

The application of support vector regression (SVR) for stream flow prediction on the Amazon basin

Melise du Toit^{1,2}, Josefine M. Wilms², Gideon J.F. Smit¹ and Willie Brink¹

¹ Department of Mathematical Sciences (Applied Mathematics), Stellenbosch University, South Africa

² Advanced Mathematical Modelling, Modelling and Digital Sciences, CSIR Stellenbosch, South Africa

Long-term forecasting of river runoff is important for climate scientists and hydrologists. By analysing the processes of a river basin characterized by measurable variables, an empirical data-driven model can be constructed. The support vector regression technique is used in this study to analyse historical stream flow occurrences and predict stream flow values for the Amazon basin. Up to twelve month predictions are made and the coefficient of determination and root-mean-square error are used for accuracy assessment. Compared to previous studies, satisfactory results are obtained. Inclusion of environmental aspects such as precipitation and evaporation are suggested for more accurate predictions.

Keywords: Support vector machine, Support vector regression, Amazon basin, Stream flow prediction

1. Introduction

Research on model-generated river runoff is essential for climate scientists and hydrologists to predict and understand future changes in river runoff that may be associated with global climate change. The hydrologic cycle is closed by returning the correct amount of water to the river mouth with the appropriate timing and position (Miller *et al.*, 1993). River engineers and scientists use these results for the study of various hydro-environmental aspects, such as the increasing international concern of riverine pollution problems and the growing flood levels of rivers (Falconer *et al.*, 2005). Furthermore, sediment transport and salinity changes within the river basin can be examined and predicted (Falconer *et al.*, 2005; Miller *et al.*, 1993).

Numerous hydrological models have been implemented by researchers to analyse the behaviour of river basins and to model river flow in such basins by mapping the natural phenomena to a simulation program (Falconer *et al.*, 2005). These models are known as physically based or process models, since they are based on the physical behaviour of the specific river basin system as well as the mathematical description of the river flow (Falconer *et al.*, 2005; Solomatine and Ostfeld, 2008). A physically based model consists of a numerical process which involves the computation of an efficient and accurate solution to equations based on the physical laws obtained for the specific system. The accuracy of a process model is tested by comparing its results to past observations, and if a desired accuracy is obtained, such a model may be

used to calculate and predict future changes in the particular system.

Even though various hydraulic and hydrologic process models have been constructed for river basin systems, limited knowledge of the required modelling processes in a system may result in an unreliable model. However, such a system may consist of a process characterized by measurable variables and contain a sufficient amount of concurrent input and output data associated with the particular process (Solomatine and Ostfeld, 2008). By analysing the relationship between the input and output data an empirical mathematical model, known as a data-driven model, can be constructed to model and predict future output variables (Solomatine and Ostfeld, 2008).

A detailed understanding of the physical processes and behaviour of a river basin system is therefore not required for the construction of a data-driven model. Instead, data-driven modelling involves a study of the relationship between the system's state variables (Solomatine *et al.*, 2008). This may allow for the improvement of physically based models.

The objective of this study is the description and implementation of an empirically based (data-driven) model for river runoff. In particular, a supervised machine learning model known as support vector regression (SVR) will be considered. This model is used to analyse the stream flow history of gauging stations in a river basin in order to determine future stream flow. The Amazon River in South America is considered for the application of this data-driven

71 model and an attempt to accurately predict stream
72 flow is made.

74 2. Instrumentation and Method

76 2.1. Study area and available data

78 Stream flow data for the Amazon basin have been
79 obtained from the Observation Service for the
80 geodynamical, hydrological and biogeochemical
81 control of erosion/alteration and material transport
82 (SO HYBAM). This association manages 20 gauging
83 stations that are distributed in the Amazon. The
84 stream flow records of three are considered for this
85 study: the Obidos station in Rio Amazonas, the
86 Manacapuru station in Rio Solimões, and the Lábrea
87 station in Rio Purus, shown in Fig. 1.



89 **Figure 1.** Study area and location of the gauging stations.

92 2.2. Support vector regression: model formulation

94 In order to forecast an outcome $y(t + \Delta t)$ at an
95 instant Δt from current time t , a regression method
96 can be constructed. The purpose of such a method is
97 to formulate a function $f(\mathbf{x})$ such that $f(\mathbf{x}) =$
98 $y(t + \Delta t)$. The function f takes an input vector
99 $\mathbf{x} = (x_1, x_2, \dots, x_m)$ of m known variables, including
100 current and past data records [$y(t), y(t - 1), \dots,$
101 $y(t - q)$], where $q \leq m$. The input vector may also
102 consist of any other available numerical variables.

104 An extension of the support vector machine (SVM),
105 formulated by Cortes and Vapnik (1995), is known as
106 the support vector regression (SVR) technique. A
107 thorough description on the construction of the SVR
108 technique, its optimization parameters (C and ϵ) and
109 its applications in the field of hydrology can be found
110 in Raghavendra and Deka (2014). An important
111 concept of the SVR method is that it attempts to find
112 a simple function that can fit all the data while
113 minimizing the sum of prediction errors above a
114 predefined margin (Callegari *et al.*, 2015).

116 For cases where the SVR model has to optimize
117 nonlinear functions, the input vector \mathbf{x} is mapped to a

118 feature space where its relationship with y is
119 linearized. This mapping function is known as a
120 kernel function. A detailed discussion on kernel
121 functions is given by Raghavendra and Deka (2014).

123 2.3. Model training and testing

125 The process of formulating a function $f(\mathbf{x})$ on a
126 given subset of the available data (known as the
127 *training set*) is known as training. During training,
128 the model is tested by fitting it to a second sample set
129 (known as the *validation set*). Finally, the trained
130 SVR model is verified by an accuracy measurement
131 on a third subset of the given samples, known as the
132 *test set* (Solomatine and Ostfeld, 2008).

134 For the Obidos gauging station, monthly stream flow
135 data from 1970 to 2000 are considered. Furthermore,
136 data from 1973 to 2003 and from 1968 to 1998 are
137 available for the Manacapuru and Lábrea stations,
138 respectively. For each station, data for the first 15
139 years are used as training sets. The following 10
140 years' data constitutes the validation set and the
141 remaining 5 years' data are used for testing.

143 2.4. Feature and kernel function selection

145 Each input vector \mathbf{x} consists of 12 antecedent stream
146 flow periods (months). The value of y represents the
147 flow in the next period. One month predictions are
148 made, where after forecasting is extended for up to
149 12 months. Evaluation is done by calculating the
150 coefficient of determination (R^2) of the predicted and
151 observed stream flow values. The purpose of R^2 is to
152 give an estimation of how well observed models are
153 replicated by the fitted model, based on the
154 percentage of total variation of outcomes interpreted
155 by the model. The R^2 percentage therefore represents
156 the percentage of variation of predicted outcomes that
157 are explained by the fitted model. Furthermore, the
158 root-mean-square error (RMSE) percentage indicates
159 residual variance between observed and forecasted
160 outcomes, and will be used for evaluation in this
161 study.

163 Linear, Polynomial (Poly) and Radial Basis (RBF)
164 kernel functions are considered. These SVR
165 formulations are expressed as follows:

167 Linear: $k(x_i, x_j) = x_i^T x_j,$

168 Poly: $k(x_i, x_j) = (\gamma x_i^T x_j + r)^d,$ and

169 RBF: $k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2); \gamma > 0.$

170

171 The mapping of features x_i and x_j to the feature
 172 space is represented by $k(x_i, x_j)$. An outline of the
 173 kernel functions and their hyperparameters are given
 174 by Granata *et al.* (2016). A built-in module in the
 175 Python programming language, known as Optunity,
 176 is used to optimize the parameters of each kernel.
 177

178 3. Results and Discussion

180 3.1. Optimal hyperparameters and kernel functions

181
 182 Historical stream flow records of the respective
 183 stations are examined. For each station, the optimal
 184 hyperparameters of the considered kernel functions
 185 are calculated in order to determine the best
 186 generalized model for the given data. The R^2 value
 187 for each optimized model is listed in Tables 1 to 3.
 188 For the training and validation sets, every kernel
 189 function provides an R^2 greater than 0.9, indicating
 190 that at least 90% of the total variation of predicted
 191 outcomes are explained by the fitted models. The
 192 RBF and polynomial kernel functions provide the
 193 best results for each station. However, the RBF
 194 kernel is less complex in comparison to polynomial
 195 kernels, since it contains fewer parameters. Further
 196 investigation is therefore done by only considering
 197 the RBF kernel.
 198
 199

200 **Table 1.** Optimized kernel-specific hyperparameters and R^2 for
 201 one month predictions of river flow at the Obidos gauging station.

OBIDOS GAUGING STATION							
Kernel	Optimal Parameters					R^2	
	C	ϵ	γ	d	r	Training set	Validation set
RBF	641	0.0303	0.067	-	-	0.983	0.967
Linear	304	0.0262	-	-	-	0.976	0.965
Poly	381	0.0307	0.1	2	0.3	0.981	0.966

202 **Table 2.** Optimized kernel-specific hyperparameters and R^2 for
 203 one month predictions of river flow at the Manacapuru gauging
 204 station.
 205

MANACAPURU GAUGING STATION							
Kernel	Optimal Parameters					R^2	
	C	ϵ	γ	d	r	Training set	Validation set
RBF	570	0.02	0.03	-	-	0.937	0.923
Linear	385	0.048	-	-	-	0.912	0.904
Poly	78	0.0056	0.1	3	0.5	0.942	0.925

206 **Table 3.** Optimized kernel-specific hyperparameters and R^2 for
 207 one month predictions of river flow at the Lábrea gauging station.
 208

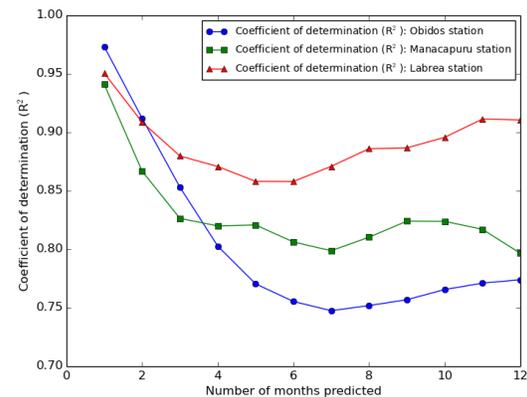
LABREA GAUGING STATION							
Kernel	Optimal Parameters					R^2	
	C	ϵ	γ	d	r	Training set	Validation set
RBF	255	0.05	1.7	-	-	0.985	0.965
Linear	84	0.015	-	-	-	0.956	0.951
Poly	675	0.0329	0.1	5	0.11	0.984	0.959

211 3.2. Extended stream flow forecasting

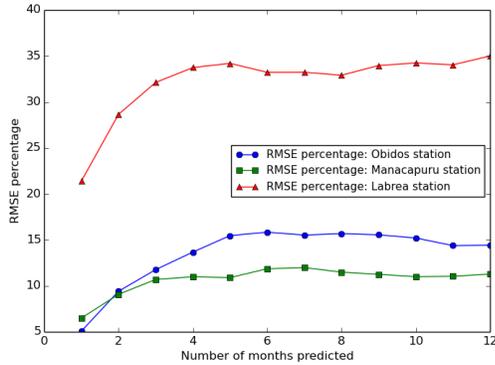
212
 213 The optimized RBF models are applied to the testing
 214 data for forecasting. At an instant (month) t , twelve
 215 antecedent observed flow values $\mathbf{x} = [y(t), y(t-1),$
 216 $\dots, y(t-11)]$ are used to predict flow $f(\mathbf{x})_{\{t+1\}}$
 217 for month $t+1$. This is known as *one month*
 218 *forecasting*. Similarly, for *two month forecasting*, an
 219 input vector $\mathbf{x} = [f(\mathbf{x})_{\{t+1\}}, y(t), \dots, y(t-10)]$ is
 220 used to predict stream flow for month $t+2$.
 221 Forecasting extending up to 12 months is done on the
 222 given test set of each station. The corresponding R^2
 223 values and $RMSE$ percentages are determined and
 224 shown in Figs. 2 and 3, respectively.
 225

226 For each gauging station the best results were
 227 obtained for one month forecasting. An R^2 of 0.973
 228 is obtained for the Obidos station, whereas R^2 values
 229 of 0.94 and 0.95 are obtained for the Manacapuru and
 230 Lábrea stations, respectively. Furthermore, the $RMSE$
 231 percentages are obtained respectively as 5.06%,
 232 6.49% and 21.38%. R^2 is a relative error of fit,
 233 whereas $RMSE$ is an absolute measure of fit. Since
 234 $RMSE$ is the square root of a variance, it can be
 235 explained as the standard deviation of the
 236 unexplained variance. This clarifies the larger $RMSE$
 237 values obtained for the Lábrea station. Compared to
 238 stream flow forecasting studies done by Veiga *et al.*
 239 (2015), Lin *et al.* (2006) and Callegari *et al.* (2015),
 240 these results are quite satisfactory.
 241

242 Extended forecasting produces less accurate results.
 243 However, it should be taken into account that
 244 predicted stream flow values were used to make
 245 future predictions. Also, stream flow is the only
 246 environmental/hydrological variable considered.
 247



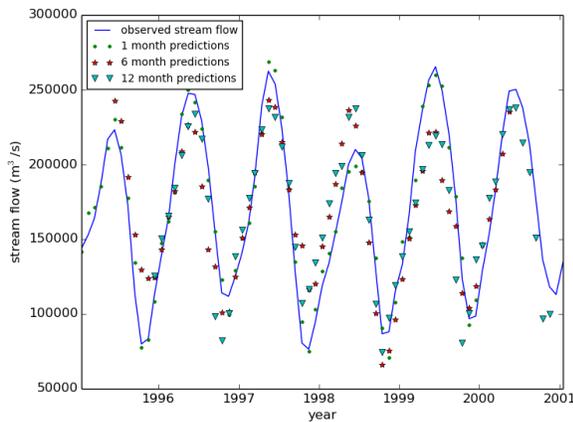
248 **Figure 2.** R^2 results for extended forecasting.
 249
 250



251
252 **Figure 3.** RMSE percentages for extended forecasting.
253

254 **3.3. Illustrations of stream flow predictions**
255

256 Figure 4 is an illustration of one, six and twelve
257 month extended stream flow forecasting compared to
258 observed stream flow. The worst predictions are
259 made at the minimum and maximum stream flow
260 occurrences, whereas good results are obtained for
261 the upward and downward flow tendencies.
262



263
264 **Figure 4:** Stream flow discharge at the Obidos station for 1, 6, and
265 12 month predictions.
266

267 **4. Conclusions**
268

269 Research on long-term forecasting of river runoff
270 predictions is important for climate scientists and
271 hydrologists, since these results are used for the study
272 of various hydro-environmental aspects. Numerous
273 physically based hydrologic models have been
274 implemented by researchers for this task, but due to
275 limited knowledge of the necessary modelling
276 processes in a river basin, inaccurate results have
277 been obtained. Therefore, by analysing the processes
278 of a river basin characterized by measurable
279 variables, an empirical data-driven model can be
280 constructed. The support vector regression (SVR)
281 machine learning technique was used in this study to
282 analyse historical stream flow occurrences in order to

283 predict stream flow values. Predictions for up to
284 twelve months were made and the coefficient of
285 determination as well as the root-mean-square error
286 were used as accuracy measurements. Satisfactory
287 results were obtained and local stream flow data
288 proved to be a trustworthy hydrological factor when
289 predicting a specific river's stream flow. Even though
290 the effects of precipitation may already be present in
291 stream flow data, an understanding of the relationship
292 between stream flow and precipitation may lead to a
293 more accurate prediction of stream flow. Explicitly
294 including precipitation and other environmental
295 aspects such as temperature and evaporation when
296 building an SVR model will therefore be addressed in
297 further studies.
298

299 **5. References**

300
301 Callegari, M., Mazzoli, P., De Gregorio, L.,
302 Notarnicola, C., Pasolli, L., Petitta, M., Pistocchi,
303 A. (2015). Seasonal river discharge forecasting
304 using support vector regression: a case study in the
305 Italian Alps. *Water*. 7: 2494-2515.
306
307 Cortes, C., Vapnik, V. (1995). Support-vector
308 networks. *Machine Learning*. 20: 273-297.
309
310 Falconer, R., Lin, B., Harpin, R. (2005).
311 Environmental modelling in river basin
312 management. *International Journal of River Basin
313 Management*. 3: 169-184.
314
315 Lin, J., Cheng, C., Chau, K. 2006. Using support
316 vector machines for long-term discharge
317 prediction. *Hydrological Sciences*. 51(4): 599-612.
318
319 Miller, J.R., Russel, G.L., Caliri, G. (1993).
320 Continental-scale river flow in climate models.
321 *Journal of Climate*. 7: 914-928.
322
323 Raghavendra, S., Deka, P.C. (2014). Support vector
324 machine applications in the field of hydrology: a
325 review. *Applied Soft Computing*. 19: 372-386.
326
327 Solomatine, D., Ostfeld, A. (2008). Data-driven
328 modelling: some past experiences and new
329 approaches. *Journal of Hydroinformatics*. 10: 3-22.
330
331 Veiga, V.B., Hassan, Q.K., He, J. (2015).
332 Development of flow forecasting models in the
333 Bow River at Calgary, Alberta, Canada. *Water*. 7:
334 99-115.