



1 COPULA AND ARMA BASED STUDY OF CONTROLLED OUTFLOW AT FARAKKA
2 BARRAGE

3 Uttam Singh; Venkappayya R. Desai; Pramod K. Sharma; and Chandra S.P. Ojha

4 Research Scholar¹; Professor²; Associate Professor³; Professor⁴

5 Email: uttamsingh426@gmail.com; venkapd@civil.iitkgp.ernet.in; drpksharma07@gmail.com; cspojha@gmail.com

6 ^{1,3,4}Department of Civil Engineering, Indian Institute of Technology Roorkee-247667

7 ²Department of Civil Engineering, Indian Institute of Technology Kharagpur-721302

8 *Corresponding Author: Email: uttamsingh426@gmail.com

9

10

11



12 COPULA AND ARMA BASED STUDY OF CONTROLLED OUTFLOW AT FARAKKA
13 BARRAGE

14 Uttam Singh¹; Venkappayya R. Desai²; Pramod K. Sharma³; and Chandra S.P. Ojha⁴

15 ^{1,3,4}Department of Civil Engineering, Indian Institute of Technology Roorkee-247667

16 ²Department of Civil Engineering, Indian Institute of Technology Kharagpur-721302

17

18 **Abstract**

19 In this study, 25 years mean monthly out flow discharge data of Farakka barrage was used (i.e.,
20 from 1949 to 1968). Farakka barrage is located between on Ganga River. Spatial and temporal
21 variation in flow rate for any particular area is very common due to various meteorological and
22 other factors existing in nature. But large variations in these factors cause extreme events (e. g.,
23 floods and droughts). Monthly outflow discharge for a particular critical month are predicted
24 using statistical models (ARMA Model and Copula Model). Different Copulas (i.e., Normal, t,
25 Frank, Clayton, Gumbel-Hoggard, Ali-Mikhail-Haq) are used for this purpose and the copula
26 model is selected based on distribution functions (Normal distribution, Lognormal distribution,
27 Extreme value type-1 distribution, Generalized Extreme value type, Gamma distribution,
28 Weibull distribution, Exponential distribution). The distribution is selected based on the Mean
29 square error (MSE), Akaike Information Criterion (AIC), and Bayesian Information Criterion
30 (BIC). The model parameters were computed using the Maximum Likelihood (ML) estimation
31 method.

32

33 **Key words:** Farakka barrage; ARMA, Copulas, Simulation; Discharge.



34 **1. Introduction**

35 An accurate flood-frequency analysis is critical for the design of many civil infrastructures such
36 as drainage system and flood proof walls. Copula word is taken from Latin language and the
37 meaning of copula is link and the concept of copula was introduced in mathematical and
38 statistical manner by Sklar (1959) in a theorem that describes a copula as a function. Afterwards,
39 many researchers such as Genest and MacKay (1986), Genest and Rivest (1993) and Nelsen
40 (1999), Favre *et al.* (2004), Genest and Favre (2007) and Salvadori and De Michele (2007) used
41 in hydrology applications. Crucial steps for copulas modeling are driving the bivariate
42 distribution of peak flow and volume, volume and duration, peak flow and duration (Zhang and
43 Singh, 2006). Archimedean copulas (Clayton, Frank, Gumbel-Hoggard, Ali-Mikhail-Haq,
44 Indpendance and Joe) can be used for bivariate modeling peak flow and volume, volume and
45 duration, peak flow and duration.

46 Dependence structure of data set is captured by copulas, thus they are used for describing the
47 dependence of o extreme output values and also useful for dependance non parametric
48 measurement. Statistical dependence among three random variables two copulas are used for
49 modeling. The Archimedean copulas are prepared by association measurement of Kendall's *tau*
50 (Osorio et al. 2009). The probability density function for the two-dimensional random variable
51 representing volume and time is given in graphic form. The graphs both represents Clayton
52 copula and Gumbel-Hougaard functions. The Gumbel-Hougaard copula was best suited for this
53 study because it shows lower value in selection criterion function. Gumbel-Hougaard copula
54 shows better matching of empirical and theoretical distribution function. The results obtained in
55 the study, risk values at extreme analyzed values of controlled discharge and flood control
56 capacity are not monotonic. It represents that simulations were completed for sets of only 10000
57 cycle elements and only 10000 cycles (Twaróg, 2016). Peak flow and hydrograph volume both



58 can be jointly studied by bivariate approach (e.g., Goel et al. 1998; Yue et al. 1999; Favre et al.
59 2004; Shiau et al. 2007). The selection of the, different criterion should be consider among the
60 candidate copula (Chowdhary et al. 2011; Requena et al. 2013). The first criterion is the
61 goodness-of-fit test which relates the ability of copula to characterize the data (Genest et al.
62 2009), The second criterion is estimation of kendall's tau return period estimation by copula. It
63 relates the adequacy of copula, for a large copula value $t \in [0, 1]$, which is based on the Kendall's
64 function $K_C(t) = P[C_\theta(u_1, u_2) \leq t]$ (Genest and Rivest, 1993). The third criterion is the estimation
65 of Akaike Information Criterion (AIC) (e.g., Zhang and Singh, 2006). A copula-based model and
66 a distributed hydro-meteorological model and a copula-based model can be studied by
67 combining extension of observed flood series (Requena, et al. 2015). Significant number of
68 researchers found in their research that Gumbel-Hougaard copula as the most suitable choice to
69 model the dependence structure relating to the peak flow discharge and the flood volume (De
70 Michele et al., 2005; Zhang and Singh, 2007, Karmakar and Simonovic, 2009 and Li et al.,
71 2013). A copula-based approach was used to derive a bivariate distribution function of two
72 constituent flood variables, with regard to a real-world case study. It was found to provide an
73 effective and straightforward strategy for inferring probability functions from multivariate
74 sample data. Powerful tests developed inside copula framework allowed to investigate the
75 empirical dependence structure in an accurate manner, especially with respect to the evaluation
76 of tail dependencies (Balistocchi, 2017). The dependence of copula model between intensity and
77 rain fall duration, both properties of marginal distribution and dependence between intensity and
78 storm duration were preserved. The Joint cumulative distribution functions represents
79 dependence between independent variables of their marginal distribution of copula (Joe, 1997
80 and Nelsen, 2006). Gaussian copula was used for generation of 1020 synthetic data sets. Among



81 the data sets, 21 data sets lies beyond the range of acceptance so these data sets were omitted. Of
82 course it is not possible to cover all input-output cases in trained models the extrapolation limit
83 are required (Hooshyaripor et al. 2014). Best copula model can be selected by coarse grid model
84 selection with supposedly known marginal parameters in which 15 families of copulas were
85 divided into 4 categories and selection with uncertain marginal parameters (Parent et al. 2013).
86 Copula is a tool for modeling multivariate distribution in which input is the marginal
87 distribution. Multivariate distribution function couples to the corresponding marginal
88 distribution. (Poulin et al., 2007; Salvadori et al., 2007). The monsoon rainfall of Assam,
89 Meghalaya and Nagaland, Manipur, Mizoram, Tripura, Gumbel–Hoggard copula model was well
90 simulates for rain fall estimation (Ghosh, 2010). Marginal distributions and correlations values
91 are used to simulate the Gaussian model. They were taken four case studies to demonstrate its
92 usefulness in the reference of determination of field significance analysis, analysis of regional
93 risk , frequency analysis and design of hydrograph derivation by QdF models. (Renard et al.
94 2007). Copulas are very good tool to model multivariate data and they are very useful in
95 financial economics as well and in the analysis of multivariate survival data. Dependent variables
96 are very useful Monte Carlo simulations for copula model. It estimates the structural
97 dependence of the data set and describe accurately for dependence of extreme out come.
98 (Muhaisen, et al. 2006). Multivariate probability distributions with arbitrary marginal can be
99 constructed in a flexible manner with the introduction of copulas (Wang et al. 2001). Major
100 issue of a copula is the compatibility with dimensions though they were successfully tested and
101 applied on several hydrological problems. (Kao and Govindaraju, 2008). Application of copula
102 in the engineering problem need moderate and minimal computational effort and accuracy of the
103 output is also satisfactory (Kao et al., 2012). For two copula approach the spatial dependence of



104 rainfall dependence in sub-basins decreases up to 18 %. To predict decrease runoff error spatial
105 rainfall dependence could be recommended for copula modeling (Razmkhah, 2016).

106 The aim of this paper is to generate the out flow discharge data at Farakka barrage using
107 Copulas. In this study, Normal Copula, T- Copula, Frank Copula, Clayton Copula, Gumbel-
108 Hoggard (GH) copula, Ali-Mikhail-Haq(AMH) copula are used and best copula is selected for
109 generation of discharge data based on copula parameters, Mean square error(MSE), Akaike
110 Information criterion(AIC), Bayesian Information criterion (BIC).

111 ARIMA model was developed to forecast monthly inflow discharge in a reservoir system
112 (Mohan et al., 1955). Criteria for model selection are residual variance(Katz et al. 1981), Akaike
113 information criteria (Akaike 1974) and Posterior probability criteria (Kashyap 1977).

114 **2. Copulas used for study**

115 Copulas are alternative methods for dealing with multivariate extremes, and these are very
116 popular in recent times. Consider a moment pair of random variables U and V , with their
117 distribution functions $F(u) = P [U \leq u]$ and $G(v) = P [V \leq v]$, respectively, and a joint distribution
118 function $H(u, v) = P [U \leq u, V \leq v]$. Each pairs having of real numbers (u, v) , associated three
119 numbers: $F(u)$, $G(v)$, and $H(u, v)$ and each numbers are lie in the interval $[0,1]$. In other words,
120 each pair of real numbers i.e. (u, v) leads to a point $\{F(u), G(v)\}$ in the unit square $[0, 1] \times [0, 1]$,
121 and this ordered pair in turn corresponds to a number $H(u, v)$ in $[0,1]$. We will show that this
122 correspondence, those values are assign in the joint distribution function to each values of
123 ordered pair in the individual distribution functions. Such functions are named as copulas.

124 A copula is used as a tool in modeling multivariate distribution in which marginal distributions
125 are input data and neglect restrictions mentioned in pervious text. Copula means couples or joins



126 multivariate distribution functions to their corresponding distribution functions of their
 127 corresponding marginal distribution functions (Poulin *et al.*, 2007; Salvadori *et al.*, 2007).
 128 Definition which is given below is given by Sklar (1959), if p-dimensional distribution function
 129 then F can be written as:

$$130 \quad F(u_1, u_2, u_3, \dots, u_p) = C(F(u_1), F(u_2), F(u_3), \dots, F(u_p)) \quad (1)$$

131 where F_1, \dots, F_p = Marginal distribution functions. If F_1, \dots, F_p are continuous then the
 132 copula C is unique and has the representation (Poulin *et al.*, 2007):

$$133 \quad C(x_1, x_2, \dots, x_p) = F(F^{-1}(x_1), F^{-1}(x_2), \dots, F^{-1}(x_p)), \quad (2)$$

$$134 \quad 0 \leq x_1, \dots, x_p \leq 1$$

135 Copula is expressed for two random variables, U and V, with their CDFs, respectively, as $F_u(u)$
 136 and $F_v(v)$, let $X = F_u(u)$ and $Y = F_v(v)$, Where, X and Y are random variables which is uniformly
 137 distributed with their values x and y. The list copulas and its equations with generating function
 138 is shown in Table 1.

139 **Table 1. : List of Copulas and its equation, generating function and relation with τ .**

S. No.	Copula	Equation	Generating function	Relation with τ
1	Normal	$C(x_1, x_2, \dots, x_p) = P[U_1 \leq F^{-1}_1(x_1), U_2 \leq F^{-1}_2(x_2), \dots, U_p \leq F^{-1}_p(x_p)]$		
2	T	$C(x_1, \dots, x_d) = F(F^{-1}_1(x_1), \dots, F^{-1}_d(x_d))$		
3	Frank	$C_\theta(x, y) = \frac{1}{\theta} \ln \left[1 + \frac{[\exp(\theta x) - 1][\exp(\theta y) - 1]}{\exp(\theta) - 1} \right]$	$\Phi(t) = \ln \left[\frac{\exp(\theta t) - 1}{\exp(\theta) - 1} \right]$	$\tau = 1 - \frac{4}{\theta} [D_1(-\theta) - 1]$
4	Clayton	$C_\theta(x, y) = [x^{-\theta} + y^{-\theta} - 1]^{-1/\theta}$	$\Phi(t) = t^{-\theta} - 1$	$\tau = \frac{\theta}{\theta + 1}$
5	Gumbel-Hoggard	$C_\theta(x, y) = \exp \{ - [(-\ln x)^\theta + (-\ln y)^\theta]^{1/\theta} \}$	$\Phi(t) = (-\ln t)^\theta$	$\tau = 1 - \theta^{-1}$
6	Ali-Mikhail-Haq	$C_\theta(x, y) = \frac{xy}{1 - \theta(1-x)(1-y)}$	$\Phi(t) = \ln \left[\frac{1 - \theta(1-t)}{t} \right]$	$\tau = \left(\frac{3\theta - 2}{\theta} \right) - 2/3 (1 - 1/\theta)^2 \ln(1 - \theta)$

140



141 Where,

142 Θ = Parameter which controlling the dependence between x and y.

143 Φ = Generator of the copulas.

144 Debye function is expressed as follows.

$$145 D_n(\beta, x) = \frac{n}{x^n} \int_0^x \frac{t^n}{(e^t - 1)^\beta} dt \quad (3)$$

$$146 D_1(1, \Theta) = \frac{1}{\Theta} \int_0^\Theta \frac{t}{e^t - 1} dt \quad (4)$$

147

148 **3. Dataset used for Copulas**

149 Mean monthly discharge at Farakka barrage data set about twenty-five years from 1949 to 1973
150 data has taken from Water Resources Information System of India at Farakka barrage project,
151 Farakka, West Bengal, India.

152 The observed data set are divided into two parts. One part contains twenty years' data (from
153 1949 to 1968) has been used for parameter estimation i.e. in model calibration, next five years'
154 data (from 1969 to 1973) has been used for model validation and testing. Parameter estimation
155 data is arranged such a way that pre-monsoon (December to May) and post monsoon (June to
156 November) data is separated and making two series of dataset for copulas.

157 **4. Selection of distribution for Copulas**

158 For modeling of controlled outflow, bivariate Copula has taken in this study. As Copula accepts
159 CDF of variables, distribution functions of two variables, should be known. The distribution
160 functions are chosen on the basis of AIC, BIC values, k-s test and probability plots. The
161 distributions that are tested to know the parent distribution of two variables are normal
162 distribution, lognormal distribution, extreme value type I distribution, generalized extreme value
163 distribution, gamma distribution, weibull and exponential distributions. We used data set for



164 different times i.e., from Dec. -May 1949 to Dec. -May 1968 (Figure 1), from Jun. -Nov. 1949 to
165 Jun. -Nov. 1968 (Figure 2), Dec. -May 1949 to Dec. -May 1968 (Figure 3), Jun. -Nov. 1949 to
166 Jun. -Nov. 1968 (Figure 4), Dec. -May 1949 to Dec. -May 1968 (Figure 5), Jun.-Nov.1949 to
167 Jun. -Nov. 1968 (Figure 6).

168 The violet colour represents the data set for different times and red colour represents normal
169 distribution, green colour represents lognormal distribution, etc as shown in Figures 1-6. Figure 1
170 represents cumulative distribution function of data points along with all distributions in pre
171 monsoon seasons (Dec.- May 1949 to Dec.-May 1968).

172 Figure 2 represents cumulative distribution function of data points along with all distributions in
173 post monsoon seasons (Jun. - Nov. 1949 to Jun. - Nov 1968). Figure 3 represents Probability
174 density function of data points along with all distributions in post monsoon seasons (Dec.- May
175 1949 to Dec.-May 1968). Figure 4 represents Probability density function of data points along
176 with all distributions in post monsoon seasons (Jun. -Nov. 1949 to Jun. -Nov. 1968). Select the
177 standard distribution which is best fit for original data sets. Violet colour of Figure 5 represents
178 the data points and other colour represents the various distributions of mean monthly discharge
179 (Dec. -May 1949 to Dec. -May 1968) Select the best fit standard probability distribution. Violet
180 colour of this Figure 6 represents the data points and other colour represents the various
181 distributions of mean monthly discharge (Jun.-Nov.1949 to Jun. -Nov. 1968). Selection of the
182 distribution function can be based the best fit for original data sets.

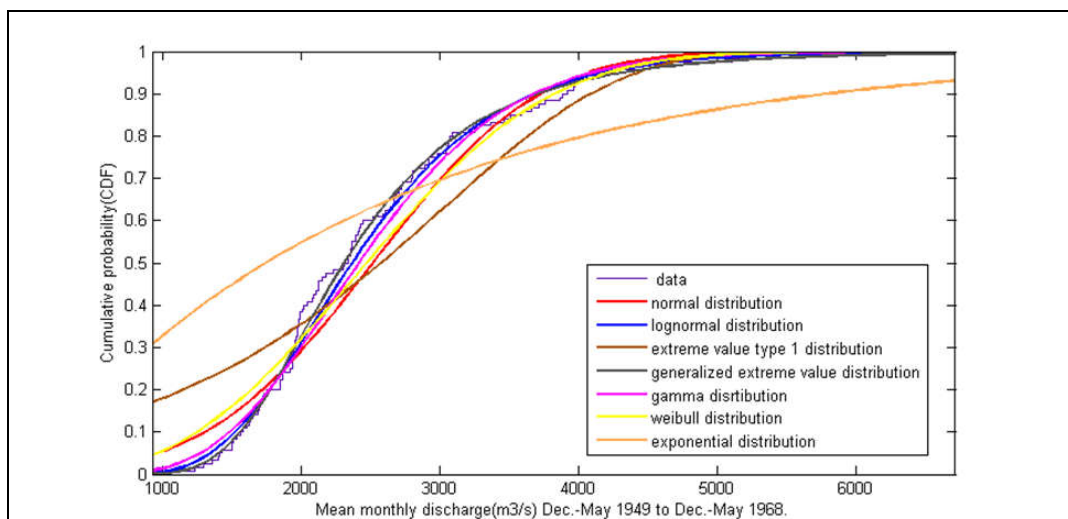


Figure 1. CDF of mean monthly discharge(m^3/s) Dec. -May 1949 to Dec. -May 1968.

183
184

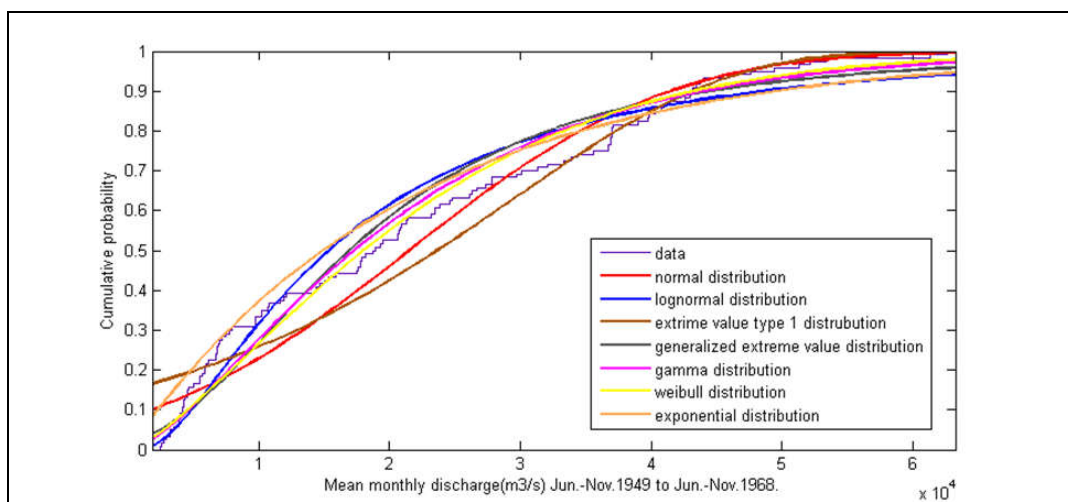
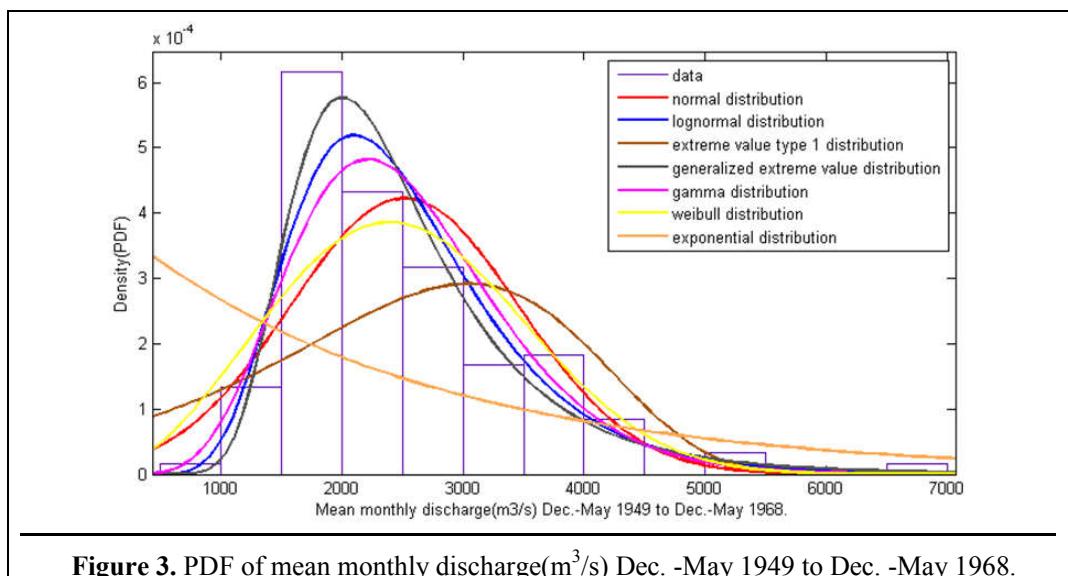
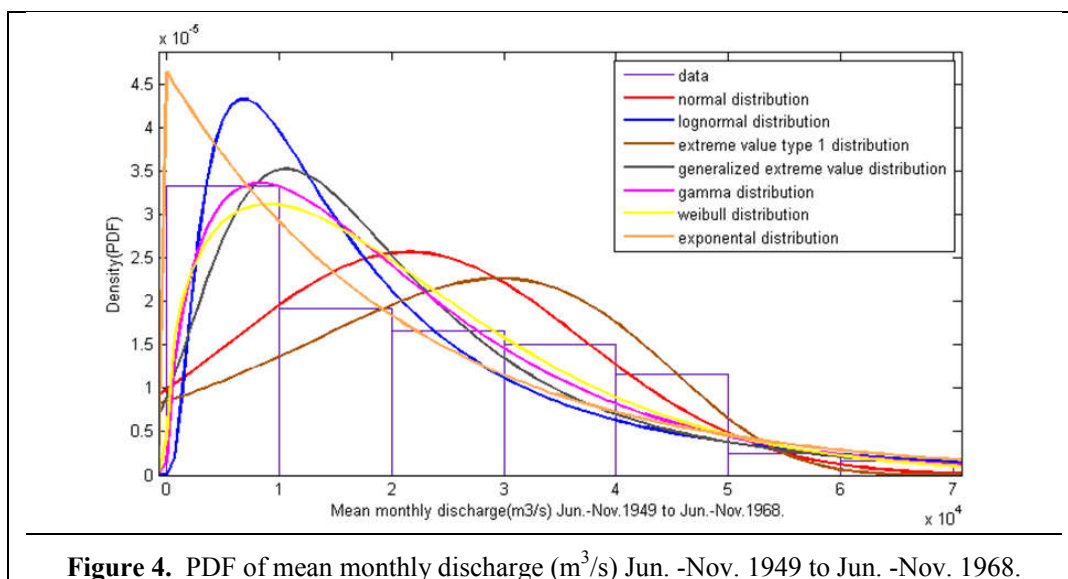


Figure 2. CDF of mean monthly discharge (m^3/s) Jun. -Nov. 1949 to Jun. -Nov. 1968.

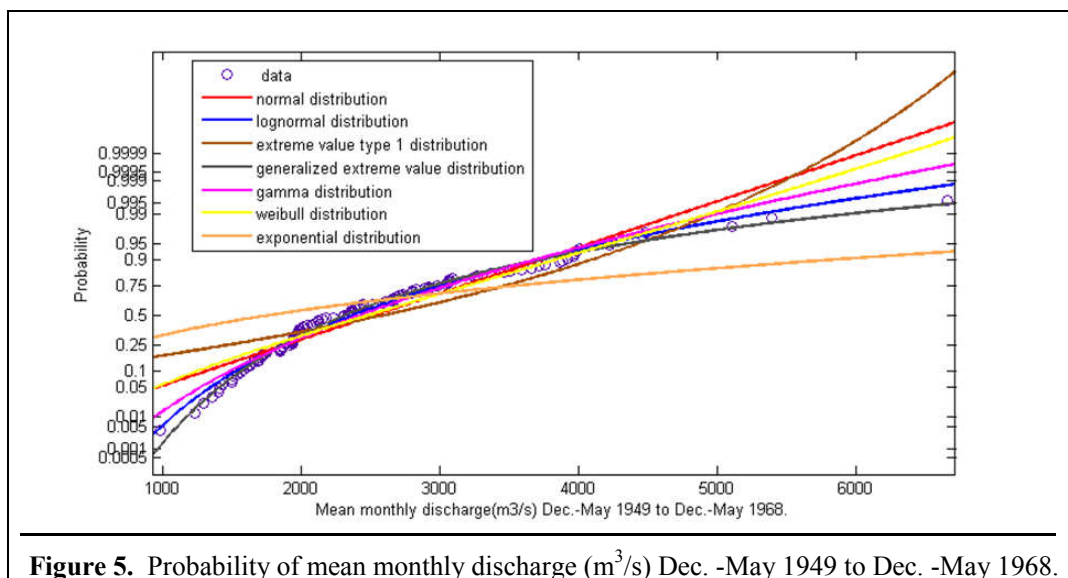
185
186
187
188
189
190
191



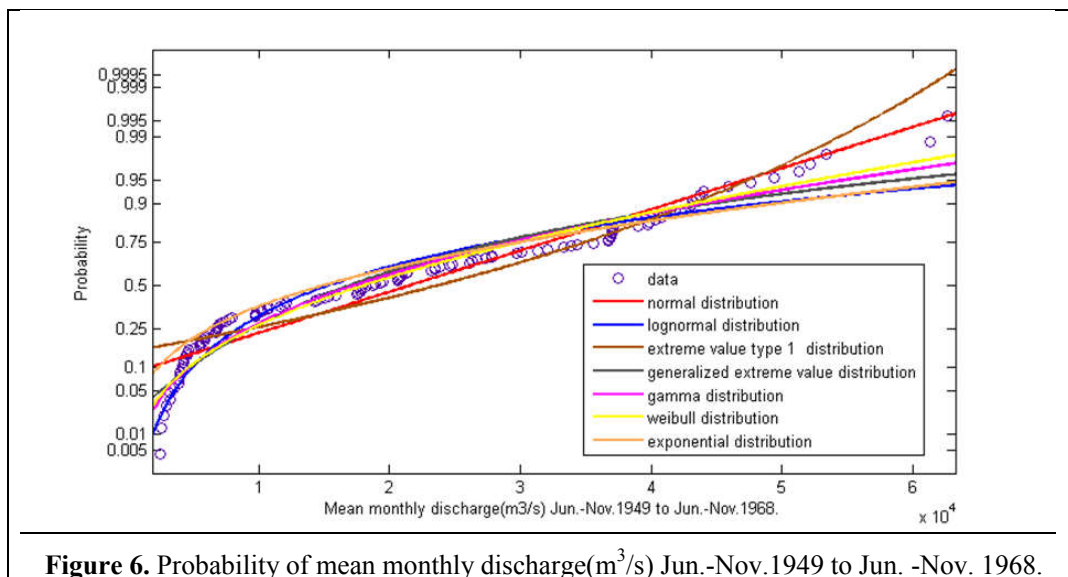
192
 193



194
 195



196
 197
 198
 199



200
 201



202 **4.1 Mean square Error (MSE) or Mean squared deviation (MSD)**

203 It is measurement of the mean of the squares of the errors or deviations i.e., the difference
204 between the estimator and what is estimated value (Table 2). MSE represents the risk function
205 corresponding to the expected value of the squared error loss or quadratic loss. The difference in
206 the MSE because of randomness. Lowest value of AIC is good for model.

$$207 \text{MSE} = \sum \frac{(e_{cdf} - p_{cdf})^2}{n} \quad (5)$$

208 Where,

209 e_{cdf} = Empirical Cumulative Density Function

210 p_{cdf} = Predicted Cumulative Density Function

211

212 **4.2 Akaike Information Criterion (AIC)**

213 For a given data set and given set of models . AIC measures relative quality of statistical
214 methods and it compute the each model's quality, relative to other models quality (Table 2).
215 Hence, AIC criteria is used for model selection and lowest value of AIC is proffered for model. .

$$216 \text{AIC} = n * \ln (\text{MSE}) + 2K + \frac{2K*(K+1)}{n-K-1} \quad (6)$$

217 Where,

218 n = Number of data points.

219 K = Number of parameters.

220



221 **4.3 Bayesian Information Criterion (BIC)**

222 It is a model selection criterion , model is selected among the finite set of model. Model with
223 lowest value of BIC is preferred (Table 2). It is mainly based on likelihood function and it having
224 approximate same conditions as Akaike information criterion (AIC).

$$225 \text{ BIC} = n \cdot \ln(\text{MSE}) + K \cdot \ln(n) \quad (7)$$

226 Where,

227 n = Number of data points.

228 K = Number of parameters.



Table 2. Statistic of distributions of data.

Data	Distribution	MSE	AIC	BIC
June-Nov.	Normal	0.19185140	-194.02150	-188.54911
Dec. -May	Normal	0.18211572	-200.27095	-194.79857
June-Nov.	Lognormal	0.19240251	-193.67728	-188.20489
Dec. -May	Lognormal	0.19197824	-193.94219	-188.46980
June-Nov.	Extreme value type 1	0.17340721	-206.15091	-200.67853
Dec. -May	Extreme value type 1	0.14660442	-226.29948	-220.82709
June-Nov.	Gen.extreme value	0.09875020	-271.61252	-263.45694
Dec. -May	Gen. extreme value	0.09914008	-271.13967	-262.98409
June-Nov.	Gamma	0.19490464	-192.12678	-186.65439
Dec.-May	Gamma	0.19024585	-195.02997	-189.55758
June-Nov.	Weibull	0.19556132	-191.72315	-186.25077
Dec.-May	Weibull	0.17272273	-206.62552	-201.15314
June-Nov.	Exponential	0.16679494	-212.88491	-210.13132
Dec.-May	Exponential	0.11907597	-253.32532	-250.57173
June-Nov.	Kernel_normal	0.16931844	-211.08299	-208.32939
June-Nov.	Kernel_box	0.16880085	-211.45037	-208.69678
June-Nov.	Kernel_triangle	0.16920205	-211.16550	-208.41191
June-Nov.	Kernel_epanechnikov	0.16901130	-211.30086	-208.54727
Dec. -May	Kernel_normal	0.18343705	-201.47214	-198.71855
Dec. -May	Kernel_box	0.18311280	-201.68445	-198.93085
Dec. -May	Kernel_triangle	0.18327642	-201.57727	-19882367
Dec. -May	Kernel_epanechnikov	0.18321604	-201.61681	-198.86322

229
 230



231 **4.4 Kolmogorov – Smirnov test**

232 The Kolmogorov–Smirnov test (K–S test or KS test) is a nonparametric test of the equality of
 233 continuous, one-dimensional probability distributions that can be used to compare a sample with
 234 a reference probability distribution (one-sample K–S test), or to compare two samples (two-
 235 sample K–S test) (Table 3). The two-sample K–S test is one of the most useful and general
 236 nonparametric methods for comparing two samples, as it is sensitive to differences in both
 237 location and shape of the empirical cumulative distribution functions of the two samples. In the
 238 Figure 7 and 8, green colour shows empirical CDF and red colour shows generalized extreme
 239 value of CDF. On the basis of Figure 7, generalized extreme value distribution is representing
 240 best fit for cumulative distribution function (Jun.-Nov.1949 to Jun. -Nov. 1968). Further, on the
 241 basis of Figure 8, generalized extreme value distribution is represents best fit for cumulative
 242 distribution function (Dec. -May 1949 to Dec. -May 1968).

243 **Table 3. k-s statistics of distributions of data.**

			K - S	Test	
Data	Distribution	H	p	k-s	cv
June-Nov.	Normal	0	0.0509	0.1222	0.1225
Dec.-May	Normal	1	0.0336	0.129	0.1225
June-Nov.	Lognormal	0	0.0691	0.117	0.1225
Dec.-May	Lognormal	0	0.3695	0.0824	0.1225
June-Nov.	Extreme value type 1	1	0.0015	0.1712	0.1225
Dec.-May	Extreme value type 1	1	0.000055	0.207	0.1225
June-Nov.	Gen. extreme value	0	0.092	0.1118	0.1225
Dec.-May	Gen. extreme value	0	0.669	0.0649	0.1225
June-Nov.	Gamma	0	0.1526	0.1021	0.1225
Dec.-May	Gamma	0	0.1784	0.0989	0.1225



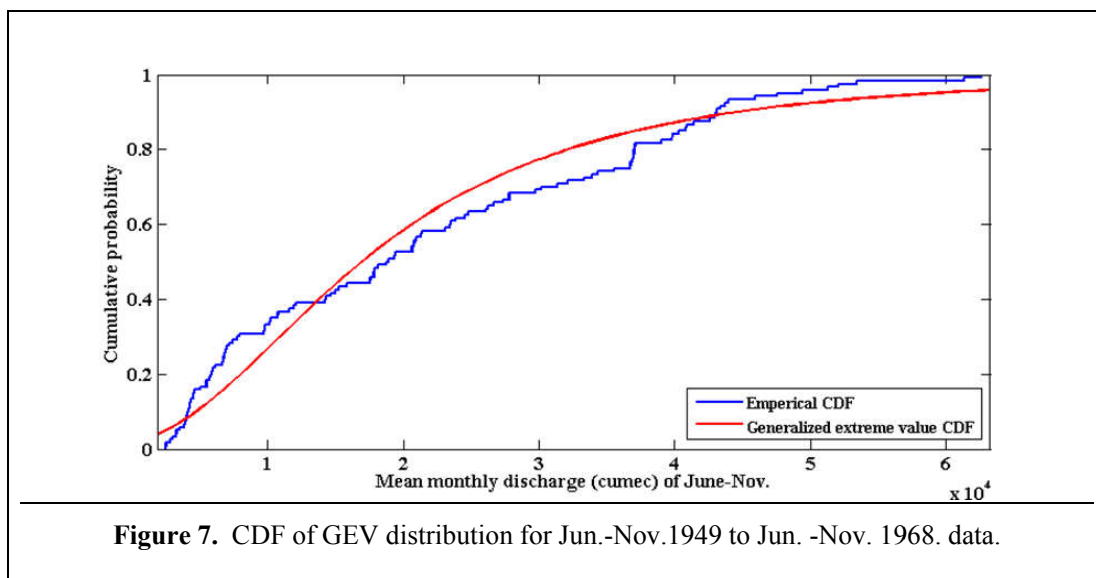
June-Nov.	Weibull	0	0.1162	0.1075	0.1225
Dec.-May	Weibull	0	0.1046	0.1095	0.1225
June-Nov.	Exponential	0	0.0748	0.1156	0.1225
Dec.-May	Exponential	1	6.16E-17	0.3939	0.1225
June-Nov.	Kernel_normal	1	0.0286	0.1315	0.1225
June-Nov.	Kernel_box	1	0.0141	0.1421	0.1225
June-Nov.	Kernel_triangle	1	0.0255	0.1333	0.1225
June-Nov.	Kernel_epanechnikov	1	0.0195	0.1374	0.1225
Dec.-May	Kernel_normal	0	0.8141	0.0567	0.1225
Dec.-May	Kernel_box	0	0.8095	0.057	0.1225
Dec.-May	Kernel_triangle	0	0.8221	0.0562	0.1225
Dec.-May	Kernel_epanechnikov	0	0.8074	0.0571	0.1225

244

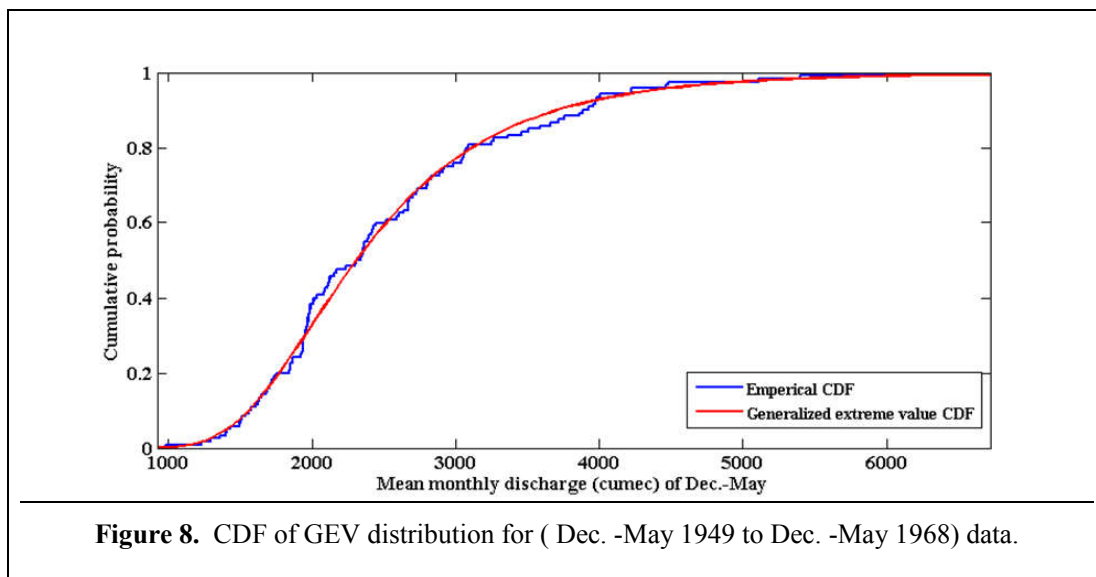
245

246

247



248



249

250

251



252 **5. Copula parameter estimation**

253 **Table 4. : Copulas and its parameter.**

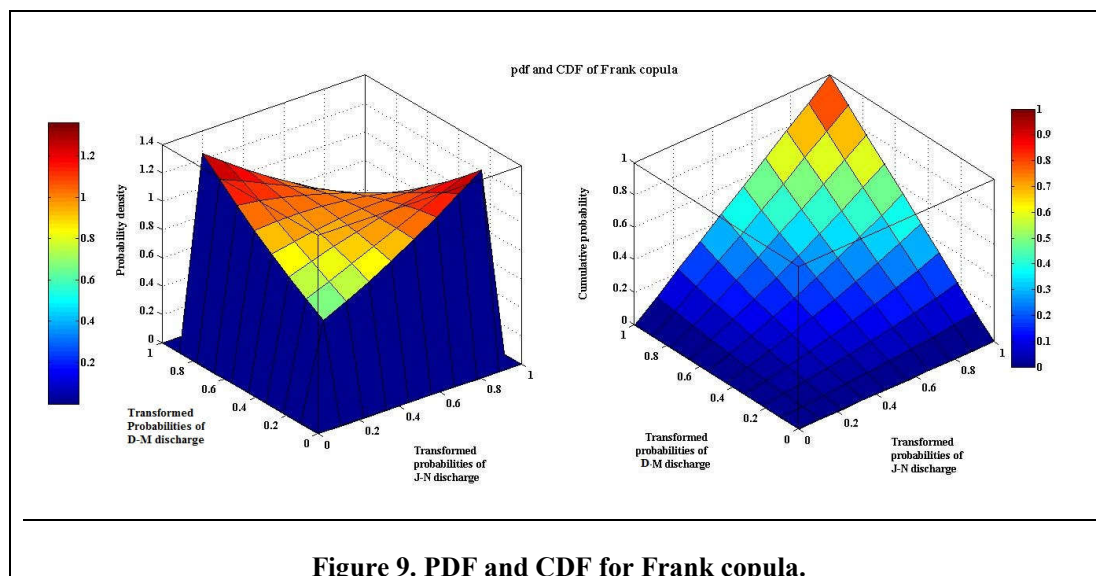
	Parameter					
Copula model	rho	nu	MLE	MSE	AIC	BIC
Gaussian	-0.2338		3.38	0.002125	-736.448	-732.69
t	-0.2344	3.79E+06	3.37	0.002126	-734.331	-731.65
Frank	-1.1424		4.557	0.00206	-740.190	-736.44
AMH	-1		2.284	0.002072	-739.477	-735.72
Clayton	1. 45E-06		0.6685	0.002209	-731.817	-728.06
GH	1		-7.2E-07	0.002201	-732.216	-728.46

254

255

256 For a best copula model MLE should be high and MSE, AIC, BIC should be minimum from the
 257 above data frank is best model for predicting the data (Table 4). Figure 9 shows the probability
 258 density variation from green to red, green having lowest probability density and green colour
 259 having maximum probability density. It also represents the probability density function and
 260 cumulative distribution function for frank copula which is best for prediction of discharge data.

261



262

263 **6. Validation test of Copula**

264 Validation test of Frank Copula is performed by comparing observed and empirical CDF in
 265 calibration and validation test. Here, observed CDF is CDF of Frank Copula and Empirical CDF
 266 is taken from some non-parametric method (Table 5). Formula of empirical CDF of copula is
 267 given below. In the Figure 10, the blue points shows data points at calibration and validation
 268 state. Blue points represents data points in calibrated and validation stage by Frank copula as
 269 shown in Figure 10.

270

Table 5. : Statistics in calibration and validation test.

Statistics	Calibration Test	Validation Test
MSE	0.00206	0.00147
R ²	0.94	0.9

271

272

273

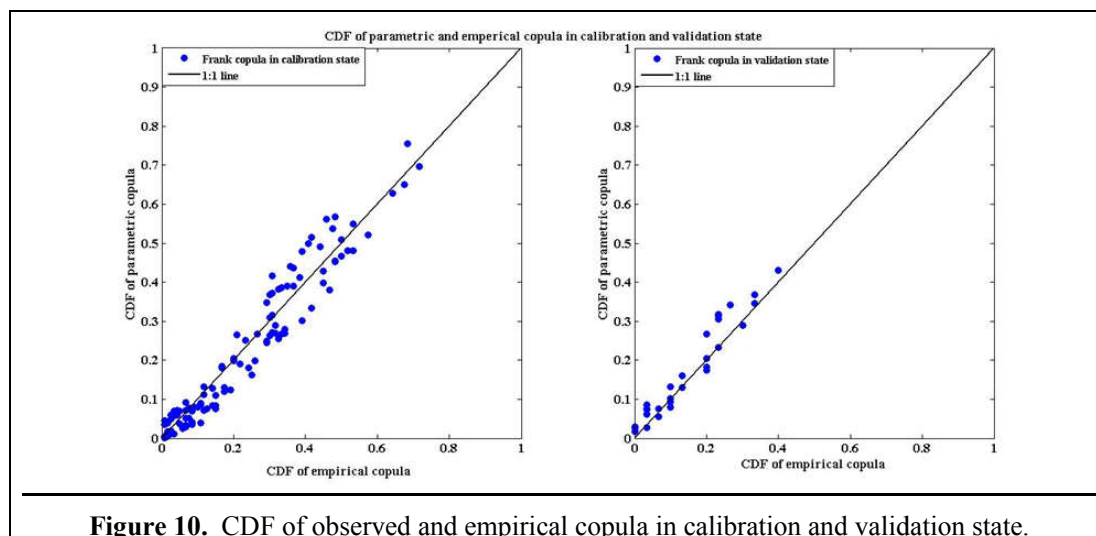


Figure 10. CDF of observed and empirical copula in calibration and validation state.

274

275

276 7. Statistical Approach for ARMA

277

278 In this approach linear type stationary ARMA models are fitted by observed discharge data

279 where stationary means the ARMA models that are generated from a time series does not

280 changing its underlying probability distribution function (pdf) from which different values of

281 time series are pulled out. In loose sense stationarity indicates time series has constant mean and

282 variance throughout the process where time series is the collection of random variables, plotted

283 corresponding of its time, follow on their own distribution (figure. 11). In ARMA model AR i.e.

284 auto regressive term indicates lag of time series value and moving average is the lag in error

285 term. Generally, ARIMA is conventional class of model where “I” integration term indicates

286 order of difference required to do the time series stationary but in this study it is done by

287 normalizing all the discharge data through its long term mean and standard deviation (figure. 12).

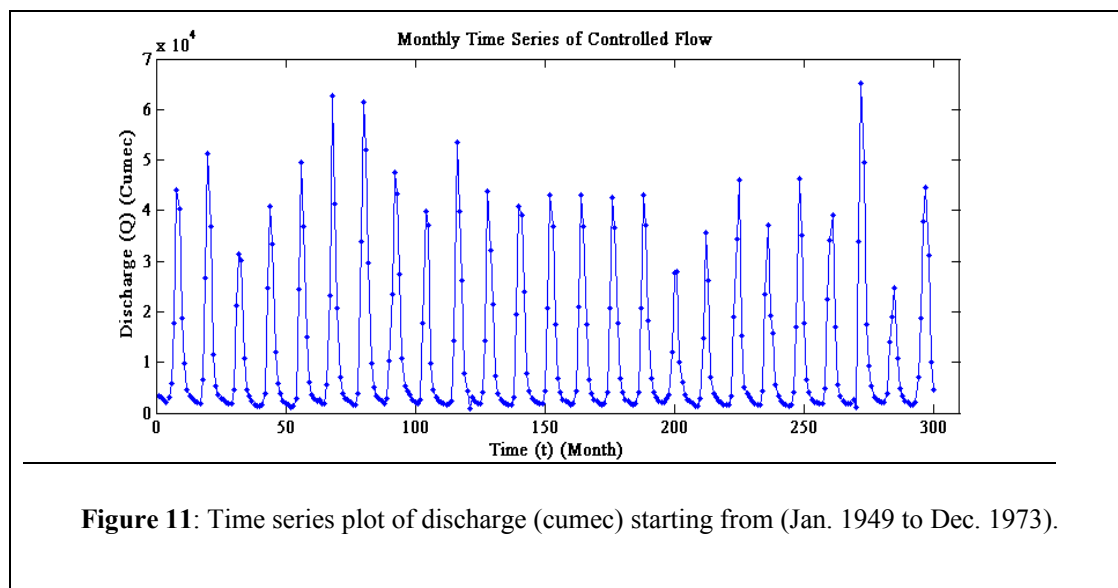
288 The mathematical form of normalization is given below.

$$289 \quad Z_i = \frac{X_i - \bar{X}}{\sigma_i} \dots\dots\dots (8)$$

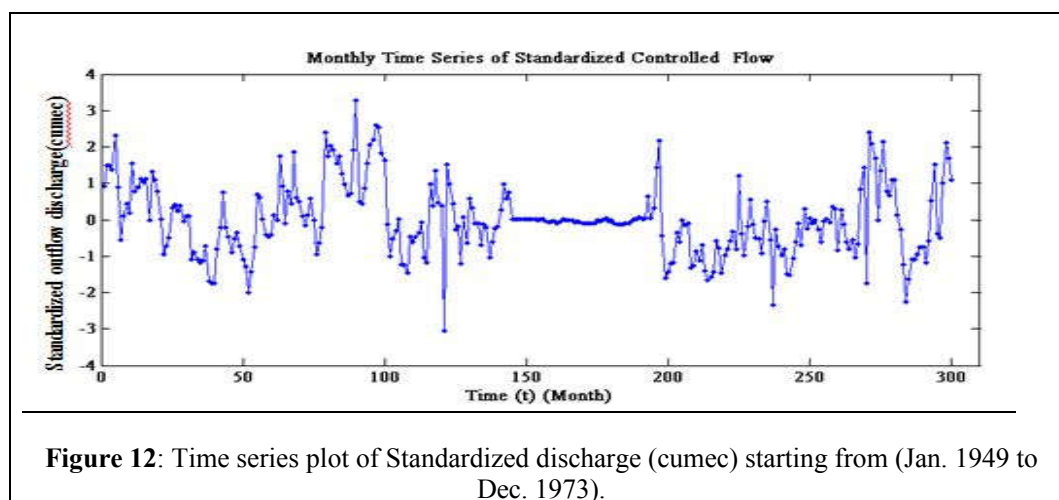


290 Where, X_i =Value of mean monthly discharge, \bar{X} =Long term average, σ_i =Long term standard
291 deviation, $i=1$ to N , N is total number of data point in monthly step. The normalization or
292 differencing in the data is not only make it stationary but also removes periodicity from the time
293 series where periodicity can be defined as correlation i.e. linear association of data with the
294 previous some lag value of data. As we are interested to only capture unknown information from
295 a process which are unknown due to noise or random term (stochastic factor in the process), so
296 deterministic part in terms of long term mean, periodicity, seasonality, trend, sudden drop or
297 jump is necessary to remove from the time series since these deterministic terms already reflects
298 known information about the process, are not required to model. Generally monthly discharge
299 time series shows periodicity and seasonality in the data set and it is necessary to remove before
300 calibrate (finding Parameter of model) to ARMA model as this type of model is developed to
301 capture unknown information from noise i.e. random process.

302 The observed data set are divided into two parts. One part contains twenty years' data (from
303 1949 to 1968) has been used for parameter estimation i.e. in model calibration, next five years'
304 data (from 1969 to 1973) has been used for model validation and testing. The mean monthly
305 discharge data used for model calibration may have serial correlation i.e. any data in particular
306 time step depends on its previous adjacent data and may follow so on. The time series plot of
307 observe discharge depicts this serial correlation, seasonality or periodicity in terms of
308 information contain in the series by showing some regularity or similar oscillation of the series.



309



310

311

312

313 8. Spectral Analysis

314 The observe time series is analyzed in frequency domain to indicate exactly in which months

315 periodicity present in the data that is only indicates by correlogram. In this frequency domain



316 analysis an assumption is taken as time series is a random sample of a process over time and is
 317 made up of oscillations of all possible frequencies. The time series is approximated by signal
 318 process contains deterministic term in wave form and noise or random term by which the
 319 information is extracted from time series and shows prominent spike in variance spectrum plot.
 320 The contributing equations for spectral analysis are given below

321

$$322 \quad X_t = \alpha_0 + \sum_{k=1}^{n-1/2, n/2} [\alpha_k \cos(2\pi f_k t) + \beta_k \sin(2\pi f_k t)] + \varepsilon_t \quad (9a)$$

$$323 \quad f_k = \frac{k}{N}; \quad P = \frac{1}{f_k}; \quad \alpha_0 = \bar{x} \quad (9b)$$

$$324 \quad \alpha_k = \frac{2}{N} \sum_{i=1}^n x_t \cos(2\pi f_k) \quad k = 1, 2, 3, \dots, M \quad (9c)$$

$$325 \quad \beta_k = \frac{2}{N} \sum_{i=1}^n x_t \sin(2\pi f_k) \quad k = 1, 2, 3, \dots, M \quad (9d)$$

326

327 Where;

328 N = Observation numbers, X_t = Observe rainfall data, t = Time step in month

329 P = Periodicity in the data, \bar{X} = Mean of the series (average monthly rainfall)

330 α_k = Cosine wave form, β_k = Sine wave form of time series.

331 M = Maximum lag typically consider $0.25N$.

332 Values of α_k and β_k in equation number 7 are valid up to $k = N/2$.

333

334 8.1. Line spectrum

335 The spike in the line spectrum confirms the presence of particular month periodicity in the data
 336 (Figure 13 and Table 6)and lime spectrum is plot between spectral density versus angular



337 frequency. It is also known as variance spectrum. Line spectrum plot is drawn by using discharge
 338 data and standardized discharge data.

$$339 \quad I_k = \frac{N}{2} [\alpha_k^2 + \beta_k^2]; \quad k = 1, 2, 3, \dots, M \quad (10a)$$

$$340 \quad \alpha_k = \frac{2}{N} \sum_{i=1}^n x_t \cos(2\pi f_k) \quad k = 1, 2, 3, \dots, M \quad (10b)$$

$$341 \quad \beta_k = \frac{2}{N} \sum_{i=1}^n x_t \sin(2\pi f_k) \quad k = 1, 2, 3, \dots, M \quad (10c)$$

342

$$343 \quad \omega_k = \frac{2\pi k}{N}; \quad k = 1, 2, \dots, M \quad (11)$$

344 Where, ω_k = Angular frequency and I_k = Spectral density.

345 N = Observation numbers,

346

347

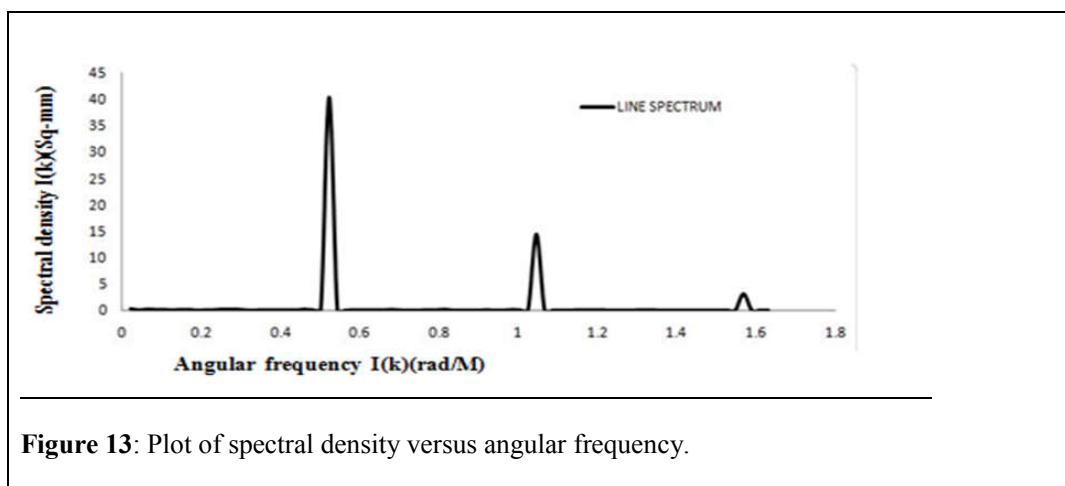


Figure 13: Plot of spectral density versus angular frequency.

348

349



350 **Table 6: Showing the spectral density and frequency data corresponding to spikes.**

Spike	Spectral density (I_k) (sq-mm)	Angular frequency (ω_k)(rad/M)	Periodicity (Months)
1	$4.05 * 10^{10}$	0.52	12
2	$4.48 * 10^9$	1.05	6
3	$4.60 * 10^9$	1.6	4

351

352 **9. Model Description**

353 Auto regressive moving average models are developed using white noise series. In the present
 354 study the information form observed time series has captured not only developing ARMA (p, q)
 355 model but also by pure AR (p) and MA (q) model. The block diagram for AR (p), MA (q) and
 356 ARMA (p, q) process are shown below.

357

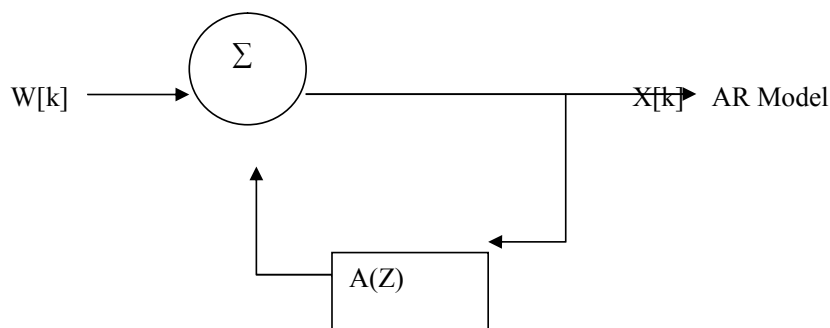
358

359

360

361

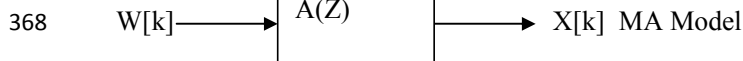
362



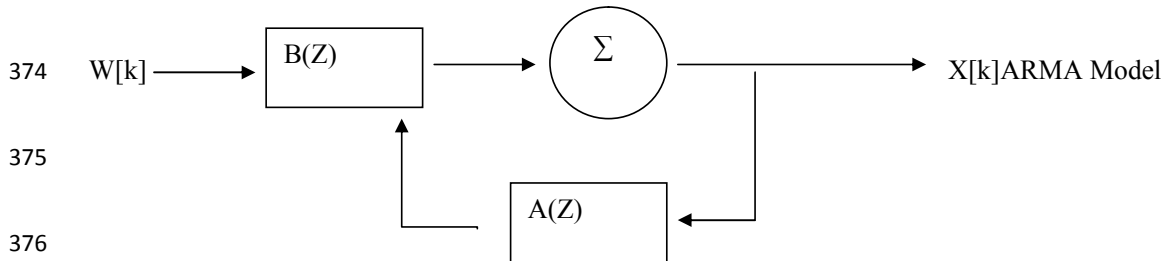


363 $X[k] = \sum a_n X[k-n] + W[k]$, where $X[k]$ = Discrete value or k^{th} sequence of random variable
 364 (discharge), a_n = AR parameter for n^{th} order sum over $n=1$ to N , N = number of data point, $X[k-$
 365 $n]$ = n^{th} lag of random variable (discharge), $A(Z)$ = AR polynomial equation, $W[k]$ = Error term
 366 associated in the model prediction.

367



370 $X[k] = \sum b_n W[k-n]$, where $X[k]$ = Discrete value or k^{th} sequence random variable (rainfall), b_n =
 371 MA parameter for n^{th} order sum over $n = 0$ to $M-1$, M = number of error point, $W[k-n]$ = n^{th} lag
 372 of white noise or error term, $B(Z)$ = MA polynomial equation, $W[k]$ = Error term associated in the
 373 model prediction.



377 $X[k] = \sum a_n X[k-n] + \sum b_n W[k-n]$, where $X[k]$ = Discrete value or k^{th} sequence of random
 378 variable (discharge), a_n = AR parameter for n^{th} order, $X[k-n]$ = n^{th} lag of random variable
 379 (discharge) sum over for $n=1$ to N , N = number of data point, $a(Z)$ = AR polynomial equation,
 380 $W[k]$ = Error term associated in the model prediction (white noise), b_n = MA parameter for n^{th}
 381 order sum over $n=0$ to $M-1$, M = number of error point, $W[k-n]$ = n^{th} lag of white noise or error
 382 term, $B(Z)$ = MA polynomial equation.



383 An ARMA (p, q) model which having autoregressive order i.e. p and moving average order i.e.

384 q can be expressed as following from of equation.

$$385 \quad X_t - \Phi_1 X_{t-1} - \Phi_2 X_{t-2} - \dots - \Phi_p X_{t-p} = \varepsilon_t + \Theta_1 \varepsilon_{t-1} + \Theta_2 \varepsilon_{t-2} + \dots + \Theta_q \varepsilon_{t-q} \quad (12)$$

$$386 \quad \Phi(L)X_t = \Theta(L) \varepsilon_t, \quad \Phi(L) = 1 - \sum_{j=1}^p \phi_j L^j \quad \text{and} \quad \Theta(L) = \sum_{j=0}^q \theta_j L^j \quad (13)$$

387 Where, Back shift operator $L^j X_t = X_{t-j}$, it shift the value for j th lag.

388 **10. Model Calibration**

389

390 Two types of model (prediction model) has developed using white noise series but model
 391 identification and parameter estimation are not done by conventional Box-Jenkins and Yule-
 392 Walker method. In this present study model identification has done by picking up some
 393 candidate ARMA model of order up to ten and five for AR and MA process as for most
 394 hydrologic cases AR parameter (table 7) and MA parameter (table 8). The model selection is
 395 based maximum likelihood estimate (MLE) criteria for prediction model. The underlying
 396 equations for MLE criteria for model selection which are used for present study has given in
 397 following form.

$$398 \quad \text{MLE criteria: } \text{MLE} = -\frac{N}{2} \ln(\sigma_i) - n_i \quad (14)$$

399 Where, N is the total data sets those are used for model calibration, σ_i is the variance of residual
 400 series where residual is the difference between observe data and corresponding to model output
 401 and n_i is the total number of parameter of a model.

402 The parameter, MLE values for candidate models are shown in table below.



403

Table 7: Showing only AR parameter for ARMA model.

MODELS	AR PARAMETERS									
	ϕ_1	ϕ_2	ϕ_3	ϕ_4	ϕ_5	ϕ_6	ϕ_7	ϕ_8	ϕ_9	ϕ_{10}
ARMA(1,0)	0.66451									
ARMA(2,0)	0.60492	0.09057								
ARMA(3,0)	0.60295	0.07795	0.02135							
ARMA(4,0)	0.60080	0.06671	-0.06864	0.15228						
ARMA(5,0)	0.58985	0.07236	-0.07197	0.10985	0.07035					
ARMA(6,0)	0.58063	0.05453	-0.05759	0.10015	-0.01679	0.14551				
ARMA(7,0)	0.58787	0.05451	-0.05024	0.09341	-0.01130	0.18096	-0.06163			
ARMA(8,0)	0.58702	0.05674	-0.05001	0.09497	-0.01299	0.18255	-0.05265	-0.01534		
ARMA(9,0)	0.58727	0.05783	-0.05312	0.09458	-0.01517	0.18504	-0.05505	-0.02823	0.02195	
ARMA(10,0)	0.58604	0.05918	-0.04949	0.08478	-0.01712	0.17845	-0.04629	-0.03598	-0.01817	0.06831
ARMA(1,1)	0.75312									
ARMA(2,1)	1.39970	-0.43496								
ARMA(3,1)	-0.05225	0.52430	0.00826							
ARMA(4,1)	1.24323	-0.31924	-0.11346	0.12318						
ARMA(5,1)	1.33450	-0.37355	-0.12367	0.15158	-0.03725					
ARMA(6,1)	0.28522	0.22950	-0.03610	0.07931	0.01672	0.16687				
ARMA(7,1)	0.74461	-0.03680	-0.05859	0.10082	-0.02648	0.18399	-0.08480			
ARMA(8,1)	-0.20826	0.52546	-0.00582	0.05756	0.05380	0.17932	0.10911	-0.08011		
ARMA(9,1)	-0.19884	0.51957	-0.00903	0.05670	0.05344	0.17970	0.09762	-0.07628	0.01668	
ARMA(10,1)	0.56699	0.07040	-0.04841	0.08390	-0.01537	0.17811	-0.04268	-0.03706	-0.01876	0.06879
ARMA(10,2)	-0.26579	-0.29836	0.50422	0.09413	0.02310	0.25209	0.10769	0.09484	-0.09784	0.05943
ARMA(10,3)	0.51149	-0.07838	0.78081	-0.33277	-0.05451	0.24153	-0.11466	-0.01912	-0.17039	0.13406
ARMA(10,4)	0.38327	0.00193	0.77221	-0.22377	-0.10340	0.23160	-0.08697	-0.02519	-0.17577	0.11709
ARMA(10,5)	-1.13146	-0.06556	0.95650	0.77413	0.07620	-0.10191	0.33892	0.13060	-0.20583	-0.18946
ARMA(1,2)	0.88700									
ARMA(1,3)	0.93403									
ARMA(1,4)	0.93193									
ARMA(1,5)	0.93484									
ARMA(2,2)	-0.02219	0.53793								
ARMA(2,3)	0.21627	0.65017								
ARMA(2,4)	0.19993	0.67908								



ARMA(2,5)	0.23180	0.66432							
ARMA(3,2)	0.64304	0.68470	-0.38491						
ARMA(3,3)	0.42446	0.65453	-0.16819						
ARMA(3,4)	-0.21319	0.15327	0.78757						
ARMA(3,5)	-0.72403	0.74153	0.76954						
ARMA(4,2)	0.69864	0.48430	-0.37004	0.11205					
ARMA(4,3)	-0.37744	1.04098	0.50283	-0.33690					
ARMA(4,4)	-0.42158	1.00669	0.53456	-0.29754					
ARMA(4,5)	0.49824	-0.28955	0.00477	0.62246					
ARMA(5,2)	-0.50475	-0.14324	0.64171	0.13615	0.05173				
ARMA(5,3)	-0.39712	-0.33977	0.55777	0.19782	-0.01648				
ARMA(5,4)	-0.66573	0.97770	0.86047	-0.19813	-0.14756				
ARMA(5,5)	-0.17019	0.61541	-0.04977	-0.06419	0.40870				
ARMA(6,2)	0.35887	-0.35880	0.29042	0.13706	-0.04217	0.24540			
ARMA(6,3)	0.42395	-0.24326	0.70088	-0.18135	-0.08643	0.20389			
ARMA(6,4)	0.39509	-0.20829	0.69730	-0.14510	-0.11156	0.19805			
ARMA(6,5)	0.08679	0.26106	-0.29711	0.48029	0.60613	-0.29824			
ARMA(7,2)	-0.17285	-0.32911	0.46577	0.11014	0.03296	0.19965	0.14650		
ARMA(7,3)	-0.00080	-0.29558	0.54497	0.01053	0.00388	0.20120	0.12022		
ARMA(7,4)	1.36622	-0.64456	0.93643	-0.80636	0.06862	0.28351	-0.21098		
ARMA(7,5)	0.07913	0.23442	-0.31250	0.48226	0.60847	-0.29135	0.02455		
ARMA(8,2)	-0.06562	0.56493	-0.08504	0.05863	0.04780	0.16703	0.08868	-0.09679	
ARMA(8,3)	-0.66504	-0.38677	0.35814	0.33139	0.05602	0.22303	0.23437	0.17652	
ARMA(8,4)	0.96323	-0.34874	0.82988	-0.53430	0.01236	0.20510	-0.08258	-0.08917	
ARMA(8,5)	-1.07677	-0.28472	0.64295	0.49656	0.08941	-0.03977	0.43053	0.33915	
ARMA(9,2)	-0.28074	-0.29950	0.52283	0.11113	0.02688	0.25629	0.13284	0.07980	-0.11560
ARMA(9,3)	-0.40831	-0.32266	0.46804	0.18178	0.04158	0.25586	0.16121	0.10094	-0.10421
ARMA(9,4)	-0.62826	0.98023	0.80696	-0.29996	-0.08852	0.21088	0.05200	-0.13123	-0.07542
ARMA(9,5)	-1.1315	0.65017	0.78081	0.48226	-0.11156	0.28351	0.08868	-0.08917	-0.1758

404

405



406 **Table 8: Showing only MA parameter and constant for ARMA model.**

MODELS	MA Parameters					Constant
	θ_1	θ_2	θ_3	θ_4	θ_5	
ARMA(1,0)						-0.00051
ARMA(2,0)						-0.00129
ARMA(3,0)						-0.00139
ARMA(4,0)						-0.00351
ARMA(5,0)						-0.00456
ARMA(6,0)						-0.00525
ARMA(7,0)						-0.00497
ARMA(8,0)						-0.00477
ARMA(9,0)						-0.00501
ARMA(10,0)						-0.00559
ARMA(1,1)	-0.16221					-0.00215
ARMA(2,1)	-0.82815					-0.00338
ARMA(3,1)	0.66709					0.00311
ARMA(4,1)	-0.66822					-0.00463
ARMA(5,1)	-0.75677					-0.00408
ARMA(6,1)	0.30166					-0.00591
ARMA(7,1)	-0.15794					-0.00470
ARMA(8,1)	0.79691					-0.00455
ARMA(9,1)	0.78870					-0.00475
ARMA(10,1)	0.01917					-0.00562



ARMA(10,2)	0.87471	0.90483				0.00162
ARMA(10,3)	0.19815	-0.77084	0.45834			-0.01064
ARMA(10,4)	0.21564	0.19305	-0.75702	-0.11842		-0.01140
ARMA(10,5)	1.78713	1.23513	-0.25609	-1.00000	-0.67359	-0.03063
ARMA(1,2)	-0.33688	-0.16404				-0.00522
ARMA(1,3)	-0.35536	-0.15329	-0.15629			-0.00608
ARMA(1,4)	-0.35436	-0.15683	-0.15869	0.01476		-0.00611
ARMA(1,5)	-0.35372	-0.15800	-0.15768	0.01995	-0.02017	-0.00608
ARMA(2,2)	0.62654	-0.03828				0.00221
ARMA(2,3)	0.35863	-0.42746	-0.26765			-0.01054
ARMA(2,4)	0.37873	-0.42551	-0.28864	-0.05125		-0.01094
ARMA(2,5)	0.36628	-0.40484	-0.26811	-0.07566	-0.09411	-0.01076
ARMA(3,2)	-0.09104	-0.63486				-0.00590
ARMA(3,3)	0.15828	-0.52132	-0.18806			-0.00865
ARMA(3,4)	0.77956	0.30867	-0.57369	-0.18270		-0.01575
ARMA(3,5)	1.35415	0.09708	-0.76904	-0.36178	-0.14700	-0.01603
ARMA(4,2)	-0.11722	-0.50226				-0.00684
ARMA(4,3)	1.01092	-0.42902	-0.79243			-0.01703
ARMA(4,4)	1.06432	-0.37089	-0.81084	-0.04031		-0.01738
ARMA(4,5)	0.08596	0.39643	0.17633	-0.41673	-0.30529	-0.00931
ARMA(5,2)	1.12739	0.92096				0.01845
ARMA(5,3)	1.02365	1.06681	0.12434			0.01720
ARMA(5,4)	1.31475	-0.15571	-1.00000	-0.28767		-0.01530



ARMA(5,5)	0.78848	-0.10724	-0.07385	0.08787	-0.31487	-0.01845
ARMA(6,2)	0.23672	0.56680				0.00669
ARMA(6,3)	0.16843	0.41609	-0.56455			-0.00732
ARMA(6,4)	0.19811	0.39653	-0.57059	-0.04465		-0.00843
ARMA(6,5)	0.54966	0.11588	0.33089	-0.21262	-0.80621	-0.01506
ARMA(7,2)	0.77331	0.87867				0.00266
ARMA(7,3)	0.59455	0.73423	-0.18396			-0.00062
ARMA(7,4)	-0.78780	0.25034	-0.96854	0.50600		-0.00029
ARMA(7,5)	0.56501	0.13511	0.35062	-0.18961	-0.79343	-0.01516
ARMA(8,2)	0.65268	-0.12361				-0.00555
ARMA(8,3)	1.26941	1.23172	0.38231			0.02268
ARMA(8,4)	-0.37589	0.20412	-0.86433	0.21865		-0.00413
ARMA(8,5)	1.78003	1.50000	0.18756	-0.54292	-0.54283	0.04868
ARMA(9,2)	0.88633	0.91781				0.00390
ARMA(9,3)	1.01837	1.02326	0.13147			0.01121
ARMA(9,4)	1.28045	-0.18837	-1.00000	-0.27298		-0.01673
ARMA(9,5)	1.35415	-0.1533	-0.757	-0.3618	-0.67359	-0.011398

407

408 **10.1 Maximum likelihood rule**

409 A likelihood value for every of the candidate models (table 9) is calculated which model
 410 represents highest likelihood value is chosen for data generation. Gaussian process, general
 411 expression of log-likelihood function for the i^{th} model is given below.

$$412 \quad L_i = \ln(p[z, \hat{\varphi}_i]) - n_i \quad (15)$$



413 It may be approximated as-

$$414 \quad L_i = -\frac{N}{2} \ln(\sigma_i) - n_1 \quad (16)$$

415 Where, L_i is the likelihood value, z represents a vector of historical series i.e. parameter vector,
 416 MA and parameters $(\theta_1, \theta_2, \dots; \phi_1, \phi_2, \dots; \sigma_i)$, σ_i represents the residual variance and n_i is the
 417 number of parameters. As the number of parameters increase, the likelihood value decreases.

418 **Table 9: Showing MLE values constant for ARMA model.**

Models	MLE Values	Models	MLE Values	Models	MLE Values
ARMA(1,0)	6.473	ARMA(10,2)	-14.515	ARMA(5,2)	-16.810
ARMA(2,0)	53.207	ARMA(10,3)	-15.273	ARMA(5,3)	-29.974
ARMA(3,0)	4.812	ARMA(10,4)	-17.592	ARMA(5,4)	-37.916
ARMA(4,0)	4.546	ARMA(10,5)	-51.252	ARMA(5,5)	-13.428
ARMA(5,0)	3.268	ARMA(1,2)	5.395	ARMA(6,2)	-39.489
ARMA(6,0)	0.671	ARMA(1,3)	2.531	ARMA(6,3)	-9.516
ARMA(7,0)	-0.729	ARMA(1,4)	1.738	ARMA(6,4)	-9.926
ARMA(8,0)	-1.772	ARMA(1,5)	0.438	ARMA(6,5)	-18.510
ARMA(9,0)	-2.787	ARMA(2,2)	2.875	ARMA(7,2)	-22.491
ARMA(10,0)	-2.796	ARMA(2,3)	-7.629	ARMA(7,3)	-17.672
ARMA(1,1)	6.031	ARMA(2,4)	-7.240	ARMA(7,4)	-20.179
ARMA(2,1)	1.747	ARMA(2,5)	-12.101	ARMA(7,5)	-18.904
ARMA(3,1)	-30.943	ARMA(3,2)	-7.923	ARMA(8,2)	-17.325
ARMA(4,1)	1.620	ARMA(3,3)	-8.838	ARMA(8,3)	-26.489
ARMA(5,1)	0.698	ARMA(3,4)	-9.603	ARMA(8,4)	-13.709



ARMA(6,1)	-0.723	ARMA(3,5)	-43.814	ARMA(8,5)	-39.369
ARMA(7,1)	-1.542	ARMA(4,2)	-5.464	ARMA(9,2)	-26.378
ARMA(8,1)	-4.142	ARMA(4,3)	-29.525	ARMA(9,3)	-26.855
ARMA(9,1)	-5.747	ARMA(4,4)	-30.252	ARMA(9,4)	-44.180
ARMA(10,1)	-4.396	ARMA(4,5)	-6.041	ARMA(9,5)	-37.972

419

420 11. Model Validation

421 In the present study ARMA(2,0) (table 9) models has selected as one time step ahead and
 422 prediction model by Maximum MLE criteria respectively. The selected model is validate to
 423 examine whether the assumptions used for selection of the model are valid.

424 11.1 Significance of residual mean

425 This test examines the validity of the assumption that the error series $e(t)$ has zero mean. A
 426 statistic $\eta(e)$ is defined as:

$$427 \quad \eta(e) = \frac{N^{1/2} \bar{e}}{\rho^{1/2}} \quad (17)$$

428 Where, \bar{e} = Estimated residual mean.

429 ρ = Estimated residual variance.

430 The statistic $\eta(e)$, approximated distribution as $t(\alpha, N-1)$, α represents the significance level at
 431 test is being carried out. If the value of $\eta(e) < t(\alpha, N-1)$, (table 10) then the mean of the residual
 432 series is not significantly different from zero (-)ve series passes the test.

433

Table 10: Showing the statistic $\eta(e)$.



434

ρ	α	N	$t(\alpha, N-1)$	$\eta(e)$
0.161472	.05	30	1.699	0.1

435

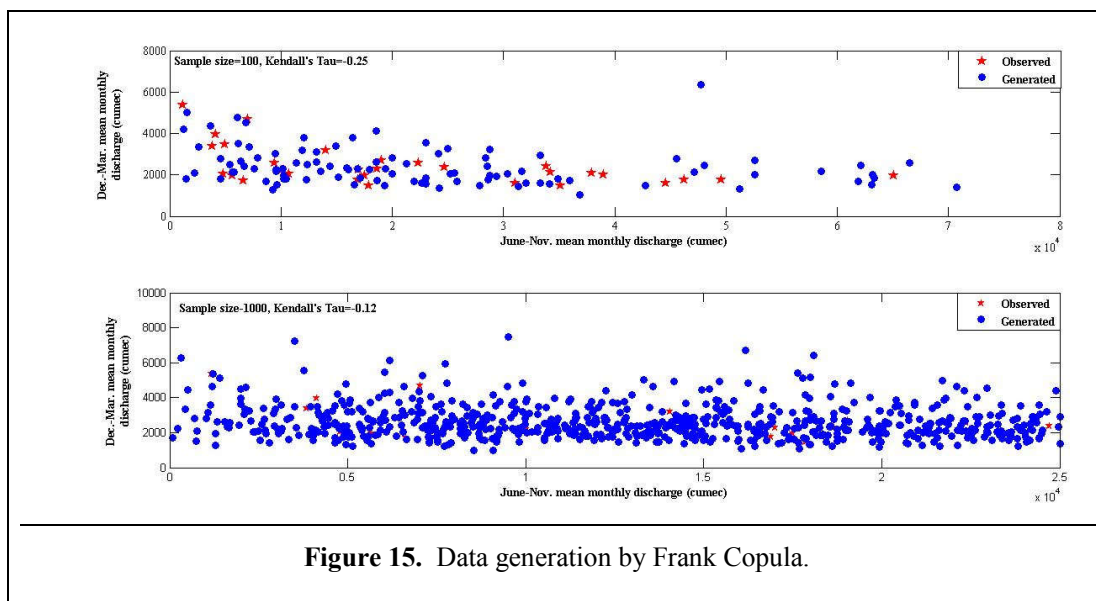
436

437 The value of $\eta(e) < t(\alpha, N-1)$ (i.e. $1.22 < 1.699$). It shows that mean of the residual series is not
438 significantly different from zero (-) series passes the test.

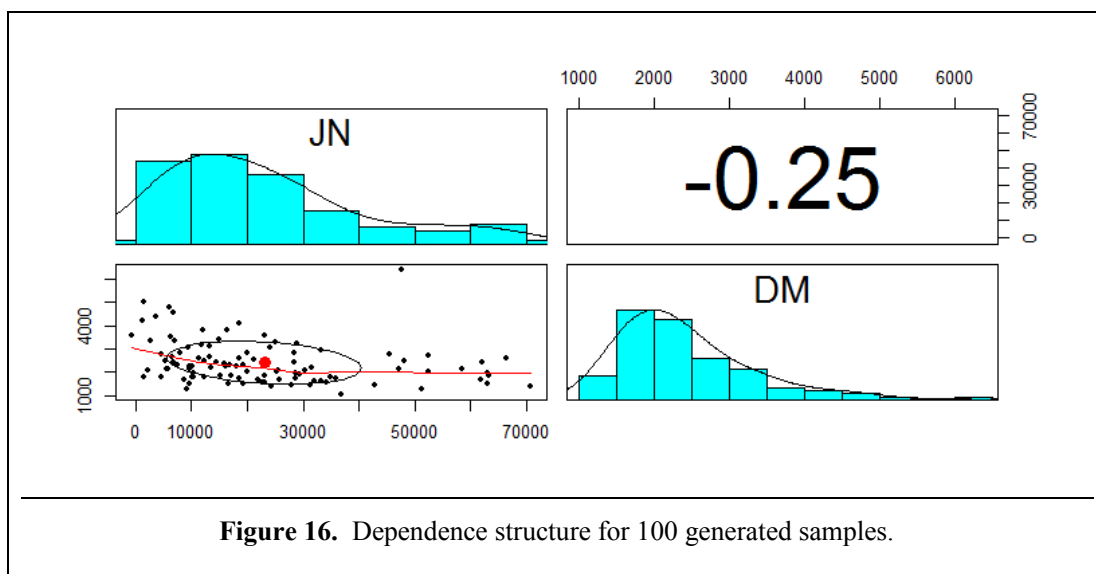
439 12. Results and Discussion

440 Outflow data for future are generated by using Frank Copula. The sample size for data
441 generation is taken 100 and 1000. These generated values are compared with observed validation
442 data set. The comparison and the individual dependence of generated samples are also shown in
443 Figure 15, which represents observed data in red color and generated data in blue, data
444 generation by copula with sample size 100 and 1000 of kendall's tau 0.25 and 0.32, respectively.
445 Figure 16 describes dependence structure for 100 generated samples i.e. Copula is a statistical
446 theory on dependence and measurement of association.

447



448
 449
 450



451
 452

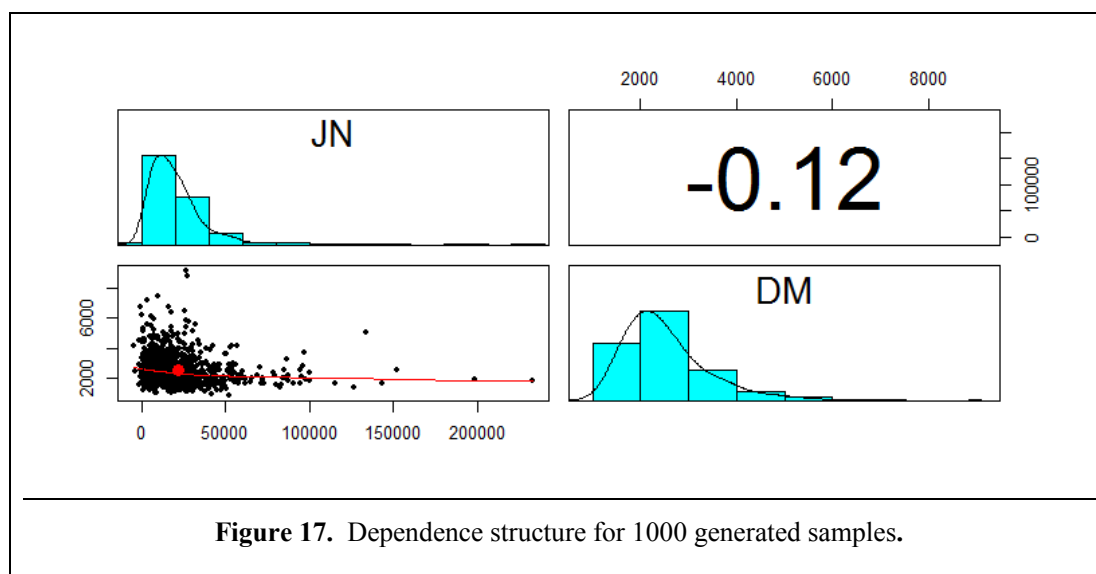


Figure 17. Dependence structure for 1000 generated samples.

453
454

455 Figure 17 describes dependence structure for 1000 generated samples i.e. Copula is a statistical
456 theory on dependence and measurement of association and teal colour shows generated data set
457 month wise. Dependence structure of a multivariate distribution is described by copula, it might
458 be appropriate to use measures of dependence which are copula-based. Linear correlation
459 coefficient can be opposed by the concordance measure spearman's rho and kendall's tau as well
460 as tail dependence and it is expressed by under laying copula.

461 Figure 20 is a time series of discharge data, blue colour represents observed data from jan. 1968
462 to dec. 1973 and red colour represents generated data from jan. 1974 to dec. 2004. The green
463 colour shows observed data set and red colour shows generated data set. When generated data set
464 is small, it shows good results because errors incorporated are less in comparison to large data
465 set generation.

466 In the study area, has analyzed that best model in ARMA (2,0) model, and Frank Copula model
467 for generating discharge data at Farakka barrage on the basis of Mean Square Error (MSE),



468 Coefficient of determination (R^2) (Figure 18 and Figure 19). Based upon all above test Frank
469 Copula is the best model for generating outflow discharge data at Farakka barrage.

470

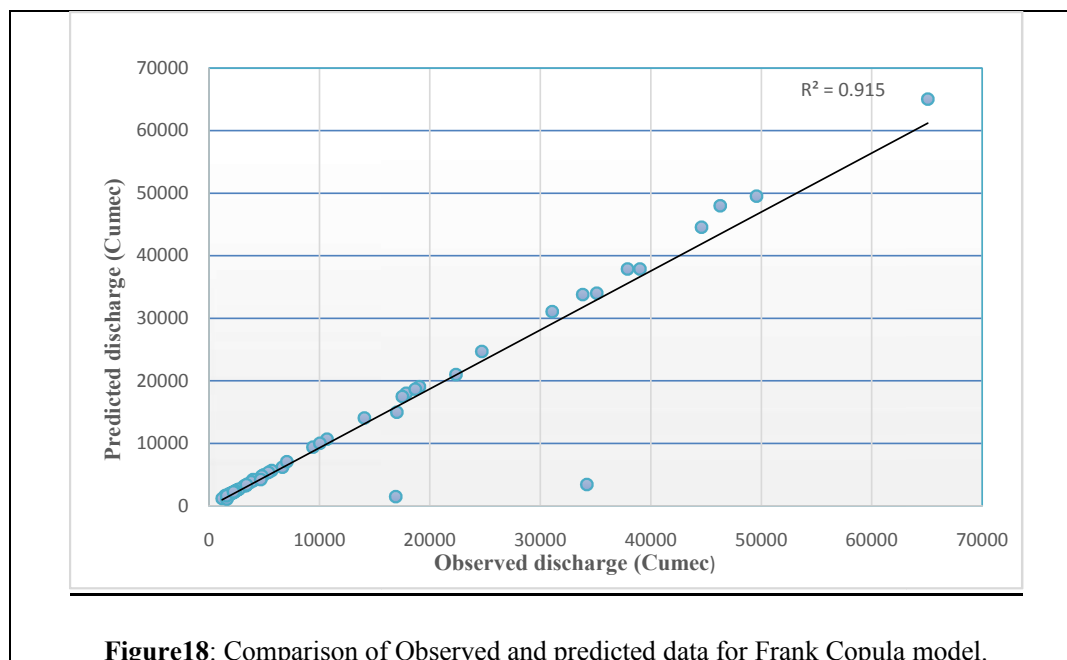


Figure18: Comparison of Observed and predicted data for Frank Copula model.

471

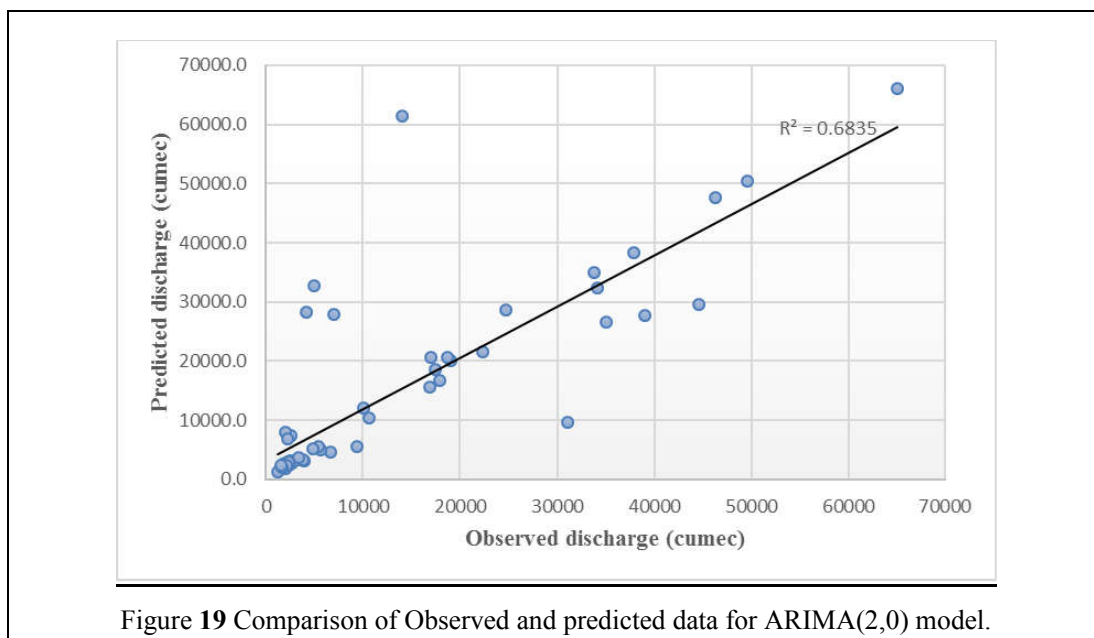
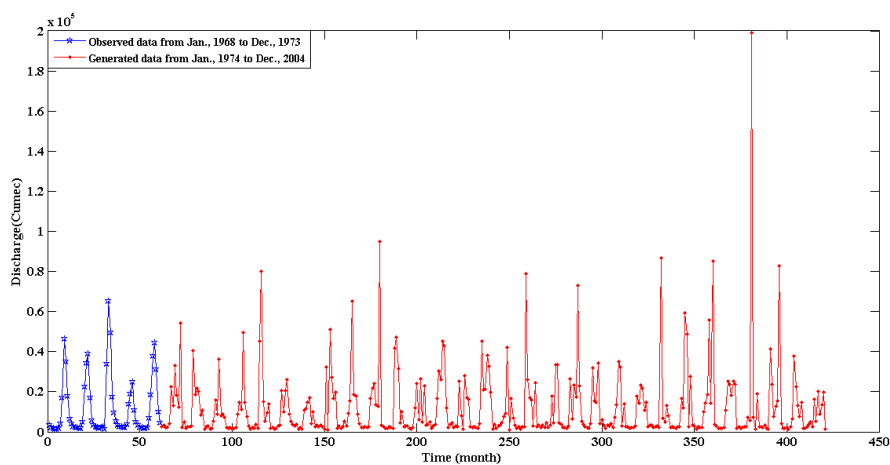


Figure 19 Comparison of Observed and predicted data for ARIMA(2,0) model.

472

473



474

475

Figure 20. Time series of observed data and generated data.

476

477



478 **13. Conclusion**

479

480 Copula based study and ARMA models are used in this study Frank copula is selected among the

481 copulas based on parameter for discharge data generation and ARMA(2,0) is selected among the

482 ARMA models. Errors incorporated in the copula model is less in comparison to the ARMA(2,0)

483 model and the value of Coefficient of determination (R^2) for Frank is close to one i.e 0.915.

484 Frank copula estimated better result over ARMA(2,0).

485 A copula based study which can be used to derive bivariate distribution function of flow rate

486 variables and it shows the real world case study. Best suited model for this study is frank copula

487 among all above copula in term of non parametric tests i.e. AIC, MSE, BIC and Kolmogrov-

488 Smirnov test. When generated data sample data set, copula shows convergence of sample data

489 set to estimated population. Copula models are an alternative approach and in this study Frank

490 Copula model is used for data generation at Farakka barrage. Bivariate series are prepared based

491 on pre monsoon and post monsoon outflow data. Moreover they are very useful in this study of

492 dependent variable. Copula is very useful for describing the dependence of extreme outcome

493 because it captures the structural dependence of data. The autocorrelation is not captured in the

494 bivariate model .

495

496

497

498

499

500

501



502 **References**

503

504 Akaike, H.: A new look at the statistical model identification. IEEE transactions on automatic
505 control, 19(6), 716-723. 1974.

506

507 Balistrocchi, M., Orlandini, S., Ranzi, R., & Bacchi, B.: Copula-Based Modeling of Flood
508 Control Reservoirs, *Water Resources Research*, 53(11), 9883-9900, 2017.

509 Chowdhary, H., Escobar, L. A., Singh, V. P.: Identification of suitable copulas for bivariate
510 frequency analysis of flood peak and flood volume data. *Hydrology Research*, 42(2–3), 193–216,
511 2011.

512 De Michele, C., Salvadori, G., Canossi, M., Petaccia, A., & Rosso, R.: Bivariate statistical
513 approach to check adequacy of dam spillway. *Journal of Hydrologic Engineering*, 10(1), 50-57,
514 2005.

515 Favre, A. C., El Adlouni, S., Perreault, L., Thiémonge, N., & Bobée, B.: Multivariate
516 hydrological frequency analysis using copulas, *Water resources research*, 40(1), 2004.

517 Genest, C., Favre, A. C., Béliveau, J., & Jacques, C.: Metaelliptical copulas and their use in
518 frequency analysis of multivariate hydrological data. *Water Resources Research*, 43(9). W09401.
519 DOI:10.1029/2006WR005275, 2007.

520 Genest, C., Rémillard, B., & Beaudoin, D.: Goodness-of-fit tests for copulas: A review and a
521 power study. *Insurance: Mathematics and economics*, 44(2), 199-213. 2009.

522 Genest, C., & Rivest, L. P.: Statistical inference procedures for bivariate Archimedean copulas, *J*
523 *of the American statistical Association*, 88(423), 1034-1043, 1993.

524 Ghosh, S.: Modelling bivariate rainfall distribution and generating bivariate correlated rainfall
525 data in neighbouring meteorological subdivisions using copula, *Hydrological Processes*, 24(24),
526 3558-3567, 2010.



- 527 Goel, N. K., Seth, S. M., & Chandra, S.: Multivariate modeling of flood flows. *Journal of*
528 *Hydraulic Engineering*, 124(2), 146-155, 1998.
- 529 Hooshyaripor, F., Tahershamsi, A., & Golian, S.: Application of copula method and neural
530 networks for predicting peak outflow from breached embankments. *Journal of Hydro-*
531 *Environment Research*, 8(3), 292-303, 2014.
- 532 Joe, H.: *Multivariate models and multivariate dependence concepts*. CRC Press, Chapman and
533 Hall, London, 1997.
- 534 Kao, S.-C., and Govindaraju, R. S.: Trivariate statistical analysis of extreme rainfall events via
535 Plackett family of copulas, *Water Resour. Res.*, 44(2), W02415, 2008.
- 536 Kao, S. C., & Chang, N. B.: Copula-based flood frequency analysis at ungauged basin
537 confluences: Nashville, Tennessee, *Journal of Hydrologic Engineering*, 17(7), 790-799, 2011.
- 538 Karmakar, S., and Simonovic, S. P.: “Bivariate flood frequency analysis. Part 2: A copula-based
539 approach with mixed marginal distributions”, *J. Flood Risk Manage.*, 2(1), 32–44, 2009.
- 540 Kashyap, R.: A Bayesian comparison of different classes of dynamic models using empirical
541 data, *IEEE Transactions on Automatic Control*, 22(5), 715-727, 1977.
- 542 Katz, R. W., & Skaggs, R. H.: On the use of autoregressive-moving average processes to model
543 meteorological time series, *Monthly Weather Review*, 109(3), 479-484, 1981.
- 544 Li, T., Guo, S., Chen, L., & Guo, J.: Bivariate flood frequency analysis with historical
545 information based on copula, *Journal of Hydrologic Engineering*, 18(8), 1018–1030, 2013.
- 546 Muhaisen, O. S., Osorio, F., & García, P. A.: Two-copula based simulation for detention basin
547 design, *Civil Engineering and Environmental Systems*, 26(4), 355-366, 2009.
- 548 Mohan, S., & Vedula, S.: Multiplicative seasonal ARIMA model for longterm forecasting of
549 inflows, *Water resources management*, 9(2), 115-126, 1995.



- 550 Osorio, F., Muhaisen, O., & García, P. A.: Copula-based simulation for the estimation of optimal
551 volume for a detention basin, *Journal of Hydrologic Engineering*, 14(12), 1378-1382,2009.
- 552 Poulin, A., Huard, D., Favre, A. C., & Pugin, S.: Importance of tail dependence in bivariate
553 frequency analysis, *Journal of Hydrologic Engineering*, 12(4), 394-403,2007.
- 554 Razmkhah, H., AkhoundAli, A. M., Radmanesh, F., & Saghafian, B.: Evaluation of rainfall
555 spatial correlation effect on rainfall-runoff modeling uncertainty, considering 2-copula, *Arabian
556 Journal of Geosciences*, 9(4), 323,2016.
- 557 Renard, B., & Lang, M.: Use of a Gaussian copula for multivariate extreme value analysis: some
558 case studies in hydrology, *Advances in Water Resources*, 30(4), 897-912,2007.
- 559 Requena, A. I., Chebana, F., & Mediero, L.: A complete procedure for multivariate index-flood
560 model application, *Journal of Hydrology*, 535, 559-580,2016.
- 561 Requena, A. I., Flores, I., Mediero, L., & Garrote, L.: Extension of observed flood series by
562 combining a distributed hydro-meteorological model and a copula-based model, *Stochastic
563 environmental research and risk assessment*, 30(5), 1363-1378,2016.
- 564 Requena AI, Mediero L, Garrote L.: A bivariate return period based on copulas for hydrologic
565 dam design: accounting for reservoir routing in risk estimation, *Hydrol Earth Syst Sc* 17:3023–
566 3038,2013.
- 567 Requena, A. I., Prosdocimi, I., Kjeldsen, T. R., & Mediero, L.: A bivariate trend analysis to
568 investigate the effect of increasing urbanisation on flood characteristics, *Hydrology
569 Research*, 48(3), 802-821,2017.
- 570 Salvadori, G., & De Michele, C.: On the use of copulas in hydrology: theory and practice,
571 *Journal of Hydrologic Engineering*, 12(4), 369-380,2007.



- 572 Salvadori, G., De Michele, C., Kottegoda, N. T., & Rosso, R.: Extremes in nature: an approach
573 using copulas (Vol. 56), Springer Science & Business Media,2007.
- 574 Shiau, J. T., Feng, S., & Nadarajah, S.: Assessment of hydrological droughts for the Yellow
575 River, China, using copulas, *Hydrological Processes*, 21(16), 2157-2163,2007.
- 576 Sklar A.: Fonction de re'partition a' n dimensions et leurs marges, vol. 8. Publications de
577 L'Institute de Statistique, Universite' de Paris: Paris; 229–231,1959.
- 578 Twaróg, B.: An assessment of risks posed by control rule parameters implemented in a flood
579 control reservoir, carried out with the application of elements of ruin theory and of bivariate
580 distribution of a random variable based on the copula function, *American Journal of*
581 *Environmental Engineering*, 6(2), 62-71,2016.
- 582 Vrugt, J. A., Diks, C. G., Gupta, H. V., Bouten, W., & Verstraten, J. M.: Improved treatment of
583 uncertainty in hydrologic modeling: Combining the strengths of global optimization and data
584 assimilation, *Water resources research*, 41(1),2005.
- 585 Wang, Q. J.: A Bayesian joint probability approach for flood record augmentation, *Water*
586 *Resour. Res.*, 37(6), 1707–1712,2001.
- 587 Zhang L, Singh V.P.: Trivariate flood frequency analysis using the Gumbel–Hougaard copula, *J.*
588 *Hydrol Eng* 12:431–439,2007.
- 589 Zhang, L. S. V. P., & Singh, V. P.: Bivariate flood frequency analysis using the copula method,
590 *Journal of hydrologic engineering*, 11(2), 150-164,2006.