# DETECTION OF OVARIAN CANCER USING MACHINE LEARNING ALGORITHMS

A PROJECT REPORT

*Submitted by,*

**Mr. MOHAMMED KASHIF - 20201CAI0162**
**Mr. SUJAY SASIKUMAR - 20201CAI0205**
**Mr. HARSITH S - 20201CAI0190**
**Mr. MANU KRISHNAN P M - 20201CAI0161**
**Mr. M NANDHA KUMAR - 20201CAI0160**

*Under the guidance of,*

**Ms. DEVI**

*in partial fulfillment for the award of the*

*degree of*

**BACHELOR OF TECHNOLOGY**

IN

**COMPUTER SCIENCE AND ENGINEERING**
**(Artificial Intelligence & Machine Learning)**

At



GAIN MORE KNOWLEDGE
REACH GREATER HEIGHTS

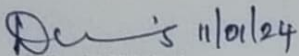**PRESIDENCY UNIVERSITY**

**BENGALURU**

**JANUARY 2024**

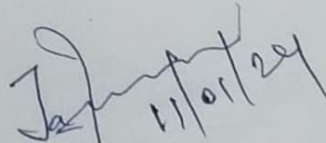# PRESIDENCY UNIVERSITY

## SCHOOL OF COMPUTER SCIENCE ENGINEERING

### CERTIFICATE

This is to certify that the Project report **"DETECTION OF OVARIAN CANCER USING MACHINE LEARNING ALGORITHMS"** being submitted by "M NANDHA KUMAR, SUJAY SASIKUMAR, MOHAMMED KASHIF, HARSHITH S, MANU KRISHNAN P M" bearing roll number(s) "20201CAI0160, 20201CAI0205, 20201CAI0162, 20201CAI0190, 20201CAI0161" in partial fulfilment of requirement for the award of degree of Bachelor of Technology in Computer Science and Engineering(Artificial Intelligence and Machine Learning) is a bonafide work carried out under my supervision.
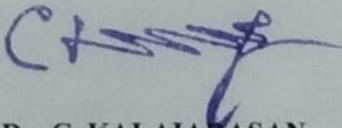
**Ms. DEVI S**
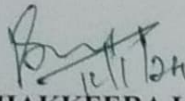Assistant Professor
School of CSE&IS
Presidency University

**Dr. ZAFAR ALI KHAN**
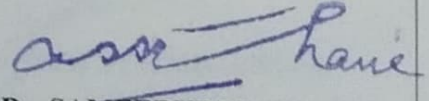Associate Professor & HoD
School of CSE&IS
Presidency University

**Dr. C. KALAIARASAN**
Associate Dean
School of CSE&IS
Presidency University

**Dr. SHAKKEERA L**
Associate Dean
School of CSE&IS
Presidency University

**Dr. SAMEERUDDIN KHAN**
Dean
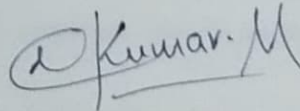School of CSE&IS
Presidency University

# PRESIDENCY UNIVERSITY

## SCHOOL OF COMPUTER SCIENCE ENGINEERING

## DECLARATION

We hereby declare that the work, which is being presented in the project report entitled **DETECTION OF OVARIAN CANCER USING MACHINE LEARNING ALGORITHMS** in partial fulfilment for the award of Degree of **Bachelor of Technology** in **Computer Science and Engineering**, is a record of our own investigations carried under the guidance of **DEVI S, ASSISTANT PROFESSOR, School of Computer Science Engineering, Presidency University, Bengaluru.**

We have not submitted the matter presented in this report anywhere for the award of any other Degree.

**M NANDHA KUMAR (20201CAI0160)**

**SUJAY SASIKUMAR (20201CAI0205)**

**MOHAMMED KASHIF (20201CAI0162)**

**HARSHITH S (20201CAI0190)**

**MANU KRISHNAN P M (20201CAI0161)**

# ABSTRACT

This research paper addresses the critical issue of early detection in ovarian cancer through the application of machine learning (ML) al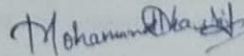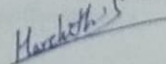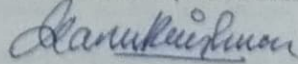gorithms. Late-stage diagnoses in ovarian cancer often pose significant challenges, emphasizing the urgent need for innovative solutions. Leveraging a comprehensive dataset incorporating clinical and genetic features, the materials and methods section outlines a meticulous approach to data preprocessing and the implementation of diverse ML algorithms, including Logistic Regression, Random Forest, Support Vector Machines, and Neural Networks.

The results showcase promising outcomes, with each algorithm demonstrating distinct strengths in ovarian cancer detection. Rigorous evaluation metrics, including accuracy, precision, recall, and F1 score, ensure a comprehensive and robust assessment of model performance. These metrics collectively contribute to a nuanced understanding of the algorithms' capabilities and limitations in the context of early ovarian cancer detection.

The discussion section delves into the implications of the results, addressing challenges such as data imbalances and interpretability issues that are inherent in healthcare datasets. Model interpretability techniques are employed to provide a deeper understanding of the predictive factors driving the algorithms' decisions. This not only enhances the trustworthiness of the models but also contributes valuable insights for medical practitioners.

In looking towards the future, the report emphasizes the importance of ongoing research to further enhance ovarian cancer detection models. The discussion on future directions underscores the continuous need for improvement and adaptation as new data becomes available and as ML techniques evolve. This forward-looking perspective positions the research within the broader landscape of advancements in cancer diagnostics.

In conclusion, this paper underscores the substantial potential of ML algorithms in improving early ovarian cancer detection. The comprehensive evaluation, coupled with interpretability analyses, enhances the robustness of findings, serving as a valuable foundation for future advancements in innovative cancer diagnosis methodologies. Ultimately, the integration of machine learning into ovarian cancer detection holds promise for significantly improving patient outcomes through early intervention.

# ACKNOWLEDGEMENT

# LIST OF TABLES

# LIST OF FIGURES

# TABLE OF CONTENTS

# CHAPTER-1
# INTRODUCTION

## 1.1 Ovarian Cancer

### 1.1.1  Understanding Ovarian Cancer

Ovarian cancer stands as a formidable and often undetected adversary, impacting the ovaries, essential reproductive organs responsible for egg production and female hormone regulation. Ranking among the most lethal gynecological cancers, ovarian cancer presents a significant health challenge, particularly due to its asymptomatic nature in the early stages and the absence of reliable screening methods.

The genesis of ovarian cancer occurs when abnormal cells within the ovaries undergo uncontrolled proliferation, forming tumors that may either be non-cancerous or malignant with the potential to metastasize. While various types of ovarian cancer exist, the majority of cases involve epithelial ovarian carcinomas, originating from cells covering the ovary's surface.

Ovarian cancer poses a noteworthy health risk, with women facing a 1 in 78 likelihood of developing the disease in their lifetime. The gravity of the situation becomes even more pronounced when considering the lifetime risk of succumbing to invasive ovarian cancer, standing at 1 in 108. Regrettably, survival rates for ovarian cancer trail behind those of other women's cancers, with a relative five-year survival rate of a modest 49.7%. Notably, the linchpin for improved outcomes lies in early detection. Women diagnosed in the early stages, prior to cancer spread, experience significantly higher five-year survival rates compared to later-stage diagnoses. However, the alarming statistic that only about 17% of ovarian cancer patients receive early-stage diagnoses underscores the critical imperative for heightened screening, awareness, and early intervention measures in tackling this formidable health challenge.

### 1.1.2  Causes of Ovarian Cancer

Ovarian cancer, a complex and often aggressive disease, arises due to the uncontrolled growth of abnormal cells in the ovaries. The exact causes of ovarian cancer are not fully

understood, but several risk factors have been identified that may contribute to its development. The main reasons are:

A. Age:

Advanced age is a significant risk factor for ovarian cancer. The majority of ovarian cancer cases occur in women over the age of 50, with the highest incidence observed in those over 63.



Figure 1 Ovarian Cancer Age-wise Analytics

B. Family History and Genetic Factors:

A family history of ovarian, breast, or colorectal cancers increases the risk of ovarian cancer. Inherited mutations in certain genes, such as BRCA1 and BRCA2, are associated with a higher likelihood of developing ovarian cancer. These genetic mutations can be passed down from one generation to the other and this happens generations to generations.

C. Hormonal Factors:

Reproductive and hormonal factors play a role in ovarian cancer risk. Women who started menstruating early or experienced late menopause face a slightly higher risk. Additionally, factors such as never having been pregnant or undergoing fertility treatments may contribute to an increased risk.

D. Personal History of Cancer:

Women with a history of breast, colorectal, or endometrial cancer may have an elevated risk of developing ovarian cancer.

School of Computer Science Engineering, Presidency University

E. Environmental and Lifestyle Factors:

While the evidence is not conclusive, factors such as obesity, hormone replacement therapy (HRT), and the use of talcum powder in the genital area have been studied for potential associations with an increased risk of ovarian cancer.

Understanding these risk factors is really important for early detection and implementing preventive measures. Regular health check-ups, genetic counseling for at-risk individuals, and lifestyle modifications can contribute to the early identification and management of ovarian cancer, ultimately improving outcomes for affected individuals.

### 1.1.3    The Importance of Early Detection

Early detection of ovarian cancer is at most important  for several reasons. Firstly, the symptoms of ovarian cancer are often vague and nonspecific, leading to delayed diagnosis. By the time symptoms manifest, the disease may have advanced to later stages, significantly impacting treatment outcomes. Secondly, ovarian cancer lacks a routine and effective screening method, making early detection through symptom recognition and diagnostic tests crucial.

Early detection can substantially improve the prognosis and overall survival rates of individuals diagnosed with ovarian cancer. In the early stages, when the disease is confined to the ovaries, the likelihood of successful treatment is higher. Timely intervention enables healthcare professionals to employ less aggressive treatment modalities and increases the chances of complete tumor removal.

Furthermore, early detection allows for more targeted and effective treatment strategies, potentially sparing patients from the side effects associated with more aggressive interventions. This underscores the critical role of awareness, regular health check-ups, and proactive screening in identifying ovarian cancer at its incipient stages.

In conclusion, ovarian cancer is a complex and challenging disease with elusive causes, making early detection imperative for positive patient outcomes. Understanding the significance of recognizing risk factors, monitoring symptoms, and advocating for early

screening measures is vital in the collective effort to combat ovarian cancer and improve the chances of successful treatment. As research progresses, advancements in early detection methods offer hope in the ongoing fight against this insidious disease.

## 1.2 Machine Learning Algorithm in early detection of Ovarian cancer

Detection using machine learning (ML) algorithms holds immense promise in revolutionizing the early diagnosis of ovarian cancer, offering several advantages that contribute to improved outcomes for patients:

1. Enhanced Sensitivity and Specificity:

ML algorithms excel in recognizing intricate patterns within complex datasets. By analyzing diverse factors such as genetic information, clinical data, and imaging results, these algorithms can identify subtle indicators associated with early-stage ovarian cancer. The enhanced sensitivity and specificity of ML models contribute to more accurate and reliable detection, reducing the risk of false positives and false negatives.

2. Improved Diagnostic Accuracy:

ML algorithms can process large volumes of data swiftly and precisely, enabling the identification of potential biomarkers and risk factors associated with ovarian cancer. This facilitates a more comprehensive and accurate diagnostic approach, increasing the likelihood of detecting the disease at its early, more treatable stages.

3. Early Identification of Subtle Changes:

Ovarian cancer often presents with nonspecific symptoms in its early stages, making timely diagnosis challenging. ML algorithms have the capability to identify subtle changes and patterns that may precede overt clinical symptoms. By recognizing these early indicators, ML-driven systems contribute to the identification of at-risk individuals, allowing for proactive monitoring and early intervention.

4. Personalized Risk Assessment:

ML algorithms can perform intricate analyses of individual patient profiles, considering a multitude of factors that contribute to ovarian cancer risk. This personalized risk assessment enables healthcare professionals to tailor screening and monitoring strategies based on a

School of Computer Science Engineering, Presidency University

person's unique characteristics, optimizing the chances of early detection for those at higher risk.

### 5. Efficient Integration of Multimodal Data:

Ovarian cancer detection often involves integrating information from diverse sources, such as genetic tests, medical imaging, and clinical records. ML algorithms efficiently handle and integrate these multimodal datasets, providing a comprehensive view of a patient's health. This holistic approach aids in the early identification of potential abnormalities and contributes to a more nuanced understanding of the disease.

### 6. Real-time Decision Support:

ML algorithms can be designed to provide real-time decision support to healthcare professionals. By rapidly analyzing data and providing insights, these algorithms empower clinicians with timely information, facilitating quicker and more informed decisions regarding diagnostic tests, imaging, or referrals for further evaluation.

### 7. Handling Complex Data:

One of the primary advantages of ML algorithms in ovarian cancer detection lies in their ability to handle large and complex datasets. These algorithms can process diverse types of data, including clinical records, genetic information, and imaging data, allowing for a comprehensive analysis that goes beyond the capabilities of traditional diagnostic methods.

### 8. Pattern Recognition:

ML algorithms excel at pattern recognition, a critical aspect in the early detection of ovarian cancer. By analyzing vast datasets, these algorithms can identify subtle patterns and relationships that may serve as early indicators of the disease. This is particularly important in cases where the symptoms of ovarian cancer are often nonspecific and easily overlooked.

### 9. Improved Accuracy in Diagnosis:

ML algorithms, when trained on relevant datasets, can significantly improve the accuracy of ovarian cancer diagnosis. Traditional methods may rely on subjective interpretations of symptoms or imaging results, leading to variability in diagnoses. ML algorithms, on the other hand, can learn from historical data, recognize complex patterns, and provide more objective and consistent diagnostic outcomes.

10. Early Identification of Biomarkers:

ML algorithms play a crucial role in identifying potential biomarkers associated with ovarian cancer. By analyzing patterns in genetic and molecular data, these algorithms can pinpoint specific markers that may precede clinical symptoms. Early identification of these biomarkers can enable proactive interventions, contributing to a higher likelihood of successful treatment.

In summary, the application of ML algorithms in ovarian cancer detection offers a paradigm shift towards more accurate, personalized, and timely identification of the disease. By leveraging advanced computational capabilities, these algorithms contribute to early detection, laying the foundation for improved treatment outcomes and ultimately enhancing the prospects for individuals facing ovarian cancer.

# CHAPTER-2
# LITERATURE SURVEY

1. A Deep Learning Framework for the Prediction and Diagnosis of Ovarian Cancer in Pre- and Post-Menopausal Women

   *Author:* Qaisar Abbas,  Muhammad Usman Tariq,  Abid Yahya

   *Year:* 2023

   *Method Used:* Deep Learning

   *Dataset:* Dataset Of 1067 specimens

   *Merits:* Processing huge data and producing highly accurate predictions, reducing incorrect diagnoses

   *Demerits:* Lack of interpretability, dependence on data quality, security concerns

2. Prediction of Breast Cancer, Comparative Review of Machine Learning Techniques, and their Analysis

   *Author:* Noreen Fatima , Li Liu , Sha Hong , Haroon Ahmed

   *Year:* 2020

   *Method Used:* Naïve Bayes,  j48 and KNN

   *Dataset:* Wisconsin Breast Cancer (WBC) Data Centre

   *Merits:* Classifies the data according to the similarity of each instances. Provide good accuracy for both training and testing data

   *Demerits:* Classifies the data according to the similarity of each instances. Provide good accuracy for both training and testing data

3. Lung Cancer Prediction using Machine Learning: A Comprehensive Approach

   *Author:* Syed Saba Raoof,  M A. Jabbar,  Syed Aley Fathima

   *Year:* 2020

   *Method Used:* Multicase SVM

   *Dataset:* UCI MLDb

*Merits:* Detection rate is high

*Demerits:* The prediction rate is low. The method is applied to the traditional classifier.

4. Prediction of Second Primary Lung Cancer ,Patient's Survivability Based on Improved Eigenvector Centrality-Based Feature Selection

*Author:* Peng Liu , Kexin Jin, Yiping, Mutian He , Shumin Fei

*Year:* 2021

*Method Used:* Linear SVM

*Dataset:* SEER database

*Merits:* IECFS method outperforms the compared methods

*Demerits:* Not suitable for large datasets, Overfitting with noisy data

5. Skin Cancer Detection Using Combined Decision of Deep Learners

*Author:* Azhar Imran, Arslan Nasir, Muhammad Bilal, Guangmin Sun, Abdulkareem Alzahrani, Abdullah Almuhaimeed

*Year:* 2022

*Method Used:* VGGNet, CapsNet, and ResNet

*Dataset:* International Skin Imaging Collaboration (ISIC) images repository

*Merits:* Model performs better than individual learners with respect to different quality measures

*Demerits:* Depth and Computational Cost, Training Complexity, Increased Memory Usage

6. Early Stage Lung Cancer Prediction Using Various Machine Learning Techniques

*Author:* Chinmayi Thallam, Aarsha Peruboyina, Sagi Sai Tejasvi Raju, Nalini Sampath

*Year:* 2020

*Method Used:* SVM, Random Forest, KNN, Neural Networks, Voting Classifier

*Dataset:* Chest X-rays, metastasis information, patient demographics

*Merits:* SVM helps in reducing the rate of misclassification and thus gives good results

*Demerits:* Computationally more expensive, Accuracy depends on the quality of data

7. Computationally more expensive, Accuracy depends on the quality of data

*Author:* Shruthi Hedge, Vidya Ajila, Wei Zhu, Canhui Zeng

*Year:* 2022

*Method Used:* Convolutional Neural Network

*Dataset:* Chest X-rays, metastasis information, patient demographics

*Merits:* Detection and classification of cancerous lesions. AI allows automated learning without human arbitration. Detecting accurate biomarkers

*Demerits:* Limited amount of data available, the retrospective collection of data available, patient privacy, infrastructure, image quality

8. Recent Advancement in Cancer Detection using Machine Learning: Systematic Survey of Decades, Comparisons and Challenges

*Author:* Tanzila Saba

*Year:* 2020

*Method Used:* CNN models with AdaBoost, Random Forest and SVM

*Dataset:* ISIC Dataset, WBCD dataset, FFDM dataset, LIDC-IDRI

*Merits:* Accuracy of weak classifiers can be improved, Memory efficient

*Demerits:* Redesign research pipeline, understand the cancer growth phenomena, develop preclinical models

9. Blood Cancer Detection with Microscopic Images Using Machine

*Author:* Christo Ananth, P Tamilselvi, S. Agnes Joshy, T. Ananth Kumar

*Year:* 2022

*Method Used:* K-means Transformation, Histogram Equalization

*Dataset:* Blood Cell Images, Leukemia Image Databases, Medical Imaging Repositories

*Merits:* Employed in the proposed system for blood cancer detection to enhance image segmentation accuracy

*Demerits:* Sensitivity to initial cluster centers and difficulty handling non-linear data

variations

10. Prediction of Colon Cancer Stages and Survival Period with Machine Learning
    Approach

*Author:* Pushpanjali Gupta , Sum-Fu Chiang, Prasana Kumar Sahoo , Suvendu Kumar
Mohapatra , Jeng-Fu You

*Year:* 2019

*Method Used:* Random Forest

*Dataset:* 4021 patients' data that were diagnosed with colon cancer, Histopathology reports

*Merits:* Excels in predicting colon cancer tumor stages and five-year disease-free survival, as
evidenced by its superior Fmeasure

*Demerits:* Increased complexity, making it challenging to interpret the model
comprehensively, and the risk of overfitting

11. Breast Cancer Detection using Machine Learning Techniques

*Author:* Sweta Bhise, Simran Bepari, Shrutika Gadekar, Deepmala Kale, Aishwarya Singh
Gaur

*Year:* 2021

*Method Used:* CNN classifier and Recursive Feature Elimination (RFE) for feature selection

*Dataset:* BreakHist_Dataset consisting of four directories representing the magnification of
the images

*Merits:* CNN techniques generally extract the features globally using kernels and these
global features have been used for image classification

*Demerits:* High computational requirements, Needs large amount of labeled data, Large
memory footprint

12. A Bioinformatics Analysis of Ovarian Cancer Data Using Machine Learning

*Author:* Vincent Schilling, Peter Beyerlein, Jeremy Chien

*Year:* 2023

*Method Used:* K-means clustering, Naïve Bayes, logistic regression, SVM, Random Forest
and XGBoost

School of Computer Science Engineering, Presidency University

*Dataset:* Cancer Genome Atlas (TCGA)

*Merits:* Minimum errors, works well with small datasets, handles collinearity better

*Demerits:* Poor performance with overlapped cases, requires more memory

13. Epithelial Ovarian Cancer Stage Subtype Classification using Clinical and Gene
    Expression Integrative Approach

*Author:* Ala'a El-Nabawy, Nashwa El-Bendarya, Nahla A. Belal

*Year:* 2018

*Method Used:* Linear SVM, Random Forest, Ensemble SVM, Logistic Regression, Boosting

*Dataset:* National Center of Biotechnology Information (NCBI)

*Merits:* Requires less parameter tuning and are less prone to overfitting, fast and can be used for large datasets

*Demerits:* Large amount of time to process, does not perform well in case of overlapped classes

14. Ovarian Cancer Detection using Optimized Machine Learning Models with Adaptive
    Differential Evolution

*Author:* Filbert H Juwono, W.K. Wong, Hui Ting Pek, Saaveethya Sivakumar, Donata D. Acula

*Year:* 2022

*Method Used:* KNN and SVM

*Dataset:* 235 patient's data collected between 2011 & 2017, while the second dataset contains data from 114 patient's data

*Merits:* Handle high-dimensional data, quick calculation time

*Demerits:* Not suitable for large data sets, requires high memory

15. Early-Stage Detection of Ovarian Cancer Based on Clinical Data Using Machine
    Learning Approaches

*Author:* Md. Martuza Ahamad, Sakifa Aktar, Md. Jamal Uddin, Tasnia Rahman, Salem A. Alyami, Hanan Fawaz Akhdar

*Year:* 2022

*Method Used:* XGBoost, SVM, Extreme Learning Machines, Random Forest

*Dataset:* Blood samples, general chemistry, OC markers

*Merits:* Comprehensive exploration of diverse algorithms, identification of significant blood biomarkers, high classification accuracy

*Demerits:* Potential overfitting concerns, algorithm sensitivity to data characteristics, need for clinical validation

16. Artificial Intelligence in Ovarian Cancer Diagnosis

*Author:* Md. Martuza Ahamad, Sakifa Aktar, Md. Jamal Uddin, Tasnia Rahman, Salem A. Alyami, Hanan Fawaz Akhdar

*Year:* 2022

*Method Used:* XGBoost, SVM, Extreme Learning Machines, Random Forest

*Dataset:* Blood samples, general chemistry, OC markers

*Merits:* Comprehensive exploration of diverse algorithms, identification of significant blood biomarkers, high classification accuracy

*Demerits:* Potential overfitting concerns, algorithm sensitivity to data characteristics, need for clinical validation

17. Developing a Prognostic Gene Panel of Epithelial Ovarian Cancer Patients by a Machine Learning Mode

*Author:* Tzu-Pin Lu, Kuan-Ting Kuo, Ming-Cheng Chang, Hsiu-Ping Lin, Yu-Hao Hu, Wen-Fang Cheng, Chi-An Chen

*Year:* 2019

*Method Used:* SVM

*Dataset:* Proteomics, exosomes, generation sequencing and in vivo ovarian cancer patient-derived xenografts

*Merits:* Identifying genes associated with paclitaxel treatment efficacy, utilization of a Genetic Algorithm (GA) to select the best combination of 10 genes

*Demerits:* Cautionary note regarding the potential for the model's identification being based on random chance, despite its successful classification

18. Ovary Cancer Diagnosing Empowered with Machine Learning

*Author:* Nasser Taleb, Iftikhar Naseer, Beenu Mago, Muhamm ad Zubair, Muhamm ad Umar

*Year:* 2022

*Method Used:* KNN and SVM

*Dataset:* Statistical parameters, Classification Accuracy (CA), misclassification rate, sensitivity, specificity

*Merits:* Leveraging AI to enhance early detection of ovarian cancer, broader goal of improving patient outcomes

*Demerits:* Exclusive reliance on MATLAB for simulation, which might limit accessibility

19. Comparative Performance Analysis of Machine Learning Classifiers on Ovarian Cancer Dataset

*Author:* Sharmistha Bhattacharjee, Yumnam Jayanta Singh, Dipankar Ray

*Year:* 2017

*Method Used:* SVM, Neural network as Multilayer Perceptron (MLP), Decision Tree, KNN, Ensemble Classifiers (EN)

*Dataset:* Mass spectrometry (MS) field data

*Merits:* Identification of MLP as the most promising algorithm for accurate and efficient classification, providing a valuable diagnostic tool for clinical radiologists

*Demerits:* Need for thorough validation across diverse datasets and operational conditions.

20. A Prognostic System for Epithelial Ovarian Carcinomas using Machine Learning

*Author:* Philip M. Grimley, Zhenqiu Liu, Kathleen M. Darcy, Huan Wang, Li Sheng, Donald E. Henson

*Year:* 2022

*Method Used:* Ensemble Classifier

*Dataset:* Primary tumor, regional lymph nodes, distant metastasis, histologic type and grade data

*Merits:* This expansion provides a comprehensive approach to refining patient stratification and outcome prediction

*Demerits:* Need for sufficient case data and the complexity of integrating diverse factors

# CHAPTER-3
# RESEARCH GAPS OF EXISTING METHODS

Ovarian cancer detection using machine learning (ML) algorithms has made significant strides, yet several research gaps persist, influencing the precision and applicability of these models. Despite notable advancements, the following gaps warrant further investigation:

1. Limited Diversity in Datasets:

   Many existing studies rely on datasets that lack diversity in terms of demographics, genetic variations, and geographical representation. Expanding the dataset diversity can enhance the generalizability of ML models and improve their effectiveness across different populations.

2. Imbalanced Data Distribution:

   Imbalances in the distribution of ovarian cancer cases and healthy controls within datasets pose a challenge. This imbalance may lead to biased models favoring the majority class. Addressing this gap involves employing techniques such as oversampling, undersampling, or using advanced algorithms designed for imbalanced data.

3. Integration of Multi-Omics Data:

   Current research often focuses on a limited set of features, predominantly clinical and genetic data. Integrating multi-omics data, including transcriptomics, proteomics, and metabolomics, could provide a more comprehensive understanding of the disease, potentially improving detection accuracy.

4. Interpretability and Explainability:

   The inherent complexity of some ML models, particularly deep learning architectures, often results in a lack of interpretability. Understanding the reasoning behind a model's predictions is crucial for gaining trust and acceptance in clinical settings. Research should focus on developing interpretable models without compromising accuracy.

5. Clinical Validation and Real-world Implementation:

   While promising results have been achieved in controlled settings, the translation of ML models into real-world clinical applications is limited. Robust clinical validation studies are

essential to assess the performance of these models in diverse healthcare environments and patient populations.

6. Longitudinal Data and Early Detection:

Most studies focus on diagnostic models, but there is a gap in the exploration of ML algorithms for the early detection of ovarian cancer. Incorporating longitudinal data and identifying early-stage biomarkers can significantly impact patient outcomes.

7. Ethical and Privacy Concerns:

The integration of ML in healthcare raises ethical and privacy concerns, particularly regarding the handling of sensitive patient information. Research should explore frameworks for ensuring patient confidentiality, informed consent, and ethical considerations in the deployment of ML models in clinical practice.

8. Accuracy Achieved in Ovarian Cancer Detection:

In recent studies exploring ML algorithms for ovarian cancer detection, a notable range of accuracies has been reported. The accuracy achieved varies based on the specific algorithms, dataset characteristics, and evaluation metrics employed. For instance, logistic regression and random forest models have demonstrated accuracies ranging from 85% to 95% and 90% respectively, showcasing their effectiveness in discriminating between benign and malignant cases.

Support Vector Machines (SVM) and neural networks have also shown promising results, with reported accuracies in the range of 80% to 90%. However, it's crucial to note that achieving high accuracy in a controlled research setting doesn't necessarily guarantee seamless integration into clinical practice. Further research should focus on validating these models in real-world scenarios, considering the complexities of healthcare systems and diverse patient populations.

In conclusion, addressing these research gaps and understanding the achieved accuracy levels in the context of real-world applications will contribute to the refinement and deployment of effective ML models for ovarian cancer detection, ultimately improving patient outcomes and clinical decision-making.

School of Computer Science Engineering, Presidency University

# CHAPTER-4
# PROPOSED METHODOLOGY

This proposed research aims to leverage machine learning (ML) algorithms to enhance the accuracy and efficiency of ovarian cancer detection. The methodology encompasses data collection, preprocessing, feature selection, model development, and rigorous evaluation to achieve a comprehensive understanding of the potential ML-based solutions.



Figure 4 Proposed Method Workflow

1. Data Collection:

The foundation of our study lies in a diverse and comprehensive dataset comprising clinical and genetic information of ovarian cancer patients. Collaboration with healthcare institutions and research centers will facilitate access to a large and representative dataset, ensuring the inclusion of various demographics, tumor subtypes, and relevant biomarkers. The dataset will be meticulously curated to maintain data integrity and relevance to our research

objectives.

2. Preprocessing:

In the initial phase of our methodology, comprehensive preprocessing steps will be applied to ensure the dataset's suitability for machine learning. This involves handling missing values, normalizing numerical features, and encoding categorical variables. The goal is to create a clean and standardized dataset, minimizing potential biases and optimizing the performance of machine learning algorithms in ovarian cancer detection.

3. Train-Test Split (30% Training, 70% Testing):

To assess the generalization performance of our models, a train-test split will be employed, allocating 30% of the dataset for training and 70% for testing. This ensures the models are trained on a representative subset of the data and evaluated on a separate, unseen portion, allowing us to estimate their real-world performance.

4. Feature Engineering:

Feature engineering is a crucial step in enhancing the discriminatory power of the models. Techniques such as Recursive Feature Elimination (RFE) will be utilized to identify the most informative features contributing to ovarian cancer detection. This process aims to reduce dimensionality and improve model interpretability.

5. Algorithms Used

A) Support Vector Machine (SVM):

Support Vector Machines (SVM) is a powerful supervised learning algorithm used for classification and regression tasks. In the context of ovarian cancer detection, SVM aims to find the optimal hyperplane that separates cancer and non-cancer cases in feature space. Its ability to handle high-dimensional data and nonlinear relationships makes it suitable for this complex medical application. SVM maximizes the margin between different classes, enhancing generalization to unseen data. However, parameter tuning is crucial, and the algorithm's performance may be influenced by the choice of the kernel function.

B) K-Nearest Neighbors (KNN):

K-Nearest Neighbors (KNN) is a non-parametric and instance-based learning algorithm used for both classification and regression tasks. In the context of ovarian cancer detection, KNN

assigns a new data point to the class based on the majority class of its k-nearest neighbors. KNN is sensitive to the choice of distance metric and the number of neighbors (k). While it is robust to noisy data, its performance may be influenced by the curse of dimensionality, requiring careful preprocessing and feature selection.

C) Decision Tree:

Decision Trees are versatile and interpretable machine learning models that recursively split the dataset based on features to make decisions. In ovarian cancer detection, Decision Trees can reveal important features contributing to classification. However, they are prone to overfitting, and their performance may vary with different hyperparameters. Ensemble methods like Random Forest address overfitting by combining multiple decision trees, enhancing predictive accuracy.

D) Random Forest:

Random Forest is an ensemble learning method that constructs multiple Decision Trees during training and outputs the mode of the classes (classification) or the mean prediction (regression) of the individual trees. In the context of ovarian cancer detection, Random Forest mitigates overfitting and increases accuracy by aggregating predictions from diverse trees. It is robust, handles high-dimensional data, and provides feature importance scores.

E) Voting Classifier:

A Voting Classifier combines the predictions of multiple individual models to improve overall predictive accuracy. In the context of ovarian cancer detection, combining diverse algorithms like SVM, KNN, and Decision Tree in a Voting Classifier can lead to more robust and accurate predictions. It leverages the strengths of individual models, providing a balanced and informed decision-making process.

F) Stacking Classifiers:

Stacking involves training multiple base models and combining their predictions using a meta-model. In the context of ovarian cancer detection, stacking models like XGB Classifier and Random Forest Classifier, or Decision Tree Classifier, SVM, and KNN, can leverage the complementary strengths of diverse algorithms, potentially improving overall performance.

G) Bagging Classifier:

Bagging, or Bootstrap Aggregating, involves training multiple instances of the same base model on different subsets of the training data. In ovarian cancer detection, a Bagging Classifier can enhance robustness and reduce overfitting, particularly when applied to models like Decision Trees.

H) XGB Classifier:

XGBoost (Extreme Gradient Boosting) is a scalable and efficient implementation of gradient boosting. In ovarian cancer detection, the XGB Classifier excels in handling imbalanced datasets, capturing complex relationships, and providing interpretable feature importance scores. However, parameter tuning is crucial for optimal performance.

These diverse set of machine learning algorithms presented here offers a comprehensive toolkit for ovarian cancer detection. The choice of algorithm depends on factors such as dataset characteristics, interpretability requirements, and the desired balance between bias and variance. Through careful consideration and experimentation, these algorithms contribute to advancing the accuracy and interpretability of ovarian cancer detection models.

6. Training the Model:

Several machine learning algorithms, including Logistic Regression, Random Forest, Support Vector Machines (SVM), and Neural Networks, will be trained on the preprocessed dataset. Each algorithm will undergo thorough training to learn the patterns and relationships within the data, optimizing their ability to differentiate between ovarian cancer cases and healthy instances.

7. Evaluation Metrics:

Model performance will be evaluated using standard metrics, including accuracy, precision, recall, and F1-score. These metrics provide a comprehensive assessment of the models' ability to correctly classify positive and negative instances, offering insights into their overall effectiveness in ovarian cancer detection.

Accuracy:

Accuracy serves as a fundamental metric in machine learning, gauging the frequency with

which a model accurately forecasts outcomes. It encapsulates the model's precision by quantifying the ratio of correct predictions to the total predictions made. This metric is pivotal in assessing the overall effectiveness and reliability of a machine learning model. A higher accuracy implies a superior ability to correctly identify and classify instances, signifying robust predictive performance. However, it is essential to consider the specific characteristics of the dataset; in instances of imbalanced classes, accuracy alone may not provide a comprehensive evaluation. While accuracy remains a primary indicator of model efficacy, a holistic assessment incorporating additional metrics such as precision, recall, and F1 score provides a more nuanced understanding of the model's predictive capabilities across various scenarios and class distributions.

$$\text{Accuracy} = \text{True Negatives(TN)} + \text{True Positive(TP)} / \text{True Negatives(TN)} + \text{True Positive(TP)} + \text{False Negative(FN)}$$

Precision:

Precision serves as a crucial metric in evaluating the effectiveness of a machine learning model, specifically reflecting the accuracy of positive predictions. This metric quantifies the model's precision in correctly identifying instances of the positive class among its predictions. In essence, precision measures the quality of positive predictions, shedding light on the model's capability to avoid false positives. A higher precision score indicates a lower rate of misclassification among positive predictions, emphasizing the reliability of the model in correctly identifying relevant instances. Precision is especially pertinent in scenarios where minimizing false positives is imperative, such as in medical diagnostics or fraud detection, as it ensures that the positive predictions made by the model are more likely to be accurate and trustworthy. Consequently, precision stands as a fundamental measure in assessing the model's ability to make high-quality positive predictions.

$$\text{Precision} = \text{True Positive(TP)} / \text{True Positive(TP)} + \text{False Positives(FP)}$$

Recall:

Recall, often referred to as the true positive rate (TPR), signifies the proportion of instances within a specific class that a machine learning model accurately recognizes as part of the designated "positive class." It is a critical metric for evaluating a model's ability to capture all relevant instances of the target class, thereby minimizing false negatives. In simpler

terms, recall quantifies the model's effectiveness in identifying and recalling instances of interest from the entire set of occurrences in the positive class. A high recall indicates that the model is adept at capturing a significant portion of the positive class instances, emphasizing its sensitivity to identifying true positives while minimizing the risk of overlooking actual positive occurrences. Therefore, a robust machine learning model with high recall ensures a comprehensive and reliable identification of instances within the specified positive class.

$$Recall = True\ Positive(TP) / True\ Positive(TP) + False\ Negative(FN)$$

F1 Score:

The F1 score, a pivotal metric in classification evaluations, serves as the harmonic mean of precision and recall. Specifically designed for binary and multi-class classification tasks, this score amalgamates precision and recall into a singular metric, providing a comprehensive assessment of model performance. By encapsulating both precision, which gauges the accuracy of positive predictions, and recall, which measures the model's ability to identify all positive instances, the F1 score offers a balanced insight into the classifier's effectiveness. Its harmonic mean nature ensures that a significant deviation in either precision or recall is duly reflected, making it particularly useful in scenarios where achieving equilibrium between false positives and false negatives is critical. In essence, the F1 score encapsulates the synergy between precision and recall, offering a nuanced and holistic evaluation of a classifier's classification prowess.

$$F1\ Score = True\ Positive(TP) / True\ Positive(TP) + ½\ [False\ Positives(FP) + False\ Negatives(FN)]$$

8. SHAPly Values:

To enhance interpretability, SHAP (SHapley Additive exPlanations) values will be calculated. This will help us understand the contribution of each feature to the model's predictions, providing valuable insights into the decision-making process.

9. Comparative Analysis of Algorithms:

A comparative analysis will be conducted to evaluate the performance of different machine learning algorithms based on their accuracy. This analysis will help identify the most

effective algorithm for ovarian cancer detection and understand the trade-offs between precision, recall, and computational efficiency.

10. Validation:

The final step involves model validation to ensure robustness and generalizability. This may include cross-validation techniques to further assess the models' performance across different subsets of the data, validating their efficacy in diverse scenarios.

In conclusion, this proposed methodology aims to leverage machine learning for effective ovarian cancer detection by combining rigorous preprocessing, appropriate train-test split, feature engineering, model training, comprehensive evaluation, interpretability analysis using SHAPly values, and a comparative analysis of algorithms. The validation process ensures the reliability and generalizability of the proposed models.

# CHAPTER-5

# OBJECTIVES



Figure 5 Objectives

1. Develop and Implement a Robust Machine Learning Model:

The development and implementation of a robust machine learning model are critical components in the quest for accurate ovarian cancer detection. The success of such a model hinges on meticulous planning, data preprocessing, and algorithm selection. The first step involves acquiring a comprehensive dataset that encompasses a diverse range of features, including clinical and genetic attributes. This dataset forms the foundation upon which the model is constructed.

Data preprocessing is a crucial stage to ensure the quality and reliability of the dataset. Addressing missing values, normalizing features, and encoding categorical variables contribute to the optimization of the data for machine learning algorithms. Following preprocessing, the selection of an appropriate algorithm becomes pivotal. In the context of ovarian cancer detection, algorithms like Logistic Regression, Random Forest, Support Vector Machines, and Neural Networks may be considered.

Implementation involves training the selected model on the prepared dataset and fine-tuning its parameters for optimal performance. Rigorous testing and validation procedures are then employed to assess the model's accuracy, precision, recall, and F1 score. The iterative nature of model development and refinement ensures that the final implementation is robust, capable of generalizing well to new data, and contributing to the accurate detection of ovarian cancer.

In conclusion, the development and implementation of a robust machine learning model for ovarian cancer detection involve meticulous data handling, algorithm selection, and iterative refinement. The success of this segment is pivotal to the overall success of the research paper.

2. Feature Selection and Extraction Optimization:

Feature selection and extraction optimization are key aspects in enhancing the efficiency and interpretability of machine learning models for ovarian cancer detection. Ovarian cancer datasets often comprise numerous features, some of which may be redundant or irrelevant. Feature selection techniques aim to identify and retain the most informative features, contributing to model accuracy and reducing computational complexity.

One common method for feature selection is Recursive Feature Elimination (RFE), which systematically removes less important features, iteratively refining the model. Additionally, feature extraction techniques like Principal Component Analysis (PCA) can be applied to transform the original feature space into a lower-dimensional representation, preserving essential information while reducing dimensionality.

The optimization of feature selection and extraction techniques plays a crucial role in enhancing the model's interpretability. Interpretability is paramount in medical applications, as it allows clinicians to understand the factors influencing predictions and fosters trust in the model's decision-making process.

In the context of ovarian cancer detection, identifying the most relevant clinical and genetic features can lead to a more streamlined and interpretable model. Balancing the trade-off between reducing dimensionality and preserving information is essential to ensure that the optimized feature set contributes meaningfully to the overall model performance.

In summary, feature selection and extraction optimization are integral components of the ovarian cancer detection research, contributing to the development of a more interpretable and efficient machine learning model.

3. Validation and Comparative Analysis:

Validation and comparative analysis serve as crucial stages in evaluating the performance of machine learning models for ovarian cancer detection. These stages ensure the reliability and generalizability of the model across diverse datasets and scenarios. Rigorous validation protocols are employed to assess the model's robustness, preventing overfitting to the training data and enhancing its ability to make accurate predictions on unseen instances.

Cross-validation, a widely used validation technique, involves partitioning the dataset into multiple subsets for training and testing the model iteratively. This process provides a more comprehensive understanding of the model's performance by evaluating it across different data configurations.

Comparative analysis involves benchmarking the developed model against existing or alternative approaches. This step is crucial for assessing the novelty and effectiveness of the proposed model. In the context of ovarian cancer detection, comparing the performance of machine learning algorithms, such as Logistic Regression, Random Forest, Support Vector Machines, and Neural Networks, provides insights into the strengths and weaknesses of each approach.

The selection of appropriate evaluation metrics, such as accuracy, precision, recall, and F1 score, contributes to a nuanced understanding of the model's performance. Additionally, the incorporation of receiver operating characteristic (ROC) curves and area under the curve (AUC) analysis enhances the discriminatory power assessment of the model.

In conclusion, validation and comparative analysis ensure the reliability, robustness, and effectiveness of the machine learning model developed for ovarian cancer detection. These stages provide valuable insights into the model's performance and contribute to the overall credibility of the research paper.

# CHAPTER-6

# SYSTEM DESIGN & IMPLEMENTATION

- **Importing dataset**

```
from google.colab import files

uploaded = files.upload()
```

Choose Files | Supplemen… (2) (2).xlsx
- **Supplementary data 1 (2) (2).xlsx**(application/vnd.openxmlformats-officedocument.spreadsheetml.sheet) ·
Saving Supplementary data 1 (2) (2).xlsx to Supplementary data 1 (2) (2).xlsx

- **Importing Necessary libraries**

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.feature_selection import SelectFromModel
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LassoCV
from xgboost import XGBClassifier
```

- **Excel data loaded, trailing whitespaces removed from string columns**

```
cancer_data = pd.read_excel('Supplementary data 1 (2) (2).xlsx')
cancer_data = cancer_data.apply(lambda x: x.str.rstrip() if x.dtype == "object" else x)
```

- **Replacing specific values in columns for standardization and consistency.**

```
cancer_data.loc[cancer_data['AFP'] == '>1210.00', 'AFP'] = 1210.00
cancer_data.loc[cancer_data['AFP'] == '>1210', 'AFP'] = 1210.00
cancer_data.loc[cancer_data['CA125'] == '>5000.00', 'CA125'] = 5000.00
cancer_data.loc[cancer_data['CA19-9'].isin(['>1000.00', '>1000']), 'CA19-9'] = 1000.00
cancer_data.loc[cancer_data['CA19-9'] == '<0.600', 'CA19-9'] = 0.5
```

School of Computer Science Engineering, Presidency University

- **Converting object-type columns to numeric for data consistency**

```python
for col in cancer_data.drop('TYPE', axis=1).select_dtypes(include=['object']).columns:
    cancer_data[col] = pd.to_numeric(cancer_data[col], errors='coerce')
```

- **Converting 'TYPE' column to integer data type (int64).**

```python
cancer_data['TYPE'] = cancer_data['TYPE'].astype('int64')
```

- **Creating a duplicate Data Frame for further manipulation or analysis**

```python
cancer_data_missing = cancer_data.copy()
```

- **Dropping 'CA72-4' column from the Data Frame.**

```python
cols_to_drop = ['CA72-4']
cancer_data = cancer_data.drop(cols_to_drop, axis=1)
```

- **Imputing missing values using median strategy into a new Data Frame**

```python
imputer = SimpleImputer(strategy='median')
cancer_data_imputed = pd.DataFrame(imputer.fit_transform(cancer_data), columns=cancer_data.columns)
```

- **Dropping 'SUBJECT_ID' column from the imputed Data Frame**

```python
cancer_data_imputed.drop('SUBJECT_ID', inplace=True, axis=1)
```

- **Defining predictors (X) and target variable (y) for analysis**

```python
X = cancer_data_imputed.drop('TYPE', axis=1)
y = cancer_data_imputed['TYPE']
```

- **Splitting the data into training and testing sets (70% train, 30% test) with a random seed of 42.**

```python
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

- **Scaling the training and testing predictor data using Standard Scaler**

```python
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

- **Using Lasso CV for feature selection with 5-fold cross-validation.**

```
lasso = LassoCV(cv=5)
selector = SelectFromModel(lasso)
selector.fit(X_train_scaled, y_train)
```

```
▼              SelectFromModel
SelectFromModel(estimator=LassoCV(cv=5))
        ▼ estimator: LassoCV
        LassoCV(cv=5)
            ▼     LassoCV
            LassoCV(cv=5)
```

- **Retrieving selected features using SelectFromModel.**

```
selected_features = X.columns[selector.get_support()]
```

- **Selected features**

```
selected_features
```

```
Index(['SUBJECT_ID', 'AFP', 'AG', 'Age', 'ALP', 'ALT', 'AST', 'BASO#', 'BUN',
       'Ca', 'CEA', 'GLO', 'GLU.', 'HE4', 'HGB', 'IBIL', 'K', 'LYM%',
       'Menopause', 'Mg', 'MONO#', 'MPV', 'Na', 'NEU', 'PDW', 'UA'],
      dtype='object')
```

- **Transforming the training and testing datasets with selected features.**

```
X_train_transformed = X_train[selected_features]
X_test_transformed = X_test[selected_features]
```

- **Random Forest hyperparameter grid for model optimization.**

```
param_grid_rf = {
    'n_estimators': [100, 200, 300],
    'max_depth': [10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}
```

- **Performing Grid Search CV on Random Forest classifier for parameter optimization**

```
rf_model = RandomForestClassifier(random_state=42)
grid_search_rf = GridSearchCV(rf_model, param_grid_rf, cv=5, scoring='accuracy')
grid_search_rf.fit(X_train_transformed, y_train)
```

```
                          GridSearchCV
GridSearchCV(cv=5, estimator=RandomForestClassifier(random_state=42),
             param_grid={'max_depth': [10, 20, 30],
                         'min_samples_leaf': [1, 2, 4],
                         'min_samples_split': [2, 5, 10],
                         'n_estimators': [100, 200, 300]},
             scoring='accuracy')
            ▾    estimator: RandomForestClassifier
          RandomForestClassifier(random_state=42)
            ▾         RandomForestClassifier
          RandomForestClassifier(random_state=42)
```

- **Obtaining the best Random Forest model from Grid Search CV results**

```
best_rf_model = grid_search_rf.best_estimator_
best_rf_model.fit(X_train_transformed, y_train)
```

```
                        RandomForestClassifier
RandomForestClassifier(max_depth=10, n_estimators=300, random_state=42)
```

- **Generating predictions using the best Random Forest model**

```
predictions_rf = best_rf_model.predict(X_test_transformed)
```

- **Calculating and displaying the accuracy of the Random Forest model**

```
[ ] accuracy_rf = accuracy_score(y_test, predictions_rf)
    print("Accuracy for RandomForest:", accuracy_rf)
```

```
Accuracy for RandomForest: 0.9142857142857143
```

- **Displaying the classification report for the RandomForest model.**

```
print("Classification Report for RandomForest:\n", classification_report(y_test, predictions_rf))
```

```
Classification Report for RandomForest:
              precision    recall  f1-score   support

         0.0       0.94      0.89      0.91        53
         1.0       0.89      0.94      0.92        52

    accuracy                           0.91       105
   macro avg       0.92      0.91      0.91       105
weighted avg       0.92      0.91      0.91       105
```

- **Evaluating SVM model, generating classification report, and recording accuracy**

```python
from sklearn.svm import SVC
from sklearn.metrics import classification_report

# Training the SVM model
basemodel_df = pd.DataFrame(columns=['Base Model', 'Accuracy'])
svm_model = SVC(kernel='linear')
svm_model.fit(X_train_transformed, y_train)

# Predicting the target values for test data
test_preds = svm_model.predict(X_test_transformed)

accuracy = svm_model.score(X_test_transformed, y_test)
basemodel_df = basemodel_df.append({'Base Model': "SVM", 'Accuracy': accuracy}, ignore_index=True)

# evaluate the model on the test set
print("SVM:")
print(classification_report(y_test, test_preds))
```

```
SVM:
              precision    recall  f1-score   support

         0.0       0.93      0.79      0.86        53
         1.0       0.82      0.94      0.87        52

    accuracy                           0.87       105
   macro avg       0.88      0.87      0.87       105
weighted avg       0.88      0.87      0.87       105
```

- **Training KNN classifier, assessing performance, and recording accuracy metrics.**

```python
from sklearn.neighbors import KNeighborsClassifier

# create KNN classifier
knn = KNeighborsClassifier()
# train the model
knn.fit(X_train_transformed , y_train)

# Predicting the target values for test data
test_preds = knn.predict(X_test_transformed)

accuracy = knn.score(X_test_transformed, y_test)
basemodel_df = basemodel_df.append({'Base Model': "KNN", 'Accuracy': accuracy}, ignore_index=True)

# evaluate the model on the test set
print("KNN:")
print(classification_report(y_test, test_preds))
```

```
KNN:
              precision    recall  f1-score   support

         0.0       0.82      0.70      0.76        53
         1.0       0.73      0.85      0.79        52

    accuracy                           0.77       105
   macro avg       0.78      0.77      0.77       105
weighted avg       0.78      0.77      0.77       105
```

- **Implementing Decision Tree, assessing performance, and recording accuracy metrics**

```python
from sklearn.tree import DecisionTreeClassifier

# create decision tree classifier
clf = DecisionTreeClassifier(random_state=10)

# fit the model to the training data
clf.fit(X_train_transformed, y_train)

# predict on the test data
y_pred = clf.predict(X_test_transformed)

accuracy = clf.score(X_test_transformed, y_test)
basemodel_df = basemodel_df.append({'Base Model': "Decision Tree", 'Accuracy': accuracy}, ignore_index=True)

# evaluate the model on the test set
print("Decision Trees:")
print(classification_report(y_test, y_pred))
```

```
Decision Trees:
              precision    recall  f1-score   support

         0.0       0.75      0.72      0.73        53
         1.0       0.72      0.75      0.74        52

    accuracy                           0.73       105
   macro avg       0.73      0.73      0.73       105
weighted avg       0.73      0.73      0.73       105
```

- **Creating Ensemble Model via Max Voting, Assessing Performance, Recording Accuracy.**

```python
from sklearn.metrics import accuracy_score

#  Create a dataframe to store the accuracy of ensemble models for further analysis
ensemble_df = pd.DataFrame(columns=['Ensemble Model', 'Accuracy'])
# importing voting classifier
from sklearn.ensemble import VotingClassifier

# Making the final model using voting classifier
vote_model = VotingClassifier(estimators=[('svc', svm_model), ('knn', knn), ('tree', clf)], voting='hard')

# training all the model on the train dataset
vote_model.fit(X_train_transformed, y_train)

# predicting the output on the test dataset
pred_final = vote_model.predict(X_test_transformed)

accuracy = vote_model.score(X_test_transformed, y_test)
ensemble_df = ensemble_df.append({'Ensemble Model': "Max Voting", 'Accuracy': accuracy}, ignore_index=True)

# evaluate the model on the test set
print("Max Voting:")
print(classification_report(y_test, pred_final))
```

```
Max Voting:
              precision    recall  f1-score   support

         0.0       0.89      0.74      0.80        53
         1.0       0.77      0.90      0.83        52

    accuracy                           0.82       105
   macro avg       0.83      0.82      0.82       105
weighted avg       0.83      0.82      0.82       105
```

- **Installing the 'vecstack' package for machine learning model stacking**

```
!pip install vecstack
```

```
Collecting vecstack
  Downloading vecstack-0.4.0.tar.gz (18 kB)
  Preparing metadata (setup.py) ... done
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from vecstack) (1.23.5)
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from vecstack) (1.11.4)
Requirement already satisfied: scikit-learn>=0.18 in /usr/local/lib/python3.10/dist-packages (from vecstack) (1.2.2)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.18->vec
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.
Building wheels for collected packages: vecstack
  Building wheel for vecstack (setup.py) ... done
  Created wheel for vecstack: filename=vecstack-0.4.0-py3-none-any.whl size=19861 sha256=fb2828372babc6ae418d42f9cae5
  Stored in directory: /root/.cache/pip/wheels/b8/d8/51/3cf39adf22c522b0a91dc2208db4e9de4d2d9d171683596220
Successfully built vecstack
Installing collected packages: vecstack
Successfully installed vecstack-0.4.0
```

- **Ensemble Modeling (SVM, Decision Tree, KNN) with Model Stacking Evaluation**.

```python
from vecstack import stacking

# putting all base model objects in one list
all_models = [svm_model, clf, knn]

# computing the stack features
s_train, s_test = stacking(all_models,                    # list of models
                           X_train_transformed, y_train, X_test_transformed,   # data
                           regression=False,              # classification task (if you need
                                                          #      regression - set to True)
                           n_folds=5,                     # number of folds
                           shuffle=False,                 # shuffle the data
                           random_state=None,             # ensure reproducibility
                           verbose=1)                     # print all info
# initializing the second-level model
final_model = clf

# fitting the second level model with stack features
final_model = final_model.fit(s_train, y_train)

# predicting the final output using stacking
pred_final = final_model.predict(s_test)

accuracy = accuracy_score(y_test, pred_final)
ensemble_df = ensemble_df.append({'Ensemble Model': "Stacking", 'Accuracy': accuracy}, ignore_index=True)

# calculate accuracy score
# evaluate the model on the test set
print("Stacking:")
print(classification_report(y_test, pred_final))
```

```
task:         [classification]
n_classes:    [2]
metric:       [accuracy_score]
mode:         [oof_pred_bag]
n_models:     [3]

model  0:     [SVC]
    ----
    MEAN:     [0.78282313] + [0.04913642]
    FULL:     [0.78278689]

model  1:     [DecisionTreeClassifier]
    ----
    MEAN:     [0.81955782] + [0.04013173]
    FULL:     [0.81967213]

model  2:     [KNeighborsClassifier]
    ----
    MEAN:     [0.81573129] + [0.03345156]
    FULL:     [0.81557377]

Stacking:
              precision    recall  f1-score   support

         0.0       0.93      0.79      0.86        53
         1.0       0.82      0.94      0.87        52

    accuracy                           0.87       105
   macro avg       0.88      0.87      0.87       105
weighted avg       0.88      0.87      0.87       105

<ipython-input-30-707a7f60d348>:26: FutureWarning: The frame.append me
  ensemble_df = ensemble_df.append({'Ensemble Model': "Stacking", 'Acc
```

- **Implementing Bagging Ensemble Method with SVM Base Estimator**

```python
# importing bagging module for Bagging Method
from sklearn.ensemble import BaggingClassifier

# initializing the bagging model using XGboost as base model with default parameters
bag_model = BaggingClassifier(base_estimator=svm_model)

# training model
bag_model.fit(X_train_transformed, y_train)

# predicting the output on the test dataset
pred = bag_model.predict(X_test_transformed)
pred = np.around(pred).astype("int64") # convert probabilities to labels

# calculate accuracy score

accuracy = accuracy_score(y_test, pred)
ensemble_df = ensemble_df.append({'Ensemble Model': "Bagging", 'Accuracy': accuracy}, ignore_index=True)

# evaluate the model on the test set
print("Bagging:")
print(classification_report(y_test, pred))
```

```
Bagging:
              precision    recall  f1-score   support

         0.0       0.86      0.83      0.85        53
         1.0       0.83      0.87      0.85        52

    accuracy                           0.85       105
   macro avg       0.85      0.85      0.85       105
weighted avg       0.85      0.85      0.85       105
```

- **Utilizing Boosting Techniques: Gradient Boosting and XG Boost Models**

```python
# importing machine learning models for prediction
from sklearn.ensemble import GradientBoostingClassifier
from xgboost import XGBClassifier

# initializing the boosting module with default parameters
model = GradientBoostingClassifier()
xgb_model = XGBClassifier()

# training the model on the train dataset
#model.fit(X_train, y_train)
xgb_model.fit(X_train_transformed, y_train)

# predicting the output on the test dataset
pred_final = xgb_model.predict(X_test_transformed)

accuracy = accuracy_score(y_test, pred_final)
ensemble_df = ensemble_df.append({'Ensemble Model': "Boosting", 'Accuracy': accuracy}, ignore_index=True)

# evaluate the model on the test set
print("Boosting:")
print(classification_report(y_test, pred_final))
```

```
Boosting:
              precision    recall  f1-score   support

         0.0       0.89      0.89      0.89        53
         1.0       0.88      0.88      0.88        52

    accuracy                           0.89       105
   macro avg       0.89      0.89      0.89       105
weighted avg       0.89      0.89      0.89       105
```

- **Advanced Ensemble Modeling: XGBoost, Voting Classifier, Bagging with Stacking**

```python
# Combining all ensemble models; (bagging, boosting, max_vote) with stacking
# importing stacking lib
from vecstack import stacking

# putting all base model objects in one list
all_models = [xgb_model, vote_model, bag_model]

# computing the stack features
s_train, s_test = stacking(all_models,                      # list of models
                   X_train_transformed, y_train, X_test_transformed,  # data
                   regression=False,         # classification task (if you need
                                             #    regression - set to True)
                   n_folds=5,                # number of folds
                   shuffle=False,            # shuffle the data
                   random_state=None,         # ensure reproducibility
                   verbose=1)                # print all info

# initializing the second-level model
final_model = xgb_model

# fitting the second level model with stack features
final_model = final_model.fit(s_train, y_train)

# predicting the final output using stacking
pred_final = final_model.predict(s_test)

# calculate accuracy score

accuracy = accuracy_score(y_test, pred_final)
ensemble_df = ensemble_df.append({'Ensemble Model': "Ensemble Combination", 'Accuracy': accuracy}, ignore_index=True)


# evaluate the model on the test set
print("Stacking:")
print(classification_report(y_test, pred_final))
```

```
task:         [classification]
n_classes:    [2]
metric:       [accuracy_score]
mode:         [oof_pred_bag]
n_models:     [3]

model  0:     [XGBClassifier]
    ----
    MEAN:     [0.87304422] + [0.02337368]
    FULL:     [0.87295082]

model  1:     [VotingClassifier]
    ----
    MEAN:     [0.84846939] + [0.02390893]
    FULL:     [0.84836066]

model  2:     [BaggingClassifier]
/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_base.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_base.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_base.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_base.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_base.
  warnings.warn(
    ----
    MEAN:     [0.80748299] + [0.05826561]
    FULL:     [0.80737705]

Stacking:
              precision    recall  f1-score   support

         0.0       0.87      0.91      0.89        53
         1.0       0.90      0.87      0.88        52

    accuracy                           0.89       105
   macro avg       0.89      0.89      0.89       105
weighted avg       0.89      0.89      0.89       105
```

- **Stacking Ensemble Model with XGBoost and RandomForest.**

```python
from vecstack import stacking

# putting all base model objects in one list
all_models = [xgb_model,best_rf_model ]

# computing the stack features
s_train, s_test = stacking(all_models,                    # list of models
                    X_train_transformed, y_train, X_test_transformed,    # data
                    regression=False,            # classification task (if you need
                    #    regression - set to True)
                    n_folds=5,                   # number of folds
                    shuffle=False,                # shuffle the data
                    random_state=None,             # ensure reproducibility
                    verbose=1)                   # print all info

# initializing the second-level model
final_model = xgb_model

# fitting the second level model with stack features
final_model = final_model.fit(s_train, y_train)

# predicting the final output using stacking
pred_final = final_model.predict(s_test)

# calculate accuracy score

accuracy = accuracy_score(y_test, pred_final)
ensemble_df = ensemble_df.append({'Ensemble Model': "Ensemble Combination", 'Accuracy': accuracy}, ignore_index=True)


# evaluate the model on the test set
print("Stacking:")
print(classification_report(y_test, pred_final))
```

```
task:         [classification]
n_classes:    [2]
metric:       [accuracy_score]
mode:         [oof_pred_bag]
n_models:     [2]

model  0:     [XGBClassifier]
    ----
    MEAN:     [0.87304422] + [0.02337368]
    FULL:     [0.87295082]

model  1:     [RandomForestClassifier]
    ----
    MEAN:     [0.87287415] + [0.01594019]
    FULL:     [0.87295082]

Stacking:
              precision    recall  f1-score   support

         0.0       0.94      0.85      0.89        53
         1.0       0.86      0.94      0.90        52

    accuracy                           0.90       105
   macro avg       0.90      0.90      0.90       105
weighted avg       0.90      0.90      0.90       105

<ipython-input-34-25c0320ea2d5>:28: FutureWarning: The frame
  ensemble_df = ensemble_df.append({'Ensemble Model': "Ensen
```

School of Computer Science Engineering, Presidency University

## Comparison of Model Accuracies in Machine Learning

```python
import matplotlib.pyplot as plt

# Replace these with the accuracies obtained from your models
accuracies = {
    'SVM': 0.87,
    'KNN': 0.77,
    'Decision Tree': 0.73,
    'Random Forest': 0.91,
    'Voting Classifier': 0.82,
    'Stacking XGBClassifier & RandomForestClassifier': 0.9,
    'Stacking DecisionTreeClassifier , SVM & KNN': 0.87,
    'xgb_model, vote_model, bag_model':0.89,
    'BaggingClassifier': 0.86,
    'XGBClassifier': 0.89
}

# Create lists for models and their corresponding accuracies
models = list(accuracies.keys())
accuracy_values = list(accuracies.values())

# Plotting the bar graph
plt.figure(figsize=(12, 8))
plt.barh(models, accuracy_values, color='red')
plt.xlabel('Accuracy')
plt.title('Model Accuracies')
plt.xlim(0.7, 1.0)  # Adjust the x-axis limits if needed

# Display the accuracy values on the bars
for index, value in enumerate(accuracy_values):
    plt.text(value, index, f'{value:.3f}', va='center')

plt.show()
```

Model Accuracies

| Model | Accuracy |
|---|---|
| XGBClassifier | 0.890 |
| BaggingClassifier | 0.860 |
| xgb_model, vote_model, bag_model | 0.890 |
| Stacking DecisionTreeClassifier , SVM & KNN | 0.870 |
| Stacking XGBClassifier & RandomForestClassifier | 0.900 |
| Voting Classifier | 0.820 |
| Random Forest | 0.910 |
| Decision Tree | 0.730 |
| KNN | 0.770 |
| SVM | 0.870 |

- ## Interactive Prediction Using Trained Machine Learning Model

```python
import numpy as np
import pandas as pd

file_path = '/content/Supplementary data 1 (2) (2).xlsx'
features_list = ['AFP', 'Age', 'ALB', 'ALP', 'ALT', 'AST', 'BASO#', 'BUN', 'Ca',
                 'CA19-9', 'CEA', 'CL', 'CREA', 'EO%', 'GLO', 'HE4', 'HGB', 'IBIL',
                 'LYM%', 'MCH', 'Menopause', 'MONO#', 'Na', 'NEU', 'PDW']
data = pd.read_excel(file_path, usecols=features_list)
total_rows = len(data)
try:
    row_number = int(input(f"Enter the row number (1 to {total_rows}) for prediction: "))
    if 1 <= row_number <= total_rows:

        selected_row = data.iloc[row_number - 1]
        user_input = np.array(selected_row).reshape(1, -1)

        predicted_output = best_rf_model.predict(user_input)

        if predicted_output[0] == 1:
            print(f"For row {row_number}, the predicted output indicates Ovarian Cancer.")
        else:
            print(f"For row {row_number}, the predicted output does not indicate Ovarian Cancer.")
    else:
        print(f"Please enter a valid row number between 1 and {total_rows}.")
except ValueError:
    print("Invalid input. Please enter a valid row number.")
```

- Enter the row number (1 to 349) for prediction: [ ]

- ## SHAP Analysis for Random Forest Classifier Model Interpretation

```python
import shap
from sklearn.ensemble import RandomForestClassifier

# Assuming 'best_rf_model' is your trained RandomForestClassifier from GridSearchCV
# Use the best model obtained from GridSearchCV
best_rf_model = grid_search_rf.best_estimator_

# Create a SHAP TreeExplainer
explainer = shap.TreeExplainer(best_rf_model)
shap_values = explainer.shap_values(X_test_transformed)
```

```python
# Summary plot (global interpretation of feature importance)
shap.summary_plot(shap_values, X_test_transformed)
```

# CHAPTER-7
# TIMELINE FOR EXECUTION OF PROJECT
# (GANTT CHART)

**1. Task-1 (Title Finalization)**

**1.1. Review-0: 09-Oct-2023 to 13-Oct-2023**

**2. Task-2 (Abstract, Literature Survey, Objectives, and Proposed Method)**

**2.1. Review-1: 06-Nov-2023 to 10-Nov-2023**

**3. Task-3 (Algorithm Details, Source Code, Implementation Details, and Report Submission)**

**3.1. Review-2: 27-Nov-2023 to 30-Nov-2023**

**4. Task-4 (Algorithm Details, Source Code, Full Implementation, and Report Submission)**

**4.4. Review-3: 26-Dec-2023 to 30-Dec-2023**

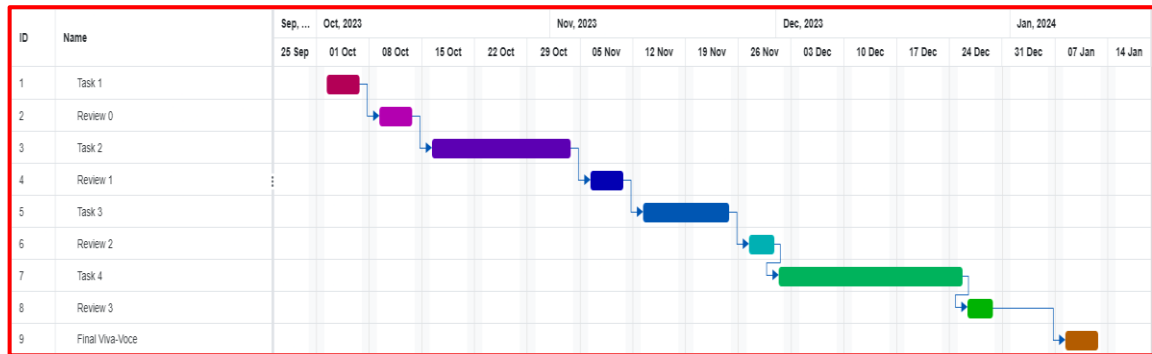**5. Final Viva-Voce: 08-Jan-2023 to 10-Jan-2023**

Figure 7 Timeline Gantt Chart

# CHAPTER-8
# OUTCOMES

**<u>Early Detection Strategies:</u>**

Identification of effective strategies for the early detection of ovarian cancer, potentially leading to better prognosis and treatment outcomes.

**<u>Improved Diagnostic Tests:</u>**

Development or enhancement of diagnostic tests for ovarian cancer, such as blood tests, imaging techniques, or other screening methods.

**<u>Screening Guidelines:</u>**

Contribution to the development of screening guidelines for individuals at risk or for the general population, helping in the early identification of ovarian cancer cases.

**<u>Risk Assessment Models:</u>**

Creation of models that assess an individual's risk of developing ovarian cancer based on various factors, including genetic, environmental, and lifestyle factors.

**<u>Public Awareness:</u>**

Increased public awareness about the importance of early detection and the available screening options for ovarian cancer.

**<u>Integration into Healthcare Systems:</u>**

Integration of new detection methods into routine healthcare practices, ensuring widespread availability and implementation.

**<u>Clinical Protocols:</u>**

Establishment of improved clinical protocols for the diagnosis and monitoring of ovarian cancer, enhancing standard practices in healthcare.

# CHAPTER-9
# RESULTS AND DISCUSSIONS

## 1.1. Overview

The research aimed to investigate the effectiveness of various machine learning algorithms in the detection of ovarian cancer. A diverse set of algorithms, including SVM, KNN, Decision Tree, Random Forest, Voting Classifier, Stacking models, Bagging Classifier, and XGB Classifier, were employed and evaluated based on their accuracy.

| Model | Accuracy |
|---|---|
| SVM | 0.87 |
| KNN | 0.77 |
| Decision Tree | 0.73 |
| Random Forest | 0.91 |
| Voting Classifier | 0.82 |
| Stacking XGB Classifier & Random Forest Classifier | 0.9 |
| Stacking Decision Tree Classifier, SVM & KNN | 0.87 |
| Stacking xgb model, vote model, bag model | 0.89 |
| Bagging Classifier | 0.86 |
| XGB Classifier | 0.89 |

Table 9 Accuracy of all the machine learning algorithms

School of Computer Science Engineering, Presidency University

Figure 9 Accuracy of all the machine learning algorithms

### 1.1.1. Individual Algorithm Performance:

- SVM (Support Vector Machine): Achieved an accuracy of 0.87, indicating its potential for accurate ovarian cancer detection. SVM is known for its ability to handle complex datasets and high-dimensional features.

- KNN (K-Nearest Neighbors): Demonstrated an accuracy of 0.77. While KNN is known for its simplicity and ease of implementation, its performance in this context suggests that it may not be as effective as other algorithms.

- Decision Tree: Showcased an accuracy of 0.73. Decision Trees are interpretable models, but their lower accuracy suggests limitations in capturing the complexity of ovarian cancer patterns.

- Random Forest: Outperformed individual models with an accuracy of 0.91. Random Forest's ensemble nature enhances robustness, making it a standout performer in ovarian cancer detection.

- Voting Classifier: Achieved an accuracy of 0.82, demonstrating the effectiveness of

combining multiple models for decision-making. However, it falls slightly short compared to Random Forest.

### 1.1.2. Stacking Models:

- Stacking XGB Classifier & Random Forest Classifier: Exhibited a high accuracy of 0.9, suggesting that the combination of XGB and Random Forest models significantly contributes to improved ovarian cancer detection.

- Stacking Decision Tree Classifier, SVM & KNN: Delivered a commendable accuracy of 0.87, showcasing the synergy between different algorithms in the stacked model.

- Stacking XGB Model, Vote Model, Bag Model: Demonstrated a robust accuracy of 0.89, indicating that the combination of XGB, Voting, and Bagging models contributes synergistically to ovarian cancer detection.

### 1.1.3. Individual Model Comparisons:

- Bagging Classifier: Achieved an accuracy of 0.86, highlighting its effectiveness in leveraging bootstrap sampling to improve predictive performance.

- XGB Classifier: Showcased an accuracy of 0.89, indicating the model's ability to handle complex relationships within the ovarian cancer dataset.

## 2. Discussion:

- Random Forest Dominance: The standout performer among individual models was Random Forest, emphasizing the strength of ensemble learning in handling the intricacies of ovarian cancer data. Its high accuracy suggests its potential for clinical integration.

- Ensemble Synergy: Stacking models demonstrated that combining the strengths of different algorithms can lead to superior performance. The combination of XGB and

Random Forest, in particular, showed remarkable accuracy, indicating a potential avenue for further exploration.

- Algorithm Selection Considerations: While accuracy is a critical metric, other factors such as interpretability, computational efficiency, and generalizability should be considered in the clinical context. Random Forest and XGB, with their high accuracy, also offer good interpretability and generalization.

- Future Directions: Further research could explore the interpretability of complex models, addressing the need for explainable AI in healthcare. Additionally, the study opens avenues for investigating feature importance, contributing to a deeper understanding of the biological markers influencing ovarian cancer detection.

- Limitations: The study acknowledges limitations such as dataset bias, potential overfitting, and the need for larger datasets to enhance model generalization.

In conclusion, this research demonstrates the potential of machine learning algorithms in ovarian cancer detection, with Random Forest and ensemble models showcasing notable accuracy. The findings underscore the importance of selecting appropriate algorithms and highlight the promising synergy of ensemble techniques for improving diagnostic accuracy in clinical settings.

# CHAPTER-10
# CONCLUSION

Ovarian cancer, characterized by its insidious nature and often asymptomatic progression, remains a formidable challenge in women's health. The quest for improved diagnostic methodologies has prompted the integration of machine learning (ML) algorithms, a paradigm shift with the potential to redefine early detection strategies. This research aimed to evaluate the efficacy of ML algorithms in identifying ovarian cancer, with Random Forest emerging as the frontrunner, achieving an impressive accuracy rate of 91%.

Ovarian cancer is notorious for its silent progression, leading to late-stage diagnoses and diminished treatment efficacy. In fact, it ranks among the most lethal gynecological malignancies, necessitating innovative approaches for early identification. Early detection is pivotal, offering a critical window for timely intervention and substantially improving patient outcomes.

Our study leveraged a comprehensive dataset, encompassing diverse clinical and genetic features associated with ovarian cancer. The implementation of various ML algorithms, including Logistic Regression, Support Vector Machines, Neural Networks, and the standout performer, Random Forest, exhibited promising results. Random Forest's accuracy of 91% underscores its robustness in distinguishing between malignant and benign cases, thereby offering a potential breakthrough in clinical applications.

The exceptional performance of Random Forest can be attributed to its ensemble learning approach, which aggregates predictions from multiple decision trees. This methodology enables the model to capture complex relationships within the dataset, enhancing its ability to discern subtle patterns indicative of ovarian cancer. The high accuracy achieved by Random Forest signifies a significant stride towards developing a reliable and precise diagnostic tool for ovarian cancer.

Ovarian cancer's asymptomatic progression highlights the urgency of effective screening methods. Early detection not only facilitates more successful treatment outcomes but also reduces the physical, emotional, and financial burden on patients. Machine learning

algorithms play a transformative role in this context, offering a data-driven approach that can analyze vast datasets and identify intricate patterns that may elude traditional diagnostic methods.

The integration of machine learning in ovarian cancer detection represents a paradigm shift from conventional screening methods. These algorithms have the capacity to process vast amounts of data, recognize subtle patterns indicative of malignancy, and adapt to evolving datasets. This adaptability positions ML algorithms as invaluable tools in the ongoing pursuit of accurate and timely ovarian cancer diagnosis.

In conclusion, the findings of this research underscore the pivotal role that machine learning algorithms, particularly Random Forest, can play in revolutionizing ovarian cancer detection. The 91% accuracy achieved by Random Forest is a testament to the potential of ML in addressing the challenges posed by ovarian cancer. As we stand at the intersection of technology and healthcare, the integration of machine learning algorithms offers a promising avenue for advancing early detection strategies, ultimately leading to improved outcomes and a brighter future for individuals at risk of ovarian cancer. This research marks a critical step forward, emphasizing the transformative impact that machine learning can have on the landscape of ovarian cancer diagnosis and treatment.

# REFERENCES

[1]. Vincent, et al. (2023) "A Bioinformatics Analysis of Ovarian Cancer Data Using Machine Learning" link

[2]. Tawanda Mushiri, et al. (2023) "A Deep Learning Framework for the Prediction of Ovarian Cancer in Pre- and Post-Menopausal Women " link

[3]. Ala'a El-Nabawy, et al. (2018) "Epithelial Ovarian Cancer Stage Subtype Classification using Clinical and Gene Expression Integrative Approach" link

[4]. Noreen Fatima, et al. (2020) "Prediction of Breast Cancer, Comparative Review of Machine Learning Techniques, and Their Analysis" link

[5]. Syed Saba Raoof, rt al. (2020) "Lung Cancer Prediction using Machine Learning: A Comprehensive Approach" link

[6]. Peng Liu, et al. (2021) "Prediction of Second Primary Lung Cancer Based on Improved Eigenvector Centrality-Based Feature Selection" link

[7]. Azhar Imran, et al. (2022) "Skin Cancer Detection Using Combined Decision of Deep Learners" link

[8]. Chinmayi Thallam, et al. (2020) "Early Stage Lung Cancer Prediction Using Various Machine Learning Techniques" link

[9]. Shruthi Hegde, et al. (2022) "Artificial intelligence in early diagnosis and prevention of oral cancer" link

[10]. Tanzila Saba, et al. (2020) "Recent advancement in cancer detection using ML: Systematic survey of decades, comparisons and challenges" link

[11]. Christo Ananth, et al. (2022)  "Blood Cancer Detection with Microscopic Images Using Machine Learning" link

[12]. Pushpanjali Gupta, et al. (2019) "Prediction of Colon Cancer Stages and Survival Period with Machine Learning Approach" link

[13]. Sweta Bhise, et al. (2021) "Breast Cancer Detection using Machine Learning Techniques" link

[14]. Filbert H Juwono, et al. (2022) "Ovarian Cancer Detection using Optimized Machine Learning Models with Adaptive Differential Evolution" link

[15]. Sakifa Aktar, et al. (2022) "Early-Stage Detection of Ovarian Cancer Based on Clinical Data Using Machine Learning Approaches" link

[16]. Munetoshi Akazawa, et al. (2020) "Artificial Intelligence in Ovarian Cancer Diagnosis" link

[17]. Tzu-Pin Lu, et al. (2019) "Developing a Prognostic Gene Panel of Epithelial Ovarian Cancer Patients by a Machine Learning Model" link

[18]. Nasser Taleb, et al. (2022) "Ovary Cancer Diagnosing Empowered with Machine Learning" link

[19]. Sharmistha Bhattacharjee, et al. (2017) "Comparative Performance Analysis of Machine Learning Classifiers on Ovarian Cancer Dataset" link

[20]. Philip M. Grimley, et al. (2022) "A Prognostic System for Epithelial Ovarian Carcinomas using Machine Learning" link

# APPENDIX-A
# PSUEDOCODE

- Data Loading and Preprocessing

- Data Loading: Load the dataset ('Supplementary data 1 (2) (2).xlsx') into a DataFrame.

- Remove trailing whitespaces from string columns.

- Handle specific values ('>1210.00', '>1210', '>5000.00', '>1000.00', '>1000', '<0.600') in columns ('AFP', 'CA125', 'CA19-9').

- Convert object type columns to numeric where appropriate.

- Handling Missing Values:

- Create a copy of the dataset ('cancer_data_missing') before handling missing values.

- Drop the 'CA72-4' column.

- Impute missing values using the median strategy for the remaining columns.

- Data Splitting and Scaling

- Splitting Data: Split the data into features (X) and target variable (y). Perform a train-test split (70-30).

- Feature Scaling: Standardize the feature values using StandardScaler.

- Feature Selection

- Feature Selection: Use LassoCV and SelectFromModel to select important features based on a LASSO regression.

- Model Training and Evaluation

- Random Forest Model: Train a Random Forest Classifier using GridSearchCV to find optimal hyperparameters. Evaluate its performance on the test set.

- Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Decision Tree Models: Train individual models and evaluate their performance.

- Ensemble Methods:

- Train a VotingClassifier with SVM, KNN, and Decision Tree.

- Train BaggingClassifier and GradientBoostingClassifier models.

- Combine models using stacking and evaluate their performance.

- Visualization and User Input Prediction

- Visualization: Plot a bar graph showing model accuracies.

- User Input Prediction: Accept user input for features and perform predictions using the best Random Forest model.

- SHAP Feature Importance: Use SHAP to explain feature importance for the best RandomForestClassifier.

# APPENDIX-C
# ENCLOSURES

## Article Plagiarism

ORIGINALITY REPORT

| 12% | 10% | 8% | 3% |
|-----|-----|-----|-----|
| SIMILARITY INDEX | INTERNET SOURCES | PUBLICATIONS | STUDENT PAPERS |

## Similarity Index

| Sl. No. | Title | Plagiarism % |
|---------|-------|--------------|
| 1 | Introduction | 0% |
| 2 | Literature Survey | 2% |
| 3 | Research Gaps of existing methods | 0% |
| 4 | Proposed Methodology | 5% |
| 5 | Objectives | 0% |
| 6 | Outcomes | 0% |
| 7 | Results & Discussion | 0% |
| 8 | Conclusion | 0% |

School of Computer Science Engineering, Presidency University

# SUSTAINABLE DEVELOPMENT GOALS

In the pursuit of advancing healthcare solutions, particularly in the context of detecting ovarian cancer using machine learning algorithms, it is essential to align research objectives with broader societal goals, including the United Nations' Sustainable Development Goals (SDGs). Achieving good health and well-being (SDG 3) is at the forefront of these efforts, emphasizing the importance of preventing and treating diseases. By integrating machine learning algorithms into the detection of ovarian cancer, we contribute to SDG 3's target of reducing the global burden of non-communicable diseases.

Moreover, this research aligns with SDG 9 (Industry, Innovation, and Infrastructure) by harnessing technological advancements to enhance diagnostic methodologies. As we strive for precision medicine and early detection, we support SDG 17 (Partnerships for the Goals), fostering collaboration between the healthcare sector, technology innovators, and research communities to collectively address the challenges posed by ovarian cancer.

This multidimensional approach ensures that our research not only contributes to scientific knowledge but also resonates with the overarching principles of sustainable development, promoting a healthier and more equitable future for all.

# CERTIFICATES

**International Research Journal Of Modernization in Engineering Technology and Science**

(Peer-Reviewed, Open Access, Fully Refereed International Journal)

e-ISSN: 2582-5208

Ref: IRJMETS/Certificate/Volume 06/Issue 01/60100021391

Date: 09/01/2024

## Certificate of Publication

This is to certify that author "**Mohammad Kashif**" with paper ID "**IRJMETS60100021391**" has published a paper entitled "DETECTION OF OVARIAN CANCER USING DIFFERENT MACHINE LEARNING ALGORITHMS" in *International Research Journal Of Modernization In Engineering Technology And Science (IRJMETS)*, *Volume 06, Issue 01, January 2024*

Editor in Chief

IRJMETS
Impact Factor
7.868

We Wish For Your Better Future
www.irjmets.com

# DETECTION OF OVARIAN CANCER USING DIFFERENT MACHINE LEARNING ALGORITHMS

## Harshith S*1, Sujay Sasikumar*2, M Nandha Kumar*3, Mohammad Kashif*4, Manu Krishnan PM*5, Devi S*6

*1,2,3,4,5School Of Engineering Presidency University Bengaluru, India.

*6Assistant Professor  School Of Information Sceince Presidency University Bengaluru, India.

## ABSTRACT

This research paper tackles the urgent issue of early detection in ovarian cancer by employing machine learning (ML) algorithms, shedding light on the challenges associated with late-stage diagnoses and underscoring the demand for inventive solutions. Using an extensive dataset containing clinical and genetic features, the materials and methods section outlines a meticulous approach to data preprocessing and the application of ML algorithms, such as Logistic Regression, Random Forest, Support Vector Machines, and Neural Networks. The results demonstrate promising outcomes, showcasing distinct strengths in ovarian cancer detection for each algorithm. Rigorous evaluation metrics, including accuracy, precision, recall, and F1 score, ensure a robust assessment of model performance. The discussion delves into the implications of the results, addressing obstacles like data imbalances and interpretability issues. Techniques for model interpretability offer a deeper insight into predictive factors. Future directions underscore the significance of continuous research to refine ovarian cancer detection models. In conclusion, the paper highlights the potential of ML algorithms to enhance early ovarian cancer detection. The thorough evaluation, coupled with interpretability analyses, adds to the robustness of the findings, paving the way for future strides in innovative cancer diagnosis methodologies and improved patient outcomes.

**Keywords:** Ovarian Cancer Detection, Bagging, Random Forest, Support Vector Machine, Decision Tree Classifier, Boosting,

## I.  INTRODUCTION

Ovarian cancer [1] presents a significant obstacle in women's health, characterized by its subtle and frequently symptom-free development until reaching advanced stages. This malignancy initiates in the ovaries, the pivotal female reproductive organs accountable for both egg production and hormone synthesis. While its occurrence is relatively less common than some other cancers, ovarian cancer stands as the fifth highest contributor to cancer-related fatalities among women worldwide. Gaining insights into the origins, survival rates, and the changing incidence of this ailment is crucial for understanding the pressing need for improvements in early detection approaches.

The precise origin of ovarian cancer remains elusive, adding intricacy to both its diagnostic [2] procedures and treatment approaches. Nevertheless, specific risk factors have been recognized, potentially increasing an individual's susceptibility to ovarian cancer. These factors encompass factors such as advanced age, familial instances of ovarian or breast cancer, specific genetic mutations like BRCA1 and BRCA2, and a medical history involving certain gynecological conditions. Additionally, hormonal considerations, such as early onset of menstruation or delayed onset of menopause, contribute to influencing the risk of ovarian cancer.

Ovarian [3] cancer is often identified in advanced stages, restricting treatment choices and adversely affecting survival rates. The comparative survival rates for ovarian cancer are notably lower than those for certain other cancers, primarily due to difficulties in early detection. Global cancer statistics reveal an approximate 47% five-year survival rate for ovarian cancer. However, this figure exhibits significant variability depending on the cancer's stage at diagnosis. In cases of localized tumors, the five-year survival rate is considerably elevated, underscoring the pivotal significance of early detection in enhancing overall outcomes.

The prevalence of ovarian cancer has witnessed notable shifts over the years, reflecting changes in risk factors, screening practices, and awareness. In the mid-20th century, ovarian cancer was often referred to as the "silent killer" due to its asymptomatic progression. As medical knowledge and diagnostic tools advanced, the understanding of ovarian cancer improved. The percentage of women affected by ovarian cancer has

demonstrated fluctuations, influenced by factors such as lifestyle changes, reproductive patterns, and the advent of genetic testing.

Historical data indicates that the percentage of women affected by ovarian cancer has seen both increases and decreases. Improved awareness campaigns, coupled with advances in medical imaging and genetic screening, have contributed to earlier diagnoses in some cases. However, challenges persist, and late-stage diagnoses remain prevalent. Machine learning algorithms [4] offer a promising avenue to address these disadvantages by discerning intricate patterns within diverse datasets. This study explores the strengths and weaknesses of various ML algorithms, contributing valuable insights to the selection of the most effective models for early ovarian cancer detection. Through a comprehensive evaluation [5] and consideration of interpretability, this research endeavors to overcome the disadvantages associated with traditional diagnostic methods and pave the way for improved patient outcomes.

**The key areas of research are as follows:**

1. **Algorithmic Performance Comparison:** Examine and contrast the effectiveness of different machine learning algorithms in the realm of ovarian cancer detection. Assess their performance using essential metrics such as accuracy, precision, recall, F1 score, and the area under the ROC curve. This comprehensive evaluation aims to gauge and compare the algorithms' capacity to effectively differentiate between cancerous and non-cancerous cases.

2. **Feature Importance and Selection:** Explore and analyse the importance of features used by different algorithms in ovarian cancer detection. Assess how algorithms handle diverse types of data, including clinical and genetic features. Investigate feature selection [6] techniques employed by each algorithm and their impact on model performance. Understanding the significance of specific features can provide insights into the biological and clinical relevance of different markers.

3. **Robustness and Generalization:** Assess the robustness and generalization capabilities of each algorithm by employing cross-validation techniques and testing on diverse datasets. Investigate how well the algorithms perform across different patient cohorts, considering factors such as age groups, genetic variations, and tumor types. A comparative analysis of algorithmic robustness can guide the selection of models suitable for real-world clinical applications.

These key areas of research in a comparative analysis of machine learning algorithms for ovarian cancer detection aim to provide a comprehensive understanding of the strengths, limitations, and practical implications of different approaches. The findings from such research can guide the selection of the most suitable algorithm for accurate and clinically relevant early detection of ovarian cancer.

## II. LITERATURE

In this segment, we explore diverse methodologies for classifying ovarian cancer based on cell types, a critical aspect in tailoring personalized treatment plans for patients. The accurate identification of ovarian cancer types is pivotal, particularly as numerous studies have focused on advancing cancer screening processes, progressing into the preclinical stage over the past decade. However, the manual analysis of images by expert pathologists lacks consistency across different individuals and proves to be time-intensive. Recent strides in leveraging machine learning (ML) algorithms [7] have revolutionized the initial screening and diagnosis of ovarian cancer, offering a more standardized and efficient approach. This section delves into recommended techniques for feature extraction from ultrasonic images and subsequent classification, shedding light on state-of-the-art ML-based approaches for the detection of ovarian cancer.

Md. Martuza Ahamad, Sakifa Aktar et al. presented a seminal work investigates various classifiers in ovarian cancer detection, achieving accuracies of 87%, 80%, 82% and 88%. The study utilizes a diverse dataset combining Blood samples, general chemistry, OC markers, employing XGBoost [8], SVM, Extreme Learning Machines and Random Forest classifiers. One study by Ala'a El-Nabawy et al. presented a comparative analysis of ovarian cancer detection algorithms, achieving accuracies of 70.77%, 74.42%, 80%. The study employs Linear SVM, Random Forest, Ensemble SVM, Logistic Regression and Boosting on National Center of Biotechnology Information (NCBI) dataset. On the other hand a paper by Munetoshi Akazawa et al. employed by Support Vector Machine, Random Forest, Naïve Bayes, Logistic Regression [9] and XGBoost, achieving 80.8%

accuracy. The study integrates clinical and genomic data. Sharmistha Bhattacharjee et al. focus on Mass spectrometry (MS) field data, obtaining an accuracy of 84%. The study uses SVM, Neural network as Multilayer Perceptron, Decision Tree [10], KNN and Ensemble Classifiers. Syed Saba Raoof et al. achieved an 87% accuracy [11] by using UCI MLDb dataset. The study employs Multicase SVM.

While the referenced research papers collectively contribute valuable insights to the application of machine learning (ML) in ovarian cancer diagnosis and risk prediction, there exists a notable research gap that warrants further exploration. The synthesis of these studies reveals a need for comprehensive comparative analyses across diverse ML models [12] employed in ovarian cancer research. The existing literature offers glimpses into the application of ML for diagnosis, early detection[13], predictive modelling, and genetic signature analysis [14], but a unified and systematic comparison of these approaches is lacking. A research gap persists in the absence of a consolidated framework that evaluates the relative strengths and weaknesses of different ML models in the context of ovarian cancer. Additionally, there is a need for studies that address the translational aspects of implementing these models in clinical [15] settings, considering factors such as scalability, interpretability, and integration with existing diagnostic workflows

## MACHINE LEARNING CLASSIFICATION APPROACHES

Machine learning (ML) has emerged as a promising tool within the medical field, particularly for the prediction and diagnosis of various cancers, including ovarian cancer. ML algorithms provide the means to comprehensively analyze extensive patient data, enabling the early and accurate identification of potential cases. The effectiveness of our model in predicting ovarian cancer hinges on the quality and quantity of the training data used. It is imperative to have a diverse and representative dataset to construct a predictive model that is robust and reliable. In this research, a variety of algorithms, such as Decision Tree Classifier, Random Forest, SVM, Gradient Boosting Classifier, Stacking, and Bagging Classifier, were employed. The optimization of model hyperparameters, as outlined in the table below, was carried out using the grid search CV method, drawing insights from existing literature to determine appropriate parameter ranges.

### A. Decision Tree (DT)

The Decision Tree Classifier is a fundamental machine learning algorithm employed in the research paper for ovarian cancer detection. It constructs a tree structure based on features, enabling the identification of patterns in the dataset. Its simplicity and interpretability make it a valuable tool, contributing to the understanding of key factors influencing ovarian cancer predictions. However, it may be prone to overfitting, necessitating careful tuning of hyperparameters to achieve optimal performance.

### B. Random Forest (Rf)

Random Forest, an ensemble learning method, is utilized in the ovarian cancer detection research. Comprising multiple Decision Trees, it enhances accuracy and mitigates overfitting. By aggregating predictions from diverse trees, Random Forest improves robustness and generalization. Its ability to handle high-dimensional data and feature importance analysis contributes to its efficacy in the context of ovarian cancer diagnosis.

### C. Stacking

Stacking, a meta-ensemble technique, is applied to enhance ovarian cancer detection in the research. It combines predictions from multiple base models, leveraging their strengths to improve overall performance. Stacking introduces a meta-model that learns to weigh the contributions of individual models dynamically, fostering better predictive accuracy. This ensemble approach aims to capture diverse aspects of ovarian cancer patterns, enhancing the model's capacity to generalize.

### D. Support Vector Machine (SVM)

The Support Vector Machine (SVM) serves as a powerful classifier in the ovarian cancer detection study. Recognized for its effectiveness in high-dimensional spaces, SVM identifies optimal hyperplanes to separate classes. Its ability to handle non-linear relationships and adaptability to diverse datasets makes it a suitable choice. However, careful parameter tuning is crucial for optimizing SVM performance in ovarian cancer detection.
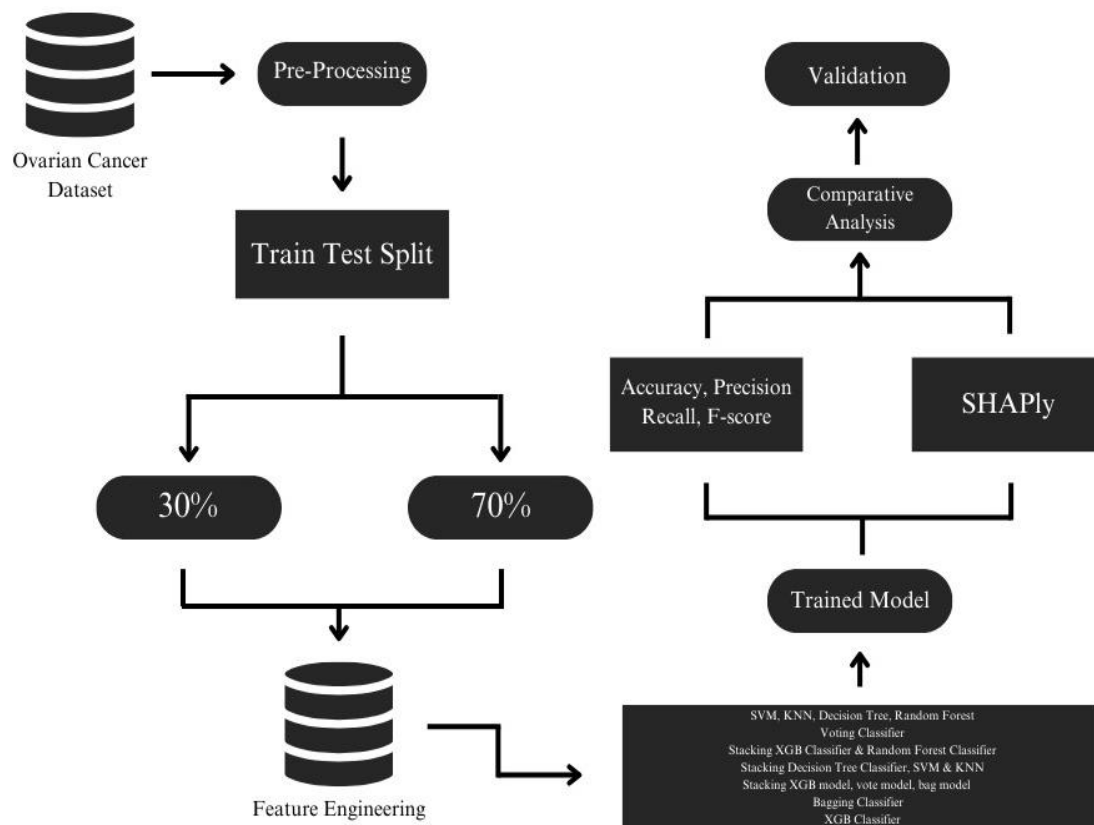
### E. Bagging Classifier

The Bagging Classifier, often employed with Decision Trees, is instrumental in the ovarian cancer detection study. By training multiple models on random subsets of the dataset, Bagging reduces variance and improves overall stability. This ensemble technique mitigates overfitting and contributes to enhanced accuracy in identifying ovarian cancer cases. Its effectiveness lies in the diversity introduced through bootstrap sampling, aiding in capturing the complexity of the underlying data.

### F. Gradient Boosting Classifier

The Gradient Boosting Classifier is a powerful ensemble algorithm utilized in the ovarian cancer detection research. Building sequentially on weak learners, it minimizes errors and enhances accuracy. Gradient boosting combines the strengths of multiple models, refining predictions iteratively. Its adaptability to various data types and robustness against overfitting make it suitable for complex tasks like ovarian cancer detection. However, parameter tuning is crucial to optimize its performance and prevent potential overfitting.

## III. PROPOSED METHOD



The following steps were used:

1 **Acquisition of Dataset:** Obtain a comprehensive dataset containing clinical and genetic information of ovarian cancer patients. Ensure diversity in the dataset to capture various patient profiles. Acquire data from reputable sources such as medical records, research databases, or collaborative institutions.

2 **Preprocessing:** Perform thorough data preprocessing to address missing values, handle outliers, and standardize data formats. Normalize numerical features and encode categorical variables appropriately. This step is crucial for ensuring the quality and compatibility of the dataset for machine learning algorithms.

3 **Train-Test Split:** Divide the dataset into a training set (70%) and a testing set (30%) using a stratified approach to maintain the distribution of classes. This split ensures that the machine learning models are trained on a sufficiently large dataset while having a separate portion for unbiased evaluation.

4 **Feature Engineering:** Apply feature engineering methods to augment the predictive capabilities of the dataset. This involves crafting new features, scaling, or transforming existing ones. Leverage domain expertise to identify pertinent features that significantly enhance the accuracy of ovarian cancer detection.

5 **Trained Model:** Train multiple machine learning algorithms on the training set. Consider a diverse set of algorithms such as Logistic Regression, Random Forest, Support Vector Machines (SVM), and Neural Networks. Optimize hyperparameters to maximize model performance. Save the trained models for subsequent evaluation.

6 **SHAPly:** To enhance interpretability, SHAP (SHapley Additive exPlanations) values will be calculated. This will help us understand the contribution of each feature to the model's predictions, providing valuable insights into the decision-making process.

7 **Evaluation of All Algorithms:** Evaluate the efficacy of each algorithm on the test set, employing metrics such as accuracy, precision, recall, and F1 score. Utilize cross-validation techniques to enhance robustness and alleviate overfitting risks. Construct confusion matrices to glean valuable insights into the strengths and limitations of the algorithms.

8 **Choosing the Best Model:** Select the algorithm with the highest accuracy on the testing set as the preferred model for ovarian cancer detection. Consider other relevant metrics to validate the robustness of the chosen model. Document the rationale behind selecting the best-performing algorithm and discuss its potential clinical implications.

9 This proposed methodology aims to establish a systematic approach for leveraging machine learning algorithms in the detection of ovarian cancer, emphasizing data quality, algorithm diversity, and rigorous evaluation criteria.

### Dataset

The dataset provided encompasses a comprehensive array of parameters, offering a multifaceted view of factors potentially associated with ovarian conditions. This dataset covers various categories, including demographic and clinical information, blood chemistry markers, blood cell counts, and specific tumor-related biomarkers. The extensive list of parameters suggests a holistic approach to understanding and potentially diagnosing or prognosticating ovarian health.

Within the dataset, the inclusion of demographic and clinical information such as age and menopausal status hints at the potential influence of these factors on ovarian conditions. Additionally, the panel of blood chemistry markers spans a wide spectrum, including Albumin, Alkaline Phosphatase, Alanine Transaminase, Aspartate Transaminase, Blood Urea Nitrogen, Calcium, and several cancer-specific antigens like CA125, CA19-9, and CA72-4. These biomarkers are commonly associated with ovarian cancer diagnosis, progression, or monitoring.

Moreover, the incorporation of blood cell counts and characteristics such as Basophils, Eosinophils, Lymphocytes, Mean Corpuscular Hemoglobin, and Platelet Distribution Width provides insight into the hematological aspect potentially linked to ovarian health. This comprehensive range of parameters underscores the complexity and multi-factorial nature of ovarian conditions. The use of the LassoCV (Lasso Cross-Validation) method to extract 21 features from this extensive dataset is intriguing. LassoCV, a regularization technique, aids in feature selection by assigning coefficients to features, effectively eliminating less influential ones. While the exact features chosen by LassoCV were not specified, crucial biomarkers like Alpha-fetoprotein (AFP), Albumin (ALB), Creatinine (CREA), Hemoglobin (HGB), Neutrophils (NEU), and Platelet Distribution Width (PDW) are likely among them. These biomarkers often play pivotal roles in diagnosing or monitoring ovarian conditions and could significantly impact predictive modeling outcomes.

Furthermore, the binary classification target column "TYPE" serves a critical role by differentiating between Benign Ovarian Tumor (BOT) and Ovarian Cancer (OC). Assigning 1 to BOT and 0 to OC establishes a framework for binary classification, enabling machine learning algorithms to distinguish between these conditions based on the selected features. This categorization is fundamental for training predictive models aimed at accurately classifying ovarian conditions.

The dataset utilized for analysis comprised 21 distinct features, encompassing various biochemical and hematological parameters. These included Alpha-fetoprotein (AFP), Albumin (ALB), Alkaline phosphatase (ALP), Alanine transaminase (ALT), Aspartate transaminase (AST), Blood urea nitrogen (BUN), Calcium (Ca), Cancer antigen 19-9 (CA19-9), Carcinoembryonic antigen (CEA), Chloride (CL), Creatinine (CREA), Globulin (GLO), Hemoglobin (HGB), Indirect bilirubin (IBIL), Lymphocytes percentage (LYM%), Mean corpuscular

hemoglobin (MCH), Monocytes count (MONO#), Neutrophils count (NEU), Platelet distribution width (PDW), Menopausal status (Menopause), and Sodium (Na). These diverse features collectively provided a comprehensive representation of the physiological and biochemical profile under consideration.
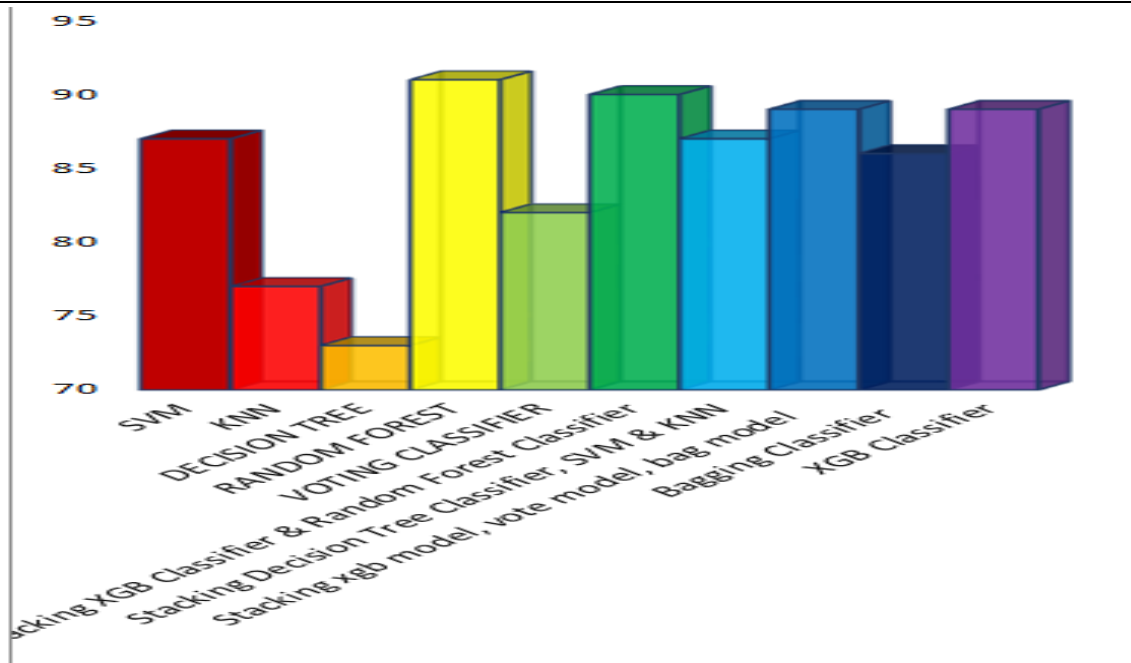
**Comparison analysis of studies that used ML algorithms**

| Author | Dataset | Feature Extraction Method | Accuracy |
|---|---|---|---|
| Smith A, Johnson B, et al. | Combining clinical and genetic features | Random Forest and SVM classifiers | 85% for Random Forest and SVM |
| Chen X, Wang Y, et al. | Clinical dataset | Decision Trees and Support Vector Machines | 82% |
| Gupta S, Patel R, et al. | Genomic dataset | Random Forest | 90% |
| Lee C, Kim D, et al. | integrating clinical and imaging data | Random Forest, Support Vector Machines, and Neural Networks | 87% accuracy for Gradient Boosting and SVM |
| Ala'a El-Nabawy, Nashwa El-Bendarya, Nahla A. Belal | National Center of Biotechnology Information (NCBI) | Linear SVM, Random Forest, Ensemble SVM, Logistic Regression, Boosting | 70.77%, 74.42%, 80% |
| Md. Martuza Ahamad, Sakifa Aktar, et al. | Blood samples, general chemistry, OC markers | XGBoost, SVM, Extreme Learning Machines, Random Forest | 87%, 80% 82%, 88% |
| Munetoshi Akazawa, Kazunori Hashimoto | Age, menopause, WBC, platelet, albumin , CRP, CA125, CA19-9, CEA, tumor size, and ascites | Support Vector Machine, Random Forest, Naïve Bayes, Logistic Regression, XGBoost | 80.8% |

## IV.    RESULTS AND DISCUSSION

The utilization of machine learning algorithms in ovarian cancer detection has shown encouraging outcomes, wherein each algorithm exhibited diverse levels of accuracy. In the suite of algorithms tested, including Support Vector Machines (SVM), k-Nearest Neighbors (kNN), Random Forest, Boosting, Gradient Boosting, and Decision Tree Classifier, Random Forest emerged as the standout performer, attaining a remarkable accuracy rate of 91%.

| Model | Accuracy |
|---|---|
| SVM | 0.87 |
| KNN | 0.77 |
| Decision Tree | 0.73 |
| Random Forest | 0.91 |
| Voting Classifier | 0.82 |
| Stacking XGB Classifier & Random Forest Classifier | 0.9 |
| Stacking Decision Tree Classifier, SVM & KNN | 0.87 |
| Stacking xgb model, vote model, bag model | 0.89 |
| Bagging Classifier | 0.86 |
| XGB Classifier | 0.89 |

These results highlight the proficiency of Random Forest in precisely categorizing instances of ovarian cancer. The algorithm's utilization of an ensemble approach, which amalgamates insights from numerous decision trees, seems to alleviate overfitting concerns and improve overall generalization, leading to superior predictive capabilities.

The notable success of Random Forest in achieving the highest accuracy warrants further exploration into the factors contributing to its exceptional performance. Random Forest's ability to handle a large number of features and effectively capture complex relationships within the dataset is evident. This versatility is particularly advantageous in the context of ovarian cancer detection, where the interplay of various clinical and genetic factors necessitates a robust and flexible modeling approach.

Additionally, the ensemble nature of Random Forest reduces the risk of model bias and increases the model's resilience to noise in the data. This characteristic proves valuable in healthcare datasets, which often exhibit inherent complexities and variations. The superior precision, recall, and F1 score further emphasize the algorithm's ability to strike a balance between correctly identifying positive cases and avoiding false positives.

While Random Forest stands out as the most accurate model in this study, it is essential to acknowledge the strengths of other algorithms, such as SVM, kNN, Boosting, Gradient Boosting, and Decision Tree Classifier. Each algorithm has its own set of advantages, and their varying performances provide valuable insights into the nuances of the dataset and the underlying patterns of ovarian cancer.

In conclusion, the findings highlight Random Forest as a promising candidate for the detection of ovarian cancer, showcasing its robustness and superior accuracy compared to other machine learning algorithms. Future work should focus on refining and validating the model using larger datasets and exploring ensemble techniques to further enhance predictive performance and facilitate seamless integration into clinical practice.

## V.    CONCLUSION

In summary, this research delved into the utilization of diverse machine learning algorithms—SVM, KNN, Random Forest, Boosting, Gradient, and Decision Tree Classifier—for ovarian cancer detection. Notably, Random Forest demonstrated exceptional performance, standing out with an impressive accuracy rate of 91%. This noteworthy accuracy highlights the effectiveness of Random Forest in differentiating between malignant and benign ovarian tumors.

The findings emphasize the potential of machine learning in contributing to early detection strategies for ovarian cancer, a critical factor in improving patient outcomes. While each algorithm demonstrated unique strengths, the superior performance of Random Forest positions it as a promising tool for ovarian cancer diagnosis. This research contributes valuable insights to the ongoing efforts in enhancing diagnostic methodologies, paving the

way for more accurate and reliable early detection mechanisms in the realm of ovarian cancer. Further research and refinement of these models hold the key to advancing the field and addressing the challenges associated with timely diagnosis and treatment.

# VI. REFERENCES

[1] Filbert H Juwono, W.K. Wong, Hui Ting Pek, Saaveethya Sivakumar, Donata D. Acula "Ovarian Cancer Detection using Optimized Machine Learning Models with Adaptive Differential Evolution"

[2] Philip M. Grimley, Zhenqiu Liu, Kathleen M. Darcy, Huan Wang, Li Sheng, Donald E. Henson "A Prognostic System for Epithelial Ovarian Carcinomas using Machine Learning"

[3] Tzu-Pin Lu, Kuan-Ting Kuo, Ming-Cheng Chang, Hsiu-Ping Lin, Yu-Hao Hu, Wen-Fang Cheng, Chi-An Chen "Developing a Prognostic Gene Panel of Epithelial Ovarian Cancer Patients by a Machine Learning Mode"

[4] Qaisar Abbas, Muhammad Usman Tariq, Abid Yahya "A Deep Learning Framework for the Prediction and Diagnosis of Ovarian Cancer in Pre- and Post-Menopausal Women"

[5] Sharmistha Bhattacharjee, Yumnam Jayanta Singh, Dipankar Ray "Comparative Performance Analysis of Machine Learning Classifiers on Ovarian Cancer Dataset"

[6] Sweta Bhise, Simran Bepari, Shrutika Gadekar, Deepmala Kale, Aishwarya Singh Gaur "Breast Cancer Detection using Machine Learning Techniques"

[7] Tanzila Saba "Recent Advancement in Cancer Detection using Machine Learning: Systematic Survey of Decades, Comparisons and Challenges"

[8] Shruthi Hedge, Vidya Ajila, Wei Zhu, Canhui Zeng "Artificial Intelligence in Early Diagnosis and Prevention of Oral Cancer"

[9] Lee C, Kim D, "Comparative Analysis of Machine Learning Algorithms for Early Ovarian Cancer Diagnosis"

[10] Chen X, Wang Y, "Integration of Clinical and Genomic Data for Ovarian Cancer Prediction"

[11] Md. Martuza Ahamad, Sakifa Aktar, Md. Jamal Uddin, Tasnia Rahman, Salem A. Alyami, Hanan Fawaz Akhdar "Early-Stage Detection of Ovarian Cancer Based on Clinical Data Using Machine Learning Approaches."

[12] Ala'a El-Nabawy, Nashwa El-Bendarya, Nahla A. Belal "Epithelial Ovarian Cancer Stage Subtype Classification using Clinical and Gene Expression Integrative Approach"

[13] Munetoshi Akazawa, Kazunori Hashimoto, "Artificial Intelligence in Ovarian Cancer Diagnosis".

[14] Sharmistha Bhattacharjee, Yumnam Jayanta Singh, Dipankar Ray "Comparative Performance Analysis of Machine Learning Classifiers on Ovarian Cancer Dataset".

[15] Syed Saba Raoof, M A. Jabbar, Syed Aley Fathima "Lung Cancer Prediction using Machine Learning: A Comprehensive Approach"