

PARAMETRIC AND SEMIPARAMETRIC ESTIMATION METHODS FOR  
BIVARIATE COPULA IN RAINFALL APPLICATION

RAHMAH BINTI MOHD LOKOMAN

A dissertation submitted in partial fulfilment of the  
requirements for the award of the degree of  
Master of Science

Faculty of Science  
Universiti Teknologi Malaysia

APRIL 2017

Specially dedicated to my beloved parents,  
Mohd Lokoman bin Kasiran  
and Normah binti Khalil,  
all my siblings, my family and all my friends.  
Thank you for your support, love, and encouragement.

## ACKNOWLEDGEMENT

In the Name of Allah Most Gracious Most Merciful. First and foremost, all praise to the Almighty Allah who allowed me and blessed me with the capability to accomplish this dissertation. Greetings and blessings to the beloved, Prophet Muhammad, His Family and Companion.

I would like to express my gratitude and a big thank you to my supervisor, Assoc. Prof. Dr. Fadhilah Binti Yusof, and my co-supervisors, Dr. Arifah Bahar for the encouragement, guidance, advice and motivation that has given to me throughout the completion of this dissertation. With their continued support and interest, I am able to finish up this dissertation report.

Besides, I also would like to express my gratitude to the staff of the Mathematical Science Department, Universiti Teknologi Malaysia for their kind assistance. Also not forgotten, the staff for both UTM library, Perpustakaan Sultanah Zanariah (PSZ) and Perpustakaan Raja Zarith Sofiah (PRZS), thanks for their assistance in supplying the related literature. This dissertation gives me an opportunity to learn more about subjects that are related to Copula.

Finally, to the most important people in my life, my beloved parents, Mohd Lokoman, and Normah and my family, I specially thank them and grateful for their prayers and supports on the completion of this dissertation.

## ABSTRACT

Hydrological phenomena such as drought, flood, and rainfall are one of the natural phenomena that often provide dependent multivariate variables. The correlation of the hydrologic dependent variables can be described by using copula. To determine a specified copula structure that fitted with the marginal variables, the copula dependence parameter needs to be estimated. This study focuses on the application of parametric and semiparametric approaches in estimating the copula dependence parameter. The performance of seven parameter estimation methods namely, maximum likelihood (ML) estimation, inference function of margins (IFM), maximization by parts (MBP), pseudo maximum likelihood (PML), the inversion of rank correlation coefficient approach based on Kendall's tau and Spearman's rho and maximum likelihood based on kernel density estimation (MLKDE) are compared in the simulation and empirical studies. The simulation and empirical studies are limited to the case of bivariate copulas. The result from the simulation study shows that the parametric approaches are inefficient when the marginal distributions are misspecified. Among the parametric approaches, MBP performs better than MLE and IFM. While, for semiparametric approaches, PML performs well and consistent for any correlation and sample size. The PML can be efficient and consistent with the parametric once the sample size is large. The empirical study is done by applying the estimation methods to identify the dependence of the daily rainfall at two rain gauge stations Station Kuala Krai and Station Ulu Sekor. The result from the empirical study is consistent with the result from the simulation study. Thus, it can be concluded that MBP is preferred when the copula and the marginal distributions are known. While, PML is preferred when the marginal distribution is unknown, where the situation is common in a real data application.

## ABSTRAK

Fenomena hidrologi seperti kemarau, banjir dan hujan adalah salah satu fenomena semula jadi yang sering membentuk pemboleh ubah multivariat bersandar. Korelasi antara pemboleh ubah hidrologi tersebut boleh digambarkan dengan menggunakan *Copula*. Untuk memadamkan struktur *Copula* tertentu yang sesuai dengan pemboleh ubah marginal, parameter bersandar *Copula* perlu dianggarkan dahulu. Kajian ini memberi tumpuan kepada penggunaan kaedah parametrik dan semiparametrik dalam menganggarkan parameter bersandar *Copula*. Prestasi tujuh kaedah penganggar iaitu anggaran kebolehdajian maksimum (MLE), fungsi taksiran marginal (IFM), kaedah pengoptimuman bahagian demi bahagian (MBP), pseudo kebolehdajian maksimum (PML), kaedah penyongsangan terhadap pekali kedudukan korelasi *Kendall tau* dan *Spearman rho* serta kebolehdajian maksimum berdasarkan anggaran ketumpatan kernel (MLKDE) telah dibandingkan dalam kajian simulasi dan kajian empirikal. Kajian simulasi dan kajian empirikal adalah terhad kepada kes *Copula* bivariat. Hasil daripada kajian simulasi menunjukkan bahawa pendekatan parametrik tidak efisien apabila taburan marginal disalah spesifikasikan. Antara kaedah parametrik, MBP mempunyai prestasi yang lebih baik daripada MLE dan IFM. Sementara itu, untuk kaedah semiparametrik, PML mempunyai prestasi yang lebih baik dan konsisten bagi setiap tahap korelasi dan saiz sampel. PML boleh menjadi efisien dan konsisten seperti parametrik apabila saiz sampel adalah besar. Kajian empirikal dilakukan dengan menggunakan kaedah anggaran untuk mengenal pasti kebersandaran hujan harian di dua stesen tolok hujan: Stesen Kuala Krai dan Stesen Ulu Sekor. Hasil kajian empirikal adalah konsisten dengan keputusan daripada hasil kajian simulasi. Secara kesimpulan, MBP lebih sesuai digunakan apabila *Copula* dan taburan marginal dapat dikenal pasti. Manakala PML lebih sesuai digunakan apabila taburan marginal tidak diketahui dan tidak dapat dikenalpasti. Situasi begini adalah biasa dalam aplikasi data sebenar.

## TABLE OF CONTENTS

<b>CHAPTER</b>	<b>TITLE</b>	<b>PAGE</b>
	<b>DECLARATION</b>	ii
	<b>DEDICATION</b>	iii
	<b>ACKNOWLEDGEMENT</b>	iv
	<b>ABSTRACT</b>	v
	<b>ABSTRAK</b>	vi
	<b>TABLE OF CONTENTS</b>	vii
	<b>LIST OF TABLES</b>	xiii
	<b>LIST OF FIGURES</b>	xvi
	<b>LIST OF SYMBOLS</b>	xviii
	<b>LIST OF ABBREVIATIONS</b>	xxi
	<b>LIST OF APPENDICES</b>	xxiii
<b>1</b>	<b>INTRODUCTION</b>	<b>1</b>
	1.1 Background of the study	1
	1.2 Problem statement	5

1.3	Research questions	6
1.4	Objectives of the study	6
1.5	Scopes of the study	7
1.6	Significance of Study	8
1.7	Research outline	9
<b>2</b>	<b>LITERATURE REVIEW</b>	<b>10</b>
2.1	Introduction	10
2.2	Copula	10
2.2.1	Definition of copula	11
2.2.2	Archimedean copulas	12
2.2.2.1	Gumbel-Hougaard family	14
2.2.2.2	Clayton family	14
2.2.2.3	Ali-Mikhail-Haq Family	15
2.2.2.4	Frank Family	15
2.2.3	Elliptical copulas	15
2.2.3.1	Gaussians copula	16
2.2.3.2	Students's $t$ copula	17
2.3	Parameter estimation methods	18
2.3.1	Parametric approaches	18
2.3.1.1	Maximum likelihood estimation	19
2.3.1.2	Inference function of margins	19

	2.3.1.3	Maximization by parts	20
	2.3.2	Semiparametric approaches	23
	2.3.2.1	Pseudo maximum likelihood	23
	2.3.2.2	Maximum likelihood based on kernel density estimation	24
	2.3.2.3	Inversion of the rank correlation coefficient approach based on Kendall's tau and Spearman's rho	26
	2.4	Previous comparative studies on copula parameter estimation methods	27
	2.5	Hydrologic analysis studies using copulas	29
	2.6	Summary	31
<b>3</b>		<b>METHODOLOGY</b>	<b>36</b>
	3.1	Introduction	36
	3.2	The research theoretical framework	36
	3.3	Copula model	38
	3.4	Parameter Estimation Methods	39
	3.4.1	Parametric approaches	40
	3.4.1.1	Maximum likelihood estimation	40
	3.4.1.2	Inference function of margins	41
	3.4.1.3	Maximization by parts	43
	3.4.2	Semiparametric Approaches	45



	3.4.2.1	Pseudo maximum likelihood	45
	3.4.2.2	Maximum likelihood based on kernel density estimation	46
	3.4.2.3	Inversion of the rank correlation coefficient based on Kendall's tau	50
	3.4.2.4	Inversion of the rank correlation coefficient based on Spearman's rho	51
3.5		Simulation study	54
	3.5.1	Case 1: No misspecification in the marginal distribution	55
	3.5.2	Case 2: Exists a misspecification in the marginal distribution	57
3.6		Empirical study	59
	3.6.1	Marginal distribution function	60
	3.6.2	Goodness of fit (GOF) test	61
3.7		Summary of the research methodology	62
<b>4</b>		<b>DATA ANALYSIS AND DISCUSSION</b>	<b>64</b>
	4.1	Introduction	64
	4.2	Simulation study	64
	4.2.1	Performance comparison of parametric and semiparametric estimation methods for copula in terms of precision based on RMSE for $\theta$	66

4.2.1.1	Case 1: No misspecification in the marginal distributions	67
4.2.1.2	Case 2: Exist a misspecification in the marginal distributions	75
4.2.2	Performance comparison of parametric and semiparametric estimation methods for copula based on the computational performance	78
4.2.2.1	Computational performance between the parametric methods	82
4.2.2.2	Computational performance between the semiparametric methods	87
4.3	Empirical study	87
4.3.1	Study area and data	88
4.3.2	The correlation level between the rainfalls data	90
4.3.3	Marginal distributions of the daily rainfall data	91
4.3.4	Joint daily rainfall data by Copula method	92
4.4	Summary	96
<b>5</b>	<b>CONCLUSION AND RECOMMENDATION</b>	<b>97</b>
5.1	Introduction	97
5.2	Conclusion	98
5.3	Recommendation for future research	101

**REFERENCES**

**102**

Appendices A - B

107-109

## LIST OF TABLES

TABLE NO.	TITLE	PAGE
2.1	Summarization of previous studies on hydrologic analysis using the copula method	32
2.2	Summarization of previous comparative studies on copula parameter estimation methods (part 1).	33
2.3	Summarization of previous comparative studies on copula parameter estimation methods (part 2).	34
3.1	The properties of the Archimedean and elliptical copulas	38
3.2	The Relationship of Spearman's rho ( $\rho$ ) and Kendall's tau ( $\tau$ ) with the copula families	53
4.1	Values of $\theta$ corresponding to Kendall's $\tau$ value	65
4.2	Rank of the seven parameter estimation method for Case 1 based on the RMSE of $\theta$ corresponding to sample size $n = 50$ .	71
4.3	Rank of the seven parameter estimation method for Case 1 based on the RMSE of $\theta$ corresponding to sample size $n = 100$ .	71

4.4	Rank of the seven parameter estimation method for Case 1 based on the RMSE of $\theta$ corresponding to sample size $n = 1000$ .	72
4.5	Rank of the seven parameter estimation method for Case 1 based on the RMSE of $\theta$ corresponding to sample size $n = 5000$ .	72
4.6	Rank of the seven parameter estimation method for Case 2 based on the RMSE of $\theta$ corresponding to sample size $n = 50$ .	76
4.7	Rank of the seven parameter estimation method for Case 2 based on the RMSE of $\theta$ corresponding to sample size $n = 100$ .	76
4.8	Rank of the seven parameter estimation method for Case 2 based on the RMSE of $\theta$ corresponding to sample size $n = 1000$ .	77
4.9	Rank of the seven parameter estimation method for Case 2 based on the RMSE of $\theta$ corresponding to sample size $n = 5000$ .	77
4.10	Average time spent (second) by each method for parameter estimation.	79
4.11	Average of the absolute value of the difference between updated estimator, $\theta_k$ and the previous estimator, $\theta_{k-1}$ for 500 repetitions.	84
4.12	Descriptive statistics of the daily rainfall for Station A and Station B	89
4.13	Kendall's tau correlation and $p$ -values for the rainfall data	90

4.14	Test of goodness-of-fit for marginal distribution based on the AIC result.	91
4.15	Estimated marginal parameters for Gamma distribution of daily rainfall data of Station A and Station B.	92
4.16	The estimators of the dependence parameter.	93
4.17	Test of goodness-of-fit for copula function based on the AIC result.	94
4.18	The copula estimator of Student's $t$ copula and the estimated AIC by the parametric and semiparametric methods.	94

## LIST OF FIGURES

FIGURE NO.	TITLE	PAGE
2.1	Kernel density estimate (KDE) a random sample. $n = 100$ from a Chi-square distribution with 5 degrees of freedom with different smoothing bandwidths, $h$ .	25
3.1	The theoretical framework for this research	37
3.2	The procedures to compare the performance of the estimation methods for Case 1.	56
3.3	The procedures to compare the performance of the estimation methods for Case 2.	58
3.4	The summary framework of this research methodology.	63
4.1	The ranking for each parameter estimation method based on the RMSE from Table 4.2.	73
4.2	The ranking for each parameter estimation method based on the RMSE from Table 4.3.	73
4.3	The ranking for each parameter estimation method based on the RMSE from Table 4.4.	74
4.4	The ranking for each parameter estimation method based on the RMSE from Table 4.5.	74

4.5	Average time spent by each method for parameter estimation of all correlation levels with sample size $n = 50$ .	80
4.6	Average time spent by each method for parameter estimation of all correlation levels with sample size $n = 100$ .	80
4.7	Average time spent by each method for parameter estimation of all correlation levels with sample size $n = 1000$ .	81
4.8	Average time spent by each method for parameter estimation of all correlation levels with sample size $n = 5000$ .	81
4.9	The average of the distance between the updated estimator and the previous estimator for 500 repetitions with $n = 50$ .	85
4.10	The average of the distance between the updated estimator and the previous estimator for 500 repetitions with $n = 100$ .	85
4.11	The average of the distance between the updated estimator and the previous estimator for 500 repetitions with $n = 1000$ .	86
4.12	The average of the distance between the updated estimator and the previous estimator for 500 repetitions with $n = 5000$ .	86
4.13	The location of the selected Station A and B with their respective neighbouring rain gauge stations.	88
4.14	Scatter plot of the daily rainfall data from Station A and B in millimeter unit (mm)	89



**LIST OF SYMBOLS**

$f_X(x)$	-	Probability Density Function for random variable $x$ .
$f_Y(y)$	-	Probability Density Function for random variable $y$ .
$F_X(x)$	-	Cumulative Distribution Function for random variable $x$ .
$F_Y(y)$	-	Cumulative Distribution Function for random variable $y$ .
$F_X^{-1}(u)$	-	Inverse of the Cumulative Distribution Function for random variable $x$ .
$F_Y^{-1}(v)$	-	Inverse of the Cumulative Distribution Function for random variable $y$ .
$\alpha$	-	Corresponding parameter of $f_X(x)$
$\beta$	-	Corresponding parameter of $f_Y(y)$
$\hat{\alpha}$	-	Estimator for parameter $\alpha$
$\hat{\beta}$	-	Estimator for parameter $\beta$

$C(u, v)$	-	Copula Function
$D_k(x)$	-	Debye Function
$\Phi_\theta$	-	The bivariate standard normal CDF
$\Phi_\theta^{-1}$	-	The inversed of univariate normal CDF
$t_\theta$	-	The bivariate standard Student's $t$ CDF
$t_\theta^{-1}$	-	The inversed of the univariate Student's $t$ CDF
$df$	-	Degree of freedom
$n$	-	Sample size
$\rho$	-	Spearman's rho
$\tau$	-	Kendall's tau
$\theta$	-	Copula Parameter
$\hat{\theta}$	-	Estimator for Copula Parameter
$\phi(t)$	-	Generator of Archimedean Copula
$\phi^{-1}(s)$	-	Inverse generator
$\ell_c(\alpha, \beta, \theta)$	-	Log-likelihood of the copula density function

$\ell_m(\alpha, \beta)$	-	Log-likelihood of marginal density functions
$a$	-	Shape parameter
$b$	-	Scale parameter
$\lambda$	-	Rate parameter
$K(y)$	-	Kernel function
$h$	-	Smoothing parameter for kernel density estimation.
$f''(y)$	-	Second derivative of $f(y)$
$o(y)$	-	Little o notation that shows the asymptotic behaviour of the given function.
$\hat{f}_{-i}(y_i)$	-	kernel estimator called as leave-one-out density estimator estimated from the data without the observation $y_i$

**LIST OF ABBREVIATIONS**

AIC	-	Akaike Information Criterion
PDF	-	Probability Density Function
CDF	-	Cumulative Distribution Function
GOF	-	Goodness-of-fit
IFM	-	Inference Function of Margins
MBP	-	Maximization by parts
ML	-	Maximum Likelihood
MLE	-	Maximum Likelihood Estimation
PML	-	Pseudo Maximum Likelihood
MLKDE	-	Maximum likelihood based on kernel density estimation.
ROT	-	Rule of thumb
LSCV	-	Least square cross validation

SJ	-	Sheater and Jones
MLKDE-rot	-	Maximum likelihood based on kernel density estimation by using bandwidth selector rule of thumb.
MLKDE-sj	-	Maximum likelihood based on kernel density estimation by using bandwidth selector Sheater and Jones plug-in.
MLKDE-lscv	-	Maximum likelihood based on kernel density estimation by using bandwidth selector least squares cross-validation.
MISE	-	Mean integrated squared error
AMISE	-	Approximate mean integrated square error
SE	-	Squared error
RMSE	-	Root mean squared error
iTAU	-	The inversion of the rank correlation coefficient Kendall's tau.
iRHO	-	The inversion of the rank correlation coefficient Spearman's rho.
NaNs	-	Undefined values
CV	-	Coefficient of variation

**LIST OF APPENDICES**

<b>APPENDIX</b>	<b>TITLE</b>	<b>PAGE</b>
A	R Code: Parametric Approaches	107
B	R Code: Semiparametric Approaches	109

## CHAPTER 1

### INTRODUCTION

#### 1.1 Background of the study

Hydrological phenomena such as drought, flood and rainfall are one of the natural phenomena that often provide dependent multivariate observations. For example, from a drought phenomenon, we can observe the drought duration, intensity, and magnitude. According to Salvadori and De Michele (2007), those random variables play an important role between each other, where such an analysis of a joint distribution between the variables can identify the characteristics of drought. Therefore, it is necessary to find the joint distribution and estimate the dependence between the variables.

Based on the traditional approach, the joint distribution has been described using bivariate or multivariate distribution functions such as bivariate gamma, bivariate normal or multivariate normal distribution. However, there are limitations to this approach which made it difficult to execute mathematically. The marginal distributions must belong to the same family of the joint distribution function and the marginal parameters may affect the dependence between the variables. In addition, Salvadori and De Michele (2007) stated that canonical Pearson's coefficient of linear correlation,  $\rho$  is usually used as the dependence parameter between variables in

hydrology process. However, the canonical Pearson's coefficient may show  $\rho = 0$  in some cases which means there is no dependency between the variables although the variables are obviously not independent. This is because it only shows a linear dependence. Thus, to overcome this problem, a flexible method called copula method is introduced.

Copula method was introduced by Sklar (1959). A copula function is a joint distribution function of a combination of two or more uniform marginal distributions. This method can overcome the limitations of the traditional approach because it allowed us to specify any distribution function to the marginal distributions and then choose any copula to construct the dependence structure of the variables. In the work of Zhang and Singh (2007), they have proved that the copula method is able to derive bivariate joint distributions of rainfall variables that have different marginal distributions and without assuming the variables to be normal or independent. Many different copula families that are able to cover a wide scope of dependence structures have been proposed and developed, for example, Archimedean, Gaussian, and Student's  $t$  copula families. Further information on copula families is discussed in Chapter 2 of this study.

In hydrologic application, the most copula families that have been used for analysis are Archimedean copula families. According to Nelsen (2006) and Zhang and Singh (2007), Archimedean copulas that usually have closed form are very popular and desirable in constructing the dependence structure of the hydrologic variables. It is because of the ease in constructing the functions and they can be applied when the variables correlation is either positive or negative.

Archimedean copulas are divided into two groups, symmetric and asymmetric copulas as mentioned by Chen et al (2013). The symmetric Archimedean family is directed by one dependence parameter,  $\theta$ . They stated that the limitation of the symmetric Archimedean copula is that it can only measure one dependence structure between two variables, where all possible pairs of variables that can be paired up will



have the same dependence structure. Thus, symmetric Archimedean copulas are only suitable for structuring the dependence of two variables, but inadequate for more than two variables. To overcome the limitation of the symmetric Archimedean copula, an asymmetric copula is constructed. The asymmetric copula is a nested form of the symmetric copula. Consequently, the asymmetric copula will be able to describe different dependence structures between two or more variables. Other than Archimedean copulas, elliptical copulas such as Gaussian and Student's  $t$  copula also have been widely used for an analysis of multivariate hydrologic variables. Elliptical copula family is implicit copulas where they do not have a closed form.

To determine a specified copula structure that fitted with the marginal variables, the parameters of the copula function need to be estimated first. There are many parameter estimation methods have been proposed and developed for estimating the dependence parameter of the copula. These methods are classified into three categories, parametric approaches, semiparametric approaches, and nonparametric approaches. For hydrological analysis, parametric and semi-parametric approaches are the most common estimation methods that have been used to estimate the copula parameter. However, the nonparametric method is very rarely used because no specific parametric forms are assumed for either the copula or the marginal distributions and the copula is estimated based on empirical distributions by simple observation on the data sets. Therefore, determining an empirical copula relies on the amount of the available observation data and the formation of an empirical copula depended on a large amount of the data which is one of the limitations in hydrologic application.

In parametric approaches, the marginal distributions are assumed to follow a parametric distribution. The parameters of interest are marginal parameters and copula dependence parameter. Parametric methods are popular because they estimate the estimator precisely. However, they have a weakness against a misspecified marginal parametric distribution. For that reason, the semiparametric approach is implemented to overcome the problem by assuming the marginal distributions to be nonparametric, which allow marginal empirical cumulative distribution functions be plugged into the

marginal functions. Thus, the copula dependence parameter is the only parameter of interest in the semiparametric approach.

In the hydrological analysis, the most common parameter estimation methods that have been used are Maximum likelihood (ML) estimation and Inference Function of Margins (IFM) for parametric approaches. Whereas, pseudo maximum likelihood (PML) and rank correlation coefficient of Kendall's tau and Spearman's rho methods have been used for semiparametric approaches. Among these five methods, Kendall's tau method is the most popular method for estimating bivariate copula probably because it has a closed form of one-to-one relationship between rank correlation, tau ( $\tau$ ) and the copula parameter,  $\theta$  which has made the estimation process become easier. Vandenberghe et al. (2010) and Chen et al. (2015) also preferred to use Kendall's tau method than ML estimation or PML because it is easier to estimate the copula parameter based on Kendall's tau rank correlation coefficient rather than finding the fitted marginal distributions and maximizing a log-likelihood function that leads to a complicated algorithm.

Other than five estimation methods mentioned above, there are two other copula estimation methods that have been developed, which are maximization by parts (MBP) under the parametric approach and maximum likelihood based on kernel density estimation (MLKDE) under semiparametric approach. Song et al. (2005) proposed MBP to overcome some loss made by IFM. Meanwhile, MLKDE has the same structure as PML, where the difference in MLKDE is the marginal distributions are estimated by kernel density estimation. There are large research of hydrological studies that use ML estimation, IFM, PML, the inversion of rank correlation coefficient approach based on Kendall's tau and Spearman's rho to estimate the copula dependence parameter. However, studies that implement MBP and MLKDE are rare to find in hydrologic application literature.

## 1.2 Problem statement

Recently, there has been an increase of interest in joining distribution functions of multivariate hydrologic observations using the copula method. Copula method is able to assess the relation between the variables without concerning a specific marginal distribution. The copula is estimated using parametric and semiparametric approaches. The most common methods that have been used in estimating the copula parameter are Maximum likelihood (ML) estimation and Inference Function of Margins (IFM) for parametric approaches and pseudo maximum likelihood (PML), and the rank correlation coefficient approach based on Kendall's tau and Spearman's rho for semiparametric approaches.

Although ML estimation, IFM, PML, and the inversion of the rank correlation coefficient approach based on Kendall's tau and Spearman's rho have been widely used for hydrologic analysis, there are limited comparative studies that focus on copula estimation methods in a hydrologic application. In addition, there are two other copula estimation methods that have been developed but rarely used in hydrologic analyses, the methods are maximization by parts (MBP) under the parametric approach and maximum likelihood based on kernel density estimation (MLKDE) under semiparametric approach. These seven estimation methods have different steps and techniques to estimate the parameter. Therefore, a comparison study is important to evaluate the performance of the estimation methods. This study is conducted to compare the precision and the performance of seven parameter estimation methods for copula in a hydrologic application.

### 1.3 Research questions

The problem statement raises several research questions. The questions are listed as follow:

- i. How to estimate the copula dependence parameter,  $\theta$  using parametric and semiparametric estimation methods?
- ii. What is the performance of parametric and semiparametric estimation in terms precision based on their value of root mean square error (RMSE)?
- iii. Which parameter estimation methods that are suitable and efficient for estimating the dependence parameter of hydrologic variables?

### 1.4 Objectives of the study

The objectives of this study are listed as follows:

- i. To estimate the copula dependence parameter,  $\theta$  using parametric and semiparametric estimation methods.
- ii. To evaluate and to compare the performance of parametric and semiparametric estimation methods for copula in terms of efficiency and precision.
- iii. To identify the estimation methods that are suitable, efficient, and precise in estimating the dependence parameter of hydrologic variables.

## 1.5 Scopes of the study

This study focuses on the application of parametric and semiparametric approaches in estimating the copula dependence parameter. The performance of seven parameter estimation methods namely, maximum likelihood (ML) estimation, inference function of margins (IFM), maximization by parts (MBP), pseudo maximum likelihood (PML), the inversion of rank correlation coefficient approach based on Kendall's tau and Spearman's rho and maximum likelihood based on kernel density estimation (MLKDE) are compared in the simulation and empirical studies. The simulation and empirical studies are limited to the case of bivariate copulas.

In the simulation study, simulation data are generated from Clayton copula as the true copula with four different values of true copula parameter dependence that are corresponding to Kendall's tau,  $\tau = 0.20, 0.50, 0.60,$  and  $0.80$ . The sample sizes of the generated data are set to  $n = 50, 100, 1000,$  and  $5000$ . 500 repetitions of data generation and estimation process are done for each combination of different data sample size,  $n$  and copula dependence level,  $\theta$ .

While, for the empirical study, rainfall data are used as the empirical data. The data are selected from two Kelantan rain gauge stations which are located in the north-east of Peninsular Malaysia. The selected rain gauge stations are Station Kuala Krai, 5522047 (Station A) and Station Ulu Sekor, 5520001 (Station B). Three types of marginal distributions are considered in fitting the hydrologic variables: Gamma, Weibull and Exponential distributions. The marginal information is used in the estimation process done by the parametric approach. For the joint distribution function, six copulas that are usually used in the hydrologic application are selected. The copulas are Gumbel, Clayton, Frank, Ali-Mikhail-Haq, Gaussian, and Student's  $t$  copulas.

## 1.6 Significance of the study

The analysis in this study involves the implementation of copula method in combining the hydrologic variables and estimating the copula dependence parameter,  $\theta$ . This study allows the characteristics of marginal distributions, parameter estimation methods, and copula families to be recognized. Furthermore, the main significance of this study is it will become as another example of application and comparative study of parametric and semiparametric estimation approaches in estimating the dependence parameter of copulas model since there are only a few previous studies that compare the parameter estimation methods for copulas in a hydrologic application. Other parameter estimation methods such as MBP and MLKDE which are rarely used in hydrologic analyses are also discovered in this study.

In the simulation study, the research process leads in developing the methodology for simulating data of marginal distributions based on the given true marginal and copula distributions and the true value of the dependence parameter. The simulation process allows the generation of  $n$  sample sizes of data and desirable repetitions of estimation process for each combination of different data sample size,  $n$  and copula dependence level,  $\theta$ .

In addition, the performance of the methods based on the measured root of mean square error (RMSE) comparison can give statistical evidence in choosing which the parameter estimation methods that are more accurate and efficient to estimate the copula dependence parameter. This is important because the copula dependence parameter will affect the precision in estimating the copula function that is fitted to the data. This study also provides the result of computational performance for the seven estimation methods.

## 1.7 Research outline

This dissertation report consists of five chapters. Chapter 1 starts with the introduction of copula and parameter estimation methods in the background of the study. Then, it is followed by the statement of the problem and the questions that arise in the problem statement. After that, the purpose or the objectives of the study and the scope that are used in the study are highlighted. Finally, the possible significance or contributions that the study can provide are also presented in this chapter.

Chapter 2 consists of the general review about copula function and the expression of some copula families that have been widely used in hydrologic analyses. Some parameter estimation methods that can be used to estimate the copula dependence parameter are reviewed and the algorithm or mathematical formulation of some estimation methods are also presented. In addition, previous comparative studies that focus on copula parameter estimation methods and hydrologic analysis studies using copulas are discussed.

Chapter 3 describes the research methodology. It consists a brief explanation about the simulation and empirical studies for comparing the performance and efficiency of the parametric and semiparametric estimation methods. The procedures that are used for the both studies are also explained. The steps that involved in the parametric and semiparametric estimation methods are also described in this chapter.

Chapter 4 presents all the results and findings of the simulation and empirical studies. In this chapter, the performance of the seven copula estimation methods is compared based on the measured root mean squared error (RMSE) and time spend of each method. Finally, in the last chapter, Chapter 5, concludes the research project based on the results and findings from the simulation and empirical studies. Some recommendations that need to be done for further research are also suggested.

## REFERENCES

- Ali M. M., Mikhail N. N., and Haq M. S. (1978). A Class of Bivariate Distributions Including the Bivariate Logistic. *Journal of Multivariate Analysis*. 8, 405-412.
- Ariff N.M., Jemain A.A., Ibrahim K., and Wan Zin W.Z. (2012). IDF relationships using bivariate copula for storm events in Peninsular Malaysia. *Journal of Hydrology*. 470-471, 158–171.
- Bouyé E., Gausse N., and Salmon M. H. (2001). Investigating Dynamic Dependence Using Copulae.
- Bowman A. W. (1984). An alternative method of cross-validation for the smoothing of kernel density estimates. *Biometrika*. 71, 353-360.
- Brahimi B., and Necir A. (2012). A semiparametric estimation of copula models based on the method of moments. *Statistical Methodology*. 9, 467-477.
- Charpentier A., Fermanian J., and Scaillet O. (2007). The estimation of copulas: theory and practice. Copulas: from theory to application in finance. Risk Books, London.
- Chen L., Singh V. P., and Guo S. (2013). Measure of Correlation between River Flows Using the Copula-Entropy Method. *Journal of Hydrologic Engineering*. 18, 1591-1606.
- Chen L., Singh V. P., Guo S., Zhou J., and Zhang J. (2015). Copula-based method for multisite monthly and daily streamflow simulation. *Journal of Hydrology*. 528, 369–384.



- Cherubini U., Luciano E., and Vecchiato W. (2004). *Copula Methods in Finance*. Wiley, Chichester, UK.
- Clayton D. G. (1978). A Model for Association in Bivariate Life Tables and Its Application in Epidemiological Studies of Familial Tendency in Chronic Disease Incidence. *Biometrika*. 65(1), 141-151.
- Dupuis D. J. (2007). Using Copulas in Hydrology: Benefits, Cautions, and Issues. *Journal of Hydrologic Engineering*. 12(4), 381-393.
- Falahi M., Karamouz M., and Nazif S. (2012). *Hydrology and hydroclimatology: principles and applications*. Taylor & Francis.
- Fermanian J. D., and Scaillet O. (2003). Nonparametric estimation of copulas for time series.
- Frank M. J. (1979). On the simultaneous associativity of  $F(x, y)$  and  $x + y - F(x, y)$ . *Aequationes Mathematicae*. 19, 194-226.
- Fu G., and Butler D. (2014). Copula-based frequency analysis of overflow and flooding in urban drainage systems. *Journal of Hydrology*. 510, 49-58.
- Genest, C., Ghoudi, K., and Rivest, L. P. (1995). A semiparametric estimation procedure of dependence parameters in multivariate families of distributions. *Biometrika*. 82(3), 543-552.
- Gumbel E. J. (1960). Distributions del valeurs extremes en plusieurs dimensions. *Publications de l'Institut de Statistique de L'Université de Paris*. 9, 171-173.
- Jaki T. and West R. W. (2008). Maximum Kernel Likelihood Estimation. *Journal of Computational and Graphical Statistics*. 17(4), 976-993.
- Joe H. and Xu J. J. (1996). The Estimation Method of Inference Functions for Margins for Multivariate Models. *Technical Report no. 166*. Department of Statistics, University of British Columbia.

- Joe H. (2005). Asymptotic Efficiency of the Two-Stage Estimation Method for Copula-Based Models. *Journal of Multivariate Analysis*. 94, 401-419
- Jou P. H., Akhoond-Ali A. M., Behnia A. and Chinipardaz R. (2008). Parametric and Nonparametric Frequency Analysis of Monthly Precipitation in Iran. *Journal of Applied Sciences*. 8, 3242-3248.
- Kim G., Silvapulle M., and Silvapulle P. (2007). Comparison of Semiparametric and Parametric Methods for Estimating Copulas. *Computational Statistics & Data Analysis*. 51, 2836-2850.
- Kim K. D., and Heo J. H. (2002). Comparative study of flood quantiles estimation by nonparametric models. *Journal of Hydrology*. 260, 176-193.
- Kim T. W., Valdés J. B., and Yoo C. (2006). Nonparametric Approach for Bivariate Drought Characterization Using Palmer Drought Index. *Journal of Hydrologic Engineering*. 11(2), 134-143.
- Kojadinovic I., and Yan J. (2010). Comparison of Three Semiparametric Methods for Estimating Dependence Parameters in Copula Models. *Insurance: Mathematics and Economics*. 47, 52–63.
- Lawless J. F., and Yilmaz Y. E., (2011). Comparison of Semiparametric Maximum Likelihood Estimation and Two-Stage Semiparametric Estimation in Copula Models. *Computational Statistics and Data Analysis*. 55, 2446–2455.
- Nelsen R.B., 2006. An Introduction to Copulas. Springer, New York.
- Parzen E. (1962). On the estimation of a probability density and the mode. *The Annals of Mathematical Statistics*. 33, 1965-1976
- Quintela-del-Rio A. (2011). On bandwidth selection for nonparametric estimation in flood frequency analysis. *Hydrological Processes*. 25, 671-678.
- Requena A. I., Mediero L., and Garrote L. (2013). A Bivariate Return Period Based on Copulas For Hydrologic Dam Design: Accounting for Reservoir Routing in Risk Estimation. *Hydrology and Earth System Sciences*. 17, 3023-3038

- Rosenblatt M. (1956). A central limit theorem and a strong mixing condition. *Proceedings of the National Academy of Sciences of the United States of America*. 42, 43-47
- Rudemo M. (1982). Empirical choice of histograms and kernel density estimators. *Scandinavian Journal of Statistics*. 9, 65-78.
- Salvadori G., De Michele C., Kottegoda N. T., and Rosso R. (2007). Extremes in Nature: An Approach Using Copulas. *Water Science and Technology Library*. 56. Springer Netherlands
- Sharma A., Lall U., and Tarboton D. G. (1998). Kernel bandwidth selection for a first order nonparametric streamflow simulation model. *Statistic Hydrology and Hydraulics*. 12, 33-52.
- Sheater S. J. (2004). Density Estimation. *Statistical Science*. 19(4), 588-597.
- Sheater S. J. and Jones M. C. (1991). A reliable data-based bandwidth selection method for kernel density estimation. *Journal of the Royal Statistical Society. Series B: Methodological*, 53: 683-690.
- Silverman B.W. (1986). Density Estimation for Statistics and Data Analysis. 1st Edition. Chapman and Hall, London.
- Sklar A. (1959). Fonctions de répartition à n dimensions et leurs marges. *Publications de l'Institut de Statistique de L'Université de Paris*. 8, 229-231.
- Song P. X. K., Fan Y., and Kalbfleisch J. D. (2005). Maximization by Parts in Likelihood Inference. *Journal of the American Statistical Association*. 100(472), 1145-1158.
- Tsukahara H. (2005). Semiparametric Estimation in Copula Models. *The Canadian Journal of Statistics*. 33, 357-375.
- Vandenberghe S., Verhoest N. E. C., and De Baets B. (2010). Fitting bivariate copulas to the dependence structure between storm characteristics: A detailed analysis based on 105 year 10 min rainfall. *Water Resources Research*. 46, W01512.

- Yee, K. C., Jamaludin S., Yusof F., and Mean F. H. (2014). Bivariate copula in fitting rainfall data. *AIP Conference Proceedings*. 1605(1), 986-990.
- Yusof F., Mean F. H., Jamaludin S., and Yusof Z. (2013). Characterization of Drought Properties with Bivariate Copula Analysis. *Water Resource Manage.* 27, 4183-4207.
- Zhang L., and Singh V. J. (2007). Bivariate rainfall frequency distributions using
- Zhang R., Czadoa C., and Min A. (2011). Efficient maximum likelihood estimation of copula based meta *t*-distributions. *Computational Statistics & Data Analysis*. 55(3), 1196-1214.