

**A MULTI-OMICS CLASSIFIER FOR PREDICTION OF ANDROGEN  
DEPRIVATION TREATMENT RESPONSE IN PROSTATE CANCER DATASET**

**ALAGBE, EMMANUEL OLUWATOBA**

**(20PBF02172)**

**B. Tech (Hons) Microbiology, Ladoke Akintola University of Technology, Ogbomosho.**

**JANUARY, 2023**

**A MULTI-OMICS CLASSIFIER FOR PREDICTION OF ANDROGEN  
DEPRIVATION TREATMENT RESPONSE IN PROSTATE CANCER DATASET**

**BY**

**ALAGBE, EMMANUEL OLUWATOBA**

**(20PBF02172)**

**B. Tech (Hons) Microbiology, Ladoke Akintola University of Technology, Ogbomosho.**

**A DISSERTATION SUBMITTED TO THE SCHOOL OF POSTGRADUATE  
STUDIES IN PARTIAL FULFILMENT OF THE REQUIREMENT FOR THE  
AWARD OF MASTER OF SCIENCE (M.Sc.) DEGREE IN BIOINFORMATICS  
DEPARTMENTMENT OF COMPUTER AND INFORMATION SCIENCES,  
COVENANT UNIVERSITY.**

**JANUARY, 2023**

## **ACCEPTANCE**

This is to attest that this dissertation is accepted in partial fulfilment of the requirements for the award of the degree of Master of Sciences in Bioinformatics in the Department of Computer and Information Sciences, College of Science and Technology, Covenant University, Ota, Nigeria.

**Miss Adefunke F. Oyinloye**

**(Secretary, School of Postgraduate Studies)**

**Signature and Date**

**Prof. Akan B. Williams**

**(Dean, School of Postgraduate Studies)**

**Signature and Date**

## **DECLARATION**

I, ALAGBE, EMMANUEL OLUWATOBA (20PBF02172) declare that this dissertation entitled “A Multi-Omics Classifier for Prediction of Treatment Response in Prostate Cancer Dataset” is a representation of my work, and is written and implemented by me under the supervision of Dr. Itunuoluwa Isewon of the Department of Computer and information sciences, Covenant University, Ota, Nigeria. I attest that this dissertation has in no way been submitted either wholly or partially to any other university or institution of higher learning for the award of a masters’ degree. All information cited from published and unpublished literature has been duly referenced.

Signature:

Date:

## **CERTIFICATION**

This is to certify that this dissertation titled “**A MULTI-OMICS CLASSIFIER FOR PREDICTION OF ANDROGEN DEPRIVATION TREATMENT RESPONSE IN PROSTATE CANCER DATASET**” is original research carried out by **ALAGBE, EMMANUEL OLUWATOBA (20PBF02172)** in the Department of Computer and Information Sciences, College of Science and Technology, Covenant University, Ota, Ogun State, Nigeria under the supervision of Dr. Itunuoluwa Isewon. We have examined and found this work acceptable as part of the requirements for the award of Master of Science (M.Sc.) in Bioinformatics.

**Dr. Itunuoluwa M. Isewon**  
Supervisor

**Signature and Date**

**Prof. Olufunke O. Oladipupo**  
Head of Department

**Signature and Date**

**Prof. Olusegun Folorunsho**  
External Examiner

**Signature and Date**

**Prof. Akan B. Williams**  
(Dean, School of Postgraduate Studies)

**Signature and Date**

## **DEDICATION**

I dedicate this project to God Almighty for His unending supply of grace, and strength given to me during my master's degree Programme. Furthermore, to biomedical and computational researchers who paved the way, providing platforms for insightful biomedical investigations.

## **ACKNOWLEDGEMENTS**

I am eternally grateful to God who has always caused me to triumph and has once again granted me the strength, wisdom, and grace to complete this Programme.

I acknowledge the Chancellor, Dr. David O. Oyedepo for his commitment to the vision of this exemplary university and for remaining a viable role model to us. I am grateful to the Vice-chancellor, Prof. Abiodun H. Adebayo for his visionary leadership and partnership with the chancellor to see the vision of the university become a reality. Furthermore, I acknowledge all members of management and the board of regents for their commitment to the university.

I appreciate the Covenant Applied Informatics and Communication Africa Centre of Excellence (CApIC-ACE) for sponsoring my studentship. I truly appreciate the head of department Prof. Olufunke O. Oladipupo for being accessible and consistently insistent on comprehensive knowledge of the concepts taught in this Programme. My profound gratitude goes to my supervisor Dr. Itunuoluwa M. Isewon for believing in me, supporting me through this Programme, for her ever-present willingness to pay the sacrifice to see me grow, and ultimately for insisting on quality research.

I am eternally grateful to my amazing parents: Pastor and Evangelist Alagbe, for their encouragement, prayers, and selfless sacrifices which continues to serve as a robust foundation for whatever achievement I make. To my wonderful sisters Ayo and Moyinoluwa Alagbe, thanks for believing in your brother. Thanks for your immense financial and emotional support.

I am appreciative of my dear partner and muse: Ogungbesan Temitope, thanks for tutoring me, inspiring me, and brainstorming with me. I'm thankful to my amazing coursemates: Adebola precious, Erika, Shekari, Itunu, Tobi, and Miracle for their amazing support as we journeyed together through the MSc programme. I acknowledge every member of the CUPGF executives, thanks for carrying the burden with me which enabled me to balance managing the fellowship and my academics. To my beloved brothers: Olumikki, Dolad, and Sogo, I'm grateful for your encouragement and your being great examples to me. I am also grateful to my friends without whom my CU story would be largely incomplete: Kaki joy, Ify, Daniel Olusegun, Daniel (Black), Lois, Tomi, Nonye, Tomi (rommie), Folo, Stephen, Samuel, Faith, Adegbesin, Akinbola David, Esther salami, Sonia, Daniella to name a few. You guys have been awesome.

Finally, I'm overflowing with gratitude to my dear friend: Abimbola (mo)Bowofoluwa Sharon, whose persistence resulted in me applying for the ACE scholarship thus starting the cascade leading up to this success. Thanks for your encouraging words and prayers, thanks for being my cheerleader from the very beginning, and thanks for continually inspiring me to be a better researcher and person in general. You are greatly cherished and acknowledged.



# TABLE OF CONTENTS

<b>CONTENTS</b>	<b>PAGES</b>
<b>ACCEPTANCE</b>	III
<b>DECLARATION</b>	IV
<b>CERTIFICATION</b>	V
<b>DEDICATION</b>	VI
<b>ACKNOWLEDGEMENTS</b>	VII
<b>TABLE OF CONTENTS</b>	IX
<b>LIST OF FIGURES</b>	XIII
<b>LIST OF TABLES</b>	XVI
<b>ABSTRACT</b>	XVII
<b>CHAPTER ONE: INTRODUCTION</b>	1
1.1 Background Information	1
1.2 Statement of the Problem	4
1.3 Research Questions	5
1.4 Aim and Objectives of the Study	5
1.5 Research Methodology	6
1.6 Significance of the Study	7
1.7 Scope of the Study	7
1.8 Limitation of the Study	7
1.9 Organization of the study	8
<b>CHAPTER TWO: LITERATURE REVIEW</b>	9
2.1 Introduction	9
2.2 Prostate Cancer	9
2.2.1 Symptoms	10
2.2.2 Risk Factors	11
2.2.4 Treatment Strategies	12
2.3 Machine Learning	13

2.3.1	Supervised Learning	15
2.3.2	Unsupervised Learning	16
2.3.3	Semi supervised Learning	16
2.3.4	Reinforcement Learning	16
2.4	Multi-Omics	18
2.4.1	Adoption of Machine Learning in Omics	20
2.5	Genomics	20
2.5.1	History of Genomics	21
2.5.2	Applications of Genomics	22
2.6	Transcriptomics	22
2.6.1	History of Transcriptomics	22
2.6.2	Applications of Transcriptomics	23
2.7	Proteomics	23
2.7.1	History of Proteomics	24
2.7.2	Applications of Proteomics	24
2.8	Precision Medicine	24
2.9	Review of Related Works	26
2.10	Summary of Findings	29
<b>CHAPTER THREE: RESEARCH METHODOLOGY</b>		<b>30</b>
3.1	Introduction	30
3.2	Datasets	30
3.3	Preprocessing	32
3.3.1	RNAseq data	32
3.3.2	miRNAseq data	33
3.3.3	CNV data	33

3.3.4	RPPA data	33
3.3.5	Integration of multiple Omics	33
3.4	Machine Learning	34
3.5	Performance Evaluation metrics	37
3.6	Pipeline of the Study	38
<b>CHAPTER FOUR: RESULT AND DISCUSSION</b>		<b>40</b>
4.1	Introduction	40
4.2	Results from Objective One	40
4.2.1	Breast Cancer: A well curated multi omics human dataset	40
4.2.2	Preprocessing	41
4.2.3	Performance evaluation (Step 1)	54
4.2.4	Performance evaluation (Step 2)	55
4.3	Results from Objective Two	57
4.3.1	Preprocessing	57
4.3.2	Performance Evaluation (Step 1)	71
4.3.3	Performance Evaluation (Step 2)	73
4.4	Results from Objective Three	74
4.4.1	Omics Comparison (Step 1)	76
4.4.2	Omics Comparison (Step 2)	78
4.5	Results from Objective Four	79
<b>CHAPTER FIVE: CONCLUSION AND RECOMMENDATION</b>		<b>81</b>
5.1	Summary	81
5.2	Conclusion	82
5.3	Contribution to knowledge	82
5.4	Recommendation	82

<b>REFERENCES</b>	<b>83</b>
<b>APENDIX A</b>	<b>116</b>
<b>APENDIX B</b>	<b>137</b>
<b>APENDIX C</b>	<b>141</b>

## LIST OF FIGURES

<b>Figures</b>	<b>Title of Figures</b>	<b>Pages</b>
1.1	Cancer induced morphological changes in the prostate gland	2
2.1	Diagrammatic representation of a machine learning workflow	15
2.2	Reinforcement learning mde of operation	17
2.3	Diagrammatic representation of different learning methods adoption in Oncology	18
2.4	Multi-omics data types and approaches used for studying disease	19
3.1	Flow diagram showing the process of integrating the omics datasets	34
3.2	Flow diagram showing the software used per step in the pipeline of the study	38
3.3	Flow diagram showing the pipeline of the study	38
4.1	MA plot showing the differentially expressed genes between LumA and Basal Breast cancer subtypes in the RNAseq dataset.	42
4.2	MA plot showing the differentially expressed genes between LumA and LumB Breast cancer subtypes in the RNAseq dataset.	43
4.3	MA plot showing the differentially expressed genes between LumA and Her2 Breast cancer subtypes in the RNAseq dataset.	44
4.4	Pie chart revealing the distribution of genes in Basal, LumB, Her2 subtypes differentially expressed from LumA subtype in the RNAseq dataset.	45
4.5	MA plot showing the differentially expressed genes between LumA and Basal Breast cancer subtypes in the miRNAseq dataset	46
4.6	MA plot showing the differentially expressed genes between LumA and LumB Breast cancer subtypes in the miRNAseq dataset.	47
4.7	MA plot showing the differentially expressed genes between LumA and Her2 Breast cancer subtypes in the miRNAseq dataset.	48
4.8	Pie chart revealing the distribution of genes in Basal, LumB, Her2 subtypes differentially expressed from LumA subtype in the miRNAseq dataset.	49
4.9	Distribution of samples in the four Breast cancer dataset.	52
4.10	Distribution of Breast cancer subtypes in the integrated multi-omics datasets.	53
4.11	Distribution of Breast cancer subtypes in the integrated multi-omics datasets after segmenting for step 1 and step 2 experiments.	53

4.12	Performance of the ten Machine learning algorithms in distinguishing LumA from other subtypes in the breast cancer dataset.	55
4.13	Performance of the ten Machine learning algorithms in classifying Basal, LumB, and Her2 subtypes in the breast cancer dataset	57
4.14	MA plot showing the differentially expressed genes between Complete response and Partial Response Prostate cancer treatment outcome groups in the RNAseq dataset	59
4.15	MA plot showing the differentially expressed genes between Complete response and Progressive disease Prostate cancer treatment outcome groups in the RNAseq dataset	60
4.16	MA plot showing the differentially expressed genes between Complete response and Stable disease Prostate cancer treatment outcome groups in the RNAseq dataset	61
4.17	Venn diagram revealing the distribution of DEGs in Partial response, Progressive disease, Stable disease treatment outcome groups in the RNAseq dataset	62
4.18	MA plot showing the differentially expressed genes between Complete response and Partial Response Prostate cancer treatment outcome groups in the miRNAseq dataset	63
4.19	MA plot showing the differentially expressed genes between Complete response and Progressive disease Prostate cancer treatment outcome groups in the miRNAseq dataset.	64
4.20	MA plot showing the differentially expressed genes between Complete response and Stable disease Prostate cancer treatment outcome groups in the miRNAseq dataset.	65
4.21	Venn diagram revealing the distribution of genes in Partial response, Progressive disease, Stable disease treatment outcome groups differentially expressed from Complete response treatment outcome group in the miRNAseq dataset.	66
4.22	Distribution of sample in the four Breast cancer dataset.	69
4.23	Distribution of prostate cancer treatment outcomes in the integrated multi-omics datasets.	70

4.24	Distribution of Breast cancer treatment outcomes in the integrated multi-omics datasets after segmenting for step 1 and step 2 experiments.	70
4.25	Performance of the ten Machine learning algorithms in distinguishing Complete response from other treatment outcome groups in the prostate cancer dataset.	72
4.26	Performance of the ten Machine learning algorithms in distinguishing Partial response, Progressive disease, Stable disease treatment outcome groups in the prostate cancer dataset.	74
4.27	Performance of the Decision tree algorithm in distinguishing Complete response from other treatment outcome groups in the prostate cancer dataset based on the different omics comparisons.	77
4.28	Performance of the Decision tree algorithm in distinguishing Partial response, Progressive disease, Stable disease treatment outcome groups in the prostate cancer dataset based on the different omics comparisons	78
4.29	Pie chart showing the distribution of the predicted treatment outcome labels	80

## LIST OF TABLES

<b>Tables</b>	<b>Title of Tables</b>	<b>Pages</b>
3.1	Dimensionality of Breast Cancer datasets	31
3.2	Dimensionality of Prostate Cancer datasets	32
4.1	Distribution of DEGs in BRCA datasets using LumA as reference	50
4.2	Dimensionality of BRCA datasets before and after preprocessing	51
4.3	Distribution of DEGs in PCa datasets using Complete response as reference	67
4.4	Dimensionality of PCa datasets better and after preprocessing	68
4.5	Dimensionality of different Omics combinations	75
4.6	Prostate cancer patient identifiers and predicted treatment outcomes	79



## ABSTRACT

Prostate cancer (PCa) is estimated to cause over 375,000 deaths and nearly 1.4 million new cases globally. Several factors contribute to PCa heterogeneity, consequently, the stage of the disease decides the strategy employed in combating the disease. The problem of missing data frequently plagues clinical research. The primary treatment outcomes in the TCGA prostate cancer phenotypic dataset had 120 (19.26%) missing values. Treatment strategies could be negatively impacted by limited care giver experience, “trial and error” approaches to treatment, and the genetic makeup of an individual. To the best of our knowledge, using Machine Learning (ML) to forecast treatment response among PCa patients had not been investigated. The aim of this study is to develop a classifier (that acts as a decision support system) from multi-omics datasets for predicting treatment response in PCa patients. RNAseq, miRNAseq, reverse phase protein array (RPPA), copy number variation (CNV) were used in the study. This study employed R programming to preprocess the data. Differential expression analysis for the RNAseq and miRNAseq conducted using the DESeq2 library. Python programming was used to implement the ML algorithms which include XGBoost, Adaboost, multilayer perceptron, decision tree, logistic regression, support vector machine, gradient boosting classifier, Random forests, naive bayes, and K -nearest neighbors. The performance metrics used include macro f1 score, macro recall, macro precision, weighted f1 score, weighted recall, weighted precision, specificity, sensitivity, accuracy, and area under the receiver operator curve. It was discovered that tree-based models were better for the task than probability and kernel-based models. This study computationally demonstrated that multi-omics strategies are generally superior to single-omics strategies, but the adoption of such strategy isn't a foolproof solution. A classifier capable of predicting treatment outcomes amongst PCa patients was built and the predicted labels for patients with missing phenotypic values in the TCGA dataset was provided.

***Keywords: Prostate cancer, Precision Oncology, Multi-omics, Machine Learning Treatment response.***