# VISVESVARAYA TECHNOLOGICAL UNIVERSITY

## 'JNANA SANGAMA' BELAGAVI-590018, KARNATAKA



### MINI-PROJECT REPORT

### ON

## "CANCER DETECTION USING MACHINE LEARNING"

**SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENT
FOR THE V SEMESTER, BE, MINI-PROJECT-BCS586**

**BACHELOR OF ENGINEERING
IN
COMPUTER SCIENCE & ENGINEERING**

### Submitted By

1. **VARSHA P N**
   **[1CG22CS119]**

2. **TRISHA M**
   **[1CG22CS116]**

3. **SUSHMITHA M R**
   **[1CG22CS114]**

4. **SHIVANI S R**
   **[1CG22CS104]**

### Under the guidance of:

**Mr. Harish T A** MTech.,
Asst. Prof., Dept. of CSE,
CIT, Gubbi, Tumakuru .

### HOD:

**Dr. Shantala C P,**Ph.D
Prof & Head, Dept. of CSE,
CIT, Gubbi, Tumakuru.



## Channabasaveshwara Institute of Technology

**(NAAC Accredited &ISO 9001:2015 Certified Institution)**
NH 206 (B.H. Road), Gubbi, Tumakuru – 572 216. Karnataka.

**(Affiliated to Visvesvaraya Technological University, Belagavi & Recognized by AICTE New Delhi)**

**2024-25**

## DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

## <u>CERTIFICATE</u>

This is to certify that the mini project work entitled **"CANCER DETECTION USING MACHINE LEARNING"** has been successfully carried out by **VARSHA P N[1CG22CS119] , TRISHA M[1CG22CS116] , SUSHMITHA M R[1CG22CS114] , SHIVANI S R[1CG22CS104]** bonafide students of **CHANNABASAVESHWARA INSTITUTE OF TECHNOLOGY, GUBBI,**

**TUMAKURU,** under our supervision and guidance and submitted in partial fulfillment for V Semester BE, Mini-project-BCS586 by **Visvesvaraya Technological University, Belagavi** during the academic year of 2024–2025. It is certified that all corrections/suggestions indicated for internal assessment have been incorporated in the report deposited in the departmental library. The mini project report has been approved as it satisfies the academic requirements for the above said degree.

**Guide:**                                                    **H.O.D:**

**Mr. Harish T A M.Tech,**                      **Dr. Shantala C P Ph.D**
Asst. Prof., Dept. of CSE,                   Prof & Head,Dept. of CSE,
CIT, Gubbi, Tumakuru.                        CIT, Gubbi, Tumakuru.

**Principal:**

**Dr. Suresh D S Ph.D**
Director & Principal
CIT, Gubbi, Tumakuru.

# ABSTRACT

This study explores the development and evaluation of predictive models for Cancer detection using a dataset sourced from Kaggle. Cancer is one of the leading causes of death worldwide . This project focuses on leveraging machine learning techniques to predict cancer diagnoses based on clinical and diagnostic data. The goal is to build a model that can effectively classify cancer cases into malignant or benign categories .The dataset used in this study contains various features, such as the size, shape, and texture of tumors, along with other biological markers that can serve as indicators of cancer. Exploratory data analysis (EDA) and Multiple machine learning algorithms are applied to the task, including Support Vector Machine (SVM), Random Forest, and XGBoost, with hyperparameter tuning via GridSearchCV to optimize model performance. Models are evaluated using several metrics such as accuracy, precision, recall, F1-score, and confusion matrices, providing a comprehensive understanding of the classifier's performance. This machine learning-based approach supports medical professionals in detecting cancer more efficiently and accurately, ultimately contributing to better patient outcomes and improved healthcare strategies.

# ACKNOWLEDGEMENT

# Table of Contents

# CHAPTER 1

# INTRODUCTION

The integration of artificial intelligence and machine learning into healthcare is revolutionizing the way diseases are diagnosed and managed. This project focuses on building a robust machine learning model for cancer detection, leveraging advanced data analysis and detective modeling techniques. The goal is to classify tumors as either **malignant** (cancerous) or **benign** (non-cancerous) using data derived from fine needle aspirate images . By analyzing and interpreting features such as the mean radius, texture, perimeter, and compactness of cells, the project aims to assist in the early and accurate diagnosis of cancer. This early detection is crucial as it significantly increases the chances of successful treatment and survival.

The dataset used in this project is highly detailed and includes both numerical and categorical features representing the physical characteristics of the cells. However, one challenge lies in the imbalanced distribution of the target classes, with benign cases being more prevalent than malignant ones. To address this, the Synthetic Minority Oversampling Technique (SMOTE) is applied to balance the dataset, ensuring that the machine learning models are trained equitably and are sensitive to both classes.

The project employs a systematic workflow that begins with data preprocessing and exploratory analysis to identify key patterns and correlations. Feature selection is conducted using the Gini Importance measure from a Random Forest model, which helps in narrowing down the most impactful detectors. Multiple machine learning algorithms, including Random Forest, Support Vector Machines (SVM), and XGBoost, are trained and optimized using cross-validation and hyperparameter tuning via GridSearchCV to ensure the best possible performance.

This project underscores the potential of machine learning in revolutionizing healthcare diagnostics. By providing a fast, reliable, and automated tool for cancer classification, it complements traditional diagnostic methods, helping medical practitioners make informed

decisions with greater confidence.

## 1.1 Objectives

- **Data Analysis :**Explore and understand the dataset, identifying key features contributing to the classification of tumors.

- **Feature Engineering:** Select the most important features to optimize model performance and reduce computational overhead.

- **Model Training:** Train multiple machine learning algorithms (Random Forest, SVM, and XGBoost) to classify tumors effectively.

- **Model Deployment:** Save and load the trained models for future predictions and deployment in real-world applications.

## 1.2 Problem Statement

Cancer, a leading global cause of mortality, involves abnormal cell growth that can spread and threaten life. Early and accurate detection is crucial for improving survival rates, but traditional methods like imaging and biopsies are often time-consuming, costly, and subject to human error. The challenge is to develop efficient, reliable, and scalable systems capable of analyzing complex biological data to detect cancer accurately. This project aims to utilize machine learning to build predictive models that assist in early diagnosis, reduce diagnostic delays, and support medical professionals in making precise and timely treatment decisions, ultimately enhancing patient outcomes.

## CHAPTER 2

# LITERATURE SURVEY

Machine learning (ML) has become a transformative tool in medical diagnostics, particularly for cancer detection. By analyzing complex and high-dimensional datasets, ML algorithms can uncover patterns and relationships that may be challenging for human experts to identify. This capability makes ML essential for tasks like classifying tumors as malignant or benign, providing timely and accurate results that can significantly improve patient outcomes.

Effective feature selection plays a critical role in cancer detection projects. High-dimensional medical datasets often include redundant or irrelevant features that can negatively affect model performance and interpretability. Studies demonstrate that techniques like Gini Importance in Random Forests or domain knowledge-driven feature prioritization are vital in isolating significant predictors such as tumor size, shape, or texture, thus enhancing the accuracy and efficiency of predictive models.

## MACHINE LEARNING ALGORITHMS USED:

### 1. Random Forest Classifier

Random Forest, an ensemble learning method, is used for its robustness in handling complex data structures and its ability to manage overfitting through bootstrap aggregation. It evaluates the importance of each feature using Gini Importance, which helps identify the most significant factors for accurate predictions. This algorithm is highly effective in high-dimensional datasets like those in cancer detection.

### 2. Support Vector Machines (SVM)

SVM is a popular algorithm for classification tasks due to its ability to work well in high-dimensional spaces. By using kernel functions, SVMs can handle non-linear relationships between features. In this project, hyperparameters like C (regularization) and gamma (kernel coefficient) are optimized through GridSearchCV to maximize performance.

## 3. XGBoost (Extreme Gradient Boosting)

XGBoost is a gradient boosting algorithm known for its speed and performance. It iteratively improves model predictions by minimizing errors, making it particularly useful for structured datasets. In this project, XGBoost is fine-tuned using parameters like n_estimators to achieve high accuracy and reliability.

## 4. Synthetic Minority Oversampling Technique (SMOTE)

To address class imbalance in the dataset, SMOTE is employed to create synthetic examples of the minority class. This technique ensures the model learns effectively from both classes, improving its ability to identify malignant cases accurately.

These algorithms, combined with hyperparameter tuning and feature selection, form the foundation of the project, enabling robust and accurate cancer detection.

# CHAPTER 3
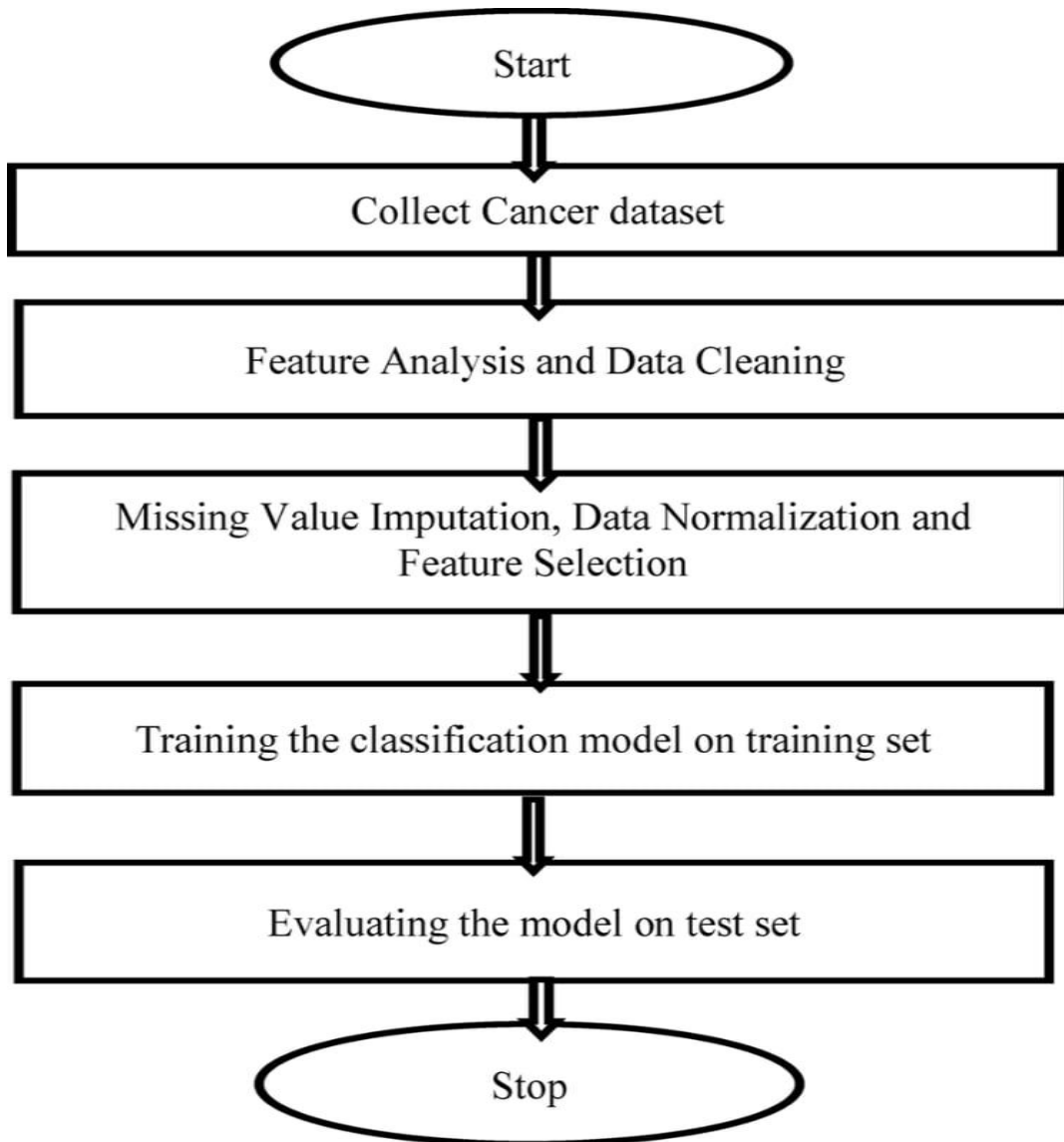
## SYSTEM DESIGN

### 3.1 ARCHITECTURE

```
                    ┌─────────────┐
                    │    Start    │
                    └─────────────┘
                           │
                           ▼
          ┌──────────────────────────────────────┐
          │         Collect Cancer dataset        │
          └──────────────────────────────────────┘
                           │
                           ▼
          ┌──────────────────────────────────────┐
          │    Feature Analysis and Data Cleaning │
          └──────────────────────────────────────┘
                           │
                           ▼
          ┌──────────────────────────────────────┐
          │  Missing Value Imputation, Data       │
          │  Normalization and Feature Selection  │
          └──────────────────────────────────────┘
                           │
                           ▼
          ┌──────────────────────────────────────┐
          │ Training the classification model on  │
          │           training set                │
          └──────────────────────────────────────┘
                           │
                           ▼
          ┌──────────────────────────────────────┐
          │     Evaluating the model on test set  │
          └──────────────────────────────────────┘
                           │
                           ▼
                    ┌─────────────┐
                    │    Stop     │
                    └─────────────┘
```

Fig 3.1. Architecture of the proposed model

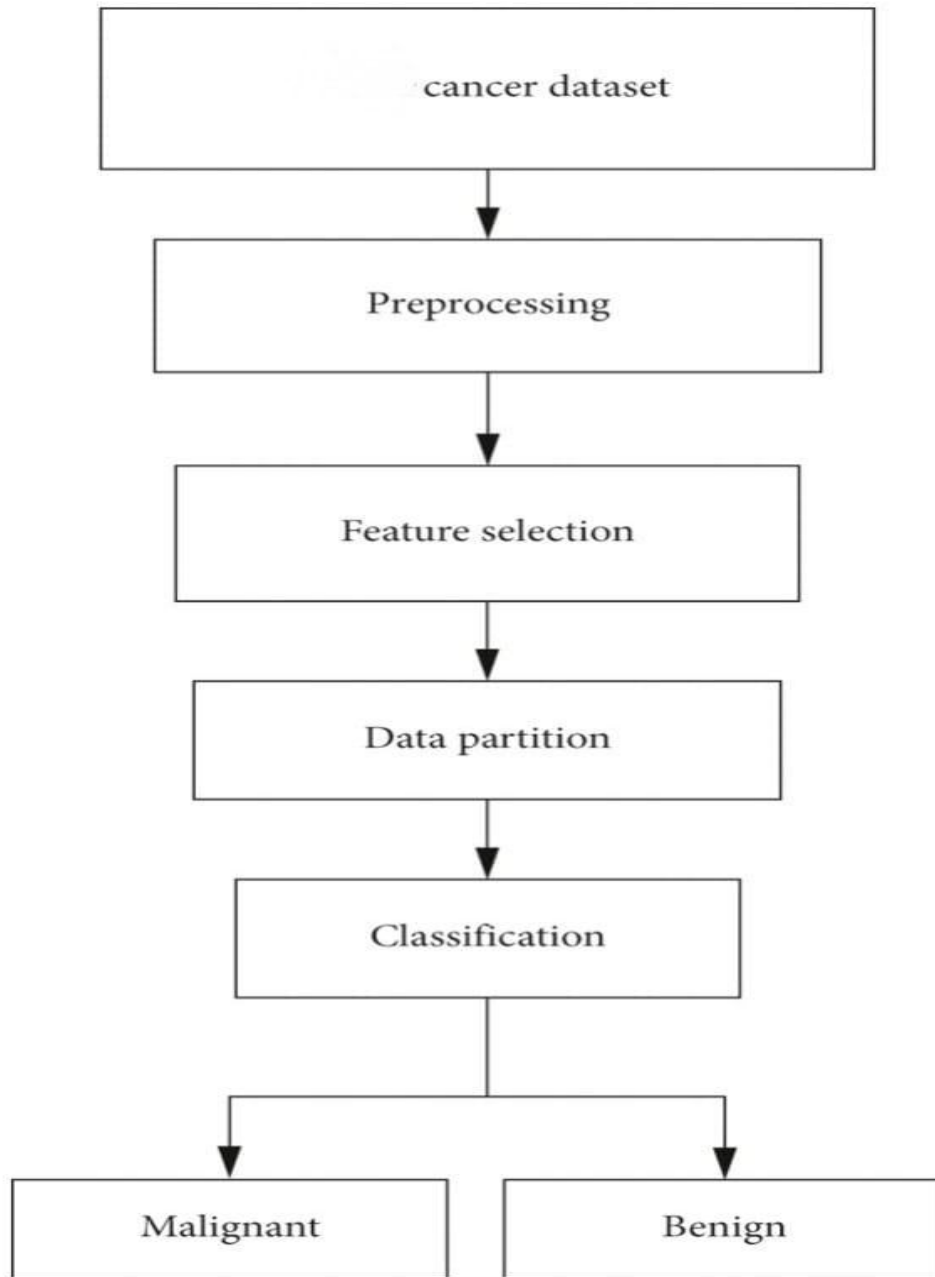## 3.2 DATAFLOW DIAGRAM



Fig 3.2. Dataflow diagram of proposed system
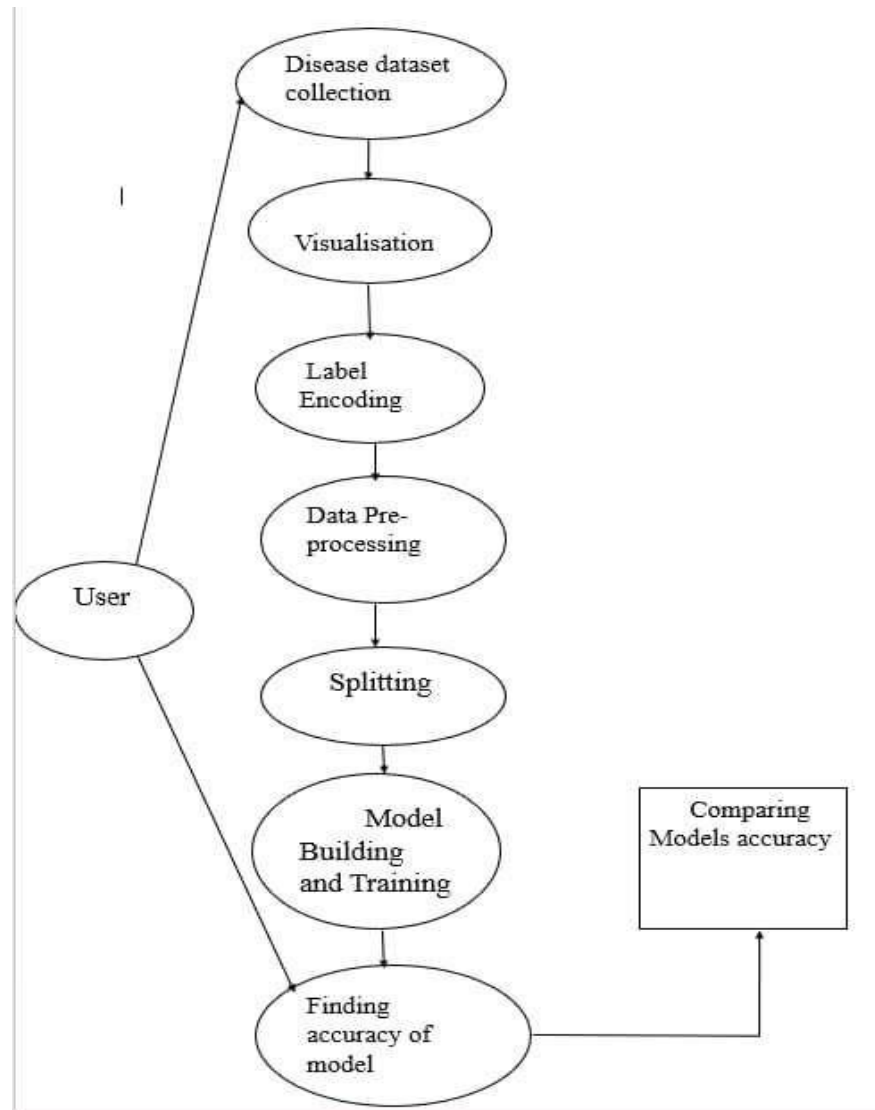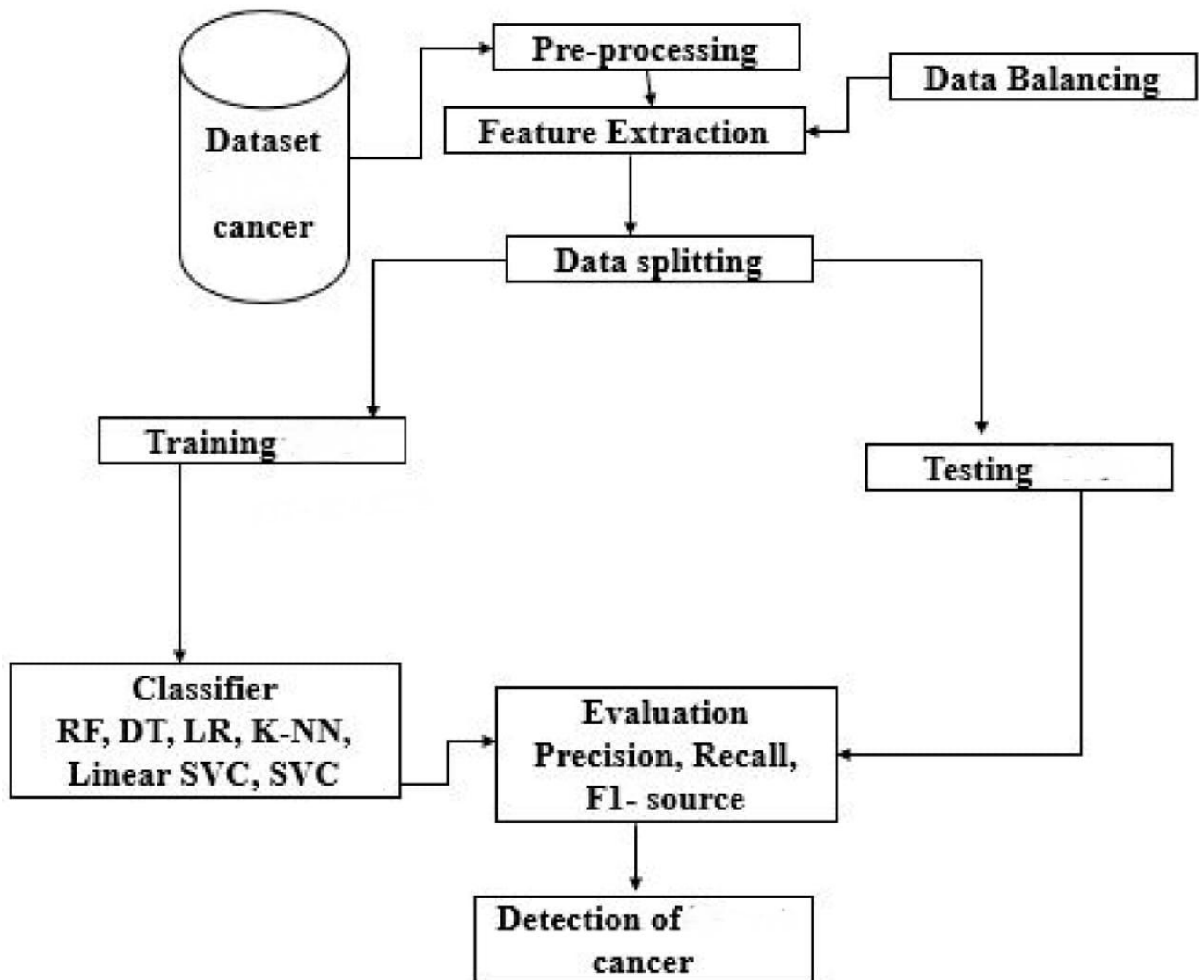
## 3.3 USE CASE DIAGRAM



Fig 3.3 use case diagram of proposed system

## 3.4 SEQUENCE DIAGRAM



3.4 Sequence diagram of proposed system

# CHAPTER 4

# METHODOLOGY

Methodology includes steps ranging from data acquisition and preprocessing to model evaluation and interpretation of results.

## 4.1 Data Collection

The dataset is downloaded from Kaggle. The Cancer Detection Model1.pynb file contains Cancer Data.csv file. The cancer dataset consists of 569 rows and 32 columns and Cancer data consist of 569 rows and 30 columns including independent column . and 1 column which is target(dependent)column(diagnosis).



Fig.4.1.1 total 32columns in data set

## 4.2 Data Preprocessing:

In the data preprocessing step, the primary focus was to prepare the dataset for analysis and modelling. The dataset consisted of 32 columns, where 32 columns represent feature selection, and the target column (labelled diagnosis) detects the presence or absence of cancer. The independent features are float64 values. The dependent column represented by M and B, where:

- M indicates the presence of a Cancer.

- B indicates the absence of a Cancer.

➢ To ensure data quality, the .isnull () function from the Pandas library in Python was employed. This function was used to check for any missing (null) values across all columns in the dataset. Upon running the check, it was confirmed that no null values were present in any of the columns, ensuring a complete and clean dataset for further analysis.

```
radius_se                0
texture_se               0
perimeter_se             0
area_se                  0
smoothness_se            0
compactness_se           0
concavity_se             0
concave_points_se        0
symmetry_se              0
fractal_dimension_se     0
radius_worst             0
texture_worst            0
perimeter_worst          0
area_worst               0
smoothness_worst         0
compactness_worst        0
concavity_worst          0
concave_points_worst     0
symmetry_worst           0
fractal_dimension_worst  0
dtype: int64
```

**Note: There is no null value in the dataset.**

Fig.4.2.1.Checking of null values in dataset

➢ Categorical Encoding: Convert categorical variables (like the diagnosis column with values 'M' and 'B') into numeric values (0 for benign and 1 for malignant) using Label Encoder or similar methods.

➢ Feature Scaling: Scale numeric features (such as radius, texture, area, etc.) to standardize the data for models that require it (e.g., SVM, k-NN).

## 4.3 Exploratory Data Analysis(EDA)

### 4.3.1 Visualizing Data

To gain deeper insights into the dataset and understand the relationships between features and the target column (diagnosis), various visualization techniques were employed:

1. **HeatMap plot**

    heatmap is a graphical representation of data that uses a system of color coding to represent different values. Heatmaps are used in various forms of analytics but are most commonly used to show relationships between two variables,one plotted on each axis and make it easier to discern areas of high and low concentration.

2. **Box Plot**

    A boxplot,is a standardized way of displaying the distribution of a data set based on its five-number summary of data points used to show distributions of numeric data values. it also organizes large amounts of data, and visualizes outlier values.

3. **Distribution Plot**

    A distribution plot is a visualization that shows the distribution of numerical data in a data set. It can be used to compare the range and distribution of data groups.Distribution plots compare observed data to expected outcomes. They display two plots for the same variable: a histogram that represents the observed data, and a density curve that represents the expected outcomes.

### 4.3.2 Class Imbalance Detection

Visualize the distribution of benign and malignant classes to check if the dataset is imbalanced. If one class dominates, techniques like SMOTE (Synthetic Minority Over-sampling Technique) can be used to balance the data.

## 4.4 Data Splitting:

### 1. Overview of Data Splitting:

➤ The dataset contains features that describe various characteristics of cell nuclei (such as radius_mean, texture_mean, area_mean, etc.) and a target variable (diagnosis) indicating whether the tumor is Malignant (M) or Benign (B).

➤ The purpose of data splitting is to divide the dataset into separate sets for training and testing to avoid overfitting and ensure the model generalizes well on unseen data.

```
Index(['radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean',
       'smoothness_mean', 'compactness_mean', 'concavity_mean',
       'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
       'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
       'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
       'fractal_dimension_se', 'radius_worst', 'texture_worst',
       'perimeter_worst', 'area_worst', 'smoothness_worst',
       'compactness_worst', 'concavity_worst', 'concave points_worst',
       'symmetry_worst', 'fractal_dimension_worst'],
      dtype='object')
```

Fig.4.4.1 Features of dataset

### 2. Train-Test Split:

➤ Training Set: This subset of the data is used to train the machine learning models, allowing the model to learn from the features and their relationship to the target variable.

➤ Testing Set: This subset is kept aside during training and is used to evaluate the performance of the trained model. The testing data is not seen by the model during training, providing an unbiased estimate of the model's performance.

➤ Typically, 80% of the data is used for training, and 20% for testing. This is known as an 80-20 split, but in some cases, a different split ratio like 70-30 is used based on dataset size and model requirements.

### 3. Cross-Validation :

➢ Cross-validation is a technique used to ensure that the model performs well across different subsets of the data. It divides the dataset into multiple folds (e.g., 10-fold cross-validation).

➢ KFold Cross-Validation is used,data is split into K number of subsets and  model is trained and tested K times, each time with a different fold as the test set and the remaining folds as the training set. This helps provide a more reliable performance estimate by averaging the results over different test sets.

### 4. Stratified Sampling :

➢ In case the dataset has an imbalanced class distribution(more Benign than Malignant cases), Stratified Sampling can be used. This ensures that both the training and testing sets have a proportionate representation of each class (Malignant and Benign) based on their occurrence in the full dataset.

➢ This prevents the model from being biased toward the majority class and helps it learn patterns for both Malignant and Benign cases.

### 5. Data Resampling :

➢ SMOTE (Synthetic Minority Over-sampling Technique): If the dataset is highly imbalanced (for example, if there are significantly more Benign cases than Malignant), SMOTE can be used to oversample the minority class (Malignant) in the training set by generating synthetic samples. This helps in balancing the dataset, allowing the model to perform better in detecting Malignant cases.

### 6. Splitting Code Implementation:

➢ The data is typically split using a function like train_test_split() from scikit-learn, which handles the random shuffling and splitting process. For cross-validation, KFold from scikit-learn can be used to generate train-test splits across different folds.

## 4.5 Feature Selection:

➢ feature selection is a crucial step to ensure the model uses the most informative attributes for accurate classification. Initially, the dataset contains several features related to the characteristics of cell nuclei, such as radius_mean, texture_mean, and area_mean. To identify the most relevant features, the project uses Random Forest's feature importance method. This technique ranks the features based on how much they contribute to reducing the Gini impurity during the tree-building process. The top-ranked features, like radius_mean and texture_mean, are then selected for model training, while less relevant or redundant features are discarded. This helps improve model performance, reduce computation time, and avoid overfitting.

## 4.6 Model used:

the project employs three machine learning models.

### 4.6.1 Random Forest Classifier

➢ Random Forest is a supervised machine learning algorithm extensively used for both classification and regression tasks. In the context of classification, such as diagnosing diseases, Random Forest excels at analyzing complex relationships in data and making robust predictions based on multiple input features like patient symptoms or test results.

➢ This algorithm operates by building an ensemble of decision trees, where each tree makes an independent prediction. The final classification is determined by aggregating the predictions through majority voting, enhancing the model's accuracy and reducing the risk of overfitting.

➢ A key feature of Random Forest is its use of randomness in both sampling the training data (bootstrap sampling) and selecting subsets of features for splitting nodes within each tree. This ensures diversity among the trees and leads to better



generalization to unseen data.

Fig 4.6.1 Random Forest

**ALGORITHM**

**Step 1:** Collect and Prepare Data

Gather and preprocess the dataset for classification. Handle missing values, remove duplicates, and encode categorical variables as needed. Split the dataset into training and testing sets.

**Step 2:** Initialize the Random Forest

Specify the hyper parameters, including the number of trees maximum tree depth, and other parameters like minimum sample ssplit. Note that the Random Forest automatically applies bootstrap sampling during training.

**Step 3:** Train the Model

Fit the Random Forest Classifier to the training data. The algorithm automatically:

- Builds multiple decision trees by sampling subsets of data (handled internally).

- Selects a random subset of features at each split within the decision trees.

**Step 4:** Evaluate the Model

Use the testing set to evaluate the model's performance. Predictions from individual trees are aggregated (majority voting for classification). Assess model accuracy and other metrics like precision, recall, F1- score.

**Step 5:** Make Predictions

Apply the trained Random Forest Classifier to new, unseen data for classification. Optionally, interpret feature importance scores to understand which features significantly influence predictions

### 4.6.2 Support Vector Machine

➢ Support Vector Machine (SVM) is a supervised machine learning algorithm commonly used for classification tasks, such as diagnosing diseases by distinguishing between different health conditions based on patient data like symptoms or test results.

➢ SVM aims to find the optimal hyperplane that separates data points into predefined classes with the maximum margin between the closest points (support vectors) of each class. This margin maximization ensures better generalization of the model to unseen data.

➢ In cases where data is not linearly separable, SVM uses kernel functions (e.g., linear, polynomial, radial basis function) to map the input data into a higher-dimensional space, making separation possible.



Fig 4.6.2 Support Vector Machine

**ALGORITHM**

**Step 1:** Collect and Prepare Data

Gather and clean the dataset for your classification task. Handle missing values, remove duplicates, and preprocess the data. Scale the features (e.g., standardization or normalization) because SVMs are sensitive to the range of feature values. Split the dataset into training and testing sets.

**Step 2:** Define the Model

Choose the SVM type based on the problem:

- Linear SVM for linearly separable data.

- Kernel SVM (e.g., RBF, polynomial, or sigmoid kernel) for non-linearly separable data. Specify hyper parameters such as C (regularization parameter) and kernel type.

**Step 3:** Train the Model

- Fit the SVM model to the training data. The algorithm:

- Finds the hyper plane that maximizes the margin (distance) between the two classes.

- Considers a small margin of error, controlled by the regularization parameter $C$, for misclassified points.
- For non-linear data, maps the input space to a higher-dimensional feature space using the selected kernel function.

**Step 4:** Evaluate the Model

Use the testing set to predict the target variable. Evaluate the model's performance using metricslike accuracy, precision, recall, F1-score. Analyze the support vectors (data points closest to the hyper plane ) to understand the classification boundaries.

**Step 5:** Make Predictions

Apply the trained SVM model to new, unseen data for classification. Interpret the decision boundaries and optionally perform hyper parameter tuning (e.g., grid search) to optimize the model for better performance.

### 4.6.3 XG Boost

➢ XG Boost is a high-performance, supervised machine learning algorithm based on the gradient boosting framework. It is widely used for both classification and regression tasks, including disease diagnosis, due to its speed and predictive power.

➢ XG Boost builds an ensemble of decision trees in a sequential manner, where each tree corrects the errors made by its predecessor. The final prediction is an aggregation of the individual trees' outputs, weighted by their learning contributions.

➢ What sets XG Boost apart is its focus on regularization (to prevent overfitting), parallel processing, and efficient handling of missing values, making it suitable for large datasets with complex feature interactions.
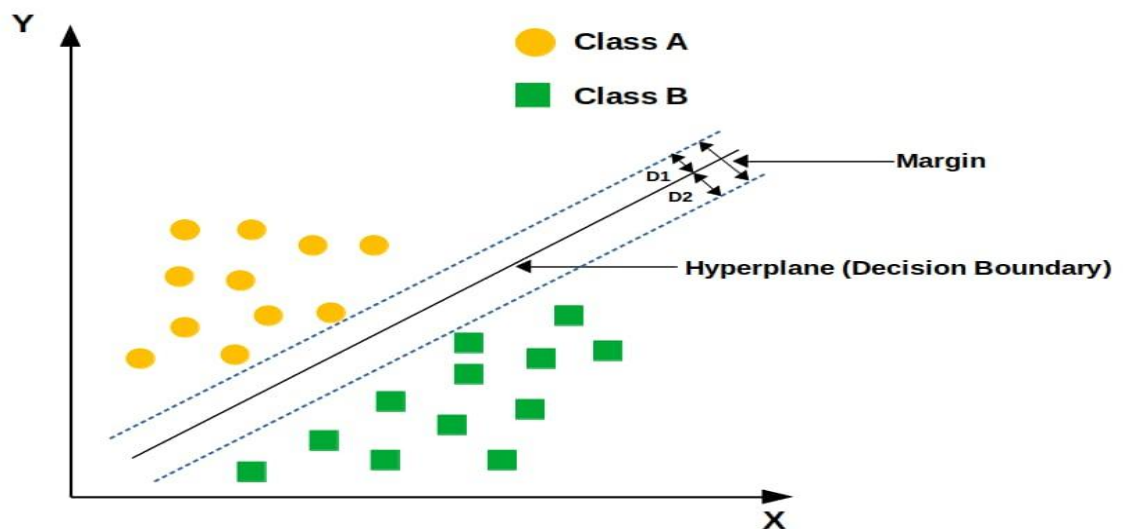


Fig 4.6.3 XG Boost

## ALGORITHM

**Step 1:** Collect and Prepare Data

Gather and clean the dataset for the classification task. Handle missing values, remove duplicates, and preprocess the features (e.g., encoding categorical variables, scaling). Split the dataset into training and testing sets.

**Step 2:** Define the Model

Initialize the XG Boost model with hyper parameters such as:

Number of boosting rounds ,The step size at each iteration , Maximum depth of each decision tree, Fraction of training data used in each boosting round, Fraction of features used for each tree, Loss function.

**Step 3:** Train the Model

Fit the XG Boost model to the training data using gradient boosting. In each boosting round:

- Train a decision tree on the residual errors from the previous tree.

- Update the model by adding the new tree's predictions weighted by the learning rate.

- Repeat the process for the specified number of boosting rounds (or until convergence).

**Step 4:** Evaluate the Model

Use the testing set to evaluate the model's performance. Predict the target variable and evaluate performance using metrics like accuracy, precision, recall, F1-score. Optionally, perform cross- validation for more robust results.

**Step 5:** Make Predictions

Apply the trained XG Boost model to new, unseen data for classification. Optionally, tune hyper parameters (using grid search or random search) to improve the model's performance.

## 4.7 Model Evaluation

**Accuracy:** This is the proportion of correct predictions out of all the predictions made by the model. However, accuracy can be misleading if the data is imbalanced ,as a model that always predicts the majority class can have a high accuracy but not be useful. Mathematically, accuracy is calculated as the ration of the number of correctly predicted instances to the total number of instances:

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total Number of Predictions}}$$

**Recall:**

- Recall, also known as sensitivity or true positive rate, measures the proportion of true positive predictions (correctly predicted positive instances) out of all actual positive instances in the dataset.

- It focuses on the model's ability to capture all positive instances and answers the question: "Of all actual positive instances, how many did the model correctly predict as positive?"

**F1-score:** This is a harmonic mean of precision and recall, and can be useful for imbalanceddatasets where both precision an drew al l need to be considered.

# CHAPTER 5

# RESULTS

| | id | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean | compactness_mean | concavity_mean | concave points_mean | ... | radius_worst | texture_worst | perimeter_worst | area_worst | smoothness_wo |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 842302 | M | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.30010 | 0.14710 | ... | 25.380 | 17.33 | 184.60 | 2019.0 | 0.162 |
| 1 | 842517 | M | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.08690 | 0.07017 | ... | 24.990 | 23.41 | 158.80 | 1956.0 | 0.123 |
| 2 | 84300903 | M | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.19740 | 0.12790 | ... | 23.570 | 25.53 | 152.50 | 1709.0 | 0.144 |
| 3 | 84348301 | M | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.24140 | 0.10520 | ... | 14.910 | 26.50 | 98.87 | 567.7 | 0.209 |
| 4 | 84358402 | M | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.19800 | 0.10430 | ... | 22.540 | 16.67 | 152.20 | 1575.0 | 0.137 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 564 | 926424 | M | 21.56 | 22.39 | 142.00 | 1479.0 | 0.11100 | 0.11590 | 0.24390 | 0.13890 | ... | 25.450 | 26.40 | 166.10 | 2027.0 | 0.141 |
| 565 | 926682 | M | 20.13 | 28.25 | 131.20 | 1261.0 | 0.09780 | 0.10340 | 0.14400 | 0.09791 | ... | 23.690 | 38.25 | 155.00 | 1731.0 | 0.116 |
| 566 | 926954 | M | 16.60 | 28.08 | 108.30 | 858.1 | 0.08455 | 0.10230 | 0.09251 | 0.05302 | ... | 18.980 | 34.12 | 126.70 | 1124.0 | 0.113 |
| 567 | 927241 | M | 20.60 | 29.33 | 140.10 | 1265.0 | 0.11780 | 0.27700 | 0.35140 | 0.15200 | ... | 25.740 | 39.42 | 184.60 | 1821.0 | 0.165 |
| 568 | 92751 | B | 7.76 | 24.54 | 47.92 | 181.0 | 0.05263 | 0.04362 | 0.00000 | 0.00000 | ... | 9.456 | 30.37 | 59.16 | 268.6 | 0.089 |

569 rows × 32 columns

Fig  5.1 Reading Data Set

| | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean | compactness_mean | concavity_mean | concave points_mean | symmetry_mean | fractal_dimension_mean | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| radius_mean | 1.000000 | 0.323782 | 0.997855 | 0.987357 | 0.170581 | 0.506124 | 0.676764 | 0.822529 | 0.147741 | -0.311631 | ... |
| texture_mean | 0.323782 | 1.000000 | 0.329533 | 0.321086 | -0.023389 | 0.236702 | 0.302418 | 0.293464 | 0.071401 | -0.076437 | ... |
| perimeter_mean | 0.997855 | 0.329533 | 1.000000 | 0.986507 | 0.207278 | 0.556936 | 0.716136 | 0.850977 | 0.183027 | -0.261477 | ... |
| area_mean | 0.987357 | 0.321086 | 0.986507 | 1.000000 | 0.177028 | 0.498502 | 0.685983 | 0.823269 | 0.151293 | -0.283110 | ... |
| smoothness_mean | 0.170581 | -0.023389 | 0.207278 | 0.177028 | 1.000000 | 0.659123 | 0.521984 | 0.553695 | 0.557775 | 0.584792 | ... |
| compactness_mean | 0.506124 | 0.236702 | 0.556936 | 0.498502 | 0.659123 | 1.000000 | 0.883121 | 0.831135 | 0.602641 | 0.565369 | ... |
| concavity_mean | 0.676764 | 0.302418 | 0.716136 | 0.685983 | 0.521984 | 0.883121 | 1.000000 | 0.921391 | 0.500667 | 0.336783 | ... |
| concave points_mean | 0.822529 | 0.293464 | 0.850977 | 0.823269 | 0.553695 | 0.831135 | 0.921391 | 1.000000 | 0.462497 | 0.166917 | ... |
| symmetry_mean | 0.147741 | 0.071401 | 0.183027 | 0.151293 | 0.557775 | 0.602641 | 0.500667 | 0.462497 | 1.000000 | 0.479921 | ... |
| fractal_dimension_mean | -0.311631 | -0.076437 | -0.261477 | -0.283110 | 0.584792 | 0.565369 | 0.336783 | 0.166917 | 0.479921 | 1.000000 | ... |

Fig 5.2 Check for Multicollinearity

```
radius_mean                      0
texture_mean                     0
perimeter_mean                   0
area_mean                        0
smoothness_mean                  0
compactness_mean                 0
concavity_mean                   0
concave points_mean              0
symmetry_mean                    0
fractal_dimension_mean           0
radius_se                        0
texture_se                       0
perimeter_se                     0
area_se                          0
smoothness_se                    0
compactness_se                   0
concavity_se                     0
concave points_se                0
symmetry_se                      0
fractal_dimension_se             0
```

Fig 5.3 Checking of null values



Fig 5.4 Feature selection chart

```
Feature Ranking:
Feature radius_mean (0.139416)
Feature texture_mean (0.137679)
Feature perimeter_mean (0.118768)
Feature area_mean (0.117989)
Feature smoothness_mean (0.077427)
Feature compactness_mean (0.057240)
Feature concavity_mean (0.048971)
Feature concave points_mean (0.045257)
Feature symmetry_mean (0.035983)
Feature fractal_dimension_mean (0.033661)
Feature radius_se (0.030638)
Feature texture_se (0.017393)
Feature perimeter_se (0.015467)
Feature area_se (0.014069)
Feature smoothness_se (0.013795)
```

Fig 5.5 Feature Ranking



Fig 5.3 Heatmap plot

Fig 5.4Box plot



Fig 5.5 Distribution plot

## Random Forest Model

```
Fitting 10 folds for each of 4 candidates, totalling 40 fits
[1 1 0 1 0 1 0 0 1 0 0 1 0 1 0 1 1 1 1 1 0 0 0 0 1 0 1 1 0 0 0 1 1 1 1 1 1
 1 1 0 0 0 1 0 1 1 1 1 0 1 1 0 1 1 0 1 0 1 0 1 1 1 1 1 0 1 1 0 1 1 0 0 0
 1 1 1 0 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 1 0 1 1 1 1 1 0 1 1 0 0 0 0 0 0 1
 0 1 0 0 1 0 1 1 1 0 0 1 1 0 0 0 0 0 1 0 1 1 0 0 1 0 0 0 0 0 0]
Best Parameters : {'n_estimators': 500}
Classification Report:                     precision    recall  f1-score   support

             0       1.00      0.97      0.99        68
             1       0.97      1.00      0.99        75

      accuracy                           0.99       143
     macro avg       0.99      0.99      0.99       143
  weighted avg       0.99      0.99      0.99       143

Accuracy Score 0.986013986013986
Confusion Matrix :
 [[66  2]
 [ 0 75]]
```
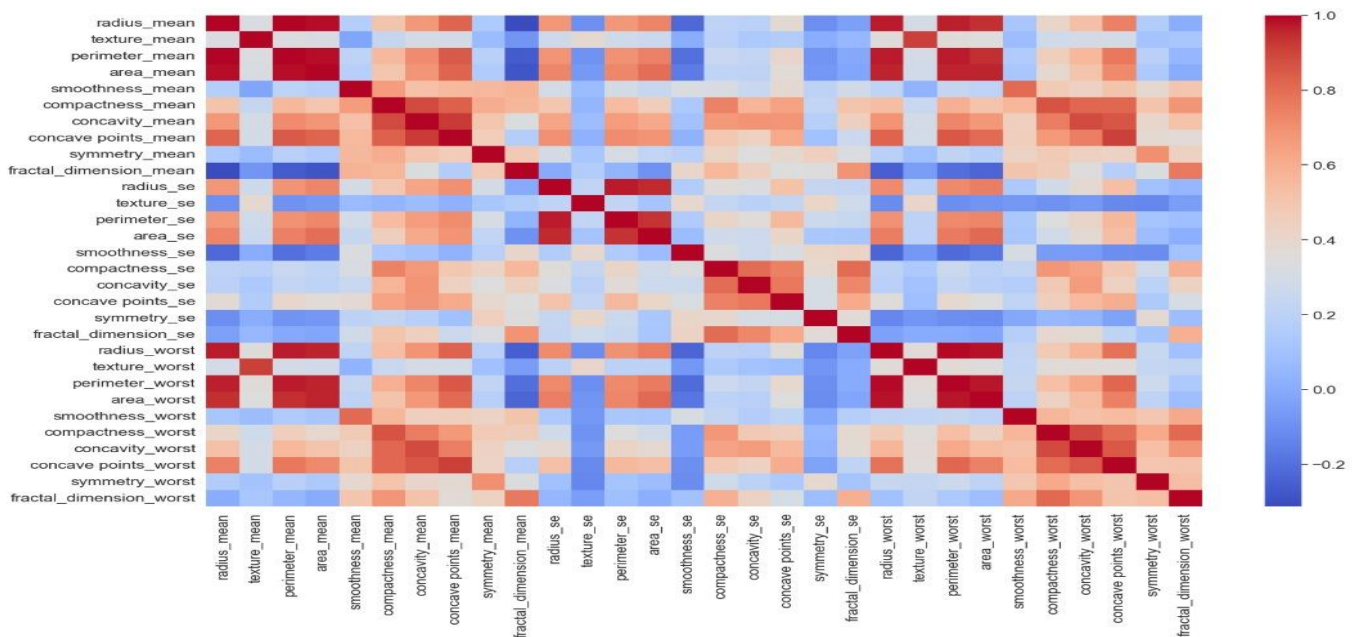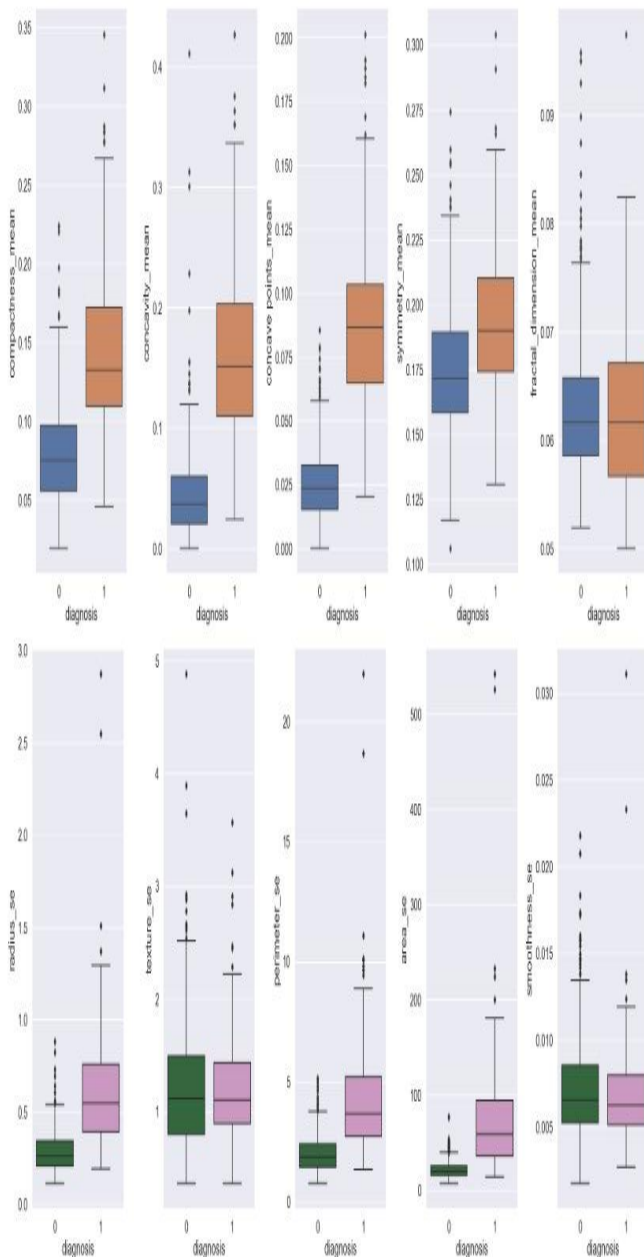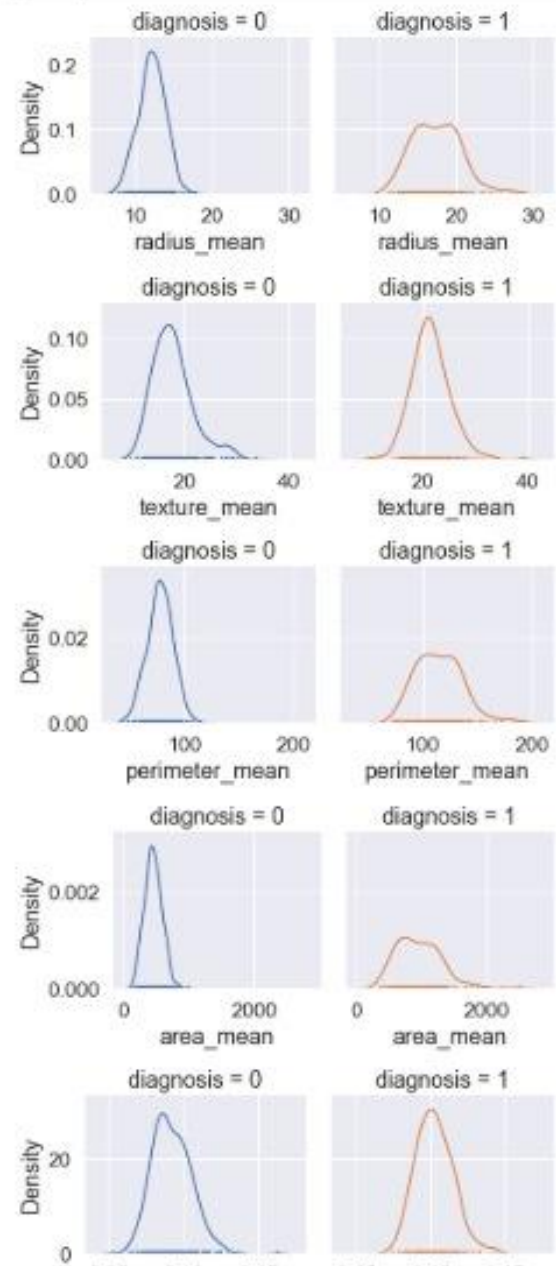
## SVC Model

```
Fitting 10 folds for each of 36 candidates, totalling 360 fits
[1 1 0 1 0 1 0 0 1 0 0 1 0 1 0 1 1 1 1 1 0 0 0 0 1 0 1 1 0 0 0 1 1 1 1 1 1
 1 1 0 0 0 1 0 1 1 1 1 0 1 1 0 1 1 0 1 0 1 0 1 1 0 1 1 1 1 1 0 1 1 0 1 1 0 0 0
 1 1 1 0 0 0 1 0 1 1 1 1 1 1 1 0 1 0 1 0 1 1 1 1 1 0 1 1 0 0 0 0 0 0 0 0 0 1
 0 1 0 0 1 0 1 1 1 0 0 1 1 0 0 0 0 0 1 0 1 1 0 0 1 0 0 0 0 0 1]
Best Parameters : {'C': 100, 'gamma': 0.0001}
Classification Report:                     precision    recall  f1-score   support

             0       0.97      0.94      0.96        68
             1       0.95      0.97      0.96        75

      accuracy                           0.96       143
     macro avg       0.96      0.96      0.96       143
  weighted avg       0.96      0.96      0.96       143

Accuracy Score 0.958041958041958
Confusion Matrix :
 [[64  4]
 [ 2 73]]
```

## XG Boost

```
Fitting 10 folds for each of 4 candidates, totalling 40 fits
[1 1 0 1 0 1 0 0 1 0 0 1 0 1 0 1 1 1 1 1 0 0 0 0 1 0 1 1 0 0 0 1 1 1 1 1 1
 1 1 0 0 0 1 0 1 1 1 1 0 1 1 0 1 1 0 1 0 1 0 1 1 0 1 1 1 1 1 0 1 1 0 1 1 0 0 0
 1 1 1 0 0 0 1 0 1 1 1 1 1 1 1 0 1 0 1 0 1 1 1 1 1 0 1 1 0 0 0 0 0 0 0 0 0 1
 0 1 0 0 1 0 1 1 1 0 0 1 1 0 0 0 0 0 1 0 1 1 0 0 1 0 0 0 0 0 0]
Best Parameters : {'n_estimators': 500}
Classification Report:                     precision    recall  f1-score   support

             0       1.00      0.97      0.99        68
             1       0.97      1.00      0.99        75

      accuracy                           0.99       143
     macro avg       0.99      0.99      0.99       143
  weighted avg       0.99      0.99      0.99       143

Accuracy Score 0.986013986013986
Confusion Matrix :
 [[66  2]
 [ 0 75]]
```

Fig 5.8 Creating of Random Forest, XG Boost and SVC Models

```
[1 0 0 1 0 0 0 0 1 1 0 0 1 0 1 1 0 0 0 0 1 1 0 0 1 1 0 1 0 0 0 0 1 0 0 0 0
 0 0 0 0 0 0 1 0 1 1 1 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 1 0 0 0 0 1 1 0 0 0 0
 0 1 0 0 1 1 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 1 0 1 0 1 1 0 0 0 0 1 1 0 0 0
 0 1 1]
```

Fig 5.9 Predicted values

## Best Parameter

```
load_model.best_params_
```

```
{'n_estimators': 100}
```

## Accuracy Score

```
print (accuracy_score (pred1,y_test))
```

```
0.9912280701754386
```

Fig 5.10 best parameter and Accuracy score

# CHAPTER 6

# CONCLUSION

In conclusion, the cancer detection project successfully demonstrates the application of machine learning techniques to classify tumors as either Malignant or Benign based on various cell feature attributes. Through careful data preprocessing, including handling missing values, feature encoding, and normalization, the project ensures that the dataset is clean and ready for model training. Feature selection, driven by Random Forest's feature importance, plays a vital role in improving model accuracy by focusing on the most significant attributes.

The project utilizes multiple machine learning models, such as Support Vector Machine (SVM), Random Forest, and XGBoost, to predict cancer types. These models, each optimized through techniques like GridSearchCV, are evaluated based on their performance, with Random Forest and XGBoost proving to be highly effective in handling the dataset's complexities. By leveraging cross-validation and addressing class imbalance through methods like SMOTE, the project ensures reliable and generalizable results.

Overall, the project highlights the potential of machine learning in medical diagnostics, particularly in early cancer detection. The approach not only provides a framework for future research in cancer prediction but also opens avenues for further optimization and real-world application in healthcare.

# REFERENCES

1. Breiman, L.(2001). "Random Forests." *Machine Learning*, 45(1), 5–32. This paper introduces the Random Forest algorithm, emphasizing its use in classification and regression tasks, which are critical for disease prediction models.

2. Chen, T., & Guestrin, C. (2016). "XGBoost: A Scalable Tree Boosting System." *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794.This study describes XGBoost, highlighting its speed and scalability, often applied in medical data analysis for its accuracy.

3. Deo, R. C. (2015). "Machine Learning in Medicine." *Circulation*, 132(20), 1920–1930. The article reviews various machine learning models, including SVM and Logistic Regression, for applications in medical diagnosis and disease prediction.

4. Sharma, A., & Virmani, J. (2022). "A Comprehensive Review on Machine Learning Algorithms for Disease Prediction." *Biomedical Signal Processing and Control*, 71, Article103144. This review compares ML algorithms, detailing their performance in healthcare applications, particularly for chronic disease prediction.

5. Sharma, A., & Virmani, J. (2022). "A Comprehensive Review on Machine Learning Algorithms for Disease Prediction." *Biomedical Signal Processing and Control*, 71, Article 103144.

6. Deo, R. C. (2015). "Machine Learning in Medicine." *Circulation*, 132(20), 1920–1930.

7. Lundberg, S. M., & Lee, S. I. (2017). "A Unified Approach to Interpreting Model Predictions."

*Advances in Neural Information Processing Systems*, 30, 4765–4774.

8. Haq, A. U., Li, J. P., Memon, M. H., Nazir, S., & Sun, N. (2018). "A Hybrid Intelligent System Framework for the Prediction of Heart Disease Using Machine Learning Algorithms." *Mobile Information Systems*, Article 3860146.

9. Zhang, Z., & Yang, L. (2021). "Machine Learning in Disease Prediction: Current Status and Future Directions." *Frontiers in Bioinformatics*, 1, Article 651159.

10. Goldstein, B. A., Navar, A. M., Pencina, M. J., & Ioannidis, J. P. A. (2017). "Opportunities and Challenges in Developing Risk Prediction Models with Machine Learning and Artificial Intelligence." *JAMA*, 320(7), 1797–1806.

11. World Health Organization. (2011). "Hemoglobin concentrations for the diagnosis of anemia and assessment of severity." *WHO Report No. WHO/NMH/NHD/MNM/11.1*.

12. Obaidy, M. A., et al. (2021). "Prevalence and Risk Factors of Anemia among Children Aged 5 months-12 years." *Mosul Journal of Nursing*, 9(1), 131-137.

13. Amin, N., & Habib, A. (2015). "Comparison of different classification techniques using WEKA for hematological data." *American Journal of Engineering Research*, 4(3), 55-61.

14. Jaiswal, M., Srivastava, A., & Siddiqui, T. J. (2019). "Machine learning algorithms for anemia disease prediction." *Recent Trends in Communication, Computing, and Electronics*. Springer Singapore.

15. Shilpa, S. A., Nagori, M., & Kshirsaga, V. (2011). "Classification of anemia using data mining techniques." *Swarm, Evolutionary, and Memetic Computing*, 113-121. Springer.

16. Sow, B., et al. (2020). "Assessing the relative importance of social determinants of health in malaria and anemia classification based on machine learning techniques." *Informatics for Health and Social Care*, 45(3), 229-241.