



This is a repository copy of *Using self-organizing maps to infill missing data in hydro-meteorological time series from the Logone catchment, Lake Chad basin.*

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/165137/>

Version: Accepted Version

Article:

Nkiaka, E., Nawaz, N.R. and Lovett, J.C. (2016) Using self-organizing maps to infill missing data in hydro-meteorological time series from the Logone catchment, Lake Chad basin. *Environmental Monitoring and Assessment*, 188 (7). 400. ISSN 0167-6369

<https://doi.org/10.1007/s10661-016-5385-1>

This is a post-peer-review, pre-copyedit version of an article published in *Environmental Monitoring and Assessment*. The final authenticated version is available online at:
<https://doi.org/10.1007/s10661-016-5385-1>

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 **Nkiaka, E., Nawaz, N.R. & Lovett, J.C. 2016. Using self-organizing maps to infill missing**
2 **data in hydro-meteorological time series from the Logone catchment, Lake Chad basin.**
3 **Environmental Monitoring and Assessment.**
4 **DOI 10.1007/s10661-016-5385-1**

5
6 **Using Self-Organizing Maps to infill missing data in hydro-meteorological time series from**
7 **the Logone catchment, Lake Chad basin**

8
9 E. Nkiaka¹, N.R. Nawaz¹ and J.C. Lovett¹

10 ¹School of Geography, University of Leeds:

11 Correspondence to: E. Nkiaka (gyenan@leeds.ac.uk)

12
13 **Abstract**

14 Hydro-meteorological data is an important asset that can enhance management of water resources.
15 But existing data often contains gaps, leading to uncertainties and so compromising their use.
16 Although many methods exist for infilling data gaps in hydro-meteorological time series, many of
17 these methods require inputs from neighbouring stations, which are often not available, while other
18 methods are computationally demanding. Computing techniques such Artificial Intelligence can
19 be used to address this challenge. Self-Organizing Maps (SOMs), which are a type of Artificial
20 Neural Network, was used for infilling gaps in a hydro-meteorological time series in a Sudano-
21 Sahel catchment. The coefficients of determination obtained were all above 0.75 and 0.65 while
22 the average topographic error was 0.008 and 0.02 for rainfall and river discharge time series
23 respectively. These results further indicate that SOMs are a robust and efficient method for infilling
24 missing gaps in hydro-meteorological time series.

25
26 **Keywords:** Artificial Neural Networks, hydro-meteorological data, infilling missing data, Logone
27 catchment, Self-Organizing Maps,

28
29 **1) Introduction**

30 Economic progress, rising standard of living, growing populations and expansion of
31 commercial agriculture in developing countries is putting increasing pressure on fresh water
32 resources (WWAP, 2015). At the same time climate extremes such as droughts and floods are
33 becoming more frequent (Coumou & Rahmstorf, 2012). Better informed water resource
34 management is needed to respond to demand and climate variability. A major requirement for
35 planning is the availability of good quality and long term hydro-meteorological data. This data
36 provides indicators of past hydro-climatic behaviour of a region/catchment and is fundamental to
37 the development of models for prediction of system behaviour (Harvey et al., 2012).

38 Existing hydro-meteorological time series used for planning and management decisions
39 often contains missing observations, particularly in developing countries. The gaps are caused by
40 many reasons, including equipment failure, destruction of equipment by natural catastrophes such
41 as floods, war and civil unrest, mishandling of observed records by personnel or loss of files
42 containing the data in a computer system (Elshorbagy et al., 2000). The presence of gaps, even if
43 there are very short, in a hydro-meteorological time series can hinder calculation of important
44 statistical parameters as data patterns maybe hidden. This can compromise their use for water

45 resources planning as it increases the level of uncertainty in the datasets (Ng and Panu, 2010;
46 Campozano et al., 2014). This problem is particularly acute in the Sudano-Sahel region where
47 rainfall is highly variable in both space and time, meteorological and flow gauging stations are
48 scarce and the available datasets are riddled with gaps.

49 Several methods exist for infilling gaps in hydro-meteorological time series. However, the
50 application of each method depends on a range of factors including the information available for
51 that station; additional datasets from neighbouring stations; the percentage of gaps present within
52 the time series to be infilled; the season within which the gaps are present; the length of the existing
53 data series; and the type of application that the infilled series will be used for (Mwale et al., 2012).
54 These infilling methods range from simple techniques such as linear interpolation, Inverse Distance
55 Weighting (IDW) and Thiessen polygons; to more complicated advanced techniques such as time
56 series models, Markov models, Global Imputation, Multiple Regression models, Artificial
57 Intelligence (Kalteh & Hjorth, 2009; Presti et al., 2010; Ismail et al., 2012; Mwale et al., 2012;
58 Campozano et al., 2014).

59 Most of the methods, require additional input data from neighbouring stations in order to
60 produce reliable results and these additional inputs are often not available. Furthermore, some of
61 the methods are time consuming and demand substantial computer power for simulation because of
62 the complicated algorithms involved (Presti et al., 2010). Some methods also require that the time
63 series be split into different seasons to obtain reliable results. Although these challenges could be
64 overcome by using numerical models (e.g. hydrological models); models also demand high data
65 inputs and cannot be applied to many stations at the same time due to parameter calibration
66 requirements which are site specific and consequently results cannot be transferred to other stations
67 even within the same catchment (Harvey et al., 2012).

68 Some of these challenges can be overcome by using computing techniques such as Artificial
69 Intelligence (AI) (Daniel et al., 2011). In this class of technique, the most promising approaches
70 include Artificial Neural Networks (ANN), Fuzzy Logic (FL) and Genetic Algorithms (GA). The
71 application of Artificial Intelligence in hydrology and water resources management is well
72 established (ASCE 2000; Kingston et al., 2008a, 2008b; Daniel et al., 2011). Among the AI class
73 of models, ANNs are probably the most popular as these use available data to learn about the
74 behaviour of a time series. In addition, they possess capabilities for modelling complex nonlinear
75 systems; do not require prior knowledge of the system process(s) under study and are robust even
76 in the presence of missing observations in the time series (Mwale et al., 2012). The main advantage
77 of ANNs over conventional methods is their ability to model physical processes without the need
78 for detailed information of the system (Daniel et al., 2011); and they have often been used for
79 infilling gaps in hydro-meteorological time series (Kalteh & Hjorth, 2009; Dastorani et al., 2010;
80 Adeloje et al., 2012; Ismail et al., 2012; Mwale et al., 2012; Mwale et al., 2014; Kim et al., 2015).

81 Within the ANN family, the Multilayer Perceptron (MLP) is one of the most widely used
82 for infilling gaps in hydro-meteorological time series (Kalteh & Hjorth, 2009; Dastorani et al.,
83 2010; Mwale et al., 2012; Mwale et al., 2014; Kim et al., 2015). Although MLP is robust for
84 performing this task, it usually demands a long time series for training; and if part of the data to be
85 used for training is missing, additional pre-processing of the time series will have to be carried out
86 to provide estimates in the input space before the training can begin (Rustum & Adeloje, 2007;
87 Mwale et al., 2012). This therefore limits application in situations where significant portions of the
88 time series to be used for training have incomplete data; or for short time series as the data may

89 not be sufficient for training. It is also computationally intensive and needs additional storage
90 memory (Kalteh et al., 2008).

91 Another member of the class of ANNs known as Self-Organizing Maps (SOMs), which is
92 a competitive and unsupervised ANN, is becoming popular for infilling gaps in hydro-
93 meteorological times series and has been shown to outperform ANNs-MLP (Kalteh & Hjorth,
94 2009; Mwale et al., 2014; Kim et al., 2015). Many studies have successfully applied SOMs for
95 infilling gaps in hydro-meteorological time series with satisfactory results (Kalteh & Hjorth, 2009;
96 Rustum & Adeloje, 2011; Adeloje et al., 2012; Mwale et al., 2012; Mwale et al., 2014; Kim et
97 al., 2015).

98 Self-Organizing Maps (SOMs) were first introduced by Kohonen, (1995, 1997). The
99 success of their application in other research disciplines led to their wide application in water
100 resources processes and systems research especially for data mining, infilling of missing data,
101 estimation and flow forecasting, clustering etc. (Kalteh et al., 2008). This is due to their ability to
102 convert nonlinear statistical relationships between high dimensional data onto a low dimensional
103 display (Ismail et al., 2012). Data points that show similar characteristics are placed closed to each
104 other or clustered together in the output space. This mapping approach does a quasi-preservation
105 of the most important topological and metric relationship of the original data (Rustum & Adeloje,
106 2007). Adeloje et al. (2012) asserted that, the ability of SOMs to cluster data together makes them
107 robust for data mining and infilling datasets with gaps and outliers as the gaps/outliers are replaced
108 by their features in the map. The SOMs algorithm generally executes assigned tasks using an
109 unsupervised and competitive learning approach to discover patterns in the data (Kalteh &
110 Berndtsson, 2007) thus, the whole process is entirely data driven. A SOM is made up of two layers:
111 a multi-dimensional input layer and an output layer. Both layers are fully connected by adjustable
112 weights and the output layer is made up of neurons arranged in a two dimensional grid of nodes
113 (Figure 1). Each neuron in the output layer of the SOM contains exactly the same set of variables
114 contained in the input vectors. Despite its wide application for infilling missing data in many
115 studies around the world, it has rarely been used Africa in general and the Sudano-Sahel region in
116 particular.

117 A Self-Organizing Maps approach was applied to infill missing data in monthly rainfall and
118 daily river discharge time series from January 1950 to December 2007 in the Logone river
119 catchment covering Cameroon, the Central Africa Republic and Chad. Infilling of missing gaps in
120 hydro-meteorological time series usually precedes most hydro-climatic studies (Kashani &
121 Dinpashoh, 2012), and this work is part of an on-going research project to assess the vulnerability
122 of this catchment to drought and flood events under anticipated increased climate variability.

123 The paper is structured as follows: Section 2 describes the data and methodology used in
124 the study. In Section 3 the results obtained are presented and discussed. Section 4 gives a general
125 summary and conclusion of the study.

126

127 **2) Methodology**

128 **2.1) Study area**

129 The Logone catchment is part of the greater Lake Chad basin. It lies between latitude 6°-12°N
130 and longitude 13°-16°E and is a transboundary catchment in the Sudano-Sahel transitional zone in
131 Central Africa with an estimated catchment area of 86,500 km² (Figure 1).

132 The Logone River has its source in Cameroon through the Mbere and Vina Rivers, which flow
 133 from the northeastern slopes of the Adamawa plateau. It is joined in Lai by the Pende River from
 134 the Central Africa Republic and flows from south to north to join the Chari River in Ndjamena
 135 (Chad) and continue flowing in a northward direction before finally emptying into Lake Chad. The
 136 climate in the catchment is characterized by high spatial variability and is dominated by seasonal
 137 changes in the tropical continental air mass (the Harmattan) and the marine equatorial air mass
 138 (monsoon) (Candela et al., 2014).

139
 140 **2.2) Data Sources**

141 Monthly gauge rainfall was obtained from SIEREM (Boyer et al., 2006) available for 18
 142 stations covering the period 1950-2000 while daily river discharge data was obtained from the
 143 Lake Chad Basin Commission (LCBC). Discharge time series are available for the stations of Lai,
 144 Bongor, Katoa and Logone Gana covering the period 1973-1998 for Lai and 1983-2007 for the
 145 rest of the stations.

146
 147 **2.3) Implementation of the SOM Algorithm**

148 A SOM algorithm is implemented in a series of steps.

149 The multi-dimensional input data is first standardized to make sure that very high or low
 150 value variables do not dominate the map. Since SOMs use Euclidian metrics to measure distances
 151 between vectors, standardization gives equal weight to all the input variables (Vesanto et al., 2000).
 152 In this analysis, data was not standardized because rainfall and river discharge time series were
 153 trained separately.

154 The input vector is then chosen at random and presented to each of the individual neurons for
 155 comparison with their weight vectors in order to identify the weight vector most similar to the
 156 presented input vector. The identification uses the Euclidian distance defined as:

$$157 \quad D_i = \sqrt{\sum_{j=1}^n m_j (x_j - w_{ij})^2}; \quad i = 1, 2, 3 \dots M \quad (1)$$

159 Where:
 160 D_i = Euclidian distance between the input vector and the weight vector i ; $x_j = j$ element of the
 161 current vector; $w_{ij} = j$ element of the weight vector I ; n = the dimension of the input vector; $m_j =$
 162 “mask”.

163 When the input vector contains missing elements, the mask is set to zero for such elements and
 164 because of this, the SOM algorithm can conveniently handle missing elements in the input vector.
 165 The neuron whose vector closely matches the input vector (i.e. with D_i minimum) is chosen as the
 166 winning node or best matching unit (BMU).

167 After finding the BMU, the weight vector of the winner neuron is adjusted so that the BMU and
 168 its adjacent neurons move closer to the input vectors in the input space, thereby increasing the
 169 agreement between the input vector and the weight vector. This adjustment is carried out using the
 170 following equation:

$$171 \quad w_t(t + 1) = w_t(t) + \alpha(t)h_{ci}[x(t) - w_t(t)] \quad (2)$$

172
 173
 174

175 Where: w_t = element of the weight vector; t = time; $\alpha(t)$ = learning rate at time t ; $h_{ci}(t)$ =
 176 neighbourhood function centred in the winner unit c at time t .

177 From here, each node in the map develops the ability to recognize input vectors that are similar to
 178 itself. This ability is referred to as self-organizing as no external information is added for this
 179 process to take place. The learning procedure continues until the SOM algorithm converges.
 180 Generally, the learning rate decreases monotonically as the number of iterations increase as shown
 181 by the following equation:
 182

$$183 \quad \alpha(t) = \alpha_0 \left(\frac{0.005}{\alpha_0} \right)^{\frac{t}{T}} \quad (3)$$

184 Where: $\alpha(t)$ = learning rate; α_0 = initial learning rate; T = training length

185 The neighbourhood function used in this analysis is Gaussian centred in the winner unit c ,
 186 calculated as:
 187
 188

$$189 \quad h_{ci}(t) = \exp \left\{ -\frac{\|r_c - r_i\|^2}{[2\sigma^2(t)]} \right\} \quad (4)$$

190 Where:

191 $h_{ci}(t)$ = neighbourhood function centred in the winner unit c at time t ; r_c and r_i = positions of nodes
 192 c and i on the SOM grid; $\sigma(t)$ = neighbourhood radius which also decreases monotonically as the
 193 number of iterations increases.
 194

195 The quality of the trained SOM is measured by the total average quantization error and total
 196 topographic error. The average quantization error is a measure of how good the map fits the input
 197 data (it measures the average distance between each data vector and its Best Matching Unit
 198 (BMU)). The smaller the quantization error, the smaller the average of the distance from the vector
 199 data to the prototypes, meaning that the data vectors are closer to its prototypes; it is a positive real
 200 number with a value close to zero indicating a good fit between the input and the map. The
 201 quantization error is calculated as:
 202

$$203 \quad q_e = \frac{1}{N} \sum_{i=1}^N \|X_i - W_{ic}\| \quad (5)$$

204 Where: q_e = quantization error; N = number of input vectors used to train the map; X_i = i th data
 205 sample or vector; W_c = prototype vector of the best matching unit for X ; $\|\cdot\|$ = denotes the Euclidian
 206 distance.

207 Topographic error measures how well the topology of the data is preserved by the map by
 208 considering the map structure. The lower the topographic error, the better the SOM preserves the
 209 topology of the data. It is a positive real number between 0 and 1 with a value close to 0 indicating
 210 good quality. It is calculated as:

$$211 \quad t_e = \frac{1}{N} \sum_{i=1}^N u(X_i) \quad (6)$$

212 Where: t_e = topographic error; N = number of input vectors used to train the map;

213 u_i = binary integer such that it is equal to 1 if the first and second BMU for X_i are not adjacent
 214 units; otherwise it is zero.

215 Since there is always a trade-off between which of the two can be minimized at the expense of the
 216 other, in this study, effort was focused on reducing the topographic error to ensure that the infilled
 217 values reflect the seasonal trend of the different time series. The coefficient of determination (R^2)
 218 was used to check the quality of the newly generated time series. R^2 gives the proportion of the
 219 variance of one variable that is predictable from the other variable and varies between 0 and 1. R^2
 220 is calculated as:

$$222 \quad R^2 = \frac{\sum_{i=1}^n [(x_i - \bar{x})(y_i - \bar{y})]}{\sum_{i=1}^n [(x_i - \bar{x})^2] \sum_{i=1}^n [(y_i - \bar{y})^2]} \quad (7)$$

223
 224 Where: x_i = the i th observed value; y_i = the i th trained value; \bar{x} = the mean of observed value; \bar{y}
 225 = the mean of the trained value; n = the number of observations.

227 **2.4) Setting of SOM algorithm parameters**

228 According to Gabrielsson & Gabrielsson, (2006), the radius of the SOM should be chosen
 229 wide enough at the beginning of the learning process so that the map can be ordered globally as
 230 the radius decreases monotonically with time. To determine the optimum number of neurons, if M
 231 is the total number of input elements, Garcia and Gonzalez (2004), propose that the number of
 232 neurons in the output can be calculated as:

$$234 \quad N = 5\sqrt{M} \quad (8)$$

235 Where: M = total number of samples and N = the number of neurons.

236 Once N is known, Garcia and Gonzalez, (2004) further propose that the number of rows and
 237 columns of N can be calculated by:

$$239 \quad \frac{l_1}{l_2} = \sqrt{\frac{e1}{e2}} \quad (9)$$

240
 241 Where l_1 and l_2 are the number of rows and columns respectively, $e1$ is the biggest eigenvalue of
 242 the training data set and $e2$ is the second biggest eigenvalue.

243 In the initialization phase of the algorithm, since the learning process involve in the
 244 computation of a feature map is a stochastic process, according to Gabrielsson & Gabrielsson,
 245 (2006) the accuracy of the map depends on the number of iterations executed by the SOM
 246 algorithm. These authors recommend that for good statistical accuracy, the number of iterations
 247 should be at least 500 times the number of network nodes. In this study, the random initialization
 248 option was used as it is recommended for hydrological applications e.g. (Kalteh et al., 2008), while
 249 the default parameters set by the SOM software for map size and lattice (rows and columns) were
 250 adopted that were exactly the same as using equations (8) and (9).

251
 252 The basic steps required to complete the infilling process consists of the following:

- 253 1) Data gathering and normalization: The data to be infilled (e.g. rainfall and discharge time
 254 series) is assembled together and standardized; these are the depleted input vectors;
- 255 2) Training: The depleted input vector (data matrix) is introduced to the iterative training
 256 procedure to form the SOM. At the beginning of the training, weight vectors must be
 257 initialized by using either a random or a linear initialization method. The process of

- 258 comparison and adjustment continues until the optimal number of iteration is reached or
259 the specified error criteria are attained.
- 260 3) Extracting information from the trained SOM: Check all the minimum Euclidian distances
261 and isolate the SOM's BMU for the depleted input vector (i.e. with missing values). The
262 BMU identified in this step is a node of trained SOM and thus has the full complement of
263 the missing values;
 - 264 4) Replacement of missing values: Replace the missing values of the input depleted vector by
265 their corresponding values in BMU identified in step 3 above.
- 266

267 **2.5) Application of SOM**

268 For the application of the SOM algorithm for infilling of missing data in this analysis, a
269 SOM toolbox developed at Helsinki University of Technology Finland
270 (www.cis.hut.fi/projects/somtoolbox/) was used in the Matlab® 2014b environment and a batch
271 training algorithm was adopted. Due to the fact both datasets (rainfall and river discharge) had
272 different time-steps, each of the datasets were trained separately. The data was presented in
273 columns with each column representing measurements from each station. The entries without data
274 were recorded as NaN (Not a Number) to meet Matlab® data entry requirements. To train all the
275 data together in a single simulation, the data entries should overlap such that there is no single
276 day/month for all the stations with no data entry.

277 The stations with the longest period of continuous missing observation were Katoa with
278 1418 consecutive days (01/04/1997-18/025/2001) approximately 4 years and Lai with 1200
279 consecutive days (31/01/1979-15/05/1982), approximately 3 years. Donomanga had the longest
280 period of missing monthly rainfall observations.

281

282 **3) Results and Discussion**

283 Initial simulation results using discharge time series produced an average topographic error
284 of 0.04 and a visual inspection of the time series was carried out to check the seasonal trends.
285 Sporadic cases of numerical instability were noticed especially in portions of the time series with
286 extensive gaps where infilling was done. In some cases, high flow values were observed in the dry
287 season and low flow values observed in the rainy season. This was not logical as periods of high
288 flows could not be followed by a single day of abrupt low flow and vice versa. These values were
289 manually deleted for all the stations and a second simulation was performed using the same initial
290 parameters. After this second simulation, these abnormalities disappeared and the average
291 topographic error reduced to 0.02. Results of the overall performance of the model are shown in
292 Table 1.

293 The results indicate that after the second simulation, the model was able to replicate with
294 high accuracy the trends and flow magnitudes (high and low) in the respective seasons as shown
295 in Figures 3 to 6. This justifies the low value of average topographic error 0.02 and the high values
296 of R^2 . From these results, the newly trained time series were used to infill missing gaps in the
297 different time series in the Logone catchment. The preservation of topology, especially for
298 discharge time series is important because seasonal variation causes high and low flows. The
299 results obtained indicate that this seasonal variation was well preserved across all the gauging
300 stations during the infilling process. In this research more emphasis was put on reducing the
301 topographic error to ensure that the infilled values reflect the seasonal variation of the time series.

302 However, a visual observation of flow hydrographs (Figures 3-6) indicate that, the possibility of
303 errors in the original river discharge time series may not be discounted especially for the Bongor
304 station, and this may have a negative impact on the overall performance of the SOM algorithm in
305 this study.

306 The results obtained for rainfall observations were similar to those obtained for discharge
307 with the lowest R^2 value of 0.76 and average topographic error of 0.008. Although some authors
308 (Kalteh & Berndtsson, 2007; Mwale et al., 2012) have proposed that to the rainfall time series
309 should be trained together according to spatial location to improve the results, this method was not
310 applied in this study because results obtained were judged to be satisfactory. Of the 18 rainfall
311 stations, 10 had R^2 values of 0.90 and above while 7 stations had R^2 values of 0.80 and above with
312 only one station (Donomanga) which has the highest percentage of missing observations having a
313 value of 0.76. However, it was noticed that the performance of the model reflected the spatial
314 location of the stations. For example, apart from Bongor CF with a R^2 of 0.80, all stations located
315 above 10°N had R^2 values above 0.90 while most stations located below this latitude had R^2 values
316 between 0.80-0.90. Since the graphs of the all the 18 rain gauge stations cannot be shown, (Figures
317 7 & 8) are used for illustration. Furthermore, it was observed that the SOM algorithm was able to
318 preserve seasonal variation when infilling missing data in rainfall time series just as it did for
319 discharge.

320 The results also indicate that, although this method is quite robust for infilling gaps in
321 hydro-meteorological time series, it cannot be used for infilling gaps in time series with extended
322 periods of missing observations as model performance starts diminishing. This is logical as in such
323 situations the model does not have sufficient data to learn from, thus cannot correctly replicate the
324 pattern in the data. For example time series of measured discharge at Katoa had 1200 consecutive
325 days of missing observations, which represent 13% of the total data entries, produced an R^2 of 0.65
326 compared to Logone Gana with 97 consecutive days of missing observations with an R^2 of 0.91.
327 This implies that time series with extended periods of missing observations should not be used as
328 the model may infill the missing observations but still fail to replicate the pattern in the data.
329 Although, as shown by Kalteh et al. (2007) and Mwale et al. (2012) this issue can be resolved for
330 rainfall time series by training such time series with data from the same spatial zone, this cannot
331 apply for discharge time series as it is influenced by other catchment characteristics and the river
332 morphology which vary along the river channel.

333 Nevertheless results obtained suggest that SOMs are suitable for infilling gaps in hydro-
334 meteorological time series in Sudano-Sahel catchments. Results obtained from this study are
335 comparable to those obtained by Mwale et al. (2012, 2014) in the Lower Shire Floodplain in
336 Malawi, Kang & Yusuf (2012) in the Kelantan and Damansara river basins in Malaysia and Kim
337 et al. (2015) in the Taehwa watershed in Korea

338 The relationship between discharges measured at various stations along the Logone River
339 is shown in Figure 9. The Unified distance matrix (U-matrix) is a graphical display used to illustrate
340 the clustering of the reference vectors in the SOM, it shows the distance between neighbouring
341 map units. The U-matrix can be seen as several component planes which are stacked together one
342 on top of the other. Component planes can either be coloured or grey shaded in a two dimensional
343 lattice. Light colours indicate areas in which the variables are close to each other in the input space,
344 while dark colours illustrate large distances between variables in the input space. Dark colours can
345 be seen as cluster separators while light colours are clusters themselves. Component planes are

346 therefore, mostly used for visualizing the correlation between the various variables in the SOM
347 since they can give information concerning the spread of values in each component (Gabrielsson
348 & Gabrielsson, 2006).

349 From Figure 9, the relationship between the discharges measured at Bongor, Katoa and
350 Logone Gana is not very discernible. To illustrate that there is no relationship between the
351 discharge time series, Figure 10 shows that the discharges measured at Katoa and Logone Gana
352 gauging stations, which are located downstream of Bongor, are paradoxically lower than discharge
353 measured at Bongor station upstream. This can partly be explained by the fact that during the rainy
354 season when the river overflows its banks, immediately after Bongor station, part of the flow is
355 diverted to fill the Maga dam and part is lost to the floodplains. During the dry season, water is
356 withdrawn from the river without control for various purposes by the inhabitants thus reducing the
357 quantity that eventually reaches Logone Gana station located downstream. This can also be
358 attributed to transmission losses as a result of infiltration to the aquifer through channel bed. Seeber
359 (2013) observed that the discharge recorded at Ndjamen flow gauging station located downstream
360 was lower than that recorded upstream at the Logone Gana station. Candela et al. (2014) reported
361 that a significant proportion of groundwater in the Lake Chad aquifer system was from the Logone
362 River through river and aquifer interactions.

363

364 **4) Conclusion**

365 The main objective of this study was to use Self-Organizing Maps (SOMs) to infill missing
366 gaps in hydro-meteorological time series in the Logone catchment using data from four river
367 discharge and 18 rain gauge stations riddled with gaps

368 The combination of artificial intelligence and human intelligence (to be able to distinguish
369 the seasonal discharge trends, patterns and magnitudes) greatly improved the overall performance
370 of the SOM algorithm in handling missing data. Other advantages of SOMs include: (i) it does not
371 require input data from neighbouring stations; (ii) unlike other ANN methodologies it does not
372 require extra datasets to train the time series; (iii) it is not computationally intensive; and (iv) it
373 does not require extra storage capacity.

374 Results obtained from this study indicate that, the SOMs algorithm is quite robust for infilling
375 gaps in hydro-meteorological time series, though it is not suitable for infilling gaps in time series
376 with extended periods of missing observations as model performance starts diminishing. This
377 methodology can be used by practitioners to enhance the planning and management of water
378 resources in areas where available records are infested with missing observations. Preservation of
379 topology through a good replication of trends and discharge magnitudes in the time series obtained
380 in this study will reduce the data input uncertainty in our future modelling studies in the catchment.

381

382 **Acknowledgements**

383 This research was supported by a Commonwealth Scholarship award to the first author. We
384 are grateful to SIEREM and the Lake Chad Basin Commission for providing the data used in this
385 research.

386

387 **References**

388 Adeloje, A. J., Rustum, R., & Kariyama, I. D. (2012). Neural computing modelling of the reference
389 crop evapotranspiration. *Environmental Modelling & Software*, 29, 61-63.

390
391 Alhoniemi, E., Himberg, J., Parhankangas, J. & Vesanto, J. (2002). SOM Toolbox - Online
392 Documentation.
393
394 ASCE Task Committee on Application of Artificial Neural Networks in Hydrology (2000).
395 Artificial Neural Networks in Hydrology. II: Hydrologic Applications. *Journal of Hydrologic*
396 *Engineering*, 5:2(124), 124-137.
397
398 Boyer, J., Dieulin C., Rouche, N., Cres, A., Servat, E., & Paturel, J. (2006). SIEREM: An
399 Environmental Information System for Water Resources, Vol. IAHS Publication 308. IAHS Press:
400 Wallingford, United Kingdom.
401
402 Campozano, L., Sanchez, E., Aviles, A., & Samaniego, E. (2012). Evaluation of infilling methods
403 for time series of daily precipitation and temperature: The case of the Ecuadorian Andes. *Maskana*,
404 5(1), 99-115.
405
406 Candela, L., Elorza, F. J., Tamoh, K., Jiménez-Martínez, J., & Aureli, A. (2014). Groundwater
407 modelling with limited data sets: the Chari– Logone area (Lake Chad Basin, Chad). *Hydrological*
408 *Processes*, 28, 3714-3727.
409
410 Coumou, D. & Rahmstorf, S. (2012). A decade of weather extremes. *Nature Climate Change*, 2,
411 491-496.
412
413 Daniel, E. B., Camp, J. V., LeBoeuf, E. J., Penrod, J. R., Dobbins, J. P. & Abkowitz, M. D. (2011).
414 Watershed modelling and its applications: A state-of-the-art review. *The Open Hydrology*
415 *Journal*, 5, 26-50.
416
417 Dastorani, M. T., Moghadamnia, A., Piri, J. & Ramirez, M. R. (2010). Application of ANN and
418 ANFIS models for reconstructing missing flow data. *Environmental Monitoring and Assessment*,
419 166(1-4), 421-34.
420
421 Elshorbagy, A. A., Panu, U. S. & Simonovic, S. P. (2000). Group-based estimation of missing
422 hydrological data: Approach and general methodology. *Hydrological Sciences Journal*, 45(6),
423 849-866.
424
425 Gabrielsson, S., & Gabrielsson, S. (2006). The use of Self-Organizing Maps in Recommender
426 Systems: A survey of the Recommender Systems field and a presentation of a State of the Art
427 Highly Interactive Visual Movie Recommender System. Master Thesis, Uppsala University.
428
429 Garcia, H., & Gonzalez, L. (2004). Self-organizing map and clustering for wastewater treatment
430 monitoring. *Engineering Applications of Artificial Intelligence*. 17(3), 215–225.
431

432 Harvey, C. L., Dixon, H. & Hannaford, F. (2012). An appraisal of the performance of data-infilling
433 methods for application to daily mean river flow records in the UK. *Hydrology Research*,
434 43(5) 618-636. DOI: 10.2166/nh.2012.110.
435

436 Ismail, S., Shabri, A., & Samsudin, R. A. (2012). Hybrid model of self-organizing maps and least
437 square support vector machine for river flow forecasting. *Hydrology and Earth System Sciences*,
438 16, 4417-4433.
439

440 Kagoda, P. A., Ndiritu, J., Ntuli, C., & Mwaka, B. (2010). Application of radial basis function
441 neural networks to short-term streamflow forecasting, *Physics and Chemistry of the Earth*, 35, 571-
442 581.
443

444 Kalteh, A. M. & Hjorth, P. (2009). Imputation of missing values in a precipitation–runoff process
445 database. *Hydrology Research*, 40(4), 420-432.
446

447 Kalteh A. M., Hjorth, P., & Berndtsson R. (2008). Review of the self-organizing map (SOM)
448 approach in water resources: Analysis, modelling and application. *Environmental Modelling &*
449 *Software*, 23, 835, – 845.
450

451 Kalteh, A. M., & Berndtsson, R. (2007). Interpolating monthly precipitation by self-organizing
452 map (SOM) and multilayer perceptron (MLP), *Hydrological Sciences Journal*, 2(2), 305-317.
453

454 Kang, H. M., & Yusof, F. (2012). Application of Self-Organizing Map (SOM) in Missing Daily
455 Rainfall Data in Malaysia. *International Journal of Computer Applications*, 48(5).
456

457 Kashani, M. H., & Dinpashoh, Y. (2012). Evaluation of efficiency of different estimation methods
458 for missing climatological data. *Stochastic Environmental Research and Risk Assessment*, 26, 59–
459 71.
460

461 Kim, M., Baek, S., Ligaray, M., Pyo, J., Park, M., & Cho, K. H. (2015). Comparative Studies of
462 Different Imputation Methods for Recovering Streamflow Observation. *Water*, 7, 6847–6860.
463

464 Kingston, G. B., Dandy, G. C. & Maier, H. R. (2008a). AI Techniques for Hydrological Modelling
465 and Water Resources Management. Part 2 - Optimization, in L. N. Robinson (editor). *Water*
466 *Resources Research Progress*, Nova Science Publishers, pp. 67-99.
467

468 Kingston, G. B., Dandy, G. C., & Maier, H. R. (2008b). AI Techniques for Hydrological Modelling
469 and Water Resources Management. Part 1 - Simulation, in L. N. Robinson (editor). *Water*
470 *Resources Research Progress*, Nova Science Publishers, pp. 15-65.
471

472 Kohonen, T. *Self-Organizing Maps*, Springer Series in Information Sciences, vol. 30, Springer,
473 Heidelberg, 1st ed., 1995; 2nd ed., 1997.
474

475 Mwale, F. D., Adeloje A. J., & Rustum R. (2014). Application of self-organising maps and multi-
476 layer perceptron-artificial neural networks for streamflow and water level forecasting in data-poor
477 catchments: the case of the Lower Shire floodplain, Malawi. *Hydrology Research*, 45(6), 838-854.
478

479 Mwale, F. D., Adeloje, A. J., & Rustum R. (2012). Infilling of missing rainfall and streamflow
480 data in the Shire River basin, Malawi – A self-organizing map approach. *Physics and Chemistry
481 of the Earth*, 50-52, 34-43.
482

483 Ng, W. W., & Panu, U. S. (2010). Infilling missing daily precipitation data at multiple sites using
484 the multivariate truncated normal distribution model for weather generation. *Water*, 8 pp.
485

486 Presti, R. L., Barca, E. & Passarella, G. A. (2010). Methodology for treating missing data applied
487 to daily rainfall data in the Candelaro River Basin (Italy). *Environmental Monitoring and
488 Assessment*, 160, 1-22.
489

490 Rustum, R., & Adeloje, A. J. (2007). Replacing Outliers and Missing Values from Activated
491 Sludge Data Using Kohonen Self-Organizing Map. *Journal of Environmental Engineering*, 133(9),
492 909-916.
493

494 Seeber, K. (2013). Consultation of the Lake Chad Basin Commission on Groundwater
495 Management. Project: Sustainable Management of the Lake Chad Basin, BGR No:05-2355.
496

497 Vesantu, J., Himberg, J., Alhoniemi, E & Parhankangas, J. (2000). SOM Toolbox for Matlab 5.
498 Report A57. Helsinki University of Technology, Helsinki, Finland.
499

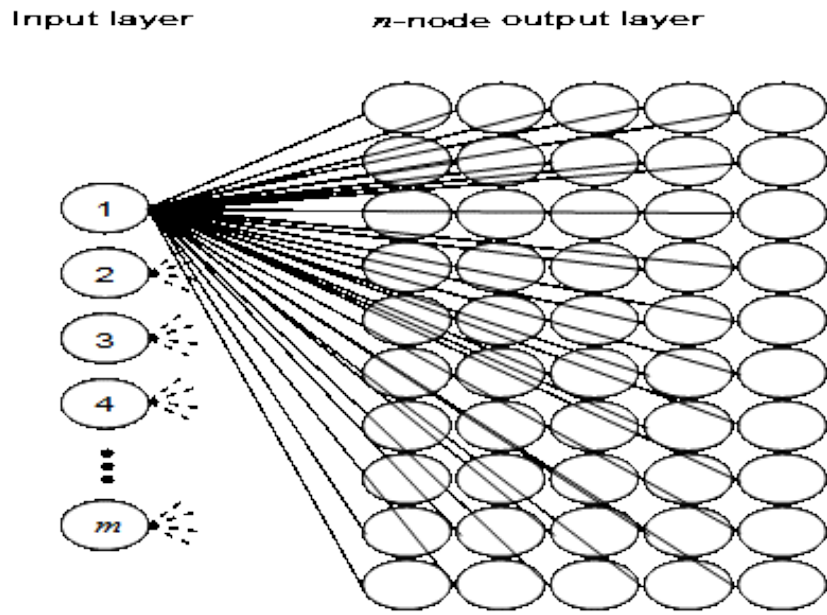
500 WWAP (United Nations World Water Assessment Programme). 2015. The United Nations World
501 Water Development Report 2015: Water for a Sustainable World. Paris, UNESCO
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518

519
520
521
522 Table 1: Station location, percentage of missing data, results of statistical evaluation and average
523 topographic error.

Flow gauging	Latitude	Longitude	Time interval	Proportion of missing data (%)	R²	Average topographic error
Lai	11.55	15.15	1973-1997	17.5	0.85	0.02
Bongor	10.83	15.08	1983-2007	19.2	0.8	
Katoa	10.27	15.42	1983-2007	26.8	0.65	
Logone Gana	9.40	16.30	1983-2007	6.45	0.91	
Rain gauge stations						
Ngaoundere	7.35	13.56	1950-2000	7.52	0.86	0.008
Baibokoum	7.73	15.68	1950-2000	8.82	0.88	
Bekao	7.92	16.07	1950-2000	5.88	0.90	
Pandzangue	8.10	15.82	1950-2000	14.2	0.81	
Donia	8.30	16.42	1950-2000	12.9	0.84	
Moundou	8.57	16.08	1950-2000	5.39	0.94	
Doba	8.65	16.85	1950-2000	4.08	0.94	
Delli	8.72	15.87	1950-2000	5.88	0.91	
Donomanga	9.23	16.92	1950-2000	16.2	0.76	
Guidari CF	9.27	16.67	1950-2000	12.3	0.85	
Goundi	9.37	17.37	1950-2000	6.05	0.91	
Kello	9.32	15.80	1950-2000	8.99	0.88	
Lai	9.40	16.30	1950-2000	5.23	0.92	
Bongor	10.27	15.40	1950-2000	10.8	0.80	
Yagoua	10.35	15.25	1950-2000	8.17	0.92	
Bouso	10.48	16.72	1950-2000	6.37	0.93	
Bailli	10.52	16.44	1950-2000	5.23	0.95	
Massenya	11.40	16.17	1950-2000	5.72	0.95	

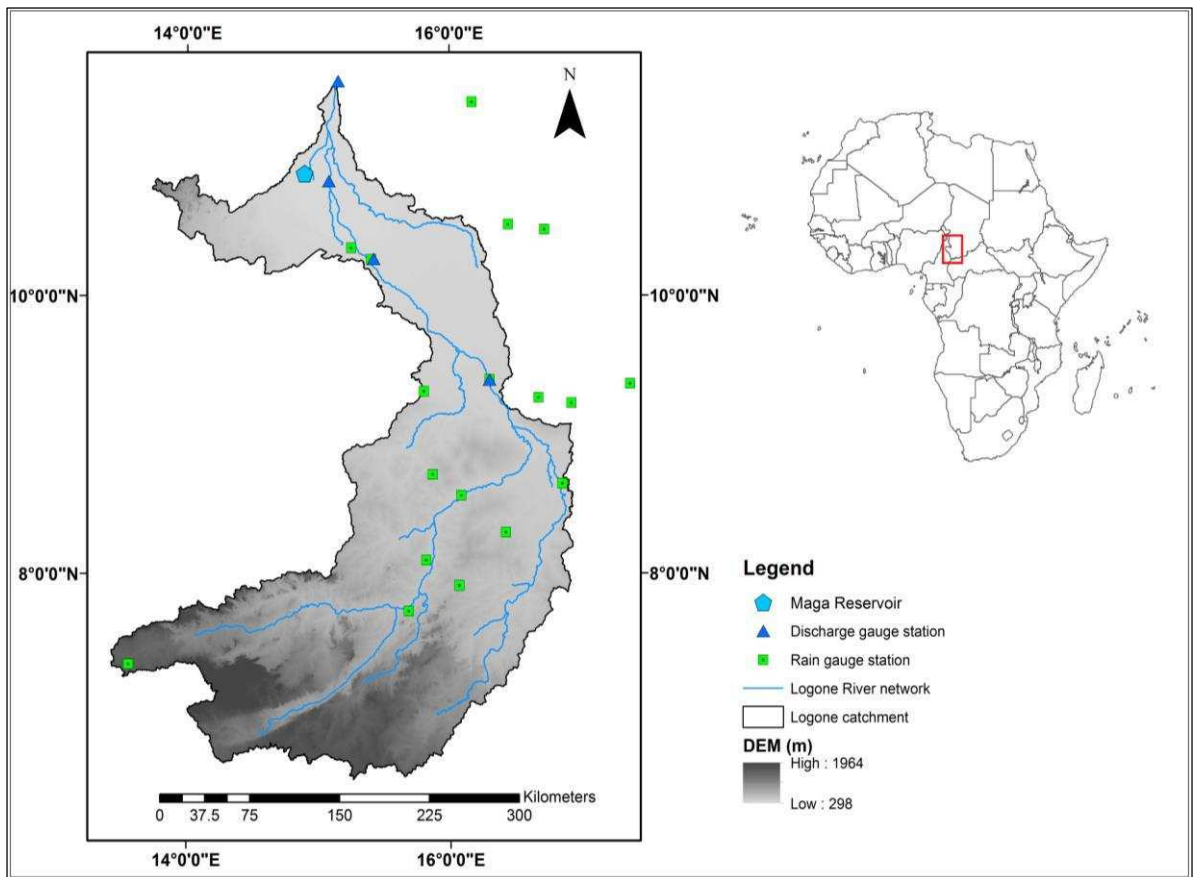
524 Latitude and Longitude in degrees

525



526
527
528

Figure 1: Architecture of an SOM (Adapted from Kagoda et al., 2010)



529
530
531

Figure 2: Map of study area showing rain and flow gauging stations

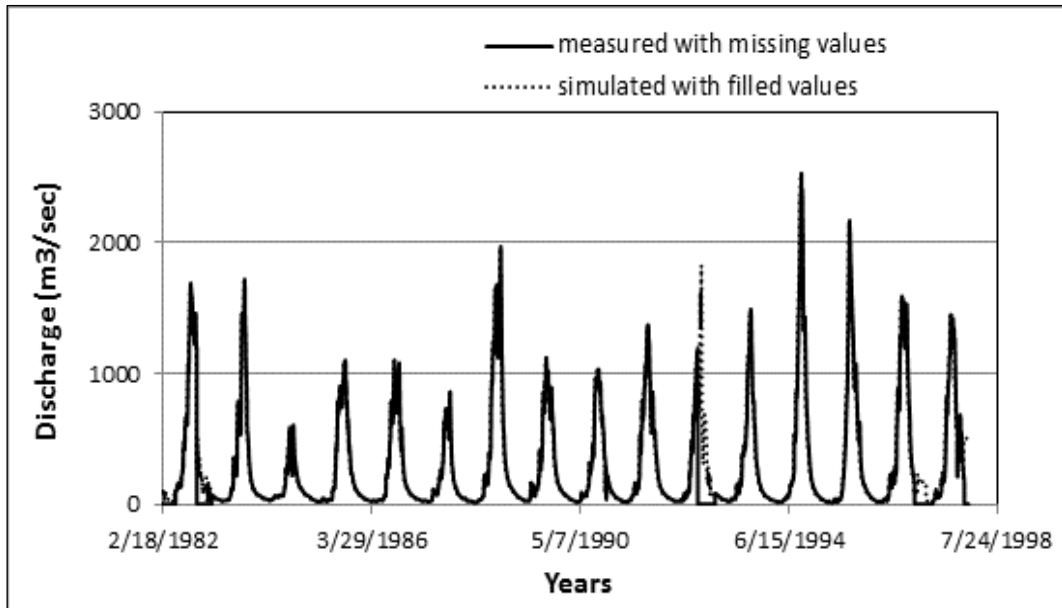


Figure 3: Observed and simulated discharge for Lai station 1973-1997

532
533
534

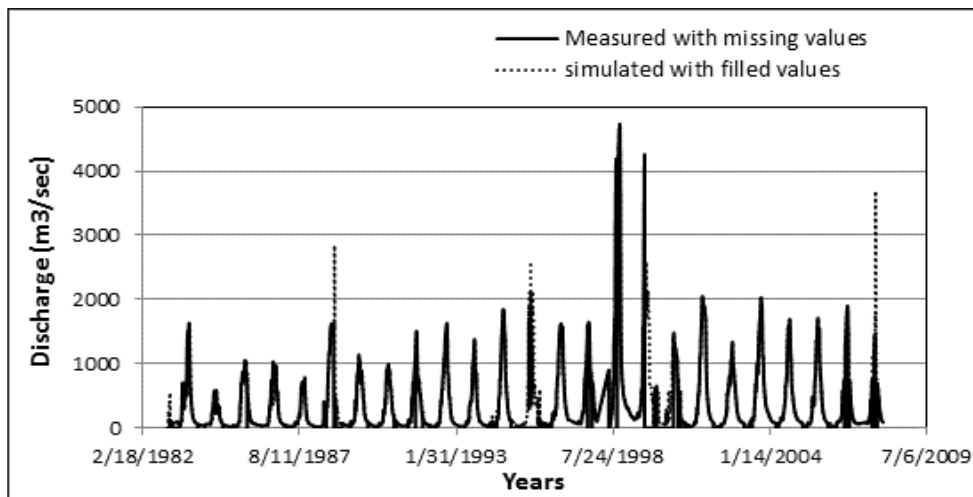


Figure 4: Observed and simulated discharge for Bongor station 1983-2007

535
536
537

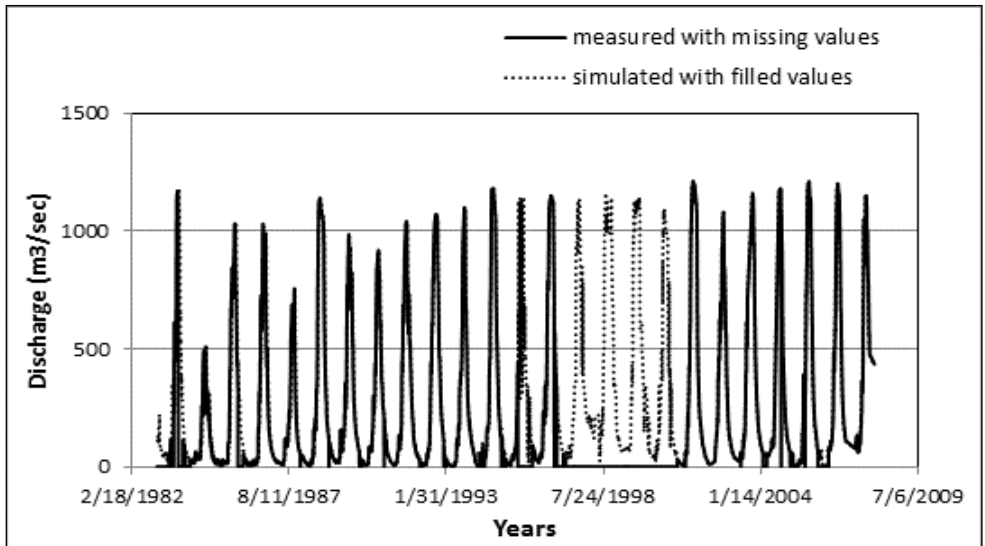


Figure 5: Observed and simulated discharge for Katoa station 1983-2007

538
539
540

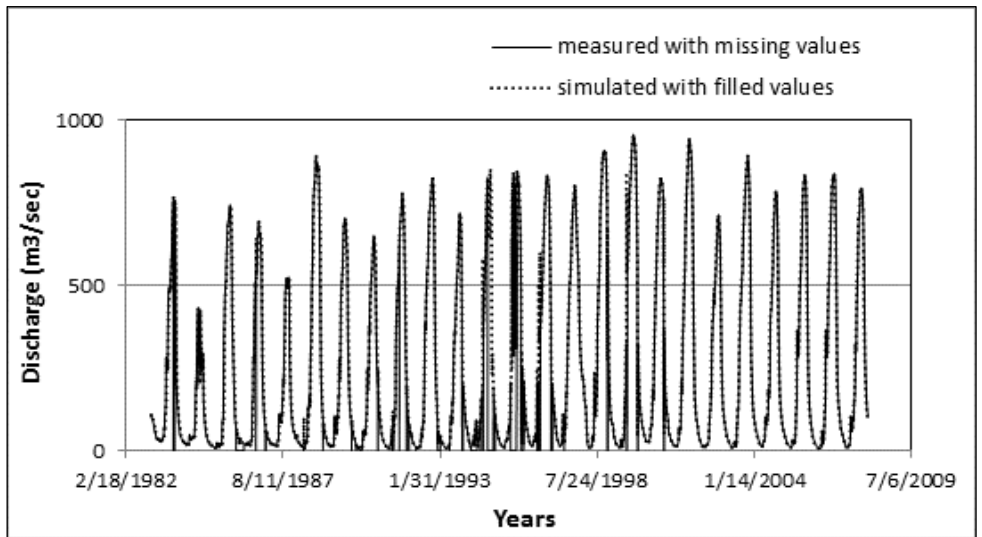


Figure 6: Observed and simulated discharge for Logone Gana station 1983-2007

541
542
543

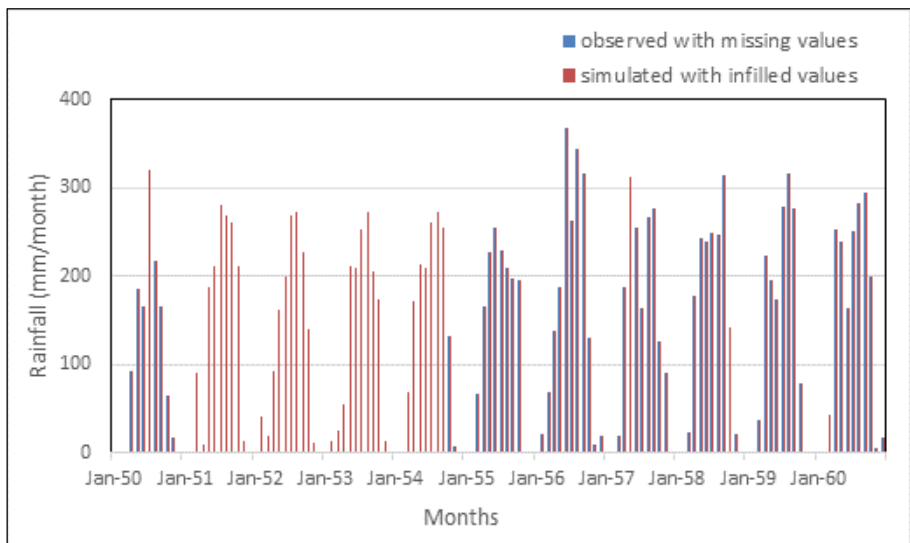
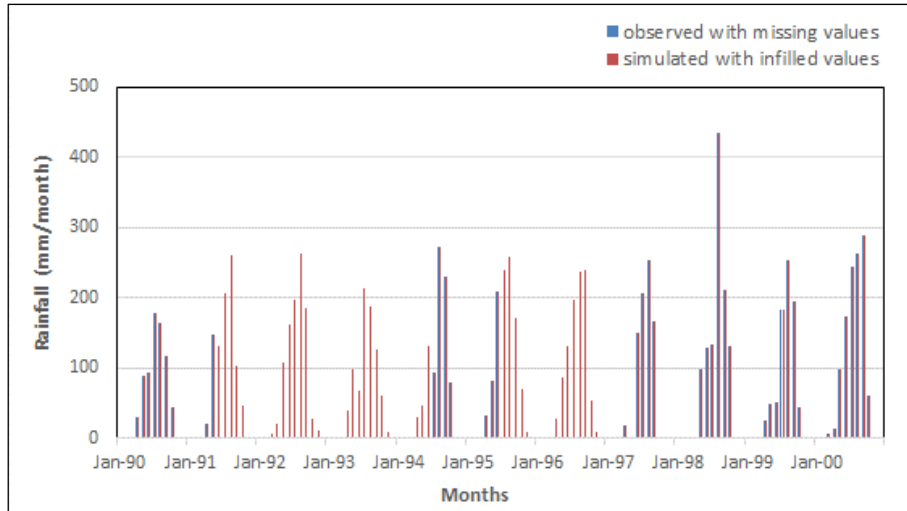


Figure 7: Observed and simulated rainfall for Ngaoundere (1950-1960)

544
545

546

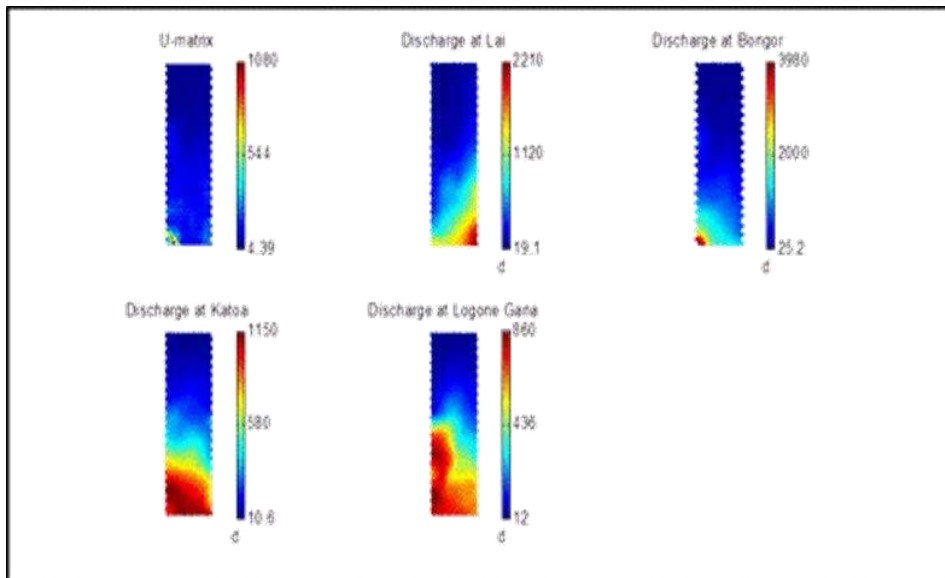


547

548

549

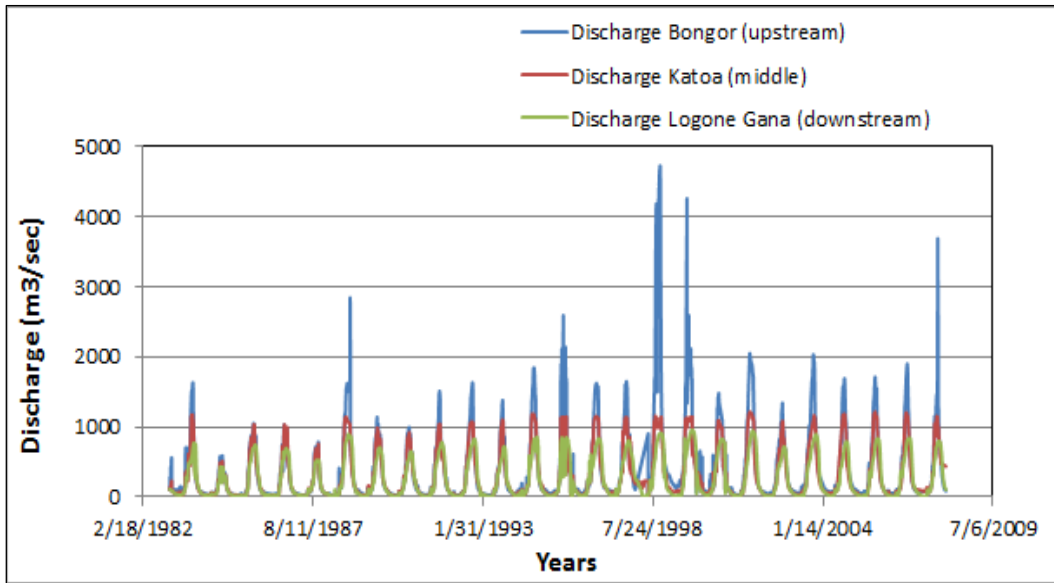
Figure 8: Observed and simulated rainfall for Kello (1990-2000)



550

551

Figure 9: Component planes for discharge at all the stations



552

553

554

Figure 10: Discharge at Bongor, Katoa and Logone Gana 1983-2007