

Addis Ababa  
University  
(Since 1950)



## **Statistical Analysis of Rainfall pattern in Dire Dawa, Eastern Ethiopia**

**REDIAT TAKELE**

**A thesis submitted to  
the Department of Statistics**

**Presented in Partial Fulfillment of the Requirement for the Degree of Master  
of Science in Statistics (Bio-Statistics)**

**Addis Ababa University**

**Addis Ababa, Ethiopia**

**December, 2012**

Addis Ababa University

School of Graduates Studies

This is to certify that the thesis prepared by Rediat Takele, entitled: Statistical Analysis of Rainfall Pattern in Dire Dawa, Eastern Ethiopia and Submitted in partial fulfillment of the Requirement for the Degree of Master of Science in Statistics (Bio-Statistics) complies with requirements of the University and meets the accepted standards with respect to originality and quality.

Signed by Examining committee:

Examiner \_\_\_\_\_ Signature \_\_\_\_\_ Date \_\_\_\_\_

Examiner \_\_\_\_\_ Signature \_\_\_\_\_ Date \_\_\_\_\_

Advisor \_\_\_\_\_ Signature \_\_\_\_\_ Date \_\_\_\_\_

---

Chair of Department or Graduate Program Coordinator

# Abstract

## **Statistical Analysis of Rainfall pattern in Dire Dawa, Eastern Ethiopia**

**REDIAT TAKELE**

**Addis Ababa University, 2012**

This paper investigated the rainfall pattern of Dire Dawa, Eastern Ethiopia. Descriptive Statistics, spectrum analysis, cross-spectral analysis as well as univariate Box-Jenkins methodology to build Seasonal ARIMA model were used with the objective to analyze rainfall pattern in Dire Dawa for the period from January, 1982-December, 2011 based on data from Dire Dawa and adjacent stations: Dengego and Haramaya. Descriptive Statistics result shows that the mean annual rainfall of Dire Dawa, Dengego and Haramaya are 611mm, 774 mm and 772mm, respectively. The amount of rainfall at Dengego and Haramaya are more or less the same on average in all seasons, and are much higher than that of Dire Dawa over the study period. The variability of annual rainfall in Dire Dawa during the last 30-year period is a bit larger than neighboring station's rainfall (Dengego and Haramaya), indicating that climate instability is high in Dire Dawa than other stations. A time series model for Dire Dawa station was adjusted, processed, diagnostically checked and lastly an ARIMA (5, 0, 0)\*(0, 1, 1)<sub>12</sub> model is established and this model is used to forecast two years monthly rainfall value. The results indicate that relatively there is a tendency of increasing trend for forecasted rainfall values. Spectrum analysis result showed that a rainfall extreme event recurs every 2.5 years in Dire Dawa. Further results indicated that rainfall pattern of Dengego and Dire Dawa are found to be more related.

## *Acknowledgment*

*First and for most, I would like to extend my unshared thanks to the almighty God for providing me the opportunity for what I have achieved and for his mercy.*

*I wish to record my sincere thanks and appreciation to my advisor, Dr. Butte Gotu, for his unreserved hospitality and abiding patience. The value of his dedication, kind and untiring guidance and warmhearted advice throughout the execution of this thesis is not only difficult to estimate but also very hard to express adequately by the usual terms of acknowledgement. Without his assistance, the production of this paper would not have been possible.*

*I am indebted to Dr. Solomon Harrar, University of Montan for his valuable comments and suggestion.*

*Finally, I would like to express my gratitude to the National meteorological Agency (NMA) of Ethiopia for giving me the data for study.*

*I am also grateful to my family especially for my father Takele Figa for your painstaking support through all the years of sacrifice I subjected you to.*

# Table of Contents

<i>LIST OF TABLE</i> .....	<i>vii</i>
<i>LIST OF FIGURE</i> .....	<i>viii</i>
<i>ACRONYMS</i> .....	<i>ix</i>
<b>1. Introduction</b> .....	<b>1</b>
1.1. Background of the study .....	1
1.2. Statement of the problem.....	4
1.3. Objective of the study .....	6
1.3.1. General objective .....	6
1.3.2. Specific Objective.....	6
1.4. Significance of the Study .....	7
<b>2. Literature Review</b> .....	<b>8</b>
2.1. Concepts and definition .....	8
2.1.1. General.....	8
2.1.2. Rainfall Characteristics.....	9
2.1.3. The main effects of Rainfall .....	9
2.2. Empirical Literature Review .....	13
<b>3. Data and Methodology</b> .....	<b>18</b>
3.1. Data source and variable of the study .....	18
3.2. Methodology .....	18
3.2.1. Frequency domain approach.....	19
3.2.1.1. Spectral analysis.....	23
3.2.1.2. Cross- spectral Analysis.....	33

3.2.2. Time domain approach .....	36
3.2.2.1. Test for Stationarity .....	44
3.2.2.2. Handling Missing Values.....	51
3.2.2.3. Building ARIMA Models .....	52
<i>4. Data Analysis.....</i>	<i>69</i>
4.1. Descriptive Analysis .....	69
4.2. Testing Stationarity.....	71
4.3. Model Building for monthly rainfall series .....	77
4.3.1. Model Identification.....	77
4.3.2. Parameter Estimation .....	79
4.3.3. Diagnostic Checking.....	82
4.3.4. Forecasting.....	87
4.3.4.1. Forecasting accuracy Evaluation .....	88
4.4. Result from spectral analysis .....	91
4.5. Result from cross-spectral analysis.....	96
<i>5. Conclusion and Limitation of the study.....</i>	<i>97</i>
5.1. Conclusion .....	97
5.2. Limitation of the study.....	98
<i>Reference .....</i>	<i>99</i>
<i>Appendices .....</i>	<i>106</i>

## **LIST OF TABLE**

<b>Table 3.1:</b> Behavior of the ACF and PACF for ARMA Models .....	53
<b>Table 3.2:</b> Behavior of the ACF and PACF for Pure SARMA Models .....	53
<b>Table 4.1:</b> Summary of rainfall amount by Station (January, 1982- December, 2011) .....	70
<b>Table 4.2:</b> Summary of ADF unit-roots test (at level and after first seasonal and regular differencing) .....	73
<b>Table 4.3:</b> Summary of Parameter Estimates and selection criteria .....	80
<b>Table 4.4:</b> Correlations of Parameter Estimates for the fitted model .....	82
<b>Table 4.5:</b> Residual white noises check with Ljung-Box test for the fitted model .....	84
<b>Table 4.6:</b> Result of Independence, Homoscedasticity and Normality tests for residual of fitted model.....	84
<b>Table 4.7:</b> Forecasting Accuracy Statistic .....	889
<b>Table 4.8:</b> Actual and fitted values of the series (January, 2011-December, 2011).....	889
<b>Table 4.9:</b> Confidence Intervals for the Smoothed Spectra of the annual and monthly rainfall Series .....	93

## **LIST OF FIGURE**

<i>Figure-1: Time plot for first seasonal differenced transformed rainfall series (Dire Dawa) .....</i>	<i>74</i>
<i>Figure-2: ACF and PACF plot for first seasonal differenced transformed rainfall series (Dire Dawa) .....</i>	<i>75</i>
<i>Figure-3: Diagnostics of the residuals from the fitted model.....</i>	<i>86</i>
<i>Figure-4: Forecast plot for total monthly rainfall of Dire Dawa.....</i>	<i>90</i>
<i>Figure-5: Raw periodogram of annual rainfall of Dire Dawa.....</i>	<i>91</i>
<i>Figure-6: Smoothed periodogram estimate of spectrum of annual rainfall of Dire Dawa. ....</i>	<i>94</i>

## ***ACRONYMS***

ACF	Autocorrelation Function
ADF	Augmented Dickey-Fuller
AIC	Akaike information criteria
AR	Autoregressive
ARIMA	Autoregressive Integrated Moving Average
ARMA	Autoregressive Moving Average
MA	Moving Average
MM	Milliliter
NMA	National Meteorological Agency
NMSA	National Meteorological service Agency
PACF	Partial Autocorrelation Function
SARIMA	Seasonal Autoregressive Integrated Moving Average
SBC	Schwartz's Bayesian Criterion

# ***1. Introduction***

## **1.1. Background of the study**

Located within the tropics, Ethiopia has great geographical diversity with high and rugged mountainous, flat-topped plateaus, deep gorges, etc. The Great Rift Valley divides the country into two parts forming the eastern and western highlands. Its altitudinal range lies between 120m below sea level and 4600m above sea level (Admasu, 1989). The differences in altitude and relief create a large variation in climate in various regions of the country. In places that are characterized as semi-arid zones, climate shows wide fluctuation from year to year and even within seasons in the year. Semi arid regions receive very small, irregular, and unreliable rainfall (Workneh, 1987).

The annual cycle of the climatology of the rainfall over tropical Africa and in particular over Ethiopia, is strongly determined by the position of the Inter Tropic Convergence Zone (ITCZ) (Griffiths, 1971). Variations in rainfall pattern throughout the country are the result of differences in elevation and seasonal changes in the atmospheric pressure systems that control the prevailing winds. The climate of Ethiopia is characterized by high rainfall variation (Yilma et al., 1994). In Ethiopia, several regions receive rainfall throughout the year, but in some regions rainfall is seasonal and low making irrigation necessary (Alemeraw and Eshetu, 2009). Rainfall is the most critical and key variable both in atmospheric and hydrological cycle. Rainfall patterns usually have spatial and temporal variability. These variabilities affect agricultural production, water supply,

transportation, environment and urban planning, thus, the entire economy of a country, and the existence of its people. Rainfall variability is assumed to be the main cause for the frequently occurring climate extreme events such as drought and flood. These natural phenomena affect badly the agricultural production and hence the economy of the nation. In regions where the year-to-year variability is high, people often suffer great calamities due to floods or droughts. Even though damage due to extremes of rainfall cannot be avoided completely, a forewarning could certainly be useful (Nicholls, 1980).

Ethiopia is one of the countries whose economy is highly dependent on rain-fed agriculture and also facing recurring cycles of flood and drought. Current climate variability is already imposing a significant challenge to Ethiopia in general and Dire Dawa in particular, by affecting food security, water and energy supply, poverty reduction and sustainable development efforts, as well as by causing natural resource degradation and natural disasters. Metrologically, Dire Dawa Administration is characterized by an arid and semi-arid climate, thus, receives low and erratic rainfall (Bekele, 1997). Prolonged droughts time and again affected the rural part of Dire Dawa. For example, the impacts of the 2004 and 2005 droughts incident (which posed food shortage to 85% of the rural population) is still fresh in the memories of many people. Recurrent floods in the past caused substantial human life and property loss in many parts of the urban kebele. For example, the August 2006 flood claimed 256 lives, displaced 2,500 families, caused direct damage estimated at ETB 100 million and indirect damage of similar magnitude (DDA, 2006).

Methods of prediction of rainfall extreme events have often been based on studies of physical effects of rainfall or on statistical studies of rainfall time series. Rainfall forecast is relevance to the agriculture sector, since it contributes significantly to the economy of countries like Ethiopia. In order to model and predict hydrologic events, one can use stochastic methods like time series methods. Numerous attempts have been made to predict behavioral pattern of rainfall using various techniques (Yevjevich, 1972; Dulluer and Kavas, 1978; Tsakiris, 1998).

Awareness about the characteristics of the rainfall over an area such as the source, quantity, variability, distribution and the frequency of rainfall is essential for the implication in utilization and associated problems. Assessing rainfall variability is practically useful in making decision, risk management and optimum usage of water resources of countries. Thus, it is important to obtain accurate rainfall forecast at various geographic levels of Ethiopia and work towards identifying periodicities in order to help policy makers improve their decisions by taking into consideration the available and future water resources. In this study, Spectrum analysis, Cross-spectral analysis as well as univariate Box-Jenkins methodology to build Seasonal ARIMA model are used for assessing the rainfall pattern in Dire Dawa based on data from Dire Dawa and adjacent stations: Dengego and Haramaya.

## **1.2. Statement of the problem**

Rainfall variability and associated droughts have historically been major causes of food shortages and famines in Ethiopia (Wood, 1977; Pankhurst and Johnson, 1988; Ketema, 1999; Bewket and Conway, 2007). Even though rainfall variability and drought are not new phenomena in Ethiopia; the frequency of occurrence of drought has been increasing during the past decades (Ketema, 1999). Floods and droughts are considered to be the two extreme conditions of variability of rainfall. Frequent and prolonged droughts have claimed the lives of millions of people. On the other hand, major flood hazards have occurred in different parts of the country in 1988, 1993, 1994, 1995, 1996, 2006 and 2010 leading to loss of life and property. The 2006 catastrophic flood led to the destruction of huge infrastructure and the death of more than 650 people and the displacement of more than 35,000 people in Dire Dawa, South Omo and West Shewa (NMA, 2007). Similar situations were experienced over Afar, Western Tigray, Gambella and over the low-lying areas of Lake Tana.

The issue of flood continues to be of growing concern in Dire Dawa administration especially to people residing in lowlands, along or near the flood courses as well as villages located at the foot of hills and mountains. Flood disasters are occurring more frequently, and are having an ever more dramatic impact on Dire Dawa in terms of the costs on lives, livelihoods and environmental resources. On the otherhand, the Administration being located in an arid and semi-arid part has often experienced drought during recent years.

Till now, except a few studies in some regions of Ethiopia (Yilma et al., 1994; Ketema, 1999; Alamerew and Eshetu, 2009), most studies are based on assessing the effect of rainfall pattern on only crop production. However, the reality shows that the impact of rainfall variability pattern is not restricted to crop production. Therefore, there is a need for studying the impact of rainfall in detail.

So far, there is no enough accurate forecasting model developed for rainfall and a detailed study in determining the characteristics of oscillation that exist in the rainfall series in Dire Dawa region of Ethiopia that aims to help decision makers for better preparations.

The key questions to be addressed in this research are:

- What is the pattern of rainfall at key locations in the Dire Dawa area?
- What is the frequency of rainfall extreme events in the last 30 years?
- Can we infer rainfall extreme events (drought and flood) in the study area?

## **1.3. Objective of the study**

### **1.3.1. General objective**

The main objective of this study is to analyze rainfall pattern in Dire Dawa using appropriate time series methods based on 30 years (January, 1982-December, 2011) data recorded at Dire Dawa station.

### **1.3.2. Specific Objective**

1. To develop a time series model for monthly rainfall of Dire Dawa.
2. To forecast the rainfall pattern in the study area.
3. To determine the characteristics of oscillations that appears in the rainfall series with a view to predicting rainfall extreme events in the study areas.
4. To determine the relationship between rainfall pattern in Dire Dawa and adjacent area (Dengego and Haramaya).

## 1.4. Significance of the Study

Knowledge of what happens to the water that reaches the earth surface will assist the study of many surface and subsurface water problems, for efficient control and management of water resources. For a country like Ethiopia, whose welfare depends very much on rain-fed agriculture, a quantitative knowledge of water requirements of the region, availability of water for plant growth and supplemental irrigation, etc. on a monthly or seasonal basis is an essential requirement for agricultural development. In this regard, increased capacity to manage future climate change and weather extremes can also reduce the magnitude of economic, social and human damage and eventually, lead to better resistance. Assessing seasonal rainfall characteristics based on past records is essential to evaluate rainfall extreme risk and to contribute to development of mitigation strategies. Therefore, a reliable rainfall forecasting and assessing behavior at station, regional and national levels is very important. The results of this research paper will hopefully be used:

- For forecasting the pattern of rainfall in the study area.
- To provide information that would be helpful for decision makers in formulating policies to mitigate the problems of rain water resources management, soil erosion, flooding and drought. In particular, the forecasting model to be developed will be a valuable instrument for the agricultural sector.
- To provide information for the early warning system in the region.
- As a basis for further study in Dire Dawa area.

## ***2. Literature Review***

### **2.1. Concepts and definition**

#### **2.1.1. General**

Weather and climate over the earth are not constant with time: they change on different time series ranging from the geological to the diurnal through annual, seasonal and intra-seasonal time scales. Such variability is an inherent characteristic of the climate. The study of climatic fluctuations involves description and investigation of causes and effects of these fluctuations in the past and their statistical interpretation. Much of the work done is about variability of the two important meteorological parameters: rainfall and temperature.

Rainfall is the resultant product of a series of complex interactions taking place within the earth-atmosphere system. Basically, it occurs when water vapor masses condense sufficiently, driven by the cooling of air masses through upward movement. Rainfall is only water that falls from the sky, whereas precipitation is any wet things that fall from the sky, which include snow, frozen rain....etc.

Water in all its forms and in all its various activities plays a crucial role in sustaining both the climate and life. It is also a major factor for planning and management of water resource project and agricultural production. Eventhough Ethiopia enjoys a fairly good amount of rainfall, wide variability in its distribution with respect to space and time are responsible for the two extremes events (floods and droughts) (Yilma et. al,1994).

### **2.1.2. Rainfall Characteristics**

Rainfall varies with latitude, elevation, topography, seasons, distance from the sea, and coastal Sea-surface temperature. The northwestern, southwestern and eastern part of Ethiopia where elevation is up to 1750m above sea level, climates is dominated by tropical rain with mean annual temperature greater than 18<sup>o</sup>c and mean annual rainfall ranging from 680mm to 1200mm. The northeastern and southeastern part of Ethiopia are dominated by dry climate with mean annual temperature ranging from 27<sup>o</sup>c to 30<sup>o</sup>c and the mean annual rainfall less than 450mm (NMA, 1996). Usually these parts of the country are characterized by strong winds, high temperature and low relative humidity. The important weather systems that cause rainfall over Ethiopia are Sub Tropical Jet (STZ), Inter Tropical Convergence Zone (ITCZ), Red Sea Convergence Zone (RSCZ), Tropical Easterly Jet (TEJ) and the Somalia Jet.

Season is defined meteorologically as a period when an air mass characterized by homogeneous weather elements such as temperature, relative humidity, wind, rainfall etc., dominate a region or part of a country (NMA, 1996). Ethiopia has three main seasons: kiremt (long rainy season), Belg (short rainy season), and Bega (dry season).

### **2.1.3. The main effects of Rainfall**

Trends in rainfall extremes have enormous implications. Extreme rainfall events cause significant damage to agriculture, ecology, and infrastructure. They also cause disruption to human activities, injury, and loss of life. Socioeconomic activities including agriculture, power generating, water supply, human health, etc. are also very sensitive to

climate variations. As a result, Ethiopia's economy is heavily dependent on rainfall for generating employment, income, and foreign currency (Alemerew and Eshetu, 2009). Thus, rainfall is considered as the most important climatic element that influences Ethiopian agriculture. The country has experienced frequent and extensive drought that caused food shortages and famine (Wood, 1997 and Ketema, 1999). The severity and frequency of occurrence of rainfall extremes events (meteorological, hydrological, and agricultural) vary for different parts of the country.

**Drought:** Many researchers now believe that the occurrence of various droughts in Africa, especially in southern Africa and the Horn, are caused by physical processes related to the occurrence of El Niño -Southern Oscillation (ENSO) events. El Niño-Southern Oscillation (ENSO) is a coupled air and ocean phenomenon with global weather implications. When past ENSO events are compared with drought and famine periods in Ethiopia, they show a remarkable association (Wolde-Georgis, 1997, Bekele, 1997). Some drought years have coincided with EN events, while others have followed it.

According to DDAEPA (2011) the trend of decreasing annual rainfall and increased rainfall variability is contributing to drought conditions in Dire Dawa Administration. The average annual rainfall patterns of Dire Dawa for the periods 1999 to 2008 and 1984 to 1991 show two important trends. First, annual average rainfall has declined from the mean value by about 8.5% and 10% respectively. Secondly, the variability of rainfall shows an overall increasing trend, suggesting greater rainfall unreliability. These rainfall patterns have led to serious drought/flood episodes throughout the Administration.

**Flood:** Floods are known as the most frequent and devastating natural disasters in both developed and developing countries (Osti et al., 2008). Between 2000 and 2008 East Africa has experienced many episodes of flooding. Almost all of these flood episodes have significantly affected large parts of Ethiopia. Being one of the largest countries in East Africa, Ethiopia's topography characteristics has made the country pretty vulnerable to floods and resulting destruction and damage to life, economic, livelihoods, infrastructure, services and health system (FDPPA, 2007).

Flooding is common in Ethiopia during the rainy season between June and September and the major type of flooding which the country is experiencing are flash flood and river floods (FDPPA, 2007).

Like other regions of Ethiopia, the issue of flood continues to be of growing concern in Dire Dawa especially to peoples residing in lowlands, along or near the flood courses as well as village located at the foot of hills and mountains. Flood disasters are occurring more frequently, and having an ever more dramatic impact on Dire Dawa in terms of the costs on lives, livelihoods and environmental resources. The topography of Dire Dawa Administration mainly consists of mountains and hills with steep slopes, valleys, and river basins. The fact is that these sloppy areas of the administration are surrounded by the mountainous areas of the neighboring wereda. Haramaya and Kersa wereda are the main areas contributing to the disaster flood event in Dire Dawa. The catchment characteristics accompanied with its large area coverage coupled with torrential rain fall during the short and long rainy season had been the main factors that contribute to the pervious flood events.

According to DDA (2006) the following are the major flood events in Dire Dawa in the last three decades that cause great loss of human lives and property:

- ❖ The flood disaster in May 1984 – 42 people were killed, and property worth 10 million birr was lost.
- ❖ The flood disaster in August 2006 – 256 people killed, 244 people unaccounted for, and 10, 000 people made homeless and 1827 households in 17 rural Kebeles were adversely affected ; property of worth 27 million birr had been lost.
- ❖ Flood disaster in April 2010- though there was no loss of human life and live animal because of the application of early warning system, property of worth 28 million birr had been lost.

**Soil Erosion:** It is the movement of soil particles from one place to another by wind or water, which is considered to be a major environmental problem. Erosion has been going on and has produced river valleys and shaped hills and mountains. Such erosion is generally slow but can cause a rapid increase in the rate at which soil is eroded (i.e. a rate faster than natural weathering of bedrock can produce new soil). This has resulted in a loss of productive soil from crop and grazing land, as well as layers of infertile soil being deposited on formerly fertile crop lands: the formation of gullies: silting of lakes and streams, and land slips.

## 2.2. Empirical Literature Review

Several studies have taken place in the analysis of pattern and distribution of rainfall in various regions of the world. Different time series methods with different objectives are employed to analyze rainfall data in various literatures. In terms of using a formal time series model to forecast, the patterns and intensity of rainfall overtime, Harvey et al., (1987) investigate how patterns of rainfall correlate with general weather conditions and frequency of the cycles of rainfall. They used rainfall data from Brazil for a particular region which often suffers from drought to assess the cyclical behavior of rainfall. They use a model that allows cyclical components to be modeled explicitly. They found that cyclical components are stochastic rather than deterministic, and the gains achieved from forecast by taking account of the cyclic component are small in the case of Brazil.

Mahsin et al. (2012) use Box-Jenkins methodology to build seasonal ARIMA model for monthly rainfall data taken for Dhaka station, Bangladesh, for the period from 1981-2010. In their paper, ARIMA (0, 0, 1) (0, 1, 1)<sub>12</sub> model was found adequate and the model is used for forecasting the monthly rainfall. Seyed et al.,(2011) use time series method to model weather parameter in Iran at Abadeh Station and recommended ARIMA(0,0,1)(1,1,1)<sub>12</sub> as the best fit for monthly rainfall data and ARIMA(2,1,0)(2,1,0)<sub>12</sub> for monthly average temperature for Abadeh station.

Al-Ansari et al. (2003) dealt with the statistical analysis of the rainfall measurements for three meteorological stations in Jordan: Amman Airport (central Jordan), Irbid (northern Jordan) and Mafraq (eastern Jordan). Normal statistical and power spectrum analyses as

well as ARIMA model were performed on the long-term annual rainfall measurements at the three stations. The result shows that possible periodicities of the order of 2.3 - 3.45, 2.5 - 3.4 and 2.44-4.1 years for Amman, Irbid and Mafraq stations, respectively, were obtained. A time series model for each station was adjusted, processed, diagnostically checked and lastly an ARIMA model for each station is established with a 95% confidence interval and the model was used to forecast 5 years annual rainfall values for Amman, Irbid and Mafraq meteorological stations. Further result indicated that there is decreasing trend for forecasted rainfall results in all stations.

Winstanley (1973a, b) reported that monsoon rains from Africa to India decreased by more than 50% from 1957 to 1970 and predicted that the future monsoon seasonal rainfall, averaged over 5 to 10 years is likely to decrease to a minimum around 2030.

Stringer (1972) reported that at least 35 quasi-periods with more than one year in length have been discovered in records of pressure, temperature, precipitation, and extreme weather conditions over many parts of the earth surface. A very common quasi-periodic oscillation is the quasi-biennial oscillation (QBO), in which the climatic events recur every 2 to 2.5 years.

Laban (1986) uses time series methods based on ARIMA and Spectral Analysis of areal annual rainfall of two homogenous region in East Africa and recommended ARMA(3,1) as the best fit for areal indice of relative wetness\dryness and dominant quasi-periodic fluctuation around 2.2-2.8 years,3-3.7 years,5-6 years and 10-13 years.

Nicholson and Entekhabi (1986) conducted a detailed power spectrum analysis of African annual rainfall series using Blackman-Tukey and Fourier methods. Their analysis revealed that quasi-periodicities were clustered in four bands at 2.2–2.4, 2.6–2.8, 3.3–3.8 and 5.0-6.3 years, common throughout equatorial and southern Africa but only weakly evident in northern Africa.

Adejuwon (2010) studied annual rainfall in Nigeria using power spectral analysis based on Benin, Sapele, Warri and Forcados Synoptic station in Edo and Delta States (formerly Mid-Western Nigeria) over 67 years and found that Benin synoptic station shows significant spectral peaks at 6.7, 4.6 and 3.7 years periodicities. The most pronounced peak at the station was found to be 3.7 years periodicity. In Sapele, the most pronounced periodicity of 5 years was observed. Although, the spectral peaks were significant at 4.6 and 3.7 years, respectively, at Warri, the most pronounced of these peaks was found to be 3.7 years. However, in the case of Forcados, a single significant spectral peak of 3.6 years cycle was prominent and it was then concluded that periodicities were evident with significant cycles of between 3 and 6 years.

Amha (2010) studied the monthly rainfall in Tigray region based on Mekelle station. He employed univariate Box-Jenkins method to analyze rainfall in the region and found that SARIMA model is suitable for forecasting future value of monthly rainfall data and used this model to forecast 12-month rainfall pattern in the study area. Further he concluded that there is no tendency of decreasing or increasing pattern of monthly rainfall over the forecast period from January 2010 to September 2011.

Assessment made by NMA (2007) reported that the mean annual rainfall is likely to increase along Dire Dawa by 3.4 % by 2030, 6.4 % by 2050 and 10.5 % by 2080 compared to the 1961-1990.

A study made by Mersha (2002) on rainfall cyclicity over selected stations in Ethiopia also shows that there appears to be cyclic tendency in the annual rainfall data ,particularly, Gode, Dire Dawa, Negelle, and Debre Zeit station. Seifu (2004) also show that the rainfall at Dire-Dawa, D/Markos and Jijiga has periodic tendency.

Alamerew and Eshetu (2009) assess local climate of Addis Ababa. They used time domain approach like ARIMA for modeling mean minimum temperature and frequency domain time series approach ,particularly, spectral analysis for rainfall data. Using spectral analysis they determined periodicity of drought in Addis Ababa and found that the periodicity of 11.24 year is dominant cycle for the annual rainfall of Addis Ababa. Based on their result, they concluded that drought recurs in Addis Ababa region between 10 to 11 years.

Haile (1988) also found that drought occur every 6-8 years in the semi-arid regions of Ethiopia including Dire Dawa Administration. The study by Tsegay (1998) based on the occurrences of drought and the frequencies of rainfall deviation from the average also suggest that drought occur every three to five years in Eastern Ethiopia and six to eight years in northern Ethiopia and every eight to ten years for the rest of the country. Many authors (including Haile, 1988; Funk et. al., 2005) suggest that Ethiopian drought is caused by El Niño-Southern Oscillation.

According to Ketema (1999), considerable variation of rainfall pattern and distribution in Dire Dawa resulted in recurrent drought. For this fact, frequency analysis has been done to determine rainfall amount for different return periods for Dire Dawa using data covering the period 1979-98. Accordingly, the study results show that the probability of drought occurrence per year was 25%. On top of this, among causes of crop loss in Dire Dawa Administration, drought and excess water account for 32% and 4%, respectively (WWDSE, 2004).

Studies made at NMA (2007) also have shown that there is a link between El Nino and LaNina phenomena and Ethiopian rainfall.

Yilma et al. (1994) assessed the statistical link between annual rainfall of Addis Ababa and sun spot number series. They used transfer function plus noise model and spectral analysis using smoothed periodogram method to determine the periodic behavior in the sunspot number series and annual rainfall series of Addis Ababa for the period 1900-1991 and found that periodogram of Addis Ababa annual rainfall series show similar feature with sunspot series periodogram, i.e both series show similar dominant periodicity of 10.22 year. Based on their result, they concluded that an „Average“ sun activity may be associated with a contemporaneous „Average“ rainfall process, but the anomalous sun activities may later induce anomalies in rainfall (Drought and Flood) through Ocean-Atmospheric phenomenon.

# ***3. Data and Methodology***

## **3.1. Data source and variable of the study**

A time series of monthly rainfall data in mm for the period January, 1982 to December, 2011 collected by the National Meteorological Agency of Ethiopia were used in the study. The data were collected from the synoptic stations of Dire Dawa, Haramaya, and Dengego. The site was chosen due to availability of relatively long series of meteorological data.

## **3.2. Methodology**

Time series is broadly defined as any series of measurements taken at different times. Although the ordering is usually through time, particularly in terms of some equally spaced time intervals, the ordering may also be taken through other dimension, such as time series. There are various objectives for studying time series. These include the understanding and description of the generating mechanism, the forecasting of future values, and optimal control of a system. The intrinsic nature of a time series is that its observations are dependent or correlated with the order of observation. Therefore, the order of observation matters when analyzing time series data. The development of climatology as a science has given rise to growing statistical applications on climatic information. For instance, time series analysis is used in order to evaluate the temporal behavior of rainfall.

In this study a univariate Box-Jenkins Methods (Box et al., 1994), in particular, Seasonal Autoregressive Integrated Moving Average (SARIMA) methods and frequency domain approach (Spectral analysis) are employed.

### 3.2.1. Frequency domain approach

This approach examines contributions of different frequencies in explaining the variance of the series. Analysis is based on the estimated spectral density function, which provides information on the properties of the time series data. The main building blocks for analysis in the frequency domain are trigonometric function of sinusoids: Sines and Cosines.

#### Spectrum

The spectrum of a time series is the distribution of variance of the series as a function of frequency. The spectrum contains no new information beyond that in the autocovariance function (ACVF), and in fact the spectrum can be computed mathematically by transformation of the (ACVF). That relationship also means that the ACVF can be expressed as a cosine transform of the spectral density function, or spectrum. That is, if

the autocovariance function ( $\gamma(h)$ ), of stationary process satisfies  $\sum_{h=-\infty}^{\infty} |\gamma(h)| < \infty$ , then it

has the representation (Shumay and Stoffar, 2010):

$$\gamma(h) = \int_{-1/2}^{1/2} e^{2\pi i \omega h} f(\omega) d_{\omega}, \quad h = 0, \pm 1, \pm 2, \dots$$

as the inverse transform of the spectral density, which has the representation

$$f(\omega) = \sum_{h=-\infty}^{\infty} \gamma(h) e^{-2\pi i \omega h}, \quad -\frac{1}{2} \leq \omega \leq \frac{1}{2}$$

The spectrum is of interest because many natural phenomena have variability that is frequency-dependent, and understanding the frequency dependence may yield information about the underlying physical mechanisms.

Spectrum can be also defined in terms of a model in which a time series  $x_1, \dots, x_n$ , where  $n$  is odd, consists of a linear combination of many sinusoids at frequencies  $(\omega_j)$  (Shumay and Stoffer, 2010):

$$x_t = \mu + \sum_{j=1}^{(n-1)/2} [\cos(2\pi\omega_j t) dA(\omega_j) + \sin(2\pi\omega_j t) dB(\omega_j)] \quad (1)$$

for  $t=1, 2, \dots, n$ . When  $n$  is even, the representation (1) can be modified by summing to  $(n/2-1)$  and adding additional component given by  $a_{n/2} \cos(\pi t)$ .

The spectrum that is being estimated in a sense then is not really the spectrum of the observed series, but the spectrum of the unknown infinitely long series from which the observed series is assumed to have come. Then we have:

$$x_t = \mu + \int_0^{\pi} \{\cos(2\pi\omega_j t) dA(\omega_j) + \sin(2\pi\omega_j t) dB(\omega_j)\}$$

The frequencies  $(\omega_j)$  are related to the size  $(n)$  by:

$$\omega_j = j/n, \quad 1 \leq j < n/2 \quad (2)$$

The wavelength, or period, of the wave is the distance from peak to peak, and is the inverse of the frequency

$$\lambda = \frac{1}{\omega} \quad (3)$$

The peaks are the high points in the wave; the troughs are the low points. The wave varies around a mean of zero. The frequencies of the sinusoids are at intervals of  $(1/n)$  and considered as Fourier frequencies, or standard frequencies. The vertical distance from zero to the peak is called the amplitude.

In developing the definition of the spectrum fulfilling, the above eqn. (1), additional assumptions must be made: the amplitudes are random variables with expected value:

$$E\{dA(\omega)\} = E\{dB(\omega)\} = 0 \quad (4)$$

$$E\{(dA(\omega))^2\} = E\{(dB(\omega))^2\} = \delta^2 \quad (5)$$

From the relationship between variance and amplitude of sinusoid, equation (5) implies that the variance associated with the standard frequency is  $\delta_j^2$ . It must also be assumed that the amplitudes associated with various standard frequencies are uncorrelated:

$$E\{dA(\omega)dB(\omega^*)\} = 0 \quad \forall_{\omega, \omega^*} \quad (6)$$

$$E\{dA(\omega)dA(\omega^*)\} = 0 \quad \omega \neq \omega^* \quad (7)$$

With the assumptions above, it can then be shown that the expected value of  $(x_t)$  and the variance of  $(x_t)$  are given by:

$$E(x_t) = \mu \quad (8)$$

$$\text{Var}(x_t) = \delta^2 = E\{(x_t - \mu)^2\} = \sum_{j=0}^{\lfloor n/2 \rfloor} \delta_j^2 \quad (9)$$

and the autocorrelation function of  $(x_t)$  is

$$\rho_k = \frac{\sum_{j=1}^{N/2} \delta_j^2 \cos(2\pi \omega_j k)}{\sum_{j=1}^{N/2} \delta_j^2} \quad (10)$$

Equation (9) shows that the variance of the series  $(x_t)$  is the sum of the variances associated with the sinusoidal components at the different standard frequencies. Thus the variance of the series can be decomposed into components at the standard frequencies, that is, the variance can be expressed as a function of frequency.

Finally, the spectrum can be defined as:

$$f(\omega_j) = \delta_j^2, 1 \leq j \leq n/2 \quad (11)$$

A plot of  $f(\omega_j)$  against frequencies  $(\omega_j)$  shows the variance contributed by the sinusoidal terms at each of the standard frequencies. From equation (9), the variance of  $(X_t)$  can then be expressed as the sum of the spectral components as:

$$\delta^2 = \sum_{j=1}^{n/2} f(\omega_j) \quad (12)$$

The variance contributed at frequency ( $\omega_j$ ) is the spectrum  $f(\omega_j)$  at that frequency. The shape of the spectral values,  $f(\omega_j)$  plotted against  $\omega_j$  indicates which frequencies are most important to the variability of the time series.

Considering that  $(\delta_j^2)$  are spectral values, equation (11) gives an important relationship between the spectrum and the autocorrelation function (ACF): the ACF is expressed as a cosine transform of the spectrum. Therefore, transformation from the time domain to the frequency domain is made by taking the Fourier transform of the time series.

### **3.2.1.1. Spectral analysis**

The main objective of spectral analysis is to estimate and study the spectrum. In other words, spectral analysis is used to estimate the contribution of a particular band of frequencies to the overall variance in terms of a time series (Ayoade, 1973). The contributions of oscillations of various wavelengths to the variable of a time series are shown by spectrum of a time series. Spectral analysis is therefore concerned with estimating the unknown spectrum of the process from the data and quantifying the relative importance of different frequency bands to the variance of the process.

Various methods have been developed to estimate the spectrum from an observed time series. In the pre-computer era, harmonic analysis (the direct method) was used and the

results were usually displayed as a plot of amplitude against frequency known as the periodogram (Tabony, 1979). Blackman and Tukey (1959) later developed indirect approach that was computer based, while Fast Fourier transform for direct method was derived by Cooley and Tukey (1965) in which the ordinary „direct“ formulae is replaced by more computationally efficient ones. However, the two most important spectral analyses are the maximum entropy spectral analysis and the power spectral analysis.

**The maximum entropy spectral analysis**, an indirect method, is a method of analyzing time series that employs autoregressive method to extract the maximum amount of information from the available data (Burroughs, 1992). Its success as a method of assessing periodicities in a time series depends mainly on the „signal-to-noise“ ratio in the time series. As observed by Burroughs (1992), meteorological series rarely meet the signal-to-noise criteria that can exploit the advantages of maximum entropy spectral analysis which is that of not adding or subtracting information from the data. This method, therefore, will not be employed in this study for assessing periodicities.

**The Power spectral analysis** involves the presentation of the square of the amplitude of the harmonics of time series as a function of the frequency of the harmonics (Burroughs, 1992). It is a non-parametric procedure in which a finite set of observations is used to estimate a function defined over the range  $(0, \Pi)$ . The plot of the power of the variance against the frequency is known as power spectrum. The area under the curve in a power spectrum is proportional to the variance.

## **Steps in spectral analysis**

In estimating the spectrum by the smoothed periodogram the following four steps are important:

1. Subtract mean (Detrend time series)
2. Compute discrete Fourier transform (DFT)
3. Compute (raw) periodogram
4. Smooth the periodogram to get the estimated spectrum

### **Step 1: Subtract mean from each observation (Detrend time series)**

The first step in estimation of the spectrum by the smoothed periodogram method is subtraction of the sample mean. This operation has no effect on the variance. The most obvious problem with not subtracting the mean is that an abrupt offset is introduced when the series is padded with zeros in a later step in the analysis. Any obvious trend should also be removed prior to spectral estimation. Trend produces a spectral peak at zero frequency, and this peak can dominate the spectrum such that other important features are obscured.

### **Step 2: Compute discrete Fourier transform (DFT)**

Discrete Fourier transforms: Say  $(X_t)$  is an arbitrary time series of length  $n$ . The time series can be expressed as the sum of sinusoids at the Fourier frequencies of the series (Bloomfield, 2000):

$$\mathbf{x}_t = A(0) + \left\{ 2 \sum_{j=0}^{n/2} [A(\omega_j) \cos 2\pi\omega_j t + B(\omega_j) \sin 2\pi\omega_j t] \right\} + 2 \{ A(\omega_{n/2}) \cos 2\pi\omega_{n/2} t \}, t = 1, 2, \dots, n \quad (13)$$

Where the summation is over Fourier frequencies and the last term in braces is included only if  $n$  is even. The coefficients in Eqn. (13) are given by

$$A(\omega_j) = \frac{2}{n} \sum_{t=0}^{n-1} \mathbf{x}_t \cos 2\pi\omega t \quad (14)$$

$$B(\omega_j) = \frac{2}{n} \sum_{t=0}^{n-1} \mathbf{x}_t \sin 2\pi\omega t$$

Equations (13) transform the time series into two series of coefficients of sinusoids. It can also more succinctly be expressed in complex notation by making use of the Euler relation

$$e^{ix} = \cos x + i \sin x \quad (15)$$

And its inverse

$$\cos x = 1/2 \{ e^{ix} + e^{-ix} \} \quad \sin x = \frac{1}{2i} \{ e^{ix} - e^{-ix} \} \quad (16)$$

Suppose  $(X_t)$  is such a real-valued time series expressed as complex numbers. The discrete Fourier transform (DFT) of  $(X_t)$  can be written in complex notation as

$$d(\omega) = \frac{1}{n} \sum_{t=0}^{n-1} x_t e^{-2\pi i \omega t} \quad (17)$$

**Step 3: Compute raw periodogram.**

The relationship (Eqn.(13)) transforms the time series into a series of coefficients at its Fourier frequencies. The discrete Fourier transform is the complex expression of these coefficients

$$d(f) = \frac{A(\omega)}{2} - i \frac{B(\omega)}{2} \quad (18)$$

where  $A$  and  $B$  are identical to the quantities defined in Eqn. (14).

The original data can be recovered from the DFT using the inverse transform

$$x_t = \sum_j d(\omega_j) e^{-2\pi i \omega_j t} \quad (19)$$

,which is the complex equivalent of equation (13).

The discrete Fourier transform has two representations. The first is in terms of its real and imaginary parts,  $\frac{A(\omega)}{2}$  and  $-\frac{B(\omega)}{2}$  given in Eqn. (18). The second is in terms of its magnitude  $R(\omega)$  and phase  $\phi(\omega)$  as:

$$d(\omega) = R(\omega) e^{-i\phi(\omega)} \quad (20)$$

The magnitude, given by

$$R(\omega) = |d(\omega)| \quad (21)$$

measures how strongly the oscillation at frequency ( $\omega$ ) is represented in the data. The strength of the periodic component is more often represented by the *periodogram* defined as:

$$I(\omega) = n[R(\omega)]^2 = n|d(\omega)|^2$$

$$\text{where: } |d(\omega)|^2 = \frac{1}{n} \left[ \left( \sum_{t=1}^n (x_t - \bar{x}) \cos(2\pi\omega t) \right)^2 + \left( \sum_{t=1}^n (x_t - \bar{x}) \sin(2\pi\omega t) \right)^2 \right] \quad (22)$$

The sine and cosine terms at the Fourier frequencies are orthogonal, and so the variance of the time series ( $x_t$ ) can be decomposed into components at the individual frequencies. For the sine and cosine transforms, the sum of squares of the original data can be expressed as

$$\sum_{t=0}^{n-1} x_t^2 = nA(0)^2 + \left\{ 2n \sum_{j=0}^{n/2} [A(\omega_j)^2 + B(\omega_j)^2] \right\} + nA(\omega_{n/2})^2 \quad (23)$$

where the last term is included only if  $n$  is even. The analog for the discrete Fourier transform in complex notation is

$$\sum_{t=0}^{n-1} |x_t|^2 = \sum_j |d(\omega_j)|^2 = \sum_j I(\omega_j) \quad (24)$$

If ( $x_t$ ) is a time series expressed as departures from its mean, the sums of squares in equations (23) and (24) are simply  $n$  times the variance.

Equation (24) therefore indicates that (a) the sum of the periodogram ordinates equals the sum of squares of departures of the time series from its mean, (b) the sum of periodogram ordinates divided by the series length equals the series variance, and (c) the periodogram ordinate at Fourier frequency ( $\omega_j$ ) is proportional to the variance accounted for by that frequency component. The periodogram at this stage is called **raw periodogram**, meaning it has not yet been smoothed.

By considering the fact that the periodogram estimates are independent and identically approximated as  $\chi^2$ -distributed with  $\nu$ -degree of freedom, then we have (Shummay and Stoffar, 2010):

$$\frac{2I(\omega_{jn})}{f(\omega_j)} \xrightarrow{d} \text{iid } \chi_\nu^2 \quad (25)$$

where  $f(\omega_j) > 0$ , for  $j = 1, \dots, n$ .

The distributional given in Eqn. (25) can be used to derive an approximate confidence interval for the spectrum in the usual way. Let  $\chi_\nu^2(\alpha)$ , denote the lower  $\alpha$  probability tail for the chi-squared distribution with  $\nu$ - degrees of freedom, that is,

$$\Pr\{\chi_\nu^2 \leq \chi_\nu^2(\alpha)\} = \alpha \quad (26)$$

Then, an approximate  $100(1 - \alpha)$  % confidence interval for the spectrum would be of the form:

$$\frac{2I(\omega_j)}{\chi_\nu^2(1-\alpha/2)} \leq f(\omega) \leq \frac{2I(\omega_j)}{\chi_\nu^2(\alpha/2)} \quad (27)$$

#### Step 4: Smoothing the raw periodogram

The periodogram is a wildly fluctuating estimate of the spectrum with high variance. For a stable estimate, the periodogram must be smoothed. Bloomfield (2000) recommends the Daniell window as a smoothing filter for generating an estimated spectrum from the periodogram. A plot of the filter weights therefore has the form of a trapezoid. The advantage of the Daniell filter over the rectangular filter for smoothing the periodogram is that the Daniell filter has less leakage, which refers to the influence of variance at non-Fourier frequencies on the spectrum. Successive smoothing by Daniell filters with different spans gives an increasingly smooth spectrum (Bloomfield, 2000).

An averaged (or smoothed) periodogram is defined as the average of the periodogram values, say,

$$\bar{f}(\omega) = \frac{1}{L} \sum_{k=-m}^m I(\omega_j + k/n) \quad (28)$$

over the frequency band  $\beta$ ,

$$\beta = \left\{ \omega^* : \omega_j - \frac{m}{n} \leq \omega^* \leq \omega_j + \frac{m}{n} \right\} \quad (29)$$

Where:  $m$  is spans of denial filter and

$$L = 2m + 1 \quad (30)$$

is an odd number, chosen such that the spectral values lie in the interval  $\beta$ .

Under the assumption that the spectral density is fairly constant in the band  $\beta$ , and in view of Eqn.(25), we can show that under appropriate conditions, for large  $n$ , the periodograms in Eqn.(28) are approximately distributed as independent  $f(\omega) \chi^2_2/2$  random variables, for  $0 < \omega < 1/2$ , as long as we keep  $L$  fairly small relative to  $n$ . Thus, under these conditions,  $L\bar{f}(\omega)$  is the sum of  $L$  approximately independent  $f(\omega) \chi^2_2/2$  random variables. It follows that, for large  $n$ ,

$$\frac{2L\bar{f}(\omega)}{f(\omega)} \sim \chi^2_{2L} \quad (31)$$

In this scenario, it seems reasonable to call the length of the interval defined by Eqn.(30),

$$\beta_\omega = L/n \quad (32)$$

the bandwidth. Bandwidth, of course, refers to the width of the frequency band used in smoothing the periodogram.

The result in Eqn.(31) can be rearranged to obtain an approximate  $100(1 - \alpha) \%$  confidence interval for the true spectrum,  $f(\omega)$  of the form (Shummay and Stoffer,2010):

$$\frac{2L\bar{f}(\omega)}{\chi^2_{2L}(1-\alpha/2)} \leq f(\omega) \leq \frac{2L\bar{f}(\omega)}{\chi^2_{2L}(\alpha/2)} \quad (33)$$

The use of the confidence intervals and the necessity for smoothing requires that we make a decision about the bandwidth( $B_\omega$ ) over which the spectrum will be essentially constant. Taking too broad a band will tend to smooth out valid peaks in the data when

the constant variance assumption is not met over the band. Taking too narrow a band will lead to confidence intervals so wide that peaks are no longer statistically significant. Thus, we note that there is a conflict here between variance properties or bandwidth stability, which can be improved by increasing  $B\omega$  and resolution, which can be improved by decreasing  $B\omega$ . A common approach is to try a number of different bandwidths, and to look qualitatively at the spectral estimators for each case. Wider Daniell filter, greater in smoothing and greater in decrease in resolution.

The problem of resolution was due to the fact that simple averaging was used in computing  $\bar{f}(\omega)$  defined in Eqn. (28). There is no particular reason to use simple averaging, and we might improve the estimator by employing a weighted average, say

$$\hat{f}(\omega) = \sum_{k=-m}^m h_k I(\omega_j + k/n) \quad (37)$$

Where the weights  $h_k > 0$  satisfy,  $\sum_{k=-m}^m h_k = 1$ . That means, the sum of periodogram weights must equal one for the spectral estimate to be an unbiased estimate of the true spectrum (Bloomfield, 2000). In particular, it seems reasonable that the resolution of the estimator will improve if we use weights that decrease as distance from the center weight  $h_0$  increases. An approximation in Eqn. (37) that seems to work well is found if we replace

$$L \text{ by } L_h = \left( \sum_{k=-m}^m h_k^2 \right)^{-1}$$

and

use the approximation:

$$\frac{2L_h \hat{f}(\omega)}{f(\omega)} \approx \chi_{2L_h}^2 \quad (38)$$

In analogy to Eqn. (33), we will define the bandwidth in this case to be

$$B_\omega = L_h / n \quad (39)$$

The result in Eqn. (38) can be rearranged to obtain an approximate  $100(1 - \alpha) \%$  confidence interval for the true spectrum,  $f(\omega)$  of the form (Shumay and Stoffar,2010):

$$\frac{2 L_h \hat{f}(\omega)}{\chi_{2L_h}^2(1-\alpha/2)} \leq f(\omega) \leq \frac{2 L_h \hat{f}(\omega)}{\chi_{2L_h}^2(\alpha/2)} \quad (40)$$

### 3.2.1.2. Cross- spectral Analysis

The cross-spectrum is similar to the power spectrum except that the covariance is substituted for the variance of the series. Cross-spectrum analysis examines the relationship between a pair of series.

Two main objective of cross-spectrum analysis are to determine the existence of correlation between two variables and to simultaneously separate each of two time series in to its harmonic component.

The cross-spectrum is generally a complex-valued function, and it is often written as (Shumay and Stoffar, 2010):

$$f_{xy}(\omega) = C_{xy}(\omega) - iq_{xy}(\omega) \quad (41)$$

where

$$C_{xy}(\omega) = \sum_{h=-\infty}^{\infty} \gamma_{xy}(h) \cos(2\pi\omega h)$$
$$q_{xy}(\omega) = \sum_{h=-\infty}^{\infty} \gamma_{xy}(h) \sin(2\pi\omega h)$$
$$\gamma_{xy}(\omega) = E(x_{t+h} - \mu_x)(y_t - \mu_y),$$

The principal tool that measures the strength of relation between two series in cross-spectral analysis is squared coherence function defined as (Shumay and Stoffar, 2010):

$$\rho_{xy}^2(\omega) = \frac{|f_{xy}(\omega)|^2}{f_{xx}(\omega)f_{yy}(\omega)}, \quad -1/2 \leq \omega \leq 1/2 \quad (42)$$

where  $f_{xx}(\omega)$  and  $f_{yy}(\omega)$  are individual spectra of the series  $x_t$  and  $y_t$  respectively.

Then spectral Matrix of a Bivariate Process is given by:

$$f(\omega) = \begin{pmatrix} f_{xx}(\omega) & f_{xy}(\omega) \\ f_{yx}(\omega) & f_{yy}(\omega) \end{pmatrix} \quad (43)$$

An estimate of the above spectral matrix is given as:

$$\hat{f}(\omega) = \sum_{k=-m}^m h_k I(\omega_j + k/n) , \text{ which is similar to the eqn. (37) defined in power spectra.}$$

Then an estimate of the squared coherence between two series,  $x_t$  and  $y_t$  is given by:

$$\hat{\rho}_{xy}^2 = \frac{|\hat{f}_{xy}(\omega)|^2}{\hat{f}_{xx}(\omega)\hat{f}_{yy}(\omega)} \quad (44)$$

We can test the hypothesis that  $\rho_{xy}^2 = 0$ , using the test statistic (Shuamay and Stoffar, 2010):

$$F = \frac{\hat{\rho}_{xy}^2(\omega)}{(1 - \hat{\rho}_{xy}^2(\omega))} (L - 1) \quad (45)$$

which has an approximate F-distribution with (2) and  $(2L - 2)$  degrees of freedom.

Solving (45) for a particular significance level  $\alpha$  leads to

$$C_\alpha = \frac{F_{2,2L-2}(\alpha)}{L - 1 + F_{2,2L-2}(\alpha)} \quad (46)$$

as the approximate value that must be exceeded for the original squared coherence to be able to reject  $\rho_{xy}^2 = 0$  at a specified frequency.

### 3.2.2. Time domain approach

This approach focuses on modeling some future value of the series as function of the current and the past. Conducting investigations using standard statistical methodologies is an essential step in the development of climatology (Polyak, 1996). In this respect, the time-domain approach of univariate time series continues to be an important topic. An intrinsic feature of the time-domain approach is that, typically, adjacent points in time are correlated and that future values are related to past and present values. Autoregressive integrated moving average (ARIMA) modeling is one of the most widely implemented methods for analyzing univariate time series data (Box and Jenkins, 1976). In order to understand the modeling procedure, it is useful to briefly introduce the following basic models.

#### Autoregressive (AR) models

Autoregressive models are based on the idea that the current value of the series,  $x_t$ , can be explained as a function of  $p$  past values,  $x_{t-1}, x_{t-2}, \dots, x_{t-p}$ , where  $p$  determines the number of steps into the past needed to forecast the current value.

An autoregressive model of order  $p$ , abbreviated AR ( $p$ ), can be written as:

$$x_t = \phi_1 x_{t-1} + \phi_2 x_{t-2} + \dots + \phi_p x_{t-p} + w_t \quad (47)$$

where  $x_t$  is stationary series,  $\phi_1, \phi_2, \dots, \phi_p$  are the parameters of the AR ( $\phi_p \neq 0$ ). Unless otherwise stated, we assume that  $w_t$  is a Gaussian white noise series with mean zero and variance  $\sigma_w^2$ . The highest order  $p$  is referred to as the order of the model.

The model in lag operators takes the following form:

$$(1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p) x_t = w_t,$$

where the lag (backshift) operator B is defined as:  $B^p x_t = x_{t-p}$ ,  $p=0,1,2,\dots$

More concisely we can express the model as:

$$\phi(B) x_t = w_t \tag{48}$$

The autoregressive operator  $\phi(B)$  is defined to be

$$\phi(B) = 1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p \tag{49}$$

The values of  $\phi$  which make the process stationary are such that the roots of  $\phi(B) = 0$  lie outside the unit circle in the complex plane (Chatfield, 1991). If all roots of  $\phi(B)$  are larger than one in absolute value, then the process is a stationary process satisfying the autoregressive equation and can be represented as:

$$x_t = \sum_{j=0}^{\infty} \psi_j w_{t-j} \tag{50}$$

The coefficients  $\psi_j$  converge to zero, such that  $\sum_{j=0}^{\infty} |\psi_j| < \infty$ . If some roots are “exactly” one in modulus, no stationary solution exists.

A plot of the ACF of a stationary AR ( $p$ ) model show a mixture of damping sine and cosine patterns and exponential decays depending on the nature of its characteristic roots.

Another characteristics feature of AR ( $p$ ) models is that the partial autocorrelation function defined as  $\text{PACF}(j) = \text{corr.}(x_t, x_{t-j} | x_{t-1}, x_{t-2}, \dots, x_{t-j+1})$  becomes “exactly” zero for values larger than  $p$  (Tsay, 2005).

## **Moving average (MA) Models**

As an alternative to the autoregressive representation in which the  $x_t$  on the left-hand side of the equation are assumed to be combined linearly, the moving average model of order  $q$ , abbreviated as MA ( $q$ ), assumes the white noise ( $w_t$ ) on the right-hand side of the defining equation are combined linearly to form the observed data.

A series  $x_t$  is said to follow a moving average process of order  $q$ , or simply MA ( $q$ ) process if

$$x_t = w_t + \theta_1 w_{t-1} + \theta_2 w_{t-2} + \dots + \theta_q w_{t-q} \quad (51)$$

where  $\theta_1, \theta_2, \dots, \theta_q$  are the MA parameters. MA( $q$ ) models immediately define stationary, every MA process of finite order is stationary (Diebold et al., 2006). In order to preserve a unique representation, usually the requirement is imposed that all roots of  $\theta(B) = 1 + \theta_1 B + \theta_2 B^2 + \dots + \theta_q B^q = 0$  are greater than one in absolute value. If all roots of  $\theta(B) = 0$  lie outside the unit circle, the MA process has an autoregressive representation of generally infinite order  $\sum_{j=0}^{\infty} \psi_j x_{t-j} = w_t$  with  $\sum_{j=0}^{\infty} |\psi_j| < \infty$ . MA process as with an infinite order autoregressive representation are said to be invertible.

A characteristic feature of MA (q) is that their ACF,  $\rho_j$ , becomes statistically insignificant after  $j=q$ . The property of the ACF should be reflected in the correlogram, which should „cut off“ after q. The PACF converges to zero geometrically.

### **Autoregressive –Moving average (ARMA)**

We now proceed with the general development of autoregressive, moving average, and mixed autoregressive moving average (ARMA), models for stationary time series. In most cases, it is best to develop a mixed autoregressive moving average model when building a stochastic model to represent a stationary time series. The order of an ARMA model is expressed in terms of both  $p$  and  $q$ . The model parameters relate to what happens in period  $t$  to both the past values and the random errors that occurred in past time periods. A general ARMA model can be written as follows:

$$x_t = \phi_1 x_{t-1} + \phi_2 x_{t-2} + \dots + \phi_p x_{t-p} + w_t + \theta_1 w_{t-1} + \theta_2 w_{t-2} + \dots + \theta_q w_{t-q} \quad (52)$$

Equation (52) of the time series model will be simplified by a backward shift operator  $B$  to obtain

$$\phi(B) x_t = \theta(B) w_t \quad (53)$$

The ARMA model is stable –i.e., it has a stationary „solution“ –if all roots of  $\phi(B)=0$  are larger than one in absolute value. The representation is unique if all roots of  $\phi(B)=0$  lie outside the unit circle and  $\phi(B)$  and  $\theta(B)$  do not have common roots. Stable ARMA

models always have an infinite order MA representation. If all roots of  $\phi(B)$  are larger than one in absolute value, it has an infinite order AR representation. The process is invertible only when the roots of  $\theta(B)$  lie outside the unit circle. Furthermore, a process is said to be causal when the roots of  $\phi(B)$  lie outside the unit circle.

To have ARMA  $(p,q)$  model, both ACF and PACF should show a pattern of decaying to zero. The autocorrelation of an ARMA  $(p, q)$  process is determined at greater lags by the AR  $(p)$  part of the process as the effect of the MA part dies out. Thus, eventually the ACF consists of mixed damped exponentials and sine terms. Similarly, the partial autocorrelation of an ARMA  $(p, q)$  process is determined at greater lags by the MA  $(q)$  part of the process. Thus, eventually the partial autocorrelation function will also consist of a mixture of damped exponentials and sine waves.

### **Autoregressive Integrated Moving Averages (ARIMA) Models**

Autoregressive integrated moving average (ARIMA) models are specific subset of univariate modeling, in which a time series is expressed in terms of past values of itself (the autoregressive component) plus current and lagged values of a „white noise“ error term (the moving average component). ARIMA models are univariate models that consist of an autoregressive polynomial, an order of integration (d), and a moving average polynomial.

A process  $(x_t)$  is said to be an autoregressive integrated moving average process, denoted by ARIMA (p, d, q) if it can be written as:

$$\phi(B) \nabla^d x_t = \theta(B) w_t \quad (54)$$

where  $\nabla^d = (1 - B)^d$  with  $\nabla^d x_t$  and  $d^{\text{th}}$  consecutive differencing (Vandale,1983)

If  $E(\nabla^d x_t) = \mu$ , we write the model as

$$\phi(B) \nabla^d x_t = \alpha + \theta(B) w_t \quad (55)$$

Where:  $\alpha$  is a parameter related to the mean of the process  $\{x_t\}$ , by  $\alpha = \mu (1 - \phi_1 - \dots - \phi_p)$  and this process is called a white noise process, that is, a sequence of uncorrelated random variables from a fixed distribution (often Gaussian) with constant mean  $E(x_t) = \mu$ , usually assumed to be “zero” and constant variance. if  $d=0$ , it is called ARMA(p,q) model while when  $d=0$  and  $q=0$ , it is referred to as autoregressive of order  $p$  model and denoted by AR ( $p$ ). When  $p=0$  and  $d=0$ , it is called Moving Average of order  $q$  model, and is denoted by MA ( $q$ ).

## Seasonal ARIMA (SARIMA)

In this section, we introduce several modifications made to the ARIMA model to account for seasonal and non-stationary behavior. Often, the dependence on the past tends to occur most strongly at multiples of some underlying seasonal lags. Seasonal ARIMA (SARIMA) is used when the time series exhibits a seasonal variation. Natural phenomena such as temperature and rainfall have strong components corresponding to seasons. Hence, the natural variability of many physical, biological, and economic processes tends to match with seasonal fluctuations. Because of this, it is appropriate to introduce autoregressive and moving average polynomials that identify with seasonal lags. The resulting pure seasonal autoregressive moving average model, say,

ARMA  $(P, Q)_S$ , then takes the form (Shumay and Stoffer, 2010):

$$\Phi_P (B^S) x_t = \Theta_Q (B^S) w_t \quad (56)$$

with the following definition of the operators

$$\Phi_P (B^S) = 1 - \Phi_{1S}B^S - \Phi_{2S}B^{2S} - \dots - \Phi_{PS}B^{Ps}$$

and (57)

$$\Theta_Q (B^S) = 1 + \Theta_{1S}B^S + \Theta_{2S}B^{2S} + \dots + \Theta_{QS}B^{Qs}$$

are the seasonal autoregressive operator and the seasonal moving average operator of orders  $P$  and  $Q$ , respectively, with seasonal period  $S$ . Analogous to the properties of non-seasonal ARMA models, the pure seasonal ARMA  $(P, Q)_S$  is causal only when the roots

of  $\Phi_P (Z^S)$  lie outside the unit circle, and it is invertible only when the roots of  $\Theta_Q (Z^S)$  lie outside the unit circle.

In general, we can combine the seasonal and non-seasonal operators into a multiplicative seasonal autoregressive moving average model, denoted by  $ARMA(p, q) \times (P, Q)_s$ , and write as the overall model as:

$$\Phi_P (B^S) \phi (B) x_t = \Theta_Q (B^S) \theta(B) w_t \quad (58)$$

A seasonal autoregressive notation ( $P$ ) and a seasonal moving average notation ( $Q$ ) will form the multiplicative seasonal autoregressive integrated moving average model, denoted by  $ARIMA (p, d, q)^*(P, D, Q)_s$ , of Box and Jenkins (1976) and is given by:

$$\Phi_P (B^S) \phi(B) \nabla_S^D x_t = \alpha + \Theta_Q (B^S) \theta(B) w_t \quad (59)$$

where  $w_t$  is the usual Gaussian white noise process. The ordinary autoregressive and moving average components are represented by polynomials  $\phi(B)$  and  $\theta(B)$  of orders  $p$  and  $q$ , respectively and the seasonal autoregressive and moving average components by  $\Phi_P (B^S)$  and  $\Theta_Q (B^S)$  of orders  $P$  and  $Q$ . The ordinary and seasonal difference components can be written as

$$\nabla^d = (1 - B)^d \quad \text{and} \quad \nabla_S^D = (1 - B^S)^D \quad (60)$$

### 3.2.2.1. Test for Stationarity

Before developing a Box-Jenkins modeling process, it is important to check whether the data under study meets basic assumptions such as stationarity. A time series is said to be stationary if there is no systematic change in mean (no trend), if there is no systematic change in variance and if periodic variations have been removed. Stationarity may be classified as strict stationary and weak stationary.

**Strict stationary:** A stochastic process is said to be strictly stationary if its statistics (e.g., mean, variance, serial correlation) are not affected by a shift in the time origin, that is, if the joint probability distribution associated with  $n$  observations  $(x_1, x_2, \dots, x_n)_t$  made at time origin  $t$ , is the same as that associated with  $n$  observations  $(x_1, x_2, \dots, x_n)_{t+k}$  made at time origin  $t+k$ . In other words,  $x(t)$  is a strictly stationary process when the two processes  $x(t)$  and  $x(t+k)$  have the same statistics for any time lag of  $k$ .

**Weak stationarity:** Weak stationarity means that only the lower order moments of the distribution function: the mean constant over time and the covariance function invariant on time, that is, it depends only on time differences (i.e., lags). This is also called stationarity in a wide sense.

If the probability distribution associated with any set of times is a multivariate normal distribution, the process is called a normal or Gaussian process. Since the multivariate normal distribution is fully described by its first and second order moments, it follows that weak stationarity and an assumption of normality imply **strict stationarity**.

Most of the probability theory of time series is concerned with stationary time series, for this reason, time series analysis often requires one to turn a non-stationary series into a stationary one so as to use this theory (Brockwell and Davis, 1996). Usually a common mechanism used to transform a series into a stationary one is differencing, which means calculating the difference among pairs of observations at some time interval. Many testing procedures for stationarity are employed in the literature. We use the following four relatively simple techniques.

***Time plot:*** The first step in the analysis of time series is usually to plot the data and obtain simple descriptive measures of the main property of the series via a visual inspection of the time series plot. This may reveal one or more of the following characteristics: Seasonality, trends either in the mean level or the variance of the series, long- term cycles, and so on. If any such patterns are present, then these are signs of non-stationarity.

***The Correlogram Test:*** One way to characterize a series with respect to its dependence over time is to plot its sample autocorrelation function (SACF). The sample partial autocorrelation function, denoted by SPACF, is similar to the SACF and can be described as the correlation between  $x_t$  and  $x_{t-s}$  (observations of the time series recorded at two moments in time  $s$  time units apart) after controlling for the common linear effects of the intermediate lags. Both functions are used in Box-Jenkins modeling as correlograms to reveal important information regarding the order of the autoregressive (AR) and moving average (MA) factors present in the generating process of the given

time series as well as to assess stationarity. Enders (1995) expresses that inspection of SACF serves as a rough indicator of whether non stationarity is present in a series. Wei (1990) states that if the sample ACF decays very slowly, it indicates that differencing is needed.

While the stationarity tests described in the above sections make use of subjective visual inspection of data plots and correlograms, formal tests were developed to help with determining stationarity. These tests, also known as unit root tests and stationarity tests are based for the most part on formal statistical tests and the difference between them lies in the stringency of the assumptions they use as well as in the form of the null and alternative hypotheses they adopt. The standard Dickey-Fuller test (DF) is based on i.i.d. errors and has as the null hypothesis the unit root. On the other hand, the Phillips-Perron test is nonparametric and allows for some heterogeneity and serial correlation in the innovations. There exist many other unit root and stationarity tests as well as generalizations and combinations of the ones mentioned above. However, in this study, Dickey-Fuller test (DF) is to be used.

**Unit Root Test:** For a univariate time series, the Unit Root test is frequently employed for testing stationarity. The test first poses the null hypothesis that the given time series has a unit root, which means that the time series is non-stationary, and tests if the null hypothesis is to be statistically rejected in favor of the alternative hypothesis that the given time series is stationary. To detect whether a given series has non stationarity, let's assume that the relationship between current value (in time  $t$ ) and last value (in time  $t-1$ ) in the time series is as following (Enders, 1995):

$$x_t = \phi x_{t-1} + w_t \quad (61)$$

where  $x_t$  is an observation value at time  $t$ ,  $w_t$  is a white noise process. This model is a first order autoregressive process. The time series  $x_t$  converges, as  $t \rightarrow \infty$ , to a stationary time series if  $|\phi| < 1$ . If  $|\phi| = 1$  or  $>1$ , the series  $x_t$  is not stationary and the variance of  $x_t$  is time dependent (Diebold et al., 2006). In other words, the series has a unit root. The Unit Root test subsequently tests the following one-sided hypothesis

$$H_0: \phi = 1 \text{ (has a unit root)}$$

$$H_1: \phi < 1 \text{ (has root outside the unit circle)}$$

The name, unit root, comes from the fact that the coefficient of  $x_{t-1}$  is unity, if the time series is non-stationary, and the Unit Root test, as the name suggests, tests if  $\phi$  is unity or not. If  $x_{t-1}$  is subtracted from the right and left sides of the above equation, we get:

$$\nabla x_t = (\phi - 1)x_{t-1} + w_t \quad (62)$$

This equation is expressed as a first order difference equation. If  $\phi$  is taken one in the equation, the effect of unit root can be removed from the actual series that has non stationarity via a first differencing. The tests above are valid only if  $w_t$  is white noise. In particular,  $w_t$  is assumed not to be autocorrelated, but would be so if there was autocorrelation in the dependent variable of the regression ( $x_t$ ) which has not been modeled. If this is the case, the test would be „oversized“, meaning that the true size of the test (the proportion of times a correct null hypothesis is incorrectly rejected) would be

higher than the nominal size used (e.g. 5%). The solution is to „augment“ the test using  $p$  lags of the dependent variable (Brooks, 2008).

**The Augmented Dickey-Fuller (ADF) Test:** To use the Augmented Dickey-Fuller test, here are the various cases of the test equation:

When the time series is flat (i.e. doesn't have a trend) and potentially slow- turning around zero, use the following test equation:

$$\nabla x_t = \Phi x_{t-1} + \theta_1 \nabla x_{t-1} + \dots + \theta_p \nabla x_{t-p} + w_t \quad (63)$$

When the time series is flat and potentially slow-turning around a non-zero value, use the following test equation:

$$\nabla x_t = \theta_0 + \Phi x_{t-1} + \theta_1 \nabla x_{t-1} + \dots + \theta_p \nabla x_{t-p} + w_t \quad (64)$$

When the time series has a trend in it (either up or down) and is potentially slow-turning around a trend line you would draw through the data, use the following test equation:

$$\nabla x_t = \theta_0 + \beta t + \Phi x_{t-1} + \theta_1 \nabla x_{t-1} + \dots + \theta_p \nabla x_{t-p} + w_t \quad (65)$$

Where:  $\nabla x_t$  is the first differenced value of  $series(x_t)$ ,  $w_t$  is the error term,

$x_{t-1}$  is the first lagged value of the series ( $x_t$ )

$\nabla x_{t-j}$  is the  $j^{\text{th}}$  lagged first differenced of values of  $x_t$ ,

$\theta_0, \beta, \Phi = \phi - 1, \theta_1, \theta_2, \dots, \theta_p$  are parameters to be estimated .

A problem now arises in determining the optimal number of lags of the dependent variable. Although several ways of choosing  $p$  have been proposed, they are all somewhat arbitrary, and are thus not presented here. Instead, the following two simple rules of thumb are suggested. First, the frequency of the data can be used to decide. So, for example, if the data are monthly, use 12 lags, if the data are quarterly, use 4 lags, and so on. Clearly, there would not be an obvious choice for the number of lags to use in a regression containing higher frequency financial data (e.g. hourly or daily). Second, an information criterion can be used to decide. So choose the number of lags that minimizes the value of an information criterion (AIC, SBC) (Brooks, 2008).

In this study, we use the test-statistic associated with the Ordinary least squares estimate of  $\Phi$ . This is called the **Dickey-Fuller test- statistic**. The Dickey-Fuller-test now estimates  $\pi = \phi + 1$  by  $\hat{\pi}$ , obtained from an ordinary regression and checks for  $\Phi = 0$  by computing the test statistic:

$$\tau = n \widehat{\Phi} = n (\hat{\pi} - 1) \quad (66)$$

Where  $n$  is the number of observations on which the regression is based.

The hypothesis of the Augmented Dickey-Fuller test for all the above three cases is

$H_0: \Phi = 0$  – (the data needs to be differenced to make it stationary)

$H_1: \Phi < 0$  – (the data is stationary and doesn't need to be differenced)

Unfortunately, the Dickey-Fuller t-statistic does not follow a standard t-distribution as the sampling distribution of this test statistic is skewed to the left with a long, left-hand-tail.

Hence, the test statistic follows the so called Dickey-Fuller distribution which cannot be explicitly given but has to be obtained by Monte-Carlo and bootstrap methods (Falk, 2006).

But note that if there is a trend term in the series, the test statistic follows the so called F-distribution.

The augmented Dickey–Fuller (ADF) statistic, used in the test, is a negative number. The more negative it is, the stronger the rejection of the hypothesis that there is a unit root at some level of confidence (Elliott et.al., 1996)

**Variance Comparisons:** The behavior of variance associated with different orders of differencing, can provide a useful means of deciding the appropriate order of differencing to achieve stationary (Hamilton, 1994). A time series that is non-stationary in mean can be made stationary by the first differencing. But, if the series is also not stationary in the rate of change of the mean (i.e. slope), stationarity can be achieved by taking the second difference, or the difference of the first difference. We should, however, bear in mind that each successive differencing will decrease the variance of the series, but at some point, higher-order differencing will have an opposite effect. When variance increases, it means that the series has been over-differenced (Nagpaul, 2005).

### 3.2.2.2. Handling Missing Values

Similar to other statistical analysis a problem frequently encountered in time series data analysis is missing observations in a data series. Missing data must also be addressed in the time series context. Special consideration in handling missing data in a time series application is that the missing data cannot simply be omitted from the data series. Missing observations must be replaced by appropriately estimated values so that the alignment of data between time periods will not be offset inappropriately. In order to replace those observations, there are several options available in the literature. As discussed in Liu and William (2001), missing data in a time series may be estimated using one of the following methods. Firstly, replace with the mean of the series. Secondly, replace with the naïve forecast. Naive model is the simplest form of a univariate forecast model. It uses the current time value for the next time, that is  $\hat{x}_{t+1} = x_t$ . Thirdly, replace with a simple trend forecast. This is accomplished by estimating the regression equation of the form,  $x_t = \alpha + \beta t$ , where  $t$  is the time for the periods prior to the missing value. Then use the equation to fit the time periods missing. Finally, replace with an average of the last two known observations that bound the missing observation.

### 3.2.2.3. Building ARIMA Models

To identify a perfect ARIMA model for a particular time series data, Box and Jenkins (1976) proposed a methodology that consists of four phases: i) Model identification; ii) Estimation of model parameters; iii) Diagnostic checking for the identified model and iv) Application of the model (i.e. forecasting).

#### ***i)* Model Identification**

The purpose of the identification stage is to determine the differencing required to achieve stationarity and also the order of both the seasonal and the non- seasonal AR and MA operators for the residual series.

There are a number of identification methods proposed in the literature. The autocorrelations function (ACF) and the partial autocorrelation functions (PACF) are the two most useful tools in any attempt at time series model identification (Granger and Newbold, 1986).

***Autocorrelation Function (ACF):*** The sample ACF ( $\Gamma_k$ ) measures the amount of linear dependence between observations in a time series that are separated by a lag  $k$ . To use the ACF in model identification, estimate  $\Gamma_k$  and then plot  $\Gamma_k$  series against lag  $k$  up to a maximum lag of about five times the seasonality interval and this should be less than to one fourth of the series under study (Hipel et al., 1977). To identify the number of non seasonal and seasonal autoregressive, and non seasonal and seasonal moving average

parameters, we examine the ACF based on the theoretical pattern for the identified parameters using Table-3.1 and 3.2 as summarized (Shumway and Stoffer, 2010):

**Table 3.1:** Behavior of the ACF and PACF for ARMA Models

	AR(P)	MA(Q)	ARMA(P,Q)
ACF	Tails off	Cuts off after lags Q	Tails off
PACF	Cuts off after lags P	Tails off	Tails off

**Table 3.2:** Behavior of the ACF and PACF for Pure SARMA Models

	AR(P) <sub>s</sub>	MA(Q) <sub>s</sub>	ARMA(P,Q) <sub>s</sub>
ACF	Tails off at lags K <sub>s</sub> , K=1, 2, ...	Cuts off after lags Q <sub>s</sub>	Tails off at lags K <sub>s</sub>
PACF	Cuts off after lags P <sub>s</sub>	Tails off at lags k <sub>s</sub>	Tails off at lag P <sub>s</sub>

Where the values at non seasonal lags  $h \neq Ks$ , for  $K = 1, 2, \dots$ , are zero. When the process is SARIMA  $(0, d, q) \times (0, D, Q)_S$  model,  $\Gamma_k$  truncates and is not significantly different from zero after lag  $q+sQ$ . If  $\Gamma_k$  spikes out at lags that are multiples of  $s$ , this implies the presence of a seasonal autoregressive component. The failure of the autocorrelation function to truncate at other lags may imply that a non seasonal autoregressive term is required.

The autocorrelation of order  $k$  is simply the correlation between  $x_t$  and  $x_{t-k}$ , i.e.

$$\rho_k = \frac{E\{(x_t - \bar{x})(x_{t-k} - \bar{x})\}}{E\{(x_t - \bar{x})^2\}} \quad (67)$$

In practice, one never knows the true autocorrelations and partial autocorrelations and at the identification stage, one has to rely on the sample autocorrelation and partial autocorrelation functions imitating the behavior of the corresponding parent quantities.

True autocorrelations ( $\rho_k$ ) can be estimated by:

$$r_k = \frac{1/n \sum_{t=k+1}^n (x_t - \bar{x})(x_{t-k} - \bar{x})}{1/n \sum_{t=1}^n (x_t - \bar{x})^2} \quad (68)$$

where  $\bar{x}$  is the sample mean of the  $x_t$ 's.

**Partial Autocorrelation Function (PACF):** Partial autocorrelation function can also be used for determining the possible order of seasonal autoregressive, non-seasonal autoregressive, moving average and seasonal moving average that should be incorporated in the model by the help of Table-3.1 and 3.2 above. When the process is a

pure SARIMA  $(p, d, 0) \times (P, D, 0)_{12}$  model,  $r_{kk}$  cuts off and is not significantly different from zero after lag  $p+SP$ . If  $r_{kk}$  dampouts at lags that are multiples of  $s$ , this suggests the incorporation of a seasonal moving average component in to the model. The failure of the partial autocorrelation function to truncate at other lags may imply that a non seasonal MA term is required (Hipel et al., 1977). To obtain an estimate for partial autocorrelations  $(\rho_{kk})$  at lag  $k$ , we can employ successive autoregressive estimation procedure. The first step is to model the  $x_t$  series by finite autoregressive models of order  $K$  given by (Box and Jenkins, 1976):

$$x_t = \rho_0 + \sum_{k=1}^K \rho_{kk} x_{t-k} \quad (69)$$

Where  $\rho_{kk}$  is the  $k^{\text{th}}$  autoregressive coefficient and  $k=1, 2, \dots, K$ . Estimate of these coefficients by ordinary least squares or maximum likelihood estimation method gives the  $k^{\text{th}}$ - sample partial autocorrelation (Hipel et al., 1977).

## **ii) Parameter Estimation**

After choosing the most appropriate model (step (i) above), the model parameters are estimated by using several estimation procedures. The estimation-stage results will be used to check: (i) parameter estimates, (ii) the appropriateness of coefficient estimates which includes the statistical significance of estimated coefficient and standard error and correlation matrix.

In maximum likelihood methods, the likelihood function is maximized in order to obtain the parameter estimates. The likelihood of a set of data is the probability of obtaining that particular set of data, given its distribution. The philosophy behind maximum likelihood estimates is to find a set of parameters which maximize the likelihood of observing the data to which the model is being fitted.

The linear optimization algorithm is used to maximize the likelihood function with respect to the parameter space (Shumway and Stofferr, 2010).

In Time series analysis, there may be several adequate models that can be used to represent a given data set, and hence, numerous criteria for model comparison have been introduced in the literature. One of them is based on the so-called information criteria. The idea is to balance the risks of under fitting (selecting an order smaller than the true order) and over fitting (selecting an order larger than the true order).

Akaike (1978) introduced a criterion called Akaike Information Criterion (AIC) in the literature. The AIC is a mathematical selection criterion of model building. When there

are several competing models to choose from, select the model that gives the minimum of the AIC defined by (Shumay and Stoffar, 2010):

$$AIC = \log \hat{\delta}_k^2 + \frac{n+2k}{n} \quad (71)$$

Where:

$\hat{\delta}_k^2 = \frac{SSE_k}{n}$  denotes the maximum likelihood estimator for the error variance and  $k$  is the number of seasonal and non-seasonal autoregressive and moving average parameters to be estimated in the model, that is, according to Wei (1990),  $k = p + q + P + Q + 1$  and  $n$  is the number of observations. The optimal order of the model is chosen by the value of  $k$ , which is a function of  $p$  and  $q$ ,  $P$  and  $Q$  so that the value of  $k$  yielding the minimum AIC specifies the best model. Wei (1990) expressed the need to select the model that has fulfilled all the diagnostic checks and has as few parameters as possible in terms of parsimony.

Schwartz (1978) suggested a Bayesian criterion called Schwartz's Bayesian Criterion (SBC) having the form:

$$SBC = \log \hat{\delta}_k^2 + \frac{k \log n}{n} \quad (72)$$

The (SBC) can also measure the parsimony of model building. A model that has the smallest SBC value among the competing models fit to time series is preferred. The optimal order of the model is chosen by the value of  $k$ , which is a function of  $p$  and  $q$ ,  $P$  and  $Q$  so that the value of  $k$  yielding the minimum SBC specifies the best model as AIC.

### **iii) Diagnostic Checking**

After fitting a provisional time series model, we can assess its adequacy in various ways. The usual approach is to extract from the data, a sequence of residuals to correspond to the underlying, last unobservable, white noise sequence, and to check that the statistical properties of these residuals are indeed consistent with white noise. Most diagnostic tests deal with the residual assumptions in order to determine whether the residuals from fitted model are independent, have a constant variance, and are normally distributed. Several diagnostic statistics and plots of the residuals can be used to examine the goodness of fit of the tentative model to the historical data.

The first approaches that can be used to evaluate the adequacy of a model are the plot of the errors over time, which can be written (Shumay and stoffar, 2010):

$$w_t = (x_t - \hat{x}_t^{t-1}) / \sqrt{p_t^{t-1}} \quad (73)$$

where  $x_t - \hat{x}_t^{t-1}$  is the one-step-ahead prediction of  $(x_t)$  based on the fitted model and  $\sqrt{p_t^{t-1}}$  is the estimated one-step-ahead error variance. If visual inspections of the errors reveal that they are randomly distributed over time, then we have a good model.

The autocorrelations function (ACF) of the series can also be used to examine whether the residual of the fitted model is white noise or not. If the ACF is significantly different from zero, this implies that there is dependence between observations (Janacek and Swift, 1993; Ferguson et al., 2000). There are different applications related to the Residual ACF

(RACF) for the independence of residuals. The first one is the correlogram drawn by plotting  $r_k(w)$  against lag  $k$ .

$$r_{ak} = \frac{\sum_{t=k+1}^n w_t w_{t-k}}{\sum_{t=1}^n w_t^2} \quad (74)$$

Under the assumption that residual follows a white noise process the standard errors of these  $(r_{ak})$  are approximately equal to  $1/\sqrt{T}$ . Thus, under the null hypothesis that residual follows a white noise process, roughly 95% of the autocorrelation coefficient  $(r_{ak})$  should fall within the range  $\pm 1.96/\sqrt{T}$ . If more than 5% of the coefficient fall outside of this range, then most likely residual does not follow a white noise process (Lehmann and Rode, 2001).

There are many statistical tests used for diagnostic checking of randomness. The Ljung-Box Q statistic, Turning point and Runs tests can be used for the diagnostic checking of residuals for independence.

**Ljung-Box Q (LBQ) Statistic:** The Ljung-Box Q or  $Q(r)$  statistic can be employed to check independence of residual instead of visual inspection of the sample autocorrelations. A test of hypothesis can be done for the model adequacy by choosing a level of significance and then comparing the value of calculated  $\chi^2$  with the  $\chi^2$ -table critical value. A useful test in these concepts is the portmanteau lack of fit test. This uses the entire residual sample with null hypothesis that  $H_0: \rho_1 = \rho_2 = \dots = \rho_k = 0$

The test statistic is calculated by the following equation (Ljung and Box, 1998):

$$Q(r) = n'(n' + 2) \sum_{j=1}^k \frac{\mu_j^2}{n-j} \quad (75)$$

where:  $n' = (n - d)$ ,  $n$  is the number of observations in the original time series,  $\mu_j^2$  is the sample autocorrelation of the residuals at lag  $j$  and  $d$  is the degree of non- seasonal differencing used to transform the original time series values into stationary time series values and  $k$  is sufficiently large integer.

This test statistic is the modified  $Q$  - statistic originally proposed by Box and Pierce (1970). Under the null hypothesis of model adequacy, Ljung and Box (1978) show that the  $Q$  - statistic approximately follows the  $\chi^2_{(K-M)}$  distribution where  $m$  is the number of parameters estimated in the model. If a model is correctly specified, residuals should be uncorrelated and  $Q(r)$  should be small (p- value should be large).

**Runs Test:** The runs test can be used to decide if a data set is from a random process. A run is defined as a series of increasing values or a series of decreasing values. The number of increasing (or decreasing) values is the length of the run. In a random data set, the probability that the  $(i + 1)^{th}$  value is larger or smaller than the  $i^{th}$  value follows a binomial distribution, which forms the basis of the runs test. The first step in the runs test is to compute the sequential differences  $(Y_i - Y_{i-1})$ . Positive values indicate an increasing value, whereas negative values indicate a decreasing value. In other terms, if  $Y_i > Y_{i-1}$  a 1 (one) is assigned for an observation and a 0 (zero) otherwise. The series then has an associated series of 1s and 0s. To determine if the number of runs is the correct number

for a series that is random, let  $T$  be the number of observations,  $T_a$  be the number above the mean,  $T_b$  be the number below the mean and  $R$  be the observed number of runs. Then, using combinatorial methods, the probability  $P(R)$  can be established and the mean and variance of  $R$  can be derived (Cromwell *et al.*, 1994). When  $T$  is relatively large ( $>20$ ) the distribution of  $R$  is approximately normal.

$$Z_N = \frac{R - E(R)}{\sqrt{V(R)}} \approx N(0,1) \quad (76)$$

where:

$$E(R) = \frac{T + 2T_a T_b}{T}$$

$$V(R) = \frac{2T_a T_b (2T_a T_b - T)}{2T(T - 1)}$$

The test of series randomness is rejected if the calculated  $Z_N$  value exceeds the selected critical value obtained from the standard normal distribution table.

**Turning Point Test:** A turning point means when the series changes from increasing to decreasing or vice versa. That is,  $X_{t-1} < X_t > X_{t+1}$  or  $X_{t-1} > X_t < X_{t+1}$ . Let  $T$  = the number of turning points in an  $n$ - period series. In order to carry out the test of white noise with this test, we must determine the distribution of the number of turning points in a series. It is known that with increasing  $n$ -the distribution of  $T$  is approximately normally distributed (Kendal and Ord, 1990). Then, the *test* statistic ( $N_T$ ) defined and approximated in Eqn.(77) should be compared with the  $Z$ -table critical value. The hypothesis of randomness should be rejected at  $\alpha$  significance level if the absolute value

of  $N_T > N_{T(1-\alpha/2)}$ , where  $N_{T(1-\alpha/2)}$  is the  $(1- \alpha /2)$  quartile of standard normal distribution (Cromwell *et al.*, 1994):

$$N_T = \frac{|T - \mu_T|}{\sqrt{Var(T)}} \approx N(0,1) \quad (77)$$

Where:  $\mu_T = (2/3)(n - 2)$   
 $Var(T) = (16n - 29) / 90$

**Test for normality of the residuals:** If a data set is distributed according to the bell shaped curve of the normal distribution, this set can be referred to as normal. Therefore, the histogram and Q-Q plot of a data set give information related to normality. There are several statistical tests used for the diagnostic checking of normality. The null hypothesis for any test of normality is that the data are normally distributed. In this study, Shapiro-Wilk Test will be used for the diagnostic checking of residuals for normality. The Shapiro-Wilk statistic “W” is proportional to the ratio of the squared slope of the normal probability plot to the usual mean square estimate (Gibbons, 1994):

$$W = \frac{(\sum_{i=1}^n a_{i,n} x_i)^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (78)$$

Where:  $a_{i,n}$  for the W- statistic are given in Table (Gibbons, 1994).

### **Test approaches for diagnostic checking of homosecdasticity**

For the diagnostic checking of residuals in terms of homoscedasticity, Goldfeld-Quandt, Breusch and Pagan and Spearman's rho tests are commonly used for time series data. But in this study, we use only Goldfeld-Quandt test and Breusch and Pagan test.

**Goldfeld-Quandt Test:** This test is very useful for determining whether a transformation of the data is needed. If there is a change in variance (heteroscedasticity) of residuals, a transformation is necessary for the data. The following is the steps for the Goldfeld-Quandt statistic (Greene, 2000):

- 1) Rank or order the residuals from the model fit to the data in ascending order,
- 2) Divide the residuals into three parts:  $n_1$  observations in the first part,  $p$  observations in the middle part, and  $n_2$  observations in the second part ( $n_1 + n_2 + p = n$ ). Usually  $p$  is taken to be one-sixth of  $n$ .

- 3) Run a regression on the first  $n_1$  observations, obtain the residuals ( $\hat{\varepsilon}_{1i}$ ), and calculate the residual variance  $s_1^2 = \sum_{i=1}^{n_1} \hat{\varepsilon}_{1i}^2 / (n_1 - 2)$ . Similarly run a regression on the second  $n_2$  observations, obtain the residuals ( $\hat{\varepsilon}_{2i}$ ), and calculate the variance

$$s_2^2 = \sum_{i=1}^{n_2} \hat{\varepsilon}_{2i}^2 / (n_2 - 2).$$

- 4) Calculate the test statistic:  $F_{cal} = \frac{S_2^2}{S_1^2}$ , iff  $S_2^2 > S_1^2$ , with  $(n_2 - k, n_1 - k)$  degree of freedom.

- 5) If  $F_{cal}$  is smaller than the critical value ( $F_\alpha(n_2 - k, n_1 - k)$ ), the residuals are assumed to be homoscedastic.

**Breusch-Pagan Test:** This involves applying ordinary least-square (OLS) to:

$$\frac{\hat{\varepsilon}_i^2}{\hat{\sigma}^2} = \gamma_0 + \gamma_1 X_{1i} + \gamma_2 X_{2i} + \dots + \gamma_K X_{Ki} + u_i$$

and calculating the regression sum of squares (RSS).

The test statistic is:

$$\chi_{\text{cal}}^2 = \frac{\text{RSS}}{2}$$

Decision rule: Reject the null hypothesis of homoscedasticity:  $H_0: \gamma_1 = \gamma_2 = \dots = \gamma_K = 0$  if:

$$\chi_{\text{cal}}^2 > \chi_{\alpha}^2(K)$$

If the constant variance and normality assumptions are not true, they are often reasonably well satisfied when the observations are transformed by a Box-Cox transformation (Shumay and Stoffar, 2010):

$$x_i^{(\lambda)} = \begin{cases} \frac{x_i^\lambda - 1}{\lambda}, \lambda \neq 0 \\ \ln x_i, \lambda = 0 \end{cases} \quad (79)$$

Choose the value of  $\lambda$  that maximizes:

$$l(\lambda) = -\frac{n}{2} \ln \left[ \frac{1}{n} \sum_{j=1}^n (x_j^{(\lambda)} - \bar{x}^{(\lambda)})^2 \right] + (\lambda - 1) \sum_{j=1}^n \ln x_j \quad (80)$$

If the selected model is inadequate, the three-step model building process is typically repeated several times until a satisfactory model is finally obtained. The final selected model can then be used for prediction purposes (Wei, 1990).

#### **iv) Forecasting**

The last step in time series modeling is forecasting. There are two kinds of forecasts: sample period forecasts and post-sample period forecasts. The former are used to develop confidence in the model and the latter to generate genuine desired forecasts. In forecasting, the goal is to predict future values of a time series,  $x_{t+m}$ ,  $m = 1, 2, \dots$  based on the data collected to the present,  $x = \{x_t, x_{t-1}, \dots, x_1\}$ . Throughout this section, we will assume  $x_t$  is stationary and the model parameters are known.

#### **Minimum mean square error forecasts**

ARIMA (p, d, q) process can be written as

$$\phi(B)\nabla^d x_t = \theta(B)w_t \quad (81)$$

Forecasting a value  $x'_{t+m}$ ,  $m = 1, 2, 3 \dots$  when we are currently standing at time  $t$  is said to be made  $m$ -step ahead forecast. Then three explicit forms of the model for the observation ( $x'_{t+m}$ ) generated by the ARIMA process may be expressed as follows (Box and Jenkins, 1976):

#### **Directly in terms of the difference equation by**

$$x_{t+m} = \phi_1 x_{t+m-1} + \dots + \phi_{p+d} x_{t+m-p-d} - \theta_1 w_{t+m-1} - \dots - \theta_q w_{t+m-q} \quad (82)$$

**Infinite weighted sum of current and previous shocks (  $w_t$  )**

$$x_{t+m} = \sum_{j=0}^{m-1} \psi_j w_{t+m-j} \quad (83)$$

where:  $\psi_0 = 1$  and  $\psi_j$ 's may be obtained by equating the coefficients in

$$\phi(B)(1+\psi_1 B + \psi_2 B^2 + \dots) = \theta(B)$$

**Infinite weighted sum of previous observations, plus a random shock**

$$x_{t+m} = \sum_{j=0}^{m-1} \pi_j x_{t+m-j} \quad (84)$$

where:  $\sum_{j=1}^{\infty} \pi_j = 1$  and  $\pi_j$ 's may be obtained by equating the coefficients in

$$\phi(B) = (1 - \pi_1 B - \pi_2 B^2 - \dots) \theta(B)$$

Standing at origin  $t$ , we can take a minimum mean square error predictor,  $\tilde{x}_{t+m}$  of  $x_{t+m}$ , which is a linear function of current and previous observations of  $x_t, x_{t-1}, \dots$ , then it will also be a linear function of current and previous shocks .

The minimum mean square error predictor ( $\tilde{x}_{t+m}$ ) for lead time  $m$  is the conditional expectation of  $x_{t+m}$ , at the origin  $t$ . From Eqn. (83), we can obtain

$$\tilde{x}_{t+m} = E(x_{t+m} | x_t, \dots, x_1) = \sum_{j=0}^{m-1} \psi_j E(w_{t+m-j}) \quad (85)$$

Then the mean-square prediction error can be written as (Shumay and Stoffar, 2010):

$$P_{t+m}^t = E(x_{t+m} - \tilde{x}_{t+m})^2 = \delta_w^2 \sum_{j=0}^{m-1} \psi_j^2 \quad (86)$$

To assess the precision of the forecasts, prediction interval can be calculated as:

$$x_{t+m}^t + C_{\frac{\alpha}{2}} \sqrt{P_{t+m}^t}, \text{ Where } C_{\frac{\alpha}{2}} \text{ is chosen to get the desired degree of confidence.}$$

### Forecasting Accuracy Measures

Once forecasts are made they can be evaluated if the actual values of the series to be forecasted are observed. There are some measurements of the accuracy of forecasts. These are root mean square error (RMSE), mean absolute error (MAE), mean absolute percentage error (MAPE) and Theil's inequality coefficient (Theil-U).

$$\text{Mean Square Error (MSE)} = \frac{\sum_{t=1}^n \hat{w}_t^2}{n};$$

$$\text{MeanAbsoluteError(MAE)} = \frac{\sum_{t=1}^n |\hat{w}_t|}{n};$$

$$\text{Root Mean Squared Error (RMSE)} = \sqrt{\frac{\sum_{t=1}^n \hat{w}_t^2}{n}};$$

$$\text{Mean Absolute Percentage Error (MAPE)} = \frac{\sum_{t=1}^n \frac{|\hat{w}_t|}{x_t}}{n};$$

$$\text{Theil's inequality coefficient (U-Statistics)} = \frac{\sqrt{\frac{1}{n} \sum_{t=1}^n (x_t - f_t)^2}}{\sqrt{\frac{1}{n} \sum_{t=1}^n f_t^2} \sqrt{\frac{1}{n} \sum_{t=1}^n x_t^2}};$$

Where:  $n$  – the number of observation,  $x_t$  and  $f_t$  is actual observation for time period  $t$  and forecast for the same period, respectively. The scaling of  $U$  is such that it will always lie between 0 and 1. If  $U = 0$ ,  $x_t = f_t$  for all forecasts and there is a perfect fit; if  $U = 1$  the predictive performance is as bad as it possibly could be. The Theil's inequality coefficient can be decomposed in to bias, variance and covariance proportions.

$$\text{Bias Proportion} = \frac{(\bar{f} - \bar{x})^2}{\frac{1}{n} \sum_{t=1}^n (x_t - f_t)^2} ;$$

$$\text{Variance Proportion} = \frac{(S_f - S_x)^2}{\frac{1}{n} \sum_{t=1}^n (x_t - f_t)^2} ;$$

$$\text{Covariance Proportion} = \frac{2(1-r) S_f S_x}{\frac{1}{n} \sum_{t=1}^n (x_t - f_t)^2} ;$$

Where:  $\bar{f}$ ,  $\bar{x}$  and  $S_f$ ,  $S_x$ , are means and standard deviations of forecast series (F) and actual series(X) respectively, and  $r$  is the correlation between F and X.

The bias and the variance proportions show how far the mean of the forecast series is from the mean of the actual series and how far the variation of the forecast is from the variation of the actual series, respectively. The covariance proportion measures the remaining asymmetric forecast errors. The sum of these three measures would be one. If the forecast is good, the bias and the variance proportions should be small.

MSE, MAE, and RMSE depend on the scale of the dependent variable. These should be used as relative measures to compare forecasts for the same series across different models; the smaller the error, the better the forecasting ability of that model according to that criterion. However, MAPE and Theil's inequality coefficient are scale invariant.

Therefore, this study used mean absolute percentage error (MAPE) and Theil's inequality coefficient (U-statistic).

## ***4. Data Analysis***

### **4.1. Descriptive Analysis**

In this section we will describe rainfall of Dire Dawa, Haramaya and Dengego station and mainly focused on analyzing monthly as well as annually rainfall data of Dire Dawa station recorded by the National Meteorological Agency of Ethiopia (NMAE).

The statistical software package used for most of the analysis is R-15.10.1. For some test and graphs, however, SAS-9.2 was used.

Table-4.1 below and Figure-1 in the Appendix, shows the summary of annual and monthly rainfall. From the table, it can be seen that the mean annual rainfall of Dire Dawa, Dengego and Haramaya are 611mm, 774 mm and 772mm, respectively. The minimum rainfall values were observed most frequently in months of dry seasons in all stations and maximum rainfall values were recorded in months of the main rainy season in all stations. In Dire Dawa, the minimum annual rainfall of 228.3mm was recorded during the year 1984 and the maximum of 719.7 mm in the year 1997. The minimum and maximum annual rainfalls of Dengego were 473.8 mm in 1985 and 1066.2mm in 1997, respectively. While the minimum and maximum annual rainfall of Haramaya was 430.3mm in 1987 and 1043.7mm in 1983, respectively. The amount of rainfall at

Dengego and Haramaya are more or less the same on average in all seasons, and is much higher than that of Dire Dawa over the study period.

**Table 4.1:** Summary of rainfall amount by Station (January, 1982-December, 2011)

Rainfall(in mm)		Dire Dawa	Dengego	Haramaya
Annual	Mean	611	774	772
	Minimum	228.3	473.8	430.3
	Maximum	719.7	1066.2	1043.7
	CV	0.23	0.18	0.19
Monthly	Mean	54.7	64.6	64.4
	Minimum	0.00	0.00	0.00
	Maximum	249.9	303.4	305.2
	CV	1.04	1.0	0.74

The coefficient of variation (CV) provides a measure of year-to-year variation in the data series. NMA (1996) documented that a series with CV less than 0.20 can be considered as less variable, while CV between 0.20 and 0.30 is moderately variable, and CV greater than 0.30 is highly variable.

The monthly rainfall data in all stations show that there is high variability since their CV is greater than 0.30. However, relatively the monthly variability of rainfall of Dire Dawa and Dengego are similar. When we see the overall variability, the high variability occurs during months of dry seasons while the low variability is observed during the rainy season in all stations. Dengego and Haramaya stations have less variability in the annual rainfall. However, the variability of annual rainfall in Dire Dawa during the last 30-year

period is a bit larger than neighboring stations. This may indicate that climate instability is high in Dire Dawa than other stations.

## 4.2. Testing Stationarity

**Visual Inspection:** Regardless of which technique is used, the first step in any time series analysis is to plot the observed values against time. A number of qualitative aspects are noticeable as you visually inspect the graph. The pattern of the time series plot in Figure-2 of the Appendix does not show any systematic upward and downward change about the mean. This indicates that the series is non-seasonally stationary.

From the autocorrelation function plot in Figure-4(Appendix), the presence of seasonality behavior and seasonally non-stationarity of the rainfall series is clear. Because it shows strong seasonal wave pattern at the multiple of seasonal intervals( $s=12$ ) and declining slowly while non seasonal lags are relatively decaying quite rapidly. This can be interpreted as a 6-month seasonal pattern that cycles between summer when there is little to no rainfall, and winter when rain is at its peak.

Therefore, from the time series plot and autocorrelation plot (Figure-2 and Figure-4 in the Appendix), we observed that our series has seasonal variation, and seasonal differencing is needed to make it stationary.

Next, we perform some formal tests of stationarity to confirm the visual inspection of non-seasonal stationarity.

**Unit root test:** The test first impose that the rainfall series has a unit root, versus the series is stationary. We test the null hypothesis using the available Augmented Dickey-Fuller (ADF) test. Such test are used in order to know whether our series are stationary or not at level or after differencing. From Figure-2 in the Appendix, we observe that the rainfall series of Dire Dawa doesn't have a trend and potentially slow- turn around zero. Therefore, we can use the test equation given in Eqn. (64). For our series, maximum lag lengths (P) of the ADF test to remove serial correlation from the residuals of the regressions based on the relationship between the current value and the past value of the series was selected as 12.

If the estimate of  $\Phi$  is nearly zero in the fitted regression eqn. (64), the rainfall series ( $x_t$ ) needs differencing, and if the estimate of  $\Phi < 0$ , then the series is already stationary (Makridakis et al., 1998). The results of ADF test for the monthly rainfall series of Dire Dawa at level (without differencing), after first regular and seasonal differencing are summarized in Table-4.2 below.

**Table 4.2:** Summary of ADF unit-roots test (at level and after first seasonal and regular differencing)

Series	ADF test statistic	1%crit.Value	5%crit.Value	10% crit.Value	Decision
Original series	-0.871	-1.442	-1.976	-2.572	Don't Reject null hypothesis
First regular differencing	-1.412	-4.157	-3.56	-3.217	Don't Reject null
First Seasonally differenced	-4.97	-3.44	-2.87	-2.57	Reject the null

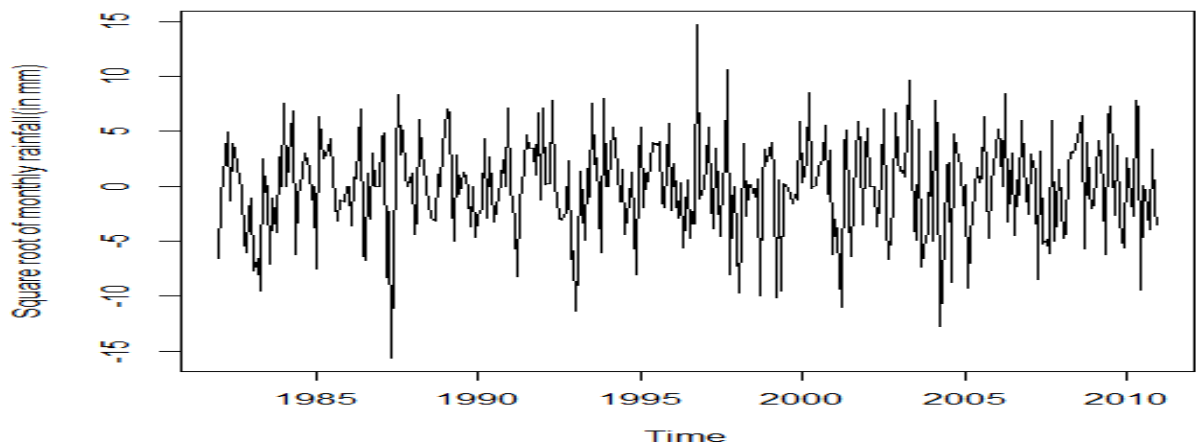
The ADF test statistic for the original and first regular differencing monthly rainfall series given in Table-4.2 are greater than the critical values at 1%, 5% and 10% significance levels. According to these results, the null hypothesis, which has a unit root, for the data sequences should not be rejected at all significance levels. This figure further confirms that original series as well as series obtained after first regular differencing are not stationary. However, as also suggested by visual analysis, after first seasonal differencing, the computed ADF test statistic are smaller than the critical values at 1%,5% and 10% significant levels. This leads to the rejection of the test that there is a unit-root problem.

**Variance Comparison:** The sample variance decreases until a stationary series has been found. Increase in the differencing order tends to increase the variance indicating over differencing. In this study the following results were obtained:

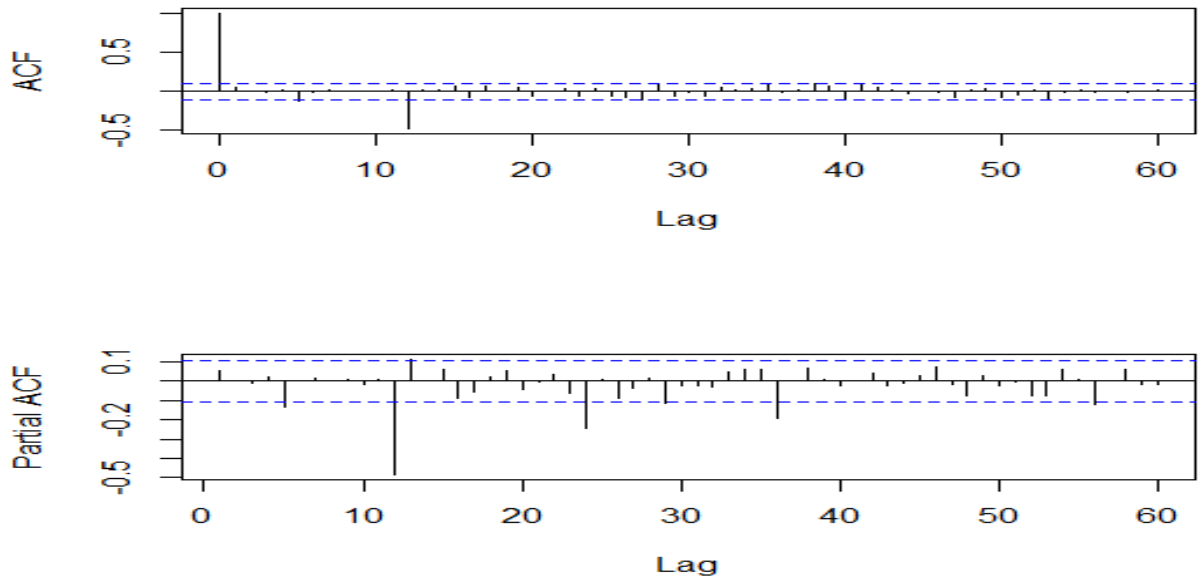
$\text{Var}(\nabla x_t)=5127.7$ ,  $\text{Var}(x_t)=3252.7$ ,  $\text{Var}(\nabla_{12}x_t)=2704.76$ , and  $\text{Var}(\nabla_{12}^2x_t)=7981.4$  values. From these it is clear that:

- i)  $\text{Var}(\nabla x_t) > \text{Var}(x_t)$
- ii)  $\text{Var}(\nabla_{12}^2x_t) > \text{Var}(x_t) > \text{Var}(\nabla_{12}x_t)$ .

These results suggest that non-seasonal first differencing ( $\nabla x_t$ ) has been over-differenced and hence the original series is non-seasonally stationary. The first seasonal differencing would rather be important.



**Figure-1:** Time plot for first seasonal differenced transformed rainfall series (Dire Dawa)



**Figure-2:** ACF and PACF plot for first seasonal differenced transformed rainfall series (Dire Dawa)

Consequently, all tests for stationarity seem to agree and suggest that the first-seasonal differencing of the series make stationarity around a constant mean, which is approximately zero and calculated its standard deviation of 52.05mm (see Figure-1).

Moreover, to determine whether stationarity has been achieved, either by trend removal or by differencing, one may examine the autocorrelation function (ACF) and partial autocorrelation function (PACF) sequence of the residual or processed series (Janacek and Swift, 1993). The sequence corresponding to a stationary process should converge quite rapidly to zero as the value of the lag increases. From this point of view, autocorrelation and partial autocorrelation plot shown in Figure-2 are in support of monthly rainfall series stationary after having first seasonal difference.

To test whether residual from fitted model come from normally distributed series, we use histogram, QQ-plot of the residual and Shapro-wilks test. The histogram and QQ-plot of the residual shown in Figure-5(a) (Appendix), show skewed rather than normally distribution. Shipiro-Wilks test also confirms this fact since it results in P-value of 7.03e-09, which is smaller than 0.05. This implies that the residuals from the fitted models do not come from a normal distribution. This therefore, suggests some kind of transformation for our monthly rainfall series to achieve normality.

By applying Box-Cox transformation defined in Eqn.(79) to the monthly rainfall series, the normality assumptions of residual will be achieved using optimum value  $\lambda = 0.5$ , which maximizes the likelihood function defined in Eqn.(80) with various iteration by the help of SAS 9.2-software. Figure-5(b) in the Appendix reveals that after this transformation the problem of non-normality seems dealt with. The Shapiro–Wilk test result (p-value=0.397) is also in support of the normality of transformed series.

Since normality is not fulfilled for the original series, three stages model building is performed with square root transformed data with the help of ACF and PACF plot shown in Figure-2 above.

## 4.3. Model Building for monthly rainfall series

Fitting a model to time series data involves plotting the data, transforming the data when appropriate, identifying the dependence orders of the model, parameter estimation, diagnostics tests, and model choice. In this section, a univariate SARIMA methodology is used to model monthly rainfall of Dire Dawa.

### 4.3.1. Model Identification

Once the degree of differencing has been determined, the autoregressive and moving-average orders are selected by examining the sample autocorrelations and sample partial autocorrelations.

To use the sample autocorrelation and sample partial autocorrelations functions for tentative model parameters identification, we consider the ACF and PACF shown in Figure-2 above. Using Table-3.1 and Table-3.2 as a guide, preliminary values of  $p$ ,  $q$ ,  $P$  and  $Q$  are chosen. Because we are dealing with estimates, it will not always be clear whether the sample ACF or PACF is tailing off or cutting off. Also, two models that are seemingly different can actually be very similar (Shumway and Stoffer, 2010). With this in mind, we should not worry about being so precise at this stage of the model fitting. At this stage, a few preliminary values of  $p$ ,  $q$ ,  $P$  and  $Q$  should be at hand, and we can start estimating the parameters.

First, concentrating on the seasonal lags, the characteristics of the ACF and PACF of our transformed series in Figure-2 tend to show a strong peak at  $h = 12$  in the autocorrelation

function, combined with peaks at  $h = 12, 24, 36$  in the partial autocorrelation function. Hence it appears that either:

- (i) the ACF is cutting off after lag 12 and the PACF is tailing off in the seasonal lags,
- (ii) the ACF and PACF are both tailing off in the seasonal lags.

Table 3.2 suggests either (i) a seasonal moving average of order  $Q = 1$ , or (ii) due to the fact that both the ACF and PACF may be tailing off at the seasonal lags, perhaps both components,  $P = 1$  and  $Q = 1$ , are needed.

To identify the between-season model, we focus the lags  $h=1, 2 \dots 11$  and identify order based on Table-3.1.

First, we set the ACF to be tailing-off and the PACF to be cut-off after lag 5, we identify  $p=5$  and  $q=0$ . Also it is possible to think of the PACF to be tailing-off and the ACF to cut-off after lag 5, leading to identify  $P=0$  and  $Q=5$ .

Fitting the four models suggested by these observations, we obtain:

$$\text{SARIMA}(0, 0, 5) \times (0, 1, 1)_{12}$$

$$\text{SARIMA}(5, 0, 0) \times (0, 1, 1)_{12}$$

$$\text{SARIMA}(0, 0, 5) \times (1, 1, 1)_{12}$$

$$\text{SARIMA}(5, 0, 0) \times (1, 1, 1)_{12}$$

### 4.3.2. Parameter Estimation

In this section, we assume we have  $n$  observations  $(x_1, \dots, x_n)$  from a causal and invertible Gaussian ARMA  $(p, q)(P, Q)$  process in which initial order of parameters,  $p$ ,  $q$ ,  $P$  and  $Q$  are known. Our goal is to estimate the value of parameter:  $\theta_1, \dots, \theta_q, \phi_1, \dots, \phi_p, \Theta_1, \dots, \Theta_Q, \Phi_1, \dots, \Phi_P$ . Estimating the parameters for Box-Jenkins models follows a non-linear estimation and parameter estimates are usually obtained by maximum likelihood method, which is asymptotic for any time series (Brockwell and Davis, 1996). Hence, we use maximum likelihood estimation method for transformed total monthly rainfall to estimate the parameters. The results are summarized in Table- 4.3 below.

**Table 4.3:** Summary of Parameter Estimates and selection criteria

Model	Parameter	Estimate	Std.Error	T-ratio	P-value	Criteria
(i)	$\mu$	-0.002	0.039	-0.05	0.959	AIC=1863.206
	$\theta_5$	0.165	0.053	3.09	0.002	SBC=1874.763
	$\Theta_{12}$	0.775	0.037	20.88	<.0001	$\delta_k^2 = 11.81$
(ii)	$\mu$	-0.002	0.039	-0.05	0.959	AIC=1860.17
	$\phi_5$	-0.181	0.053	-3.41	0.0006	SBC=1870.73
	$\Theta_{12}$	0.777	0.037	20.99	<.0001	$\delta_k^2 = 10.71$
(iii)	$\mu$	-0.002	0.039	-0.05	0.959	AIC=1865.19
	$\theta_5$	0.165	0.053	3.09	0.002	
	$\Theta_{12}$	0.772	0.037	20.88	<.0001	SBC=1879.61
	$\Phi_{12}$	-0.007	0.071	0.09	0.925	$\delta_k^2 = 12.19$
(iv)	$\mu$	-0.002	0.039	-0.05	0.959	AIC=1864.13
	$\phi_5$	-0.181	0.053	-3.41	0.0006	
	$\Theta_{12}$	0.775	0.037	20.88	<.0001	SBC=1878.54
	$\Phi_{12}$	-0.011	0.071	-0.16	0.877	$\delta_k^2 = 11.85$

Table-4.3 above displays the list of the parameters for each temporally entertained model. For each model parameter, the table presents the estimated value, standard error,  $t$ -value, AIC, SBC and variance ( $\hat{\delta}_k^2$ ) for the estimate. As indicated by McDowall et al., (1980), parameters must differ significantly from zero and all significant parameters must be included in the model. The T-ratios ( $t_{cal}$ ) for associated parameter estimate were

compared with the critical value of obtained from the t-distribution. The result indicated that the seasonal and non-seasonal moving averages as well as non-seasonal autoregressive parameters are all significant since their p-values is smaller than 0.05 and should be retained in the model. However, the constant ( $\mu$ ), non-seasonal autoregressive parameters ( $\phi_1, \phi_2, \phi_3, \phi_4$ ) and seasonal autoregressive parameters ( $\phi_{12}$ ) in all selected models are insignificant since their p-values is greater than 0.05. Therefore, it can be conclude that the parameters is statistically close to zero and that it should be omitted from the model.

The AIC, SBC and variance of the estimate ( $\hat{\delta}_k^2$ ) in Table-4.3 are computed according to Eq. (71 and 72) and are used to compare selected models fit best to the monthly rainfall series. The model with the smaller information criteria is said to fit the data better. Since SARIMA (5, 0, 0)\*(0, 1, 1)<sub>12</sub> model has lower AIC, SBC and variance of estimate than other model, it fits the transformed monthly rainfall series of Dire Dawa better, and can be further analyzed.

The correlations of the parameter estimates are shown in Table-4.4, which can be used to assess the extent to which collinearity may have influenced the results. If two parameter estimates are very highly correlated, you might consider dropping one of them from the model, but in our case correlation between seasonal moving average and non-seasonal autoregressive parameter is -0.079, which indicate that less correlation is observed between parameter estimates.

**Table 4.4:** Correlations of Parameter Estimates for the fitted model

Parameter	$\hat{\sigma}^2$	$\theta_{12}$	$\phi_5$
$\hat{\sigma}^2$	1.000	0.032	-0.017
$\theta_{12}$	0.032	1.000	-0.079
$\phi_5$	-0.017	-0.079	1.000

### 4.3.3. Diagnostic Checking

In this section, we will assess how well the selected model fit the actual rainfall data. If the model fits the data well, the residuals of the fitted model are random (Chatfield, 1991). In ARIMA modeling, the selection of a best model to data is directly related to how well the residual analysis is performed (Kadri et al., 2005). Therefore, several diagnostic statistics and plots of the residuals can be used to examine the goodness of fit of the selected model to the data.

The first step in this stage is plotting the standardized innovations (or residuals) of monthly rainfall of Dire Dawa. If the model fits well, the standardized residuals should behave as an identically and independently distributed sequence with mean zero and variance one (Shumay and Stoffar, 2010). The time plot should be inspected for any clear departures from this assumption. Inspection of the time plot of the standardized residuals for square root transformed rainfall series in Figure-3 below, which is obtained by computing Eqn. (73), shows no clear patterns (trend or seasonality behavior).

Under the assumption that residuals follow a white noise process, the standard errors of the residual ACF, in our case are approximately equal to  $1/\sqrt{360}$ . Thus, under the test that residual follows a white noise process, roughly 95% of the residual autocorrelations should fall within the range of  $\pm 1.96/\sqrt{360}$ . It is clear, as shown in Figure- 3 below that there is no pattern in the residuals autocorrelation plot for the selected model, which means there is no autocorrelation coefficient which lies outside the two standard errors significantly for the fitted models. Therefore, this indicates that residual for the fitted models are not significantly different from a white noise.

There are many statistical tests used for diagnostic checking of randomness. In this study, the Ljung-Box Q-statistic, Turning-point test and runs tests are used for checking independence of residual.

**The Ljung-Box Q-statistic** can be employed to check independence; here we can perform a general test that takes into consideration, the magnitudes of residual ACF as a group instead of individual visual inspection of the sample autocorrelations. For example, it may be the case that, individually, they are small in magnitude, say; each one is just slightly less than  $2/\sqrt{n}$  in magnitude, but, collectively, the values are large. A test of this hypothesis can be done for the model adequacy by choosing a level of significance and then comparing the value of calculated chi-square with the critical value. The Q-statistic for each group of six lags are computed using Eq. (75) and the results are summarized in Table-4.5 and Figure-3 below.

**Table 4.5:** Residual white noises check with Ljung-Box test for the fitted model

To Lag	Q- Statistic	DF	P-Value	-----Autocorrelations-----					
6	7.60	4	0.1072	0.037	0.033	-0.021	-0.017	0.008	-0.032
12	12.24	10	0.2695	-0.064	-0.034	-0.023	0.040	0.073	-0.006
18	17.67	16	0.3435	0.095	0.024	-0.007	-0.063	-0.025	-0.025
24	22.56	22	0.4267	0.014	-0.074	-0.033	0.069	0.030	0.024
30	27.21	28	0.5069	-0.021	-0.036	-0.084	0.042	-0.038	-0.012
36	32.48	34	0.5420	-0.062	-0.008	0.016	0.070	0.066	0.012
42	39.07	40	0.5120	0.055	0.079	0.043	-0.055	0.042	0.030
48	40.61	46	0.6970	0.019	-0.024	-0.027	-0.007	-0.042	0.018
54	45.06	52	0.7411	0.059	-0.039	0.009	-0.006	-0.073	-0.020
60	47.08	58	0.8469	-0.000	-0.019	-0.029	-0.055	-0.023	-0.001
66	50.60	64	0.8885	0.026	0.012	0.065	0.029	0.036	-0.032
72	57.42	70	0.8594	-0.019	0.051	-0.044	-0.078	-0.013	-0.067

From Table-4.5 and Figure-3, we can observe that the p-value is greater than 0.05 for all lags, which implies that the white noise hypothesis is not rejected.

**Table 4.6:** Result of Independence, Homoscedasticity and Normality tests for residual of fitted model

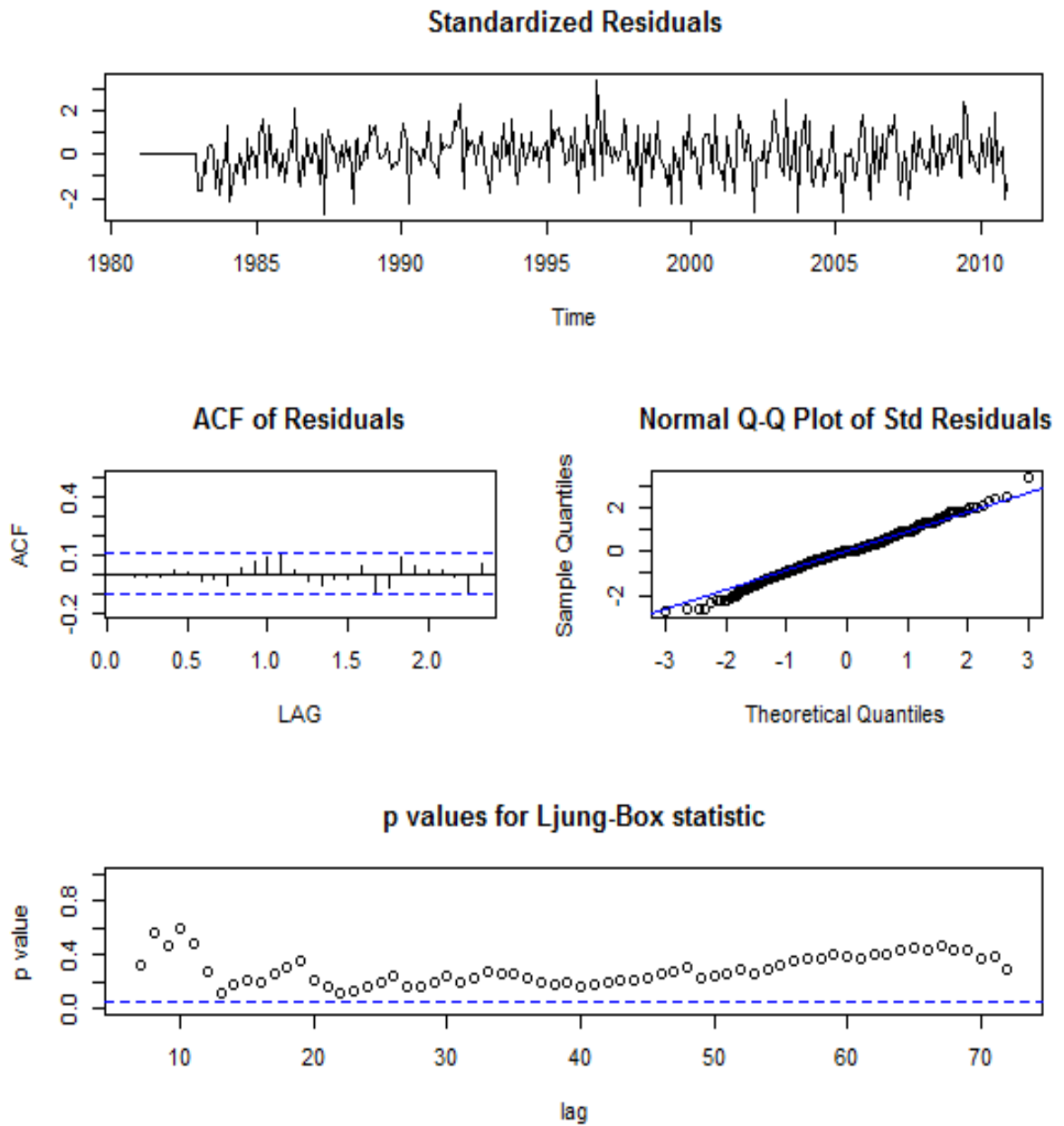
	Runs test	Turning point test	Goldfeld-Quandt Test	Breusch and Pagan test	Shapiro-Wilk test
Test statistic	-4.67	0.07	0.79	0.18	0.887
Tabulated value	$\pm 1.96$	$\pm 1.96$	1.00	1.00	0.996
P-value	0.297	0.354	0.061	0.072	0.397

Two alternative methods, namely runs and turning point tests, are applied to check the independence assumption of residuals of the fitted model. The results of these tests were presented in Table 4.6 above. The fitted model was found to be consistent with the independence assumption for both test methods since their test statistic is smaller than the critical values. These results cope up with plot of the autocorrelation of residual in Figure-3 below and Table 4.5.

For the selected best models, the results related to the normality of residuals using Shapiro–Wilk tests are given in Table 4.6. Since p-value is greater than 0.05, indicating the residuals of the fitted models is normally distributed. In addition to these tests, Figure 5(b) (Appendix) shows the histograms and QQ-plot of the residuals. As expected, the curves significantly reflect a normal distribution.

Test statistics value of the Goldfeld-Quandt(G-Q) Test and Breusch and Pagan(B-P) test, for the homoscedasticity of the residuals are also presented in Table 4.6. All calculated values are found to be smaller than the respective critical values, which indicating that the residual variance is constant.

Therefore, the hypothesis that the residuals are white noise cannot be rejected indicating that the fitted model is adequate. That is, SARIMA (5,0,0) x (0,1,1)<sub>12</sub> model is adequate for modeling the square root transformed monthly rainfall series of Dire Dawa.



**Figure-3:** Diagnostics of the residuals from the fitted model

#### 4.3.4. Forecasting

Since the model diagnostic tests show that all the parameter estimates are significant and the residual series is white noise, the estimation and diagnostic checking stages of the modeling process are complete. We can now proceed to forecasting the rainfall series with fitted SARIMA (5, 0, 0)\*(0, 1, 1)<sub>12</sub> model. Forecasting refers to the process of predicting future rainfall values from a known time series. In this study, forecasting is performed as follows:

According to Eqn.(58),the SARIMA (5, 0, 0) x (0, 1, 1)<sub>12</sub> model can be written as

$$(1-\phi_5 B^5)(1-B^{12})y_t=(1+\theta_{12}B^{12})e_t \quad (87)$$

This equation can also be multiplied out and rewritten in a form that is used in forecasting as shown in Eq. (88) below.

$$y_t=y_{t-12}+\Phi_5(y_{t-5}-y_{t-17})+e_t+\theta_{12}e_{t-12} \quad (88)$$

$$\text{Where: } B^5 y_t=y_{t-5} \text{ and } (1-B^{12}) y_t=y_t-y_{t-12}$$

The above equation can be re-expressed as:

$$y_{t+m}=y_{t+m-12}+\Phi_5(y_{t+m-5}-y_{t+m-17})+e_{t+m}+\theta_{12}e_{t+m-12} \quad (89)$$

After substituting the estimated parameter values in Eq. (89) above, we obtain the following equation:

$$\hat{y}_{t+m}=\hat{y}_{t+m-12}-0.181(\hat{y}_{t+m-5}-\hat{y}_{t+m-17})+\hat{e}_{t+m}+0.777\hat{e}_{t+m-12} \quad (90)$$

Forecasts are to be made at the origin,  $t = \text{Dec}, 2011$  for lead times  $m = 1, \dots, 24$  for the total monthly rainfall series in the coming 24-month. For example, the one-step ahead forecast at the origin,  $t = \text{Dec}, 2011$ , which give us forecast of actual series for the month of January, 2012 are given by:

$$\hat{x}_{Dec,2011+1} = x_{Jan,2012} = \hat{x}_{Dec,2011+1} - 0.033(x_{Dec,2011+4} - \hat{x}_{Dec,2011+6}) + \hat{w}_{Dec,2011+1} + 0.604\hat{w}_{Dec,2011+1}$$

This way we find 24 month-step ahead forecast and prediction interval for total monthly rainfall series of Dire Dawa, which is presented in Figure-4 below.

#### **4.3.4.1. Forecasting accuracy Evaluation**

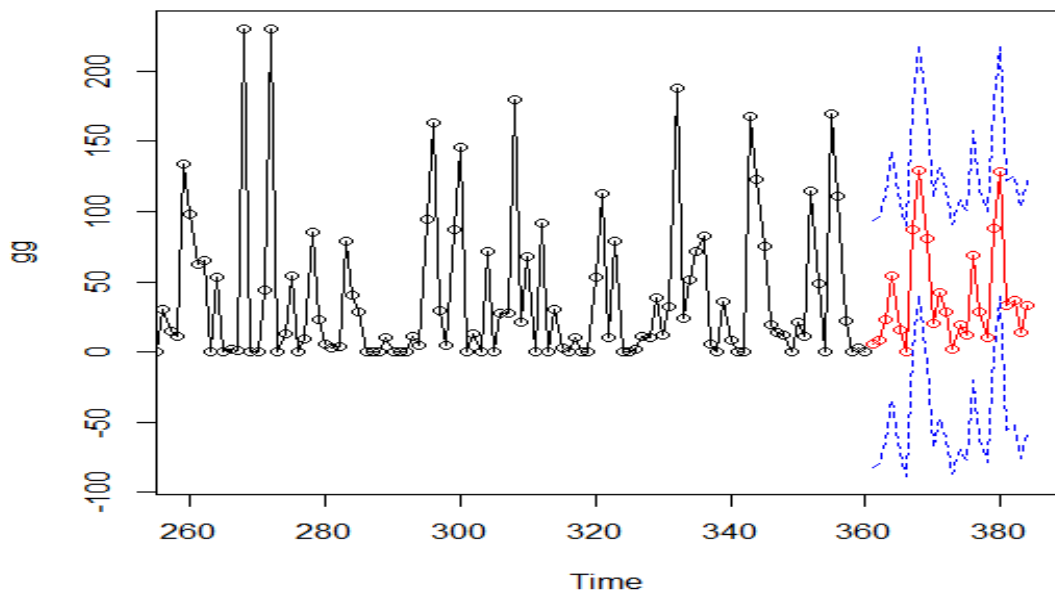
If the fitted SARIMA  $(5, 0, 0)*(0, 1, 1)_{12}$  model has to perform well in forecasting, the forecast error will be relatively small. The accuracy of forecasts is usually measured using root mean square error (RMSE), mean absolute error (MAE), Mean absolute percentage error (MAPE) and Theil's inequality coefficient (Theil-U). The result shows that the Mean Absolute Percentage Error (MAPE) turn out to be 3.56%, which is relatively less than 4% and Theil's inequality coefficient (U-statistic) turn out to be 0.018, which is relatively close to zero. Besides this result, the bias and variance proportion are also very small, which are 0.047 and 0.001, respectively. Thus, measures indicate that the forecasting inaccuracy is low (see Table-4.7). In addition, it is summarized in Table 4.8 that, the actual and forecasted values of monthly rainfall series from Jan, 2011 to Dec, 2011, more or less supports the value of measures.

**Table 4.7:** Forecasting Accuracy Statistic

Forecast Sample: January, 2012 to December, 2013	
Accuracy Measures	Variable
	Monthly rainfall
Root Mean Squared Error	46.52
Mean Absolute Error	23.02
Mean Absolute Percent Error	3.562
Theil Inequality Coefficient	0.018
Bias proportion	0.047
Variance proportion	0.001

**Table 4.8:** Actual and fitted values of the series (January, 2011-December, 2011)

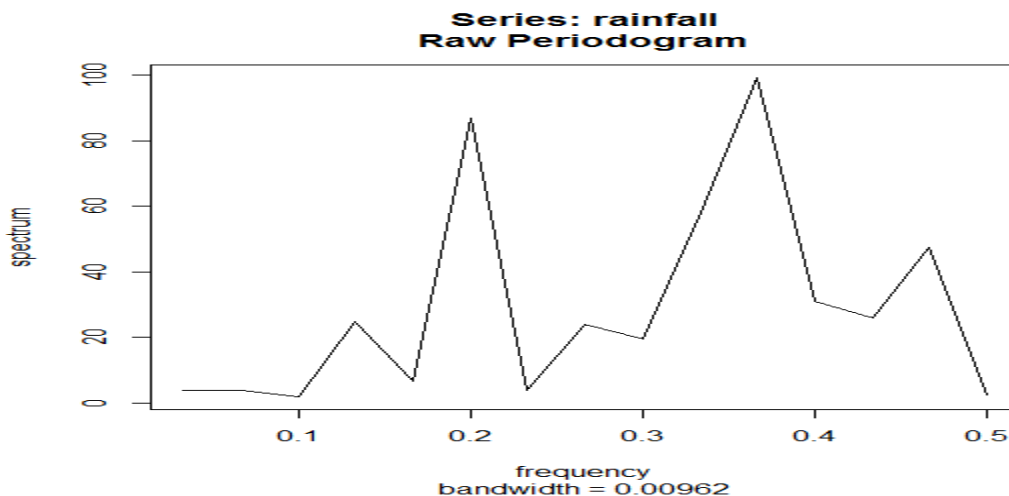
Date	Actual	Forecast	Residual
Jan, 2011	0	3.2	-3.2
Feb, 2011	20.8	18.9	1.9
Mar, 2011	10.7	9.7	1.0
Apr, 2011	115.1	120.3	-5.2
May, 2011	48.9	53.1	-4.2
Jun, 2011	0	1.2	-1.2
Jul, 2011	169.8	156.3	13.5
Aug, 2011	111.1	99.8	11.3
Sep, 2011	22.4	38.7	-16.3
Oct, 2011	0	0.7	-0.7
Nov, 2011	3	7.3	-4.3
Dec, 2011	0	6.2	-6.0



**Figure-4:** Forecast plot for total monthly rainfall of Dire Dawa.

## 4.4. Result from spectral analysis

Quite often, hydrologic phenomenon depicts cyclic and stochastic processes. The periodogram method has dominated hydrology for years because of underlying periodicities in hydrologic processes (Yevjevich, 1972). The power spectral Analysis is applied for annual and monthly rainfall series of Dire Dawa. The first step in estimation of the spectrum by the smoothed periodogram method is subtraction of the sample mean and removing any obvious trend. Figure- 7 of the Appendix show that rainfall series expressed as departures from its mean and raw periodogram of this series can be seen in Figure 5 below. Each point of raw periodogram represents the variance of the rainfall series contributed by a frequency range centered at the point.



**Figure-5:** Raw periodogram of annual rainfall of Dire Dawa

Since raw periodogram is a wildly fluctuating estimate of the spectrum with high variance, as shown in Figure-5 above, the periodogram must be smoothed to ensure stable estimate. Since proper amount of smoothing is somewhat subjective and depends

on the characteristics of the data, generally, it is a good idea to try several bandwidths that seem to be compatible with the general overall shape of the spectrum, as suggested by the periodogram (Shumay and Stoffar,2010).

Considering the tradeoff in smoothness, stability and resolution in selecting widths of Daniell filters, which leads to span of length,  $m= 3$  and  $m=15$  as a suggested value of smoothing for our annual and monthly rainfall, respectively. The result is displayed in Figure-6 and (Figure-8 in the Appendix).

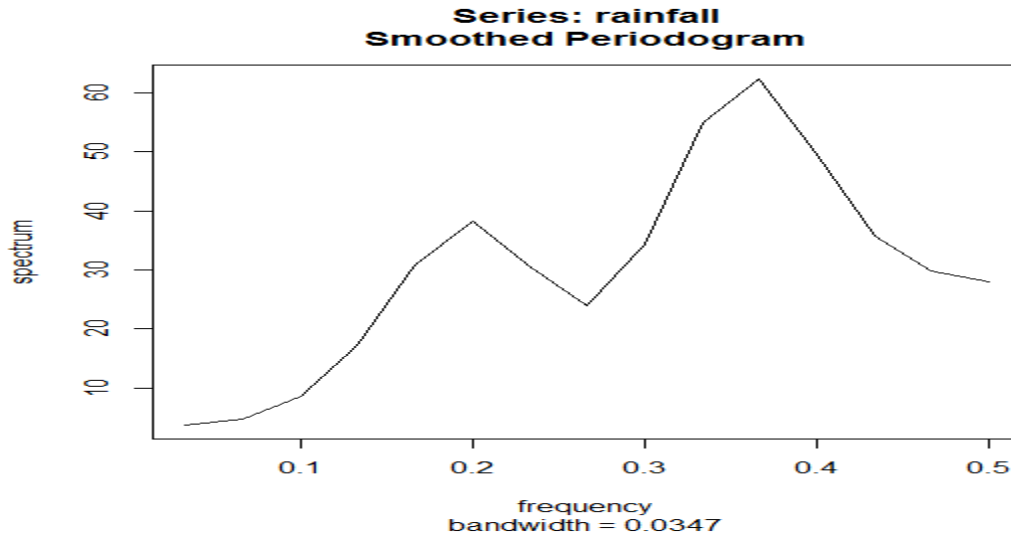
The width of the center mark on the 95% confidence interval indicator indicates the band width (Venables and Ripley, 1999) cited also in (Alemrew and Eshetu, 2009). The bandwidth, using Eqn. (32) is equal to 0.1 cycles per year and 0.0417 cycles per month for the spectral estimator. This bandwidth means we are assuming a relatively constant spectrum over about 20% and 8.3% for monthly and annual rainfall with entire frequency interval (0, 0.5), respectively.

By using degrees of freedom ( $d_f$ ) =  $2L=6$  and 30 for annual and monthly rainfall respectively, we obtain  $\chi_6^2(0.025) =1.24$  and  $\chi_6^2(0.975) =14.45$  and  $\chi_{30}^2(0.025) =46.98$  and  $\chi_{30}^2(0.975) =16.79$ . Substituting these into Eqn.(33),we can construct approximate  $100(1-\alpha) \%$  confidence intervals of spectral density for the frequency bands identified as having the maximum power, as shown in Table 4.9 below.

**Table 4.9:** Confidence Intervals for the Smoothed Spectra of the annual and monthly rainfall Series

Series	Frequency	Period	power	Lower	Upper
Annual rainfall	0.37	2.5years	62.3	39.34	255.75
	0.2	5years	38.2	16.76	156.74
Monthly rainfall	0.083	12month	9.66	6.17	17.25
	0.167	6month	5.51	3.52	9.85
	0.250	4month	7.61	4.86	13.59

If the lower confidence limit for the spectral value is greater than the baseline level at some predetermined level of significance, we may claim that frequency value as a statistically significant peak (Shumay and Stofar, 2010). As we observed from Table-4.9 above, an approximate 95% confidence interval for the spectrum  $f_s(0.37)$  is [39.34, 255.75], which is again too wide to be of much use, but we do notice, that the lower value 39.34 is higher than any other periodogram ordinate shown in Figure-6, so it is safe to say that this value is statistically significant. Similarly, an approximate 95% confidence interval for the spectrum  $f_s(0.083)$  is [6.17, 17.25].



**Figure-6:** Smoothed periodogram estimate of spectrum of annual rainfall of Dire Dawa.

The smoothed spectra shown in Figure-6 provide a sensible compromise between the noisy version, shown in Figure- 5, and a more heavily smoothed spectrum, which might lose some of the peaks.

As presented in Figure-6 above, the peak with an average period of 2.5years, contributing the highest percent to total variance of the annual rainfall series, corresponds to the existence of "strong" peak at a frequency band centered at 0.37 cycles per year. This indicates that the most dominant cyclical component whose periodicity is 2.5 years is observed in annual rainfall, which can be related to Quasi-biennial oscillation (2-3 yr cycles). There is also a subsidiary peak corresponding to a cycle of 5 years with associated frequency band centered around 0.2 cycles per year. These can also be attributed to the El Nino phenomenon.

If monthly data are used their spectrum will exhibit narrow and sharp peaks at seasonal frequencies, as shown in Figure-8 of Appendix. For such type of data, commonly occurred periodic oscillation: Semi-Annual Oscillation (SAO), Annual Oscillation (AO), Quasi-Biennial Oscillation (QBO) and El-Nino/Southern Oscillation (ENSO) are approximately obtained by 4-7 months, 10-14 months, 22-32 months, and 40-66 months, respectively (Sinta and Hariadi, 2003).

The peak with an average period of 12 month contributes a higher percent to total variance of the monthly rainfall series. This indicates that the most dominant peak observed in monthly rainfall of Dire Dawa is annual Oscillation (AO). There is also a subsidiary peak corresponding to a periodicity of 4 month and 6 month with associated frequency band centered around 0.25 and 0.167 cycles per month, respectively. These can also be attributed to the Semi-Annual Oscillation (SAO).

It has been well documented that the inter-annual variability of rainfall in Ethiopia are strongly related to the ENSO phenomena (Haile, 1988; Funk et. al., 2005, NMA, 2007, Bekele, 1997). This oscillation was used for estimating possible future rainfall extreme events conditions (Alemrew and Eshetu, 2009, Yilma et al., 1994). The possibility of using prominent peaks in the spectrum to predict the long-range behavior of the rainfall is attractive (Sinta and Hariadi, 2003). Oduro and Adukpo (2006) predict that a rainfall extremes event in Ghana recurs every 5.6 years. Stringer (1972) also estimated that the climatic events recur every 2 to 2.5 year in Africa. Alemrew and Eshetu (2009) inferred that drought recurs in Addis Ababa between 10 to 11 years.

## 4.5. Result from cross-spectral analysis

In this study, the cross-spectrum analysis are used to provides a means of determining the contributions of fluctuations in various frequency bands to covariance quantities such as the fluctuations in rainfall series at different nearby stations. Coherence square of cross-spectral analysis is the principal tool that will be used in our analysis of pairs of rainfall series, which will isolate those frequencies that are important at both stations.

Figure-9 in the Appendix shows the squared coherence between rainfall series of Dire Dawa with Haramaya and Dengego stations with  $L_h = 15$  and degree of freedom (df) =  $2(15) = 30$  according Eqn.(38). And  $F_{2,30}(0.05)=3.34$  at 5% significance level . Hence, the hypothesis of no coherence is rejected for the values of estimated square coherence (Eqn.44) that exceed  $C_{0.05} = 0.16$ (see eqn.46).

As it can be shown in the same Figure that, at lower seasonal frequency, Dire Dawa rainfall characteristics are significantly associated with both Haramaya and Dengego rainfall, that is, strongly coherent. However, the degree of relatedness is a bit higher in Dengego than Haramaya rainfall to the rainfall of Dire Dawa.

## ***5. Conclusion and Limitation of the study***

### **5.1. Conclusion**

In this study, 30 years annual and monthly rainfall records were analyzed, using data from Dire Dawa, Dengego and Haramaya weather stations in the Eastern Ethiopia, mainly to study the rainfall pattern of Dire Dawa. Univariate Box-Jenkins methodology and Spectral analysis were used to examine the modes of variation of rainfall data.

Based on the overall results of the research, the following conclusions could be drawn:

- A time series model for monthly rainfall series of Dire Dawa was adjusted, processed, diagnostically checked and lastly SARIMA model is established with a 95% prediction interval that can adequately be used to forecast 2 years monthly rainfall values. Further results reveal that there is a tendency of relatively increasing pattern of monthly rainfall over the forecast period from January 2012 to December 2013.
- Predominant cycles with periodicities of 2.5 years were found in the annual rainfall of Dire Dawa. On other hand, Annual oscillations dominate monthly rainfall of the region. In general, it can be inferred that extreme rainfall events recurs every 2-3 year in Dire Dawa.
- Rainfall pattern of Dengego is found to be more related to variability pattern of rainfall in Dire Dawa as compared to Haramaya.

## **5.2. Limitation of the study**

Many studies investigated that sunspot numbers as an indicator for the various aspects of world weather or climate. Moreover, exploring their statistical connection may lead to better understanding and identification of some temporal and spatial patterns of rainfall. However, the data is not included because of the unavailability of data.

## Reference

- Adane, A., (2009). Hydrological Drought Analysis - Occurrence, Severity, Risks: The Case of Wabi Shebele River Basin, Ethiopia. PHD. dissertation.
- Adejuwon, J.O., (2010). A spectral analysis of rainfall in Edo and Delta States (formerly Mid-Western Region). *International Journal of Climatology*, **Vol.31**: 2365–2370pp, Nigeria.
- Admasu, G., (1989). Regional flood frequency Analysis. Technical Report. Royal Institute of Technology Stockholm, Sweden.
- Akaike, H., (1978). A Bayesian analysis of the minimum AIC procedure.
- Alamerew, B. and Eshetu, W., (2009). Assessment of Local Climate in Addis Ababa, *Journal of Ethiopian Statistical Association*, **Vol. 18**, 55-57 PP.
- Al-Ansari, A., Al-Shamali B. and Shatnawi A., (2006). Statistical Analysis of rainfall Records at three Major Meteorological Stations in Jordan, Al-Mararah University special publications, **Vol.12**.
- Amha, G., (2010). Modelling and forecasting monthly rainfall in Tigray region. Case of Mekelle station. (Unpublished M.Sc. Thesis), Department of Statistics, Addis Ababa University, Ethiopia.
- Ayoade, J.O., (1973). Annual rainfall trends and periodicities in Nigeria. *Nigeria Geographical Journal*, **Vol.16**, 167–176pp.
- Bayazit, M., (1981). Statistical Methods in Hydrology. *Istanbul Technical University Press*, No. 1197, Istanbul.
- Bekele, F., (1997). Ethiopian use of ENSO information in its seasonal forecast. *Internet Journal for African Studies (IJAS)*, **Vol.1**.
- Bewket, W., and Conway, D., (2007). A note on the temporal and spatial variability of rainfall in the drought-prone Amhara region of Ethiopia. *Int. J. Climatology*. **Vol.27**, 1467-1477pp.
- Blackman, R.B., and Tukey, J.W., (1959). The measurement of power spectra, from the point of view of communications engineering, New York, Dover.
- Bloomfield, P., (2000). Fourier analysis of time series: an Introduction, 2<sup>nd</sup> ed.: New York, John Wiley & Sons, Inc.

- Box, G. E. P., and Jenkins, G.M., (1976). *Time Series Analysis: Forecasting and Control*; Holden day.
- Box, G.E.P. and Pierce, D.A., (1970). Distributions of residual autocorrelations in autoregressive integrated moving average models. *J. American Stat. Assoc.*, **Vol.72**, 397-402pp.
- Box, G. E. P. and Cox, D. R.,(1964). „An analysis of transformation“.
- Breusch, T. and Pagan, A.,(1979). “A Simple Test of Heteroscedasticity and Random Coefficient Variation”. *Econometrica*, 47, 1287-1294.
- Brockwell, P.J., Davis, R.A. and Rockwell, P.J.,(2000). *Introduction to time series and forecasting*, 2<sup>nd</sup> edition, Springer texts in statistics.
- Brockwell, P.J. and Davis, R. A.,(1996). *Introduction to Time Series and Forecasting*, 2<sup>nd</sup> ed. , Springer, New York.
- Brooks, C., (2008). *Introductory Econometrics for Finance*, 3<sup>rd</sup> ed., *Cambridge University Press*, UK.
- Burroughs ,W.J., (1992). *Weather Cycles, Real or Imaginary?*, *Cambridge University Press*: Cambridge.
- Chatfield, C., (1991). *The Analysis of Time Series- An Introduction*, 4<sup>th</sup> ed., Chapman and Hall: London.
- Cooley, J.W., and Tukey, J.W., (1965). An algorithm for the machine computation of complex Fourier series: *Math. Comput.*, **Vol.19**, 297-301pp.
- Cromwell, J.B., Labys, W.C. and Terraza, M., (1994). *Univariate Tests for Time Series Models*. A Sage Publications, Series/number: 07-99, 96 pp., London.
- Desalegn, R.,(1991). *Famine and survival strategies: a case study from northeast Ethiopia*. Nordiska Afrikainstitutet, Uppsala.
- De Gooije, G. and Hyndman, J.R.,(2006). 25 Years time series forecasting. *International Journal of forecasting*, **Vol.22**.
- Dickey, D.A. and Fuller, W.A.,(1981). Likelihood Ratio Statistics for Autoregressive Time Series with a Unit Root", *Econometrica*, **Vol.49**, 1057-1072pp.

- Diebold, F.X. , Kilian, L. and Nerlove, M.,(2006) .Time Series Analysis. Working Paper No. 06-011, University of Maryland, College Park.
- Dire Dawa Administration (DDA)., (2006). Integrated Dvelopment Management study, Dire Dawa, Ethiopia.
- Dire Dawa Environmental Protection Authority (DDAEP),.(2011). Dire Dawa Administration program of adaptation to climate change.
- Douben, K. (2006). Characteristics of River Floods and Flooding: A global overview. *Irrigation and Drainage* ,55: S9-S21.
- Dulleur, J. W and Kavas, M.L., (1978). Stochastic Models for Monthly Rainfall Forecasting and Synthetic Generation; *Journal of applied Meteorology*, **Vol.17**, 1528-1535 pp.
- Enders, W., (1995). Applied Econometric Time Series, John Wiley and Sons, Inc., New York.
- Elliott, G., Rothenberg, T. J. and Stock, J.H., (1996). "Efficient Tests for an Autoregressive Unit Root", *Econometrica*, **Vol. 64**, 813–836pp.
- Falk, M.,(2005). A First Course on Time Series Analysis -Examples with SAS; University of Wurzburg.
- Federal Disaster Prevention and Preparedness Agency (FDPPA), (2007). Regional summary of Multi Agency Flood Impact assessment of 2006. Addis Ababa, Ethiopia.
- Ferguson, T.S., Genest, C. and Hallin, M.,(2000). "Kendal's Tau for Serial Dependence", *The Canadian Journal of Statistics*, **Vol.28**, 587-604pp.
- Fuller, W.A. (1995). Introduction to Statistical Time Series, 2<sup>nd</sup> ed. New York:Wiley.
- Gibbons, R.D. ,(1994). Statistical Methods for Groundwater Monitoring. John Wiley & Sons, New York.
- Granger,K. and Newbold,J., (1986). Forecasting Economic Time Series; *Academic Press*,USA.
- Greene, W.H.,(2000) Econometric Analysis, Prentice Hall International, Inc., New Jersey, USA.

- Griffiths, C.G.,(1971). The variation with height of the top brightness of precipitating convective cloud. Report of WMO technical report, N0 .237/2000, Switzerland, zananghug.
- Haile, T.,(1988).Causes and characters of drought in Ethiopia. *Ethiopian Journal of Agricultural Sciences*,**Vol. 10**, 85 – 97pp.
- Hamilton, J., (1994). Time Series Analysis. Princeton University Press, New Jersey.
- Harvey, R., Andrew C. and Souza, R.C., (1987).Assessing and Modeling the Cyclical Behavior of Rainfall in North-East Brazil, *Journal of Climate and Applied Meteorology*, **Vol.26**, 1339-1344pp.
- Hipel, K.W. and McLeod, A.I.,(1994). Time Series Modelling of Water Resources and Environmental Systems. Development in Water Science, **Vol. 45**, Elsevier Scientific Publishing Company, Amsterdam, Netherlands.
- Hipel, K.W., McLeod, A.I. and Lennox, W.C., (1977). Advances in Box-Jenkins Modeling, Model Construction", *Journal of Water Resources Research*, **Vol.13**, 567-575pp.
- Janacek, G. and Swift, L.,(1993.) Time Series Forecasting, Simulation, Application, Ellis Horwood, New York,USA.
- Kadri,Y., Ahmet, K. and Fazli, O.,(2005). Testing the Residuals of an ARIMA Model on the Cekerek Stream Watershed in Turkey. *Turkish Journal of Environmental Engineering*,**Vol.29**.
- Kendall, M. and Ord, J.K.,(1990). Time Series, 3<sup>rd</sup> ed. Edward Arnold:156–180.
- Ketema, T.,(1999). Test of homogeneity, frequency analysis of rainfall data and estimate of drought probabilities in Dire Dawa, Eastern Ethiopia. *Ethiopian Journal of Natural Resources*,**Vol. 1**: 125-136pp.
- Laban, A.J.and Ogallo,H., (1986).Stochastic modelling of regional annual rainfall anomalies in East Africa. *Journal of Applied Statistics*, **Vol.13**.
- La'zaro, R., Rodrigo, F.S., Gutie'rrez, L., Domingo, F. and Puigdefa'bregas, J., (2001). "Analysis of a 30-year rainfall record (1967-1997) in semi-arid SE Spain for implications on vegetation". *Journal of Arid Environments*,**Vol. 48**:373-395pp.

- Lehmann, A. and Rode, M.,(2001).Long-Term Behaviour and Cross-Correlation Water Quality Analysis of the River Elbe, Germany". *Journal of Water Resources*, **Vol.35**, 2153-2160.
- Liu, L.M. and William, J.,( 2001). Data mining on time series: An illustrative using fast food restaurant franchise data; scientific, Computing Associate Corp.
- Ljung, G.M. and Box, G.E.P., (1978).”On a Measure of Lack of Fit in Time Series Models", *Biometrika*, **Vol.65**,297-303pp.
- Makridakis, S., Wheelwright, S.C. and Hyndman, R.J., (1998). Forecasting methods and Application; New York: John Wiley & Sons.
- Mersha E., (2002). Determination of Rainfall Cyclicity over Selected Location in Ethiopia. *Ethiopian J. Water Science and Technology*, **Vol.5**.
- Mahsin, M.D. Yesmin, A. and Monira, B.,(2012). Modeling Rainfall in Dhaka Division of Bangladesh Using Time Series Analysis. *Journal of Mathematical Modelling and Application*, **Vol. 1**, No.5, 67-73pp.
- Nicholls, N., (1980). Long-range weather forecasting: value, status and prospects, *Review Geophysical Space Physics*, 18: 771–788pp.
- Nicholson, S. E. and Entekhabi, D., (1986). The quasi-periodic behaviour of rainfall variability in Africa and its relationship to the Southern Oscillation. *Journal of Climate and Applied Meteorology* ,**Vol.34**: 311–348pp.
- NMA., (2007).National Adaptation Programme of Action of Ethiopia (NAPA).Final draft report. National Meteorological Agency, Addis Ababa,Ethiopia.
- NMA.,(1996). Climate & agroclimate resources of Ethiopia. NMSA Meteorological Research Report Series, **Vol.1**, Addis Ababa.
- Oduro, A. K. and Adukpo, D.C.,(2006). Spectral Characteristics of the Annual Mean Rainfall Series in Ghana, *West Africa Journal of Applied Ecology*, **Vol. 9**.
- Osti, R., Tanaka S., and Tokioka, T., (2008). Flood Hazard Mapping in Developing Countries: problems and prospects. *Disaster Prevention and Management*,17(1):104-113
- Pankhurst, R. and Johnson, D.H.,(1988). The great drought and famine of 1888-92 in northeast Africa. In *The ecology of survival: case studies from northeast African history* (pp47-72), Lester Crook Academic Publishing, London.

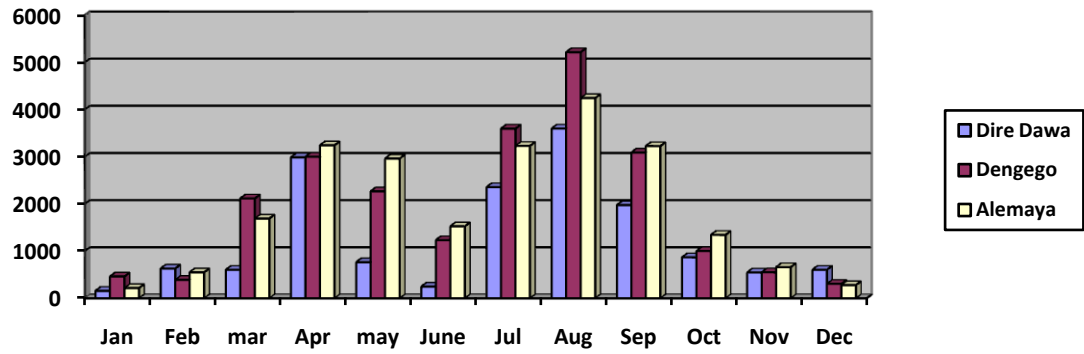
- Pankratiz , A., (1983). Forecasting with Univariate Box-Jenkins: Concepts and Cases; John Wiley & Sons, Inc. New York.
- Schwartz, M.,(1978). Estimating the dimension of a model, Annual Statistics.
- Seifu, A.,(2004). Rainfall Variation and its Effect on Crop Production in Ethiopia.AAU, unpublished M.SC thesis, Ethiopia.
- Seyed, A., Shamsnia, M.,Naeem,S. and Ali, L.,(2011).Modelling weather parameter using stochastic methods(ARIMA Model)(Case Study:Abadeh station,Iran).2011 international conference on Environment and Industrial innovation IPCBEE,12.
- Shapiro, S.S. and Wilk, M.B.,(1965). “An Analysis of Variance Test for Normality (complete samples).” *Biometrika*. **Vol.52**, 591-611pp.
- Shumway, R.H. and Stoffer, D.S.,(2010). Time series analysis and its applications with R Examples, Springer, 3<sup>rd</sup> Ed.
- Sinta, B. S. and Hariadi, T.E.,(2003). The Spectrum Analysis of Rainfall in Indonesia. *Indonesian Journal of Physics*, **Vol.14**, No.3.
- Tamiru, F., (2009). Impact Assessment of global climate change on Some Components of Hydrometeorology in Ethiopia; Kochi University of Technology (unpublished).
- Tesfaye, H.,(1988). .Causes and Characters of Drought in Ethiopia. *Ethiopian Journal of Agricultural Sciences*, **Vol.10**, 85-97pp.
- Tabony, R. C., (1979). Spectral and filter analysis of long-period rainfall records in England and Wales. *Meteor. Magazine*, 108: 97–118.
- Tsakiris, G.,(1998). Stochastic Modeling of Rainfall Occurrences in Continuous Time; *Hydrological Science Journal*, **Vol. 21**, Athens, Greece.
- Tsay,.S.R,(2005). Analysis of Financial Time Series; 2<sup>nd</sup> ed. John Wiley & Sons, Inc.,Hoboken.
- Tsegay, W.,(1998). El Niño and Drought Early Warning in Ethiopia. In Using Science Against Famine: Food Security Early Warning System and El Niño, *Internet Journal of African Studies*.
- Wei, W.W.S., (1990). Time Series Analysis. Addison-Wesley Publishing Company, Inc., 478 pp., New York-USA.

- Winstanley, D.,(1973a). Recent rainfall trends in Africa, the Middle East and India. *Nature*. 243: 464–465pp.
- Winstanley, D.,(1973b). Rainfall Patterns and General Atmospheric circulation. *Nature*. 245: 190–194pp.
- Wing, H. C., Gabriel ,B. and Ashbindu ,S.,(2008). Trends and Spatial Distribution of Annual and Seasonal Rainfall in Ethiopia; *International Journal of Climatology*; Published online in Wiley Inter Science [www.interscience.wiley.com](http://www.interscience.wiley.com).
- Wolde-Georgis, T.,(1997). The Use of El Nino Information as Drought Early Warning in Ethiopia. *Internet Journal of African Studies*, **Vol.1**, No.2. Case studies.
- Wood A.,(1977). A preliminary chronology of Ethiopian droughts. In Drought in Africa,vol. 2 (pp.68-73), Dalby D, Church RJH, Bezzaz F (eds.). *International African Institute*, London.
- Workineh ,D., (1987). Some aspects of meteorological drought in Ethiopia., Drought and Hunger in Africa: denying famine a future, M.H. Glantz (ed.). Cambridge University Press.
- Yevjevich, V. (1972). Stochastic Process in Hydrology; Water Resources Publications, Collin.
- Yilma, S., Dameree, G. R.and Delleur, J.W.,(1994).Sunspot number as a possible indicator of annual rainfall at Addis Ababa, Ethiopia. *International Journal of climatology*,**Vol.14**, 911-923pp.
- Young, PC, Jakeman, AJ, Post DA. ,1997. Recent Advances in Data-Based Modelling and Analysis of Hydrological Systems. *Water Science and Technology*, **Vol.36**: 99-116pp

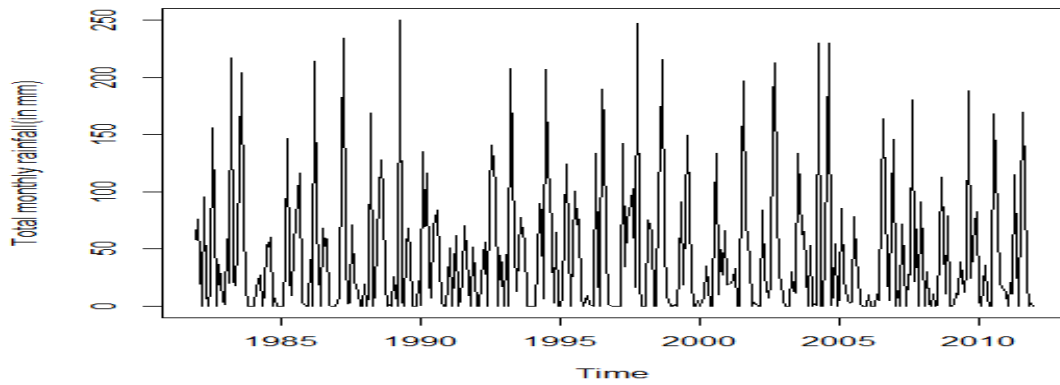
# Appendices

## List of Figure

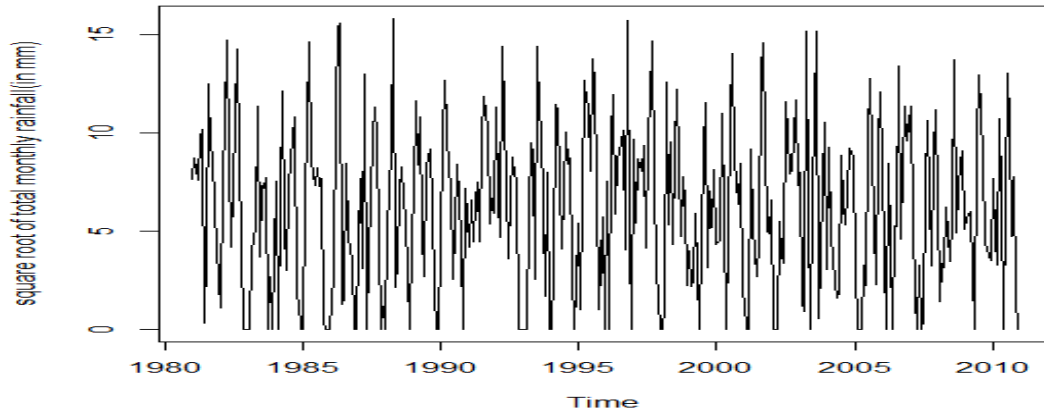
**Figure-1:** Plot of total monthly rainfall in mm



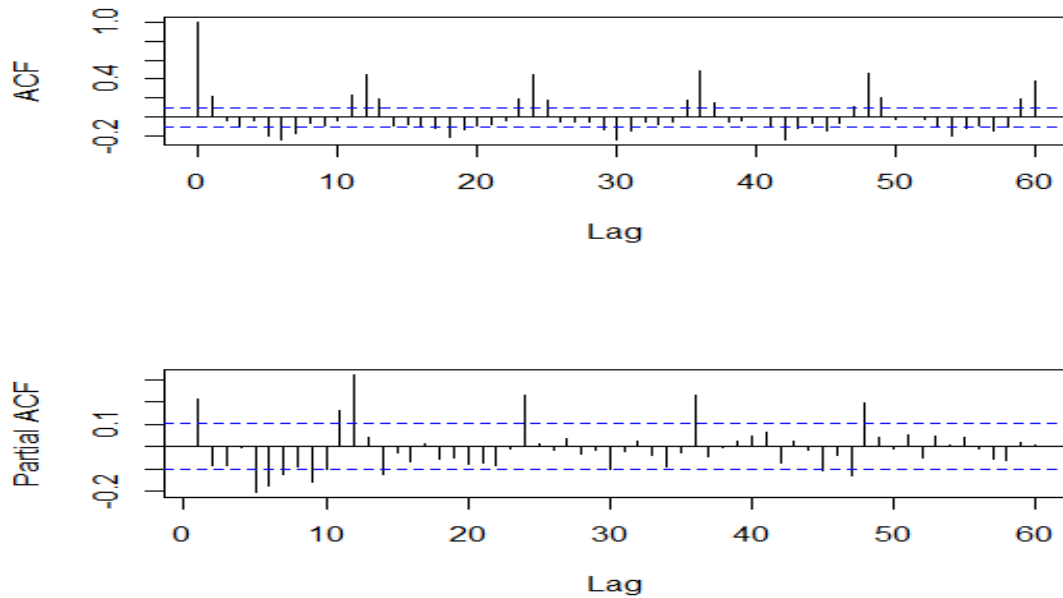
**Figure-2:** Time plot for total monthly rainfall in mm (Dire Dawa)



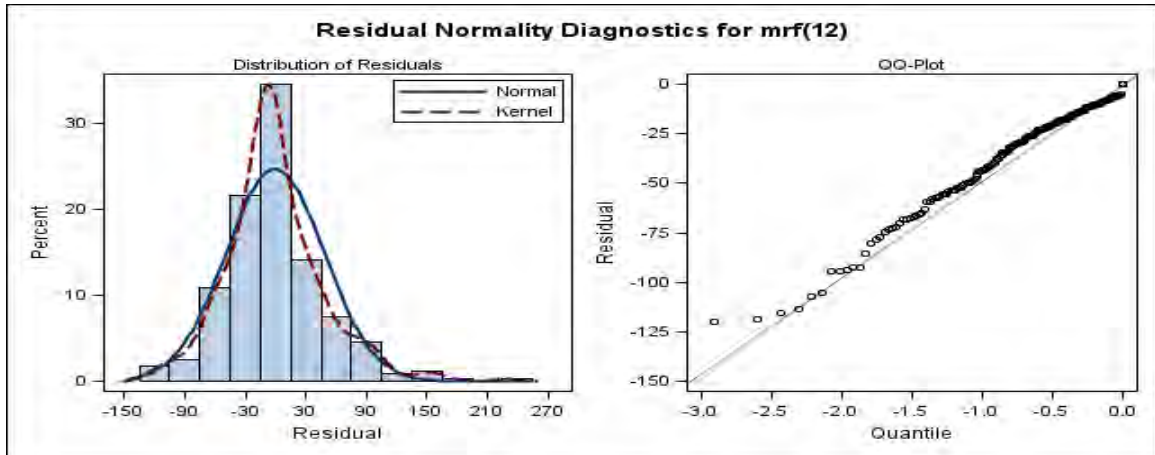
**Figure-3:** Time plot for total monthly rainfall of square root transformed rainfall series (Dire Dawa)



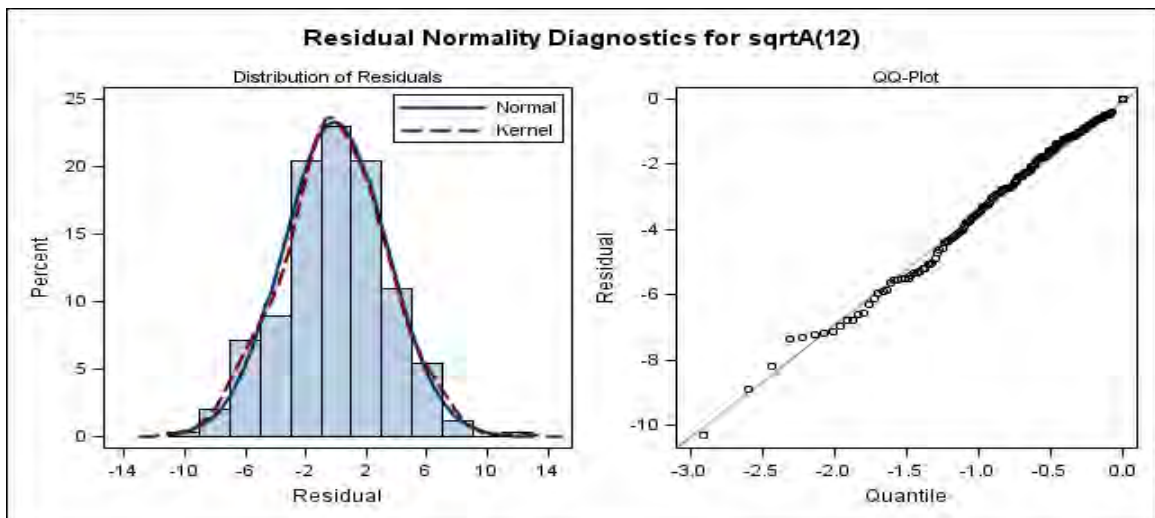
**Figure-4:** Autocorrelation and Partial Autocorrelation plots for the rainfall series.



**Figure-5:** Normality Diagnostics plot for: (a) untransformed series model, (b) square root transformed series

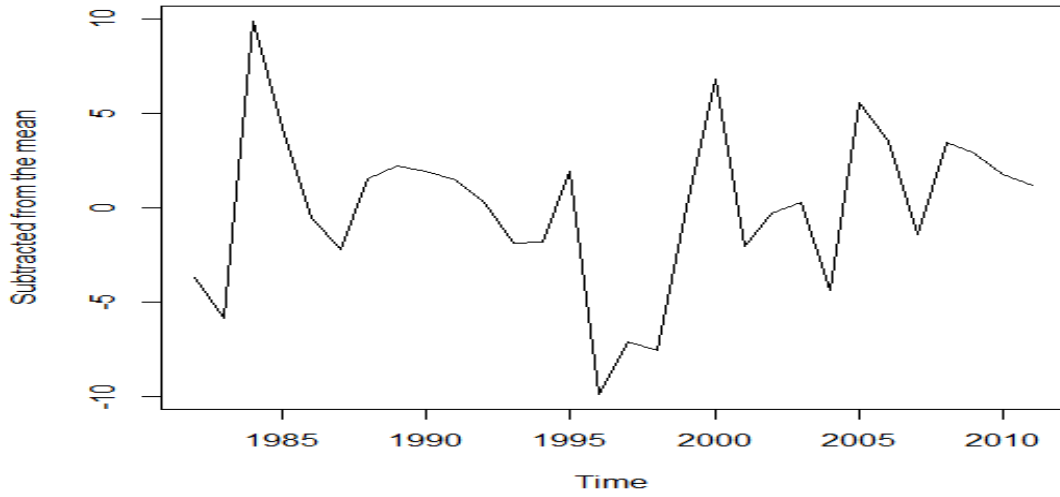


(a)

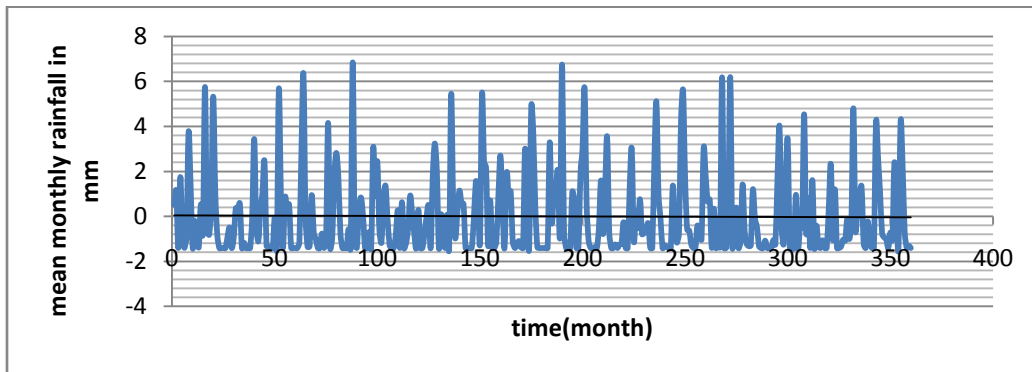


(b)

**Figure-6:** Rainfall series plot as departure from mean: (a) Annual and (b) Monthly

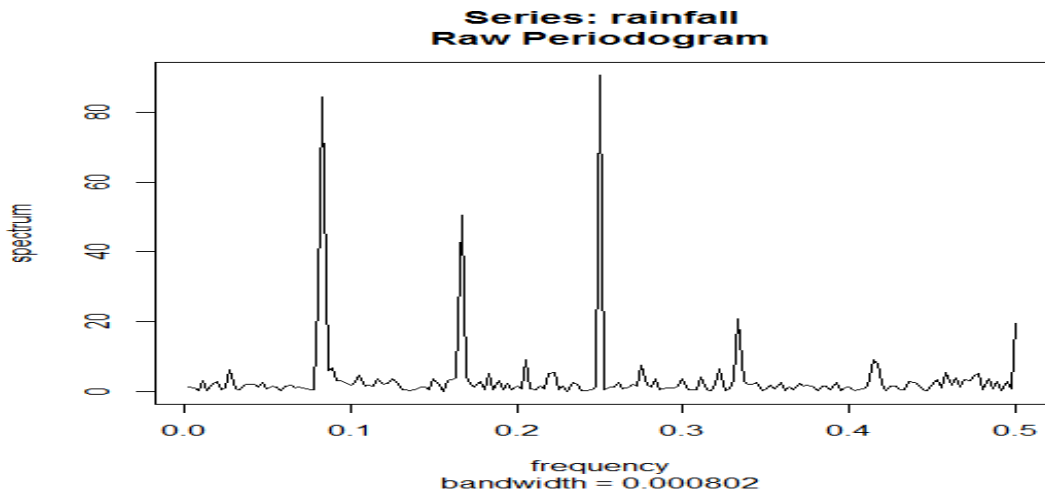


(a)

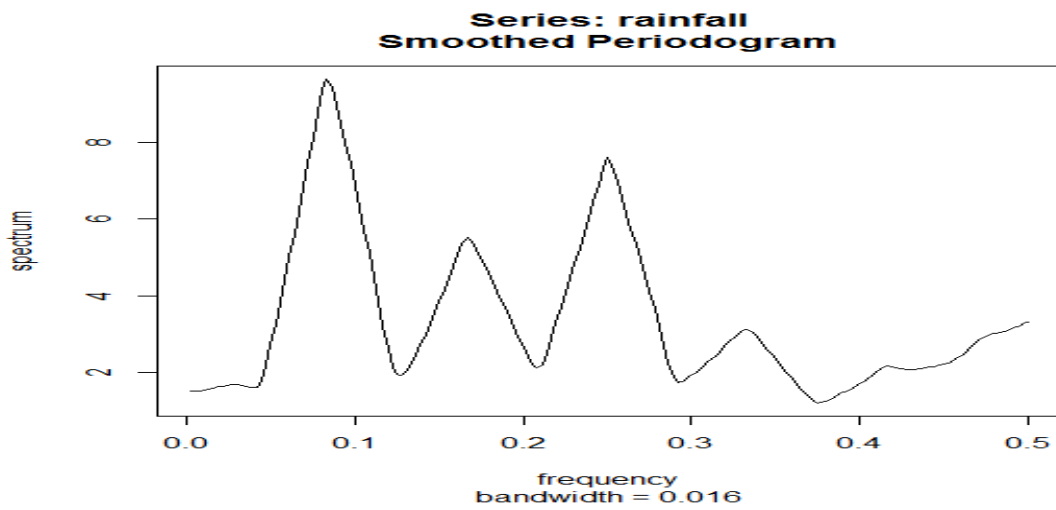


(b)

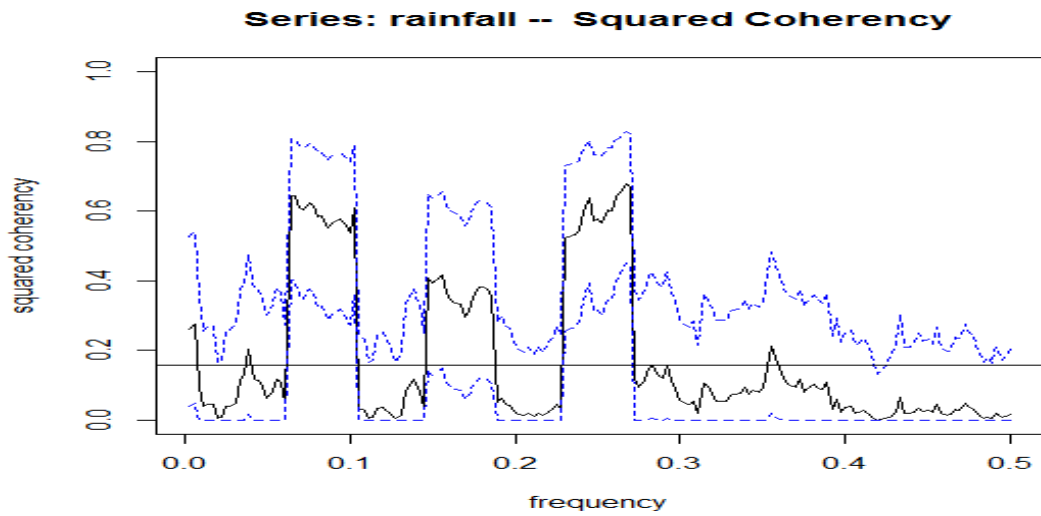
**Figure-7:** Raw periodogram of monthly rainfall of Dire Dawa



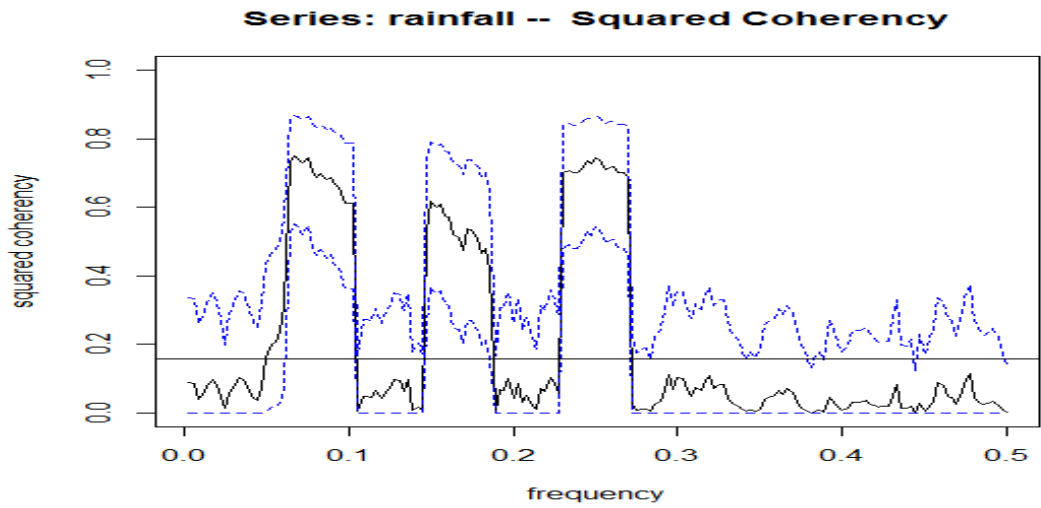
**Figure-8:** Smoothed periodogram estimate of spectrum of monthly rainfall of Dire Dawa.



**Figure-9:** Squared coherency between rainfall series of: (a) Dire Dawa and Haramaya station; (b) Dire Dawa and Dengego stations.



(a)



(b)

