

**Judith Osuna Ramírez**

**Use of artificial intelligence in the prediction of  
response in patients with COVID-19 pneumonia  
treated with low-dose anti-inflammatory  
radiotherapy**

**Final Degree Project**

**directed by Dr. Victor Hernandez Masgrau  
and Dr. Meritxell Arenas Prat**

**Biomedical Engineering Degree**



**UNIVERSITAT ROVIRA I VIRGILI**

**Tarragona**

**2024**

*Dedicated*

*to my parents, Ramón and Rosario,  
who have supported and encouraged me throughout  
my entire academic journey.*

*To my brother Javier, to Marina,  
and to my nephew Oliver, who are all so special to me.*

*To my university schoolmates, especially Gabriela and Òscar,  
who have become significant individuals in my life.*

*To my lifelong friend Alba, who has been by my side since childhood.*

*To Dani, my best chance.*

*To my supervisors Víctor and Meritxell,  
for believing in me.*

*To Alberto, the hospital IT specialist  
who helped me to get started in the world of radiomic data.*

*To Paula, my work partner-in-crime.*

*To Marta and Raquel,  
for their unwavering support and assistance in every possible way.*

*To Hatem, the professor who has given me the most confidence  
and support in all the technical aspects of this work.*

*Lastly, to Amsterdam, for providing me with the most wonderful views  
during most of the creation of this project.*

# Abstract

This project presents a multi-omic predictive model designed to predict the response to low-dose pulmonary radiotherapy in patients suffering from pneumonia caused by COVID-19. Clinical, metabolomic, and radiomic characteristics were acquired from 50 patients with varying responses to the same treatment. Using supervised machine learning models, the best classifier for predicting the response to this treatment was studied.

The results highlight that predictive models for binary classification achieve improved results when dealing with imbalanced response data and employing advanced feature selection methods.

This approach offers substantial advancements in predicting responses to a certain treatment, marking a significant step towards precision and personalized medicine through the use of Artificial Intelligence as a clinical support system.

**Keywords:** COVID-19; low-dose radiotherapy; Machine Learning; multi-omics.

# Resumen

Este proyecto presenta un modelo predictivo multi-ómico diseñado para predecir la respuesta a la radioterapia pulmonar de baja dosis en pacientes que sufren neumonía causada por COVID-19. Se adquirieron características clínicas, metabolómicas y radiómicas de 50 pacientes con respuestas variables al mismo tratamiento. Mediante modelos supervisados de aprendizaje automático, se estudió cuál era el mejor clasificador para predecir la respuesta a este tratamiento.

Los resultados ponen de manifiesto que los modelos predictivos de clasificación binaria obtienen mejores resultados cuando se tratan los datos de respuesta desequilibrados y se emplean métodos avanzados de selección de características.

Esta aproximación ofrece avances sustanciales en la predicción de respuestas a un determinado tratamiento, marcando un paso significativo hacia la medicina de precisión y personalizada mediante el uso de la Inteligencia Artificial como sistema de apoyo clínico.

**Palabras clave:** COVID-19; radioterapia; aprendizaje supervisado; multi-ómica.

# Resum

Aquest projecte presenta un model predictiu multi-òmic dissenyat per a predir la resposta a la radioteràpia pulmonar de baixa dosi en pacients que sofreixen pneumònia causada per COVID-19. Es van adquirir característiques clíniques, metabolòmiques i radiòmiques de 50 pacients amb respostes variables al mateix tractament. Mitjançant models supervisats d'aprenentatge automàtic, es va estudiar quin era el millor classificador per a predir la resposta a aquest tractament.

Els resultats posen de manifest que els models predictius de classificació binària obtenen millors resultats quan es tracten les dades de resposta desequilibrats i es fan servir mètodes avançats de selecció de característiques.

Aquesta aproximació ofereix avanços substancials en la predicció de respostes a un determinat tractament, marcant un pas significatiu cap a la medicina de precisió i personalitzada mitjançant l'ús de la Intel·ligència Artificial com a sistema de suport clínic.

**Paraules clau:** COVID-19; radioteràpia; aprenentatge supervisat; multi-òmica.



# Table of Contents

1. Introduction . . . . .	1
1.1. Background . . . . .	1
1.2. IPACOVID clinical trial . . . . .	2
1.3. Application of Artificial Intelligence in analysing patient response to LD-RT . . . . .	2
1.4. Multiomic data . . . . .	3
1.4.1. Clinical data . . . . .	3
1.4.2. Metabolomic data . . . . .	3
1.4.3. Radiomic data . . . . .	3
1.5. Clinical applicability of predictive models . . . . .	4
2. Objectives . . . . .	5
3. Materials and methods . . . . .	6
3.1. Data collection . . . . .	6
3.2. Data extraction . . . . .	8
3.3. Data preparation . . . . .	16
3.4. Data augmentation . . . . .	19
3.5. Data preprocessing . . . . .	22
3.6. Predictive models . . . . .	24
3.7. Model validation . . . . .	27
4. Results . . . . .	28
5. Discussion . . . . .	30
6. Conclusions . . . . .	32
References . . . . .	35
Annexes . . . . .	36
I. Feature Importance Ranking . . . . .	36
II. Confusion Matrices . . . . .	38
III. Radiomic features . . . . .	42
IV. Codes . . . . .	43

# Abbreviations

<b>AI</b>	Artificial Intelligence
<b>AUC</b>	Area Under the Curve
<b>bpm</b>	breaths per minute
<b>CM</b>	Confusion Matrix
<b>CT</b>	Computed Tomography
<b>CTV</b>	Clinical Target Volume
<b>CV</b>	Cross-Validation
<b>COVID</b>	COronaVirus Disease 2019
<b>DICOM</b>	Digital Imaging and Communications in Medicine
<b>DT</b>	Decision Tree
<b>EHR</b>	Electronic Health Record
<b>FiO<sub>2</sub></b>	Fraction of Inspired Oxygen
<b>FN</b>	False Negatives
<b>FP</b>	False Positives
<b>GDS</b>	Global Deterioration Scale
<b>GLCM</b>	Gray-Level Co-Occurrence Matrix
<b>GLDM</b>	Gray-Level Dependence Matrix
<b>GLRLM</b>	Gray-Level Run-Length Matrix
<b>GLSZM</b>	Gray-Level Size Zone Matrix
<b>GOT</b>	Glutamic Oxaloacetic Transaminase
<b>GPT</b>	Glutamic Pyruvic Transaminase
<b>HUSJR</b>	Hospital Universitari Sant Joan de Reus
<b>KDE</b>	Kernel Density Estimate
<b>LBP</b>	Local Binary Patterns
<b>LD</b>	Low Dose
<b>LR</b>	Logistic Regression
<b>ML</b>	Machine Learning
<b>MIQ</b>	Mutual Information Quotient

<b>mRMR</b>	minimum Redundancy Maximum Relevance
<b>NGTDM</b>	Neighborhood Gray-Tone Difference Matrix
<b>NIFTI</b>	Neuroimaging Informatics Technology Initiative
<b>PaO<sub>2</sub></b>	Partial Pressure of Oxygen
<b>PM</b>	Predictive Model
<b>PTV</b>	Planning Target Volume
<b>RF</b>	Random Forest
<b>ROC</b>	Receiver Operating Characteristic
<b>ROI</b>	Region Of Interest
<b>rpm</b>	respirations per minute
<b>RT</b>	Radiotherapy
<b>SMOTE</b>	Synthetic Minority Oversampling TEchnique
<b>SoC</b>	Standard of Care
<b>SpO<sub>2</sub></b>	Peripheral Capillary Oxygen Saturation
<b>TN</b>	True Negatives
<b>TP</b>	True Positives
<b>VIF</b>	Variance Inflation Factor
<b>WHO</b>	World Health Organization
<b>YAML</b>	YAML Ain't Markup Language

# 1. Introduction

## 1.1 Background

Throughout the COVID-19 pandemic, the world faced an unprecedented challenge. The outbreak began in December of 2019 in a city called Wuhan, in China [1]. Thereafter, as illustrated in the Figure 1 provided by the World Health Organization (WHO), the coronavirus spread to every corner of the globe, with some countries being more severely affected than others, but no nation was completely spared. Since the WHO declared the global pandemic situation on the 11th of March 2020, over 775 million cases have been reported worldwide [2].

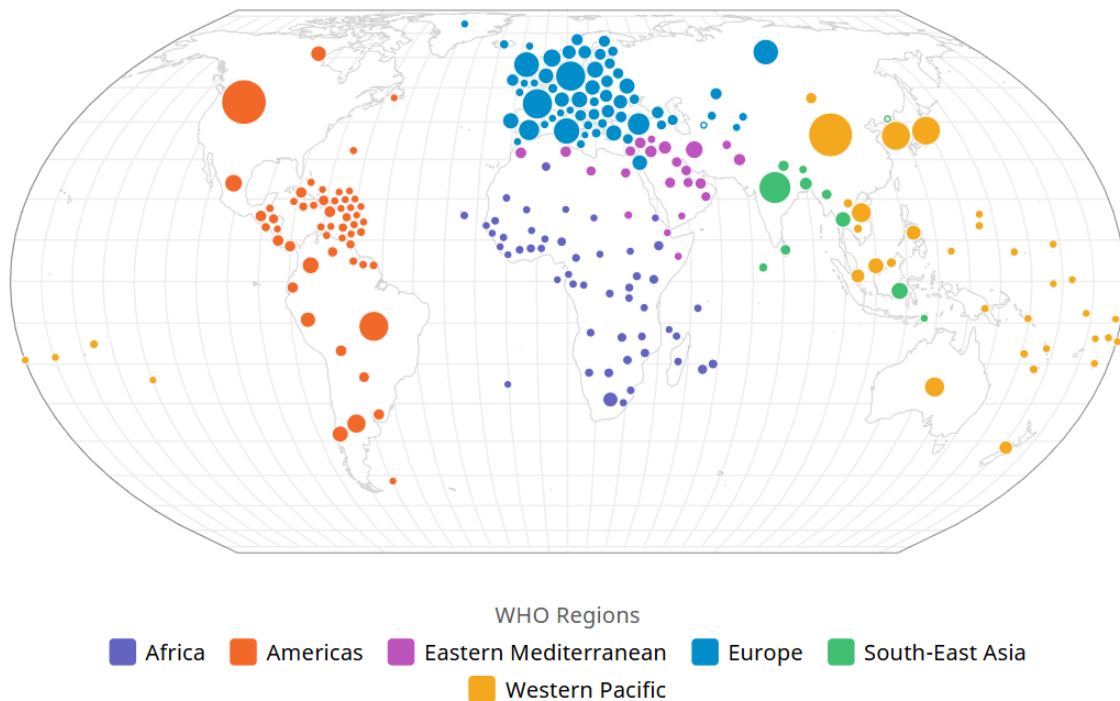


Figure 1: Total cumulative number of COVID-19 cases reported to the World Health Organization [2].

Some people were asymptomatic carriers of the virus, while others experienced mild symptoms similar to a common cold. However, many people, especially the elderly and those with pre-existing medical conditions, required prolonged hospitalization due to severe complications such as bilateral pneumonia. For these patients, the standard treatments administered included antiviral and anti-inflammatory drugs, as well as mechanical ventilation. In cases where patients were not candidates for invasive mechanical ventilation and pharmacological treatments were not sufficient to bring about recovery, an innovative approach offered new hope. This treatment consisted of a clinical trial, named IPACOVID, of total lung Low Dose Radiotherapy (LD-RT) to improve patients' oxygenation [3].

## 1.2 IPACOVID clinical trial

IPACOVID clinical trial is a Phase 3 non-randomized ambispective study involving 100 participants, comprising 50 control patients and 50 cases. The mean age of both cohorts was 85 years, and they were matched by sex, comorbidities, and SAFI values. The control cohort received only the standard of care (which includes antiviral or anti-inflammatory drugs). The experimental cohort received the standard of care plus a single radiation dose of 0.5 Gy to the entire thorax [3].

As demonstrated by *Arenas et al.*[4], radiotherapy administered at low doses (LD-RT) modulates the inflammatory response with efficacy, producing an anti-inflammatory effect. This means that we can apply LD-RT to treat benign diseases such as osteoarthritis, humeral epi-condylitis, scapular–humeral peri-arthritis, or heel spurs.

Therefore, the clinical trial hypothesized that applying LD-RT could improve patients' respiratory function by reducing lung inflammation caused by COVID-19 pneumonia, thereby decreasing the mortality rate. The trial results showed that patients in the experimental cohort who received LD-RT in addition to Standard of Care (SoC) experienced improvements in respiratory parameters, shorter hospitalization periods, reduced 1-month mortality rates, and prolonged overall survival compared to those who received SoC alone.

## 1.3 Application of Artificial Intelligence in analysing patient response to LD-RT

In this project, we focused on analysing the responses to the LD-RT pulmonary treatment in the experimental cohort previously mentioned. While the vast majority of those fifty patients responded positively and recovered with the help of the medication, there were some for whom this treatment was ineffective. Our goal was to study and identify the relationships between variables in patients who responded positively and those who did not. This is helpful to better determine which patients are likely to be suitable candidates for this treatment in the future.

To achieve this, we conducted a hybrid study that integrated clinical, metabolomic, and radiomic data to develop a highly accurate and effective predictive model (PM) for assessing responses to LD-RT treatment. After the data was collected and cleaned, we applied several supervised ML algorithms to identify the one that yields the best results.

## 1.4 Multiomic data

All data used in this study was collected from the fifty patients participating in the IPACOVID clinical trial conducted at Hospital Universitari Sant Joan de Reus, Spain.

While many researchers and practitioners have focused on enhancing the quality of models (such as the architecture and automated feature selection) to improve the results, there are limited efforts towards improving the data quality. Nonetheless, the performance of a ML model is upper bounded by the quality of the data used [5].

As indicated above, this project combines data from various sources to gather the maximum amount of information about each patient. The following sections will explain each type of data and provide further details.

### 1.4.1 Clinical data

Clinical data is derived from the patient's historical information, some of which is stored in the hospital's EHR (Electronic Health Record). In this context, it encompasses personal details about the patient, including pre-existing conditions, functional assessments and respiratory parameters [6].

### 1.4.2 Metabolomic data

"The suffix omics refers specifically to the whole of something. Therefore, the omics sciences are those sciences that allow us to study a large number of molecules, that are involved in the functioning of an organism." [7]

In this context, metabolomic data pertains to metabolites and cellular components present in the blood that have been shown to be important in COVID-19 patients [6], as determined through blood analysis.

### 1.4.3 Radiomic data

Radiomics has been defined as an area of study focused on extracting quantitative metrics, known as radiomic features, from medical images. These features reflect various aspects of tissue and lesion characteristics, such as heterogeneity and shape. When used alone or together with demographic, histologic, genomic, or proteomic information, radiomic features can aid in addressing clinical challenges [8].

In this study, our medical images were simulation CT scans stored in DICOM format and carried out before the LD-RT treatment.

## **1.5 Clinical applicability of predictive models**

One of the primary advantages of Artificial Intelligence (AI) in the medical field is its capacity to process and analyse large volumes of clinical data with a speed and accuracy that far surpass human capabilities. This enables the identification of patterns and trends that might otherwise go unnoticed. In this specific context, AI can evaluate a multitude of clinical, metabolomic, and imaging variables to more accurately predict which patients will respond favorably to the treatment. This not only optimizes medical resources but can also significantly improve patient outcomes.

Nevertheless, this advanced tool also presents several challenges. For an AI model to be installed in a clinical setting, it must demonstrate that it ensures data security and privacy, is free from bias, and has proven efficacy [9]. In other words, it must be a reliable model. Furthermore, the implementation of this technology needs that healthcare professionals possess a thorough understanding of its functionality and limitations. It is essential that they do not place blind trust in its efficacy and remember that the final decision should always be theirs. AI should be regarded as a decision-support tool, not a replacement for clinical judgment.

## 2. Objectives

The main objective of this project is to develop a PM capable of determining the response to the LD-RT treatment in patients with COVID-19 pneumonia.

To achieve this, clinical and metabolomic data were collected, while radiomic features were extracted from CT images. To address the issue of data imbalance—where positive responses were more frequent than negative ones—data augmentation techniques were employed. Subsequently, two different techniques were evaluated for preprocessing and selecting non-correlated variables for model training. Finally, multiple supervised ML models were developed and compared using various evaluation metrics to assess their efficacy.

The goal is to identify the best combination of preprocessing techniques and model training methods to create an effective algorithm for predicting which future COVID-19 patients are likely to benefit from LD-RT treatment.

### 3. Materials and methods

The steps used to prepare the data to the PM are illustrated in Figure 2 and described in the following sections.

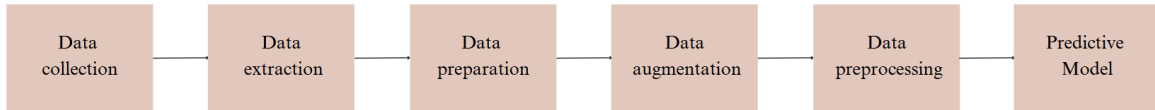


Figure 2: Diagram showing the methodology workflow.

#### 3.1 Data collection

Both the clinical and metabolomic databases used in the present study were collected from the IPACOVID prospective clinical trial [3] and described by *Piqué et al.* [6].

Radiomic databases were not available in those datasets and were created from scratch from the CT basal images of the 50 patients.

Before proceeding, we needed to process each of the fifty CT images individually to isolate the target area that outlines the lung volume. In radiotherapy, target volumes typically include a margin around the lungs to ensure accurate radiation delivery. For our purposes, though, we are only interested in extracting and analyzing the specific lung information.

The steps to obtain the desired CT images are the following:

1. **Lung Segmentation Wizard:** The software used in the Radiotherapy department of HUSJR to prescribe treatments to patients is called *Eclipse* [10] and included a tool designed to automatically delineate the lungs. Nevertheless, this tool proved ineffective for our images, as COVID-19 patients had mucus and fluid in their lungs, resulting in abnormal lung structures. This necessitated an additional processing step.
2. **Manual adjustment of the segmentation:** With the drawing tools, we could paint and erase specific areas using the desired brush size, correcting any parts that the segmentation tool either missed or incorrectly highlighted. This permits to manually contour and refine each section of the CT scan.
3. **Post processing Tool:** This tool allowed us to remove parts of the drawing that were smaller than a specified threshold ( $0.5 \text{ cm}^2$ , in this case) and retain the ‘n’ largest segments, in this instance, 2. This process ensures that the resulting structure accurately represents the lungs by eliminating any accidental or irrelevant portions.

Once everything had been verified, we could proceed to download the resulting anonymized image, which included the structures: Body, PTV, CTV, and, ultimately, our target structure, z\_lungs, as shown in Figure 3; where PTV (Planning Target Volume) is the volume including the CTV with an additional margin to account for uncertainties in treatment delivery and CTV (Clinical Target Volume) is the volume encompassing the tumor or treatment area, along with any surrounding tissues that may require treatment.

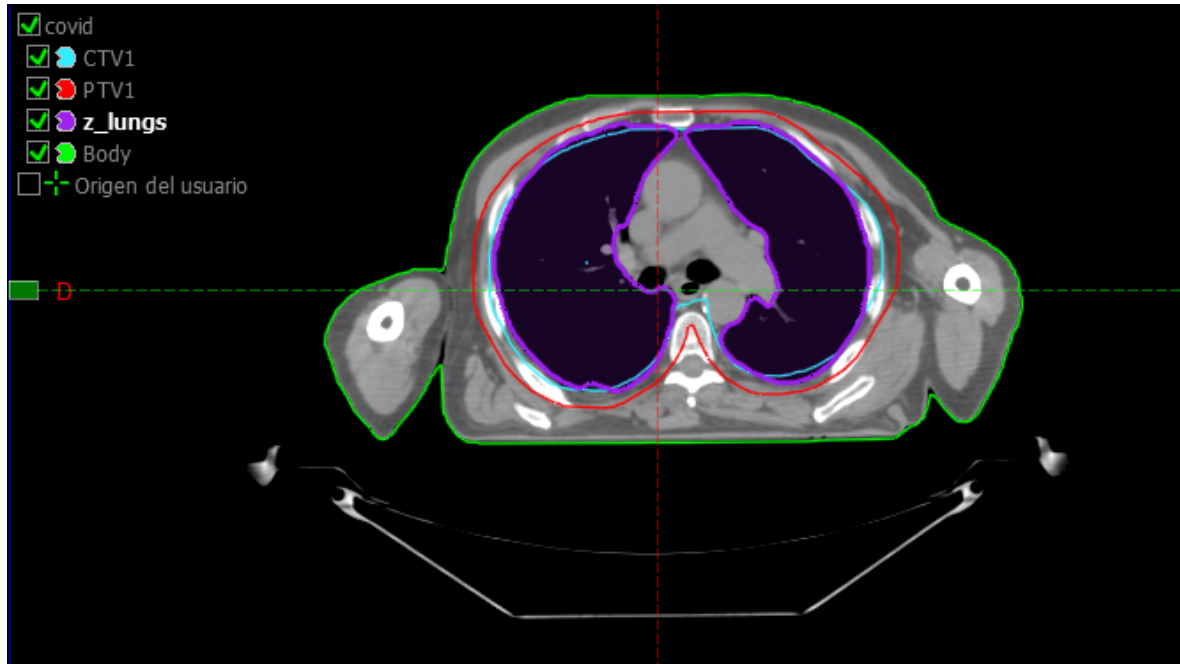


Figure 3: Computed Tomography Structures

The resulting 'z\_lungs' structure from the CT scan above is represented in 3D in Figure 4.

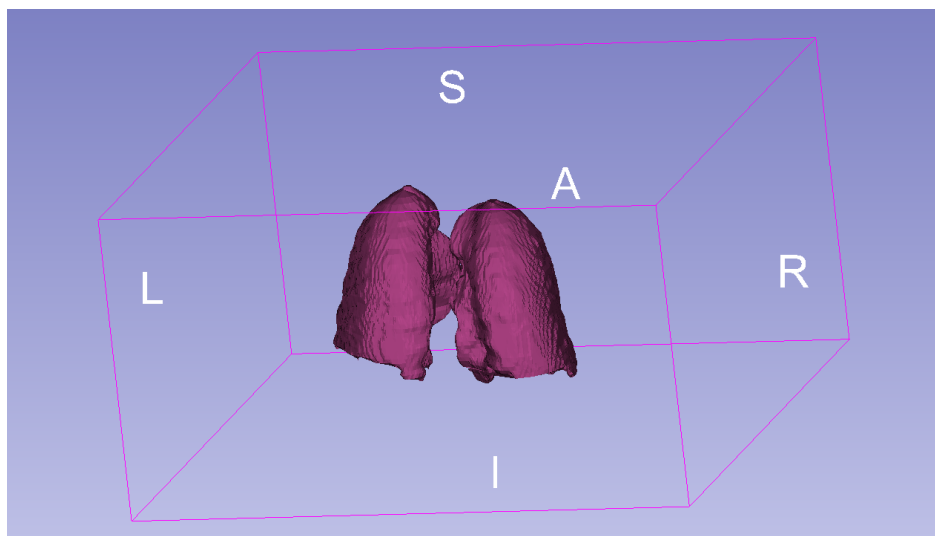


Figure 4: Our target structure: 'z\_lungs'

## 3.2 Data extraction

At this point, we needed to select the appropriate clinical and metabolomic features for this algorithm, as those used in the article from *Piqué et al.* [6] also took into account clinical and metabolomic data after performing the CT scan. In our case, this information is redundant since the prediction focuses on finding differences between patients before treatment, specifically at the time of the CT scan.

### *Clinical data*

Starting with the clinical data, we were left with the following features selected as appropriate for the model, given that they focus on data obtained before the LD-RT treatment and reflect the relevant information for response prediction at the time of the CT scan.

Table 3.1 provides a brief description of the clinical variables selected, detailing the units used in the database, as well as the normal value ranges or their coding into binary and categorical variables.

Variable	Units	Description	Normal Range/ Codification
<b>Age</b>	Years	The patient's age at the time of the CT scan.	Elderly people
<b>Sex</b>	Unitless	The gender of the patient.	0 (Male) or 1 (Female)
<b>Neurological diseases</b>	Unitless	Presence of previous neurological diseases.	1 (Presence) or 0 (Absence)
<b>Cardiovascular diseases</b>	Unitless	Presence of previous cardiovascular diseases.	1 (Presence) or 0 (Absence)
<b>Respiratory diseases</b>	Unitless	Presence of previous respiratory diseases.	1 (Presence) or 0 (Absence)
<b>Other diseases</b>	Unitless	Presence of other diseases not classified in the previous categories.	1 (Presence) or 0 (Absence)
<b>Barthel scale</b>	Scale	Assessment of the patient's ability to perform activities of daily living.	From 0 (total dependence) to 5 (total independence)

Variable	Units	Description	Normal Range/ Codification
<b>GDS</b>	Scale	Measurement to assess the level of cognitive decline in individuals with dementia.	From 0 (no depression) to 6 (severe depression)
<b>Days symptoms</b>	Days	Number of days since the onset of symptoms up to the time of the CT scan.	Variable based on clinical context
<b>Days hospitalization</b>	Days	Number of days the patient remained admitted in the hospital from their admission to their discharge.	Variable based on clinical context
<b>SAFI</b>	Ratio	SpO <sub>2</sub> /FiO <sub>2</sub> .	Greater than 400
<b>PAFI</b>	Ratio	PaO <sub>2</sub> /FiO <sub>2</sub> .	Greater than 300
<b>Pulmonary affectation</b>	Scale	Measurement of pulmonary affectation.	From 1 (5-25%) to 4 (more than 75%)
<b>PO<sub>2</sub> venous</b>	mmHg	Partial pressure of oxygen in venous blood.	40-50 mmHg
<b>Oxygen saturation</b>	Percentage	Percentage of oxygen saturation in the blood.	95-100%
<b>Liters oxygen</b>	Liters	Amount of oxygen administered to the patient.	Variable based on patient needs
<b>Cardiac frequency</b>	bpm	Number of heartbeats per minute.	60-100 bpm
<b>Respiratory frequency</b>	rpm	Number of breaths per minute.	12-20 rpm
<b>CURB-65 scale</b>	Score	Score for assessing the severity of pneumonia and the need for hospitalization.	From 0 (low risk) to 5 (high risk)

Table 3.1: Description of Clinical Variables.  
GDS: Global Deterioration Scale.

The variable nomenclature has been coded according to the guidelines provided by *Piqué et al.*[6]. Additionally, some descriptions and annotated values have been verified using [11], [12], and [13].

Therewith, a thorough description of these features is provided in Table 3.2, including basic statistics such as the mean values, dispersion or standard deviation, as well as maximum and minimum measures. Additionally, an analysis of the total number of missing values is

performed, allowing us to understand the extent of missing data and decide on management strategies such as imputation or removal. Finally, the types of variables are analyzed to identify any categorical variables within the data that need to be processed for model training.

	Mean	Std	Min	Max	Missing values (%)	Type
<b>Age</b>	84,7	6,58	60	96	0	Integer
<b>Sex</b>	0,5	0,50	0	1	0	Binary
<b>Neurological diseases</b>	0,3	0,45	0	1	0	Binary
<b>Cardiovascular diseases</b>	0,8	0,39	0	1	0	Binary
<b>Respiratory diseases</b>	0,4	0,49	0	1	0	Binary
<b>Other diseases</b>	0,9	0,33	0	1	0	Binary
<b>Barthel scale</b>	3,5	1,13	1	5	0	Categorical
<b>GDS</b>	2,1	1,45	1	6	0	Categorical
<b>Days symptoms</b>	5,1	1,64	2	7	0	Integer
<b>Days hospitalization</b>	18,6	13,90	0	57	0	Integer
<b>SAFI</b>	280,9	93,02	90	471	0	Integer
<b>PAFI</b>	262,3	99,05	59	628	4	Integer
<b>Pulmonary affectation</b>	3,1	0,76	1	4	2	Categorical
<b>PO<sub>2</sub> venous</b>	41,9	12,88	16	82,5	32	Float
<b>Oxygen saturation</b>	93,2	2,98	86	99	2	Integer
<b>Liters oxygen</b>	6,9	4,57	2	15	4	Integer
<b>Cardiac frequency</b>	101,6	37,23	61	329	2	Integer
<b>Respiratory frequency</b>	20,5	3,73	15	30	4	Integer
<b>CURB-65 scale</b>	2,9	0,68	2	4	0	Integer

Table 3.2: Clinical data description.  
Std: Standard Deviation;  
GDS: Global Deterioration Scale.

Considering that the values pertain to severely ill patients, we had to adopt a more flexible approach when identifying outliers, as all variables might be affected by COVID-19. In this context, we did not observe any values that were clearly out of range, and therefore, there was insufficient reason to exclude them based on extreme values alone. Nonetheless, we noted that one variable, *PO<sub>2</sub> venous*, contained a higher percentage of missing values (32%). This variable, along with the three categorical variables (*Barthel scale*, *GDS*, *Pulmonary affectation*), needed to be addressed in the data preparation section to ensure robustness in our analysis. Additionally, we observed a balanced distribution of gender among the observed patients, which suggests that our model is unlikely to introduce gender bias.

**Metabolomic data**

Continuing with the metabolomic data, the same description of these data is provided in the Table 3.3 below.

<b>Variable</b>	<b>Units</b>	<b>Description</b>	<b>Normal Range</b>
<b>IL6</b>	pg/mL	A pro-inflammatory cytokine involved in immune response and inflammation.	10 pg/mL
<b>Ferritin</b>	ng/mL	A protein that stores iron in the body; its level reflects iron reserves.	- Men: 30-300 ng/mL - Women: 30-200 ng/mL
<b>PCR</b>	mg/L	A marker of acute inflammation in the body.	10 mg/L
<b>LDH</b>	U/L	An enzyme involved in cellular metabolism; can indicate tissue damage or disease.	140-280 U/L
<b>CD4</b>	cells/ $\mu$ L	A type of T cell crucial for immune response; used to monitor immune status.	500-1,500 cells/ $\mu$ L
<b>CD8</b>	cells/ $\mu$ L	T cells that help eliminate infected and tumor cells.	200-1,000 cells/ $\mu$ L
<b>DimerD</b>	ng/mL	A fragment of fibrin breakdown; used to assess the presence of thrombosis.	500 ng/mL
<b>Creatinin</b>	mg/dL	A waste product from muscle metabolism, used to assess kidney function.	- Men: 0.7-1.3 mg/dL - Women: 0.6-1.1 mg/dL
<b>GOT</b>	U/L	A liver enzyme that helps evaluate liver function and damage.	10-40 U/L
<b>GPT</b>	U/L	A liver enzyme that also assesses liver function and damage.	7-56 U/L
<b>Hemoglobin</b>	g/dL	A protein in red blood cells that carries oxygen.	- Men: 14-17 g/dL - Women: 12-16 g/dL
<b>Leukocytes</b>	cells/ $\mu$ L	White blood cells that help fight infections.	4,000-11,000 cells/ $\mu$ L
<b>Lymphocytes</b>	cells/ $\mu$ L	A type of white blood cell essential for the immune response.	1,000-4,000 cells/ $\mu$ L
<b>Platelets</b>	cells/ $\mu$ L	Blood cells involved in clotting.	150,000-450,000 cells/ $\mu$ L

Table 3.3: Description of metabolomic variables. U: Units; GOT: Glutamic Oxaloacetic Transaminase; GPT: Glutamic Pyruvic Transaminase.

Using the information provided in [14], Table 3.3 has been completed.

The subsequent table presents a comprehensive description of the metabolomic data used, including the key statistical measures, percentage of missing values and the type of each variable 3.4.

	Mean	Std	Min	Max	Missing values (%)	Type
<b>IL6</b>	64.4	114.22	0	627	4	Float
<b>Ferritin</b>	1349.6	1836.44	42.7	8836	2	Float
<b>PCR</b>	7.9	6.79	0.3	32.8	0	Float
<b>LDH</b>	315.4	137.30	168	800	2	Float
<b>CD4</b>	240.7	151.91	0	827	8	Float
<b>CD8</b>	162.3	223.15	0	1523	8	Float
<b>DimerD</b>	2196.4	4074.35	220	26680	0	Float
<b>Viral_load_baseline</b>	1.7	1.13	0.4	4.89	52	Float
<b>RNA_total_concentration</b>	383.7	46.30	259.2	487.2	8	Float
<b>Glucose</b>	160.3	81.92	43	390	0	Float
<b>Creatinin</b>	885.2	6250.36	0.15	44198	0	Float
<b>GOT</b>	31.3	22.29	8	137	0	Integer
<b>GPT</b>	24.3	20.83	6	97	0	Integer
<b>Hemoglobin</b>	36178.3	17117.35	10	44450	0	Integer
<b>Leukocytes</b>	9407.4	5396.68	2460	27910	0	Integer
<b>Lymphocytes</b>	1062.2	1617.70	200	9400	0	Integer
<b>Platelets</b>	222634.0	95107.09	10700	583000	0	Integer

Table 3.4: Metabolomic data description.

Std: Standard Deviation; GOT: Glutamic Oxaloacetic Transaminase;

GPT: Glutamic Pyruvic Transaminase.

Table 3.4 provides a detailed statistical overview of the metabolomic data. In this table, outliers were evident in several variables, specifically *Ferritin*, *DimerD*, *Creatinin*, and *Hemoglobin*, as their values significantly exceed the normal range established in Table 3.3. Furthermore, *Viral load baseline* exhibited a substantial proportion of missing values, with 52% of its data being unknown. Additionally, the analysis revealed that all variables were numerical, with no categorical variables requiring preprocessing before the predictive modeling phase.

### ***Radiomic data***

Lastly, the extraction of radiomic features from CT images was a more complex process than the previous data extraction and consisted of several steps. Instead of using the 3D Slicer software [15], which was employed in previous works [16] and provided limited feature extraction with much information left undetected, a direct extraction was performed using the Python *PyRadiomics* library [17].

In the following paragraphs, the steps for extracting radiomic features will be explained.

#### **1. Images format transformation from DICOM to NIfTI**

To perform radiomic analysis with tools like PyRadiomics, converting images from DICOM (Digital Imaging and Communications in Medicine) to NIfTI (Neuroimaging Informatics Technology Initiative) format was essential. NIfTI, designed for neurological imaging, simplifies handling 3D and 4D images by consolidating data into one or two files, facilitating data management and integration with segmentation masks. In contrast, DICOM's complexity, involving multiple files and extensive metadata, can complicate processing and reduce efficiency. Therefore, converting to NIfTI enhances compatibility with processing tools, improves performance, and supports more efficient analysis, especially in neuroimaging [18].

In this work, the thoracic image were converted to NIfTI format. Additionally, the contour created in the previous process, named '*z\_lungs*', which exclusively delineated the lung regions, was also converted to NIfTI format to serve as the mask. This resulted in two files for each patient: '*Image.nii*' and '*Mask.nii*', where '*Mask.nii*' represented the segmentation of the lung regions.

The code for this transformation is provided in Annexes: Code 1.

#### **2. Configuration file for feature extraction**

When extracting features from medical or scientific images, a YAML (YAML Ain't Markup Language) configuration file is crucial. YAML is a human-readable and simple data format designed to be easy for humans to understand and write while being easily interpretable by machines.

In this particular project, the YAML file outlined the application of various image filters, including *original images*, *wavelet transforms*, as well as *square*, *square root*, *logarithm*, *exponential*, *gradient*, and *local binary patterns (LBP 2D and 3D)*. The configuration also detailed the types of features to be extracted, covering *shape-based*

*metrics, first-order statistics, and texture-related features such as the gray level co-occurrence matrix, gray level run length matrix, gray level size zone matrix, gray level dependence matrix, and neighboring gray tone difference matrix.* This setup ensured a comprehensive analysis by applying multiple image transformations and extracting a broad range of features to fully capture different aspects of the image data.

The YAML configuration file is shown in Annexes: Code 2.

### 3. Feature extraction using PyRadiomics

	Number of features
<b>Original</b>	107
<b>Wavelet</b>	744
<b>Square</b>	93
<b>Square Root</b>	93
<b>Logarithm</b>	93
<b>Exponential</b>	93
<b>Gradient</b>	93
<b>LBP_2D</b>	93
<b>LBP_3D</b>	279
<b>Total radiomic features:</b>	<b>1688</b>

Table 3.5: Radiomic Feature Filters and their corresponding number of extracted features. LBP: Local Binary Patterns.

Using the *PyRadiomics* library, 1688 features were extracted from the CT images to quantify characteristics relevant to predicting treatment responses.

Shape-based features were applied exclusively to the original CT images. These features describe geometric properties such as volume, surface area, and sphericity. Since they are focused solely on the size and form of the segmented region, shape-based features are limited in number and do not capture variations in pixel intensity or texture.

In contrast, texture features extracted using various filters were more numerous and provided a more comprehensive analysis. These features capture patterns and distributions of gray levels, analyzing the relationships and variations in pixel intensities within the segmented region.

Among these, the wavelet filter significantly increased the number of extracted features by decomposing the image into eight sub-bands, each representing different levels of detail and orientation—specifically, low-low, low-high, high-low, and high-high frequencies in both horizontal and vertical directions. This decomposition allowed for texture features to be extracted from each sub-band, effectively multiplying the number of features by eight.

Similarly, the LBP 3D filter contributed to the large number of features extracted, yielding a total of 279 features. This is attributed to its extensive exploration of local texture patterns across various spatial scales in three dimensions.

The code implemented to extract radiomic features is appended in Annexes, in Code 3.

### 3.3 Data preparation

As mentioned by *Shrivastava et al.*, ‘If the data is inadequate, or contains extraneous and irrelevant information, Machine Learning (ML) algorithms may produce less accurate and less understandable results. Thus, data preparation is an important step in the ML process.’ [19]

1. **Outliers:** The first step was to analyse our clinical and metabolomic data for any obvious outliers. Outliers may arise from human error during manual data entry or due to instrumental errors. In both cases, the most appropriate approach is to manually remove these data points when it is certain they are incorrect and treat them as null values [20]. This ensures that they do not adversely affect the accuracy of our model. For radiomic features, identifying outliers was not applicable, as there are no established criteria for defining outliers in this type of data. Therefore, our outlier analysis was focused solely on clinical and metabolomic data.
2. **Columns with fixed values:** Variables that remain constant across patients do not contribute to differentiating them. Therefore, removing these columns was an important step to take before proceeding with preprocessing. This ensures that these columns are not considered from the outset and do not affect feature selection.
3. **Missing values:** When working with data, it is crucial to properly handle missing values, as they can negatively impact the performance of a PM. A straightforward option is to replace them with the mean (for numerical variables) or the mode (for categorical variables) of the column. However, this method does not account for relationships between variables and can introduce bias if many values are missing. A more advanced approach is the *Iterative Imputer* function from the Scikit-learn library [21]. This technique, as its name suggests, iteratively predicts missing values by treating the column with missing values as the target variable in each iteration while the other columns, whether they have missing values or not, serve as predictors to train a regression model. Finally, through successive iterations, this process refines the imputations based on the relationships learned from the other columns, leading to more accurate and realistic imputations.

Additionally, we have used a threshold of 20% of missing values to consider that a column should be removed, as they can severely impact data quality and introduce bias, compromising the effectiveness of the model. This is exemplified by the *Viral load baseline* variable, as shown in Figure 5 and quantitatively observed in Table 3.4, where a significant proportion of the data was missing, indicated by white areas.

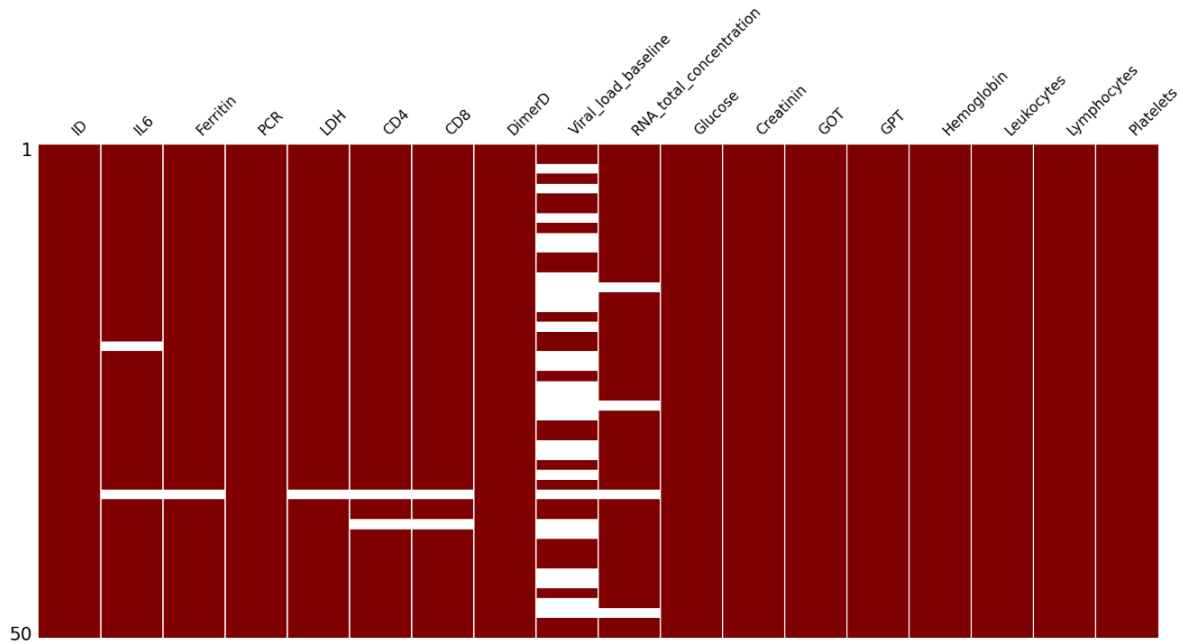


Figure 5: Missing values representation in Metabolomic data

Moreover, it was important to reassess the percentage of missing values after removing outliers, as this can potentially increase the proportion of missing data beyond the 20% threshold.

4. **Categorical features:** Transforming categorical variables is a crucial step in preparing data for ML models, as many of these models require numerical inputs and cannot directly process categorical data. One commonly used technique for this transformation is *one-hot encoding* from the Scikit-learn library [22]. This technique converts a categorical variable into multiple binary columns, with each column corresponding to a possible category. In each record, one of these columns will take the value 1 to indicate the presence of that specific category, while the other columns are set to 0. By applying this transformation, categorical variables (in Table 6a) were effectively converted into numerical data (in Table 6b), which enhances the models' ability to interpret and process the different categories present in the dataset.

GDS	
<b>Patient 1</b>	2
<b>Patient 2</b>	5

(a) Categorical feature

	GDS 1	GDS 2	GDS 3	GDS 4	GDS 5	GDS 6
<b>Patient 1</b>	0	1	0	0	0	0
<b>Patient 2</b>	0	0	0	0	1	0

(b) Categorical feature transformed binary columns

Figure 6: Two tables showing the categorical features transformation steps.

5. **Normalization:** Data normalization is crucial to ensure that all features contribute equally to the ML model, preventing features with larger scales from dominating the learning process. In this case, the *MinMaxScaler* from the Scikit-learn library [23] method was used, which transforms each feature to a defined range, typically between 0 and 1, according to the formula shown in Equation 3.1. This normalization ensures that all features are on a comparable scale, preventing biases and facilitating a balanced contribution of each feature during model training.

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}} \quad (3.1)$$

Data preparation codes are added in Annexes: Codes 4 (for clinical data) and 5 (for metabolomic data). Radiomic data only needed to be normalized as values where free from outliers, missing values and categorical features.

### 3.4 Data augmentation

As shown in Figure 7, in our dataset, we had 34 patients who exhibited a positive response to the treatment, meaning they recovered, demonstrating the treatment's efficacy. However, the remaining 16 patients, out of a total of 50, showed no signs of improvement within a three-day period and subsequently passed away during this time frame. As we can see, our dataset was significantly imbalanced between the classes. If we had trained the model with this data, we would have obtained a model biased towards the majority class, as demonstrated in previous research [24], which highlights the importance of balancing imbalanced data to improve classification accuracy in ML.

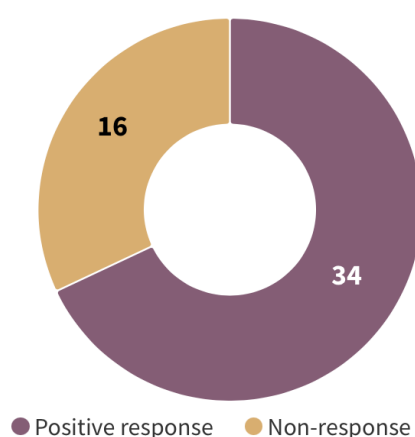


Figure 7: Distribution of treatment response between classes.

One common approach to addressing imbalanced data classes is to reduce the size of the majority class to match that of the minority class. Nevertheless, given the limited number of patients in our dataset, this method was not ideal, as it would result in a significant loss of valuable information.

An alternative and more suitable strategy involved augmenting the minority class by generating synthetic data based on the existing information. This approach allowed us to maintain the integrity of the dataset while mitigating the effects of class imbalance, leading to a more balanced and representative dataset for model training.

#### 1. Clinical and metabolomic data augmentation:

To balance clinical and metabolomic data, which are numerical in nature, a well-known oversampling method called *SMOTE* (*Synthetic Minority Oversampling TEchnique*) was employed as it has been demonstrated that in low-dimensional data, the use of this technique is beneficial [25].

SMOTE operates by randomly selecting nearby neighbours for each sample in the minority class, using a distance metric such as Euclidean distance. The number of selected neighbors is an adjustable parameter that influences the diversity of the generated synthetic samples. To create a new sample, a nearby neighbor is chosen, and the difference between the features of the original sample and those of the neighbor is calculated. This difference is then multiplied by a random number between 0 and 1, and the resulting value is added to the original features, thereby generating a new synthetic sample that lies somewhere between the original sample and its neighbour. This process is repeated until the minority class reaches the same size as the majority class [26].

When augmenting data, we needed to observe that the synthetic data followed the same Gaussian distribution as the original data. This is important because it ensures that the synthetic data accurately represents the same characteristics and maintains the same variability as the original dataset.

Figure 8 illustrates the Gaussian distribution observed for the metabolomic variable *Creatinine*.

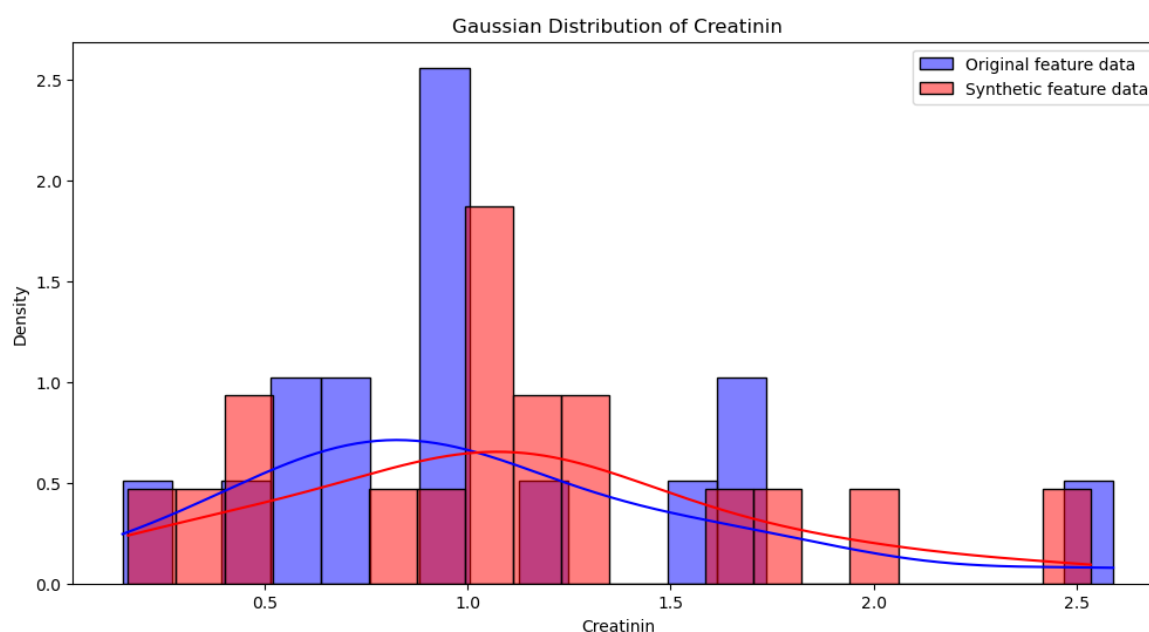


Figure 8: Gaussian distribution with Creatinin feature.

The codes used to augment clinical and metabolomic data are included in Annexes: Code 6.

## 2. Image augmentation:

In the context of image data augmentation, we needed to generate synthetic images rather than merely creating synthetic features from the original features extracted, as this approach generally yields better model performance. To achieve this, we utilized the *PyTorch* library to apply various transformations to the original images (see Figure 9a), thereby creating new, synthetic variations (see Figure 9b).

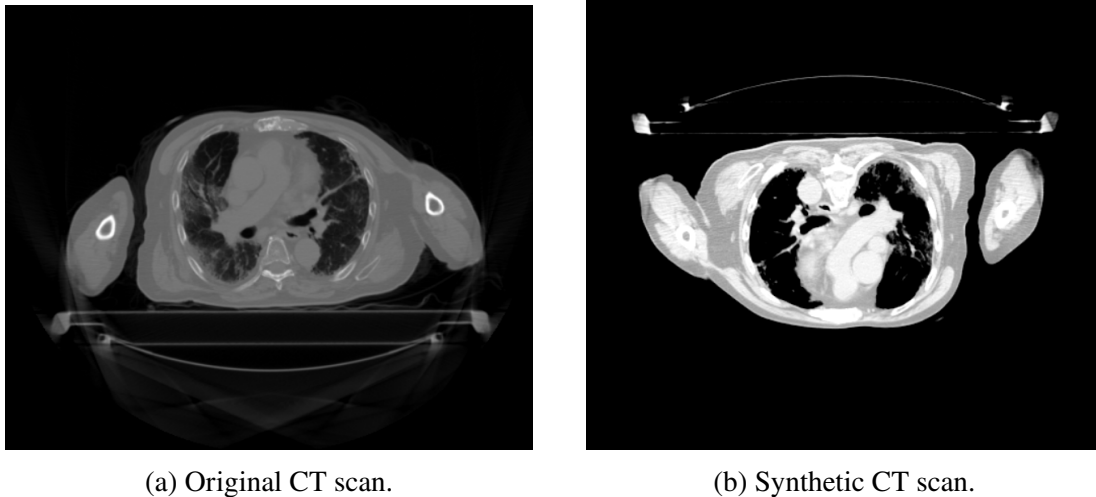


Figure 9: Original and synthetic Computed Tomography scans comparison.

The transformation process involved several random alterations to the images, including the addition of *noise*, changes in *brightness*, and *flipping*.

Additionally, it was imperative to adjust the mask associated with each image in accordance with the flipping transformations applied to the synthetic images. This adjustment ensures that the features extracted from the transformed images matches information from the lungs. In Figure 10 an example of a mask image produced through these transformations is provided below (see Figure 10a), alongside its original counterpart (see Figure 10a) for comparison.



Figure 10: Original and synthetic mask comparison.

The codes performed to obtain synthetic images are appended in Annexes in Code 7.

### 3.5 Data preprocessing

Selecting only the most relevant features decreases the model's complexity, which in turn speeds up training and makes the model more interpretable. Additionally, it helps eliminate noise from irrelevant features, confirming that the model focuses on the most informative data.

Two different methods for selecting the necessary features were applied and compared:

- **Manual selection:** The feature extraction process was conducted manually, involving the deliberate selection of individual methods to filter and refine the dataset. This approach required careful consideration and choice of each technique to guarantee the most relevant features were identified and retained for further analysis.

1. **Correlation evaluation:** Manually selecting features involved evaluating the correlation between features and the target variable to identify and remove those with weaker relationships. This approach helps to reduce the risk of biased predictions that may arise due to feature correlation. In this process, we used Pearson's correlation coefficient  $r$  given by the following Equation 3.2:

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) \left(\sum_{i=1}^n (Y_i - \bar{Y})^2\right)}} \quad (3.2)$$

Where:

- $X$  and  $Y$  are two features in the dataset.
- $\bar{X}$  and  $\bar{Y}$  are the means of  $X$  and  $Y$ .

2. **Variance Inflation Factor (VIF) calculation:** Additionally, calculating the VIF (see Equation 3.3) was essential for statistically measuring the extent of multicollinearity among predictor variables. Features with a VIF exceeding a threshold of 5 indicate excessive correlation with other features and should be removed to improve model interpretability [27].

$$\text{VIF}_i = \frac{1}{1 - R_i^2} \quad (3.3)$$

Where:

- $R_i^2$  is the value from regressing the  $i$ -th variable on all the other variables.

This approach ensures that the remaining features contribute valuable and distinct information to the model.

The code for manual feature selection is included in the Annexes: Code 8.

- **Automatic selection:** In parallel with the manual feature selection process, an automatic feature selection method was implemented using the mRMR (minimum Redundancy Maximum Relevance) algorithm, specifically utilizing the MIQ (Mutual Information Quotient) variant. The mRMR method filters features considering both the relevance for predicting the outcome variable and the redundancy within the selected features, performing less memory consumption and time required to achieve computation efficiency [28].

The code for automatic feature selection is shown in Annexes: Code 9.

### 3.6 Predictive models

Once the data was cleaned, preprocessed, and features were selected, the next step was to implement various predictive models (PMs) based on the input data and the three classifiers used. The steps to be followed in this phase are illustrated in the Figure 11:

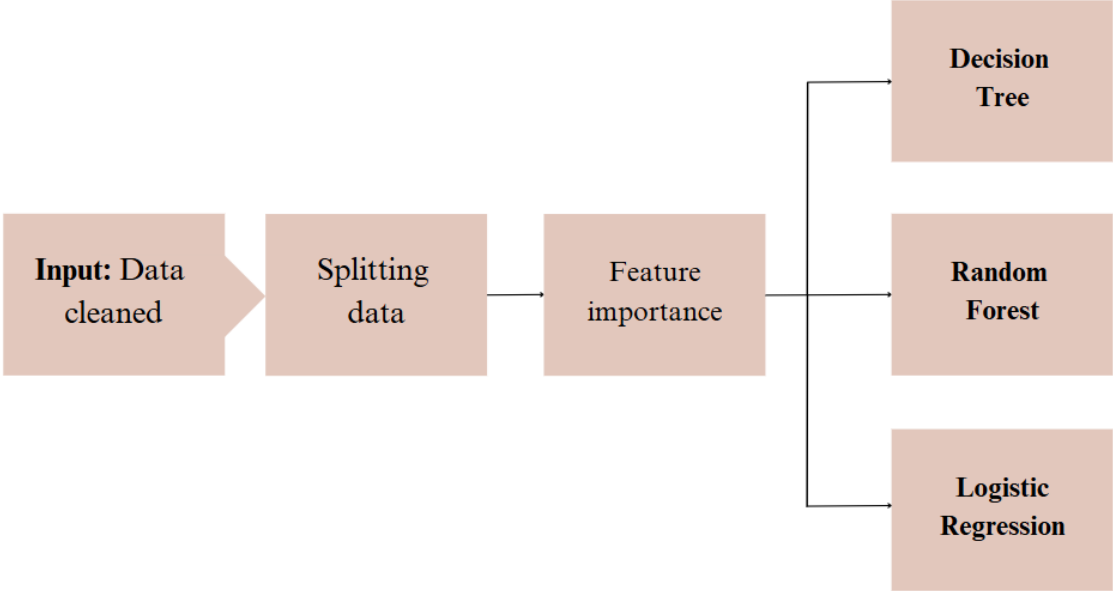


Figure 11: Predictive Model workflow representation.

#### *Splitting data*

Before training the models, the dataset was split into two subsets: 90% for training and 10% for testing. Due to the limited amount of data, *Cross-Validation (CV)* was implemented within the 90% training subset to maximize data utilization and provide a more reliable estimate of the model’s generalization capability.

CV involves dividing the training set into multiple smaller subsets called *folders*. This process was repeated several times (in this case, 5 folds), where the model was trained on a subset of the folders and validated on the remaining fold (which corresponded to 20% of the training data). Each fold was used as the validation set once throughout the process.

Additionally, the separate testing set was used to evaluate the performance of the best PM on unseen data.

A visual representation of this process is provided in Figure 12:

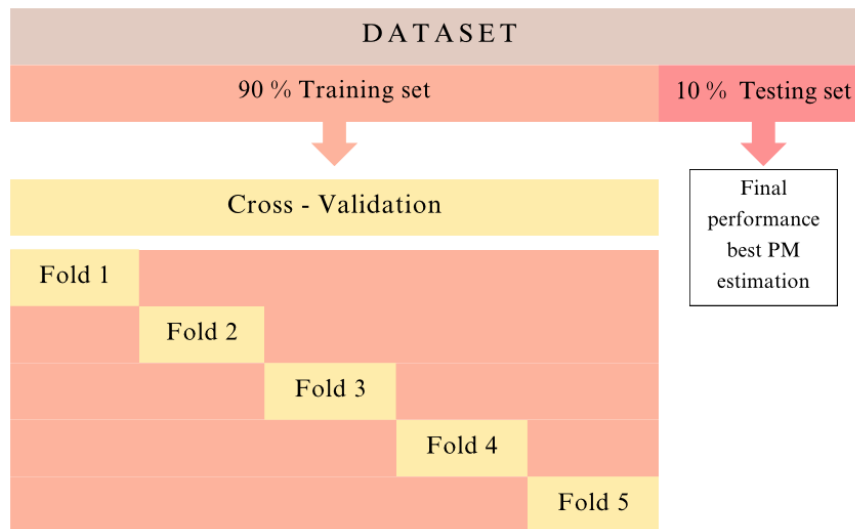


Figure 12: Splitting data process representation.  
PM: Predictive Model.

### ***Feature Importance***

To enhance the accuracy and efficiency of the models, feature importance weighting was conducted using a forest of trees (Random Forest). This model is particularly effective for classification tasks and offers valuable insights into the significance of each feature in predicting the target variable. By evaluating the weights assigned to each feature [29], we could identify the most impactful ones, 9 in this case, allowing us to reduce the dimensionality of the dataset and potentially improve model performance.

### ***Model Implementation***

In this analysis, three widely recognized classification models were employed due to their effectiveness in supervised learning tasks: *Random Forest (RF)*, *Logistic Regression (LR)*, and *Decision Trees (DT)*. These models are well-known for their robust performance in classification tasks, making them suitable candidates for our comparative study. The models were evaluated and compared using various performance metrics to determine which one provided the best results for our specific application.

- **DT** is a classification model that works by splitting the data into subsets based on feature values that best separate the classes. At each node, the algorithm chooses a feature and a threshold that maximizes class separation (using criteria like Gini impurity or entropy). This process continues until the data is sufficiently pure (mostly one class) at the leaf nodes, where the final class prediction is made. DTs are easy to interpret but can overfit the data if the tree becomes too complex [30].

- **RF** is an ensemble classification method that builds multiple DTs using random subsets of the data and features. Each tree in the forest makes its own classification, and the final output is determined by majority voting among all the trees. This approach reduces the variance and overfitting associated with individual decision trees, leading to more accurate and stable predictions, though it sacrifices some interpretability [31].
- **LR** is a linear classification model that predicts the probability that a given input belongs to a specific class. It does this by fitting a linear equation to the input features and applying a logistic (sigmoid) function to convert the output into a probability between 0 and 1. Based on this probability, the model assigns the input to one of the two classes. LR is simple and effective for linearly separable data but may struggle with more complex, non-linear relationships [32].

Finally, the model that demonstrated the best performance was selected to make the final predictions on the test set using OPTUNA as an hyperparameter tuning optimization [33], aiming to maximize specificity and generalization capability for future applications.

All the PM codes are provided in Annexes: Codes 10 (RF PM code), 12 (DT PM code), 11 (LR PM code) and 13 (Best PM testing).

### 3.7 Model validation

The models' performance was evaluated using the following metrics:

1. **Sensitivity:** Also referred to as, recall, this metric indicates the proportion of True Positives (TP) that the model has correctly identified. It is calculated as:

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (3.4)$$

2. **Specificity:** Specificity measures the proportion of True Negatives (TN) that the model has correctly identified. The formula is:

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (3.5)$$

3. **Confusion Matrix (CM):** The confusion matrix is a table that visually summarizes the performance of the classification model by showing TP, False Positives (FP), TN, and False Negatives (FN). An example of a CM is in the Table 3.6:

	Positive Prediction	Negative Prediction
Actual Positive	TP	FN
Actual Negative	FP	TN

Table 3.6: Confusion Matrix Representation

## 4. Results

In this section, the efficacy on predicting the response of COVID-19 patients treated with LD-RT was evaluated. For this, data augmentation techniques, various feature selection methods and three different classification algorithms were implemented.

### *Original Data Results*

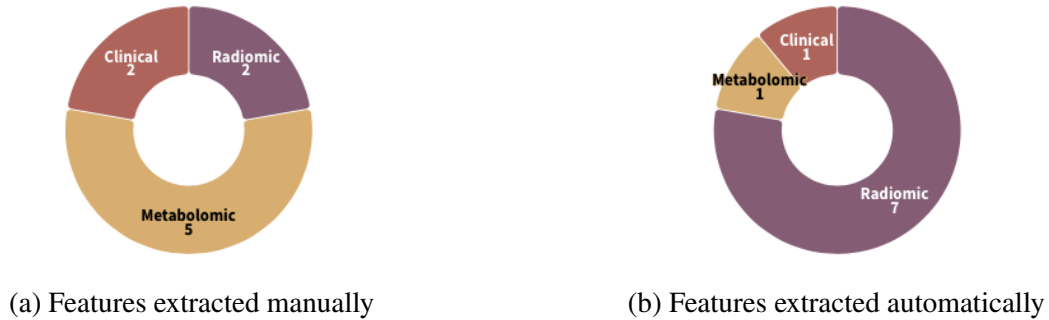


Figure 11: Feature Importance Types in Original Data.

For the original data, which was imbalanced and lacked data augmentation techniques, the number of manually selected features from each group based on their importance is represented in Figure 11a. Additionally, for the automatic feature selection, the selected data type representation is visualized in Figure 11b. The features, ordered by their importance in predicting the target variable, are shown in Annexes in supplementary Figures 15 and 16, respectively.

The results obtained using only the original imbalanced data are presented as follows in Table 4.1, distinguishing between manual and automatic feature selection:

	Features extracted manually		Features extracted automatically	
	Sensitivity	Specificity	Sensitivity	Specificity
<b>Decision Tree</b>	0.65	0.43	0.71	0.50
<b>Random Forest</b>	0.87	0.57	0.90	0.43
<b>Logistic Regression</b>	0.87	0.57	1.00	0.00

Table 4.1: Imbalanced data results

We observed that sensitivity, which is the rate of true positives, was generally high, with the RF classifier using automatically extracted features performing the best in this metric. However, specificity remained very low across all results obtained.

### Augmented Data Results

Next, we examined the results obtained with the augmented data, using a similar approach as before.



(a) Features extracted manually

(b) Features extracted automatically

Figure 12: Feature Importance Types in Augmented Data.

The data type of the important features identified were, in the case of manually selected features, the represented in Figure 12a and for the automatically selected features (see Figure 12b). The corresponding features importance, ordered by their weight in predicting the target variable, are shown in Annexes in supplementary Figures 17 and 18, respectively.

All radiomic features selected are explained in Annexes: Table 8.

The results obtained with this data input are as follows:

	Features extracted manually		Features extracted automatically	
	Sensitivity	Specificity	Sensitivity	Specificity
<b>Decision Tree</b>	0.68	0.77	0.74	0.77
<b>Random Forest</b>	0.84	0.77	0.81	0.83
<b>Logistic Regression</b>	0.68	0.70	0.67	0.60

Table 4.2: Augmented data results

With data augmentation to address the class imbalance issue, better results were achieved compared to not addressing the class disparity. The best performance, as before, was obtained with the RF classifier using automatically selected features, achieving a sensitivity of 0.81 and a specificity of 0.83.

### Testing Set Results

As mentioned earlier, the best-performing model underwent hyperparameter tuning to maximize its effectiveness, utilizing 10% of the data as a testing set to evaluate performance on unseen data. This model achieved sensitivity and specificity values of 1.00 and 0.50, respectively.

All CM from the three classifiers are represented in Annexes in Figures 27 to 31.

## 5. Discussion

The current study introduces a novel approach by integrating clinical, metabolomic, and radiomic data to predict responses to LD-RT treatment in patients with COVID-19 pneumonia. Although predictive modeling using such diverse data sources is not entirely new, our study provides unique insights by focusing on the specific application of these techniques to predict LD-RT response of COVID-19 pneumonia patients.

A critical element of this study was the use of the PyRadiomics library for extracting radiomic features from CT images. While 3D Slicer also implements PyRadiomics, it lacks the extensive image filtering capabilities available through the library itself. The direct use of PyRadiomics enabled us to extract a greater number of features through multiple image transformations, revealing patterns that would otherwise remain undetectable. This advanced feature extraction process significantly contributed to enabling automation and scalability as well as it allowed to identifying impactful features and enhancing model performance.

Our analysis revealed that radiomic features were predominantly selected across various models, largely due to their higher quantity compared to clinical and metabolomic features. Conversely, clinical features were less emphasized in most PM, suggesting that, despite their potential relevance, they did not contribute as significantly to the model's performance in this context. This suggests that although clinical variables are currently the primary factors used to plan RT treatment and predict patient response, other data, such as metabolomic and radiomic features, should also be considered.

Our findings reveal a notable trade-off between sensitivity and specificity across different models. Specifically, the RF classifier, when applied to the original imbalanced data with automatically extracted features, achieved a high sensitivity of 0.97 but exhibited low specificity at 0.50. Given its low specificity, this model was deemed unsuitable as we aim to accurately identify TP and avoid administering treatment to patients who are unlikely to benefit from it. To address the class imbalance, a data augmentation methodology was applied, which improved the results: the RF classifier achieved a sensitivity of 0.81 and a specificity of 0.83. This enhancement underscores the effectiveness of data augmentation in improving model performance and ensuring better identification of negative cases.

In both cases, the best-performing model utilized automatically selected features using the mRMR method, which also demonstrated a lower emphasis on clinical features than the manual method. Moreover, this approach not only yielded general superior results in terms of sensitivity and specificity but also significantly reduced the time required for feature filtering, compared to manual feature selection.

When evaluated on the test set, the RF classifier achieved a sensitivity of 1.00 but a specificity of only 0.50, highlighting the impact of the limited amount of training data on the model's ability to generalize to unseen data. The small dataset size was a significant constraint in our study, which is a common challenge in ML applications, particularly in specialized healthcare studies.

One notable limitation of the current study is the restricted number of data samples, which impacts the model's performance and generalizability. Despite this limitation, we have demonstrated that predictive models can be useful and have significant potential. The high sensitivity observed indicates that the model has the capability to correctly identify responses to LD-RT treatment in a substantial proportion of cases. However, the low specificity suggests that improvements are needed to reduce FP.

At this juncture, future work should undoubtedly focus on expanding the dataset by collaborating with additional hospital centers willing to administer LD-RT to their most severely affected COVID-19 patients. This expansion is essential to enhance the results and develop a reliable and unbiased predictive model. Once this significant step is achieved, the next phase involves implementing the model in a practical clinical setting. Specifically, this entails integrating the model as a clinical support system for radiotherapy departments through a web-based implementation hosted on the hospital's server. This will provide doctors with a valuable tool for improving patient care and treatment decision-making.

## 6. Conclusions

This Final Degree Project emphasizes the intricate and sensitive nature of PMs within health data analysis. It has become evident that even minor oversights or inconsistencies in the dataset can significantly alter the results of these models. For PM to achieve their full potential, they require input data that is not only reliable and clean but also well-balanced and sufficiently voluminous. The quality and quantity of the data directly impact the efficacy and accuracy of the model's predictions.

The primary objective of this project was to develop a PM capable of determining the response to LD-RT treatment in patients with COVID-19 pneumonia. Although the results were constrained in the test set due to the low dimensionality of the data, reliable outcomes were achieved in the training set. Nonetheless, they did not achieve the level of specificity in the test set, necessary for implementation in clinical practice.

Despite these challenges, the insights provided by PMs can be invaluable, depending on the context and situation. The ability to foresee potential outcomes before they occur represents a transformative advancement in medicine. However, it is crucial to acknowledge that such models are not infallible and that some degree of error is inevitable. For this reason, they should never replace clinical decisions but rather support them.

To fully realize the potential of predictive modeling in healthcare, there is an urgent need to advance the integration of technological and biomedical engineering expertise into healthcare environments. Professionals from diverse fields, including telecommunications engineers, computer scientists, and physicists, should collaborate closely with healthcare practitioners. This interdisciplinary approach is essential for driving innovations and breakthroughs in medical data analysis and for paving the way toward a more effective and efficient future in healthcare. Moreover, leveraging these predictive models can lead to significant cost savings by optimizing resource allocation, reducing unnecessary treatments, and improving patient outcomes.

In conclusion, while the road to integrating predictive modeling into healthcare is fraught with challenges, it holds the potential of revolutionizing healthcare. By ensuring high-quality data, fostering interdisciplinary collaboration, investing in technological advancements, AI can become an essential component of future medical practice.

## References

- [1] N. Zhu, D. Zhang, W. Wang, *et al.*, “A novel coronavirus from patients with pneumonia in china, 2019,” *New England journal of medicine*, vol. 382, no. 8, pp. 727–733, 2020.
- [2] W. H. Organization, *Total cumulative number of covid-19 cases reported to who*, <https://data.who.int/dashboards/covid19/cases>, Accessed on the 10th of April, 2020, 2024.
- [3] M. Arenas, B. Piqué, L. Torres-Royo, *et al.*, “Treatment of covid-19 pneumonia with low-dose radiotherapy plus standard of care versus standard of care alone in frail patients: The seor-gicor ipacovid comparative cohort trial,” *Strahlentherapie und Onkologie*, vol. 199, no. 9, pp. 847–856, 2023.
- [4] M. Arenas, S. Sabater, V. Hernández, *et al.*, “Anti-inflammatory effects of low-dose radiotherapy,” *Strahlentherapie und Onkologie*, pp. 1–7, 2012.
- [5] A. Jain, H. Patel, L. Nagalapatti, *et al.*, “Overview and importance of data quality for machine learning tasks,” in *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 3561–3562.
- [6] B. Piqué, K. Peña, F. Riu, *et al.*, “Sars-cov-2 serum viral load and prognostic markers proposal for covid-19 pneumonia in low-dose radiation therapy treated patients,” *Journal of clinical medicine*, vol. 12, no. 3, p. 798, 2023.
- [7] K. B. Méndez-Rodríguez, M. J. Santoyo-Treviño, K. Saldaña-Villanueva, M. Rodríguez-Aguilar, R. Flores-Ramírez, and F. J. Pérez-Vázquez, “Metabólica como nueva herramienta para el diagnóstico oportuno en enfermedades no transmisibles,” *Revista de Salud Ambiental*, vol. 19, no. 2, pp. 109–115, 2019.
- [8] M. E. Mayerhoefer, A. Materka, G. Langs, *et al.*, “Introduction to radiomics,” *Journal of Nuclear Medicine*, vol. 61, no. 4, pp. 488–495, 2020.
- [9] K. Lekadir, A. Feragen, A. J. Fofanah, *et al.*, *Future-ai: International consensus guideline for trustworthy and deployable artificial intelligence in healthcare*, 2024. arXiv: 2309.12325 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/2309.12325>.
- [10] Varian Medical Systems, *Eclipse treatment planning system*, [Accessed: 4-Aug-2024], 2024. [Online]. Available: <https://www.varian.com/products/radiotherapy/treatment-planning/eclipse>.
- [11] S. Venegas, M. Cortés, L. Flores, *et al.*, “Correlación de spo<sub>2</sub>/fio<sub>2</sub> versus pao<sub>2</sub>/fio<sub>2</sub> para monitoreo de la oxigenación en pacientes con trauma de tórax,” *Med Crit*, vol. 32, no. 4, pp. 201–207, 2018. DOI: 10.35366/TI184E.
- [12] Lab Tests Guide, *Po<sub>2</sub> normal values*, [Accessed: 4-Aug-2024], 2024. [Online]. Available: <https://www.labtestsguide.com/po2#po2-normal-values>.

- [13] M. Moser, M. Lehofer, G. Hildebrandt, M. Voica, S. Egner, and T. Kenner, “Phase- and frequency coordination of cardiac and respiratory function,” *Biological Rhythm Research*, vol. 26, no. 1, pp. 100–111, 1995.
- [14] MSD Manual, *Pruebas de sangre: Valores normales*, [Accessed: 4-Aug-2024], 2024. [Online]. Available: <https://www.msmanuals.com/es/professional/recursos/valores-normales-de-laboratorio/pruebas-de-sangre-valores-normales>.
- [15] 3D Slicer, *Documentation: Radiomics extension*, [Accessed: 4-Aug-2024], 2024. [Online]. Available: <https://www.slicer.org/wiki/Documentation/4.8/Extensions/Radiomics>.
- [16] J. P. F. Oliver, *Development and evaluation of machine learning models for COVID-19 related research*, Final Degree Project, Universitat Rovira i Virgili, 2023. [Online]. Available: <https://hdl.handle.net/20.500.11797/TFG7026>.
- [17] J. J. Van Griethuysen, A. Fedorov, C. Parmar, *et al.*, “Computational radiomics system to decode the radiographic phenotype,” *Cancer research*, vol. 77, no. 21, e104–e107, 2017.
- [18] X. Li, P. S. Morgan, J. Ashburner, J. Smith, and C. Rorden, “The first step for neuroimaging data analysis: Dicom to nifti conversion,” *Journal of neuroscience methods*, vol. 264, pp. 47–56, 2016.
- [19] H. Shrivastava and S. Sridharan, “Conception of data preprocessing and partitioning procedure for machine learning algorithm,” *International Journal of Recent Advances in Engineering & Technology (IJRAET)*, vol. 1, no. 3, pp. 2347–2812, 2013.
- [20] J.-F. Beaumont and L.-P. Rivest, “Dealing with outliers in survey data,” in *Handbook of statistics*, vol. 29, Elsevier, 2009, pp. 247–279.
- [21] S.-l. developers, *Sklearn.impute.iterativeimputer*, Accessed: 2024-08-20, 2024. [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.impute.IterativeImputer.html>.
- [22] Scikit-learn Developers, *Sklearn.preprocessing.onehotencoder*, Accessed: 2024-08-20, 2024. [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.OneHotEncoder.html>.
- [23] Scikit-learn, *Sklearn.preprocessing.minmaxscaler*, Accessed: 2024-08-03, 2024. [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html>.
- [24] V. S. Spelman and R. Porkodi, “A review on handling imbalanced data,” in *2018 international conference on current trends towards converging technologies (ICCTCT)*, IEEE, 2018, pp. 1–11.

- [25] R. Blagus and L. Lusa, “Smote for high-dimensional class-imbalanced data,” *BMC bioinformatics*, vol. 14, pp. 1–16, 2013.
- [26] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, “Smote: Synthetic minority over-sampling technique,” *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [27] J. I. Daoud, “Multicollinearity and regression analysis,” in *Journal of Physics: Conference Series*, IOP Publishing, vol. 949, 2017, p. 012 009.
- [28] Z. Zhao, R. Anand, and M. Wang, “Maximum relevance and minimum redundancy feature selection methods for a marketing machine learning platform,” in *2019 IEEE international conference on data science and advanced analytics (DSAA)*, IEEE, 2019, pp. 442–452.
- [29] Scikit-Learn, *Example of feature importance with random forest*, [Accessed: 4-Aug-2024], 2024. [Online]. Available: [https://scikit-learn.org/stable/auto\\_examples/ensemble/plot\\_forest\\_importances.html](https://scikit-learn.org/stable/auto_examples/ensemble/plot_forest_importances.html).
- [30] Scikit-learn, *Decisiontreeclassifier*, Accessed: 2024-08-25, 2024. [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>.
- [31] Scikit-learn, *Randomforestclassifier*, Accessed: 2024-08-25, 2024. [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>.
- [32] Scikit-learn, *Logisticregression*, Accessed: 2024-08-25, 2024. [Online]. Available: [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html).
- [33] O. Authors, *Optuna: A hyperparameter optimization framework*, Accessed: 2024-08-03, 2024. [Online]. Available: <https://optuna.readthedocs.io/en/v4.0.0/>.
- [34] J. Osuna, *Radiomics-ai*, Accessed: 2024-08-04, 2023. [Online]. Available: <https://github.com/JudithOsuna/Radiomics-AI>.

# Annexes

## I. Feature Importance Ranking

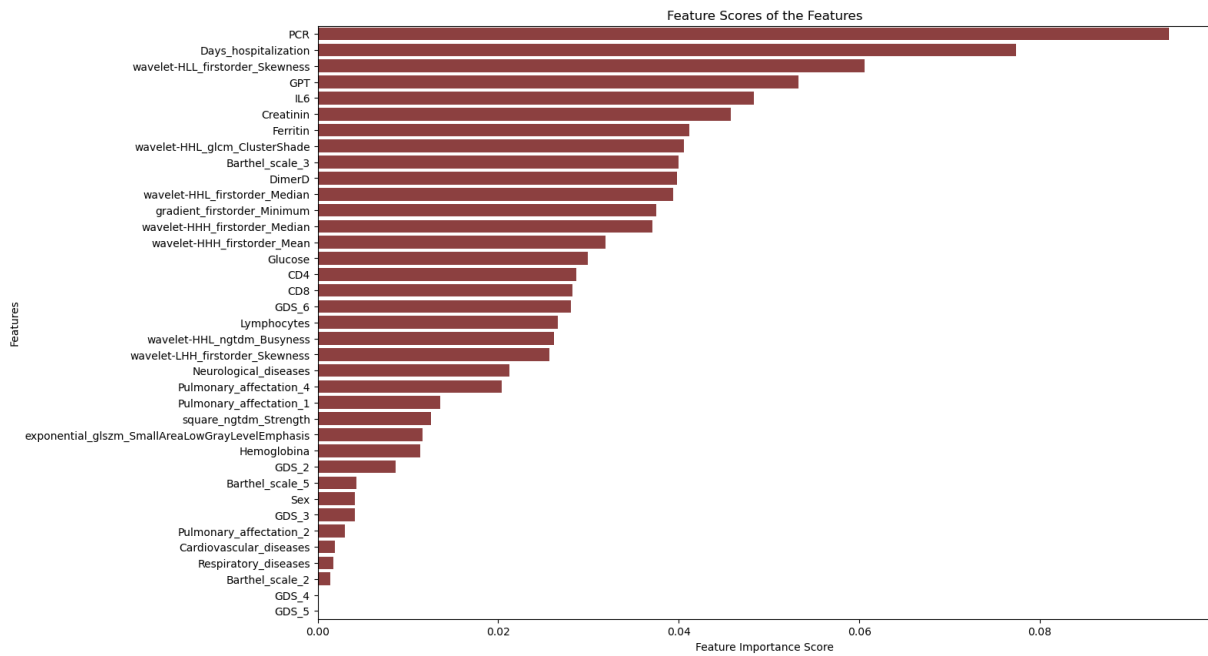


Figure 15: Feature importance with original and manually selected data

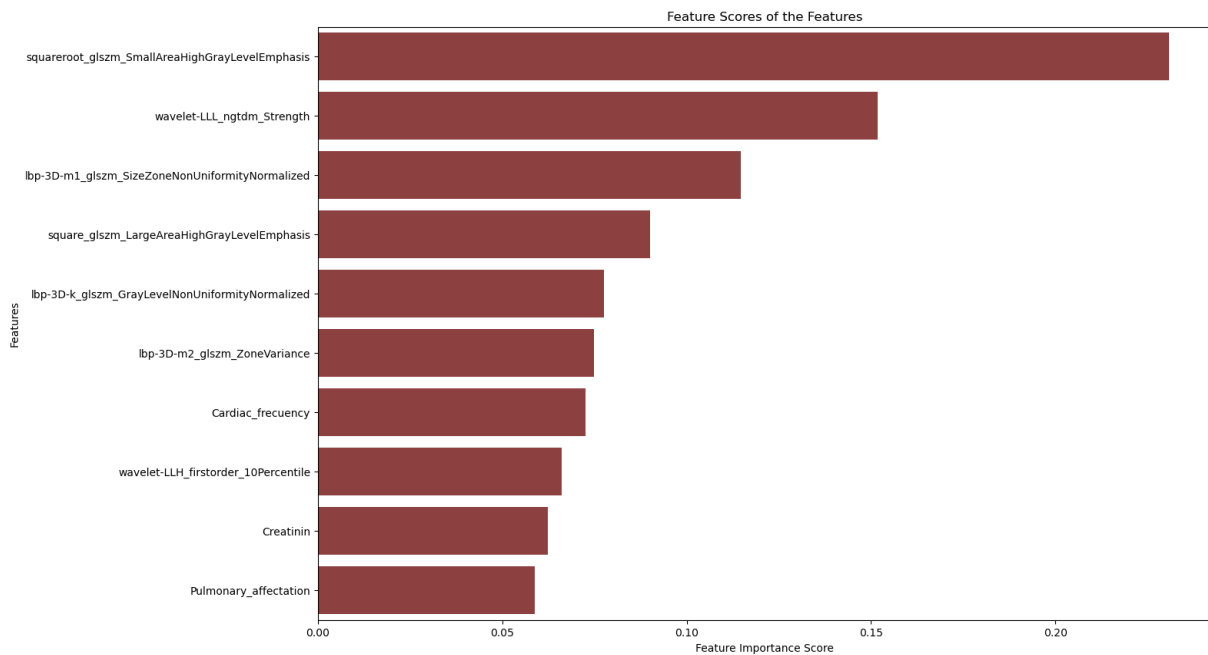


Figure 16: Feature importance with original and automatically selected data

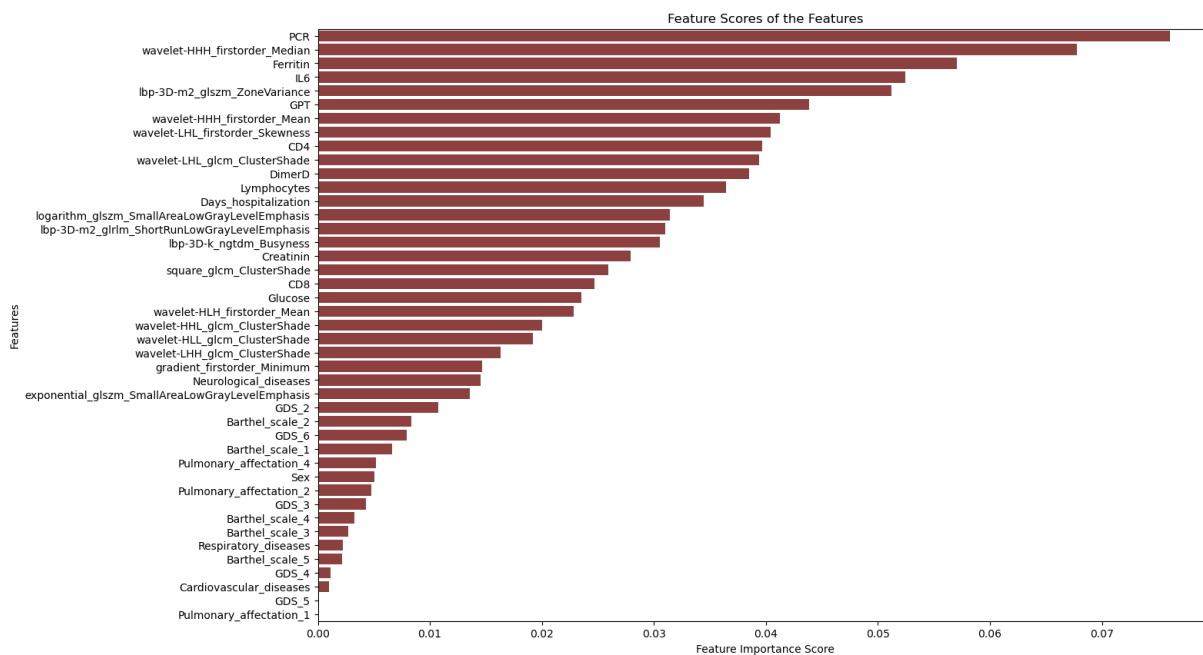


Figure 17: Feature importance with augmented and manually selected data

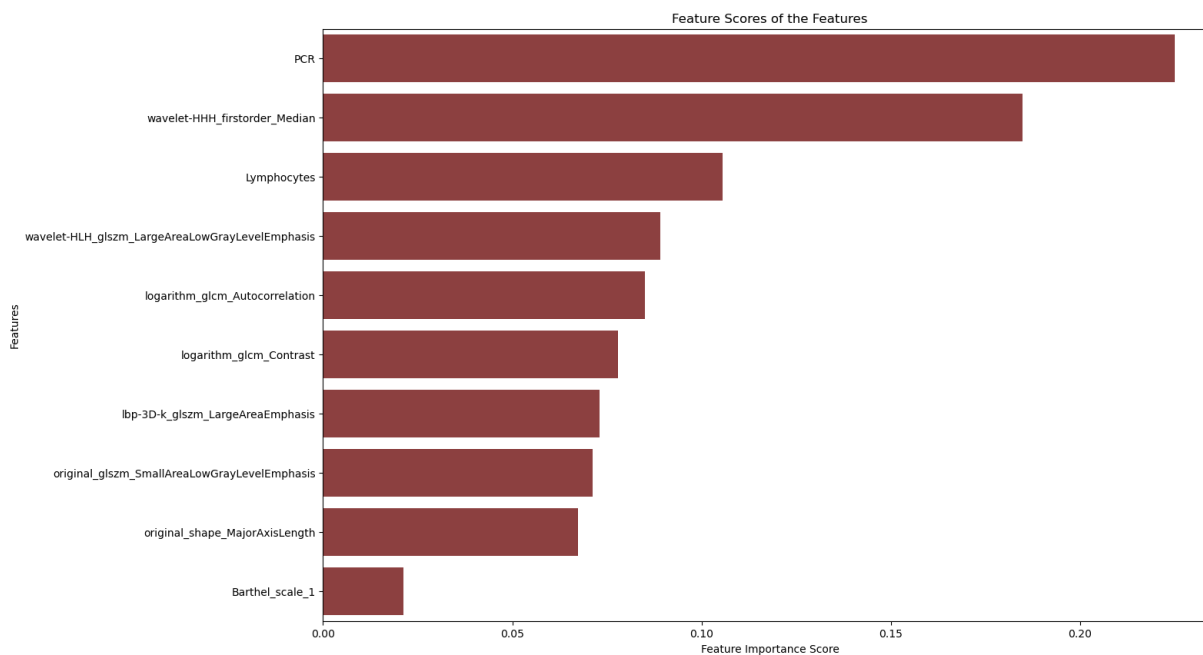


Figure 18: Feature importance with augmented and automatically selected data

## II. Confusion Matrices

### Decision Tree CM

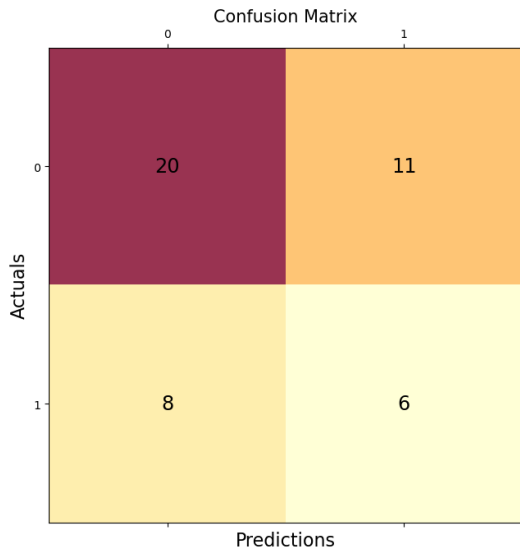


Figure 19: Confusion Matrix original data extracted manually

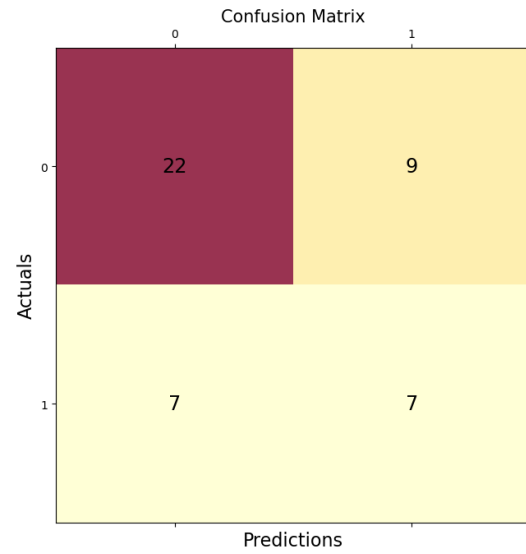


Figure 20: Confusion Matrix original data extracted automatically

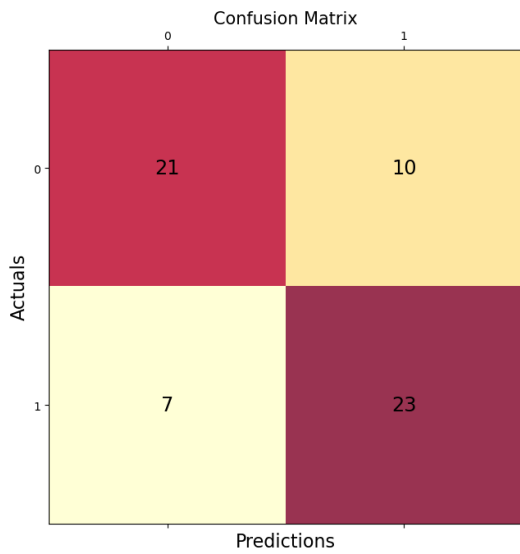


Figure 21: Confusion Matrix augmented data extracted manually

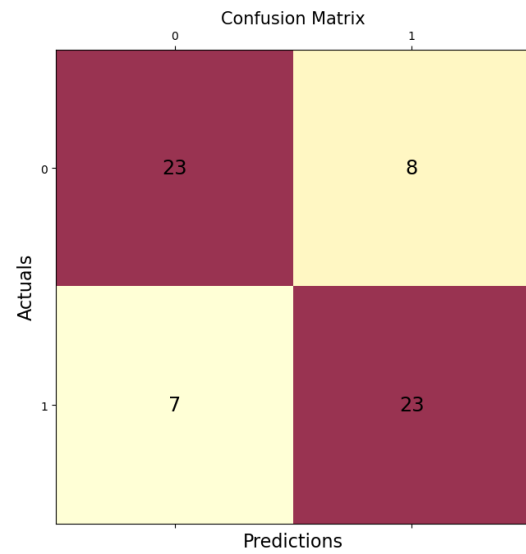


Figure 22: Confusion Matrix augmented data extracted automatically

### Random Forest CM

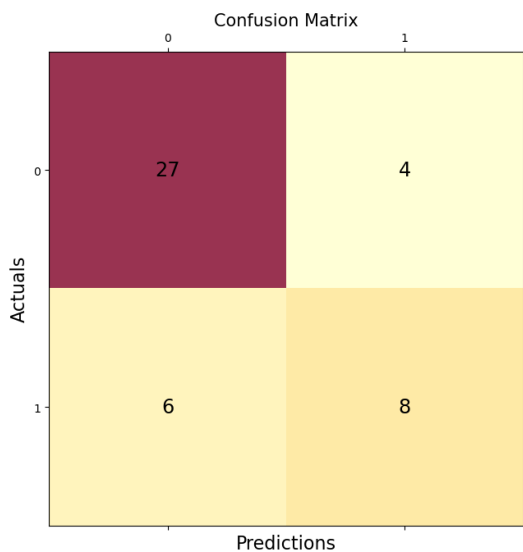


Figure 23: Confusion Matrix original data extracted manually

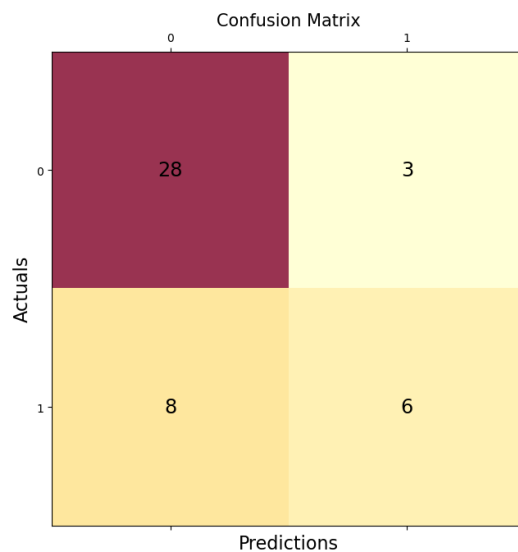


Figure 24: Confusion Matrix original data extracted automatically

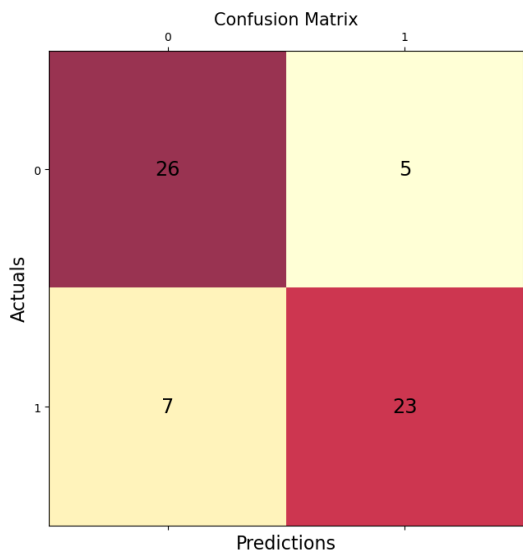


Figure 25: Confusion Matrix augmented data extracted manually

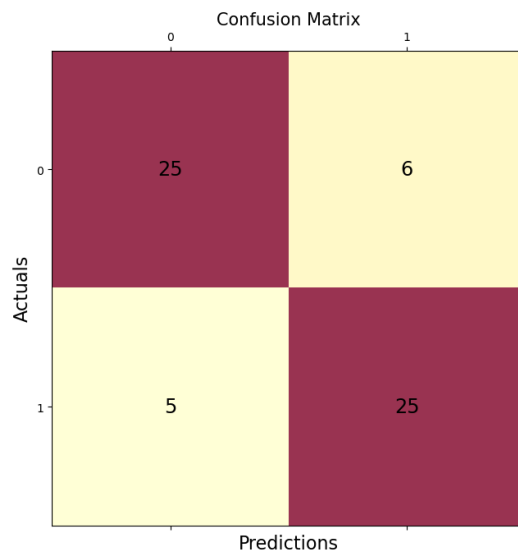


Figure 26: Confusion Matrix augmented data extracted automatically

## Logistic Regression CM

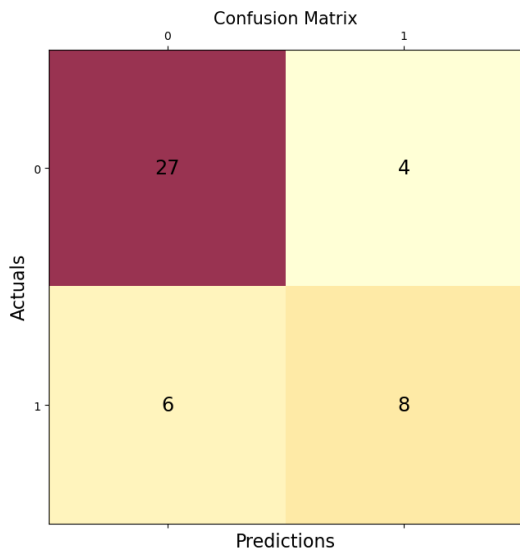


Figure 27: Confusion Matrix original data extracted manually

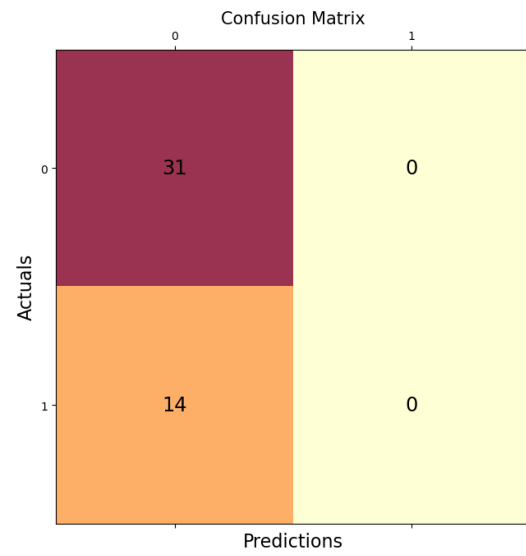


Figure 28: Confusion Matrix original data extracted automatically

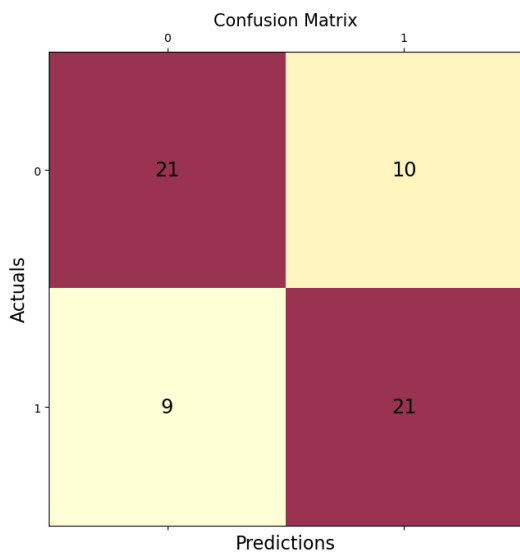


Figure 29: Confusion Matrix augmented data extracted manually

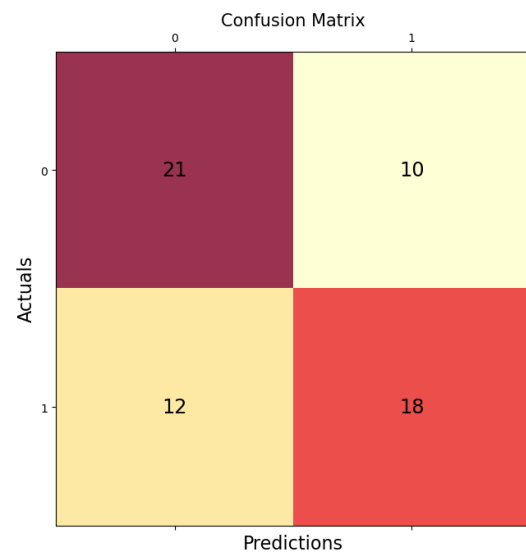


Figure 30: Confusion Matrix augmented data extracted automatically

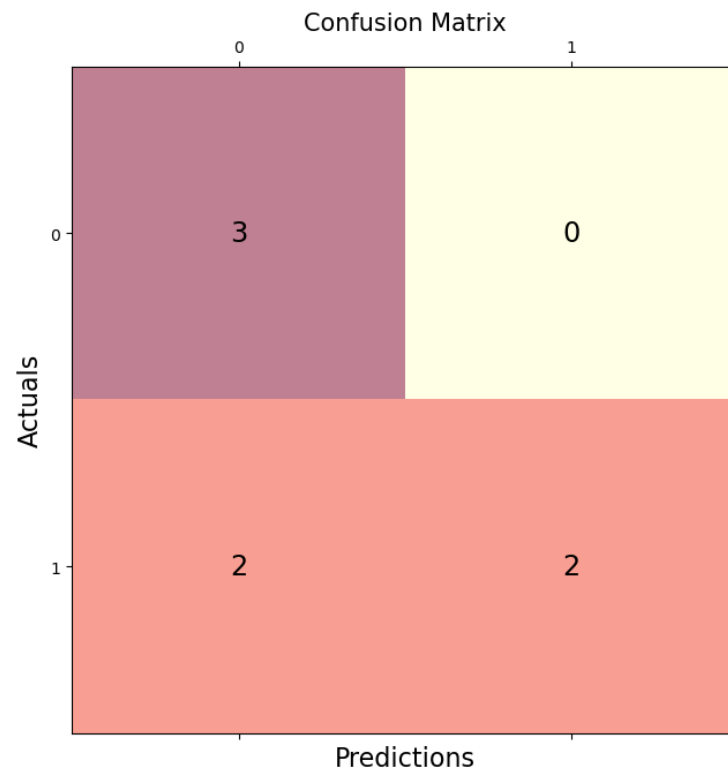
**Validation set CM**

Figure 31: Confusion Matrix validation set

### III. Radiomic features

Feature	Description
Wavelet HLL first order skewness	Measures the skewness of pixel intensity distribution in the HLL wavelet component.
Wavelet HHL glcm clustershade	Quantifies the shading effect of pixel clusters in the HHL wavelet component using GLCM.
Squareroot glszm small area high gray level emphasis	Emphasizes high gray levels in small areas of the GLSZM matrix.
Wavelet LLL ngtdm strength	Measures the strength of pixel gray-tone differences in the LLL wavelet component using NGTD.
Lbp 3d m1 glszm size zone non uniformity normalized	Assesses non-uniformity of zone sizes in the GLSZM using LBP 3D method M1.
Square glszm large area high gray level emphasis	Emphasizes large areas with high gray levels in the GLSZM matrix.
Lbp 3d k glszm gray level non uniformity normalized	Evaluates non-uniformity of gray levels in GLSZM with LBP 3D method K.
Lbp 3d m3 glszm zone variance	Measures variance of zones in GLSZM using LBP 3D method M3.
Wavelet hhh first order median	Computes the median of pixel intensities in the HHH wavelet component.
Lbp 3d m2 glszm zone variance	Evaluates variance of zones in GLSZM using LBP 3D method M2.
Wavelet hhh first order mean	Computes the mean of pixel intensities in the HHH wavelet component.
Wavelet lhl first order skewness	Measures the skewness of pixel intensity distribution in the LHL wavelet component.
Wavelet hhh first order median	(Repeated) Computes the median of pixel intensities in the HHH wavelet component.
Wavelet hlh glszm large area low gray level emphasis	Emphasizes large areas with low gray levels in GLSZM using HLH wavelet component.
Logarithm glcm autocorrelation	Measures pixel intensity correlation in the GLCM matrix using logarithmic transformation.
Logarithm glcm contrast	Quantifies the contrast of pixel intensities in the GLCM matrix using logarithmic transformation.
Lbp 3d k glszm large area emphasis	Emphasizes large areas in GLSZM using LBP 3D method K.
Original glszm small area low gray level emphasis	Emphasizes low gray levels in small areas of the GLSZM matrix.
Original shape major axis length	Measures the length of the major axis in the shape of the ROI.

Table 8: Radiomic feature selected description

## IV. Codes

All codes used in this work are provided below and are available in the public repository: <https://github.com/JudithOsuna/Radiomics-AI> [34].

### From DICOM to Nifti

```

1 from DicomRTTool.ReaderWriter import DicomReaderWriter
2 import os
3 import SimpleITK as sitk
4
5 # Main directory containing DICOM data
6 main_directory = "C:/Users/jor14/Desktop/Not_treatment_worked"
7 # Directory where the output Nifti files will be saved
8 output_directory = "C:/Users/jor14/Desktop/Not_treatment_worked_nii"
9
10 # Instantiate the DicomReaderWriter with a description and argument to get the maximum
11 Dicom_reader = DicomReaderWriter(description='Images', arg_max=True)
12
13 # Iterate through each sub-directory in the main directory
14 for subdir, dirs, files in os.walk(main_directory):
15     for dir in dirs:
16         current_path = os.path.join(subdir, dir)
17         print(f"Processing {current_path}...")
18
19         # Process DICOM data in the current directory
20         Dicom_reader.walk_through_folders(current_path)
21         all_rois = Dicom_reader.return_rois(print_rois=False)
22         # Set the contour names and their associations
23         Dicom_reader.set_contour_names_and_associations (contour_names=['z_lungs'])
24
25         # Determine which indexes have all the required ROIs
26         indexes = Dicom_reader.which_indexes_have_all_rois()
27
28         if indexes:
29             # Select the last index that contains all the required ROIs
30             pt_indx = indexes[-1]
31             Dicom_reader.set_index(pt_indx)
32             # Retrieve images and masks for the selected index
33             Dicom_reader.get_images_and_mask()
34
35             # Prepare the output subdirectory based on the current directory's name
36             subfolder_name = os.path.basename(current_path)
37             output_subdir = os.path.join(output_directory, subfolder_name)
38             print(f"Saving images in {output_subdir}")
39             # Create the output subdirectory if it does not exist
40             if not os.path.exists(output_subdir):
41                 os.makedirs(output_subdir)
42
43             # Retrieve the image and mask arrays
44             image = Dicom_reader.ArrayDicom
45             mask = Dicom_reader.mask
46             # Retrieve the SimpleITK image handles
47             dicom_sitk_handle = Dicom_reader.dicom_handle
48             mask_sitk_handle = Dicom_reader.annotation_handle
49
50             # Save the images and masks as Nifti files

```

```

51     sitk.WriteImage(dicom_sitk_handle, os.path.join(output_subdir, 'Image.nii'))
52     sitk.WriteImage(mask_sitk_handle, os.path.join(output_subdir, 'Mask.nii'))
53
54     print(f"Saved Image and Mask in {output_subdir}")
55     else:
56         print(f"No suitable indexes found in {current_path}. Skipping...")
57
58 print("Processing completed.")

```

Code 1: Code for transforming DICOM images to NIfTI format

## File 'parameters.yaml'

```

1 # This is an example of a parameters file
2 # It is written according to the YAML-convention (www.yaml.org) and is checked by the code
  for consistency.
3 # Three types of parameters are possible and reflected in the structure of the document:
4 #
5 # Parameter category:
6 #   Setting Name: <value>
7 #
8 # The three parameter categories are:
9 # - setting: Setting to use for preprocessing and class specific settings. if no <value>
  is specified, the value for
10 #   this setting is set to None.
11 # - featureClass: Feature class to enable, <value> is list of strings representing enabled
  features. If no <value> is
12 #   specified or <value> is an empty list ('[]'), all features for this class are enabled.
13 # - imageType: image types to calculate features on. <value> is custom kwarg settings
  (dictionary). if <value> is an
14 #   empty dictionary ('{}'), no custom settings are added for this input image.
15 #
16 # Some parameters have a limited list of possible values. Where this is the case, possible
  values are listed in the
17 # package documentation
18
19 # Settings to use, possible settings are listed in the documentation (section "Customizing
  the extraction").
20 setting: {correctMask: True}
21
22 # Image types to use: "Original" for unfiltered image, for possible filters, see
  documentation.
23 imageType:
24   Original: {} # for dictionaries / mappings, None values are not allowed, '{}' is
    interpreted as an empty dictionary
25   Wavelet: {}
26   Square: {}
27   SquareRoot: {}
28   Logarithm: {}
29   Exponential: {}
30   Gradient: {}
31   LBP2D: {}
32   LBP3D: {}
33 # Featureclasses, from which features must be calculated. If a featureclass is not
  mentioned, no features are calculated
34 # for that class. Otherwise, the specified features are calculated, or, if none are
  specified, all are calculated (excluding redundant/deprecated features).

```

```

35 featureClass:
36     # redundant Compactness 1, Compactness 2 an Spherical Disproportion features are
      disabled by default, they can be
37     # enabled by specifying individual feature names (as is done for glcm) and including
      them in the list.
38     shape:
39     firstorder: # specifying an empty list has the same effect as specifying nothing.
40     glcm:
41     glrlm: # for lists none values are allowed, in this case, all features are enabled
42     glszm:
43     gldm: # contains deprecated features, but as no individual features are specified, the
      deprecated features are not enabled
44     ngtdm:

```

Code 2: Parameters file needed to run the radiomic features extraction code

## Radiomic features extraction

```

1 from radiomics import featureextractor
2 import os
3 import pandas as pd
4 import warnings
5
6 # Suppress warnings for cleaner output
7 warnings.filterwarnings('ignore')
8
9 def radiomic_features_extraction(dataDir, params):
10     # Initialize an empty DataFrame to store the results
11     df = pd.DataFrame()
12
13     # Count the number of subdirectories in the main directory
14     count = 0
15     for path in os.listdir(dataDir):
16         if os.path.isdir(os.path.join(dataDir, path)):
17             count += 1
18
19     # Iterate through each subdirectory to extract radiomic features
20     for i in range(1, count + 1):
21         # Define the paths for the image and mask files
22         imagePath = os.path.join(dataDir, str(i).zfill(2) + "/" + "Image.nii")
23         labelPath = os.path.join(dataDir, str(i).zfill(2) + "/" + "Mask.nii")
24
25         print(imagePath)
26         print(labelPath)
27
28         # Initialize the feature extractor with the provided parameters
29         extractor = featureextractor.RadiomicsFeatureExtractor(params)
30         extractor.enableAllFeatures() # Enable all available features
31
32         try:
33             # Execute feature extraction and store the results
34             result = extractor.execute(imagePath, labelPath)
35
36             # Append the results to the DataFrame
37             df = df._append(result, ignore_index=True)
38         except Exception as e:
39             # Print any errors that occur during feature extraction

```

```

40         print(f"Error processing case {i}: {e}")
41
42     # Drop columns containing "diagnostics" in their names
43     for i in df.columns.values:
44         if "diagnostics" in i:
45             df.drop([i], axis=1, inplace=True)
46
47     return df
48
49 # Define directories and parameter file
50 dataDir = "C:/Users/jor14/Desktop/Treatment_worked_Nifti/"
51 params = "C:/Users/jor14/Desktop/parameters.yaml"
52
53 # Extract features from the first dataset
54 df_worked = radiomic_features_extraction(dataDir, params)
55
56 # Display the first 34 rows of the extracted features
57 df_worked.head(34)
58
59 # Define directories and parameter file for the second dataset
60 dataDir = "C:/Users/jor14/Desktop/Not_treatment_worked_Nifti/"
61 params = "C:/Users/jor14/Desktop/parameters.yaml"
62
63 # Extract features from the second dataset
64 df_not_worked = radiomic_features_extraction(dataDir, params)
65
66 # Display the first 16 rows of the extracted features
67 df_not_worked.head(16)
68
69 print("TIME TO CONCATENATE")
70
71 # Assign a result label to each DataFrame
72 df_worked = df_worked.assign(Result=1)
73 df_not_worked = df_not_worked.assign(Result=0)
74 # Concatenate the two DataFrames
75 radiomic_features = pd.concat([df_worked, df_not_worked])
76 # Reset the index of the concatenated DataFrame
77 radiomic_features.index = range(radiomic_features.shape[0])
78
79 # Display the first few rows of the combined DataFrame
80 radiomic_features.head()
81
82 # Save the combined DataFrame to an Excel file
83 radiomic_features.to_excel('radiomic_features_excel.xlsx', index=False)
84 print("DONE")

```

Code 3: Code for extracting radiomic features from Nifti images

## Clinical data preparation

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.experimental import enable_iterative_imputer
4 from sklearn.impute import IterativeImputer
5
6 # Load the data from the Excel file
7 df_clinical = pd.read_excel('clinical_features_excel_id.xlsx')

```

```

8
9 # Rename 'Gender' to 'Sex' and convert 'MALE' to 1 and others to 0
10 df_clinical.rename(columns={'Gender': 'Sex'}, inplace=True)
11 df_clinical['Sex'] = np.where(df_clinical['Sex'] == 'MALE', 1, 0)
12
13 # Separate the 'ID' column
14 ids = df_clinical['ID'] # Save the IDs separately
15 df_without_id = df_clinical.drop(columns=['ID'])
16
17 # Apply imputation to fill missing values
18 imputer = IterativeImputer(random_state=0)
19 imputed_data = imputer.fit_transform(df_without_id)
20 df_imputed_no_id = pd.DataFrame(imputed_data, columns=df_without_id.columns)
21
22 # Exclude 'PO2_venous' from integer conversion
23 column_to_exclude = 'PO2_venous'
24 df_imputed_no_id[df_imputed_no_id.columns.difference([column_to_exclude])] =
25     df_imputed_no_id[
26     df_imputed_no_id.columns.difference([column_to_exclude])
27 ].apply(np.floor).astype(int)
28
29 # Add the 'ID' column back to the imputed DataFrame
30 df_imputed_with_id = df_imputed_no_id.copy()
31 df_imputed_with_id['ID'] = ids.values
32
33 # List of columns to transform with one-hot encoding
34 columns_to_transform = ['GDS', 'Barthel_scale', 'Pulmonary_affectation']
35
36 # Create dummies/one-hot encoding for the specified columns
37 df_one_hot_encoded = pd.get_dummies(df_imputed_with_id, columns=columns_to_transform,
38     prefix=columns_to_transform)
39
40 # Replace any boolean values with integers (False = 0, True = 1)
41 df_one_hot_encoded.replace({False: 0, True: 1}, inplace=True)
42
43 # Merge the encoded DataFrame with the results DataFrame based on 'ID'
44 df_combined_clinical = pd.merge(df_one_hot_encoded, dfresult, on='ID', how='outer')
45
46 # Display the first few rows of the combined DataFrame
47 print(df_combined_clinical.head())
48
49 # Save the resulting DataFrame with 'ID' to a new Excel file
50 df_combined_clinical.to_excel('clinical_ID.xlsx', index=False)
51
52 print("File 'clinical_ID.xlsx' created successfully.")

```

Code 4: Code for filling in missing values and processing categorical characteristics of clinical data

## Metabolomic data preparation

```

1 import pandas as pd
2 from sklearn.experimental import enable_iterative_imputer
3 from sklearn.impute import IterativeImputer
4
5 # Load the data from the Excel file
6 df_metabolomics = pd.read_excel('metabolomic_features_id.xlsx')

```

```

7
8 # Separate the 'ID' column
9 ids = df_metabolomics['ID'] # Save the IDs separately
10 df_without_id = df_metabolomics.drop(columns=['ID'])
11
12 # Apply imputation to fill missing values
13 imputer = IterativeImputer(random_state=0)
14 imputed_data = imputer.fit_transform(df_without_id)
15 df_imputed_no_id = pd.DataFrame(imputed_data, columns=df_without_id.columns)
16
17 # Add the 'ID' column back to the imputed DataFrame
18 df_imputed_with_id = df_imputed_no_id.copy()
19 df_imputed_with_id['ID'] = ids.values
20
21 # Display the first few rows of the imputed DataFrame
22 print(df_imputed_with_id.head())
23
24 # Save the imputed DataFrame with 'ID' to a new Excel file
25 df_imputed_with_id.to_excel('metabolomic_ID.xlsx', index=False)
26
27 print("File 'metabolomic_ID.xlsx' created successfully.")

```

Code 5: Code to fill in missing values in metabolomics data

## Clinical and Metabolomic data augmentation

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from imblearn.over_sampling import SMOTE
5
6 # Load the data from the Excel file
7 df = pd.read_excel('metabolomic.xlsx')
8
9 # Separate features (X) and the target variable (y)
10 X = df.drop('Result', axis=1) # Replace 'Result' with the name of your target column
11 y = df['Result']
12
13 # Apply SMOTE to balance the data
14 smote = SMOTE(sampling_strategy='auto', random_state=42)
15 X_resampled, y_resampled = smote.fit_resample(X, y)
16
17 # Combine the resampled features and target variable into a DataFrame
18 df_resampled = pd.concat([pd.DataFrame(X_resampled, columns=X.columns),
19                            pd.DataFrame(y_resampled, columns=['Result'])], axis=1)
20
21 # Save the balanced DataFrame to a new Excel file
22 df_resampled.to_excel('metabolomic_augmented.xlsx', index=False)
23
24 print("The data has been balanced and saved in 'metabolomic_augmented.xlsx'")
25
26 # Plot the distribution of a selected feature
27 feature = 'LDH' # Replace 'LDH' with the column name you want to plot
28
29 # Original data (Result=0)
30 original_data = df[df['Result'] == 0][feature]

```

```

31 # Synthetic data (Result=0, after SMOTE)
32 synthetic_data = df_resampled[(df_resampled['Result'] == 0) &
    (~df_resampled.index.isin(df.index))][feature]
33
34 # Plotting the distributions
35 plt.figure(figsize=(12, 6))
36
37 # Histogram and KDE of original data
38 sns.histplot(original_data, kde=True, color='blue', stat="density", label='Original
    (Result=0)', bins=20)
39
40 # Histogram and KDE of synthetic data
41 sns.histplot(synthetic_data, kde=True, color='red', stat="density", label='Synthetic
    (SMOTE)', bins=20)
42
43 plt.title(f'Gaussian Distribution of {feature} - Original vs Synthetic Data')
44 plt.xlabel(feature)
45 plt.ylabel('Density')
46 plt.legend()
47 plt.show()

```

Code 6: Code to increase the unbalanced negative response on clinical and metabolomic data

## CT images augmentation

```

1 import nibabel as nib
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import os
5 import shutil
6 import torch
7 from skimage import exposure
8
9
10 # Directory containing the numbered subfolders
11 base_dir = 'C:/Users/jor14/Desktop/Not_treatment_worked_Nifti/' # Change this to your
    base directory path
12
13 # Create a folder to save synthetic images and masks
14 output_dir = 'synthetic_images_real'
15 if os.path.exists(output_dir):
16     shutil.rmtree(output_dir) # Remove the existing directory if it already exists
17 os.makedirs(output_dir)
18
19 # Function to add random noise to an image
20 def add_random_noise(image_np, noise_level=0.1):
21     noise = np.random.normal(loc=0, scale=noise_level, size=image_np.shape)
22     return image_np + noise
23
24
25 # Function to adjust brightness and contrast of an image
26 def adjust_brightness_contrast(image_np, brightness_factor=1.0, contrast_factor=1.0):
27     # Adjust brightness
28     image_np = image_np * brightness_factor
29     # Adjust contrast
30     p2, p98 = np.percentile(image_np, (2, 98))

```

```

31 image_np = exposure.rescale_intensity(image_np, in_range=(p2, p98)) * contrast_factor
32 return np.clip(image_np, 0, 255)
33
34 # Function to apply random flipping to an image and its mask
35 def apply_random_flipping(image_np, mask_np):
36     if np.random.rand() > 0.5:
37         image_np = np.flip(image_np, axis=0)
38         mask_np = np.flip(mask_np, axis=0)
39     if np.random.rand() < 0.5:
40         image_np = np.flip(image_np, axis=1)
41         mask_np = np.flip(mask_np, axis=1)
42     return image_np, mask_np
43
44 # Function to generate synthetic images from an original image and mask
45 def generate_synthetic_image(image_np, mask_np):
46     # Add random noise
47     noisy_image = add_random_noise(image_np)
48
49     # Adjust brightness and contrast
50     brightness_factor = np.random.uniform(0.8, 1.2) # Random brightness factor
51     contrast_factor = np.random.uniform(0.8, 1.2) # Random contrast factor
52     final_image = adjust_brightness_contrast(noisy_image, brightness_factor,
53         contrast_factor)
54
55     # Apply random flipping
56     final_image, final_mask = apply_random_flipping(final_image, mask_np)
57
58     return final_image, final_mask
59
60 num_original_images = 16 # Number of original images to generate
61 num_synthetic_images = 18 # Number of synthetic images to generate
62
63 # Store original images and masks
64 original_images = []
65 original_masks = []
66
67 # Load 16 original images and their masks
68 for i in range(1, 17):
69     folder_path = os.path.join(base_dir, f'{i:02d}')
70     image_path = os.path.join(folder_path, 'Image.nii')
71     mask_path = os.path.join(folder_path, 'Mask.nii')
72
73     # Load image and mask
74     image = nib.load(image_path)
75     mask = nib.load(mask_path)
76
77     # Convert data to numpy arrays
78     image_data_np = image.get_fdata()
79     mask_data_np = mask.get_fdata()
80
81     # Convert numpy arrays to PyTorch tensors
82     image_data_tensor = torch.from_numpy(image_data_np.astype(np.float32))
83     mask_data_tensor = torch.from_numpy(mask_data_np.astype(np.float32))
84
85     # Store tensors and affine matrices in lists
86     original_images.append((image_data_tensor, image.affine))
87     original_masks.append((mask_data_tensor, mask.affine))
88
89 # Generate and save synthetic images and masks

```

```

89 for i in range(1, num_synthetic_images + 1):
90     # Select a random original image to generate a synthetic image
91     idx = (i - 1) % num_original_images
92     original_image_tensor, affine_image = original_images[idx]
93     original_mask_tensor, affine_mask = original_masks[idx]
94
95     # Convert tensors to numpy arrays
96     original_image_np = original_image_tensor.numpy()
97     original_mask_np = original_mask_tensor.numpy()
98
99     # Generate synthetic data
100    synthetic_image_np, synthetic_mask_np = generate_synthetic_image(original_image_np,
101        original_mask_np)
102
103    # Create a subfolder for each synthetic image
104    output_subdir = os.path.join(output_dir, f'{i:02d}')
105    os.makedirs(output_subdir)
106
107    # Create NIfTI images for synthetic data
108    synthetic_image = nib.Nifti1Image(synthetic_image_np, affine_image)
109    synthetic_mask = nib.Nifti1Image(synthetic_mask_np, affine_mask)
110
111    # Save synthetic images and masks in the corresponding subfolder
112    nib.save(synthetic_image, os.path.join(output_subdir, 'Image.nii'))
113    nib.save(synthetic_mask, os.path.join(output_subdir, 'Mask.nii'))
114
115    # Compare intensity distributions of original and synthetic images
116    original_sample = original_images[0][0].numpy().flatten()
117    synthetic_sample = synthetic_data_np.flatten()
118
119    plt.figure(figsize=(12, 6))
120    plt.hist(original_sample, bins=50, alpha=0.5, label='Original')
121    plt.hist(synthetic_sample, bins=50, alpha=0.5, label='Synthetic')
122    plt.legend(loc='upper right')
123    plt.title('Intensity Distribution')
124    plt.xlabel('Intensity')
125    plt.ylabel('Frequency')
126    plt.show()

```

Code 7: Code for augmenting unbalanced Computed Tomography Images

## Preprocessing all data manually

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 # For data imputation and scaling
8 from sklearn.experimental import enable_iterative_imputer
9 from sklearn.impute import IterativeImputer
10 from sklearn.preprocessing import StandardScaler
11
12 # For calculating Variance Inflation Factor (VIF)
13 from statsmodels.stats.outliers_influence import variance_inflation_factor
14

```

```

15 import warnings
16 warnings.filterwarnings("ignore")
17
18 # Load the dataset
19 df_radiomic = pd.read_excel('clinical_nulls.xlsx')
20
21 # Remove columns with constant values
22 df_radiomic = df_radiomic.loc[:, (df_radiomic != df_radiomic.iloc[0]).any()]
23
24 # List to store names of columns with object data type
25 object_columns = []
26
27 # Check the data type of each column
28 for col in df_radiomic.columns:
29     print(col, ': ', df_radiomic[col].dtype)
30     if df_radiomic[col].dtype == 'object':
31         object_columns.append(col)
32
33 # Check if there are any object type columns and print them
34 if object_columns:
35     print("\nThe following variables are of type object in the DataFrame:")
36     for col in object_columns:
37         print(col)
38 else:
39     print("\nThere are no object-type variables in the DataFrame.")
40
41 # Extract the independent variables (X) and dependent variable (y)
42 X = df_radiomic.drop('Result', axis=1)
43 y = df_radiomic.Result
44
45 # Display the distribution of the target variable
46 print(y.value_counts())
47
48 # Calculate the correlation matrix
49 corr_df = df_radiomic.corr(method='pearson')
50 print(corr_df.head())
51
52 # Plot the correlation matrix
53 sns.set(rc={'figure.figsize':(11.7, 8.27)})
54 plt.figure(2)
55 sns.heatmap(corr_df)
56 plt.show(block=False)
57
58 # Identify highly correlated variables
59 highly_correlated_vars = set()
60 threshold = 0.75
61 target_variable = "Result"
62
63 for i in range(len(corr_df.columns)):
64     for j in range(i):
65         if abs(corr_df.iloc[i, j]) > threshold:
66             var1 = corr_df.columns[i]
67             var2 = corr_df.columns[j]
68             corr_with_target_var1 = abs(df_radiomic[var1].corr(df_radiomic
69 [target_variable]))
70             corr_with_target_var2 = abs(df_radiomic[var2].corr(df_radiomic
71 [target_variable]))
72             if corr_with_target_var1 < corr_with_target_var2:
73                 highly_correlated_vars.add(var1)

```

```

74         else:
75             highly_correlated_vars.add(var2)
76
77 # Drop highly correlated variables
78 df_radiomic.drop(columns=highly_correlated_vars, inplace=True)
79 X.drop(columns=highly_correlated_vars, inplace=True)
80
81 # Recalculate the correlation matrix if needed
82 correlation_matrix_updated = df_radiomic.corr(method='pearson')
83
84 # Function to calculate and filter by VIF
85 def calculate_vif(x, thresh=5):
86     output = x.copy()
87     while True:
88         vif = [variance_inflation_factor(output.values, i) for i in range(output.shape[1])]
89         max_vif = max(vif)
90         if max_vif > thresh:
91             max_index = vif.index(max_vif)
92             output = output.drop(output.columns[max_index], axis=1)
93         else:
94             break
95     return output
96
97 # Select features based on VIF
98 selected_features = calculate_vif(X)
99 print(selected_features.head())
100
101 # Get the names of the selected columns
102 selected_features_columns = selected_features.columns
103
104 # Update df_radiomic to only contain the selected columns and the 'Result' column
105 df_radiomic_cleaned = df_radiomic[list(selected_features_columns) + ['Result']]
106
107 # Initialize the standard scaler
108 scaler = StandardScaler()
109
110 # Select the columns to standardize (all except 'Result')
111 columns_to_scale = df_radiomic_cleaned.drop(columns='Result')
112
113 # Apply standardization to the selected columns
114 df_radiomic_cleaned[columns_to_scale.columns] = scaler.fit_transform(columns_to_scale)
115
116 # Save the cleaned and standardized DataFrame to an Excel file
117 df_radiomic_cleaned.to_excel('MetabolicData.xlsx', index=False)

```

Code 8: Code to prepare manually all data for the Predictive Model

## Preprocessing all data automatically

```

1 import pandas as pd
2 import pymrmr
3 from sklearn.preprocessing import MinMaxScaler
4
5 # Load the dataset from an Excel file
6 data = pd.read_excel('C:/Users/jor14/Desktop/radiomics&metabolomic&clinical.xlsx')
7
8 # Assuming the 'Result' column is the target variable

```

```

9 X = data.drop(columns=['Result']) # Features (all columns except 'Result')
10 y = data['Result'] # Target variable
11
12 # Scale the features to a range between 0 and 1
13 scaler = MinMaxScaler()
14 X_scaled = scaler.fit_transform(X) # Fit and transform the features
15
16 # Concatenate the scaled features and the target variable for use in pymrmr
17 data_for_mrmr = pd.concat([y, pd.DataFrame(X_scaled, columns=X.columns)], axis=1)
18
19 # Select the top 10 features using the mRMR method with the 'MIQ' criterion
20 selected_features = pymrmr.mRMR(data_for_mrmr, 'MIQ', 10) # 'MIQ' stands for Mutual
    Information Quotient
21
22 # Add the target variable to the list of selected features
23 selected_features_with_target = ['Result'] + selected_features
24
25 # Create a new DataFrame with only the selected features and the target variable
26 selected_features_df = data[selected_features_with_target]
27
28 # Save the resulting DataFrame to a new Excel file
29 selected_features_df.to_excel('C:/Users/jor14/Desktop/ALLpreprocessed_auto.xlsx',
    index=False)
30
31 # Display the first few rows of the resulting DataFrame
32 print("DataFrame with the selected features and the target variable:")
33 print(selected_features_df.head())

```

Code 9: Code to select features automatically for the Predictive Model

## Predictive Model with RF

```

1 # Import necessary libraries
2 import pandas as pd # For data manipulation
3 import numpy as np # For numerical operations
4 import matplotlib.pyplot as plt # For plotting graphs
5 import seaborn as sns # For data visualization
6
7 import warnings # To ignore warnings
8 warnings.filterwarnings("ignore")
9
10 from sklearn.model_selection import train_test_split, cross_val_predict, StratifiedKFold
    # For data splitting and cross-validation
11 from sklearn.ensemble import RandomForestClassifier # For Random Forest model
12 from sklearn.metrics import confusion_matrix, recall_score, f1_score, precision_score,
    accuracy_score, roc_curve, auc # For model evaluation metrics
13 from sklearn.model_selection import learning_curve # For plotting learning curves
14
15 # Load the dataset
16 df = pd.read_excel('C:/Users/jor14/Desktop/allpreprocessed_auto.xlsx')
17
18 # Separate the dataset into features (X) and target variable (y)
19 X = df.drop('Result', axis=1)
20 y = df['Result']
21 print(X.shape)
22

```

```

23 # Split the dataset into training and testing sets with stratification to maintain class
    proportions
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1, train_size=0.9,
    stratify=y, random_state=42)
25
26 # Convert target variables to arrays for compatibility with scikit-learn
27 y_train = y_train.values
28 y_test = y_test.values
29
30 # Convert feature datasets to arrays for compatibility with scikit-learn
31 X_train = X_train.values
32 X_test = X_test.values
33
34 # Extract the feature names
35 feature_names = X.columns
36 print(feature_names)
37
38 # Initialize and train the Random Forest Classifier
39 clf = RandomForestClassifier(n_estimators=100, random_state=42)
40 clf.fit(X_train, y_train)
41
42 # Calculate feature importance scores
43 feature_scores = pd.Series(clf.feature_importances_,
    index=X.columns).sort_values(ascending=False)
44
45 # Select the top n features based on their importance scores
46 n = 9
47 selected_features = feature_scores.index.values[:n]
48 print(selected_features)
49
50 # Get the indices of the selected features
51 cols_index = [X.columns.tolist().index(col) for col in selected_features]
52
53 # Plot the feature importance scores
54 f, ax = plt.subplots(figsize=(15, 10))
55 sns.barplot(x=feature_scores.index, y=feature_scores.values, ax=ax, color='maroon',
    alpha=0.8)
56 ax.set_title("Feature scores of the features")
57 ax.set_yticklabels(feature_scores.index)
58 ax.set_xlabel("Feature importance score")
59 ax.set_ylabel("Features")
60 plt.show()
61
62 # Reduce the training and testing datasets to the selected features
63 X_train = X_train[:, cols_index]
64 X_test = X_test[:, cols_index]
65
66 #-----
67 # MODIFY THIS PART OF THE CODE FOR LR AND DT PM
68 #-----
69 # Implement cross-validation with stratified folds
70 cv = StratifiedKFold(n_splits=5)
71 pred_cv = cross_val_predict(clf, X_train, y_train, cv=cv)
72 #-----
73
74 # Calculate ROC curve and AUC
75 fpr, tpr, thresholds = roc_curve(y_train, pred_cv)
76 roc_auc = auc(fpr, tpr)
77

```

```

78 # Function to print the confusion matrix in a readable format
79 def print_confusion_matrix(y_true, y_pred):
80     cm = confusion_matrix(y_true, y_pred)
81     print('True positive = ', cm[1][1])
82     print('False positive = ', cm[0][1])
83     print('False negative = ', cm[1][0])
84     print('True negative = ', cm[0][0])
85
86 # Generate and print the confusion matrix
87 conf_matrix = confusion_matrix(y_train, pred_cv)
88 conf_matrix = np.flipud(conf_matrix)
89 conf_matrix = np.fliplr(conf_matrix)
90
91 # Print various model evaluation metrics
92 print('Accuracy: {:.2f}'.format(accuracy_score(y_train, pred_cv)))
93 print('Precision: {:.2f}'.format(precision_score(y_train, pred_cv)))
94 print('Recall: {:.2f}'.format(recall_score(y_train, pred_cv)))
95 print('F1 Score: {:.2f}'.format(f1_score(y_train, pred_cv)))
96 print('AUC: {:.2f}'.format(auc(fpr, tpr)))
97 print('Confusion Matrix:')
98 print(conf_matrix)
99 print_confusion_matrix(y_train, pred_cv)
100
101 # Extract confusion matrix components for specificity calculation
102 tn, fp, fn, tp = confusion_matrix(y_train, pred_cv).ravel()
103
104 # Calculate specificity
105 specificity = tn / (tn + fp)
106 print(f'Specificity: {specificity:.2f}')
107
108 # Plot ROC curve
109 plt.figure()
110 plt.plot(fpr, tpr, color='maroon', lw=2, label='AUC = %0.2f' % roc_auc)
111 plt.plot([0, 1], [0, 1], color='gray', lw=2, linestyle='--')
112 plt.xlim([0.0, 1.0])
113 plt.ylim([0.0, 1.05])
114 plt.xlabel('False Positive Rate')
115 plt.ylabel('True Positive Rate')
116 plt.title('ROC Curve')
117 plt.legend(loc="lower right")
118 plt.show()
119
120 # Plot the confusion matrix as a heatmap
121 fig, px = plt.subplots(figsize=(7.5, 7.5))
122 px.matshow(conf_matrix, cmap=plt.cm.YlOrRd, alpha=0.8)
123 for m in range(conf_matrix.shape[0]):
124     for n in range(conf_matrix.shape[1]):
125         px.text(x=m, y=n, s=conf_matrix[n, m], va='center', ha='center', size='xx-large')
126
127 # Set labels and title for the confusion matrix plot
128 plt.xlabel('Predictions', fontsize=16)
129 plt.ylabel('Actuals', fontsize=16)
130 plt.title('Confusion Matrix', fontsize=15)
131 plt.show()
132
133 # Generate learning curves
134 train_sizes, train_scores, valid_scores = learning_curve(clf, X, y, cv=5, scoring='f1',
135     train_sizes=np.linspace(0.1, 1.0, 10))

```

```

136 # Calculate the mean and standard deviation of the training and validation scores
137 train_mean = np.mean(train_scores, axis=1)
138 train_std = np.std(train_scores, axis=1)
139 valid_mean = np.mean(valid_scores, axis=1)
140 valid_std = np.std(valid_scores, axis=1)
141
142 # Plot the learning curves
143 plt.figure(figsize=(10, 6))
144 plt.plot(train_sizes, train_mean, 'o-', color='purple', label='Training')
145 plt.plot(train_sizes, valid_mean, 'o-', color='red', label='Cross-Validation')
146
147 # Fill between the standard deviation bands for training and validation curves
148 plt.fill_between(train_sizes, train_mean - train_std, train_mean + train_std,
149                 color='purple', alpha=0.1)
149 plt.fill_between(train_sizes, valid_mean - valid_std, valid_mean + valid_std, color='red',
150                 alpha=0.1)
151
152 # Set plot title and labels
153 plt.title('Learning Curves')
154 plt.xlabel('Training Set Size')
155 plt.ylabel('F1 Score')
156 plt.legend(loc='best')
157 plt.show()

```

Code 10: Code to train and evaluate the Random Forest classifier

## Predictive Model with LR

```

1 #-----
2 # Import necessary libraries
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.model_selection import StratifiedKFold, cross_val_predict
5
6 # Initialize the Logistic Regression model
7 lr = LogisticRegression(random_state=42)
8
9 # Perform Stratified K-Fold Cross-Validation and get predictions
10 cv = StratifiedKFold(n_splits=5)
11 pred_cv = cross_val_predict(lr, X_train, y_train, cv=cv, method='predict')
12 #-----

```

Code 11: Code to modify the Predictive Model with Logistic Regression

## Predictive Model with DT

```

1 #-----
2 # Import necessary libraries
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.model_selection import StratifiedKFold, cross_val_predict
5
6 # Initialize the Decision Tree model
7 dt = DecisionTreeClassifier(random_state=42)
8
9 # Perform Stratified K-Fold Cross-Validation and get predictions

```

```

10 cv = StratifiedKFold(n_splits=5)
11 pred_cv = cross_val_predict(dt, X_train, y_train, cv=cv, method='predict')
12 #-----

```

Code 12: Code to modify the Predictive Model with Decision Tree

## Best PM with OPTUNA and testing

```

1
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import warnings
7
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.model_selection import train_test_split, cross_val_score, cross_val_predict,
   StratifiedKFold
10 from sklearn.metrics import confusion_matrix, recall_score, f1_score, precision_score,
   accuracy_score, roc_curve, auc
11 import optuna
12
13 # Suppress warnings
14 warnings.filterwarnings("ignore")
15
16 # Load the dataset
17 df = pd.read_excel('C:/Users/jor14/Desktop/ALLaugmented_featureauto.xlsx')
18
19 # Extract features (independent variables) and target (dependent variable)
20 X = df.drop('Result', axis=1)
21 y = df['Result']
22
23 # Split the data into training and testing sets
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1, train_size=0.9,
   stratify=y, random_state=42)
25
26 # Convert variables to arrays
27 y_train = y_train.values
28 X_test = X_test.values
29 X_train = X_train.values
30 y_test = y_test.values
31
32 # Get the feature names
33 feature_names = X.columns
34 print(feature_names)
35
36 # Initialize and fit the RandomForest model
37 clf = RandomForestClassifier(n_estimators=100, random_state=42)
38 clf.fit(X_train, y_train)
39
40 # Calculate feature importance scores
41 feature_importance_scores = pd.Series(clf.feature_importances_,
   index=X.columns).sort_values(ascending=False)
42 n=10
43
44 # Display feature importance scores
45 print(feature_importance_scores)

```

```

46
47 # Select the top n features
48 selected_features = feature_importance_scores.index.values[:n]
49 print(selected_features)
50
51 # Get the indices of the selected features
52 cols_index = [X.columns.tolist().index(col) for col in selected_features]
53
54 # Plot the feature importance scores
55 plt.figure(figsize=(15, 10))
56 sns.barplot(x=feature_importance_scores.values, y=feature_importance_scores.index,
57             color='maroon', alpha=0.8)
58
59 # Set plot title and labels
60 plt.title("Feature Importance Scores")
61 plt.xlabel("Importance Score")
62 plt.ylabel("Features")
63
64 # Show the plot
65 plt.show()
66
67 # Reduce the training and testing sets to the selected features
68 X_train = X_train[:, cols_index]
69 X_test = X_test[:, cols_index]
70
71 # Objective function for hyperparameter tuning with Optuna
72 def objective(trial):
73     # Define the hyperparameters to tune
74     n_estimators = trial.suggest_int('n_estimators', 10, 200)
75     max_depth = trial.suggest_int('max_depth', 2, 32)
76     min_samples_split = trial.suggest_int('min_samples_split', 2, 10)
77     min_samples_leaf = trial.suggest_int('min_samples_leaf', 1, 4)
78     bootstrap = trial.suggest_categorical('bootstrap', [True, False])
79
80     # Create the model with the suggested hyperparameters
81     model = RandomForestClassifier(
82         n_estimators=n_estimators,
83         max_depth=max_depth,
84         min_samples_split=min_samples_split,
85         min_samples_leaf=min_samples_leaf,
86         bootstrap=bootstrap,
87         random_state=42
88     )
89
90     # Evaluate the model using cross-validation
91     cv = StratifiedKFold(n_splits=5)
92     scores = cross_val_score(model, X, y, cv=cv, scoring='f1')
93
94     return scores.mean()
95
96 # Create and optimize the study with Optuna
97 study = optuna.create_study(direction='maximize')
98 study.optimize(objective, n_trials=50)
99
100 # Print the best parameters and best score found by Optuna
101 print(f"Best Parameters: {study.best_params}")
102 print(f"Best Score (F1): {study.best_value}")
103
104 # Train the best model with the selected features and parameters

```

```

104 best_params = study.best_params
105 best_model = RandomForestClassifier(**best_params, random_state=42)
106 best_model.fit(X_train, y_train)
107
108 # Evaluate the model on the training set using cross-validation predictions
109 y_probs = cross_val_predict(best_model, X_train, y_train, cv=StratifiedKFold(n_splits=5),
110                             method="predict")
111 fpr, tpr, thresholds = roc_curve(y_train, y_probs)
112 roc_auc = auc(fpr, tpr)
113
114 # Define function to print confusion matrix details
115 def print_confusion_matrix(y_true, y_pred):
116     cm = confusion_matrix(y_true, y_pred)
117     print('True Positive = ', cm[1][1])
118     print('False Positive = ', cm[0][1])
119     print('False Negative = ', cm[1][0])
120     print('True Negative = ', cm[0][0])
121
122 # Print model evaluation metrics
123 print('Accuracy: {:.2f}'.format(accuracy_score(y_train, y_probs)))
124 print('Precision: {:.2f}'.format(precision_score(y_train, y_probs)))
125 print('Recall: {:.2f}'.format(recall_score(y_train, y_probs)))
126 print('F1 Score: {:.2f}'.format(f1_score(y_train, y_probs)))
127 print('AUC: {:.2f}'.format(roc_auc))
128
129 # Convert probabilities to binary predictions
130 y_pred = (y_probs >= 0.5).astype(int)
131
132 # Confusion matrix for the training set
133 cm_train = confusion_matrix(y_train, y_pred)
134 print('Confusion Matrix:')
135 print(cm_train)
136 print_confusion_matrix(y_train, y_pred)
137
138 # Plot ROC curve
139 plt.figure()
140 plt.plot(fpr, tpr, color='maroon', lw=2, label='AUC = %0.2f' % roc_auc)
141 plt.plot([0, 1], [0, 1], color='gray', lw=2, linestyle='--')
142 plt.xlim([0.0, 1.0])
143 plt.ylim([0.0, 1.05])
144 plt.xlabel('False Positive Rate')
145 plt.ylabel('True Positive Rate')
146 plt.title('ROC Curve')
147 plt.legend(loc="lower right")
148 plt.show()
149
150 # Plot the confusion matrix
151 fig, ax = plt.subplots(figsize=(7.5, 7.5))
152 ax.matshow(cm_train, cmap=plt.cm.YlOrRd, alpha=0.8)
153 for i in range(cm_train.shape[0]):
154     for j in range(cm_train.shape[1]):
155         ax.text(x=i, y=j, s=cm_train[j, i], va='center', ha='center', size='xx-large')
156
157 # Set labels and title
158 plt.xlabel('Predictions', fontsize=16)
159 plt.ylabel('Actuals', fontsize=16)
160 plt.title('Confusion Matrix', fontsize=15)
161 plt.show()

```

```

162 # Plot learning curves
163 train_sizes, train_scores, valid_scores = learning_curve(best_model, X_train, y_train,
    cv=5, scoring='f1', train_sizes=np.linspace(0.1, 1.0, 10))
164
165 # Calculate mean and standard deviation of the scores
166 train_mean = np.mean(train_scores, axis=1)
167 train_std = np.std(train_scores, axis=1)
168 valid_mean = np.mean(valid_scores, axis=1)
169 valid_std = np.std(valid_scores, axis=1)
170
171 # Plot the learning curves
172 plt.figure(figsize=(10, 6))
173 plt.plot(train_sizes, train_mean, 'o-', color='purple', label='Training')
174 plt.plot(train_sizes, valid_mean, 'o-', color='red', label='Cross-Validation')
175
176 # Fill between the standard deviation bands
177 plt.fill_between(train_sizes, train_mean - train_std, train_mean + train_std,
    color='purple', alpha=0.1)
178 plt.fill_between(train_sizes, valid_mean - valid_std, valid_mean + valid_std, color='red',
    alpha=0.1)
179
180 # Set plot title and labels
181 plt.title('Learning Curves')
182 plt.xlabel('Training Set Size')
183 plt.ylabel('F1 Score')
184 plt.legend(loc='best')
185 plt.show()
186
187 # Evaluate the model on the test set
188 clf = RandomForestClassifier(n_estimators=100, random_state=42)
189 clf.fit(X_train, y_train)
190 pred_test = clf.predict(X_test)
191
192 fpr_test, tpr_test, thresholds_test = roc_curve(y_test, pred_test)
193 roc_auc_test = auc(fpr_test, tpr_test)
194
195 # Print test set evaluation metrics
196 print('Accuracy: {:.2f}'.format(accuracy_score(y_test, pred_test)))
197 print('Precision: {:.2f}'.format(precision_score(y_test, pred_test)))
198 print('Recall: {:.2f}'.format(recall_score(y_test, pred_test)))
199 print('F1 Score: {:.2f}'.format(f1_score(y_test, pred_test)))
200 print('AUC: {:.2f}'.format(roc_auc_test))
201
202 # Confusion matrix for the test set
203 cm_test = confusion_matrix(y_test, pred_test)
204 print('Confusion Matrix:')
205 print(cm_test)
206 print_confusion_matrix(y_test, pred_test)
207
208 # Plot ROC curve for the test set
209 plt.figure()
210 plt.plot(fpr_test, tpr_test, color='maroon', lw=2, label='AUC = %0.2f' % roc_auc_test)
211 plt.plot([0, 1], [0, 1], color='gray', lw=2, linestyle='--')
212 plt.xlim([0.0, 1.0])
213 plt.ylim([0.0, 1.05])
214 plt.xlabel('False Positive Rate')
215 plt.ylabel('True Positive Rate')
216 plt.title('ROC Curve')
217 plt.legend(loc="lower right")

```

```
218 plt.show()
219
220 # Plot the confusion matrix for the test set
221 fig, ax = plt.subplots(figsize=(7.5, 7.5))
222 ax.matshow(cm_test, cmap=plt.cm.YlOrRd, alpha=0.5)
223 for i in range(cm_test.shape[0]):
224     for j in range(cm_test.shape[1]):
225         ax.text(x=i, y=j, s=cm_test[j, i], va='center', ha='center', size='xx-large')
226
227 # Set labels and title
228 plt.xlabel('Predictions', fontsize=16)
229 plt.ylabel('Actuals', fontsize=16)
230 plt.title('Confusion Matrix', fontsize=15)
231 plt.show()
```

Code 13: Code to test the best Predictive Model