

**PERFORMANCE EVALUATION OF ADABOOST AND XGBOOST
ALGORITHMS FOR BREAST CANCER CLASSIFICATION
USING THE WISCONSIN BREAST CANCER DATASET**

Thesis

By

**ARDILA MANDA PUTRI
2217031004**



**FACULTY OF MATHEMATICS AND NATURAL SCIENCES
LAMPUNG UNIVERSITY
BANDAR LAMPUNG**

2026

ABSTRACT

PERFORMANCE EVALUATION OF ADABOOST AND XGBOOST ALGORITHMS FOR BREAST CANCER CLASSIFICATION USING THE WISCONSIN BREAST CANCER DATASET

By

Ardila Manda Putri

Breast cancer is one of the most common types of cancer among women worldwide and remains a leading cause of cancer-related deaths. Early detection plays a crucial role in improving patient survival rates. The advancement of machine learning techniques enables the development of classification models that can assist in the diagnosis of breast cancer more quickly and accurately. This study aims to compare the performance of the AdaBoost and XGBoost algorithms in classifying breast cancer using the Wisconsin Breast Cancer Dataset (WBCD). The research process includes data preprocessing, normalization, and handling class imbalance using the Synthetic Minority Over-sampling Technique (SMOTE). The models were trained using the AdaBoost and XGBoost algorithms and evaluated using the 5-Fold Cross Validation method. Model performance was assessed using several evaluation metrics, including accuracy, precision, recall, and F1-score. The results show that the AdaBoost model achieved an accuracy of 96.2%, precision of 94.7%, recall of 94.4%, and an F1-score of 94.5%. Meanwhile, the XGBoost model achieved an accuracy of 96.1%, precision of 93.3%, recall of 95.5%, and an F1-score of 94.3%. These findings indicate that both algorithms demonstrate excellent classification performance. AdaBoost performs better in terms of accuracy, precision, and F1-score, while XGBoost shows a higher recall in detecting breast cancer cases.

Keywords: machine learning, breast cancer classification, AdaBoost, XGBoost.

ABSTRAK

EVALUASI KINERJA ALGORITMA ADABOOST DAN XGBOOST UNTUK KLASIFIKASI KANKER PAYUDARA MENGGUNAKAN WISCONSIN BREAST CANCER DATASET

Oleh

Ardila Manda Putri

Kanker payudara merupakan salah satu jenis kanker dengan tingkat kejadian yang tinggi pada perempuan di dunia dan menjadi penyebab utama kematian akibat kanker. Deteksi dini sangat penting untuk meningkatkan peluang kesembuhan pasien. Perkembangan metode machine learning memungkinkan pengembangan model klasifikasi yang dapat membantu proses diagnosis kanker payudara secara lebih cepat dan akurat. Penelitian ini bertujuan untuk membandingkan kinerja algoritma AdaBoost dan XGBoost dalam melakukan klasifikasi kanker payudara menggunakan Wisconsin Breast Cancer Dataset (WBCD). Tahapan penelitian meliputi proses preprocessing data, normalisasi, serta penanganan ketidakseimbangan kelas menggunakan Synthetic Minority Over-sampling Technique (SMOTE). Model kemudian dilatih menggunakan algoritma AdaBoost dan XGBoost dengan metode evaluasi 5-Fold Cross Validation. Kinerja model dievaluasi menggunakan metrik accuracy, precision, recall, dan F1-score. Hasil penelitian menunjukkan bahwa model AdaBoost memperoleh accuracy sebesar 96,2%, precision 94,7%, recall 94,4%, dan F1-score 94,5%. Sementara itu, model XGBoost memperoleh accuracy sebesar 96,1%, precision 93,3%, recall 95,5%, dan F1-score 94,3%. Hasil tersebut menunjukkan bahwa kedua algoritma memiliki performa klasifikasi yang sangat baik, dengan AdaBoost unggul pada accuracy, precision, dan F1-score, sedangkan XGBoost memiliki nilai recall yang lebih tinggi dalam mendeteksi kasus kanker payudara.

Kata-kata kunci: *machine learning*, klasifikasi kanker payudara, AdaBoost, XGBoost.

**PERFORMANCE EVALUATION OF ADABOOST AND XGBOOST
ALGORITHMS FOR BREAST CANCER CLASSIFICATION
USING THE WISCONSIN BREAST CANCER DATASET**

ARDILA MANDA PUTRI

Thesis

**In a Partial Fulfillment of The Requirements for
BACHELOR OF MATHEMATICS**

In the

Department of Mathematics

Faculty Of Mathematics And Natural Sciences



**FACULTY OF MATHEMATICS AND NATURAL SCIENCES
LAMPUNG UNIVERSITY
BANDAR LAMPUNG**

2026

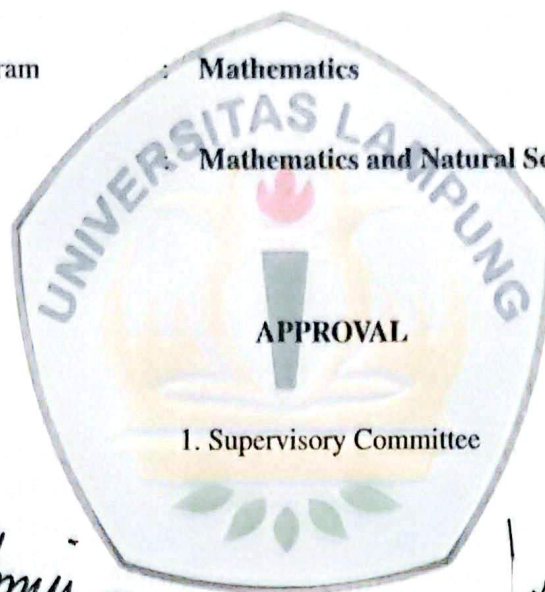
Thesis Title : PERFORMANCE EVALUATION OF ADABOOST AND XGBOOST ALGORITHMS FOR BREAST CANCER CLASSIFICATION USING THE WISCONSIN BREAST CANCER DATASET

Student Name : Ardila Manda Putri

ID Number Of Student : 2217031004

Study Program : Mathematics

Faculty : Mathematics and Natural Sciences



1. Supervisory Committee

Dr. Khoirin Nisa, S.Si, M.Si.
NIP 197407262000032001

Misgiyati, S.Pd, M.Si.
NIP 198509282023212032

2. Associate Dean for Academic Affairs and Collaboration
Faculty of Mathematics and Natural Sciences Lampung University

Mulyono, S.Si., M.Si., Ph.D.
NIP. 197406112000031002

VALIDATED BY

1. Examination Committee

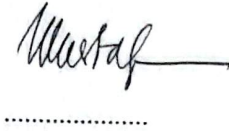
Head : Dr. Khoirin Nisa, S.Si, M.Si.



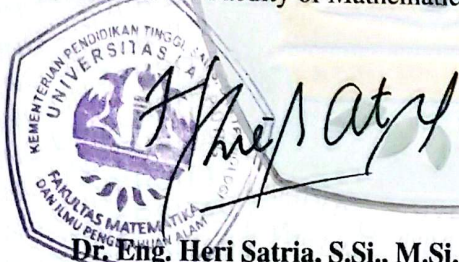
Secretary : Misgiyati, S.Pd, M.Si.



Examiner
Non-Supervisor : Prof. Mustofa, M.A., Ph.D.



2. Dean of the Faculty of Mathematics and Natural Sciences



Dr. Eng. Heri Satria, S.Si., M.Si.
NIP. 197110012005011002

Thesis Defense Date: 19 May 2026

STUDENT THESIS STATEMENT

The undersigned below:

Name : **Ardila Manda Putri**
ID Number Of Student : **2217031004**
Major : **Mathematics**
Thesis Title : **Performance Evaluation of AdaBoost and XGBoost Algorithms for Breast Cancer Classification Using the Wisconsin Breast Cancer Dataset**

Hereby declare that this thesis is the result of my own work and all writings contained in this thesis have followed the rules of scientific writing at the Lampung University.

Bandar Lampung, 19 May 2026

Author,



Ardila Manda Putri

BIOGRAPHY

Ardila Manda Putri was born in West Lampung Regency, Lampung Province, on August 12, 2004. She is the third of four children of Mr. Suroso and Mrs. Susi Patrawati. She has two older brothers, Agus Nendy and Edho Apriliyansyah, and one younger brother, Sheva Lukiyanto.

She began her early education at Gelora Mekar in the academic year 2008-2010. elementary education at SDN 1 Karang Agung from 2010 to 2016. She then continued her junior high school education at SMPN 1 Way Tenong from 2016 to 2019, and later pursued her senior high school education at SMAN 1 Way Tenong from 2019 to 2022.

In 2022, Ardila Manda Putri continued her undergraduate studies in the Department of Mathematics, Faculty of Mathematics and Natural Sciences, University of Lampung through the SNMPTN admission pathway. During her time as a student, the author frequently participated in the annual Dies Natalis events of the Department of Mathematics, serving as a Member of the Mathematics Competition Division in 2023, Coordinator of the Mathematics Competition Division in 2024, and Steering Committee of the Mathematics Competition Division in 2025.

From December 2024 to January 2025, Ardila Manda Putri carried out Practical Work (KP) or an internship at the Regional Development Planning, Research, and Innovation Agency of Bandar Lampung City as a form of self-development and application of the knowledge gained during her studies. Furthermore, from August to November 2025, the author also participated in an internship program at the Central Bureau of Statistics of Bandar Lampung City. Then, from January to February 2026, the author conducted a Community Service Program (KKN) in Gotong Royong Subdistrict, Central Tanjung Karang District, Bandar Lampung City.

WORD OF INSPIRATION

“Allah will not burden a person but according to his ability.”

(Q.S Al-Baqarah: 286)

“With difficulty there is ease, so when you have finished a task, do it with all diligence.”

(Q.S Al-Inshirah: 6)

“Verily, Allah will not change the condition of a people until they change their own condition.”

(Q.S Ar-Rad: 11)

“Tell yourself, tomorrow we may arrive, tomorrow we may achieve it.”

(Baskara Putra - Besok Mungkin Kita Sampai)

“All your rises and falls are a natural part of life. Dreams and questions will be answered in time. Give yourself a limit for grieving—be sad only as much as needed. Celebrate your feelings as a human being.”

(Baskara Putra - Mata Air)

DEDICATION

By expressing my deepest gratitude to Allah SWT for all His mercy, blessings, and grace, the author sincerely dedicates this humble work with great respect and love to:

Thank you for the love, sacrifices, guidance, and prayers that Father and Mother gave during your lifetime. Although you are no longer here, your love and advice will always live in the author's heart. Every step and achievement accomplished by the author is inseparable from the struggles and values of life that you instilled. May this work become one form of devotion and pride for Father and Mother in His presence.

Thank you to the extended family who have always given prayers, care, support, and encouragement to the author in completing this work.

Thank you to the lecturers who have provided guidance, direction, motivation, and valuable knowledge throughout the process of completing this work.

Thank you to my dear friends and loved ones who have always been present in giving support, encouragement, and prayers throughout every process the author has gone through.

To my beloved almamater, Lampung University.

ACKNOWLEDGEMENT

Praise be to Allah SWT and blessings and greetings be upon the beloved prophet Muhammad SAW. With His mercy and grace, the author was able to complete this thesis entitled “Performance Evaluation of AdaBoost and XGBoost Algorithms for Breast Cancer Classification using the Wisconsin Breast Cancer Dataset”.

The completion of this final work cannot be separated from the support, guidance, advice, and prayers from various parties. Therefore, on this occasion the author would like to express her gratitude to:

1. Mrs. Dr. Khoirin Nisa, S.Si, M.Si., as the first supervisor for her patience and willingness to provide guidance, assistance, motivation, input, and support to the author in completing this thesis.
2. Mrs. Misgiyati, S.Pd., M.Si., as the second supervisor for her patience and willingness to provide direction, support and helpful advice to the author in completing this thesis.
3. Prof. Mustofa, M.A., Ph.D., as a thesis examiner who has provided constructive criticism and suggestions during the thesis preparation process.
4. Mr. Dr. Aang Nuryaman, S.Si., M.Si., as the Head of the Department of Mathematics, Faculty of Mathematics and Natural Sciences, Lampung University.
5. Prof. Dra. Wamiliana, M.A., Ph.D., as Academic Supervisor, who has provided guidance and support during the study period.
6. Mr. Dr. Eng. Heri Satria, S.Si., M.Si as the Dean of the Faculty of Mathematics and Natural Sciences, Lampung University.
7. All lecturers, staff, and employees of the Department of Mathematics, Faculty of Mathematics and Natural Sciences, Lampung University who have helped during the lecture period.

8. To the late Mr. Suroso, whom I called Father, and the late Mrs. Susi Patrawati, whom I called Mother. You are the two souls I miss the most, who empowered me to rise whenever I felt like giving up. Thank you for everything you have given. Though in the end, I must struggle and persevere on my own without your presence by my side, I offer my prayers as a substitute for your presence. May every step I take become a continuous blessing for both of you in eternity. May Allah SWT grant you the most beautiful place by His side. Thank you for being the reason I never gave up on life.
9. The three beloved brothers Nendy, Edho and Sheva, thanks to them for always giving enthusiasm, prayers, and great affection to the author and hopefully always compact and provide happiness for the family so that it makes motivation to complete this thesis.
10. To my sister in law, Mesa. Thank you for being a part of our family. I feel like I haven't just gained an in-law, but a sibling of my own. Thank you for all your support and kindness all this time.
11. To my beloved nieces and nephew, Cinta, Abiyan, and Qinnar. Thank you for being my source of happiness and for cheering me up throughout this journey. May you grow up to be good-hearted and responsible children.
12. For my dear friends who are like family: Nazwa, Hani, Erisa, Nadhia, and Kinasih. My deepest gratitude for being an incredible support system during the challenging process of completing this thesis. Thank you for being a significant part of my journey. I am forever grateful for the laughter, the shared coffee, and the constant encouragement that helped me persevere.
13. To my three best friends since high school, Alya, Elisa and Novita. Thank you for your support from afar and for all the care you've shown me throughout this time. This thesis is finally complete thanks to your prayers and encouragement.
14. For Selvani, who is both a friend and a sister to me. Thank you for being my constant support and for always lending an ear to my every concern. I pray that your future journey is made easy and that you are blessed with endless joy.
15. All parties who have helped the author in completing this thesis.
16. Finally, I want to thank myself. Thank you for the resilience you showed when life felt lonely without the presence of parents. Thank you for embracing

independence, for surviving the struggles alone, and for pushing forward to this very moment. Thank you for not surrendering to your situation and for successfully proving that you can stand firm and thrive on your own

May Allah SWT reward all the kindness that has been given in the best possible way. Hopefully, this thesis can provide benefits for the readers. The author realizes that this thesis writing is far from perfect. Therefore, criticism and suggestions are highly expected to improve this work.

Bandar Lampung, 19 May 2026

Ardila Manda Putri

TABLE OF CONTENTS

TABLE OF CONTENT	xiii
LIST OF TABLE	xv
LIST OF FIGURE	xvi
I. INTRODUCTION	1
1.1 Background and Issues	1
1.2 Research Problems	3
1.3 Research Objectives	3
1.4 Research Benefits	3
II. LITERATURE REVIEW	4
2.1 Machine Learning	4
2.2 Classification Analysis	6
2.3 Ensemble Learning	8
2.4 AdaBoost	11
2.5 Gradient Boosting	13
2.6 Extreme Gradient Boosting (XGBoost)	15
2.7 Model Evaluation	17
2.8 Breast Cancer	19
III. RESEARCH METHODS	21
3.1 Time and Place of Research	21
3.2 Research Data	21
3.3 Research Methods	21
IV. RESULTS AND DISCUSSION	24
4.1 Breast Cancer Data Characteristics	24
4.2 Data Preprocessing	27
4.2.1 Checking for Missing Values	27
4.2.2 Checking for Duplicate Data	28
4.2.3 Data Type Checking	30
4.3 Data Imbalance	30
4.4 Data Normalization	31

4.5	Splitting Data with K-Fold Cross Validation	32
4.6	AdaBoost and XGBoost Models	38
4.7	Model Evaluation Results	47
V.	CONCLUSION	76
5.1	Conclusion	76
5.2	Recommendation	76
	LITERATURE	77
	APPENDIX	82

LIST OF TABLE

1	The first 10 datasets used	25
2	Results of preprocessing missing values	28
3	Results of preprocessing duplicate data	29
4	Results of preprocessing type data	30
5	Splitting on the 1st Fold	32
6	Splitting on the 2nd Fold	33
7	Splitting on the 3rd Fold	35
8	Splitting on the 4th Fold	36
9	Splitting on the 5th Fold	37
10	Example dataset AdaBoost	38
11	Example dataset XGBoost	42
12	The predictions for the first stump	43
13	Average Evaluation Results of AdaBoost and XGBoost Models	72

LIST OF FIGURE

1	Various types of machine learning techniques (Source: Sarker, 2021)	6
2	Framework of ML predictive model (Source: Sarker, 2021)	8
3	Confusion Matrix (Source: Salvador, 2024)	18
4	Flowchart Research Methodology	23
5	Countplot class distribution in the dataset	31
6	Countplot class distribution after SMOTE – Fold 1	33
7	Countplot class distribution after SMOTE – Fold 2	34
8	Countplot class distribution after SMOTE – Fold 3	35
9	Countplot class distribution after SMOTE – Fold 4	36
10	Countplot class distribution after SMOTE – Fold 5	37
11	AdaBoost Confusion Matrix – Fold 1	48
12	XGBoost Confusion Matrix – Fold 1	50
13	AdaBoost Confusion Matrix – Fold 2	53
14	XGBoost Confusion Matrix – Fold 2	55
15	AdaBoost Confusion Matrix – Fold 3	58
16	XGBoost Confusion Matrix – Fold 3	60
17	AdaBoost Confusion Matrix – Fold 4	63
18	XGBoost Confusion Matrix – Fold 4	65
19	AdaBoost Confusion Matrix – Fold 5	68
20	XGBoost Confusion Matrix – Fold 5	70

I. INTRODUCTION

1.1 Background and Issues

The advancement of artificial intelligence technology, especially in the domain of machine learning, has produced various algorithms that can improve the accuracy of data classification. One widely used approach is ensemble learning, a method that combines several simple models (weak learners) to form a stronger model with better performance. Among ensemble techniques, boosting is one of the most popular due to its ability to gradually correct classification errors through a weighting process (Freund & Schapire, 1997). Two notable boosting algorithms are AdaBoost (Adaptive Boosting) and XGBoost (Extreme Gradient Boosting). AdaBoost works by giving greater weight to data that was misclassified in the previous iteration, so that the next model focuses more on that data. This algorithm is simple, effective, and suitable for use with weak learners such as decision stumps or shallow decision trees. However, its weakness is that it is sensitive to outliers and noise, because problematic instances can be given high weights, thereby reducing the performance of the model (Asri, et al., 2016).

XGBoost, on the other hand, is an advanced implementation of gradient boosting equipped with computational optimization, regularization, and feature pruning to reduce overfitting. XGBoost is known for its speed, memory efficiency, and flexibility in hyperparameter settings (Chen & Guestrin, 2016). However, its high complexity requires the selection of appropriate parameters to avoid overfitting, especially in small datasets (Mathew, 2023). The selection of these two methods in this study was not without reason. First, boosting has been proven to improve accuracy, recall, and F1-score in various classification datasets. Second, both AdaBoost and XGBoost are algorithms that consistently rank at the top in competitions and scientific studies, especially in the medical domain (Asri, et al., 2016). Third, a direct comparison of

the two on the same dataset can provide insights into their advantages, disadvantages, and potential applications in the real world.

One important domain of machine learning application is healthcare, particularly breast cancer classification. Breast cancer is among the most common cancers in women globally and a primary cause of death. Early detection of breast cancer is crucial to improving patient survival rates, requiring accurate, fast, and reliable diagnostic methods (Xiao, et al., 2019). In research related to machine learning, the Wisconsin Breast Cancer Dataset (WBCD) is often used as a benchmark. This dataset contains numerical features extracted from breast tissue samples labeled as benign or malignant (Freund & Schapire, 1997). WBCD is an ideal dataset for testing the performance of classification algorithms. The evaluation of the AdaBoost and XGBoost algorithms on this dataset is expected to illustrate the effectiveness of boosting in distinguishing malignant and benign cancer cells. A number of previous studies have shown that boosting, particularly XGBoost, often performs better than other algorithms such as Logistic Regression, Support Vector Machine, or Random Forest (Asri, et al., 2016). AdaBoost has also been reported to improve accuracy and sensitivity, although in some studies its performance is slightly below that of XGBoost. Recent research has even combined boosting with feature selection techniques to improve accuracy to over 95 % on WBCD (Asri, et al., 2016).

This shows the high relevance of both algorithms in supporting breast cancer diagnosis. However, there are various challenges in applying boosting to breast cancer classification. First, the risk of overfitting remains a problem, especially in complex algorithms such as XGBoost, if the parameters are not tuned appropriately (Chen & Guestrin, 2016). Second, the issue of interpretability makes the prediction results difficult for medical personnel to understand, thereby reducing the level of confidence in clinical decision-making. To address these issues, various studies have proposed approaches such as SMOTE to balance the data, the use of regularization and cross-validation to prevent overfitting (Mathew, 2023). These innovations reinforce the importance of research comparing AdaBoost and XGBoost, especially in the context of breast cancer classification.

Therefore, the Wisconsin Breast Cancer Dataset is used in this study to assess how well the AdaBoost and XGBoost algorithms perform for classifying breast cancer. It is anticipated that the study's findings will advance scientific knowledge of how well

both algorithms function, as well as provide recommendations for the application of accurate, robust, and reliable machine learning in supporting early detection of breast cancer.

1.2 Research Problems

The research problems can be formulated as follows:

1. How well does the AdaBoost and XGBoost algorithm perform in classifying breast cancer using the Wisconsin Breast Cancer Dataset (WBCD)?
2. How do the performances of AdaBoost and XGBoost compare in terms of evaluation metrics such as accuracy, precision, recall, and F1-score?

1.3 Research Objectives

The objectives of this study are as follows:

1. To examine the performance of the AdaBoost and XGBoost algorithm in classifying breast cancer using the WBCD.
2. To compare the performance outcomes of AdaBoost and XGBoost through evaluation indicators such as accuracy, precision, recall, and F1-score.

1.4 Research Benefits

The benefits of this study are:

1. Contributes to the scientific understanding of machine learning, especially regarding the effectiveness of boosting methods in medical dataset classification.
2. Serves as a reference for future studies exploring ensemble learning or disease classification using machine learning algorithms.
3. Provides insights for healthcare practitioners and medical researchers on the potential use of AdaBoost and XGBoost to support early breast cancer detection.

II. LITERATURE REVIEW

2.1 Machine Learning

Machine Learning (ML) is a sector of artificial intelligence that concentrates on creating algorithms that can learn autonomously from data to generate predictions or decisions (Tian, et al., 2024). In recent years, artificial intelligence (AI), especially machine learning (ML), has experienced rapid growth in the realms of data analysis and computing, enabling applications to operate intelligently (Tian et al., 2024). ML typically equips systems with the capability to learn and improve from experience automatically, without needing explicit programming, and is commonly regarded as one of the most significant contemporary technologies in the fourth industrial revolution (4IR or Industry 4.0) (Tian, et al., 2024).

In the early days of AI research, the emphasis was on using explicitly defined statements in formal languages that computers could utilize for reasoning through logical inference rules. This approach is referred to as the knowledge base method (Janiesch, et al., 2021). Nevertheless, this paradigm has numerous limitations, as humans often find it challenging to articulate all the implicit knowledge necessary for executing complex tasks. Machine learning addresses these shortcomings. Broadly speaking, machine learning refers to a method whereby a computer program enhances its performance through experience in relation to specific tasks and performance metrics. This method's goal is to automate the process of creating analytical models that can perform cognitive tasks like natural language translation and object detection. This is achieved by using algorithms that learn iteratively from particular training data, allowing computers to discover complex patterns and hidden insights without explicit programming. Machine learning is quite effective, especially when it comes to high-dimensional data tasks like grouping, regression, and classification. It can help produce consistent and trustworthy decision-making by learning from previous analyses and finding patterns in large data sets. Consequently,

machine learning algorithms have been effectively utilized across various fields, including fraud detection, credit assessment, analysis of the next-best offer, as well as in speech and image recognition and natural language processing (NLP) (Janiesch, et al., 2021).

Machine Learning algorithms are mainly divided into four categories: Supervised learning, Unsupervised learning, Semi-supervised learning, and Reinforcement learning (Sarker, 2021).

1. Supervised

Using sample pairs of inputs and outputs, supervised learning often uses machine learning to find a function that links input to output. To infer a function, it uses a selection of training cases and labeled training data. Supervised learning takes a task-oriented approach and is carried out when specified goals are identified as being achievable from a given set of inputs. The two most common supervised tasks are "regression," which seeks to fit the data, and "classification," which distinguishes the data. One example of supervised learning is text categorization, which is the process of identifying the class label or sentiment of a text, such a tweet or a product review.

2. Unsupervised

Unsupervised learning is a data-driven method that analyzes unlabeled datasets without the need for human intervention. This approach is frequently used to organize results, identify important patterns and structures, find generative features, and carry out exploratory research. Clustering, density estimation, feature learning, dimensionality reduction, association rule discovery, and anomaly detection are the main tasks related to unsupervised learning.

3. Semi-supervised

Using both labeled and unlabeled data, semi-supervised learning can be defined as a hybrid of the supervised and unsupervised methods discussed above. It is therefore in the middle between learning "without supervision" and learning "with supervision." Semi-supervised learning is beneficial in real-world situations where labeled data may be limited in certain settings while unlabeled data is plentiful. More accurate predictions than those derived

only from the labeled data utilized in the model are the main goal of a semi-supervised learning model. Machine translation, fraud detection, data labeling, and text categorization are a few domains that use semi-supervised learning.

4. Reinforcement

Reinforcement learning is a type of machine learning that uses an environment-driven method to enable software agents and machines to independently determine the optimal course of action in a given context or environment to improve their performance. This approach, which is based on the concepts of rewards and penalties, seeks to use knowledge gathered from interactions with the environment to inform choices that either raise rewards or lower risks. It is a potent tool for training AI models, which can facilitate greater automation or improve the operational effectiveness of intricate systems like robotics, self-driving cars, and supply chain and manufacturing logistics. Nevertheless, it is not recommended to use it for easier or more straightforward issues.

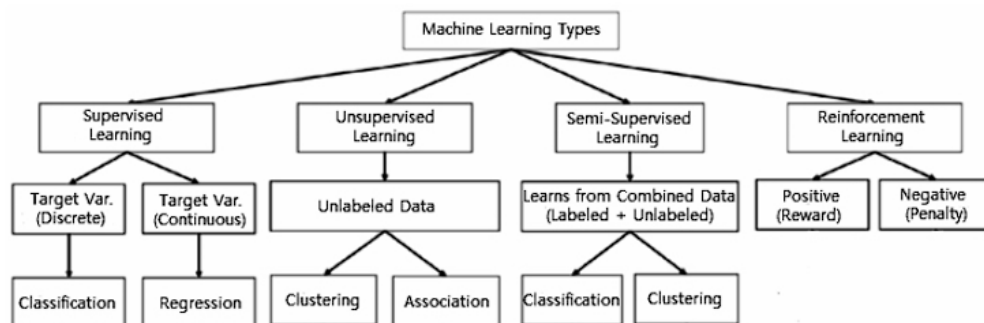


Figure 1 Various types of machine learning techniques (Source: Sarker, 2021)

2.2 Classification Analysis

Classification is considered a supervised learning approach in machine learning, which deals with predictive modeling where a class label is assigned to a particular example. To explain it mathematically, it creates a function (f) that converts input variables (X) into output variables (Y), which stand for labels, categories, or targets. This procedure can be used for both structured and unstructured data to identify the class of particular data points. For instance, identifying whether an email is "spam"

or "not spam" in email service providers is an example of a classification problem (Sarker, 2021).

1. Binary Classification

Binary classification tasks involve categorizing data into two distinct class labels, such as "true or false" or "yes or no." In these tasks, one class may represent the normal condition, while the other signifies an abnormality. For example, in a medical test scenario, "cancer not detected" is the normal state, while "cancer detected" represents the abnormal condition. Likewise, in the context of email services, "spam" and "not spam" are examples of binary classification.

2. Multiclass classification

This phrase usually refers to classification tasks with more than two class labels. The idea of normal versus abnormal outcomes does not apply to multiclass classification tasks, in contrast to binary classification tasks. Rather, examples are classified as being a part of a specific class. One example of a multiclass classification problem is the categorization of several network attack types in the NSL-KDD dataset, where attack categories are separated into four class labels: R2L (Root to Local Attack), U2R (User to Root Attack), DoS (Denial of Service Attack), and probing attack.

3. Multi-label Classification

In machine learning, multi-label classification is a critical aspect where a single example can be linked to multiple classes or labels. As in multi-level text classification, this idea goes beyond multiclass classification, where the classes taken into consideration are arranged hierarchically, enabling each sample to belong to multiple classes at each level of the hierarchy. Articles on Google News, for instance, may be grouped under "technology," "city name," or "latest news," among other headings. Unlike classic classification challenges where class labels are exclusive, multi-label classification uses advanced machine learning algorithms that can predict many non-exclusive classes or labels.

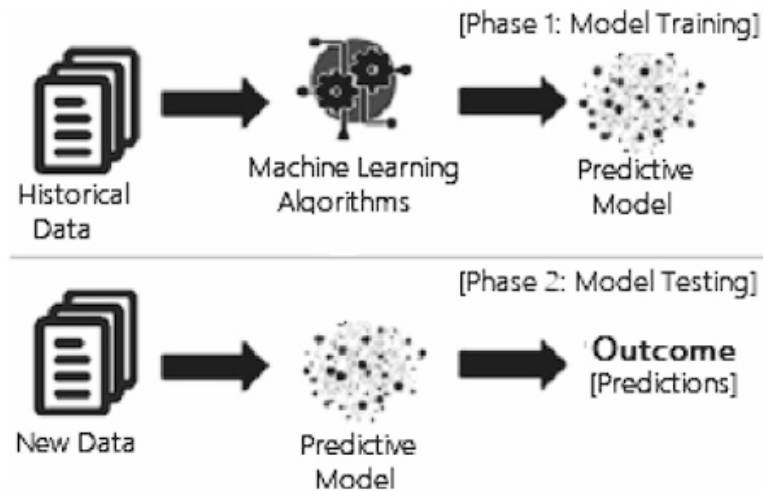


Figure 2 Framework of ML predictive model (Source: Sarker, 2021)

2.3 Ensemble Learning

Ensemble learning is a method that integrates two or more machine learning algorithms to achieve better performance than when the algorithms are applied separately. To create a more accurate overall prediction, the predictions of each individual learner are combined using a combination rule rather than relying on a single model. Ensemble approaches can generally be separated into sequential and parallel ensembles. The parallel approaches train various base classifiers independently and merge their predictions using a combiner (Mienye & Sun, 2022). Along with its extension, the random forest algorithm, bagging is a well-known parallel ensemble approach. Parallel ensemble algorithms leverage the simultaneous creation of base learners to promote diversity among the members of the ensemble (Liu, et al., 2016). Conversely, sequential ensembles do not train the base models independently. They are trained in a stepwise manner, so that each model in every iteration learns to remedy the mistakes made by the preceding model (Mienye & Sun, 2022).

The accuracy and adaptability of each learner are the main factors that determine how effective ensemble learning techniques are. If a machine learning model exhibits significant generalization abilities on data that has never been seen before, it is considered exact. Conversely, ML models exhibit diversity when their errors on new instances differ (Bian, et al., 2020). Hence, diversity is regarded as the

distinction among the base learners within an ensemble. Unlike accuracy, measuring diversity lacks a universally accepted guideline (Alshdaifat, et al., 2021). When building ensemble classifiers, it is also difficult to achieve diversity among the basic models. The base learners in many ensembles are created using overlapping sections of the same training data, which results in linked models and makes it challenging to promote variety. Various ensemble strategies attempt to promote diversity either heuristically or implicitly. For example, bagging fosters diversity through the subsampling of the training data, while boosting does so by reweighting the training data (Liu, et al., 2016). Below are some methods of ensemble learning.

1. Boosting

A machine learning approach called "boosting" turns weak learners into effective classifiers. It is an ensemble meta-algorithm that aims to reduce variance and bias. A weak learner is defined as a classifier that performs slightly better than random chance, whereas strong learners achieve high accuracy and form the backbone of the boosting ensemble algorithms (Sun, et al., 2021). The concept of the boosting algorithm was initially presented by Schapire (1990) in response to Kearns and Valiant's question about whether multiple weak learners could be combined to create a strong learner. Schapire's (1990) research profoundly influenced the fields of machine learning and statistics, leading to the creation of various boosting algorithms, such as AdaBoost and XGBoost (Sun, et al., 2021).

The fundamental concept of boosting is the iterative application of the base learning algorithm to modified input data. Therefore, it enhances base learners with a high bias and low variance, such as decision stumps (a decision tree with one internal node) (Mienye, et al., 2020). The base learner focuses on samples that are misclassified because they are given more weight. As a result, if the base classifier is biased against particular samples, the algorithm corrects the bias by giving those examples greater weight. However, this iterative learning approach makes boosting unsuitable for learning noisy data because the weight given to noisy samples is usually much greater than the weights given to the other samples, thereby forcing the algorithm to focus excessively on the noisy samples, resulting in overfitting (Mienye, et al., 2020). Nevertheless, one of the most effective algorithms in applied machine learning is boosting-based ensemble approaches. Boosting learning includes

AdaBoost, Gradient Boosting, XGBoost, LightGBM and CatBoost.

2. Bagging

Bootstrap aggregating, commonly known as bagging, was introduced by Breiman in 1994 to improve the classification accuracy of machine learning models by amalgamating predictions from randomly created training datasets (Mishra, et al., 2022). In this approach, diversity is achieved by generating bootstrapped versions of the original data, wherein different subsets of the input data are chosen at random, with replacement, from the initial training dataset. As a result, the different training sets are considered diverse and utilized to train multiple base learners using the same machine learning algorithm. Essentially, the bagging technique entails dividing the training data for each base learner through random sampling, thereby producing b distinct subsets that are employed to train b base learners (Mishra, et al., 2022).

Base learners perform better when bagging is used, especially when the learning method is unstable. Instead of addressing bias, its main goal is to lessen the variance among the ensemble members. Consequently, bagging achieves the best results when the ensemble members exhibit high variance and low bias (Alelyani, 2021). Bagging's ability to reduce variance without raising bias is a major advantage. Additional advantages of bagging include its ability to use the bootstrapping technique to produce variation in the input data. For large datasets, bagging requires less computational time compared to many machine learning algorithms because it trains the model using a smaller sample size (Alelyani, 2021). Bagging has the disadvantage of increasing model accuracy without taking interpretability into account.

3. Stacking

Stacked generalization, commonly known as Stacking, is an ensemble learning method that institutes a distinct machine learning algorithm to merge the predictions from multiple ensemble members. It specifically entails constructing models with several base algorithms, referred to as level-0 models, and employing a meta-learning algorithm that trains an additional model to integrate the predictions from these base models (Mienye, et al., 2020). The

fundamental concept of stacking is that the level-0 base learners are trained on the training dataset and are then evaluated on out-of-sample or unseen data; the predicted labels from these base models on the unseen data, along with the actual labels, create input-output pairs for a new dataset intended for meta-learner training (Liang, et al., 2021). Meta-learning represents the aspect of machine learning wherein algorithms utilize the outputs of other ML algorithms to generate more precise predictions based on the outputs of the base classifiers (Hospedales, et al., 2022). Because it combines the strengths of multiple powerful classifiers, the stacking strategy is very successful in producing classifications that outperform the individual models within the ensemble.

Moreover, stacking employs various base algorithms along with the same dataset to create models that are diverse and tackle the predictive modeling challenge in different ways. In contrast to bagging, which primarily utilizes decision tree models trained on portions of the input data, the stacked models leverage multiple algorithms trained on the same dataset (Miguel-Hurtado, et al., 2016). Additionally, stacking uses a single model to determine the optimal strategy to combine the predictions from the base learners, in contrast to boosting, which trains models sequentially to improve the predictions of earlier models. However, in regression tasks, linear regression is predominantly used as the meta-classifier (Miguel-Hurtado, et al., 2016).

2.4 AdaBoost

AdaBoost is the first boosting method created by Freund and Schapire. It utilizes weak classifiers to form a more robust and reliable classifier. As noted by Tavish, the process begins with the base learner, or weak learner, giving equal importance to every observation. After assigning weights to all observations, the weak learner is then applied for predictions (Shahri, et al., 2021). Observations that the weak learner incorrectly labels are given increased weights, and the subsequent base learner takes over for prediction. This method continues until it reaches the limit of the base learning algorithm, represented by T_i , with the iteration marked as t . In the concluding stage, the outputs from the weak learners are merged to create a more powerful learner that improves prediction accuracy (Shahri, et al., 2021).

The learning process of AdaBoost consists of training an initial classifier by utilizing a basic algorithm, which is often a decision tree. Based on the classifier's predictions, the sample weights are changed, and the updated samples are then used to train the subsequent classifier. As a result, the samples that were misclassified receive higher weights, while the samples that were correctly classified are given lower weights, making sure that the upcoming classifiers focus more on the misclassified samples (Wang, et al., 2019). The various base learners are incorporated one after another and assigned weights to create a powerful classifier. Given m labelled training instances $S = \{(x_1, y_1), \dots, (x_i, y_i), \dots, (x_m, y_m)\}$, where y_i is the target label of sample x_i , and $y_i \in Y = \{-1, +1\}$, the weight D_1 of the sample x_i and the weight update D_{t+1} are computed as:

$$D_1(i) = \frac{1}{n}, i = 1, 2, \dots, n \quad (2.4.1)$$

and

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \exp(-\alpha_t y_i h_t(x_i)), i = 1, 2, \dots, n \quad (2.4.2)$$

The function $h_t(x)$ serves as the base classifier, with $t = 1, \dots, T$ representing the number of iterations. The notation Z_t denotes the normalization factor, while α_t refers to the weight assigned to the classifier $h_t(x)$. The weight α_t indicates the level of contribution of the classifier $h_t(x)$ in forming the final decision. In cases where $h_t(x)$ misclassifies an instance, that data point is assigned a higher weight in the subsequent iteration ($t+1$). Moreover, Z_t is determined in such a way that D_{t+1} remains a valid distribution. The values of Z_t and α_t are obtained through the following equations:

$$Z_t = \sum_{i=1}^n D_t(i) \exp(-\alpha_t y_i h_t(x_i)) \quad (2.3)$$

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right) \quad (2.4)$$

where ϵ_t symbolizes the classifier's error rate, which is acquired using:

$$\epsilon_t = P[h_t(x_i) \neq y] = \sum_{i=1}^n D_t(i) \mathbb{I}[h_t(x_i) \neq y_i] \quad (2.5)$$

Following the completion of the specified number of iterations, the final strong classifier is calculated using:

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right) \quad (2.6)$$

AdaBoost is simple to use and requires minimal hyperparameter adjustment. Moreover, AdaBoost's adaptability allows it to use a variety of algorithms as its foundational learner. As a result, it can improve performance by using a method that is appropriate for a particular application. However, because AdaBoost uses an iterative learning approach, it is susceptible to outliers and noisy data, which can lead to overfitting (Mienye, et al., 2020).

2.5 Gradient Boosting

A machine learning technique called gradient boosting uses the boosting approach to create strong ensembles. Known as gradient boosted decision trees (GBDT), it mainly uses decision trees as the fundamental learner to produce a powerful ensemble classifier. The gradient boosting method was initially presented by Breiman, who observed that boosting can be viewed as an optimization method for a suitable loss function (Mienye, et al., 2020). The training process of this algorithm consists of progressively training new models to build a strong classifier. It is created in a gradual fashion, akin to other boosting methods, but its main concept is to produce base learners that are closely aligned with the negative gradient of the loss function associated with the complete ensemble (Mienye, et al., 2020).

Suppose we are given a training dataset $S = \{x_i, y_i\}_{i=1}^N$. The gradient boosting technique seeks to construct an approximation $\hat{F}(x)$ of the target function $F^*(x)$, which maps the predictor variables x to the response variables y , by minimizing the loss function $L(y, F(x))$. In essence, GBDT forms an additive approximation of $F(x)$ through a weighted summation of functions:

$$F_m(x) = F_{m-1}(x) + \rho_m h_m(x) \quad (2.7)$$

where ρ_m represents the weight of the m^{th} function, $h_m(x)$. These functions are the

decision tree models in the ensemble. The algorithm performs the approximation iteratively. Meanwhile, a constant approximation of $F^*(x)$, is achieved using:

$$F_0(x) = \arg \min_a \sum_{i=1}^n L(y_i, a) \quad (2.8)$$

successive base learners aim to minimize

$$(\rho_m, h_m(x)) = \arg \min_{\rho, h} \sum_{i=1}^N L(y_i, F_{m-1}(x_i) + \rho h(x_i)) \quad (2.9)$$

Instead of directly carrying out the optimization task, each h_m can be interpreted as a greedy step in the gradient descent optimization of $F^*(x)$ (Bentéjac, et al., 2021). Consequently, each h_m is trained on a new dataset $D = \{x_i, r_{mi}\}_{i=1}^N$, where r_{mi} denotes the residual errors, defined as the difference between the prediction of an individual base classifier and the actual label (Mienye, et al., 2020). These residual errors, also referred to as pseudo-residuals, are computed as:

$$r_{mi} = \left[\frac{\delta L(y_i, F(x))}{\delta F(x)} \right]_{F(x)=F_{m-1}} \quad (2.10)$$

The coefficient ρ_m is determined through line search optimization. However, the technique may experience overfitting if the iterative process is not sufficiently regularized. For certain loss functions, such as the quadratic loss, when h_m fits the residuals perfectly, the residuals in the next iteration become zero, causing the process to stop prematurely (Zhang, et al., 2019). In addition, various regularization hyperparameters have been explored to optimize the additive learning process of GBDT. A common approach to regularization in GBDT is shrinkage, which reduces each gradient descent step as follows.

$$F_m(x) = F_{m-1}(x) + v\rho_m h_m(x) \quad (2.11)$$

Where v is typically set to 0.1 (Jiang, et al., 2020). Like other boosting techniques, gradient boosting's capacity to identify intricate data patterns by fixing previous models' mistakes is one of its key advantages. Nevertheless, models built with this method may overfit and capture noise when the dataset contains noisy data (Zhang,

et al., 2019). Gradient boosting is considered particularly effective for problems involving small datasets.

2.6 Extreme Gradient Boosting (XGBoost)

XGBoost is an extension of gradient boosting, an algorithm that can find optimal solutions to regression and classification problems based on Gradient Boosting Decision Trees (Andriansyah & Fridayanthie, 2023). This algorithm's fundamental idea is to lower the loss function (a mechanism for evaluating the model) by continually adjusting the learning parameters. XGBoost improves performance and lowers model complexity to prevent overfitting by using a more structured model to construct a regression tree structure. Boosting is an ensemble learning algorithm that assigns different weights to the training data distribution in each iteration (Sunata, et al., 2020). Each iteration of the boost adds weight to incorrectly classified samples and reduces weight to correctly classified samples. Boosting combines multiple weak classifiers to create a robust classifier. The traditional Gradient Boosting technique gave rise to XGBoost, a boosting variation. Gradient Boosting was changed into XGBoost to increase generalization performance, scalability, and speed. Data structuring is the first step in using XGBoost. Since XGBoost only works with numeric vectors, all categorical data types must be transformed into numeric forms. One Hot Encoding can be used to complete this conversion. The following step involves cleaning the data and performing feature engineering (Shahri, et al., 2021). The following equation states that the general function can be used to produce the estimated model.

$$\hat{y}_i^{(t)} = \sum_{k=1}^t f_k(x_i) = \hat{y}_i^{(t-1)} + f_t(x_i) \quad (2.12)$$

where,

$$\begin{aligned} \hat{y}_i^{(t)} &= \text{predictions at step } t \\ f_t(x_i) &= \text{the learner at step } t \\ x_i &= \text{the input variable} \\ \hat{y}_i^{(t-1)} &= \text{predictions at step } t - 1 \end{aligned}$$

As the equation below illustrates, the goal is to prevent overfitting in XGBoost:

$$F_{obj}(\theta) = \sum_{i=1}^n l(\hat{y}_i, y_i) + \sum_{k=1}^K \Omega(f_k) \quad (2.13)$$

where,

$l(\hat{y}_v, y_i)$ = a loss function that measure the difference between the prediction \hat{y}_v and the target y_i

$\Omega(f_k) = \gamma T + \frac{1}{2} \lambda \sum_{k=1}^K w_k^2$, a regularized term that penalizes complex models

γ = the min. loss needed to further partition the leaf node

T = the number of leaves in the tree

λ = a regularized parameter to scale the penalty

γT = spanning the tree pruning

w = weight of leaves

Subsequently, the loss function in XGBoost undergoes a second-order expansion, and the final objective is formally expressed in the equation (Krawczyk, 2016).

$$J^{(t)} \approx \sum_{i=1}^n \left[g_i w_{q(x_i)} + \frac{1}{2} (h_i w_{q(x_i)}^2) \right] + \gamma T + \frac{1}{2} \lambda \sum_{k=1}^K w_k^2 \quad (2.14)$$

where,

h_i = the second derivative of the loss function

g_i = the first derivative

The optimal weight of leaf j , w_j can be identified using the following equation,

$$w_j = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda} \quad (2.15)$$

The cumulative loss associated with every leaf node can be expressed through the

subsequent loss function in the equation

$$J^{(t)} \approx \sum_{j=1}^T \left[\left(\sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) (w_j^2) \right] + \gamma T \quad (2.16)$$

Where I_j refers to all data samples in the leaf node j . Accordingly, changes in model performance can be determined from the objective function once a node split occurs in the decision tree (Zhang, et al., 2018).

2.7 Model Evaluation

In order to assess the effectiveness of different machine learning algorithms in predicting and classifying breast cancer, several performance metrics were computed and compared to provide a comprehensive evaluation of the models (Prastyo, et al., 2020). The performance of each algorithm was then determined based on the average values of these metrics (Salvador, 2024), which include

1. Confusion Matrix

An actual and expected classification table that includes true positives, true negatives, false positives, and false negatives is called a confusion matrix. It offers insights into model performance by highlighting specific types of classification errors. However, in high-dimensional problems with many classes, the confusion matrix may be less informative and requires careful interpretation. To address this, metrics such as accuracy, precision, and recall summarize its information into more interpretable forms.

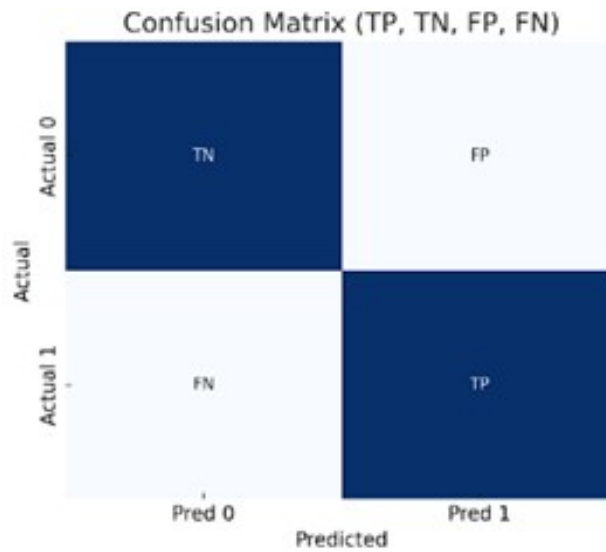


Figure 3 Confusion Matrix (Source: Salvador, 2024)

2. Accuracy Score

The proportion of accurately predicted occurrences to all instances in the dataset is known as the accuracy score. It is computed by dividing the total number of forecasts by the sum of true positives and true negatives.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Where TP represents true positive (i.e., breast cancer patients who are correctly identified as having malignant cancer), TN represents true negative (i.e., breast cancer patients who are correctly identified as having benign cancer), FP represents false positive (i.e., cancer patients who are incorrectly identified as having malignant cancer), and FN represents false negative (i.e., cancer patients who are incorrectly identified as having benign cancer).

3. Precision Score

The percentage of real positive predictions compared to all instances projected as positive is known as precision. It is computed as:

$$\text{Precision} = \frac{TP}{TP + FP}$$

4. Recall

Recall measures the percentage of real positives that the model successfully detects, it is sometimes referred to as sensitivity or true positive rate. It is

defined as:

$$\text{Recall} = \frac{TP}{TP + FN}$$

5. F-1 Score

The F1 score provides a fair assessment by taking into account both metrics. It is the harmonic mean of precision and recall. It is given as:

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

2.8 Breast Cancer

Breast cancer is the most common cause of cancer death in women according to the World Health Organization (Pinheiro & Becker, 2025). It's the most commonly diagnosed cancer in the world and estimate 2.3 million new cases and 685,000 deaths occurred in 2020 (Arnold, et al., 2022). In Brazil, breast cancer is also one of the most common cancers in women, for each year of the 2023 2025 triennium, 73,610 new cases were estimated (Pinheiro & Becker, 2025).

The three main parts of the breast are ducts, connective tissue, and lobules; cancer usually starts in the ducts or lobules. Cells that proliferate abnormally and invade neighboring tissues are known as cancerous tumors. Mammography, ultrasound, MRI, biopsy, clinical breast examination, and genetic testing are a few of the current diagnostic techniques. Breast cancer tumors are classified as benign or malignant, with treatment choice based on grade, stage, and molecular subtype of BC, options are surgery, radiation therapy, chemotherapy, neoadjuvant chemotherapy, and adjuvant chemotherapy (Burguin, et al., 2021).

Studies have shown that increased breast cancer rates are linked to lifestyle risk factors (drinking, being overweight, being inactive), reproductive and hormonal risk factors (early menarche, later menopause, higher age at first childbearing, fewer children, reduced breastfeeding, menopausal hormone therapy, oral contraceptives), and increased mammography (Shang & Xu, 2022). Since 2003, breast cancer cases and deaths among women in the United States and Australia have continued to rise, while the incidence has increased slightly, the mortality rate has decreased (Siegel, et al., 2020).

The incidence of breast cancer is increasing in all three nations, despite the fact that the morbidity and death rates for female breast cancer in the US and Australia are far higher than those in China. The mortality rate of female breast cancer in the United States and Australia declined slightly, while breast cancer mortality in China was growing (Smolarz, et al., 2022). In transitioned countries, the incidence of breast cancer is 88% greater than in transitioning countries (55.9 and 29.7 per 100,000, respectively).

However, compared to women in advanced countries, death rates for women in developing nations are 17% higher (15.0 and 12.8 per 100,000, respectively) (Shang & Xu, 2022). With a 15% mortality-to-incidence ratio, breast cancer makes up around 30% of all female malignancies worldwide. The global incidence ranges from 27 per 100,000 people in Africa and East Asia to 97 per 100,000 people in North America, reflecting the link between breast cancer incidence and economic development, as well as related social and lifestyle factors (Siegel, et al., 2020). Mortality from breast cancer in women can be better reduced by increasing access to aggressive prevention, early detection, and early treatment (Smolarz, et al., 2022).

III. RESEARCH METHODS

3.1 Time and Place of Research

This research was conducted in the even semester of the 2025/2026 academic year at the Department of Mathematics, Faculty of Mathematics and Natural Sciences, University of Lampung.

3.2 Research Data

The data used in this study is secondary data, namely the Wisconsin Breast Cancer Dataset, which can be accessed via the UCI Machine Learning Repository website at the following link: <https://archive.ics.uci.edu/dataset/15/breast+cancer+wisconsin+original>.

3.3 Research Methods

The methodology of this study consists of several stages as follows:

1. Data Analysis
2. Data Preprocessing
 - a. Handling missing values (if any)
 - b. Normalizing or standardizing numerical features to ensure they are on a consistent scale
 - c. Using the Synthetic Minority Oversampling Technique (SMOTE), which creates synthetic samples for the minority class in order to attain a more balanced distribution, to address class imbalance

3. Splitting Data

- a. Both AdaBoost and XGBoost models are trained using the training dataset
- b. Fold CV is applied to validate the consistency of model performance

4. Model Development

- a. AdaBoost Algorithm: Implemented with decision stumps or shallow decision trees as weak learners. The algorithm iteratively adjusts the weights to focus more on misclassified instances
- b. XGBoost Algorithm: Implemented with gradient boosting decision trees. This model incorporates regularization, pruning, and hyperparameter tuning to optimize performance and reduce overfitting

5. Performance Evaluation

- a. The trained models are tested on the testing dataset
- b. Performance is evaluated using classification metrics, including:
 - Confusion Matrix: Provides a detailed overview of prediction results by displaying True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN). This allows a more comprehensive analysis of classification errors
 - Accuracy: The proportion of correct predictions among all cases
 - Precision: The percentage of all anticipated positive cases that were accurately predicted
 - Recall (Sensitivity): The percentage of true positive instances that were accurately recognized
 - F1-Score: The precision and recall harmonic mean
- c. The performance of AdaBoost and XGBoost is then compared based on these metrics

6. Result Analysis

- a. The strengths and weaknesses of AdaBoost and XGBoost in breast cancer classification are analyzed
- b. The discussion emphasizes the practical implications of applying boosting algorithms in medical diagnosis

7. Drawing conclusions regarding the comparative performance of AdaBoost and XGBoost

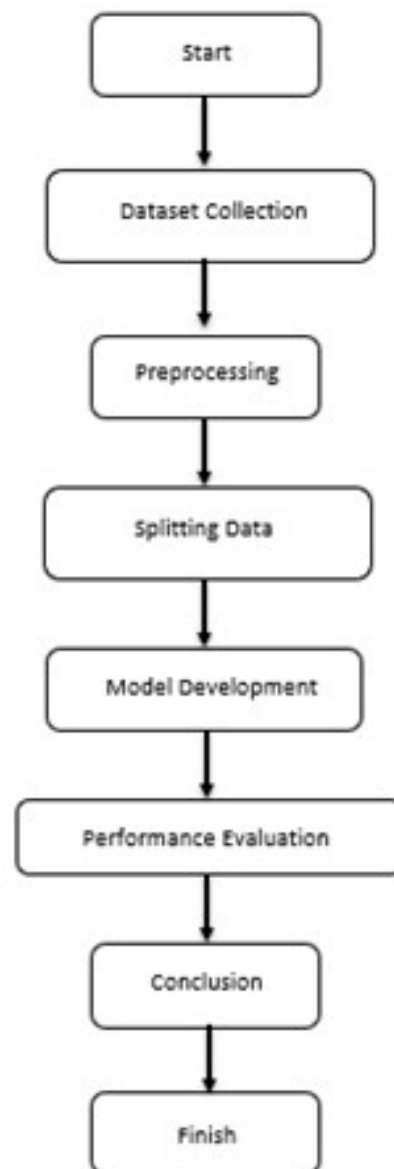


Figure 4 Flowchart Research Methodology

V. CONCLUSION

5.1 Conclusion

Based on the results obtained, several points can be concluded as follows:

- 1 The results show that the AdaBoost and XGBoost algorithms perform very well in classifying breast cancer using the Wisconsin Breast Cancer Dataset (WBCD), with an accuracy value above 96% and high precision, recall, and F1-score values, making both suitable for use as breast cancer classification models.
- 2 Based on the evaluation results, AdaBoost excels in accuracy, precision, and F1-score metrics, while XGBoost excels in recall metrics. The difference in performance between the two models is relatively small, so the choice of algorithm can be adjusted according to priority needs, namely prediction accuracy or completeness of breast cancer case detection.

5.2 Recommendation

Based on the evaluation results, it is recommended to pay more attention to data balancing techniques and hyperparameter optimization so that the model can be more sensitive to minority classes. In addition, further research can explore the use of other models or more sophisticated methods to improve performance in cancer detection, especially in dealing with data imbalance issues that may affect the accuracy and sensitivity of the model.

LITERATURE

- Alelyani, S.(2021). Stable bagging feature selection on medical data. *Journal of Big Data*, **8**(1).
- Alshdaifat, E., Al-hassan, M., & Aloqaily, A. (2021). Effective heterogeneous ensemble classification: An alternative approach for selecting base classifiers. *ICT Express*, **7**(3): 342-349.
- Andriansyah, D., & Fridayanthie E. W. (2023). Optimization of Support Vector Machine and XGBoost Methods Using Feature Selection to Improve Classification Performance. *Journal of Informatics and Telecommunication Engineering*, **6**(2): 484–493.
- Arnold, M., Morgan, E., Rungay, H., Mafra, A., Singh, D., Laversanne, M., Vignat, J., Gralow, J. R., Cardoso, F., Siesling, S., & Soerjomataram, I. (2022). Current and future burden of breast cancer: Global statistics for 2020 and 2040. *Breast*, **24**: 15-23.
- Asri, H., Mousannif, H., Al Moatassime, H., & Noel, T. (2016). Using Machine Learning Algorithms for Breast Cancer Risk Prediction and Diagnosis. *Procedia Computer Science*, **5**: 1064–1069.
- Bentéjac, C., Csörgő, A., & Martínez-Muñoz, G. (2021). A comparative analysis of gradient boosting algorithms. *Artificial Intelligence Review*, **54**(3): 1937–1967.
- Bian, Y., Wang, Y., Yao, Y., & Chen, H. (2020). Ensemble Pruning Based on Objection Maximization with a General Distributed Framework. *IEEE Transactions on Neural Networks and Learning Systems*, **31**(9): 3766–3774.

- Burguin, A., Diorio, C., & Durocher, F. (2021). Breast cancer treatments: Updates and new challenges. *Journal of Personalized Medicine*, **11**(8).
- Chen, T., & Guestrin, C. (2016). A scalable tree boosting system. *REKAYASA: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*: 785–794.
- Freund, Y., & Schapire, R. E. (1997). *Journal of Computer and System Sciences*, **55**.
- Hospedales, T., Antoniou, A., Micaelli, P., & Storkey, A. (2022). Meta-Learning in Neural Networks: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **44**(9): 5149–5169.
- Janiesch, C., Zschech, P., & Heinrich, K. (2021). *Machine learning and deep learning*.
- Jiang, J., Wang, R., Wang, M., Gao, K., Nguyen, D. D., & Wei, G. W. (2020). Boosting Tree-Assisted Multitask Deep Learning for Small Scientific Datasets. *Journal of Chemical Information and Modeling*, **60**(3): 1235–1244.
- Krawczyk, B. (2016). Learning from imbalanced data: open challenges and future directions. *Progress in Artificial Intelligence*, **5**(4): 221–232.
- Liang, M., Chang, T., An, B., Duan, X., Du, L., Wang, X., Miao, J., Xu, L., Gao, X., Zhang, L., Li, J., & Gao, H. (2021). A Stacking Ensemble Learning Framework for Genomic Prediction. *Frontiers in Genetics*, **12**.
- Liu, H., Gegoy, A., & Cocea, M. (2016). Ensemble Learning Approaches. *Studies in Big Data*, **13**: 63–73.
- Mathew, T. E. (2023). Breast Cancer Classification Using an Extreme Gradient Boosting Model with F-Score Feature Selection Technique. *Journal of Advances in Information Technology*, **14**(2): 363–372.

- Mienye, I. D., & Sun, Y. (2022). A Survey of Ensemble Learning: Concepts, Algorithms, Applications, and Prospects. *IEEE Access* , **10**: 99129-99149.
- Mienye, I., D., Sun, Y., & Wang, Z. (2020). Improved Predictive Sparse Decomposition Method with Densenet for Prediction of Lung Cancer. *International Journal of Computing*, **19**(4).
- Miguel-Hurtado, O., Guest, R., Stevenage, S. V., Neil, G. J., & Black, S. (2016). Comparing machine learning classifiers and linear/logistic regression to explore the relationship between hand dimensions and demographic characteristics. *Plos One*, Plos **11**(11).
- Mishra, S., Shaw, K., Mishra, D., Patil, S., Kotecha, K., Kumar, S., & Bajaj, S. (2022). Improving the Accuracy of Ensemble Machine Learning Classification Models Using a Novel Bit-Fusion Algorithm for Healthcare AI Systems. *Frontiers in Public Health*, **10**.
- Pinheiro, J. M. H., & Becker, M. (2025). Breast Cancer Classification Using Gradient Boosting Algorithms Focusing on Reducing the False Negative and SHAP for Explainability. *Inteligencia Artificial*, **28**(75): 63-80.
- Prastyo, P. H., Paramartha, I. G. Y., Pakpahan, M. S. M., & Ardiyanto, I. (2020). Predicting Breast Cancer: A Comparative Analysis of Machine Learning Algorithms, **3**(1): 455-459.
- Salvador, E. L. (2024). *Use of Boosting Algorithms in Household-Level Poverty Measurement: A Machine Learning Approach to Predict and Classify Household Wealth Quintiles in the Philippines*.
- Sarker, I. H. (2021). Machine Learning: Algorithms, Real-World Applications and Research Directions. In *SN Computer Science*, **2**(3). Springer.
- Schapire, R. E. (1990). *The The Strength of Weak Learnability*, **5**.

- Shahri, N. H. N. B. M., Lai, S. B. S., Mohamad, M. B., Rahman, H. A. B. A., & Rambli, A. Bin. (2021). Comparing the performance of adaboost, xgboost, and logistic regression for imbalanced data. *Mathematics and Statistics*, **9**(3): 379–385.
- Shang, C., & Xu, D. (2022). Epidemiology of Breast Cancer. *In Oncologie*, **24**(4): 649–669.
- Siegel, R. L., Miller, K. D., & Jemal, A. (2020). Cancer statistics, 2020. *CA: A Cancer Journal for Clinicians*, **70**(1): 7–30.
- Smolarz, B., Zadrożna Nowak, A., & Romanowicz, H. (2022). Breast Cancer—Epidemiology, Classification, Pathogenesis and Treatment (Review of Literature). *In Cancers*, **14**(10).
- Sun, Y., Li, Z., Li, X., & Zhang, J. (2021). Classifier Selection and Ensemble Model for Multi-class Imbalance Learning in Education Grants Prediction. *Applied Artificial Intelligence*, **35**(4): 290-303.
- Sunata, H., Azrullah, F. J., & Rianto, Y. (2020). Komparasi Tujuh Algoritma Identifikasi Fraud ATM Pada PT. Bank Central Asia Tbk, **7**(3).
- Tian, T., Wu, H., Liu, X., & Hu, Q. (2024). D3AT-LSTM: An Efficient Model for Spatiotemporal Temperature Prediction Based on Attention Mechanisms. *Electronics (Switzerland)*, **13**(20).
- Wang, F., Li, Z., He, F., Wang, R., Yu, W., & Nie, F. (2019). Feature Learning Viewpoint of Adaboost and a New Algorithm. *IEEE Access*, **14**: 149890-149899.
- Xiao, Y., Xia, J., Li, L., Ke, Y., Cheng, J., Xie, Y., Chu, W., Cheung, P., Kim, J. H., Colditz, G. A., Tamimi, R. M., & Su, X. (2019). Associations between dietary patterns and the risk of breast cancer: A systematic review and meta-analysis of observational studies. *In Breast Cancer Research*, **21**(1).

Zhang, B., Ren, J., Cheng, Y., Wang, B., & Wei, Z. (2019). Health data driven on continuous blood pressure prediction based on gradient boosting decision tree algorithm. *IEEE Access*, **7**: 32423-32433.

Zhang, D., Qian, L., Mao, B., Huang, C., Huang, B., & Si, Y. (2018). A Data-Driven Design for Fault Detection of Wind Turbines Using Random Forests and XGboost. *IEEE Access*, **6**: 21020-21031.