

Nonlinear regression in environmental sciences by support vector machines combined with evolutionary strategy

Aranildo R. Lima^{*a}, Alex J. Cannon^a, William W. Hsieh^a

^a*Department of Earth and Ocean Sciences, University of British Columbia, Vancouver, BC V6T 1Z4, Canada*

Abstract

A hybrid algorithm combining support vector regression with evolutionary strategy (SVR-ES) is proposed for predictive models in the environmental sciences. SVR-ES uses uncorrelated mutation with p step sizes to find the optimal SVR hyper-parameters. Three environmental forecast datasets used in the WCCI-2006 contest – surface air temperature, precipitation and sulphur dioxide concentration – were tested. We used multiple linear regression (MLR) as benchmark and a variety of machine learning techniques including bootstrap-aggregated ensemble artificial neural network (ANN), SVR-ES, SVR with hyper-parameters given by the Cherkassky-Ma estimate, the M5 regression tree, and random forest (RF). We also tested all techniques using stepwise linear regression (SLR) first to screen out irrelevant predictors. We concluded that SVR-ES is an attractive approach because it tends to outperform the other techniques and can also be implemented in an almost automatic way. The Cherkassky-Ma estimate is a useful approach for minimizing the mean absolute error and saving computational time related to the hyper-parameter search. The ANN and RF are also good options to outperform multiple linear regression (MLR). Finally, the use of SLR for predictor selection can dramatically reduce computational time and often help to enhance accuracy.

Keywords: Support Vector Machine, Evolutionary Strategy, Artificial Neural Network, Atmospheric Forecasts, Hybrid Systems, Machine Learning

^{*}Corresponding author. Department of Earth and Ocean Sciences, University of British Columbia, 6339 Stores Road, Vancouver, BC V6T 1Z4, Canada, Tel.: +1-604-822-5691; Fax.: +1-604-822-6088, E-mail: arodrigu@eos.ubc.ca

1. Introduction

The main idea of machine learning is that computer algorithms are capable of automatically distilling knowledge from data. From this knowledge they can construct models capable of making predictions from novel data in the future. However, environmental modeling problems are typically very noisy (Cawley et al., 2007), hence it is not easy to build successful predictive models.

Due to data modeling complexity, even for a particular machine learning method, there is often more than one way to build the model (Zhang, 2007). To build a successful predictive model, the correct adjustment of the model hyper-parameters is necessary. For example, in artificial neural network (ANN) models, the number of hidden processing units, the choice of activation functions and the regularization parameter all need to be specified (Haykin, 1998; Hsieh, 2009).

Similarly, in support vector machines (SVM) for regression (SVR), typically two or three hyper-parameters have to be tuned, such as the cost of constraint violation (C), the insensitive-loss (ϵ) and, if the Gaussian function is used as the kernel function, the width of the Gaussian (γ). In theory, the establishment of these hyper-parameters requires an optimal search in full state space.

Much effort has been spent on improving the efficiency of the SVR hyper-parameter search (Cherkassky and Ma, 2004; Friedrichs and Igel, 2005; Fan et al., 2005; Huang and Wang, 2006; Lin et al., 2008). The most common approach is the simple grid search (Fan et al., 2005; Ortiz-García et al., 2009; Pino-Mejías et al., 2010; Zeng et al., 2011), an exhaustive approach that is computationally expensive. Furthermore, the hyper-parameters are varied by fixed step-sizes through a wide range of values, limiting the search to discrete values.

Some of the most promising approaches in tuning the SVM hyper-parameters are based on evolutionary algorithms (EA), e.g. genetic algorithms (GA) (Pai and Hong, 2005; Huang and Wang, 2006; Tripathi et al., 2006), particle swarm optimization (PSO) (Lin et al., 2008), and evolutionary strategies (ES) (Friedrichs and Igel, 2005). However, these approaches need a set of adjustable parameters called EA parameters (Smit and Eiben, 2009) and typically there are more EA parameters to adjust than the number of hyper-parameters. For example, GA has population size, mutation rate, crossover

32 rate, number of generations, etc., and PSO has population size, acceleration coefficients,
33 inertia weight, number of generations, etc. Kramer et al. (2007) showed that the choice
34 of the EA parameters has a decisive impact in the final results and can undervalue the
35 algorithm performance. In other words, to use a GA it is necessary to adjust at least four
36 parameters, which is more than the three SVR hyper-parameters when using a Gaussian
37 kernel.

38 Another important issue is the data quality upon which machine learning meth-
39 ods operate. Indeed, machine learning algorithms may produce less accurate and less
40 understandable results if the data are inadequate or contain extraneous and irrelevant
41 information (Hall and Smith, 1996). It is also possible to make the modeling process less
42 time consuming and sometimes more accurate by removing predictors that are irrelevant
43 or redundant with respect to the task to be learned.

44 In this paper, our main goals are: (i) reduce the number of hyper-parameters which
45 require estimation; (ii) use an accurate initialization of the hyper-parameter search; and
46 (iii) discard irrelevant and redundant predictors. We propose a hybrid algorithm called
47 SVR-ES which uses a simple evolutionary strategy called “uncorrelated mutation with p
48 step sizes” (Eiben and Smith, 2003) to find the optimal SVR hyper-parameters. We also
49 combine the SVR-ES with stepwise linear regression (SLR) (Draper and Smith, 1998) to
50 screen out irrelevant predictors.

51 Three environmental forecast problems used in the WCCI-2006 contest – surface air
52 temperature (TEMP), precipitation (PRECIP) and sulphur dioxide concentration (SO2)
53 – are tested (Cawley et al., 2007). These three datasets contain different amounts of non-
54 linearity and noise. Several other machine learning techniques successfully used in the
55 environmental forecast problems are considered, including bootstrap-aggregated ensem-
56 ble ANN (Cannon and Lord, 2000; Krasnopolsky, 2007), SVR using the Cherkassky-Ma
57 hyper-parameter estimates (Cherkassky and Ma, 2004), the M5 regression tree (Quinlan,
58 1992; Solomatine and Xue, 2004; Haupt et al., 2009) and random forest (RF) (Breiman,
59 2001; Pino-Mejías et al., 2010). We also use SLR with these techniques to prescreen and
60 reduce the number of predictors.

61 Section 2 describes the data sets used in our study. The forecasting methods are
62 presented in Sections 3 and 4. Results and discussion of the experiments are given in

63 Section 5, followed by summary and conclusion in Section 6.

64 **2. Data Description**

65 Environmental data normally contain properties that are difficult to model by com-
66 mon regression techniques. For example, response variables may be strictly non-negative,
67 highly skewed, non-Gaussian distributed and heteroscedastic (Cawley et al., 2007). Some
68 examples are the modeling of SO₂ air pollution (Kurt et al., 2008) or statistical down-
69 scaling of temperature and precipitation (Schoof and Pryor, 2001).

70 We used as benchmark three datasets which were originally used at the WCCI-2006
71 Predictive Uncertainty in Environmental Modeling Challenge. These benchmarks, char-
72 acterized by a non-Gaussian, heteroscedastic variance structure (Cawley et al., 2007),
73 are freely available from the challenge website ([http://theoval.cmp.uea.ac.uk/~gcc/
74 competition/](http://theoval.cmp.uea.ac.uk/~gcc/competition/)).

75 Predicting precipitation accurately is still one of the most difficult tasks in meteorol-
76 ogy (Fritsch et al., 1998; Kuligowski, 1998; Strangeways, 2007). Factors responsible for
77 the difficulty in predicting precipitation are e.g. the chaotic nature of the atmosphere and
78 the complexity of the processes involved in its creation (Fritsch et al., 1998), seasonal vari-
79 ations (Wallace and Hobbs, 2006), non-stationary statistical behavior (Von Storch and
80 Zwiers, 2001), difficulties in precipitation measurements including problems with rain
81 gauges, radar and satellites (Strangeways, 2007) and the limited temporal and spatial
82 scales of global circulation models (GCMs) (Kuligowski, 1998).

83 To forecast precipitation at regional scale, GCM outputs cannot be used directly due
84 to coarse spatial resolution. For example, GCMs do not provide information on the
85 spatial structure of temperature and precipitation in areas of complex topography and
86 land use distribution. To use the output of a GCM for this task, statistical downscaling
87 models are often used (Hashmi et al., 2009).

88 In the PRECIP benchmark the input variables are large-scale circulation information,
89 such as might be obtained from a GCM, and the target is daily precipitation data recorded
90 at Newton Rigg, a relatively wet station located in the northwest of the United Kingdom.
91 This benchmark is an example where the response variable is skewed to the right. The
92 PRECIP dataset contains 106 predictors and 10,546 daily patterns.

93 With similar large-scale circulation features to those used for the PRECIP benchmark,
94 the TEMP benchmark is another statistical downscaling problem, in this case one in
95 which the target is the more normally distributed daily maximum temperature at the
96 Writtle station in the southeast of the United Kingdom. The TEMP dataset has 106
97 predictors and 10,675 daily patterns.

98 Air pollution is a major environmental problem in cities. Local emissions and topo-
99 graphic factors, allied with meteorological conditions such as high atmospheric pressure
100 and temperature inversions, cause poor dispersion of atmospheric pollutants due the
101 stagnant conditions. Human health problems (both short term and chronic problems)
102 can occur when pollutant concentrations exceed the allowable limits. Therefore predict-
103 ing the concentration of pollutants such as sulfur dioxide (SO₂) is crucial to providing
104 proper actions and control strategies in extreme situations (Kurt et al., 2008).

105 The target of the SO₂ dataset is the SO₂ concentration in urban Belfast. In order to
106 forecast the SO₂ concentration twenty-four hours in advance, meteorological conditions
107 and current SO₂ levels were used as input variables. Similar to PRECIP, SO₂ is a right
108 skewed variable with a heavy tail. The SO₂ dataset has 27 predictors and 22,956 hourly
109 patterns, more than double the number of patterns in PRECIP and TEMP.

110 Irrelevant predictors can unnecessarily increase the time needed for learning a suf-
111 ficiently accurate forecast model and in many cases also increase the size of the search
112 space. To reduce the number of predictor variables, we use the well-known prescreening
113 technique of SLR (Hocking, 1976; Venables and Ripley, 2002; Hastie et al., 2009).

114 3. Support Vector Regression

115 Support vector machines (SVM) have been widely used in the environmental sciences
116 for classification and regression problems (Pino-Mejías et al., 2010; Tripathi et al., 2006).
117 SVM were originally designed for nonlinear classification problems, then extended to
118 nonlinear regression problems (SVR) (Muller et al., 1997).

119 Suppose we are given the training data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ with n patterns. After
120 mapping the input pattern \mathbf{x} into a higher dimensional feature space using a nonlinear
121 mapping function Φ , the nonlinear regression problem between \mathbf{x} and y can be converted

122 to a linear regression problem between $\Phi(\mathbf{x})$ and y , i.e.

$$f(\mathbf{x}; \mathbf{w}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b, \quad (1)$$

123 where $\langle \cdot, \cdot \rangle$ denotes the inner product, and \mathbf{w} and b are the regression coefficients obtained
 124 by minimizing the error between f and the observed values of y . Instead of the commonly
 125 used mean squared error (MSE) norm, SVR uses the ϵ -insensitive error norm to measure
 126 the error between f and y ,

$$|f(\mathbf{x}; \mathbf{w}) - y|_\epsilon = \begin{cases} 0, & \text{if } |f(\mathbf{x}; \mathbf{w}) - y| < \epsilon \\ |f(\mathbf{x}; \mathbf{w}) - y| - \epsilon, & \text{otherwise,} \end{cases} \quad (2)$$

127 i.e., small errors ($|f - y| < \epsilon$) are ignored, whereas for large errors, the error norm
 128 approximates the mean absolute error (MAE). A key issue is that an error norm based
 129 on MAE is more robust to outliers in the data than the MSE. Using pattern data $(\mathbf{x}_i,$
 130 $y_i)$, the \mathbf{w} and b coefficients are estimated by minimizing the objective function:

$$J = \frac{C}{n} \sum_{i=1}^n |f(\mathbf{x}_i, \mathbf{w}) - y_i|_\epsilon + \frac{1}{2} \|\mathbf{w}\|^2, \quad (3)$$

131 where C (which controls the regularization) and ϵ are the hyper-parameters.

132 The global minimum solution to the linear regression problem (1) can be achieved
 133 without iterative nonlinear optimization, hence local minima in the objective function
 134 are not a problem. However, the linear regression problem can be very expensive or im-
 135 possible to compute without truncating infinite dimensional vectors. This occurs because
 136 $\Phi(\mathbf{x})$ may be a very high (or even infinite) dimensional vector. To counter this drawback
 137 a kernel trick is used in which the inner product $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$ in the solution algorithm
 138 is replaced by a kernel function $K(\mathbf{x}, \mathbf{x}')$, which does not involve the difficult handling
 139 of $\Phi(\mathbf{x})$. The minimization of (3) uses the method of Lagrange multipliers, and the final
 140 regression estimate can be expressed in the form (Bishop, 2006):

$$f(\mathbf{x}) = \sum_i K(\mathbf{x}, \mathbf{x}_i) + b, \quad (4)$$

141 where the summation is only over a subset of the given data \mathbf{x}_i called the support vectors.

142 3.1. Evolutionary Strategies

143 Evolutionary algorithms (EA) mimic nature's way of evolving successful organisms
 144 (individuals) (Haupt et al., 2009). An Evolutionary Strategy (ES) is a particular class

145 of EA that is normally used for continuous parameter optimization. An attractive point
 146 of the ES is the self-adaptation of some strategy parameters. The strategy parameters
 147 are the set of adjustable parameters used by the algorithm. Self-adaptivity means that
 148 some EA parameters are varied during a run in a specific manner: the parameters are
 149 included in the chromosomes and co-evolve with the solutions (Eiben and Smith, 2003),
 150 i.e., the algorithm is capable of adapting itself autonomously.

151 *3.1.1. Mutation and Self-adaptation*

152 Mutation is the name given to a genetic operation which uses only one parent and
 153 creates one offspring by applying some type of randomized change to the representation.
 154 Let \mathbf{G} be a chromosome defined by,

$$\mathbf{G} = (g_1, g_2, \dots, g_p), \quad (5)$$

155 where g_i ($i = 1, 2, \dots, p$) are the solution parameters. Mutations are realized by adding
 156 some Δg_i to each g_i , where Δg_i are values from a Gaussian distribution $N(0, \sigma)$, with
 157 zero mean and standard deviation σ , i.e.,

$$g'_i = g_i + N(0, \sigma), \quad i = 1, 2, \dots, p. \quad (6)$$

158 The self-adaptation consists of including the step size σ in the chromosomes so it also
 159 undergoes variation and selection, i.e. mutations are realized by replacing $(g_1, g_2, \dots, g_p; \sigma)$
 160 by $(g'_1, g'_2, \dots, g'_p; \sigma')$, where σ' is the mutated value of σ , also called the mutation step
 161 size.

162 *3.1.2. Uncorrelated Mutation with p Step Sizes*

163 An attractive feature of the “mutation with p step sizes” method is its ability to treat
 164 each dimension differently, i.e., using different step sizes for different dimensions. The
 165 chromosome given in (5) is extended to p step sizes, resulting in

$$\mathbf{G} = (g_1, g_2, \dots, g_p; \sigma_1, \sigma_2, \dots, \sigma_p), \quad (7)$$

166 and the mutation rules are given by:

$$\sigma'_i = \sigma_i e^{\tau' N(0,1) + \tau N_i(0,1)}, \quad (8)$$

167

$$g'_i = g_i + \sigma_i N_i(0, 1), \quad i = 1, 2, \dots, p, \quad (9)$$

168 where $\tau' \propto 1/\sqrt{2p}$, and $\tau \propto 1/\sqrt{2\sqrt{p}}$. The common base mutation $e^{\tau'N(0,1)}$ allows an
 169 overall change of the mutability. The flexibility to use different mutation strategies in
 170 different directions is provided by the coordinate-specific $e^{\tau N_i(0,1)}$ (Eiben and Smith,
 171 2003).

172 3.2. Hyper-parameter Optimization and SVR-ES

173 The performance of an SVR model is highly dependent on the choice of the kernel
 174 function and the hyper-parameters (Hastie et al., 2009; Hsieh, 2009). The use of a
 175 suitable nonlinear kernel function in SVR allows it to be fully nonlinear, while the use
 176 of a linear kernel function restricts SVR to a linear model. The SVR with the linear
 177 kernel is a linear regression model but with the robust ϵ -insensitive error norm instead
 178 of the non-robust MSE norm as used in multiple linear regression (MLR). In this study,
 179 we used the radial basis function (RBF) kernel (also called the Gaussian kernel) given by
 180 $K(\mathbf{x}, \mathbf{x}_i) = \exp[-\|\mathbf{x} - \mathbf{x}_i\|^2/(2\sigma_K^2)]$, with the hyper-parameter $\gamma = 1/(2\sigma_K^2)$ controlling
 181 the width σ_K of the Gaussian function.

182 Cherkassky and Ma (2004) proposed how to estimate values for the hyper-parameters
 183 C and ϵ in SVR. The estimated regularization parameter C is $\max(|\bar{y} + 3\sigma_y|, |\bar{y} - 3\sigma_y|)$
 184 where \bar{y} and σ_y are the mean and the standard deviation of the y values of the training
 185 data; and $\epsilon \sim 3 \sigma_N \sqrt{\ln n/n}$, where n is the size of data set and σ_N is the standard
 186 deviation of the noise. The noise level is estimated by using the k -nearest-neighbors
 187 regression method, where \hat{y} , the value of a new point, is estimated from the average of
 188 the k nearest points (Hastie et al., 2009). Hence the estimated noise is:

$$\hat{\sigma}_N^2 = \frac{n^{1/5}k}{n(n^{1/5}k - 1)} \sum_{m=1}^n (y_m - \hat{y}_m)^2, \quad (10)$$

189 where y_m and \hat{y}_m are the observed and estimated output values, respectively, for case m .

190 Cherkassky and Ma (2004) suggested that with d predictor variables the RBF width
 191 σ_K should behave as $\sigma_K^d \sim (0.1 - 0.5)$. Hence, no precise estimate of γ was given. From
 192 now on we will call the SVR with the Cherkassky-Ma hyper-parameters just SVR for
 193 brevity.

194 In SVR-ES, the ES is initialized from a starting solution \mathbf{h} containing the Cherkassky-
195 Ma hyper-parameters. The local search algorithm searches the candidate solutions in
196 $N(\mathbf{h})$, a set of neighbors, for an \mathbf{h}' that has a better fitness function than \mathbf{h} . Basically,
197 the fitness function assigns a fitness value to each point in the parameter space (adaptive
198 landscape (Eiben and Schoenauer, 2002)), where this value can be seen as a measure
199 of how good a solution, represented by that point in the landscape, is to the given
200 problem (Hordijk, 1996). In this way, if a solution exists (i.e. has better fitness function)
201 it is accepted as the new incumbent solution, and the search proceeds by examining
202 the candidate solution in $N(\mathbf{h}')$. In our particular case, we minimize the MSE ($f(\mathbf{h}) =$
203 $(\text{MSE})^{-1}$). The necessary steps for implementing the SVR-ES algorithm are shown in
204 Figure 1.

205 Eventually, this process will lead to the identification of a local optimum (Eiben and
206 Smith, 2003). However, as this was done for a single individual instead of a popula-
207 tion (which is necessary for a global search), this approach is highly dependent on its
208 initial condition (initial individual). Assuming that the initial hyper-parameters from
209 Cherkassky-Ma are reasonable, our local search with ES is adequate.

210 In particular, the use of SVR-ES is attractive because it can be implemented in an
211 almost automatic way, uses variable steps to find the optimal hyper-parameters, is self-
212 adaptive, and, by using $\tau' = 1/\sqrt{2p}$ and $\tau = 1/\sqrt{2\sqrt{p}}$, the only EA algorithm parameter
213 that remains to be adjusted is the number of iterations.

214 4. Artificial Neural Network and Regression Trees

215 For comparison with the SVR and SVR-ES models, an artificial neural network
216 (ANN) and two regression tree algorithms are also applied to the WCCI-2006 datasets.

217 An ANN is a biologically inspired mathematical model which is composed of a large
218 number of highly interconnected perceptrons divided in layers (input, hidden, output),
219 but working in union to solve a problem. Training of the ANN model involves adjusting
220 the parameters iteratively so the MSE between the model output \hat{y} and target y is mini-
221 mized. Overfitting is a known problem with ANN. If p , the number of model parameters,
222 is too large, the ANN may overfit the data, producing a spuriously good fit that does
223 not lead to better predictions for new cases. This motivates the use of model comparison

224 criteria, such as the corrected Akaike information criterion (AIC_c), which penalizes the
225 MSE as a function p . The AIC_c , a bias-corrected version of the original AIC, is given by

$$AIC_c = n \ln(\text{MSE}) + 2p + 2(p+1)(p+2)/(n-p-2), \quad (11)$$

226 where n is the number of effective observations (Faraway and Chatfield, 1998). Models
227 with increasing numbers of hidden neurons n_h are trained and the model that minimizes
228 AIC_c is chosen as the optimum model. The ANN architecture and training algorithm used
229 in this study are based on Cannon and Lord (2000), with one hidden layer, early stopping
230 to avoid overfitting, an ensemble (or committee) of ANN models to deal with local minima
231 in the objective function, and with bootstrap aggregation (bagging) (Breiman, 1996).

232 Conceptually simple yet powerful, regression tree analysis, applicable to data sets with
233 both a large number of patterns and a large number of variables, is extremely resistant
234 to outliers. Among the decision trees used in machine learning, the classification and
235 regression tree (CART) model, first introduced by Breiman et al. (1984), is the most
236 commonly used. A random forest (RF) (Breiman, 2001) is a bootstrap ensemble of many
237 decision trees (CART). Each tree is grown over a bootstrap sample from the training data
238 set using a randomly selected subset of the available predictors at each decision branch.

239 Like CART, M5 builds a tree-based model, however the tree constructed by M5 can
240 have multivariate linear models at its leaves - the model trees are thus analogous to
241 piecewise linear functions. The advantage of M5 over CART is that M5 model trees are
242 generally much smaller and more accurate in some problem domains (Quinlan, 1992).

243 5. Experimental Results

244 To demonstrate the practical use of the forecasting methods for environmental prob-
245 lems, experiments were performed on the datasets outlined in Section 2. In order to
246 compare linear versus non-linear approaches, we also performed experiments with MLR.
247 Linear models are simple, fast, and often provide adequate and interpretable descriptions
248 of how the inputs affect the output. In particular for prediction purposes they can some-
249 times outperform fancier nonlinear models (Hastie et al., 2009). All the forecast models
250 can be built using the free R software environment for statistical computing (R Develop-
251 ment Core Team, 2011). For RF we used the R package randomForest (Liaw and Wiener,

252 2002). For SVR we used the R package e1071 (Dimitriadou et al., 2011). We developed
 253 the SVR-ES using the R package e1071 and the R native libraries; SVR-ES code is freely
 254 available from the project website: <http://forai.r-forge.r-project.org/>. For M5
 255 we used the package RWeka (Hornik et al., 2009). ANN uses the monmlp package. For
 256 MLR and SLR we used the package stats.

257 Data were standardized (zero mean and unit variance) and divided into a training set
 258 (75% of the data) and an independent test set (last 25% of the data) which is used to
 259 test the trained model. For the SVR we used a 5-fold cross validation within the training
 260 set to train and validate the model. For SVR-ES we subdivided the training set into two
 261 parts leaving the final 20% to validate the model. The RF and ANN perform their own
 262 split-sample validation via the out-of-bootstrap samples in the bootstrap aggregation.

263 The ANN setup was as follows: the range of initial random weights was between
 264 $[-0.5 : 0.5]$, 30 ensemble members were used for bagging, and 5000 was used as the
 265 maximum number of iterations. For RF, we used 500 as the number of generated trees
 266 and 5, the default value, as the minimum size of terminal nodes. For M5 we used the
 267 default values of the package RWeka. For the SRV-ES, C and ϵ were initialized by the
 268 Cherkassky-Ma guidelines and γ by 0.001.

269 For SVR, the C and ϵ values used were from Cherkassky-Ma, but for γ we did a grid
 270 search using the range suggested by Cherkassky-Ma (see section 3.2) and the extended
 271 range $[2^{-10}, 2^4]$ suggested by Lin and Lin (2003). Table 1 shows that using γ from
 272 Cherkassky-Ma yielded poorer results than those from Lin and Lin. Henceforth, SVR
 273 will refer to the model using a γ search following Lin and Lin.

274 In order to compute the relative accuracy between the MLR and the non-linear meth-
 275 ods we calculated the skill score (SS) (Hsieh, 2009) of the MAE and MSE. The SS is
 276 defined by:

$$SS = \frac{A - A_{\text{ref}}}{A_{\text{perfect}} - A_{\text{ref}}}, \quad (12)$$

277 where A is a particular measure of accuracy, A_{perfect} is the value of A for a set of perfect
 278 forecasts (in our case MAE and MSE equal to zero), and A_{ref} is the value of A computed
 279 over the set of reference forecasts (in our case the MAE or MSE value of the MLR model
 280 using all predictors). Positives SS means better performance than the reference model
 281 and negative values, worse performance.

282 According to the skill scores for PRECIP, the SVR (with the extended γ range)
283 achieved the best MAE results (Figure 2) and the SVR-ES the best MSE results (Fig-
284 ure 3). That the two SVR methods perform better than the other methods relative to
285 the MAE could be expected due to the ϵ -insensitive error norm of the SVR being more
286 similar to the MAE. The SLR reduced the number of predictors from 106 to 59, cutting
287 approximately 45% of the predictors. Basically, the SLR did not have an impact on the
288 MLR accuracy. For the ANN, the MAE SS is also essentially the same with and without
289 prescreening by SLR. However, an improvement is noted in the MSE when the SLR was
290 applied. The combination of SLR with M5 and RF tended to diminish the skill. As
291 M5 and RF do their own variable selection, it makes sense that SLR is not helpful. On
292 other hand, SVR-ES had a minor improvement with fewer predictors. All the non-linear
293 methods had better MAE results than the MLR (Figure 2). However the M5 had poor
294 MSE performance when compared with the MLR (around 10% worse) (Figure 3). The
295 SVR-ES had good results when compared with MLR, with skill scores exceeding 20% in
296 MAE and over 10% in MSE.

297 Figures 4 and 5 show that for TEMP the SVR-ES achieved the best results in both
298 MAE and MSE. Excluding M5, all the non-linear methods again had better MAE SS
299 than the MLR. The SLR reduced the 106 predictors to 60, cutting approximately 44%
300 of the predictors. Again the strongest improvement was ANN combined with SLR. SLR
301 was also beneficial to SVR-ES, SVR, and RF but, as in PRECIP, detrimental to M5.

302 For the SO2 dataset, Figures 6 and 7 show RF as the best performer. However, SVR-
303 ES kept its good performance in both MSE and MAE. Figure 6 shows again that MLR is
304 the worst in terms of MAE performance. SVR-ES has MAE SS performance around 10%
305 better than the MLR, and RF and SVR around 15% better. SLR reduced the number of
306 predictors for SO2 from 27 to 21, cutting approximately 23% of the predictors. Although
307 SLR did not have major impact on the forecast accuracy for the SO2 dataset, the results
308 are still desirable since the reduction of predictors reduces computing time.

309 To clarify why SVR performs better than SVR-ES in Figures 2 and 6, we performed
310 a small experiment on the PRECIP and SO2 datasets (with all predictors used) by
311 varying the fitness function in ES. Instead of maximizing only $(\text{MSE})^{-1}$, we also tested
312 using $(\text{MAE})^{-1}$ as fitness function, and calculated the skill scores using the MLR values

313 as reference. Figure 8 compares the skill score of MAE and MSE with $f_1 = (\text{MSE})^{-1}$ as
314 fitness function, and $f_2 = (\text{MAE})^{-1}$ as fitness function.

315 For the PRECIP dataset (Figure 8a), with f_1 as fitness function, SVR-ES has worse
316 MAE SS than SVR, but the two are comparable when f_2 is used instead. However, for
317 the MSE SS, SVR-ES did better than SVR, regardless of whether f_1 or f_2 was used.
318 Similarly, for the SO2 dataset (Figure 8b), SVR-ES with f_1 has worse MAE SS than
319 SVR, but slightly better MAE SS than SVR with f_2 . For the MSE SS, SVR-ES (with
320 either f_1 or f_2) did better than SVR, though the difference between f_1 and f_2 is more
321 pronounced than in the PRECIP dataset. More guidelines on why an error measure can
322 be minimized while a different error measure can remain unchanged or diminish can be
323 found at Jachner et al. (2007) and Lima et al. (2010). In general, it is not possible for
324 a single error measure to perform best on all criteria (e.g. cost, reliability, sensitivity,
325 resistance to outliers, relationship to decision making, etc.) (Armstrong and Collopy,
326 1992), hence the choice of the fitness function is dependent on the desired goal.

327 The heart of the learning problem is generalization (Witten et al., 2011), i.e., whether
328 the methods can retain satisfactory performance on new datasets. Based on the previous
329 results (Figures 3, 4 and 5), M5 is not recommended. On the other hand, ANN, SVR,
330 SVR-ES, and RF have better performance than the MLR when considering both MAE
331 and MSE over the three tested datasets.

332 It is difficult to evaluate the computing time of two methods under every parameter
333 setting because different values of the SVR hyper-parameters affect the training time (Fan
334 et al., 2005). To illustrate the time reduction provided by SLR, a small experiment
335 using the PRECIP dataset was performed. Varying the size of the training dataset, we
336 measured the cpu time used to train the SVR model (with γ fixed at 0.001) and forecast
337 the test set, using (i) all predictors and (ii) only the predictors selected by SLR. For 1000,
338 3000, 5000 and 7000 points in the training set, using all predictors required cpu times
339 of 10, 135, 371 and 727 s, respectively, while using only the predictors selected by SLR
340 required 5, 60, 212 and 469 s, respectively, thus confirming that screening of predictors
341 by SLR saves computing time.

342 To provide some guidelines and explanation for SVR results, Figures 9 and 10 show
343 changes related to the range of the γ value, with γ varying between $[2^{-10}, 2^4]$. The

344 MAE and MSE values were rescaled to range between $[0, 1]$ in order to compare the
345 sensitivity between the different datasets. First, as discussed in Cherkassky-Ma, indeed
346 there is dependency related to the γ values and the dataset dimensionality. However, this
347 dependency is not present over the full tested range. For example, in Figures 9 and 10,
348 the solid lines with triangles correspond to the full dimensionality (106 predictors) and
349 the dashed lines with triangles correspond to a reduced dimensionality (60 predictors)
350 of the TEMP dataset; within the range of $[2^{-10}, 2^{-6}]$ there is no difference between
351 the dimensionality and the values of MAE/MSE. The same behavior occurs in the SO2
352 dataset over the range of $[2^{-10}, 2^{-4}]$ (solid and dashed lines with x symbols). Another
353 interesting point is that the best γ values are independent of the dimensionality, i.e.,
354 solid lines and dashed lines converge to the same minimal point. This means that the
355 dependency of γ is not totally related to the dataset dimensionality but with the dataset
356 characteristics, as can be seen in the range of $[2^{-10}, 2^{-6}]$. However, in this range, the
357 MAE and MSE results for the PRECIP dataset had different optimal γ values, though
358 this difference is only one unit ($\log_2 \gamma$) as the grid search used discrete steps of this step
359 size.

360 Finally, to show the differences between various approaches to setting SVR hyper-
361 parameters, we performed an experiment on the PRECIP dataset again, using all pre-
362 dictors and the same test set as before, but only the last 500 points of the training set.
363 We tested the modified GA proposed by Leung et al. (2003) with different parameters
364 settings (GA1, GA2 and GA3), SVR using the procedure recommended by Cherkassky
365 and Ma (2004) but with the extended range $\gamma = [2^{-10}, 2^4]$ as suggested by Lin and Lin
366 (2003), SVR-ES with 200 generations (which is less than the 500 used in GA1 and GA2)
367 and SVR using 3-D grid search with the range of the three hyper-parameters recom-
368 mended by Fan et al. (2005). For GA1, we used 500 as the number of generations, 0.1
369 as the mutation rate, 0.9 as the crossover weight and 50 as the size of the population
370 (which was initialized randomly). For GA2, the corresponding values were 500, 0.2, 0.8
371 and 20, respectively, while for GA3, the values were 200, 0.3, 0.5 and 100, respectively.
372 The reference model was again the MLR (trained on the 500 points).

373 SVR-ES and SVR had the best results according to the MSE and MAE SS (Figure 11).
374 As shown by Kramer et al. (2007), changing the GA parameters can improve/worsen the

375 final results. Similar to Figure 2, SVR had better performance in the MAE SS than
376 SVR-ES, but this can be reversed by using $(MAE)^{-1}$ as the fitness function (Figure 8).
377 The 3-D grid search had similar performance as GA1. In terms of computing efforts,
378 SVR required 15 evaluations (for the 15 different γ values used in the search), SVR-
379 ES 200 evaluations, 3-D grid search about 1500, and the three GA models 1400–3500
380 evaluations.

381 6. Summary and Conclusion

382 In summary, we used three environmental datasets to test five non-linear forecast-
383 ing methods (four well-known and one that we proposed). Except for M5, the nonlin-
384 ear methods (ANN, SVR, SVR-ES and RF) generally outperformed the linear method
385 (MLR). Prescreening of predictors by SLR is generally beneficial for the nonlinear models
386 (except for M5), as it reduces the computing time and may increase the forecast skills
387 (especially for ANN). As explained in Witten et al. (2011), there is no universal best
388 learning method. SVR-ES had very good accuracy when the datasets were TEMP and
389 PRECIP, while RF had the best accuracy when the dataset was SO2. During pollution
390 events, the SO₂ concentration spikes much more dramatically than the high values ob-
391 served in the temperature and precipitation data, hence the architecture of RF may be
392 particularly suited for predicting variables with a heavy-tailed distribution.

393 The best overall method tends to be SVR-ES. SVR (with the Cherkassky-Ma esti-
394 mates for C and ϵ and an extended grid search in γ) worked well in terms of the MAE
395 skill score, and provided satisfactory performance in terms of the MSE. It used a rela-
396 tively modest amount of computing time for the hyper-parameter search. When using
397 the fitness function of $(MSE)^{-1}$ in the ES, SVR-ES may sometimes underperform SVR
398 in terms of the MAE skill score, but changing the fitness function to $(MAE)^{-1}$ appears
399 to eliminate this problem.

400 As recommended by Eiben and Smith (2003), a common approach is to start with
401 simple models; if the results are not of good enough quality, then move to more com-
402 plex models. The SVR-ES accuracy can be improved using more complex ES such as
403 correlated mutations (Eiben and Smith, 2003) or covariance matrix adaptation evolution
404 strategy (CMA-ES) (Friedrichs and Igel, 2005). Multi-processors can also be used in the

405 computation, i.e. given a chromosome \mathbf{G} (equation 7), k independent mutations can be
406 made among ρ processors (Verhoeven and Aarts, 1995).

407 References

- 408 Armstrong, J.S., Collopy, F., 1992. Error measures for generalizing about forecasting methods: Empirical
409 comparisons. *International Journal of Forecasting* 8, 69–80.
- 410 Bishop, C.M., 2006. *Pattern Recognition and Machine Learning (Information Science and Statistics)*.
411 Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- 412 Breiman, L., 1996. Bagging predictors. *Machine Learning* 24, 123–140. 10.1007/BF00058655.
- 413 Breiman, L., 2001. Random forests. *Machine Learning* 45, 5–32. 10.1023/A:1010933404324.
- 414 Breiman, L., Friedman, J., Olshen, R., Stone, C., 1984. *Classification and Regression Trees*. Wadsworth
415 Inc.
- 416 Cannon, A.J., Lord, E.R., 2000. Forecasting summertime surface-level ozone concentrations in the lower
417 fraser valley of british columbia: An ensemble neural network approach. *J. Air & Waste Manage* 50,
418 pp. 322–339.
- 419 Cawley, G.C., Janacek, G.J., Haylock, M.R., Dorling, S.R., 2007. Predictive uncertainty in environmental
420 modelling. *Neural Networks* 20, 537 – 549. *Computational Intelligence in Earth and Environmental*
421 *Sciences*.
- 422 Cherkassky, V., Ma, Y., 2004. Practical selection of svm parameters and noise estimation for svm
423 regression. *Neural Netw.* 17, 113–126.
- 424 Dimitriadou, E., Hornik, K., Leisch, F., Meyer, D., , Weingessel, A., 2011. e1071: Misc Functions of the
425 Department of Statistics (e1071), TU Wien. R package version 1.5-26.
- 426 Draper, N.R., Smith, H., 1998. *Applied Regression Analysis (Wiley Series in Probability and Statistics)*.
427 Wiley series in probability and mathematical statistics. Applied probability and statistics, John Wiley
428 & Sons Inc. 2 sub edition.
- 429 Eiben, A.E., Schoenauer, M., 2002. Evolutionary computing. *Information Processing Letters* 82, 1–6.
- 430 Eiben, A.E., Smith, J.E., 2003. *Introduction to Evolutionary Computing*. Natural Computing Series,
431 Springer, Berlin.
- 432 Fan, R.E., Chen, P.H., Lin, C.J., 2005. Working set selection using second order information for training
433 support vector machines. *J. Mach. Learn. Res.* 6, 1889–1918.
- 434 Faraway, J., Chatfield, C., 1998. Time series forecasting with neural networks: a comparative study
435 using the air line data. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 47,
436 123–140. 10.1111/1467-9876.00109.
- 437 Friedrichs, F., Igel, C., 2005. Evolutionary tuning of multiple svm parameters. *Neurocomputing* 64,
438 107–117. *Trends in Neurocomputing: 12th European Symposium on Artificial Neural Networks 2004*.
- 439 Fritsch, J.M., Houze, R.A., Adler, R., Bluestein, H., Bosart, L., Brown, J., Carr, F., Davis, C., Johnson,
440 R.H., Junker, N., et al., 1998. Quantitative precipitation forecasting: Report of the eighth prospectus

- 441 development team, us weather research program. *Bulletin of the American Meteorological Society* 79,
442 285–299.
- 443 Hall, M., Smith, L., 1996. Practical feature subset selection for machine learning, in: McDonald, C. (Ed.),
444 *Computer Science '98 Proceedings of the 21st Australasian Computer Science Conference ACSC'98*,
445 Springer. pp. 181–191.
- 446 Hashmi, M.Z., Shamseldin, A.Y., Melville, B.W., 2009. Statistical downscaling of precipitation: state-
447 of-the-art and application of bayesian multi-model approach for uncertainty assessment. *Hydrology*
448 *and Earth System Sciences Discussions* 6, 6535–6579.
- 449 Hastie, T., Tibshirani, R., Friedman, J., 2009. *The Elements of Statistical Learning: Data Mining,*
450 *Inference, and Prediction, Second Edition.* Springer Series in Statistics, Springer. 2nd ed. 2009. corr.
451 3rd printing 5th printing. edition.
- 452 Haupt, S., Pasini, A., Marzban, C., 2009. *Artificial intelligence methods in the environmental sciences.*
453 Springer.
- 454 Haykin, S., 1998. *Neural Networks - A Comprehensive Foundation.* Pearson Education. second edition.
- 455 Hocking, R.R., 1976. A biometrics invited paper. the analysis and selection of variables in linear regres-
456 sion. *Biometrics* 32, pp. 1–49.
- 457 Hordijk, W., 1996. A measure of landscapes. *Evolutionary Computation* 4, 335–360.
- 458 Hornik, K., Buchta, C., Zeileis, A., 2009. Open-source machine learning: R meets Weka. *Computational*
459 *Statistics* 24, 225–232.
- 460 Hsieh, W.W., 2009. *Machine Learning Methods in the Environmental Sciences: Neural Networks and*
461 *Kernels.* Cambridge University Press, New York, NY, USA.
- 462 Huang, C.L., Wang, C.J., 2006. A ga-based feature selection and parameters optimization for support
463 vector machines. *Expert Systems with Applications* 31, 231–240.
- 464 Jachner, S., van den Boogaart, K.G., Petzoldt, T., 2007. Statistical methods for the qualitative as-
465 sessment of dynamic models with time delay (r package qualv). *Journal of Statistical Software* 22,
466 1–30.
- 467 Kramer, O., Gloger, B., Goebels, A., 2007. An experimental analysis of evolution strategies and particle
468 swarm optimizers using design of experiments, in: *Proceedings of the 9th annual conference on Genetic*
469 *and evolutionary computation*, ACM, New York, NY, USA. pp. 674–681.
- 470 Krasnopolsky, V.M., 2007. Neural network emulations for complex multidimensional geophysical map-
471 pings: Applications of neural network techniques to atmospheric and oceanic satellite retrievals and
472 numerical modeling. *Reviews of Geophysics* 45.
- 473 Kuligowski, Robert J., A.P.B., 1998. Localized precipitation forecasts from a numerical weather predic-
474 tion model using artificial neural networks. *Wea. Forecasting* 13, 1194–1204.
- 475 Kurt, A., Gulbagci, B., Karaca, F., Alagha, O., 2008. An online air pollution forecasting system using
476 neural networks. *Environment International* 34, 592–598.
- 477 Leung, F., Lam, H., Ling, S., Tam, P., 2003. Tuning of the structure and parameters of a neural network
478 using an improved genetic algorithm. *Neural Networks, IEEE Transactions on* 12, 79–88.
- 479 Liaw, A., Wiener, M., 2002. Classification and regression by randomforest. *R News* 2, 18–22.

480 Lima, A.R., Silva, D.A., Mattos Neto, P.S., Ferreira, T.A., 2010. An experimental study of fitness
481 function and time series forecasting using artificial neural networks, in: Proceedings of the 12th
482 annual conference companion on Genetic and evolutionary computation, ACM, New York, NY, USA.
483 pp. 2015–2018.

484 Lin, K.M., Lin, C.J., 2003. A study on reduced support vector machines. *Neural Networks, IEEE*
485 *Transactions on* 14, 1449 – 1459.

486 Lin, S.W., Ying, K.C., Chen, S.C., Lee, Z.J., 2008. Particle swarm optimization for parameter deter-
487 mination and feature selection of support vector machines. *Expert Systems with Applications* 35,
488 1817–1824.

489 Muller, K., Smola, A., Rtsch, G., Schlkopf, B., Kohlmorgen, J., Vapnik, V., 1997. Predicting time
490 series with support vector machines, in: Gerstner, W., Germond, A., Hasler, M., Nicoud, J.D. (Eds.),
491 *Artificial Neural Networks ICANN'97*. Springer Berlin / Heidelberg. volume 1327 of *Lecture Notes*
492 *in Computer Science*, pp. 999–1004. 10.1007/BFb0020283.

493 Ortiz-García, E.G., Salcedo-Sanz, S., Pérez-Bellido, A.M., Portilla-Figueras, J.A., 2009. Improving the
494 training time of support vector regression algorithms through novel hyper-parameters search space
495 reductions. *Neurocomputing* 72, 3683–3691.

496 Pai, P.F., Hong, W.C., 2005. Forecasting regional electricity load based on recurrent support vector
497 machines with genetic algorithms. *Electric Power Systems Research* 74, 417–425.

498 Pino-Mejías, R., de-la Vega, M.D.C., Anaya-Romero, M., Pascual-Acosta, A., Jordán-López, A.,
499 Bellinfante-Crocci, N., 2010. Predicting the potential habitat of oaks with data mining models and
500 the r system. *Environmental Modelling & Software* 25, 826 – 836.

501 Quinlan, J.R., 1992. Learning with continuous classes, in: Adams, A., Sterling, L. (Eds.), Proceedings
502 of the 5th Australian joint Conference on Artificial Intelligence, World Scientific. pp. 343–348.

503 R Development Core Team, 2011. R: A Language and Environment for Statistical Computing. R
504 Foundation for Statistical Computing. Vienna, Austria. ISBN 3-900051-07-0.

505 Schoof, J., Pryor, S., 2001. Downscaling temperature and precipitation: a comparison of regression-based
506 methods and artificial neural networks. *International Journal of Climatology* 21, 773–790.

507 Smit, S., Eiben, A., 2009. Comparing parameter tuning methods for evolutionary algorithms, in: Evo-
508 lutionary Computation, 2009. CEC '09. IEEE Congress on, pp. 399 –406.

509 Solomatine, D.P., Xue, Y., 2004. M5 model trees and neural networks: Application to flood forecasting
510 in the upper reach of the huai river in china 9, 491–501.

511 Strangeways, I., 2007. Precipitation: theory, measurement and distribution. Cambridge University Press.

512 Tripathi, S., Srinivas, V., Nanjundiah, R.S., 2006. Downscaling of precipitation for climate change
513 scenarios: A support vector machine approach. *Journal of Hydrology* 330, 621–640.

514 Venables, W., Ripley, B., 2002. Modern applied statistics with S. Statistics and computing, Springer.

515 Verhoeven, M., Aarts, E., 1995. Parallel local search. *Journal of Heuristics* 1, 43–65.
516 10.1007/BF02430365.

517 Von Storch, H., Zwiers, F., 2001. Statistical Analysis in Climate Research. Cambridge University Press.

518 Wallace, J., Hobbs, P., 2006. Atmospheric Science: An Introductory Survey. *International Geophysics*

519 Series, Elsevier Academic Press.

520 Witten, I.H., Frank, E., Hall, M.A., 2011. Data Mining: Practical Machine Learning Tools and Tech-
521 niques. Morgan Kaufmann, Burlington, MA. 3 edition.

522 Zeng, Z., Hsieh, W.W., Burrows, W.R., Giles, A., Shabbar, A., 2011. Surface wind speed prediction
523 in the canadian arctic using non-linear machine learning methods. Atmosphere-Ocean 49, 22–31.
524 <http://www.tandfonline.com/doi/pdf/10.1080/07055900.2010.549102>.

525 Zhang, G.P., 2007. Avoiding pitfalls in neural network research. IEEE Transactions on Systems, Man,
526 and Cybernetics, Part C 37, 3–16.

527 **List of Tables**

528 1 Comparison of MAE and MSE values from the test set using SVR with
529 the range of γ suggested by Cherkassky and Ma (2004) and by Lin and
530 Lin (2003) for the PRECIP, TEMP, and SO2 datasets, using either all
531 predictors (ALL) or predictors selected by SLR. 20

Table 1: Comparison of MAE and MSE values from the test set using SVR with the range of γ suggested by Cherkassky and Ma (2004) and by Lin and Lin (2003) for the PRECIP, TEMP, and SO2 datasets, using either all predictors (ALL) or predictors selected by SLR.

	MAE				MSE			
	Cherkassky-Ma		Lin and Lin		Cherkassky-Ma		Lin and Lin	
	ALL	SLR	ALL	SLR	ALL	SLR	ALL	SLR
PRECIP	0.0534	0.0530	0.0379	0.0380	0.00951	0.00944	0.00572	0.00572
TEMP	0.2520	0.2415	0.0622	0.0618	0.09199	0.08543	0.00642	0.00630
SO2	0.0277	0.0263	0.0244	0.0244	0.00386	0.00371	0.00325	0.00326

532 **List of Figures**

533	1	Procedure of the SVR-ES	22
534	2	Skill score of MAE for the PRECIP test dataset with all predictors used	
535		(dark bar) and with predictors selected by SLR (light bar). The reference	
536		forecast, MLR with all predictors, has simply 0 for the skill score.	22
537	3	Skill score of MSE for the PRECIP test dataset.	23
538	4	Skill score of MAE for the TEMP test dataset.	23
539	5	Skill score of MSE for the TEMP test dataset.	24
540	6	Skill score of MAE for the SO2 test dataset.	24
541	7	Skill score of MSE for the SO2 test dataset.	25
542	8	Skill score of MSE (dark bar) and MAE (light bar) for (a) PRECIP and (b)	
543		SO2 test datasets, where f_1 denotes SVR-ES with $f_1=(MSE)^{-1}$ as fitness	
544		function and f_2 denotes SVR-ES with $f_2=(MAE)^{-1}$	26
545	9	Results of MAE as γ varies between $[2^{-10}, 2^4]$, for the PRECIP test dataset	
546		with all predictors used (solid line with circles), with predictors selected	
547		by SLR (dashed line with circles), for the TEMP test dataset with all	
548		predictors used (solid line with triangles) and with predictors selected	
549		by SLR (dashed line with triangles), for the SO2 test dataset with all	
550		predictors used (solid line with x) and with predictors selected by SLR	
551		(dashed line with x).	27
552	10	Results of MSE as γ varies between $[2^{-10}, 2^4]$	28
553	11	Results of MSE skill score (left side) and MAE skill score (right side) for	
554		the PRECIP test dataset using the last 500 points of the training set and	
555		all the predictors to train the models. The models are respectively: the	
556		modified GA proposed by Leung et al. (2003) with different parameters	
557		settings (GA1, GA2 and GA3), SVR with the procedure recommended	
558		by Cherkassky and Ma (2004) and the extended range $\gamma = [2^{-10}, 2^4]$ sug-	
559		gested by Lin and Lin (2003), SVR-ES with 200 generations and 3-D	
560		grid search with the hyper-parameters' range recommended by Fan et al.	
561		(2005).	29

Figure 1: Procedure of the SVR-ES

```

begin
   $\tau \rightarrow 0$ ; //  $\tau$ : Iteration number;
  if Screenout irrelevant Predictors = TRUE then
    | Perform stepwise regression;
  end
  Initialize  $\mathbf{G}(\tau)$ ; //  $\mathbf{G}(\tau)$ : Chromosome according to Equation (7);
  Evaluate  $f(\mathbf{G}(\tau))$ ; //  $f(\mathbf{G}(\tau))$ : Fitness function value ;
  while not termination condition do
     $\tau \rightarrow \tau + 1$ ;
    Perform mutation operation according to Eqns. (8) and (9) to generate new chromosome  $\mathbf{G}'$  ;
    Evaluate  $f(\mathbf{G}'(\tau))$ ;
    if  $f(\mathbf{G}') > f(\mathbf{G})$  then
      |  $f(\mathbf{G}) \leftarrow f(\mathbf{G}')$ ;
    end
  end
end
end

```

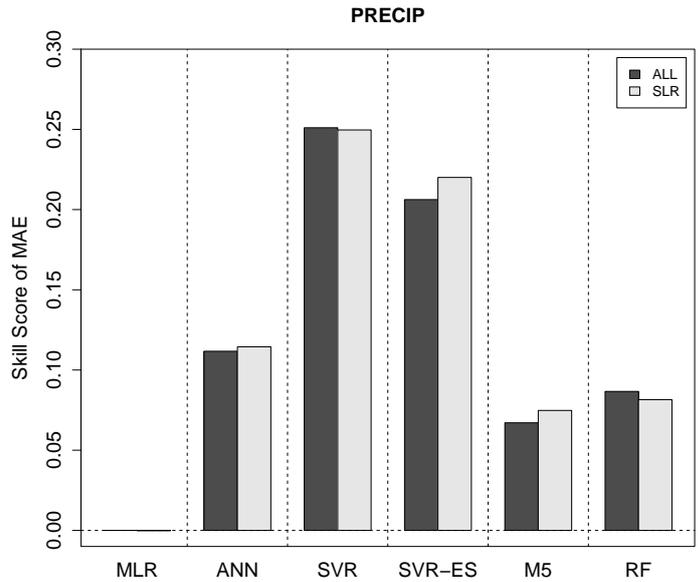


Figure 2: Skill score of MAE for the PRECIP test dataset with all predictors used (dark bar) and with predictors selected by SLR (light bar). The reference forecast, MLR with all predictors, has simply 0 for the skill score.

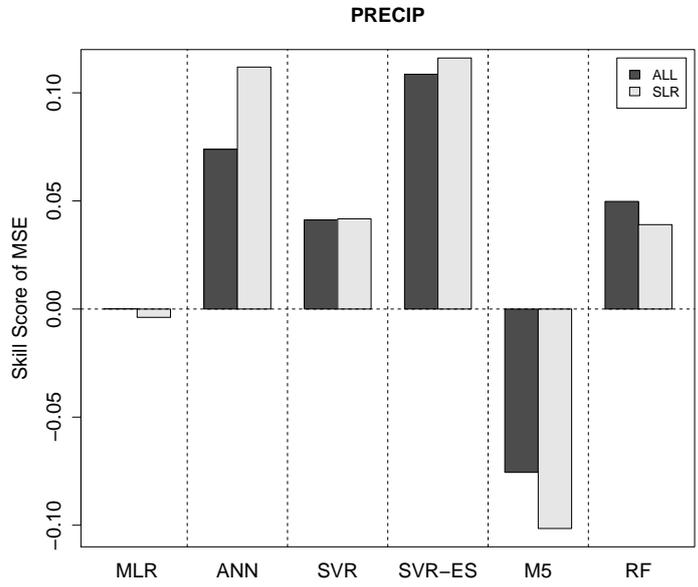


Figure 3: Skill score of MSE for the PRECIP test dataset.

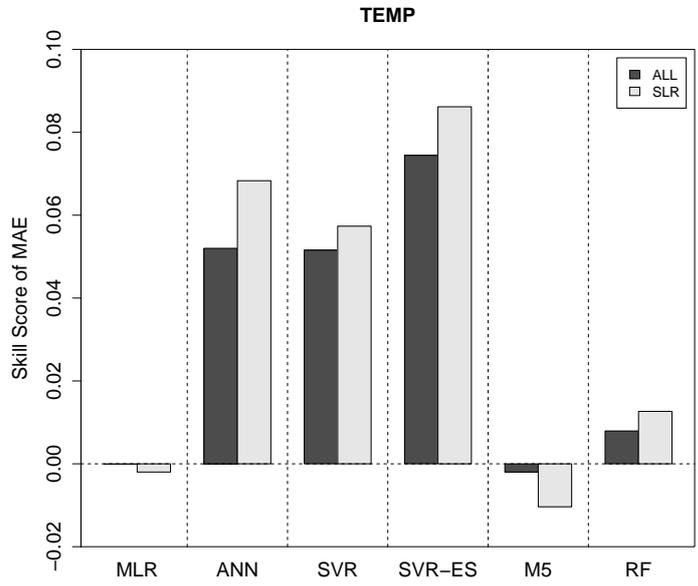


Figure 4: Skill score of MAE for the TEMP test dataset.

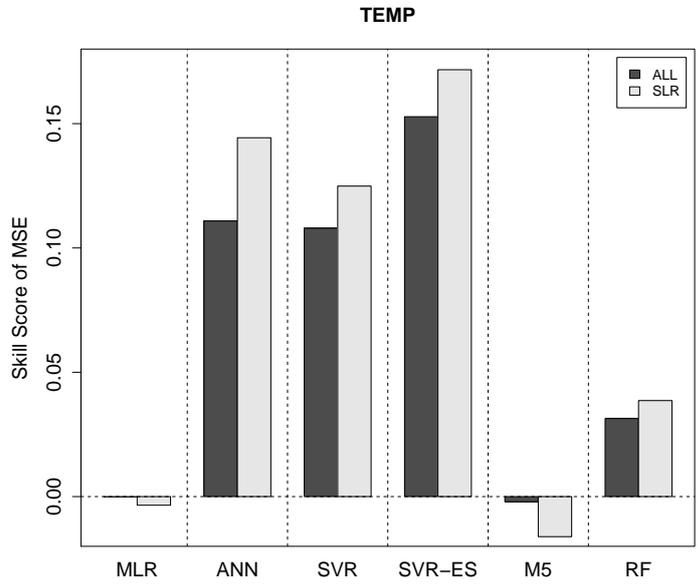


Figure 5: Skill score of MSE for the TEMP test dataset.

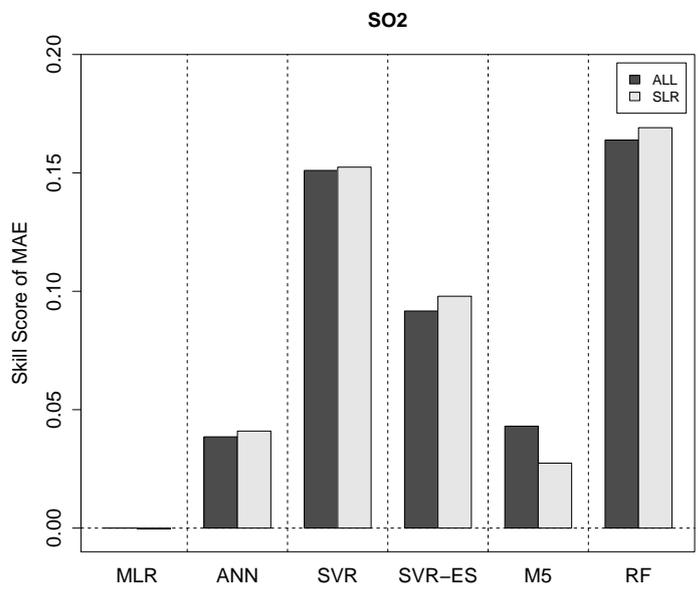


Figure 6: Skill score of MAE for the SO2 test dataset.

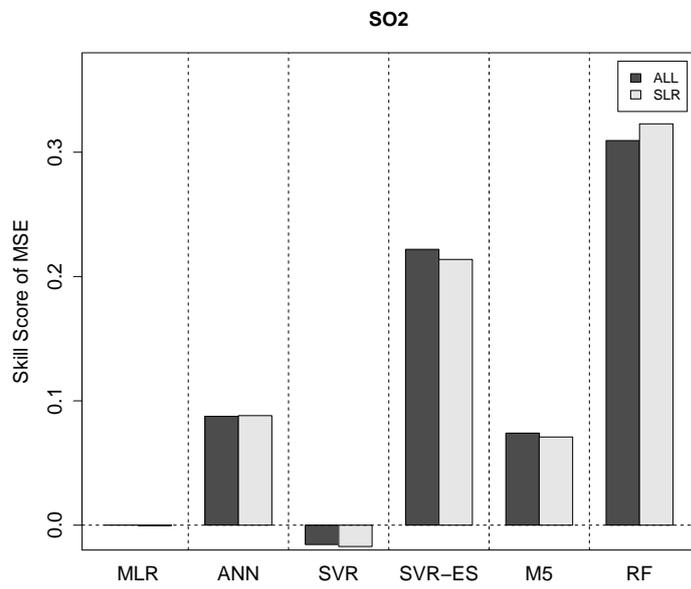


Figure 7: Skill score of MSE for the SO2 test dataset.

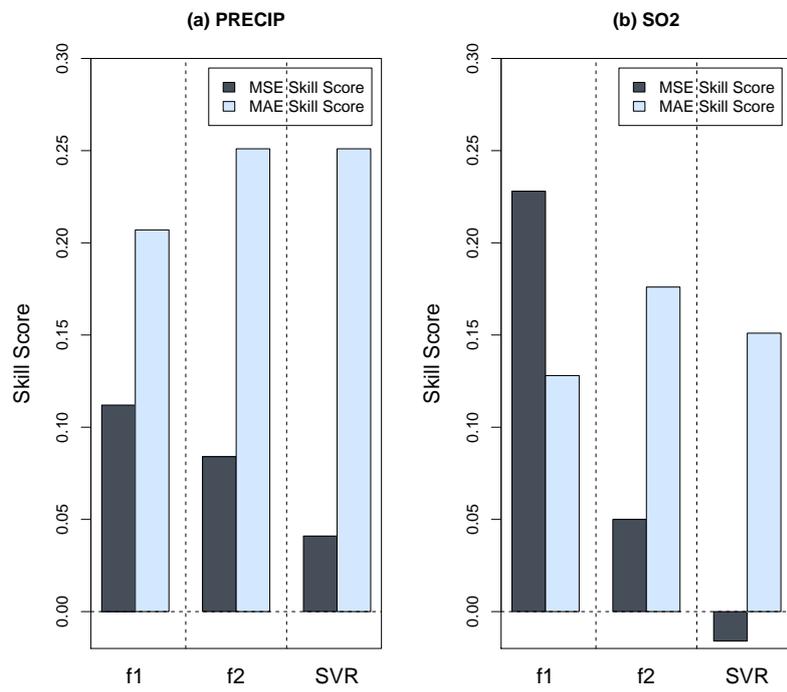


Figure 8: Skill score of MSE (dark bar) and MAE (light bar) for (a) PRECIP and (b) SO2 test datasets, where f1 denotes SVR-ES with $f_1=(\text{MSE})^{-1}$ as fitness function and f2 denotes SVR-ES with $f_2=(\text{MAE})^{-1}$.

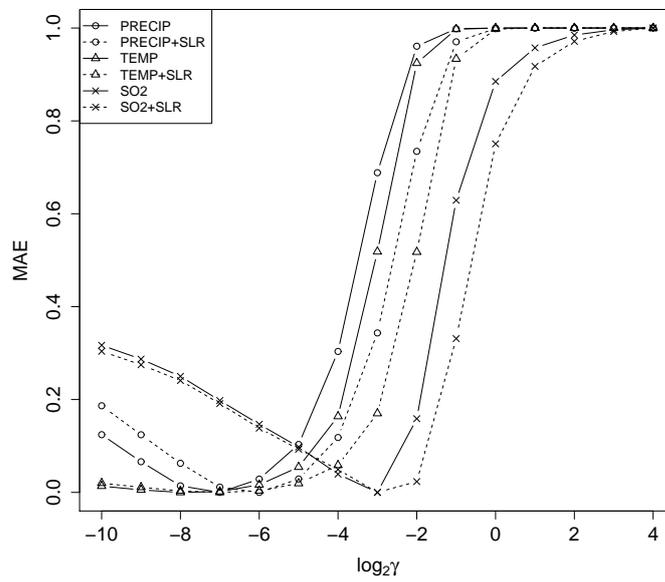


Figure 9: Results of MAE as γ varies between $[2^{-10}, 2^4]$, for the PRECIP test dataset with all predictors used (solid line with circles), with predictors selected by SLR (dashed line with circles), for the TEMP test dataset with all predictors used (solid line with triangles) and with predictors selected by SLR (dashed line with triangles), for the SO2 test dataset with all predictors used (solid line with x) and with predictors selected by SLR (dashed line with x).

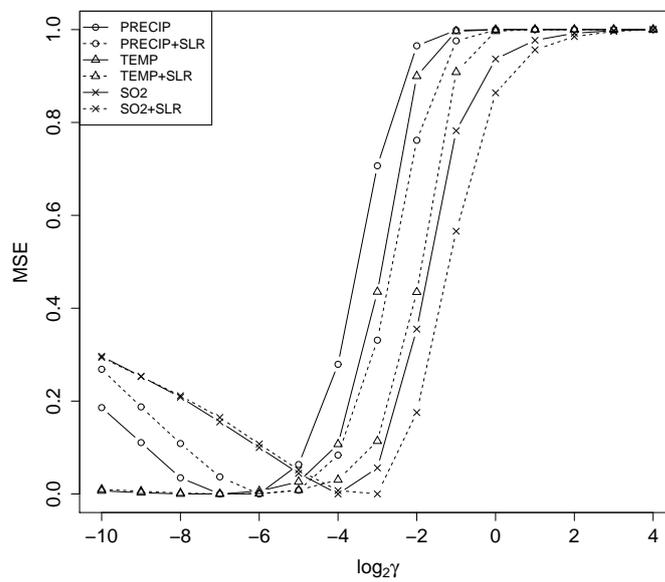


Figure 10: Results of MSE as γ varies between $[2^{-10}, 2^4]$.

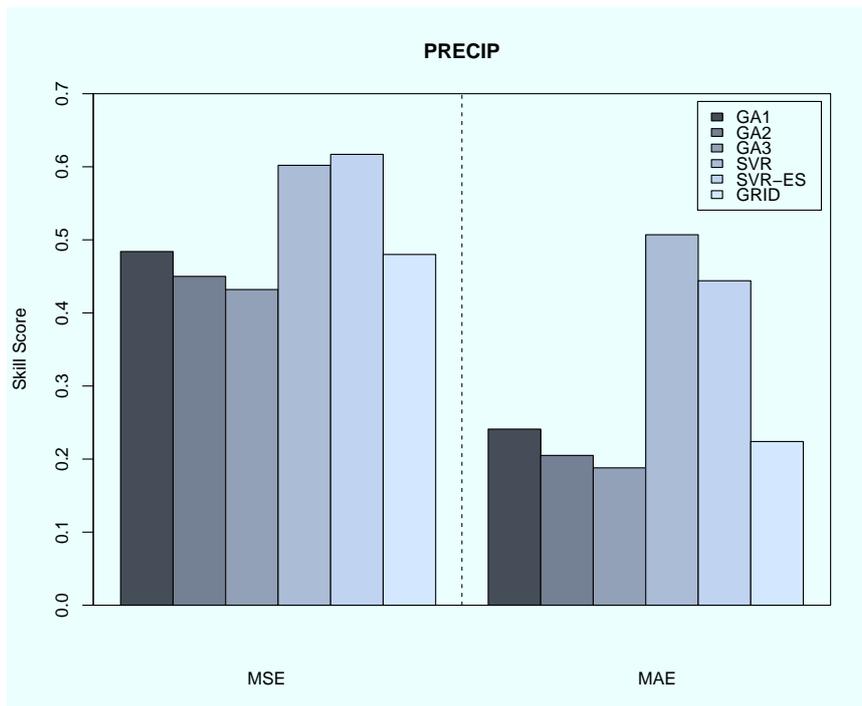


Figure 11: Results of MSE skill score (left side) and MAE skill score (right side) for the PRECIP test dataset using the last 500 points of the training set and all the predictors to train the models. The models are respectively: the modified GA proposed by Leung et al. (2003) with different parameters settings (GA1, GA2 and GA3), SVR with the procedure recommended by Cherkassky and Ma (2004) and the extended range $\gamma = [2^{-10}, 2^4]$ suggested by Lin and Lin (2003), SVR-ES with 200 generations and 3-D grid search with the hyper-parameters' range recommended by Fan et al. (2005).