



CSE422 : ARTIFICIAL INTELLIGENCE

Project Report

Project Title :Lung Cancer Prediction

Name	ID
Omor Bin Amjad Chowdhury	23241085
Asif Imtiaz Chowdhury	21241050
Rayhan Sharif Sadif	21101063

Submitted To:

Sumaiya Akter
Lecturer
Brac University

Swattic Ghose
Lecturer
Brac University

Table of Contents

Section No	Content	Page No
1	Introduction	4
2	Dataset Description	5
3	Dataset pre-processing	6
4	Feature scaling	8-9
5	Dataset splitting	9-10
6	Model Training and Testing	11-12
7	Model Selection	13-17
8	Conclusion	18

INTRODUCTION

Lung Cancer Prediction is an artificial intelligence (AI) project that aims to predict the risk of lung cancer in individuals based on various factors and datasets. The project's objectives, motivation, and problems it aims to solve can be summarized as follows:

Objectives: The project aims to develop an AI tool that can accurately predict the risk of lung cancer for individuals with or without a significant smoking history, based on analyses of low-dose computed tomography (LDCT) scans from patients in the U.S. and Taiwan. This tool can help improve early detection, diagnosis, and treatment of lung cancer.

Motivation: Lung cancer is the leading cause of cancer death in the United States and around the world. The development of an AI tool that can accurately predict lung cancer risk can help reduce the number of lung cancer deaths by identifying individuals at high risk and providing timely interventions.

Problems: The project faces several challenges, such as the heterogeneity of lung cancer and the need for large and diverse datasets to train and validate the AI models. Additionally, AI tools currently cannot replace physicians in clinical trials and have limitations in terms of decision-making.

Solutions: AI has the potential to help treat lung cancer in various aspects, including detection, diagnosis, decision-making, and prognosis prediction. By integrating and analyzing large datasets, AI models can strengthen various aspects of lung cancer therapy, such as early detection, auxiliary diagnosis, prognosis prediction, and immunotherapy practice.

In summary, our project aims to develop an AI tool that can accurately predict lung cancer risk, helping to improve early detection, diagnosis, and treatment. The project is motivated by the high mortality rate associated with lung cancer and aims to address the challenges posed by the heterogeneity of lung cancer and the need for large datasets.

Dataset Description:

Dataset Link: <https://www.kaggle.com/datasets/rishidamarla/cancer-patients-data>

Task Type: The project focuses on a regression task, specifically predicting Lung Cancer Prediction of each person.

Number of Data Points: 1000

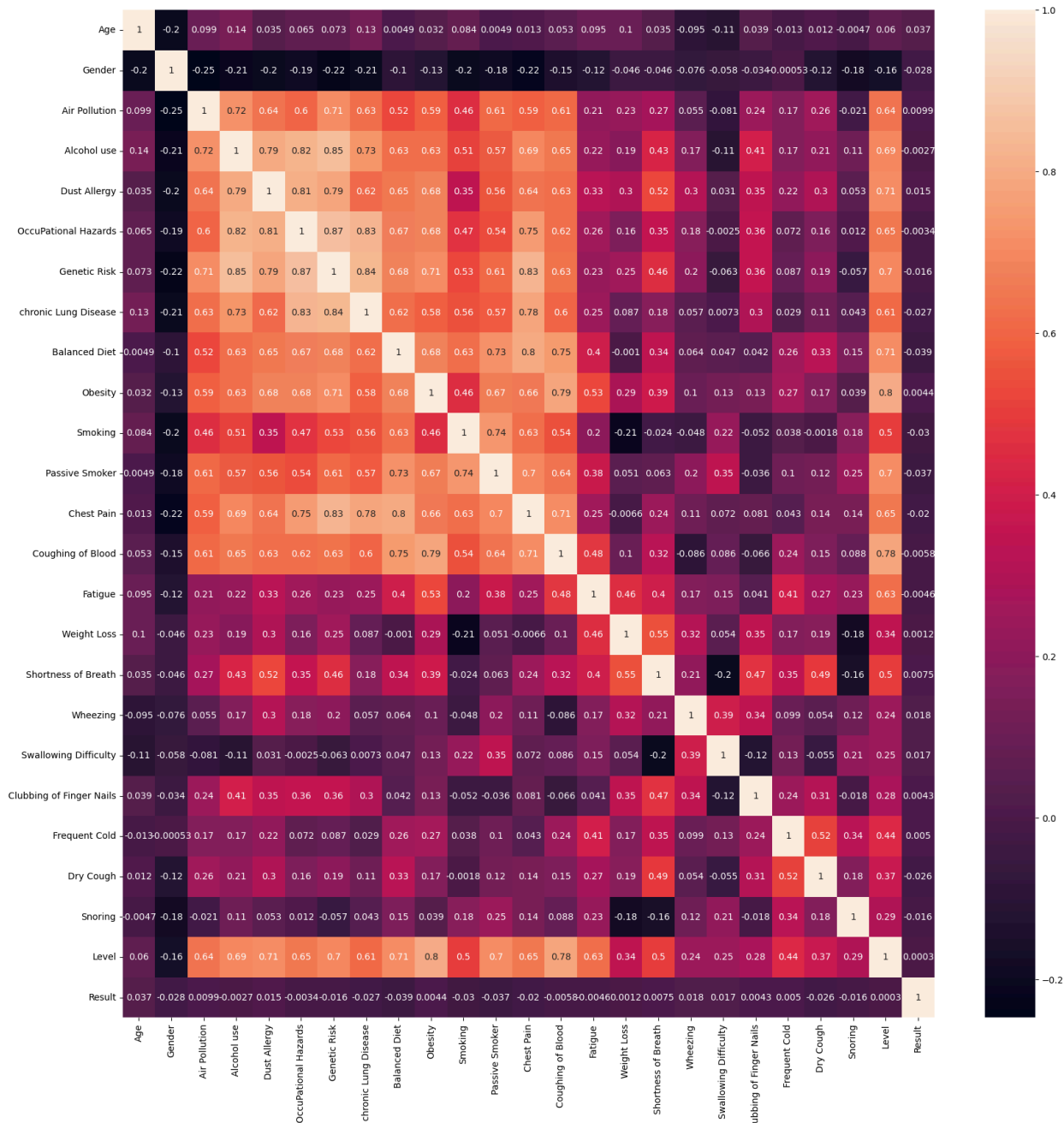
No of Features: 28 Features

Qualitative: Gender ,Level and Age

Quantitative: All others

Category: Classification problem as the risk level falls in one of three categories: Low, Medium or High.

Correlation of Features:



Dataset Pre-processing:

Several steps were taken during the preprocessing phase to prepare the dataset for effective machine learning modeling. These steps are essential for ensuring data quality and usability, as well as improving model performance. The following preprocessing procedures were carried out:

- Duplicate columns like "Smoking" and "Swallowing Difficulty" were removed.
- Columns like "Patient id" which do not contribute to the results were dropped.
- Missing values were handled through imputation using the "SimpleImputer" from `sklearn`. The mean strategy was employed.
- Categorical features like "Gender" and "Level" were encoded with appropriate values.
- Less informative features, such as "Genetic Risk" and "Occupational Hazard", were dropped using correlation analysis, visualized using Seaborn.

Feature Scaling:

Data Loading and Exploration:

It includes data exploration using functions like `info()`, `head()`, and visualization using Seaborn and Matplotlib for creating count plots, strip plots, histograms, and heatmap.

Data Preprocessing:

The dataset undergoes preprocessing steps, including handling missing values (`SimpleImputer`), mapping categorical values to numerical ones, and dropping unnecessary columns.

Data Scaling:

Min-Max scaling is applied to the features using `MinMaxScaler` from `scikit-learn`.

Model Training and Evaluation:

Machine learning models such as Decision Tree Classifier, Gaussian Naive Bayes, and Random Forest Classifier are trained and evaluated.

The models are trained both with and without feature scaling to observe the impact on performance.

Hyperparameter Tuning:

Grid Search Cross-Validation (`GridSearchCV`) is employed for hyperparameter tuning in the Support Vector Classifier (SVC).

Model Evaluation Metrics:

The code calculates and prints various evaluation metrics, including accuracy, confusion matrix, and classification report, to assess the models' performance.

Data Visualization:

Seaborn's heatmap is utilized to visualize confusion matrices, aiding in understanding the model's performance.

Machine Learning Models:

Decision Tree Classifier, Gaussian Naive Bayes, Random Forest Classifier, and Support Vector Classifier are used as machine learning models.

These features collectively contribute to the development, training, and evaluation of machine learning models for lung cancer prediction.

Dataset Splitting:

The splitting of the data into training and testing sets is an important step in preparing the dataset for machine learning modeling. This split enables the model to be evaluated on unseen data, providing an estimate of its performance on real-world data. The following are the steps and details for splitting the dataset in this project:

Defining Features and Target Variable:

The features (X) and the target variable (y) were defined. 'X' includes all the columns except for the target variable

Splitting Ratio:

The dataset was divided into training (70%) and test (30%) sets, with two variants for each: one scaled using "MinMaxScaler" and the other unscaled. The data to be split was selected randomly.. A 70-30 split is commonly used in machine learning as it provides a substantial amount of data for learning while retaining enough data to test and validate the models effectively.

Random Split with Fixed Random State:

The `train_test_split` function from scikit-learn was used for this purpose. A `random_state` was set to ensure reproducibility of the results.

Result of the Split:

The training set consisted of samples (features and target), and the testing set comprised samples.

This methodical approach to dataset splitting ensures that the models are trained on a diverse set of data points and are tested on different samples to gauge their generalization ability effectively.

Model Training and Evaluation:

The project included training and testing three distinct machine learning models: decision trees, naive bayes, random forest, support vector machine and grid search. Each model was chosen for its distinct characteristics and ability to address regression issues. The following explains the training and evaluation process for each model:

Decision Tree:

Training: Decision Trees are trained by recursively partitioning the dataset based on feature values. The algorithm selects the best feature at each node to split the data, aiming to maximize information gain or Gini impurity reduction.

Testing: During testing, new data traverses the tree, following the learned decision rules at each node until it reaches a leaf. The predicted class at the leaf is assigned to the input data.

Naive Bayes:

Training: Naive Bayes is a probabilistic model based on Bayes' theorem. It assumes that features are conditionally independent given the class. The model calculates class probabilities and conditional probabilities of features given the class during training.

Testing: In testing, the model uses Bayes' theorem to calculate the probability of each class given the input features. The class with the highest probability is assigned to the input data.

Random Forest:

Training: Random Forest is an ensemble model that builds multiple decision trees during training. Each tree is trained on a random subset of the data and a random subset of features to promote diversity. The final prediction is made by averaging or voting among the individual trees.

Testing: During testing, the input data passes through each tree in the forest, and the final prediction is determined by aggregating the predictions of all trees.

Support Vector Machine (SVM):

Training: SVM aims to find a hyperplane that best separates data points of different classes. During training, it identifies the optimal hyperplane by maximizing the margin between classes. SVM can handle linear and non-linear separation using different kernels.

Testing: In testing, the input data are mapped into the feature space, and their position relative to the hyperplane is used to determine the class label.

Grid Search:

Training and Testing: Grid Search is not a model but a hyperparameter tuning technique. It systematically searches through a predefined hyperparameter grid to find the combination that yields the best model performance. It involves training and testing the model with various hyperparameter values, allowing for optimization.

Evaluation Process:

Each model was trained on the training set and then used to make predictions on the test set.

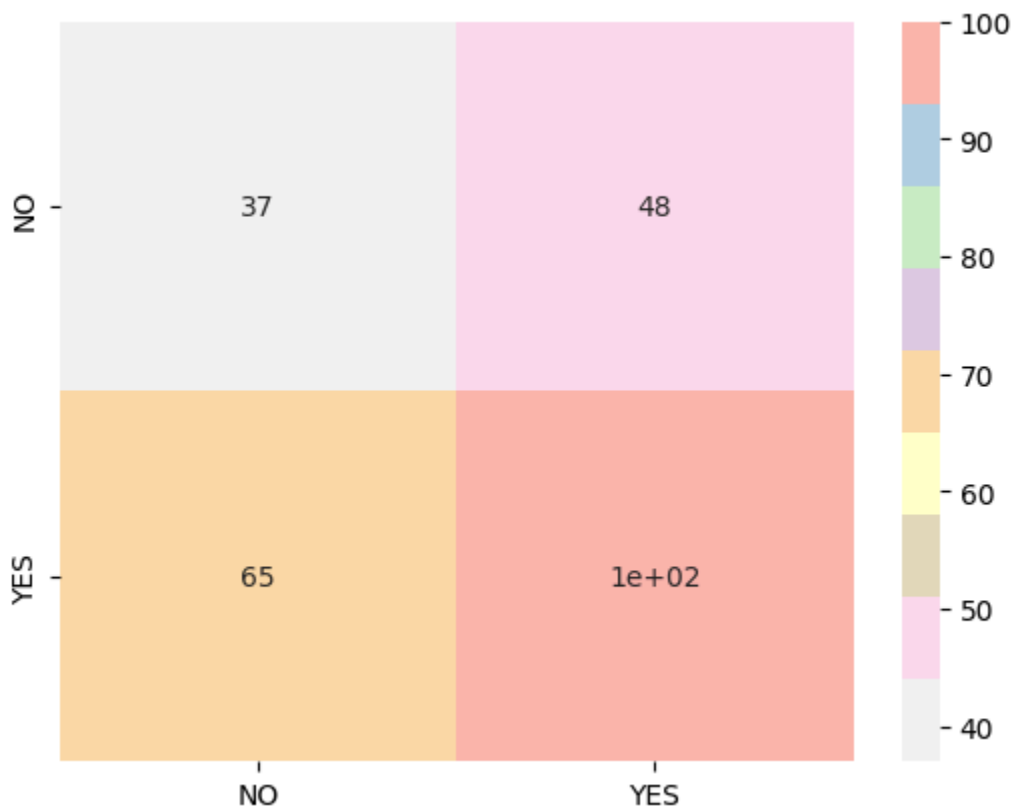
The evaluation focused on training and testing before and after scaling. This allows for an in-depth comparison of how each model performs in predicting and understanding the strengths and weaknesses of each approach in the context of the given dataset.

Model Selection and Comparison Analysis

The project involved a thorough comparison and analysis of three machine learning models: Decision Tree, Naive Bayes, Random Forest. The initial performance of these models was evaluated based on Without Scaling Training and testing accuracy of the model and with Scaling Training and testing accuracy of the model.

The results were as follows:

Decision Tree:



Without Scaling Training accuracy of the model is 0.7453333333333333

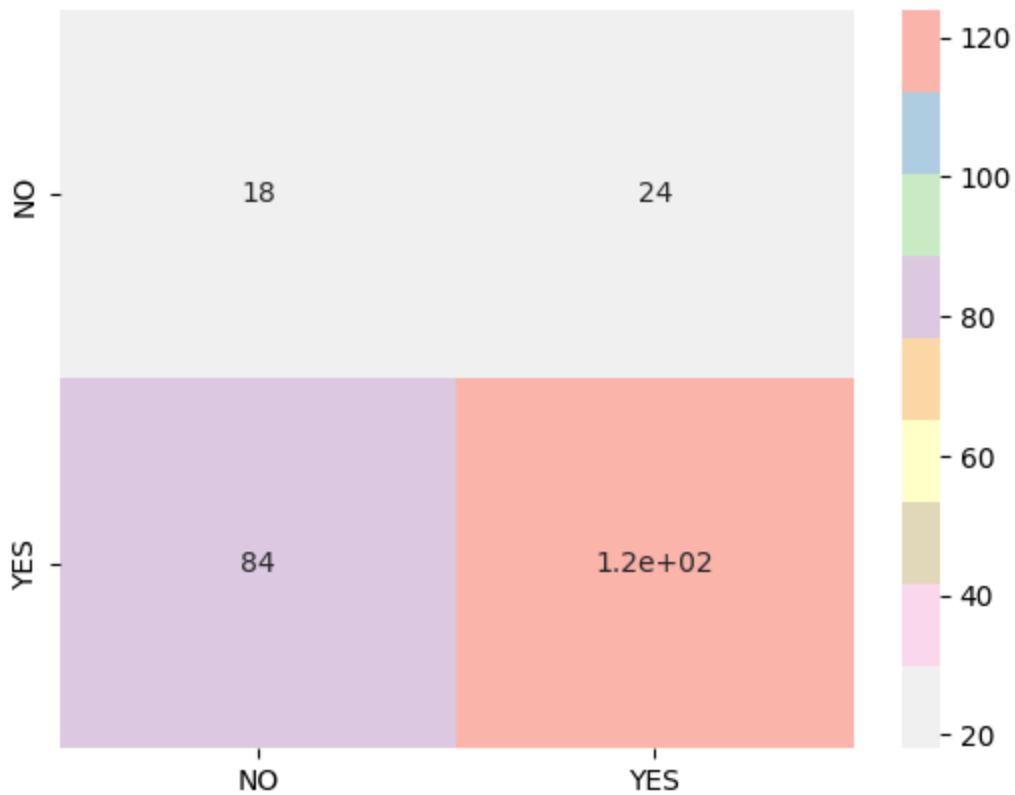
Without Scaling Testing accuracy of the model is 0.492

With Scaling Training accuracy of the model is 0.7453333333333333

With Scaling Testing accuracy of the model is 0.548

[0 148]		precision	recall	f1-score	support
0		0.44	0.36	0.40	102
1		0.61	0.68	0.64	148
accuracy				0.55	250
macro avg		0.52	0.52	0.52	250
weighted avg		0.54	0.55	0.54	250

Naive Bayes:



Without Scaling Training accuracy of the model is 0.60

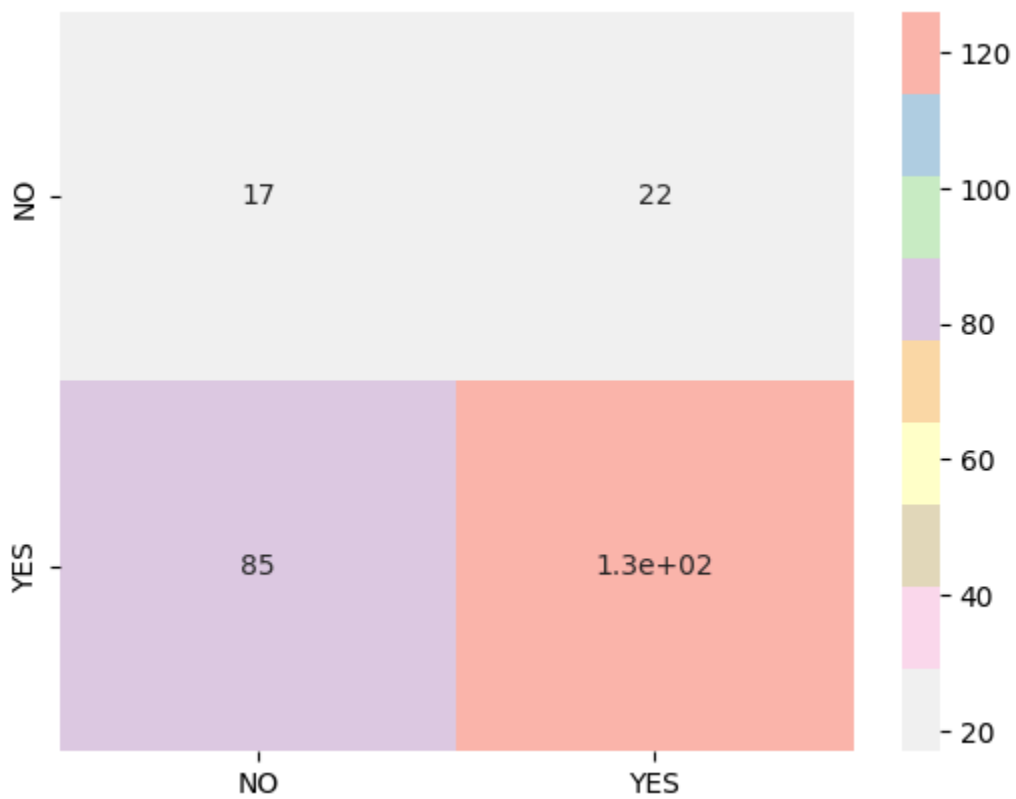
Without Scaling Testing accuracy of the model is 0.57

With Scaling Training accuracy of the model is 0.60

With Scaling Testing accuracy of the model is 0.57

[84 124]					
	precision	recall	f1-score	support	
0	0.43	0.18	0.25	102	
1	0.60	0.84	0.70	148	
accuracy			0.57	250	
macro avg	0.51	0.51	0.47	250	
weighted avg	0.53	0.57	0.51	250	

Random Forest:



Without Scaling Training accuracy of the model is 0.74

Without Scaling Testing accuracy of the model is 0.52

With Scaling Training accuracy of the model is 0.75

With Scaling Testing accuracy of the model is 0.57

	precision	recall	f1-score	support
0	0.44	0.17	0.24	102
1	0.60	0.85	0.70	148
accuracy			0.57	250
macro avg	0.52	0.51	0.47	250
weighted avg	0.53	0.57	0.51	250

Outcome from Model Comparison:

Based on these metrics, the Random Forest model significantly outperformed the other two. The random forest model proved to be the most robust and consistent performer, both before and after scaling. Naive bayes show no change before and after scaling and the decision tree improves its testing accuracy but there is no change in its training accuracy.

Conclusion

This project, which centered on lung cancer prediction. It successfully navigated through the various stages of a machine learning project, from data preprocessing to model training and evaluation. The journey resulted in insightful findings that not only shed light on the capabilities of various machine learning models, but also highlighted the importance of careful feature selection and model evaluation.

Finally, this project not only met its goal of predicting product selling prices, but it also provided a thorough understanding of the complexities involved in machine learning projects, from data preparation to model evaluation and selection. These discoveries pave the way for further analysis and continuous improvement in the field of predictive modeling.