# TRANSCRIPTIONAL SIGNATURE FOR TUBERCULOSIS PREDICTION USING ENSEMBLE LEARNING TECHNIQUE

## BY

## OKEZIE, ADAUGO FIONA

### (18PCG01762)

### B.Sc Computer Science, University of Nigeria, Nsukka, Enugu State

A DISSERTATION SUBMITTED TO THE SCHOOL OF POSTGRADUATE STUDIES IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE AWARD OF MASTER OF SCIENCE (M.Sc.) DEGREE IN COMPUTER SCIENCE IN THE DEPARTMENT OF COMPUTER AND INFORMATION SCIENCES, COLLEGE OF SCIENCE AND TECHNOLOGY, COVENANT UNIVERSITY.

### DECEMBER, 2020

# ACCEPTANCE

This is to attest that this dissertation is accepted in partial fulfillment of the requirements for the award of Master of Science (M. Sc.) degree in Computer Science in the Department of Computer and Information Science, College of Science and Technology, Covenant University, Ota, Ogun State, Nigeria.


**Mr. John A. Philip**                                      _____

(Secretary, School of Postgraduate Studies)                      Signature and Date


**Prof. Akan B. Williams**                                 _____

(Dean, School of Postgraduate Studies)                           Signature and Date

# DECLARATION

**I, OKEZIE, ADAUGO FIONA** with matriculation number **18PCG01762,** herby declare that this dissertation **TRANSCRIPTIONAL SIGNATURE FOR TUBERCULOSIS PREDICTION USING ENSEMBLE LEARNING TECHNIQUE** was carried out by me under the supervision of Prof. Victor C. Osamor. This project is an original study in the Department of Computer and Information Sciences, College of Science and Technology, Covenant University, Ota, Nigeria. All scholarly information used in this study is fully acknowledged..

**OKEZIE, ADAUGO FIONA** _____

**Signature and Date**

# CERTIFICATION

This is to certify that the dissertation titled "**TRANSCRIPTIONAL SIGNATURE FOR TUBERCULOSIS PREDICTION USING ENSEMBLE LEARNING TECHNIQUE**" is an original research work carried out by **OKEZIE, ADAUGO FIONA** with matriculation number **18PCG01762** under the supervision of Prof Victor .C. Osamor in the Department of Computer and Information Sciences, College of Science and Technology, Covenant University, Ota, Ogun State, Nigeria.


**Prof Victor C. Osamor**                   _____

**(Supervisor)**                                   Signature and Date



**Dr. Olufunke O. Oladipupo**               _____

**(Head of Department)**                          Signature and Date



**Prof. Olumide B. Longe**                  _____

**(External Examiner)**                           Signature and Date



**Prof. Akan B. Williams**                  _____

**Dean, School of Postgraduate Studies**        Signature and Date

# DEDICATION

This project is dedicated to God Almighty for the grace and strength to complete the program and to my family for their love and support throughout my program.

# ACKNOWLEDGEMENTS

I want to thank the almighty God for his sufficient grace to start this program and complete it successfully in due time. I appreciate him for the wisdom, knowledge, and understanding he gave me throughout the program.

I acknowledge and appreciate my parents Engr and Mrs. Ifeanyi Okezie, for their continuous support and love for me throughout my program. Their motivation and encouragement were a steady source of strength during the period. I also acknowledge the support I received from my siblings in all aspects.

Most importantly, to my supervisor Prof. V.C Osamor who has been an ideal thesis and project supervisor for me. His soothing advice, insightful criticism, and patient encouragement aided the writing of this thesis in innumerable ways. He has been more than just a supervisor to me throughout my attachment to him. I want to appreciate him and acknowledge his contributions.

In the department of Computer and Information Science, my gratitude and appreciation go to my HOD, Dr. Olufunke O. Oladipupo and to Prof. A.A. Azeta, for being patient and supportive throughout my entire program. My thanks also go to the entire Management of Covenant University for seeing that the vision has been fulfilled towards accomplishing great excellence and a thorough academic program.

I am grateful to the PG coordinator, Dr. Aderonke Oni, and all the faculty members of the CIS department for their productive remarks during the numerous presentations we had during the cause of this work. I acknowledge Dr. I. Odunayo for his significant contribution, support, and guidance during the project work and throughout my program. I also appreciate Dr. O. Emebo for been supportive throughout the course of my study. I am grateful.

Finally, my appreciation goes to my wonderful friends and collegues, who were with me throughout the program and encouraged me in one way or another. I appreciate them for all efforts made to ensure the program was stress-free.

I appreciate you all, and God bless you.

# TABLE OF CONTENT

# LIST OF FIGURES

# LIST OF TABLES

# ABSTRACT

Tuberculosis has the most considerable death rate among diseases caused by a single micro-organism type. The disease is a significant issue for most third world countries due to poor diagnosis and treatment potentials. Early diagnosis of tuberculosis is the most effective way of managing the disease in patients and reducing the mortality rate caused. Despite the several methods that exist in diagnosing tuberculosis, the limitations ranging from the cost in carrying out the test to the time taken to obtain the results have hindered early diagnosis of the disease. There is, therefore, the need to research alternative diagnostic methods that can aid diagnosis. Better testing with non-sputum samples, like blood, is now desirable and sustainable for diagnoses. Currently, blood transcriptional signatures (genes) are being considered since blood is easily accessible and can tell the state of the body at any point in time, and results can be gotten promptly. Lots of research will discover relevant transcriptional signatures to aid tuberculosis detection. These signatures can be easily observed from analyzing blood samples to know genes triggered in the body after tuberculosis-causing bacteria have infected it. This project work aims to develop a predictive model that would help in the diagnosis of TB and also identify relevant signatures that are affiliated with tuberculosis. The method used to carry out this research involved analyzing tuberculosis gene expression data obtained from GEO (Transcript Expression Omnibus) database and identifying relevant genes used to develop a classification model to aid tuberculosis diagnosis. A classifier combination of K.Nearest Neighbor, Bayes, and Support Vector Machine was used to develop the classification model. The weighted voting ensemble technique was used to improve the classification model's performance. The transcriptional signature obtained from the research includes "C4orf41", "GNPAT", "DHX15", "AGGF1", "ANKRD17", "TM2D1", "VAMP4". while the performance accuracy of the ensemble classifier was 0.95 which showed a better performance than the single classifiers which had 0.94, 0.92 and 0.87 obtained from KNN, SVM and NB respectively. The research clearly shows that the identified signatures can help in the early diagnosis of tuberculosis. The developed model can also assist health practitioners in the timely diagnosis of tuberculosis, which would reduce the mortality rate caused by the disease, especially in developing countries.

**Keywords:** Ensemble Learning, Weighted Voting Method, Tuberculosis, Machine learning Diagnosis, Predictive Model, Biomarkers