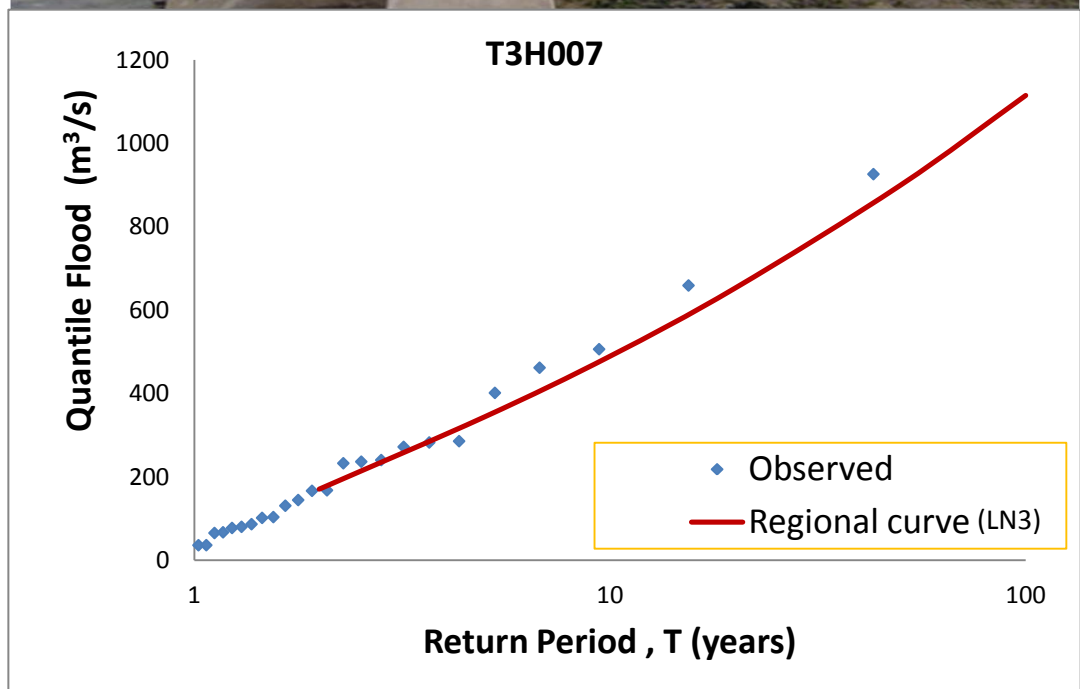


# Regional Flood Frequency Analysis in Southern Africa

Alem Tadesse Haile



**UNIVERSITY OF OSLO**

**FACULTY OF MATHEMATICS AND NATURAL SCIENCES**



# Regional Flood Frequency Analysis in Southern Africa

Alem Tadesse Haile



Master Thesis in Geosciences

Discipline: Hydrology

Department of Geosciences

Faculty of Mathematics and Natural Sciences

UNIVERSITY OF OSLO

**1. September 2011**

© Alem Tadesse Haile, 2011

Supervisors: Lars Gottschalk (UiO) and Lena M. Tallaksen (UiO)

This work is published digitally through DUO – Digitale Utgivelser ved UiO

<http://www.duo.uio.no>

It is also catalogued in BIBSYS (<http://www.bibsys.no/english>)

All rights reserved. No part of this publication may be reproduced or transmitted, in any form or by any means, without permission.

Cover photo: the photograph of Mzimvubu River at Ku-Makuhola gauging site (T3H007): from <http://www.dwaf.gov.za/Hydrology/CGI-BIN/HIS/CGIHis.exe/Photo?Station=T3H007> (Retrieved, 28/08/2011) and the lower graph indicates the distribution of the AMS of the river (from 1972-2008).

## Acknowledgments

First and foremost, I would like to thank my supervisors Prof. Lars Gottschalk and Prof. Lena M. Tallaksen for their consistent supervision, crucial advices, invaluable suggestions and unreserved cooperation throughout the project work and the preparation of the thesis. Their support and encouragement from the beginning up to the completion of this study is kindly appreciated. I also express my deep gratitude and appreciation to Prof. Chong-Yu Xu, for his kind supports in data collections and communications during the project work.

I am grateful to all my lecturers in the section of Physical Geography, Hydrology and Geomatics for their academic supports and constructive feedbacks. Many thanks goes to my study group mates most especially, Nils Charles Prieur and Søren Nykjær Boje in the Hydrology discipline who make up a significant diversity from different context (open minded, introducing with different cultures, supporting in any faced difficulties that enabled me to develop alternative models of thinking, an open minded culture and staying confident without feeling alone). Your team spirit, academic and social input will be always remaining an extraordinary and forever missed. I would always be indebted to you all.

I have a great pleasure in thanking my family and all friends for their being with me and for their continuous encouragement. Especially Hadgu Girmay, Tesfamariam Birhane, Kibrom Araya, Libaragachew Demile, Samai Sanon and others, your discussions and friendship were valuable academically, spiritually and socially throughout my studies and it will forever be appreciated and unforgettable. I am always proud and honorable because you all are my kind friends.

**“አባት ዘይብሉ አይገዛካ ደቂዓዲካ ዘብሉ አይዓዲካ ከም ዚበሃል እዩ እሞ ኒዚገበርኩምላይ ኩሉ ካብ ልቢ የመስግን”**

Last but not least, I would like to pass my great acknowledgments to all Geosciences staffs, students and library workers for their cooperation in giving necessary information and technical supports. I am also thankful to Quota scheme and staff workers that they gave me the chance for M.Sc. study at university Oslo, Norway, and the Norwegian government (lånekassen) for financial supports during my studies.



## Abstract

Extreme floods are natural disasters often associated with losses of life, and severe impact to agricultural production and infrastructures. However, efficient estimations of the magnitude of such extreme events with their non-exceedance probabilities, either for design or risk management planning purposes, are often limited by the data availability (i.e., both in quality and quantity). In this study, regional frequency analyses of annual maximum series (AMS) of flood events from unregulated rivers of southern Africa were conducted. This includes preliminary data analysis (data screening and outlier analysis), sensitivity analysis, identification of homogenous regions and suitable regional distribution models for the regions, development of regional growth curves and regression models to estimate the quantile floods for ungauged catchments. The study area comprises five countries (459 stations): Namibia, Malawi, Zambia, Zimbabwe and South Africa. The AMS derived for each station were examined for validity, dependency, and the existence of outliers. After thorough examinations of the AM flood events, 122 gauging sites were selected for further analysis. The study area was divided into nine possible homogenous regions based on the geographical grouping method together with the heterogeneity tests. The AMS from Namibia, Zimbabwe, Zambia and Malawi were grouped as regions *R1*, *R2*, *R3* and *R4*, respectively, while the South African catchments were further classified into five possibly homogenous regions. The choice of an appropriate regional flood distribution model was performed based on L-moment approaches together with the index flood procedures and goodness-of-fit (GOF) tests. The Generalized Pareto (GPA), Pearson type III (PE3), Three-parameter lognormal (LN3) and the Generalized Extreme Value (GEV) distributions were found to be suitable models for AMS of floods in southern Africa catchments. Regional flood frequency curves were constructed based on the best regional distribution for the nine regions and design floods estimated for the return periods of 2-500 years. Based on assessments the accuracy of the derived quantile-quantile plots, it was concluded that the performance of this regional approaches was satisfactory and also confirmed when validated against sites not included in the regional analysis.

**Keywords** Outliers; index flood; regional flood frequency analysis; L-moment; Southern Africa





# I. TABLE OF CONTENTS

I. TABLE OF CONTENTS.....	i
II. LIST OF TABLES .....	iv
III. LIST OF FIGURES .....	vi
1. INTRODUCTION .....	1
1.1 General background.....	1
1.2 Objective of the study.....	4
1.3 Target Group.....	4
1.4 Limitation of the study .....	5
1.5 Structure of the Thesis .....	5
2. STUDY AREA AND DATA .....	7
2.1 Study area description .....	7
2.1.1 Geography of Southern Africa .....	8
2.1.2 Climate and Vegetation .....	10
2.1.3 Hydro climatology.....	12
2.2 Data collection.....	15
2.2.1 Data source .....	15
2.2.2 Selection of flood data: Annual Maximum Series (AMS).....	16
2.2.3 Site characteristics.....	17
3. THEORY AND METHODOLOGY.....	19
3.1 Background.....	19
3.1.1 Methods of RFFA.....	19
3.1.2 Procedures of RFFA.....	21
3.2 Exploratory data analysis.....	21
3.2.1 Data screening .....	22
3.2.2 Empirical distribution.....	26
3.2.3 Outlier detection and treatments .....	28
3.2.4 Index flood method .....	32
3.3 Regionalization.....	34
3.3.1 Delineation of homogeneous regions .....	34

3.3.2	Homogeneity Test .....	35
3.4	Choice of regional flood frequency distribution.....	36
3.4.1	Theoretical distribution functions .....	36
3.4.2	Fitting the regional data to empirical distribution .....	39
3.4.3	L-moment and L- moment ratio diagram .....	40
3.4.4	Goodness-of-fit (GOF) test .....	42
3.5	Regional flood frequency curve .....	44
3.6	Evaluation the performance of frequency distributions .....	45
3.6.1	Quantile-quantile (qq) plots .....	46
3.6.2	Growth curve verification .....	46
3.7	Regional estimation for ungauged catchments .....	47
4.	RESULT AND ANALYSIS .....	49
4.1	Exploratory data analysis.....	49
4.1.1	Data Screening .....	49
4.1.2	Autocorrelation and Spatial Correlation.....	50
4.1.3	Empirical distribution.....	52
4.1.4	Outlier Analysis.....	53
4.1.5	At-Site flood characteristics .....	57
4.1.6	Choice of the Index Flood .....	57
4.2	Identification of homogenous regions .....	59
4.2.1	Delineation of homogenous regions.....	59
4.2.2	Heterogeneity test.....	60
4.2.3	Regionalization of outliers .....	61
4.2.4	Outlier sensitivity analysis .....	64
4.3	Identification of regional flood frequency distribution .....	65
4.3.1	The L-moment ratio diagram .....	65
4.3.2	Goodness-of- fit (GOF) measures .....	68
4.4	Regional flood frequency curves .....	72
4.5	Performance evaluation using simulation.....	77
4.5.1	Quantile-Quantile (qq) plots.....	77
4.5.2	Verifications of the regional flood frequency curves.....	78

4.6	Regional estimation for ungauged catchments .....	81
5.	DISCUSSION .....	82
5.1	Data and outlier analysis.....	82
5.2	Regional Homogeneity .....	85
5.3	Regional flood frequency distribution.....	87
5.4	Regional flood frequency curve .....	89
5.5	Performance evaluation of empirical distributions.....	91
5.6	Estimation of design floods from ungauged catchments.....	92
6.	CONCLUSION AND RECOMMENDATIONS .....	94
7.	LIST OF REFERENCES .....	97
	APPENDICES.....	100
A.	Selected stations .....	100
1.	Stations used for regional flood frequency analysis.....	100
2.	Stations used for model verifications.....	103
B.	At- site statistical behaviors of annual maximum floods .....	104
C.	The quantile flood of the normalized observed versus simulated values .....	107
D.	The regional Regression of at-site median values .....	110
E.	Theatrical distributions and their relationships .....	112

## II. LIST OF TABLES

Table 2.1 Annual rainfall, evapotranspiration and surface runoff for Southern Africa countries: from Pallet et al. (1997).....	13
Table 2.2 Summary of the daily runoff data available from five countries of Southern Africa .....	15
Table 2.3 Summary of the AMS data selected from five countries of Southern Africa .....	17
Table 3.1 Some of the popular probability plotting models: <i>where 'r' is rank of the observed value to be plotted, n is number of observation.</i> .....	27
Table 3.2 Summary of the employed distributions and their parameter .....	39
Table 4.1 Summary of the outlier analysis in Malawi gauging sites .....	56
Table 4.2 The sensitivity analysis of the index floods to the largest observations in stations which contains one or two large outliers.....	58
Table 4.3 The result of Hosking and Wallis' (1997) homogeneity test and the summery of their regional statistics of Southern Africa Catchments.....	61
Table 4.4 Summery of regional outliers of Southern African floods.....	62
Table 4.5 Unaccepted outliers.....	63
Table 4.6 The comparisons of the relative differences of the at-site and regional weighted average statistics after suspected outliers were removed for the sample series. $\theta_1$ and $\theta_2$ are parameters values estimated from the sample series before and after the suspected outliers were excluded from the series, respectively. ....	64
Table 4.7 The result of Anderson - Darling goodness-of-fit test at 10 % level of significant (Viglione, 2010) .....	69
Table 4.8 Hosking and Wallis (1997) goodness-of-fit test statistics for regional frequency distribution. ....	70
Table 4.9 Regional weighted average L-moments for the grouped regions of southern Africa .....	73
Table 4.10 Summary of the regional growth curves which includes the non-exceedance probability, the best distribution models for respective regions and their parameter values and quantile floods for a range of recurrence intervals.....	74
Table 4.11 Selected stations for model verifications and their index floods .....	78
Table 4.12 Derived regression models to predict the median values from catchment characteristics in Southern Africa. ....	81
Table 1 Site characteristics of Namibia catchments (from 1969-2004).....	100
Table 2 Site characteristics of Malawi rivers (from 1954-1990) .....	100
Table 3 Site characteristics of Zambia Rivers (from 1970-2004).....	100

Table 4 Site characteristics Zimbabwe catchments (From 1957-1990).....	100
Table 5 Site characteristics of South Africa Rivers (From 1969-2008).....	101
Table 6 Selected stations for model verifications .....	103
Table 7 at site statistical characteristics for Namibia Catchments .....	104
Table 8 At-site statistical characteristics for Zimbabwe Catchments .....	104
Table 9 At-site statistical characteristics of Zambia Catchments .....	104
Table 10 At-site statistical characteristics for Malawi Catchments .....	105
Table 11 At-site statistical characteristics for South Africa Catchments.....	105
Table 12 Theoretical distribution functions and their moments. Taken from Geo4310 lecture notes by Gottschalk and Krasovskaia (2001) and Hosking and Wallis (1997): $x =$ observed values, $m =$ mean value, $\sigma =$ standard deviation, $C_s =$ coefficient of variance $\alpha =$ scale parameter, $\mu =$ location parameter and $k =$ shape parameter.....	112
Table 13 Theoretical relationships of L-moments and the inverse of the some cumulative distribution function (Gottschalk and Krasovskaia, 2001; Hosking and Wallis, 1997).	113

### III. LIST OF FIGURES

<i>Figure 2.1 Geographic locations of the southern African countries. Taken from Southern Africa FRIEND.....</i>	<i>7</i>
<i>Figure 2.2 Typical regions in the Okavango Delta, with free canals and lakes, swamps and islands: Taken from Wikipedia, the free encyclopedia. ....</i>	<i>8</i>
<i>Figure 2.3 The Kalahari Desert (shown in maroon) &amp; Kalahari Basin (orange): Taken From Wikipedia, the free encyclopedia .....</i>	<i>9</i>
<i>Figure 2.4 Köppen-Geiger climate classification map of Southern Africa. Adapted from Peel et al. (2007) .....</i>	<i>11</i>
<i>Figure 2.5 Regional distribution of precipitation: Taken from webpage of The Kunene River Awareness Kit.....</i>	<i>12</i>
<i>Figure 2.6 Distribution of regional surface water drainages. Taken from the website of The Kunene River Awareness Kit.....</i>	<i>14</i>
<i>Figure 2.7 The spatial distributions of Stations used for this analysis .....</i>	<i>18</i>
<i>Figure 3.1 Generalized Extreme Value (GEV) distributions: Gumbel (EV1, <math>k=0</math>), Frechet (EV2, <math>k=-0.04</math>) and Weibull (EV3, <math>k=0.04</math>) distribution (Gottschalk and Krasovskaia, 2001).....</i>	<i>38</i>
<i>Figure 4.1 Plotting observed data series from Zambia station ‘1591470’ (1970-2004): a) daily average runoff series; and (b) AMS. ....</i>	<i>50</i>
<i>Figure 4.2 the auto-correlation plots for station ‘1591001’: the left is for the daily time series for two years time lag and the right is for the AMS for the time lag of 35 years at 95% confidence interval (the dotted lines at the right plot). ....</i>	<i>51</i>
<i>Figure 4.3 An example of Gumbel Probability Plotting for the normalized annual maximum floods from three stations in South Africa: (a) station ‘K9H001’ which contains lower bounded observation; (b) station ‘X2H006’ with upper bounded normalized series; (c) station ‘U2H048’ contains annual maximum flood series with outlier; and (d) station ‘U2H048’ presents annual maximum flood series after an outlier has been removed. ...</i>	<i>52</i>
<i>Figure 4.4 Examples of visual inspections of the AMS containing suspected outliers for station B8H010 in South Africa: (a) Time series plotting of the annual maximum series (AMF); (b) Gumbel plotting positions; and (c) Histogram of the annual extreme events. ....</i>	<i>54</i>
<i>Figure 4.5 Delineation of southern Africa catchments into hydrologically homogenous region. The further classifications of South African drainage areas are shown in the right side of the map. The abbreviation NA indicates the countries or regions which have no available data. ....</i>	<i>60</i>

<i>Figure 4.6 L-moment ratio diagram for the annual maximum floods from Malawi gauging sites. The diagram shows the influence of a single outlier in station '1992100' (Table 4.5) in case of fitting theoretical distributions to the regional data. ....</i>	<i>63</i>
<i>Figure 4.7 L-moment diagrams showing the relationships between the theoretical distribution curves and the regional data from five countries of Southern Africa: the name of regions is labeled under the pictures from (a-i). ....</i>	<i>67</i>
<i>Figure 4.8 Regional flood frequency curves for 9 regions in Southern Africa: the title of each curve indicates the name of the regions. The curves were developed from best fitted distribution of respective regions in Table 4.10. ....</i>	<i>76</i>
<i>Figure 4.9 Examples of quantile-quantile plots of the normalized empirical discharge against the simulated values from the best fitted distributions: a) Pearson type III (PE3) for Region R1-Nambia; and b) Generalized Pareto distribution (GPA) for ZA_R5 .....</i>	<i>77</i>
<i>Figure 4.10 shows the comparison of the probability plots of the quantile floods between the observed series (Doted) and estimated values from the best fitted of regional frequency curves (solid line) .....</i>	<i>80</i>
<i>Figure 1 plotting the normalized quantile values of the observed against randomly simulated using best fitted regional distribution .....</i>	<i>109</i>
<i>Figure 2 The regional regression coefficients showing the relationships between the index flood (median) and catchments area. ....</i>	<i>111</i>





# 1. INTRODUCTION

## 1.1 General background

Extreme events, such as floods are among the catastrophic natural events that cause severe consequences for human society. In many countries of the world, floods are causing damages to properties and agricultural lands that result in huge economic and life losses for the affected areas. For example, in Southern Africa (the study area), it is often reported in WebPages such as UN news center<sup>1</sup> that every country of the region are on alert for potentially disastrous flooding. The UN Office for the Coordination of Humanitarian Affairs (OCHA) in January 27, 2011 warned that floods in Southern Africa could be severe and lead to food shortages. Five countries (Botswana, Mozambique, Namibia, Zambia and Zimbabwe) have also recently forecasted serious flooding phenomenon that could affect tens of thousands of people, and damage infrastructure, crops and homes (UNNC, 2011).

How frequently a flood event of a given magnitude may be expected to occur is of great important, because almost every activities on a particular flooded areas might be controlled by it (Hosking and Wallis, 1997). The frequency of floods with various risks of exceedance, are therefore needed for a wide range of engineering problems, planning for weather-related emergencies, reservoir management, pollution control, and insurance risk calculations (Gottschalk and Krasovskaia, 2001; Kjeldsen et al., 2002; Saf, 2008). Estimation must be fairly accurate not only aimed at the preventing of catastrophes, but also at avoiding excessive costs in case of overestimating the flood magnitude, or excessive damage while underestimating the flood potential.

Flood frequency analysis is a hydrologic field dealing with estimation of a flood magnitude corresponding to any required return period of occurrence. Based on experience, people have some idea as to how often floods of a given size occur at given places. Hydrologists have been attempting to formalize these ideas by establishing networks of gauging stations and analyze the recorded information (Hipei, 1994). In hydrological events, there are numerous and unpredictable sources of uncertainties about the physical processes (Hosking and Wallis, 1997). Thus, stochastic models (such as flood frequency analysis) are very important and desirable to estimate how often a specified event will occur on average in a particular area.

---

<sup>1</sup> <http://www.un.org/apps/news/story.asp?NewsID=37347&Cr=flood&Cr1>

This is due to the fact that statistical methods are acknowledging the existence of uncertainties and enable its effects to be quantified by confidence intervals.

The frequency analysis of extreme events from a single site is well established and might be easier than at the regional level. However, it is most often case that many related samples having the same statistical behavior may available at different measuring sites. A more appropriate estimation could then be to analyze all the data samples together than using only individual series. This approach is known as regional frequency analysis (RFA) (Hosking and Wallis, 1997). Regional flood frequency analysis (henceforth RFFA) may be practiced in a joint use of at-site and regional data. The method assumes that the extreme events at several sites in a region may have similar statistical characteristics (Cunnane, 1989). The author also suggested that, though the assumption of homogeneity of the regions is a gross simplification, the method is convenient and effective. The advantages of regional approaches are also frequently illustrated in the literature (Farquharson et al., 1992; Gottschalk and Krasovskaia, 2001; Hosking and Wallis, 1997; Kachroo et al., 2000; Kjeldsen et al., 2002; Mkhandi and Kachroo, 1997; Mkhandi et al., 2000; Rosbjerg, 2007; Saf, 2008; Saf et al., 2008; Shu and Ouarda, 2008; Wiltshire, 1986). These studies suggested that RFFA is more reliable estimation of design floods for two fundamental reasons: (1) due to short and uneven record lengths, the regional data of homogenous regions have smaller standard error than those estimated at individual station data only; and (2) it has the ability to estimate design floods for the homogenous regions and allow estimation from gauged sites to ungauged sites.

Nowadays, hydrologists have been using the advanced method of regional flood frequency analysis which compromises the use of L-moments together with the index-flood method (Hosking and Wallis, 1997; Saf, 2008) . For example, the methodology has been successfully applied in Southern Africa flood studies such as RFFA studies for South Africa and Botswana (Farquharson et al., 1992); Southern Africa (Mkhandi and Kachroo, 1997); and South Africa (Kjeldsen et al., 2002).

A RFFA is based on the recorded observations from sites in homogenous region and then a single form distribution is fitted to the pooled data (NERC, 1975). For flood modeling, a range frequency distributions have been suggested, but none has been accepted as universal distribution (Mkhandi and Kachroo, 1997). For example, the survey by Cunnane (1989) suggested that the RFFA studies conducted in a number of countries aimed at selecting a “best” national distribution for annual maximum series (AMS) recommended different

distributions such as Log Pearson Type Three (LP3) distribution for USA (USWRC, 1981), Generalized Extreme Value (GEV) for UK flood studies (NERC, 1975), LP3 for Australia (Institution of Engineers, 1977), Two Component Extreme Value (TCEV) for Italy (Rossi et al., 1984) and Two Parameter Log-Normal (LN2) distribution for Canada flood studies (Spence, 1973).

In addition to the studies above, very few studies on RFFA for southern Africa are documented in the literature. The most notable works are the technical document of RFFA for Southern Africa (Mkhandi and Kachroo, 1997) and flood frequency analysis for Southern Africa catchments (Mkhandi et al., 2000). The authors found that the L-moment diagram to be an appropriate analytical tool for the identification of a suitable frequency distribution together with goodness-of-fit tests. For the delineated Southern Africa homogenous regions (Kachroo et al., 2000), the Pearson Type III (PE3) with Probability Weighted Moments (PWM) and/or LP3 with maximum likelihood (ML) methods of parameter estimators were recommended as the appropriate flood model. However, these studies may not being sufficient when scaling down to the country levels. For example, the later study by Kjeldsen et al. (2002) concluded that the appropriate regional flood frequency distribution for South Africa particularly in the KwaZulu-Natal province could be the GPA distributions.

Having the above extreme value theory and pervious outcomes as a motivation, this study is aimed at extracting information as much as possible from the available runoff data series and previous studies. Thus, the outcomes of this work can be provided additional inputs for the improvement of the flood hydrology in Southern Africa. The study had attempted to improve some inputs used for the analysis such as the record lengths of runoff data; and choices on the procedures of data analysis. The work has been accomplished by implementing different methods/inputs of RFFA that can be useful in designing flood problems corresponding to specified exceedance probability or simply risk. As a result, the flood models and flood magnitudes corresponding to required recurrence intervals were furnished for the catchments of the Southern Africa. The spatial and temporal variability of flood events with respect to regional climatic variables and catchment structure were also identified.

## **1.2 Objective of the study**

The main objective of this study was to analyze flood frequency distribution for homogeneous regions in Southern Africa which may serve as a basic input to improve the design and economic appraisal of civil engineering structures, and to have optimum land use planning and /or decreasing risk due to flood damages.

The specific tasks that helped to achieve the overall objective of the study were:

- Exploratory data analysis ( Data screening and Outlier analysis)
- Grouping gauging sites into hydrologically homogeneous regions
- Identify an appropriate theoretical distribution of flood flows in Southern Africa
- Develop regional frequency curves for the delineated regions and
- Regional estimation for ungauged catchments

## **1.3 Target Group**

This master thesis contributes to the NUFU - Water Sciences project; Water resources and hydrological extremes theme. The overall goal of the project is to improve human welfare by efficient utilizations of the inadequate resources of the community through improved access and availability of healthy and safe water (NUFU, 2010) i.e.,

- Based on an inventory of existing data and earlier work, identify emerging tasks within flood and drought research addressing the need of the regions
- Identify the variability of hydrological extremes, flood and drought, with respect to the available information such as regional climatic variables and catchment characteristics
- Develop maps that show the spatial behavior of extreme hydrological events using a combination of data sources with high resolution satellite data.

The project has been implementing by conducting basic and operational research to address and relate health to water quality, availability, climate change and poverty through postgraduate research and training at Master and PhD levels. Currently, this project has been implementing in Malawi. It is expected that the communities in the selected study sites in Malawi will attain improved health and welfare and the innovations replicated to other parts of Southern Africa including Botswana, Malawi and South Africa (NUFU, 2010).

## 1.4 Limitation of the study

The main constraint of the study was collecting sufficient runoff data information both in quality and quantity. The region of Southern Africa has 12 countries including Madagascar (see for details in section 2.1). However, for half of the region which includes 6 countries (Lesotho, Swaziland, Tanzania, Mozambique, Botswana and Madagascar), the author couldn't find sufficient runoff data for RFFA. The runoff stations in these countries have insignificant discharge (zeros and nearly zero values) and very short length of records (from 3 to 6 years length). Besides this, the data series available from the stations (even the stations used for analysis) have also a lot of information gaps (in some stations it is more than 10 years). The reason might be due to the following three sources: (1) frequent and sustainable dry seasons, i.e., most of the seasons are dry that the record indicates a lot of zeros and nearly zero values (especially in Botswana), (2) suitability of the sites for measuring and (3) political and economical problem.

Since the study area is located in the arid and semi-arid zones, many of the problems were associated with estimating floods such as the difficulties of measuring flood flows and the variability of flood events (Farquharson et al., 1992). The Authors also illustrated that the difficulty of establishing a reasonable rating curve-particularly at high flow levels is the worst problems in this area. This may arise due to the access for gauging near the peak of a short flood, the long periods without flow, and the instability of the channel control and cross-section area owing to the scouring effect of floods. Hence, the uncertainty of the data should be considered during that analysis.

## 1.5 Structure of the Thesis

The thesis has six sections and the outline of each section is presented as follows:

**Section 1: Introduction-** introduces the general backgrounds and relevant previous findings, presents the objective and motivates the thesis; introduce the contribution of this work done and the aims of the projects; and the outline of the subsequent sections. Some term definitions, importance and applications of flood frequency analysis are also introduced in this section.

**Section 2: Study area and data-** presents the study area and the detailed activities accomplished during data collection and preparation for analysis.

**Section 3: theory and methodology of the study-**describes the procedures and the methods of the study and gives the theoretical background of each method used. It starts with the procedures of data screening and examinations, reviews the necessity and application of the methods used for RFFA and statistical test which were used during the analysis. It reviews the types of models, the procedures and their approaches.

**Section 4: Result and analysis-** presents the details of the main outputs of the research. It presents the analysis of the data behavior, outlier detection and treatments, groups of homogenous regions, choice of best fitted regional distribution, development of regional flood frequency curves and the quantile flood of the rivers and other result for example, evaluation the performance of distribution functions, the regional parameter values and regional L-moments and L-moment diagrams, the sensitivity analysis were presented and analyzed under this section.

**Section 5: Discussion-**discussed the methods, results and the choice with respect to their theoretical backgrounds.

**Section 6: Conclusion and recommendations -** the conclusions reached in the research are presented in this section. In addition, the recommendations for the future researchers that should be focused on are given in this section.

Finally, the references and appendices are presented in the last pages of the thesis.

## 2. STUDY AREA AND DATA

### 2.1 Study area description

Southern Africa is a region located in the southernmost of the African continent which covers total Area of 6,938,014 km<sup>2</sup>. The region comprises the countries: Angola, Botswana, Lesotho, Malawi, Mozambique, Namibia, South Africa, Swaziland, Zambia, Tanzania and Zimbabwe. The entire land of Southern Africa is varied, ranging from forests and grasslands to deserts. The region has both low-lying coastal areas, and mountains. The natural resources of the region is the world's largest source of elements such as platinum and the platinum group elements like chromium, and cobalt, as well as uranium, gold, titanium, iron and diamonds (Wikipedia, 2011).



*Figure 2.1 Geographic locations of the southern African countries. Taken from Southern Africa FRIEND<sup>2</sup>*

---

<sup>2</sup> <http://www.ru.ac.za/static/institutes/iwr/friend/?request=institutes/iwr/friend>

### 2.1.1 Geography of Southern Africa

Southern Africa is located in the southern part of the African continent and is bordering to: east-coastal plains of Mozambique and Tanzania with Indian Ocean; south-coastal areas South Africa with Southern ocean; west-Angola and Namibia with Atlantic Ocean; and north-the inlands of Democratic Congo and Kenya countries. The Geography of southern Africa consists of a series of undulating plateaus that cover most of South Africa, Namibia, and Botswana and extend into central Angola. Contiguous with this are uplands in Zambia and Zimbabwe. The Coastal Mountains and escarpments which flank the high ground are also found in northern Mozambique, South Africa, Namibia, Angola, and along the Mozambique-Zimbabwe border. Southern Zimbabwe and much of South Africa are within a region of scrublands and grasslands known as the Veld<sup>3</sup>. To the southeast of the Veld is the Drakensberg range-the main mountain range of Southern Africa. The Drakensberg rises to more than 3,475 meters and extends roughly northeast to southwest for 1,125 km parallel to the southeastern coast of South Africa. This includes the region's highest mountain-Lesotho's mount Ntlenyana with an elevation 3,482 m.a.s.l (meters above mean sea level) (SouthernAfrica, 2011).



*Figure 2.2 Typical regions in the Okavango Delta, with free canals and lakes, swamps and islands: Taken from Wikipedia, the free encyclopedia<sup>4</sup>.*

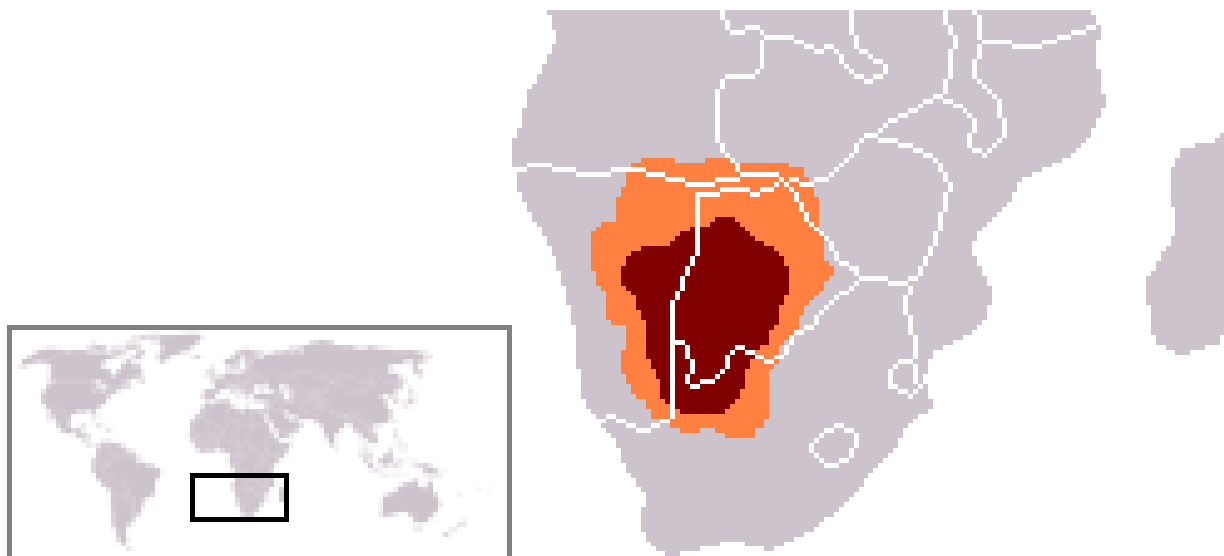
<sup>3</sup> (Afrikaans: "field") -it is a name given to various types of open country in [Southern Africa](#) that is used for pasturage and farmland.

<sup>4</sup> <http://upload.wikimedia.org/wikipedia/commons/6/61/Okavango11.jpg>



The region contains a lot of unique geographical and geomorphologic features such as the Okavango Delta (or Okavango Swamp), in Botswana which is the largest inland delta; the third largest desert called Kalahari; and the largest salt pans of the Makgadikgadi Pan in Botswana and Etosha Pan in Namibia. The pan is all that remains of the formerly huge lake Makgadikgadi, which once covered an area larger than Switzerland, but dried up several thousand years ago (Wikipedia, 2011).

The Kalahari Desert is the largest desert of the region which extends 900,000 km<sup>2</sup> from the arid to semi-arid sandy area in Southern Africa, covering much of Botswana and parts of Namibia and South Africa (see Fig. 2.3) (Wikipedia, 2011). Fig. 2.3 shows the extent of the desert with the orange color indicates the surrounding Kalahari Basin which covers over 2,500,000 km<sup>2</sup>. As it can be seen from the figure the drainage of the desert is extending farther into Botswana, Namibia and South Africa, and encroaching into parts of Angola, Zambia and Zimbabwe. It forms the central depression of the Southern African plateaus. Its elevation rises to the great escarpment, which flanks the plateau almost unbroken line from the Zambezi River to Angola.



*Figure 2.3 The Kalahari Desert (shown in maroon) & Kalahari Basin (orange): Taken From Wikipedia, the free encyclopedia<sup>5</sup>*

The second largest desert in the region is the Namib Desert which extending 1,900 km from Namibia, Angola, along the entire coast of Namibia to the Olifants River in South Africa. It is almost rainless area, 80 –130 km wide over most of its length. It is mainly a smooth platform

---

<sup>5</sup> <http://upload.wikimedia.org/wikipedia/commons/b/bc/LocationKalahari.PNG>

of bedrock of various types and ages. In the southern half, the platform is covered with sand. The eastern part, the inner Namib, supports large numbers of ruminant mammals like antelope. The shore area is densely populated by marine birds, including Flamingos, Pelicans, and Penguins (SouthernAfrica, 2011).

### **2.1.2 Climate and Vegetation**

The driving elements of the hydrologic cycle are the temporal and spatial distribution of water, the intensity of precipitation, temperatures and many other physical and chemical processes that shape the landscape. Climate is perhaps the most important driver with respect to determining the amount, distribution and the availability of water in the environment. It is known that climate can be commonly defined as the weather averaged over a period of around 30 years of a particular region and mainly affected by the latitude, topography, altitude, ice or snow cover, as well as nearby water bodies and their currents.

The Southern African climates are seasonal, ranging from arid to semi-arid and from temperate to tropical. According to Peel et al. (2007), the climate of the region can be broadly divided into two Köppen climate Groups:

- I) Class B** - Dry climates including the southwestern countries bordering the Kalahari Desert including the Angola, Botswana, Zimbabwe, Namibia and South Africa countries with climates ranging from semi-arid and sub-humid in the east to hyper-arid in the west parts.
- II) Class C** - Moist mid-latitude climates with mild winters which include the eastern countries: Tanzania, Malawi, Mozambique, Swaziland, Lesotho and the Indian Ocean island countries, with climatic conditions ranging from Dry to Moist Subtropical Mid-Latitude conditions.

The region is located between the Atlantic and Indian Oceans on the west and east, respectively. These are high pressure zones and played impotent role in the region's climate. Angola and Namibia on the west coast are influenced by the cold Benguela current from the Atlantic Ocean, which produces a drier climate. By contrast, the east coast is influenced by the southward-flowing Mozambique current, which brings warm water and humid air from the Equator and creates a humid, warm climate (KRAK, 2011).

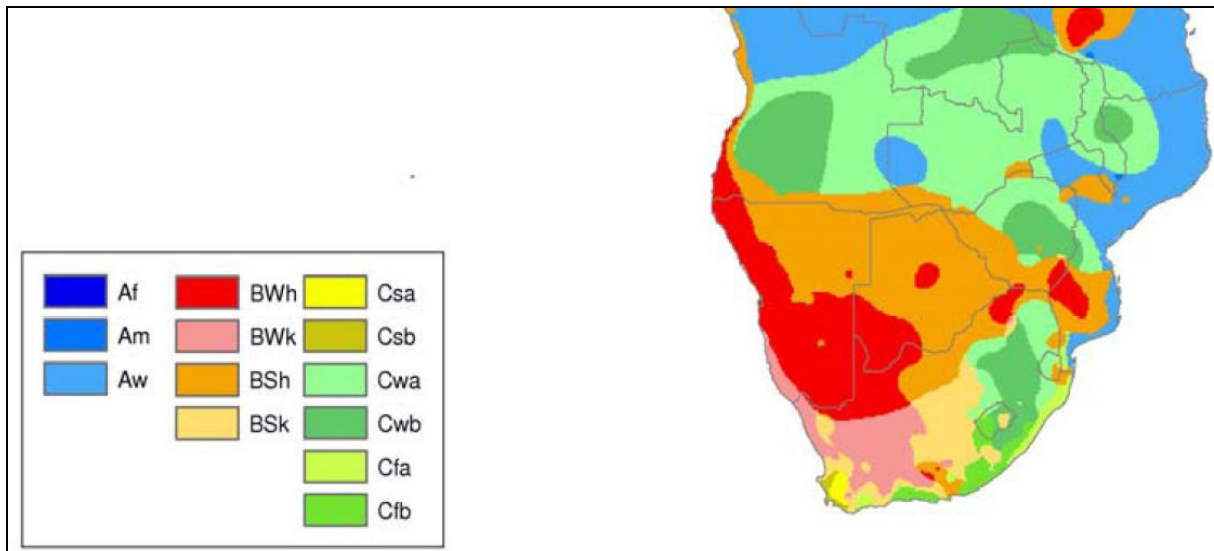


Figure 2.4 Köppen-Geiger climate classification map of Southern Africa. Adapted from Peel et al. (2007)

The region has two distinct seasons – a wet season roughly from November to April and a dry season from May to October. It is prone to frequent droughts and uneven rainfall distribution. There is a strong rainfall gradient from east to west parts interior of southern Africa. In Swaziland and Lesotho to the east, both altitude and exposure to moist air coming off the Indian Ocean produce the heaviest and most reliable rainfall. The total rainfall of the region gradually decreases westward, so that much of the central and western regions are semi-desert with low and variable rainfall over the whole of this interior region, rainfall mainly occurs in the summer season in the form of thunderstorms. There are also large daily and seasonal temperature ranges as a result of the effects of altitude and “continental” position (the lack of ocean influences). Winters are usually dry and sunny while summers are wet and hot (KRAK, 2011).

The seasonality of the climate is therefore the main control of the hydrological regime on plant growth of the region. On the favor of this seasonal climate, there are mainly four types of vegetation: savanna woodlands (known as *miombo* forest) in the north, a series of dry woodlands to the south of arid and semi-arid grassland, scrubland, and bush land in the Namib and Kalahari deserts and their environs, and Mediterranean vegetations along the southern coast (SouthernAfrica, 2011).

### 2.1.3 Hydro climatology

The hydro climatology of southern Africa described for this work includes the precipitation, evapotranspiration, surface water distribution and the drainage of the rivers.

#### A. Precipitation and Evapotranspiration

The region has variable precipitation levels ranging from low (< 250 mm/yr) over large parts, to relatively high (> 1200 mm/yr), which tends to be concentrated in the north of the Southern African Development Community (SADC) with some smaller areas along the south-east coast (KRAK, 2011). Most rain falls in the summer months which are most commonly from December to March with the exception of the Western Cape of South Africa, which has a temperate climate. Rainfall of the region is highly variable in intensity and distribution, particularly high degree in the drier regions (Pallett et al., 1997). When rain falls, it is often periodic, arriving in short intense rainstorm during warm weather.

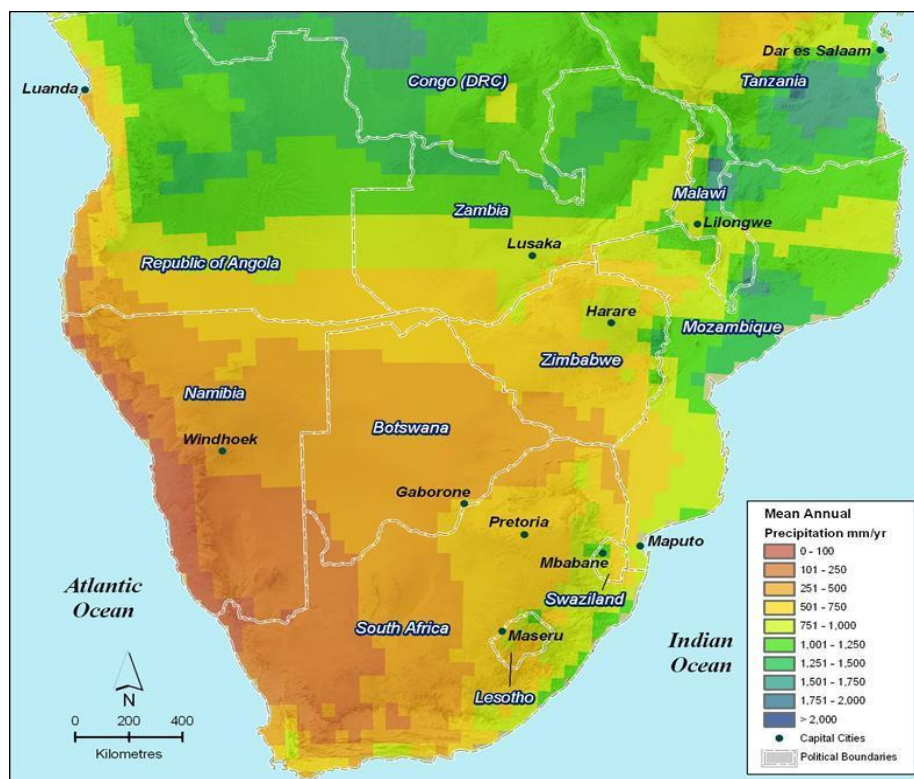


Figure 2.5 Regional distribution of precipitation: Taken from webpage of The Kunene River Awareness Kit<sup>6</sup>.

<sup>6</sup> <http://www.kunenerak.org/en/river/hydrology/hydrology+of+southern+africa.aspx>

Rain falling in intense downpours often runs off into river channels as it falls faster than can be absorbed into the soil. Table 2.1 summarizes the rainfall, evapotranspiration and surface runoff statistics for the region.

Table 2.1 Annual rainfall, evapotranspiration and surface runoff for Southern Africa countries: from Pallet et al. (1997)

Country	Rainfall range	Average Rainfall		Potential evapotranspiration		Total surface runoff	
		mm	Mm	mm	Mm	mm	10 <sup>3</sup> m <sup>3</sup>
Angola	25-1600	800	997	1300-2600	104	130.0	
Botswana	250-650	400	233	2600-3700	0.6	0.35	
Lesotho	500-2000	700	21	1800-2100	136	4.13	
Malawi	700-2800	1000	119	1800-2000	60	7.06	
Mozambique	350-2000	1100	879	1100-2000	275	220.0	
Namibia	10-700	250	206	2600-3700	1.5	1.24	
South Africa	50-3000	500	612	1100-3000	39	47.45	
Swaziland	500-1500	800	14	2000-2200	111	1.94	
Tanzania	300-1600	750	709	1100-2000	78	74.0	
Zambia	700-1200	800	602	2000-2500	133	100.0	
Zimbabwe	350-1000	700	273	2000-2600	34	13.1	
Total			4665			599.27	

The Southern Africa has extremely high water losses from evaporation and evapotranspiration, with only a small percentage of rainfall reaching aquifers through groundwater recharge or surface water through run-off (Pallett et al., 1997). For example, it can be seen from Table 2.1 that in all countries of Southern Africa, the annual potential evapotranspiration is higher than the annual precipitation.

## B. Surface water and Drainages

The surface resources are unevenly distributed across the region with Namibia and in particular Botswana has very sparse surface water resources. As shown in Fig 2.6, many of the water channels across the region, especially those in areas of low rainfall, high temperatures and high rates of evaporation are not permanent rivers, only flow after the intense rainfall events that characterize precipitation in the region. However, the courtiers of South Africa, Zambia, Mozambique and Angola contain relatively good surface runoffs.



*Figure 2.6 Distribution of regional surface water drainages. Taken from the website of The Kunene River Awareness Kit<sup>7</sup>*

The region is generally drained eastward towards the Indian Ocean, a pattern exemplified by the largest rivers, the Zambezi and Limpopo. The Zambezi is the longest river in the region, and its catchment includes much of Angola, Zambia, and Zimbabwe. The only major river flowing into the Atlantic Ocean and passing through both desert areas and connecting three countries is the Orange River. This river rises in the Lesotho Highlands as the Sinqu River, flows west as the Orange across South Africa, and finally to Atlantic Ocean. It passes the southern edge of Kalahari Desert and winds through the Nimbi Desert before draining into the Atlantic Ocean in South Africa, which serves as a border between South Africa and Namibia. It is about 2,100 km long and drains parts of South Africa, Lesotho, and Namibia (SouthernAfrica, 2011). There is also one river called Okavango Rives, which permanently flows to the northwest of Okavango Delta. This river forms important marshes that are rich in wildlife (Wikipedia, 2011).

<sup>7</sup> <http://www.kunenerak.org/en/river/hydrology/hydrology+of+southern+africa.aspx>

## 2.2 Data collection

The regional flood frequency study aimed to include data from all Southern Africa countries such as Angola, Botswana, Lesotho, Malawi, Mozambique, Namibia, South Africa, Swaziland, Tanzania, Zambia and Zimbabwe (Fig. 2.1). However, due to the difficulties involved in obtaining data, the data provided for this study were only from five countries: Malawi, Namibia, South Africa, Zambia and Zimbabwe. From these countries, daily average runoff data from 459 gauging stations with average record length of 35-40 years were collected. The size of the gauged catchments ranges from 72.8 to 850,530 km<sup>2</sup>. A country-wise breakdown of the data is presented in Table 2.2.

Table 2.2 Summary of the daily average runoff data available from five countries of Southern Africa

S.N	Country	Area (km <sup>2</sup> )	No. stations	Catchment area (km <sup>2</sup> )	Data source
1	South Africa	1,221,037	342	119-850530	Webpage Of DWA
2	Zambia	752,618	55	110-284538	GRDC
3	Malawi	118,484	23	72.8-149500	Glad (2010)
4	Namibia	825,418	30	3810-334000	GRDC
5	Zimbabwe	390,757	9	277-5307	SADC- project
Total			<b>459</b>		

### 2.2.1 Data source

The critical issue during the data collection was to find sufficient data of good quality. An attempt, i.e., both officially and personally was made to find the required streamflow information from different data sources. However, in most of the data sources it was impossible to find the sufficient information even for the countries which are listed in Table 2.2. Hence, after considerable efforts, data from a total of 459 stations which contain mean daily runoff data were collected from four different sources. The sources are: (1) for Zimbabwe catchments, nine stations were available from the SADC- project; (2) for Namibia and Zambia, 85 stations were obtained from GRDC (The Global Runoff Data Centre, 56068



Koblenz, Germany), (3) for South Africa, 342 stations were downloaded from the webpage<sup>8</sup> of Department of Water Affairs , South Africa; and (4) for Malawi, 23 stations were collected from Glad (2010). Glad (2010) discussed that the daily average runoff data for Malawi catchments are provided by The Ministry of Irrigation and Water Development in Malawi, and The FRIEND program (Flow Regimes from Experimental and Network Data).

### **2.2.2 Selection of flood data: Annual Maximum Series (AMS)**

The point of departure in design stochastic models (such as RFFA) is having one (or several) observation series. In RFFA, our concern is to analyze the flood characteristics based on the extreme events of the pooled daily time series. The extreme value theory can therefore provide a theoretical basis for selecting the required extreme series. The popular methods used for extreme event selection are the Peak Over Threshold (POT) method - all values higher than a predefined threshold level is chosen (Lang et al., 1999) and Annual Maximum Series (AMS) – is a typical example of block maxima method of extreme value theorem (Engeland, 2005). The block maxima method selects maximum extreme events for each block. In most RFFA, a block is considering as a year, thus the highest daily flow data within a year is chosen (i.e., AMS) (Rootzen and Tajvidi, 2006)..

The choice of the methods depends on the behavior of the data available and use of flood models. For this work, however, the AMS was adapted in agreement with the discussion of Cunnane (1989), that the choice of the AM series was not based on any objective manner rather based on the following advantages; the method is widely accepted, convenient to apply, consistent, and less sensitive to outliers and subjectivity.

For all the stations listed in Table 2.2, the AMS data were selected and later subjected for exploratory data analysis in order to choose representative stations for the study area (the details of the methods will be discussed in section 3.2.1). Finally, the total number of stations was reduced from 459 to 122 (112 for RFFA and 10 for model validation) and the national break down of the stations is presented in Table 2.3.

---

<sup>8</sup> <http://www.dwaf.gov.za/Hydrology/CGI-BIN/HIS/CGIHis.exe/Station>



Table 2.3 Summary of the AMS data selected from five countries of Southern Africa

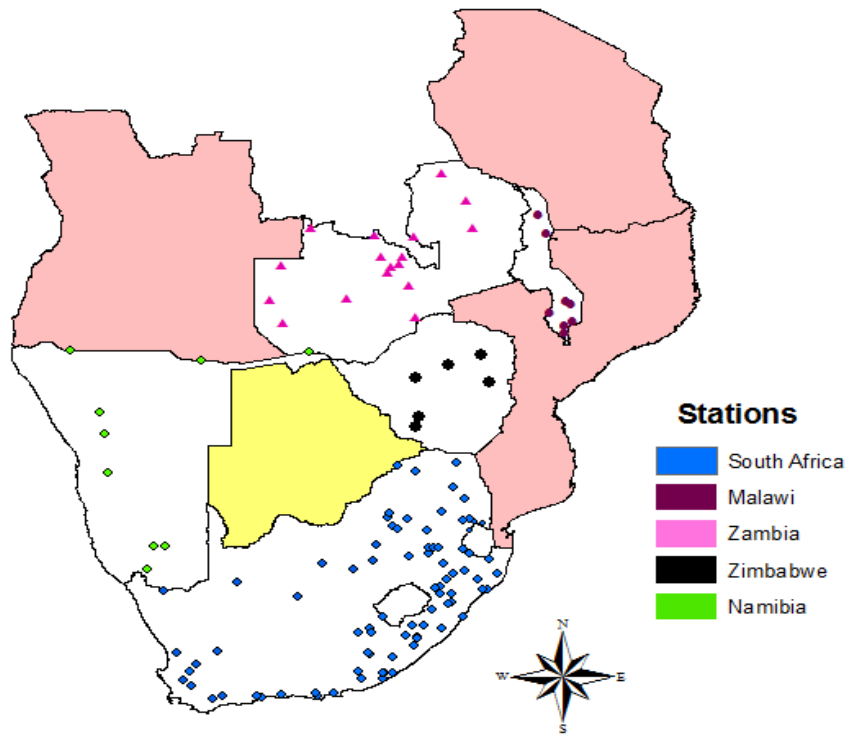
S.N	Country	N. stations <sup>9</sup>	Record period	No. of years
1	South Africa	83 (8)	1969- 2008	40
2	Zambia	17 (2)	1970-2004	35
3	Malawi	8	1957-1990	37
4	Namibia	8	1969-2004	36
5	Zimbabwe	7	1954-1990	37
Total		<b>122</b>		

However, during the whole process of choosing stations and selecting extreme observations from the daily time series, there was no any control to assess/quantify the uncertainties that may arise from the data available. Some of the stations that were collected from different sources had some redundancy though the data for the some station did not match. Hence, the data from different sources were used as inputs for the confirmation of the simple observation on reliability of the data. This phenomenon was due case for some stations of South Africa and then, finally, the stations were selected by cross checking with the official data from the webpage of the water affairs of South Africa.

### 2.2.3 Site characteristics

The site characteristics of the selected stations for this study are presented in the Appendices A (Tables: 1-5). These Tables include the code of the stations, name of river and their gauging sites, the locations (both latitude and longitude in degree), the catchment area coverage (km<sup>2</sup>) for each selected stations for this analysis. Fig. 2.7 shows the locations of the stations that were used in the detailed analysis. The GIS tool so called ArcGIS - ArcMap software was used to plot the location of the stations for each country. Even though the stations in Malawi, Namibia and Zimbabwe are limited in number and insufficient to represent the flood situations in each of the countries, the overall distribution of the stations within the region is satisfactory.

<sup>9</sup> The stations in the brackets were used for validation



*Figure 2.7 The spatial distributions of Stations used for this analysis*

### 3. THEORY AND METHODOLOGY

#### 3.1 Background

Regional flood frequency analysis (RFFA) is an approach to estimate the quantile floods,  $Q_T$  (i.e., the flood magnitude of  $Q$  corresponding to a given recurrence interval  $T$ ) for any site in a region. The magnitude of  $Q_T$  is expected to be expressed in terms of flood data recorded at all gauging sites in the region. However, RFFA has also the ability to include sites which did not have sufficient data available or ungauged catchments in the region. Some RFFA methods, mainly index flood method (section 3.2.4), assumes that a region is a set of catchments in which its flood frequency and parameter behavior is homogeneous in some quantifiable manner. RFFA take advantage of this homogeneity to produce quantile estimates which, in most cases, are more trustworthy than those obtainable from at-site data alone (Cunnane, 1989; Hosking and Wallis, 1997; Mkhanda et al., 2000). Because observed event are short and most likely uneven, this assumption can play substantial roles in reducing errors during quantile estimation and extrapolate the estimations beyond the recorded return periods. A preference that can be made from this discussion is that, estimated quantiles from regional data analysis could be more trustworthy than those estimated from individual series (Cunnane, 1988).

##### 3.1.1 Methods of RFFA

Since the early 1960, around 12 methods have been developed and the details are briefly summarized in the literature by Cunnane (1988 and 1989). The author also illustrated that the development of these models were based on the definitions and notations briefed in the following paragraphs. Most of the methods were based on use of annual maximum (AM) series while a few are based on peaks over a threshold (POT).

Let  $\{Q_{ij}, j = 1, 2, \dots, n_i, i = 1, 2, \dots, M\}$  be annual maxima at  $M$  gauging sites with total AM observations of  $n_i$  in site  $i$  and a total of  $L = \sum_{i=1}^M n_i$  station years of record in a region. For any site  $i$ , the usual assumption is that  $(Q_{ij}, j = 1, 2, \dots, n_i)$  is a random sample from the same parent population. Most RFFA used the normalized series of the data. The scale factor which

is known by index flood ( $\mu_i$ ) of station  $i$  is the median or the mean of the series. This parameter is used for normalized the series in to dimensionless variate and/or up scaling the regional estimation to at-site quantile flood estimation (see section 3.2.4 for details).

The dimensionless of the data is in the form  $X_{ij} = Q_{ij}/\mu_i$ , which is referred to the Index Flood Method. When  $\mu = \bar{Q}$ , the corresponding variate is  $X_{ij} = Q_{ij}/\bar{Q}_i$ , with the properties  $E(X) = 1$ ,  $\sigma_x = Cv(Q)$ , and the third order moment,  $g_x = g_Q$ . The variate  $X$  is a ratio of two random variables rather than a single scaled random variable though the assumptions of such distinctions are usually ignored in practice (Cunnane, 1988). The fundamental assumption of regional homogeneity or RFFA is that the quantile 'X' is assumed to have a common form of distribution with identical parameter values at all sites in the region.

Using the above definitions and notations, Cunnane (1988) suggested that all RFFA methods can be used in two form of quantile flood estimations as presented below:

#### **a) At-site/regional RFFA quantile estimation**

Among the popular methods of RFFA, the Dalrymple method (Dalrymple, 1960) is a common method of At-site/Regional quantile estimation and was also adapted for this thesis. This has been applying in different RFFA studies across the world (Cunnane, 1989). It is a regional averaging index flood method based on equal records length,  $n$ , from unregulated rivers which have been tested for homogeneity (Cunnane, 1988). The results obtained from this method are in a form of standardized variant  $X$  where its  $X-T$  relation is assumed to hold at all sites  $i$  in the region, with  $Q_T = X_T * \mu_i$  (where  $\mu_i$  is the index flood of at-site annual maximum floods).

#### **b) Regional only flood quantile estimation - ungauged catchment**

In case of ungauged catchment, there is no sample available from which the at-site index value can be estimated. Nevertheless, one of the merits of regional analysis is to solve such problems reasonably. That is, once the regional frequency analysis is done, the normalized quantile flood with the index flood approaches may be used to estimate the quantiles for ungauged catchments. The index values ( $\mu$ ) of the ungauged catchments, however, can be estimated using a relation between  $\mu$  and catchment characteristics, obtained using multiple

regression from the available neighboring data sets (Cunnane, 1988; Ellouze and Abida, 2008; Noto and La Loggia, 2009; Rosbjerg, 2007).

### **3.1.2 Procedures of RFFA**

In this thesis, the analysis of regional flood frequencies were carried out using different packages of R-software (which includes the *lmomco*<sup>10</sup>, *bootstrap*<sup>11</sup>, *LmomRFA*<sup>12</sup> and *nsRFA*<sup>13</sup> packages) and the index flood method together with L- moment approaches. A fundamental assumption of the index flood method is that the normalized data at different sites in a region follow the same distribution and should satisfy the conditions of independent and identically distributed (i.i.d) (Cunnane, 1989; Hosking and Wallis, 1997).

For data series available at large number of sites, the quantile estimation is required at each river station. According to Hosking and Wallis ( 1997), the regional flood frequency analysis using index-flood procedures together with L-moments were derived using the following steps (the details will follow in the next sections):

- i. Extrapolator data analysis (data screening and outlier analysis)
- ii. Develop homogeneous regions
- iii. Fitting the regional data to an appropriate frequency distribution
- iv. Estimation the parameter values for the fitted distribution
- v. Estimation of quantile floods of the regions
- vi. Develop regional flood frequency curve
- vii. Develop regional flood frequency analysis for ungauged catchments

## **3.2 Exploratory data analysis**

Exploratory data analysis is a method which employs some statistical tools that provides conceptual and computational tools for discovering patterns to further hypothesis development and refinement. It is an approach for data analysis that utilizes a variety of techniques to maximize insight into a dataset, extract important variables, detect outliers and irregularity of the observations, test underlying assumptions, and develop robust models

---

<sup>10</sup> <http://cran.r-project.org/web/packages/lmomco/index.html>

<sup>11</sup> <http://cran.r-project.org/web/packages/bootstrap/index.html>

<sup>12</sup> <http://cran.r-project.org/web/packages/lmomRFA/index.html>

<sup>13</sup> <http://cran.r-project.org/web/packages/nsRFA/>

(Behrens, 1997). This method of analysis is carried out by visual and graphical data exploration i.e., using different statistical and graphical tools such as plotting time series of parent data, probability plots of the extremes, histograms, and autocorrelation functions and mean excess functions (Embrechts et al., 1997).

In regional flood frequency analysis, there are several factors that may influence the certainty of the analysis. For example, related data may exist at several sites or different population may exist at a single site. Hosking and Wallis (1997) suggested procedures that can be used to make sure that the observation series are representative of the real process, i.e., (i) checking each site's data separately which may contain outliers and repeated values; (ii) checking for trends and independency in the data; and (iii) checking for inter-site dependency.

Moreover, it is important to check whether the sites/data fulfills the requirements of the analysis. In this work, the exploratory data analysis was accomplished based on two steps i.e., (1) data screening and (2) outlier analysis.

### **3.2.1 Data screening**

Data screening is first task of an exploratory data analysis which employed methods that can filtered the unwanted observation from the data series as well as the sites from the analysis (Hosking and Wallis, 1997; Kachroo et al., 2000). For this work, the following methods of data screening were performed;

#### **i) Looking at the data series**

It was performed by visual inspections of some simple statistical methods of data analysis such as the time series plotting, probability plotting, histogram, and autocorrelation plots for both time series data, but more focused for AMS. In the first step, all the data were examined at their time series plots of the sample. The main criteria that were used to select stations were based on length of record period (above 15 years), continuous (no consecutive gap) and common record period.

Therefore, once the above method of data screening was carried out, stations contain the following conditions were excluded from subsequent step of data analysis.

- i) Stations which have short record length ( i.e., <15 years)

- ii) Stations which consist a lot of NO data in the series ( i.e. contains more consecutive gaps)
- iii) If a station contains insignificant magnitude of observed series
- iv) Rivers which reflect not natural phenomenon i.e., if stations have repeated values for long period of time and/or some constant fluctuations. This could be in catchments which are under control at somewhere upstream or in stations that the gauging instruments are not able to measure high magnitude floods.

## **ii) Checking for independent and identically distribution**

By principle, it is known that flood frequency analysis is carried out when the at-site data are independent (without serial correlation and trends) and identically distributed (from the same population), i.e., when the conditions of independently, identical distribution (i.i.d.) are satisfied (Gottschalk and Krasovskaia, 2001; Hosking and Wallis, 1997; Kjeldsen et al., 2002). This provides that the extreme events might appear randomly and all might have the same frequency distribution. However, due to the complicity of the flood environment, it may be expected that the extreme events may not satisfy the conditions of i.i.d and/or stationary (Engeland, 2005).

The presence of temporal dependency implies repetition of information given by previous values i.e., correlated with time. Various studies were carried out to investigate the effect of the presence of dependence in annual maximum series on parameter estimation. For instance, the review by Mkhandi et al.(2000) illustrated that the presence of dependence in data leads to biased quantile estimates and larger standard error than when independence and the correct model form is assumed.

Another requirement of RFFA is that the AMS at different stations in a homogeneous region should be spatially independent. Stations which have significant spatial correlation implies that a lower degree of additional regional information can be obtained by considering both stations in the estimation of regional parameters (Mkhandi and Kachroo, 1997; Mkhandi et al., 2000). That is, the presence of two stations which are significantly correlated may be considered as providing redundant information.

Therefore, the **serial and cross correlations** i.e., the dependence of the observations within a given site and across stations were examined by computing the autocorrelation and spatial correlation coefficient, respectively.

**Autocorrelation coefficient**- is a normalized measure of the linear correlation among successive values in a time series. The use of the autocorrelation function in characterizing the behavior of a time series lies in its ability to determine the degree of dependence present in the values. For a random process, a descriptor of the random structure of the process needs to be added and the autocovariance function (acf) determines this structure as an acceptable approximation (Gottschalk, 2005). The covariance  $\beta(t, t')$  of the state of a random process between two different points in time  $X(t)$  and  $X(t')$  defines this autocovariance function of ( $t$  and  $t'$ ):

$$\beta(t, t') = \beta_x(t, t') = Cov[X(t), X(t')] = E[X(t).X(t')] - m(t)m(t') \quad (3.1)$$

Similarly, the autocorrelation function ( $\rho(t, t')$ ) is defined as autocovariance divided by consecutive standard deviations ( $\sigma$ ) of time ( $t$  and  $t'$ ) by:

$$\rho(t, t') = \rho_x(t, t') = \frac{\beta(t, t')}{\sigma(t)\sigma(t')} \quad (3.2)$$

which is the correlation coefficient between  $X(t)$  and  $X(t')$ .

For sample size of  $n$  observations, the sample autocorrelation were estimated by calculating the sample covariance  $\hat{\beta}(k)$  first and then correlation coefficients  $r(k)$  as follows:

$$\hat{\beta}(k) = \hat{\beta}(k \Delta t) = \frac{1}{n - k} \sum_{j=1}^{n-k} x(t_j)x(t_{j+k}) - \bar{x}^2, \quad k = 0 \dots K \quad (3.3)$$

$$r(k) = r(k \Delta t) = \frac{\hat{\beta}(k)}{s_x^2}, \quad k = 0 \dots K \quad (3.4)$$

where,  $\bar{x}^2$  and  $s_x^2$  are the square mean and variance of the sample series, respectively and  $k$  is the time lag in terms of the interval  $\Delta t$  between observations in time up to  $K$ , which is the maximum lag. The correlation coefficients between two consecutive observations of the sample series were plotted and the degree of dependence was rejected at 5% significant level.



**Spatial Correlation coefficient-** calculates the dependency of the AMS between nearby stations. In applied situation (Gottschalk, 2005), the first and second-order sample moments were determined from the observations  $x(u_i, t_k)$  in  $M$  stations at points  $u_i$ ,  $i = 1, 2 \dots M$  stations at  $k$  points of time,  $t_k$ ,  $k = 1, 2, \dots, n$ . As a first step, the sample means  $\bar{x}$  were calculated for each of  $M$  stations as;

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x(u_i, t_k), \quad i = 1, 2 \dots M \quad (3.5)$$

The variance,  $\hat{\beta}_{ii}$  of the observations which measures the deviation of individual observations from the expected value can be obtained as;

$$\hat{\beta}_{ii} = s_x^2 = \frac{1}{n} \sum_{k=1}^n (x(u_i, t_k))^2 - \bar{x}^2, \quad i = 1, 2 \dots M \quad (3.6)$$

The pair wise covariance and correlation is also estimated using;

$$\hat{\beta}_{ij} = (s_i * s_j) = \frac{1}{n} \sum_{k=1}^n x(u_i, t_k)x(u_j, t_k) - \bar{x}_i \bar{x}_j, \quad i, j = 1, 2 \dots M \quad (3.7)$$

where  $s_i$  and  $s_j$  are the standard deviations of random variables of  $X_i$  and  $X_j$ , respectively. The pair wise correlation coefficients were calculated as;

$$r_{ij} = \frac{\hat{\beta}_{ij}}{s_i * s_j}, \quad i, j = 1, \dots, M \quad (3.8)$$

Thus, one of the two stations which reflect strong pair correlation coefficient was excluded from the analysis.

### 3.2.2 Empirical distribution

It is a statistical method used to ensure that all the observations are valid representations of the hydrological characteristics under considerations (Haan, 2002). Haan (2002) briefed that after the data have been accepted as valid, basic statistical moments of the data should be computed and should be plotted as a probability plots. Determining the empirical distribution of a given extreme event is referred to determine the probability of the assign data- it could be probability density function,  $f(x)$  or cumulative probability function,  $F(x)$ .

Let  $X$  is a random variable, taking values that are real number. The relative frequency with which these values occurred defines the frequency distribution or probability distribution of  $X$  and is specified by the cumulative distribution functions;

$$F(x) = p_r(X \leq x) \quad (3.9)$$

where,  $p_r(A)$  denotes for the probability of the event  $A$ .  $F(x)$  is an increasing function of  $x$ , and  $0 \leq F(x) \leq 1$  (Haan, 2002).

Probability plotting of hydrological data requires that individual observations or data points should be independent of each other and representative of the same population (Haan, 2002). Generally, a sample will not contain the smallest or largest value of the unknown population (Cunnane, 1988; Gottschalk and Krasovskaia, 2001; Haan, 2002; Hosking and Wallis, 1997). Thus, plotting positions of '0 and 1' should be avoided from the sample series unless one has additional information on the population limits.

Gottschalk and Krasovskaia (2001) also discussed that for a given sample size  $n$ , ranked in ascending order, if the distribution  $F(x)$  is assumed to be a uniform distribution giving a probability  $p_r$  related to the  $r^{th}$  ordered value of  $x_{(r)}$ :  $p_r = F(x_{(r)})$  for the interval  $(0,1)$ , then the frequency distribution is:

$$f_r(p) = \frac{n!}{(r-1)!(n-r)!} p^{r-1} [1-p]^{n-r} \quad (3.10)$$

The mean and variance of this distribution are therefore obtained as:

$$E[p_r] = \frac{r}{n+1}; \quad Var[p_r] = \frac{r(n-r+1)}{(n+1)^2(n+2)} \quad (3.11)$$

where,  $E[p_r]$  is an expression often used for plotting the empirical distribution. In hydrology it is usually called the Weibull's plotting position (Weibull, 1939). The popular and alternative methods of plotting positions are given in Table 3.1

Table 3.1 Some of the popular probability plotting models: where 'r' is rank of the observed value to be plotted, n is number of observation.

S.N	Plotting- Position Models	Non-exceedance probably
1	Weibull (Weibull, 1939)	$p_r = \frac{r}{n+1}$
2	Gringorten (Gringorten, 1963)	$p_r = \frac{r - 0.44}{n + 0.12}$

Using the above methods, the graphical presentation of the relationship between observed values  $x$  and cumulative distribution function  $F(x)$  in arithmetic scale is not usually simple when extreme values are of interest (Embrechts et al., 1997; Gottschalk and Krasovskaia, 2001). However, hydrologists have been using a modified graph which is known by Reduced Gumbel Variate ( $y$ ) for presenting the relationship between the probability and observed values. The Probability plots are designed for particular theoretical distributions (i.e., extreme value type I (EV1) (Gumbel, 1958)) by transforming the scale of the probability axis so that a given distribution is represented by a straight line. This provides the reduced form of  $F(x)$  from  $n$  observations, and it is linearly related to the observed values,  $x$  (Gottschalk and Krasovskaia, 2001).

Therefore, for this analysis, the reduced Gumbel probability plotting was adapted. Gottschalk and Krasovskaia (2001) and Cunnane (1989) suggested that the method is recommended probability plotting for extreme value analysis and can be approximately expressed by Gringorten plotting position. It also has a theoretical background for graphical representation of empirical distributions.

Suppose  $x_{i1}, x_{i2} \dots x_{in}$  are extreme floods of  $n$  observations from  $i$ -site, ranked in ascending order. The probability  $p_r$  related to the  $r^{th}$  ordered value  $x_{(r)}$  is the cumulative distribution function of the theoretical distribution (i.e.,  $p_r(x_{(r)}) = F(x_{(r)})$  for the interval  $(0,1)$ ). The Gumbel (EV1) distribution of these extreme events is;

$$F(x) = \exp \left[ -\exp \left\{ -\frac{(x_{(r)} - u)}{\alpha} \right\} \right] \quad (3.12)$$

where,  $u$  and  $\alpha$  are the location and scale parameters of the Gumbel distribution, respectively. The cumulative distribution function  $F$  is reduced into the Gumbel variate,  $y$  which related through;

$$y = -\ln \left[ -\ln \left( F(x_{(r)}) \right) \right] \quad (3.13)$$

when  $F(x_{(r)})$  is approximated by Gringorten probability plotting, the reduced Gumbel plotting position  $y$  for the  $r^{th}$  observation was estimated using:

$$y_{(r)} = -\ln \left[ -\ln \left( \frac{r - 0.44}{n + 0.12} \right) \right] \quad (3.14)$$

As a result, the probably plot curve for every station was presented in the form of,  $x_{(r)}$ :

$$x_{(r)} = u + \alpha y_{(r)} \quad (3.15)$$

where  $x_{(r)}$  is the normalized  $r^{th}$  events (y-axis) and  $\alpha$  is the linear correlation of the observed series, and  $y_{(r)}$  is the reduced Gumbel variate. The curves were closely examined whether the AMS in every station represents the random process and comes from single population of the sample series.

### 3.2.3 Outlier detection and treatments

Outliers are observation values which may not be representative of the sample i.e., they might have apparently different frequency distribution. When the empirical distribution of the extreme values of the observed data is plotted, this can be located strongly deviate from the rest of the dataset (Gottschalk and Kundzewicz, 1995; Haan, 2002; Kottegoda, 1984). The existence of outliers means that there is an existences of extremes of extreme events, which can be the reason for many of the problems raised in the regional analysis of hydrological data (Gottschalk and Kundzewicz, 1995). For example, it is entirely possible that a 100 year event is contained in 20 years record. If this is the case, assigning a normal plotting position

(i.e.,  $p_r = 1/21$ ) to this value would not be representative of its true return period. Thus, this value should be detected and treated as an outlier.

Although no observations are fully trustworthy, it will be considered that the data under some circumstances might be reasonably representative of the sample (Kottegoda, 1984). This means that, the uncertainties introduced during measurement and recording are supposed to be eliminated by establishing prior treatments to the data. Kottegoda (1984) discussed that no one can define exactly “what an outlier is?”. This is because the detection and treatments or decisions to be considered as unacceptable observation depends on the person’s choice or judgment. This choice may depend on the experience, personal judgment, type of the data and the robustness of the model for the analysis and interpretation of output should account for such pitfalls.

There are some principal statistical test reviewed and recommended in the literatures by Kottegoda (1984), Gottschalk and Kundzewicz (1995) and Haan (2002) which are useful to identify and eventually eliminate outliers. Among these, the most relevant methods that are easily applicable and can be used directly to identify whether the largest observations of the sample are outliers, were employed in this analysis. The outliers from all series were detected using the following methods, simultaneously:

- i. Observed suspected outliers by visual inspection of the AM series
- ii. Identified the symmetry of the distribution of the extreme floods from the skewness coefficients
- iii. Detected the outliers in both tails of the AM floods (AMF) frequency distributions using Bulletin 17B method (USWRC, 1981) and
- iv. Tested the significant of the suspected outliers of the AMS using student test statistics

**Looking at the Extreme events** – was performed based on prior knowledge of the series at each gauging sites. It is visual inspection of observed values under some statistical methods such as reduced Gumbel variate – plotting position, histogram and plotting of the full time series of annual maximum flood (AMF).

As defined earlier, an outlier is an observation which is positioning apart from the rest of the sample groups. Hence, these visual inspection methods can give as the impression of the existence of outliers in the observation series.

**Skewness coefficient ( $C_s$ )** - for a random variable  $X$ , it is the third moment which measures the asymmetry of the probability distribution of the observed values. The skewness coefficient is defined as;

$$C_s = E \left[ \left( \frac{X - \mu}{\sigma} \right)^3 \right] = \frac{E[(X - \mu)^3]}{[E(X - \mu)^2]^{3/2}} \quad (3.16)$$

where,  $\mu$  and  $\sigma$  are the mean and standard deviation of the random variable  $X$ , respectively. Therefore, the skewness coefficient ( $g_x$ ) for a sample of  $n$  observations is given;

$$g_x = \frac{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^3}{\left[ \frac{1}{n} \sum_{j=1}^n ((x_j - \bar{x})^2) \right]^{3/2}} \quad (3.17)$$

The skewness values can be positive, negative or nearly zero. The value indicates the asymmetry of the frequency distribution. For example, if  $C_s$  is approximately zero, then it indicates symmetrical distribution of the data. However, in this case (i.e., outlier analysis for maxima events), a distribution with positive and large skewness coefficient indicates heavy tail to the right. This gives the impression that the series of data contains one or more outlying observations in the sample (Gottschalk and Kundzewicz, 1995). Gottschalk and Kundzewicz (1995) also illustrated that “when the observed data contains a single outlier in the series, the skewness coefficient can be changed by a factor of two”.

Therefore, the sample coefficient of skewness was used as indicator whether the observed series has one or more outliers. However, this doesn't mean that the coefficient test whether a specific observation is an outlier rather merely indicates the presence of outliers.

**Bulletin 17B method** - was introduced by the United States Water Resources Council (1981), and suggested to be useful as a tool for outlier detection in the book by Haan (2002). It is a method that can detect outliers from both sides of the frequency distribution curve of the extreme events- upper and lower tails. This method is preferable to other methods because

outliers at both sides (minimum and maximum) of the series can be investigated based on the sample characteristics (i.e., the mean, standard deviation and sample size)

Suppose, a sample  $x_{ij}, j= 1, 2 \dots n_i$  observations of random variable  $X_i$  are ranked in ascending order, the threshold levels for high and low outliers ( i.e.,  $x_H$  and  $x_L$  , respectively) of the sample series were defined as;

$$x_H = \bar{x} + K_n s_x \quad (3.18)$$

$$x_L = \bar{x} - K_n s_x \quad (3.19)$$

where, the  $\bar{x}$  and  $s_x$  are the mean and standard deviation of the observed values, respectively; and  $K_n$  is approximated from the logarithmic value of the number of the observations,  $n$  as follows;

$$K_n \approx 1.055 + 0.981 \log_{10} n \quad (3.20)$$

Because the data used for this analysis were extreme maxima with expected positive skewness coefficient; the threshold level were needed only for higher outliers and estimated through equation. 3.18. The observations which were larger than the threshold level of the at-site sample series were detect as outliers.

**Student deviation test (t-test)** - suppose the observation,  $x_{ij}, j= 1, 2 \dots n$  belong to a normal distribution with the same mean. The null hypothesis is then rejected if the largest observation at 95% confidence interval significantly deviates from the expected value of the samples (Gottschalk and Kundzewicz, 1995; Kottegoda, 1984).

$$t = \frac{(x_{(n)} - \bar{x})}{s_x} \quad (3.21)$$

where,  $x_{(n)}$  is the largest value among the total of  $n$  observations and  $\bar{x}$  and  $s_x$  are the mean and standard deviation of the observed sample, respectively.

**Outlier treatments** - the outliers identified by the above procedures were treated in agreement with the recommendations of Cunnane (1989) as follows:

- Because outliers are rare and do not appear in every sample, a test of any hypothesis about their frequency of occurrence must be done on a regional base.
- Outliers can be excluded from the observations only if it is certain that the AM floods can be adequately modeled by a single known distribution form.
- If the AMS are regarded as true observations, but come from two very different sub-populations, then the outliers must be retained.
- If the retained outliers are few in number and an efficient method of parameter estimation like PWM or ML is used, their influence is very small and should be used only in regional estimation procedures.

Therefore, the regional distributions of the stations which contained one or more significant outliers were identified first and consequently, for every outlying observation: the region, station, normalized value, the expected exceedance probability and return periods were estimated regionally.

### **3.2.4 Index flood method**

This is a procedure that assumes the data at different sites in a region follows the same distribution except for scaling factor (i.e., index flood). The procedures are a convenient way of pooling summary statistics from different data samples (Hosking and Wallis, 1997). The index flood might be the mean ( $\bar{Q}$ ) or the median ( $\tilde{Q}$ ) of annual maximum floods (Chebana and Ouarda, 2009; Farquharson et al., 1992; Gottschalk and Krasovskaia, 2001; Hosking and Wallis, 1997; NERC, 1975). For gauged sites, the  $\bar{Q}$  or  $\tilde{Q}$  values can be estimated in a straight forward manner from the observations. However, if the record is too short (< 15 years) and for ungauged sites, estimation of the index flood could be difficult. Thus, an adjustment for climatic and catchment variations might be necessary (section 3.6).

It is obvious that an appropriate choice of representative scale factor may increase the performances of the estimation parameter values of the candidate distribution, especially, for the sites which have outliers in the sample of the floods. As mentioned earlier, for the catchments which have measured flood values, the index flood should be the mean or the median of the extreme events. However, there is no clear guideline for the choice of the value



of index flood, i.e., under which circumstances does the mean be the index flood, otherwise the median of the sample.

The at-site mean annual maximum flood is often used as the index flood for RFFA (Farquharson et al., 1992; Kjeldsen et al., 2002; Noto and La Loggia, 2009; Rosbjerg, 2007; Stedinger and Lu, 1995b; Yang et al., 2010). Nevertheless, Viglione et al. (2007) discussed and suggested the advantages of using the sample median as the index value. When the parent distributions are skewed and used for flood frequency analysis, the median is preferable to mean. This is also discussed in the literature by Noto and La Loggia (2009) that the sample median can be used instead of mean in case of outlying observations in the sample.

In this work, almost all observations in the sample reflect positive skewness coefficients. Hence, to make sure that the index floods was used in a better way, the index flood was chosen by performing sensitivity analysis of the sample mean and median to large observations.

Suppose, a homogeneous region with  $M$  sites, each site  $i$  having a sample size  $n_i$ , and an observed annual flood series  $Q_{ij}$ ,  $j = 1, 2, \dots, n_i$ . The index flood  $\mu_i$  is the mean ( $\bar{Q}$ ) or the median ( $\tilde{Q}$ ) of the AMS. The better index flood was chosen by comparing the relative differences of the mean and median after the largest observation were removed from the series. The relative difference was calculated as:

$$\text{Relative difference of Mean} = \frac{(\bar{Q}_1 - \bar{Q}_2)}{\bar{Q}_1} * 100, \text{ and,} \quad 3.22$$

$$\text{Relative difference of Median} = \frac{(\tilde{Q}_1 - \tilde{Q}_2)}{\tilde{Q}_1} * 100 \quad 3.23$$

where, the numbers 1 and 2 indicate for the parameters values with and without the largest observation in the series, respectively. The less sensitive index flood ( $\mu_i$ ) was therefore used to normalize the sample series (i.e., the dimensionless variate  $X_{ij}$  which is assumed to have the same form of distribution at every site  $i$  in the region) as:

$$X_{ij} = Q_{ij}/\mu_i \quad 3.24$$

where  $X_{ij}$  is the normalized value of the  $j^{\text{th}}$ -observation in  $i$ -site and the at-site  $i$  quantile flood ( $Q_{T(i)}$ ) will be estimated by up scaling the regional quantile ( $X_T$ ) for the return period ( $T$ ) as;

$$Q_{T(i)} = \mu_i * X_T \quad 3.25$$

### 3.3 Regionalization

Regionalization, in the context of flood frequency analysis, refers to grouping of basins into homogeneous regions and selection of appropriate frequency distributions for the identified regions. A more specific definition of a homogeneous region is that the region consists of sites having the same standardized frequency distributional form and parameters (Burn, 1988; Chebana and Ouarda, 2008; Cunnane, 1988; Gottschalk and Krasovskaia, 2001; Hosking and Wallis, 1993; Hosking and Wallis, 1997; Kachroo et al., 2000; Tveito, 1993; Wiltshire, 1985, 1986). These studies summarized that to accomplish grouping of stations or basins in to hydrological homogeneous regions, two basic steps should be conducted. The first step is delineating the regions using different methods with catchment, environment and climate information and the second step is applying heterogeneity tests (i.e., evaluating if the regions contain statistically similar sites or not).

#### 3.3.1 Delineation of homogeneous regions

Due to the complexity in understanding the factors that have direct and indirect effect on the generation of flood, there are no simple guidelines for identifying homogeneous regions (Kachroo et al., 2000). Meanwhile, experience, prior information and personal judgments can provide possible guidelines to delineate regions with similar hydrological features.

There were several attempts made by different authors to identify hydrologically homogeneous regions and their emphasis were either on geographical considerations or on hydrological characteristics or a combination of both (Kachroo et al., 2000). For example, Hosking and Wallis (1997) discussed some of the methods such as geographical convenience-based on the administrative areas; subjective partitioning - defines region subjectively by

inspection of the site characteristics; objective partitioning – regions formed by assigning sites to one of two groups depending on some threshold values; and clustering method, which is the standard method of statistical multivariate analysis for dividing a data set into regions. It is obvious that catchments might not necessarily have exactly the same behavior given the limited sample size, and the dynamic and infinite factors influencing flood generation. Hence, it is not necessarily to group sites to exactly satisfy the homogeneity tests (Hosking and Wallis, 1997). That is, an approximate homogeneity could be sufficient to ensure that the regional flood frequency analysis is preferred to at-site analysis.

For this analysis, the geographical regionalization approach was carried out. It is more convenient, because it may divide a region into different regions based on variation of soil, climate and topography with latitude and longitude (Cunnane, 1989). It is also advantageous if ungauged/poorly gauged catchments are assumed to be assigned in the identified region. However, geographical proximity of two stations is not guaranteed that they have similar form of flood frequency distributions (Cunnane, 1989; Hosking and Wallis, 1997).

All the above backgrounds were considered and the regionalization was accomplished using the following procedures together with the available previous grouping information such as Mkhanda and Kachroo (1997) and Karchroo et al.(2000).

- a) Geographic information such as drainage characteristics and geographically continuous catchments were used to identify likely homogenous regions.
- b) Each region identified in procedure (a) where checked for heterogeneity by its statistical data behavior.
- c) The final numbers of regions were found by testing the homogeneity measure at each grouping procedures- which started with one group, if not homogeneous then continue to separate to two or three groups etc.

### **3.3.2 Homogeneity Test**

The homogeneity test was based on the heterogeneity measured ( $H$ ) suggested by Hosking and Wallis (1997). The assessment of the regional heterogeneity is obtained by comparing the L-moments (i.e., particularly the at-site  $L-Cv$ ) of observed samples in the region. This is

performed after the normalized regional series are fitted to kappa distribution and later compared with 500 times Monte Carlo simulated values (Hosking and Wallis, 1997) using:

$$H = \frac{(V - \mu_v)}{\sigma_v} \quad (3.26)$$

where,  $V$  is the weighted standard deviation of the at-site observed  $L-Cv$  values,  $\mu_v$  and  $\sigma_v$  are the mean and standard deviation of simulated  $L-Cv$  values, respectively. Suppose that a region has  $M$  sites, for the sample and simulated regions,  $V$  is then calculated:

$$V = \left\{ \frac{\sum_{i=1}^M n_i (t^i - t^R)^2}{\sum_{i=1}^M n_i} \right\}^{1/2} \quad (3.27)$$

where,  $n_i$  is record length at site  $i$ ,  $t^i$  is the sample  $L-Cv$  at site  $i$  and  $t^R$  is the regional average sample  $L-Cv$ .

The regions were declared in agreement with criteria established by Hosking and Wallis (1993 and 1997). That is, a regions was regard as acceptably homogeneous if  $H < 1$ ; possibly heterogeneous if  $1 < H < 2$ ; and definitely heterogeneous if  $H > 2$ .

### 3.4 Choice of regional flood frequency distribution

#### 3.4.1 Theoretical distribution functions

In flood event analysis, the annual maximum flow  $Q_T$  corresponding to a given recurrence intervals  $T$ , can be estimated from the annual flood series using varies theoretical distributions. If  $T$  is large compared with the record length of the series  $n$  and the chosen model is inappropriate, the error of the  $T$ -year estimated flood can be very large and consequently, the associated design losses could be considerable. Thus, an acceptable design procedures is essentially required to choose a model that minimize such uncertainties (Rossi et al., 1984).

It has been suggested that, the theoretical distribution functions used in hydrology for RFFA are as a rule borrowed directly from the Probability Theory. For example, the two-parameter Gumbel (EV1) and Exponential (EXP); and three-parameter Generalized Logistic (GLO),

Generalized Extreme Value (GEV), Generalized Pareto (GPA), Pearson type III (P3) and Lognormal (LN3) distributions etc are among the commonly employed distributions in recent RFFA (Gottschalk and Krasovskaia, 2001).

The theoretical distributions that are of a special importance related to the family of extreme value distributions are the EV1, GEV and GPA distributions. It is well known that the GEV distribution can be Gumbel distribution for  $k \approx 0$ , the Frechet distribution (EV2,  $k < 0$ ) which is unbounded towards the extreme maxima, and the Weibull distribution (EV3,  $k > 0$ ), bounded towards the extreme maxima (for example, see Fig. 3.1) (Gottschalk and Krasovskaia, 2001). In recent hydrology, it has been reported that the GPA and GEV distributions are commonly applied as better distributions to predicate extreme events. They are also interrelated. The GPA appears as a limiting form for extreme values over a given threshold, Peak over Threshold (POT). Whereas, the GEV distribution is the limiting form for extremes events selected as the largest value over a certain time interval; say a year, Annual Maximum Series (AMS). The relationship between the GEV and GPA distributions are discussed in the literature (Engeland, 2005; Gottschalk and Krasovskaia, 2001). If the number of events observed over a time interval follows the Poisson distribution and the intensity of POTs per time interval (year) is  $\lambda$ , the general expression for the distribution of annual maxima  $Z$  is:

$$F_Z(z) = \exp\{-\lambda(1 - F_X(z))\} \quad (3.28)$$

where,  $F_X(x)$  is the distribution for the POT values. If they are distributed in accordance with GPA, their distribution is;

$$F(X) = \exp \left[ -\lambda \left\{ 1 - k \frac{(x - u)}{\alpha} \right\}^{\frac{1}{k}} \right] \quad (3.29)$$

This can be identified as a reduced version of the GEV distribution. It can also be written in the form of an ordinary GEV with changed location,  $u' = u + \alpha/k (1 - \lambda^{-k})$  and scale  $\alpha' = \alpha \lambda^{-k}$  parameters (for  $k=0$  the corresponding expressions are:  $u' = u + \alpha \ln(\lambda)$  and  $\alpha' = \alpha$ ).

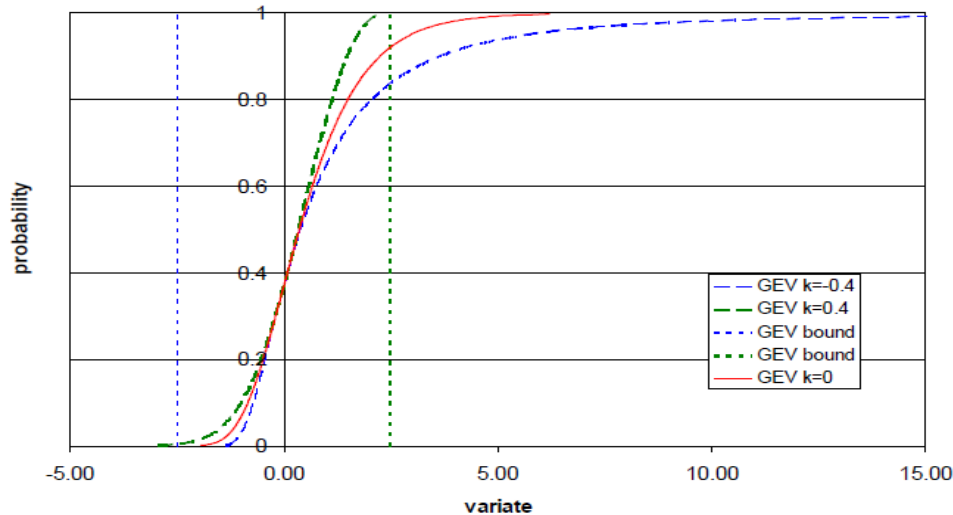


Figure 3.1 Generalized Extreme Value (GEV) distributions: Gumbel (EV1,  $k=0$ ), Frechet (EV2,  $k=-0.04$ ) and Weibull (EV3,  $k=0.04$ ) distribution (Gottschalk and Krasovskaia, 2001).

In RFFA, a single frequency distribution (one of the above distributions) is fitted to data of several sites in a homogenous region. In most cases, regions can be slightly heterogeneous. This implies that there is no single distribution that may fit exactly for all sites in the region. Therefore, the main aim of RFFA is not to find the exact frequency distribution of the region rather it is to find a distribution that will yield as accurate as possible quantile estimates for each site (Hosking and Wallis, 1997).

The problem of the choice of an appropriate theoretical distribution would not have arisen if a certain distribution function was chosen based on the hydro climatic premises. Unfortunately, our prior knowledge of hydrological processes could not help as much as possible. The extreme value theory can be instead useful but only in a very broad meaning. It does not give an answer for the question “what type of an extreme value distribution is going to be chosen?” This problem in fact might be solved based on the generalized extreme value (GEV) distribution if the shape parameter, ‘ $k$ ’ of this distribution can be estimated accurately (Gottschalk and Krasovskaia, 2001).

Including the recently popular methods of GEV and GPA, many flood frequency distributions have been practiced for flood modeling, but none has been accepted as universal (Mkhandi and Kachroo, 1997). Hence, seven distributions were considered for the evaluation of the possible distributions that can represent the average frequency distribution of the regional data

in Southern Africa. The distributions and their parameters are presented in Table 3.2 (the details are presented in Appendices E (Table 12) and Hosking and Wallis (1997) pages 191-209).

Table 3.2 Summary of the employed distributions and their parameter

Two-parameter distributions:	Three-parameter distribution:
Location ( $\mu$ ), scale ( $\alpha$ )	Location ( $\mu/\zeta^*$ ), scale ( $\alpha/\beta^*$ ) and shape ( $k/\alpha^*$ )
<ul style="list-style-type: none"> <li>• Gumbel (EV1)</li> <li>• Exponential (EXP)</li> </ul>	<ul style="list-style-type: none"> <li>• Lognormal (LN3)</li> <li>• Person type III (PE3)</li> <li>• Generalized Pareto (GPA)</li> <li>• Generalized Logistic (GLO)</li> <li>• Generalized extreme Value (GEV)</li> </ul>

*Note: the parameters with asterisk are for the Pearson type III (PE3) distribution and in exponential distribution the  $\mu$  parameter denotes for the lower end point of the distribution.*

The considerations of these extreme event models were also based on previous studies in Southern Africa. For example, the RFFA for South Africa (Kjeldsen et al., 2002); and the RFFA for Southern Africa (Mkhandi and Kachroo, 1997) considered all the above distributions. However, using the same methods of RFFA and AMS dataset, the suggested flood models were different. Mkhandi and Kachroo (1997) suggested that the PE3/PWM or LP3/LM distributions as an appropriate distribution for all catchments in Southern Africa while the Kjeldsen et al. (2002) recommended that the GPA distribution as suitable frequency distribution for Kwazulu-natal province catchments in South Africa.

### 3.4.2 Fitting the regional data to empirical distribution

After the regions were identified and confirmed as homogenous, the next step was to choose a stochastic model which can represent the regional flood characteristics. The fitting of the theoretical distributions to the regional observations was carried out with the following procedures:

- Normalization of the observed series with respect to their index flood

- Calculation of the weighted at-site L-moments and then weighted average regional L-moments using PWMs
- Development of the L-moment ratio diagram and plotting the at-site L-moments together with theoretical distribution
- choosing candidate distributions from the diagrams
- Finally statistical confirmation test of the best fitted distribution

The above procedures were compiled into three principal steps: firstly by establishing the L-moment ratio diagrams (L-skewness vs. L-Kurtosis), secondly by conducting goodness-of-fit test, and after developing regional growth curve ( i.e., the regional average weighted curve of the at-site samples) lastly by model performance evaluation.

### **3.4.3 L-moment and L- moment ratio diagram**

L-moments ( $\lambda_n$ ,  $n=1, 2, \dots$ ) are related to the expected order statistics and have come to replace the use of ordinary moments in hydrological analysis. The main purpose of estimating L-moments and probability weighted moments (PWM) is similar to ordinary product moments, but increase the certainty of the models because the natural estimator of  $\lambda_n$  is based on a linear combination of the ordered of the observed data values (Kjeldsen et al., 2002). The L-moments approach covers the characterization of probability distributions, the summary of observed data samples, the fitting of probability distributions to data, and testing the hypothesis about the distributional form (Hosking and Wallis, 1997).

According to Hosking and Wallis (1997), ordinary moments are not always satisfactory because of two major reasons: (1) they do not always reveal easily interpreted information about the shape of a distribution, and (2) parameter estimates of distributions fitted by the moments are often less accurate than those obtained by other methods, such as the PWMs method of parameter estimator. L-moments have the theoretical advantages over conventional moments for being able to characterize a wider range of distributions and, when estimated from a sample, they are more robust to the presence of outliers in the data.

The “L” in L-moments gives attention to the linearity in forming the moments by linear combinations of the probability-weighted moments as given by:



$$\beta_r = E\{X[F_X(x)]^r\} \quad (3.30)$$

where  $\beta_r$  is the  $r^{\text{th}}$  order of PWM and  $F_X(x)$  is the cumulative distribution function of a stochastic variable,  $X$ . Let a site  $i$  has sample size of  $n$  observations, arranged in ascending order (i.e.,  $x_{1:n} \leq x_{2:n} \leq \dots \leq x_{n:n}$ ). The sample L-moments were estimated from the unbiased sample estimator ( $\beta_r$ ) of the probability weighted moments as (Hosking and Wallis, 1997):

$$\beta_r = n^{-1} \sum_{j=r+1}^n \left[ \frac{(j-1)(j-2) \dots (j-r)}{(n-1)(n-2) \dots (n-r)} \right] x_{j:n} \quad (3.31)$$

The unbiased sample estimators of the first four PWMs were calculated through equation (3.31) and then used to calculate the first four sample L-moments as follows:

$$\lambda_1 = \beta_0 \quad (3.32)$$

$$\lambda_2 = 2\beta_1 - \beta_0 \quad (3.33)$$

$$\lambda_3 = 6\beta_2 - 6\beta_1 + \beta_0 \quad (3.34)$$

$$\lambda_4 = 20\beta_3 - 30\beta_2 + 12\beta_1 - \beta_0 \quad (3.35)$$

It is known that the first order L-moment  $\lambda_1$  is the expected value of the normalized values of the AMS. Finally, the L-moment ratios,  $\tau$  were calculated from L-moment values as:

$$\tau_2 = \frac{\lambda_2}{\lambda_1}, \quad \tau_3 = \frac{\lambda_3}{\lambda_2}, \quad \tau_4 = \frac{\lambda_4}{\lambda_2} \quad (3.36)$$

where,  $\tau_2, \tau_3$ , and  $\tau_4$  are L-coefficient of variance ( $L-Cv$ ), L-skewness ( $L-Cs$ ) and L-kurtosis ( $L-Kurt$ ), respectively.

The weighted averages of the regional L- moments for the delineated regions were estimated from the above relationships. Suppose, there are  $M$  sites in a region with sample size  $n_1, n_2, \dots, n_M$  and the sample L-moment ratios at site  $i$  are denoted by  $\tau^i, \tau_3^i, \tau_4^i$  etc. The regional weighted average L-moment ratios are then given as (Gottschalk and Krasovskaia, 2001; Yang et al., 2010):

$$\bar{\tau}_r = \frac{\sum_{i=1}^M n_i \tau_r^i}{\sum_{i=1}^M n_i} \quad r = 2, 3, 4 \dots \quad (3.37)$$

Using the regional and theoretical L-moment relationships (see Appendices E: Table 13), the L-moment ratio diagram (i.e., *L-Cs* vs. *L-kurt*) was prepared. For every region, the weighted at-site L-skewness and L-kurtosis coefficients was plotted on the diagram along with the theoretical curves of the candidate distributions listed in Table 3.2. The choice of best fitted distribution in the L-moment analysis was then performed by comparing L-moment ratios samples with the theoretical values. This was carried out in agreement with the suggestions, for example, (Ben-Zvi, 2010; Ellouze and Abida, 2008; Kjeldsen et al., 2002; Rosbjerg, 2007; Yang et al., 2010) that if the points corresponding to the regional averages are located near the curve corresponding to a given distribution, the nearby distribution was practically a choice for the parent distribution of the region.

#### 3.4.4 Goodness-of-fit (GOF) test

For a given region that contains sites with similar statistical distribution and parameter values, the main aim of this test is to examine whether the candidate distribution fits to a data set better than the others. There are several methods available for testing the goodness-of-fit of theoretical distribution for extreme events both at-site and/or regional average data. For example, the graphical (like histogram and probability plotting) and statistical tests such as chi square test ( $\chi^2$ - test), Kolmogorov-Smirnov statistical methods discussed in the literature by Gottschalk and Krasovskaia (2001) and Hosking and Wallis (1997); the  $Z^{DIST}$ -test suggested by Hosking and Wallis (1997) and the Anderson-Darling goodness-of-fit test which were used by Viglione et al. (2007) and Viglione (2010).

In recent RFFA, the following two methods are popularly used and also discussed here as potential methods for regional analysis.

**The Anderson-Darling goodness-of-fit test** – can be used in RFFA studies to assess the fitness of the candidate regional frequency distributions. This method is based on statistical frequency distribution behavior of the observed and Monte-Carlo simulation values. The model is used in Viglione et al. (2007) and also well documented in the R-software, ‘nsRFA’

package by Viglione (2010). The procedure computes the statistical summary of the observed data and tests the degree of the data fitting to the expected distribution function.

For a given sample  $x_{ij}$  ( $j = 1, 2 \dots n$ , from station  $i=1, 2 \dots M$ ), extracted from a real distribution  $F_R(x)$ , the test was used to check the null hypothesis,  $H_0: F_R(x) = F(x, \theta)$ , where  $F(x, \theta)$  is the hypothetical distribution and  $\theta$  is an array of parameters estimated from the average value of sample  $x_j$ . Thus, the test measures the departure between the hypothetical distribution  $F(x, \theta)$  and the cumulative frequency function  $F_n(x)$  defined as follows:

$$\begin{aligned} F_n(x) &= 0, & x < x_{(1)} \\ F_n(x) &= j/n, & x_{(j)} \leq x < x_{(j+1)} \\ F_n(x) &= 1, & x_{(n)} \leq x \end{aligned} \quad (3.38)$$

where,  $x_{(j)}$  is the  $j^{\text{th}}$  element of the ordered sample (in increasing order). The test statistic is:

$$Q^2 = n \int_x [F_n(x) - F(x, \theta)]^2 \psi(x) dF(x) \quad (3.39)$$

where  $\psi(x)$ , in the case of the Anderson-Darling test is:

$$\psi(x) = [F(x, \theta)(1 - F(x, \theta))]^{-1} \quad (3.40)$$

In practice, the statistic is calculated as:

$$A^2 = -n - \frac{1}{n} \sum_{j=1}^n \{(2j-1) \ln[F(x(j), \theta)] + (2n+1-2j) \ln[1 - F(x(j), \theta)]\} \quad (3.41)$$

The obtained statistic  $A^2$  might be expected to be confronted with the population of the  $A^2$ 's that one obtain if samples effectively belongs to the  $F(x, \theta)$  candidate distribution model (Viglione, 2010).

The null hypothesis was rejected when the probability of the Anderson-Darling statistics  $A^2$  was greater than a probability at level of significance,  $\alpha = 10\%$  (for example, if the  $p_r(A^2)$  is greater than 0.9) (Viglione, 2010).

**$Z^{DIST}$  - goodness-of-fit test** - a method suggested by Hosking and Wallis (1997). It verifies the selected distribution from the L-moments ratio diagram by comparing the observed regional L-kurtosis to the theoretical values of various candidate distributions. This has been recently a popular method in flood frequency analysis. For example, the RFFA studies in South Africa (Kjeldsen et al., 2002); in Tunisia (Ellouze and Abida, 2008); in Pearl River Delta in South China (Yang et al., 2010) etc, were adequately applied the model for selecting best distribution. The statistics test is given by:

$$Z^{DIST} = \frac{(\tau_4^{DIST} - t_4^R + \beta_4)}{\sigma_4} \quad (3.42)$$

where,  $t_4^R$  is the regional average L-kurtosis of the observed dataset in the homogeneous region and  $\tau_4^{DIST}$  is the theoretical L-kurtosis, and  $\sigma_4$  is the standard deviation of  $t_4^R$  obtained from 1000 times Monte Carlo simulations using  $DIST$  distribution ( i.e., from the candidate regional distribution). The test declared that a particular distribution is considered acceptable at the 90% confidence interval if  $|Z^{DIST}| \leq 1.64$ .

For this work, both the above methods were adapted for the goodness- of-fit test statistics at 90% confidence interval. This is because both methods have advantages and disadvantages. In agreement with the suggestions by Hosking and Wallis (1997), the  $Z^{DIST}$ -goodness-of-fit test was used to choose the best fitted regional distribution (i.e., the lowest  $|Z^{DIST}|$  value was chosen as the best fitted distribution), but the model was not able to test two parameter distributions. Meanwhile, the Anderson-darling test is able to test the fitness of all types of extreme event distributions for rejection, but it is less applicable in the selection of an appropriate regional distribution. Hence, both methods were used in order to fulfill the gaps one by the other.

### 3.5 Regional flood frequency curve

In every RFFA, the main goal of the analysis is to develop regional frequency curve that can represent the average weighted distribution of the homogenous regions. It is the final procedure of flood frequency analysis to estimate the normalized regional quantile floods

( $X_T$ ); flood frequency curve ( $X_T$  vs.  $T$ ); and at-site flood quantiles,  $Q_T$  for a give return period,  $T$ .

For a given region, the model parameters derived from the best fitted distribution to the observed data are the most essential one. Because, these values are used to compute standardized quantile estimates,  $X_T$  for the return periods  $T$ , and then used to construct regional frequency curves for the homogenous region (i.e., a curve showing  $X_T$  against return period,  $T$ ) (Hosking, 1990; Hosking and Wallis, 1997; Kachroo et al., 2000; Mkhandi et al., 2000; Rosbjerg, 2007; Stedinger and Lu, 1995a; Yang et al., 2010). As mentioned earlier, this curve is assumed to be valid for all sites in the region.

The regional growth curves for southern Africa were constructed by performing the following steps;

- The parameter values such as shape ( $k$ ), location ( $\alpha$ ) and scale ( $\mu$ ) parameters for the best fitted distributions were estimated using the regional and theoretical relationships (see the details in Appendices E: Table 13, *Hosking (1990)*, *Hosking and Wallis (1997)* and *Viglione (2010)*).
- The model parameters estimated for a given region were then used to compute the standardized quantile estimates for the return periods  $T$ , where  $T= 2, 5, 10, 20, 50, 100, 200$  and  $500$  years.
- As a result, the regional frequency curve (i.e.,  $X_T$  vs.  $T$ ) for each region was developed
- The at-site quantile floods  $Q_T$  can be then up scaled from the regional quantile flood ( $X_T$ ) through equation 3.25.

### **3.6 Evaluation the performance of frequency distributions**

The results which have been obtained from statistical analysis are essentially uncertain, and to be trustful, methods of uncertainty assessments should be applied (Hosking and Wallis, 1997). The authors also illustrate that “the assessment of the accuracy of the estimates should therefore take into account the possibility of heterogeneity in the region, misspecification of the frequency distribution and statistical dependence between observations at different sites, to an existent that is consistent with the data”. Hence, for this analysis, two methods of

uncertainty assessments were performed. These are the plotting position i.e., quantile-quantile plot using Monte Carlo simulations; and regional growth curve verification.

### **3.6.1 Quantile-quantile (qq) plots**

The performance of the best distribution model identified for the respective regions were evaluated by comparing observed values with simulated values. The argument was that the values that obtained by randomly simulated after 1000 times of Monte Carlo simulations should be matched to the particular characteristics of the data (i.e., the intersection of the values should be closed to the line 1:1). The best frequency distribution was subjected to randomly simulate the same size as observed series. Thus, the quantiles of the normalized streamflows and simulated values are plotted on one graph that represents on the *x-axis* and *y-axis*, respectively.

### **3.6.2 Growth curve verification**

The developed regional curves or models have to be validated mostly by an independent dataset. Identification and selection of stations for validation are required which depend highly on data availability. However, in this study, the insufficient number of stations in all the countries was a main constraint to choose stations for model validation. Consequently, a total of 10 stations from two countries having more than 17 stations were chosen. Of these, 8 gauging sites were from South Africa and 2 were from Zambia (Appendices A: Table 6). Similar to other stations, these stations were subjected for data screening.

In order to obtain independent verifications, these series were withdrawn from the analysis (i.e., all the stations were used neither for homogeneity measures nor for derivation of regional frequency distribution). Hence, the at-site frequency curve was established by up scaling the regional curve (i.e., by multiplying the regional curve by the index flood of the specified stations through equation 3.25) and the observed AMS data was then used for validating the curves for a particular region.

However, the comparison of the regional flood frequency distribution against the at-site observed data wouldn't mean that they can be used to discriminate the regional distribution curves. This is because, the at-site data is representing only one of an infinite number of

relations of the real underlying population or in the inverse that the regional curve is the average of the numbered of at-site statistics in the region. However, the probability plots that shown the observed values together with the simulated from the regional values may reveal such as systematic regional bias in the estimation of the quantile events (Kjeldsen et al. (2002).

### 3.7 Regional estimation for ungauged catchments

In case of ungauged catchment, the main problem is that estimating the appropriate index flood without any given data or insufficient data. In this case, the mean or median values of the ungauged catchments can be estimated using a relation between the index flood ( $\bar{Q}$  or  $\tilde{Q}$ ) and catchment characteristics, obtained using multiple regression from the available neighboring data sets (Cunnane, 1988; Ellouze and Abida, 2008; Noto and La Loggia, 2009; Rosbjerg, 2007).

Different authors used different relationships because the relationships are related to the available catchment information. For example, the relationships that presented by Cunnane (1988) to estimate the index flood for ungauged catchment are:

$$\bar{Q} = cA^a(MAP)^rS^s \quad (3.43)$$

Where,  $A$  is the catchment area,  $MAP$  is the mean annual precipitation,  $S$  is slope of the catchment, and the  $c, a, r$  and  $s$  are the regression coefficients obtaining from the relationships of homogenous catchments.

Another example is the regression equation established by (Rosbjerg, 2007). This relates the index flood,  $\mu$  with the catchment characteristics such area ( $A$ ) and the mean annual precipitation ( $MAP$ ) as follows;

$$\mu = a(A.MAP)^b \quad (3.44)$$

where,  $a$  and  $b$  represents the regression coefficients.

This index flood-catchment characteristics relationship was adapted in this analysis. For each of the countries or regions identified in section 3.4, the index flood was supposed to relate with the catchments characteristics such as Area, mean annual precipitation (MAP), topography etc by means of multiple regression. However, for this analysis, the available catchment information was only the catchment area. Thus, a simple regression model was developed for those defined regions. Thus, these regression equations could be later used to estimate the index flood for ungauged catchments.



## **4. RESULT AND ANALYSIS**

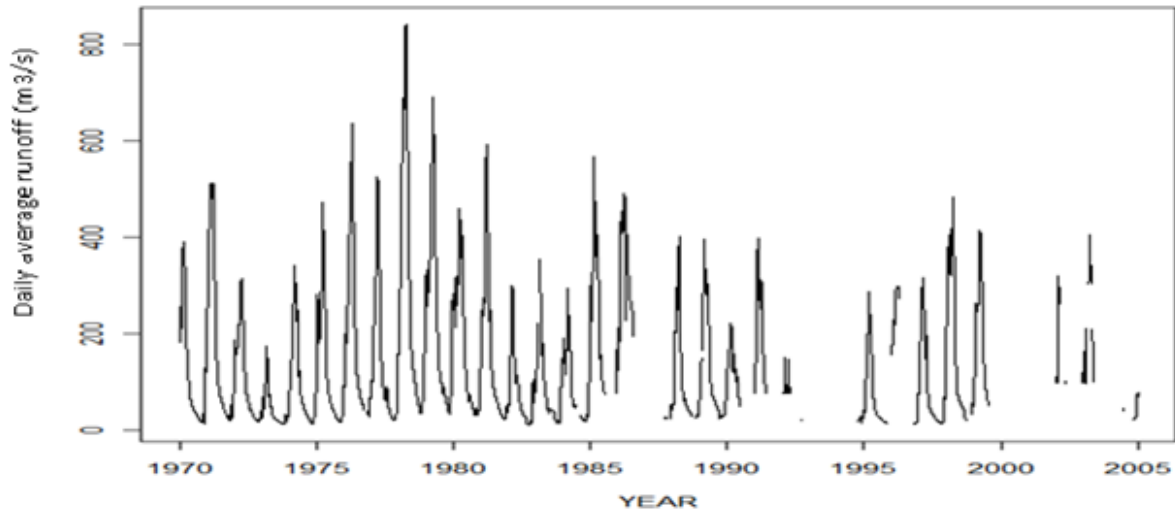
### **4.1 Exploratory data analysis**

This section presents the results from the data screening and later data analysis. A pooled daily average runoff data from 459 gauging stations were collected and AMS data from each station were generated. The stations containing AM series were then subjected to screening under different criteria and the outputs are presented in the following sections.

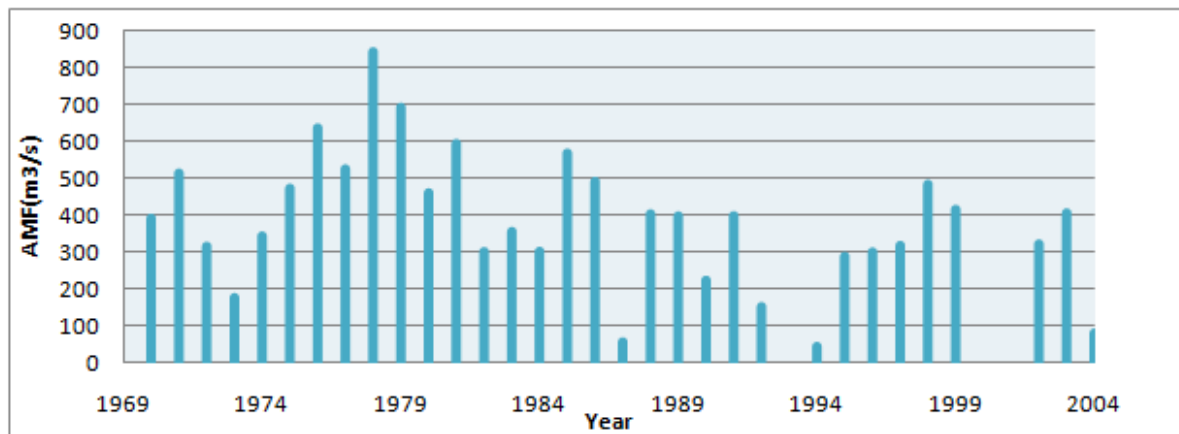
#### **4.1.1 Data Screening**

After all the sample series that were collected from all stations and different sources were examined at their data record lengths, the AMS collected from more than 50% of the stations were with short record length (< 15 years), non-continuous (with consecutive gaps for years) and recorded at different periods. This indicates that more than 230 stations were excluded in this section.

Following these preliminary data screening, the AMS of the rest 229 stations were examined with time series plots (Fig 4.1) to observe the significance of the observations for analysis in terms of quality and quantity. Fig 4.1 shows an example of the observed daily average runoff series and the generated AM floods of Kafue River at Ndubeni, Zambia. Even though the time series of the station looks inadequate (especially the daily series) to represent the flood information, this is a typical example for the stations which were used for this analysis. As it can be seen from Fig. 4.1, most of the stations in both time series (i.e., the daily data and AMS) contain a lot of gaps even with no data for more than a year. Besides this, we can clearly observe from the graph, that the daily average runoff curves show unequal and non-uniform changes over time in the distribution. It is difficult to conclude the data behavior with gaps of information (i.e., no available data for some periods); however, the curves reflect that the nature of the data is non-stationary with time.



a)



b)

Figure 4.1 Plotting observed data series from Zambia station '1591470' (1970-2004): a) daily average runoff series; and (b) AMS.

#### 4.1.2 Autocorrelation and Spatial Correlation

The independence of the time series was tested by computing the first order serial autocorrelation coefficient and then checked it for significance at the 95% confidence interval (Fig. 4.2). As it can be seen from the diagrams, the daily average runoff series is strongly correlated with pronouncing periodic fluctuations. However, except for lag one correlation, the AMS that generated from this parent population are not significantly correlated. This is a

typical example of autocorrelation analysis in this study, which illustrates that though the daily runoff sounds periodic dependence, the AMS serial correlations plots reflect insignificance dependence with time at the 95% confidence intervals.

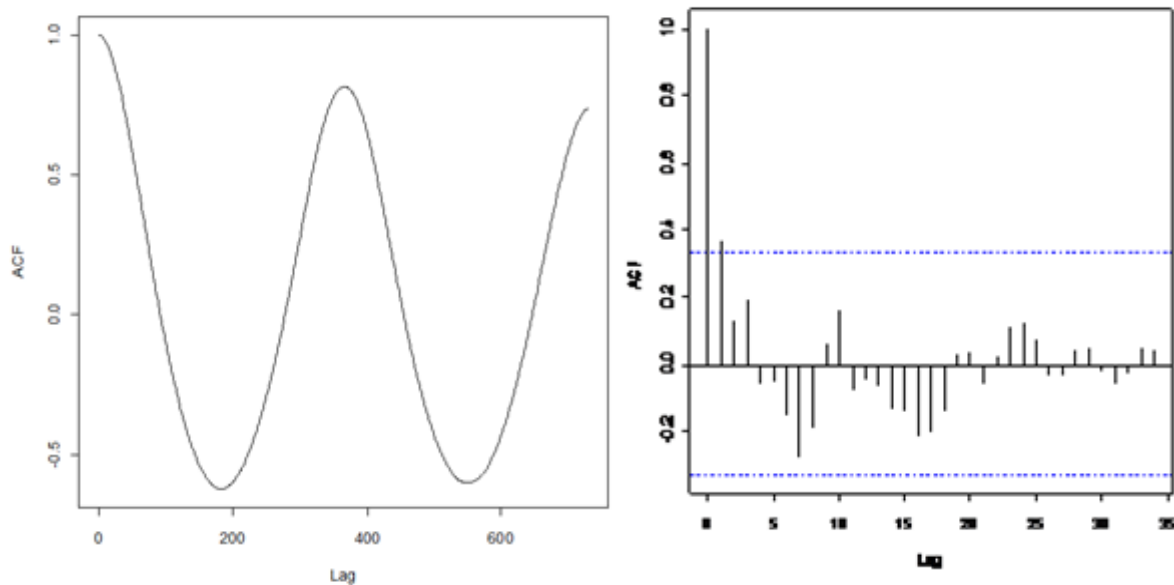


Figure 4.2 the auto-correlation plots for station '1591001': the left is for the daily time series for two years time lag and the right is for the AMS for the time lag of 35 years at 95% confidence interval (the dotted lines at the right plot).

The result obtained from the cross-correlation coefficients also illustrated that the relationships of AMS in across stations were insignificant in all catchments collected from Malawi, Zambia (except 4 stations), Zimbabwe, and Namibia countries. However, the pair correlation coefficients of South Africa catchments showed that good spatial correlations between stations and 64 stations had significant correlations with correlation coefficient  $> 0.89$ . Of these, 34 of the stations were in the analysis while the others were excluded.

### 4.1.3 Empirical distribution

The graphical representation of the relationship between the observed values and their recurrence probabilities were plotted using the Gumbel plotting positions. The horizontal axis of the probability plot shows reduced variate in linear scale and the random variable values,  $x$  were plotted on vertical axis (Fig. 4.3). This was carried out after the streamflows were normalized by their index flood.

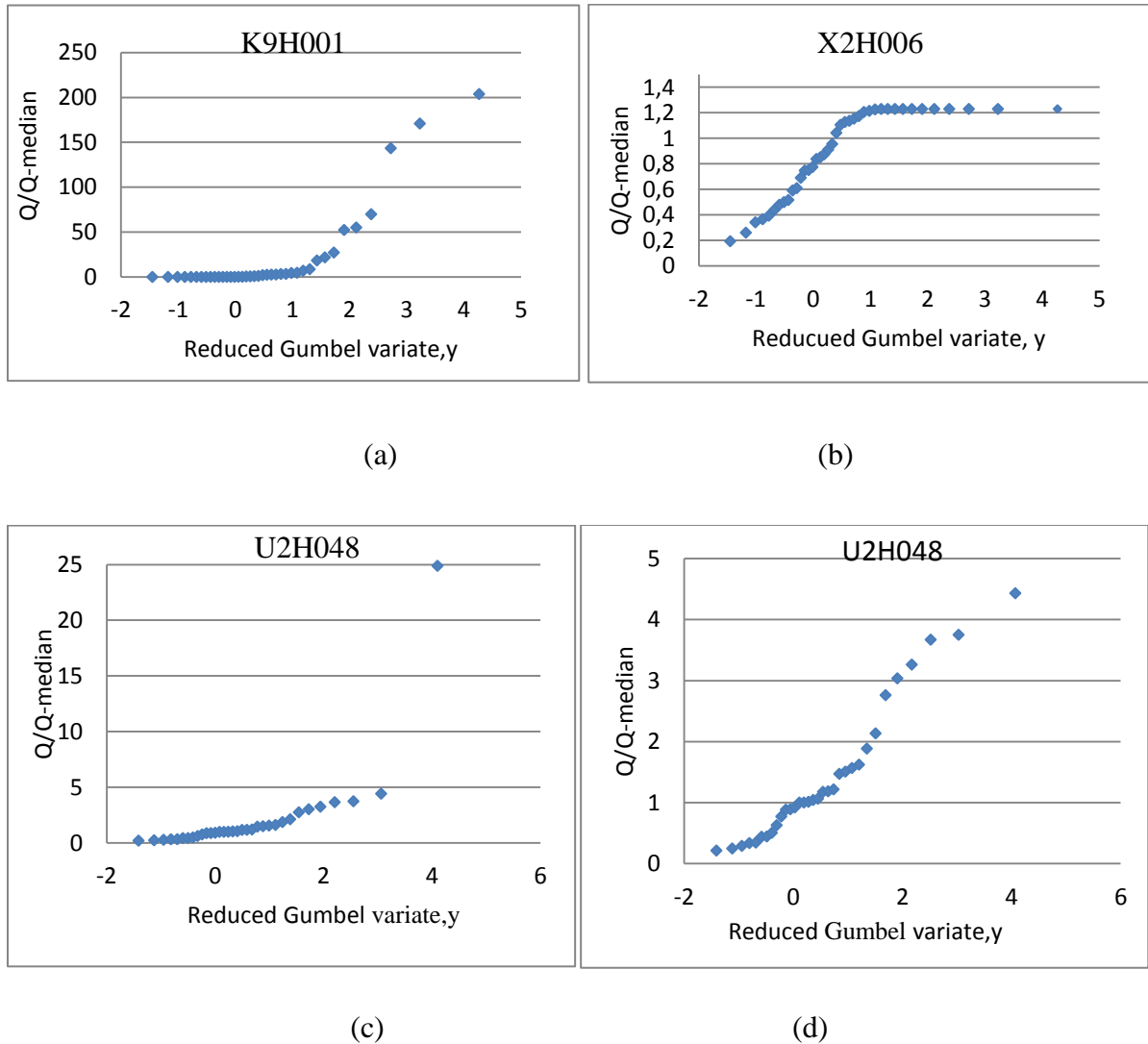


Figure 4.3 An example of Gumbel Probability Plotting for the normalized annual maximum floods from three stations in South Africa: (a) station 'K9H001' which contains lower bounded observation; (b) station 'X2H006' with upper bounded normalized series; (c) station 'U2H048' contains annual maximum flood series with outlier; and (d) station 'U2H048' presents annual maximum flood series after an outlier has been removed.

Fig.4.3 (a) shows that the distribution is bounded to the lower part of the curve (that the series gives repeated values of  $x_l$  or nearly values) with high values to the upper side (i.e., surprisingly the normalized values extended up to 200). On the other way, Fig.4.3 (b) shows that the frequency of the observations is bounded to the upper side of the curve (that the series gives repeated values of  $x_n$ ). Both the curves reflect that the flood events from these stations did not follow the assumptions of random distribution.

Whereas, Figs. 4.3 (c) and (d) for station ‘U2H048’ reflect that the influence of one outlier on probability plot of the observations series. Fig. 4.3 (d) shows the probability of the observed values after one outlier had been removed from the series in (c). The largest normalized value was diminished from 24.27 to 4.5 (i.e., almost four times lesser than).

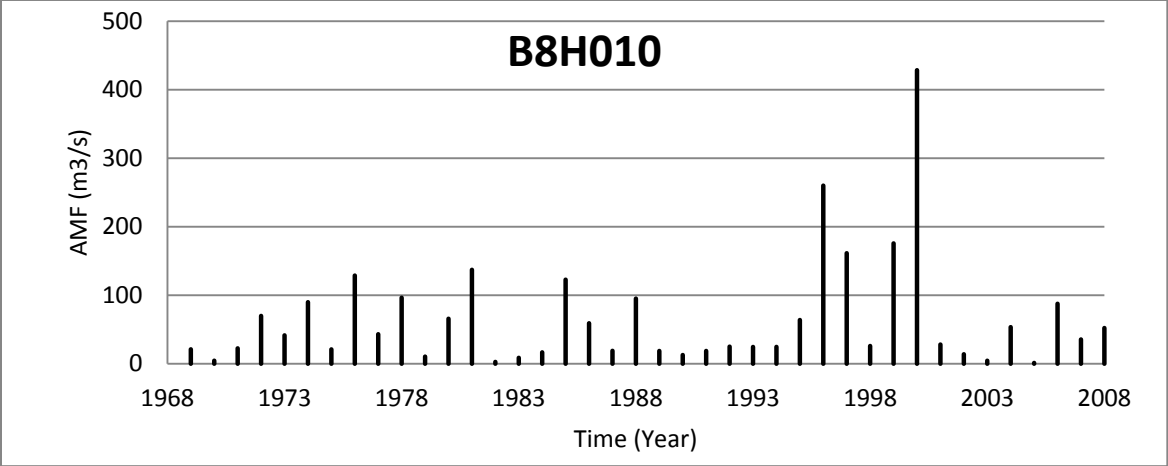
After the empirical distribution of all stations was thoroughly examined, the AMS from 73 stations didn’t have smooth change of curves and bounded into either or both sides of the frequency distribution curve (for example, Figs. 4.3 (a) and (b)). The flood distribution from these stations were considered as erroneous and excluded from the analysis. Thus, a total of 122 stations were screened out for this analysis. However, stations which contain likely outliers were subjected to a further outlier analysis (see section 4.1.4).

#### **4.1.4 Outlier Analysis**

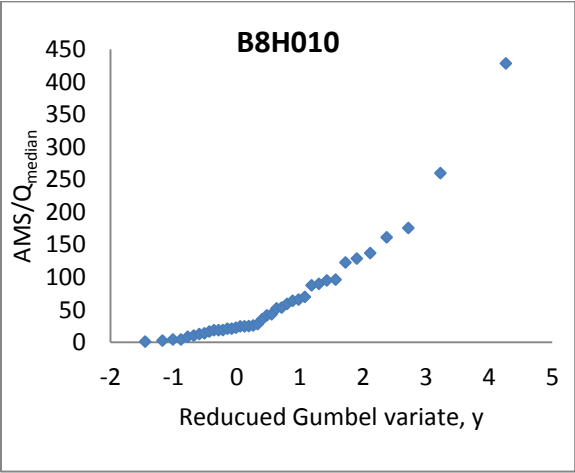
Table 4.1 presents the summary of outlier’s analysis in Malawi catchments. This is a typical example used to show the outlier analysis conducted in this thesis. This has done in two phases; initially possible outliers were indentified based on the graphical inspection from the plotting time series, probability plots and histograms (Fig 4.4). All observations which are located at a distant apparently from the rest of the series were preliminary selected as suspected outliers. For instance, station ‘B8H010’ illustrates the presence of very far positioned large observation with a flood of 428.47m<sup>3</sup>/s in 2000. Following the outlier identification, the degree of the significance was checked using the skewness coefficient, the test for threshold value from Bulletin 17B test and test for significance (Table 4.1).

It is known that the coefficient of skewness (i. e.,  $|Cs| > 0$ ) tells us the existence of skweness in the series. Hence, from Malawi catchments, it can be seen that stations 1992100, 1992900,

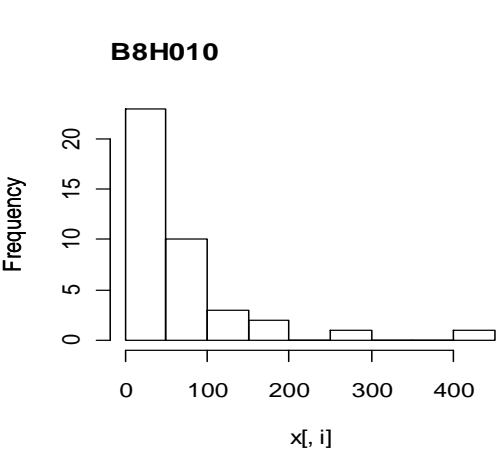
1992950 and 1992690 disclosed strong positive skewness with coefficient  $>1$  (especially, station '1992100' with  $C_s = 5.2$ ) while the other four stations contain reasonable sample series (i.e., their coefficient values are less than one). These results therefore strongly supported for the outlier analysis using Bulletin 17B test. When the observed series in stations which reflect strong skewness coefficient were examined under the tests, all were shown one or two big observations in the series and were treated as special observations.



a)



b)



c)

Figure 4.4 Examples of visual inspections of the AMS containing suspected outliers for station B8H010 in South Africa: (a) Time series plotting of the annual maximum series (AMF); (b) Gumbel plotting positions; and (c) Histogram of the annual extreme events.

However, the outlier analysis using the student test and Bulletin 17B test were slightly disagreed. Even though the outliers obtained from Bulletin 17B test were significant under the t-statistical test at the 95% confidence interval, the total number of observations which were

detected as outliers were not the same. The total number of outlier throughout the region obtained from Bulletin 17B test and t-test were 5 and 10, respectively. The final decision was taken by considering the references that supplied from the statistical values, the visual inspections and personal judgments. As a result, for example, from Malawi catchments three observations have been considered as real outliers. These procedures were applied for all stations in Southern Africa provided for this study.

Initially, 87 of the 122 stations have experienced at least one flood which is much larger than the second highest annual maximum flood. From these stations 95 floods have been detected as suspected outliers and more than 85 of the outliers were found in South Africa stations. This implies that almost all the AMS of every station in South Africa have at least one extreme event. As Kjeldsen et al., (2002) explained that the annual maximum floods of 1974, 1976, 1978, 1884, 1996 and 2000 were among the largest records of 40 years floods in south Africa. During this period almost all catchments of the country have uniformly high floods and most of the outliers were recorded during this time.

Finally, the numbers of stations that contain outliers were reduced from 87 to 53 (of which 47 stations were in South Africa) and the total numbers of the outliers were filtered from 95 to 62. However, in agreement with the recommendation of Cunnane (1989) on treatments of outliers in AMS, the analysis was conducted regionally. Hence, for every outlier: the station, region, normalized value, significance and their sensitivities were estimated from the regional data characteristics and will be discussed later in sections 4.2.

Table 4.1 Summary of the outlier analysis in Malawi gauging sites

Station		1992100	1992200	1992850	1992400	1992700	1992900	1992950	1992690
Top three observations	1 <sup>obs</sup>	639.8	398.2	571	338	1073	2142.03	850	523
	2 <sup>obs</sup>	184.3	373	498.37	291.1	954	2005.49	795.52	349
	3 <sup>obs</sup>	61.13	350	488	287	952	1568	631.5	247
Mean	$\bar{x}$	53.11	219.94	242.29	166.09	620.19	1015.47	266.1	134.45
Standard deviation	Stdv.	107.96	90.89	126.02	78.65	191.58	365.25	211.58	102.87
Coefficient of skewness	$C_s$	5.2	-0.04	0.81	0.25	0.22	1.47	1.43	1.93
Bulletin 17B method	$x_H$	329.3	452.4	564.6	367.2	1110.1	1949.6	807.2	397.5
Number outliers <sup>14</sup>		1	0	0	0	0	2	1	1
Student–test	t <sub>1</sub>	5.434	1.961	2.608	2.186	2.364	3.084	2.76	3.777
	t <sub>2</sub>	1.212	1.684	2.032	1.59	1.742	2.711	2.502	2.086
	t <sub>3</sub>	0.074	1.431	1.95	1.537	1.732	1.513	1.727	1.094
Degree of Significant ( $\alpha=0.05$ )	1 <sup>obs</sup>	F	P	F	F	F	F	F	F
	2 <sup>obs</sup>	P	P	P	P	P	F	F	F
	3 <sup>obs</sup>	P	P	P	P	P	P	P	P
Number of significant outliers <sup>15</sup>		1	0	1	1	1	2	2	2

*NB: from every station in the study area, the largest three observations have been chosen for student t-test. The letters in the last rows of the table indicates the result of significant outliers i.e., ‘P’ indicates insignificant outliers while the ‘F’ represents significantly outlying observations.*

<sup>14</sup> The total number of outliers in each station from Bulletin 17B test

<sup>15</sup> The total number of outliers in each station from student- test



#### **4.1.5 At-Site flood characteristics**

The at-site statistics for all stations which includes the first four ordinary moments and the first three L-moment ratios are presented in Appendices B (Tables 7- 11). As discussed earlier in section 4.1.4, almost all the stations were with one or two flood events that are much larger than the rest of the AMS in the sample. This was clearly observed on the third order moment of the at-site statistics (i.e. the skewness coefficient). Though positive skewness was expected from the frequency of extreme maxima, the magnitude of the coefficients was relatively high in almost all stations of Southern Africa except the Zambia and Malawi AM series. However, except for exceptional stations the AM floods of Zambia and Malawi are almost symmetrically distributed. The Zambia flood statistics also included four stream gauging sites with negative skewness coefficients though the values were small. The magnitudes of the coefficient of variance also illustrated the proportion of the flood series. The  $C_v$  values for Zambia and Malawi stations were slightly lower when compared to others, i.e., their average  $C_v$  values are in the range 0.4-0.6. Whereas, the AMF from other countries reflected with average  $C_v$  values close to 1.0 and beyond in South Africa sites. The different levels of variability in the observed flood samples might be attributed to varying hydrological phenomena that generated the flood events over the different regions. This reflects that the AM series collected for this study were well- behaved heterogeneous and most of the stations that showed large  $C_s$  could be also due to the presence of one or more outliers.

#### **4.1.6 Choice of the Index Flood**

Table 4.2 presents the degree of the sensitivity of sample median versus sample mean when outlying observations were with and without the sample series in the flood frequency analysis. The result shows that the sample mean of the index flood is more sensitive to outliers than the median. After the largest observation were excluded from the series, the relative deviation of the sample mean is changed with the range 6.18-33.16% while the median is varied in the range 0.09-3.04 %, which is around 30% less sensitive than the mean.

This supports the suggestions by Hampel (1974) which claims that the sensitivity of the sample median is known to be less than the sample mean in case outliers exist in the sample and this phenomenon was more likely found in samples from highly skewed distributions.

Table 4.2 The sensitivity analysis of the index floods to the largest observations in stations which contains one or two large outliers.

Index Flood ( $\mu$ )	Stations contain outliers							
	W2H005	W2H006	W2H009	W5H005	V1H001	V2H004	V3H002	V7H020
The values of index flood with outliers, $\Theta_1$								
Mean ( $\bar{Q}$ )	161.8	141.1	46.3	28.9	433.2	103.4	46.4	95.3
Median ( $\tilde{Q}$ )	111.8	92.3	23.9	16.4	258.3	73.1	29.7	65.4
The values of index flood after removing one large observation from the series, $\Theta_2$								
Mean ( $\bar{Q}1$ )	149.0	130.4	30.9	25.4	406.4	91.0	43.0	85.4
Median ( $\tilde{Q}1$ )	108.4	91.2	23.8	16.1	251.0	72.7	28.8	63.8
The relative difference in $\% = \left(\frac{\theta_1 - \theta_2}{\theta_1}\right) * 100$								
Median ( $\tilde{Q}$ )	3.04	1.21	0.09	1.81	2.83	0.55	3.02	2.46
Mean ( $\bar{Q}$ )	7.92	7.55	33.16	12.23	6.18	12.00	7.21	10.31

Therefore, the median of the observed series (which has the quantile probability of  $p_r = 0.5$ ) of each site was considered as index flood for all catchments in Southern Africa. This agrees with suggestions' by Viglione et al., (2007) and Noto and La Loggia (2009) that if the series of the sample is skewed to the right the median might be better index flood than mean. The median of each station is therefore estimated and presented in Appendix A (Tables: 1-5). Henceforth, the flood of Southern Africa, i.e., the observed AMS at each site was normalized by their midpoint of the AM series and in return the regional curves will be up scaled to at-site quantile flood by multiplying by the median values through equations 3.24 and 3.25, respectively.

Though the median was preferred as better index flood, quantifying the uncertainty of the median was another difficult task. Nevertheless, an attempt was made to solve this problem using Jackknifing method of standard error estimation. The method is adapted in different studies since 1958 to estimate the standard error. For example, details are discussed by Efron and Gong (1983) that the standard error of the mean can be easily computed from the sample values. However, the trouble with this is that there is no obvious ways to extend to estimators other than sample mean, for example, the sample median. The jackknifing could be therefore an alternative way of making this extension. The values were estimated by create sample data sets from the data leaving out one data point at a time and take the median of these sample sets.

The jackknife estimate of standard error of the median  $SE_{jack}$  is given by:

$$SE_{jack} = \left[ \frac{(n-1)}{n} \sum_{j=1}^n (\tilde{Q}_j - \tilde{Q}_j mean)^2 \right]^{1/2} \quad (4.1)$$

where  $\tilde{Q}_j$  is jackknife median value and  $\tilde{Q}_j mean$  is average of jackknife median values from the sample set. The relative standard errors of the sample median for each station was estimated using R-bootstrapping package<sup>16</sup> and presented in Appendices A (last columns of Tables:7-11).

*SE<sub>jack</sub>*

## 4.2 Identification of homogenous regions

### 4.2.1 Delineation of homogenous regions

The grouping of sites in to homogeneous regions was carried out by applying the hierarchical geographic regionalization method. The procedures considered that the entire region should be geographically continuous and at every grouping step, the AM series from different sites in the region should satisfy the Hosking and Wallis' (1997) Homogeneity test. At the initial step, i.e., when all the stations were considered as a region, the computed heterogeneity measure was very large with the magnitude of  $H = 33.2$ . This value provides that the statistical characteristics of the sites were strongly different and further divisions of sites in to homogenous regions were needed.

Therefore, the first step was to separate sites into their respective countries and test for homogeneity. Since the gauging stations in the countries of Namibia, Zimbabwe, Zambia and Malawi were few in number and covered large areas between sites, each country was considered as a region (Fig. 4.5). In fact, the stations at every country level were still not sufficient that the total stations provided for this analysis were varied from 7-15 (Table 4.3). Thus, the AM series from these stations may not necessarily satisfy the condition of regional homogeneity and consequently, all failed to the Hosking and Wails (1997) heterogeneity test (Table 4.3). Especially, the heterogeneity measure for Namibia catchments was very big value, i.e., with  $H = 7.5$ . However, these four countries were considered as acceptable regions in agreement with Hosking and Wallis (1997) suggestion that “even though a region is

<sup>16</sup> <http://cran.r-project.org/web/packages/bootstrap/index.html>

moderately heterogeneous, regional analysis will still yield much more accurate quantile estimates than at-site analysis”.

Since around 75% of the total stations used for this study were collected from South African catchments, the further classification of catchments into statistically homogenous regions was implemented only in this country. The choice of which sites do belong to their appropriate regions was performed using previous studies by Mkhundi and Kachroo (1997), geographical consistency of the regions and at-site L-moments (i.e., basically the at-site  $L-Cv$ ). As a result, 9 regions (5 of them were possibly heterogeneous regions while the other 4 regions were defiantly heterogeneous) were formulated and the results are presented in Fig. 4.5.

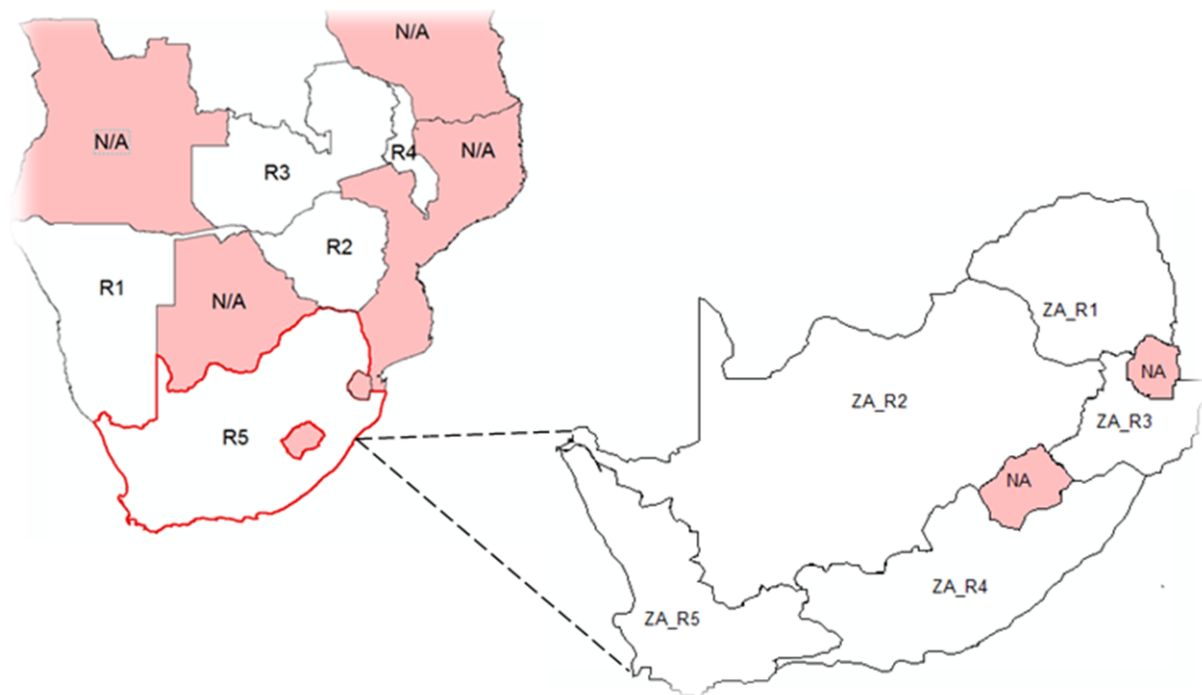


Figure 4.5 Delineation of southern Africa catchments into hydrologically homogenous region. The further classifications of South African drainage areas are shown in the right side of the map. The abbreviation NA indicates the countries or regions which have no available data.

#### 4.2.2 Heterogeneity test

The results summarized in Table 4.4 show that regions  $R1$  and  $R4$  reflect definitely heterogeneity with  $H$  values 7.5 and 5.43, respectively. For regions  $R2$  and  $R3$ , the heterogeneity measure suggested that the regions were moderately heterogeneous (i.e., with  $H$

values 2.31 and 2.35, respectively) while the other five regions of South Africa are possibly heterogeneous regions with  $H$  values in the range of 1.41-1.91.

When the sites containing all AM series including the outliers listed in Table 4.4 were used for heterogeneity test, it was difficult to form a region with sites which have the same statistical behaviors in its nearby stations. That is, the existence of outliers in the sample series was the main problem during homogenization of the regions. However, after the outliers listed in Table 4.4 were excluded from the series, the regions that are listed in Table 4.3 were formed with possible regions.

Table 4.3 The result of Hosking and Wallis' (1997) homogeneity test and the summary of their regional statistics of Southern Africa Catchments

S.N	Region	Drainage basins	NO stat.	$H$	Regional L-moments			Degree of Heterogeneity
					L-Cv	L-Cs	L-kur	
*	<b>All stations</b>	All southern Africa	112	33.2	0.45	0.32	0.19	Strongly Heterogeneous
1	<i>R1</i> (Namibia)	All Namibia	8	7.5	0.391	0.259	0.132	Definitely Heterogeneous
2	<i>R2</i> (Zimbabwe)	All Zimbabwe	7	2.31	0.487	0.280	0.117	Moderately Heterogeneous
3	<i>R3</i> (Zambia)	All Zambia	15	2.35	0.268	0.037	0.114	Moderately Heterogeneous
4	<i>R4</i> (Malawi)	All Malawi	8	5.43	0.287	0.193	0.154	Definitely Heterogeneous
5	<i>ZA_R1</i>	<i>A5-A9, B &amp; X</i>	14	1.43	0.508	0.372	0.186	Possibly Heterogeneous
6	<i>ZA_R2</i>	<i>A1-A4, C &amp; D3-D8</i>	17	1.91	0.550	0.439	0.223	Possibly Heterogeneous
7	<i>ZA_R3</i>	<i>V &amp; W</i>	13	1.41	0.430	0.350	0.196	Possibly Heterogeneous
8	<i>ZA_R4</i>	<i>D1 &amp; 2, R, S, T, U, L and Q</i>	21	1.88	0.397	0.298	0.166	Possibly Heterogeneous
9	<i>ZA_R5</i>	<i>E, G, and H</i>	8	1.82	0.379	0.263	0.105	Possibly Heterogeneous
	Total		112					

*Note: \*Denotes when all stations supposed to be one region and the letters A-X in the drainage basin 'column' denotes the drainage regions of South Africa catchments. The specification of the drainages is described on the website of South Africa, Department of water affairs, hydrology section.<sup>17</sup>*

### 4.2.3 Regionalization of outliers

<sup>17</sup> <http://www.dwaf.gov.za/hydrology/cgi-bin/his/maps/Drainage%20Regions.htm>

After all the procedures in section 4.1.4 were performed and the results were analyzed, the next step was to decide “how to handle these extreme of extreme events”. Table 4.4 presents the summary of the regionally significant outliers from all the regions which were defined in sections 4.2.1 and 4.2.2. For every outlier: the station, region, normalized values were calculated from the characteristics of the regional data. As it can be seen from Table 4.4, the outliers were almost located in all regions of the study area except in Zimbabwe catchments. The total number of outliers in every region varied from 2-7 and their normalized values were extended from 4.63 to 26.36.

Table 4.4 Summary of regional outliers of Southern African floods

Region	Station	Year	Normalized Outliers
<i>R4</i>	1992100	1987	21.77
	1992900	1984	4.65
	1992690	1978	4.63
<i>R2</i>	63533035	1972	11.58
<i>R1</i>	1258200	1985	9.02
	1259110	1989	8.38
<i>ZA_R1</i>	A5H006	2000	10.78
	B8H010	2000	13.48
	X1H001	1974	10.12
	X1H014	1984	14.87
	X2H022	2000	13.31
<i>ZA_R2</i>	C3H003	1975	10.82
	D3H008	1974	14.93
	D8H003	1988	13.69
	D7H005	1988	15.01
	A2H006	1978	16.22
	A2H013	1976	11.59
	A2H021	1996	13.77
<i>ZA_R5</i>	H6H008	1986	9.50
<i>ZA_R4</i>	R3H003	1970	15.88
	S3H004	1974	14.14
	U2H048	1976	24.27
	T4H001	1975	15.95
<i>ZA_R3</i>	W2H009	1981	26.36
	W5H005	1984	10.20
	V2H004	1984	8.04
	V7H020	1975	7.31

Though majority of the outliers were very large in magnitude compare to the at-site index flood, three observations were recognized as unacceptable observations. The magnitude of these outliers is highly inconsistent with the regional observed values and their return period

is more than 25,000 years. The reason was not well-investigated but all were excluded from the analysis and treated as NO data. Table 4.5 presents the list of these outliers and their probability (recurrence interval), location and year of occurrence.

Table 4.5 Unaccepted outliers

Region	Station	Year	Normalized Outliers	Exceedance probability, $p_r(X \geq x)$	Return period, T (year)
R1	1992100	1987	21.95	$1.21 \cdot 10^{-12}$	$8.3 \cdot 10^{11}$
ZA_R3	W2H009	1981	26.37	$3.58 \cdot 10^{-5}$	277778
ZA_R4	U2H048	1976	24.89	$2.86 \cdot 10^{-5}$	34960

For example, the effect of such outliers can be seen boldly from the diagram shown in Fig. 4.6. The Figure shows the L-moment ratio diagrams for the regional flood distribution of Malawi catchments (left) with and (right) without outlier. When we see at a particular station which indicated by an arrow (Fig. 4.6), the at-site average L-moments were almost decreased by half in case of no outliers in the observation series. This shows how a single outlier in station ‘1992100’, Malawi (Table 4.5) affects the center of the distribution both the at-site and regional L-moments and consequently the choice of frequency distributions.

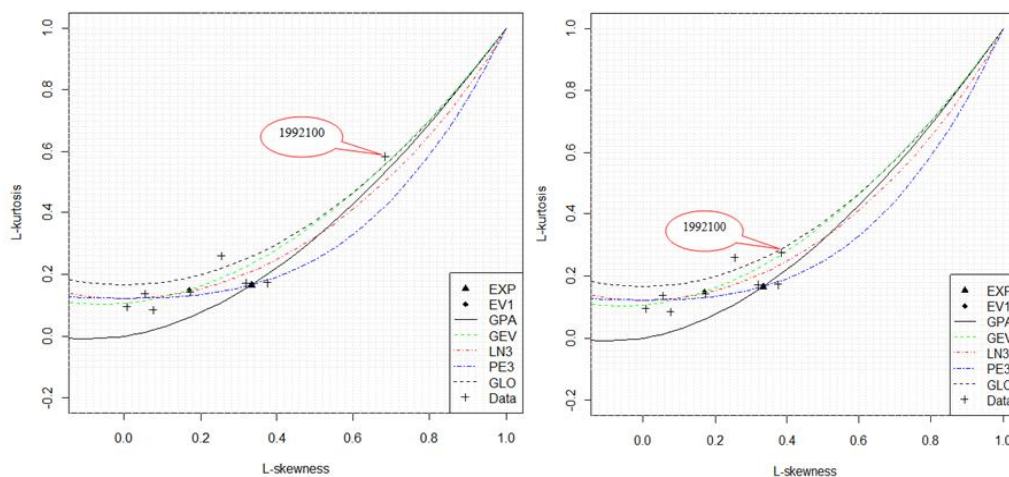


Figure 4.6 L-moment ratio diagram for the annual maximum floods from Malawi gauging sites. The diagram shows the influence of a single outlier in station ‘1992100’ (Table 4.5) in case of fitting theoretical distributions to the regional data.

## 4.2.4 Outlier sensitivity analysis

Table 4.6 presents an example that shows the influence of the at-site and regional statistics by the existence of outliers in a particular region (i.e., stations from region *ZA\_RI*). It is obvious that the existence of extreme events may more affect the kurtosis and skewness parts of the flood hydrographs. The analysis illustrated that though the at-site sample statistics were more sensitive to large observations in the series, the regional weighted average statistics that generated from the at-site observations was less sensitive.

Table 4.6 The comparisons of the relative differences of the at-site and regional weighted average statistics after suspected outliers were removed for the sample series.  $\theta_1$  and  $\theta_2$  are parameters values estimated from the sample series before and after the suspected outliers were excluded from the series, respectively.

### I. At-site L-moments

At-site statistics including the largest observations ( $\theta_1$ )

Station	$\lambda_l$	<i>L-Cv</i>	<i>L-Cs</i>	<i>L-kurt.</i>
A5H006	1.86	0.57	0.43	0.27
B8H010	2.04	0.56	0.46	0.27
X1H001	1.65	0.52	0.45	0.27
X1H014	2.24	0.62	0.52	0.31
X2H022	1.80	0.56	0.48	0.30

At-site statistics after the largest observations have been removed from each station ( $\theta_2$ )

A5H006	1.65	0.53	0.35	0.20
B8H010	1.97	0.52	0.38	0.18
X1H001	1.43	0.47	0.36	0.17
X1H014	1.98	0.59	0.47	0.27
X2H022	1.65	0.50	0.36	0.17

Relative difference (%) =  $\frac{(\theta_1 - \theta_2)}{\theta_1} * 100$

A5H006	11.75	7.02	18.51	26.36
B8H010	3.13	8.06	18.44	34.55
X1H001	12.97	9.17	19.41	35.47
X1H014	11.54	5.27	8.87	13.53
X2H022	8.65	10.62	24.03	43.80

### II. Regional L-moments (*ZA\_RI*)

With outlier, $\theta_1$	1.65	0.51	0.39	0.21
Without outlier, $\theta_2$	1.58	0.49	0.36	0.18
Relative difference (%) = $\frac{(\theta_1 - \theta_2)}{\theta_1} * 100$	4.26	3.42	8.24	15.93

After the suspected outliers were removed from respective stations, the at-site statistics relatively varied from 3.13% to 43.8% whereas, the regional statistics reduced from 3.4% (*L-*



$C_v$ ) to 15.93% ( $L\text{-Kurt}$ ). This implies that the use of these extreme events as random observations might be relatively less influenced the regional analysis.

Therefore, except the observations which are listed in Table 4.5, all the outliers in Table 4.4 were kept in the data sample. This was merely because of three reasons in agreement with the recommendations by Cunnane (1989) and applications by Kjeldsen et al., (2002): (1) the observations were real and occurred randomly (i.e., most of the outliers were recorded when there were high floods throughout the regions); (2) the frequency of outliers in the regional data were very few in number (from 2-7 observations); and (3) the regional statistics were less sensitive to suspected outliers. Hence, all large observations in every region were accepted as random variables which can play substantial role in the analysis of the flood situations in Southern Africa. Thus, henceforth, the regional flood frequency analysis was carried out by weighting all the observations in the samples series.

### **4.3 Identification of regional flood frequency distribution**

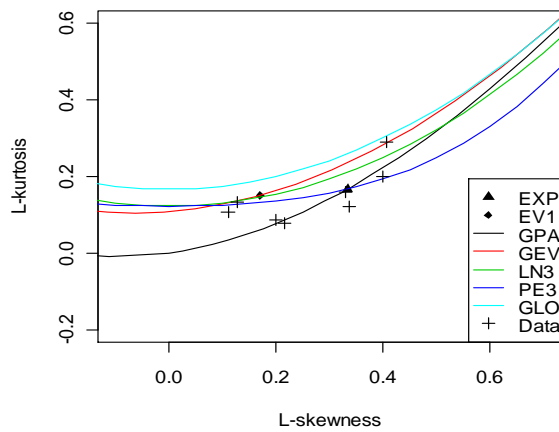
The identification of an appropriate regional flood models for each of the grouped regions was accomplished based on the L-moment ratio diagrams, a goodness-of-fit tests and later evaluated their performance using quantile plots and model validations.

#### **4.3.1 The L-moment ratio diagram**

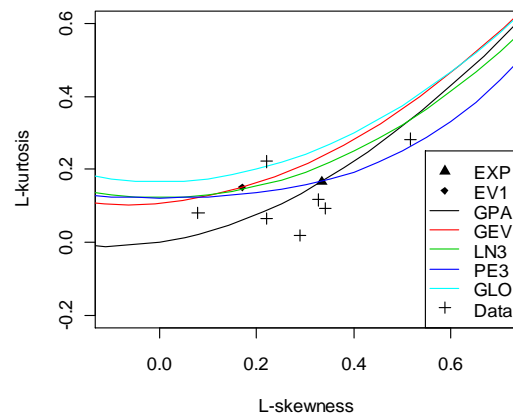
The choice of an appropriate regional distribution was performed initially by comparing the L-moment ratios diagram of the sample L-skewness versus L- kurtosis to the theoretical values. Fig. 4.10 presents the relationships between the population  $L\text{-}C_s$  and  $L\text{-}Kurt$  for a range of distributions, commonly employed in flood frequency analysis. These include the Gumbel (EV1), Exponential (EXP), Three-parameter Lognormal (LN3), Generalized Pareto (GPA), and Generalized extreme Value (GEV), Pearson Type III (PE3) and Generalized Logistic (GLO) distributions.

From the diagrams shown in Fig. 4.10, the visual observation indicates that the GPA distribution could be the best distribution for seven regions which includes regions  $R1$ ,  $R2$ ,

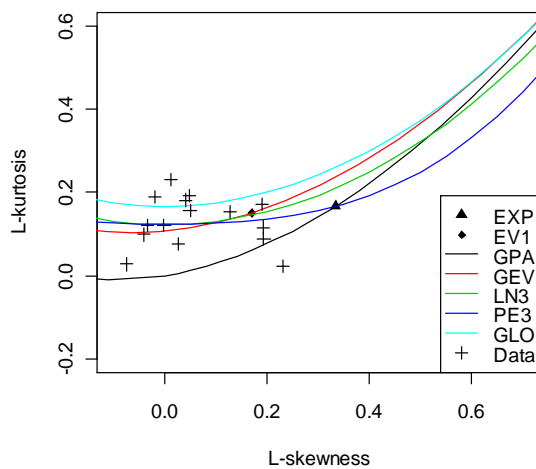
$ZA\_R1$ ,  $ZA\_R2$ ,  $ZA\_R3$ ,  $ZA\_R4$  and  $ZA\_R5$ . For regions  $R1$ ,  $ZA\_R1$ ,  $ZA\_R2$  and  $ZA\_R3$  the PE3 distribution should be also considered as possible regional distributions. The same for Regions  $ZA\_R3$  and  $ZA\_R4$  that the LN3 distribution could be adopted as regional distribution. However, for the other two regions (i.e.,  $R3$  and  $R4$ ), the choice of the best distribution is not easily based on the L-moment diagram. For example, for  $R3$ , it is difficult to conclude simply from the graph as all the LN3, PE3 and GEV might be considered as possible regional distribution models. This is also the same case for Malawi catchments, i.e., the L-moment diagram in Fig.4.10 (d)-the EXP, GEV, LN3 and PE3 could be the candidate regional frequency distributions.



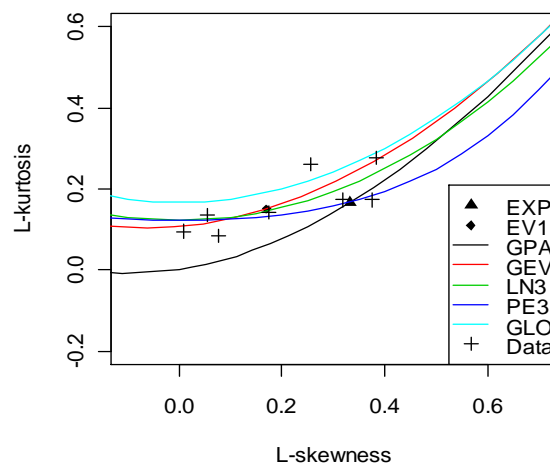
a)  $R1$ -Namibia (NA)



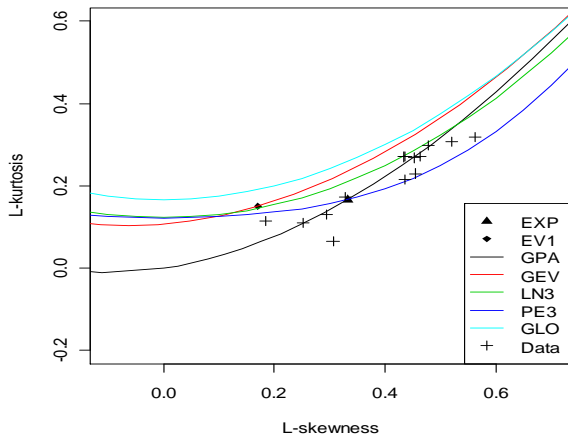
b)  $R2$ -Zimbabwe (ZIM)



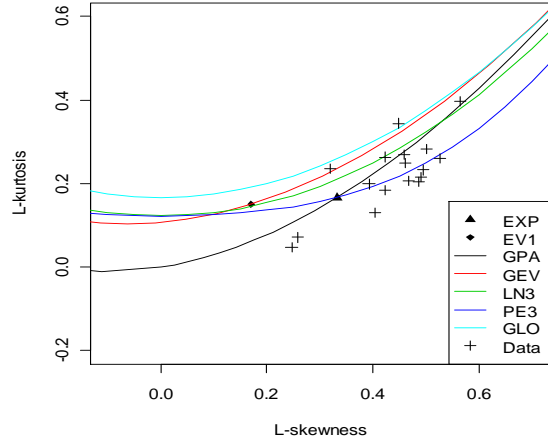
c)  $R3$ -Zambia (ZM)



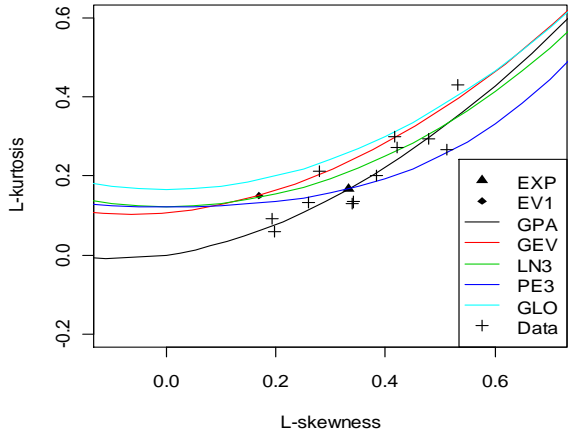
d)  $R4$ -Malawi (MW)



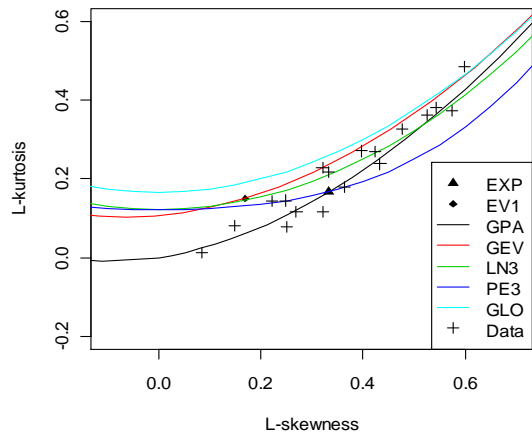
e) ZA\_R1



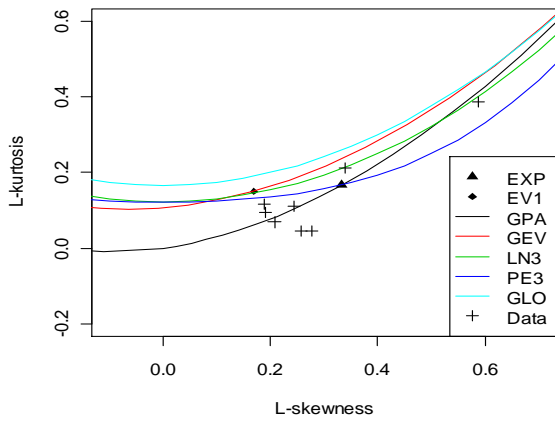
f) ZA\_R2



g) ZA\_R3



h) ZA\_R4



i) ZA\_R5

Figure 4.7 L-moment diagrams showing the relationships between the theoretical distribution curves and the regional data from five countries of Southern Africa: the name of regions is labeled under the pictures from (a-i).

This indicates that, it is difficult to say a given distribution is a best representative of the flood distribution without being doing some conformations. Hence, the suggestions from the L-moment diagram were supported using goodness-of-fit tests presented in section 4.3.2. However, though all the considered distributions were subjected to goodness-of-fit tests, it is clear from the L-moment diagram that GLO distribution couldn't be a model for any region in Southern Africa.

In addition, it can be seen from the diagrams that the weighted statistics of some stations are located a bit far from the groups. For example, Fig. 4.7 (i)-region *ZA\_R6* shows that there is one station which has different statistical behaviors than the other stations in the group. This is station 'H6H009' containing two big floods in the series (see Table 4.4 and at-site statistics-Appendix B (Table 11)). As Kjeldsen et al. (2002) discussed, the significance difference between the average statistical properties of the stations for example, station 'H6H009' and the others in the group, was mainly due to acceptance of infrequent outliers

### **4.3.2 Goodness-of-fit (GOF) measures**

The choices of distributions from the L-moment diagram might not guarantee that the distribution is the real representative of flood statistics in the given region. For this reason, a confirmation of the candidate distributions is needed using the so called Goodness-of-fit test. Hence, two statistical tests namely the Anderson-Darling test and Hosking and Wallis (1993 and 1997)  $Z^{\text{DIST}}$  test were employed.

Table 4.7 summarizes the results from the Anderson-Darling goodness-of-fit test corresponding to the theoretical distributions for 9 regions. The results of the GOF test illustrates that in almost all of the regions, the PE3 and/or GPA were the distributions accepted at 90 % confidence intervals i.e., all the values with asterisk (\*) indicate that the distribution is accepted as regional distribution.

When the results were compared with the diagrams shown in Fig. 4.7, the L-moment diagrams illustrate that except for regions *R3* and *R4*, the GPA should be an appropriate regional distribution whereas, the GOF test indicates that the PE3 could be an appropriate

flood model for all regions ( i.e., the distribution was accepted at 90% confidence interval in all regions). Though the GOF suggests that the PE3 could be a regional model for all regions, the overall suggestions from the test were reasonably agreed with the diagrams. However, there are also some disagreements. For example, this was observed in R4-Malawi flood flows. In the L-moment diagram, the regional data are distributed around the GEV, LN3 and PE3 theoretical curves while the test statistics indicates that the GPA, EXP and PE3 are among the accepted distributions. The same in R3-Zambia that the regional L-moments ratios are close to GEV, LN3 and PE3 curves, whereas the GOF test suggests that the regional data were fitted only to the PE3 distribution.

Table 4.7 The result of Anderson - Darling goodness-of-fit test at 10 % level of significant (Viglione, 2010)

Distriution	Regions								
	<i>R1</i>	<i>R2</i>	<i>R3</i>	<i>R4</i>	<i>ZA_R1</i>	<i>ZA_R2</i>	<i>ZA_R3</i>	<i>ZA_R4</i>	<i>ZA_R5</i>
EV1	1	1	1	1	1	1	1	1	1
LN3	0.981	0.998	1	0.9967	0.8207*	0.955	0.3203*	0.998	0.997
EXP	1	0.992	0.981	0.789*	0.969	1	0.963	1	0.97
GPA	0.95	0.208*	0.974	0.643*	0.167*	0.925	0.898*	0.606*	0.813*
GEV	1	0.916	1	0.998	0.997	0.998	0.991	0.996	1
PE3	0.541*	0.622*	0.541*	0.758*	0.868*	0.649*	0.899*	0.632*	0.734*
GLO	1	0.982	1	1	1	1	0.999	1	1

*\*Denotes regional frequency distributions passed the test statistics*

Therefore, in order to have further confirmations and choose appropriate regional distribution, the  $Z^{DIST}$  goodness-of-fit test was implemented. Table 4.8 presents the result of  $Z^{DIST}$  goodness-of-fit test for the considered three parameter regional distributions. A distribution with goodness-of-fit test value  $|Z^{DIST}| \leq 1.64$  was considered as acceptable regional distribution at 90% confidence level. For a given region, a distribution was chosen as giving the best fit among the chosen candidates if the  $|Z^{DIST}|$  is small value. All distributions that already passed the statistical tests are marked with the asterisks (\*) and ranked with numbers **1, 2, 3** and **4**.

Table 4.8 Hosking and Wallis (1997) goodness-of-fit test statistics for regional frequency distribution.

Distribution	Regions								
	<i>R1</i>	<i>R2</i>	<i>R3</i>	<i>R4</i>	<i>ZA_R1</i>	<i>ZA_R2</i>	<i>ZA_R3</i>	<i>ZA_R4</i>	<i>ZA_R5</i>
GLO	2.84	4.04	2.52	1.19 *(4)	3.05	4.16	2.15	2.25	4.59
GEV	1.72	3.00	-0.88 *(3)	-0.17 *(1)	2.32	3.51	1.36 *(3)	1.37*(2)	3.31
LN3	1.14 *(3)	2.29	-0.33 *(1)	-0.53 *(2)	0.94 *(2)	1.87	0.27 *(1)	0.09*(1)	2.41
PE3	0.08 *(1)	1.04 *(2)	-0.44 *(2)	-1.24 *(3)	-1.44 *(3)	-0.95 *(1)	-1.60 *(4)	-2.12	0.83 *(2)
GPA	-1.12 *(2)	0.23 *(1)	-7.44	-3.36	-0.20 *(1)	1.00 *(2)	-1.13 *(2)	-1.44*(3)	-0.13 *(1)

*Note: The z-values with asterisk (\*) denote that the empirical distributions are accepted as regional flood models and the numbers in the brackets refer to the rank between the possible regional distributions.*

When the GOF results from both methods were compared, the appropriate regional distributions that obtained from both methods were more or less in good agreement except for regions *R3*, *R4* and *ZA\_R4*. The regional distributions that passed under both statistical tests were similar in all regions. However, for example, for region *ZA\_R4*, the Anderson-Darling test suggests that the GPA and PE3 could be the possible regional distributions whereas, the result from  $Z^{DIST}$ -goodness of test reflects that the LN3, GEV and GPA could be accepted regional distributions. This means that the regional distribution which was ranked as best distribution in the  $Z^{DIST}$ -goodness-of-fit test shouldn't be considered under the Anderson-Darling test. Hence, the selection of best distribution was determined based on the average weighted regional L-moments presented in Table 4.9. The regional weighted average of *L-Cs* versus *L-Kurt.* values were more close to LN3, then to GEV and GPA distributions than to PE3 distribution. In addition, for regions *R1*, *R3* and *R4*, the possible regional frequencies distributions were better identified by the  $Z^{DIST}$ -statistical test. That is, the possible regional distributions were better agreed to the L-moment diagrams than the one obtained from the Anderson-Darling test.

When the results from the L-moment ratio diagram (Fig. 4.7) were compared to the results obtained from the goodness-of-fit tests (Tables 4.7 and 4.8), good correspondences were found in most of the regions. Both the fitting criteria suggest possible regional distribution functions i.e., the distributions which had chosen from the visual inspection of the regional data were confirmed and could be considered as appropriate regional distributions for southern Africa flood studies.

Therefore, from the results summarized in both Figure 4.7-L-moment diagrams and Tables 4.7 and 4.8-the GOF tests of the empirical distribution models, the following suggestions were forwarded:

- a) Pearson type III distribution provides best fit for 2 regions: regions *R1*-Nambia and *ZA\_R2*-South Africa catchments.
- b) In three regions which include Zimbabwe-*R2* and both *ZA\_R1* and *ZA-R5* regions of South Africa, the analysis recommends GPA distribution.
- c) LN3 distribution was best fitted model for three regions; the *R3*-Zambia and both *ZA-R3* and *ZA\_R4* regions in South Africa.
- d) Even though the GEV in the L-moment ratio diagram was a candidate in *R3*-Zambia and *R4*-Malawi, the model under the GOF test was best fitted only to flood events in

Malawi catchments and considered as an acceptable distribution for floods  $R3$ ,  $R4$ ,  $ZA\_R3$  and  $ZA\_R4$  regions.

However, the Gumbel (EV1), Exponential (EXP) and Generalized logistic (GLO) distributions couldn't be part of the models for Southern Africa flood studies i.e., none of these models passed under the statistical test for any region of the study. Even though the GLO under the Z-test and EXP under Anderson-Darling test were passed for  $R4$ -Malawi catchments, they could not be as good as the others. For example in the  $Z^{DIST}$  - test, the GLO was ranked on the fourth place and EXP was rejected.

#### 4.4 Regional flood frequency curves

As discussed in section (4.3.1), using the index flood, the annual maximum floods at each gauging station were reduced to dimensionless form before fitting the individual series to regional frequency distribution. Here, the regional flood frequency curves were derived by combined all the dimensionless curves from each of the stations of the entire regions.

- For the best regional distribution, the parameter values (the location ( $\mu$ ), scale ( $\alpha$ ) and shape ( $k$ ) were estimated from the regional data (i.e., from the theoretical relationship of the regional distributions and sample L-moments) (Table 4.10).
- Using these parameter values and the inverse functions of the best fitted distributions the regional quantile floods,  $X_T$  for the return periods,  $T=2, 5, 10, 20, 50, 100, 200$  and  $500$  years were computed (Table 4.10).
- Finally, the regional flood frequency curves for the best fitted distributions were constructed and plotted together with the normalized regional data (i.e., the normalized streamflows/ regional quantile  $X_T$  versus return period  $T$  (Fig.4.7).

**Regional L-moments-** the weighted averages of regional L-moments for the delineated regions of Southern Africa were estimated from the combination of weighted at-site L-moments. Once again, the values are estimated after the at-site flood events have been normalized by their respective median, i.e., by the index value. The first two L-moments and the first three L-moment ratios for each region considered in the study area are given in Table 4.9. It can be seen from Table 4.9 that the weighted expected values of the regional normalized floods are greater than one in all regions and reached 2.035 in region  $ZA\_R2$ . This



implies that the median value was smaller than the mean value in almost all sites i.e., all the sample series were positively skewed and in some of the regions (for example, in regions ZA\_R1, ZA\_R2), very large floods and also large variability between flood events were experienced in the sample series.

Table 4.9 Regional weighted average L-moments for the grouped regions of southern Africa

Regions	Regions L-moments				
	$\lambda_1$	$\lambda_2$	$L-Cv$	$L-Cs$	$L-kurt.$
<i>R1</i>	1.357	0.559	0.412	0.259	0.132
<i>R2</i>	1.587	0.817	0.487	0.280	0.117
<i>R3</i>	1.013	0.275	0.267	0.0367	0.114
<i>R4</i>	1.106	0.325	0.287	0.193	0.154
<i>ZA_R1</i>	1.642	0.868	0.516	0.391	0.211
<i>ZA_R2</i>	2.01	1.112	0.548	0.436	0.224
<i>ZA_R3</i>	1.476	0.660	0.437	0.374	0.234
<i>ZA_R4</i>	1.413	0.631	0.431	0.362	0.219
<i>ZA_R5</i>	1.358	0.549	0.392	0.280	0.123

The regional flood frequency curves were then developed for 9 of the regions considered in this study based on regional L-moments. Table 4.10 provides the summary of the construction which includes the best fitted regional distribution, their estimated regional parameter values, and regional quantile floods for the specified non-exceedance probabilities (return period ranged from 2-500years).

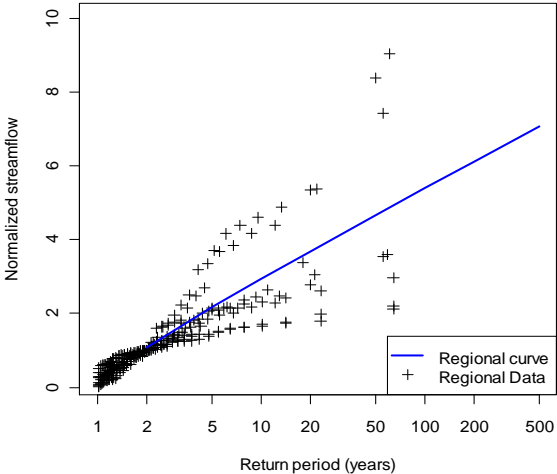
Table 4.10 Summary of the regional growth curves which includes the non-exceedance probability, the best distribution models for respective regions and their parameter values and quantile floods for a range of recurrence intervals.

Return period, $T$ (year)		2	5	10	20	50	100	200	500			
Non-exceedance probability, $F(x) = 1 - 1/T$		0.5	0.8	0.9	0.95	0.98	0.99	0.995	0.998			
		Parameter values			Dimensionless regional quantile floods, $X_T$							
Regions	Regional Distribution	$\mu$ ( $\xi^*$ )	$\alpha$ ( $\beta^*$ )	$K(\alpha^*)$	$X_2$	$X_5$	$X_{10}$	$X_{20}$	$X_{50}$	$X_{100}$	$X_{200}$	$X_{500}$
<i>R1</i>	PE3	-0.0128	0.831	1.648	1.09	2.17	2.94	3.69	4.67	5.40	6.12	7.07
<i>R2</i>	GPA	-0.149	1.952	0.124	1.15	2.70	3.76	4.73	5.90	6.69	7.43	8.30
<i>R3</i>	LN3	0.995	0.486	-0.0751	0.995	1.42	1.65	1.85	2.07	2.23	2.38	2.56
<i>R4</i>	GEV	0.828	0.453	-0.0365	0.995	1.53	1.89	2.25	2.73	3.10	3.47	3.99
<i>ZA_R1</i>	GPA	0.0129	1.425	-0.125	1.05	2.55	3.82	5.19	7.20	8.88	10.71	13.39
<i>ZA_R2</i>	PE3	0.565	3.308	0.167	1.11	3.25	5.09	7.04	9.70	11.77	13.87	16.67
<i>ZA_R3</i>	LN3	1.05	0.899	-0.793	1.06	2.14	3.06	4.10	5.70	7.10	8.68	11.06
<i>ZA_R4</i>	LN3	1.022	0.875	-0.767	1.02	2.06	2.93	3.91	5.39	6.68	8.11	10.26
<i>ZA_R5</i>	GPA	0.191	1.312	0.124	1.06	2.10	2.81	3.47	4.26	4.79	5.28	5.87

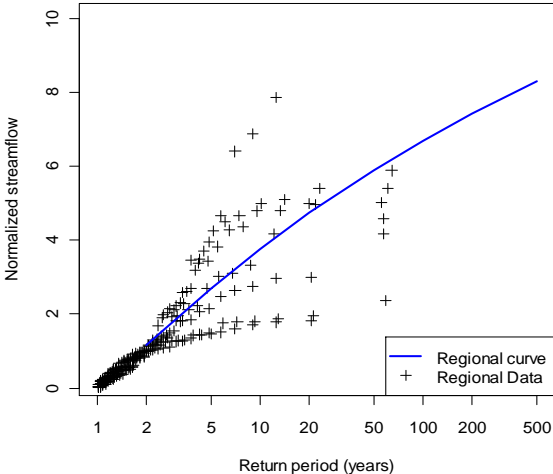
Note: The parameters  $\mu$ ,  $\alpha$  and  $k$  represent for location, scale and shape parameter values of a given distribution, respectively. In case of PE3, these parameters have another form which estimated primarily from other parameters (moments). The parameters are instead represents by the symbols in brackets with the asterisk(\*) i.e.,  $\xi$  (location),  $\beta$  (scale) and  $\alpha$  (shape) see brief discussions and relationships by Hosking and Wallis (1997),. These parameters are estimated directly from the relationship of the regional L-moment and some empirical coefficients (moments) such as the  $\mu$ -gamma mean,  $\sigma$ -gamma standard deviation and  $\gamma$  - the third moment as follows;

$$\text{If } \gamma \neq 0, \text{ let } \alpha = 4/\gamma^2, \quad \beta = 1/2\sigma|\gamma|, \text{ and } \xi = \mu - 2\sigma/\gamma$$

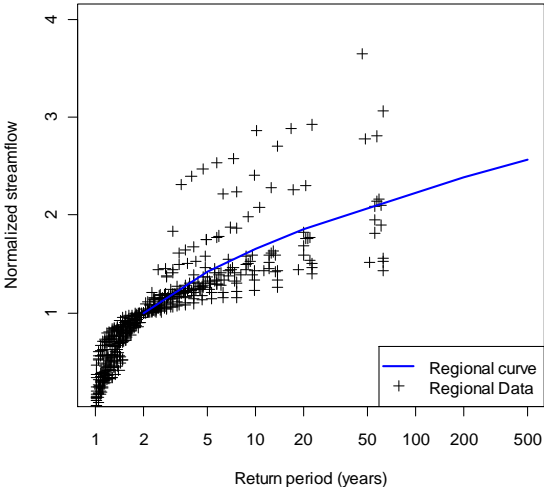
Once Table 4.10 has been developed, the quantile results together with the recurrence interval  $T$  were used to construct the regional frequency curves. Hence, the derived frequency curves were plotted together with regional normalized data and given in Fig. 4.8. The curves illustrate the regional relationships between the flood magnitudes and recurrence intervals. Though the difference between the observed and growth curves increased with recurrence interval, the curves generally reflect good consistence with trend of observed regional data. It can also be seen that the difference is high in regions which have had statistically heterogonous catchments (for example, (a) Namibia catchments) and all curves underestimated regionally large observations.



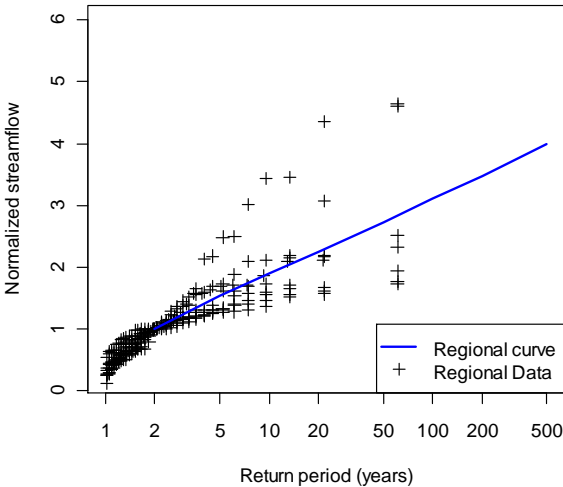
a) R1-Namibia



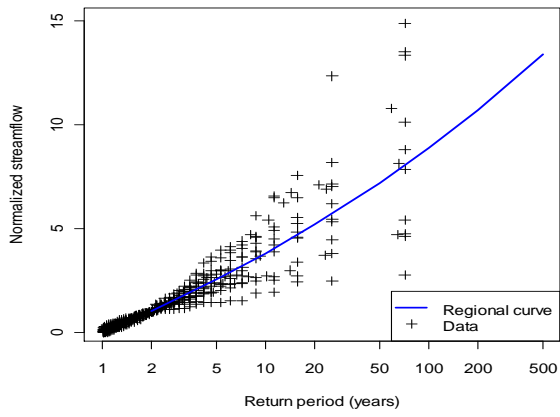
b) R2-Zimbabwe



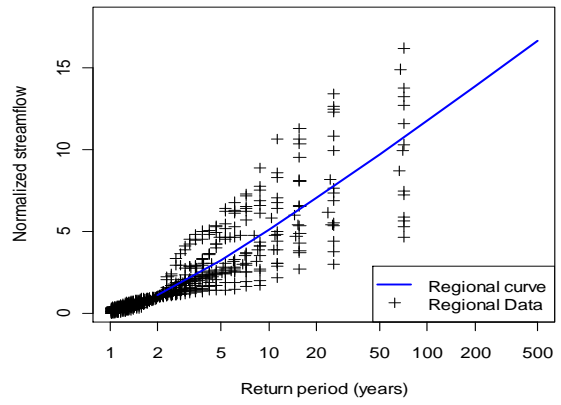
c) R3-Zambia



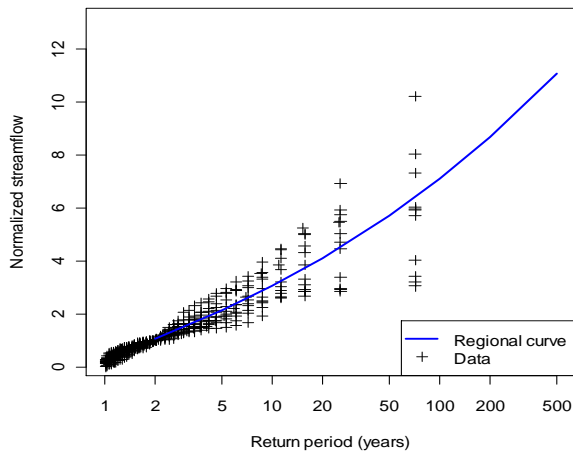
d) R4-Malawi



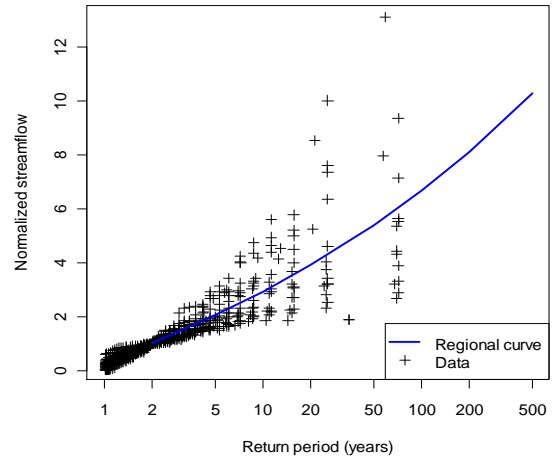
e) ZA\_R1



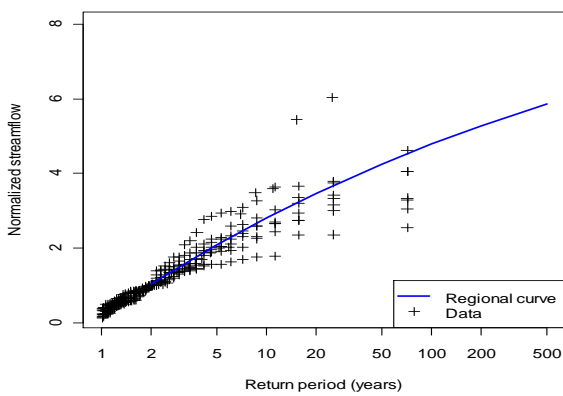
f) ZA\_R2



g) ZA\_R3



h) ZA\_R4



j) ZA\_R5

Figure 4.8 Regional flood frequency curves for 9 regions in Southern Africa: the title of each curve indicates the name of the regions. The curves were developed from best fitted distribution of respective regions in Table 4.10.

## 4.5 Performance evaluation using simulation

The performance of the chosen distribution as a best fitted regional model was assessed using two methods (1) Plotting-position i.e., quantile-quantile plot; and (2) Regional growth curve verification.

### 4.5.1 Quantile-Quantile (qq) plots

For 9 regions of the study, the quantile-quantile plots of the normalized observed floods versus the simulated values that generated from the best fitted regional frequency distributions were developed and furnished in Appendices C (Fig.1 (a-i)). Fig. 4.9 below also shows examples of the quantile-quantile plots for two regions. The figures in Appendices C illustrate that almost all the qq plots were well fitted to the line 1:1. This implies that the frequency distributions that were chosen as a best distribution could be an appropriate regional flood models for all southern Africa catchments. However, for example, Fig. 4.9 (a) shows that the regional flood frequency model for *RI*- Namibia underestimated the large quantiles of the region. Even though the model simulated values were in agreement with the prediction of small observations, the coordinators of the largest three q-q values deviated from the line 1:1. These situations happened in all quantile-quantile plots generated from all the best regional distributions for this study (Appendices C: Fig. 1)

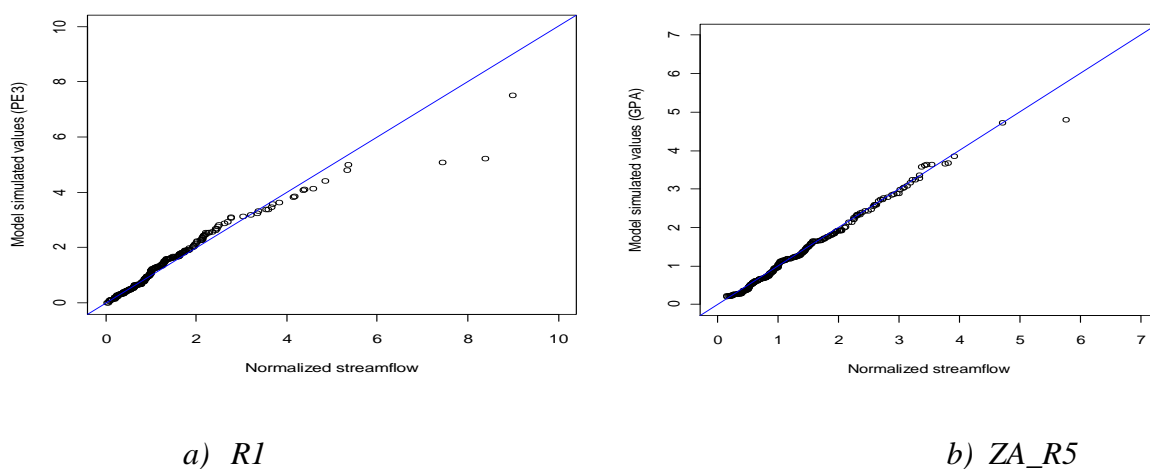


Figure 4.9 Examples of quantile-quantile plots of the normalized empirical discharge against the simulated values from the best fitted distributions: a) Pearson type III (PE3) for Region RI-Nambria; and b) Generalized Pareto distribution (GPA) for ZA\_R5

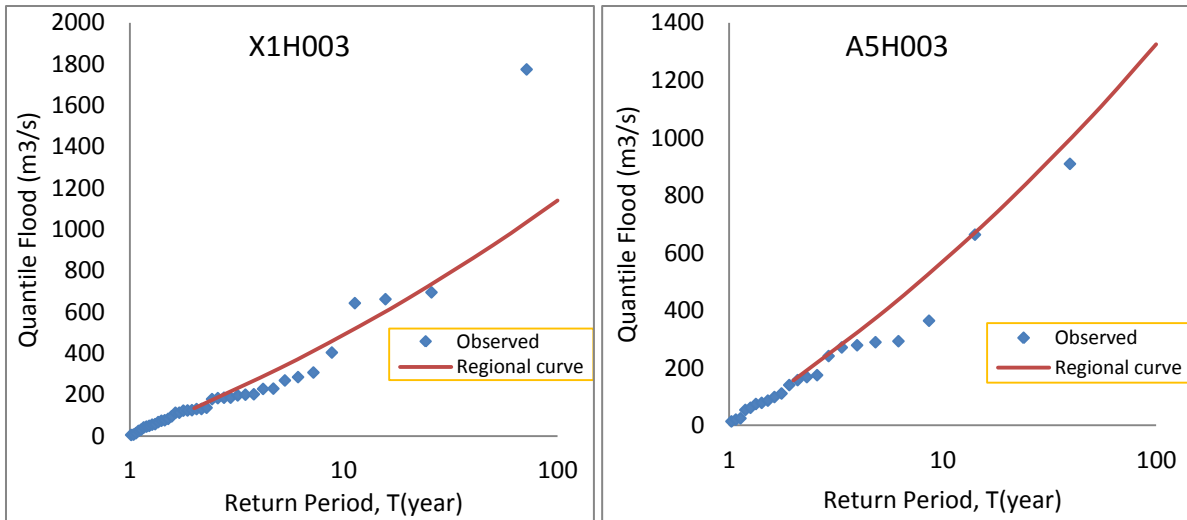
## 4.5.2 Verifications of the regional flood frequency curves

The ability of the selected regional frequency distribution was also evaluated by comparing the difference between at-site observed values and regional estimations. For the validation of the selected models, 10 stations were used. Although the aim was to verify all the regional curves, the numbers of the stations were not sufficient in all regions. Hence, the regional curves from five regions were verified using two stations from each region. Of these, four regions were from South Africa and the other was region *R3-Zambia* (Table 4.11).

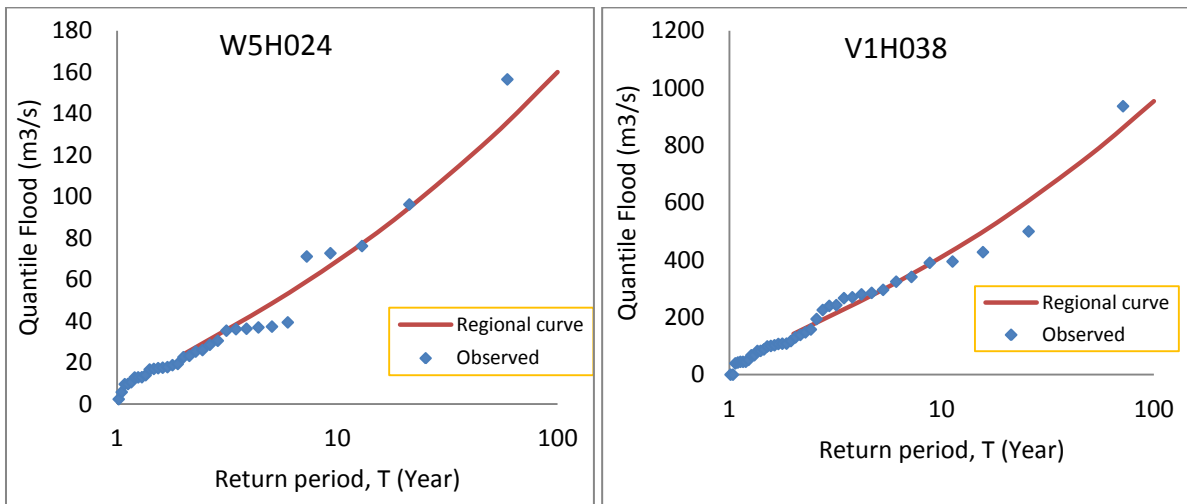
Table 4.11 Selected stations for model verifications and their index floods

Station code	River @ Station	Region	Index flood (m <sup>3</sup> /s)
1591404	Kafue@Kafue Hook Bridge	<i>R3</i>	881.0
1591474	Kafue@ Kafironda		194.0
A5h003	Limpopo River @ Botswana	<i>ZA_R1</i>	149.18
<u>X1h003</u>	Komati River @ Tonga		128.44
<u>C2h008</u>	Vaal River @ Woodlands	<i>ZA_R2</i>	172.9
D8h008	Orange River @ Pella Mission		239.3
<u>V1H038</u>	Klip River @ Ladysmithdorpsgronde	<i>ZA_R3</i>	134.34
<u>W5H024</u>	Mpuluzi River @ Dumbarton		22.35
<u>T3h007</u>	Mzimvubu River @ Ku-Makhola	<i>ZA_R4</i>	166.88
<u>U2h012</u>	Sterk River @ Grootboek		14.44

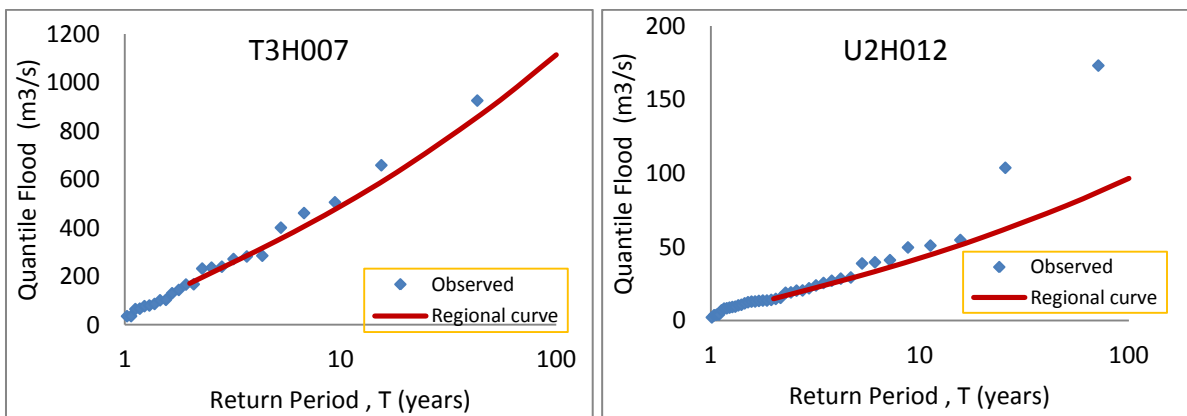
The regional curves were constructed only from the best fitted distribution, i.e., mainly for *R3-LN3*, *ZA\_R1-GPA*, *ZA\_R2-PE3*, *ZA\_R3-LN3* and *ZA\_R4-LN3* distributions. Fig. 4.10 (a-e) show the  $Q_T - T$  relationship for each of the 10 stations collected from five of these regions and calculated through equation 3.25. Due to the fact that the series had records with the range from 22 - 40 years, the comparison was done for the quantile floods up to 100 years return periods. For all regions, the estimated  $Q_T - T$  relationships are in good agreement with the observed flood events. However, except for ‘*A5H003*’, the chosen distributions underestimated the largest observation of the at-site sample series. This shows that the curves have consistence with the suggestions from Fig.4.8 that all the flood models underestimated the largest observations of the sample series.



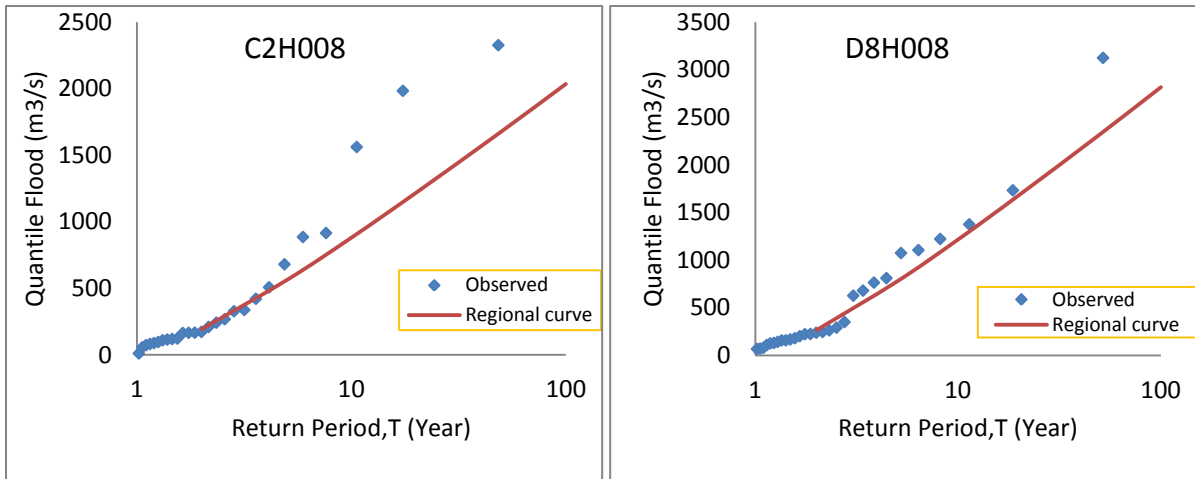
a) Region ZA\_R1



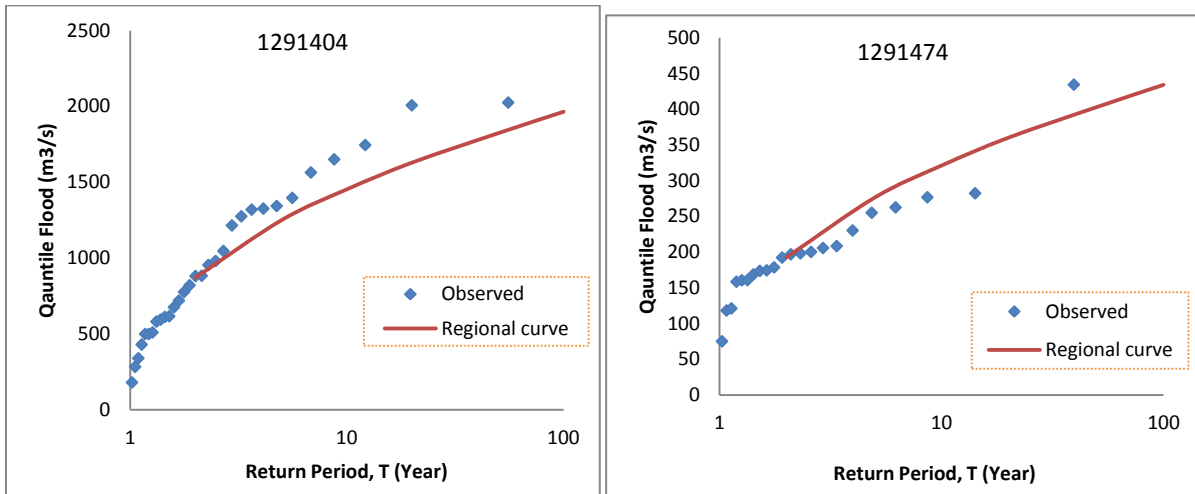
b) Region ZA\_R3



c) Region ZA\_R4



d) Region ZA\_R2



e) Region R3-Zambia

Figure 4.10 shows the comparison of the probability plots of the quantile floods between the observed series (Dotted) and estimated values from the best fitted of regional frequency curves (solid line)



## 4.6 Regional estimation for ungauged catchments

The regional estimations for ungauged catchments were derived from the relationships of the median and catchment area of the neighboring gauged catchments. The relationships between the median (the index flood) and the area of the catchments of the region are presented in Appendices D (Fig.2) and also summarized in Table 4.12. The regression equations were developed after an attempt had been made to fit the regional data to both exponential and logarithmic equations. As a result, a model which gives best  $R^2$  value was chosen as regional regression models to predict the median values of the ungauged /low record series catchments.

Table 4.12 Derived regression models to predict the median values from catchment characteristics in Southern Africa.

Regions	Regression model , where, $\tilde{Q}$ = index flood-median ( $m^3/s$ ) , $A$ = catchment area ( $km^2$ )	$R^2$
<i>R1</i>	$\tilde{Q} = 0.0371(A)^{0.78}$	0.5117
<i>R2</i>	$\tilde{Q} = 0.1265(A)^{0.8681}$	0.8696
<i>R3</i>	$\tilde{Q} = 0.4249(A)^{0.6662}$	0.9381
<i>R4</i>	$\tilde{Q} = 89.786*\ln(A)-409.02$	0.7022
<i>ZA_R1</i>	$\tilde{Q} = 14.755*\ln(A)-49.338$	0.3664
<i>ZA_R2</i>	$\tilde{Q} = 52.664*\ln(A)-340.28$	0.7683
<i>ZA_R3</i>	$\tilde{Q} = 66.461*\ln(A)-395.91$	0.5218
<i>ZA_R4</i>	$\tilde{Q} = 0.6089(A)^{0.6639}$	0.5927
<i>ZA_R5</i>	$\tilde{Q} = 42.282*\ln(A) -187.1$	0.889

## 5. DISCUSSION

### 5.1 Data and outlier analysis

In arid and semi-arid zones such as southern Africa, rivers which may have no flow for periods of time and sometimes extend to a number of years. Floods of such region raise mainly from intense convective thunderstorms of very limited areal extent and thus affect catchments randomly with little spatial pattern (Cunnane, 1989).

For this study, the daily average streamflow data were collected from five countries (459 stations) in Southern Africa. The departure of meaningful regional flood frequency analysis is that the data available should be continuous, long recorded period, non-regulated, independently, identical distribution (Cunnane, 1989). When the AM series were selected from 459 stations of parent distributions, more than 50% of the gauging sites reflected insignificant magnitude of observation series (AM series of zero and nearly zero values for most of the recorded period), with long gaps of information, and short record length (< 15 years). This illustrated that the AM series revealed that they were poor in quality and shouldn't be used in further analysis. When the rest of the AM flood samples were examined for dependency (i.e., the dependency between consecutive time series and across gauging sites) (section 4.1.2), a total of 68 stations had AM observations with strong cross-correlation coefficients. It was expected that due to inadequate stations collected from every corner of the countries, the spatial correlation of the AMS from Malawi, Namibia, Zimbabwe and Zambia stations might be insignificant, thus their correlation coefficients were < 0.4. However, the pair correlation coefficients of South Africa catchments showed good correlations and these 64 stations were from this country with correlation coefficient > 0.89 and the other 4 stations were from Zambia. Of these, 34 stations were used in the analysis while the others were excluded. This is in good agreement with conclusions of Mkhanda et al. (2000) that the AMS from South Africa had good correlation because of the dense network of gauging stations in the country. None of the AMS autocorrelation plots show significant correlation at the  $\alpha = 0.05$ , and thus it was accepted that the AM series were independently, identical distribution.

The curves from **empirical distribution** (for example, Fig.4.3) were also used to examine the frequency of the outliers (if available), shape of the curve and choice of theoretical distribution functions, etc. Fig.4.3 (a) shows that the distribution is bounded to the lower part of the curve (the series gives repeated values of  $x_l$  or nearly values) with high values to the

upper side (i.e., surprisingly the normalized values extended up to 200). This implies that the values from this particular river might not be natural, i.e., at somewhere upstream, the channel might be regulated. On the other way, Fig.4.3 (b) shows that the frequency of the observations is bounded to the upper side of the curve (the series gives repeated values of  $x_n$ ). It could also be the same reason that the upper limit of floods is controlled at somewhere upstream and/or the gauging instrument might not able to measure the floods, in case, above the maximum level.

It was also observed that the AMS of some stations were with suspected outliers. For example, Fig. 4.3 (c) (station 'U2H048') illustrates that the probability plot of the observations is with one outlier, i.e., the normalized values raised up to 24.27 times larger than the index flood. Whereas, Fig. 4.3 (d) shows the probability of the observed values after one outlier were removed from the series in (c). As shown in the plots, the largest normalized value was decreased almost by 80% i.e., the largest normalized value changes from 24.27 to 4.5.

The presence of such stations containing outliers, upper and lower bounded observations obviously affected the choice of a representative regional distribution. This was clearly observed when the observations were tried to fit to the theoretical frequency curves under L-moment ratio diagram (Fig. 4.7). As a result, AM floods from 73 stations were found experiencing not random phenomenon (i.e. upper or lower bounded series) and all were excluded from further analysis.

Following a thorough data screening, the AMS from 122 stations (112 for regional analysis and 10 for regional growth curve verification) were selected. The flood series of these stations were assumed that they are independent observations at the 95% confidence interval, independent across stations, from non-regulated rivers with record period > 15 years and few gaps of information. The gaps of these few runoff information were also treated as NA (NO data) values.

**Outliers-** the AMS at each site was closely examined whether the extreme events come from a single population (i.e., outlier analysis). It is suggested that the existence of outliers can be the reason for many of the problems raised in the regional analysis of hydrological data (Gottschalk and Kundzewicz, 1995). The significance of the outliers and their influence of the regional estimations were estimated based on detection techniques and then treatments (section 4.1.4 and 4.2.3). However, in consistent with the suggestion of Cunnane (1989),

every hypothesis was carried out at regional base. Here, the difficult task was investigating how many outliers of different degree of severity have actually occurred in the region dataset. As a result, a total of 29 outliers which were relatively more than 10 times higher than the at-site index flood and inconsistent with the rest of the frequency of the normalized regional data were selected as regional outliers (Table 4.4). The normalized values of the detected outliers were in the range between 4 and 27 and most of the outliers were recorded when there was high runoff throughout the study regions. Especially, the most extreme floods obtained from the sub-tropical hurricanes in South Africa (Kjeldsen et al., 2002), floods of 1974, 1976, 1988, 1996 and 2000 are among the periods which records large observation.

Though majority of the outliers were very large in magnitude compare to the at-site index flood, their frequency might be consistent with the theory of regional extreme values. However, there were three observations (Table 4.5) with normalized values extended from 21.95 to 26.37 and their recurrence intervals lies above 27778 years. These AM floods were recognized as unacceptable observations. The reason was not well-investigated but all were extremely deviated from the rest of the regional data. The values were excluded from the sample series and treated as NO data.

The treatment of the outliers was performed in agreements with the recommendations Cunnane (1989). The outliers were regionalized as shown in Table 4.4 and the influence of the suspected outliers on at-site and regional statistics was summarized in Table 4.6. The outliers were considered as random observation because most of the observations have common record time and were recorded during high floods throughout the region. The total numbers of outliers recognized as regionally large observations were with a range from 2-7 observations per region (too few to form parent distribution). In addition, the sensitivity analysis illustrated that the regional weighted average statistics was less sensitive to existence or removal of one or two outliers from the series. That is, after the suspected outliers were removed from respective stations, the at-site statistics relatively varied from 3.13% to 43.8% whereas, the regional statistics reduced from 3.4% ( $L-Cv$ ) to 15.93% ( $L-Kurt$ ). This is because the analysis adapted the PWMs methods of parameter estimation.

Therefore, except for three observations listed in Table 4.5, all outliers were accepted as random variables, and allowed to have in distribution selection and parameter estimation. This was in agreement with the recommendation of Cunnane (1989) that “if the AM floods come from two different sub-populations and are regarded as true observations, then the outliers

must be retained and treated as random and unbiased". It was also experienced by other studies in south Africa such as Kjeldsen et al. (2002).

## 5.2 Regional Homogeneity

Regional flood estimations methods are based on the premise that the normalized values of the floods have the same distribution at every site in the chosen region (i.e., the at-site statistics such as  $C_v$  and  $C_s$  of the normalized values are considered to be constant across the region) (Cunnane, 1989; Hosking and Wallis, 1997). In this work, the identification of statistically homogenous regions were based on delineating geographically homogenous regions and later evaluating for their heterogeneity using Hosking and Wallis (1997) homogeneity measures ( $H$ ).

The delineation of sites into statistically homogenous regions therefore adapted the context of geographically continues regions. This was because geographical regions are more continent if the region is used to assign ungauged and poorly gauged catchments, and may be easy to divided a country or the sturdy drainages by its variation of soil, climate and topography with latitude and longitude (Cunnane, 1989; Hosking and Wallis, 1997). But this doesn't mean that geographically neighboring catchments could necessarily fulfill the homogeneity assumption. Based on these assumptions, all the AMS collected from 112 stations were delineated into 9 regions based on hierarchical grouping procedures and summarized in section 4.2. At the initial stage, after all the stations failed the heterogeneity measure ( $H$ ), all gauging sites were grouped into belonging countries. Hence, the AMS from Namibia, Malawi, Zimbabwe and Zambia were grouped into four regions which represent every country as one region (with  $H$  values 7.5, 5.43, 2.35 and 2.31, respectively). The catchments from Zambia and Zimbabwe were considered as moderately heterogeneous; and the Namibia and Malawi as definitely heterogeneous regions. The regional coefficient of variances of these regions was in the range 0.287-0.487, which all met the recommendations suggested by Cunnane (1989). Cunnane (1989) stated that in regions relatively low  $C_v$  ( $< 0.6$ ), a small degree of heterogeneity doesn't cancel out the benefit of using an at-site/regional estimation method. Hence, all AMS from these countries were considered as statistically accepted regions. Meanwhile, the catchments of South Africa were more than 75 % of the total stations provided for this analysis. The sites were therefore grouped into five regions (see for details Table 4.3) slightly in agreement with

the delineations of Mkhandi and Kachroo (1997). The AMS of the sample sires from all these regions were with heterogeneity measure  $H$  ranged between 1 and 2 (i.e., possibly heterogeneous regions). Though the stations in South Africa were classified into more regions, they still have not formed exactly homogeneous regions. This could be due to the following possibilities: naturally the regions can have slightly different climate and geographical locations, hydrological regimes (such as storage and water table levels), insufficient stations and datasets, highly variance with the flood events etc. The regions were accepted as valid regions in agreement with recommendations by Hosking and Wallis (1997) and Cunnane (1989) that the heterogeneity measure of the regions could be realistically representative of the complex system if the regions are with slight heterogeneous.

However, the sizes of the stations in all regions were considered inadequate for any meaningful regionalization (especially in Malawi, Zambia, Namibia and Zimbabwe catchments, and even in South Africa catchments). In addition to the inadequate data available and short record periods, the most extreme events in the sample series resulted from infrequent sub-tropical cyclones (Kjeldsen et al., 2002) were another main problem during homogenization of the regions. Especially in South Africa catchments, floods of 2000, 1996, and 1976 and in some regions floods of 1974 were among the main floods that make difference during the selection of homogenous regions. In agreement with the recommendations of Cunnane (1989) on treatments of outliers in AMS, the regions that are listed in Table 4.3 were therefore developed, after the suspected outliers were excluded from the sample series.

### 5.3 Regional flood frequency distribution

The identification of an appropriate regional flood frequency distribution for each of the homogenized regions was based on the L-moment diagram and Goodness-of-fit tests. The L-moment diagram indicates that the GPA could be the most suitable regional distribution for Namibia (*R1*), Zimbabwe (*R2*) and all South African catchments. It also suggests that the PE3 distribution could be stochastically an alternative flood frequency model for AM series of catchments in these regions and the LN3 distribution for *ZA\_R3* and *ZA\_R4* catchments. However, it was not easy to distinguish the most likely regional distribution for Zambia and Malawi catchments as the GEV, LN3 or PE3 distributions could be preferably an appropriate regional distribution.

For every region, the GOF tests, however, suggested the acceptable regional distribution and their rank of performances. From the tests analysis, the PE3 distribution was found to be an appropriate regional distribution for all regions of Southern Africa. It was the best regional frequency distribution only for two regions such as Namibia and *ZA\_R1* catchments and accepted regional frequency distribution for other 6 regions (Table 4.8). Whereas, the Andersen-Darling goodness-of-fit test (Table 4.7) recommends that the PE3 distribution could be a valid regional flood frequency model for all catchments in Southern. The performance of PE3, followed by GPA (for 7 regions) and then by LN3 (for 6 regions) distributions, while the GEV distribution recommended in recent studies performed as a fourth place (i.e., only best fitted to Malawi catchments and acceptable distribution to other three regions). All the candidate distributions which were selected from the L-moment diagram were confirmed by the GOF tests. Hence, the results summarized in Table 4.8 were therefore recognized as a test that described the performance of the candidate distributions very well. As a result, this study suggested that the GPA-for regions Zimbabwe and two regions of South Africa (*ZA\_R1* and *ZA\_R5*), and LN3 for regions *ZA\_R3* and *ZA\_R4* as the most suitable regional flood frequency distribution. For region *R4* –Malawi flood events, the GEV distribution was chosen as a best fitted regional flood model. The possible regional distributions for every region considered in this study were ranked in Table 4.8. However, from both L-moment diagram and the test statistics, the GLO and the two-parameter distributions (i.e., EV1 and EXP) performed very poor i.e., none of these distributions could be used for flood modeling in southern Africa.

This study also suggested that the statistical flood frequency for every country carried in this study could be characterized by PE3 for Namibia, GPA for Zimbabwe, LN3 for Zambia, GEV for Malawi and the GPA for South Africa flood events- using PWMs method of parameter estimators. In fact, these regions are characterized by relatively different physiographic and climatic conditions. This implies that, the flood situations of every catchment should be modeled based on the best regional flood frequency distribution. Though the empirical models are with the expected uncertainties (i.e., at 90 % confidence interval), these distributions might reasonably represent the flood phenomenon in the regionalized catchments.

In Southern Africa, few flood frequency studies, such as the RFFA for southern Africa carried by Mkhanda and Kochroo (1997); and Mkhanda et al. (2000) have concluded that the best flood frequency distribution procedures for the region could be the LP3/MOM and/or P3/PWM. In contrast to this conclusion, the RFFA in south Africa by Kjeldsen et al., (2002) concluded that the suitable regional distribution based on the L-moments diagram are the GPA and GLO distributions. The literature by Kjeldsen et al.(2002) also discussed that, different studies in South Africa particularly in Kwazulu-natal province suggested different distribution function. This indicated that despite the generalization of flood frequency distribution at regional (Southern Africa) level, the scaling down frequency analysis could give different results and probably better estimations. This could mainly depend on the data inputs.

However, though the delineation of the regions was slightly different, the summary of the frequency distribution models were generally in agreement with the conclusions by Mkhanda and Kachroo (1997) i.e., the three-parameters of PE3 with PWMs could be regional distribution for all flood events of Southern Africa. This study also supports the flood studies by (Kachroo et al., 2000; Kjeldsen et al., 2002; Mkhanda and Kachroo, 1997) that the three parameter distributions are more capable of fitting flood data in southern Africa. The two-parameter distributions (i.e., EV1 and EXP) performed very poor i.e., any of the distributions shouldn't be considered for regional frequency analysis for all Southern Africa floods. Though the regions (for Kwazulu-natal province catchments-South Africa) were slightly modified the identification of the regional distribution from this study was also slightly in agreement with the finding by Kjeldsen et al. (2002). The study area in this work includes part of ZA\_R3 and ZA\_R4 regions, but in both regions the regional flood data were fitted to LN3 distribution. The GEV and GPA distributions should also be considered as the flood events



models of the region. The results are in good agreement because the GPA distribution was suggested as regional distribution in both studies. Based on the GOf test, the LN3 was also chosen as best distribution in half of the regions of the previous study ( i.e., strongly in agreement with this study) , but, in this study, the choice of the GLO distribution was not considered at all.

#### **5.4 Regional flood frequency curve**

The regional frequency curves are simulated values which developed to represent the average values of the at-site flood frequencies. A regional curve is essentially a frequency distribution of  $Q/\tilde{Q}$ . It associates a return period  $T$  with  $Q/\tilde{Q}$  and this relationship is assumed to be valid for all catchments in the region, or alternatively represent the mean of the different relationships for the different catchments in the region (NERC, 1975). Hence, the curves (Fig.4.8) represent for all regions reflect good agreement with regional normalized flood events. Despite the large deviation of the regionally largest observations from the curve with increasing return period, the agreement is relatively good in flood events of homogenous regions, for example, from South Africa regions ( $ZA\_R1$ ,  $ZA\_R2$ ,  $ZA\_R3$ ,  $ZA\_R4$  and  $ZA\_R5$ ). However, because the flood statistics from the first four regions ( $R1$ - $R4$ ) were ranged from moderately heterogeneous ( $R1$  and  $R4$ ) to strongly heterogenous ( $R2$  and  $R3$ ) (see section 4.2); the curves were somewhat disagree with the observed samples, especially with increase return period.

It can also be seen that none of the regional frequency curves were able to represent the regional large observations. That is, as discussed in section 5.1, there were some regional observations which were large relative to the index flood, but since they are few in numbers, the values were included in the analysis. This supports the suggestions by Cunnane (1989) that many of the existing statistical flood frequency estimation models underestimated the frequency of very large floods. Hence, all the regional curves (Fig. 4.8) developed in this analysis reflect underestimation for the quantiles of large observations in the region.

The regional frequency curves (constructed from 9 regions) reflect that all curves have different flood characteristics. This could be due to the fact that the flood in different regions has different flood statistics. However, and it is also obvious that the degree of the flood regime variation depends on the meteorological and catchment characteristics that generating

the flood events (Hosking and Wallis, 1993; Kachroo et al., 2000; NERC, 1975; Tveito, 1993). Despite the effort of homogenization of the stations, examination of the regional curves also showed that some regions have almost the same regional frequency pattern. These regions were in South Africa especially regions *ZA\_R3* and *ZA\_R4*. The AM floods in the catchments of both regions could be modeled using LN3 distribution and the gradient of the regional curves also showed almost the same pattern. However, examination of the regional homogeneity showed that both regions have different flood statistics and this might happen due to the acceptance of outliers in both regions.

The gradient of the regional frequency curve describes the probability of extreme floods. i.e., a curve that shows large gradient reflects large variability of the regional floods or there are floods which occurred rarely and vice versa for the gently sloped curves. Therefore, as shown in Fig. 4.8, the regional curves derived from AM floods of Zambia, Malawi and one region of South Africa (*ZA\_R5*) catchments were gently sloped. The remaining 6 regions ranged from moderate (*R1* and *R2*) to steeper (*ZA\_R3* and *ZA\_R4*), and the steepest slopes were in *ZA\_R1* and *ZA\_R2* regions. Especially, the slope of the curve for floods *ZA\_R2* shows the steepest curve of the study area. This can be observed, for example, the 500 years flood of this region reaches up to 17 times of the index flood. This implies that these regions which have steeply sloped curves were established from highly variable floods.

This can also be confirmed by the results from the weighted average regional L-moment ratios (i.e., regional coefficient of variance, Table 4.9). The regions which show small values of *L-Cv* have gentle slopes and the large *L-Cv* values have high gradient. For instance, the steep flood frequency curves such as *ZA\_R1* ( $\tau_2 = 0.516$ ) and *ZA\_R2* ( $\tau_2 = 0.548$ ) were the curves among the regions which reflect a high variation of the flood regimes and included large observation in the stations of the entire regions. They are regions which represent the inland zones of the South Africa catchments. Thus, they have high and erratic rainfall, large area and most of the time dry climate (see Fig. 2.4 and 2.5) especially *ZA\_R2*. As Table 4.4 shows, region *ZA\_R2* includes the semi-arid and semi-desert drainages of South Africa such as the Limpopo (*A1-A4*), Olifats (*C*) and the desert area of Orange (*D3-D8*) drainages.

The regional growth curves of the coastal zones which includes *ZA\_R4* ( $\tau_2 = 0.437$ ) and *ZA\_R5* ( $\tau_2 = 0.392$ ) and partially *ZA\_R3* ( $\tau_2 = 0.431$ ) (see Table 4.3 and Fig.4.9 for the details of the drainages) of South Africa has small gradient which indicates the little variability in the flood size and/or frequent of floods in the entire region. In these regions, the

catchments are relatively small, but supplies relatively high floods than the other catchments in South Africa. This might be due to the fact that coastal areas may have more frequent precipitation.

The regional growth curves of regions *R4-Malwi* and *R3-Zambia* are among the gentlest slopes and the regional coefficients of variance are 0.287 and 0.267, respectively. These might indicate that the tropical humid climate and the drainage characteristics could be the main factors that create such small variation of the flood events. This low relative steepness may be explained by the relatively low variability in rainfall in the region. The flood study report (NERC, 1975) discussed that the effect of catchment characteristics such as large area coverage between stations, high ground water levels, and low soil storages can produce gentle regional frequency curves. The *R1-Nambia* and *R2-Zimbawie* are relatively steep as compare to the Malawi and Zambia frequency curves. The reason may be due to the heterogeneity in the flood data used in this analysis. The AM floods were pooled from different parties of the countries and were also inadequate in numbers (i.e., 8 and 7 stations, respectively). In addition, since some of the regions of Namibia were collected from semi-desert dry climate of the coastal areas (see Fig. 2.2 and 2.5), the rainfall from these regions could have relatively high variability as compare to coastal zones of the other regions (Southern Africa).

## **5.5 Performance evaluation of empirical distributions**

There are several components that contribute to the errors in regional quantile estimation such as errors arising from estimation of the index flood, misspecification of the regional flood frequency distribution, heterogeneity in the region, number of sites in a region, record length and outliers at each sites and other unexpected uncertainties (Hosking and Wallis, 1997).

Sections 4.5.1 and 4.5.2 assessed the performance of the chosen regionally best fitted distributions. The quantile-quantile plots - the normalized values of observed and simulated were plotted to see the ability of selected model that were chosen based on the observed regional data. As shown in Appendix C (Fig. 1), the diagrams indicated that the simulated values were in good agreement with the normalized observed, i.e., the predictive ability of the selected regional models were very good except for the largest observations. However, all the regional models underestimated the largest observed events i.e., the simulated versus

observed values were located below the line 1:1. This was observed in all quantile-quantile plots shown in Appendices C (Fig.1) and agreed with the regional flood frequency curves.

**Model Validation-** the comparison of the regional flood frequency curve against the at-site observed flood quantiles. In consistent with review by Kjeldsen et al. (2002), this was not used to discriminate the regional distribution curves rather to observe if there was any systematic regional bias in the estimation of the quantile events.

Fig.4.10 presents the comparison diagrams of the at-site quantile floods of 10 stations from five regions, i.e., the at-site quantile floods derived from the regional curves vs. at-site quantile observed floods. In all regions, the diagrams reflect that the regional curves were reasonably in agreement with the quantiles of observed values. Since the regional curves were derived from the average values of slightly different samples in the region, it is expected to observe some discrepancies in between the quantiles. However, high differences were observed in regions which were not homogenous like *R3-Zambia* (Fig.4.10 (*e*)) and in regions which contain highly variable flood events such as *ZA\_R2*. Fig.4.10 (*e*), station '1291404' shows clear deviations between the at-site observed and simulated quantiles. Though both the curves were followed the same curve pattern, the diagrams reflect that the relative difference between the quantile increases with return period.

The choice of at least two stations per region might help to see whether the best regional models underestimated or overestimated the frequency of the largest observations. Hence, the results from the diagrams show that all the models underestimated the large observations in a site. Therefore, it can be again concluded that from both the regional curve verifications and quantile-quantile plots, the efficiency of the chosen flood models in simulating largest quantiles is not that much worthy. This might be due to the facts that the regional data is heterogeneous, the acceptance of some existed outliers in the region and inadequate stations for the analysis.

## 5.6 Estimation of design floods from ungauged catchments

From Table 4.13 and Appendices D (Fig.2), it can be observed that the development of regression models for estimation of the index flood gives reasonably good fit. The correlation coefficient ( $R^2$ ) obtained for each of the derived regression model was above 50% except for

region *ZA\_RI* ( $R^2=0.3664$ ) of southern Africa. However, the efficiency of the models might be considered that the  $R^2$  values were developed from few stations. This might be largely attributed to the limited available information and disproportion of catchment size. For example, in region *ZA\_RI*, for station 'A5H006', the median value that was estimated from the catchment area of 98,240 km<sup>2</sup> is 115m<sup>3</sup>/s which is relatively very small.

When the result from this work were compared with previous studies by Mkhandi and Kachroo (1997), the regression equations from this study provide better  $R^2$  values. For example, the regression equations for Zambia, Malawi, Namibia and Zimbabwe suggested by Mkhandi and Kachroo (1997), were developed from multiple regression and the index flood was the mean of AMS with  $R^2$  value 0.78, 0.234, 0.265 and 0.732 respectively. Meanwhile, the regression values from this work are with  $R^2$  value 0.938, 0.702, 0.512 and 0.869, respectively.

However, because the floods of given catchments are not only characterized by the catchment area, the uncertainty that generates from these regression models might be very high, especially on the logarithmic equations. Therefore, by taking this into consideration, the results presented in Table 4.13 or Appendices D (Fig.2 (a-i)) might be served as optional models in estimating the index flood for ungauged catchments. However, this study suggested that the index flood models should be established as function of several catchment characteristics and incorporating all the necessary catchment information.

## 6. CONCLUSION AND RECOMMENDATIONS

The main aim of the RFFA conducted in this study was to extract sufficient information from rarely available at-site flood events, i.e., to adequately estimate the frequency of these flood events. These are certainly cases for the extreme events which are of interest in hydraulic structure safety, emergencies, human and resource risk managements and other resources utilizations. Coping with floods in an efficient manner and reducing damages necessitate efficient methods for the estimation of design floods; and later developments of flood management plans such as flood inundation zones and sufficient methods of flood forecasting techniques.

The data collected for this analysis was daily average runoff from 459 stations of five countries. After the AMS data were selected and subjected for preliminary data analysis, AMS data from 122 (112 for RFFA and 10 stations for model verification) were screened. The AM flood events generated from the daily average were positively skewed, highly variable and includes large observations (outliers) in the series. However, after the outliers have been detected and treated, except for three observations, all the AMS collected from the above stations were considered as random observations and used in the regional flood frequency analyses.

Using the site characteristics mainly the continuity of the geographical locations of the catchments, the gauging sites were grouped into 9 regions. Four countries which include stations of Zambia-*R3*, Namibia-*R1*, Malawi-*R4* and *R2*-Zimbabwe were considered each as one region. However, due to existence of more stations, the South African gauging sites were grouped in to five regions: namely the *ZA\_R1*, *ZA\_R2*, *ZA\_R3*, *ZA\_R4*, and *ZA\_R5*. When these regions were examined with respect to heterogeneity measures, the AMF obtained from the first four regions/countries failed the homogeneity test. Regions *R1* and *R4*- definitely heterogeneous while the other two regions of *R2* and *R3* -moderately heterogeneous and all the regions of South Africa are grouped as possibly heterogeneous regions.

The identification is performed using the hierarchical grouping method based on the geographical locations together with the Hosking and Wallis (1997) heterogeneity test. However, the study believed that the available information were not sufficient to form meaningful homogenous regions. It is therefore recommended that further studies should

investigate appropriate techniques, and consider number of sites and all necessary site characteristics for effective delineation of homogenous regions.

The identification of an appropriate regional distribution for each of the 9 regions is done based on the L-moment approaches: the L-moment ratio diagrams and GOF test of the L-moments and later evaluated their performances by their probability plots. For this study seven types of theoretical distribution (such as the EV1, EXP, LN3, GEV, GPA, PE3 and GOL) have been employed. The results obtained from this analysis are generally in agreement with the suggestions of previous study of Southern Africa by Kachroo et al.(2000). Both the L-moment diagrams and the GOF tests indicated that three parameter distributions are more capable of modeling the AM flood events in southern Africa. That is, the most suitable distributions for all regions of Southern Africa could be the PE3 and/or GPA distributions with PWMs method of parameter estimations. The GOF tests statistics at 10% level of significance, however, recommend the possible regional distributions and their prediction abilities (ranks). As a result, the PE3 distribution is best fitted to the regional AM floods of *R1-Nambia* and *ZA\_R2* catchment; and GPA is best fitted empirical distribution to three regions: *R2-zimbabwe*, and two regions of southern Africa namely regions *ZA\_R1* and *ZA\_R5*. For regions such as: *R3*, *ZA\_R3* and *ZA\_R4* the three-parameter log-normal distribution provides best regional frequency distribution and the GEV distribution was chosen as the best fitted to the AM floods of Malawi catchments. However, none of the EV1, EXP and GOL flood frequency models was able to model the flood events of the Southern Africa catchments.

Therefore, based on this result, it can be concluded that the PE3, GPA, LN3 and GEV emerged as underlying regional distributions, but none of the other three distributions could be considered as regional flood models in southern Africa. From the present study a particular flood frequency model for each country's flood events were suggested that the PE3 for Namibia, GPA for Zimbabwe, and South Africa, LN3 for Zambia and GEV for Malawi flood events could be an appropriate empirical distributions.

For every region, from the relations of the best regional distribution and at-site AM floods, the regional flood frequency curve for the return periods  $T$  of 2, 5, 10, 20, 50,100, 200 and 500 years has been developed. Since some of the regional AM floods are large even compare to the normalized regional data, the curves revealed that the suggested best regional flood models underestimated the magnitudes of these extreme events. It is known that the slope of

the regional frequency curves represents variability of the flood events in the region. Hence, the slopes of the growth curves derived in this study reflect good correspondence with the location and climate phenomenon of the regions /catchments in the study region. Thus, the curves are relatively categorized as gentle for regions *R3*, *R4*, and *ZA\_R5* and moderate for regions *R2*, *R3* and the coastal areas of south Africa i.e., regions *ZA\_R3* and *ZA\_R4*, while, the arid and semi- arid regions and desert areas of south Africa ( *ZA\_R1* and *ZA\_R2*) as steeply sloped curves.

The performance of the chosen distributions as the best regional flood models and the curves that developed from these distributions were assessed by employing the quantile-quantile plots and the model verification techniques. The results from both methods displayed that the flood frequencies of the regions are well addressed by the chosen distributions except for large observation (outliers). However, the relative difference between the observed and simulated quantile floods increased with return period. In addition, as concluded in the above paragraph, both qq plots and the curve verifications suggested that the best fitted regional distributions are not able to estimate the largest regional observations.

Finally, an attempt has been made to develop a regression regional models that can able to estimate the quantile floods of the ungauged catchments from the regional relationships of the gauged catchment characteristics. The results of the regression have shown that there is good correlation between the areas and the median values of the AM series of southern Africa catchments. For each of the recommended regions, the regression models are given with their corresponding  $R^2$  values so that one can judge the quality of the data and the likely uncertainty of the estimation using the regression models. However, since the regression is done only based on the index flood-catchment area relationships; the study recommends as an optional use of these relationships and if more data are available to carry out further analysis.



## 7. LIST OF REFERENCES

- Behrens, J. T. (1997). Principles and procedures of exploratory data analysis. *Psychological Methods*, 2(2), 131-160.
- Burn, D. H. (1988). Delineation of groups for regional flood frequency-analysis. *Journal of Hydrology*, 104(1-4), 345-361.
- Chebana, F., & Ouarda, T. (2008). Depth and homogeneity in regional flood frequency analysis. *Water Resources Research*, 44(11), W11422
- Chebana, F., & Ouarda, T. (2009). Index flood-based multivariate regional frequency analysis. *Water Resources Research*, 45, W10435.
- Cunnane, C. (1988). Methods and merits of regional flood frequency-analysis. *Journal of Hydrology*, 100(1-3), 269-290.
- Cunnane, C. (1989). *Statistical distribution for flood frequency analysis*. World Meteorological Organization, Operational Hydrology Report No.33, WMO-NO.718. Geneva, Switzerland.
- Dalrymple, T. (1960). *Flood-frequency analysis*. U.S. Geological Survey Water Supply Paper 1543-A, Washington.
- Efron, B., & Gong, G. (1983). A Leisurely Look at the Bootstrap, the Jackknife, and Cross-Validation. *The American Statistician*, 37(1), 36-48.
- Ellouze, M., & Abida, H. (2008). Regional flood frequency analysis in Tunisia: Identification of regional distributions. *Water Resources Management*, 22(8), 943-957.
- Embrechts, P., Klüppelberg, C., & Mikosch, T. (1997). *Modelling extremal events for insurance and finance*. Berlin: Springer.
- Engeland, K. (2005). *A short introduction to extreme value theory*. Lecture notes in Stochastic Hydrology. Department of Geosciences, University of Oslo.
- Farquharson, F. A. K., Meigh, J. R., & Sutcliffe, J. V. (1992). Regional flood frequency-analysis in arid and semiarid areas. *Journal of Hydrology*, 138(3-4), 487-501.
- Glad, P. A. (2010). *Meteorological and hydrological conditions leading to severe regional drought in Malawi*. MSc., The University of Oslo, Oslo. Retrieved from <http://urn.nb.no/URN:NBN:no-26053>
- Gottschalk, L. (2005). Methods of analyzing variability. In M. G. Anderson (Ed.), *Encyclopedia of hydrological sciences* (Vol. 1, pp. 95-122). Chichester: John Wiley & Sons.
- Gottschalk, L., & Krasovskaia, I. (2001). *Regional flood frequency analysis, - a theoretical background*. Lecture notes in Stochastic Hydrology. Department of Geophysics. University of Oslo.
- Gottschalk, L., & Kundzewicz, Z. (1995). Analysis of outliers in Norwegian flood data. In Z. Kundzewicz (Ed.), *New Uncertainty Concepts in Hydrology and Water Resources* (pp. 245-251): Cambridge University Press.
- Gringorten, I. I. (1963). A plotting rule for extreme probability paper. *Journal of Geographical Research*, 68(3), 813-814.
- Gumbel, E. J. (1958). *Statistics of extremes*. New York: Columbia University Press.
- Haan, C. T. (2002). *Statistical methods in hydrology*. Ames, Iowa: Iowa State Press.
- Hampel, F. R. (1974). The Influence Curve and Its Role in Robust Estimation. *Journal of the American Statistical Association*, 69(346), 383-393.
- Hipei, K. W. (1994). *Stochastic and statistical methods in hydrology and environmental engineering*. Dordrecht: Kluwer Academic.

- Hosking, J. R. M. (1990). L-Moments: Analysis and Estimation of Distributions Using Linear Combinations of Order Statistics. *Journal of the Royal Statistical Society. Series B (Methodological)*, 52(1), 105-124.
- Hosking, J. R. M., & Wallis, J. R. (1993). Some statistics useful in regional frequency-analysis. *Water Resources Research*, 29(2), 271-281.
- Hosking, J. R. M., & Wallis, J. R. (1997). *Regional frequency analysis: an approach based on L-moments*. Cambridge: Cambridge University Press.
- Institution of Engineers, A. (1977). *Australian rainfall and runoff: flood analysis and design : S.I. metric, 1977*: Institution of Engineers, Australia.
- Kachroo, R. K., Mkhandi, S. H., & Parida, B. P. (2000). Flood frequency analysis of southern Africa: I. Delineation of homogeneous regions. *Hydrological Sciences Journal-Journal Des Sciences Hydrologiques*, 45(3), 437-447.
- Kjeldsen, T. R., Smithers, J. C., & Schulze, R. E. (2002). Regional flood frequency analysis in the KwaZulu-Natal province, South Africa, using the index-flood method. *Journal of Hydrology*, 255(1-4), 194-211.
- Kottegoda, N. T. (1984). Investigation of outliers in annual maximum flow series. *Journal of Hydrology*, 72(1-2), 105-137.
- KRAK. (2011). The Kunene River Awareness Kit (RAK) , Southern Africa. *The River Basin* Retrieved 26/07/2011, from <http://www.kunenerak.org/en/river.aspx>
- Lang, M., Ouarda, T. B. M. J., & Bobée, B. (1999). Towards operational guidelines for over-threshold modeling. *Journal of Hydrology*, 225(3-4), 103-117.
- Mkhandi, S. H., & Kachroo, R. K. (1997). *Regional flood frequency analysis for Southren Africa*. Southern Africa FRIEND, Technical Documents in Hydrology No.15 UNESCO, Paris, France.
- Mkhandi, S. H., Kachroo, R. K., & Gunasekara, T. A. G. (2000). Flood frequency analysis of southern Africa: II. Identification of regional distributions. *Hydrological Sciences Journal-Journal Des Sciences Hydrologiques*, 45(3), 449-464.
- NERC. (1975). Flood studies report. *Hydrological studies* (Vol. I). London, UK.
- Noto, L. V., & La Loggia, G. (2009). Use of L-moments approach for regional flood frequency analysis in Sicily, Italy. *Water Resources Management*, 23(11), 2207-2229.
- NUFU. (2010). NUFU-Water Sciences Project. *Theme 3: Water Rsources and Hydrological Extremes* Retrieved 01/24/2011, from <<http://www.geo.uio.no/watersciences/themes-of-project>>
- Pallett, J., Heyns, P., Marais, C., & Larsson, V. (1997). *Sharing water in southern Africa*. Windhoek: Desert Research Foundation of Namibia.
- Peel, M. C., Finlayson, B. L., & McMahon, T. A. (2007). Updated world map of the Köppen-Geiger climate classification. *Hydrol. Earth Syst. Sci.*, 11(5), 1633-1644.
- Rootzen, H., & Tajvidi, N. (2006). Multivariate generalized Pareto distributions. *Bernoulli*, 12(5), 917-930.
- Rosbjerg, D. (2007). Regional flood frequency analysis. *Extreme Hydrological Events: New Concepts for Security*, 78, 151-171.
- Rossi, F., Fiorentino, M., & Versace, P. (1984). Two-Component Extreme Value Distribution for Flood Frequency Analysis. *Water Resources Research*, 20(7), 847-856.
- Saf, B. (2008). Application of index procedures to flood frequency analysis in Turkey. *Journal of the American Water Resources Association*, 44(1), 37-47.
- Saf, B., Dikbas, F., & Yasar, M. (2008). Regional flood frequency analysis of the lower west Mediterranean subregion of Turkey with L-moments. *Fresenius Environmental Bulletin*, 17(4), 427-433.
- Shu, C., & Ouarda, T. (2008). Regional flood frequency analysis at ungauged sites using the adaptive neuro-fuzzy inference system. *Journal of Hydrology*, 349(1-2), 31-43.

- SouthernAfrica. (2011). Southern Africa. *Encyclopædia Britannica Online* Retrieved 2/3/2011, from <http://www.britannica.com/EBchecked/topic/556618/Southern-Africa>
- Spence, E. S. (1973). Theoretical Frequency Distributions for the Analysis of Plains Streamflow. *Canadian Journal of Earth Sciences*, 10(2), 130-139.
- Stedinger, J., & Lu, L. (1995). Appraisal of regional and index flood quantile estimators. *Stochastic Hydrology and Hydraulics*, 9(1), 49-75.
- Stedinger, J. R., & Lu, L. H. (1995). Appraisal of regional and index flood quantile estimators. *Stochastic Hydrology and Hydraulics*, 9(1), 49-75.
- The Global Runoff Data Centre (GRDC), 56068 Koblenz, Germany.
- Tveito, O. E. (1993). *A regional flood frequency analysis of Norwegian catchments*. Institute report series NO.86, Institute of geophysics, University of Oslo, Oslo.
- UNNC. (2011). Heavy rainfall triggers flood alert in southern Africa. *UN News Center* Retrieved 15/08/2011, from <http://www.un.org/apps/news/story.asp?NewsID=37347&Cr=flood&Cr1>
- USWRC. (1981). *Guidelines for determining flood flow frequency*. Hydrology Subcommittee Bulletin 17B. Reston, Virginia.
- Viglione, A. (2010). Non-supervised Regional Frequency Analysis. *R-software, CRAN-nsRFA package, Version 0.7-0* Retrieved 05/03/2011, from [http://www.idrologia.polito.it/~alviglio/index\\_en.htm](http://www.idrologia.polito.it/~alviglio/index_en.htm)
- Viglione, A., Laio, F., & Claps, P. (2007). A comparison of homogeneity tests for regional frequency analysis. *Water Resources Research*, 43(3), W03428.
- Weibull, W. A. (1939). *A statistical theory of the strength of materials*. Ingeniorsvetenskapsakademien, Handlingar 151, Stockholm.
- Wikipedia. (2011). Southern Africa. *Wikipedia, the free encyclopedia* Retrieved 02/03/2011, from [http://en.wikipedia.org/wiki/Southern\\_Africa](http://en.wikipedia.org/wiki/Southern_Africa)
- Wiltshire, S. E. (1985). Grouping basins for regional flood frequency-analysis. *Hydrological Sciences Journal-Journal Des Sciences Hydrologiques*, 30(1), 151-159.
- Wiltshire, S. E. (1986). Regional flood frequency-analysis .1. Homogeneity statistics. *Hydrological Sciences Journal-Journal Des Sciences Hydrologiques*, 31(3), 321-333.
- Yang, T., Xu, C. Y., Shao, Q. X., & Chen, X. (2010). Regional flood frequency and spatial patterns analysis in the Pearl River Delta region using L-moments approach. *Stochastic Environmental Research and Risk Assessment*, 24(2), 165-182.

## APPENDICES

### A. Selected stations

#### 1. Stations used for regional flood frequency analysis

Table 1 Site characteristics of Namibia catchments (from 1969-2004)

S.N	Grdc_No	Nat_Id	River	Station	Area (km <sup>2</sup> )	Lat.( <sup>o</sup> )	Long. ( <sup>o</sup> )	Index flood (m3/s)
1	1255100	2811M01	Kunene River	Ruacana	89600	-17.4	14.2	591.7
2	1257100	2511M01	Okavango	Rundu	97300	-17.9	19.75	452.5
3	1258200	2962M03	Ugab River	Vingerklip	14200	-20.4	15.47	21.2
4	1258300	2971M02	Omaruru River	Etemba	3810	-21.43	15.67	38.0
5	1258501	2991M01	Kuiseb River	Schlesien	6520	-23.27	15.8	27.1
6	1291200	2400M01	Kwando River	Kongola	170000	-17.683	23.2667	33.5
7	1291100	2300M01	Zambezi	Katima Mulilo	334000	-17.47	24.3	3689.5
8	1259110	0497M03	Loewen	Altdorn	7000	-26.8	18.22	69.2

Table 2 Site characteristics of Malawi rivers (from 1954-1990)

S.N	Station code <sup>18</sup>	River @ <sup>19</sup> station	Area (km <sup>2</sup> )	Lat. ( <sup>o</sup> )	Long. ( <sup>o</sup> )	Index Flood (m3/s)
1	1992100	Domasi @ Domasi	72.8	-15.23	35.38	28.9
2	1992200	luweya @ Zayuka	2320	-11.80	34.37	231.0
3	1992850	Ruo @ M1 Roadbridge	193	-16.50	35.40	227.6
4	1992690	RiviRivi @ Balaka	748	-15.34	34.93	174.2
5	1992900	shire @ chiromo	149500	-16.55	35.13	610.5
6	1992700	shire @ Liwonde	130200	-15.07	35.20	923.9
7	1992400	south rukuru@Phewzi	11132	-10.90	34.05	182.9
8	1992950	Thuchila@chonde	1440	-16.00	35.50	113.0

Table 3 Site characteristics of Zambia Rivers (from 1970-2004)

S.N	GRDC Code	National code	River	Station	Area (km <sup>2</sup> )	Lat. ( <sup>o</sup> )	Long.( <sup>o</sup> )	Index Flood (m3/s)
1	1591001	2400	Zambezi	Senanga	284538	-16.12	23.25	2049.8
2	1591100	1145	Makonde	Chivata village	3354	-13.33	23.15	51.0
3	1591237	1425	Luakela	Sachibondo	632	-11.53	24.42	30.3
4	1591401	4977	Kafue	Kasaka	150971	-15.82	28.22	1440.0
5	1591404	4669	Kafue	Kafue hook bridge	96239	-14.93	25.92	881.9
6	1591406	4280	Kafue	Machiya ferry	23065	-13.65	27.62	412.4
7	1591441	4302	Luswishi	Lwendo	2668	-12.92	27.35	102.0
8	1591470	4260	Kafue	Ndubeni	18509	-13.40	27.82	394.1
9	1591471	4200	Kafue	Mpatamato	12001	-13.25	28.13	323.0
10	1591480	4200	Kafue	Wusakili	9088	-12.88	28.25	230.9
11	1591490	4015	Muchindamu	Muchindamu	110	-11.87	27.13	11.1
12	1591720	5815	Mulungushi	Great north road bridge	1448	-14.30	28.55	41.8
13	1591820	2250	Luanginga	Kalabo	34621	-14.97	22.68	165.1
14	1593100	6670	Luapula	Chembe ferry	123072	-11.97	28.75	947.4
15	1593740	6350	Lukulu	Kasama/luwingu road bridge	6504	-10.18	30.97	197.2

Table 4 Site characteristics Zimbabwe catchments (From 1957-1990)

<sup>18</sup> The code of the stations was given randomly by the author  
<sup>19</sup> The location of the gauging stations

S.N	FRIEND Code	Local code	River	Location	Area (km <sup>2</sup> )	Lat. (°)	Long. (°)	Index Flood (M3/S)
1	63351113	A13	Gweru	Gweru River Causew.	4201	-18.72	28.80	112.8
2	63341012	C12	Mupfure	Twyford Weir	5307	-18.12	30.22	197.6
3	63315506	D6	Shawanhowe	Mutoko Road Brg.	1194	-17.63	31.60	84.15
4	63535231	B31	Thuli	Thuli Gorge	4090	-21.08	28.83	141.5
5	63535215	B15	Lumane	Insindi Weir	277	-20.60	29.60	15.5
6	63422219	E19	Macheke	Condo U/S	3383	-18.92	31.95	122.6
7	63533035		Thuli	Ntalali Causeway	5880	-21.32	28.95	450.6

Table 5 Site characteristics of South Africa Rivers (From 1969-2008)

S.NO	Station code	River@ <sup>20</sup> station	Area (km <sup>2</sup> )	Long (°)	Lat. (°)	Index Flood (m3/s)
1	A2H006	Pienaars River @ Klipdrift	1028	28.47	-26.18	17.8
2	A2H012	Krokodil River @ Kalkheuwel	2551	28	-26.43	68.0
3	A2H013	Magalies River @ Scheerpoort	1171	27.86	-26.39	13.5
4	A2H021	Pienaars River @ Buffelspoort	7483	27.74	-25.83	26.0
5	A2H023	Jukskei River @ Nietgedacht	686	28.08	-26.22	70.9
6	A5H006	Limpopo River @ Botswana	98240	28.04	-23.02	115.0
7	A6H029	Mogalakwena River @ Glen Alpine	11292	28.85	-23.23	38.4
8	A9H004	Mutale River @ Tengwe	320	30.59	-23.03	80.8
9	B4H003	Steelpoort River @ Buffelskloof	2240	29.91	-25.76	34.2
10	B6H004	Blyde River @ Chester	2241	30.92	-25.01	93.5
11	B7H010	Ngwabitsi River @ Harmony	318	30.44	-24.13	15.2
12	B8H010	Letsitele River @ Mohlabas Location	477	30.43	-24.48	31.8
13	B1H005	Olifants River @ Wolwekrans	3256	29.29	-26.39	112.4
14	C1H005	Leeu Spruit @ Welbedacht	341	29.4	-27.11	17.5
15	C1H006	Blesbok Spruit @ Rietvley	1094	29.62	-27.32	147.5
16	C2H001	Mooi River @ Witrand	3595	27.15	-27.29	7.5
17	C3H003	Harts River @ Taung	10990	24.86	-27.98	23.0
18	C4H004	Vet River @ Fizantkraal	16153	26.19	-28.03	179.0
19	C6H003	Vals River @ Mooifontein	7765	26.6	-27.4	136.0
20	C8H001	Wilge River @ Frankfort	15673	28.56	-27.57	256.2
21	C1H007	Vaal River @ Goedgeluk	4686	29.78	-27.3	81.3
22	D1H003	Orange River @ Aliwal-North	37075	26.86	-31.37	1485.5
23	D1H006	Kornet Spruit @ Maghaleen	2969	27.41	-30.73	254.3
24	D1H011	Kraai River @ Roodewal	8688	26.96	-31.65	224.4
25	D3H008	Orange River @ Marksdrift	99316	23.81	-29.87	332.2
26	D8H003	Orange River @ Vioolsdrif	850530	17.77	-29.23	393.2
27	D5H003	Fish River @ Hardeheuwel	1509	20.43	-32.21	21.41
28	D7H005	Orange River @ Upington	364560	21.29	-28.93	377.3
29	E2H003	Doring River @ Melkboom	24044	18.71	-32.6	229.3
30	G1H013	Berg River @ Drieheuvelds	2934	18.98	-33.97	204.3
31	H1H003	Bree River @ Ceres Toeken Geb.	657	19.32	-34.2	56.7
32	H6H009	Riviersonderend @ Reenen	2007	20.24	-34.8	101.6
33	H7H013	Buffeljags River @ Eenzaamheid	602	20.68	-34.17	99.1
34	K3H003	Maalgate River @ Knoetze Kama	145	22.35	-34.4	28.2
35	K2H002	Great-Brak River @ Wolvedans	131	22.27	-34.73	18.1
36	K5H002	Knysna River @ Milwood Forest Res.	133	23.14	-34.35	19.8
37	L7H006	Groot River @ Grootrivierspoort	29560	24.63	-34.58	56.9
38	Q9H002	Koonap River @ Adelaide	1245	26.41	-33.51	17.2
39	Q9H012	Great Fish River @ Brandt Legte	23067	26.55	-33.98	69.9
40	R1H015	Keiskamma River @ Farm 7	2530	27.46	-33.31	85.7
41	R3H003	Nahoon River @ Farm 305	473	27.89	-33.21	53.3
42	R2H005	Buffalo River @ King Williams Town	411	27.53	-33.39	38.9
43	S3H004	Black-Kei River @ Cathcart's Gift	1413	26.85	-32.07	14.8
44	S3H006	Klaas Smits River @ Weltevreden	2170	26.81	-32.29	14.6
45	S5H002	Tsomo River @ Wyk Maduma	2359	27.87	-32.69	52.0
46	R2H006	Mgqakwebe River @ Msenge Ridge	119	27.41	-33.36	15.7
47	U1H005	Mkomazi River @ Lot 93 1821	1744	29.95	-30.36	221.0
48	U2H006	Karkloof River @ Shafton	339	30.38	-30.27	23.0

@<sup>20</sup> The location of the gauging stations

49	U2H048	Mgeni River @ Midmar	928	30.23	-30.06	29.8
50	V1H001	Tugela River @ Tugela Drift	4176	29.86	-28.89	258.3
51	V1H010	Little Tugela River @ Winterton	782	29.65	-28.91	119.3
52	V2H004	Mooi River @ Doornkloof	1546	30.36	-29.36	73.1
53	V3H002	Buffels River @ Schurvepoort	1518	30.03	-27.76	29.7
54	V6H004	Sondags River @ Kleinfontein	658	30.13	-28.69	55.5
55	V7H020	Boesmans River @ Wagendrift	744	29.9	-29.5	65.4
56	T3H005	Tina River @ Mahlunqulu	2597	28.9	-31.93	148.4
57	T3H006	Tsitsa River @ Xonkonxa	4268	28.87	-31.52	295.7
58	T3H009	Mooi River @ Maclear	307	28.39	-31.37	57.4
59	T4H001	Mtamvuna River @ Gundrift	715	29.94	-30.79	41.2
60	T5H003	Polela River @ Coxhill	140	29.59	-30.57	26.2
61	T5H004	Mzimkulu River @ Fp 1609030	545	29.51	-30.4	59.7
62	W1H009	Mhlatuze River @ Riverview	2408	31.86	-29.62	55.0
63	W2H005	White Mfolozi River @ Overvloed	3939	31.44	-28.7	111.8
64	W2H006	Black Mfolozi River @ Native Res 12	1648	31.7	-28.24	92.3
65	W2H009	White Mfolozi River @ Doornhoek	432	30.89	-28.74	23.9
66	W4H006	Phongolo River @ M'Hlali	6846	31.93	-28.15	240.4
67	W5H005	Hlelo River @ Ishlelo	804	30.83	-27.49	16.4
68	W5H022	Assegai River @ Zandbank	2313	31.09	-27.96	64.2
69	X1H001	Komati River @ Hooggenoeg	5499	31.13	-26.2	83.0
70	X1H014	Mlumati River @ Lomati	1119	31.65	-26.1	59.0
71	X2H005	Nels River @ Boschrand	642	31.11	-26.37	15.9
72	X2H016	Krokodil River @ Tenbosch	10365	32.29	-25.49	132.7
73	X2H015	Elands River @ Lindenau	1554	30.84	-25.91	56.6
74	X2H022	Kaap River @ Dolton	1639	31.32	-26.13	36.3

## 2. Stations used for model verifications

Table 6 Selected stations for model verifications

Station code	River @ Station	Country	Area(km <sup>2</sup> )	Latitude (°)	Longitude (°)	Data Available	Record Length (Years)
1591404	Kafue@Kafue Hook Bridge	Zambia	96239	-14.93	25.92	1973-2005	30
1591474	Kafue@ Kafironda	Zambia	Na	-13.633	27.59	1969-1991	22
A5H003	Limpopo River @ Botswana	South Africa	98160	-22.96	28.04	1959-1980	22
X1H003	Komati River @ Tonga	South Africa	8614	-25.82	31.92	1969-2008	40
C2H008	Vaal River @ Woodlands	South Africa	47214	-26.76	27.68	1969- 1996	27
D8H008	Orange River @ Pella Mission	South Africa	821850	-29.08	19.17	1979-2008	29
V1H038	Klip River @ Ladysmithdorpsgronde	South Africa	1644	-28.66	29.76	1971- 2008	38
W5H024	Mpuluzi River @ Dumbarton	South Africa	1446	-26.42	30.95	1976-2008	33
T3H007	Mzimvubu River @ Ku-Makhola	South Africa	6906	-30.95	29.11	1972-2008	24
U2H012	Sterk River @ Groothoek	South Africa	438	-29.48	30.53	1969-2008	34

## B. At- site statistical behaviors of annual maximum floods

Table 7 at site statistical characteristics for Namibia Catchments

S.N.	Station	$\bar{x}$	$Cv$	$Cs$	$C_{KURT}$	$\lambda_2$	$\tau_2$	$\tau_3$	$\tau_4$	$SE_{jack}$
1	1255100	686.2	0.654	0.765	2.844	251.3	0.366	0.215	0.077	233
2	1257100	493.3	0.378	0.496	3.063	107.5	0.218	0.127	0.133	19.8
3	1258200	39.8	1.046	1.895	7.435	20.5	0.516	0.400	0.200	0.9
4	1258300	68.2	1.018	1.367	4.968	36.4	0.534	0.337	0.121	20.9
5	1258501	32.0	0.713	0.880	3.826	12.7	0.397	0.201	0.086	2.2
6	1291200	44.9	0.570	1.321	4.554	13.6	0.302	0.329	0.158	1.4
7	1291100	3784.3	0.473	0.462	2.958	1030.0	0.272	0.110	0.106	675
8	1259110	94.3	1.188	3.224	16.563	48.8	0.517	0.405	0.290	4.3

Table 8 At-site statistical characteristics for Zimbabwe Catchments

S.N.	Station	$\bar{x}$	$Cv$	$Cs$	$C_{KURT}$	$\lambda_2$	$\tau_2$	$\tau_3$	$\tau_4$	$SE_{jack}$
1	63315506	107.0	0.841	0.838	3.389	50.0	0.467	0.220	0.067	66.2
2	63341012	358.6	0.980	0.823	2.444	190.4	0.531	0.289	0.018	192.4
3	63351113	114.5	0.585	0.285	2.688	38.6	0.337	0.079	0.080	10.32
4	63535215	27.4	0.974	1.097	3.163	14.1	0.514	0.341	0.101	0.29
5	63535231	208.4	0.984	1.197	3.912	109.0	0.523	0.327	0.117	31.71
6	63533035	1220.9	1.368	2.444	9.917	749.1	0.614	0.516	0.283	147.6
7	63422219	127.5	0.797	2.517	13.529	48.9	0.384	0.220	0.221	7.82

Table 9 At-site statistical characteristics of Zambia Catchments

S.N.	Station	$\bar{x}$	$Cv$	$Cs$	$C_{KURT}$	$\lambda_2$	$\tau_2$	$\tau_3$	$\tau_4$	$SE_{jack}$
1	1591001	2046.3	0.271	-0.319	2.66	328.6	0.161	-0.040	0.099	43.4
2	1591100	46.1	0.502	0.048	2.49	13.4	0.292	0.026	0.075	9.01
3	1591237	32.4	0.359	0.401	3.38	6.7	0.207	0.127	0.154	1.20
4	1591401	1412.9	0.263	0.104	3.45	212.5	0.150	0.011	0.230	36.1
5	1591404	1006.6	0.558	0.791	3.44	316.3	0.314	0.193	0.116	5.2
6	1591406	419.0	0.434	0.230	3.66	103.6	0.247	0.041	0.182	11.7
7	1591441	96.1	0.465	0.032	3.47	25.3	0.264	-0.033	0.120	25.9
8	1591470	382.9	0.471	0.212	3.75	102.0	0.266	0.049	0.193	24.3
9	1591471	287.3	0.488	-0.379	2.25	81.6	0.284	-0.074	0.027	80.8
10	1591480	214.9	0.476	0.309	3.57	58.3	0.271	0.049	0.156	16.6
11	1591490	14.9	0.878	0.707	2.60	7.3	0.488	0.230	0.024	1.74
12	1591720	45.2	0.593	0.884	4.03	14.9	0.330	0.189	0.172	0.82
13	1591820	154.9	0.508	-0.062	2.79	45.4	0.293	-0.002	0.120	8.95
14	1593100	1091.5	0.727	0.722	2.97	445.9	0.409	0.193	0.087	191
15	1593740	196.3	0.249	-0.253	4.03	28.2	0.144	-0.019	0.190	10.6



Table 10 At-site statistical characteristics for Malawi Catchments

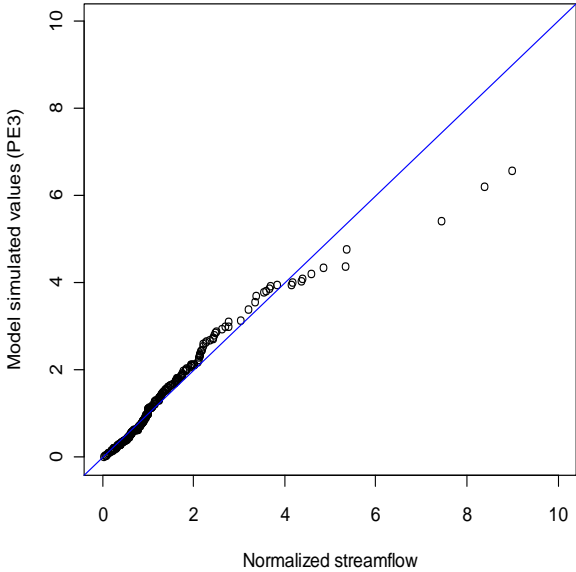
S.N.	Station	$\bar{x}$	$Cv$	$Cs$	$C_{KURT}$	$\lambda_2$	$\tau_2$	$\tau_3$	$\tau_4$	$SE_{jack}$
1	1992100	35.4	0.870	3.759	21.18	12.8	0.363	0.384	0.276	2.40
2	1992200	219.9	0.413	0.042	2.51	53.1	0.242	0.008	0.095	4.38
3	1992850	242.3	0.520	0.805	3.63	70.8	0.292	0.173	0.142	25.38
4	1992400	166.1	0.474	0.247	2.61	45.5	0.274	0.076	0.085	6.80
5	1992700	620.2	0.309	0.223	3.15	112.0	0.181	0.055	0.137	19.84
6	1992900	1015.5	0.360	1.474	5.96	194.7	0.192	0.255	0.260	6.17
7	1992950	266.1	0.795	1.428	4.54	109.1	0.410	0.375	0.173	21.90
8	1992690	134.4	0.765	1.930	8.46	52.1	0.388	0.319	0.173	2.44

Table 11 At-site statistical characteristics for South Africa Catchments

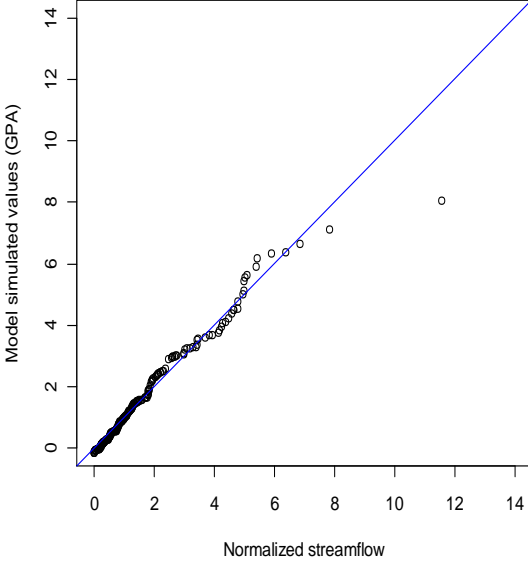
S.N.	Station	$\bar{x}$	$Cv$	$Cs$	$C_{KURT}$	$\lambda_2$	$\tau_2$	$\tau_3$	$\tau_4$	$SE_{jack}$
1	A2H006	48.6	1.395	2.044	7.08	31.0	0.637	0.527	0.260	0.63
2	A2H012	96.5	0.892	2.150	8.09	40.9	0.424	0.423	0.264	0.55
3	A2H013	30.0	1.336	1.869	6.06	18.8	0.626	0.496	0.232	7.27
4	A2H021	76.0	1.262	1.711	5.22	45.8	0.603	0.486	0.205	16.67
5	A2H023	87.7	0.795	2.586	13.07	33.4	0.381	0.320	0.235	2.40
6	A5H006	214.4	1.235	2.426	10.09	122.5	0.571	0.433	0.271	73.38
7	A6H029	68.0	0.977	0.950	2.74	35.7	0.525	0.306	0.066	3.73
8	A9H004	102.9	0.879	1.147	4.54	49.0	0.476	0.253	0.111	0.44
9	B4H003	45.7	0.787	1.636	6.05	18.5	0.406	0.329	0.173	4.92
10	B6H004	116.7	0.848	1.393	5.29	52.3	0.448	0.294	0.131	21.27
11	B7H010	45.5	1.607	2.422	8.48	31.0	0.681	0.563	0.319	1.79
12	B8H010	64.8	1.257	2.810	13.19	36.4	0.562	0.463	0.271	22.78
13	B1H005	195.2	1.189	1.828	5.98	112.2	0.575	0.435	0.216	27.45
14	C1H005	22.0	1.289	2.740	10.91	12.0	0.547	0.448	0.344	1.48
15	C1H006	216.9	0.907	0.930	3.21	107.7	0.496	0.259	0.071	42.16
16	C2H001	13.7	1.360	2.397	8.83	8.3	0.606	0.501	0.282	2.50
17	C3H003	41.9	1.259	3.108	15.89	23.1	0.551	0.460	0.270	2.19
18	C4H004	238.9	0.965	1.022	4.06	124.8	0.522	0.247	0.047	107.32
19	C6H003	221.3	1.211	1.903	6.17	126.9	0.573	0.461	0.249	17.28
20	C8H001	493.3	1.239	1.917	6.52	288.9	0.586	0.490	0.215	131.58
21	C1H007	168.6	1.039	1.409	4.51	88.7	0.526	0.405	0.131	24.51
22	D1H003	1873.1	0.747	0.976	3.58	766.0	0.409	0.250	0.078	254.65
23	D1H006	261.1	0.520	0.214	2.03	78.6	0.301	0.085	0.013	93.61
24	D1H011	348.2	1.041	2.246	7.92	164.0	0.471	0.476	0.326	20.87
25	D3H008	633.5	1.386	3.776	19.81	334.8	0.529	0.564	0.397	46.22
26	D8H003	1137.0	1.266	3.138	17.03	634.3	0.558	0.423	0.183	13.58
27	D5H003	33.1	1.200	2.205	9.59	19.2	0.581	0.393	0.199	11.32
28	D7H005	1308.7	1.260	2.438	10.87	760.1	0.581	0.468	0.207	58.15
29	E2H003	276.4	0.710	1.015	3.67	107.7	0.390	0.245	0.111	65.84
30	G1H013	271.5	0.638	0.723	2.76	97.3	0.358	0.208	0.069	17.60
31	H1H003	64.0	0.490	0.854	3.56	17.6	0.275	0.189	0.117	2.37
32	H6H009	189.3	1.423	3.187	14.48	106.8	0.564	0.587	0.387	6.04
33	H7H013	139.8	0.713	0.906	3.12	54.5	0.390	0.278	0.046	22.81
34	K3H003	36.4	0.766	1.627	5.93	14.3	0.393	0.340	0.212	0.49
35	K2H002	23.0	0.626	0.681	2.82	8.1	0.352	0.191	0.094	0.03
36	K5H002	29.9	0.785	0.790	2.47	12.9	0.432	0.257	0.046	2.43
37	L7H006	144.5	1.469	1.975	6.10	95.6	0.662	0.533	0.272	32.87
38	Q9H002	46.0	1.359	1.551	4.26	29.7	0.644	0.488	0.188	2.56
39	Q9H012	329.9	1.641	2.136	7.03	232.9	0.706	0.604	0.305	15.16

40	R1H015	100.0	0.837	1.981	8.93	41.5	0.414	0.222	0.143	27.00
41	R3H003	100.3	1.701	3.327	15.86	67.7	0.676	0.557	0.371	6.94
42	R2H005	63.5	1.026	2.261	9.35	30.9	0.487	0.424	0.269	11.07
43	S3H004	27.0	1.349	3.720	19.76	14.3	0.531	0.526	0.363	8.33
44	S3H006	26.1	1.152	2.073	7.81	14.4	0.551	0.433	0.238	4.87
45	S5H002	48.8	0.839	2.105	9.70	19.9	0.408	0.208	0.180	0.05
46	R2H006	16.2	0.672	0.574	2.88	6.2	0.382	0.147	0.078	2.59
47	U1H005	250.4	0.543	1.083	3.97	74.7	0.298	0.249	0.145	10.05
48	U2H006	31.8	1.141	4.215	24.52	13.8	0.434	0.495	0.394	2.31
49	U2H048	53.7	2.175	5.545	36.39	34.4	0.641	0.600	0.486	0.06
50	V1H001	432.8	0.901	1.299	4.19	204.6	0.473	0.341	0.134	45.03
51	V1H010	146.8	0.677	1.024	3.48	54.2	0.369	0.259	0.133	9.49
52	V2H004	103.4	1.086	3.402	15.36	43.1	0.417	0.531	0.430	2.47
53	V3H002	46.4	1.002	1.282	4.08	24.4	0.527	0.339	0.131	6.16
54	V6H004	68.4	0.699	1.472	5.59	25.1	0.367	0.280	0.212	9.62
55	V7H020	95.3	0.908	2.686	12.61	39.8	0.418	0.421	0.273	9.96
56	T3H005	186.6	0.883	1.891	7.64	82.8	0.444	0.334	0.217	29.89
57	T3H006	325.0	0.701	1.142	4.14	123.2	0.379	0.270	0.117	76.55
58	T3H009	90.5	0.848	1.243	4.26	40.5	0.448	0.321	0.117	1.57
59	T4H001	86.2	1.418	3.318	16.01	48.8	0.566	0.575	0.374	14.86
60	T5H003	35.8	0.769	2.015	8.68	13.7	0.383	0.365	0.179	12.39
61	T5H004	74.1	0.517	1.961	6.97	18.8	0.254	0.396	0.271	10.66
62	W1H009	141.5	1.351	2.643	11.88	84.9	0.600	0.513	0.266	7.71
63	W2H005	161.8	0.890	2.352	8.96	66.5	0.411	0.418	0.300	20.96
64	W2H006	141.1	0.863	1.674	5.93	61.3	0.435	0.385	0.202	6.87
65	W2H009	46.3	2.174	5.406	34.81	29.5	0.638	0.631	0.524	0.98
66	W4H006	307.4	0.605	0.689	2.82	104.6	0.340	0.194	0.093	14.87
67	W5H005	28.9	1.045	2.939	14.09	13.3	0.460	0.479	0.294	1.82
68	W5H022	81.0	0.686	0.668	2.62	31.4	0.387	0.197	0.060	17.51
69	X1H001	136.6	1.148	2.769	13.08	70.8	0.518	0.453	0.268	2.69
70	X1H014	132.4	1.421	2.645	10.77	82.0	0.619	0.520	0.308	12.14
71	X2H005	25.3	1.079	2.276	8.75	12.8	0.504	0.437	0.271	0.1
72	X2H016	258.5	1.085	1.857	6.06	134.4	0.520	0.455	0.230	59.51
73	X2H015	56.9	0.662	0.865	3.67	20.9	0.367	0.184	0.114	17.75
74	X2H022	65.4	1.308	3.342	17.58	36.4	0.556	0.478	0.298	19.37

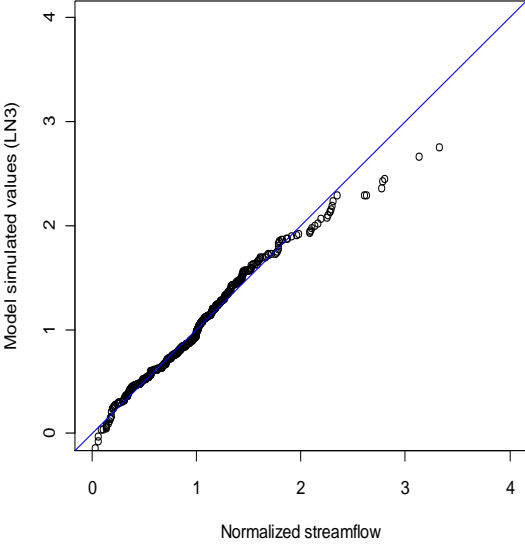
**C. The quantile flood of the normalized observed versus simulated values**



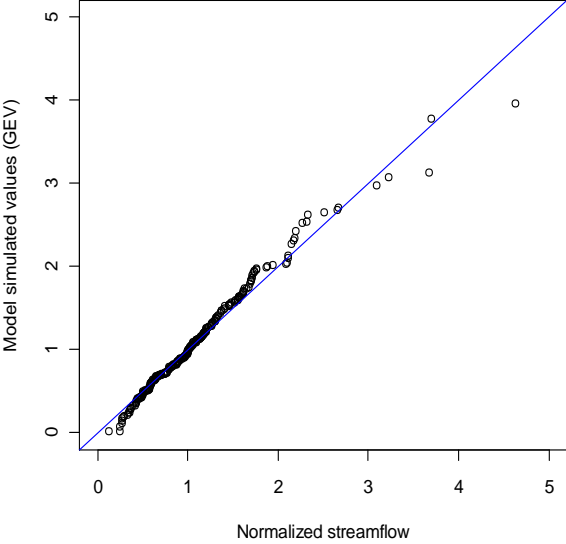
*a) R1*



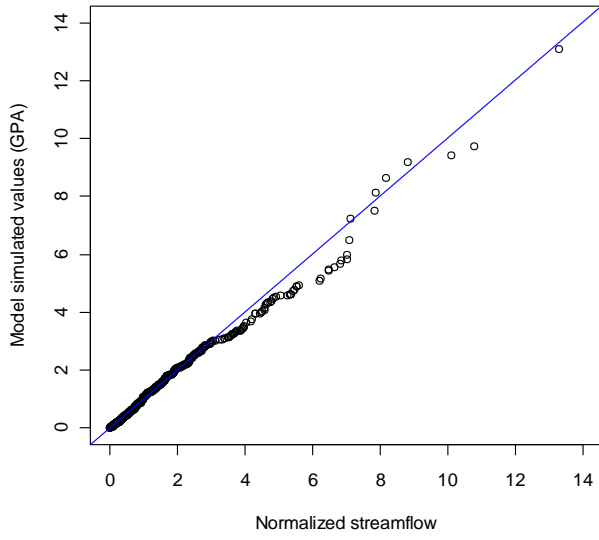
*b) R2*



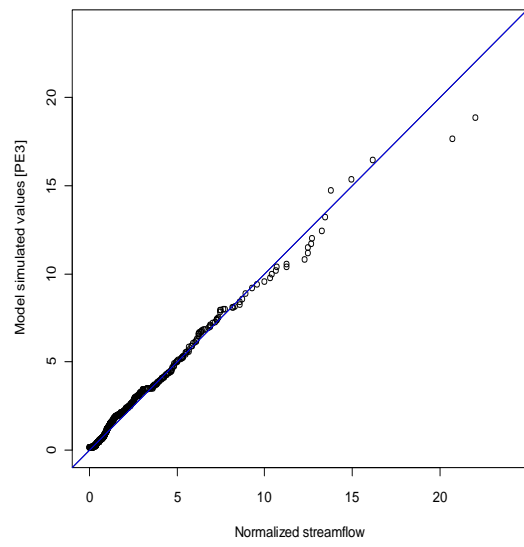
*c) R3*



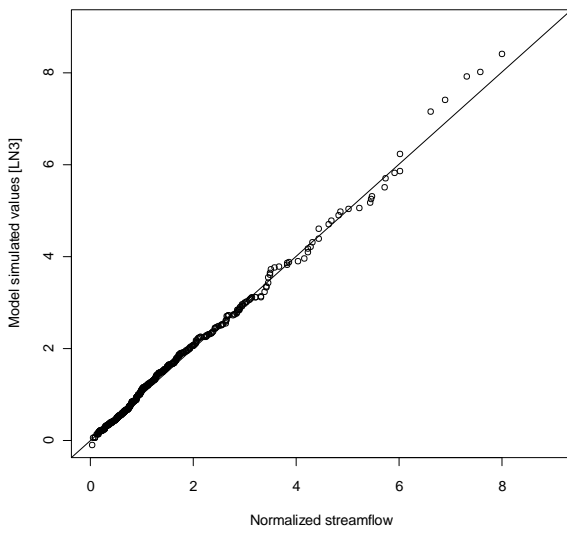
*d) R4*



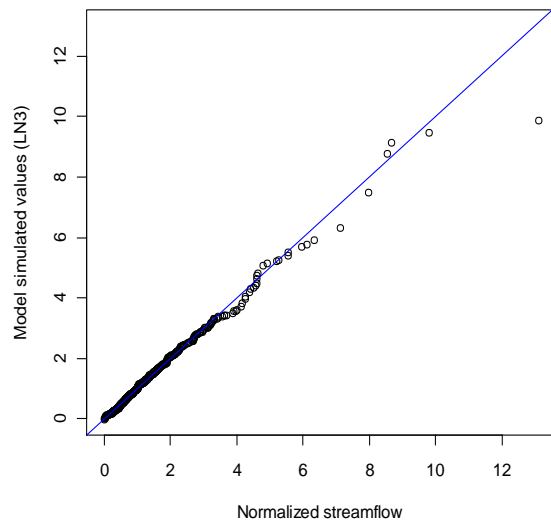
e) ZA\_R1



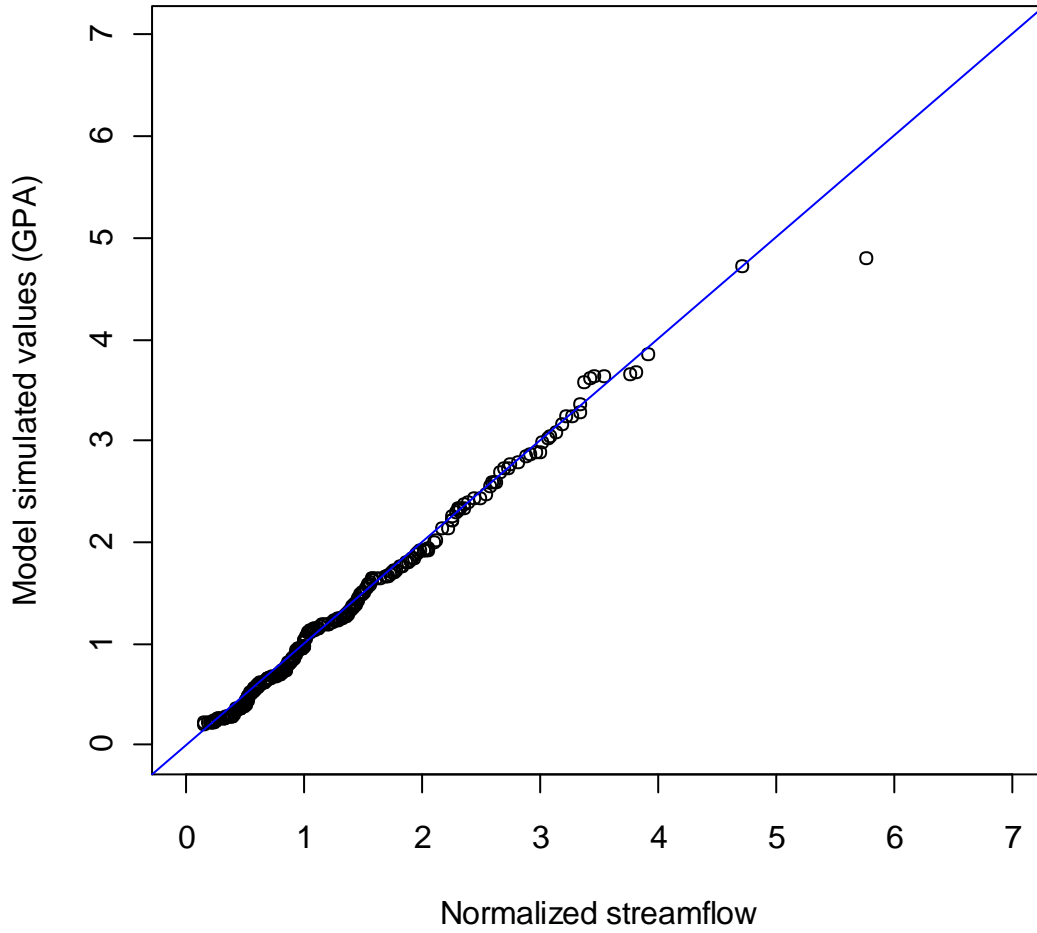
f) ZA\_R2



a) ZA\_R3



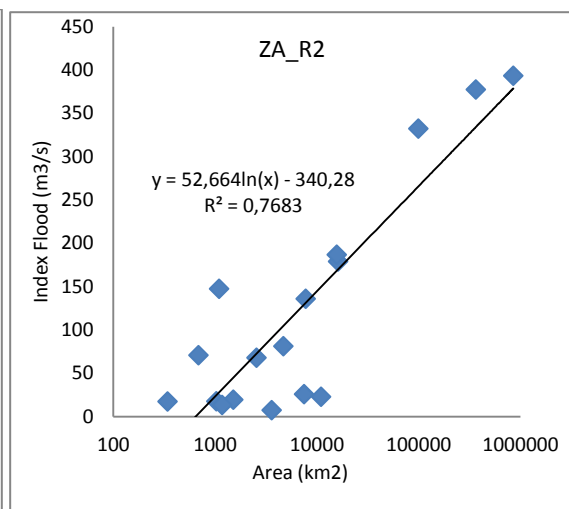
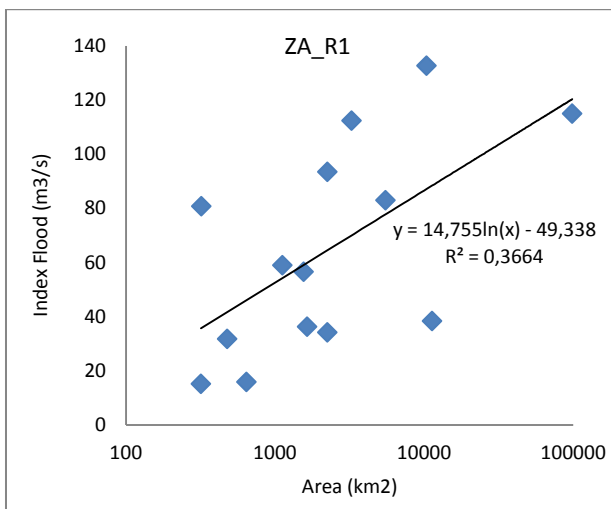
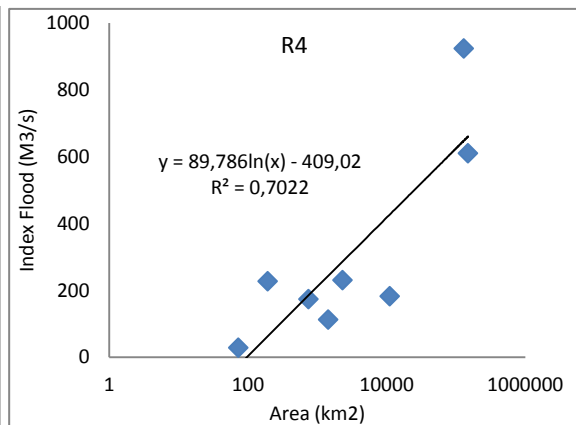
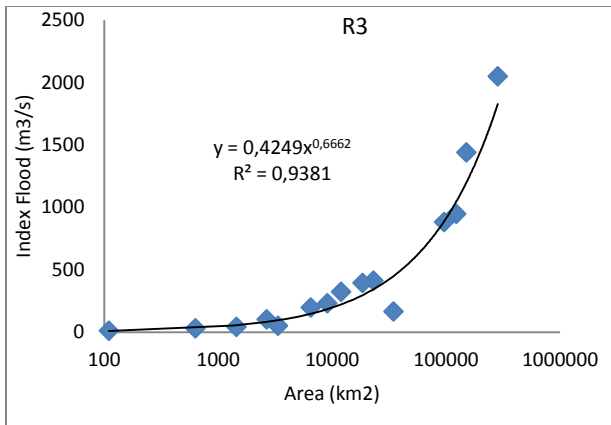
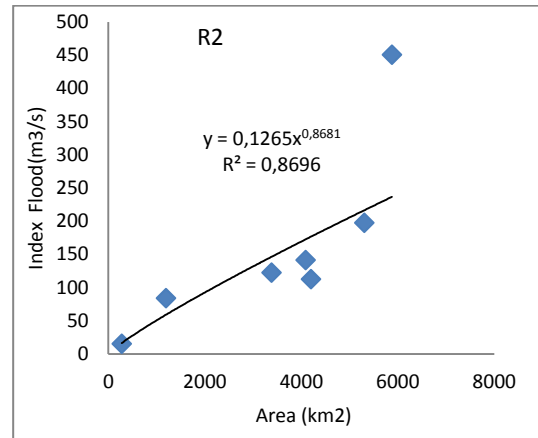
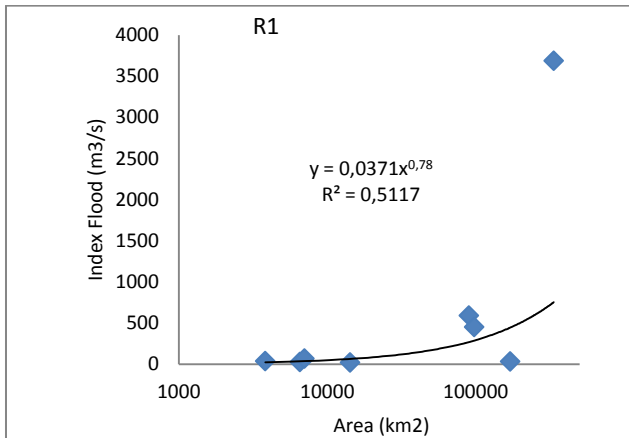
h) ZA\_R4



i) ZA\_R5

*Figure 1 plotting the normalized quantile values of the observed against randomly simulated using best fitted regional distribution*

## D. The regional Regression of at-site median values



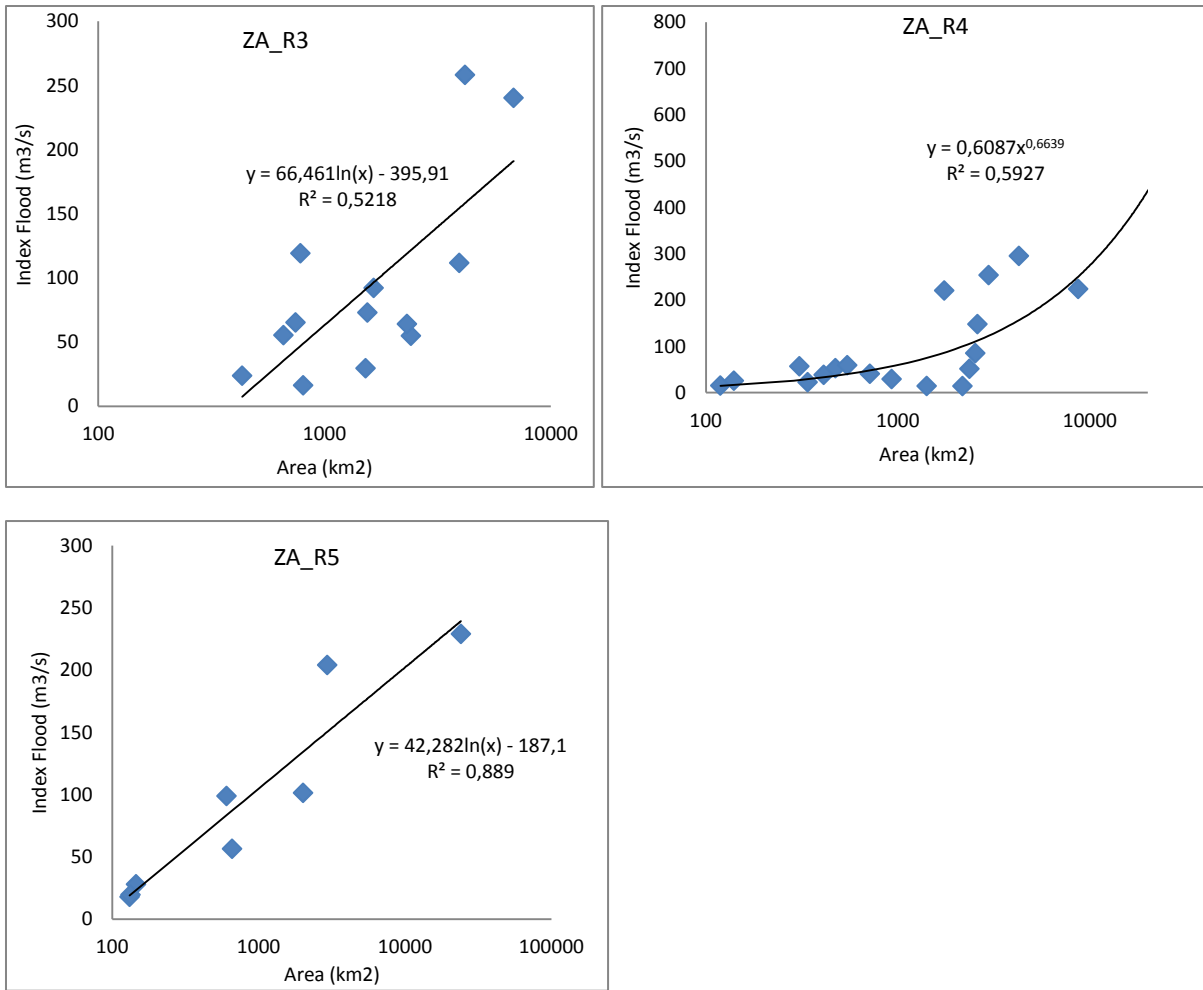


Figure 2 The regional regression coefficients showing the relationships between the index flood (median) and catchments area.

## E. Theoretical distributions and their relationships

Table 12 Theoretical distribution functions and their moments. Taken from Geo4310 lecture notes by Gottschalk and Krasovskaia (2001) and Hosking and Wallis (1997):  $x$  = observed values,  $m$  = mean value,  $\sigma$  = standard deviation,  $C_s$  = coefficient of variance  $\alpha$  = scale parameter,  $\mu$  = location parameter and  $k$  = shape parameter

DISTRIBUTION	MOMENTS
<b>Normal:</b> $F(x) = \int_{-\infty}^x \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-m}{\sigma}\right)^2\right\}$	$m$ $\sigma$
<b>Uniform:</b> $F(x) = (x-a)/(b-a)$	$m = \frac{1}{2}(a+b)$ $\sigma = \frac{1}{\sqrt{12}}(b-a)$
<b>Lognormal:</b> $F(x) = \int_0^x \frac{1}{x\sigma_n\sqrt{\pi}} \exp\left\{-\frac{1}{2}\left(\frac{(\ln x - m_n)^2}{\sigma_n^2}\right)\right\} dx$	$m = \exp(m_n + \sigma_n^2/2)$ $\sigma = m\sqrt{\exp(m_n + \sigma_n^2/2)}$ $C_s = (\exp(3\sigma_n^2) - 3\exp(\sigma_n^2) + 2) / (\exp(\sigma_n^2) - 1)^{3/2}$
<b>Generalised Normal:</b> $F(x) = \Phi\left[-k^{-1} \ln\{1 - k(x-u)/\alpha\}\right]$	$m = u + \frac{\alpha}{k}(1 - \exp(k^2/2))$
<b>Gamma:</b> $F(x) = \int_0^x \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} dx$	$m = \alpha\beta$ $\sigma = \beta\sqrt{\alpha}$ $C_s = 2/\sqrt{\alpha}$
<b>Pearson type III:</b> $F(x) = \int_0^x \frac{1}{\beta^\alpha \Gamma(\alpha)} (x-\gamma)^{\alpha-1} e^{-(x-\gamma)/\beta} dx$	$m = \alpha\beta + \gamma$ $\sigma = \beta\sqrt{\alpha}$ $C_s = 2/\sqrt{\alpha}$
<b>Gumbel (EV1):</b> $F(x) = \exp\left(-\exp\left(\frac{x-u}{\alpha}\right)\right)$	$m = u + 0.5772\alpha$ $\sigma = \frac{\pi\alpha}{\sqrt{6}}$ $C_s = 1.14$
<b>Generalised Extreme Value (GEV):</b> $F(x) = \exp\left\{-\left[1 - \frac{k(x-u)}{\alpha}\right]^{\frac{1}{k}}\right\}$	$m = u + \frac{\alpha}{k}[1 - \Gamma(1+k)]$ $\sigma = \alpha\sqrt{[\Gamma(1+2k) - \Gamma^2(1+k)]}$ $C_s = \frac{\Gamma(1+3k)/\Gamma^3(1+k) - 3\Gamma(1+2k)/\Gamma^2(1+k) + 2}{[\Gamma(1+2k)/\Gamma^2(1+k) - 1]}$
<b>Generalised Pareto (GPA):</b> $F(x) = 1 - \left[1 - \frac{k(x-u)}{\alpha}\right]^{\frac{1}{k}}$	$m = u + \alpha/(1+k)$ $\sigma = \alpha/\sqrt{(1+2k)(1+k)^2}$ $C_s = 2(1-k)\sqrt{1+2k}/(1+3k)$



Table 13 Theoretical relationships of L-moments and the inverse of the some cumulative distribution function (Gottschalk and Krasovskaia, 2001; Hosking and Wallis, 1997)

DISTRIBUTION AND ITS INVERSE	L-MOMENTS
<b>Normal:</b> $x = m + \sigma \Phi^{-1}[F]$	$\lambda_1 = m$ $\lambda_2 = \sigma / \sqrt{\pi}$ $\tau_3 = 0$ $\tau_4 = 0.1226$
<b>Gamma:</b> The inverse of $F(x) = \int_0^x \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} dx$	$\lambda_1 = \alpha\beta$ $\lambda_2 = \frac{\beta}{\sqrt{\pi}} \frac{\Gamma(\alpha + \frac{1}{2})}{\Gamma(\alpha)}$ .. $\tau_3 = 6I_{1/3}(\alpha, 2\alpha) - 3$
<b>Pearson type III:</b> $x(F)$ not explicitly defined	$\lambda_1 = \alpha\beta + \gamma$ $\lambda_2 = \frac{\beta}{\sqrt{\pi}} \frac{\Gamma(\alpha + \frac{1}{2})}{\Gamma(\alpha)}$ $\tau_3 = 6I_{1/3}(\alpha, 2\alpha) - 3$
<b>Gumbel (EV1):</b> $x = u - \alpha \ln[-\ln(F)]$	$\lambda_1 = u + 0.5772\alpha$ $\lambda_2 = \alpha \ln(2)$ $\tau_3 = 0.1669; \quad \tau_4 = 0.1504$
<b>Generalised Extreme Value (GEV):</b> $x = u + \frac{\alpha}{k} \left[ 1 - (-\ln(F))^k \right]$	$\lambda_1 = u + \frac{\alpha}{k} [1 - \Gamma(1+k)]$ $\lambda_2 = \frac{\alpha}{k} (1 - 2^{-k}) \Gamma(1+k)$ $\tau_3 = \frac{3 \cdot 2^{-k} - 2 \cdot 3^{-k} - 1}{1 - 2^{-k}}$ $\tau_4 = (1 - 6 \cdot 2^{-k} + 10 \cdot 3^{-k} - 5 \cdot 4^{-k}) / (1 - 2^{-k})$
<b>Generalized Pareto (GPA):</b> $x = u + \alpha \left[ 1 - (1-F)^k \right]$	$\lambda_1 = u + \alpha / (1+k)$ $\lambda_2 = \alpha / (1+k) (2+k)$ $\tau_3 = (1-k) / (3+k)$ $\tau_4 = (1-k)(2-k) / (3+k)(4+k)$
<b>Lognormal:</b> $y = \exp(m_n + \sigma_n \Phi^{-1}[F])$	$\lambda_1 = \exp(m_n + \sigma_n^2 / 2)$ $\lambda_2 = \exp(m_n + \sigma_n^2 / 2) \operatorname{erf}(\sigma_n / 2)$ . $\tau_3 = 6\pi^{-\frac{1}{2}} / \operatorname{erf}(\sigma_n / 2) \int_0^{\sigma_n/2} \operatorname{erf}(x/\sqrt{3}) \exp(-x^2) dx$
<b>Exponential (EXP):</b> $x(F) = u - \alpha \ln(1-F)$	$\lambda_1 = u + \alpha$ $\lambda_2 = 1/2\alpha$ $\tau_3 = 1/3$ $\tau_4 = 1/6$
<b>Generalized logistic (GLO):</b> $x(F) = \begin{cases} u + \frac{\alpha \left[ 1 - \left( \frac{1-F}{F} \right)^k \right]}{k}, & K \neq 0 \\ u - \alpha \ln\{(1-F)/F\}, & K = 0 \end{cases}$	$\lambda_1 = u + \alpha(1/k - \pi / \operatorname{sinc}k\pi)$ $\lambda_2 = \alpha k \pi / \operatorname{sinc}k\pi$ $\tau_3 = -k$ $\tau_4 = (1 + 5k^2) / 6$