**A PRELIMENERY REPORT ON**


# CANCER CLASSIFICATION AND DETECTION


SUBMITTED TO THE VISHWAKARMA INSTITUTE OF INFORMATION
TECHNOLOGY, PUNE
IN THE PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE AWARD OF THE DEGREE

OF


# BACHELOR OF TECHNOLOGY (AI & DS ENGINEERING)


**SUBMITTED BY**

| | |
|---|---|
| **APURVA BELSARE** | **Exam Seat No.:22110511** |
| **ADITI DEVADE** | **Exam Seat No.:22110354** |
| **PRADNYA KEDARI** | **Exam Seat No.:22110926** |
| **AARYA UMBARE** | **Exam Seat No.:22111339** |

# DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND DATA SCIENCE


**BRACT'S**
**VISHWAKARMA INSTITUTE OF INFORMATION TECHNOLOGY**

**SURVEY NO. 3/4, KONDHWA (BUDRUK), PUNE – 411048,
MAHARASHTRA (INDIA).
BRACT'S**

**VISHWAKARMA INSTITUTE OF INFORMATION TECHNOLOGY, PUNE**
**2024 -2025**



# CERTIFICATE

This is to certify that the project report entitles

**CANCER CLASSIFICATION AND DETECTION**

Submitted by

| | |
|---|---|
| **APURVA BELSARE** | **Exam Seat No.:22110511** |
| **ADITI DEVADE** | **Exam Seat No.:22110354** |
| **PRADNYA KEDARI** | **Exam Seat No.:22110926** |
| **AARYA UMBARE** | **Exam Seat No.:22111339** |

is a bonafide student of this institute and the work has been carried out by him/her under the supervision of **Prof. Yashwant Ingle** and it is approved for the partial fulfillment of the requirement of VISHWAKARMA INSTITUTE OF INFORMATION TECHNOLOGY, for the award of the degree of **Bachelor of Technology** (Artificial Intelligence and Data Science).

**Prof. Yashwant Ingle**                                          **Prof. Santosh Kumar**
Guide                                                                          Head,
Department of AI & DS                                       Department of AI & DS

**Dr. Vivek Deshpande**
Director,
BRACT's Vishwakarma Institute of Information Technology, Pune-48

Place : Pune
Date :

ii

# ACKNOWLEDGEMENT

**AARYA UMBARE**      **22111339**
**APURVA BELSARE**    **22110511**
**ADITI DEVADE**      **22110354**
**PRADNYA KEDARI**    **22110926**

# ABSTRACT

Ovarian cancer, often diagnosed in its advanced stages due to its ambiguous symptoms, ovarian cancer is frequently detected in its advanced stages, making early identification and effective treatment extremely difficult. In order to transform the diagnosis of ovarian cancer, this initiative uses cutting-edge data science techniques to create machine learning models. The study uses methods like Principal Component Analysis (PCA) for dimensionality reduction and Synthetic Minority Over-Sampling Technique (SMOTE) to address class imbalance after evaluating a large dataset of gene expression levels and tumor subtypes.

K-Nearest Neighbors (KNN), XGBoost, Decision Trees, Logistic Regression, and Support Vector Machines (SVM) were among the machine learning algorithms that were trained and thoroughly assessed for classification performance, accuracy, and F1 scores. KNN was found to be the most successful algorithm for this challenge based on the findings. Bar graphs and scatter plots were among the visualizations that shed light on model performance and gene expression patterns. This study emphasizes how data-driven methods can help with precision medicine, increase early diagnosis, and eventually lead to better patient outcomes. The models and insights that have been established provide a strong basis for future use in personalized healthcare solutions and ovarian cancer detection.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# 01.INTRODUCTION

## 1.1 OVERVIEW

Ovarian cancer is a type of cancer that originates in the ovaries, the reproductive organs in women responsible for producing eggs and hormones. It is a significant health concern, with various cell types actively dividing in the ovaries, leading to the development of tumors. This form of cancer is challenging to detect in its preliminary stages due to a lack of specific symptoms. As a result, it often goes unnoticed until it reaches an advanced stage, making effective treatment more complex.

Epithelial ovarian cancer is the most common type, and risk factors include age, family history, and inherited gene mutations like BRCA1 and BRCA2. Symptoms may include abdominal bloating, pelvic pain, difficulty eating, and frequent urination.Treatment typically involves surgical removal of the tumor, and in some cases, the ovaries, and surrounding tissues.

The prognosis for ovarian cancer depends on factors such as the stage at diagnosis, the type of cancer, and the individual's overall health. Early detection is crucial for better outcomes, but currently, there are challenges in developing effective screening methods. Research efforts are ongoing to enhance our understanding of ovarian cancer, improve early detection methods, and develop more effective treatments. Increased awareness and support for ovarian cancer research are essential in addressing this complex and often late-diagnosed disease.

## 1.2 MOTIVATION

Selecting ovarian cancer as a focus for a data science project provides an opportunity to significantly contribute to advancements in healthcare. Several key motivations support this choice:

- Fostering Early Detection: A primary goal is to tackle the issue of late-stage ovarian cancer diagnoses. Leveraging data science techniques enables the analysis of diverse datasets to identify subtle patterns that indicate early-stage ovarian tumors, ultimately improving patient outcomes.

- Utilizing Diverse Datasets: This project emphasizes the wealth of available data, including patient demographics, genetic profiles, imaging studies, and clinical records. By applying advanced data science methodologies, the project aims to uncover meaningful relationships across these data types, offering a comprehensive understanding of ovarian tumors.

- Harnessing Machine Learning for Diagnosis: Machine learning has proven effective in identifying disease patterns within extensive datasets. Utilizing these algorithms for ovarian cancer detection can enhance diagnostic precision, assisting healthcare professionals in making more accurate and timely decisions.

- Empowering Clinical Decision-Making: Beyond building models, the project seeks to apply insights to clinical practice by decoding complex data relationships. The ultimate aim is to provide actionable insights that enable healthcare providers to make better-informed decisions regarding patient care and treatment strategies.

- Contributing to Scientific Progress: Ovarian cancer continues to be a challenging area for research. Data-driven approaches can generate valuable insights, uncovering novel correlations and advancing scientific knowledge. This project has the potential to foster new research opportunities and contribute to the collective understanding of ovarian cancer.

In conclusion, this project aligns with the broader objectives of improving early detection, advancing personalized medicine, and addressing one of the most complex and late-diagnosed diseases, thus creating meaningful impacts in both research and clinical settings.

## 1.3 PROBLEM DEFINITION AND OBJECTIVE

### 1.3.1  Problem Definition

Ovarian cancer is often diagnosed at an advanced stage due to the lack of effective early detection mechanisms. Because there are insufficiently efficient early detection methods, ovarian cancer is frequently discovered at an advanced stage. This late-stage detection poses a serious healthcare dilemma because it drastically lowers survival rates. By using data science and machine learning approaches to find patterns and markers of early-stage ovarian cancer, the initiative aims to address this problem.

### 1.3.2  Objective

The primary objective of this project is to develop a machine learning model that can analyze diverse datasets, including patient demographics, genetic information, and clinical records, to enhance the accuracy of early detection. This project's main goal is to create a machine learning model that can evaluate a variety of datasets, such as clinical records, genetic data, and patient demographics, in order to improve the precision of early detection. Additionally, the project seeks to give medical professionals useful information to help them make well-informed decisions and increase scientific knowledge about ovarian cancer by using data-driven investigation. Help achieve better healthcare results and the more general objectives of individualized treatment.

## 1.4 PROJECT SCOPE AND LIMITATION

### 1.4.1 Project Scope

The scope of this project revolves around leveraging advanced machine learning techniques to classify ovarian cancer subtypes based on gene expression data. By integrating comprehensive data preprocessing, dimensionality reduction using PCA, and oversampling to address class imbalances, the project ensures robust data handling. Multiple machine learning algorithms, including KNN, SVM, XGBoost, and Decision Trees, are evaluated to identify the most accurate and reliable model for subtype prediction. The project extends beyond classification, providing insights into significant genes contributing to cancer subtypes, potentially uncovering biomarkers for early diagnosis. The ultimate goal is to develop a scalable, data-driven tool that aids in the precise classification of ovarian tumors, enhancing diagnostic accuracy and enabling personalized therapeutic interventions. This initiative contributes to the broader field of bioinformatics, with potential applications in early cancer detection and precision medicine.

### 1.4.2 Limitations

- Data Constraints:

The initiative depends on the provision of diverse and high-quality datasets. The model's generalizability may be impacted by biased or incomplete datasets.

- Medical Data Complexity:

In feature engineering and model training, integrating diverse data types—such as genetic, imaging, and clinical data—presents difficulties.

- Generalizability:

Due to variations in patient demographics or healthcare systems, models created on particular datasets may not function well on unobserved data.

- Resource Requirements:

Scalability may be constrained by the high computational requirements for processing big datasets and training intricate models.

- Interpretability:

It's still difficult to explain machine learning model predictions to medical professionals, particularly when the models are black-box.

## 1.5 METHODOLOGIES OF PROBLEM SOLVING

The methodology consists of the below major steps:

1.Data Collection

In the data collection phase, the main objective is to gather relevant datasets essential for the analysis or problem at hand. This involves acquiring data from various sources, which can include databases, APIs, online repositories, or other means. The goal is to assemble a comprehensive and representative data set aligning with the project's objectives.

2.Exploratory Data Analysis (EDA)

EDA is a crucial step where the collected data is analyzed and visualized to gain insights and understand the underlying patterns. During EDA, various statistical and graphical methods are employed to explore the characteristics of the dataset. This includes summary statistics, distribution plots, correlation analysis, and data visualization techniques. EDA helps identify trends, outliers, and potential relationships between variables, laying the foundation for informed decision-making.

3.Feature Engineering and Preprocessing

Feature engineering involves transforming or creating new features from the existing ones to enhance the predictive power of the model. This step includes handling missing data, encoding categorical variables, scaling numerical features, and creating interaction terms. Feature preprocessing ensures that the data is in a suitable format for modeling. Techniques such as normalization, standardization,

and handling outliers are applied to prepare the dataset for the machine learning algorithms.

4.Model Building and Evaluation

In the decisive step, machine learning models are developed based on the preprocessed dataset. This involves selecting appropriate algorithms that align with the nature of the problem, such as regression for predicting continuous outcomes or classification for predicting categories. Models are trained on a subset of the data and evaluated using another subset to assess their performance. Evaluation metrics, such as accuracy, precision, recall, F1-score, or regression metrics, provide insights into how well the model generalizes to unseen data. Model tuning and optimization may be performed to enhance performance.

# 02.LITERATURE SURVEY

Multi-Modal Evolutionary Deep Learning Model for Ovarian Cancer Diagnosis by Rania M. Ghoniem, Abeer D. Algarni, Basel Refky, and Ahmed A. Ewees presents a hybrid deep learning system integrating histopathological images and gene expression data for ovarian cancer diagnosis. The model employs Convolutional Neural Networks (CNN) for image analysis and Long Short-Term Memory (LSTM) networks for gene expression data, optimized with the Antlion Optimization (ALO) method. This integration demonstrates the potential for multi-modal data to enhance cancer detection accuracy.

Diagnosing Ovarian Cancer on MRI: A Preliminary Study Comparing Deep Learning and Radiologist Assessments by Tsukasa Saida, Kensaku Mori, Sodai Hoshial, and colleagues compares the diagnostic performance of deep learning models and expert radiologists using MRI data. Employing a deep convolutional neural network, the study reveals that the model achieves diagnostic accuracy comparable to that of skilled radiologists, highlighting the potential of deep learning as a supplementary diagnostic tool.

Early-Stage Detection of Ovarian Cancer Based on Clinical Data Using Machine Learning Approaches by Md. Mehedi Hasan Al Zubair, Shamsun Nahar, Md. Junaid Iqbal, and co-authors explores the use of machine learning to analyze clinical data for early ovarian cancer detection. The study applies multiple algorithms to identify significant biomarkers, achieving high accuracy in distinguishing between benign and malignant tumors, underscoring the role of machine learning in early diagnosis.

Prediction of Ovarian Cancer Using Artificial Intelligence Tools by Seyed Mohammad Ayyoubzadeh, Marjan Ahmadi, and others examines the application of artificial intelligence techniques to predict ovarian cancer. Using tumor markers and blood test results, the study evaluates various models, with the Random Forest achieving the highest accuracy (86%). This work emphasizes the potential of AI in improving diagnosis through precise and cost-effective methods.

Improved Prediction of Ovarian Cancer Using Ensemble Classifier and Shapley Explainable AI by Nihal Abuzinadah, Sarath Kumar Posa, and co-authors introduces a stacked ensemble model for ovarian cancer prediction, integrating boosting and bagging classifiers. Achieving 96.87% accuracy on the Soochow University dataset, the model uses Shapley Explainable AI to enhance interpretability, providing insights into the decision-making process and aiding clinical adoption.

Artificial Intelligence in Ultrasound Diagnoses of Ovarian Cancer: A Systematic Review and Meta-Analysis by Sian Mitchell, Manolis Nikolopoulos, and collaborators systematically reviews AI applications in ovarian cancer diagnosis via ultrasound imaging. Highlighting the effectiveness of convolutional neural networks (CNNs), the study concludes that AI models often outperform traditional methods, promising significant improvements in early detection.

A Deep Learning Framework for the Prediction and Diagnosis of Ovarian Cancer in Pre- and Post-Menopausal Women by Blessed Ziyambe, Abid Yahya, and co-authors proposes a deep learning model using CNNs to analyze medical imaging for ovarian cancer diagnosis. Considering physiological differences

between pre- and post-menopausal women, the study demonstrates high accuracy in distinguishing between benign and malignant conditions, supporting personalized diagnostic strategies.

Ovarian Cancer Beyond Imaging: Integration of AI and Multiomics Biomarkers by Sepideh Hatamikia, Stephanie Nougaret, and others explores AI-driven multiomics approaches for ovarian cancer diagnosis. Combining clinical, imaging, genomic, and transcriptomic data, the study highlights the superior predictive capabilities of integrative models over single-modality approaches. This comprehensive method provides valuable insights for personalized treatments and improved outcomes.

# 03.SYSTEM REQUIREMENTS AND SPECIFICATION

## 3.1 Assumptions and Dependencies

The ovarian cancer diagnosis system assumes the availability of a comprehensive dataset containing gene expression data and clinical features relevant to the study. It also depends on the accessibility of machine learning libraries such as Scikit-learn and TensorFlow for model development. The functionality relies on the compatibility of the selected  software configurations.

## 3.2 Functional Requirements

### 3.2.1 System Feature 1: Data Processing and Analysis

The system should enable the ingestion and preprocessing of gene expression datasets. This includes loading datasets, performing exploratory data analysis (EDA), detecting anomalies, and normalizing the data for machine learning models.

### 3.2.2 System Feature 2: Model Training and Prediction

The system must support the training and evaluation of multiple machine learning models, including K-Nearest Neighbors, Random Forest, and XGBoost, to classify ovarian cancer into subtypes. It should allow hyperparameter tuning and provide performance metrics such as accuracy, precision, recall, and F1-score for model comparison.

## 3.3 External Interface Requirements

### 3.3.1 User Interfaces

A user-friendly interface will allow researchers and clinicians to view results. Visualizations of key insights, such as feature importance and classification outcomes, should be included.

### 3.3.2 Software Interfaces

The system interfaces with Python libraries like pandas, NumPy, Scikit-learn, and Matplotlib. It also uses google colab for coding and visualization.

## 3.4 Nonfunctional Requirements

### 3.4.1 Performance Requirements

The system must provide highly accurate predictions of ovarian cancer subtypes based on gene expression data, with a target accuracy of at least 75% or higher for the optimal machine learning model. It should demonstrate robust performance across all cancer subtypes, including challenging categories such as clear cell and mucinous tumors. Additionally, the system must be efficient, processing and predicting categories within a reasonable time frame to ensure practicality

### 3.4.2 Safety Requirements

Data integrity must be maintained throughout preprocessing and analysis to ensure accurate predictions. Backup mechanisms should be in place to prevent data loss.

### 3.4.3 Software Quality Attributes

The system must be reliable, with minimal downtime, and provide accurate results. It should also be maintainable, allowing easy updates and enhancements to algorithms or datasets.

## 3.5 System Requirements

### 3.5.1 Database Requirements

The system requires a database capable of storing gene expression data, model parameters, and analysis results. This could include relational databases like MySQL or SQLite.

### 3.5.2 Software Requirements (Platform Choice)

The system should run on platforms supporting Python, such as Windows, macOS, or Linux. Essential software includes Python 3.8 or later, Google colab or Jupyter Notebook, and relevant libraries.

## 3.6 Analysis Models: SDLC Model to Be Applied

The system development follows the Iterative Model of the Software Development Life Cycle (SDLC). This approach allows for progressive refinement of requirements and implementation through cycles of design, development, testing, and feedback. It ensures flexibility in accommodating changes and enhances the system's robustness and reliability.

# 04. SYSTEM DESIGN

## 4.1 SYSTEM ARCHITECTURE

**System Architecture for Ovarian Cancer Detection**

**Data Collection**

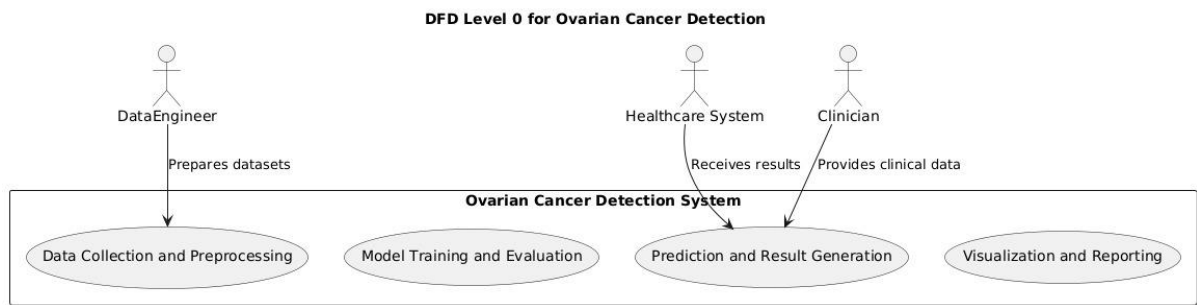Clinical Data Sources | Genetic Data Sources | Imaging Data Sources

**Data Preprocessing**

Data Cleaning

Feature Engineering

**Machine Learning Pipeline**

Model Training

Model Evaluation

**Deployment**

Model Integration

API for Clinical Use

## 4.2 DATA FLOW DIAGRAMS

**DFD Level 0 for Ovarian Cancer Detection**

DataEngineer

Healthcare System    Clinician

Prepares datasets

Receives results    Provides clinical data

**Ovarian Cancer Detection System**

Data Collection and Preprocessing | Model Training and Evaluation | Prediction and Result Generation | Visualization and Reporting

Department of Artificial Intelligence and Data Science, 2024-25

**DFD Level 1 for Ovarian Cancer Detection**

DataEngineer

Provides raw data

**Ovarian Cancer Detection System**

Data Collection

Data Preprocessing

Exploratory Data Analysis

Model Training

Model Evaluation

Clinician

Provides clinical data

Healthcare System

Receives predictions

Prediction Generation

Clinician

**DFD Level 2 for Ovarian Cancer Detection**

Test Model → Evaluate Performance → Store Evaluation Results

Select Features → Train Model → Store Trained Model

Perform Statistical Analysis → Visualize Data Patterns

Clean Data → Handle Missing Values → Normalize Data → Store Processed Data

Clinician

Requests predictions → Generate Predictions → Store Predictions

Views results → Display to Clinician

DataEngineer

Provides clinical datasets → Data Collection

**Ovarian Cancer Detection System**

Data Collection

Model Training

Exploratory Data Analysis

Prediction Generation

Data Preprocessing

Model Evaluation

DataCollectionRaw

Store Raw Data

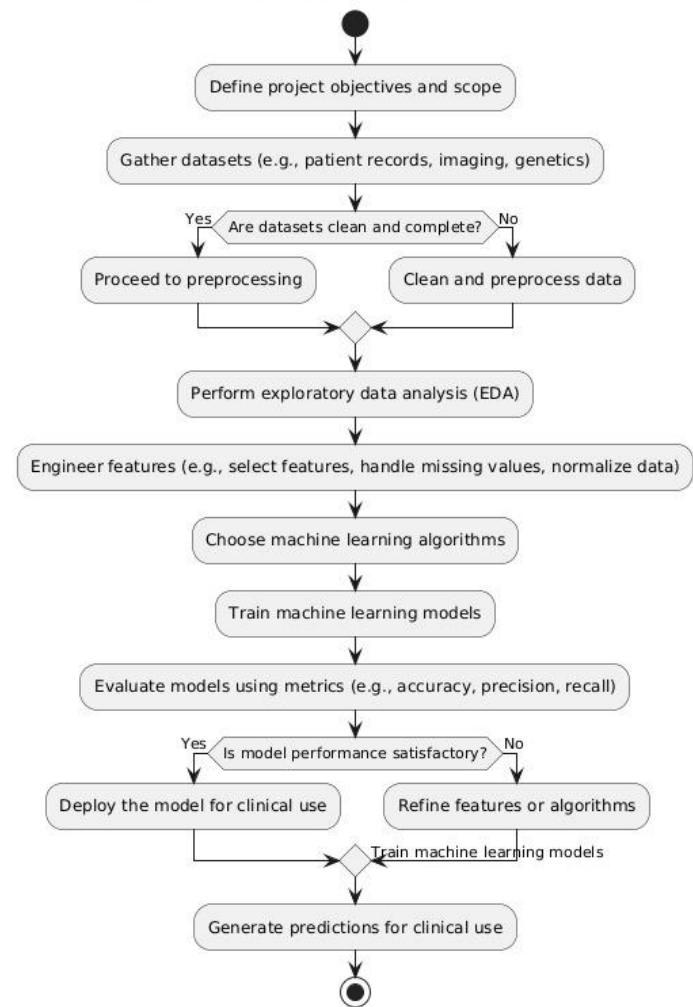Department of Artificial Intelligence and Data Science, 2024-25
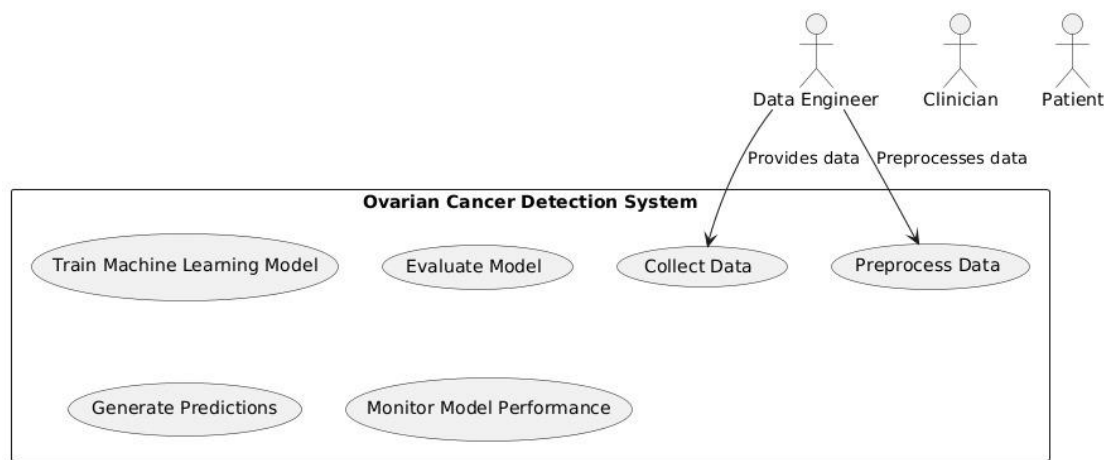
## 4.3 ENTITY RELATIONSHIP DIAGRAM

## 4.4 UML Diagrams

**Activity Diagram for Ovarian Cancer Detection Project**
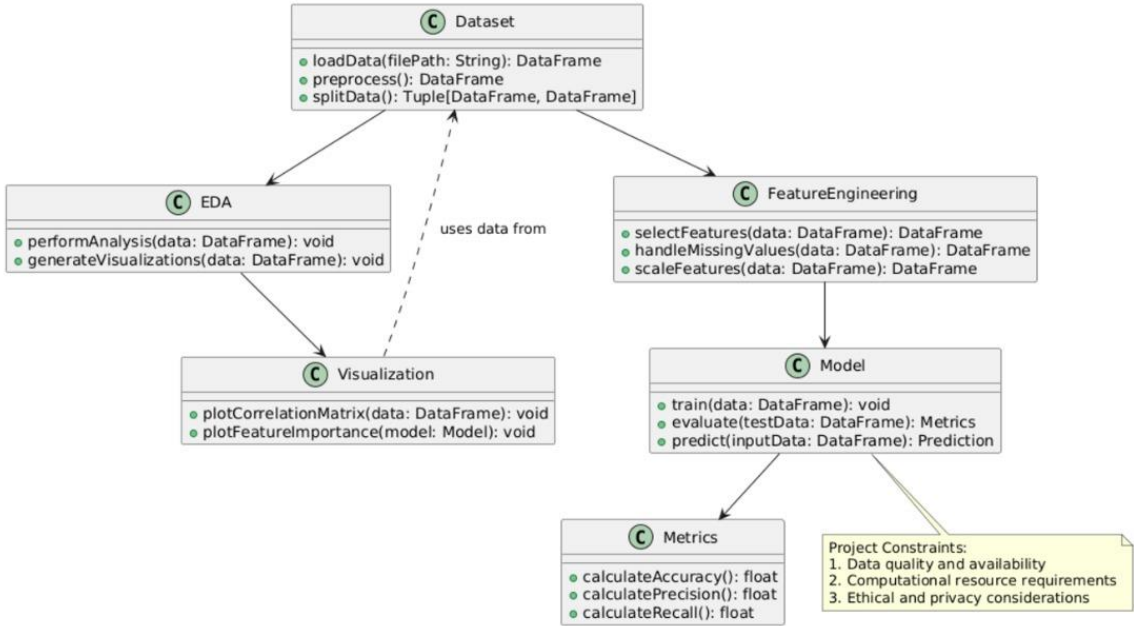


**Use Case Diagram for Ovarian Cancer Detection**



Department of Artificial Intelligence and Data Science, 2024-25

## Ovarian Cancer Classification Process



## Ovarian Cancer Detection Project - Class Diagram



Department of Artificial Intelligence and Data Science, 2024-25

# 05. PROJECT PLAN

## 5.1 Project Estimate

### 5.1.1 Reconciled Estimates

The project requires an estimated timeline of 4 months for completion, broken into distinct phases: data collection and preprocessing, exploratory data analysis (EDA) and visualization, model development and evaluation.

### 5.1.2 Project Resources

Key software resources include Python libraries such as NumPy, Pandas, Scikit-learn, and TensorFlow for data manipulation, modeling, and analysis. Jupyter Notebook serves as the primary environment for developing and documenting workflows, while visualization tools like Matplotlib and Seaborn are employed to create clear and insightful graphical representations of the data. These tools collectively provide a robust framework for building, evaluating, and interpreting predictive models for ovarian cancer diagnosis.

## 5.2 Risk Management

### 5.2.1 Risk Identification

Key risks include:
- Limited availability or quality of gene expression data.
- Model overfitting due to small or imbalanced datasets.
- Computational challenges with large datasets or complex models.

5.2.2 Risk Analysis

The project faces several risks, including data quality issues, such as inconsistent or noisy data, which could result in inaccurate conclusions. Performance challenges may arise if the models underperform, potentially affecting the reliability of predictions. Additionally, scalability concerns may emerge as the volume of data increases, potentially exceeding the computational capacity available for processing.

5.2.3 Overview of Risk Mitigation, Monitoring, and Management

To address these risks, oversampling techniques and cross-validation methods are used to handle imbalanced datasets effectively. Regular monitoring of model performance is conducted using robust evaluation metrics to ensure accuracy and reliability. Frequent reviews and continuous testing are integral to identifying and resolving issues early, ensuring smooth project progression.

**5.3 Project Schedule**

5.3.1 Project Task Set

The project encompasses several critical tasks to achieve its objectives effectively. These tasks include data acquisition, where relevant gene expression datasets are gathered from reliable sources. This is followed by preprocessing, which involves cleaning, normalizing, and transforming the data to ensure it is suitable for analysis. Exploratory analysis is then performed to understand patterns and relationships within the data. Feature selection is a crucial step to identify significant genes contributing to ovarian cancer diagnosis. The selected features are used for model training, where various machine learning algorithms are applied and optimized. Finally, the models are evaluated based on their accuracy and reliability, and all findings are documented .

5.3.2 Task Network

The project follows a structured task network with defined dependencies to maintain a logical and efficient workflow. For instance, data preprocessing must be completed before initiating feature selection, as clean and structured data is essential for identifying relevant features. Feature selection, in turn, feeds into the model training process, as only significant features are used to build predictive models. Model evaluation, which assesses the performance of trained models, depends on the completion of these preceding tasks.

5.3.3 Timeline Chart

| Week | Task |
|---|---|
| Week 1-2 | Data Collection and Initial Data Exploration |
| Week 3-4 | Data Preprocessing and Cleaning |
| Week 5-6 | Data Analysis and Feature Selection |
| Week 7-8 | Exploratory Data Analysis (EDA) |
| Week 9-10 | Model Selection and Training |
| Week 11-12 | Model Tuning and Optimization |
| Week 13-14 | Model Evaluation and Metrics Calculation |

**5.4 Team Organization**

5.4.1 Team Structure

Internal Guide: Prof. Yashwant Ingle
Our team:

Apurva Belsare.

Pradnya Kedari

Aditi Devade

Aarya Umbhare

## 5.4.2 Management Reporting and Communication

Regular team meetings and progress reports are conducted weekly. A shared online workspace ensures seamless communication and document sharing. Additionally, periodic reviews are conducted to assess the alignment of the project with its goals and identify areas for improvement.

# 06.PROJECT IMPLEMENTATION

## 6.1 Overview of Project Modules

The project is organized into distinct modules, each addressing a specific aspect of the ovarian cancer classification process. The Data Acquisition and Preprocessing Module handles the integration of gene expression datasets, ensuring data quality and preparation for analysis. The Exploratory Data Analysis (EDA) Module offers insights into the dataset's structure and characteristics through visualization techniques. The Feature Engineering Module focuses on dimensionality reduction and oversampling to handle imbalanced datasets. The Model Development Module implements machine learning algorithms for classification, while the Evaluation Module assesses the model's performance using metrics like accuracy, precision, recall, and F1-score.

## 6.2 Tools and Technologies Used

The project leverages Python programming for implementation, utilizing libraries such as NumPy and Pandas for data manipulation, Matplotlib and Seaborn for visualization, and Scikit-learn for machine learning. Dimensionality reduction is performed using Principal Component Analysis (PCA), and imbalanced datasets are addressed through oversampling techniques like SMOTE. Development and experimentation are conducted in Jupyter Notebook, ensuring a seamless, interactive environment for analysis and coding.

## 6.3 Algorithm Details

### 6.3.1 Algorithm 1: K-Nearest Neighbors (KNN)

The KNN algorithm classifies data points based on their proximity to other labeled points. It calculates the distance between data points using metrics like Euclidean

distance and assigns the most frequent class among the nearest neighbors. For this project, parameters such as n_neighbors and weights were tuned to achieve an accuracy of 78.19%.

### 6.3.2 Algorithm 2: XGBoost

XGBoost is an ensemble learning algorithm based on gradient boosting, designed for high-performance classification and regression tasks. It builds a series of decision trees iteratively, with each tree correcting the errors of its predecessor. For this project, parameters such as n_estimators, max_depth, and learning_rate were optimized, resulting in an accuracy of 73.04%.

### 6.3.3 Algorithm 3: Random Forest

Random Forest is another ensemble learning method that constructs multiple decision trees and merges their outputs for better accuracy and stability. This project tuned parameters like n_estimators and max_depth, achieving an accuracy of 68.92%.

### 6.3.4 Algorithm 4: Logistic Regression

Logistic Regression is a simple yet effective algorithm that uses a logistic function to model the probability of class membership. It is particularly useful for binary classification tasks. In this project, feature scaling and regularization techniques like L2 regularization were applied to improve performance. Logistic Regression achieved an accuracy of 65.41%.

### 6.3.5 Algorithm 5: Support Vector Machine (SVM)

SVM is a powerful algorithm that finds an optimal hyperplane to separate data points in a high-dimensional space. The RBF kernel was used in this project to handle non-linear relationships between features. Parameters like C (regularization) and gamma (kernel coefficient) were fine-tuned, resulting in an accuracy of 69.25%.

# 07.SOFTWARE TESTING

7.1 Type of Testing

To ensure the reliability and accuracy of the ovarian cancer classification system, the following types of testing were conducted:

1. Unit Testing: Individual components of the project, such as data preprocessing, feature extraction, and model training, were tested to validate their functionality and correctness.

2. Integration Testing: Modules such as data preprocessing, model training, and evaluation were tested in combination to ensure smooth data flow and proper integration of components.

3. Performance Testing: The system was tested on large datasets to evaluate its computational efficiency, response time, and scalability.

4. Accuracy Testing: The classification models were tested for accuracy, precision, recall, and F1-score to ensure reliable predictions.

5. Stress Testing: The system was tested under heavy computational loads to evaluate its robustness and error-handling capabilities.

6. Regression Testing: Repeated testing was performed after updates or optimizations to ensure that changes did not introduce new errors.

## 7.2 Test Cases & Test Results

| Test Case ID | Test Case Description | Input | Expected Output | Actual Output | Status |
|---|---|---|---|---|---|
| TC-01 | Validate data preprocessing module | Raw gene expression dataset | Cleaned and normalized dataset | Cleaned and normalized dataset | Pass |

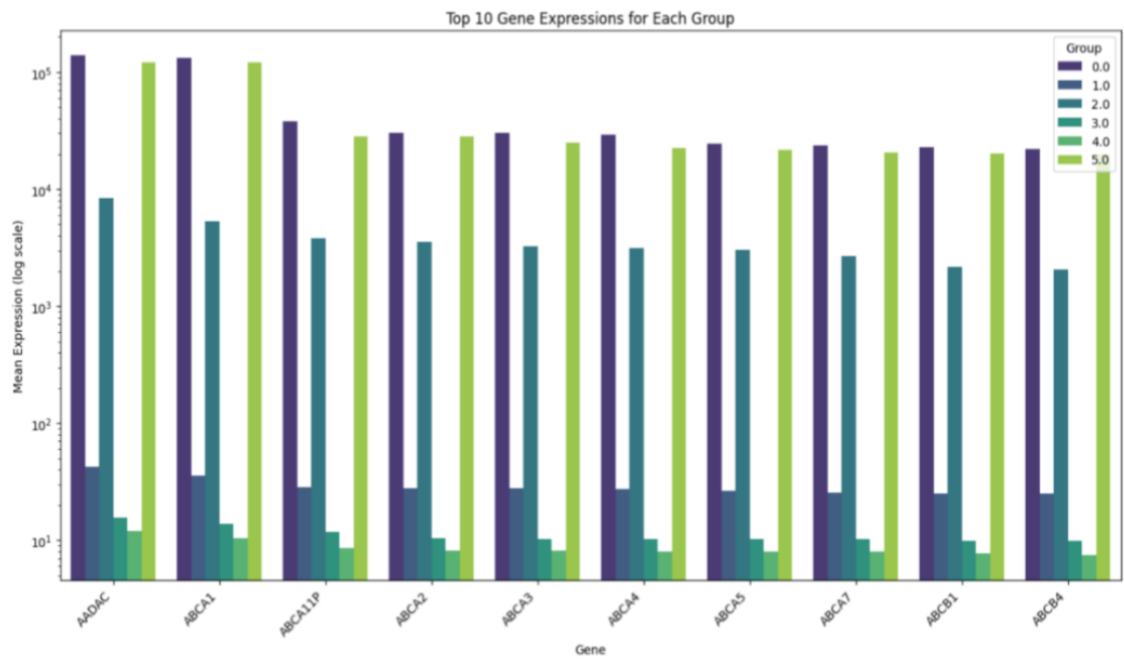| Test Case ID | Test Case Description | Input | Expected Output | Actual Output | Status |
|---|---|---|---|---|---|
| TC-02 | Check feature selection process | Processed dataset | Reduced feature set | Reduced feature set | Pass |
| TC-03 | Validate model training for KNN | Training data | Trained KNN model | Trained KNN model | Pass |
| TC-04 | Test prediction for Random Forest | Sample test data | Predicted categories | Predicted categories | Pass |
| TC-05 | Ensure accuracy evaluation | Model predictions and labels | Accuracy metric | Correct accuracy value | Pass |
| TC-06 | Stress test for scalability | Large dataset | Successful execution | Successful execution | Pass |
| TC-07 | Validate user interface response | User input (gene data) | Predicted cancer category | Predicted cancer category | Pass |

.

# 08.RESULTS

## 8.1 Outcome

An interactive bar chart was created to compare the accuracies of various machine learning models, including Logistic Regression, KNN, SVM, Decision Tree, and XGBoost. Each bar in the chart represents the accuracy of a specific model, with accuracy values displayed directly on the bars for easy reference. This visual overview highlights the performance of each model, making it simple to identify which ones are more effective for the given task.
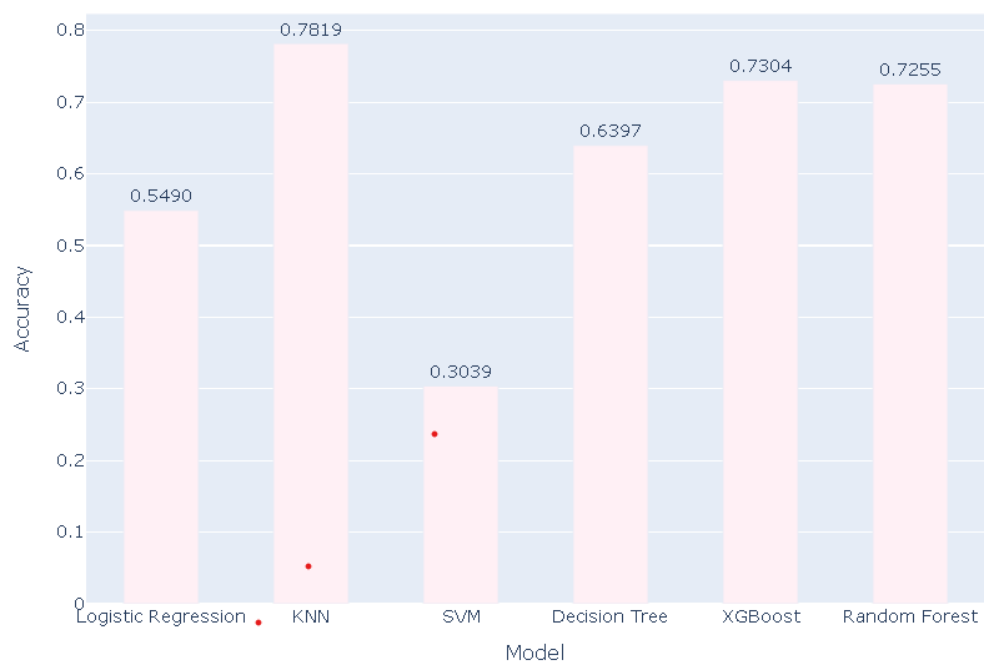
After extensive research and testing, KNN was determined to be the most optimal algorithm for this dataset. This decision was based on its superior test accuracy, F1 scores in the classification report, and robust performance reflected in the confusion matrix. The following table provides a detailed comparison of the models, offering insights into the criteria used to identify the best-performing algorithm.

## 8.2 Screen Shots

Top 10 Gene Expressions for Each Group



Model Comparison - Accuracy

# 09.CONCLUSION

## 9.1 Conclusion

This project aimed at classifying cancer subtypes through an extensive analysis of gene expression data. We meticulously collected, preprocessed, and explored diverse datasets, ensuring data accuracy. Applying Principal Component Analysis (PCA) facilitated dimensionality reduction, and machine learning models, including XGBoost, Decision Trees, Logistic Regression, K-Nearest Neighbors, and Support Vector Machines, were rigorously evaluated.

The models underwent thorough assessment, considering accuracy metrics, classification reports, and confusion matrices. Visualization techniques provided additional insights.

A final model comparison highlighted the strengths of each algorithm, aiding in the selection of the most effective approach for cancer subtype classification. This comprehensive journey ensures a robust foundation for future applications in leveraging gene expression data for cancer classification.

## 9.2 Future Scope

The proposed ML-based ovarian cancer detection system has significant potential for expansion and refinement. To increase accuracy even more, future advancements will use cutting-edge AI methods like deep neural networks and ensemble learning. Comprehensive predictive modeling and improved diagnosis accuracy can be achieved

by integrating data from other sources, such as proteomic, genomic, and imaging data.

Furthermore, adding worldwide and varied samples to the dataset will guarantee the model's resilience across a range of demographics. Healthcare practitioners will be able to make decisions more quickly with real-time deployment in clinical settings that includes automated reporting capabilities. Lastly, implementing explainable AI techniques will improve forecast transparency and foster patient and clinician trust.

## 9.3 Applications

The outcomes of this project have broad applications in the medical field, particularly in early detection and personalized treatment planning for ovarian cancer. The models developed can assist oncologists in identifying high-risk patients and determining precise treatment strategies. The methodology can also be adapted for studying other cancer types or diseases that involve complex gene interactions. Furthermore, the integration of such predictive models into healthcare systems could streamline diagnostic workflows, improve patient outcomes, and reduce the overall burden on healthcare infrastructure.

# APPENDIX

## 25% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

### Match Groups

**96 Not Cited or Quoted 25%**
Matches with neither in-text citation nor quotation marks

**0 Missing Quotations 0%**
Matches that are still very similar to source material

**0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation

**0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

### Top Sources

14%  ⊕ Internet sources

17%  📖 Publications

18%  👤 Submitted works (Student Papers)

### Integrity Flags

**0 Integrity Flags for Review**

No suspicious text manipulations found.

Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.

# REFERENCES

- Multi-Modal Evolutionary Deep Learning Model for Ovarian Cancer Diagnosis
  Rania M. Ghoniem, Abeer D. Algarni, Basel Refky, Ahmed A. Ewees.
- Diagnosing Ovarian Cancer on MRI: A Preliminary Study Comparing Deep Learning and Radiologist Assessments
  Tsukasa Saida, Kensaku Mori, Sodai Hoshial, Masafumi Sakai, Aiko Urushibara, Toshitaka Ishiguro, Manabu Minami, Toyomi Satoh, Takahito Nakajima.
- Early-Stage Detection of Ovarian Cancer Based on Clinical Data Using Machine Learning Approaches
  Md. Mehedi Hasan Al Zubair, Shamsun Nahar, Md. Junaid Iqbal, Md. Asif Hossain, Md. Mahfuzur Rahman, Md. Shafiqul Islam, Md. Rashedul Islam, Md. Abdur Rahman.
- Prediction of Ovarian Cancer Using Artificial Intelligence Tools
  Seyed Mohammad Ayyoubzadeh, Marjan Ahmadi, Alireza Banaye Yazdipour, Fatemeh Ghorbani-Bidkorpeh, Mahnaz Ahmadi.
- Improved Prediction of Ovarian Cancer Using Ensemble Classifier and Shapley Explainable AI
  Nihal Abuzinadah, Sarath Kumar Posa, Aisha Ahmed Alarfaj, Ebtisam Abdullah Alabdulqader, Muhammad Umer, Tai-Hoon Kim, Shtwai Alsubai, Imran Ashraf.
- Artificial Intelligence in Ultrasound Diagnoses of Ovarian Cancer: A Systematic Review and Meta-Analysis
  Sian Mitchell, Manolis Nikolopoulos, Alaa El-Zarka, Dhurgham Al-Karawi, Shakir Al-Zaidi, Avi Ghai, Jonathan E. Gaughran, Ahmad Sayasneh.
- A Deep Learning Framework for the Prediction and Diagnosis of Ovarian Cancer in Pre- and Post-Menopausal Women
  Blessed Ziyambe, Abid Yahya, Tawanda Mushiri, Muhammad Usman Tariq, Qaisar Abbas, Muhammad Babar, Mubarak Albathan, Ayyaz Hussain, Sohail Jabbar.
- Ovarian Cancer Beyond Imaging: Integration of AI and Multiomics Biomarkers
  Sepideh Hatamikia, Stephanie Nougaret, Camilla Panico, Giacomo Avesani, Camilla Nero, Luca Boldrini, Evis Sala, Ramona Woitek.