

AFRICAN INSTITUTE FOR MATHEMATICAL SCIENCES
(AIMS RWANDA, KIGALI)

Name: Kamau Gladys Muthoni
Course: Research Methods in Climate Science

Assignment Number: 1
Date: February 4, 2021

Contents

1	Task 1: Cluster Analysis	2
1.1	Introduction	2
1.2	Methodology	3
1.3	Results and Interpretation	3
1.3.1	Climate Viariables	3
1.3.2	Years (1960 -2010)	5
1.4	Conclusion	6
2	Task 2: Principal Component Analysis (PCA)	7
2.1	Introduction	7
2.2	Methodology	7
2.3	Results and Interpretation	7
2.4	Conclusion	10
3	Task 3: Time-Series Analysis	11
3.1	Introducttion	11
3.2	Methodology	11
3.3	Results and Interpretation	11
3.3.1	Plot of PRED and TMPD time series	11
3.3.2	Detrended TMPD	12
3.3.3	Correlogram for PRED and DTMPD	12
3.3.4	Cross-correlation of PRED and DTMPD	13
3.4	Conclusion	14
4	Task 4: Spectral and Wavelet Analysis	15
4.1	Introduction	15
4.2	Methodology	15
4.3	Results and Interpretation	15
4.3.1	Spectral analysis	15
4.3.2	Wavelet analysis	16
4.4	Conclusion	18

1 Task 1: Cluster Analysis

1.1 Introduction

The data choosen for this study was extracted from Climate Dataset (CRU) over Turkana City in Kenya with coordinates 3.3122° N, 35.5658° E. The study is based on the month of April which was recorded with the highest amount of precipitation over the span of 51 years from 1960 to 2010. The aim of the study is to identify the characteristics of three Cluster Analysis algorithms. CA groups climate variables into groups of similar type. The importance of this study is to learn how different methods of hierarchical cluster analysis work [7].

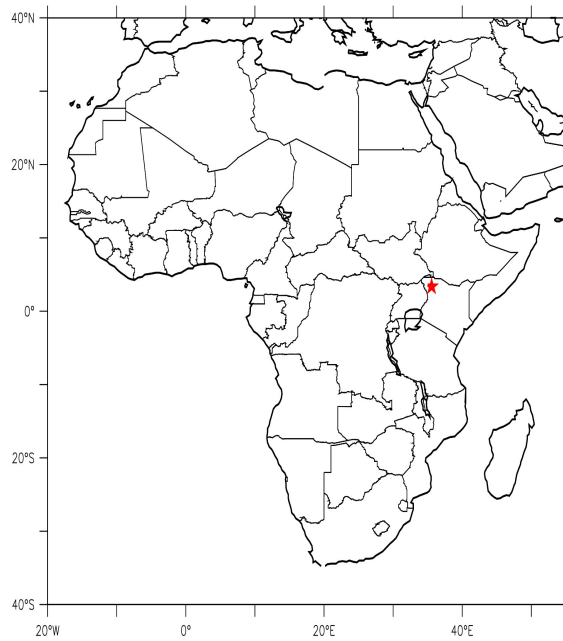


Figure 1: Map of Africa

There's no clear history of analysing climate data of Turkana using CA. This study therefore groups the various climate variables over Turkana city.

The extracted data had 51 rows and 10 columns with the following key variables:

- CLDD - Cloud Cover Days

- PRED - Precipitation Days
- WETD - Wet Days
- VAPD - Vapour Pressure
- TMND - Minimum Temperature
- TMPD - Mean Temperature
- TMXD - Maximum Temperature
- PETD - Potential Evapotranspiration
- DTRD - Diurnal Temperature Range

1.2 Methodology

After extracting the needed data using microsoft excel it was subjected to CA analysis methods. This study focused on 3 hierarchical CA methods; single linkage, average linkage and ward algorithm [7]. Vertical cluster tree was the main tool used to visualise the data. R was used as the programming software and language to up with the cluster trees. Clustering was done by grouping the data matrix to explore the grouping of climate variable and years (1960 -2010) over the city.

1.3 Results and Interpretation

The data was analysed using Hierarchical cluster method which comprises of single linkage, average linkage and ward algorithm. The study explores the grouping of climate variables and years.

1.3.1 Climate Viariables

i) Single Linkage

Starts with two closest data points and link. Next two closest data points and link until all data points were linked as shown in figure 2.

Cutting the cluster tree at height 9 creates two groups:

- CLDD, PRED, WETD - common feature is wetness.
- DTRD, PETD, VAPD, TMND, TMPD, TMXD - common feature is temparature.

PRED and WETD are more similar to each other compared to CLDD indicated by the short linkage distance. TMPD, TMXD and TMPD are also more similar to each other than VAPD, PETD and DTRD.

ii) Average Linkage

Similar to single linkage but based on the average linkage distance. Grouped the data as shown in figure 3.

Cutting the cluster tree at the height of almost 8 created 3 groups:

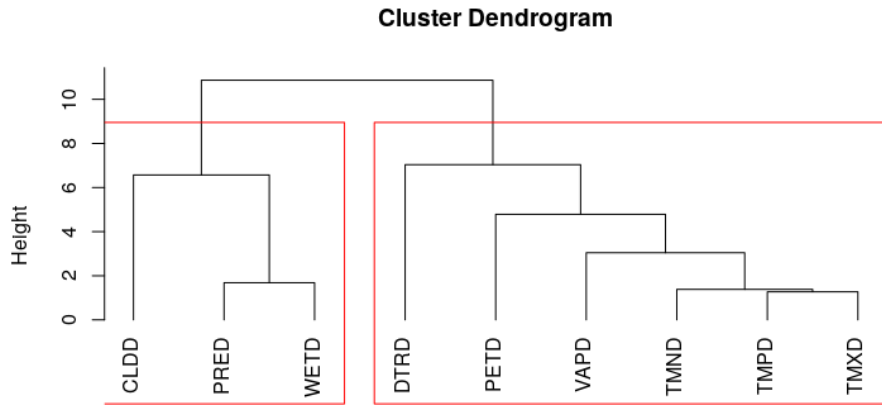


Figure 2: Single linkage

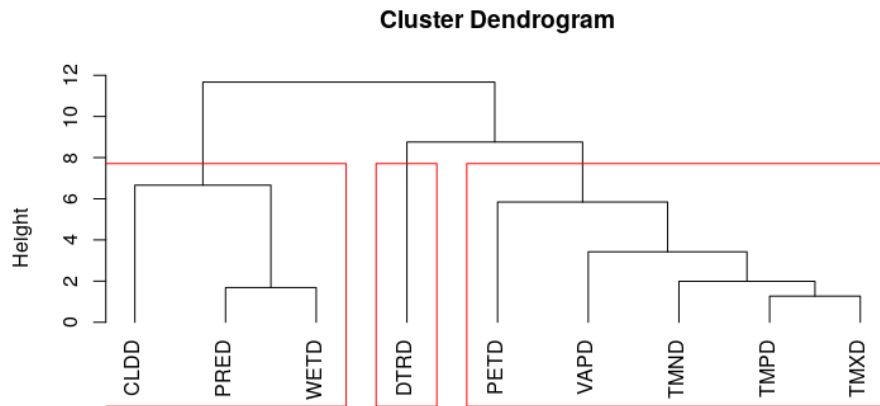


Figure 3: Average linkage

- CLDD, PRED, WETD - common feature is wetness.
- DTRD
- PETD, VAPD, TMND, TMPD, TMXD - common feature is temperature.

DTRD is an outlier. PRED and WETD are more similar to each other compared to CLDD which joins the cluster later. TMPD and TMXD are also more similar to each other more compared to PETD, VAPD and TMND.

iii) Ward Algorithm

Based on variance analysis. It seeks to minimize within group variance and maximize between group variance. Grouped the data as shown in figure 4.

Cutting the cluster tree at the height of almost 10 created 3 groups:

- CLDD, PRED, WETD - common feature is wetness.
- PETD, DTRD - potential evapotranspiration causes diurnal temperature change.
- VAPD, TMND, TMPD, TMXD - common feature is temperature.

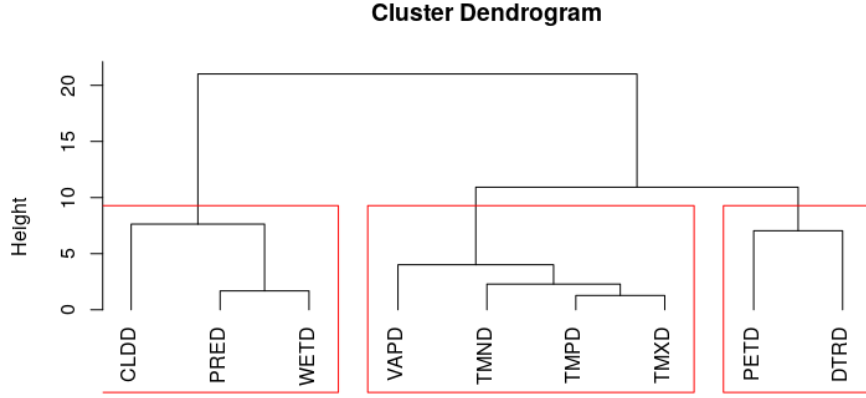


Figure 4: Ward algorithm

DTRD and PETD are of similar type. PRED and WETD are more similar to each other compared to CLDD which joins the cluster later. TMPD and TMXD are also more similar to each other more compared to VAPD and TMND.

1.3.2 Years (1960 -2010)

i) Single Linkage

Grouped the data points of years as shown in figure 5. Cutting the cluster tree at the

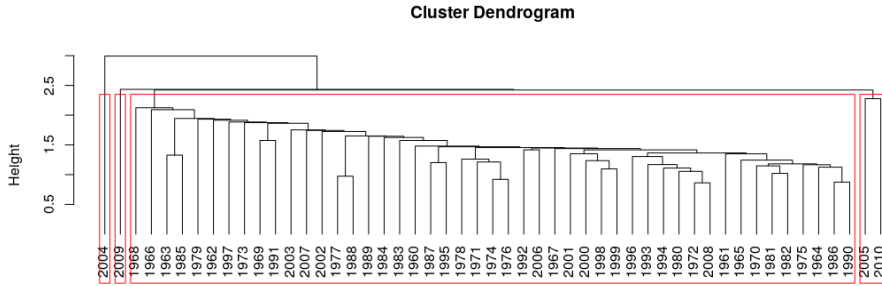


Figure 5: Single linkage

height of almost 2.5 groups the data points into 2 clusters with 2004 and 2009 as outliers. The largest group captures most of the data points indicating that all the years in this group share a common feature. It is also evident from the cluster tree that 1977 and 1988, 1974 and 1976, 1972 and 2008, 1986 and 1990 are more similar to each other since they are closest to each other.

ii) Average Linkage

Data points of years were grouped as shown in figure 6.

Cutting the cluster tree at the height of 3.5 groups the data points into 6 clusters with 2004 as an outliers. The 2 largest group captures most of the data points indicating that all the years in these groups share a common feature respectively.

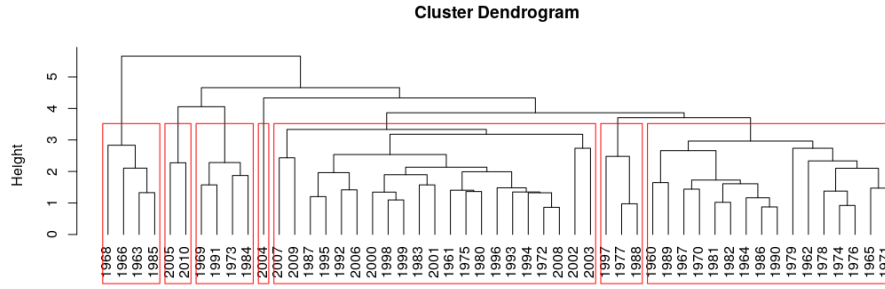


Figure 6: Average linkage

iii) Ward Algorithm

Grouped the data points of years as shown in figure 7. Cutting the cluster tree at the height of 7 groups the data points into 6 clusters.

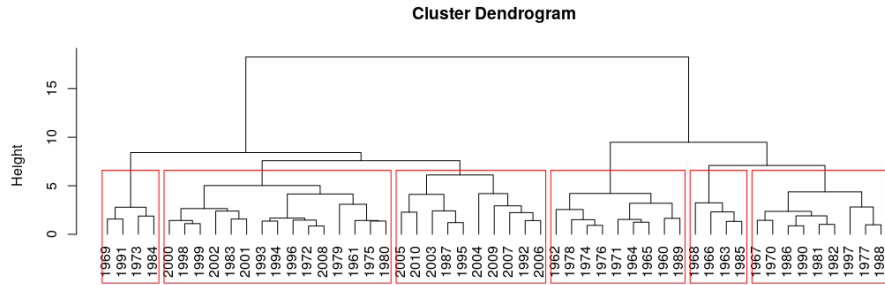


Figure 7: Ward algorithm

1.4 Conclusion

From the three hierarchical cluster analysis methods used in the 2 cases of variables its clear that single linkage and average linkages identify the outliers but the groups are not well defined. Ward algorithm clearly defines the groups but outliers are not well identified. It is also noted that cluster trees for climate variables and year variables are different. The years grouped together indicate that the climate feature was common to those years.

2 Task 2: Principal Component Analysis (PCA)

2.1 Introduction

PCA is a procedure that focuses on dimensionality reduction. It extracts fewer and independent underlying dimensions around which the data variance is organised. It also identifies the main processes that explain the largest percentage variance of the dataset [3]. The aim of the task is to study the characteristics of PCA and compare them with that of CA. The main importance of this task is to learn how PCA works. PCA has been used on various studies like climate variability and change on vulnerability and adaptation among Turkana pastoralists, establishing vegetable cover change over time and assessing the impact of climate change on food security of communities in Turkana. This study will try to get a deeper insight of Turkana climate data and explain the main processes that the city has been experiencing over a span of 51 years.

2.2 Methodology

The study used two methods of dimensionality reduction; rotated and non-rotated. These two methods had different outputs but rotated was more efficient in reducing the number of dimensions. It allows changing of the factor analysis to identify new patterns of the factor structure while unrotated PCA tries to illustrate the maximum variance value with minimal number of factors [3]. Rotated PCA helps to extract meaningful data that accurately represents the original dataset. The study looked at component loadings, variation and component scores to select the most important principal factors and better understand the dataset. The results were compared to ward algorithm to better understand PCA. R software was used for programming.

2.3 Results and Interpretation

i) Variation (Rotated and Unrotated PCA)

A total of 9 principle components were obtained as shown in table 1. The standard deviation measures variability across each principle component. Proportion of variance is the total variance percentage explained by each principle component in the original data set. PC1 explains 50% of the total variance in the original dataset. PC2 explains 23% of the total variance. In cumulative proportion PC1, PC2 and PC4 explains 85% of the total variance. As a result PC1, PC2 and PC4 were selected to explain the dataset thus dimensionality reduction.

Both unrotated and rotated differ in principle components that best describe the data set. Rotated method was preferred to unrotated as it reduces the dimensions to PC1, PC2 and PC4 which explain a variance of 44%, 23% and 18% respectively.

	Unrotated			Rotated		
	Std dev	Proportion of Variance	Cumulative Proportion	Std dev	Proportion of Variance	Cumulative Proportion
PC1	2.276	0.576	0.576	3.947	0.453	0.453
PC2	1.419	0.224	0.800	2.082	0.261	0.714
PC3	1.048	0.122	0.922	1.622	0.213	0.927
PC4	0.722	0.058	0.979	1.152	0.053	0.980
PC5	0.386	0.017	0.996	0.156	0.012	0.992
PC6	0.163	0.005	0.997	0.027	0.005	0.997
PC7	0.093	0.001	1	0.014	0.003	1.000
PC8	0.025	0	1	0.001	0	1.000
PC8	0.025	0	1	0.001	0	1.000
PC9	2.635	0	1	0	0	1.00

Table 1: PCA Summary

Selection of the most important principle components can also be selected using the scree plot show in 8 with a default setting that the varinace should be atleast 1.

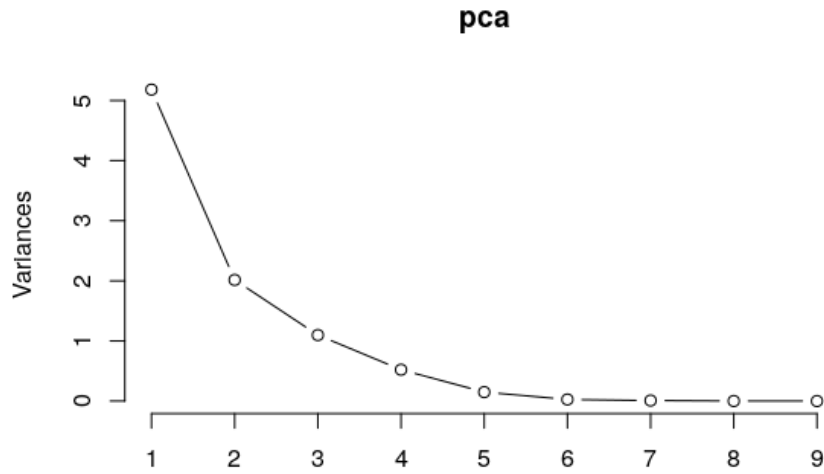


Figure 8: Scree plot

ii) Component Loadings (Rotated)

These are the correlation coefficients of the linear combination of the original variables from which principle components are constructed. Rotated PCA determines the loadings cutoff for the importnat variables that contribute highest to the factors as shown in red in table 2.

TMPD, TMND, TMXD and VAPD all increase in PC1. Due to the high temperatures in the dataset, this process is drought. In PC2 both PRED and WETD decreases. In drought seaoons there's no rain and the land is dry. In PC3 as CLDD decreases,DTRD increases. When there's no cloud during the night the temperatures tend to decrease and increase during the day. This results to increase in dirunal range temperatures.

Variables	PC1	PC2	PC3
PRED	-0.162010	-0.953727	-0.171107
PETD	0.546160	0.396223	0.678831
CLDD	-0.141978	-0.508158	-0.714389
DTRD	0.011618	0.029256	0.893948
TMPD	0.965685	0.142101	0.194406
TMND	0.983804	0.137589	0.028010
TMXD	0.917354	0.138291	0.349367
VAPD	0.935681	0.154074	-0.057063
WETD	-0.167719	-0.952899	-0.155570

Table 2: Principle component loadings

iii) Component Scores (Rotated)

Figure 9 figuratively describes percentage variance of each principle component on the original data and the component scores. The score shows how each process varies in different years. PC1 most active in 2010, most inactive in 1968 but was dormant in 1975. PC2 was most active in 1966, most inactive in 1969 but was dormant in 1960,1961 and 2002. PC4 most active in 2005, most inactive in 1963 with no dormance in any year.

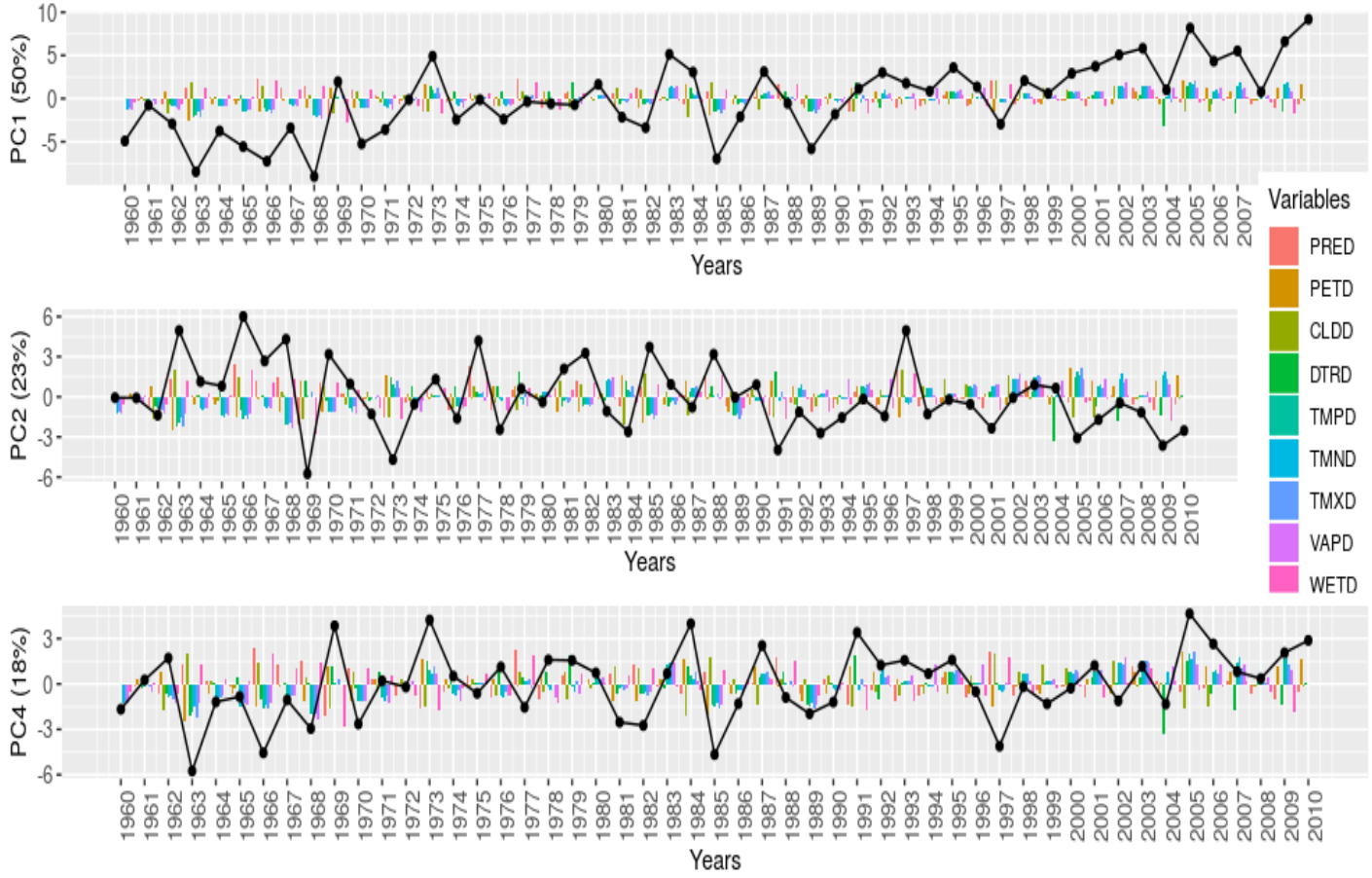


Figure 9: Component scores plot

2.4 Conclusion

PCA analysis identifies the processes that cause the climate variables to vary. PCA gives a deeper understanding of the dataset as it also reduces dimensionality which can be helpful for better climate predictions and data analysis. Both loadings and scores play an important role in analysing data using PCA. Ward algorithm defines the groups well but you can barely deduce how they are grouped and the processes involved.

3 Task 3: Time-Series Analysis

3.1 Introduction

Time series analysis is a sequence of data(continuous) that follow non-random orders. The data must be equally spaced time intervals. To deploy time series analysis the dataset must satisfy these two conditions. Time series analysis comprises of various methods that analyzes time series data for the purpose of extracting meaningful climate statistics and other features of the data [2]. The aim of this task was to explore the time series patterns in precipitation (PRED) and temperature (TMPD) data and the relationships between the two variables. This is important for the purpose of learning and analysing the patterns of temperature and precipitation in Turkana city over a span of 51years from 1960-2010. As a result it helps to study the past and predict the future.

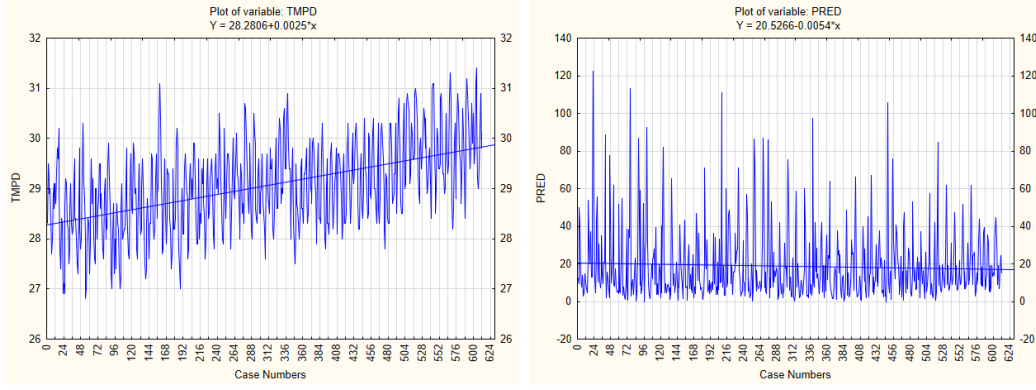
3.2 Methodology

There are various methods of time series analysis which can be classified into frequency domain (Fourier and Wavelet analysis) and time domain (autocorrelation and cross-correlation) [2]. This task focused on time domain and time series to identify if there was trend in the data. Dataset was re-extracted again using excel so as to include the month variable since the data was indexed monthly from January 1960 to December 2010. Using statistica TMPD and PRED values were subjected to time series to identify the patterns. At some point it was required to detrend the temperature variable for better analysis of the data. Detrending of the data was done using microsoft Excel.

3.3 Results and Interpretation

3.3.1 Plot of PRED and TMPD time series

It was evident from figure 10a that temperature variable has deterministic and global trend; change in mean over time(climate change). Line of best fit showed the temperature values were increasing with 0.0025°C every month from 1960 to 2010. Precipitation variable presented a very small trend as shown in figure 10b which was neglected because it was insignificant. This implied the slope of the best line of fit was insignificant but we can see the pattern of highest amount of rainfall recorded at 120mm/m was recieved in 1960 and this has reduced over the years. The minimum rainfall recorded over the span of 51years was 0mm/m. The variability in both plots indicate presence of seasonality.



(a) Trended TMPD

(b) TimeSeries PRED

Figure 10: Time series plots

3.3.2 Detrended TMPD

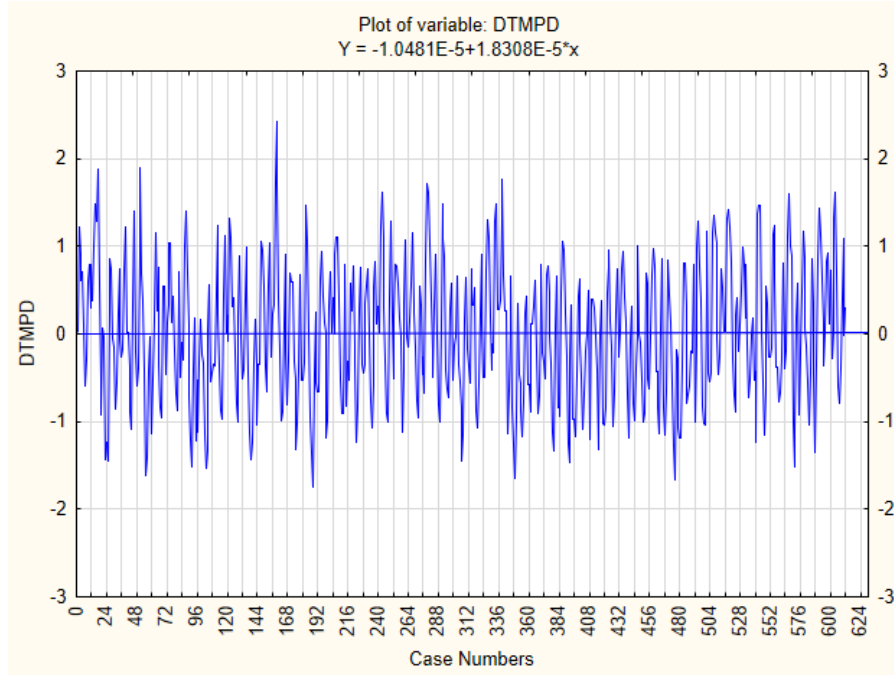


Figure 11: Detrended TMPD

There are many methods of detrending data but linear regression was used in this case. Detrending was done to remove distortion in TMPD variable to allow better analysis of its variability. Plot of the detrended data gave a clearer visual of the cyclic patterns that is increase and decrease of temperature with time compared to trended data.

3.3.3 Correlogram for PRED and DTMPD

Correlogram is synonym for autocorrelation. It is the correlation between time series data and lagged version of itself. It is used to find repeating patterns [2]. As shown in figure 12a DTMPD autocorrelation shows a significant sinusoidal cycle. With a lag of 1 month, it shows that temperature data is correlated with the past data and it varies with time. It also shows a

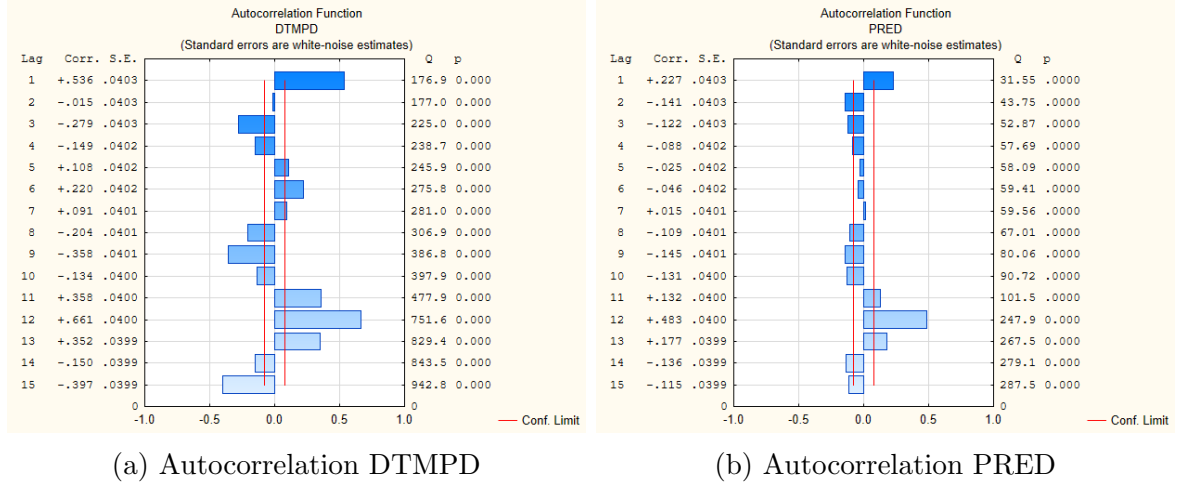


Figure 12: Autocorrelation

repetitive annual pattern at lag 12, sub-seasonal pattern at lag 6 and seasonal pattern at lag 3. The red line indicates the approximate 95% confidence interval.

Figure 12b shows a weak sinusoidal autocorrelation of precipitation. At lag 1 correlation is at approximately 0.25 and 0.5 at lag 12. Although the autocorrelation is not strong, seasonality is evident.

3.3.4 Cross-correlation of PRED and DTMPD

Cross-correlation describes the degree of correlation between two different time series. It is useful in determining whether changes in one time series has effect on the other time series [5]. Figure 13 shows temperature is significantly correlated to precipitation. At lag -13 and 11 the correlation is 0.5. The sinusidal cycle shows as temperature increases precipitation decreases and vice versa.

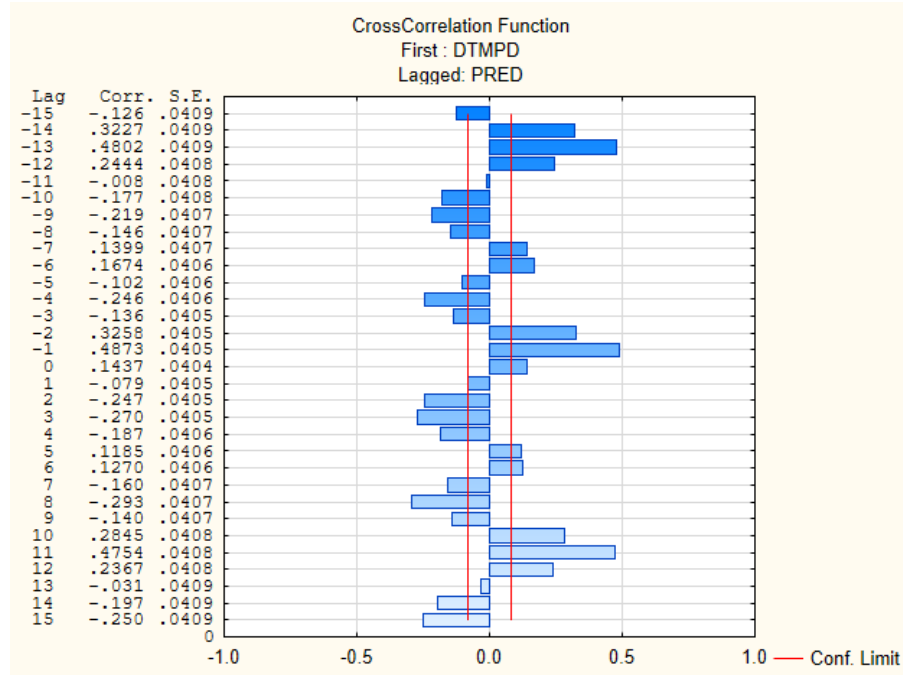


Figure 13

3.4 Conclusion

Autocorrelation, cross-correlation and time series played a vital role in identifying the underlying patterns of how temperature and precipitation varied with time. To get a deeper and better insight of the data more analysis needs to be done. Both autocorrelation and cross-correlation are based on correlation.

4 Task 4: Spectral and Wavelet Analysis

4.1 Introduction

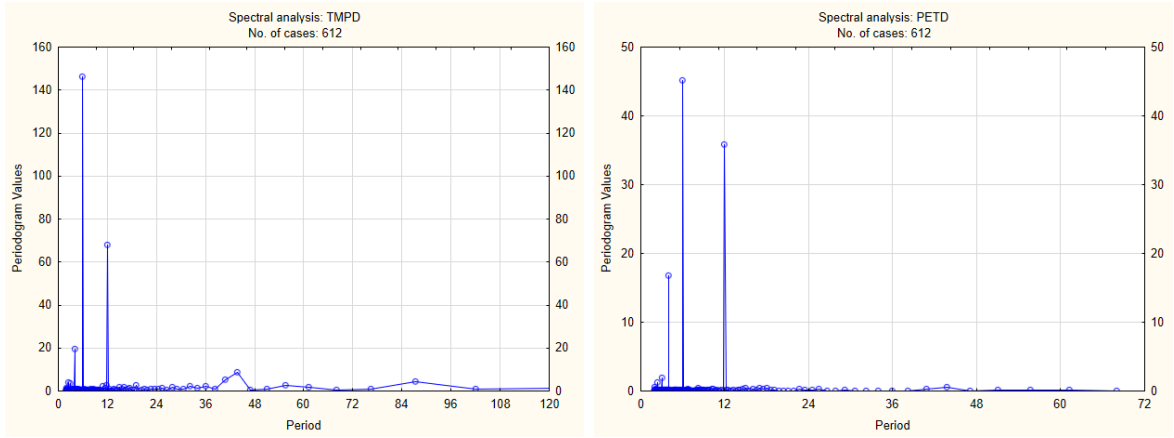
Spectral analysis is based on variance and it is useful in identifying seasonal fluctuations of different lengths [4]. From the cyclic pattern identified in fourier analysis, wavelet analysis explains the behaviour and change of the cycles and amplitude with time. The aim of this task is to explore and compare spectral analysis and wavelet analysis using the precipitation data. The study was important for learning and analysing how the cycles of precipitation varied within the span of 51 years from 1960-2010 in Turakana City. This task will also be an extension of the analysis done in task 3.

4.2 Methodology

Using statistica the data was subjected to fourier analysis to identify the seasonal fluctuations. R software was also used to come up with the wavelet images to better explain how the seasonality varies with time.

4.3 Results and Interpretation

4.3.1 Spectral analysis



(a) Spectral analysis TMPD

(b) Spectral analysis PETD

Figure 14: Spectral analysis

There are 3 cycles in temperature and precipitation spectral analysis. Figure 14a shows sub-seasonal cycle is the most dominating with a periodgram of approximately 150. It is followed by annual cycle then the quarterly cycle. This implies that there are 3 main seasons in the area; after 3, 6 and 12 months. In the temperature spectral plot, there are other seasons that repeat after 4 and 8 years. They are called ENSO and solar cycle seasons respectively.

Figure 14a shows sub-seasonal cycle is the most dominating with a periodgram of approximately 45. It is followed by annual cycle then quarterly cycle. There are only 3 precipitation seasons in the area. They repeat after 3, 6 and 12 months.

4.3.2 Wavelet analysis

Wavelet analysis explains the behaviour and change of the cycle and amplitude with time. There are different methods to use when performing wavelet analysis but the best is morlet as it is more sensitive to the time series data compared to the others [1].

Figure 15a shows a band at 0.50(6 months) of the dorminant season with high temperatures which is persistent across all the years. At period 1 high temperature season was also recorded between 1970 and 1980 and later after around 25years between 2000 and 2010. Generally the band at period 1 shows seasonal cycle that comes after the subseasonal cycle and is persistent across all the years except in 1990. Between the period of 2 and 4 the ENSO season that happens after 4 years is also visible [6].

Figure 15b shows a cycle between the period of 0.25 and 0.50 that happened between the years 1960-1970 and 1980-1990 and later decreases with time. At period 1 seasonal cycle is seen between the year 1965 and 1990. This shows the 3 seasons as identified by spectral analysis.

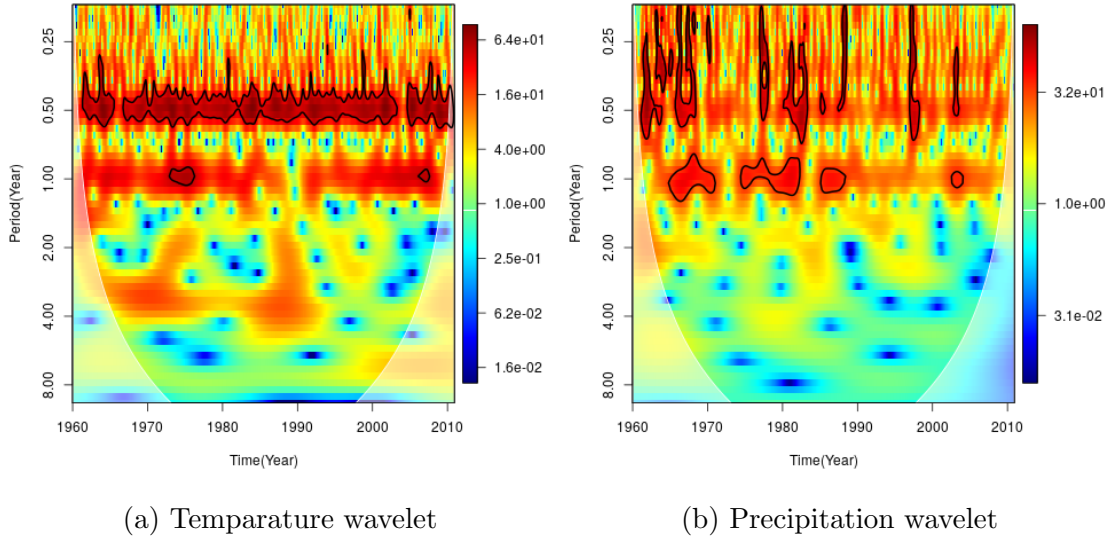


Figure 15: Wavelet analysis

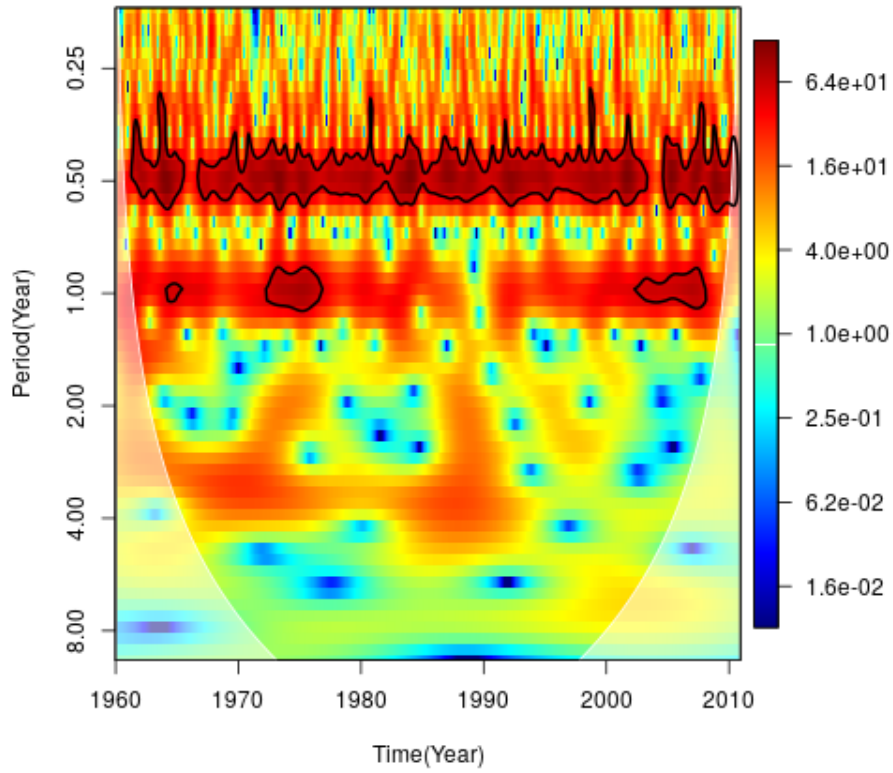


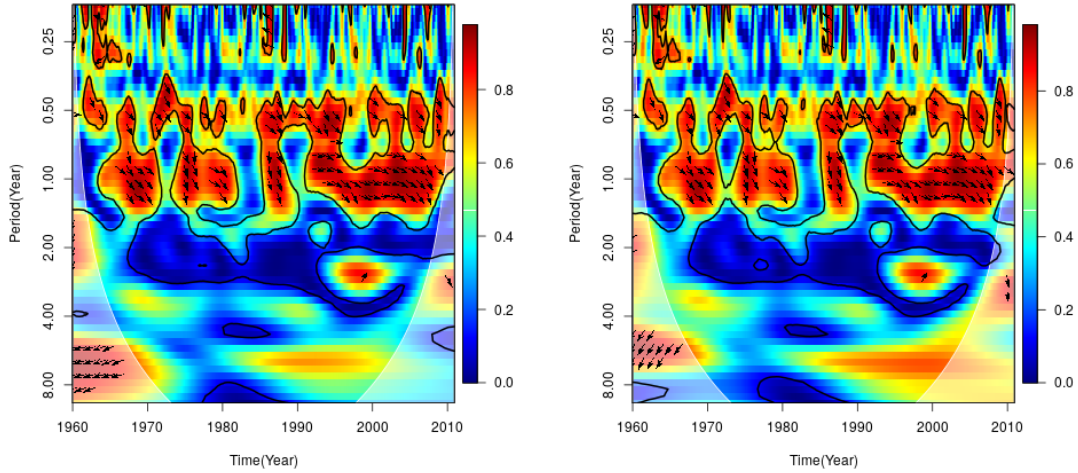
Figure 16: Detrended TMPD wavelet

Figure 16 is similar to figure 15a but here the plot is of the detrended temperature variable. It is more clear without trend but displays 3 seasonal cycles as discussed.

Wavelet coherence considers two time series data and tries to explain when they couple and decouple.

Figure 17a shows that temperature and precipitation couple at the band of period 0.50 and 1. They decouple at the band of 1 between year 1980 and 1990. At period strong coupling is evident between year 1960 and 1970. There is also a faint coupling the year between 1990 and 2000 falling in the period of 2-4 and 4-8. The arrow facing down indicates that precipitation leads temperature. This implies as temperature decreases precipitation increases.

Figure 17b is similar to figure 17a but coupling is done between detrended temperature variable and precipitation. A study of Paul and DOG methods of doing wavelet analysis was also done.



(a) Wavelet coherence

(b) Detrended wavelet coherence

They are not sensitive to the time series data set because important information and patterns were omitted. This made morlet the best and most suitable for most time series wavelet analysis.

4.4 Conclusion

Wavelet analysis is more powerful than spectral analysis at it explains the variation with time. It is clear that Turkana City has 3 main repeating seasons and the climate is changing over the area due to global warming. Temperature recorded continuously increases with time. Due to the high temperatures and low precipitation over the span of 51years, Turkana is a dry city.

References

- [1] Ingrid Daubechies. The wavelet transform, time-frequency localization and signal analysis. *IEEE transactions on information theory*, 36(5):961–1005, 1990.
- [2] James Douglas Hamilton. *Time series analysis*, volume 2. Princeton university press Princeton, NJ, 1994.
- [3] Ian T Jolliffe. Principal component analysis: a beginner’s guide—i. introduction and application. *Weather*, 45(10):375–382, 1990.
- [4] Lambert H Koopmans. *The spectral analysis of time series*. Elsevier, 1995.
- [5] Boris Podobnik and H Eugene Stanley. Detrended cross-correlation analysis: a new method for analyzing two nonstationary time series. *Physical review letters*, 100(8):084102, 2008.
- [6] Christopher Torrence and Gilbert P Compo. A practical guide to wavelet analysis. *Bulletin of the American Meteorological society*, 79(1):61–78, 1998.
- [7] Jeroen K Vermunt and Jay Magidson. Latent class cluster analysis. *Applied latent class analysis*, 11:89–106, 2002.