

## CERTIFICATION

This is to certify that **ALADELE THEOPHILUS EXCELLENT** with matriculation number **16CD004588** has carried out this research project in partial fulfillment of the requirement for the award of Bachelor of Science Degree in Computer Science in the College of Pure and Applied Sciences, Landmark University.

---

**Dr. Marion O. Adebisi**

(Project Supervisor)

---

**Date**

---

**Dr. Oladayo G. Atanda**

(Co - Supervisor)

---

**Date**

---

**Prof. Oluwakemi C. Abikoye**

(External Examiner)

---

**Date**

---

**Dr. Marion O. Adebisi**

(Head of Department)

---

**Date**

## **DEDICATION**

I dedicate this project to the Almighty God, for his guidance, wisdom, grace, strength, and patience granted to me during this period of this project. I also dedicate this project to my family, most especially my parents, Mr and Mrs. Aladele for their constant love and support for me.

## **ACKNOWLEDGEMENT**

I acknowledge God Almighty for His eternal love, strength, protection, provision, and guidance up till this very moment. I really appreciate my supervisors, most especially Dr. Marion O. Adebisi and Dr Oladayo Atanda, for their contribution to the success of this project. I also want to thank Dr Achebe also for her input and assistance also. I appreciate the department of Computer Science, Landmark University, for its thriving educational environment and my project members' assistance. Finally, I appreciate my parents, Mr and Mrs Aladele for their endless support in every area and their constant encouragement and funding of my educational studies. I am forever grateful.

## **TABLE OF CONTENTS**

<b>TABLE OF CONTENT</b>	<b>PAGE</b>
DECLARATION	ii
CERTIFICATION	iii
DEDICATION	iv
ACKNOWLEDGEMENT	v
LIST OF FIGURES	viii
LIST OF TABLES	ix
ABSTRACT	x
<b>CHAPTER ONE</b>	<b>1</b>
<b>INTRODUCTION</b>	<b>1</b>
1.1    Background of the Study	1
1.2    Statement of the Problem	2
1.3    Aim and Objectives of the Study	3
1.4    Significance of the Study	4
1.5    Scope of the Study	4
1.6    Limitations of the Study	5
1.7    Justification of the Study	5
1.8    Arrangement of Chapters	6
<b>CHAPTER TWO</b>	<b>7</b>
<b>LITERATURE REVIEW</b>	<b>7</b>
2.1    Overview of Principal Component Analysis (PCA)	7
2.2    Historical review of Principal Component Analysis (PCA)	8

2.3	The Rise of High-dimensional Biological Datasets	9
2.4	Dimensionality Reduction in Biological Data Analysis	10
2.5	Review of Related Works	11
<b>CHAPTER THREE</b>		<b>17</b>
<b>METHODOLOGY</b>		<b>17</b>
3.1	System Overview	17
3.2	Dataset Selection and Description	18
3.3	System Flowchart and Use Case Diagram	20
3.4	Data Preprocessing	23
<b>CHAPTER FOUR</b>		<b>26</b>
<b>SYSTEM IMPLEMENTATION AND DOCUMENTATION</b>		<b>26</b>
4.1	Development Tools	26
4.2	System Requirements	26
4.3	Comparison and Analyzing of result.	30
4.4	Discussion	32
<b>CHAPTER FIVE</b>		<b>35</b>
<b>SUMMARY, CONCLUSION AND RECOMMENDATION</b>		<b>35</b>
5.1	Summary	35
5.2	Conclusion	35
5.3	Recommendations	36
<b>REFERENCES</b>		<b>37</b>
<b>APPENDIX</b>		<b>41</b>

## LIST OF FIGURES

Figure. 3.1 Flowchart of PCA Implementation .....	22
Figure. 3.2 Use case diagram of PCA Implementation .....	23
Figure. 4.1 A biplot Diagram.....	29
Figure. 4.2 A Heat Map Diagram .....	30
Figure. 4.3 Scree Plot Diagram.....	31
Figure. 4.4 Variables – PCA Diagram .....	32

## **LIST OF TABLES**

Table 2.1 Review of related works .....	13
Table 4. 1 Tabular view of Principal component.....	28

## ABSTRACT

This research delves into the pivotal realm of Principal Component Analysis (PCA) in the context of biological data analysis. As high-dimensional datasets continue to inundate genomics and proteomics, the need for robust dimensionality reduction techniques becomes increasingly apparent. The exponential growth of biological data has led to a dire need for tools that can distill the signal from the noise. Traditional approaches struggle to handle the complexity, and the challenge lies in identifying which variables are crucial for understanding the biological context.

The aim of this study is to harness the power of PCA in reducing dimensionality while retaining essential biological information, thus providing a vital tool for researchers navigating the sea of biological data. Employing R/RStudio and dimension reduction packages, we conducted PCA on a vast dataset. Data preprocessing and normalization ensured robust analysis. Our analysis extends to the biological interpretation of the principal components, drawing connections between gene expression patterns and underlying biological phenomena.

The PCA results unveil the intricate relationships within the data, highlighting significant genes, pathways, or clusters of samples that have remained latent in the high-dimensional space. Notable principal components offer a glimpse into the biological significance of the dataset. This research illuminates the potential of PCA as a dimensionality reduction tool for biological data. By unveiling hidden patterns, it offers insights that can steer future experiments, direct research hypotheses, and guide further exploration. In a data-rich world, PCA emerges as a beacon of clarity, promising to enhance our understanding of the biological universe.



# CHAPTER ONE

## INTRODUCTION

### 1.1 Background of the Study

The exponential growth of data across a range of industries, from biology to banking, has highlighted the significance of effective data analysis techniques. The difficulty of interpreting and displaying datasets increases as they become more complicated and include more variables or features (Smith, 2018). The "curse of dimensionality," a phenomenon where data becomes sparse and conventional analytic methods become useless or even deceptive, might result from this expansion in data dimensionality, which also raises computational costs (Jones et al., 2019). Dimensionality reduction stands out as an essential remedy in this situation. Among the methods used for this, Principal Component Analysis (PCA) stands out as one of the most popular. By using a statistical technique called principal component analysis (PCA), data can be represented in a reduced-dimensional space without suffering a major loss of information by pinpointing the dataset's axes (or "principal components") that optimize variance (Wang & Zhang, 2020). In addition to assisting with data visualization, this procedure also prepares data for use in other machine learning activities, potentially enhancing model accuracy and effectiveness.

The R programming language, celebrated for its robust statistical and data analysis capabilities, provides comprehensive support for PCA. Moreover, with the advent of the dimension reduction packages, R's prowess extends further into the realm of biological data analysis. Given the expanding datasets in fields like genomics and the concomitant need for efficient data interpretation tools, the relevance of PCA, especially when combined with

platforms like R and RStudio, has never been higher. This study delves into the application of PCA for dimension reduction using the R and RStudio, a conjunction that promises effective data interpretation and valuable insights (Liu & Smith, 2021).

The onset of the digital age has brought with it an unprecedented influx of data. This 'data deluge' presents both opportunities and challenges. While larger datasets can offer richer insights, they also become harder to manage and interpret, especially when they span across multiple dimensions. This complexity is further amplified in fields like genomics, proteomics, and other high-throughput biological studies, where thousands of variables might be recorded for a single sample (Shah & Murthi, 2021).

Dimensionality reduction is not just a matter of convenience or computational efficiency; it's a necessity for extracting meaningful patterns from such high-dimensional data. PCA, being a linear method, offers an advantage by projecting the original data onto a lower-dimensional space, preserving as much variability as possible. This reduction facilitates clearer visualizations, more efficient storage, and improved computational performances for subsequent analysis, such as clustering or classification.

## **1.2 Statement of the Problem**

In the contemporary landscape of research, data is both an asset and a challenge. The quantity and complexity of data generated, especially in fields like genomics and bioinformatics, have grown exponentially (Ahmad et al., 2022). This vast amount of high-dimensional data poses significant hurdles in its management, interpretation, and visualization. Traditional data analysis methods often struggle to cope with this magnitude and complexity, leading to inefficiencies, increased computational costs, and at times, misleading conclusions.

Principal Component Analysis (PCA) has emerged as a powerful tool for dimensionality reduction, aiding in both visualization and data preprocessing for further analyses. (Hajibabaei et al., 2023). However, while the theoretical underpinnings of PCA are well-established, its application to specific types of complex biological datasets remains fraught with challenges. For instance, datasets in genomics often carry unique structures, correlations, and inherent noise that standard PCA implementations might not handle optimally (Malik et al., 2020).

Moreover, while the R programming language is a go-to tool for many researchers, its potential in handling PCA for intricate biological datasets, especially with the aid of the dimension reduction packages, remains underexplored. (Sepulveda et al., 2020). There is a gap in comprehensive, user-friendly guides and methodologies that meld the strengths of R and RStudio, and PCA, tailored for high-throughput biological data.

This study aims to address the aforementioned challenges by developing and elucidating a robust methodology for conducting PCA on biological datasets using the R/RStudio dimension packages. The objective is to bridge the current knowledge gap, providing researchers with a streamlined approach to harness the power of PCA in R, optimized for the unique intricacies of biological data.

### **1.3 Aim and Objectives of the Study**

The aim of this study is to use PCA algorithm in R/RStudio and use dimension Packages for dimension reduction of a dataset.

Objectives:

1. Design a step-by-step procedure for PCA in R using the dimension reduction packages.
2. To implement the designed system.