

A method for regularization of evolutionary polynomial regressions

Francisco Coelho*

*Departamento de Informática, Universidade de Évora, Rua Romão Ramalho 58, 7000-671
Évora, Portugal*

João Pedro Neto*

*University of Lisboa, Faculty of Sciences, BioISI– Biosystems & Integrative Sciences
Institute, 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 with a regularization term to control polynomial complexity. The article also 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 needs less computation time than plain EPR.

Keywords: evolutionary polynomial regression, regularization, feature extraction

*Corresponding author: Tel.: +351-919-006-379

Email addresses: `fc@di.uevora.pt` (Francisco Coelho), `jpn@di.fc.ul.pt` (João Pedro Neto)

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 other areas. The kernel trick in Support Vector Machines (SVM) ([1, 2, 3]) alleviates this problem by allowing special nonlinear transformations of the feature-space. The condition such transformations must meet is known as the *kernel trick*, $k(x, x') = \langle \varphi(x), \varphi(x') \rangle$, where φ is the feature-space transformation and $\langle \cdot, \cdot \rangle$ denotes inner product. The “trick” consists on computing the kernel $k(x, x')$ while avoiding the computation of the inner product and the transformations $\varphi(x), \varphi(x')$. A special case of polynomial transformation, the *polynomial kernel*, $k(x, x') = \langle x, x' \rangle^d$ is commonly used in regression and classification tasks with SVMs. However the kernel trick doesn’t apply to general polynomial transformations.

Polynomials, one of the most studied subjects in mathematics, generalize linear functions and define, perhaps, the simplest and most used nonlinear models. For example, Polynomial Neural Networks [4] generalize the linear component of neural networks with a polynomial function. Applications include colorimetric calibration [5], explicit formulæ for turbulent pipe flows [6], computational linguistics [7] and more recently analytical techniques for cultural heritage materials [8], liquid epoxy moulding process [9], B-spline surface reconstruction [10], product design [11] or forecasting cyanotoxins presence in water reservoirs [12]. These examples not only illustrate the wide spectrum of applications but, additionally, each one uses, at some point, Genetic algorithms (GA).

Evolutionary algorithms, including GA, were, arguably, one of the hottest topics of research in the recent decades and with good reason since they outline an optimization scheme easy to conceptualize and with very broad application. If a nonlinear (or otherwise) model requires parameterization, GAs provide a

31 simple and often effective approach to search for global optimal parameters
32 (although no guarantee of global optimality can be given). Related research
33 abound and spans from the 1950s seminal work of Nils Aall Barricelli [13] in
34 the Institute for Advanced Study of Princeton to today’s principal area of study
35 for thousands of researchers, covered in hundreds of conferences, workshops and
36 other meetings. Perhaps the key impulse to GAs came from John Holland’s
37 work and his book “Adaptation in Natural and Artificial Systems” [14].

38 One interesting variation of genetic algorithms, named *genetic programming*
39 by John Koza [15], proposes the use of GAs to search the syntactic structure of
40 complex functions. Syntactic structure search is also keen to the central ideas of
41 deep learning [16, 17], a subarea of machine learning actually producing quite
42 promising results (*e.g.* in [18]). It is also related to the work presented in
43 this paper in the sense that, unlike linear models that have a simple structure,
44 $y = \sum_i \beta_i x_i$, nonlinear (in particular polynomial) models pose an additional
45 structure search problem.

46 The idea of using GAs to find a polynomial regression is not new [19, 20, 21]
47 but still generates original research [22, 23]. The modern formulation of the use
48 of GA to find polynomial models is known as Evolutionary Polynomial Regres-
49 sion (EPR) and systematization can be traced back to the work of Davidson,
50 Savic and Walters [24]. Further developments include multi-objective optimiza-
51 tions [25].

52 Use of regularization/penalty functions is common practice in machine learn-
53 ing in general and has some applications in GA[26]. This paper describes an
54 extension of the general EPR method to find a regularized polynomial regres-
55 sion of a given dataset. Herein optimal regression results from a cost function
56 that accounts for both the root-mean-square (error) together with a novel reg-
57 ularization factor that penalizes over-fitting by polynomial complexity.

58 The next section describes the method’s details and is followed by a presen-
59 tation of some performance results. The last section draws some conclusions
60 and points future research tasks.

61 2. Genetic Algorithms for Polynomials

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

A usual representation of polynomials is through expressions of the form

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

66 where each $q_i = \prod_j x_j^{\alpha_{ij}}$ is a monomial, the exponents $\alpha_{ij} \in \mathbb{N}_0$ are non-
67 negative integers and the coefficients $\theta_i \in \mathbb{R}$ are real valued. For example
68 $p(x_1, x_2, x_3) = 2x_1 + x_2x_3 + \frac{1}{2}x_1^2x_3$ has monomials $q_1 = x_1, q_2 = x_2x_3$ and
69 $q_3 = x_1^2x_3$, exponents $\alpha_{1,1} = 1, \alpha_{2,2} = 1, \alpha_{2,3} = 1, \alpha_{3,1} = 2, \alpha_{3,3} = 1$ and all
70 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 rep-
resentation of the monomials is

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

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

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

- 80 1. For a given set of monomials $\mathcal{Q} = \{q_1, \dots, q_k\}$ find the regression coeffi-
81 cients $\Theta = \{\theta_1, \dots, \theta_k\}$ 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 $\mathcal{Q} = \{q_1, \dots, q_k\}$ a set of k monomial
 86 expressions over X_1, \dots, X_m . Define the hypothesis¹

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

87 and let the error (as “cost”) be

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

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

91 This is a simple linear regression problem obtained by expanding \mathcal{D} with
 92 columns that replicate the monomials in \mathcal{Q} . The resulting dataset, $\mathcal{D} \cup \mathcal{Q}(\mathcal{D})$,
 93 adds the monomial transformations in \mathcal{Q} to the original dataset \mathcal{D} . An alter-
 94 native formulation would just replace \mathcal{D} by $\mathcal{Q}(\mathcal{D})$. It turns out that the first
 95 formulation is a special case of the second (by including the variables in the
 96 monomial set) and has the potential for better error performance because it
 97 uses more features.

98 The second problem is treated in the GA setting: Let \mathcal{D} be a dataset as
 99 above and \mathcal{P} a set of polynomials. For each polynomial $p \in \mathcal{P}$ let \mathcal{Q}_p be the set
 100 of monomials in p (without the coefficients) and define the minimization based
 101 fitness

$$\phi_p = \min_{\Theta} J_{\text{fit}}(\Theta; \mathcal{Q}_p, \mathcal{D}) \quad (3)$$

¹The 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 ϵ 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
         $pop \leftarrow \text{ITERATEGA}(pop)$ 
         $pop \leftarrow \text{SORT}(pop, key = J)$        $\triangleright$  Sort population by regression error
         $err \leftarrow J(\text{FIRST}(pop))$ 
    end while
    return  $\text{FIRST}(pop)$ 
end function

```

102 by solving the first problem. With a fitness of every instance, the GA genetic
 103 operators (usually mutation and crossover) evolve the population \mathcal{P} until a
 104 reasonable approximation of a minimum is found. The properties of GAs and
 105 linear regression entail that Algorithm 1 converges to a polynomial that is a
 106 minimum of the fitness function, encapsulated in the error function J_{fit} .

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

110 2.1. Polynomial Encoding

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

117 Let $\{q_1, \dots, q_k\}$ be a set of monomials over the variables X_1, \dots, X_m . The

118 encoding of that set using d bits per exponent is a binary list such that

- 119 1. the initial segment of k bits defines the active state of each monomial;
- 120 2. the remaining bits are split into k segments of size $m \times d$, each representing
- 121 a monomial;
- 122 3. the bits in each monomial segment are split into m sub-segments of size
- 123 d . The j^{th} sub-segment is the binary representation of the degree of the
- 124 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 $\{x_1^3x_3, x_3^7, x_1x_2\}$ (with $m = 3$) has matrix representation

$$\begin{bmatrix} 3 & 0 & 1 \\ 0 & 0 & 7 \\ 1 & 1 & 0 \\ 6 & 2 & 5 \end{bmatrix} = \begin{bmatrix} 011 & 000 & 001 \\ 000 & 000 & 111 \\ 001 & 001 & 000 \\ 110 & 010 & 101 \end{bmatrix}_{(2)}$$

125 where the right matrix is in binary form using $d = 3$ bits. An example of this

126 set of monomials with the forth monomial, $x_1^6x_2^2x_3^5$, inactive would be

127 1110; 011, 000, 001; 000,000,111; 001,001,000; 110,010,101

128 where, for reading purposes, semicolons separate segments and commas separate

129 variables. The first $k = 4$ bits inform that the first, second and third monomials

130 are active while the fourth is not.

131 While each valid encoding represents a set of monomials the map is not

132 bijective: each set of monomials has multiple encodings, for example by changing

133 d or the order of monomial segments. However, considering the EPR task, this

134 is a minor issue and a bijective map would add computational complexity and

135 negative impact to the algorithm's performance.

136 There is one final remark concerning this encoding method. As it is, the

137 activation segment can become all zeros, representing the empty set of monomi-

138 als. This situation can be avoided with a simple hack: Given an encoding, the

139 first monomial is always considered active, thus restricting the syntactic form
 140 of encodings to binary strings starting with 1. In practice, this means that the
 141 implementation of the encoding can omit the first bit.

142 2.2. Cost Function

143 The polynomial regression error considered so far accounts for the ability
 144 to predict the transformed testset. A known problem of using a cost function
 145 based only in the dataset error (and of polynomial regressions in general) is the
 146 tendency to overfit training data. Excessive variance of the estimation method
 147 can be reduced by regularizing the error function with a penalty factor. Thus,
 148 to reduce polynomial complexity and variance by regularizing the size of the
 149 monomial set the error function from equation 2 is multiplied by a factor λ^k

$$J_{\text{reg}}(\Theta, \lambda; \mathcal{Q}, \mathcal{D}) = \lambda^k J_{\text{fit}}(\Theta; \mathcal{Q}, \mathcal{D}) \quad (4)$$

150 where k is the number of monomials in the polynomial. When $\lambda > 1$ polyno-
 151 mials with more monomials are penalized. The regularized extension of EPR is
 152 denoted by Evolutionary Polynomial Regression with Regularization (EPRR).

153 A simple exploration on the effect of the regularization parameter is depicted
 154 in Figure 1 where it is possible to observe that penalizing polynomial complexity
 155 (i.e. $\lambda < 1$) achieves better results than the opposite ($\lambda > 1$). This observation
 156 motivates further inquire, done in section 3 where we compare the regularized
 157 versions with $\lambda < 1$ and the non-regularized algorithm against several other
 158 regression methods.

159 2.3. Genetic Algorithm Parameterization

160 In general GAs offer many possibilities with respect to the choice of genetic
 161 operators and respective application rates, population evolution, *etc.* The re-
 162 sults found here were obtained using the package `genalg` [27] with standard
 163 operators (crossover and mutation) and population evolution defined by muta-
 164 tion rate of 5% and 20% elitism between generations.

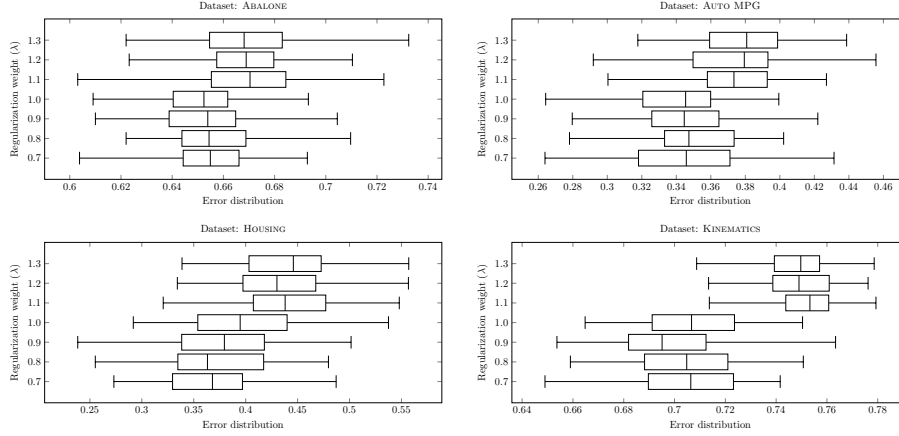


Figure 1: Error distribution by regularization exponent for common datasets. The box plots summarize error values of 75 simulations (60 iterations on a population size of 50) for each value of λ . Performance of the non-regularized EPR is plotted in the line $\lambda = 1$. Performance for $\lambda < 1$ is better than $\lambda > 1$ in all datasets.

3. Experimental Results

Here is described the experiment setup used to gather and summarize the empirical evidence that supports this comparative study of EPR and EPRR. Evaluation is focused in error distribution and, besides EPR and EPRR, also uses several common regression methods and datasets easily accessible in R, the free software environment for statistical computing and graphics [28]. Datasets and code used to produce the results and plots in this paper are available at <https://github.com/jpneto/GenAlgPoly>. 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 [29] with linear kernel, Regression Trees [30] and Conditional Inference Trees [31, 32, 33]. No ensemble regression method was used in this comparison since these methods average between several models while EPRR and the

dataset	method	quantile 25%	error mean	quantile 75%
Abalone	EPRR $\lambda = 0.7$	0.6392	0.6555	0.6677
	EPRR $\lambda = 0.8$	0.6408	0.6543	0.6636
	EPRR $\lambda = 0.9$	0.6481	0.6581	0.6707
	EPRR $\lambda = 1.0$	0.6542	0.6715	0.6816
	Linear Regression	0.6803	0.6927	0.7078
	SVM (linear kernel)	0.6916	0.7044	0.7205
	Regression Trees	0.7423	0.7520	0.7621
	Cond. Inference Trees	0.7031	0.7126	0.7264
Auto-Mpg	EPRR $\lambda = 0.7$	0.3635	0.3916	0.4147
	EPRR $\lambda = 0.8$	0.3629	0.3956	0.4228
	EPRR $\lambda = 0.9$	0.3646	0.4130	0.4215
	EPRR $\lambda = 1.0$	0.3691	0.3999	0.4057
	Linear Regression	0.4071	0.4284	0.4473
	SVM (linear kernel)	0.4116	0.4358	0.4613
	Regression Trees	0.4216	0.4501	0.4785
	Cond. Inference Trees	0.4063	0.4372	0.4663
Housing	EPRR $\lambda = 0.7$	0.4412	0.6650	0.5739
	EPRR $\lambda = 0.8$	0.4241	0.5274	0.6016
	EPRR $\lambda = 0.9$	0.4354	0.5717	0.6462
	EPRR $\lambda = 1.0$	0.4477	0.5417	0.5995
	Linear Regression	0.4898	0.5313	0.5649
	SVM (linear kernel)	0.4831	0.5469	0.6017
	Regression Trees	0.4845	0.5232	0.5720
	Cond. Inference Trees	0.4676	0.5080	0.5413
Kinematics	EPRR $\lambda = 0.7$	0.6600	0.6660	0.6720
	EPRR $\lambda = 0.8$	0.6617	0.6694	0.6751
	EPRR $\lambda = 0.9$	0.6636	0.6714	0.6761
	EPRR $\lambda = 1.0$	0.7568	0.7558	0.7739
	Linear Regression	0.7672	0.7759	0.7849
	SVM (linear kernel)	0.8074	0.8136	0.8247
	Regression Trees	0.5673	0.6021	0.5803
	Cond. Inference Trees	0.7558	0.7620	0.7686

Table 1: Tabular summary results for different regression methods on common datasets.

selected methods produce a single model.

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 75 times, each time with a different shuffling of the samples in the train and test sets. The datasets have been scaled to prevent problems with different magnitude attributes. 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

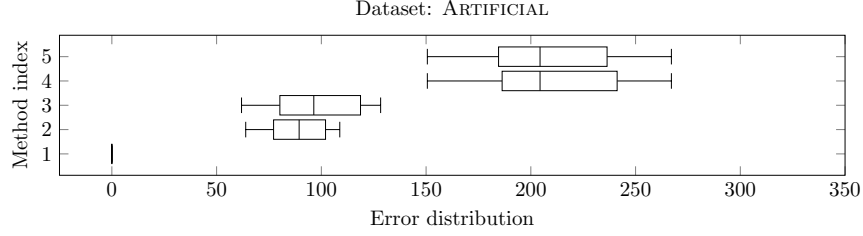


Figure 2: Testing polynomial discovery. The dataset is generated from a polynomial expression and, as shown, EPRR finds the exact generator structure: in line 1, the error box is centered in 0 and has width 0. The regression methods depicted are: 1. EPRR, 2. Linear Regression, 3. SVM, 4. Regression Trees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

193 a population with size $n = 100$ and evolves for 50 generations. For the remaining
 194 datasets the population has size $n = 300$ and evolves for 100 generations.

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

200 HOUSING concerns the task of predicting housing values in areas of Boston.
 201 There are $n = 506$ observations of $m = 13$ continuous attributes and
 202 one dependent variable, the median value of owner-occupied homes in
 203 thousands of USD;

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

207 AUTO MPG gathers fuel consumption in miles per gallon, based on two dis-
 208 crete and five continuous attributes ($m = 7$). There are $n = 398$ observa-
 209 tions;

210 KINEMATICS results from a realistic simulation of the forward kinematics of an

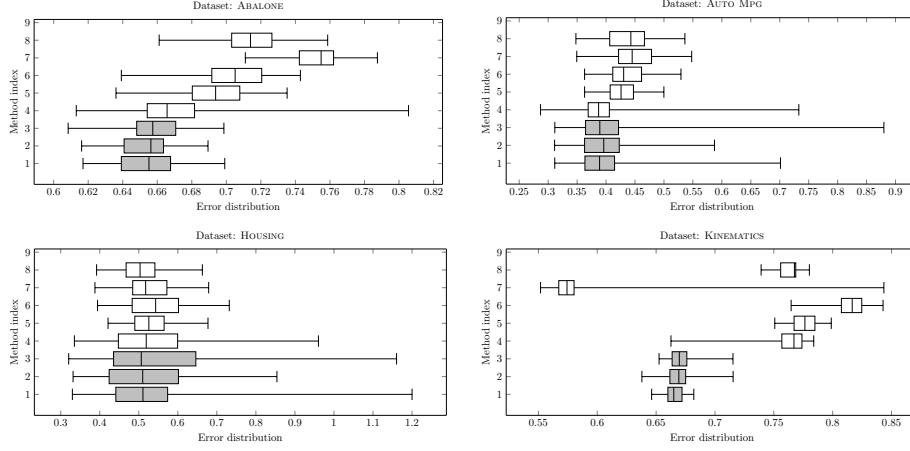


Figure 3: Graphical summary results for different regression methods on common datasets. The regression methods depicted in these figures are: 1. EPRR, $\lambda = 0.7$; 2. EPRR, $\lambda = 0.8$; 3. EPRR, $\lambda = 0.9$; 4. EPR (i.e. $\lambda = 1.0$); 5. Linear Regression; 6. SVM (linear kernel); 7. Regression Trees; 8. Conditional Inference Trees.

211 8 link robot arm. The task is to predict the distance of the end-effector
 212 from a target using $m = 8$ continuous attributes. There are $n = 8192$
 213 observations;

214 3.2. Convergence speed

215 Since this work is oriented to the EPRR model error it is necessary to assess
 216 how this depends on the number of generations of the GA. As illustrated in
 217 Figure 4, the error quickly drops during the initial 100 generations. Then, it
 218 proceeds slower achieving better solutions only with marginal error reduction.
 219 The EPR takes more than 300 generations to achieve similar error results.

220 4. Conclusion and Future Work

221 Of the regression methods considered EPRR achieves the best results in three
 222 out of four datasets, the exception being Regression Trees in the Kinematics
 223 dataset.

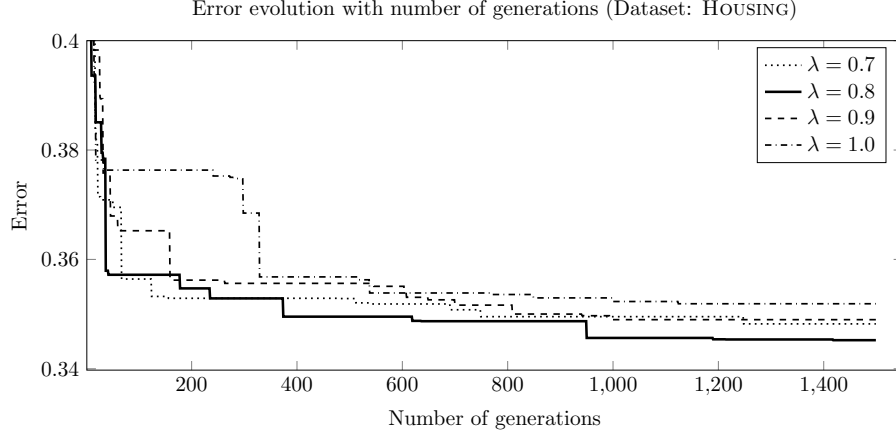


Figure 4: Learning curve: Error progress for the HOUSING dataset during a single execution of the genetic algorithm. The figure shows the fitness evolution for different regularization values. The population of each run consists of 200 polynomials.

224 Comparing EPR and EPRR — the main article’s topic — the regularized
 225 version achieves better results at ABALONE and especially KINEMATICS. A
 226 Bayesian Estimation was computed, which estimates the difference in means
 227 between two groups and yields a probability distribution over the difference.
 228 On the Abalone dataset errors are improved wrt EPR in a difference in means,
 229 for all tested λ values less than 1, resulted in a 95% HDI (Highest Density
 230 Interval) equal to $[0.01, 0.02]$ which achieves statistical significance. On the
 231 Auto-MPG dataset the 95% HDI is $[-0.01, 0.008]$ which includes zero. On the
 232 Housing dataset the 95% HDI is $[-0.02, 0.05]$, for $\lambda = 0.7$, which includes zero
 233 (the other values of λ had similar statistics).

234 For complexity considerations EPR and EPRR demand some processing
 235 time. On a quad-core computer, processing the KINEMATICS dataset (with near
 236 8K observations) takes approximately 5 minutes. Probably processing time can
 237 be reduced by one to two orders in magnitude if the algorithm is implemented
 238 with computational speed in mind. However, speed optimization is not the focus
 239 of this article.

240 Parameters like mutation chance or the amount of elitism can be tested for

241 tuning purposes. However, these type of tests need a low-level, fast implemen-
242 tation of EPR and are postponed to future investigation.

243 Acknowledgements

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

247 João Pedro Neto's work is supported by centre grant (to BioISI, Centre
248 Reference: UID/MULTI/04046/2013), from FCT/MCTES/PIDDAC, Portugal.

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

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

- 254 [1] B. Schölkopf, A. Smola, K.-R. Müller, Kernel principal component analysis,
255 in: Artificial Neural Networks ICANN'97, Springer, 1997, pp. 583–588.
- 256 [2] Z. Liang, Y. Lee, Eigen-analysis of nonlinear pca with polynomial kernels.
- 257 [3] Y. Bao, Z. Hu, T. Xiong, A pso and pattern search based memetic algorithm
258 for svms parameters optimization, *Neurocomputing* 117 (2013) 98–106.
- 259 [4] S. Dehuri, B. B. Misra, A. Ghosh, S.-B. Cho, A condensed polynomial neu-
260 ral network for classification using swarm intelligence, *Applied Soft Com-*
261 *puting* 11 (3) (2011) 3106–3113.
- 262 [5] L. Mendes, P. d. Carvalho, Adaptive polynomial regression for colorimetric
263 scanner calibration using genetic algorithms, in: *Intelligent Signal Process-*
264 *ing*, 2005 IEEE International Workshop on, IEEE, 2005, pp. 22–27.
- 265 [6] J. Davidson, D. Savic, G. Walters, Method for the identification of explicit
266 polynomial formulae for the friction in turbulent pipe flow., *Journal of*
267 *Hydroinformatics* 1 (1999) 115–126.

- [7] L. Sánchez, J. Otero, I. Couso, Obtaining linguistic fuzzy rule-based regression models from imprecise data with multiobjective genetic algorithms, *Soft Computing* 13 (5) (2009) 467–479.
- [8] L. Cséfalvayová, M. Pelikan, I. Kralj Cigić, J. Kolar, M. Strlič, Use of genetic algorithms with multivariate regression for determination of gelatine in historic papers based on FT-IR and NIR spectral data, *Talanta* 82 (5) (2010) 1784–1790.
- [9] K. Y. Chan, T. S. Dillon, C. K. Kwong, Modeling of a liquid epoxy molding process using a particle swarm optimization-based fuzzy regression approach, *Industrial Informatics, IEEE Transactions on* 7 (1) (2011) 148–158.
- [10] A. Gálvez, A. Iglesias, J. Puig-Pey, Iterative two-step genetic-algorithm-based method for efficient polynomial b-spline surface reconstruction, *Information Sciences* 182 (1) (2012) 56–76.
- [11] K. Y. Chan, C. Kwong, T. S. Dillon, Development of product design models using fuzzy regression based genetic programming, in: *Computational Intelligence Techniques for New Product Design*, Springer, 2012, pp. 111–128.
- [12] P. J. García Nieto, J. Alonso Fernández, F. de Cos Juez, F. Sánchez Lasheras, C. Díaz Muñoz, Hybrid modelling based on support vector regression with genetic algorithms in forecasting the cyanotoxins presence in the trasona reservoir (northern spain), *Environmental research*.
- [13] N. A. Barricelli, Numerical testing of evolution theories. part i: Theoretical introduction and basic tests, *Acta Biotheoretica* 16 (1-2) (1962) 69–98.
- [14] J. H. Holland, *Adaptation in natural and artificial systems: An introductory analysis with applications to biology, control, and artificial intelligence.*, U Michigan Press, 1975.

- 295 [15] J. R. Koza, Genetic Programming: vol. 1, On the programming of comput-
296 ers by means of natural selection, Vol. 1, MIT press, 1992.
- 297 [16] Y. Bengio, Learning deep architectures for AI, Foundations and trends in
298 Machine Learning 2 (1) (2009) 1–127.
- 299 [17] Y. Bengio, A. Courville, P. Vincent, Representation learning: A review and
300 new perspectives.
- 301 [18] D. Tarlow, I. Sutskever, R. S. Zemel, Stochastic k-neighborhood selection
302 for supervised and unsupervised learning, Journal of Machine Learning
303 Research.
- 304 [19] K. Maertens, J. De Baerdemaeker, R. Babuška, Genetic polynomial regres-
305 sion as input selection algorithm for non-linear identification, Soft Com-
306 puting 10 (9) (2006) 785–795.
- 307 [20] T.-L. Yu, W.-K. Lin, Optimal sampling of genetic algorithms on polynomial
308 regression, in: Proceedings of the 10th annual conference on Genetic and
309 evolutionary computation, ACM, 2008, pp. 1089–1096.
- 310 [21] C.-H. Wu, G.-H. Tzeng, R.-H. Lin, A novel hybrid genetic algorithm for
311 kernel function and parameter optimization in support vector regression,
312 Expert Systems with Applications 36 (3) (2009) 4725–4735.
- 313 [22] M. Hofwing, N. Strömberg, M. Tapankov, Optimal polynomial regression
314 models by using a genetic algorithm, in: Proceedings of the Second Inter-
315 national Conference on Soft Computing Technology in Civil, Structural and
316 Environmental Engineering Conference, (Crete, Greece), 2011009, 2011.
- 317 [23] B. Cetisli, H. Kalkan, Polynomial curve fitting with varying real powers,
318 Electronics and Electrical Engineering 112 (6) (2011) 117–122.
- 319 [24] J. Davidson, D. A. Savic, G. A. Walters, Symbolic and numerical regression:
320 Experiments and applications, Information Sciences 150 (1) (2003) 95–117.

- [25] O. Giustolisi, D. Savic, Advances in data-driven analyses and modelling using epr-moga, *Journal of Hydroinformatics* 11 (3-4) (2009) 225–236.
- [26] R. Gupta, A. Bhunia, D. Roy, A GA based penalty function technique for solving constrained redundancy allocation problem of series system with interval valued reliability of components, *Journal of Computational and Applied Mathematics* 232 (2) (2009) 275 – 284. doi:<http://dx.doi.org/10.1016/j.cam.2009.06.008>.
URL <http://www.sciencedirect.com/science/article/pii/S0377042709003549>
- [27] E. Willighagen, *genalg: R based genetic algorithm* (2012).
- [28] R Core Team, *R: A language and environment for statistical computing*.
URL <http://www.R-project.org/>
- [29] D. Meyer, E. Dimitriadou, K. Hornik, A. Weingessel, F. Leisch, *e1071: Misc functions of the department of statistics (e1071), tu wienR package version 1.6-1*.
URL <http://CRAN.R-project.org/package=e1071>
- [30] T. Therneau, B. Atkinson, B. Ripley, *rpart: Recursive partitioningR package version 4.1-1*.
- [31] T. Hothorn, K. Hornik, A. Zeileis, Unbiased recursive partitioning: A conditional inference framework, *Journal of Computational and Graphical Statistics* 15 (3) (2006) 651–674.
- [32] C. Strobl, A.-L. Boulesteix, T. Kneib, T. Augustin, A. Zeileis, Conditional variable importance for random forests, *BMC Bioinformatics* 9 (307).
URL <http://www.biomedcentral.com/1471-2105/9/307>
- [33] C. Strobl, A.-L. Boulesteix, A. Zeileis, T. Hothorn, Bias in random forest variable importance measures: Illustrations, sources and a solution, *BMC Bioinformatics* 8 (25).
URL <http://www.biomedcentral.com/1471-2105/8/25>