
27 otherwise) model requires parameterization, GAs provide a simple and often effective  
28 approach to search for locally optimal parameters. Related research abound and spans  
29 from the 1950s seminal work of Nils Aall Barricelli [2] in the Institute for Advanced  
30 Study of Princeton to today’s principal area of study for thousands of researchers, cov-  
31 ered in hundreds of conferences, workshops and other meetings. Perhaps the key im-  
32 pulse to GAs come from John Holland’s work and his book “Adaptation in Natural and  
33 Artificial Systems” [15].

34 One interesting variation of genetic algorithms, named *genetic programming* by  
35 John Koza [17], proposes the use of GAs to search the syntactic structure of complex  
36 functions. Syntactic structure search is also keen to the central ideas of deep learning  
37 [3, 4], a subarea of machine learning actually producing quite promising results (*e.g.* in

[27]). It is also related to the work presented in this paper in the sense that, unlike linear models that have a simple structure,  $y = \sum_i \beta_i x_i$ , nonlinear (in particular polynomial) models pose an additional structure search problem.

The idea of using GAs to find a polynomial regression is not new [19, 31, 30] but still generates original research [14, 5]. The modern formulation of the use of GA to find polynomial models is known as Evolutionary Polynomial Regression (EPR) and systematization can be traced back to the work of Davidson, Savic and Walters [10]. Further developments include multi-objective optimizations [13].

This paper describes an extension of the general EPR method to find a regularized polynomial regression of a given dataset. The optimal regression results from a cost function that accounts for both the root-mean-square (error) and a regularization factor to avoid overfit by penalysing polynomial complexity.

The next section describes the method's details and is followed by a presentation of some performance results. The last section draws some conclusions and points future research tasks.

## 2. Genetic Algorithms for Polynomials

This section starts with a brief introduction and outline of the evolutionary polynomial regression algorithm, EPR, and proceeds into core details as the encoding used to represent individual polynomial instances in the GA populations and the regularization of the cost function.

An usual representation of polynomials is through expressions of the form

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

where each  $q_i = \prod_j x_j^{\alpha_{ij}}$  is a monomial, the exponents  $\alpha_{ij} \in \mathbb{N}_0$  are non-negative integers and the coefficients  $\theta_i \in \mathbb{R}$  are real valued. 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  $q_3 = x_1^2x_3$ , exponents  $\alpha_{1,1} = 1, \alpha_{2,2} = 1, \alpha_{2,3} = 1, \alpha_{3,1} = 2, \alpha_{3,3} = 1$  and all other  $\alpha_{ij} = 0$  and coefficients  $\theta_1 = 2, \theta_2 = 1$  and  $\theta_3 = 1/2$ .

The exponents alone can be organized into a matrix  $[\alpha_{ij}]$  that defines the monomial structure of the polynomial. For the example above the matrix representation of the monomials is

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 2 & 0 & 1 \end{bmatrix} \sim \begin{bmatrix} x_1 \\ x_2 x_3 \\ x_1^2 x_3 \end{bmatrix}$$

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

This partial representation of polynomials makes the problem of structure search very clear: except for the trivial cases, the number of possible monomials given  $n$  variables and a maximum joint degree  $d$  grows exponentially with either  $n$  or  $d$ . But more importantly, by separating the set of monomials from the coefficients, the polynomial regression problem can be naturally split into two subproblems:

1. For a given set of monomials  $Q = \{q_1, \dots, q_k\}$  find the regression coefficients  $\Theta = \{\theta_1, \dots, \theta_k\}$  that minimize the error on a given dataset;
2. Find the fittest set of monomials, *i.e.* the polynomial that minimizes the error on the same dataset;

More precisely, concerning the first problem, let  $\mathcal{D}$  be a dataset with  $n$  observations of variables  $Y, X_1, \dots, X_m$  and  $Q = \{q_1, \dots, q_k\}$  a set of  $k$  monomial expressions over  $X_1, \dots, X_m$ . Define the hypothesis<sup>1</sup>

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

and let the error (as “cost”) be

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

---

<sup>1</sup>The expression “ $q|_{X=x}$ ” reads “ $q$  with all instances of  $X$  replaced by  $x$ .”

---

**Algorithm 1** This EPR algorithm uses linear regression for the calculation of the error  $J$  and the space of polynomials is searched in the GAs iteration step. At exit the error of the fittest instance is bounded by  $\epsilon$  or the maximum number of allowed iterations.

---

```

function EPR( $D, pop_0, \epsilon, maxiter$ )
     $pop \leftarrow pop_0; err \leftarrow 1.0 + \epsilon$ 
    while  $err > \epsilon \wedge iterations < maxiter$  do
         $pop \leftarrow \text{ITERATEGA}(pop)$ 
         $pop \leftarrow \text{SORT}(pop, key = J)$  ▷ Sort population by regression error
         $err \leftarrow J(\text{FIRST}(pop))$ 
    end while
    return  $\text{FIRST}(pop)$ 
end function

```

---

79 the usual root-mean-square (error) function. Now the first problem can be stated as:  
80 *Given a dataset  $\mathcal{D}$  and a set of monomials  $Q$  find parameters  $\Theta$  that minimize the cost*  
81  *$J_{fit}(\Theta; Q, \mathcal{D})$ .*

82 This is a simple linear regression problem obtained by expanding  $\mathcal{D}$  with columns  
83 that replicate the monomials in  $Q$ . The resulting dataset,  $\mathcal{D} \cup Q(\mathcal{D})$ , adds the monomial  
84 transformations in  $Q$  to the original dataset  $\mathcal{D}$ . An alternative formulation would just  
85 replace  $\mathcal{D}$  by  $Q(\mathcal{D})$ . It turns out that the first formulation is a special case of the second  
86 (by including the variables in the monomial set) and has better error performance —  
87 what is not surprising because it uses more features.

88 The second problem is treated in the GA setting: Let  $\mathcal{D}$  be a dataset as above and  
89  $\mathcal{P}$  a set of polynomials. For each polynomial  $p \in \mathcal{P}$  let  $Q_p$  be the set of monomials in  
90  $p$  (without the coefficients) and define the (anti) fitness

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

91 by solving the first problem. With a fitness of every instance, the GA genetic operators  
92 (usually mutation and crossover) evolve the population  $\mathcal{P}$  until a reasonable approxi-  
93 mation of a local minimum is found. The properties of GAs and linear regression entail  
94 that Algorithm 1 converges to a polynomial that is a local minimum of the fitness func-

tion, encapsulated in the error function  $J_{\text{fit}}$ .

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

### 2.1. Polynomial Encoding

The specific encoding (representation) of a set of monomials is an important aspect in the implementation of EPR. The choice described below permits active and inactive monomials for regression purposes. The active (or inactive) state of a monomial might change through mutation or crossover. This simple mechanism enhances variation in the complexity of polynomial expressions by evolutionary operations.

Let  $\{q_1, \dots, q_k\}$  be a set of monomials over the variables  $X_1, \dots, X_m$ . The encoding of that set using  $d$  bits per exponent is a binary list such that

1. the initial segment of  $k$  bits defines the active state of each monomial;
2. the remaining bits are split into  $k$  segments of size  $m \times d$ , each representing a monomial;
3. the bits in each monomial segment are split into  $m$  sub-segments of size  $d$ . The  $j^{\text{th}}$  sub-segment is the binary representation of the degree of the variable  $X_j$  in the enclosing monomial segment;

This encoding can also be viewed as the flattening of the binary exponents in the matrix representation prefixed by the activation segment. The set  $\{x_1^3 x_3, x_3^7, x_1 x_2\}$  (with  $m = 3$ ) has matrix representation

$$\begin{bmatrix} 3 & 0 & 1 \\ 0 & 0 & 7 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 011 & 000 & 001 \\ 000 & 000 & 111 \\ 001 & 001 & 000 \end{bmatrix}_{(2)}$$

where the right matrix is in binary form using  $d = 3$  bits. An encoding of this set of monomials with the extra monomial  $x_1^6 x_2^2 x_3^5$  inactive, setting  $k = 4$ , would be

1110; 011, 000, 001; 000,000,111; 001,001,000; 110,010,101

116 where, for reading purposes, semicolons separate segments and commas separate vari-  
 117 ables. The first  $k = 4$  bits inform that the first, second and third monomials are active  
 118 while the fourth is not.

119 While each valid encoding represents a set of monomials the map is not bijective:  
 120 each set of monomials has multiple encodings, for example by changing  $d$  or the or-  
 121 der of monomial segments. However, considering the EPR task, this is a minor issue  
 122 and a bijective map would add computational complexity and negative impact to the  
 123 algorithm’s performance.

124 There is one final remark concerning this encoding method. As it is, the activation  
 125 segment can become all zeros, representing the empty set of monomials. This situation  
 126 can be avoided with a simple hack: Given an encoding, the first monomial is always  
 127 considered active, thus restricting the syntactic form of encodings to binary strings  
 128 starting with 1. In practice, this means that the implementation of the encoding can  
 129 omit the first bit.

## 130 2.2. Cost Function

131 The polynomial regression error considered so far accounts for the ability to predict  
 132 the transformed testset. A known problem of using a cost function based only in the  
 133 dataset error (and of polynomial regressions in general) is the tendency to overfit train-  
 134 ing data. Excessive variance of the estimation method can be reduced by regularizing  
 135 the error function with a penalty factor. Thus, to reduce polynomial complexity and  
 136 variance by regularizing the size of the monomial set the error function from equation  
 137 1 is multiplied by a factor  $\lambda^k$

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

138 where  $k$  is the number of monomials in the polynomial. When  $\lambda > 1$  polynomials  
 139 with more monomials are penalized. The regularized extension of EPR is denoted by  
 140 Evolutionary Polynomial Regression with Regularization (EPRR).

141 A simple exploration on the effect of the regularization parameter is depicted in  
 142 Figure 1 where it is possible to observe that the typical inflection point lies around  
 143  $\lambda = 0.8$ . This value, favoring “larger” polynomials is justified by the balance of the

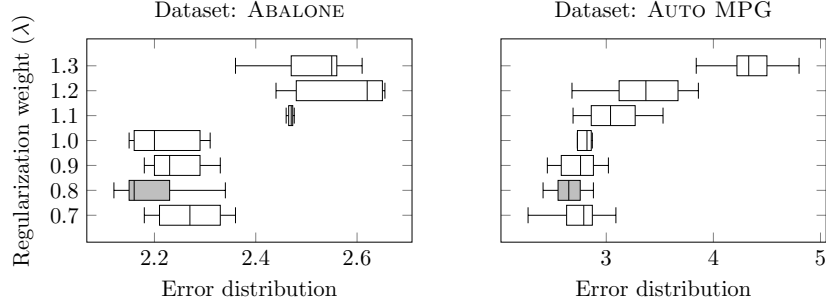


Figure 1: Error distribution by regularization exponent for two common datasets. The box plots summarize error values of ten simulations for each value of  $\lambda$ . The smallest overall error, in grey, is achieved in both datasets when  $\lambda = 0.8$ . Performance of the non-regularized EPR is plotted in the line  $\lambda = 1$ .

144 data's non-linearity and polynomial complexity: below  $\lambda = 0.8$ , even penalized, larger  
 145 monomial sets achieve better error performance than smaller ones while above that  
 146 value the size of the monomial set is excessive. Within this tension the overall error  
 147 results reduced when compared to the non-regularized EPR version.

### 148 2.3. Genetic Algorithm Parameterization

149 In general GAs offer many possibilities with respect to the choice of genetic op-  
 150 erators and respective application rates, population evolution, *etc.* The results found  
 151 here where obtained using the package `genalg` [29] with default parameters, standard  
 152 operators (crossover and mutation) and population evolution with 20% elitism between  
 153 generations.

## 154 3. Experimental Results

155 Here is described the experiment setup used to gather and summarize the empir-  
 156 ical evidence that supports this comparative study of EPR and EPRR. Evaluation is  
 157 focused in error distribution and, besides EPR and EPRR, also uses several common  
 158 regression methods and datasets easily accessible in R, the free software environment



for statistical computing and graphics [22]<sup>2</sup>. A small consideration on the convergence speed concludes this section.

### 3.1. Regression Methods and Datasets

The EPRR method is ranked against several well-known learning algorithms for regression, namely: non-regularized EPR, Linear Regression, Support Vector Machines [21], Regression Trees [28] and Conditional Inference Trees [16, 26, 25]. To achieve better error results the SVM and Regression Tree parameters are tuned in each dataset.

The performance of each method is evaluated on several common datasets. From each dataset 70% of the observations are reserved for training purposes and the remaining observations used to estimate the error. To enhance the robustness of results this process is repeated 25 times, each time with a different shuffling of the samples in the train and test sets. Some datasets with attribute values of different magnitudes have a pre-processing scaling transformation. The box plots in figures 2 and 3 resume the test set error distributions over these different runs.

One of the used datasets, ARTIFICIAL, has a special role: it is used to test if EPRR is able to discover a polynomial model. The idea of this test is to generate a polynomial dependent variable and measure the EPRR error after fitting the dataset. The genetic algorithm parameterization for this dataset uses a population with size  $n = 100$  and evolves for 50 generations. For the remaining datasets the population has size  $n = 300$  and evolves for 100 generations.

ARTIFICIAL is a polynomial dataset with four numeric features,  $x_1, \dots, x_4$ , where  $x_1, x_3$  are outcomes from Poisson random variables, and  $x_2, x_4$  from Normal random variables. The dependent variable is given by the polynomial expression  $y = x_2x_4^2 + x_1^2x_3 + 5$ . The dataset includes  $n = 50$  observations;

HOUSING concerns the task of predicting housing values in areas of Boston. There

---

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

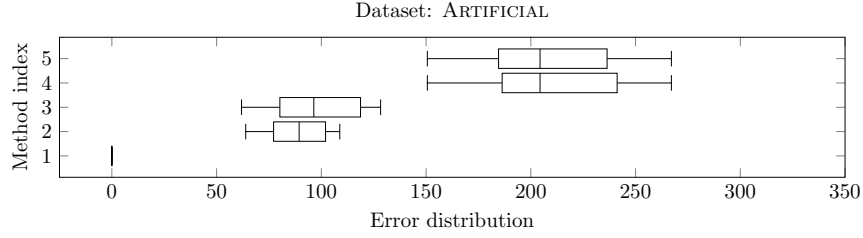


Figure 2: Testing polynomial discovery. The dataset is generated from a polynomial expression and, as shown, EPRR finds the exact generator structure: in line 1, the error box is centered in 0 and has width 0. The regression methods depicted are: 1. EPRR, 2. Linear Regression, 3. SVM, 4. Regression Trees and 5. Conditional Inference Trees

are  $n = 506$  observations of  $m = 13$  continuous attributes and one dependent variable, the median value of owner-occupied homes in thousands of USD;

ABALONE is used to predict the age of a abalone shell using  $m = 8$  numeric attributes concerning several physical measurements. There are  $n = 4177$  observations;

AUTO MPG gathers fuel consumption in miles per gallon, based on two discrete and five continuous attributes ( $m = 7$ ). There are  $n = 398$  observations;

KINEMATICS results from a realistic simulation of the forward kinematics of an 8 link robot arm. The task is to predict the distance of the end-effector from a target using  $m = 8$  continuous attributes. There are  $n = 8192$  observations;

### 3.2. Convergence speed

Since this work is oriented to the error of the EPRR model it is necessary to assess how this depends on the number of generations of the GA. As illustrated in Figure 4, the error quickly drops during the initial 50 to 100 generations. Then, it proceeds slower achieving better solutions only with marginal error reduction.

## 4. Conclusion and Future Work

Of the regression methods considered SVM achieves the best results in three out of four datasets. However SVM and Conditional Inference Trees are pre-trained, having

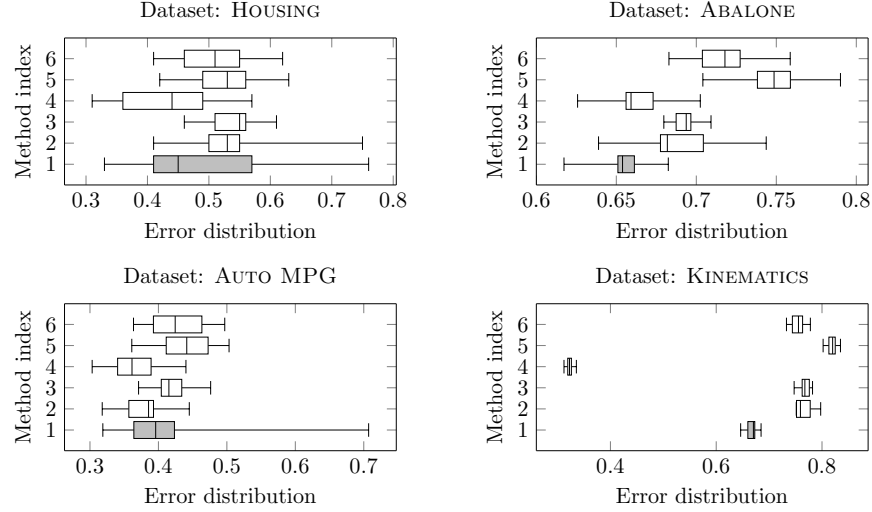


Figure 3: Summary results for different regression methods on diverse datasets. Although EPRR not always achieves the smallest expected error, performance is on-par with more sophisticated methods. The regression methods depicted in these figures are: 1. EPRR; 2. EPR; 3. Linear Regression; 4. SVM; 5. Regression Trees; 6. Conditional Inference Trees;

parameters tuned for each particular dataset unlike EPRR, that runs with the same parameterization on all datasets. Even so it is the best estimator for the ABALONE dataset and in the remaining datasets it outperforms most of the other estimators.

Comparing EPR and EPRR — the main article’s topic — the regularized version achieves much better results at ABALONE and especially KINEMATICS. On the HOUSING dataset errors are improved wrt EPR in a difference in means, resulted in a 95% HDI (Highest Density Interval) equal to  $[0.001, 0.119]$  which, while borderline, achieves statistical significance. Only in the AUTO MPG dataset EPR achieves better results, even if not that different from EPRR.

For complexity considerations EPR and EPRR demand some processing time. On a quad-core computer, processing the KINEMATICS dataset (with near 8K observations) takes approximately 5 minutes. Probably processing time can be reduced by one to two orders in magnitude if the algorithm is implemented with computational speed in mind. However, speed optimization is not the focus of this article.

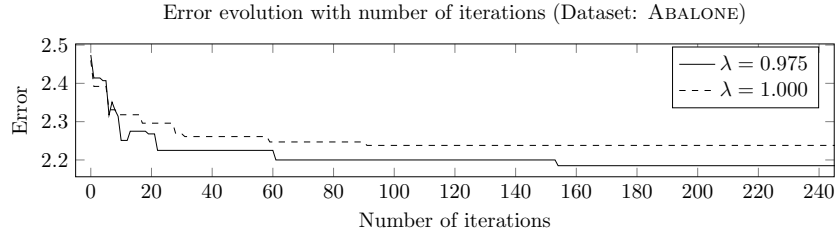


Figure 4: Learning curve: Error progress for the ABALONE dataset during a single execution of the genetic algorithm. The figure shows the fitness evolution for two different regularization values. The population for both consists of 200 polynomials. The error values seem to stabilize around iteration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

A cross-validation procedure can be implemented to refine the appropriate parameter values to achieve better errors. Namely, the regularization parameter,  $\lambda$ , can be tested with several values, instead of being fixed at 0.8. Other parameters like mutation chance or the amount of elitism can also be tested. However, these type of tests need a low-level, fast implementation of EPR and are postponed to future investigation.

## Acknowledgements

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

Datasets used herein are selected from Luís Torgo's data repository, <http://www.dcc.fc.up.pt/~ltorgo/Regression/DataSets.html>. Most can also be found in the UCI ML repository at <http://archive.ics.uci.edu/ml/>.

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

## References

- [1] Yukun Bao, Zhongyi Hu, and Tao Xiong. A pso and pattern search based memetic algorithm for svms parameters optimization. *Neurocomputing*, 117:98–106, 2013.

- 234 [2] Nils Aall Barricelli. Numerical testing of evolution theories. part i: Theoretical  
235 introduction and basic tests. *Acta Biotheoretica*, 16(1-2):69–98, 1962.
- 236 [3] Yoshua Bengio. Learning deep architectures for AI. *Foundations and trends in*  
237 *Machine Learning*, 2(1):1–127, 2009.
- 238 [4] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning:  
239 A review and new perspectives. 2013.
- 240 [5] B Cetisli and H Kalkan. Polynomial curve fitting with varying real powers. *Elec-*  
241 *tronics and Electrical Engineering*, 112(6):117–122, 2011.
- 242 [6] Kit Yan Chan, Tharam S Dillon, and Che Kit Kwong. Modeling of a liquid epoxy  
243 molding process using a particle swarm optimization-based fuzzy regression ap-  
244 proach. *Industrial Informatics, IEEE Transactions on*, 7(1):148–158, 2011.
- 245 [7] Kit Yan Chan, CK Kwong, and Tharam S Dillon. Development of product design  
246 models using fuzzy regression based genetic programming. In *Computational*  
247 *Intelligence Techniques for New Product Design*, pages 111–128. Springer, 2012.
- 248 [8] L Cséfalvayová, M Pelikan, I Kralj Cigić, J Kolar, and M Strlič. Use of genetic  
249 algorithms with multivariate regression for determination of gelatine in historic  
250 papers based on FT-IR and NIR spectral data. *Talanta*, 82(5):1784–1790, 2010.
- 251 [9] J Davidson, D Savic, and G Walters. Method for the identification of explicit  
252 polynomial formulae for the friction in turbulent pipe flow. *Journal of Hydroin-*  
253 *formatics*, 1:115–126, 1999.
- 254 [10] JW Davidson, Dragan A Savic, and Godfrey A Walters. Symbolic and numerical  
255 regression: Experiments and applications. *Information Sciences*, 150(1):95–117,  
256 2003.
- 257 [11] Akemi Gálvez, Andrés Iglesias, and Jaime Puig-Pey. Iterative two-step genetic-  
258 algorithm-based method for efficient polynomial b-spline surface reconstruction.  
259 *Information Sciences*, 182(1):56–76, 2012.

- 260 [12] Paulino José García Nieto, JR Alonso Fernández, FJ de Cos Juez, Fernando  
261 Sánchez Lasheras, and C Díaz Muñoz. Hybrid modelling based on support vector  
262 regression with genetic algorithms in forecasting the cyanotoxins presence in the  
263 trasona reservoir (northern spain). *Environmental research*, 2013.
- 264 [13] O Giustolisi and D Savic. Advances in data-driven analyses and modelling using  
265 epr-moga. *Journal of Hydroinformatics*, 11(3-4):225–236, 2009.
- 266 [14] Magnus Hofwing, Niclas Strömberg, and Martin Tapankov. Optimal polynomial  
267 regression models by using a genetic algorithm. In *Proceedings of the Second  
268 International Conference on Soft Computing Technology in Civil, Structural and  
269 Environmental Engineering Conference, (Crete, Greece), 2011009*, 2011.
- 270 [15] John H Holland. *Adaptation in natural and artificial systems: An introductory  
271 analysis with applications to biology, control, and artificial intelligence*. U Michi-  
272 gan Press, 1975.
- 273 [16] Torsten Hothorn, Kurt Hornik, and Achim Zeileis. Unbiased recursive partition-  
274 ing: A conditional inference framework. *Journal of Computational and Graphi-  
275 cal Statistics*, 15(3):651–674, 2006.
- 276 [17] John R Koza. *Genetic Programming: vol. 1, On the programming of computers  
277 by means of natural selection*, volume 1. MIT press, 1992.
- 278 [18] Zhiyu Liang and Yoonkyung Lee. Eigen-analysis of nonlinear pca with polyno-  
279 mial kernels. 2012.
- 280 [19] Koen Maertens, Josse De Baerdemaeker, and R Babuška. Genetic polynomial  
281 regression as input selection algorithm for non-linear identification. *Soft Com-  
282 puting*, 10(9):785–795, 2006.
- 283 [20] L Mendes and P d Carvalho. Adaptive polynomial regression for colorimetric  
284 scanner calibration using genetic algorithms. In *Intelligent Signal Processing,  
285 2005 IEEE International Workshop on*, pages 22–27. IEEE, 2005.

- 286 [21] David Meyer, Evgenia Dimitriadou, Kurt Hornik, Andreas Weingessel, and  
287 Friedrich Leisch. e1071: Misc functions of the department of statistics (e1071),  
288 tu wien. 2012. R package version 1.6-1.
- 289 [22] R Core Team. R: A language and environment for statistical computing. 2013.
- 290 [23] Luciano Sánchez, José Otero, and Inés Couso. Obtaining linguistic fuzzy rule-  
291 based regression models from imprecise data with multiobjective genetic algo-  
292 rithms. *Soft Computing*, 13(5):467–479, 2009.
- 293 [24] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel prin-  
294 cipal component analysis. In *Artificial Neural Networks ICANN'97*, pages 583–  
295 588. Springer, 1997.
- 296 [25] Carolin Strobl, Anne-Laure Boulesteix, Thomas Kneib, Thomas Augustin, and  
297 Achim Zeileis. Conditional variable importance for random forests. *BMC Bioin-*  
298 *formatics*, 9(307), 2008.
- 299 [26] Carolin Strobl, Anne-Laure Boulesteix, Achim Zeileis, and Torsten Hothorn. Bias  
300 in random forest variable importance measures: Illustrations, sources and a solu-  
301 tion. *BMC Bioinformatics*, 8(25), 2007.
- 302 [27] Daniel Tarlow, Ilya Sutskever, and Richard S Zemel. Stochastic k-neighborhood  
303 selection for supervised and unsupervised learning. *Journal of Machine Learning*  
304 *Research*, 2013.
- 305 [28] Terry Therneau, Beth Atkinson, and Brian Ripley. rpart: Recursive partitioning.  
306 2013. R package version 4.1-1.
- 307 [29] Egon Willighagen. genalg: R based genetic algorithm. 2012.
- 308 [30] Chih-Hung Wu, Gwo-Hshiung Tzeng, and Rong-Ho Lin. A novel hybrid ge-  
309 netic algorithm for kernel function and parameter optimization in support vector  
310 regression. *Expert Systems with Applications*, 36(3):4725–4735, 2009.

- 311 [31] Tian-Li Yu and Wei-Kai Lin. Optimal sampling of genetic algorithms on poly-  
312 nomial regression. In *Proceedings of the 10th annual conference on Genetic and*  
313 *evolutionary computation*, pages 1089–1096. ACM, 2008.