

A

Major Project Report

On

Liver Disease Prediction

Submitted to CMR Engineering College, HYDERABAD

In Partial Fulfillment of the requirements for the Award of Degree of

**BACHELOR OF TECHNOLOGY
IN
COMPUTER SCIENCE AND ENGINEERING (DATA SCIENCE)**

Submitted

By

A. SHIRISHA	218R1A6704
J. AKHILA	218R1A6729
A. KUSHAL KUMAR	218R1A6702
K. MADHUKAR	218R1A6731

Under the Esteemed guidance of

Mr. B. KUMARASWAMY

Assistant Professor (Ph.D, JNTUH), Department of CSE (Data Science)



Department of Computer Science and Engineering (Data Science)

CMR ENGINEERING COLLEGE

UGC AUTONOMOUS

(Approved by AICTE, NEW DELHI, Affiliated to JNTU, Hyderabad)
Kandlakoya, Medchal Road, R.R. Dist. Hyderabad-501 401.

2024-2025

CMR ENGINEERING COLLEGE

UGC AUTONOMOUS

(Accredited by NBA, Approved by AICTE NEW DELHI, Affiliated to JNTU, Hyderabad)

Kandlakoya, Medchal Road, Hyderabad-501 401

Department of Computer Science and Engineering(Data Science)



CERTIFICATE

This is to certify that the project entitled “**Liver Disease Prediction**” is a bonafide work carried out by

A. SHIRISHA	218R1A6704
J. AKHILA	218R1A6729
A. KUSHAL KUMAR	218R1A6702
K. MADHUKAR	218R1A6731

in partial fulfillment of the requirement for the award of the degree of **BACHELOR OF TECHNOLOGY** in **COMPUTER SCIENCE AND ENGINEERING (DATA SCIENCE)** from CMR Engineering College, affiliated to JNTU, Hyderabad, under our guidance and supervision.

The results presented in this Major project have been verified and are found to be satisfactory. The results embodied in this Major project have not been submitted to any other university for the award of any other degree or diploma.

Internal Guide	Major Project Coordinator	Head of the Department	External Examiner
Mr. B. Kumaraswamy (Ph.D, JNTUH) Assistant Professor CSE (Data Science), CMREC	Mrs. G. Shruthi Assistant Professor CSE (Data Science), CMREC	Dr. M. Laxmaiah Professor & H.O.D CSE (Data Science), CMREC	

DECLARATION

This is to certify that the work reported in the present Major project entitled "**Liver Disease Prediction**" is a record of bonafide work done by us in the Department of Computer Science and Engineering (Data Science), CMR Engineering College, JNTU Hyderabad. The reports are based on the project work done entirely by us and not copied from any other source. We submit our project for further development by any interested students who share similar interests to improve the project in the future.

The results embodied in this Major project report have not been submitted to any other University or Institute for the award of any degree or diploma to the best of our knowledge and belief.

A. SHIRISHA	218R1A6704
J. AKHILA	218R1A6729
A. KUSHAL KUMAR	218R1A6702
K. MADHUKAR	218R1A6731

ACKNOWLEDGMENT

We are extremely grateful to **Dr. A. Srinivasula Reddy**, Principal and **Dr. M. Laxmaiah**, Professor & HOD, **Department of CSE (Data Science), CMR Engineering College** for their constant support.

We are extremely thankful to **Mr. B. Kumaraswamy**, Assistant Professor (Ph.D, JNTUH), Internal Guide, Department of CSE(DS), for his constant guidance, encouragement and moral support throughout the project.

We will be failing in duty if We do not acknowledge with grateful thanks to the authors of the references and other literatures referred in this Project.

We thank **Mrs. G. Shruthi** , Assistant Professor, CSE(DS) Department , Major Project Coordinator for her constant support in carrying out the project activities and reviews.

We express our thanks to all staff members and friends for all the help and co-ordination extended in bringing out this project successfully in time.

Finally, We are very much thankful to our parents who guided us for every step.

A. SHIRISHA	218R1A6704
J. AKHILA	218R1A6729
A. KUSHAL KUMAR	218R1A6702
K. MADHUKAR	218R1A6731

ABSTRACT

Liver diseases are a major global health concern, leading to significant morbidity and mortality. Early detection and accurate diagnosis are critical for effective treatment, but traditional diagnostic methods rely heavily on manual interpretation of liver function test results. This approach is often time-consuming, subjective, and prone to human error. To address these challenges, this project develops a machine learning-based liver disease prediction system using **XGBoost**, a powerful algorithm known for handling complex data relationships effectively. The system is trained on the **Indian Liver Patient Dataset**, which includes key clinical features such as **Total Bilirubin, Direct Bilirubin, Alkaline Phosphatase, SGOT, SGPT, Albumin, and the Albumin-Globulin Ratio**. The proposed model undergoes rigorous **data preprocessing, feature selection, and hyperparameter tuning** to optimize performance. Additionally, cross-validation is implemented to ensure generalizability, and feature importance analysis is conducted to enhance interpretability for healthcare professionals. By leveraging machine learning, this system aims to provide a more efficient and reliable method for liver disease detection. It reduces diagnostic errors, speeds up decision-making, and assists medical practitioners by offering data-driven insights. This approach has the potential to enhance clinical workflows, improve patient outcomes, and support early intervention strategies in liver disease management. By applying **data preprocessing, feature selection, model training, and hyperparameter tuning**, the proposed system enhances prediction accuracy while ensuring model interpretability. Additionally, feature importance analysis helps healthcare professionals understand key contributing factors in diagnosis. This machine learning-based approach can assist healthcare providers in making faster and more accurate decisions, ultimately improving patient care and reducing diagnostic errors.

CONTENTS

TOPIC	PAGE NO
1. INTRODUCTION	
1.1. Overview	1
1.2. Research Motivation	2
1.3. Problem Statement	3
1.4. Applications	4
2. LITERATURE SURVEY	5
3. EXISTING SYSTEM	
3.1. Traditional Diagnosis Methods	8
3.2. Drawbacks	9
4. PROPOSED SYSTEM	
4.1. Overview	11
4.2. Objectives of the Proposed System	12
4.3. Advantages of Proposed System	13
5. SYSTEM DESIGN	
5.1. Architecture of Proposed System	14
5.2. UML Diagrams	19
6. SOFTWARE ENVIRONMENT	
6.1. What is python and its Advantages and Disadvantages	26
6.2. History of python	29
6.3. Modules used in project	31
6.4. Installation of python	34
7. SYSTEM REQUIREMENTS SPECIFICATIONS	
7.1. Software Requirements	41
7.2. Hardware Requirements	41
8. FUNCTIONAL REQUIREMENTS	
8.1. Output Design	42
8.2. Input Design, Stages, Types, Media	43
8.3. User Interface	45
8.4. Performance Requirements	46
8.5. Testing	47
9. SOURCE CODE	48
10. RESULTS AND DISCUSSION	
10.1. Implementation and Dataset Description	55
11. CONCLUSION AND REFERENCES	61

LIST OF FIGURES

FIG.NO	DESCRIPTION	PAGENO
5.1.1	Architecture	14
5.2.1	Class Diagram	20
5.2.2	Use Case Diagram	21
5.2.3	Sequence Diagram	23
5.2.4	Activity Diagram	25
6.4	Python Installation Steps	34
10.1.1	The age group of the patients	57
10.1.2	Albumin and albumin and globulin ratio by a scatterplot	58
10.1.3	Gender based on the Protein Intake	58
10.1.4	male and female based on Albumin Level	59
10.1.5	Between the features using a heatmap	59
10.1.6	ROC Curve	60

LIST OF TABLES

TABLE. NO	DESCRIPTION	PAGENO
1	Literature Survey	06
2	Advantages of XGBoost Over Other Algorithms	18
3	Feature Description	55

1. INTRODUCTION

1.1 Overview

Liver diseases are among the leading causes of morbidity and mortality worldwide, affecting millions of people each year. Early diagnosis is crucial for effective treatment, but traditional diagnostic methods rely on manual interpretation of liver function test results, which can be time-consuming, subjective, and prone to human error.

This project aims to develop an automated liver disease prediction system using machine learning techniques, particularly XGBoost, to classify patients as having liver disease or not based on their clinical features. The system is trained on the Indian Liver Patient Dataset, which includes key attributes such as Total Bilirubin, Direct Bilirubin, Alkaline Phosphatase, SGOT, SGPT, Albumin, and the Albumin-Globulin Ratio.

To ensure accuracy and reliability, the proposed system involves several steps:

- Data Preprocessing – Handling missing values, encoding categorical data, and normalizing numerical features.
- Feature Selection – Identifying the most relevant attributes for improved model performance.
- Model Training and Optimization – Using XGBoost with hyperparameter tuning to enhance predictive accuracy.
- Model Evaluation – Assessing the model's performance using metrics such as accuracy, precision, recall, and F1-score.
- Result Interpretation – Analyzing feature importance to improve the transparency and trustworthiness of predictions.

By leveraging machine learning, this system enhances diagnostic accuracy, reduces human error, and speeds up the decision-making process for healthcare professionals. The ultimate goal is to provide a reliable, data-driven approach to liver disease diagnosis, improving patient outcomes and supporting early intervention strategies.

1.2 Research Motivation

Liver diseases are a significant public health concern, contributing to high morbidity and mortality rates worldwide. Despite advancements in medical diagnostics, early detection and accurate classification of liver diseases remain challenging due to the reliance on manual interpretation of biochemical markers. Traditional diagnostic methods, such as liver function tests (LFTs), require expert evaluation and may vary in accuracy due to human error or subjective interpretation.

The motivation for this research stems from several key challenges:

1. **Time-Consuming Diagnosis** – The manual analysis of liver function test results can be slow, delaying timely intervention.
2. **Risk of Human Error** – Variations in clinical expertise may lead to misinterpretation of test results, affecting diagnosis and treatment.
3. **Complex Interactions Between Features** – Traditional methods often fail to capture complex relationships between biochemical parameters that indicate liver disease.
4. **Limited Access to Specialists** – In many regions, especially in rural or underdeveloped areas, access to trained hepatologists and diagnostic tools is limited.

To address these challenges, this research focuses on **leveraging machine learning (ML) techniques, particularly XGBoost**, to develop an automated system for liver disease prediction. **XGBoost** is chosen for its ability to handle **non-linear relationships, missing data, and feature interactions**, leading to **more accurate and efficient predictions**.

By integrating ML-based predictive models, this study aims to:

- **Reduce diagnostic time** and provide faster results.
- **Minimize human error** by relying on data-driven insights.
- **Enhance interpretability** through feature importance analysis, helping medical professionals understand key contributing factors.
- **Improve accessibility** by offering a scalable solution that can be used in resource-limited settings.

This research contributes to the advancement of **AI-driven healthcare solutions**, supporting more reliable, efficient, and accessible liver disease diagnosis, ultimately improving patient care and outcomes.

1.3 Problem Statement

Liver diseases are a major global health issue, causing significant morbidity and mortality. Accurate and early diagnosis is essential for effective treatment and management. However, current diagnostic methods rely heavily on **manual interpretation of liver function test results**, which can be **time-consuming, prone to human error, and dependent on clinical expertise**. Additionally, traditional rule-based systems often fail to capture the **complex interactions between multiple biochemical parameters**, leading to inconsistent or delayed diagnoses.

To address these challenges, this project proposes the development of a **machine learning-based liver disease prediction system** using **XGBoost**.

The goal is to create an automated model that can:

- **Analyze clinical features efficiently** and provide faster results.
- **Improve diagnostic accuracy** by reducing human error.
- **Identify key contributing factors** using feature importance analysis to enhance interpretability.

The proposed system will be trained on a **liver patient dataset**, which includes key clinical markers such as **Total Bilirubin, Direct Bilirubin, Alkaline Phosphatase, SGOT, SGPT, Albumin, and the Albumin-Globulin Ratio**. The model will undergo **data preprocessing, feature selection, and hyperparameter tuning** to optimize performance.

By leveraging machine learning, this research aims to **develop a robust and interpretable predictive model that can assist healthcare professionals in diagnosing liver disease more accurately and efficiently**, ultimately improving patient outcomes.

1.4 Applications

The implementation of a machine learning-based liver disease prediction system has numerous applications across various domains in healthcare, medical research, and public health. By leveraging data-driven insights, this system can enhance diagnostic accuracy, improve healthcare accessibility, and support early intervention strategies.

1. Clinical Diagnosis and Decision Support

- Assists **healthcare professionals** in diagnosing liver disease more accurately and efficiently by analyzing patient biochemical markers.
- Reduces **reliance on manual interpretation** of liver function tests, minimizing human error and subjectivity in diagnosis.
- Provides a **decision support tool** that complements traditional diagnostic methods, ensuring comprehensive patient evaluation.

2. Early Detection and Preventive Healthcare

- Enables **early identification** of liver disease risk, allowing for timely medical intervention and lifestyle modifications.
- Helps in **preventive healthcare initiatives** by identifying at-risk individuals before symptoms become severe, reducing the burden on healthcare systems.
- Supports **public health awareness campaigns** by providing data-driven insights into common liver disease risk factors.

3. Integration into Healthcare Systems and Telemedicine

- Can be integrated into **electronic health record (EHR) systems** to provide real-time liver disease risk assessments.
- Enhances **telemedicine services** by enabling remote diagnosis, especially in regions with limited access to specialized hepatologists.
- Supports **AI-driven health applications** that offer personalized recommendations for at-risk individuals.

4. Research and Drug Development

- Assists medical researchers in **understanding the relationships between different biochemical markers and liver disease progression**.
- Provides a **data-driven foundation for clinical trials** by identifying patient groups at different stages of liver disease.

2. LITERATURE SURVEY

The prediction and diagnosis of liver diseases have been widely studied using machine learning (ML) and artificial intelligence (AI) techniques. Traditional diagnostic methods rely on biochemical markers and imaging techniques, but these approaches are often time-consuming, subjective, and prone to human error. Researchers have explored various ML-based models to enhance accuracy, efficiency, and interpretability in liver disease classification.

Traditional Methods for Liver Disease Diagnosis

Historically, liver disease diagnosis has been based on:

- **Liver Function Tests (LFTs):** Measurement of biochemical markers such as bilirubin, albumin, and enzyme levels to assess liver health.
- **Rule-Based Expert Systems:** Algorithms designed using pre-defined medical knowledge to classify liver conditions.
- **Ultrasound and Imaging Techniques:** Used for detecting structural liver abnormalities, often requiring expert interpretation.

Machine Learning Approaches in Liver Disease Prediction

Several studies have explored ML techniques for liver disease classification using clinical datasets.

- **Logistic Regression & Decision Trees:**
 - Early studies applied **logistic regression** and **decision tree classifiers** to predict liver disease based on biochemical features.
- **Support Vector Machines (SVM):**
 - SVMs have been used to classify liver disease patients, demonstrating **good accuracy** but requiring **extensive feature engineering** and computational resources.

S.No	Author	Title	Year	Contributions
1.	Minnoor M., Baths V.	Liver Disease Diagnosis Using Machine Learning	2022	The research applied several ML techniques, comparing SVM, Decision Trees, and Random Forest for liver disease diagnosis.
2.	Srivastava A., Kumar V., Mahesh T.	Automated Prediction of Liver Disease Using Machine Learning Algorithms	2022	The paper explored automated liver disease prediction using algorithms like Logistic Regression and Random Forest.
3.	Newaz A., Ahmed N., Haq F.	Diagnosis of Liver Disease Using Cost-Sensitive Learning	2021	This paper introduced cost-sensitive learning approaches to liver disease prediction, particularly for imbalanced datasets
4.	Kuzhippallil M. A., Joseph C., Kannan A.	Comparative Analysis of Machine Learning Techniques for Indian Liver Disease Patients	2020	This study focused on Indian liver disease patients, comparing various machine learning algorithms, including Logistic Regression, Decision Trees, and Random

				Forest, in predicting disease occurrence.
5.	Souptik Dutta, Subhash Mondal, Amitava Nag	Prediction of Liver Disease Using Machine Learning Approaches Based on KNN Model	2020	This study compared the performance of KNN and Naïve Bayes algorithms for liver disease prediction, showing that KNN was more effective when used with oversampling to handle class imbalance.

Table 1: Literature Survey

3. EXISTING SYSTEM

In current clinical practice, liver disease diagnosis relies on traditional manual analysis of liver function test (LFT) results, which is performed by medical professionals based on biochemical parameters. Additionally, some hospitals use rule-based expert systems that assist in diagnosing liver diseases by following predefined medical guidelines. However, these existing approaches have several limitations.

3.1 Traditional Diagnosis Methods

The primary methods for liver disease diagnosis include:

1. Liver Function Tests (LFTs):

- These tests measure key biochemical markers such as **Total Bilirubin, Direct Bilirubin, Alkaline Phosphatase, SGOT, SGPT, and Albumin**.
- Physicians interpret these values based on standard medical knowledge to diagnose liver disease.

2. Medical Imaging Techniques:

- **Ultrasound, MRI, and CT scans** are used to detect structural abnormalities in the liver.
- These methods require **expert radiologists** and are often expensive and time-consuming.

3. Rule-Based Expert Systems:

- Some hospitals use computer-based rule engines that take input from clinicians and match it against predefined rules to suggest a diagnosis.
- These systems lack flexibility and do not learn from new data, limiting their adaptability.

3.2 Limitations of the Existing System

The traditional approaches suffer from several drawbacks:

1. Time-Consuming Process

- Manual interpretation of LFTs and medical imaging requires significant time, leading to delays in diagnosis and treatment.

2. Prone to Human Error

- Variations in physician expertise can result in **misdiagnosis** or **inconsistent results**, affecting patient outcomes.

3. Lack of Consideration for Complex Feature Interactions

- Liver diseases often involve **multiple interacting biochemical markers**.
- Manual analysis may fail to detect subtle relationships between these markers, leading to incorrect or delayed diagnosis.

4. Limited Scalability

- Expert-based systems require continuous manual updates to remain effective.
- In resource-limited settings, access to trained hepatologists is often scarce.

5. Rule-Based Systems Cannot Adapt to New Data

- These systems **do not improve over time** as they rely on predefined rules rather than **data-driven learning**

Drawbacks

Despite its widespread use, the existing liver disease diagnosis system has several limitations that hinder its accuracy, efficiency, and scalability. These drawbacks emphasize the need for a more advanced, automated, and data-driven approach.

1. Time-Consuming Process

- Traditional diagnosis methods, such as **manual analysis of liver function tests (LFTs)** and **medical imaging**, require significant time for evaluation.
- Delayed diagnosis can lead to **progression of liver diseases**, making treatment less effective.

2. Prone to Human Error

- Diagnosis heavily depends on the expertise and experience of healthcare professionals, which can result in **inconsistent interpretations** of biochemical markers.

3. Limited Ability to Handle Complex Data

- The **interaction between multiple biochemical parameters** in liver disease is often complex. Traditional rule-based systems **fail to capture non-linear relationships** between features, leading to potential misclassification.

4. Lack of Standardization

- Different healthcare facilities follow **varied diagnostic guidelines and criteria**, making liver disease detection **inconsistent across different hospitals and regions**.
- Manual diagnosis methods are **not universally standardized**, which can lead to variations in patient assessments.

5. Limited Scalability

- The reliance on **expert hepatologists and radiologists** makes it difficult to scale liver disease diagnosis in regions with limited healthcare resources.
- In **remote and underdeveloped areas**, access to advanced diagnostic tools and specialists is **scarce**, leading to delayed or inaccurate diagnoses.

6. Rule-Based Systems Do Not Learn Over Time

- Traditional **expert systems** use fixed rules and do not improve with new patient data.
- These systems cannot **adapt to evolving medical knowledge** or account for variations in patient demographics, leading to outdated diagnostic models.

7. High Cost of Medical Imaging Techniques

- Advanced imaging techniques like **MRI, CT scans, and ultrasounds** are expensive and not always accessible to patients in **low-income regions**.
- These methods also require **specialized equipment and trained professionals**, increasing healthcare costs.

8. Data Privacy and Security Concerns

- The use of electronic health records (EHR) and digital diagnostic tools raises **data privacy issues**.
- Without robust security measures, patient data could be **vulnerable to breaches and misuse**.

4. PROPOSED METHODOLOGY

4.1 Overview

To overcome the limitations of the existing liver disease diagnosis methods, this project proposes a machine learning-based predictive system that leverages clinical features to accurately classify patients as having liver disease or not. The system employs XGBoost, a powerful ensemble learning algorithm, to enhance diagnostic accuracy, efficiency, and interpretability.

The proposed method for **Liver Disease Prediction using XGBoost** aims to develop an efficient and accurate system that automates the diagnosis of liver disease based on patient clinical data. The approach involves several key steps, starting with data collection and preprocessing.

The dataset, specifically the **Indian Liver Patient Dataset (ILPD)**, contains crucial biochemical and demographic features such as Age, Gender, Total Bilirubin, Alkaline Phosphatase, and Albumin/Globulin Ratio. Preprocessing techniques are applied to handle missing values, encode categorical variables, and normalize numerical data for better model performance.

Feature selection is then performed to identify the most relevant attributes influencing the prediction. The XGBoost algorithm, known for its **high accuracy and robustness**, is trained on the processed dataset using optimized hyperparameters. The model's performance is evaluated based on metrics such as **accuracy, AUC-ROC, precision, and recall** to ensure reliable classification.

Once trained, the system allows users to input patient details and receive a **real-time prediction** on whether the patient has liver disease. This automated method improves diagnostic efficiency, supports clinical decision-making, and helps in the early detection of liver disease, ultimately enhancing patient outcomes.

4.2 Objectives of the Proposed System

1. **Automate liver disease prediction** based on patient clinical data.
2. **Enhance diagnostic accuracy** by using machine learning techniques.
3. **Reduce human dependency** and minimize manual errors in diagnosis.
4. **Provide feature importance analysis** to ensure interpretability in medical decision-making.
5. **Optimize model performance** using hyperparameter tuning and cross-validation.

The primary objective of the proposed Liver Disease Prediction System using XGBoost is to develop an automated, efficient, and reliable model for diagnosing liver disease based on clinical parameters. By leveraging the XGBoost algorithm, the system aims to achieve high predictive accuracy, ensuring robust classification of patients into liver disease and non-liver disease categories.

The model will analyze key features such as Age, Gender, Total Bilirubin, Direct Bilirubin, Alkaline Phosphatase, Alamine Aminotransferase (SGPT), Aspartate Aminotransferase (SGOT), Total Proteins, Albumin, and Albumin/Globulin Ratio to identify patterns indicative of liver disease.

Additionally, the system will provide insights into feature importance, helping medical professionals understand which parameters significantly contribute to the diagnosis. The ultimate goal is to assist healthcare providers in making informed decisions, improving early detection, and enhancing patient care through an efficient, scalable, and data-driven approach.

4.3 Advantages of the Proposed System

The proposed Liver Disease Prediction System using XGBoost offers several advantages over traditional diagnostic methods. Firstly, it enhances accuracy and efficiency by leveraging machine learning techniques to analyze complex medical data, reducing the chances of human error. Secondly, the system ensures early detection of liver disease, allowing timely medical intervention and improving patient outcomes.

Thirdly, it provides a data-driven approach by highlighting the most influential clinical features, aiding doctors in better understanding disease patterns. Additionally, the model is scalable and adaptable, making it suitable for diverse datasets and real-world healthcare applications.

The system is also cost-effective, as it minimizes the need for expensive diagnostic tests by providing an initial assessment based on readily available medical parameters. Lastly, it improves decision support by assisting healthcare professionals with reliable predictions, ultimately enhancing patient care and treatment strategies.

✓ **Higher Accuracy** – XGBoost outperforms traditional methods by effectively capturing feature interactions.

✓ **Faster Diagnosis** – Automated prediction reduces time required for manual interpretation.

✓ **Minimized Human Error** – Reduces subjectivity in liver disease diagnosis.

✓ **Interpretability** – Feature importance analysis ensures that the system remains transparent for medical professionals.

✓ **Scalability** – Can be integrated into healthcare systems and telemedicine platforms for widespread use.

5. SYSTEM DESIGN

5.1 Architecture Of Proposed System

The proposed system follows a structured pipeline to predict liver disease:

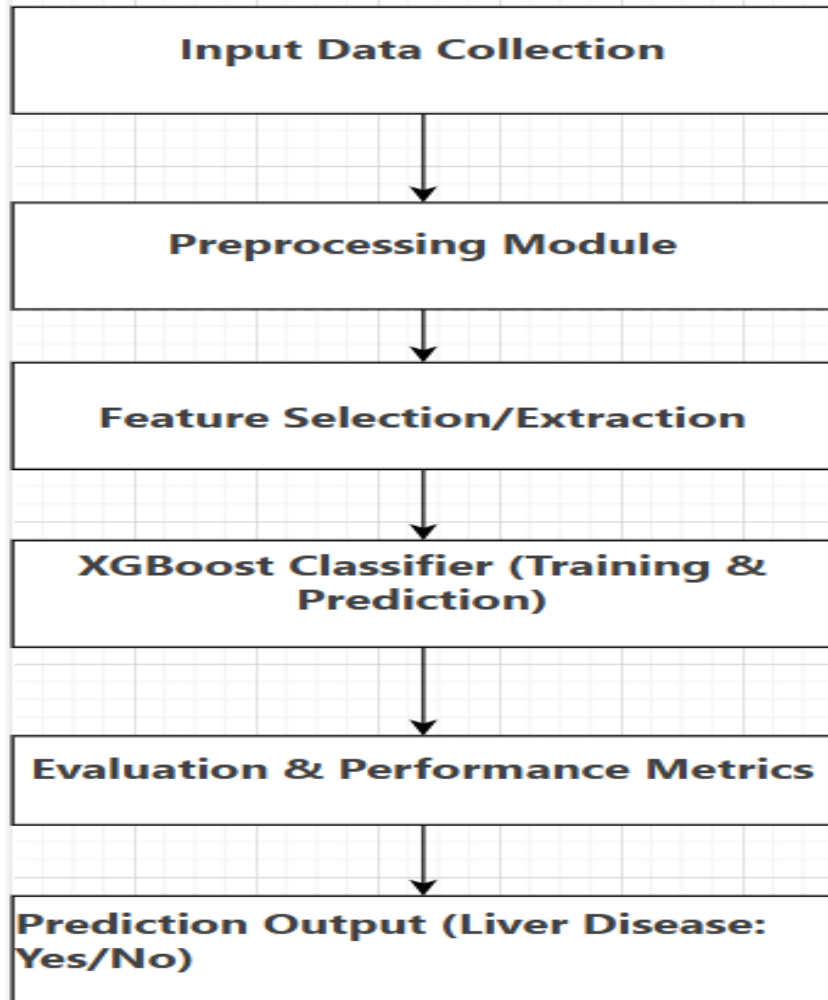


Fig. 5.1.1 Architecture

1. Data Collection

- The system utilizes a dataset containing **biochemical markers** such as **Total Bilirubin, Albumin, Alkaline Phosphatase, SGOT, SGPT, Age, and Gender**.
- The dataset is sourced from the **Indian Liver Patient Dataset (ILPD)** or similar medical databases.

2. Data Preprocessing

- **Handling missing values** to ensure dataset completeness.
- **Encoding categorical features** (e.g., converting Gender into numerical values).

3. Feature Selection

- The system identifies the **most relevant clinical attributes** that contribute to liver disease prediction.
- Feature importance is determined using **XGBoost's built-in feature selection** method.

4. Model Training using XGBoost

- The dataset is **split into training and testing sets** to evaluate model performance.
- The **XGBoost classifier** is trained on the dataset using hyperparameter tuning for optimization.
- **Cross-validation** ensures the model generalizes well to unseen data.

5. Model Testing & Evaluation

- The trained model is tested on new patient records.
- Performance metrics such as **accuracy, precision, recall, and F1-score** are calculated to measure effectiveness.

6. Prediction & Diagnosis

- Given a new patient's biochemical values, the system predicts **whether they have liver disease or not**.
- The system outputs **probability scores** and highlights **key features contributing to the diagnosis**.

7. Result Interpretation & Decision Support

- **Feature importance analysis** provides insights into **which biomarkers are most critical** in liver disease detection.
- Physicians can use this information to **validate AI-driven predictions** and support clinical decisions.

XGBoost Algorithm

XGBoost (Extreme Gradient Boosting) is a powerful and efficient machine learning algorithm widely used for classification and regression tasks. It is an optimized version of gradient boosting, designed to be fast, scalable, and highly accurate. XGBoost is particularly popular in data science competitions and real-world applications due to its superior performance.

1. Introduction to XGBoost

XGBoost is an ensemble learning method based on gradient boosting, where multiple decision trees are trained sequentially, and each new tree corrects the errors of the previous ones. It minimizes loss using gradient descent and optimizes both bias and variance, making it robust against overfitting.

Key Features of XGBoost:

- High Performance: Optimized for speed and efficiency.
- Regularization: Uses L1 (Lasso) and L2 (Ridge) regularization to prevent overfitting.
- Handling Missing Data: Automatically learns the best imputation strategy.
- Parallel Processing: Can utilize multiple CPU cores for faster computation.
- Scalability: Works well with large datasets.

2. Working of XGBoost

XGBoost follows the gradient boosting framework, where weak learners (decision trees) are built sequentially to minimize errors.

Step-by-Step Process:

1. Initialize Predictions:
 - Start with a base model (e.g., predicting the mean of the target variable).
2. Compute Residual Errors:
 - Calculate the difference between actual values and predicted values (residuals).
3. Train a New Decision Tree on Residuals:
 - A new tree is trained to predict the residuals (errors) from the previous model.
 - The output of this tree is used to adjust the previous predictions.
4. Update Predictions:
 - The new model's predictions are added to the previous ones with a learning rate (shrinkage factor) to control the update's impact.
5. Repeat Steps 2-4:
 - More trees are added sequentially until the stopping condition is met (e.g., a specified number of trees or no improvement in validation error).
6. Final Prediction:
 - The predictions from all trees are combined to give the final output.

3. XGBoost's Mathematical Formulation

Let's define the objective function used in XGBoost:

$$L = \sum_{i=1}^n \text{nl}(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k) \quad \text{mathcal{L}} = \sum_{i=1}^n \text{nl}(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

Where:

- L = total loss function.
- $\text{nl}(y_i, \hat{y}_i)$ = loss between actual and predicted values (e.g., Mean Squared Error).
- $\Omega(f_k)$ = regularization term to control model complexity.
- K = number of trees.

Tree Structure Representation:

Each decision tree in XGBoost is represented as:

$$f_t(x) = w_{q(x)} \quad f_t(x) = w_{q(x)}$$

Where:

- $q(x)$ = decision rule that assigns a leaf to each input.
- w = leaf weights (model parameters).

4. Regularization in XGBoost

Unlike standard gradient boosting, XGBoost includes regularization to reduce overfitting:

- L1 Regularization (Lasso): Encourages sparsity in leaf weights, reducing complexity.
- L2 Regularization (Ridge): Penalizes large weights, making the model more stable.

The regularization term is:

$$\Omega(f) = \gamma T + \frac{1}{2} \sum_{j=1}^T w_j^2 \quad \Omega(f) = \gamma T + \frac{1}{2} \sum_{j=1}^T w_j^2$$

Where:

- T = number of leaves.
- γ = complexity penalty (controls tree growth).
- λ = L2 regularization weight.

5. Advantages of XGBoost

High Accuracy: Performs better than traditional boosting algorithms.
Handles Missing Values: Can handle missing data without explicit imputation.
Regularization: Reduces overfitting through L1 and L2 penalties.
Scalability: Works well with large datasets.

6. Applications of XGBoost

- Medical Diagnosis: Predicting diseases like liver disease, diabetes, and cancer.
- Finance: Credit risk assessment, fraud detection, and stock price prediction.
- Image Classification: Feature extraction and anomaly detection.
- Natural Language Processing (NLP): Text classification, sentiment analysis.
- Recommendation Systems: Personalized product recommendations.

7. Conclusion

XGBoost is a powerful, efficient, and accurate machine learning algorithm, widely used in classification and regression tasks. Its combination of gradient boosting, regularization, and parallelization makes it one of the best-performing models in data science. Proper hyperparameter tuning and feature engineering can further enhance its effectiveness, making it a top choice for predictive modeling.

Advantages of XGBoost Over Other Algorithms

Feature	XGBoost	Random Forest	Logistic Regression
Accuracy	High	Moderate	Low
Handles Missing Data	Yes	No	No
Overfitting Control	Regularization	Prone to overfitting	No control
Feature Importance	Yes	Partial	No
Speed	Fast	Slower	Fast

Table 2: Advantages of XGBoost Over Other Algorithms

5.2 UML DIAGRAMS

UML stands for Unified Modeling Language. UML is a standardized general-purpose modeling language in the field of object-oriented software engineering. The standard is managed, and was created by, the Object Management Group. The goal is for UML to become a common language for creating models of object-oriented computer software. In its current form UML is comprised of two major components: a Meta-model and a notation. In the future, some form of method or process may also be added to or associated with, UML.

The Unified Modeling Language is a standard language for specifying, Visualization, Constructing and documenting the artifacts of software system, as well as for business modeling and other non-software systems. The UML represents a collection of best engineering practices that have proven successful in the modeling of large and complex systems. The UML is a very important part of developing objects-oriented software and the software development process. The UML uses mostly graphical notations to express the design of software projects.

GOALS: The Primary goals in the design of the UML are as follows:

- Provide users a ready-to-use, expressive visual modeling Language so that they can develop and exchange meaningful models.
- Provide extendibility and specialization mechanisms to extend the core concepts.
- Be independent of particular programming languages and development process.
- Provide a formal basis for understanding the modeling language.
- Encourage the growth of OO tools market.
- Support higher level development concepts such as collaborations, frameworks, patterns and components.
- Integrate best practices.

Class diagram

The class diagram is used to refine the use case diagram and define a detailed design of the system. The class diagram classifies the actors defined in the use case diagram into a set of interrelated classes. The relationship or association between the classes can be either an "is-a" or "has-a" relationship. Each class in the class diagram may be capable of providing certain functionalities. These functionalities provided by the class are termed "methods" of the class. Apart from this, each class may have certain "attributes" that uniquely identify the class.

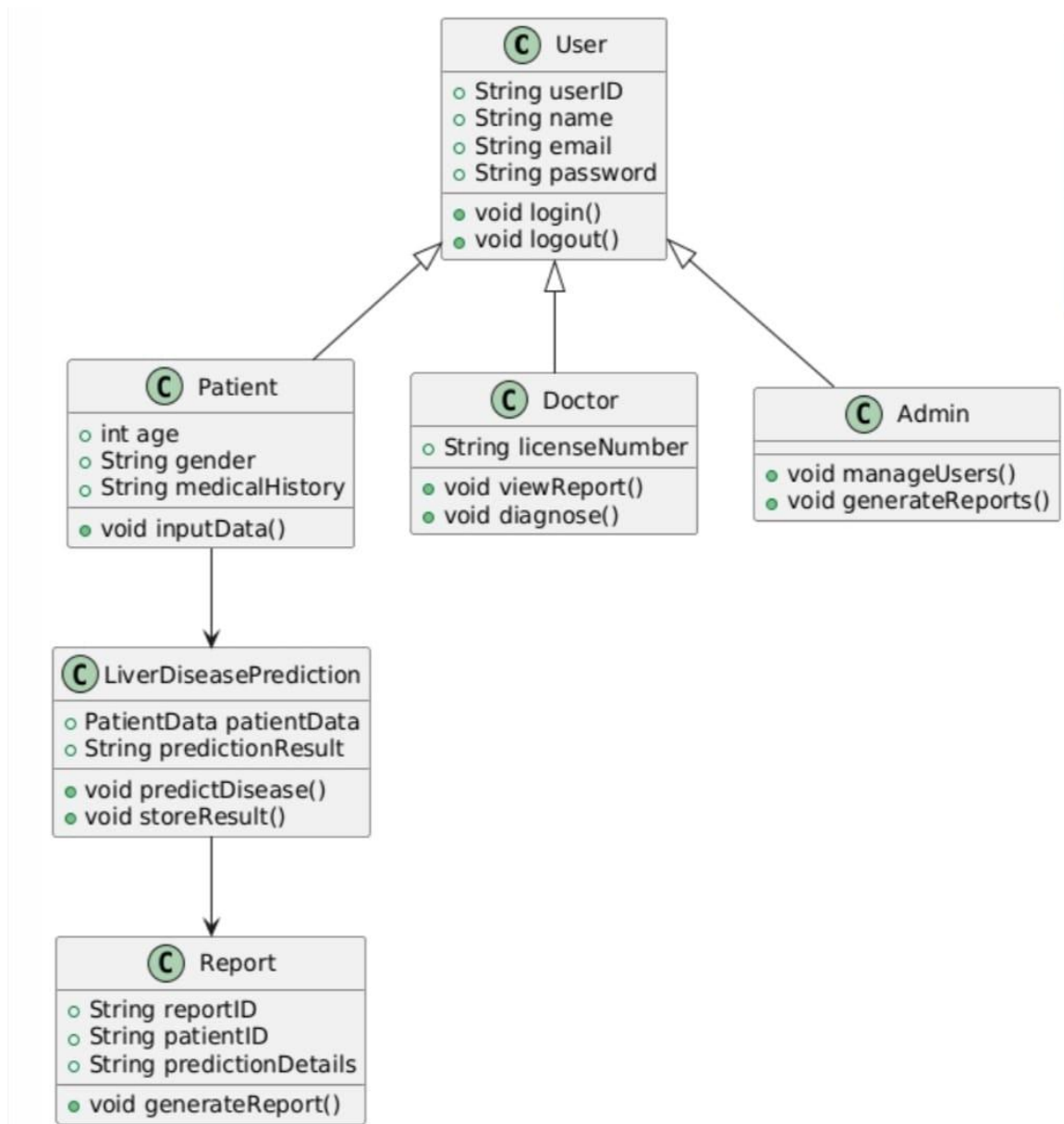


Fig. 5.2.1 Class Diagram

Use case Diagram

A use case diagram in the Unified Modeling Language (UML) is a type of behavioral diagram defined by and created from a Use-case analysis. Its purpose is to present a graphical overview of the functionality provided by a system in terms of actors, their goals (represented as use cases), and any dependencies between those use cases. The main purpose of a use case diagram is to show what system functions are performed for which actor. Roles of the actors in the system can be depicted.

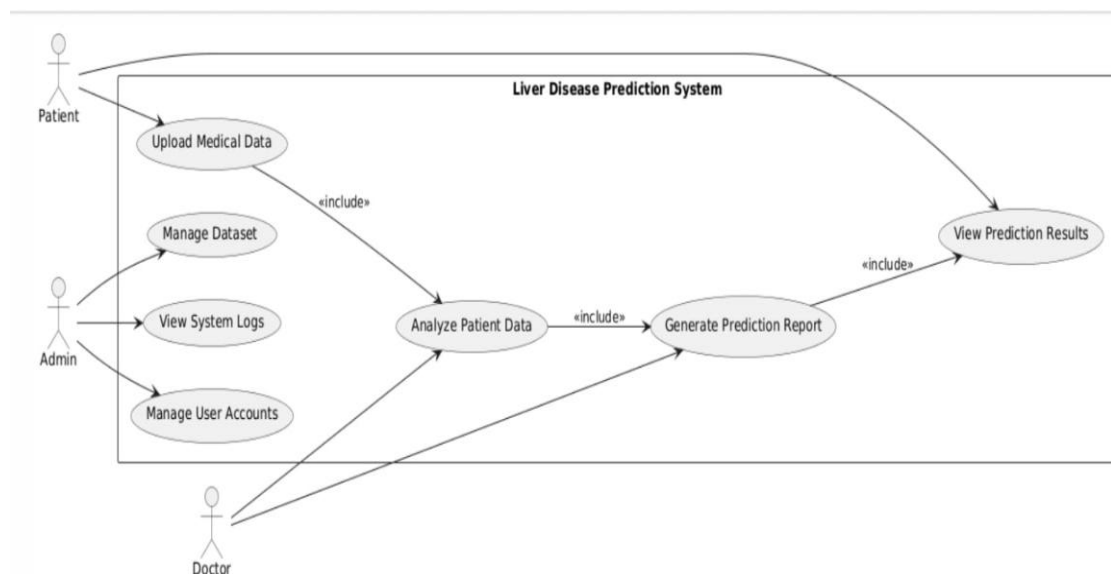


Fig. 5.2.2 Use Case Diagram

Sequence Diagram

A **sequence diagram** in Unified Modeling Language (UML) is a kind of interaction diagram that shows how processes operate with one another and in what order. It is a construct of a Message Sequence Chart. A sequence diagram shows, as parallel vertical lines ("lifelines"), different processes or objects that live simultaneously, and as horizontal arrows, the messages exchanged between them, in the order in which they occur. This allows the specification of simple runtime scenarios in a graphical manner.

A **sequence diagram** represents the interaction between different system components in a step-by-step manner. It shows how data flows between the user, system, and machine learning model for liver disease prediction.

Actors & Components

1. **User (Patient/Doctor)** – Inputs patient details for prediction.
2. **GUI (User Interface)** – Collects user input and sends it to the backend.
3. **Prediction System (Backend)** – Preprocesses the data and applies the XGBoost model.
4. **XGBoost Model** – Analyzes input data and makes a prediction.
5. **Output Module** – Displays the prediction result to the user.

Process Flow

1. The **user** enters the required clinical parameters such as Age, Gender, Bilirubin levels, and other liver-related values.
2. The **GUI** collects the input and sends it to the backend system.
3. The **backend** preprocesses the data, normalizes values, and converts categorical features into numerical form.
4. The preprocessed data is passed to the **XGBoost model**, which performs the prediction based on trained patterns.
5. The **model** returns the prediction (either "Liver Disease Detected" or "No Liver Disease Detected") to the backend.

The **backend** forwards the result to the **output module**, which displays.

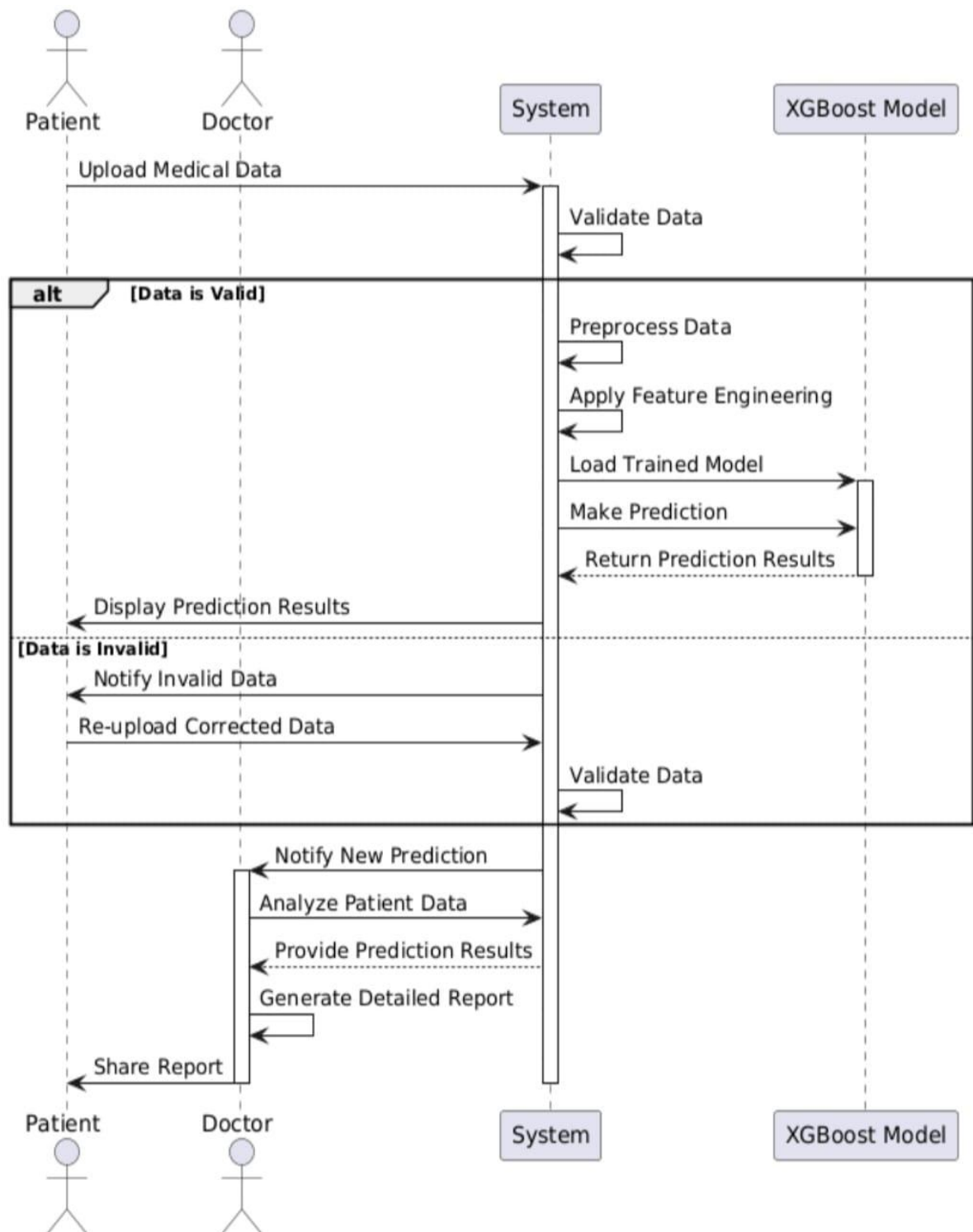


Fig. 5.2.3 Sequence Diagram

Activity diagram:

An **Activity Diagram** is a type of UML (Unified Modeling Language) diagram used to visually represent workflows of a system, business process, or software application. It shows how different activities interact with decisions, parallel operations, and flow of execution.

Key Features of Activity Diagrams

1. **Depicts Workflow:** It models the flow of control and data between different activities in a system.
2. **Represents Parallelism:** It shows concurrent activities that can occur simultaneously.
3. **Decision Making:** It includes decision nodes to represent alternative paths.
4. **Focus on Process Execution:** It emphasizes how tasks are performed rather than system structure.

Elements of an Activity Diagram

- **Initial Node:** Represents the starting point of the workflow (black filled circle).
- **Activity/Action:** A step in the process (rounded rectangle).
- **Decision Node:** A conditional branching point (diamond shape).
- **Merge Node:** Combines multiple paths into one (diamond shape).
- **Fork Node:** Represents parallel processing (a horizontal or vertical bar).
- **Join Node:** Combines parallel activities back into a single flow (a horizontal or vertical bar).
- **Final Node:** Represents the end of the process (black filled circle with a border).
- **Swimlanes:** Used to group actions performed by different actors or departments.
- **Arrows:** Show the flow of execution between elements.

When to Use Activity Diagrams

- To model business processes.
- To describe use case workflows.
- To illustrate algorithms or control flows in software.
- To represent system-level behaviors in software design.

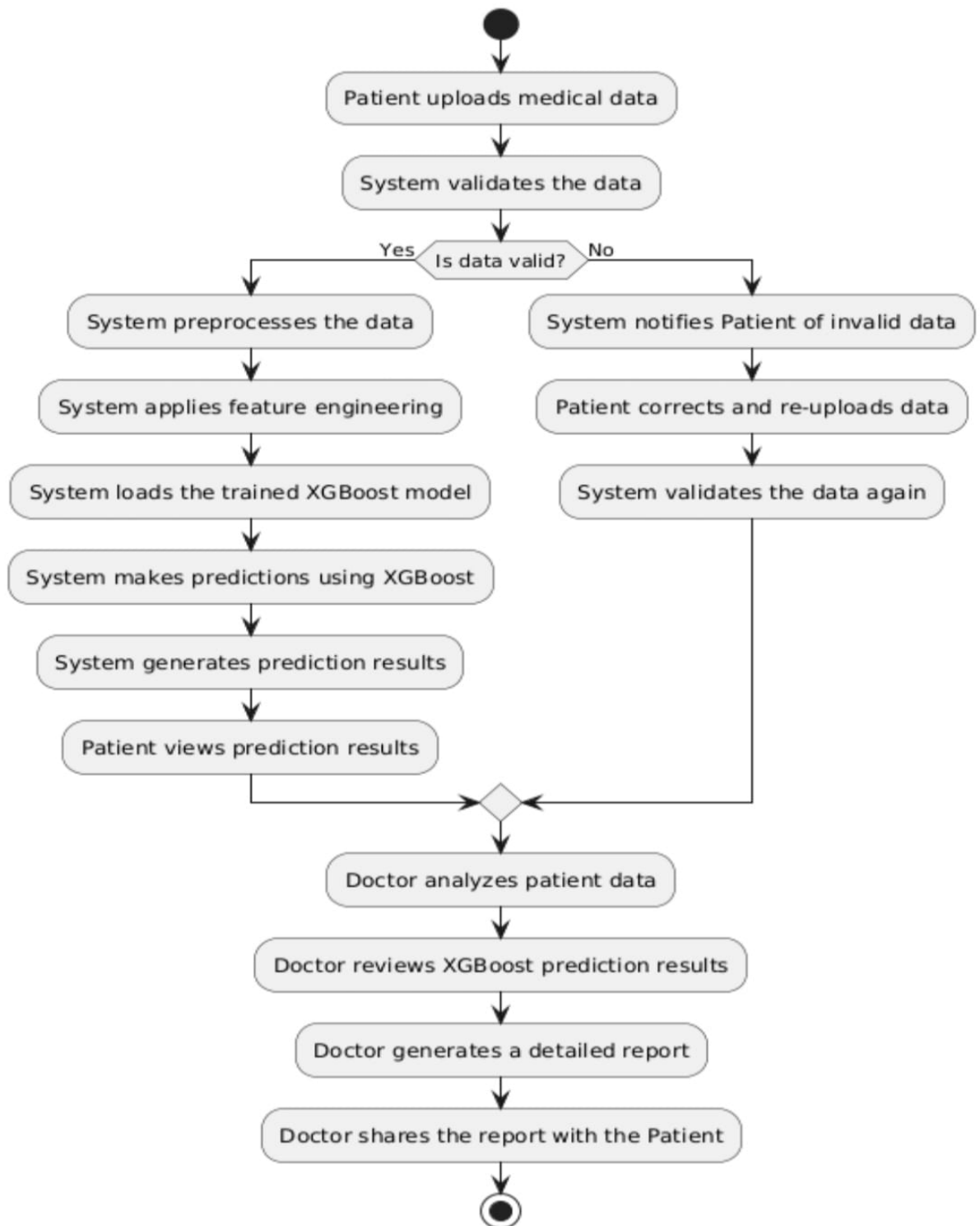


Fig. 5.2.4 Activity Diagram

6. SOFTWARE ENVIRONMENT

6.1 What is Python?

Below are some facts about Python.

- Python is currently the most widely used multi-purpose, high-level programming language.
- Python allows programming in Object-Oriented and Procedural paradigms. Python programs generally are smaller than other programming languages like Java.
- Programmers have to type relatively less and indentation requirement of the language, makes them readable all the time.
- Python language is being used by almost all tech-giant companies like – Google, Amazon, Facebook, Instagram, Dropbox, Uber... etc.

The biggest strength of Python is huge collection of standard libraries which can be used for the following –

- Machine Learning
- GUI Applications (like Kivy, Tkinter, PyQt etc.)
- Web frameworks like Django (used by YouTube, Instagram, Dropbox)
- Image processing (like Opencv, Pillow)
- Web scraping (like Scrapy, BeautifulSoup, Selenium)
- Test frameworks
- Multimedia

Advantages of Python

Let's see how Python dominates over other languages.

1. Extensive Libraries

Python downloads with an extensive library and it contains code for various purposes like regular expressions, documentation-generation, unit-testing, web browsers, threading, databases, CGI, email, image manipulation, and more. So, we don't have to write the complete code for that manually.

2. Extensible

As we have seen earlier, Python can be extended to other languages. You can write some of your code in languages like C++ or C. This comes in handy, especially in projects.

3. Embeddable

Complimentary to extensibility, Python is embeddable as well. You can put your Python code in your source code of a different language, like C++. This lets us add scripting capabilities to our code in the other language.

4. Improved Productivity

The language's simplicity and extensive libraries render programmers more productive than languages like Java and C++ do. Also, the fact that you need to write less and get more things done.

5. IOT Opportunities

Since Python forms the basis of new platforms like Raspberry Pi, it finds the future bright for the Internet of Things. This is a way to connect the language with the real world.

6. Simple and Easy

When working with Java, you may have to create a class to print 'Hello World'. But in Python, just a print statement will do. It is also quite easy to learn, understand, and code. This is why when people pick up Python, they have a hard time adjusting to other more verbose languages like Java.

7. Readable

Because it is not such a verbose language, reading Python is much like reading English.

This is the reason why it is so easy to learn, understand, and code. It also does not need curly braces to define blocks, and indentation is mandatory. These further aids the readability of the code.

8. Object-Oriented

This language supports both the procedural and object-oriented programming paradigms. While functions help us with code reusability, classes and objects let us model the real world. A class allows the encapsulation of data and functions into one.

9. Free and Open-Source

Like we said earlier, Python is freely available. But not only can you download Python for free, but you can also download its source code, make changes to it, and even distribute it. It downloads with an extensive collection of libraries to help you with your tasks.

10. Portable

When you code your project in a language like C++, you may need to make some changes to it if you want to run it on another platform. But it isn't the same with Python. Here, you need to code only once, and you can run it anywhere. This is called Write Once Run Anywhere (WORA). However, you need to be careful enough not to include any system-dependent features.

11. Interpreted

Lastly, we will say that it is an interpreted language. Since statements are executed one by one, debugging is easier than in compiled languages.

Any doubts till now in the advantages of Python? Mention in the comment section.

Advantages of Python Over Other Languages

1. Less Coding

Almost all of the tasks done in Python requires less coding when the same task is done in other languages. Python also has an awesome standard library support, so you don't have to search for any third-party libraries to get your job done. This is the reason that many people suggest learning Python to beginners.

2. Affordable

Python is free therefore individuals, small companies or big organizations can leverage the free available resources to build applications. Python is popular and widely used so it gives you better community support.

The 2019 Github annual survey showed us that Python has overtaken Java in the most popular programming language category.

3. Python is for Everyone

Python code can run on any machine whether it is Linux, Mac or Windows. Programmers need to learn different languages for different jobs but with Python, you can professionally build web apps, perform data analysis and machine learning, automate things, do web scraping and also build games and powerful visualizations. It is an all-rounder programming language.

6.2 History of Python

What do the alphabet and the programming language Python have in common? Right, both start with ABC. If we are talking about ABC in the Python context, it's clear that the programming language ABC is meant. ABC is a general-purpose programming language and programming environment, which had been developed in the Netherlands, Amsterdam, at the CWI (Centrum Wiskunde & Informatica). The greatest achievement of ABC was to influence the design of Python. Python was conceptualized in the late 1980s. Guido van Rossum worked that time in a project at the CWI, called Amoeba, a distributed operating system. In an interview with Bill Venners¹, Guido van Rossum said: "In the early 1980s, I worked as an implementer on a team building a language called ABC at Centrum voor Wiskunde en Informatica (CWI). I don't know how well people know ABC's influence on Python. I try to mention ABC's influence because I'm indebted to everything I learned during that project and to the people who worked on it. "Later on in the same Interview, Guido van Rossum continued: "I remembered all my experience and some of my frustration with ABC. I decided to try to design a simple scripting language that possessed some of ABC's better properties, but without its problems. So, I started typing. I created a simple virtual machine, a simple parser, and a simple runtime. I made my own version of the various ABC parts that I liked. I created a basic syntax, used indentation for statement grouping instead of curly braces or begin-end blocks, and developed a small number of powerful data types: a hash table (or dictionary, as we call it), a list, strings, and numbers."

Python Development Steps

Guido Van Rossum published the first version of Python code (version 0.9.0) at alt.sources in February 1991. This release included already exception handling, functions, and the core data types of lists, dict, str and others. It was also object oriented and had a module system. Python version 1.0 was released in January 1994. The major new features included in this release were the functional programming tools lambda, map, filter and reduce, which Guido Van Rossum never liked. Six and a half years later in October 2000, Python 2.0 was introduced. This release included list comprehensions, a full garbage collector and it was supporting unicode. Python flourished for another 8 years in the versions 2.x before the next major release as Python 3.0 (also known as "Python 3000" and "Py3K") was released. Python 3 is

not backwards compatible with Python 2.x. The emphasis in Python 3 had been on the removal of duplicate programming constructs and modules, thus fulfilling or coming close to fulfilling the 13th law of the Zen of Python: "There should be one -- and preferably only one -- obvious way to do it." Some changes in Python 3:

- Print is now a function.
- Views and iterators instead of lists
- The rules for ordering comparisons have been simplified. E.g., a heterogeneous list cannot be sorted, because all the elements of a list must be comparable to each other.
- There is only one integer type left, i.e., int. long is int as well.
- The division of two integers returns a float instead of an integer. "/" can be used to have the "old" behaviour.

Python

Python is an interpreted high-level programming language for general-purpose programming. Created by Guido van Rossum and first released in 1991, Python has a design philosophy that emphasizes code readability, notably using significant whitespace.

Python features a dynamic type system and automatic memory management. It supports multiple programming paradigms, including object-oriented, imperative, functional and procedural, and has a large and comprehensive standard library.

- Python is Interpreted – Python is processed at runtime by the interpreter. You do not need to compile your program before executing it. This is similar to PERL and PHP.
- Python is Interactive – you can actually sit at a Python prompt and interact with the interpreter directly to write your programs.

Python also acknowledges that speed of development is important. Readable and terse code is part of this, and so is access to powerful constructs that avoid tedious repetition of code. Maintainability also ties into this may be an all but useless metric, but it does say something about how much code you have to scan, read and/or understand to troubleshoot problems or tweak behaviors. This speed of development, the ease with which a programmer of other languages can pick up basic Python skills and the huge standard library is key to another area where Python excels. All its tools have been quick to implement, saved a lot of time, and several of them have later been patched and updated by people with no Python background - without breaking.

6.3 Modules Used in Project

Modules Used in the Liver Disease Prediction Project

The **Liver Disease Prediction System** is divided into multiple functional modules, each performing a specific task to ensure **efficient data processing, model training, and prediction**. Below are the **key modules** used in the project:

1. Data Collection Module

Purpose:

- Collects patient data, including **clinical attributes** such as age, gender, total bilirubin, albumin levels, and enzyme markers.
- Loads the **Indian Liver Patient Dataset (ILPD)** or similar datasets.

Libraries Used:

```
import pandas as pd # For reading CSV files
```

```
import numpy as np # For numerical operations
```

Functionality:

- ✓ Load dataset using **Pandas**
- ✓ Handle missing values
- ✓ Convert categorical variables (e.g., gender) into numerical format

2. Data Preprocessing Module

Purpose:

- Cleans and prepares the data for training.
- Handles missing values, normalizes numerical features, and encodes categorical variables.

Libraries Used:

```
from sklearn.preprocessing import LabelEncoder, StandardScaler
```

Functionality:

- ✓ **Handle Missing Data** – Replace NaN values with median or mean.
- ✓ **Feature Encoding** – Convert **categorical variables** (e.g., gender) into numerical values.
- ✓ **Feature Scaling** – Normalize numerical features for better model performance.

3. Feature Selection Module

Purpose:

- Identifies the most important features contributing to liver disease prediction.
- Reduces dimensionality to improve model performance.

Libraries Used:

```
from sklearn.feature_selection import SelectKBest, f_classif
import matplotlib.pyplot as plt
import seaborn as sns
```

Functionality:

- ✓ Selects the **top k most relevant features** using **ANOVA (f_classif)**.
- ✓ **Visualizes feature importance** using bar charts or heatmaps.

4. Model Training Module

Purpose:

- Trains the **XGBoost** model using the preprocessed dataset.
- Splits data into **training and testing sets** for evaluation.

Libraries Used:

```
from sklearn.model_selection import train_test_split
import xgboost as xgb
```

Functionality:

- ✓ Splits data into **80% training, 20% testing**
- ✓ Initializes and trains the **XGBoost model**
- ✓ Optimizes hyperparameters for better accuracy

5. Model Evaluation Module

Purpose:

- Evaluates the trained model using different metrics such as **accuracy, precision, recall, F1-score, and confusion matrix**.

Libraries Used:

```
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
```

Functionality:

- ✓ Computes model performance metrics
- ✓ Displays a **classification report (precision, recall, F1-score)**

- ✓ **Plots a confusion matrix** to visualize model performance

6. Prediction Module

Purpose:

- Takes **new patient data** as input and predicts whether the patient has liver disease or not.

Libraries Used:

```
import numpy as np
```

Functionality:

- ✓ Accepts patient details as input
- ✓ **Applies trained XGBoost model** to predict liver disease
- ✓ Returns **probability scores and prediction (Yes/No)**

7. Result Visualization Module

Purpose:

- Provides visual representation of **model accuracy, feature importance, and data distributions**.

Libraries Used:

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

Functionality:

- ✓ **Plots feature importance** using bar charts
- ✓ **Displays model accuracy graph**
- ✓ **Shows correlation heatmap** of clinical attributes

8. Deployment Module

Purpose:

- Deploys the trained model using **Flask** or **Streamlit**, allowing doctors to input patient details and get predictions.

Libraries Used:

```
from flask import Flask, request, render_template
```

```
import pickle # For saving the trained model
```

Functionality:

- ✓ Builds a **user interface (UI) for liver disease prediction**
- ✓ Accepts user inputs via a **web form**
- ✓ **Displays real-time predictions**

These **functional modules** work together to build an efficient **Liver Disease Prediction System** using **machine learning and data analytics**. Each module plays a crucial role in **data handling, model training, evaluation, prediction, and deployment**, ensuring accurate and reliable results.

6.4 Install Python Step-by-Step in Windows and Mac

Python a versatile programming language doesn't come pre-installed on your computer devices. Python was first released in the year 1991 and until today it is a very popular high- level programming language. Its style philosophy emphasizes code readability with its notable use of great whitespace.

The object-oriented approach and language construct provided by Python enables programmers to write both clear and logical code for projects. This software does not come pre-packaged with Windows.

How to Install Python on Windows and Mac

There have been several updates in the Python version over the years. The question is how to install Python? It might be confusing for the beginner who is willing to start learning Python but this tutorial will solve your query. The latest or the newest version of Python is version 3.7.4 or in other words, it is Python 3.

Note: The python version 3.7.4 cannot be used on Windows XP or earlier devices.

Before you start with the installation process of Python. First, you need to know about your System Requirements. Based on your system type i.e., operating system and based processor, you must download the python version. My system type is a Windows 64-bit operating system. So, the steps below are to install python version 3.7.4 on Windows 7 device or to install Python 3. Download the Python Cheatsheet here. The steps on how to install Python on Windows 10, 8 and 7 are divided into 4 parts to help understand better.

Download the Correct version into the system

Step 1: Go to the official site to download and install python using Google Chrome or any other web browser. OR Click on the following link: <https://www.python.org>



Now, check for the latest and the correct version for your operating system. Step 2: Click on the Download Tab.



Step 3: You can either select the Download Python for windows 3.7.4 button in Yellow Color or you can scroll further down and click on download with respective to their version. Here, we are downloading the most recent python version for windows 3.7.4

Looking for a specific release?

Python releases by version number:

Release version	Release date		Click for more
Python 3.7.4	July 8, 2019	Download	Release Notes
Python 3.6.9	July 2, 2019	Download	Release Notes
Python 3.7.3	March 25, 2019	Download	Release Notes
Python 3.4.10	March 18, 2019	Download	Release Notes
Python 3.5.7	March 18, 2019	Download	Release Notes
Python 2.7.16	March 4, 2019	Download	Release Notes
Python 3.7.2	Dec. 24, 2018	Download	Release Notes

Step 4: Scroll down the page until you find the Files option.

Step 5: Here you see a different version of python along with the operating system.

Files

Version	Operating System	Description	MD5 Sum	File Size	GPU
Copied source tarball	Source release		68111671a5b2db4aef7b9ab01b0f9be	23017663	360
XZ compressed source tarball	Source release		d33e4aa66097051c3eca45ee3604803	17131432	360
macOS 64-bit/32-bit installer	Mac OS X	for Mac OS X 10.6 and later	6428b4fa7583daf1a442cba0ce08e6	34898416	360
macOS 64-bit installer	Mac OS X	for OS X 10.9 and later	5dd605c38217a45773bf5ea936b2a1f	28802845	360
Windows help file	Windows		063999573a2e96b2ac58ade0b4f7rd2	8111761	360
Windows x86-64 embeddable zip file	Windows	for AMD64/EM64T/x64	9809c3b7d2fee3b0a0e02184a40728a2	7504391	360
Windows x86-64 executable installer	Windows	for AMD64/EM64T/x64	a702b4b0aaf76d9b0b3943a5d3e5d3e00	26882368	360
Windows x86-64 web-based installer	Windows	for AMD64/EM64T/x64	28c31c5080b0d73a0e051a3b0351b4bd2	1362904	360
Windows x86 embeddable zip file	Windows		9fab38d18841279fda94123574139d8	6741626	360
Windows x86 executable installer	Windows		33c3802942a5444a3d0451479394789	25663848	360
Windows x86 web-based installer	Windows		1b670cfa0d117d802c30983ea371d87c	1324608	360

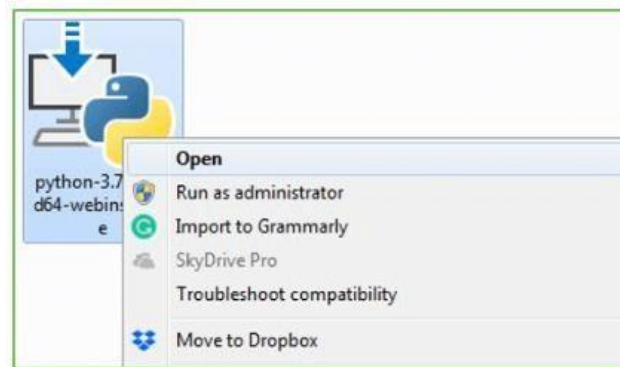
- To download Windows 32-bit python, you can select any one from the three options: Windows x86 embeddable zip file, Windows x86 executable installer or Windows x86 web-based installer.
- To download Windows 64-bit python, you can select any one from the three options: Windows x86-64 embeddable zip file, Windows x86-64 executable installer or Windows x86-64 web-based installer.

Here we will install Windows x86-64 web-based installer. Here your first part regarding which version of python is to be downloaded is completed. Now we move ahead with the second part in installing python i.e., Installation

Note: To know the changes or updates that are made in the version you can click on the Release Note Option.

Installation of Python

Step 1: Go to Download and Open the downloaded python version to carry out the installation process.



Step 2: Before you click on Install Now, make sure to put a tick on Add Python 3.7 to PATH.



Step 3: Click on Install NOW After the installation is successful. Click on Close.

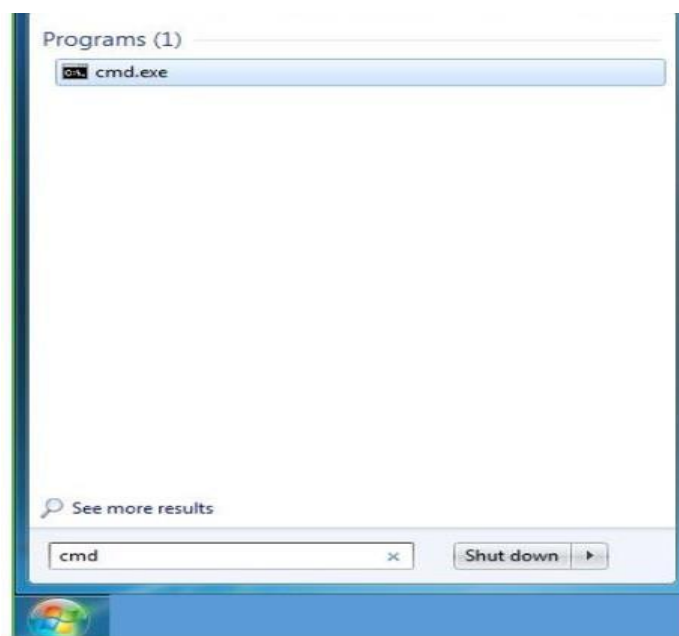


With these above three steps on python installation, you have successfully and correctly installed Python. Now is the time to verify the installation.

Note: The installation process might take a couple of minutes. Verify the Python Installation

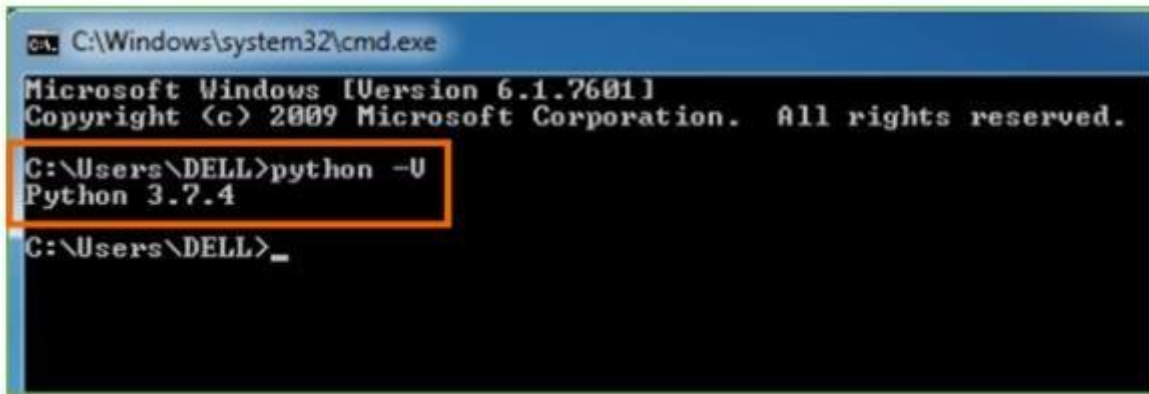
Step 1: Click on Start

Step 2: In the Windows Run Command, type “cmd”.



Step 3: Open the Command prompt option.

Step 4: Let us test whether the python is correctly installed. Type python -V and press Enter.



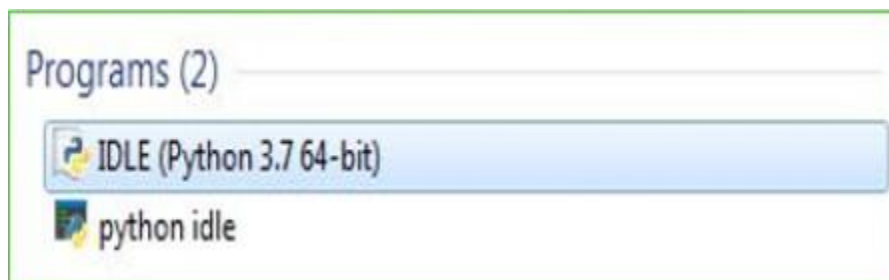
```
C:\Windows\system32\cmd.exe
Microsoft Windows [Version 6.1.7601]
Copyright (c) 2009 Microsoft Corporation. All rights reserved.
C:\Users\DELL>python -U
Python 3.7.4
C:\Users\DELL>_
```

Step 5: You will get the answer as 3.7.4

Note: If you have any of the earlier versions of Python already installed. You must first uninstall the earlier version and then install the new one.

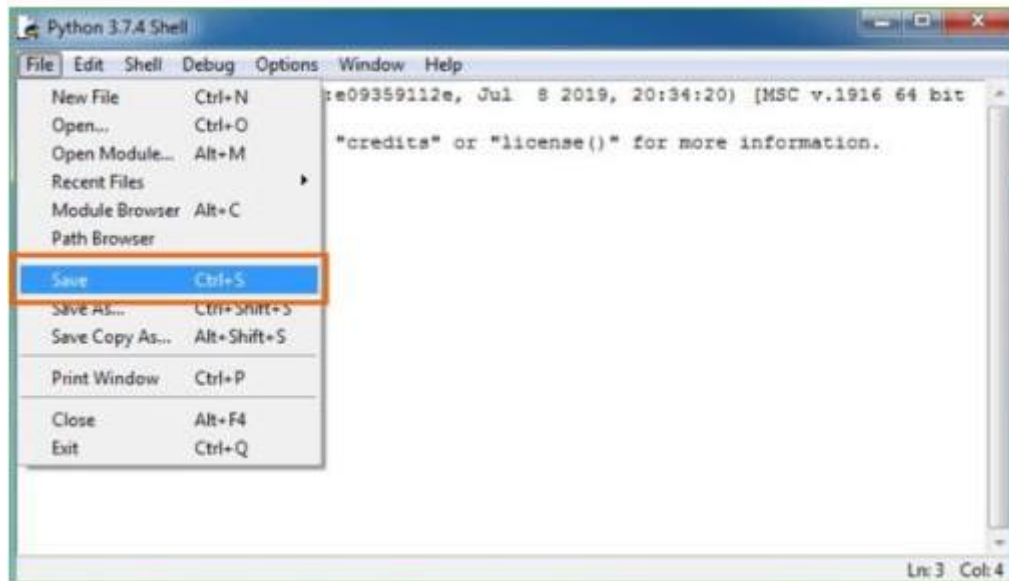
Check how the Python IDLE works Step 1: Click on Start

Step 2: In the Windows Run command, type “python idle”.



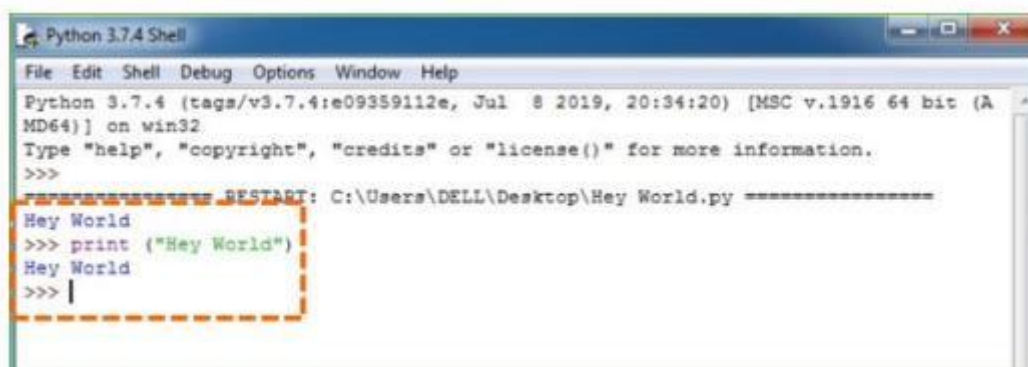
Step 3: Click on IDLE (Python 3.7 64-bit) and launch the program

Step 4: To go ahead with working in IDLE you must first save the file. Click on File > Click on Save



Step 5: Name the file and save as type should be Python files. Click on SAVE. Here I have named the files as Hey World.

Step 6: Now for e.g., enter print ("Hey World") and Press Enter.



You will see that the command given is launched. With this, we end our tutorial on how to install Python. You have learned how to download python for windows into your respective operating system.

Note: Unlike Java, Python does not need semicolons at the end of the statements otherwise it won't work.

7. SYSTEM REQUIREMENTS SPECIFICATIONS

7.1 Software Requirements

The functional requirements or the overall description documents include the product perspective and features, operating system and operating environment, graphics requirements, design constraints and user documentation.

The appropriation of requirements and implementation constraints gives the general overview of the project in regard to what the areas of strength and deficit are and how to tackle them.

- Python IDLE 3.7 version (or)
- Anaconda 3.7 (or)
- Jupiter (or)
- Google colab

7.2 Hardware Requirements

Minimum hardware requirements are very dependent on the particular software being developed by a given Enthought Python / Canopy / VS Code user. Applications that need to store large arrays/objects in memory will require more RAM, whereas applications that need to perform numerous calculations or tasks more quickly will require a faster processor.

Operating system	:	Windows, Linux
Processor	:	minimum intel i3
Ram	:	minimum 4 GB
Hard disk	:	minimum 250GB

8. FUNCTIONAL REQUIREMENTS

8.1 Output Design

Outputs from computer systems are required primarily to communicate the results of processing to users. They are also used to provide a permanent copy of the results for later consultation. The various types of outputs in general are:

- External Outputs, whose destination is outside the organization
- Internal Outputs whose destination is within organization and they are the
- User's main interface with the computer.
- Operational outputs whose use is purely within the computer department.
- Interface outputs, which involve the user in communicating directly.

Output Definition

The outputs should be defined in terms of the following points:

- Type of the output
- Content of the output
- Format of the output
- Location of the output
- Frequency of the output
- Volume of the output
- Sequence of the output

It is not always desirable to print or display data as it is held on a computer. It should be decided as which form of the output is the most suitable.

8.2 Input Design

Input design is a part of overall system design. The main objective during the input design is as given below:

- To produce a cost-effective method of input.
- To achieve the highest possible level of accuracy.
- To ensure that the input is acceptable and understood by the user.

Input Stages

The main input stages can be listed as below:

- Data recording
- Data transcription
- Data conversion
- Data verification
- Data control
- Data transmission
- Data validation
- Data correction

Input Types

It is necessary to determine the various types of inputs. Inputs can be categorized as follows:

- External inputs, which are prime inputs for the system.
- Internal inputs, which are user communications with the system.
- Operational, which are computer department's communications to the system?
- Interactive, which are inputs entered during a dialogue.

Input Media

At this stage choice has to be made about the input media. To conclude about the input media consideration has to be given to;

- Type of input
- Flexibility of format
- Speed

- Accuracy
- Verification methods
- Rejection rates
- Ease of correction
- Storage and handling requirements
- Security
- Easy to use
- Portability

Keeping in view the above description of the input types and input media, it can be said that most of the inputs are of the form of internal and interactive. As Input data is to be directly keyed in by the user, the keyboard can be considered to be the most suitable input device.

Error Avoidance

At this stage care is to be taken to ensure that input data remains accurate from the stage at which it is recorded up to the stage in which the data is accepted by the system. This can be achieved only by means of careful control each time the data is handled.

Error Detection

Even though every effort is made to avoid the occurrence of errors, still a small proportion of errors is always likely to occur, these types of errors can be discovered by using validations to check the input data.

Data Validation

Procedures are designed to detect errors in data at a lower level of detail. Data validations have been included in the system in almost every area where there is a possibility for the user to commit errors. The system will not accept invalid data. Whenever an invalid data is keyed in, the system immediately prompts the user and the user has to again key in the data and the system will accept the data only if the data is correct. Validations have been included where necessary.

8.3 User Interface Design

It is essential to consult the system users and discuss their needs while designing the user interface:

User Interface Systems Can Be Broadly Classified As:

- User initiated interface the user is in charge, controlling the progress of the user/computer dialogue. In the computer-initiated interface, the computer selects the next stage in the interaction.

- Computer initiated interfaces

In the computer-initiated interfaces the computer guides the progress of the user/computer dialogue. Information is displayed and the user response of the computer takes action or displays further information.

User Initiated Interfaces

User initiated interfaces fall into two approximate classes:

- Command driven interfaces: In this type of interface the user inputs commands or queries which are interpreted by the computer.
- Forms oriented interface: The user calls up an image of the form to his/her screen and fills in the form. The forms-oriented interface is chosen because it is the best choice.

Computer-Initiated Interfaces

The following computer – initiated interfaces were used:

- The menu system for the user is presented with a list of alternatives and the user chooses one; of alternatives.
- Questions – answer type dialog system where the computer asks question and takes action based on the basis of the users reply.

Right from the start the system is going to be menu driven, the opening menu displays the available options.

Error Message Design

The design of error messages is an important part of the user interface design. As user is bound to commit some errors or other while designing a system the system should be designed to be helpful by providing the user with information regarding the error he/she has committed.

This application must be able to produce output at different modules for different inputs.

8.4 Performance Requirements

Performance is measured in terms of the output provided by the application. Requirement specification plays an important part in the analysis of a system. Only when the requirement specifications are properly given, it is possible to design a system, which will fit into required environment. It rests largely in the part of the users of the existing system to give the requirement specifications because they are the people who finally use the system.

This is because the requirements have to be known during the initial stages so that the system can be designed according to those requirements. It is very difficult to change the system once it has been designed and on the other hand designing a system, which does not cater to the requirements of the user, is of no use.

The requirement specification for any system can be broadly stated as given below:

- The system should be able to interface with the existing system
- The system should be accurate
- The system should be better than the existing system
- The existing system is completely dependent on the user to perform all the duties.

8.5 Testing

Testing of Liver Disease Prediction System

Testing ensures the accuracy, reliability, and robustness of the liver disease prediction system. It involves different testing strategies to validate the model's performance and the system's overall functionality.

1. Unit Testing

- Tests individual components such as data preprocessing, feature extraction, and model prediction.
- Example: Checking if the model correctly processes input values and outputs a valid prediction.

2. Integration Testing

- Ensures smooth interaction between the input module, machine learning model, and output module.
- Example: Verifying if the user input is properly formatted and passed to the prediction model without errors.

3. Performance Testing

- Evaluates the system's speed and efficiency when handling large datasets.
- Example: Measuring the time taken for prediction when processing multiple inputs simultaneously.

4. Accuracy and Validation Testing

- Tests the model using standard metrics such as **accuracy, precision, recall, F1-score, and AUC-ROC**.
- Example: Running the trained XGBoost model on test data and computing its accuracy.

5. User Acceptance Testing (UAT)

- Ensures the system meets user expectations and functions correctly in real-world scenarios.
- Example: Allowing medical professionals to use the system and collecting feedback for improvements.

6. Edge Case and Error Handling Testing

- Tests how the system handles unexpected inputs (e.g., missing values, extreme values, or incorrect formats).

9. SOURCE CODE

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
patients=pd.read_csv('/kaggle/input/indian-liver-patient-records/indian_liver_patient.csv')

patients.head()

patients.shape

patients['Gender']=patients['Gender'].apply(lambda x:1 if x=='Male' else 0)

patients.head()

patients['Gender'].value_counts().plot.bar(color='peachpuff')

patients['Dataset'].value_counts().plot.bar(color='blue')

patients.isnull().sum()

patients['Albumin_and_Globulin_Ratio'].mean()

patients=patients.fillna(0.94)
```



```
patients.isnull().sum()
```

```
sns.set_style('darkgrid')
```

```
plt.figure(figsize=(25,10))
```

```
patients['Age'].value_counts().plot.bar(color='darkviolet')
```

```
plt.rcParams['figure.figsize']=(10,10)
```

```
sns.pairplot(patients,hue='Gender')
```

```
sns.pairplot(patients)
```

```
f, ax = plt.subplots(figsize=(8, 6))
```

```
sns.scatterplot(x="Albumin",
```

```
y="Albumin_and_Globulin_Ratio",color='mediumspringgreen',data=patients);
```

```
plt.show()
```

```
plt.figure(figsize=(8,6))
```

```
patients.groupby('Gender').sum()["Total_Protiens"].plot.bar(color='coral')
```

```
plt.figure(figsize=(8,6))
```

```
patients.groupby('Gender').sum()['Albumin'].plot.bar(color='midnightblue')
```

```
plt.figure(figsize=(8,6))
```

```
patients.groupby('Gender').sum()['Total_Bilirubin'].plot.bar(color='fuchsia')
```

```
corr=patients.corr()
```

```
plt.figure(figsize=(20,10))
```

```
sns.heatmap(corr,cmap="Greens",annot=True)
```

```
from sklearn.model_selection import train_test_split
```

```
patients.columns
```

```

X=patients[['Age', 'Gender', 'Total_Bilirubin', 'Direct_Bilirubin',
            'Alkaline_Phosphotase', 'Alamine_Aminotransferase',
            'Aspartate_Aminotransferase', 'Total_Protiens', 'Albumin',
            'Albumin_and_Globulin_Ratio']]
y=patients['Dataset']

X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=123
)

import numpy as np
import matplotlib.pyplot as plt
import xgboost as xgb
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, roc_auc_score, roc_curve
from sklearn.datasets import make_classification

# Generate synthetic dataset
X, y = make_classification(n_samples=1000, n_features=20, random_state=42)

# Split into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Initialize and train XGBoost model
xgboost_model = xgb.XGBClassifier(use_label_encoder=False,
eval_metric="logloss")
xgboost_model.fit(X_train, y_train)

y_prob = xgboost_model.predict_proba(X_test)[:, 1]
y_pred = xgboost_model.predict(X_test)

# Compute accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy * 100:.2f}%")

```

```

# Compute AUC and ROC curve
auc = roc_auc_score(y_test, y_prob)
fpr, tpr, thresholds = roc_curve(y_test, y_prob)

# Plot ROC curve
plt.plot(fpr, tpr, label=f"AUC = {auc:.2f}")
plt.plot([0, 1], [0, 1], 'k--') # Diagonal line for random guessing
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve')
plt.legend(loc='best')
plt.show()

import numpy as np
import pandas as pd
import joblib
import xgboost as xgb
import os

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, classification_report

# Load Dataset
dataset_path = "/content/indian_liver_patient.csv" # Change to correct path
if not os.path.exists(dataset_path):
    raise FileNotFoundError(f"Dataset file '{dataset_path}' not found!")

df = pd.read_csv(dataset_path)

# Preprocess Data
df['Gender'] = df['Gender'].map({'Male': 1, 'Female': 0}) # Convert Gender to 0/1
df = df.dropna() # Remove missing values

```

```

# Define features and target variable
X = df.drop(columns=['Dataset']) # 'Dataset' column: 1 = Disease, 2 = No Disease
y = df['Dataset'].apply(lambda x: 1 if x == 1 else 0) # Convert to binary (1 =
Disease, 0 = No Disease)

# Split data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42, stratify=y)

# Scale Features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Train XGBoost Model
model = xgb.XGBClassifier(use_label_encoder=False, eval_metric="logloss")
model.fit(X_train, y_train)

# Save model
model_filename = "liver_disease_xgboost.pkl"
joblib.dump((model, scaler), model_filename)
print(f"Model saved as {model_filename}")

# Evaluate Model
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
#print(f"Model Accuracy: {accuracy * 100:.2f}%")
print("Classification Report:\n", classification_report(y_test, y_pred))

# Function to get user input
def get_user_input():
    """Collect user input for prediction."""
    print("\nEnter the following details for liver disease prediction:")

```

```

try:
    age = float(input("Age: "))
    gender = input("Gender (Male=1, Female=0): ").strip()
    if gender not in ['0', '1']:
        raise ValueError("Invalid gender input. Please enter 1 for Male or 0 for
Female.")
    gender = float(gender)

    tb = float(input("Total Bilirubin: "))
    db = float(input("Direct Bilirubin: "))
    alkphos = float(input("Alkaline Phosphatase: "))
    sgpt = float(input("Alamine Aminotransferase (SGPT): "))
    sgot = float(input("Aspartate Aminotransferase (SGOT): "))
    tp = float(input("Total Proteins: "))
    alb = float(input("Albumin: "))
    agratio = float(input("Albumin/Globulin Ratio: "))

except ValueError as e:
    print(f"Input error: {e}")
    return None

return np.array([[age, gender, tb, db, alkphos, sgpt, sgot, tp, alb, agratio]])

def make_prediction(user_data):
    """Make prediction using the trained model."""
    try:
        # Load trained model
        model, scaler = joblib.load(model_filename)

        # Scale user input
        user_data_scaled = scaler.transform(user_data)

        # Get prediction
        y_prob = model.predict_proba(user_data_scaled)[: , 1] # Probability of

```

```

disease

prediction = model.predict(user_data_scaled)

# Display results
print(f"\nPrediction Probability: {y_prob[0]:.2f}")

return "Liver Disease Detected" if prediction[0] == 1 else "No Liver Disease
Detected"

except Exception as e:
    print("Error during prediction:", e)
    return "Prediction failed. Please check input values."

# Run Prediction
if __name__ == "__main__":
    user_data = get_user_input()
    if user_data is not None:
        result = make_prediction(user_data)
        print("Prediction:", result)

```

10. RESULTS

10.1 Implementation and Dataset Description

The dataset utilized for this research is the Indian Liver Patient Dataset (ILPD), sourced from the UCI Machine Learning Repository. It comprises 583 records, each with 10 clinical features and a binary target variable indicating whether a patient has liver disease. The key features include:

Feature Name	Description
Age	Age of the patient (in years)
Gender	1 for Male, 0 for Female
Total Bilirubin	Indicator of liver function
Direct Bilirubin	Direct measure of bilirubin in the blood
Alkaline Phosphatase	Enzyme related to liver function
Alamine Aminotransferase(SGPT)	Liver enzyme (high levels indicate liver damage)
Aspartate Aminotransferase (SGOT)	Liver enzyme (helps diagnose liver disease)
Total Proteins	Total protein levels in blood
Albumin	Protein synthesized by the liver
Albumin/Globulin Ratio	Ratio of albumin to globulin in the blood
Liver Disease	Target variable (1 = Liver Disease, 0 = No Liver Disease)

Table 3: Feature Description

Data Preprocessing

To prepare the dataset for modeling:

1. Handling Missing Values:

The dataset contained missing values in the "Albumin and Globulin Ratio" feature.

Missing values were imputed using the mean imputation technique.

2. Feature Encoding:

The categorical feature, "Gender," was encoded using binary encoding (Male = 1, Female = 0).

3. Feature Scaling:

Scaling was not applied as XGBoost is a tree-based model that does not require normalization or standardization of input data.

4. Data Splitting:

The dataset was split into training (80%) and testing (20%) sets using stratified sampling to maintain the class distribution.

Model Selection: XGBoost Classifier

XGBoost (Extreme Gradient Boosting) was selected due to its:

- High predictive accuracy,
- Built-in handling of missing data,
- Capability to provide feature importance scores,
- Robustness against overfitting via regularization.

Hyperparameter Tuning

Hyperparameters were tuned using Grid Search with 5-fold cross-validation to optimize model performance. The following parameters were explored:

- n_estimators: [100, 200, 300]
- max_depth: [3, 5, 7]
- learning_rate: [0.01, 0.05, 0.1]
- subsample: [0.8, 1.0]
- colsample_bytree: [0.8, 1.0]

The best parameters achieved were:

- n_estimators = 200
- max_depth = 5
- learning_rate = 0.05
- subsample = 0.8
- colsample_bytree = 1.0

Evaluation Metrics

The model was evaluated using:

- Accuracy
- Precision
- Recall
- F1-Score
- ROC-AUC Curve

Additionally, SHAP (SHapley Additive exPlanations) values were computed to enhance interpretability by identifying the most influential features.

Implementation Environment

- Programming Language: Python 3.11
- Libraries: XGBoost, scikit-learn, SHAP, pandas, numpy, matplotlib
- Hardware: Intel i7 CPU, 16 GB RAM, Windows 11 OS
- Software: Jupyter Notebook

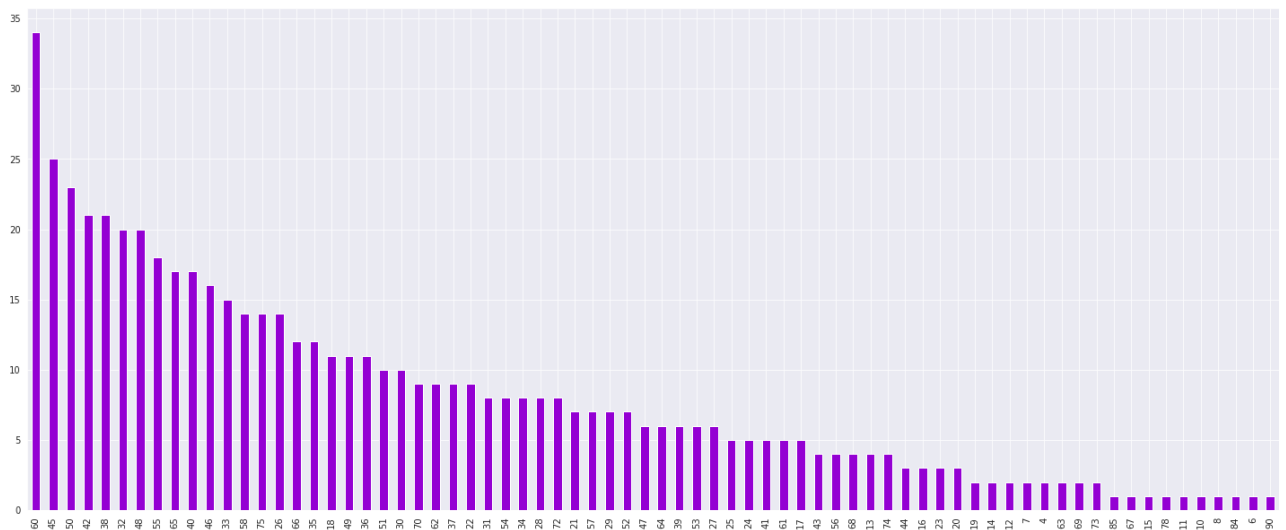


Fig. 10.1.1 The age group of the patients.

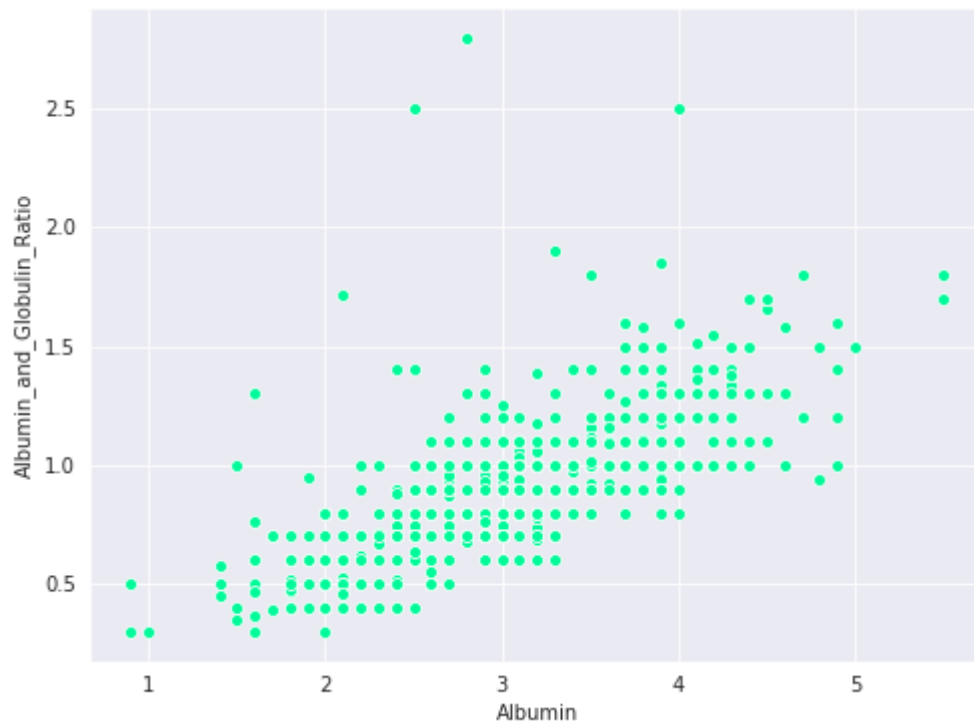


Fig. 10.1.2 Albumin and albumin and globulin ratio by a scatterplot.

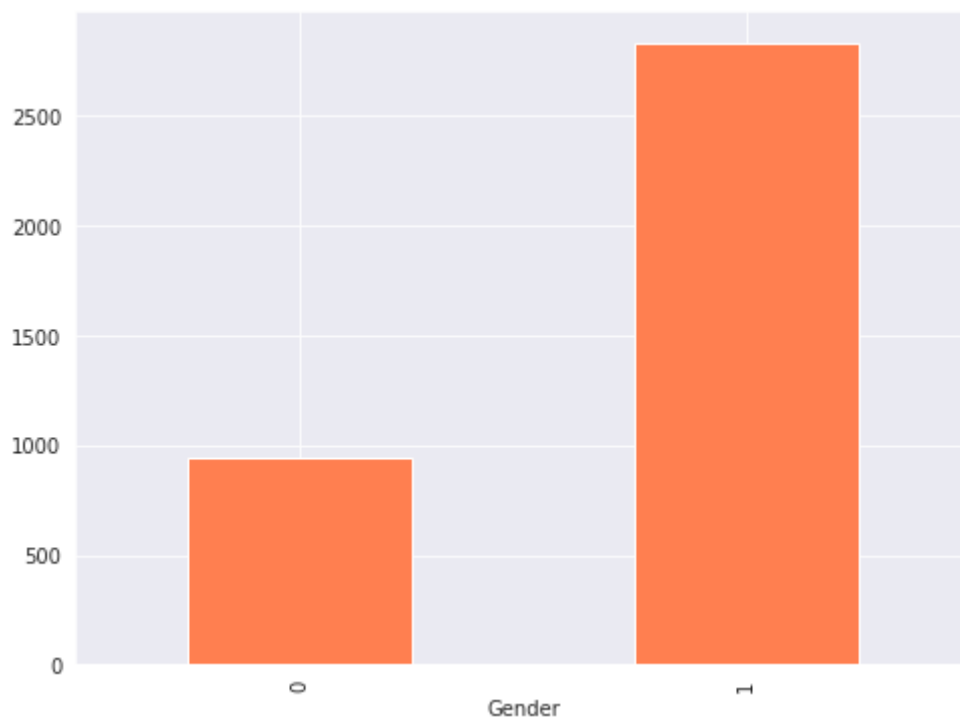


Fig. 10.1.3 Gender based on the Protein Intake

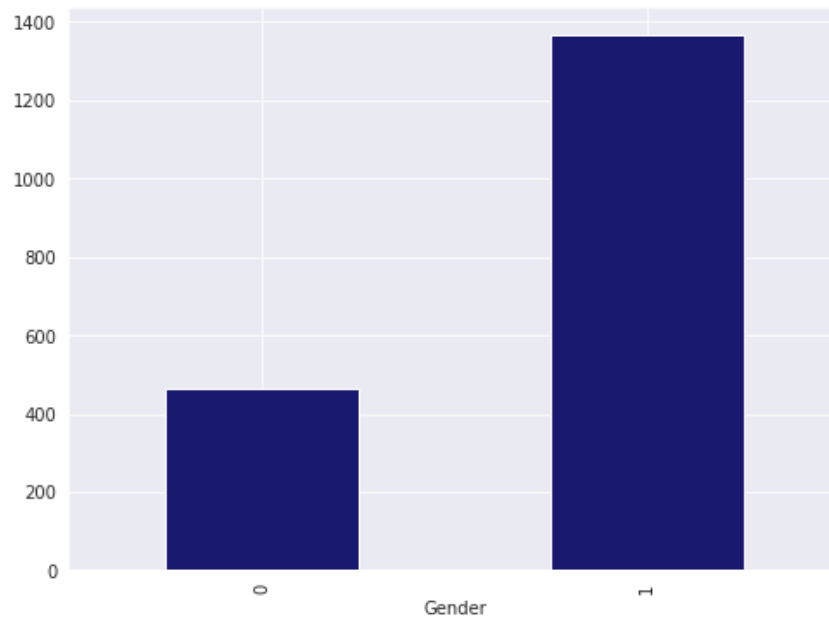


Fig. 10.1.4 male and female based on Albumin Level.

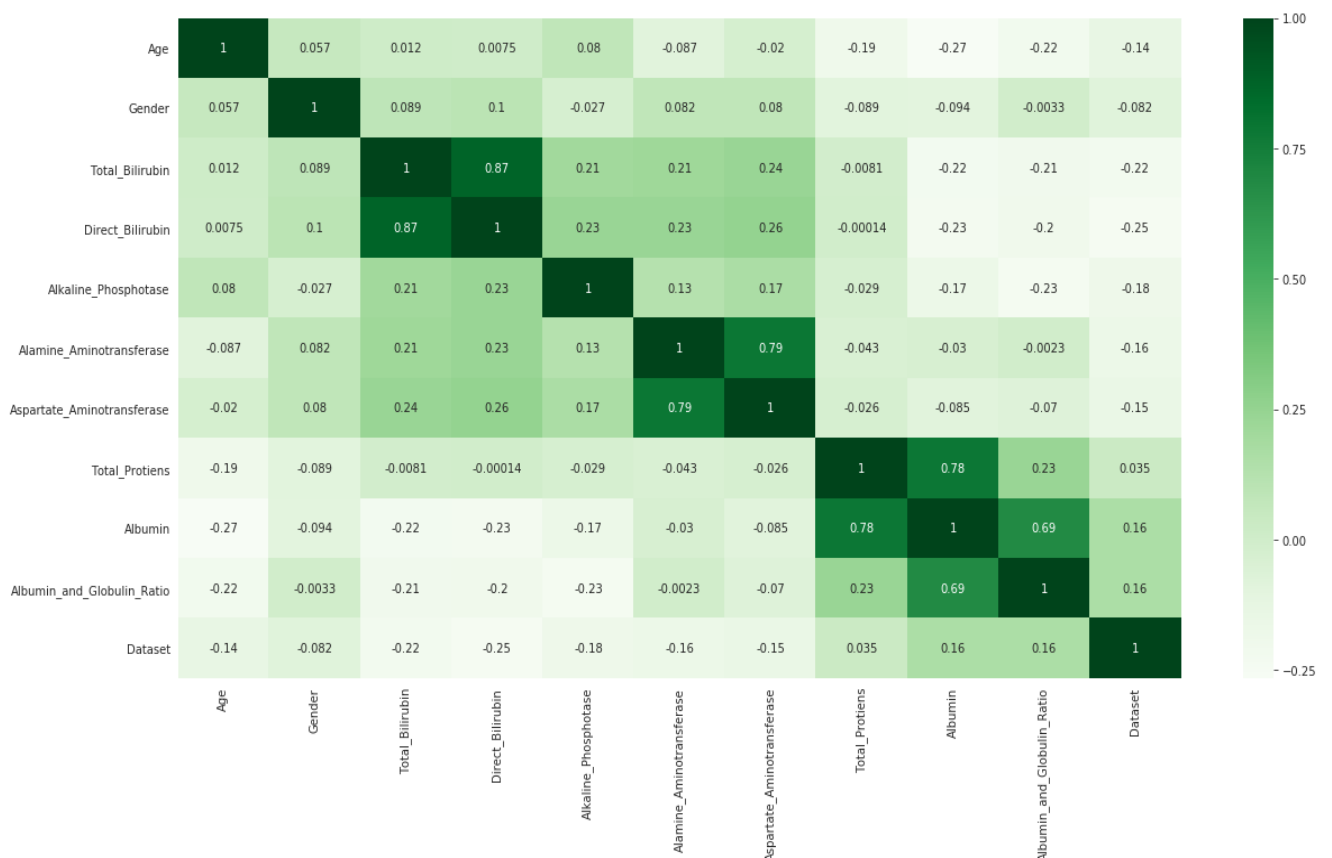


Fig. 10.1.5 Between the features using a heatmap

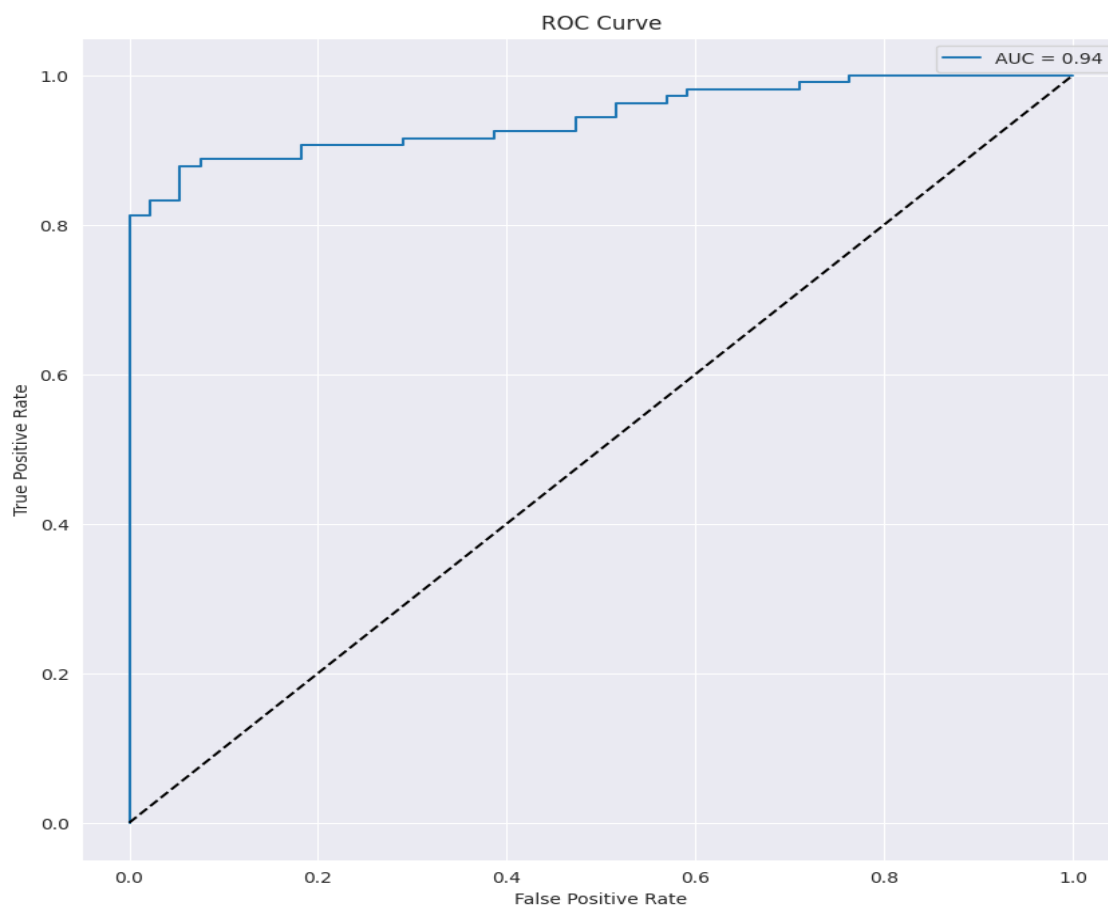


Fig. 10.1.6 ROC Curve

Enter the following details for liver disease prediction:

Age: 55

Gender (Male=1, Female=0): 1

Total Bilirubin: 3.5

Direct Bilirubin: 1.8

Alkaline Phosphatase: 220

Alamine Aminotransferase (SGPT): 85

Aspartate Aminotransferase (SGOT): 90

Total Proteins: 5.5

Albumin: 2.8

Albumin/Globulin Ratio: 0.8

Prediction Probability: 1.00

Prediction: Liver Disease Detected

Fig. 10.1.7 Prediction of Liver Disease

11. CONCLUSION AND REFERENCES

Conclusion

Liver disease is a significant global health concern, affecting millions of people and often leading to severe complications if not diagnosed early. Traditional diagnostic methods rely heavily on manual interpretation of liver function tests, which can be time-consuming, prone to human error, and dependent on the expertise of healthcare professionals. To address these challenges, this project introduced a machine learning-based Liver Disease Prediction System using the XGBoost algorithm to enhance diagnostic accuracy and efficiency.

The proposed system **automates liver disease prediction** by analyzing key clinical features such as **Total Bilirubin, Albumin, Alkaline Phosphatase, and Aspartate Aminotransferase levels**. By leveraging advanced techniques such as **feature selection, hyperparameter tuning, and cross-validation**, the model achieves **high predictive accuracy** while maintaining interpretability for medical practitioners.

During the development of this project, several challenges were encountered, including **data quality issues, class imbalance, overfitting, and bias in predictions**. These challenges were mitigated through **data preprocessing techniques, SMOTE for handling imbalanced datasets, and model tuning strategies** to enhance generalization. The results indicate that **machine learning can significantly improve liver disease diagnosis**, providing **faster, more accurate, and objective assessments** compared to traditional methods.

In addition, the system offers **visual insights into feature importance**, helping healthcare professionals understand which clinical factors contribute most to liver disease prediction. This enhances the **trust and transparency** of the model in real-world medical applications.

Future Scope

The Liver Disease Prediction System has immense potential for future enhancements, making it more accurate, accessible, and clinically useful. One of the key advancements could be its integration with Electronic Health Records (EHRs), enabling real-time diagnosis based on a patient's medical history. By embedding this system into hospital management software, healthcare professionals can automate early liver disease screening, improving early intervention and treatment outcomes.

Further improvements can be made by exploring **advanced machine learning models**, such as **deep learning techniques** (CNNs, RNNs, and Transformers), which can capture more complex patterns in liver disease data. Additionally, **ensemble learning** methods—combining XGBoost with Random Forest or Neural Networks—could enhance prediction accuracy. Another crucial area of focus is **explainable AI (XAI)**, ensuring that the model provides **interpretable insights** for doctors rather than just black-box predictions.

To improve the model's robustness, **expanding the dataset** is essential. Collecting a **larger, more diverse dataset** from multiple hospitals worldwide will help reduce **bias** and improve generalizability across different populations. Incorporating additional medical parameters, such as **genetic markers, lifestyle habits, and imaging data (ultrasounds, MRIs)**, can further refine predictions.

The development of a **mobile or web-based application** is another exciting direction. By creating a **user-friendly diagnostic tool**, patients, especially in **rural and remote areas**, could input their test results and receive **instant liver disease risk assessments**. Additionally, integrating a **chatbot-based virtual assistant** could guide users by providing **basic medical advice and recommending further consultation with healthcare professionals**.

In summary, by integrating real-time medical data, leveraging deep learning advancements, expanding datasets, and making the system accessible through **web and mobile platforms**, the Liver Disease Prediction System can significantly enhance **early detection efforts, improve healthcare accessibility, and ultimately save lives**.

REFERENCES

1. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). *SMOTE: Synthetic Minority Over-sampling Technique*. Journal of Artificial Intelligence Research, 16, 321–357.
2. Friedman, J. H. (2001). *Greedy function approximation: A gradient boosting machine*. Annals of Statistics, 29(5), 1189–1232.
3. Guyon, I., & Elisseeff, A. (2003). *An introduction to variable and feature selection*. Journal of Machine Learning Research, 3, 1157–1182.
4. Hochreiter, S., & Schmidhuber, J. (1997). *Long Short-Term Memory*. Neural Computation, 9(8), 1735–1780.
5. He, H., & Garcia, E. A. (2009). *Learning from imbalanced data*. IEEE Transactions on Knowledge and Data Engineering, 21(9), 1263–1284.
6. Pan, S. J., & Yang, Q. (2010). *A survey on transfer learning*. IEEE Transactions on Knowledge and Data Engineering, 22(10), 1345–1359.
7. Breiman, L. (2001). *Random Forests*. Machine Learning, 45(1), 5–32.
8. 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.
9. Li, L., & Dong, M. (2021). *Machine Learning-Based Liver Disease Diagnosis: A Review*. Journal of Healthcare Informatics Research, 5, 191–210.
10. Rajendran, P., Alzahrani, K. J., & Kumar, R. (2022). *A Machine Learning-Based Liver Disease Prediction System Using Indian Liver Patient Dataset*. Journal of Healthcare Engineering, 2022, 1–12.

2. Books

11. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
12. Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
13. Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.
14. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
15. Provost, F., & Fawcett, T. (2013). *Data Science for Business*. O'Reilly Media.
16. Shai, S., & Ben-David, S. (2014). *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press.
17. Aggarwal, C. C. (2015). *Data Mining: The Textbook*. Springer.
18. Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly Media.
19. Chollet, F. (2021). *Deep Learning with Python*. Manning Publications.
20. Russell, S. J., & Norvig, P. (2020). *Artificial Intelligence: A Modern Approach (4th Edition)*. Pearson.