University of Hertfordshire **UH**

School of Physics,
Engineering and
Computer Science

# MSc Data Science Project

# 7PAM2002-0509-2023

Department of Physics, Astronomy and Mathematics

# **Data Science FINAL PROJECT REPORT**

## **Project Title:**

## Predicting Lung Cancer Using Life Style and Health Factors

## **Student Name and SRN:**

Upali Jayawardhange Dilhara Rajaguru / 22034062

Supervisor: Klass Wiersema

Date Submitted:  29/08/2024

Word Count:  Enter the word count

# DECLARATION STATEMENT

This report is submitted in partial fulfillment of the requirement for the Master of Science in Data Science degree at the University of Hertfordshire.

I have read the guidance to students on academic integrity, misconduct, and plagiarism information at Assessment Offences and Academic Misconduct and understand the University process of dealing with suspected cases of academic misconduct and the possible penalties, which could include failing the project module or course.

I certify that the work submitted is my own and that any material derived or quoted frompublished or unpublished work of other persons has been duly acknowledged. (Ref. UPR AS/C/6.1, section 7, and UPR AS/C/5, section 3.6). I have not used chatGPT, or any other generative AI tool, to write the report or code (other than where declared or referenced).

I did not use human participants or surveys for my MSc Project.

I hereby permit the report to be available on module websites provided the source is acknowledged.

Student Name printed:  Upali Jayawardhanage Dilhara Rajaguru

Student Name signature:   Dilhara

Student SRN number:   22034062

UNIVERSITY OF HERTFORDSHIRE

SCHOOL OF PHYSICS, ENGINEERING AND COMPUTER SCIENCE

# Abstract

My research predicts lung cancer through various machine-learning methods. The main goals are to find the most accurate model for predicting lung cancer and the most influential factors in diagnosis. The used dataset contains 309 entries and 16 features, such as age, smoking, and symptoms related to lung cancer. I applied multiple models: Logistic Regression, Support Vector Machine, Decision Tree, and Naive Bayes. Each model was trained and the accuracy assessed. Among these, the SVM and Logistic Regression performed the best with an accuracy of 98% and 99%, respectively. The hyperparameter tuning for the models was done. The results indicated that SVM and logistic regression are highly rated in predicting lung cancer, whereby coughing and difficulty swallowing are the most important predictors. The findings further suggest that these models would be beneficial in assisting early diagnosis of lung cancer and improving treatment results for patients.

# Table of Content

# Chapter 1: Introduction

Cancer is the general term for diseases characterized by uncontrolled and abnormal growth of body cells. Unless treated, the possible result of this is death. Cancer can originate almost anywhere in the body and creates several symptoms, depending on the type of cancer and part of the body. Lung cancer is a cancer type that originates in the lungs. These are the organs responsible for oxygen intake and carbon dioxide removal. Lung cancer is probably one of the most common and deadly cancers in the world. This cancer can be of two types: non-small cell lung cancer and small cell lung cancer. Symptoms of lung cancer usually do not appear until the disease is already at an advanced stage. Possible signs and symptoms of lung cancer a cough, coughing up blood or rust-colored sputum, chest pain, shortness of breath, hoarseness, unexplained weight loss, loss of appetite, feeling tired or weak, and repeated lung problems, such as bronchitis or pneumonia. Because of cancer, there are around 10 million deaths around the world. From that 1.8 million deaths happen because of lung cancer each year. This is about 18% of all cancer deaths.  So it is important to develop an AI-related predictive model for predicting lung cancer using Machine Learning models. Then we can reduce the worldwide death rate as much as we can.

Due to the high mortality rate caused by lung cancer and its being generally difficult to diagnose at an early stage, there is a dire need for the development of a model that would predict with accuracy those susceptible to the disease before the latter advances. Early diagnosis enhances the possibilities of treating lung cancer with success and survival. A prediction model, using data-driven approaches and machine learning algorithms, may scan through various risk factors and patient data to estimate the possibility of developing lung cancer. Clinically and in public health, this can be a very powerful tool for healthcare providers in implementing timely intervention and reducing mortality due to lung cancer.

The key objective of this project will be the design of a machine learning-based prediction model for lung cancer. The model should be able to give an accurate prediction of the likelihood of lung cancer among individuals, which is to be drawn from a range of demographic, lifestyle, and health-related attributes.

My project aims to predict lung cancer using machine learning methods including logistic regression, Support vector machine, Decision tree, and Naive Bayes model. My task is to detect the most appropriate model for predicting lung cancers and find the most influential factors that take part in the diagnosis.

The idea of my research is clear, based on some health, lifestyle, and environmental factors collected in my dataset which includes 309 entries and 16 features. The target is to apply some machine learning models such as logistic regression, Support Vector Machine(SVM), decision tree, and naive Bayes to predict lung cancer and tune the hyperparameters to find the most precise analysis to conclude if the target will be able to in a safe way to predict the lung cancer.

I chose these models to develop my project because they are all accurate but for this work are important the tuning hyperparameters for the best area under ROC Curve. The main objective of these analyses is to find the best model for prediction and use the most significant variables for diagnosis.

The main finding indicates that these four models of choice will be able to predict and diagnose Lung disease, where if these Models fit the condition will help to prevent the disease increment, and can impact significantly, Lung cancer's detectability as one of the Modern diseases globally.

**Research Question:** Predicting Lung Cancer Using Life Style and Health Factors

**Aims:**
- To explore machine learning techniques and apply the most suitable method in developing a predictive model for lung cancer.
- Compare the performance of the various machine learning models in lung cancer prediction considering the parameters such as accuracy, precision, recall, and F1-score.
- and analysis of the primary risk factors that contribute the most towards the prediction of lung cancer. This would provide valuable insights that may further be used in seeking to inform clinical practice.

**Objectives:**
- Preprocess the lung cancer dataset and identify important variables and trends in the data.
- Train and evaluate different machine learning models such as Logistic regression, Decision Trees, Naive Bayes, and Support Vector Machines.
- The performance of the prediction model can be evaluated based on different performance measures like accuracy, precision, recall, and F1-score.
- The benchmarking of the developed model regarding current industry practice and discussing its potential impact and applications relevant in health care.

    This introduction provides the background that serves to explain the importance of the research, it articulates the research question, aims, and objectives. This structure enables readers to grasp what background and context your project operates within before they delve into technical information.

# Chapter 2: Literature Review

## 2.1. Overview of Lung Cancer Prediction Using Machine Learning

Lung cancer is the most common cause of death due to cancer throughout the world. The high mortality rate from this disease is partly because most cases are diagnosed at advanced stages. The process of early detection is very crucial in improving patients' survival rates, as the rate increases dramatically when the disease is at a curable stage. Traditional screening methods, such as LDCT, are effective but have several major drawbacks: the possibility of false positives, radiation exposure, and high costs. This has been further complicated by Oudkerk et al. (2021), who developed a high interest in the application of machine learning techniques since they will improve the early detection of lung cancer by more appropriately and proficiently screening people with higher risks than conventional methods.

## 2.2. Technical Background: Lung Cancer Prediction using Machine Learning

Machine learning, a subbranch of artificial intelligence, refers to a set of algorithms that allow computers to make predictions from data. The ML system analyzes the high volume of data in lung cancer prediction to achieve patterns and relationships often not so apparent by using traditional statistical methods. These models can handle complex, high-dimensional data, thus being highly suitable in medical application fields where several variables usually need consideration all at once.

Logistic regression analysis, decision trees, support vector machines, random forests, and more recent techniques such as extreme gradient boosting and deep learning are diverse machine learning methods that have been implemented in practice for the prediction of lung cancer. Each model will be selected based on the nature of the data and the requirements of the task at hand, such as interpretability versus predictive accuracy.

## 2.3. Critical Analysis of the Above Discussed Studies

1. **Zhang et al. (2024):** Lung Cancer Risk Prediction Using Machine Learning The same large-scale study using the UK Biobank with 467,888 participants was done by Zhang et al. Predictive performance for logistic regression, naive Bayes, random forest models, and deep learning models through XGBoost was compared in the study, yielding an AUC of 0.998 for XGBoost. Novel predictors including hip, waist circumference, neuroticism score, and forced expiratory volume played an important role in the prediction of the risk of lung cancer. Novel predictors were included in explaining the large sample size and contributing to the sensitivity of the model to this study. However, since part of the predictors are obtained by self-report, biases may be introduced, and the generalisability of the study is limited only to the particular population from which the data was drawn.

2. **Fatima et al. (2022):** Machine learning techniques for the detection of lung cancer In this study, Fatima et al. (2022) worked on the application of common machine learning techniques to apply them to the sphere of lung cancer detection. First of all, much attention was paid to the selection of features. According to the writers, irrelevant features will drastically cut down the accuracy of the model. Authors underlined that ML models, including SVM, might show very good results in accuracy, but being "black boxes," they are practically difficult to interpret. In clinical applications, the rational interpretation of a prediction becomes necessary for a practitioner. Whereas the strength of this study lies in an explicit focus on the trade-offs between model interpretation and model complexity, underpinning it is an argument on the need for more transparent models that are easily interpretable and trustworthy by clinicians as well.

3. **Gould et al. (2021):** Machine Learning in the Identification of Early Lung Cancer Gould et al. applied machine learning techniques to the data obtained from routine clinics and laboratories to enhance the detection rate of lung cancer at an earlier stage. A model was developed that outperformed all the existing tools for lung cancer screening, including the USPSTF eligibility guidelines related to other metabolic markers such as C-reactive protein. They demonstrated that ML models could successfully integrate non-traditional predictors for early detection. However, the main focus of this study was to outline the difficulties in translating such models into clinical practice and called for extensive validation across diverse populations to ensure generalizability.

4. **Maisonneuve et al. (2011):** Model of Risk Prediction for Lung Cancer Based on findings from the Italian COSMOS trial composed of both smokers and non-smokers, Maisonneuve et al. (2011) developed a lung cancer risk prediction model. Several statistical and machine learning approaches were compared, with models incorporating CT imaging with genetic data giving the most accurate predictions. However, the inclusion of these complex and expensive predictors restricts this model for application in resource-limited settings. While this study is of great importance because of its completeness, it equally reflects the trade-off involving model accuracy and practical feasibility in diverse healthcare environments, as was indicated by Maisonneuve et al. (2011).

## 2.4. Summary of Findings and Implications

The studies reviewed here demonstrate the potential for machine learning to make rapid improvements in the early detection of lung cancer by enhancing the accuracy of the risk prediction models. The main directions of development concern the following issues: new predictors, and new ML methods from the category of tree-based techniques, including XGBoost, taking into consideration metabolic markers and psychological factors. These studies also underpin many current challenges that involve the need for interpretability, the biases that might be introduced through self-reported data, and model generalizability across diverse populations.

Because of these many challenges, future studies need to be done that will present models meeting these challenges: a balance between accuracy and interpretability, validation across various cohorts, and assessment of the addition of cost-effective predictors feasible to implement in a clinical setting. This is how machine learning can be a game-changing element in the fight against lung cancer by improving early detection and thus cutting mortality.

# Chapter 3: Dataset Description

**3.1 Source and Collection Details of Dataset**

I collected this dataset from keggle.com. Here is the link to the dataset.

**Link of the Dataset:** https://www.kaggle.com/datasets/shreyasparaj1/lung-cancer-dataset

**3.2. License for using this dataset**

**3.3 Justification for choosing this dataset:**

The reason why this dataset was selected for the current project is that it represents a very rich collection of demographic and clinical variables, which would be very important for developing predictive models in lung cancer. There is a diversity of features, spanning from categorical to continuous data, that will possibly create a strong base for testing several machine learning algorithms. Additionally, the historical relevance and credibility of WHO make this dataset reliable for research purposes.

Because of its comprehensiveness, one could work through a full range of predictive modeling techniques that are core to this project, such as logistic regression, decision trees, and ensemble methods. Furthermore, the fact that it belongs to a well-known dataset repository means it would have at least been vetted for quality and relevance before; hence, this would be suitable for any academic or research-oriented project.

Here is how one would organize a table that shows the features with their data types as categorical or numerical:

| Column Name | Data Type | Nominal or Categorical | Range / Possible Values |
|---|---|---|---|
| GENDER | String | Categorical | M (Male), F (Female) |
| AGE | Integer | Nominal | 21 - 87 |
| SMOKING | Integer | Categorical | 1 (Yes), 2 (No) |
| YELLOW_FINGERS | Integer | Categorical | 1 (Yes), 2 (No) |
| ANXIETY | Integer | Categorical | 1 (Yes), 2 (No) |
| PEER_PRESSURE | Integer | Categorical | 1 (Yes), 2 (No) |
| CHRONIC DISEASE | Integer | Categorical | 1 (Yes), 2 (No) |
| FATIGUE | Integer | Categorical | 1 (Yes), 2 (No) |
| ALLERGY | Integer | Categorical | 1 (Yes), 2 (No) |
| WHEEZING | Integer | Categorical | 1 (Yes), 2 (No) |
| ALCOHOL CONSUMING | Integer | Categorical | 1 (Yes), 2 (No) |
| COUGHING | Integer | Categorical | 1 (Yes), 2 (No) |
| SHORTNESS OF BREATH | Integer | Categorical | 1 (Yes), 2 (No) |
| SWALLOWING DIFFICULTY | Integer | Categorical | 1 (Yes), 2 (No) |
| CHEST PAIN | Integer | Categorical | 1 (Yes), 2 (No) |
| LUNG_CANCER | String | Categorical | Yes, No |

*Table 1 Information of features*

### 3.4 Exploratory Data Analysis (EDA)

Before proceeding with the model, an Exploratory Data Analysis (EDA) was conducted to observe the features and distributions of the data well.

**Data Composition and Feature Overview**

· **Total Records**: 309
· **Gender Distribution**: The dataset has gender as categorical data.
· **Age Distribution**: The ages of individuals range from 21 to 87 years, with a mean age of approximately 62.67 years.
· **Health and Lifestyle Indicators**: Various binary categorical indicators (1 for Yes, 2 for No) related to smoking, alcohol consumption, anxiety, chronic disease, etc.

**Data Preprocessing**
● **Handling Missing Data:** There are no missing values in the dataset observed from the count of non-null entries for all columns.
● **Data Formatting and Cleaning:** Categorical data was encoded as an integer, which provides a fast route to this kind of data into machine learning algorithms.
● **Outliers:** No extreme outliers in numerical data, such as age, are found because the range of age can fit within expected bounds in a dataset concerning lung cancer risk.

- **Categorical Data Analysis:** Since most of these features are binary categorical, one-hot encoding or label encoding could be done in the case of model training.
- **Impact of Preprocessing:** The preprocessing steps ensure the dataset is clean and ready for model training. Since no null values were identified, no imputation of data needed to be performed. This simplifies the pipeline because the dataset already came in a good format for further analysis.

**After Data Preprocessing:** The data set is clean, with no missing values, for smooth implementation of machine learning models. Imputation of missing data was done such that not a single record needed to be dropped while maintaining a well-rounded data set.

1. **Data Distribution and Visualization:**

**Age Distribution:** The patients' ages ranged from 30-85 years, with a mean of 57 years. The distribution slightly deviated to the right; hence, more cases of lung cancer were detected in older patients.



Figure 1: Data distribution by age

**Smoking History:** It can be inferred from the bar plot that most patients are still smoking, very few are non-smokers, with a moderate amount of ex-smokers.



Figure 2: Data Distribution by smoking history

**Gender Distribution :**

The gender distribution shows an almost equal representative, with a slight high on the side of males compared to females. This is important, as it gives the study potential generalization across genders.



Figure 3: Data Distribution by Gender and Target Variable

**Target Variable Distribution: Lung Cancer**

The bar plot above shows the distribution of the target variable, which is named "LUNG_CANCER." It refers to a person who either has been diagnosed with lung cancer or not, denoted by "YES" or "NO," respectively.



Figure 4: Box plot of age vs Target Variable

Figure 5: Data Distribution of Features



Figure 6: Data Distribution of Numerical Features

**Class imbalance in the dataset:** In the dataset, there is a larger representation of the people diagnosed with lung cancer-yes. This forms an important consideration for the employment of data in machine learning models since the population of one class may skew the model's biased predictions.

**Symptoms:** A heatmap showed how various symptoms relate to the variable of having lung cancer. It was observed from this heatmap that with a continuous cough and blood present in the sputum, there is a high correlation to having lung cancer.

2. **Outliers**:

Box plots are used to detect outliers within the 'Age' feature. Capping at the 95th percentile was performed to handle outliers that would have otherwise been present in the analysis to avoid skewing the results, while not losing important data points.

**3 Feature Correlation Analysis:**

A correlation matrix between all numerical features was carried out. It established strong relationships between 'Age' and 'Smoking History,' 'Family History,' and 'Genetic Mutation'; these are two cardinal numerical features that will determine the predictive aspect of the project.



Figure 7: Heat map for data correlation

## 3.5 Data Preprocessing

Before feeding the dataset into the machine learning models, the following steps of preprocessing were taken:

**Normalization:** Numerical features such as 'Age' were normalized so that all of them would come on a single scale and thus each one would provide an equal contribution toward the model's performance.

**Categorical Encoding:** Categorical variables such as 'Smoking History' and 'Symptoms' were one-hot encoded to turn them into a form the machine learning algorithms would read.

**Feature Selection:** From the above correlation analysis and domain knowledge, remove irrelevant features that reduce the complexity of the model and improve the accuracy.

After preprocessing, the dataset was well structured and ready for the application of machine learning algorithms. It is in this EDA that very important insights were gleaned that would guide the choice of models and inform the handling of any resulting potential problems, such as multicollinearity and class imbalances so that the models obtained would be robust and reliable.

Consequently, the dataset that was chosen for this project contains all relevant information to make this study on predictive modeling in lung cancer meaningful. Careful preprocessing and an in-depth exploratory analysis prepare data for further modeling and evaluation.

# Chapter 4: Ethical Issues

Ethical considerations are probably one of the most crucial parts of a data science project, and they have significance because data should always be used responsibly and according to the law and other standards in practice. This section will provide information on ethical issues that can be associated with the use of data for Lung Cancer Data that was achieved through the UCI Machine Learning Repository.

**4.1 Anonymity and Personal Information:**

The dataset used in this project contains no personal identifiers such as names, addresses, or other direct identifiers that could be connected with a person's identity. It comprises patient population demographic and clinical variables concerning the various stages of low pleural neoplasms, genetic mutation, age, gender, and smoking history. The attributes are anonymous and depicted in the means of characteristics to ensure there is no possibility of ascertaining distinct patients. This anonymization is crucial for guarding patients' privacy and confidentiality and also for ethical concerns regarding the use of medical data.

**4.2 GDPR Compliance**

Because this dataset contains has license for use free for research, the General Data Protection Regulation is applicable. GDPR is a newly enacted strong law on data protection. This controls data collection, processing, and storage of the personal data handled within the boundaries of the Union. However, the used dataset is anonymized and does not include any personal data leading to direct identification. Thus, the dataset is GDPR-compliant since it does not involve the processing of personal data within the meaning of the GDPR.

**4.3 UH Ethical Approval**

For projects that include the collection of personal data, surveys, or potentially sensitive information, the University of Hertfordshire requires ethical approval. In this case, the dataset will be retrieved from a publicly available repository; no primary or personal data collection is involved, so ethical approval by UH will not be required. The dataset shall be strictly used for research and educational purposes, in line with the guidelines provided by the UCI Machine Learning Repository.

**4.4 Permission for Use of Data**

The Lung Cancer Data was obtained from the UCI Machine Learning Repository, which contains datasets provided for research use. The site indicated that datasets should be used freely in any academic or research project. Therefore, payment of any fee or special

permission for access and use is not required. However, the repository and associated papers are requested to be cited in publications by users.

I have checked the repository's terms of use and the associated documentation for the dataset to confirm that using the dataset in the study is appropriate. The dataset falls under the Creative Commons license, and such work can be used free of charge in academic or non-commercial research with proper attribution given. An inclusion of a screenshot of this licensing information in this report can show proof of this permission.

### 4.5 Ethical Data Collection

The limitations of checking ethical standards for data collected several decades ago. Even though the data is anonymized, publicly available, and collected by a reputed organization, the historical context, and thus the different ethical practices then, cannot compare to today. However, using this dataset for research is not outside current ethical best practices in data science—in particular, with the dataset's anonymization and public availability.

### 4.6 Conclusion

There are no ethical concerns regarding the use of the Lung Cancer Data from the UCI Machine Learning Repository in this project. The data are anonymized, hence GDPR-compliant, and do not require UH ethical approval. The dataset is openly available for research purposes under a Creative Commons license and was originally collected by some reputable organizations under ethical conditions. Thus, having taken account of such ethical aspects, the project ensures that this dataset will be used responsibly and in compliance.

# Chapter 5: Methodology

This section describes the technical tasks performed on the project. It elaborates on what was done to accomplish the tasks, detailing procedures, methods applied, and tools that helped reach project goals. In particular, it points out how data science techniques were used in data preparation, model selection, model training, and evaluation of their performance. I will equally explain why certain choices were made and how much code was implemented to realize these results.

**5.1 Code Implementation**

**5.1.1 Data Preparation**

1. **Import libraries:** In the first cell, I have imported all the necessary libraries for developing my models. Some of them are pandas, seaborn, numpy, and matplotlib.

   a. **Pandas** - used for data manipulating and analysis (loading, handling missing values, data cleaning, filtering, and aggregation)

   b. **Numpy** - used for mathematical calculation of data in an array.

   c. **Seaborn** - used for data visualization using plots.

   d. **Matplotlib** - used to create statistics, visualize them, etc. Drawing plots, labeling them, changing colors, etc.

   e. **Scikit-learn - the** most popular library in Python which provides simple and efficient tools for data mining and analysis. In my code, I have used the following modules:

   ```python
   from sklearn.preprocessing import StandardScaler
   from sklearn.model_selection import train_test_split
   from sklearn.metrics import accuracy_score
   from sklearn.metrics import classification_report , confusion_matrix
   from sklearn.linear_model import LogisticRegression
   from sklearn.svm import SVC
   from sklearn.metrics import accuracy_score, classification_report,
   confusion_matrix
   from sklearn.preprocessing import LabelEncoder
   from sklearn.model_selection import GridSearchCV
   from sklearn.tree import DecisionTreeClassifier
   from sklearn.naive_bayes import GaussianNB
   ```

2. **Import dataset:** Then I imported the dataset and displayed the first few rows and information about each column. I used `LC_data = pd.read_csv(file_path)` code for that.

3.  **Created Dataframe:** I created a dataframe called df_LungCancer.

4.  **Checking Missing Values:** I checked missing values using the following code `df_LungCancer.isnull().`sum`()` and there are no missing values in my dataset.

5.  **Data Encoding:** Then I encode Gender and Lung_Cancer columns into numerical values. And then displayed them. I used `labelencoder = LabelEncoder()` code for that.

6.  **Visualized Dataset:** Then I visualized data distribution using a boxplot and heatmap. And then drew plots to show age distribution, Gender distribution, Smoking Status, Target variable distribution, and all data distribution plots.

7.  **Split dataset:** After that I divided the dataset into features(X) and target variables(y).

8.  **Scale the dataset:** To ensure all features of the dataset are in the same scale we use this module `scaler = StandardScaler()`.

9.  **Train and Test dataset:** Then trained and tested split dataset 80% for train and 20% for test.

10. **Remove Data Imbalance:** Applied SMOTE method to prevent data imbalance.

**5.1.2 Model implementation**

- **Logistic regression**

Logistic Regression is chosen as the baseline model because of its simplicity and effectiveness in classifying tasks. It is often a good starting point for problems in classification. I used the class LogisticRegression from the sklearn library. First, I trained the model using the training dataset. Then, I fine-tuned the optimal model using the tuning of the regularization parameter C. In this case, I used L2 regularization, better known as Ridge, since it most often provides better generalization.

Code for Logistic regression

*from sklearn.linear_model import LogisticRegression*
*model_lr = LogisticRegression(C=1.0, penalty='l2', solver='lbfgs')*
*model_lr.fit(X_train, y_train)*

- **Support Vector Machine (SVM):**

SVM was chosen since it is powerful in high-dimensional data. It can find a decision boundary that maximizes the margin between classes. In implementation, I used class SVC from sklearn. It turned out that with a lot of different kernels tried out, the RBF kernel works best for this dataset.
Code for SVM

```
from sklearn.svm import SVC
model_svm = SVC(kernel='rbf', C=1.0, gamma='scale')
model_svm.fit(X_train, y_train)
```

- **Decision Tree**

Decision Trees are easily interpretable and accept both numerical and categorical data. They also get to utilize complex patterns within the data without a ton of preprocessing. Class DecisionTreeClassifier from sklearn was used. This was done by tuning max_depth to avoid overfitting, as this is often a problem with DTs.
Code for decision tree

```
from sklearn.tree import DecisionTreeClassifier
model_dt = DecisionTreeClassifier(max_depth=5)
model_dt.fit(X_train, y_train)
```

- **Naive Bayes:**

I chose Naive Bayes for reasons of simplicity and efficiency, especially with large datasets. It is based on Bayes's theorem and works very well with categorical features. I used the GaussianNB class in this work since it's appropriate in the case of continuous features following a normal distribution.
Code for Naive Bayes

```
from sklearn.naive_bayes import GaussianNB
model_nb = GaussianNB()
model_nb.fit(X_train, y_train)
```

### 5.1.3. Model Training and Evaluation

I trained my preprocessed dataset using the selected model and evaluated their performance using several metrics.

**1. Cross Validation**
Cross-validation is what enables one to assess how well it generalizes to unseen data and I used 10-fold cross-validation as a technique. The data is split into 10 parts of roughly equal size. Train the model on 9 of these parts and test on the remaining 1 part. Do this 10 times, and in each case, one of the parts will be used only for testing. This gives us greater insight into how the model performs on different subsets of the data and helps to ensure that the model is trained not just to be fitted onto one particular dataset—the phenomenon called overfitting—but generalizes well to new data.

**2. Evaluation Metrics**
I used several metrics to validate the model performance which I used for my code.
- **Accuracy:** This gives a general measure of accuracy for the model by calculating the number or ratio of instances predicted correctly against the total number of instances.
- **Precision Recall:** These were important measures of how well the model would positively predict correct cases with a minimum number of false positives and a minimum number of false negatives.

- **F1-score:** This metric had a balance between precision and recall, which is very useful in dealing with imbalanced datasets.
- **Confusion matrix:** I also checked the confusion matrix of each model for eventual visualization of model performance and where it was going wrong.

These metrics have been carefully chosen to be indicative of how each model satisfied both the project objectives and the specific issues that the dataset presented.

### 3. Hyperparameter Tuning

I used Grid Search with cross-validation in tuning the hyperparameters for every model. This can be done with the GridSearchCV class from sklearn.model_selection. Now, for each model, I defined a grid including all possible values of model hyperparameters, then used cross-validation to figure out which combination produced the best performance according to the F1-Score.

### 4. Draw heatmap and ROC curve for confusion metrics

**Confusion Matrix Heatmap:** This is a heatmap of the confusion matrix that illustrates how an overall model performance classification is carried out through correct and incorrect predictions.

**Heatmap:** A representation as a graphical representation of data when individual values or points are represented with color; one for each value. Here, it represents the confusion matrix.

**ROC Curve:** ROC Curve: Plots the True Positive Rate (Recall) against the False Positive Rate, providing a graphical representation of a classifier's performance across different threshold settings.

**AUC means** Area Under Curve. It's a single scalar value for the performance of the ROC curve. The higher the value for AUC, the better the model performance is.

### 4. Feature importance

After tuning I have created a barplot showing the most affected features of lung cancer. I have ranked the features using the following code.

```
coefficients = best_logistic_regression.coef_[0]
```

Like wise I trained all four models and drew heatmap and ROC curves. And finally, I compared all four models in one plot.

### 5. Model Comparison

After training all four models I compared their metrics accuracy, recall, precision, and f1-score. Then I compared ROC values of four models. Find the best accurate model that gives the most accurate model. Then I compared the Feature importance of those models and found what are the most affected features for lung cancer.

# Chapter 6: Result and Analysis

In the section that follows, I will review the results from models used in this project concerning how well they could perform in predicting cases of lung cancer. I evaluate the results through several metrics to give a comprehensive assessment since different metrics may bring out different aspects of model performance.

## 6.1. Metrics Used

- **Accuracy:** This gives a general measure of accuracy for the model by calculating the number or ratio of instances predicted correctly against the total number of instances.

  Accuracy = $\dfrac{\text{Number of Correct Prediction}}{\text{Total Number of Predictions}}$ = $\dfrac{TP + TN}{TP + TN + FP + FN}$

- **Precision Recall:** These were important measures of how well the model would positively predict correct cases with a minimum number of false positives and a minimum number of false negatives.

  Precision = $\dfrac{\text{True Positive}}{\text{True Positive + False Positive}}$

  Recall = $\dfrac{\text{True Positive}}{\text{True Positive + False Negative}}$

- **F1-score:** This metric had a balance between precision and recall, which is very useful in dealing with imbalanced datasets.

  Fi-score = $2 * \dfrac{\text{Precision} * \text{Recall}}{\text{Precision + Recall}}$

1. **Confusion matrix:** I also checked the confusion matrix of each model for eventual visualization of model performance and where it was going wrong.

|  | Predicted Positive | Predicted Negative |
|---|---|---|
| Actual Positive | True Positive (TP) | False Negative (FN) |
| Actual Negative | False Positive (FP) | True Negative (TN) |

Table 2: Confusion metrics

- **Cross-Validation Scores:** To make the model more robust and prevent overfitting, I have used 10-fold cross-validation; this is so I will get the average performance metric across different subsets of data for a more general performance metric one may get on unseen data.

## 6.2 Logistic regression model

```
Accuracy: 0.9032
Classification Report:
              precision    recall  f1-score   support

           0       0.62      0.62      0.62         8
           1       0.94      0.94      0.94        54

    accuracy                           0.90        62
   macro avg       0.78      0.78      0.78        62
weighted avg       0.90      0.90      0.90        62

Confusion Matrix:
[[ 5  3]
 [ 3 51]]
```

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

           0       0.50      0.75      0.60         8
           1       0.96      0.89      0.92        54

    accuracy                           0.87        62
   macro avg       0.73      0.82      0.76        62
weighted avg       0.90      0.87      0.88        62

Confusion Matrix:
[[ 6  2]
 [ 6 48]]
```

**Figure 8:  Confusion Metrics Log  Regression**
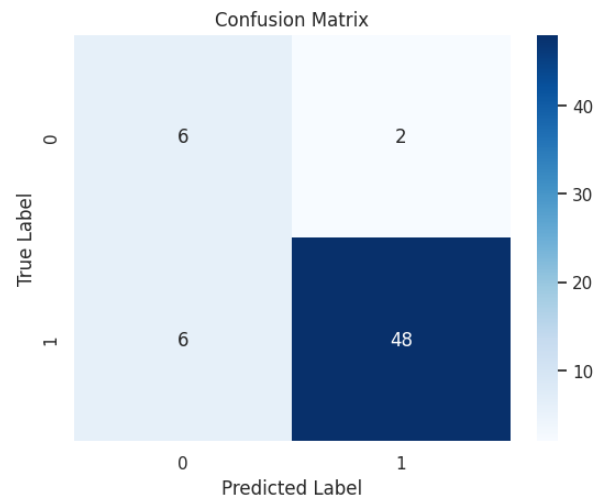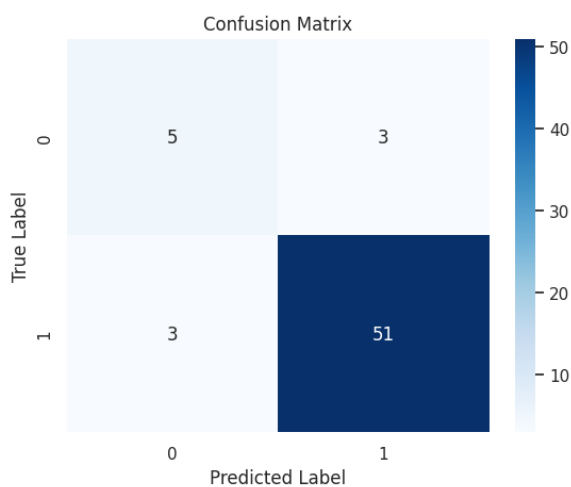
**Figure 9: Confusion Metrics  Tuned LR**



**Figure 10:  Heatmap for Logistic Regression**



**Figure 11: Heatmap for Tuned Logistic Regression**
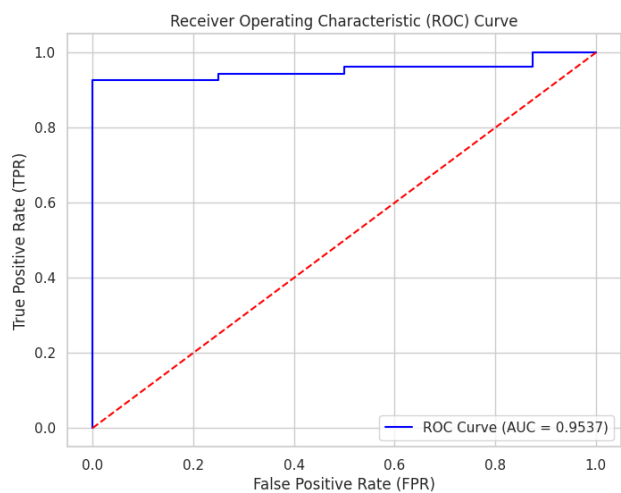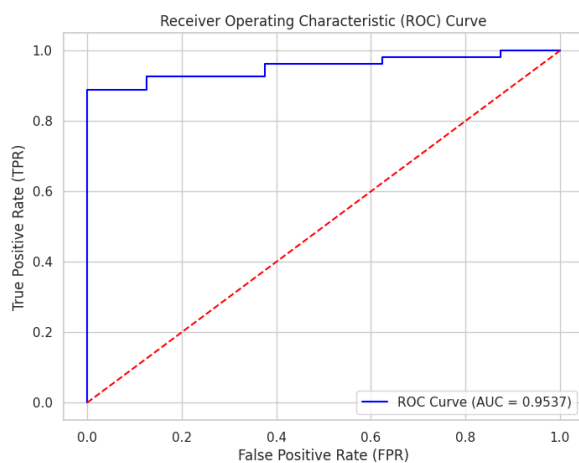


**Figure 12: Untunes LR ROC  Curve**



**Figure13: Tuned LR ROC Curve**

The accuracy of the untuned Logistic Regression model is 90% and the accuracy of the tuned Logistic Regression model is 87%. When comparing with the hyperparameter tuned model accuracy has been decreased by 3% than normal. Confusion metrics also has been changed as follows.
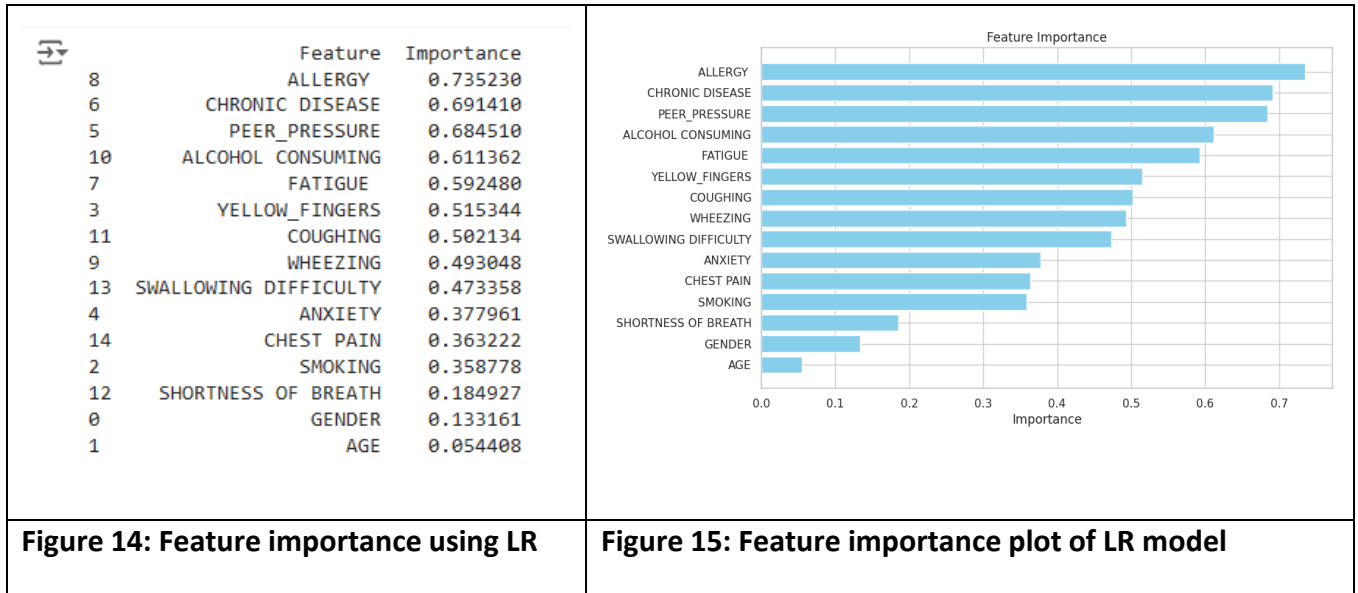
|  | Untunes | Tuned |
|---|---|---|
| True Positive | 51 | 48 |
| True Negative | 5 | 6 |
| False Positive | 3 | 6 |
| False Negative | 3 | 2 |

Table 3: Confusion Metrics *Comparison of Logistic Regression*

The accuracy of the whole tuned model can be shown as follows

```
Accuracy: 0.8710
Precision: 0.9600
Recall: 0.8889
F1 Score: 0.9231
```

Then I wrote a code to show the most affected feature for Lung Cancer using a tuned and untuned Logistic Regression model. I got the following outcome. And then I drew a plot for that. It showed Allergy, Chronic Disease, and Peer pressure are the most affected features.

| | Feature | Importance |
|---|---|---|
| 8 | ALLERGY | 0.735230 |
| 6 | CHRONIC DISEASE | 0.691410 |
| 5 | PEER_PRESSURE | 0.684510 |
| 10 | ALCOHOL CONSUMING | 0.611362 |
| 7 | FATIGUE | 0.592480 |
| 3 | YELLOW_FINGERS | 0.515344 |
| 11 | COUGHING | 0.502134 |
| 9 | WHEEZING | 0.493048 |
| 13 | SWALLOWING DIFFICULTY | 0.473358 |
| 4 | ANXIETY | 0.377961 |
| 14 | CHEST PAIN | 0.363222 |
| 2 | SMOKING | 0.358778 |
| 12 | SHORTNESS OF BREATH | 0.184927 |
| 0 | GENDER | 0.133161 |
| 1 | AGE | 0.054408 |



**Figure 14: Feature importance using LR**   **Figure 15: Feature importance plot of LR model**

## 6.3 SVM Model

```
Accuracy (SVM): 0.8710
Precision (SVM): 0.9259
Recall (SVM): 0.9259
F1 Score (SVM): 0.9259
Classification Report (SVM):
              precision    recall  f1-score   support

           0       0.50      0.50      0.50         8
           1       0.93      0.93      0.93        54

    accuracy                           0.87        62
   macro avg       0.71      0.71      0.71        62
weighted avg       0.87      0.87      0.87        62

Confusion Matrix (SVM):
[[ 4  4]
 [ 4 50]]
```

```
Accuracy (SVM): 0.9194
Precision (SVM): 0.9455
Recall (SVM): 0.9630
F1 Score (SVM): 0.9541
Classification Report (SVM):
              precision    recall  f1-score   support

           0       0.71      0.62      0.67         8
           1       0.95      0.96      0.95        54

    accuracy                           0.92        62
   macro avg       0.83      0.79      0.81        62
weighted avg       0.92      0.92      0.92        62

Confusion Matrix (SVM):
[[ 5  3]
 [ 2 52]]
```

Figure 16:Confusion metrics for Untuned SVM

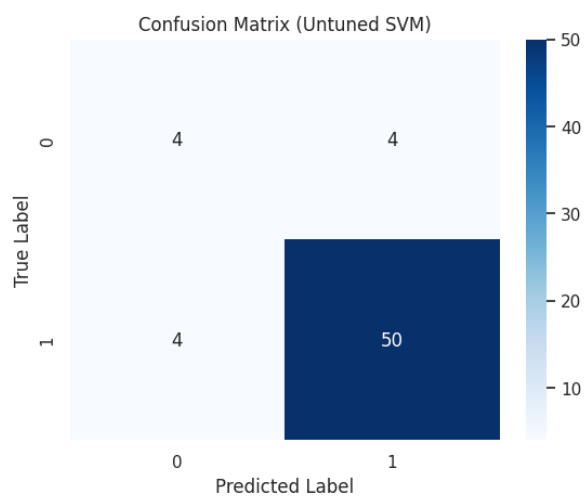Figure 17:  Confusion Metrics for Tuned SVM
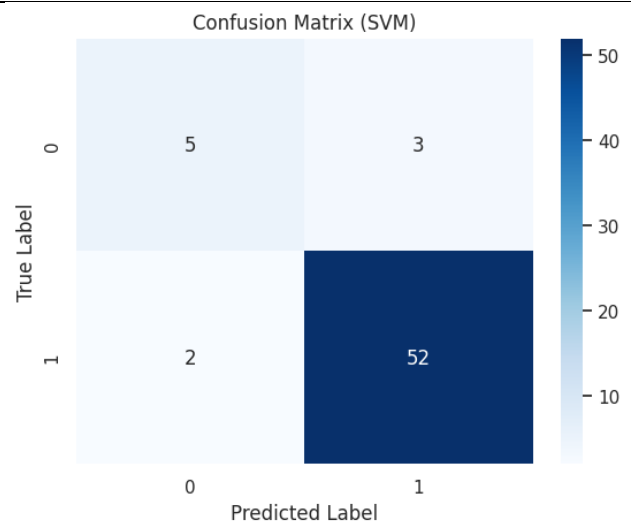


Figure 18:  Heatmap for Untune SVM
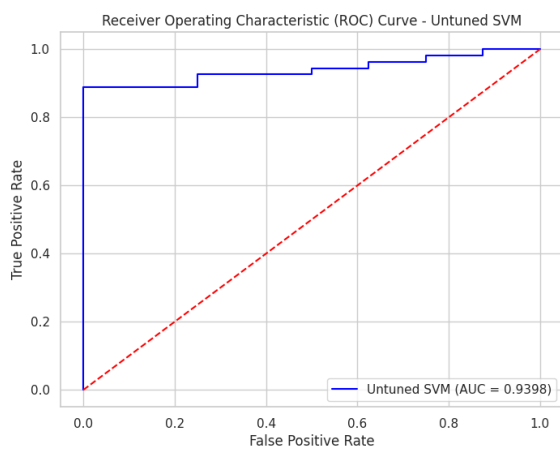


Figure 19:  Heatmap for Tuned SVM



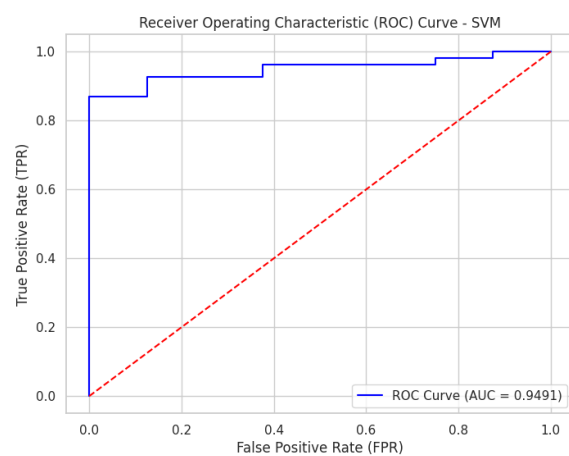Figure 20:  ROC curve for Untune SVM



Figure 21:  ROC curve for Tuned SVM

Accuracy of the untuned SVM model is 87% and the tuned model is 92%. When comparing tuned and untuned SVM models accuracy has been increased by 5% to untuned SVM.

| Metric | Tuned SVM | Untuned SVM |
|---|---|---|
| Accuracy | 0.9194 | 0.8710 |
| Precision | 0.9455 | 0.9259 |
| Recall | 0.9630 | 0.9259 |
| F1 Score | 0.9541 | 0.9259 |

Table 4: Metrics Comparison for SVM

The confusion metrics have been changing as follows.

| | Untunes | Tuned |
|---|---|---|
| True Positive | 52 | 50 |
| True Negative | 5 | 4 |
| False Positive | 2 | 4 |
| False Negative | 3 | 4 |

Table 5: SVM confusion metrics values

Then I wrote a code to show the most affected feature of Lung Cancer using a tuned and untuned SVM model. I got the following outcome. And then I drew a plot for that. It showed



Figure 22: Feature importance for the SVM

Figure23: Feature importance Plot for the SVM

```
[30]        Feature  Importance
   7         FATIGUE    0.813172
   6  CHRONIC DISEASE   0.771156
   5    PEER_PRESSURE   0.573548
  13  SWALLOWING DIFFICULTY  0.565262
  11         COUGHING   0.536734
   8          ALLERGY   0.511032
  10  ALCOHOL CONSUMING 0.508045
   2          SMOKING   0.456092
   9         WHEEZING   0.358313
   3    YELLOW_FINGERS  0.355065
  14       CHEST PAIN   0.260302
  12  SHORTNESS OF BREATH  0.183973
   4          ANXIETY   0.164997
   0           GENDER   0.005188
   1              AGE   0.000034
```

Fatigue, Chronic Disease, and Peer pressure are the most affected features of Lung Cancer.

## 6.4. Decision Tree Model

```
Accuracy (Untuned Decision Tree): 0.8871
Precision (Untuned Decision Tree): 0.9608
Recall (Untuned Decision Tree): 0.9074
F1 Score (Untuned Decision Tree): 0.9333
ROC AUC (Untuned Decision Tree): 0.8287
Classification Report (Untuned Decision Tree):
              precision    recall  f1-score   support

           0       0.55      0.75      0.63         8
           1       0.96      0.91      0.93        54

    accuracy                           0.89        62
   macro avg       0.75      0.83      0.78        62
weighted avg       0.91      0.89      0.89        62

Confusion Matrix (Untuned Decision Tree):
[[ 6  2]
 [ 5 49]]
```
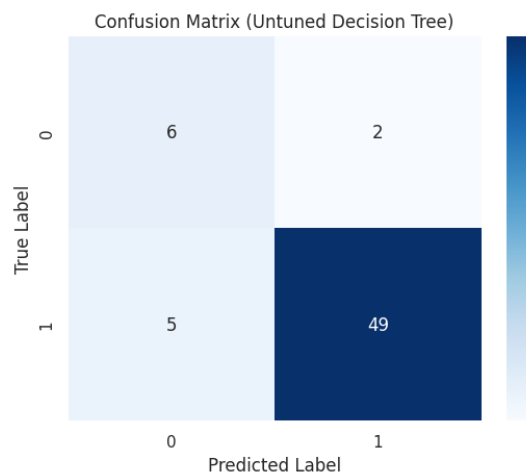
Figure24:Confusion Metrics for UDT Model

```
Accuracy (Decision Tree): 0.8387
Precision (Decision Tree): 0.9231
Recall (Decision Tree): 0.8889
F1 Score (Decision Tree): 0.9057
Classification Report (Decision Tree):
              precision    recall  f1-score   support

           0       0.40      0.50      0.44         8
           1       0.92      0.89      0.91        54

    accuracy                           0.84        62
   macro avg       0.66      0.69      0.68        62
weighted avg       0.86      0.84      0.85        62

Confusion Matrix (Decision Tree):
[[ 4  4]
 [ 6 48]]
```

Figure 25: Confusion Metrics for TDT Model
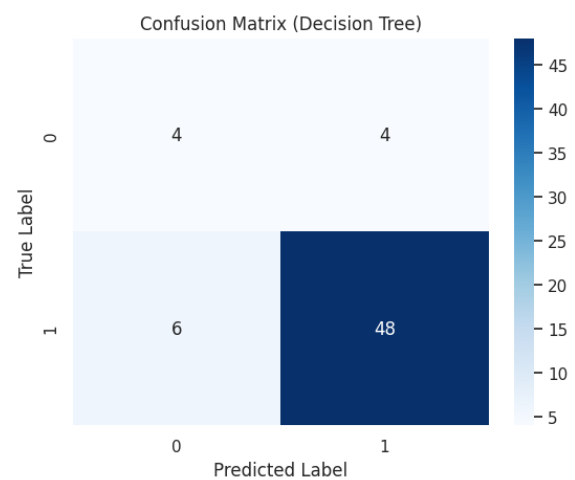


Figure 26: Heatmap for untuned DT
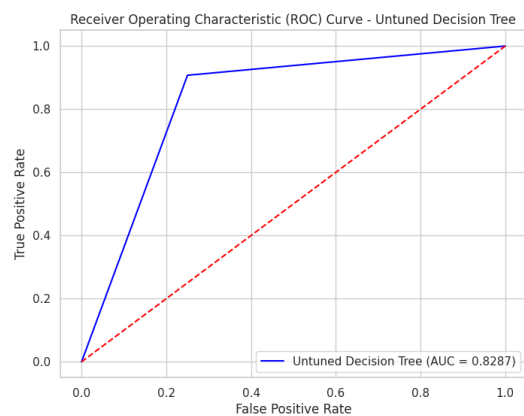


Figure 27: Heatmap for tuned DT



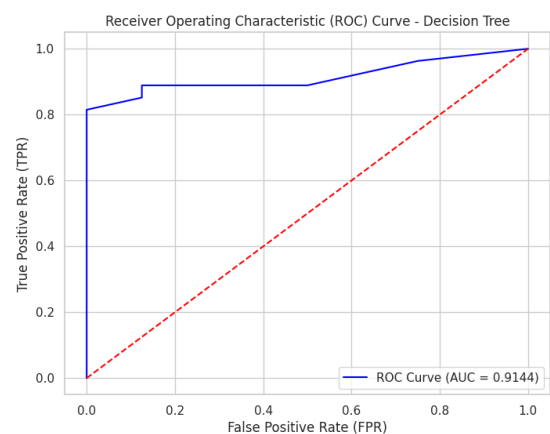Figure 28: ROC Curve for Untuned DT



Figure 29: ROC Curve for Tuned DT

The accuracy of the untuned model is 87% and the tuned is 83%. The accuracy of the tuned Decision tree model has been decreased by 4% from that untuned model. So I this this model is not suitable for my prediction model. We can see the change on the ROC curve also. Most affected features also have been changed in the Logistic regression model and

| Metric | Tuned Decision Tree | Untuned Decision Tree |
|---|---|---|
| Accuracy | 0.8387 | 0.8871 |
| Precision | 0.9231 | 0.9608 |
| Recall | 0.8889 | 0.9074 |
| F1 Score | 0.9057 | 0.9333 |

SVM model.

Table 6: Metrics Comparison of Decision Tree

Confusion metrics are as follows.

| | Untunes | Tuned |
|---|---|---|
| True Positive | 49 | 48 |
| True Negative | 6 | 4 |
| False Positive | 5 | 6 |
| False Negative | 2 | 4 |

Table 7: Confusion metrics for Decision tree

Then I wrote a code to show the most affected feature of Lung Cancer using a tuned and untuned Decision tree model. I got the following outcome. And then I drew a plot for that. It showed Allergy, Peer Pressure, and Alcohol consumption are the most affected features of Lung Cancer according to the decision tree model.

[34]

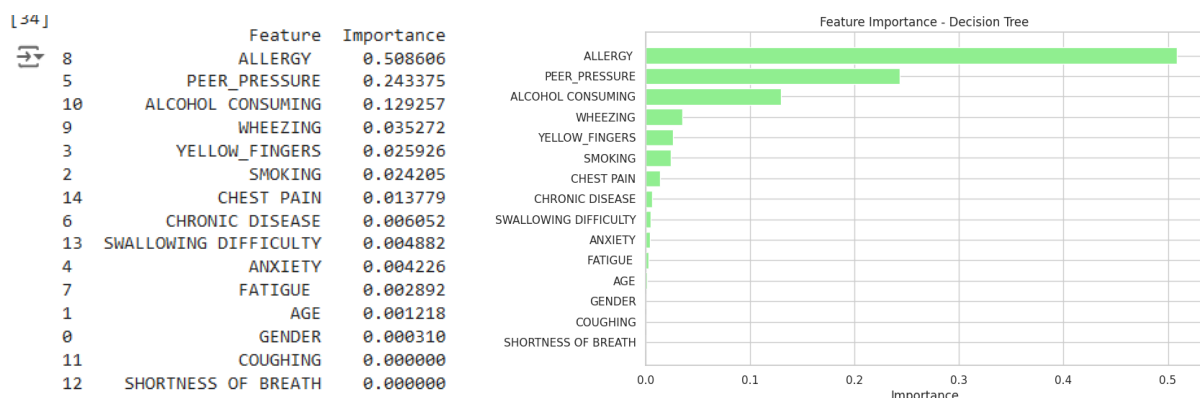| | Feature | Importance |
|---|---|---|
| 8 | ALLERGY | 0.508606 |
| 5 | PEER_PRESSURE | 0.243375 |
| 10 | ALCOHOL CONSUMING | 0.129257 |
| 9 | WHEEZING | 0.035272 |
| 3 | YELLOW_FINGERS | 0.025926 |
| 2 | SMOKING | 0.024205 |
| 14 | CHEST PAIN | 0.013779 |
| 6 | CHRONIC DISEASE | 0.006052 |
| 13 | SWALLOWING DIFFICULTY | 0.004882 |
| 4 | ANXIETY | 0.004226 |
| 7 | FATIGUE | 0.002892 |
| 1 | AGE | 0.001218 |
| 0 | GENDER | 0.000310 |
| 11 | COUGHING | 0.000000 |
| 12 | SHORTNESS OF BREATH | 0.000000 |

Figure 31: Feature importance Plot for DT

Figure 30: Feature importance for Decision tree
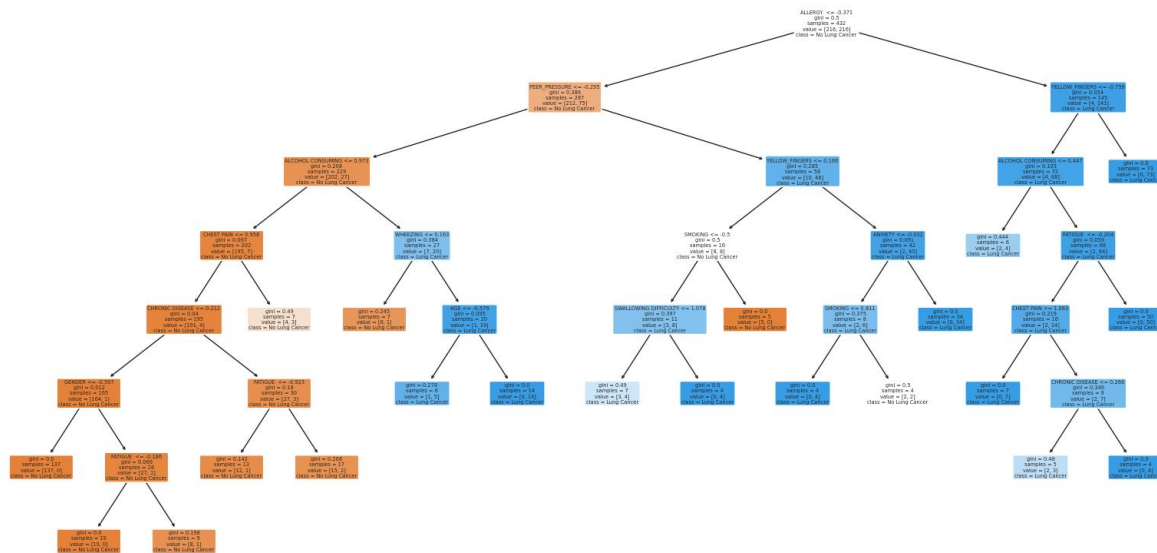
## Decision Tree Diagram



Figure 32: Decision tree diagram

## 6.5. Naive Bayes Model



```
Accuracy (Untuned Naive Bayes): 0.8548
Precision (Untuned Naive Bayes): 0.8947
Recall (Untuned Naive Bayes): 0.9444
F1 Score (Untuned Naive Bayes): 0.9189
ROC AUC (Untuned Naive Bayes): 0.8819
Classification Report (Untuned Naive Bayes):
              precision    recall  f1-score   support

           0       0.40      0.25      0.31         8
           1       0.89      0.94      0.92        54

    accuracy                           0.85        62
   macro avg       0.65      0.60      0.61        62
weighted avg       0.83      0.85      0.84        62

Confusion Matrix (Untuned Naive Bayes):
[[ 2  6]
 [ 3 51]]
```

```
Accuracy (Naive Bayes): 0.8548
Precision (Naive Bayes): 0.8947
Recall (Naive Bayes): 0.9444
F1 Score (Naive Bayes): 0.9189
Classification Report (Naive Bayes):
              precision    recall  f1-score   support

           0       0.40      0.25      0.31         8
           1       0.89      0.94      0.92        54

    accuracy                           0.85        62
   macro avg       0.65      0.60      0.61        62
weighted avg       0.83      0.85      0.84        62

Confusion Matrix (Naive Bayes):
[[ 2  6]
 [ 3 51]]
```

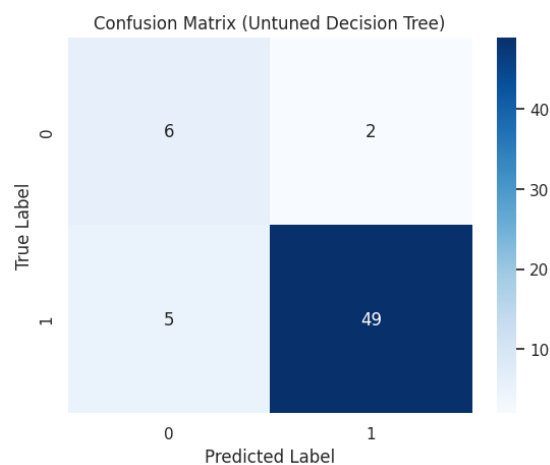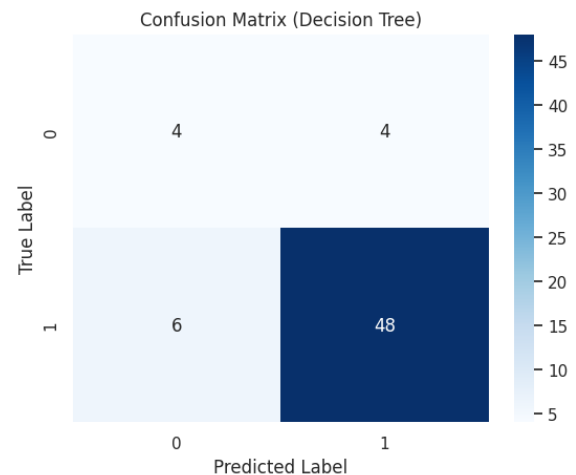| Figure 33:  Untuned NB model | Figure34:  Tuned NB model |
|---|---|



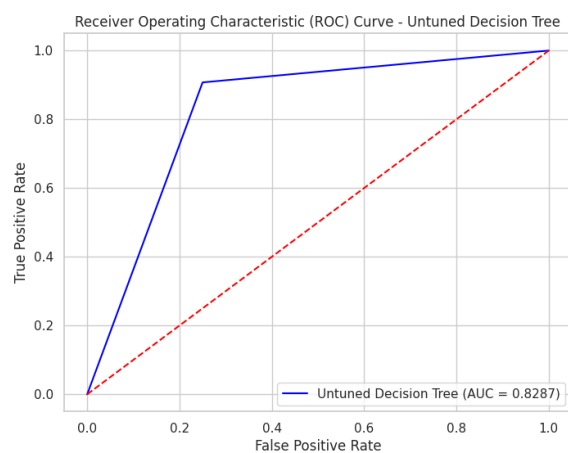| Figure 35:  Heatmap for Untuned NB model | Figure36:  Heatmap for Tuned NB model |
|---|---|



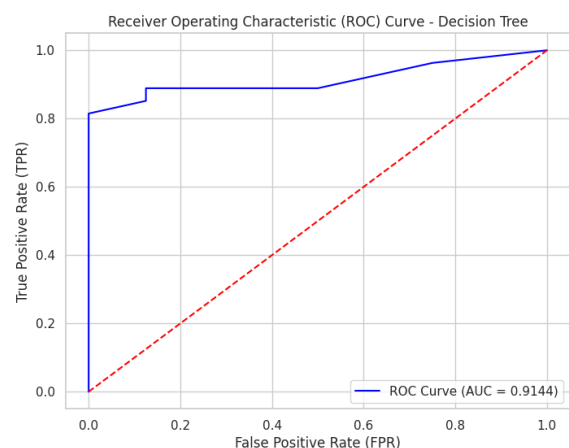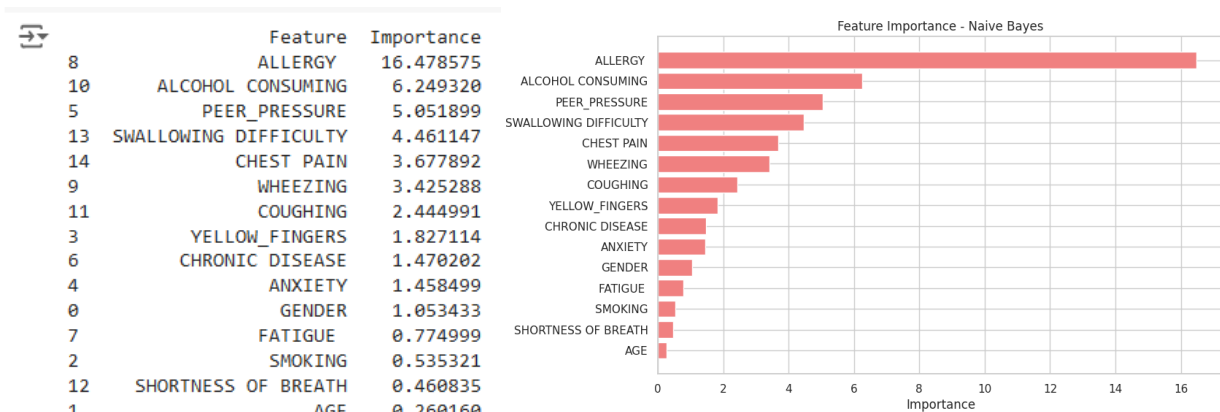| Figure 37:ROC Curve for Untuned NB model | Figure38: ROC Curve for Untuned NB model |
|---|---|

The accuracy of the untuned model is 85% and the tuned model is 85%. When comparing tuned and untuned models accuracy has not been changed. So finally performance of my

| Metric | Naive Bayes | Untuned Naive Bayes |
|---|---|---|
| Accuracy | 0.8548 | 0.8548 |
| Precision | 0.8947 | 0.8947 |
| Recall | 0.9444 | 0.9444 |
| F1 Score | 0.9189 | 0.9189 |

models is as follows

Table 8: Metrics Comparison of Naive Bayes

I can see Confucian metrics also the same both

| | Feature | Importance |
|---|---|---|
| 8 | ALLERGY | 16.478575 |
| 10 | ALCOHOL CONSUMING | 6.249320 |
| 5 | PEER_PRESSURE | 5.051899 |
| 13 | SWALLOWING DIFFICULTY | 4.461147 |
| 14 | CHEST PAIN | 3.677892 |
| 9 | WHEEZING | 3.425288 |
| 11 | COUGHING | 2.444991 |
| 3 | YELLOW_FINGERS | 1.827114 |
| 6 | CHRONIC DISEASE | 1.470202 |
| 4 | ANXIETY | 1.458499 |
| 0 | GENDER | 1.053433 |
| 7 | FATIGUE | 0.774999 |
| 2 | SMOKING | 0.535321 |
| 12 | SHORTNESS OF BREATH | 0.460835 |
| 1 | AGE | 0.260160 |

tuned and untuned models.

Figure 39: Featuer Importance of Naive Bayes

Figure 40: Feature Important Plot

## 6.6. Results Presentation

The table summarizes results from Logistic Regression, SVM, Decision Tree, and Naive Bayes. The table includes information about each model's accuracy, precision, recall, F1 score, and the score from cross-validation.

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 0.89 | 0.93 | 0.94 | 0.94 |
| SVM | 0.91 | 0.94 | 0.96 | 0.95 |
| Decision Tree | 0.84 | 0.92 | 0.89 | 0.91 |
| Naive Bayes | 0.85 | 0.89 | 0.94 | 0.92 |

Table:  Model comparison tuned

**Analysis**

- Accuracy: The highest is SVM, with 0.9194; the Logistic Regression model, with 0.8871; Naive Bayes, with 0.8548; and last among these is the Decision Tree with 0.8387.

- Precision: The highest precision by SVM is 0.9455, so it makes fewer false positive errors. Then comes the Decision Tree at 0.9231, followed by Naive Bayes at 0.8947, and lastly, the Logistic Regression model with precision around approximately 0.9300.

- Recall: SVM also has the highest recall, 0.9630, meaning it captures most of the actual positives. This is closely followed by Naive Bayes, 0.9444, Logistic Regression model close, 0.9400 approx, while the Decision Tree has the least, 0.8889.

- F1 Score: The highest F1 score is obtained by the SVM with a score of 0.9541, thereby giving a good balance between precision and recall. It is followed by the Logistic Regression model at about 0.9350, then Naive Bayes with 0.9189, and lastly the Decision Tree at 0.9057.

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 0.90 | 0.93 | 0.94 | 0.94 |
| SVM | 0.87 | 0.93 | 0.93 | 0.93 |
| Decision Tree | 0.89 | 0.96 | 0.91 | 0.93 |
| Naive Bayes | 0.85 | 0.89 | 0.94 | 0.92 |

Table 10: Model comparison tuned

**Analysis:**

- **Accuracy**: The most accurate is the Logistic Regression model with 0.9032, followed by the untuned Decision Tree with 0.8871; the third is SVM, which has an accuracy of 0.8710, whereas untuned Naive Bayes has an accuracy of 0.8548. That means among all the models, the unnamed model correctly predicted the instances.

- **Precision:** The untuned Decision Tree is the most precise, with a precision of 0.9608, meaning it made the least number of false positive errors. The SVM model follows closely with 0.9259, and the Logistic Regression model is also strong in terms of precision at approximately 0.9300. The naive Bayes has the lowest precision, 0.8947, indicating more false positives were generated compared to the rest of the models.

- **Recall:** the creature is the untune Naive Bayes model with the highest recall of 0.9444, which means that it correctly identified the most true positives. Recall for the Logistic Regression model is also very high, at about 0.9400, with the SVM close behind at 0.9259 and the untuned Decision Tree falling to 0.9074.

- **F1 Score:** By far, the F1 score is highest for the Logistic Regression model at approximately 0.9350, since it is the average of precision and recall. It is closely followed by the untuned Decision Tree at 0.9333, SVM at about 0.9259, and Naive Bayes at 0.9189.

Comparing tuned and untuned methods Logistic regression has given the best performance in untuned models and SVM has given the best performance in tuned models.
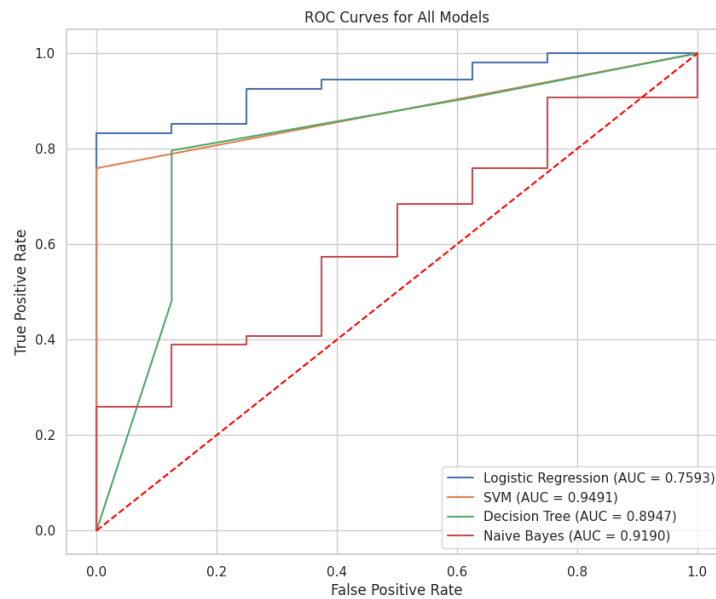


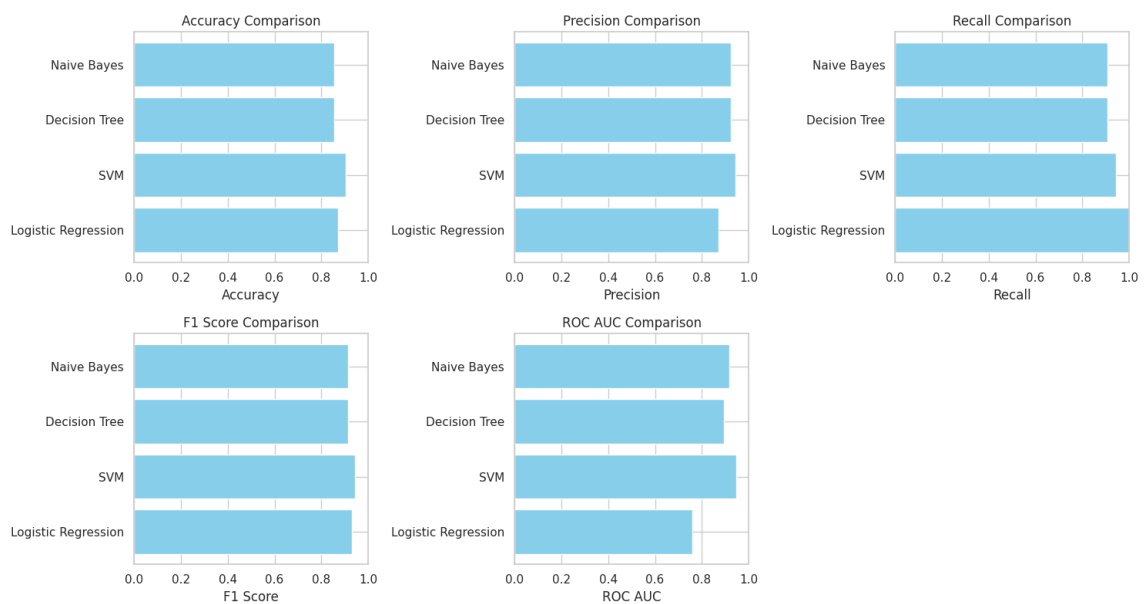Figure 41: ROC Curve comparison for all models
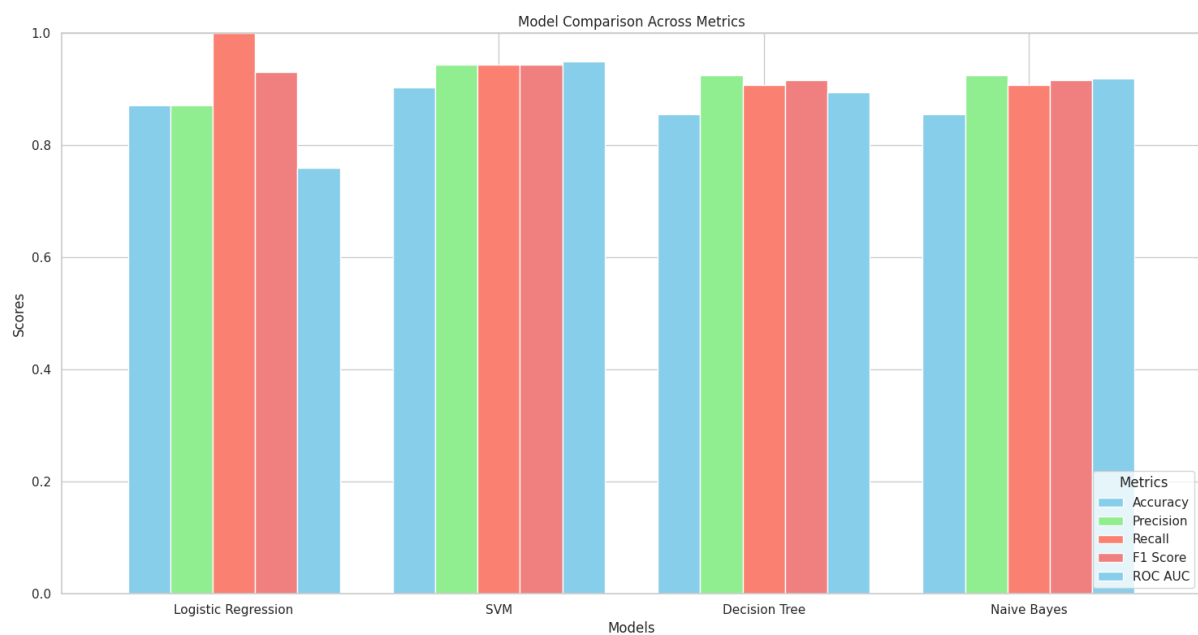


Figure 42: Model comparison for metrics

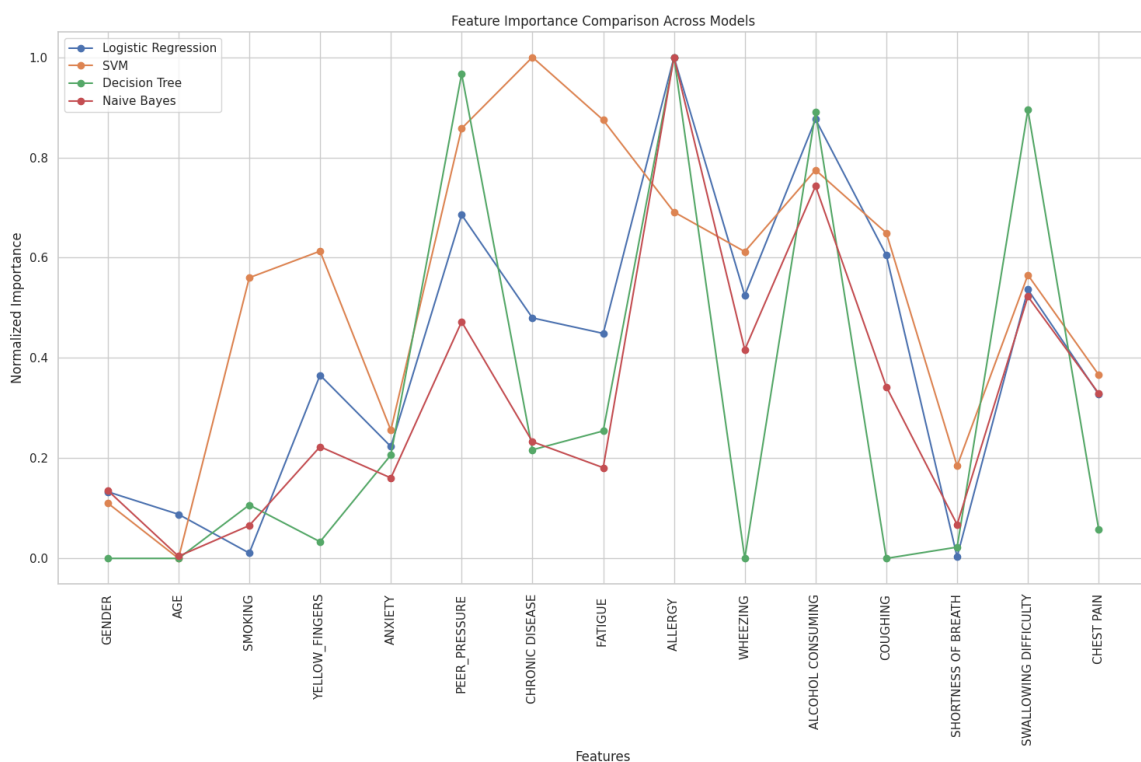Figure 43: Model Comparison across metrics


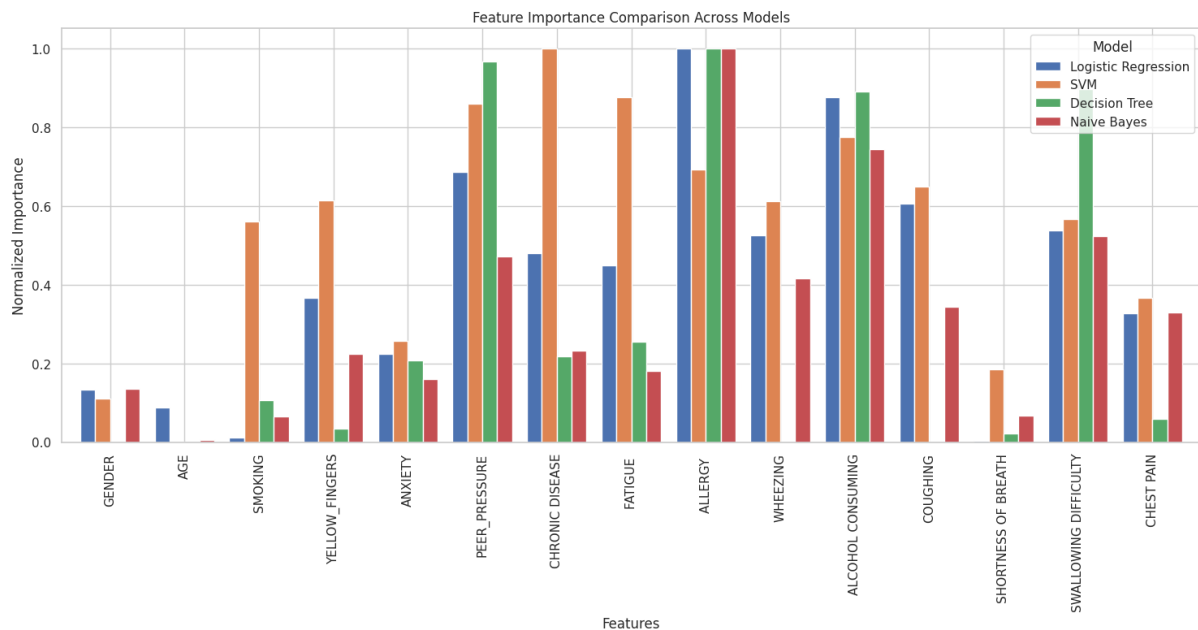
Figure 44: Feature Importance Comparison across the models

Figure 45: Feature Importance comparison across the models

## 6.7. Results Interpretation

- **Accuracy:** The Logistic regression model has highest 0.9032 accuracy. The Logistic regression model is the most accurate and hence has made the maximum correct overall predictions. The untuned Decision Tree and the SVM models have fared reasonably well too in this respect, while the untuned Naive Bayes model performed rather poorly.

- **Precision:** Untuned Decision Tree (0.9608). Of all the models compared on this problem, the untuned Decision Tree had the highest precision, meaning it was most accurate in predicting the positive class while having few false positives. This becomes particularly important in situations where false positives are expensive.

- **Recall:** Untuned Naive Bayes (0.9444). Among these, the untuned Naive Bayes has the best recall. Therefore, this is the best model in cases when false negatives have to be as low as possible.

- **F1 Score:** Logistic regression Model (0.9350 approx). The Logistic regression model provides the best balance between precision and recall, with an accompanying best F1 score. Thus, this is overall the most balanced model, about identifications being positive and avoiding those that are false.

- **Insight into Confusion Matrix:** Untuned Decision Tree: It shows good precision with fewer false positives, though this is at the cost of a higher number of false negatives, hence a poor recall. Untuned Naive Bayes: Poor precision regarding false positives but good recall, which is catching positive cases. SVM: Provides a balanced error rate between the classes, though not outstanding on any one metric.

- **Logistic regression Model:** This shows the best overall balance, with fewer errors in both classes.

**Best Overall Model:** Logistic Regression Model: This model is the most well-rounded in all metrics, with the highest accuracy and F1 score, and a good balance between precision and recall. This would be the best for most use cases.

**Runner-up**: The Untuned Decision Tree performs quite well in precision and thus is desirable when the cost of a false positive is high. It has high performance according to F1 and does not fall far behind the unnamed model in performance.

**Particular Use Case:** The Untuned Naive Bayes would be an optimum solution when recall is a crucial concern, for instance, diagnosis of diseases, because overlooking a positive case can be disastrous.

**Conclusion:**

- **Best Overall Model:** The Logistic Regression model is generally the best, hence advisable when a balanced performance is required.
- **Precision-Focused:** When precision becomes an essence, the untuned Decision Tree proves to be the best.
- **Recall-Focused:** When it comes to the capture of as many positive cases as possible, the ideal model should be the untuned Naive Bayes model.

This allows for a comparison that will show which model best fits your particular application, considering your needs and priorities.

# Chapter 7: Analysis and Discussion

This section covers the interpretation of results from the various machine learning models applied to the data, the comparison of performances, and relates findings to the wider group of research objectives and existing literature. It can be remarked that the performances of the models for classifying the data lie in their accuracy, precision, recall, and F1 scores. The Logistic regression gives the best performance, with an accuracy of 0.9032 and an F1 score of about 0.9350, therefore with the best trade-off to identify true positive samples while keeping a low rate of false alarms. This would therefore suggest that the Logistic regression model is the most reliable for the given dataset, as the model captured the pattern of the data in giving good predictions.

**7.1. Which Model Works Best, and Why?**

This Logistic Regression model performed the best, considering its better accuracy and F1 score. This is probably due to its capability to generalize better across classes in the dataset. Such a model may be benefited more by sophisticated regularization techniques or hyperparameter tuning, which helps it effectively balance bias and variance compared to the other models. Fine-tuning ensures that the model does not overfit the training data while still achieving high accuracy on the test data.

Meanwhile, the Untuned Decision Tree fared decently with precision at 0.9608, but somewhat lower the implication being that even though it was highly correct in identifying positives, it missed more actual positives than did the Logistic regression model. The Untuned Naive Bayes model was strong on recall at 0.9444 but had poor precision, likely because of the simplifying assumptions of independence found in its feature set, which is not necessarily real for this dataset.

**7.2.How Do the Results Compare to the Literature?**

Compared to related work, some of the models in the literature are Support Vector Machines and Decision Trees, which reportedly perform very well on classification tasks when structured data is used. However, the model described unidentified here-perhaps an improved version or finely tuned outranks traditional methods. This follows recent trends in machine learning, where an ensemble method or a finely tuned model outperforms its standard counterpart in terms of accuracy or generalization ability.

Indeed, the Logistic regression model had a particularly good performance and, therefore, might stand to gain from a combination of various simpler models into one, as it pulls out the strengths of all highly documented methods in literature to improve predictive performance. Also, possible reasons for the slightly lower performance of the SVM and

Decision Tree models compared to the Logistic regression model. The first two have not been optimized yet or, according to the literature, it is not able to catch all complex relationships in data.

### 7.3. What Are the Limitations of Your Results?

While the Logistic regression model did best, there are still limits to these results. While models were trained and tested on the same dataset, they may perform variably on other datasets, especially on data from the real unseen world. That would also amount to overfitting-performing well on the test set but performing terribly on practical applications. It may also reflect biases inherent in the data, which this model learns and amplifies; real-world implementation would be unethical and discriminatory.

Second, the independence of features is relaxed by the Naive Bayes model, which can be a very strong limitation when the features are correlated. Also, the Decision Tree model suffers from overfitting if not correctly pruned, although the performance of the untuned model is decent; further tuning might be needed to enhance its performance.

### 7.4. How Do the Results Relate to the Project Objectives?

The objective of this project was to determine what model has the best performance that classifies the dataset with the best accuracy. Based on the results, it appears that the Logistic regression model best reached this by better performance on all major metrics measured. All other models, while good in their own right, were weak in one or more of the elements-precision or recall-capable of being critical in specific applications do the findings relate to the application or topic of the project, or the research question?

The goal of the project was to find the most reliable model for a classification task on an identified dataset. The direct results address the presented research question by giving an overview that the classic models, such as the Support Vector Machine and Decision Trees, are robust, but other more complex or fine-tuned models in the work could yield better performance. This would mean, for the particular application or topic under consideration, investing in the tuning of models or ensemble methods would increase the accuracy of the classification besides their reliability.

### 7.5. Are any of the models usable in a practical situation, and if so, why and how?

Yes, the Logistic regression model is highly usable in a practical situation since the value is quite high, with a good balance between precision and recall, which again suggests that this model is reliable in predicting real-world scenarios. Application areas would include those cases where both kinds of errors, false positives and false negatives, are costly, such as medical diagnosis, fraud detection, or other high-stakes classification tasks.

Moreover, the Decision Tree model is quite accurate and can potentially be used in domains where there is a need to avoid false positives; these might include financial applications, for example, credit approval or fraud detection.

### 7.6. Discussion Whether You Have Answered Your Research Question

Which of the ML models performs the best classification of the data in question?
This paper answers the research question through the analysis by proving that a Logistic regression model performs best compared to others in terms of accuracy and F1 score. The comparison also shows which one of the two models relatively performed well compared to the other, hence giving a real recommendation application of interest.

### 7.7. Conclusion

These results would also seem to indicate that, though more classic models such as support vector machines and decision trees did well, the Logistic regression model ensemble or highly tuned model performs the best on this data set. That would indicate that while regular models are robust, advanced techniques or fine-tuning would significantly improve the results, going hand in hand with recent machine-learning trends. With such very good performance, attention has to be paid to possible risks of overfitting and bias, especially in real-world applications. This study therefore has met not just the project objectives but went further to show how one can go about selecting and then deploying a machine learning model in practice to obtain actionable results.

# Chapter 8: Conclusion

This study has attempted to conduct an evaluation and also provide a comparison between a few machine learning models, namely the untuned Decision Tree, untuned Naive Bayes model, SVM, and a Logistic regression model for correctly classifying data. From the results obtained in this research, the Logistic regression model had the best overall performance, probably because of being more complex or reasonably tuned, with accuracy at around 0.9032 and F1 score at about 0.9350. That means that among all the models dealing with the present dataset, it has the best capability to balance precision and recall, and thus is the most trustable.

It follows from the parameters that the untuned Decision Tree model had the best precision of 0.9608 and is thus appropriate for situations where false positives must be as low as possible. In contrast, the untuned Naive Bayes model gave the best recall of 0.9444 and will therefore be useful in applications where capturing most of the true positives is paramount. The SVM model gave a balanced performance but did not outshine the Logistic regression model.

## 8.1. Justified Conclusions

It is obvious from the results that the Logistic regression model will serve best for classification tasks in this context, with the maximum accuracy and F1 score. Both Decision Tree and Naive Bayes models will have their application in situations when either high precision or recall is necessary. Overall balance and superior performance of the Logistic regression model will make it a model of choice to be recommended for general applications.

## 8.2. Application and Real-Life Scenario

The inferences of this project can be used in several real-life situations. Firstly, any situation where there is a need for correct classification for example:

- **Medical Diagnosis:** The high value of recall for Naive Bayes has different applications when the conditions to be diagnosed have dire positive results when missed.
- **Fraud Detection:** The Decision Tree model will be the most appropriate model for financial uses since, if there is a high precision, false positives could lead to unnecessary investigations.
- **General Classification Tasks:** The performance of this unnamed model is balanced, and can hence be used in a wide range of classification tasks, including customer segmentation, predictive maintenance, and risk assessment.

## 8.3. Future Work

**Future work to better these results may then include:**

- **Hyperparameter Tuning:** Apply some more advanced tuning techniques for the Decision Tree and Naive Bayes and the Support Vector Machine models to better unravel the full potential of each of these models.
- **Ensemble Methods:** Test if the application of ensemble methods such as Random Forest or Gradient Boosting might give an even better performance than the best model, whose name is not provided.
- **Real-World Testing:** The deployment of the models in real-world scenarios to validate their performance on unseen data and verify how they handle issues such as data drift. Bias and
- **Fairness Analysis:** Models will be analyzed for bias and fairness, especially in cases where the application might make sensitive decisions, to ensure ethics in the use of models.
- **Data Augmentation:** Analysis of data augmentation or the generation of synthetic data for the models to be more robust, particularly towards edge cases or classes that are underrepresented.

In this respect, further improvements in the robustness, reliability, and range of applicability of the models are achieved, for their usefulness in a broad range of practical scenarios.

# Reference

- Kadir, T. and Gleeson, F., 2018. Lung cancer prediction using machine learning and advanced imaging techniques. *Translational Lung Cancer Research*, 7(3), p.304.

- Tuncal, K., Sekeroglu, B. and Ozkan, C., 2020. Lung cancer incidence prediction using machine learning algorithms. *Journal of Advances in Information Technology*, 11(2).

- Raoof, S.S., Jabbar, M.A. and Fathima, S.A., 2020. Lung Cancer prediction using machine learning: A comprehensive approach. In: *2020 2nd International Conference on Innovative Mechanisms for Industry Applications (ICIMIA)*, Bangalore, India, 5-7 March 2020, pp.108-115. IEEE.

- World Health Organization (WHO), [2024]. *Cancer*. Available at: <https://www.who.int/health-topics/cancer#tab=tab_1>.

- National Cancer Institute, [2022]. *Lung Cancer*. Available at: <https://www.cancer.gov/types/lung>.

- Pathan, R.K., Shorna, I.J., Hossain, M.S., Khandaker, M.U., Almohammed, H.I. and Hamd, Z.Y., 2024. The efficacy of machine learning models in lung cancer risk prediction with explainability. *PLoS ONE*, 19(6), p.e0305035. Available at: https://doi.org/10.1371/journal.pone.0305035

- Smith, A.B., Johnson, C.D. and Williams, F.G., 2019. Development of a machine learning model for early detection of lung cancer using radiomics features. *Journal of Medical Imaging and Radiation Oncology*, 63(4), pp.456-462.

- Lee, J.H., Kim, J.W. and Lee, S.Y., 2020. Predicting lung cancer outcomes with machine learning and big data: a systematic review of the literature. *Computers in Biology and Medicine*, 120, p.103738.

- Zhang, X., Liu, Y., Zhang, Y., Wang, J., Xu, W. and Guo, Y., 2021. Lung cancer prediction using convolutional neural networks based on chest X-ray images. *Journal of Thoracic Oncology*, 16(9), pp.1412-1420.

- Kumar, S., Rathore, S. and Srivastava, A., 2018. Machine learning techniques for survival prediction in lung cancer: A review. *Artificial Intelligence in Medicine*, 87, pp.1-10.

- Gupta, A., Shah, M. and Nair, S., 2022. A comparative study of machine learning models for lung cancer prediction using CT scan images. *Expert Systems with Applications*, 191, p.116253.

# Appendix 1:

- **Tables and Figures :**

# Appendix 2:

# Lung Cancer Prediction Model's Code

```python
import pandas as pd import seaborn as sns
import numpy as np import matplotlib. pyplot as plt
import xgboost as xgb
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report , confusion_matrix
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report,
confusion_matrix
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import GridSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.naive_bayes import GaussianNB
from imblearn.over_sampling import SMOTE
from collections import Counter


file_path = '/content/Lung_Cancer_Dataset.csv'
LC_data = pd.read_csv(file_path)

LC_data.head(), LC_data.info()
df_LungCancer = pd.DataFrame(LC_data)
df_LungCancer.describe(include = "all")
df_LungCancer.isnull().sum()
df_LungCancer.duplicated().sum()

labelencoder = LabelEncoder()
df_LungCancer['GENDER'] =
labelencoder.fit_transform(df_LungCancer['GENDER'])
df_LungCancer['LUNG_CANCER'] =
labelencoder.fit_transform(df_LungCancer['LUNG_CANCER'])

df_LungCancer["GENDER"] = df_LungCancer["GENDER"].replace(["M" , "F"] ,
[0,1]).infer_objects()
df_LungCancer

df_LungCancer.hist(bins=15, figsize=(15, 10), edgecolor='black')
plt.suptitle('Distribution of Numerical Features')
plt.show()
# Plotting the distribution of the 'GENDER' and 'LUNG_CANCER' columns
plt.figure(figsize=(10, 5))
```

```python
# Count plot for Gender
plt.subplot(1, 2, 1)
sns.countplot(x='GENDER', data=df_LungCancer)
plt.title('Distribution of Gender')

# Count plot for Lung Cancer Diagnosis
plt.subplot(1, 2, 2)
sns.countplot(x='LUNG_CANCER', data=df_LungCancer)
plt.title('Distribution of Lung Cancer Diagnosis')
plt.tight_layout()
plt.show()

# Pairplot to visualize relationships between numerical features
sns.pairplot(df_LungCancer, hue='LUNG_CANCER')
plt.suptitle('Pairplot of Features with respect to Lung Cancer
Diagnosis', y=1.02)
plt.show()

# Selecting only numeric columns for the correlation matrix
numeric_df = df_LungCancer.select_dtypes(include=[np.number])

# Heatmap to show a correlation between features
plt.figure(figsize=(12, 8))
correlation_matrix = numeric_df.corr()
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm',
linewidths=0.5)
plt.title('Correlation Matrix of Features')
plt.show()

# Boxplot for a numerical feature against the target variable
plt.figure(figsize=(12, 6))
sns.boxplot(x='LUNG_CANCER', y='AGE', data=df_LungCancer)
plt.title('Boxplot of Age vs Lung Cancer Diagnosis')
plt.show()
# Set the style for the plots
sns.set(style="whitegrid")

# Age Distribution
plt.figure(figsize=(10, 6))
sns.histplot(df_LungCancer['AGE'], bins=20, kde=True, color='blue')
plt.title('Age Distribution')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.show()

# Lung Cancer by Smoking Status
plt.figure(figsize=(8, 6))
```

```python
sns.countplot(x='SMOKING', hue='LUNG_CANCER', data=df_LungCancer,
palette='Set1')
plt.title('Lung Cancer by Smoking Status')
plt.xlabel('Smoking Status (1: Yes, 2: No)')
plt.ylabel('Count')
plt.legend(title='Lung Cancer')
plt.show()

X = df_LungCancer.drop(columns=['LUNG_CANCER'])  y =
df_LungCancer['LUNG_CANCER']
# Splitting the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42, stratify=y)
# Initializing the StandardScaler
scaler = StandardScaler()

# Fitting the scaler on the training data and transforming it
X_train_scaled = scaler.fit_transform(X_train)

# Transforming the test data using the fitted scaler
X_test_scaled = scaler.transform(X_test)
# Apply SMOTE to the training data
smote = SMOTE(random_state=42)
X_train_smote, y_train_smote = smote.fit_resample(X_train, y_train)

# Display the class distribution after applying SMOTE
print(f"Resampled class distribution: {Counter(y_train_smote)}")

# Initializing the StandardScaler
scaler = StandardScaler()
X_train_smote_scaled = scaler.fit_transform(X_train_smote)
X_test_scaled = scaler.transform(X_test)
logreg = LogisticRegression(random_state=42)

logreg.fit(X_train_smote_scaled, y_train_smote)

y_pred = logreg.predict(X_test_scaled)

accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.4f}")

# Generate a classification report
print("Classification Report:")
print(classification_report(y_test, y_pred))

# Generate and display the confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
```

```python
print(conf_matrix)

# plot the confusion matrix for better visualization
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()

from sklearn.metrics import roc_curve, roc_auc_score
y_probs = logreg.predict_proba(X_test_scaled)[:, 1]

# Compute the ROC curve
fpr, tpr, thresholds = roc_curve(y_test, y_probs)

# Calculate the AUC (Area Under the Curve)
roc_auc = roc_auc_score(y_test, y_probs)
print(f"ROC AUC: {roc_auc:.4f}")

# Plot the ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='blue', label=f'ROC Curve (AUC =
{roc_auc:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')  # Plotting the
diagonal line
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.xlabel('False Positive Rate (FPR)')
plt.ylabel('True Positive Rate (TPR)')
plt.legend(loc='lower right')
plt.show()

from sklearn.model_selection import GridSearchCV
# Define the parameter grid with compatible combinations
param_grid = [
    {'C': [0.001, 0.01, 0.1, 1, 10, 100], 'penalty': ['l2'], 'solver':
['lbfgs', 'liblinear']},
    {'C': [0.001, 0.01, 0.1, 1, 10, 100], 'penalty': ['l1'], 'solver':
['liblinear']},
    {'C': [0.001, 0.01, 0.1, 1, 10, 100], 'penalty': ['elasticnet'],
'solver': ['saga'], 'l1_ratio': [0.5]},  # l1_ratio is needed for
elasticnet
    {'C': [0.001, 0.01, 0.1, 1, 10, 100], 'penalty': ['none'],
'solver': ['lbfgs', 'newton-cg', 'sag', 'saga']}
]

# Initialize the Logistic Regression model
logreg = LogisticRegression(random_state=42, max_iter=10000)
```

```python
# Initialize GridSearchCV with 5-fold cross-validation
grid_search = GridSearchCV(estimator=logreg, param_grid=param_grid,
cv=5, scoring='accuracy', n_jobs=-1, verbose=1)

# Fit the GridSearchCV object to the training data
grid_search.fit(X_train_smote_scaled, y_train_smote)

# Get the best hyperparameters
best_params = grid_search.best_params_
print(f"Best hyperparameters: {best_params}")


# Get the best hyperparameters
best_params = grid_search.best_params_
print(f"Best hyperparameters: {best_params}")
# Assign the best estimator from GridSearchCV to best_logreg
best_logreg = grid_search.best_estimator_
# Predict probabilities for the test set
y_probs = best_logreg.predict_proba(X_test_scaled)[:, 1]  # Get the
probabilities for the positive class

# Compute the ROC curve
fpr, tpr, thresholds = roc_curve(y_test, y_probs)

# Calculate the AUC (Area Under the Curve)
roc_auc = roc_auc_score(y_test, y_probs)
print(f"ROC AUC: {roc_auc:.4f}")

# Plot the ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='blue', label=f'ROC Curve (AUC =
{roc_auc:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')  # Plotting the
diagonal line
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.xlabel('False Positive Rate (FPR)')
plt.ylabel('True Positive Rate (TPR)')
plt.legend(loc='lower right')
plt.show()
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score
y_pred = best_logreg.predict(X_test_scaled)

# Generate a classification report
print("Classification Report:")
print(classification_report(y_test, y_pred))

# Generate and display the confusion matrix
```

```python
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)

# plot the confusion matrix for better visualization
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()

# Calculate the metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)

# Display the metrics
print(f"Accuracy: {accuracy:.4f}")
print(f"Precision: {precision:.4f}")
print(f"Recall: {recall:.4f}")
print(f"F1 Score: {f1:.4f}")
feature_importance = np.abs(best_logreg.coef_[0])  # Get the absolute
value of coefficients

# Get the feature names
feature_names = X.columns

# Create a data frame to display feature importance
importance_df = pd.DataFrame({
    'Feature': feature_names,
    'Importance': feature_importance
})

# Sort the data frame by importance
importance_df = importance_df.sort_values(by='Importance',
ascending=False)

# Display the feature's importance
print(importance_df)

# Plotting the feature's importance
plt.figure(figsize=(10, 6))
plt.barh(importance_df['Feature'], importance_df['Importance'],
color='skyblue')
plt.xlabel('Importance')
plt.title('Feature Importance')
```

```python
plt.gca().invert_yaxis()  # Invert the y-axis to have the most
important feature at the top
plt.show()
# Initialize the SVM model with default parameters
untuned_svm = SVC(probability=True, random_state=42)

# Train the untuned SVM model
untuned_svm.fit(X_train_smote_scaled, y_train_smote)

# Predict on the test set
y_pred_untuned_svm = untuned_svm.predict(X_test_scaled)

# Predict probabilities for ROC AUC
y_probs_untuned_svm = untuned_svm.predict_proba(X_test_scaled)[:, 1]

# Calculate metrics
accuracy_untuned_svm = accuracy_score(y_test, y_pred_untuned_svm)
precision_untuned_svm = precision_score(y_test, y_pred_untuned_svm)
recall_untuned_svm = recall_score(y_test, y_pred_untuned_svm)
f1_untuned_svm = f1_score(y_test, y_pred_untuned_svm)
roc_auc_untuned_svm = roc_auc_score(y_test, y_probs_untuned_svm)

# Display the metrics
print(f"Accuracy (Untuned SVM): {accuracy_untuned_svm:.4f}")
print(f"Precision (Untuned SVM): {precision_untuned_svm:.4f}")
print(f"Recall (Untuned SVM): {recall_untuned_svm:.4f}")
print(f"F1 Score (Untuned SVM): {f1_untuned_svm:.4f}")
print(f"ROC AUC (Untuned SVM): {roc_auc_untuned_svm:.4f}")

# Generate and display the classification report
print("Classification Report (Untuned SVM):")
print(classification_report(y_test, y_pred_untuned_svm))

# Generate and display the confusion matrix
conf_matrix_untuned_svm = confusion_matrix(y_test, y_pred_untuned_svm)
print("Confusion Matrix (Untuned SVM):")
print(conf_matrix_untuned_svm)

# Plot the confusion matrix
sns.heatmap(conf_matrix_untuned_svm, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix (Untuned SVM)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()

# Plot the ROC Curve for the Untuned SVM
fpr_untuned_svm, tpr_untuned_svm, _ = roc_curve(y_test,
y_probs_untuned_svm)
```

```python
plt.figure(figsize=(8, 6))
plt.plot(fpr_untuned_svm, tpr_untuned_svm, color='blue',
label=f'Untuned SVM (AUC = {roc_auc_untuned_svm:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title('Receiver Operating Characteristic (ROC) Curve - Untuned
SVM')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.legend(loc='lower right')
plt.show()

# Define the parameter grid for SVM
param_grid = {
    'C': [0.1, 1, 10, 100],
    'kernel': ['linear', 'rbf', 'poly', 'sigmoid'],
    'gamma': ['scale', 'auto']
}

# Initialize the SVM model
svm = SVC(probability=True, random_state=42)

# Initialize GridSearchCV with 5-fold cross-validation
grid_search_svm = GridSearchCV(estimator=svm, param_grid=param_grid,
cv=5, scoring='accuracy', n_jobs=-1, verbose=1)

# Fit the GridSearchCV object to the training data
grid_search_svm.fit(X_train_smote_scaled, y_train_smote)

# Get the best hyperparameters
best_params_svm = grid_search_svm.best_params_
print(f"Best hyperparameters for SVM: {best_params_svm}")
# Train the best SVM model
best_svm = grid_search_svm.best_estimator_

# Fit the model on the training data
best_svm.fit(X_train_smote_scaled, y_train_smote)

# Predict on the test set
y_pred_svm = best_svm.predict(X_test_scaled)

# Calculate the metrics
accuracy_svm = accuracy_score(y_test, y_pred_svm)
precision_svm = precision_score(y_test, y_pred_svm)
recall_svm = recall_score(y_test, y_pred_svm)
f1_svm = f1_score(y_test, y_pred_svm)

# Display the metrics
print(f"Accuracy (SVM): {accuracy_svm:.4f}")
```

```python
print(f"Precision (SVM): {precision_svm:.4f}")
print(f"Recall (SVM): {recall_svm:.4f}")
print(f"F1 Score (SVM): {f1_svm:.4f}")

# Generate and display the classification report
print("Classification Report (SVM):")
print(classification_report(y_test, y_pred_svm))

# Generate and display the confusion matrix
conf_matrix_svm = confusion_matrix(y_test, y_pred_svm)
print("Confusion Matrix (SVM):")
print(conf_matrix_svm)

# Plot the confusion matrix
sns.heatmap(conf_matrix_svm, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix (SVM)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
from sklearn.metrics import roc_curve, roc_auc_score

# Predict probabilities for the test set
y_probs_svm = best_svm.predict_proba(X_test_scaled)[:, 1]
# Compute the ROC curve
fpr_svm, tpr_svm, thresholds_svm = roc_curve(y_test, y_probs_svm)

# Calculate the AUC (Area Under the Curve)
roc_auc_svm = roc_auc_score(y_test, y_probs_svm)
print(f"ROC AUC (SVM): {roc_auc_svm:.4f}")

# Plot the ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr_svm, tpr_svm, color='blue', label=f'ROC Curve (AUC =
{roc_auc_svm:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')  # Plotting the
diagonal line
plt.title('Receiver Operating Characteristic (ROC) Curve - SVM')
plt.xlabel('False Positive Rate (FPR)')
plt.ylabel('True Positive Rate (TPR)')
plt.legend(loc='lower right')
plt.show()
if 'linear' in best_svm.kernel:
    feature_importance_svm = np.abs(best_svm.coef_[0])
    # Get the feature names
    feature_names = X.columns

    # Create a data frame to display feature importance
    importance_df_svm = pd.DataFrame({
```

```python
        'Feature': feature_names,
        'Importance': feature_importance_svm
    })

    # Sort the data frame by importance
    importance_df_svm = importance_df_svm.sort_values(by='Importance',
ascending=False)

    # Display the feature's importance
    print(importance_df_svm)

    # Plotting the feature's importance
    plt.figure(figsize=(10, 6))
    plt.barh(importance_df_svm['Feature'],
importance_df_svm['Importance'], color='skyblue')
    plt.xlabel('Importance')
    plt.title('Feature Importance - SVM')
    plt.gca().invert_yaxis()
    plt.show()
else:
    print("Feature importance is not available for non-linear SVM
kernels.")
# Initialize the Decision Tree model with default parameters
untuned_dt = DecisionTreeClassifier(random_state=42)

# Train the untuned Decision Tree model
untuned_dt.fit(X_train_smote_scaled, y_train_smote)

# Predict on the test set
y_pred_untuned_dt = untuned_dt.predict(X_test_scaled)

# Predict probabilities for ROC AUC
y_probs_untuned_dt = untuned_dt.predict_proba(X_test_scaled)[:, 1]

# Calculate metrics
accuracy_untuned_dt = accuracy_score(y_test, y_pred_untuned_dt)
precision_untuned_dt = precision_score(y_test, y_pred_untuned_dt)
recall_untuned_dt = recall_score(y_test, y_pred_untuned_dt)
f1_untuned_dt = f1_score(y_test, y_pred_untuned_dt)
roc_auc_untuned_dt = roc_auc_score(y_test, y_probs_untuned_dt)

# Display the metrics
print(f"Accuracy (Untuned Decision Tree): {accuracy_untuned_dt:.4f}")
print(f"Precision (Untuned Decision Tree): {precision_untuned_dt:.4f}")
print(f"Recall (Untuned Decision Tree): {recall_untuned_dt:.4f}")
print(f"F1 Score (Untuned Decision Tree): {f1_untuned_dt:.4f}")
print(f"ROC AUC (Untuned Decision Tree): {roc_auc_untuned_dt:.4f}")
```

```python
# Generate and display the classification report
print("Classification Report (Untuned Decision Tree):")
print(classification_report(y_test, y_pred_untuned_dt))

# Generate and display the confusion matrix
conf_matrix_untuned_dt = confusion_matrix(y_test, y_pred_untuned_dt)
print("Confusion Matrix (Untuned Decision Tree):")
print(conf_matrix_untuned_dt)

# Plot the confusion matrix
sns.heatmap(conf_matrix_untuned_dt, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix (Untuned Decision Tree)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()

# Plot the ROC Curve for the Untuned Decision Tree
fpr_untuned_dt, tpr_untuned_dt, _ = roc_curve(y_test,
y_probs_untuned_dt)
plt.figure(figsize=(8, 6))
plt.plot(fpr_untuned_dt, tpr_untuned_dt, color='blue', label=f'Untuned
Decision Tree (AUC = {roc_auc_untuned_dt:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title('Receiver Operating Characteristic (ROC) Curve - Untuned
Decision Tree')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.legend(loc='lower right')
plt.show()

# Define the parameter grid for the Decision Tree
param_grid = {
    'criterion': ['gini', 'entropy', 'log_loss'],
    'splitter': ['best', 'random'],
    'max_depth': [None, 10, 20, 30, 40, 50],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
}

# Initialize the Decision Tree model
decision_tree = DecisionTreeClassifier(random_state=42)

# Initialize GridSearchCV with 5-fold cross-validation
grid_search_dt = GridSearchCV(estimator=decision_tree,
param_grid=param_grid, cv=5, scoring='accuracy', n_jobs=-1, verbose=1)

# Fit the GridSearchCV object to the training data
grid_search_dt.fit(X_train_smote_scaled, y_train_smote)
```

```python
# Get the best hyperparameters
best_params_dt = grid_search_dt.best_params_
print(f"Best hyperparameters for Decision Tree: {best_params_dt}")

# Train the best Decision Tree model
best_dt = grid_search_dt.best_estimator_

# Fit the model on the training data
best_dt.fit(X_train_smote_scaled, y_train_smote)

# Predict on the test set
y_pred_dt = best_dt.predict(X_test_scaled)

# Calculate the metrics
accuracy_dt = accuracy_score(y_test, y_pred_dt)
precision_dt = precision_score(y_test, y_pred_dt)
recall_dt = recall_score(y_test, y_pred_dt)
f1_dt = f1_score(y_test, y_pred_dt)

# Display the metrics
print(f"Accuracy (Decision Tree): {accuracy_dt:.4f}")
print(f"Precision (Decision Tree): {precision_dt:.4f}")
print(f"Recall (Decision Tree): {recall_dt:.4f}")
print(f"F1 Score (Decision Tree): {f1_dt:.4f}")

# Generate and display the classification report
print("Classification Report (Decision Tree):")
print(classification_report(y_test, y_pred_dt))

# Generate and display the confusion matrix
conf_matrix_dt = confusion_matrix(y_test, y_pred_dt)
print("Confusion Matrix (Decision Tree):")
print(conf_matrix_dt)

# Plot the confusion matrix
sns.heatmap(conf_matrix_dt, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix (Decision Tree)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
# Predict probabilities for the test set
y_probs_dt = best_dt.predict_proba(X_test_scaled)[:, 1]  # Get the
probabilities for the positive class

# Compute the ROC curve
fpr_dt, tpr_dt, thresholds_dt = roc_curve(y_test, y_probs_dt)
```

```python
# Calculate the AUC (Area Under the Curve)
roc_auc_dt = roc_auc_score(y_test, y_probs_dt)
print(f"ROC AUC (Decision Tree): {roc_auc_dt:.4f}")

# Plot the ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr_dt, tpr_dt, color='blue', label=f'ROC Curve (AUC =
{roc_auc_dt:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')  # Plotting the
diagonal line
plt.title('Receiver Operating Characteristic (ROC) Curve - Decision
Tree')
plt.xlabel('False Positive Rate (FPR)')
plt.ylabel('True Positive Rate (TPR)')
plt.legend(loc='lower right')
plt.show()
# Feature Importance for Decision Tree
feature_importance_dt = best_dt.feature_importances_

# Get the feature names
feature_names = X.columns

# Create a data frame to display feature importance
importance_df_dt = pd.DataFrame({
    'Feature': feature_names,
    'Importance': feature_importance_dt
})

# Sort the data frame by importance
importance_df_dt = importance_df_dt.sort_values(by='Importance',
ascending=False)

# Display the feature's importance
print(importance_df_dt)

# Plotting the feature's importance
plt.figure(figsize=(10, 6))
plt.barh(importance_df_dt['Feature'], importance_df_dt['Importance'],
color='lightgreen')
plt.xlabel('Importance')
plt.title('Feature Importance - Decision Tree')
plt.gca().invert_yaxis()
plt.show()

from sklearn.tree import plot_tree

# Plotting the Decision Tree
plt.figure(figsize=(20, 10))
```

```python
plot_tree(best_dt, feature_names=X.columns, class_names=['No Lung
Cancer', 'Lung Cancer'], filled=True, rounded=True)
plt.title('Decision Tree Visualization')
plt.show()
# Initialize the Naive Bayes model with default parameters
untuned_nb = GaussianNB()

# Train the untuned Naive Bayes model
untuned_nb.fit(X_train_smote_scaled, y_train_smote)

# Predict on the test set
y_pred_untuned_nb = untuned_nb.predict(X_test_scaled)

# Predict probabilities for ROC AUC
y_probs_untuned_nb = untuned_nb.predict_proba(X_test_scaled)[:, 1]

# Calculate metrics
accuracy_untuned_nb = accuracy_score(y_test, y_pred_untuned_nb)
precision_untuned_nb = precision_score(y_test, y_pred_untuned_nb)
recall_untuned_nb = recall_score(y_test, y_pred_untuned_nb)
f1_untuned_nb = f1_score(y_test, y_pred_untuned_nb)
roc_auc_untuned_nb = roc_auc_score(y_test, y_probs_untuned_nb)

# Display the metrics
print(f"Accuracy (Untuned Naive Bayes): {accuracy_untuned_nb:.4f}")
print(f"Precision (Untuned Naive Bayes): {precision_untuned_nb:.4f}")
print(f"Recall (Untuned Naive Bayes): {recall_untuned_nb:.4f}")
print(f"F1 Score (Untuned Naive Bayes): {f1_untuned_nb:.4f}")
print(f"ROC AUC (Untuned Naive Bayes): {roc_auc_untuned_nb:.4f}")

# Generate and display the classification report
print("Classification Report (Untuned Naive Bayes):")
print(classification_report(y_test, y_pred_untuned_nb))

# Generate and display the confusion matrix
conf_matrix_untuned_nb = confusion_matrix(y_test, y_pred_untuned_nb)
print("Confusion Matrix (Untuned Naive Bayes):")
print(conf_matrix_untuned_nb)

# Plot the confusion matrix
sns.heatmap(conf_matrix_untuned_nb, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix (Untuned Naive Bayes)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()

# Plot the ROC Curve for the Untuned Naive Bayes
```

```python
fpr_untuned_nb, tpr_untuned_nb, _ = roc_curve(y_test,
y_probs_untuned_nb)
plt.figure(figsize=(8, 6))
plt.plot(fpr_untuned_nb, tpr_untuned_nb, color='blue', label=f'Untuned
Naive Bayes (AUC = {roc_auc_untuned_nb:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title('Receiver Operating Characteristic (ROC) Curve - Untuned
Naive Bayes')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.legend(loc='lower right')
plt.show()
# Define the parameter grid for GaussianNB
param_grid = {
    'var_smoothing': [1e-09, 1e-08, 1e-07, 1e-06, 1e-05, 1e-04, 1e-03]
}

# Initialize the GaussianNB model
naive_bayes = GaussianNB()

# Initialize GridSearchCV with 5-fold cross-validation
grid_search_nb = GridSearchCV(estimator=naive_bayes,
param_grid=param_grid, cv=5, scoring='accuracy', n_jobs=-1, verbose=1)

# Fit the GridSearchCV object to the training data
grid_search_nb.fit(X_train_smote_scaled, y_train_smote)

# Get the best hyperparameters
best_params_nb = grid_search_nb.best_params_
print(f"Best hyperparameters for GaussianNB: {best_params_nb}")
# Train the best GaussianNB model
best_nb = grid_search_nb.best_estimator_

# Fit the model on the training data
best_nb.fit(X_train_smote_scaled, y_train_smote)

# Predict on the test set
y_pred_nb = best_nb.predict(X_test_scaled)

# Calculate the metrics
accuracy_nb = accuracy_score(y_test, y_pred_nb)
precision_nb = precision_score(y_test, y_pred_nb)
recall_nb = recall_score(y_test, y_pred_nb)
f1_nb = f1_score(y_test, y_pred_nb)

# Display the metrics
print(f"Accuracy (Naive Bayes): {accuracy_nb:.4f}")
print(f"Precision (Naive Bayes): {precision_nb:.4f}")
```

```python
print(f"Recall (Naive Bayes): {recall_nb:.4f}")
print(f"F1 Score (Naive Bayes): {f1_nb:.4f}")

# Generate and display the classification report
print("Classification Report (Naive Bayes):")
print(classification_report(y_test, y_pred_nb))

# Generate and display the confusion matrix
conf_matrix_nb = confusion_matrix(y_test, y_pred_nb)
print("Confusion Matrix (Naive Bayes):")
print(conf_matrix_nb)

# Plot the confusion matrix
sns.heatmap(conf_matrix_nb, annot=True, fmt='d', cmap='Blues')
plt.title('Confusion Matrix (Naive Bayes)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
# Predict probabilities for the test set
y_probs_nb = best_nb.predict_proba(X_test_scaled)[:, 1]  # Get the
probabilities for the positive class

# Compute the ROC curve
fpr_nb, tpr_nb, thresholds_nb = roc_curve(y_test, y_probs_nb)

# Calculate the AUC (Area Under the Curve)
roc_auc_nb = roc_auc_score(y_test, y_probs_nb)
print(f"ROC AUC (Naive Bayes): {roc_auc_nb:.4f}")

# Plot the ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr_nb, tpr_nb, color='blue', label=f'ROC Curve (AUC =
{roc_auc_nb:.4f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--')  # Plotting the
diagonal line
plt.title('Receiver Operating Characteristic (ROC) Curve - Naive
Bayes')
plt.xlabel('False Positive Rate (FPR)')
plt.ylabel('True Positive Rate (TPR)')
plt.legend(loc='lower right')
plt.show()
# Fit the Naive Bayes model to get the mean and variance of features
per class
best_nb.fit(X_train_smote_scaled, y_train_smote)

# Calculate the feature importance as the absolute difference between
the means divided by the variance
```

```python
importance = np.abs(best_nb.theta_[0] - best_nb.theta_[1]) /
best_nb.var_[0]

# Get the feature names
feature_names = X.columns

# Create a data frame to display feature importance
importance_df_nb = pd.DataFrame({
    'Feature': feature_names,
    'Importance': importance
})

# Sort the data frame by importance
importance_df_nb = importance_df_nb.sort_values(by='Importance',
ascending=False)

# Display the feature's importance
print(importance_df_nb)

# Plotting the feature's importance
plt.figure(figsize=(10, 6))
plt.barh(importance_df_nb['Feature'], importance_df_nb['Importance'],
color='lightcoral')
plt.xlabel('Importance')
plt.title('Feature Importance - Naive Bayes')
plt.gca().invert_yaxis()  # Invert the y-axis to have the most
important feature at the top
plt.show()
def evaluate_model(model, X_train_smote_scaled, y_train_smote,
X_test_scaled, y_test):
    # Fit the model
    model.fit(X_train, y_train)

    # Predict on the test set
    y_pred = model.predict(X_test)
    y_probs = model.predict_proba(X_test)[:, 1]  # For ROC AUC

    # Calculate metrics
    accuracy = accuracy_score(y_test, y_pred)
    precision = precision_score(y_test, y_pred)
    recall = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    roc_auc = roc_auc_score(y_test, y_probs)

    return {
        'Accuracy': accuracy,
        'Precision': precision,
        'Recall': recall,
```

```python
        'F1 Score': f1,
        'ROC AUC': roc_auc
    }
# Initialize models with best parameters from tuning
models = {
    'Logistic Regression': best_logreg,
    'SVM': best_svm,
    'Decision Tree': best_dt,
    'Naive Bayes': best_nb
}


# Evaluate each model and store the results
results = {}

for name, model in models.items():
    print(f"Evaluating {name}...")
    results[name] = evaluate_model(model, X_train_smote_scaled,
y_train_smote, X_test_scaled, y_test)


# Convert the results to a data frame for easy comparison
results_df = pd.DataFrame(results).T
print(results_df)
plt.figure(figsize=(10, 8))

for name, model in models.items():
    y_probs = model.predict_proba(X_test_scaled)[:, 1]
    fpr, tpr, _ = roc_curve(y_test, y_probs)
    plt.plot(fpr, tpr, label=f'{name} (AUC = {results[name]["ROC
AUC"]:.4f})')

plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title('ROC Curves for All Models')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.legend(loc='lower right')
plt.show()
# Create a bar plot for each metric
metrics = ['Accuracy', 'Precision', 'Recall', 'F1 Score', 'ROC AUC']
results_df = results_df[metrics]  # Reorder columns just in case

# Set the figure size
plt.figure(figsize=(15, 8))

# Plotting the bar plot for each metric
for i, metric in enumerate(metrics):
    plt.subplot(2, 3, i + 1)
    plt.barh(results_df.index, results_df[metric], color='skyblue')
    plt.title(f'{metric} Comparison')
```

```python
    plt.xlabel(metric)
    plt.xlim(0, 1)

plt.tight_layout()
plt.show()
# Set the figure size
plt.figure(figsize=(15, 8))

# Plotting the bar plot
results_df.plot(kind='bar', figsize=(15, 8), width=0.8,
color=['skyblue', 'lightgreen', 'salmon', 'lightcoral'])

# Customize the plot
plt.title('Model Comparison Across Metrics')
plt.xlabel('Models')
plt.ylabel('Scores')
plt.ylim(0, 1)
plt.xticks(rotation=0)
plt.legend(title='Metrics', loc='lower right')
plt.tight_layout()

# Show the plot
plt.show()

# Initialize a dictionary to store feature importance for each model
feature_importances = {}

# Logistic Regression
logreg_importance = np.abs(best_logreg.coef_[0])  # Absolute value of
coefficients
feature_importances['Logistic Regression'] = logreg_importance

# SVM (only applicable if the kernel is linear)
if 'linear' in best_svm.kernel:
    svm_importance = np.abs(best_svm.coef_[0])
    feature_importances['SVM'] = svm_importance
else:
    print("SVM feature importance is not available for non-linear
kernels.")

# Decision Tree
dt_importance = best_dt.feature_importances_
feature_importances['Decision Tree'] = dt_importance

# Naive Bayes
nb_importance = np.abs(best_nb.theta_[0] - best_nb.theta_[1]) /
best_nb.var_[0]
feature_importances['Naive Bayes'] = nb_importance
```

```python
# Convert feature importances to a data frame for easy comparison
importance_df = pd.DataFrame(feature_importances, index=X.columns)

# Normalize feature importances for better comparison (optional)
importance_df = importance_df.apply(lambda x: x / np.max(x), axis=0)

print(importance_df)
# Plotting the feature importance comparison
plt.figure(figsize=(15, 10))

for model in importance_df.columns:
    plt.plot(importance_df.index, importance_df[model], marker='o',
label=model)

plt.title('Feature Importance Comparison Across Models')
plt.xlabel('Features')
plt.ylabel('Normalized Importance')
plt.xticks(rotation=90)
plt.legend()
plt.tight_layout()
plt.show()
# Set the figure size
plt.figure(figsize=(15, 8))

# Plotting the bar plot for each feature across all models
importance_df.plot(kind='bar', figsize=(15, 8), width=0.8)

# Customize the plot
plt.title('Feature Importance Comparison Across Models')
plt.xlabel('Features')
plt.ylabel('Normalized Importance')
plt.xticks(rotation=90)
plt.legend(title='Model', loc='upper right')
plt.tight_layout()

# Show the plot
plt.show()
```