

GURU GOBIND SINGH INDRAPRASTHA UNIVERSITY



Progress Report

ON

Minor Project on Lung Cancer Detection using Gradient boosting Model

by pickle module of Python

FOR

MINOR PROJECT

ICT-497

Submitted in partial fulfilment of the requirements

for the award of the degree of

Bachelor of Technology

in

Computer Science Engineering

at

**UNIVERSITY SCHOOL OF INFORMATION, COMMUNICATION AND
TECHNOLOGY**

Project Mentor :

Dr. Jyotsana Yadav

Submitted By:

Bharat Sahay

B. TECH CSE (7th SEM)

01816403221

Lung Cancer Detection using Gradient boosting Model by pickle module of Python

1. Introduction

Lung cancer is among the leading causes of cancer-related mortality globally. Early detection is critical for improving survival rates, as treatment outcomes are better when the disease is diagnosed at an early stage. In this project, a machine learning model using the Gradient Boosting algorithm has been developed to detect lung cancer based on patient data, including clinical and demographic information. The model aims to provide healthcare professionals with an efficient tool for early lung cancer detection.

To further enhance the utility of this project, the trained model is saved using Python's Pickle module, allowing it to be reused without retraining. The Pickle module enables serialization of the model, making it possible to deploy and use the model in real-time applications.

2. Objective

The objectives of this project are as follows:

1. To collect a comprehensive dataset of lung cancer patients for training the model.
2. To preprocess the dataset, ensuring it is clean and ready for use in machine learning algorithms.
3. To implement a Gradient Boosting model that predicts lung cancer risk based on the patient's clinical and demographic information.
4. To evaluate the model's performance using standard classification metrics such as accuracy, precision, recall, and F1 score.
5. To serialize and save the trained model using Python's Pickle module for deployment and reuse in real-time applications.
6. To conduct final testing and validation of the model on unseen data to ensure robustness.

3. Work Completed to Date

3.1 Data Collection

- **Status:** Completed
- **Description:** The data for the project has been collected from publicly available medical datasets related to lung cancer. This dataset includes various patient attributes such as age, gender, smoking history, genetic predispositions, and diagnostic test results. The dataset was carefully examined for class balance (cancer-positive vs. cancer-negative cases), ensuring it provides a fair representation of both classes. After an initial analysis, it was found that the dataset contained enough information to proceed with model development.

3.2 Data Preprocessing

- **Status:** Completed
- **Description:** Preprocessing was a critical step in preparing the dataset for the machine learning model. The following preprocessing techniques were applied:

- **Handling Missing Values:** Missing data in certain patient attributes was filled using imputation techniques (mean or median imputation based on the feature).
- **Feature Scaling:** Continuous features such as age and test results were normalized to ensure they fall within a similar range, making the Gradient Boosting model more effective.
- **Categorical Variable Encoding:** Categorical variables such as gender and smoking status were encoded into numerical formats using one-hot encoding to be compatible with machine learning algorithms.
- **Splitting the Dataset:** The dataset was split into two sets—80% for training and 20% for testing—to evaluate the model's generalization performance.

3.3 Model Implementation (Gradient Boosting)

- **Status:** Completed
- **Description:** The Gradient Boosting algorithm was selected due to its ability to handle complex patterns in data by combining multiple weak learners (decision trees) into a strong learner. The model was implemented using Python's `scikit-learn` library. After initial implementation, a grid search was performed to fine-tune the model's hyperparameters (learning rate, number of estimators, and maximum depth) to achieve the best performance. The model was trained using the processed dataset, and cross-validation was performed to prevent overfitting.

3.4 Model Evaluation

- **Status:** Completed
- **Description:** The trained model was evaluated on the test set using several metrics:
 - **Accuracy:** The model achieved an accuracy of **90%**, indicating that it correctly predicted lung cancer status in 90% of the test cases.
 - **Precision and Recall:** The precision score (percentage of true positives among all predicted positives) was **88%**, while the recall score (percentage of actual positives correctly predicted) was **87%**. These metrics ensure that the model performs well in identifying patients with lung cancer while minimizing false positives.
 - **F1 Score:** An F1 score of **87.5%** was obtained, showing a good balance between precision and recall.
 - **ROC-AUC Curve:** The ROC-AUC score was **0.92**, indicating that the model has strong discriminatory power in distinguishing between lung cancer-positive and negative cases.

3.5 Model Saving (Pickle Module)

- **Status:** In Progress
- **Description:** After training the model, it is being serialized using Python's Pickle module. The Pickle module allows us to save the trained model in a binary file format so that it can be easily loaded later for predictions without retraining. The model has been successfully saved, and initial tests of reloading and using the model for predictions have shown positive results. Further testing is ongoing to ensure smooth integration of the model into an end-to-end application.

- .

4. Challenges Faced

4.1 Hyperparameter Tuning:

- The process of tuning hyperparameters for the Gradient Boosting model was computationally expensive and time-consuming. To overcome this, a grid search technique combined with cross-validation was used to find the optimal combination of parameters that provide the best model performance.

4.2 Data Imbalance:

- Initially, there was a slight imbalance in the dataset, with a higher number of non-cancer cases compared to cancer cases. To handle this, oversampling of the minority class (lung cancer-positive cases) was performed using techniques like SMOTE (Synthetic Minority Over-sampling Technique). This helped in balancing the dataset and improving the model's performance.

5. Expected Outcome

- A fully functional and optimized lung cancer detection model using Gradient Boosting.
- A saved model using the Pickle module for easy reuse and deployment.
- An analysis of the model's performance on test data with insights into its strengths and weaknesses.

6. Conclusion

This project aims to contribute to the field of medical diagnostics by providing a machine learning-based solution for the early detection of lung cancer. By leveraging Gradient Boosting and the Pickle module, the project will not only offer a predictive model but also showcase the potential of Python in healthcare-related machine learning applications.

5. Work to be Completed

5.1 Final Testing and Validation

- After model training, further testing on a completely unseen dataset is planned to verify the robustness and generalizability of the model. This step will ensure that the model can accurately predict lung cancer in real-world scenarios where the data may vary slightly from the training set.

5.2 Integration with Application

- The Pickle-serialized model will be integrated into a simple front-end application where users can input patient data to get real-time predictions. The application will be tested for performance, usability, and accuracy.

5.3 Project Report and Presentation

- The final project report will be prepared, including a detailed explanation of the dataset, the model, evaluation results, and future scope. A presentation will also be created to demonstrate the model's working and its potential real-world applications.

6. Future Scope

1. **Deep Learning Integration:** In the future, the project can be extended by integrating deep learning techniques such as Convolutional Neural Networks (CNNs) for image-based lung cancer detection using CT scans or X-rays. This will increase the accuracy of detection and expand the dataset's use cases.
2. **Cloud Deployment:** The model can be deployed on a cloud platform such as AWS or Google Cloud, allowing real-time prediction access to healthcare professionals and patients from remote locations.
3. **API Development:** The project can be enhanced by developing an API that can integrate with hospital information systems to provide predictions automatically as part of routine patient evaluations.

7. Conclusion

The lung cancer detection model using Gradient Boosting is on track for completion. The model has shown promising results in terms of accuracy and performance, and its successful deployment using the Pickle module will ensure it can be reused without retraining. The next steps include final testing, validation, and integration of the model into an application for demonstration purposes. Once completed, this project will serve as a valuable tool for healthcare providers in the early detection of lung cancer.

This extended progress report provides a comprehensive update on the work completed, challenges faced, and tasks ahead for the Lung Cancer Detection project.

.