**JRBM**

# Early Stopping Technique Using a Ga for Calibration in a Storm Runoff Model

SCHOLARONE™
Manuscripts

# EARLY STOPPING TECHNIQUE USING A GENETIC ALGORITHM

# FOR CALIBRATION OF AN URBAN RUNOFF MODEL

PHUONG CU THI[a], JAMES E BALL[b*], NGOC HUNG DAO[c]

[a] *Faculty of Hydrology and Water Resources, Thuy Loi University, Tay Son Street, Hanoi, Vietnam. E-mail: phuongcuthi@yahoo.com*

[b] *School of Civil and Environmental Engineering, University of Technology Sydney, Broadway, NSW, 2007, Australia. E-mail: james.ball@uts.edu.au*

[c] *Faculty of Physical Geography, National Education University, Hanoi, Vietnam. E-mail: daongochung69@gmail.com*

**ABSTRACT**

Identifying suitable parameter sets for use in catchment modelling remains a critical issue in hydrology. This paper describes an early stopping technique (EST) for use during calibration of a multi-parameter urban catchment modelling system. The proposed method takes advantage of MODE and lower confidence limit (LCL) functions in statistical analysis of spanning set of objective function values. The paper also introduces a monitoring process and regularization techniques to avoid under/overfitting during the calibration and to enhance generalisation performance. The methodology is assessed using SWMM and linked with a Genetic Algorithm for calibration of a Powells Creek catchment model in Sydney, Australia. Results demonstrate that the statistical spanning set analysis approach overcomes issues of poor interpretation and deterioration in the model's generalisation properties. By stopping early, the calibration process avoided overfitting; this was indicated by too closely fitting to the calibration dataset and a failure to fit to the monitoring dataset.

Keywords: catchment model, early stopping, SWMM, flood estimation, calibration

## INTRODUCTION

Management of storm runoff in urban areas remains a serious problem for many water facility managers. The flow data needed for this management generally are obtained either from model studies or field monitoring. Model-based studies have become very popular for catchments where there are no or insufficient flow data for management purposes. However, the reliability of model predictions depends on how well a model replicates the catchment response; in other words, the reliability of model predictions depends on the selection of suitable model parameter values.

The process of selecting parameter values is referred to as model calibration and involves the minimizing of an objective function over a search space; the search space represents the ranges of possible values for the parameters. A number of studies (Kuczera et al., 2006; Jin et al., 2010; Refsgaard and Storm, 1996; Cu and Ball, 2016) have discussed the potential for errors in the input and field data to influence the selection of parameter values resulting in over or under-fitting of the model.

Underfitting occurs when a model cannot adequately capture the structure of the data and therefore, suitable parameter sets are unable to be found. Overfitting, on the other hand, occurs when parameter set(s) are fitted too closely, or exactly, to a particular dataset and, therefore, fail to fit to data from other events, or fails to reliably predict future observations (Choi and Ball, 2002a).

There are alternative approaches to catchment simulation; these alternatives can be classified as regression based (pattern recognition) and process based. The issue of over and under-fitting of simulation models has received greatest attention in those models based on Artificial Neural Networks (ANN) which are a form of regression model (e.g. Coulibaly et al. 2000 and Piotrowski and Napiorkowski, 2013). While a number of different approaches have been attempted, Piotrowski and Napiorkowski (2013) reported that an Early Stopping Technique (EST) can be considered a good strategy to avoid over-fitting a catchment model.

53    The idea of an early stopping technique was first introduced by Nelson and Illingworth

54    (1991), who wanted to avoid the problem of overfitting in large feed-forward ANNs.

55    Applications of this technique appear in systems affected by noise or fitting signals (Zhang and

56    Yu, 2005; Mahsereci et al., 2017).  Furthermore, the technique has become one of the more

57    common approaches applied to deal with overfitting issues in ANN models due to the good

58    model performance and the quicker convergence (Piotrowski and Nappiorkowski, 2013).

59    Finally, Coulibaly et al. (2000) states that, in an ANN application to a complex catchment

60    system, EST makes it possible to avoid overfitting since the calibration stops as soon as the

61    criteria are met.

62        While ANN based approaches have been found to provide a useful approach to

63    modelling of catchment response to storm events, the more commonly used approaches would

64    be classified as process based.  Similar to the ANN based approaches, over-fitting of process

65    models during calibration is possible.  However, the focus during calibration of these catchment

66    models has been parameter uncertainty (Beven and Binley, 1992; Refsgaard and Storm, 1996;

67    Kuczera et al., 2006; Fang and Ball, 2007; Cu and Ball, 2016) rather than assessing whether the

68    model has been over-fitted.  These studies have demonstrated that, rather than a single

69    parameter value, or set of parameter values, being suitable for replication of catchment

70    response, there are multiple alternative parameter values, or sets of parameter values capable

71    of replicating recorded catchment response to one or more storm events.  As the number of

72    parameters in a distributed catchment model usually is greater than the available data for their

73    estimation, the potential for over-fitting of the parameter value pdfs has not been addressed in

74    these studies.  This poses the question of generalization of the pdfs of the parameter values.

75        Choi and Ball (2002b), however, considered potential over-fitting of a process-based

76    catchment model through application of an EST in the search for a single optimal parameter set

77    of values.  In their study, optimal was defined as the set of parameter values that had the best

78    generalization; in other words, the set of parameter values that resulted in the best replication

79 over a range of events. An EST was employed to define when calibration of parameter values

80 should cease; validation of parameter values was attempted then using independent events. Of

81 the validation events considered, it was found that use of the EST resulted in the minimum

82 difference between the predicted and field data occurred 81% of the time at an *early stop point*

83 rather than after the specified number of calibration iterations had been completed. In other

84 words, the EST approach avoided over-fitting the parameter values in 81% of the models tested.

85 The study of Choi and Ball (2002b), however, did not consider parameter uncertainty. Hence,

86 there is a need to investigate if an EST approach should be used when estimating parameter

87 uncertainty; in other words, can the parameter value pdfs be enhanced using an EST during

88 calibration. The investigation reported herein aims at addressing this question. An EST

89 approach is applied during calibration, using a Genetic Algorithm, of an urban flood model

90 applied to the Powells Creek catchment in Sydney, Australia.

91

92 **CASE STUDY CATCHMENT DETAILS**

93 ***General***

94 Description of the case study catchment can be divided into three components, namely the

95 catchment, the modelling software, and the available data. Each of these components will be

96 presented in the following sections.

97

98 ***Catchment***

99 The Powells Creek catchment (see Figure 1), sometimes referred to as the Strathfield

100 catchment, is an 8.41km$^2$ catchment situated 10km west of Sydney's central business district.

101 The drainage network comprises a closed piped system that opens out to a lined channel and

102 then into the Parramatta River. The main open channel was established in 1892 and the closed

103 pipe system was established in the 1920's. From a topographic perspective, the catchment is

104    classified as having gentle slopes between 4% and 6% with a maximum elevation of 40m AHD;

105    the minimum elevation is governed by the tidal regime of the Parramatta River.

106    Between 1958 and 2005, a gauging station on the main Powells Creek Stormwater

107    Channel was operated by UNSW.   The catchment area draining to this gauging station was

108    2.3km$^2$ of the total 8.41km$^2$.  Details of the data collected are described later.

109    **Insert Figure 1 Here**

110

111    ***Model Software***

112    There are numerous alternative software systems suitable for process-based modelling of

113    existing and potential urban catchments.   After considering these alternatives, SWMM

114    (Rossman, 2005) was used for simulation of the Powells Creek catchment.  This model has

115    received extensive application; examples of applications are presented by Sun et al. (2014) and

116    Shahed Behrouz et al. (2020).

117    For purposes of modelling, the catchment upstream of the gauging station was divided

118    into 38 subcatchments based on topographic and drainage system characteristics.  As discussed

119    by Choi and Ball (2002a), these parameters can be categorised arbitrarily as:

120    • Measured parameters.  These are parameters that are physically measured such as pipe

121    diameters, catchment areas, rainfall depth, etc.; and

122    • Inferred parameters.  These are parameters that are not measured and are determined from

123    the application of a model.  Examples of inferred parameters are Manning's roughness for

124    catchment surfaces or channels, depression storage, catchment or subcatchment

125    imperviousness, and the infiltration loss parameters.

126    In discussing inferred parameters (i.e. those parameters more likely to be adjusted

127    during calibration) for a typical application of SWMM, Fang and Ball (2007) noted that there

128    are 10 parameters.  These parameters are the subcatchment width, the impervious percentage

129    of the subcatchment, the percentage of the impervious area with zero depression storage,

130  depression storage of the impervious and previous areas, the Manning's roughness for

131  impervious and previous areas, and three infiltration parameters for Horton's infiltration

132  equation.

133

134  *Available Data*

135  As previously noted, a gauging station was in operation during the period 1958 to 2005.  In

136  addition to the flow data, continuous rainfall data was collected at two locations within the

137  gauged portion of the catchment; these locations were at the centroid of the gauged catchment

138  and at the flow gauging station.  While this rainfall data was collected for the same period as

139  the flow data, only rainfall data for the period 1981 to 1998 from the flow gauging station was

140  available for this study.

141       In addition to the flow data, continuous rainfall data was collected at two locations

142  within the gauged portion of the catchment.  While this rainfall data was collected for the same

143  period as the flow data, only rainfall data for the period 1981 to 1998 was available for this

144  study.

145       Details of the events used for the calibration, monitoring, and validation are shown in

146  Table 1 and Figure 2.

147                      **Insert Figure 2 here**

148                      **Insert Table 1 here**

149

150  **METHODOLOGY**

151  **Objective Function**

152  Calibration of a catchment model involves selection of parameter values and the testing of these

153  values using an objective function to assess the model performance.  In this study, a Modified

154  Nash Sutcliffe Efficiency (MNSE) was used as the objective function.  Following Podger

155  (2004), the MNSE can be expressed as

$$MNSE = 1 - \left[ \frac{\sum_n^i \left( (Y_n^{obs})^\lambda - (Y_n^{sim})^\lambda \right)^2}{\sum_n^i \left( (Y_n^{obs})^\lambda - (Y^{mean})^\lambda \right)^2} \right] \tag{1}$$

156    where: $Y_n^{obs}$ and $Y_n^{sim}$ are the observed and simulated flows at the nth increment of the

157    hydrograph, and $Y^{mean}$ is the mean observed flow over the n increments. As prediction of peak

158    flow is the major concern when modelling flood events, a value of $\lambda = 2$ was used to give more

159    weight to high flows as recommended by Podger (2004).

160

**Uncertainty in Modelling Parameters**

162    Numerous previous studies (Beven and Binley, 1992; Refsgaard and Storm, 1996; Kuczera et

163    al., 2006; Fang and Ball, 2007; Cu and Ball, 2016) have investigated parameter uncertainty.

164    These studies have shown that there is a range of possible values for a given parameter and that,

165    for each potential value for a parameter, there is an associated probability.

166         A common problem encountered when modelling individual events in a catchment is

167    the need to ensure generality of selected parameter values.   For this study, this concept is

168    expanded to generality of parameter values and their associated probabilities.

169

**Search algorithm**

171    A Genetic Algorithm (GA) served as the search algorithm to estimate the range and likelihood

172    of parameter values necessary for operation of SWMM.  GAs were first developed by Holland

173    (1975) and are a class of metaheuristic search algorithms that make use of evolutionary ideas

174    to generate solutions for optimization problems.  The principle of "survival" is accomplished

175    by evaluating each candidate's fitness through an appropriate objective function and a biased

176    random selection procedure of individuals for "reproduction", wherein higher rated candidates

177    are more likely to be selected.  This reproduction is undertaken using stochastic transformations

178    inspired by natural evolution, such as inheritance, mutation, selection, and crossover (Goldberg

179    1989).

180            Linking a GA with a catchment modelling system for calibration means that each

181    individual set of parameters is represented by an individual in the generated population of

182    chromosomes. In this case, the chromosomes consist of 342 genes representing the 342 model

183    parameters.  To use the GA, a population of chromosomes will be generated from a random

184    combination of values within the search space; for the generation of these parent chromosomes,

185    the values were assumed to have uniform likelihoods.  Based on the likelihood that use of a

186    chromosome will result in replication of the catchment response, individual chromosomes were

187    selected from the current population to be parents and reproduction; a child chromosome is

188    accepted for the next generation only if the objective function (MNSE) is better than the parent.

189            The reproduction of children is conducted by changing chromosomes based on the

190    following genetic features:

191    (i)      real-value coding - the chromosomes consist of parameter values;

192    (ii)     tournament selection with replacement for chromosomes (Miller et al., 1995); and

193    (iii)    cross-over (with elitism) operations with a random level at uniform probabilities (Fang

194            and Ball, 2007).

195            GAs can be used to produce an infinite number of generations.  However, the

196    under/overfitting problem exists, and hence there is a need to define when production of

197    generations should cease.  An EST approach is used for this purpose.

198            For testing the EST approach in calibration of a model for the Powells Creek catchment,

199    a linkage between a GA and SWMM was developed.  This linkage is illustrated in Figure 3.

200                          **Insert Figure 3 here**

201            Using this linkage, the steps in the calibration process are:

202    1. An initial population of 600 parameter sets (referred to as the parent data sets) are

203            developed.  These parameter sets are uniformly distributed across the search space.

204   2.   These 600 parameter sets are used with SWMM to generate 600 MNSE for the calibration

205        and monitoring datasets.

206   3.   The 300 datasets that provided the best MNSE were selected to produce 300 new datasets,

207        known as children, for the next generation.

208   4.   The 300 children datasets were used with SWMM to generate the MNSE for the new

209        datasets.  This provided a total of 600 parameter sets and MNSE at the new generation.

210   5.   Steps 3 and 4 are repeated until the calibration process ceases.

211   6.   For the purposes of this study, the maximum number of generations considered was 50.

212

213   **Early Stopping Technique**

214   The basic concept of the EST is to split the available data for calibration and validation into

215   three parts: (1) a calibration (or training) set, used to determine the catchment modelling

216   parameters; (2) a monitoring (or testing) set, serving to estimate the generalisation of the

217   catchment modelling system performance and to decide when calibration must cease; and (3) a

218   validation set, used to verify the effectiveness of the stopping criteria and to estimate expected

219   performance in the future.  Typical error curves of calibration and monitoring data are shown

220   in Figure 4.

221                              **Insert Figure 4 here**

222     During the initial phase of the calibration process, the objective function values for both

223     the calibration and monitoring data will decrease. When the objective function values

224     calculated for calibration data decrease, while those for the monitoring data increase, the

225     calibration process has reached a critical point and should terminate. Further calibration is

226     likely to result in the parameters being fitted to errors present in the data, rather than to the

227     reduction of prediction errors arising from inappropriate parameter values, in other words –

228     overfitting occurs.

229

230     **Early Stopping Technique Application**

231     A novel aspect of the research reported herein is analysis of the objective function for each

232     generation as part of an EST. To achieve this aim, regularization techniques were applied to

233     the vector of MNSE values to find convergence points of the generalization process and,

234     consequently, identify stopping points for the calibration process. Use of two statistical

235     characteristics, i.e. the mode and the confidence limits of the MNSE distribution, as a

236     regularization technique is proposed. The use of these statistical characteristics as part of an

237     EST process follows.

238     Use of the statistical characteristics in an EST process requires the following definitions:

239    • The possible values of MNSE are the space (R) which a set of real numbers;

240    • A sample of the population in each generation of the GA is defined by I (I = {1,2,3,… }), a

241     set of integer numbers;

242    • The number of a dataset (flood event) is defined by J (J = {1,2,3,… }), a set of integer

243     numbers; and

244    • The generation in the GA is defined by K  (K = {1,2,3,… }), a set of integers.

245     Using these definitions, the MNSE is a vector of values comprising j components at

246  generation k.  The statistical characteristics of the vector MNSE are used in the EST process.

247     The calibration procedure can be evaluated using a span function of MNSE generated

248  from j datasets.   The Span function is a combination of all members from all subspaces which

249  meet objective function constraints and can be expressed by Eq. (2) (modified from Zhang

250  2005):

$$\text{Span(S)} = \{\text{MNSE}^{i,j,k}: \text{MNSE}^{i,j,k} \in R , i \in I, j \in J, k \in K\} \qquad (2)$$

251  where i is the population number in a GA generation, j is the number of datasets, and k is the

252  number of iterations (generations).   Convergence of the calibration process occurs if a

253  concave/convex function (A) of Span S for k generations exists. The early stopping point can

254  be identified by the Suprema/Infimum function A (sup/inf A) of span S and can be expressed

255  in the following form:

$$\text{Stopping point} = \text{sup/inf } A(f); f \in \text{Span(S)} \qquad (3)$$

256

257     The aim of this study is to find a concave function (A) which indicates convergence of

258  the calibration metric (MNSE).  To find this convergence point, two statistical functions were

259  tested, namely MODE and lower confidence limit (LCL) of the vector MNSE.

260

261 **RESULTS AND DISCUSSION**

262 **Objective function analysis**

263 At each generation, the value of the objective function (MNSE) was calculated for each

264 chromosome in the population (300 chromosomes) for 2 events. The two statistical attributes

265 MODE and LCL were determined for analysis of the MNSE distributions to identify suitable

266 stopping points. A moving average method was employed also.

267

268 **Method 1: Trend analysis of MODE values**

269 At each generation k and for j events, MNSEs were fitted by a distribution and MODEs were

270 calculated. The Span X, therefore, consists of MODEs at k generations for the calibration and

271 monitoring events. Presented in Figure 5 are MODEs over 50 generations for the 2 selected

272 events. As can be seen from this figure, MNSE strongly fluctuates during the calibration

273 process. Overall MNSE increased during the first 2 steps from 0.85 to more than 0.90, followed

274 by a downward trend for both validation and monitoring events. However, individual values

275 are largely scattered around these trends. This indicates the absence of convergence in the

276 calibration process and hence stopping points were undefined.

<p align="center">**INSERT FIGURE 5 HERE**</p>

278 To identify a convergent point a moving average method was applied with step sizes of

279 3, 4 and 5. Figure 5 shows evolution of the calibration metric over the generations during the

280 calibration process. During calibration, the maximum function value is obtained virtually from

281 the starting point (after the first step), and meanwhile, the monitoring process produces 3 stages

282 according to different reduction rates. Nevertheless, this process resulted in poor interpretation

283 and hence the early stopping point was undefined.

284

285 **Method 2: Analysis of confidence limits**

286 ==Analysis of the confidence limits is aimed at improving estimation of the stopping point==

287 ==compared to the first method.  At each generation, the 90% LCL of MNSE (i.e. 95% of the==

288 ==MSNE were greater than this value) in the spanning set S for each event was calculated (see==

289 ==Figure 6).==

290 **INSERT FIGURE 6 HERE**

291 Figure 6 illustrates the changes in LCL for the 3 datasets (calibration, validation and

292 monitoring).  The maximum value of the LCL of MNSE was 0.906 occurring at the 15th step.

293 This confirmed that the best stopping point was the 15th generation which resulted in 95% of

294 the population having a MNSE equal to or exceeding 0.906.  The MODE values at step 15 were

295 0.90 and 0.98 in calibration and monitoring, respectively.  The distribution of MNSE at the 15th

296 generation is illustrated in Figure 7.

297 **INSERT FIGURE 7 HERE**

298

299 **Parameter analysis**

300 This section critically examines how the calibration process performed by evaluating the

301 statistical properties of the SWMM model parameter values at different stages.  Based on an

302 analysis of the objective function shown in Figure 6, the fitting process can be divided into 3

303 stages, namely:

304 • Stage 1: From the beginning to generation 15.

305 • Stage 2: From generation 15 (stopping point) to generation 31.

306 • Stage 3: From generation 31, where there is a dramatic drop in the monitoring process

307     performance, to the end of the fitting process (generation 50).

308 Subcatchment weighted average values over the catchment were calculated for each

309 parameter category. These values were fitted by a normal distribution using the closed-form

310 Maximum Likelihood Estimation of Pandey and Nguyen (1999). Statistical parameters (mean

311    and standard deviation) of each parameter at 3 points (first; generation 15; and the last

312    generation) are shown in Table 2 and are illustrated in Figure 8.

313                                    **INSERT FIGURE 8 HERE**

314                                    **INSERT TABLE 2 HERE**

315          Statistical differences in parameter sets at these four critical points were tested by 2

316    criteria; the first test was a Welch Two Sample t-test (Welch 1951), and the second test was an

317    F-Test (Lomax 2007).  The Welch Two Sample t-test is a test of difference in mean between

318    two normal distributed data sets based on assumptions of unequal variance and application of

319    the Welch distribution modification with confidence limits of 90%. The F-test seeks to compare

320    the two variances.  Both these tests use p-values to accept or to reject the null hypothesis.  If

321    the p-value is less than the chosen alpha level, then the null hypothesis is rejected and there is

322    evidence that the data tested do not originate from the same population.  Commonly, an alpha

323    value of 0.05 is accepted (Welch 1951).  An alpha value of 0.05 indicates there is less than 5

324    chances out of a hundred that a sample came from a population where that was not true.  Hence,

325    a p-value of more than 0.05 means that two data sets can be assumed to be similar. Results of

326    the t-test and F-test for each stage are shown in Tables 3 and 4.

327                                    **INSERT TABLE 3 HERE**

328                                    **INSERT TABLE 4 HERE**

329          From inspection of Figure 8 and Tables 3 and 4, it can be seen that many parameters

330    attained stability after several generations. The mean and standard deviations during the 3

331    stages were similar indicated by the p-values in the t-test and the F-test were more than 0.05

332    (Tables 3 and 4).   There was a significant shift in the distributions of two parameter

333    classifications, namely the weighted average catchment width and percentage of impervious

334    area.   These two parameter classifications influence the catchment lag and rainfall loss

335    (impervious areas in SWMM are assumed to have no continuing losses) respectively.  These

336    two parameter categories warrant further analysis.  This analysis is presented in terms of the

337    stages previously noted.

338

*Stage 1: From the starting point to the 15th generation*

340    During this stage of the calibration process, the mean values were moved and the standard

341    deviations were narrowed (less scattered values) (Figure 8b and 8c). This is illustrated by

342    differences in mean and standard deviations between the distributions resulted from the starting

343    point and the 15th generation. The p-values in the t-test and F-test were less than 0.05 (Tables

344    3 and 4) and ratios of variances were 1.99 and 8.13 for catchment width and catchment

345    impervious area, respectively. This proves that the change in parameter values during the

346    calibration process resulted in improvement of the model's performance. This was indicated by

347    robustness in objective function values (MNSE).  The MODE of the MNSE increased from

348    0.839 to more than 0.95 (Figure 6a). The LCL rose from 0.737 to more than 0.906 for the

349    calibration, validation, and monitoring events, respectively (see Figure 7).

350

*Stage 2: From the 15th generation to the 31st generation*

352    As the calibration process continued there was a shift in the parameter means. Use of an F-test

353    showed significant differences in variance; this is illustrated by a narrowing in the parameter

354    distributions (Figure 8c, 8c).  The p-values in the F-test were less than 0.05 and the ratios of

355    variances increased to 4.01 for the impervious area parameter, 2.21 and 2.36 for width and

356    pervious area depression storage, respectively. This indicated that changes had occurred in

357    these sensitive parameter distributions.

358            However, there was a slight improvement in the model's calibration performance. The

359    MODE fluctuated at around 0.9 (Figure 6a) and a decline in the LCL was observed (Figure 7).

360    This signal indicated that the calibration process had started to fit the noise/error of the outputs

361    rather than the model structure. As a result, despite the calibration process slightly improving

362    there was, however, a decline in validation and monitoring performance.

363

364    *Stage 3: From the 31st generation to the end*

365    In this stage, the parameter distributions were similar as illustrated by p-values in the t-test and

366    F-test more than 0.05 (Tables 5 and 6). Furthermore, the distributions did have a similar shape

367    (Figure 8b, 8c). The variances were slightly different with the ratio being around 1.0. This

368    similarity suggested that the fitting process in fact presented poor performance in terms of

369    refining the model parameter values (no significant change in parameters). However, it

370    continued to fit the flow errors/noise. As a result, the process only slightly improved the

371    calibration and validation. It, however, failed to fit the monitoring dataset. In this case, the

372    MODE in the monitoring event dramatically dropped to 0.71, and the LCL during monitoring

373    fell to 0.723. This can be deemed an example of poor performance.

374                          **INSERT TABLE 5 HERE**

375                          **INSERT TABLE 6 HERE**

376          From the above analysis, it can be concluded that the first stage was the best calibration

377    stage while the best point at which to stop calibration was the 15th generation.

378
379    **CONCLUSION**

380    The focus of this study was the testing of an EST applied during the calibration process of a

381    catchment modelling system.  Of particular concern was the feasibility of using an EST during

382    the estimation of the parameter uncertainty.  The EST approach was applied successfully to

383    modelling an urban catchment for flood estimation.  Furthermore, the EST approach overcame

384    the problem of overfitting in parameter estimation.

385          The EST approach applied in this study was based on statistical analysis of the objective

386    function (MNSE) at each generation of the GA.  The two statistical characteristics considered

387    were the Mode and LCL.  It was found that the LCL allowed easier identification of a stopping

388    point and the point where values of the parameters started losing generalisation through over-

389    fitting.

390    While the approach was found to be viable, the analysis undertaken required a number

391    of generations of the GA to be completed.  Hence, the stopping points could be identified by

392    analysis of the MNSE only after the simulation process had been completed; the analyses

393    undertaken relied on post-processing data from the GA and not processing the data during the

394    GA.  Further development of the approach is needed to mitigate this issue.

395

402

**REFERENCES**

404    Beven, K. and Binley, A., (1992), The future of distributed models: Model calibration and

405    uncertainty prediction, Hydrological Processes, 6:279-298. doi:10.1002/hyp.3360060305

406    Choi, K.-S. and Ball, J.E., (2002a), Parameter estimation for urban runoff modelling, *Urban

407    Water*, 4(1):31-41.

408    Choi, K.-S. and Ball, J.E., (2002b), A Generic Calibration Approach: Monitoring the

409    Calibration, *In Proc. Hydrology and Water Resources Symposium 2002*, Melbourne,

410    Australia, Pub Institution of Engineers Australia, Barton, A.C.T 650-657.

411    Coulibaly, P, Anctil, F, and Bobee, B, (2000), Daily reservoir inflow forecasting using artificial

412    neural networks with stopped training approach, *Journal of Hydrology*, 230:244–257.

413    Cu, P. and Ball, J.E., (2016), The influence of the calibration metric on design flood estimation

414        using continuous simulation, *International Journal of River Basin Management* 15(1):9-

415        20.

416    Fang, T. and Ball, J.E., (2007), Evaluation of spatially variable control parameters in a complex

417        catchment modelling system: A genetic algorithm application, *Journal of*

418        *Hydroinformatics*, 9:163-173.

419    Goldberg, D.E., (1989), *Genetic algorithms in Search, Optimisation and Machine Learning*,

420        Addison-Wesley Publishing Co. Inc., Reading, Mass., USA

421    Holland, J. H., (1975), Adaptation in Natural and Artificial Systems, University of Michigan

422        Press, Ann Arbor, MI, USA

423    Jin, X., Xu, C.-Y., Zhang, Q. and Singh, V.P., (2010), Parameter and modeling uncertainty

424        simulated by GLUE and a formal Bayesian method for a conceptual hydrological model,

425        *Journal of Hydrology*, 383:147-155.

426    Kuczera, G., Kavetski, D., Franks, S. and Thyer, M., (2006), Towards a Bayesian total error

427        analysis of conceptual rainfall-runoff models: Characterising model error using storm-

428        dependent parameters, *Journal of Hydrology*, 331:161-177.

429    Mahsereci M., Balles L., Lassner C., and Hennig P., (2017), Early Stopping without a

430        Validation Set, CoRR abs/1703.09580(1703.09580).

431    Nelson, M.C. and Illingworth, W.T., (1991), *A Practical Guide to Neural Nets*, Addison-

432        Wesley, Reading, MA, USA

433    Piotrowski, A. P. and Napiorkowski, J.J., (2013), A comparison of methods to avoid overfitting

434        in neural networks training in the case of catchment runoff modelling, *Journal of*

435        *Hydrology*, 476:97-111.

436    Podger, G., (2004), Rainfall Runoff Library User Guide, Unpublished Report, CRC for

437        Catchment Hydrology, Australia.

438    Refsgaard, J.C. and Storm, B., (1996), Construction, calibration and validation of hydrological

439         models, In Abbott, M.B., Refsgaard, J.C. (Eds.), *Distributed Hydrological Modelling*,

440         Kluwer Academic, 41-54.

441    Rossman, L., (2015), Storm Water Management Model User's Manual Version 5.1 – Manual,

442         US EPA Office of Research and Development, Washington, DC, *EPA/600/R-14/413*

443         (NTIS EPA/600/R-14/413b).

444    Shahed Behrouz, M., Zhu, Z., Matott, L.S. and Rabideau, A.J., (2020), A new tool for automatic

445         calibration of the Storm Water Management Model (SWMM), *Journal of Hydrology*, 581,

446         doi.org/10.1016/j.jhydrol.2019.124436

447    Sun, N., Hall, M., Hong, B. and Zhang, L.J., (2014), Impact of SWMM Catchment

448         Discretization: Case Study in Syracuse, New York, *Journal of Hydrologic Engineering*,

449         19(1):223-234.

450    Welch, B. L., (1951), On the Comparison of Several Mean Values: An Alternative Approach,

451         *Biometrika*, 38:330–336. doi:10.2307/2332579.

452    Zhang, T. and B. Yu., (2005), Boosting with Early Stopping: Convergence and Consistency,

453         *Annuals of Statistics*, 33(4):1538-1579.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

| No. | Comment | Response |
|---|---|---|
| **Reviewer 1** | | |
| 1 | Line 1, Title, it is better to use "Genetic Algorithm" instead of "GA" | Agreed |
| 2 | Line 12 "this paper descried" better "describes". Same line 15 "introduces" | Agreed |
| 3 | Figure 1. Replace y label "No of Function evaluations" by "Number (or No.) of function evaluations | Figure caption corrected to provide correct citation. Figure is copied from source. |
| 4 | Figure 2. Improve this figure, it looks like it was cut from the reference not properly | No longer relevant |
| 5 | Please check your Introduction, you have a mix of past and present tense in the way you cite references. | The abstract has been revised and corrected. |
| 6 | Line 112, Please introduce MNIST dataset for those who are not familiar with it | No longer relevant |
| 7 | Line 117, which three techniques? | Clarified in new version of manuscript |
| 8 | Line 122, you already said this in the previous line 119. | No longer relevant |
| 9 | Line 131, what is the most relevant part of Nguyen's work (2017)? | Nguyen reference deleted |
| 10 | Lines 133-138, A table should serve better to introduce the model's parameters. The table has to shows parameter, a brief description, and range of values that parameters can take. | Rewritten to cite references where the suggested approach has been used |
| 11 | Lines 145-147 Please introduce these events and mention the most relevant of them. | Agreed. Section rewritten. |
| 12 | The Methodology is not clear. The description of Methodology section needs considerable improvement. The steps followed should be mentioned; perhaps a scheme could help you to do this. Why do you only use Modifier Nash Sutcliffe coefficient (MNSE) as the objective function? | Methodology clarified in a rewrite of this section of the manuscript. |
| 13 | Please improve the Application section, this is not a technical report. Any other researcher interested in replicating your work must have all the elements to do it. | Methodology clarified in a rewrite of this section of the manuscript. |
| 14 | In the results and discussion, I think the following elements should be discussed: -Among the methods to avoid overfitting, early stopping is the most used in practice, according to Piotrowski and Napiorkowski (2013). Although the manuscript does not mention or address other techniques, what are the main advantages of the proposed technique over others? -What are the main drawbacks of applying the technique presented in the paper? -What are the main difference and advantages of the technique proposed in the paper with the called optimized approximation algorithm introduced in Liu et al. (2008)? | Results and discussion sections have been modified to clarify the important aspect. |
| 15 | In your references, the most recent is from 2018 (only one), being the topic that you | Literature review has been revised, 2020 references now included. |

| | | | |
|---|---|---|---|
| | | present in the paper of high relevance, there are no more recent references in this regard? | |
| | 16 | Line 386, What was the purpose of your research? | Clarified in the rewritten introduction |
| | 17 | Line 387, To whom or what does "Its" refers? | Clarified in the rewritten introduction. |
| | 18 | Lines 389, What methods? Did you make a comparison between different methods? Why do you mention this? | Clarified in the new manuscript |
| **Reviewer 2** | | | |
| | 1 | Language has to be thoroughly revised. At the moment it does not real well. The literature review also misses some key papers on hydrological model calibration, which are essential to your discussion. Also, the figures require significant rework for clarity. | Language has been checked and major changes made. Figures have been reworked for clarity. |
| | 2 | Model calibration is not an issue about performance, but rather to estimate the un-observable parameters of a system. This is conceptual difference between regression models and conceptual/physically-based models. The issue of overfitting refers to models whose complexity (namely number of parameters) is larger than what the available data supports. This is the reason why there are hardly reported "overfitting" issues in conceptual model calibration, as you accurately report on your literature review. | While the Reviewer makes a valid point, the practice of catchment modelling commonly requires the use of complex models where the relationship between the model process and real process lies only in the terminology. For example. A 2D model of flow on a road surface will require different values of Manning's 'n' depending on the alignment of the grid with the road – the model and the real world energy losses must be similar if the flow propagation is to be reliably simulated. However, the flow lengths will differ between the real world and the model depending on the grid orientation leading to the need for differing values of the same parameter name. Consequently, there is a need to investigate this practical problem. |
| | 3 | I think the literature review should be more profound, by having more references to model calibration issues, as data-driven models (as is a big chunk of your literature) does have over/under fitting issues as parameters are added in function of the available dataset. In conceptually and physically-based hydrological models, the number of parameters depend on the process and its approximation and not on the data availability, thus the same model does not increase in complexity depending on the available data. The issue with conceptual/physically-based type of models is usually in the parameter identifiability (due to equifinality), as its number is usually fixed. In other words, the number of model parameters in a conceptual model does not change once more data is available. | I disagree with the reviewer and particularly with respect to the definition of a model. The model in the context of this manuscript comprises the modeller's concept of hydrologic and hydraulic flow paths, the available data, and the processes embedded in the software. In this context, a change in data does force a change in the model – for example, a new DEM may change the subcatchments due to greater resolution. As the parameters are spatially distributed and functions of land-use, etc, a change in DEM may change parameter values.<br><br>However, the Reviewer is correct in that equifinality is an issue in identification of parameter values. In this paper, this issue is addressed through consideration ensembles of parameter values. As the desire is to use the ensemble for prediction, there is a need to ensure generalisation of the adopted parameter values is not comprised due to systematic and random errors in the model system. |
| | 4 | In addition, you are using data from only 1 storm for calibrating the model. In my experience, this is quite little data to draw conclusions, considering the size of the | The Reviewer is correct that more event data would enhance the paper. However, the problem being dealt with is one of flood prediction. Hence, a greater number of events |

| | | problem you are dealing (over 300 parameters). I would recommend using a case study with more data availability. | would require a longer period of monitoring. As described in the paper, only a 20 year period of data was available for the study.  Hence, in that period the expectation of a 1%AEP event is less than 0.2 and for a 5%AEP is approx. 0.6. As a result, meaningful events are rare. A similar result would exist for any other urban catchment – considering a longer period has the potential for non-stationarity in catchment development to influence model parameters. |
| --- | --- | --- | --- |
| | 5 | Also, the methodology can be clearer, as it is really hard to follow. For example, in l212 I is defined as the population number in each generation of GA, while in l221 i the number of GA generations and, later in l299 each generations is defined by k. Similarly, you present the index J, which indicates the number of datasets, while I suppose you should stick to calibration set to generate the new parameterisation, validation to tests its performance, and testing (or monitoring if I understand correctly) to evaluate the performance on unseen data (which should not be used to retro-fit the model as it mimics "operational conditions"). | Agreed.  Rewritten to clarify this point. |
| | 6 | l14 There is not indication of what is MODE or LCL | MODE and LCL are defined in the revised version – note that these two characteristics are standard statistical parameters. |
| | 7 | l31-32 I don't think the only reason for models is lack of data. Actually you do require data to build models | Agreed.  Rewritten to clarify. |
| | 8 | l151 why in particular using MNSE instead of simply NSE? Is it the intention to get the peak value right, or is it important to understand the dynamics of the system as well? If so, perhaps a MINMAX performance function would be more informative | MNSE was used to bias the metric to higher flows rather than the more frequent low flows that were not of interest in this study.  This has been clarified in the text. |
| | 9 | l175-196 Please avoid using metaphors (such as 'survival', 'reproduction', 'chromosomes', 'genes') to describe GA | These terms are standard in the GA literature. However, the use of these terms has been minimised. |
| | 10 | l188 Are there 342 model parameters using only 3 storms for the model calibration? | Yes, a common problem in practice. |
| | 11 | l197 it is named a metaheuristic | corrected |
| | 12 | l227-231 Is quite hard to read but seems to be essential in describing the methodology. I have two interpretations of this 1) you will find a set of parameters that yield a value of 1, thus indicating the "convergence" of the calibration process. 2) Find a concave function "A" that contains said value (MNSE = 1). In both cases, having a value of 1 in the calibration set is not overfitting already? | This section has been clarified in the revised text. |
| | 13 | l260 Why an initial population of 600? | Previous knowledge – see Fang and Ball (2007). |
| | 14 | l267 How did the parent generate the children? | Details are beyond the scope of the current paper but some clarification has been added. |
| | 15 | l269 How did the children and parents were combined? | Details are beyond the scope of the current paper but some clarification has been added. |

| 16 | l273 Why to stop at 50 generations? | The point was identification of an EST, 50 generations were considered sufficient for that purpose. |
|---|---|---|
| 17 | Table 2 Why to calibrate percentage of impervious area? This is an observable parameter | I disagree.  The model assumes effective impervious area – 100% of rainfall becomes runoff from the effective impervious area, not the total impervious area.  This has been discussed extensively in the urban drainage literature.  Furthermore, impervious areas need not be effective impervious areas as some runoff from impervious areas may flow onto pervious areas and not contribute to the subcatchment runoff. |
| 18 | Figures require considerable rework | Figures have been reworked. |
| 19 | Figure 5 Why to use NSE instead of MNSE? | MNSE was used as it is not a traditional NSE approach. |
| 20 | The notation of the equations can be greatly simplified | This has been considered with additional text to clarify the equations. |

### *List of Tables*

Table 1: Details of selected events for testing the model

| Events | Start date | End date | Peak flow (m$^3$/s) |
|---|---|---|---|
| Calibration | 13:05 23/10/1985 | 17:35 23/10/1985 | 11.89 |
| Monitoring | 08:02 07/10/1997 | 15:00 17/10/1997 | 5.93 |
| Validation | 08:02 29/06/1997 | 13:57 29/06/1997 | 6.59 |

Table 2: Statistical parameters of the model parameters

| Parameter | 1st generation | | 15th generation | | 50th generation | |
|---|---|---|---|---|---|---|
| | Mean | STD | Mean | STD | Mean | STD |
| Width (m) | 739.7 | 65.00 | 566.4 | 46.03 | 439.9 | 28.16 |
| Percentage of impervious area (%) | 43.7 | 0.72 | 44.4 | 0.38 | 44.6 | 0.39 |
| Average Catchment Roughness | 0.039 | 0.0036 | 0.039 | 0.0032 | 0.040 | 0.0032 |
| Imp. area depression storage (mm) | 3.04 | 0.28 | 3.15 | 0.29 | 3.40 | 0.22 |
| Pervious area depression storage (mm) | 28.9 | 2.74 | 29.3 | 2.87 | 28.4 | 2.16 |
| Max. Infiltration rate (mm/hr.) | 362.8 | 26.18 | 367.1 | 28.10 | 345.3 | 22.12 |
| Min. Infiltration rate (mm/hr) | 21.3 | 1.55 | 21.2 | 1.56 | 20.3 | 1.43 |
| Infiltration decay rate | 7.03 | 0.41 | 6.29 | 0.38 | 6.44 | 0.40 |

Table 3: t-test for the differences in model parameter distribution at selected generations

| Parameters | p-value | | |
|---|---|---|---|
| | 1st -15th generation | 15th -31st generation | 31st-50th generation |
| Width (m) | 2.2E-16 | 2.20E-16 | 2.20E-16 |
| Percentage of impervious area (%) | 4.54E-06 | 2.20E-16 | 0.256 |
| Average Catchment Manning | 0.185 | 1.21E-05 | 0.882 |
| Imp. depression storage (mm) | 2.46E-06 | 7.03E-14 | 7.75E-05 |
| Pervious area depression storage (mm) | 0.073 | 2.20E-16 | 2.20E-16 |
| Max. Infiltration rate (mm/hr.) | 0.055 | 0.000125 | 8.39E-10 |
| Min. Infiltration rate (mm/hr) | 0.172 | 1.38E-05 | 0.004 |
| Infiltration decay rate | 3.384E-13 | 8.66E-12 | 0.150 |

Table 4: F-test for the differences in model parameter distribution at selected generations

| Parameters | Between 1st -15th generation | | Between 15th -31st generation | | Between 31st -50th generation | |
|---|---|---|---|---|---|---|
| | p -value | Ratio | p-value | Ratio | p-value | Ratio |
| Width (m) | 3.56E-09 | 1.99 | 1.25E-11 | 2.21 | 0.104 | 1.21 |
| Percentage of impervious area (%) | 2.20E-16 | 8.13 | 2.20E-16 | 4.01 | 0.188 | 0.86 |
| Average Catchment Manning | 0.066 | 1.24 | 0.048 | 0.80 | 0.045 | 1.26 |
| Impervious area depression storage (mm) | 0.758 | 0.96 | 0.136 | 1.19 | 0.002 | 1.43 |
| Pervious area depression storage (mm) | 0.429 | 0.91 | 2.80E-13 | 2.36 | 0.012 | 0.75 |
| Maximum Infiltration rate (mm/hr.) | 0.222 | 0.87 | 0.981 | 1.00 | 3.56E-05 | 1.62 |
| Minimum infiltration rate (mm/hr.) | 0.924 | 0.99 | 0.843 | 0.98 | 0.087 | 1.22 |
| Decay rate of infiltration in Horton's equation | 0.265 | 1.14 | 0.877 | 0.98 | 0.610 | 0.94 |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**LIST OF FIGURES**

1307x925mm (96 x 96 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

```
                    ┌──────────────────────┐
                    │        Start          │
                    └──────────────────────┘
                               │
                               ▼
      ┌──────────────────────────────────────────────┐
      │   Generate initial population (600 chromosomes) │
      └──────────────────────────────────────────────┘
                               │
                               ▼
      ┌──────────────────────────────────────────────┐
      │        Formulate input data files for SWMM      │ ◄───────┐
      └──────────────────────────────────────────────┘          │
                               │                                  │
                               ▼                                  │
                    ┌──────────────────────┐                      │
                    │      Run SWMM         │                      │
                    └──────────────────────┘                      │
                               │                                  │
                               ▼                                  │
      ┌──────────────────────────────────────────────┐          │
      │  Evaluate NSE to select 300 the best chromosomes│          │
      │            and set as parents                   │          │
      └──────────────────────────────────────────────┘          │
                               │                                  │
                               ▼                                  │
      ┌──────────────────────────────────────────────┐          │
      │ Generate 300 children and combine with 300      │          │
      │  parents to produce 600 chromosomes             │          │
      └──────────────────────────────────────────────┘          │
                               │                                  │
                               ▼                                  │
                    ┌──────────────────────┐        ┌─────────┐   │
                    │   Generation k = 50    │ ──────▶│   No    │───┘
                    └──────────────────────┘        └─────────┘
                               │
                               ▼
                    ┌──────────────────────┐
                    │         Yes           │
                    └──────────────────────┘
                               │
                               ▼
                    ┌──────────────────────┐
                    │        Stop           │
                    └──────────────────────┘
```

**Early Stopping Technique**

Early stopping point at minimum of monitoring data

Monitoring

Calibration

Objective Function Value

No of Function Evaluations

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

(a)

(b)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



a. Roughness

b. Average catchment width (m)

c. Percentage of impervious area

d. Depression storage - impervious area

e. Maximum infiltration capacity

f. Decay rate

Minimug. m infiltration capacity