

**UCLV**  
Universidad Central  
"Marta Abreu" de Las Villas



**FIE**  
Facultad de  
Ingeniería Eléctrica

Departamento de Control Automático

## TRABAJO DE DIPLOMA

Título: Utilización de características radiómicas en la clasificación maligno-benigno de nódulos pulmonares a partir de TC

Autor: Rachel Abreu Llanes

Tutores: DrC. Marlen Pérez Díaz, DrC Déborah Galpert Cañizares

**UCLV**  
Universidad Central  
"Marta Abreu" de Las Villas



**FIE**  
Facultad de  
Ingeniería Eléctrica

Automatic Control Department

## TRABAJO DE DIPLOMA

Title: Use of radiomic features in the malignant-benign classification of pulmonary nodules from CT scans

Author: Rachel Abreu Llanes

Thesis Director: DrC. Marlen Pérez Díaz, DrC Déborah Galpert

Santa Clara  
Copyright©UCLV

Este documento es Propiedad Patrimonial de la Universidad Central “Marta Abreu” de Las Villas, y se encuentra depositado en los fondos de la Biblioteca Universitaria “Chiqui Gómez Lubian” subordinada a la Dirección de Información Científico Técnica de la mencionada casa de altos estudios.

Se autoriza su utilización bajo la licencia siguiente:

**Atribución- No Comercial- Compartir Igual**



Para cualquier información contacte con:

Dirección de Información Científico Técnica. Universidad Central “Marta Abreu” de Las Villas. Carretera a Camajuaní. Km 5½. Santa Clara. Villa Clara. Cuba. CP. 54 830

Teléfonos.: +53 01 42281503-1419



## ACTA DE CONFORMIDAD PARA ESTUDIANTES DE PREGRADO

Universidad Central "Marta Abreu" de Las Villas

Por una parte: Rachel Abreu Llanes estudiante de la carrera de: Ingeniería en Automática en la facultad de: Ingeniería Eléctrica, en lo adelante **EL ESTUDIANTE**. Con número de identidad permanente: 01050372016. Y por otra parte DrC Iván Santana Ching Jefe del Departamento Docente de: Control Automático en la ya mencionada facultad, en lo adelante **EL JEFE DE DEPARTAMENTO**, y DrC Marlen Pérez Díaz y DrC Déborah Galpert Cañizares profesor(es) encargado(s) de tuturar el Trabajo de Diploma **DEL ESTUDIANTE**, en lo adelante **EL TUTOR**.

Reconocen que:

- I. **EL ESTUDIANTE** se le ha aprobado como tema de investigación para su Trabajo de Diploma el titulado "Utilización de características radiómicas en la clasificación maligno- benigno de nódulos pulmonares a partir de TC"
- II. **EL ESTUDIANTE** no divulgará información concerniente a la investigación, tanto durante el desarrollo como tras la culminación de esta sin la debida autorización **DEL TUTOR** o **EL JEFE DE DEPARTAMENTO**.
- III. Que el Trabajo de Diploma fruto de la labor investigativa de **EL ESTUDIANTE** y la asesoría de **EL TUTOR**, resulta de **TITULARIDAD EXCLUSIVA** de la Universidad Central "Marta Abreu" de las Villas.
- IV. **EL ESTUDIANTE** una vez aprobada su tesis para la defensa, depositará una copia electrónica de la misma en el Repositorio Digital Institucional de la Universidad Central "Marta Abreu" de Las Villas.
- V. A partir de la defensa y aprobación del Trabajo de Diploma, la publicación total, parcial o la elaboración de cualquier obra que se derive de esta investigación por parte de **EL ESTUDIANTE**, contará con la coautoría de **EL TUTOR** y viceversa, resultando de referencia obligada esta obra en cualquier otra que se elabore. El incumplimiento de esta cláusula, puede llevar consigo el inicio de procesos de plagio. Todo lo anterior de acuerdo a la normativa de Derecho de Autor vigente en Cuba.

Y para que así conste se firma la presente en la Universidad Central "Marta Abreu" de Las Villas, a los 3 días del mes de diciembre del año 2023.

EL ESTUDIANTE

JEFE DE DEPARTAMENTO

TUTOR

TUTOR

## **PENSAMIENTO**

*“La perseverancia es el puente que conecta los sueños con la realidad.  
No se trata solo de la fuerza de voluntad, sino de la convicción profunda de que cada  
paso hacia adelante nos acerca un poco más a nuestros objetivos”*

*Anónimo*

## **DEDICATORIA**

A mi familia, en especial a mi mamá,  
por ser mi luz, mi guía, mi todo...

## AGRADECIMIENTOS

Mis agradecimientos infinitos y eterna admiración a la profe Marlen, que ha ejercido de tutora, de madre y psicóloga durante todo este tiempo, gracias por tantas atenciones, por la ayuda, la paciencia y por compartir conmigo sus conocimientos y experiencias.

A los radiólogos de los Hospitales de Santa Clara y Cienfuegos por sus enseñanzas, en especial a la doctora Zenia, por dedicarme su tiempo.

A los profesores que han formado parte del proyecto, por estar siempre al pendiente y dispuestos a ayudar.

A Dari, por su tiempo porque siempre está cuando lo necesito.

A Anay, Ernesto, Ili, Melisa, Hanlert y Darian que aportaron su granito de arena para que esto fuera posible.

A mi familia, por el apoyo, por confiar en mí, especialmente a mi mamá por el amor, su sacrificio y entrega, por ser mi guía y hacerme la persona que soy. Por impulsarme a nunca dejar de aprender.

A mi tía Madelyn, por ser mi segunda madre, porque no importa donde esté siempre puedo contar con ella.

A mi hermana que a pesar de hablar poco y estar lejos siempre la llevo presente y no para de hacerme sentir orgullosa.

A Juanca por ser como un padre para mí y apoyarme siempre.

A Luisma por su paciencia, el apoyo incondicional, por los ánimos y la confianza en mí, por tantos atardeceres y por siempre hacerme sonreír.

A los amigos que me ha dado la universidad, que más que amigos los considero familia, por todos los buenos momentos y por los que no han sido tan buenos y también han estado ahí, en especial a mis niñas de error 503 por acogerme desde el inicio y hacer de este tiempo todo lo especial que pudo ser.

Los llevo a todos en el corazón.

Mil gracias...

## RESUMEN

El cáncer constituye una de las causas principales de muerte en la sociedad actual. La tomografía computarizada es el medio que más se utiliza para diagnosticar y evaluar las respuestas a los tratamientos en esta enfermedad. Su interpretación para clasificar nódulos pulmonares es difícil debido a la complejidad y diversidad de los mismos. Los sistemas de diagnóstico asistido por computadora han demostrado su efectividad para apoyar a los radiólogos en la toma de decisiones. La radiómica tiene como objetivo extraer características de imágenes digitales automáticamente y desarrollar modelos, para predecir fenotipos de lesiones de una manera no invasiva. En este estudio se desarrolla un sistema basado en radiómica donde se analizan 275 nódulos de la base de datos LIDC-IDRI, se extraen 102 características radiómicas a través de *PyRadiomics*. Los clasificadores *Support Vector Machine* y *Random Forest* son entrenados y validados con dos vectores de características. Los mejores modelos de cada clasificador alcanzaron valores de sensibilidad, especificidad, precisión y exactitud por encima del 90%, demostrando la efectividad de este enfoque para la clasificación.

**Palabras clave:** radiómica, clasificación de nódulos pulmonares, PyRadiomics, Support Vector Machine, Random Forest

## **ABSTRACT**

*Cancer is one of the main causes of death at the present. Computed tomography is the most commonly used means to diagnose and evaluate responses to treatment in this disease. Its interpretation to classify pulmonary nodules is difficult due to their complexity and diversity. Computer-aided diagnosis systems have proven to be effective in supporting radiologists in their decision making. Radiomics aims to extract features from digital images automatically and develop models, to predict lesion phenotypes in a noninvasive way. In this study, a radiomics-based system is developed where 275 nodules from the LIDC-IDRI database are analyzed, 102 radiomic features are extracted through PyRadiomics. Support Vector Machine and Random Forest classifiers are trained and validated with two feature vectors. The best models of each classifier achieved sensitivity, specificity, precision and accuracy values above 90% demonstrating the effectiveness of this approach for classification.*

**Key words:** *radiomics, lung nodule classification, PyRadiomics, Support Vector Machine, Random Forest*

## GLOSARIO DE SIGLAS Y TÉRMINOS

**BD:** Base de datos

**CAD:** (del inglés, *Computer Aided Diagnosis*). Sistema de diagnóstico asistido por computadora.

**CNN:** (del inglés, *Convolutional Neural Network*). Red neuronal de convolución

**DL:** (del inglés, *Deep Learning*). Aprendizaje profundo

**FN:** Falso negativo

**FP:** Falso positivo

**IA:** Inteligencia artificial

**ML:** (del inglés, *Machine Learning*). Aprendizaje automático

**RF:** (del inglés, *Random Forest*). Bosque Aleatorio

**ROI:** (del inglés, *Region of Interest*). Región de interés.

**SVM:** (del inglés, *Support Vector Machine*). Máquina de Soporte Vectorial

**TC:** Tomografía computarizada

**VN:** Verdadero negativo

**VP:** Verdadero positivo

## TABLA DE CONTENIDOS

<b>INTRODUCCIÓN .....</b>	<b>1</b>
<b>CAPÍTULO 1. MARCO TEÓRICO .....</b>	<b>8</b>
<b>1.1 Tomografía computarizada .....</b>	<b>8</b>
1.1.1 <b>Calidad de imagen en TC .....</b>	<b>9</b>
<b>1.2 Cáncer de pulmón .....</b>	<b>11</b>
1.2.1 <b>Nódulos pulmonares.....</b>	<b>11</b>
<b>1.3 Sistemas CAD como soporte al estudio del cáncer de pulmón .....</b>	<b>13</b>
<b>1.4 Algoritmos de ML en imagenología médica.....</b>	<b>14</b>
1.4.1 <b>Algoritmo de bosque aleatorio .....</b>	<b>14</b>
1.4.2 <b>Redes neuronales convolucionales .....</b>	<b>15</b>
1.4.3 <b>Máquinas de soporte vectorial .....</b>	<b>16</b>
<b>1.5 Radiómica.....</b>	<b>16</b>
1.5.1 <b>Procedimiento de extracción de características radiómicas .....</b>	<b>18</b>
1.5.2 <b>Características radiómicas y su importancia.....</b>	<b>22</b>
1.5.3 <b>Aplicaciones de la radiómica en el estudio de los nódulos de pulmón .....</b>	<b>24</b>
<b>1.6 Conclusiones del capítulo.....</b>	<b>26</b>
<b>CAPÍTULO 2. MATERIALES Y MÉTODOS .....</b>	<b>28</b>
<b>2.1 Descripción general de la implementación del sistema CAD.....</b>	<b>28</b>
<b>2.2 Funciones y bibliotecas más utilizadas.....</b>	<b>30</b>
<b>2.3 Segmentación de los nódulos .....</b>	<b>30</b>
<b>2.4 Reconstrucción, etiquetado y preparación de los datos.....</b>	<b>35</b>
2.4.1 <b>Visualización de las anotaciones de malignidad mediante Pylidc .....</b>	<b>36</b>
<b>2.5 Extracción de características con PyRadiomics .....</b>	<b>38</b>
<b>2.6 Reducción de características.....</b>	<b>39</b>
<b>2.7 Implementación de los clasificadores .....</b>	<b>40</b>
2.7.1 <b>Implementación de máquinas de soporte vectorial (SVM).....</b>	<b>41</b>
2.7.2 <b>Implementación de bosques aleatorios (RF) .....</b>	<b>42</b>
<b>2.8 Entrenamiento y validación.....</b>	<b>42</b>
<b>2.9 Base de datos utilizada .....</b>	<b>43</b>
<b>2.10 Hardware y software utilizado.....</b>	<b>44</b>
<b>CAPÍTULO 3. RESULTADOS Y DISCUSIÓN .....</b>	<b>45</b>
<b>3.1 Análisis de las características radiómicas extraídas .....</b>	<b>45</b>

<b>3.2</b>	<b>Disminución de la dimensionalidad del problema .....</b>	<b>45</b>
<b>3.3</b>	<b>Resultados de los clasificadores con etiquetado binario y no binario .....</b>	<b>48</b>
<b>3.3.1</b>	<b>Valoración del mejor desempeño.....</b>	<b>53</b>
<b>3.4</b>	<b>Discusión .....</b>	<b>55</b>
<b>3.5</b>	<b>Análisis económico y medioambiental.....</b>	<b>57</b>
<b>3.6</b>	<b>Conclusiones del capítulo.....</b>	<b>58</b>
	<b>CONCLUSIONES y RECOMENDACIONES .....</b>	<b>60</b>
	<b>REFERENCIAS BIBLIOGRÁFICAS .....</b>	<b>61</b>

## INTRODUCCIÓN

El cáncer constituye una de las principales preocupaciones de la sociedad actual debido a su impacto en la salud y calidad de vida de las personas. Las variaciones de la enfermedad en diferentes partes del mundo, el papel de las instalaciones médicas disponibles y otros factores socioeconómicos, han influido en la gestión adecuada de esta enfermedad [1].

Según la Organización Mundial de la Salud, el cáncer es una de las principales causas de muerte en el mundo. Las intervenciones de control del cáncer más específicas y la inversión en la mejora de la detección precoz y el tratamiento, facilitan la reducción de la mortalidad por cáncer [2].

En el caso particular de Cuba, el cáncer también constituye una de las causas principales de muerte. La población cubana está compuesta por 11 147 405 habitantes, de los cuales el 21,6 % tiene 60 años o más de edad [3]. Según la Oficina Nacional de Estadísticas e Información de la República de Cuba (ONEI), el cáncer ha sido la segunda causa de muerte para todas las edades, más elevada durante los últimos 5 años, con 26 791 casos en el 2021 [4]. La tasa de mortalidad más elevada por tipo de cáncer en ambos sexos, según el Anuario estadístico de salud, emitido por el Ministerio de Salud Pública de Cuba (MINSAP), corresponde a los tumores malignos de tráquea, bronquios y pulmón [3]. Este tipo de cáncer normalmente comienza como un nódulo pulmonar pequeño.

Según el glosario de términos de imagen torácica, propuesto por la *Fleischner Society* [5], un nódulo pulmonar se define como una opacidad aproximadamente redondeada, que mide hasta 3 cm de diámetro. Los nódulos pulmonares pueden ser solitarios o múltiples. Es muy probable que un nódulo pulmonar en crecimiento sea maligno [6].

La detección de los nódulos pulmonares es un indicador radiológico de primer orden para el diagnóstico precoz. Diagnosticar los nódulos pulmonares en una fase temprana es crucial para mejorar las tasas de supervivencia de los pacientes [7].

Las imágenes médicas constituyen una de las fuentes de mayor importancia, por cuanto ofrecen un apoyo integral del acto médico: el diagnóstico y el seguimiento. Algunas formas de obtención de imágenes médicas son las imágenes por resonancia magnética, ultrasonido, radiografía y tomografía computarizada (TC o CT, del inglés *Computed Tomography*) [8]. Su principal objetivo es generar información de gran importancia para la caracterización de

la fisiología y/o anatomía de diversos órganos o partes del cuerpo y para optimizar la toma de decisiones, evitando costosos tratamientos y aumentando la tasa de éxito en intervenciones quirúrgicas.

La radiografía es la primera opción de modalidad para la detección de enfermedades pulmonares. Sin embargo, debido a los efectos de proyección, las radiografías de tórax solo pueden ser adecuadas para detectar nódulos grandes en comparación con la TC. Por lo tanto, por sí solas no son óptimas para la caracterización temprana del cáncer de pulmón [9]. Algunos de los problemas asociados a la detección de cáncer de pulmón con esta técnica, son la poca visibilidad debido a su pequeño tamaño, márgenes mal definidos y estructuras superpuestas [10].

La TC es el medio diagnóstico que más se utiliza en Cuba y en el mundo para diagnosticar y evaluar las respuestas a los tratamientos en el carcinoma de pulmón de células no pequeñas (CPCNP) [11]. Las imágenes generadas por los equipos médicos de esta modalidad poseen entre sus ventajas, la gran resolución espacial (cantidad de píxeles por cm) y la densidad o profundidad (niveles de grises que se pueden representar) que posee la imagen [12]. Permite explorar todo el parénquima pulmonar con mayor contraste y mejor detalle, ya que se analiza el tejido por cortes, visualizando en cada uno solo su contenido y quitando lo que está por encima y por debajo de este [13], lo cual incrementa la capacidad de detección.

A pesar de lo antes mencionado, la tecnología de TC, en dependencia del espesor del corte programado, genera una gran cantidad de imágenes que deben ser revisadas por especialistas en cada estudio. Debido a la variabilidad en forma, textura y características que presentan los nódulos pulmonares, es complicado para un radiólogo detectar, a simple vista, la malignidad de una lesión, lo que conlleva a una solución invasiva: realizar la biopsia del tejido. En este contexto, el aprendizaje automático a partir de cortes de TC puede ser de gran utilidad.

El hecho de que los sistemas modernos de TC suponen un riesgo relativamente bajo para la salud humana, en términos de radiación impartida al paciente, ha generalizado su uso como herramienta de diagnóstico médico [14]. Por tanto, si se apoya la TC con una herramienta informática, como el diagnóstico asistido por ordenador, se puede contribuir a realizar la clasificación de nódulos pulmonares de forma temprana y con ello mejorar la tasa de supervivencia por cáncer de pulmón, lo que implicaría un importante impacto clínico y social [15].

Para ayudar a los radiólogos en su tarea, para la detección de diversas patologías con alta eficiencia y bajas tasas de falsos positivos y negativos, en etapas más tempranas de la evolución de la enfermedades [16],[17], se han desarrollado los sistemas de diagnóstico asistido por computadora (CAD, del inglés *Computed Aided Diagnosis*).

Los CAD constan por lo general de cinco etapas: adquisición de imágenes, preprocesamiento, segmentación, extracción de características y clasificación [7]. Algunos solo se dedican a detectar y otros llegan hasta la separación en clases benigno-maligno. Los primeros sistemas se construyeron solo a base de procesamiento digital de imágenes [18]. En la actualidad, estos sistemas CAD han alcanzado mayores beneficios, gracias al desarrollo de la inteligencia artificial aplicada a las imágenes médicas. En la última década por ejemplo, han aumentado las investigaciones relacionadas con el aprendizaje automático, (ML, del inglés *Machine Learning*) y un tipo especial de este, el aprendizaje profundo (DL, del inglés *Deep Learning*) [19].

El ML es una rama de la inteligencia artificial que persigue desarrollar en las máquinas la habilidad de aprender a partir de un conjunto de datos, que pueden ser imágenes. Los algoritmos de ML deben procesar grandes cantidades de imágenes etiquetadas, lo cual es un problema en Medicina, donde por lo general, no existe tal variabilidad [20].

Hay muchas arquitecturas y enfoques que han demostrado ser particularmente efectivos para diversas tareas de diagnóstico y para extraer características de ciertos tipos de datos. Entre estos podemos citar las redes neuronales recurrentes (RNN, del inglés *Recurrent Neural Networks*) para series temporales y las redes neuronales convolucionales (CNN, del inglés *Convolutional Neural Networks*) para tareas de visión computarizada [21].

La red neuronal convolucional es el tipo de arquitectura de aprendizaje más profundo que existe. Tiene como ventajas la amplia gama de aplicaciones, alta precisión y rápida velocidad de análisis. Su precisión ha superado a la del ser humano en algunos aspectos [22], [23]. Por esta razón, en Radiología se han propuesto y estudiado varias aplicaciones clínicas con CNN para la clasificación, detección y segmentación de lesiones [23].

En la actualidad la Medicina avanza hacia la denominada "Medicina de precisión". Esta está destinada a transformar la investigación clínica, biomédica y el propio cuidado de los pacientes, ofreciendo nuevas oportunidades de mejora a los Sistemas de salud públicos. La esencia de este tipo de Medicina es tomar decisiones personalizadas para la prevención, el

diagnóstico y el tratamiento de enfermedades, sobre la base de datos de pacientes individuales, recopilados a través de mediciones de alta precisión y extracción e integración eficiente de la información. Surge así "la ómica" [24].

La incorporación del término ómica a los campos de la ciencia tiene su origen en la ciencia básica. Sin embargo, ahora es un sufijo muy utilizado en la investigación de Medicina clínica, para denotar el concepto de examinar grandes volúmenes de datos complejos, para identificar características o resultados precisos y extraer información valiosa [24], [25].

La radiómica en particular, es hoy una revolución en la tecnología tradicional de imágenes identificables visualmente y constituye una nueva rama. Tiene como objetivo extraer características de imágenes digitales automáticamente y desarrollar modelos, para predecir fenotipos de lesiones de una manera no invasiva [24], [26].

El término radiómica se utiliza en dos contextos diversos: para referirse al estudio de la variación genética asociada con la respuesta a la radiación, o para referirse a la correlación existente entre las características de las imágenes del cáncer y su expresión génica. Por lo tanto, asume que las características de la imagen están relacionadas con las firmas de genes que posee [27].

En el ámbito de los nódulos pulmonares y el cáncer de pulmón, el objetivo de este método es extraer características cuantitativas automatizadas de las imágenes, que puedan predecir la naturaleza de los nódulos y los tumores de forma no invasiva. Por esta razón, tiene el potencial de revolucionar el diagnóstico, la vigilancia y la planificación del tratamiento, permitiendo una gestión personalizada de forma no invasiva y rentable [25].

Actualmente, el problema de detección y clasificación de nódulos presenta dos tendencias de investigación a nivel mundial:

- 1) Con métodos de aprendizaje profundo [28], que consisten en construir la estructura de una red profunda convolucional, entrenar el modelo y, a continuación, utilizar el modelo entrenado para clasificar los datos; un procedimiento que obviamente involucra un volumen de datos muy elevado.
- 2) Con métodos tradicionales de aprendizaje automático [29], (que por lo general utilizan menos datos que las CNN). Estos se basan en la extracción de la región de interés (ROI) y su vector de características, para servir como entrada a clasificadores. Obviamente, tanto la

selección de las características, como la de los clasificadores, tendrán un impacto en los resultados [30].

En la presente investigación se asume la segunda tendencia. En este caso, los clasificadores estarán alimentados por características radiómicas, a extraer en ROI, a partir de cortes de TC. Ambas tendencias, aunque implican múltiples procedimientos técnicos, se basan en características de alta dimensión, que son extraídas de las imágenes para tareas de diagnóstico y predicción. En el primer caso, la red las selecciona por sí solas y en el segundo deben ser escogidas y calculadas previo a la clasificación [31].

Para que la radiómica o los algoritmos de aprendizaje profundo se conviertan en una herramienta clínica válida, su rendimiento debe validarse mediante pruebas clínicas. En este sentido, los esfuerzos de investigación deben abordar el impacto clínico de ambas tendencias [32].

En relación con lo explicado anteriormente, el tema de la presente investigación se relaciona con la clasificación de nódulos pulmonares en benignos o malignos, a partir de características de radiómica y algoritmos de ML. El sistema a desarrollar debe ser capaz de analizar tomografías computarizadas de tórax. Será empleado por radiólogos como herramienta de apoyo a la identificación y/o verificación de lesiones, para agilizar y mejorar la precisión de sus diagnósticos.

### **Justificación de la investigación**

Actualmente existe una directiva del estado cubano para el diseño y desarrollo de arquitectura informática, basada en la Medicina personalizada y las ómicas, en su sentido más amplio [33]. Se ha definido lo anterior entre las prioridades para mejorar la atención a la salud del pueblo. Específicamente, el cáncer, requiere especial atención, por lo que se estimula la investigación y la innovación.

Debido a la exactitud y alto contraste que brinda la TC, es una de las técnicas imagenológicas más empleadas para la detección de nódulos pulmonares. Sin embargo, la complejidad y diversidad de estos, hacen que la tarea de clasificación benigno-maligno aún sea un reto y un campo abierto a la investigación. Implementar un sistema de diagnóstico asistido por ordenador que clasifique las lesiones, a partir de las imágenes de TC, sería de gran ayuda para los especialistas en Radiología, cuyo *gold standar* hoy, para esta tarea, lo es solo a partir de un método invasivo para el paciente (la biopsia del nódulo detectado). En particular

utilizar características radiómicas como entrada a clasificadores de ML, podría reducir las tasas de falsos positivos y negativos en los diagnósticos de cáncer pulmonar, y además, constituiría un paso de avance hacia la Medicina de precisión personalizada en Cuba.

### **Objeto y campo de la investigación**

El objeto de estudio son las imágenes tomográficas con nódulos pulmonares. El campo de investigación son los algoritmos de aprendizaje automático a partir de características radiómicas para la clasificación de tales nódulos.

### **Problema de investigación**

Las tomografías computarizadas de tórax resultan en ocasiones de difícil interpretación para expertos radiólogos y neumólogos, en la tarea de clasificar los nódulos pulmonares detectados, en malignos o benignos, debido a la amplia diversidad de estos y a la variabilidad de la enfermedad entre pacientes.

### **Hipótesis**

Con la ayuda de clasificadores de aprendizaje automático, a partir de características de radiómica, es posible obtener un sistema de diagnóstico asistido por ordenador, eficaz y eficiente computacionalmente, que permita la clasificación de nódulos pulmonares con niveles de exactitud, sensibilidad y especificidad por encima del 85 %.

### **Objetivo general**

Diseñar un sistema CAD que clasifique en maligno-benigno los nódulos de pulmón a partir de una data anotada de TC, utilizando características de radiómica y clasificadores de ML.

### **Objetivos específicos**

1. Seleccionar un grupo de características radiómicas que, utilizando clasificadores de ML, permita la clasificación de nódulos pulmonares, de una forma eficiente computacionalmente.
2. Entrenar los clasificadores para el 80 % de una data anotada, de modo que clasifique la data interna con alta precisión.
3. Evaluar en términos de eficacia y eficiencia computacional los resultados obtenidos en la clasificación, a partir de los resultados de la prueba de validación.
4. Comparar los resultados contra los obtenidos anteriormente por el método de fractales.

---

El Trabajo de diploma consta de 70 páginas, con 7 tablas y 21 figuras incluidas, así como 8 anexos. Se revisaron 120 fuentes bibliográficas.

## CAPÍTULO 1. MARCO TEÓRICO

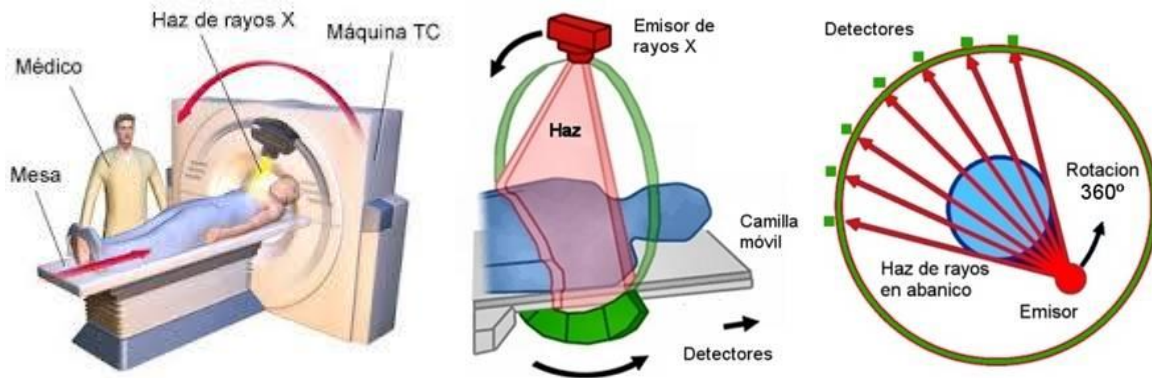
En el presente capítulo se abordan varios elementos conceptuales, tecnológicos y de revisión bibliográfica del tema a tratar. Además, se describen algunos algoritmos de *Machine Learning* y el actual papel de la radiómica como herramienta útil para el análisis de las imágenes médicas y su valor predictivo en la Medicina.

### 1.1 Tomografía computarizada

Las imágenes médicas son representaciones gráficas bidimensionales o tridimensionales de una estructura, región, órgano o tejido del cuerpo humano. Hoy en día, son una parte integral del examen de diagnóstico de un paciente en muchas situaciones clínicas. Se ha demostrado que las técnicas de imágenes médicas son eficaces para detectar nódulos pulmonares. De los equipos de obtención de imágenes, la TC es actualmente uno de los más eficaces para detectar estos nódulos, por su sensibilidad y nivel de detalle de la enfermedad pulmonar [34].

La tomografía computarizada utiliza rayos X para obtener información estructural del cuerpo humano, proporcionando vistas detalladas de los tejidos blandos, incluidos los vasos sanguíneos, el tejido muscular y los órganos [35], [36]. Puede revelar los detalles internos en tres dimensiones del organismo de forma no invasiva.

Durante el procedimiento de TC, el paciente es colocado en una mesa que se mueve lentamente a través de un anillo llamado tomógrafo. El tomógrafo emite rayos X hacia el cuerpo del paciente desde varios ángulos, mientras que un detector en el otro lado del anillo, mide la cantidad de radiación que pasa a través del cuerpo, obteniéndose una serie de radiografías bidimensionales del cuerpo. La computadora utiliza esta información para crear una imagen transversal o "corte" del cuerpo en cada ángulo, que se pueden combinar para crear una imagen 3D en escala de grises, de la estructura interna del cuerpo [14], [37]. En la Figura 1.1 se observa lo explicado anteriormente.



**Figura 1.1** Esquema de adquisición de datos en un tomógrafo [38]

Físicamente se trata de una corriente eléctrica que calienta el filamento de un cátodo, el cual desprende electrones dentro de un tubo al vacío. Mientras más alta es esta corriente (mA), mayor número de electrones son desprendidos. Estos electrones se aceleran dentro del tubo de rayos X, gracias a una diferencia de potencial respecto al ánodo (kV). Al impactar en este, arrancan los rayos X (radiación electromagnética), que se dirigen en forma de haz al cuerpo del paciente, para ser atenuados de forma diferente por los diferentes tejidos del mismo. Se genera así una imagen con contrastes entre tejidos, que son los diversos tonos de grises. Como el proceso se repite para diversos ángulos, para cada uno se llena una fila de datos (proyecciones) en una matriz de adquisición, que luego hay que reconstruir aplicando la Transformada de Radón primero y la de Fourier después [39]. De esta forma se generan los cortes 2D (imágenes), uno por cada proyección, y la composición de todos, imagen 3D o tomografía [40].

### 1.1.1 Calidad de imagen en TC

Tener buena calidad de imagen es imprescindible para poder diagnosticar. Desde el punto de vista clínico significa que el radiólogo pueda apreciar los detalles contenidos en una Norma médica, existente para cada tipo de estudio [41], a partir de los cortes 2D o de la imagen reconstruida en 3D. Sin embargo, desde el punto de vista de la Física y la Ingeniería, calidad de imagen implica tener fundamentalmente altos valores de resolución espacial, bajos niveles de ruido cuántico y alto contraste imagen [42]. Estos se obtienen de medir niveles de gris sobre regiones de interés en la imagen digital. Dichos niveles se relacionan directamente con la cantidad de rayos X que impactaron y fueron absorbidos en cada vóxel del paciente. Los

valores medidos en cada píxel de la imagen, o en el grupo de píxeles que conforman cada región de interés (ROI), son introducidos en ecuaciones matemáticas que permiten tener un valor objetivo de cada parámetro de calidad de imagen.

Así, por ejemplo, la resolución espacial es la capacidad del escáner de TC para reproducir fielmente los detalles de un objeto. Depende del tamaño de la matriz de reconstrucción, ancho del detector, espesor de corte, distancia del objeto al detector y tamaño del punto focal. En la TC, la función de transferencia de modulación se utiliza para calcular esta característica [42], [43], que no es más que calcular la Transformada de Fourier de la función de dispersión de punto, línea o escalón, del contenido de la ROI escogida para esto sobre la imagen.

El ruido cuántico por su parte, en TC es el grado de incertidumbre en la medición de la atenuación del haz de rayos X que atraviesa al paciente [42], [43]. Depende de la cantidad de fotones de rayos X que caen sobre cada porción de los detectores. Por tanto, es la fluctuación estadística o desviación típica de los números de CT en una ROI [44]. Como la estadística responde a la distribución de Poisson, en primera aproximación, la estimación del ruido cuántico sería la raíz cuadrada de los conteos obtenidos en una ROI.

El ruido está determinado por varios factores, como la filtración del haz, el tamaño de campo de visión (FOV, *del inglés field of view*), el tamaño del píxel, el espesor del corte, el algoritmo de reconstrucción y la dosis que recibe el paciente, determinada por la corriente que alimenta el tubo de rayos X. Mientras más alto es el mA, mayor estadística de conteos en una ROI y menor el valor del ruido asociado [43].

Finalmente, el contraste de imagen en TC es la capacidad de diferenciar pequeños cambios en la atenuación lineal de los rayos X entre diferentes tejidos o entre tejidos y lesiones. Esta es la principal ventaja de la TC sobre las radiografías convencionales [45].

La resolución de contraste de un sistema de imagen determina el detalle de contraste que se puede reproducir visualmente cuando hay una pequeña diferencia de la densidad en la región circundante, mostrando objetos más sutiles en la imagen [46]. La resolución de contraste es altamente degradada por el ruido [37]. Para identificar de forma fiable una estructura, la relación señal a ruido debe ser superior a 5:1 [42].

## 1.2 Cáncer de pulmón

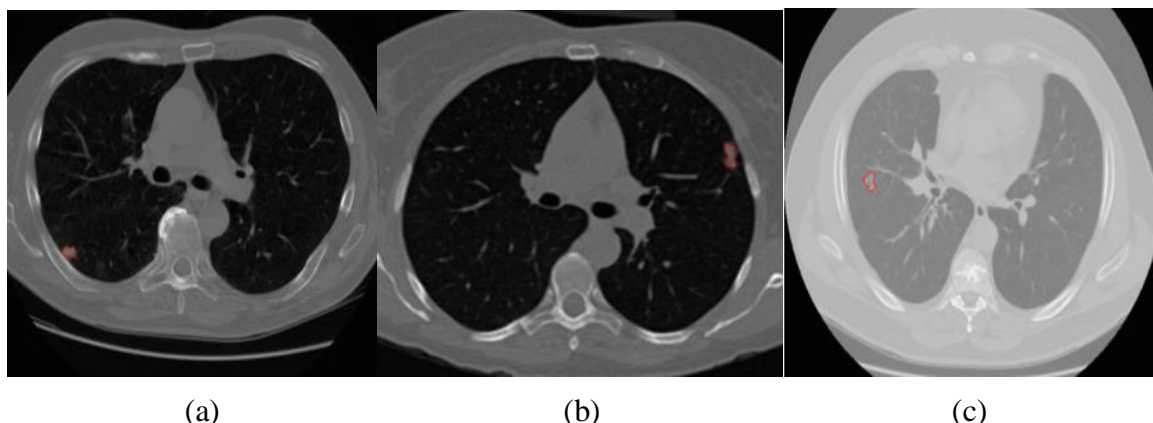
El cáncer de pulmón es provocado cuando las células tumorales, que son capaces de invadir los tejidos sanos, se asientan en los pulmones. Es originado cuando los mecanismos de control de la división celular se alteran y las células se dividen descontroladamente. Las anomalías aparecen primero como nódulos. Si un nódulo pulmonar no se detecta a tiempo, las células cancerosas pueden extenderse a otras zonas del cuerpo, provocando un fenómeno conocido como metástasis, que en muchos casos es mortal [47], [48].

Existen fundamentalmente tres clasificaciones de cáncer de pulmón: cáncer de pulmón de células pequeñas (CPCP), cáncer de pulmón de células no pequeñas (CPNP) y tumor carcinoide de pulmón. El primero de estos tipos constituye alrededor del 10% al 15% de los cánceres de pulmón, pero es también conocido como el cáncer más agresivo, debido a su rápido crecimiento en comparación con otros tipos de cáncer de pulmón. Sin embargo, alrededor del 85% del cáncer de pulmón es del tipo CPNP. El tumor carcinoide de pulmón es un tipo de cáncer poco común y tiende a crecer más lentamente en comparación con los demás [49].

Los términos "célula pequeña" y "célula grande" son simplemente términos descriptivos de cómo aparecen las células cancerosas bajo un microscopio. Estas características de las células cancerosas ayudan al médico a determinar el tipo de cáncer, el grado de anormalidad en las células y dónde se originó [49].

### 1.2.1 Nódulos pulmonares

Un nódulo pulmonar se define como una lesión pulmonar menor de 3 cm de diámetro, que está completamente rodeada de parénquima pulmonar, sin otras anomalías [6], [50]. Las lesiones redondeadas que miden más de 3 cm de diámetro se denominan masas pulmonares y deben considerarse indicativas de cáncer de pulmón, hasta que se demuestre lo contrario histológicamente [6]. En la Figura 1.2 se muestran ejemplos de nódulos pulmonares en imágenes 2D de TC.



**Figura 1.2** Imágenes de pulmón con nódulos interiores delimitados en rojo (a) nódulos muy pequeños, (b) y (c) nódulos más grandes. Extraído de [26].

Los nódulos pulmonares solitarios (NPS) se pueden identificar por la forma, calcificación, esfericidad, lobulación, espiculación, textura y localización de su estructura interna [51]. Los primeros pueden ser de tejido blando, líquido, grasa y aire. Los nódulos con esfericidad son lineales, ovoides y redondos. Los lobulados pueden ser claramente lobulados o no lobulados. Los espiculados se identifican por ser no espiculados o con marcada espiculación. Según la textura, los nódulos se caracterizan en sólidos no calcificados, calcificados y vidrio deslustrado. Dependiendo de la ubicación, pueden ser centrales o periféricos [12].

El análisis de la forma de los NPS puede ayudar a los radiólogos a distinguir entre nódulos benignos y malignos. Los pequeños nódulos que tienen bordes bien definidos son característicos de las lesiones benignas, pero no se limitan a ellas. Por el contrario, las lesiones malignas se caracterizan por un borde lobulado o un borde irregular o fisurado con distorsión de los vasos adyacentes. Las calcificaciones difusas, nodulares centrales, laminadas, o en palomita de maíz, sugieren benignidad. En cambio, en los nódulos malignos se describen calcificaciones excéntricas o moteadas [52].

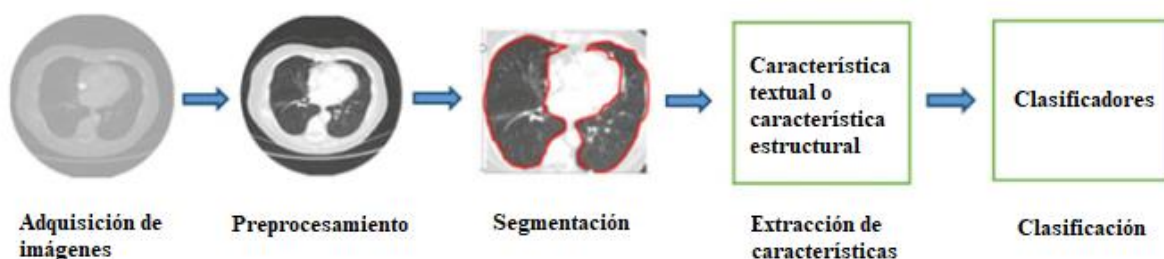
Según [7] debido a que el tamaño del nódulo está relacionado con la malignidad, la medición precisa del diámetro del nódulo es importante para el diagnóstico. La base de datos del Programa de acción temprana contra el cáncer de pulmón (por sus siglas en inglés ELCAP) sugiere que la incidencia de malignidad es del 1 % para nódulos de menos de 5 mm de diámetro, del 24 % para nódulos de 6 mm a 10 mm, del 33 % para nódulos de 11 a 20 mm y del 80 % para nódulos mayores de 20 mm.

### 1.3 Sistemas CAD como soporte al estudio del cáncer de pulmón

Los sistemas de diagnóstico asistido por computadora (CAD) son sistemas potentes desarrollados para la detección y caracterización de diversas lesiones en el campo del diagnóstico del cáncer de pulmón. El objetivo principal de tales sistemas es ayudar al radiólogo en las distintas etapas del análisis y proporcionar una segunda opinión para las decisiones finales [35], [47], [53], [54]. A lo largo de su historia, han evolucionado de ser un conjunto de reglas, a incorporar el aprendizaje automático y el aprendizaje profundo.

Los sistemas CAD también están diseñados para superar los errores de percepción y reducir los falsos negativos. Varios estudios han demostrado que agregar un sistema CAD al proceso de diagnóstico, puede mejorar la eficiencia en la detección de lesiones, al reducir la variabilidad entre observadores [55]. Además, los sistemas CAD brindan soporte cuantitativo para la toma de decisiones clínicas, como son: las recomendaciones de biopsias [56], ayudar a realizar controles de diagnóstico, reducir las biopsias falso positivas innecesarias, entre otras [57]. También se pueden utilizar para distinguir entre tumores malignos y benignos [58]. Los sistemas CAD actuales de tomografía computarizada buscan densidades pulmonares con características físicas específicas, que representan nódulos pulmonares. Por ello, estas aplicaciones para la detección de nódulos pulmonares se han convertido en un área activa de investigación.

El sistema de detección de cáncer de pulmón consta principalmente de cinco etapas: (i) adquisición de imágenes, (ii) preprocesamiento, (iii) segmentación, (iv) extracción de características y (v) clasificación [59]. Para comprender todo el proceso, se requiere tener una comprensión básica de todas las etapas involucradas en el sistema de detección [49]. La Figura 1.3 muestra un Sistema CAD con sus diferentes etapas para la detección de nódulos pulmonares.



**Figura 1.3** Etapas de un sistema de detección de nódulos pulmonares. Tomado de [47].

Para la implementación de sistemas CAD basados en imágenes médicas, se requiere disponer de altos volúmenes de estas. Para esto se pueden utilizar bases de datos públicas. Estas bases de datos son una herramienta importante para el desarrollo de algoritmos de detección de nódulos pulmonares automatizados, que pueden ayudar a los radiólogos a identificarlos de manera eficiente y precisa [47].

Durante el preprocesamiento se realiza la segmentación de la imagen. Se basa en diferentes características de la imagen y tiene como objetivo agrupar píxeles en regiones de imagen significativas [59], [60]. Una vez preprocesada y segmentada la imagen, se procede a la segmentación de los nódulos pulmonares. Esta es una etapa muy crítica e importante en el sistema de detección de cáncer de pulmón, porque el espacio de características posteriores solo se extrae del volumen segmentado. Algunos de los métodos utilizados para la segmentación pueden ser: crecimiento de regiones, umbralización y segmentación de cuencas [49].

La extracción de características es parte fundamental de cualquier sistema CAD. Esta etapa se hace necesaria cuando los datos son demasiado grandes [49]. Se utilizan de hecho, en el proceso de clasificación de nódulos. Dicha clasificación es el componente final y vital de un sistema CAD e implica la diferenciación de nódulos benignos y malignos [7]. En la actualidad estas últimas etapas se desarrollan con ayuda del ML.

## **1.4 Algoritmos de ML en imagenología médica**

El aprendizaje automático es una rama de la inteligencia artificial (IA) que permite la extracción de patrones significativos de las imágenes. En el contexto de las imágenes médicas, la idea de tener una computadora que realice tareas repetitivas de manera consistente e incansable es extremadamente atractiva para todos. Dentro de los más aceptados por la comunidad científica radiológica están las redes neuronales convolucionales, los algoritmos de bosque aleatorio y las máquinas de soporte vectorial [61].

### **1.4.1 Algoritmo de bosque aleatorio**

Un modelo basado en árbol implica la partición recursiva del conjunto de datos dado en dos

grupos, en función de un determinado criterio, hasta que se cumpla una condición de parada predeterminada. En la parte inferior de los árboles de decisión se encuentran los llamados nodos de hojas u hojas [62].

El modelo de bosque aleatorio (RF, del inglés *Random Forest*) es un algoritmo de aprendizaje basado en árboles de conjunto; es decir, el algoritmo promedia las predicciones sobre muchos árboles individuales. Los árboles individuales se construyen sobre muestras de arranque en lugar de sobre la muestra original. Esto se denomina agregación (*bootstrap*) o simplemente embolsado, y reduce el sobreajuste [62].

El basamento del RF es que la aleatorización con muchos árboles de decisión puede mejorar la precisión de la clasificación general, aumentando el grado de selección de características que tienen el mayor efecto en la clasificación [60]. Los bosques aleatorios, como algoritmos de aprendizaje de conjunto, son adecuados para conjuntos de datos medianos y grandes. Es importante tener en cuenta que, en el algoritmo de RF, el tamaño del subconjunto de variables predictoras es crucial para controlar la profundidad final de los árboles. Por lo tanto, es un parámetro que debe ajustarse durante la selección del modelo [62].

#### **1.4.2 Redes neuronales convolucionales**

Para tareas de visión computarizadas son empleadas las CNN [63], que son algoritmos de aprendizaje profundo que toman imágenes de entrada y las convolucionan con filtros o núcleos para extraer características [64]. La CNN tiene un excelente desempeño en problemas de aprendizaje automático, especialmente las aplicaciones que se ocupan de los datos de imágenes [65], [66].

Las redes convolucionales operan sobre tensores 3D, llamados mapas de características, con dos ejes espaciales (altura y anchura) y un eje de profundidad. La operación de convolución extrae parches de su mapa de características de entrada y aplica la misma transformación a todos estos parches, convirtiéndolos en vectores 1D, en forma de profundidad. Estos vectores son ensamblados, produciendo un mapa de características de salida. Este mapa sigue siendo un tensor 3D. Por tanto, tiene una anchura y una altura, pero su profundidad puede ser arbitraria, porque la profundidad de salida es un parámetro de la capa. Los distintos canales de ese eje de profundidad representan filtros. Estos filtros codifican aspectos específicos de

los datos de entrada, para que ciertas características se vuelvan dominantes. Cabe destacar que cada posición espacial en el mapa de características de salida corresponde a la misma posición en el mapa de características de entrada [67].

El problema fundamental de este método radica en que para evitar su sobreajuste requiere de un número muy alto de datos de entrenamiento, que a menudo sobrepasa las disponibilidades en Medicina [67], [68].

### 1.4.3 Máquinas de soporte vectorial

La máquina de soporte vectorial (SVM, del inglés *Support Vector Machine*) es un modelo que se utiliza para la regresión y la clasificación. SVM realiza la clasificación definiendo un nivel de separación entre clases distintivas. En la práctica, SVM se entrena con datos etiquetados, lo que también es conocido como aprendizaje supervisado, donde el algoritmo crea un hiperplano óptimo para clasificar los datos de prueba en etiquetas dadas [69].

Utilizando la teoría del aprendizaje estadístico, SVM busca una hipótesis regularizada que se ajuste bien a los datos existentes sin sobreajustarse. SVM tiene muy pocos parámetros libres y se puede optimizar utilizando la teoría de la generalización, sin una validación separada durante el entrenamiento. La principal limitación de la SVM es que el resultado es dependiente de la elección del *kernel*. Sin embargo, tiene como ventaja que para conjuntos de entrenamiento grandes, generalmente escoge una pequeña cantidad de vectores de soporte, lo que minimiza los requisitos computacionales durante las pruebas [70].

## 1.5 Radiómica

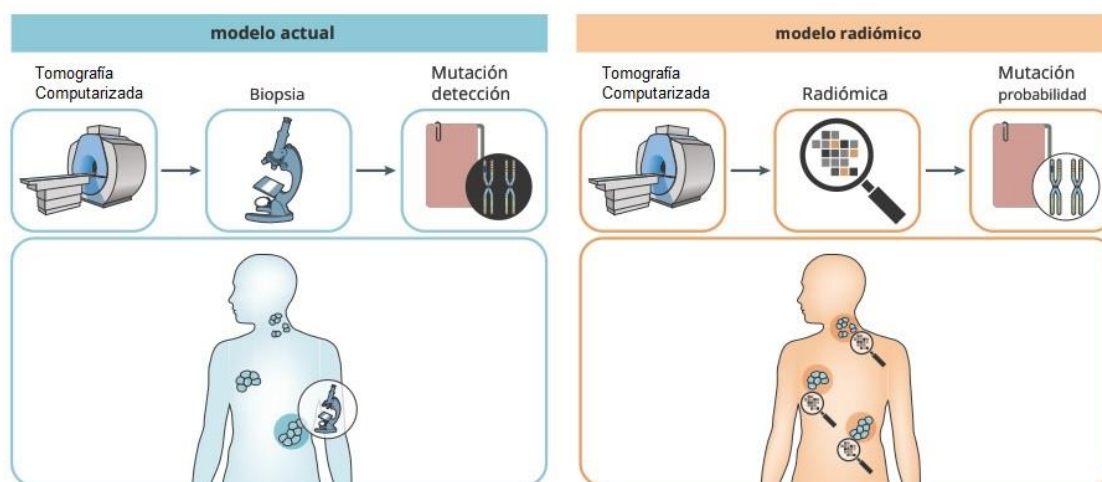
La TC permite la extracción semicuantitativa de las características anatómicas y morfológicas bidimensionales de los tumores, así como detalles fisiopatológicos tales como variaciones genéticas y funciones celulares, lo que facilita la elección individualizada del tratamiento [24]. Sin embargo, no puede predecir la heterogeneidad del tumor [44]. Por lo tanto, existe una necesidad urgente de desarrollar técnicas de imagen más sistemáticas y completas.

En 2012, el científico holandés P. Lambin, propuso por primera vez el término "radiómica" y lo definió como: extraer una gran cantidad de características a partir de imágenes

radiográficas con un enfoque de alto rendimiento, para explorar posibles vínculos con la biología y los resultados clínicos [27].

En la actualidad la radiómica es un campo emergente de investigación, centrado en el desarrollo de nuevos biomarcadores, basados en análisis de imágenes radiológicas a partir de datos. Se basa en la hipótesis de que las imágenes médicas reflejan características fisiopatológicas subyacentes y, por tanto, los análisis cuantitativos pueden ser útiles para describir la biología del volumen de imágenes. La extracción automatizada de una gran cantidad de características, puede ser informativa para el diagnóstico de enfermedades, el pronóstico y la respuesta al tratamiento [71].

La radiómica se diseñó para descifrar la heterogeneidad inherente, las características genéticas y otros fenotipos de una lesión, que permita mejorar el tratamiento del paciente. La TC es la modalidad de imagen más común para el análisis radiómico y permite una fácil comparación entre instituciones [24]. Por esta razón implica un nuevo modelo de análisis de las imágenes. La Figura 1.4 compara el modelo radiómico con el tradicional.



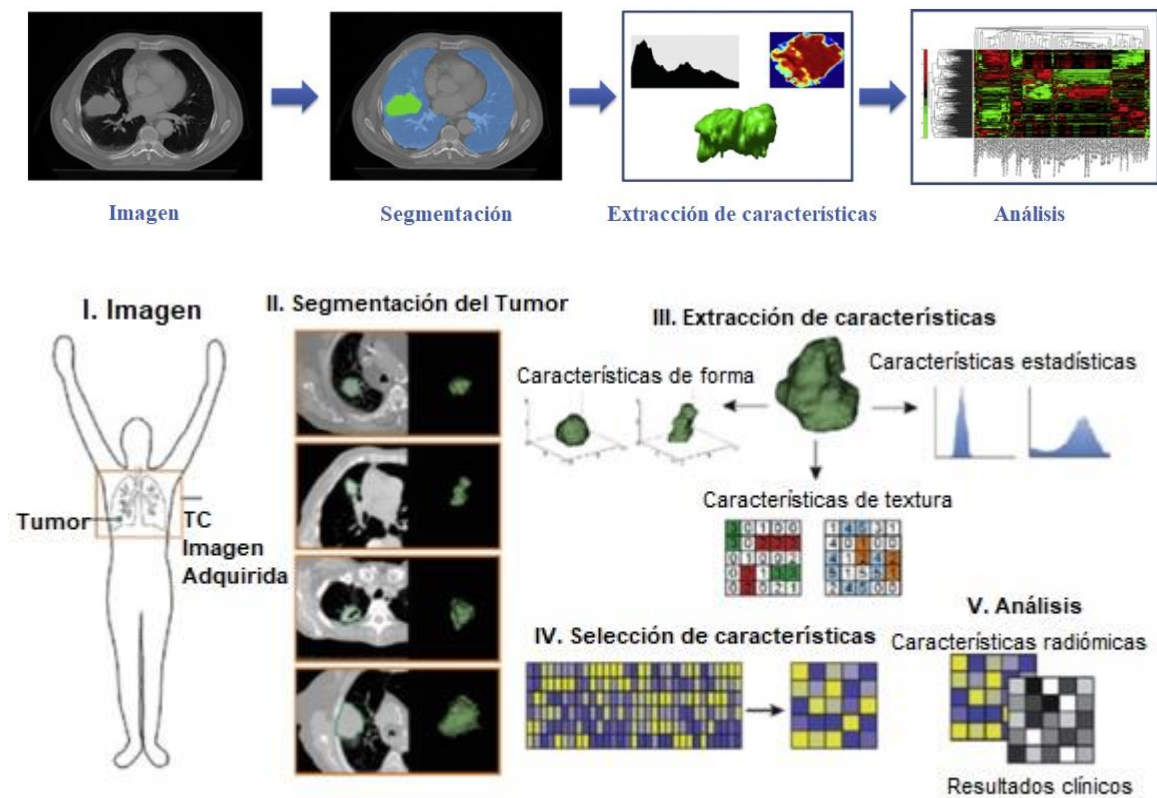
**Figura 1.4** Comparación del modelo actual tradicional con el modelo radiómico [72].

Se espera que la radiómica sea fundamental para la Medicina de precisión. Su esencia es tomar decisiones personalizadas para la prevención, el diagnóstico y el tratamiento de enfermedades, sobre la base de datos de pacientes individuales. Estos datos se recopilan a través de mediciones de alta precisión y extracción e integración eficientes de la información [24].

El objetivo final de la radiómica es construir sistemas confiables para ayudar a los médicos en lugar de reemplazarlos. Estos sistemas deben apoyar la toma de decisiones clínicas más rápidas [73]. Es un campo de investigación que está recibiendo un gran interés en los últimos años, ya que las imágenes radiológicas son mucho más que figuras anatómicas triviales [74].

### 1.5.1 Procedimiento de extracción de características radiómicas

El procedimiento radiómico se puede dividir en varios procesos con entradas y salidas definidas: (a) adquisición y reconstrucción de imágenes, (b) segmentación y representación de imágenes, (c) extracción y calificación de características, (d) bases de datos e intercambio de datos y (e) análisis informáticos [49], [75]. La Figura 1.5 describe el procedimiento.



**Figura 1.5** Flujo de trabajo de la radiómica. Extraído de [27], [49].

Veamos cada paso por separado dentro del procedimiento:

(a) En las imágenes clínicas de TC, existe gran variabilidad en los parámetros de imagen, como la resolución espacial (tamaño de píxel o tamaño de matriz de adquisición y espesor de

corte), la posición del paciente y las variaciones introducidas por el algoritmo de reconstrucción sobre los niveles de ruido y contraste imagen, entre otros [76]. Una lesión escaneada y luego procesada con dos algoritmos o métodos de reconstrucción diferentes, puede mostrar texturas significativamente diferentes, por lo que en la actualidad se hacen esfuerzos para que coincidan al menos los parámetros originales y protocolos de reconstrucción entre escáneres de diferentes compañías [77].

En la práctica de la radiómica, para contribuir a resolver el problema anterior, solo se capturan los volúmenes de interés (VOI). La filosofía básica detrás del diseño del procedimiento es recopilar la mayor cantidad de datos posible en la interfaz de usuario y utilizar la minería de datos para identificar características con el valor predictivo más alto [78].

Es común en los trabajos científicos emplear imágenes de varios repositorios disponibles en línea, los cuales han proporcionado a los investigadores parte de la gran cantidad de datos que requieren las técnicas de ML para que los modelos se entrenen y funcionen con precisión. Debido a la disponibilidad de bases de datos bien seleccionadas y validadas, la Radiología torácica permanece a la vanguardia del desarrollo de estas nuevas técnicas [79]. Para el análisis de TC de tórax, uno de los repositorios más utilizados es el del Consorcio de bases de datos de imágenes pulmonares y la Iniciativa de recursos de bases de datos de imágenes (LIDC-IRDI) [80].

(b) La segmentación de imágenes en volúmenes de interés, como tumor, tejido normal y otras estructuras anatómicas, es un paso esencial en los análisis computacionales que se requieren en radiómica [75], y según [78] es el más crítico y desafiante. Un método de segmentación ideal debe tener cuatro características clave: automatización, precisión, reproducibilidad y consistencia. Sin embargo, no existe un método de segmentación universal, adecuado para todo tipo de imágenes médicas. Incluso si el mismo algoritmo se ejecuta varias veces con diferentes inicializaciones, los resultados pueden ser variables [75]. Por tanto, el método de segmentación para radiómica debe ser lo más automático posible y con la menor cantidad de interacción del operador. Debe ser eficiente en el tiempo y proporcionar límites precisos y reproducibles [75].

Para obtener suficiente información con alta precisión, el grosor de corte de una tomografía computarizada de tórax no debe ser superior a 1,5 mm, y el número de cortes por paciente

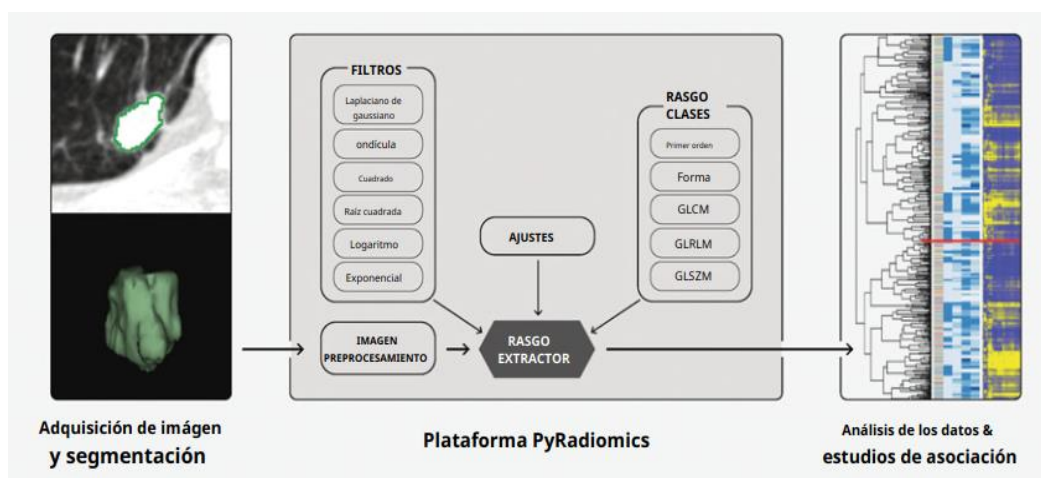
puede ser superior a 300. Por lo tanto, es muy importante tener un algoritmo de segmentación automática y reproducible que pueda lidiar con la variabilidad anterior [24].

(c) A partir de la segmentación, el foco de la radiómica es la extracción de datos de características de alta dimensión, para describir cuantitativamente los atributos de los VOI [78]. Una vez que se definen las regiones tumorales, se pueden distinguir las características de la imagen. Estas características, incluyen la forma, la intensidad, la estructura y la ondulación de la lesión, junto a su ubicación. Para la extracción de características, la información debe ser informativa pero no redundante. Este proceso es necesario para la reproducibilidad [24].

La elección del método de selección de características también depende del algoritmo de aprendizaje automático elegido para la clasificación. Una vez que se seleccionan las características radiómicas a utilizar como entrada al modelo, muchos algoritmos de ML están disponibles para entrenar modelos de clasificación o regresión. Estos algoritmos pueden clasificarse en términos generales como soluciones supervisadas (por ejemplo, regresión logística, máquina de vectores de soporte, bosque aleatorio o red neuronal artificial tradicional) y sin supervisión (por ejemplo, CNN, agrupamiento y autocodificador) [81].

(d) Para resolver el problema de la falta de definiciones de algoritmos estandarizados y procesamiento de imágenes, que dificulta seriamente la reproducibilidad y comparabilidad de los resultados, se implementó en Python una plataforma integral de código abierto denominada “PyRadiomics”. La misma permite el procesamiento y la extracción de características radiómicas a partir de imágenes médicas, utilizando un panel de algoritmos de características rígidamente codificadas. PyRadiomics proporciona una plataforma de análisis flexible, con una interfaz frontal simple y conveniente en 3D Slicer (una plataforma gratuita de código abierto para imágenes médicas) y una interfaz de fondo que permite la automatización del procesamiento de datos, la definición de características y el manejo de lotes de imágenes [82].

PyRadiomics puede extraer características radiómicas de imágenes de diferentes modalidades, utilizando cuatro pasos principales: (i) carga y preprocesamiento de imagen y mapas de segmentación; ii) aplicación de filtros habilitados; iii) cálculo de características utilizando las diferentes clases de características; y (iv) devolución de resultados [82]. La Figura 1.6 muestra este paso.



**Figura 1.6** Utilización de PyRadiomics, extraído de [82].

Las características calculadas se almacenan y devuelven en un diccionario ordenado. Cada característica se identifica con un nombre exclusivo, que consta del filtro utilizado, la clase de característica y el nombre de la característica. Además de las características calculadas, este diccionario también contiene información adicional sobre la extracción, incluida la versión actual, los filtros aplicados, la configuración y el espaciado de la imagen original [82].

(e) Lo anterior genera un volumen muy alto de datos. Cuando se dispone de tales conjuntos, se utiliza la minería de datos, que es el proceso de buscar patrones en grandes conjuntos de datos. Este proceso puede utilizar IA o enfoques estadísticos. Por otro lado, incluyen enfoques de aprendizaje automático tanto supervisados como no supervisados; pero como se ha mencionado, el poder del modelo de clasificación a emplear, depende completamente del tamaño y la calidad de los datos dentro de la base de datos a utilizar. La calidad no solo depende de las condiciones de adquisición de la imagen, sino también de la disponibilidad y confiabilidad de covariables [78].

El análisis de datos consta de dos pasos distintos: el primero es el modelado, en el que se crea un modelo de clasificación y/o regresión. En el segundo paso se utiliza el modelo para predecir resultados. Construir el modelo implica determinar el tipo de clasificador o regresor a utilizar y alimentar el modelo con un conjunto de casos preclasificados, es decir, matrices de pares característica-etiqueta, donde la etiqueta representa el estado clínico del sujeto

correspondiente. A este proceso de presentar el modelo con casos preclasificados se le denomina entrenamiento [83].

En cuanto al clasificador a utilizar, hay varias opciones disponibles que van desde soluciones conceptualmente sencillas, como el análisis discriminante lineal, K-vecinos más cercanos (KNN) y *Naive Bayes*, hasta otras más complejas, como las descritas CNN, RF y SVM. Para la regresión, Cox y regresión logística son las opciones más comunes [83].

### 1.5.2 Características radiómicas y su importancia

Las características radiómicas se pueden subdividir en estadísticas, incluidas las basadas en histogramas y basadas en texturas, basadas en modelos, basadas en transformación y basadas en la forma [84],[85].

Las características estadísticas de primer orden, describen la distribución de los valores individuales de los vóxeles, sin tener en cuenta las relaciones espaciales. Se trata de propiedades basadas en histogramas que informan los valores, medio, máximo y mínimo de las intensidades de los vóxeles de la imagen, así como su asimetría y curtosis. La asimetría refleja los desplazamientos de la curva de distribución de datos hacia la izquierda (sesgo negativo, por debajo de la media) o hacia la derecha (sesgo positivo, por encima de la media). La curtosis refleja la cola de una distribución de datos, en relación con una distribución gaussiana, debido a valores atípicos. Otras características incluyen la aleatoriedad (entropía) y la uniformidad del histograma [84], [86].

Por el contrario, las características estadísticas de segundo orden, incluyen las denominadas características texturales, que se obtienen calculando interrelaciones estadísticas entre vóxeles vecinos. Proporcionan una medida de la distribución espacial de la intensidad de los vóxeles, y por tanto, de la heterogeneidad intralesional. Estas características pueden derivarse de la matriz de concurrencia de niveles de gris (GLCM). La misma cuantifica la incidencia de vóxeles con las mismas intensidades, a una distancia predeterminada a lo largo de una dirección fija. También se puede obtener la matriz de longitud de recorrido de nivel de grises (GLRLM). Esta cuantifica los vóxeles consecutivos con la misma intensidad a lo largo de direcciones fijas [86].

Además, se puede analizar la matriz de zona de tamaño del nivel de gris (GLSZM). Esta se basa en áreas de vóxeles vecinos con el mismo nivel de gris. Otra posibilidad es el cálculo

de la matriz de zona de distancia del nivel de gris (GLDZM). Esta evalúa las zonas de vóxeles vecinos interconectados con el mismo nivel de gris. En este caso requiere que estén a la misma distancia del borde de la región de interés. También se puede calcular la matriz de diferencia de tonos de grises de vecindario (NGTDM), que cuantifica la suma de las diferencias entre el nivel de gris de un píxel o un vóxel y el nivel de gris medio de sus vóxeles vecinos, dentro de una distancia predefinida [84]. La Figura 1.7 muestra un ejemplo de extracción de matrices asociadas a algunas de estas características radiómicas.

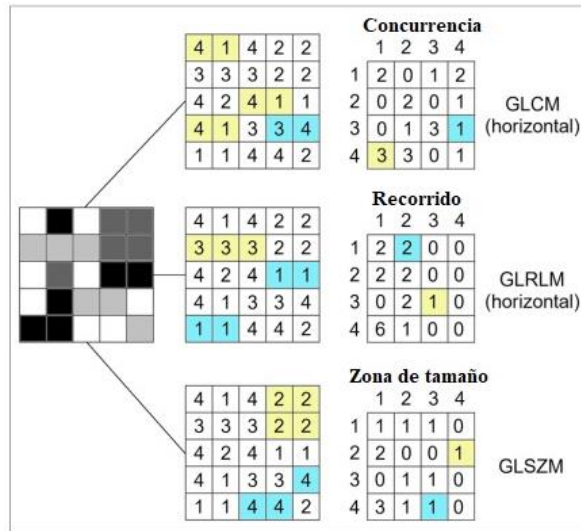


Figura 1.7 Cálculo de características de textura radiómica, tomada de [83].

Los análisis basados en modelos tienen como objetivo interpretar la información espacial del nivel de gris para caracterizar objetos o formas. Se calcula un modelo parametrizado de generación de texturas y se ajusta a la ROI. Sus parámetros estimados se utilizan como características radiómicas [84].

Los métodos basados en transformadas, incluidas las transformadas Wavelet, de Fourier, Gabor y Haar, analizan los patrones de nivel de gris en un espacio diferente [83]. La transformada wavelet describe el concepto de descomponer los datos de imagen en diferentes componentes de frecuencia llamadas ondículas (wavelets) y utilizar estos datos para extraer características relacionadas con la textura y la intensidad de la imagen [25]. Por tanto, la transformación wavelet se puede utilizar no solo para la generación de características

radiómicas, sino también para la segmentación de imágenes o como un paso de preprocesamiento para el análisis de texturas [84].

Las características basadas en formas describen las propiedades geométricas de las ROI. Muchas características basadas en forma son conceptualmente mucho más simples que otras características radiómicas, como los diámetros 2D y 3D, los ejes y sus proporciones. Las características incluyen compacidad y esfericidad, que describen cómo la forma de una ROI difiere de la de un círculo (para análisis 2D) o una esfera (para análisis 3D), y la densidad, que se basa en la construcción de un cuadro delimitador orientado mínimo (o rectángulo para análisis 2D) que encierra la ROI [84].

### **1.5.3 Aplicaciones de la radiómica en el estudio de los nódulos de pulmón**

La aplicación de la radiómica requiere una combinación de especialidades, como Radiología, la Bioinformática y la Ingeniería biomédica. Aplica características de imagen para que actúen como entidades básicas y utiliza varios algoritmos y modelos como herramientas para analizar y transformar estas entidades en información de diagnóstico útil, legible y confiable. Su aparición abre una nueva forma de entender la etiología, patología y progresión de la enfermedad. La información precisa y detallada extraída de las imágenes médicas digitales se correlaciona con las enfermedades y puede utilizarse ampliamente en el diagnóstico, el seguimiento y la planificación para el tratamiento de enfermedades malignas [87].

En el campo de los nódulos pulmonares, el uso de la radiómica se está expandiendo. Se han realizado trabajos a partir de imágenes de TC que permitan la clasificación de los nódulos pulmonares, para predecir la malignidad de los mismos. Entre los estudios precedentes, aplicando características de este tipo, se pueden mencionar: en el año 2016 en [88], a partir de imágenes del *National Lung Screening Trial*, se extrajeron 219 características radiómicas y se empleó un clasificador RF, obteniendo una precisión del 80% y sensibilidad y especificidad de 51.7% y 92.9% respectivamente, con un área bajo la curva (AUC) de 0.83. En otro estudio [89], se utilizaron datos de la BD LIDC-IRDI, con anotaciones de los radiólogos, para extraer 583 características radiómicas, empleando también RF como algoritmo clasificador, obteniéndose los resultados que se reflejan en la Tabla 1.1.

Otras investigaciones, también reflejadas en la Tabla 1.1 como [90] y [91], ambas con datos de LIDC-IRDI, la primera con una muestra de 593 imágenes de pacientes, donde se

extrajeron 150 características radiómicas y se utilizó SVM para clasificar, obtuvo buenos resultados de clasificación. En el segundo caso se analizaron imágenes de 97 pacientes con el fin de obtener un secuenciador radiómico, a partir de una CNN, igualmente con resultados prometedores.

Otros autores [92] y [93], emplearon SVM como algoritmo clasificador de ML, el primero extrajo de 72 imágenes con nódulos pulmonares, 103 características radiómicas, para alimentar el modelo de clasificación y lo combinó con el operador de selección y contracción mínimo absoluto (LASSO); mientras que el segundo analizó 75 nódulos de 72 pacientes y extrajo 750 características radiómicas de cada nódulo. Los resultados de las diversas métricas de ambos estuvieron por encima del 84 % de clasificación exitosa.

Por otra parte, en [94] se recogieron datos de 294 casos, de los cuales extrajeron 385 características radiómicas, las cuales se redujeron por el método LASSO, con resultados por encima del 80 % de buena clasificación.

En [69] se hace una comparación de los métodos de *Naives Bayes* y SVM, entrenados con diversas características, obteniendo muy buenos resultados.

Como se aprecia, la radiómica está demostrando ser eficaz no solo en la detección y estadificación del cáncer de pulmón, sino también en la predicción de la respuesta del paciente a la terapia aplicada [83]. Nuevos estudios multicéntricos son sin embargo necesarios, con el fin de buscar modelos cada vez más eficaces y con poder de generalización, que puedan ser utilizados en rutina clínica como apoyo al diagnóstico médico.

**Tabla 1.1** Estudios recientes empleando radiómica

Estudio	Descripción	Resultados
Hawkins 2016 [88]	Datos del <i>National Lung Screening Trial</i> (NLST) 219 características radiómicas Clasificador RF	AUC: 0,83 Precisión: 80% Sensibilidad: 51,7% Especificidad: 92,9%
Ma 2016 [16]	Datos del LIDC con anotaciones de radiólogos 583 características radiómicas Clasificador RF	Precisión: 82,7% Sensibilidad: 80% Especificidad: 85,5%
Wang 2016 [90]	593 pacientes LIDC-IDRI 150 características radiómicas SVM	Precisión: 76,1% Sensibilidad: 74,6% Especificidad: 78,9%
Kumar 2017 [91]	97 pacientes LIDC-IDRI Descubrimiento de un secuenciador radiómico a partir de CNN	Precisión: 77,52% Sensibilidad: 79,06% Especificidad: 76,11%
Choi 2018 [92]	72 nódulos de LIDC-IDRI 103 características radiómicas SVM-LASSO	AUC: 0,89 Precisión: 84,6%
Chen 2018 [93]	72 pacientes con 75 nódulos 750 características radiómicas SVM	Precisión: 84% Sensibilidad: 92,85% Especificidad: 72,73%
Mao 2019 [94]	294 casos 385 características radiómicas LASSO para reducir características y Regresión para clasificación	AUC: 0,97 Precisión: 89,8% Sensibilidad: 81% Especificidad: 92,2%
Shakir 2020 [69]	Imágenes de LIDC y LUNGx 105 características radiómicas usando PyRadiomics. SVM y Naives Bayes entrenado con 2, 4, 8, 16 y 20 características	Para el mejor modelo (SVM con 8 características) Sensibilidad: 81% Especificidad: 92,2%

## 1.6 Conclusiones del capítulo

La constante evolución de las imágenes de TC y las posibilidades que esta brinda han contribuido a su utilidad para detectar nódulos pulmonares. Con el avance del ML han aparecido nuevos y prometedores métodos de extracción de características para la correcta

---

clasificación de lesiones, que a simple vista no son obvias de clasificar. En particular, la radiómica ha emergido como un método robusto para detectar o predecir enfermedades. En el contexto del cáncer de pulmón, la radiómica se está convirtiendo en un método efectivo, sin embargo, al generar gran cantidad de datos, que complementan la información de las imágenes de TC, se hace necesario el empleo de la inteligencia artificial para su análisis. Insertar la radiómica dentro de un sistema CAD, es una valiosa herramienta de apoyo a los especialistas médicos, avocados a la Medicina personalizada del futuro inmediato.

## CAPÍTULO 2. MATERIALES Y MÉTODOS

En este capítulo se explica la implementación de un sistema CAD, para su utilización como medio de clasificación de malignidad de nódulos pulmonares, a partir de imágenes de tomografía computarizada de tórax. Se realiza una descripción de cada uno de los pasos a realizar en las etapas del sistema propuesto. También se explica el proceso de entrenamiento, validación y prueba de varios clasificadores y se caracterizan los materiales de software y hardware utilizados para procesar la información. Además, se describe la base de datos empleada y el manejo de esta a partir de sus anotaciones.

### 2.1 Descripción general de la implementación del sistema CAD

Para implementar el sistema CAD para clasificar maligno - benigno, se partió de imágenes anotadas de TC de pulmón, de la base de datos LIDC-IDRI [80]. Las mismas presentan una resolución de 512 x 512 píxeles y se encuentran en el formato estándar en imagenología digital y comunicaciones en Medicina (DICOM, del inglés *Digital Imaging and Communications in Medicine*). Aprovechando las potencialidades de este formato, que incluye múltiples metadatos y alta resolución espacial, y las anotaciones de la base de datos, se realizó un análisis de tomografías de tórax de 100 pacientes. Se escogieron, los casos LIDC-IDRI-0001 hasta el LIDC-IDRI-0101, (con la excepción de LIDC-IDRI-0028, LIDC-IDRI-0032, LIDC-IDRI-0062, LIDC-IDRI-0071, LIDC-IDRI-0100, que no presentan lesiones anotadas por los radiólogos). Se obtuvo así un conjunto de entrenamiento de 1165 imágenes (cortes de TC en 2D) que contienen un total de 275 nódulos pulmonares.

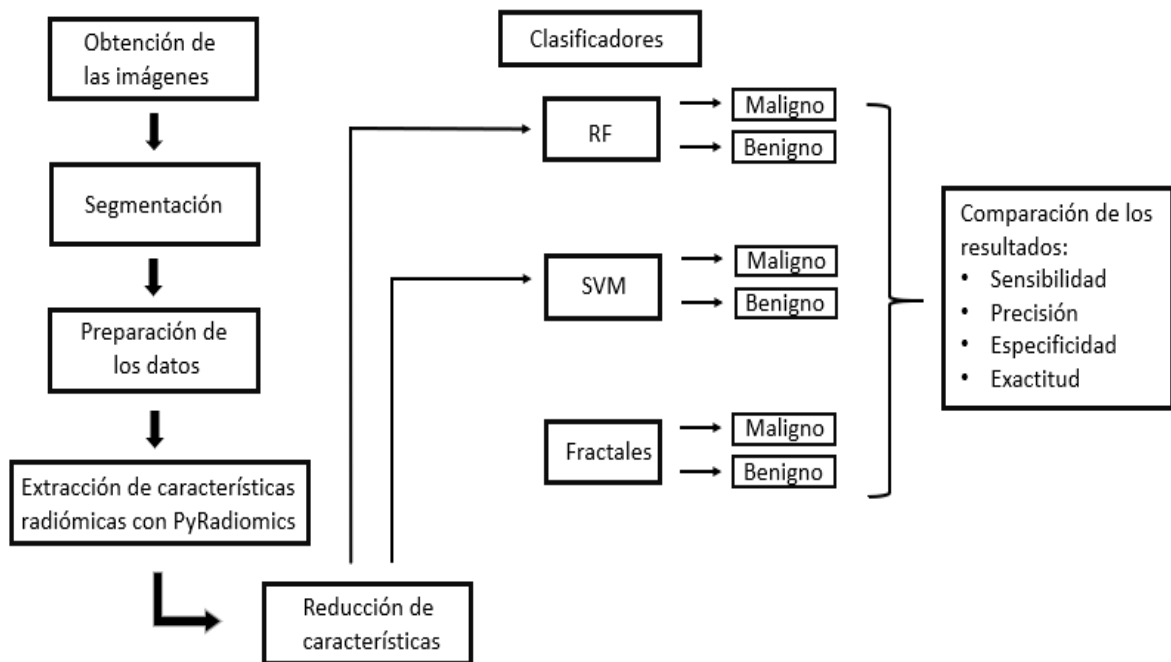
De cada TC con nódulos, se realizó un proceso de segmentación en 2D, para separar la lesión del resto de la imagen, partiendo de todos los cortes 2D donde cada nódulo era visible. Una vez segmentada la lesión, se realizó la reconstrucción en 3D de la misma, a partir del conjunto de cortes 2D en que cada una era visible y se introdujo este resultado en la plataforma PyRadiomics [82].

Haciendo uso de PyRadiomics, se obtuvieron varias características radiómicas [69]. Este número fue reducido a partir de calcular la correlación de Pearson de cada una, con el grado de malignidad reportado en la anotación de la base de datos. El objetivo de este paso fue

determinar aquellas características que eran significativamente importantes para la clasificación. Se desecharon valores de correlación con  $R < 0.6$  y correlaciones estadísticas no significativas ( $p \geq 0.05$ ).

Las características finalmente seleccionadas conformaron un vector, que constituyó la entrada a modelos de RF y SVM [60]. Estos fueron entrenados con 220 nódulos seleccionados aleatoriamente por el software sobre Python, de todos los nódulos que se estudiaron, para obtener una clasificación de las lesiones en maligno o benigno.

Durante la etapa de validación se escogieron también al azar los restantes 55 nódulos del conjunto por el software sobre Python, de la misma BD, para evaluar internamente el desempeño de ambos clasificadores y comparar resultados. Para evaluar el desempeño se calcularon un conjunto de métricas. Los resultados obtenidos se compararon contra los obtenidos por el método de fractales [95], para la misma base de datos, así como respecto a otros estudios de la literatura científica. En la Figura 2.1 se presenta un esquema con todos los pasos hasta aquí descritos.



**Figura 2.1** Descripción general del proceso

## 2.2 Funciones y bibliotecas más utilizadas

Las bibliotecas de Python más utilizadas fueron:

- *Pydicom* para la manipulación de archivos DICOM
- *SimpleITK* para el procesamiento y análisis de imágenes médicas
- *OpenCV* para el procesamiento de imágenes y la visión computacional
- *Matplotlib.pyplot* para la visualización de datos y gráficos
- *Numpy* para la manipulación de matrices y arreglos, es decir, la computación científica del lenguaje
- *Pandas* para el análisis y manipulación de datos tabulares
- *Scikit-Learn* para los algoritmos de aprendizaje automático
- *Radiomics* para llamar al extractor de características

De las funciones más utilizadas encontramos:

- *.dcmread()* y *.dcmwrite()* funciones de *pydicom* para leer y guardar archivos DICOM
- *.RreadImage()*, *.GetSpacing()*, *.GetSize()*, *.GetOrigin()* y *.GetDirection()* de *SimpleITK* para leer una imagen y obtener su espaciado, su tamaño, su orientación y dirección.
- *.imshow()*, *.waitKey()* y *destroyAllWindows()* de *opencv* para mostrar la imagen, esperar a que se presione una tecla y cerrar la ventana respectivamente
- *.subplots()* para mostrar gráficos donde se observen las coordenadas
- *.sum()*, *.mean()*, *.round()* funciones de *numpy* para hacer cálculos básicos como suma, media y redondeo
- *.read\_csv()* y *.to\_csv()* funciones de *pandas* para leer y guardar archivos CSV
- *.fit()* y *.predict()* funciones de *Scikit-Learn* para entrenar y validar los modelos de clasificación
- *.featureextractor.RadiomicsFeatureExtractor()* de *radiomics* para trabajar con el extractor de características

## 2.3 Segmentación de los nódulos

Los nódulos objeto de estudio, son visibles en un número de cortes variables entre tomografías de diversos pacientes. Esto depende del tamaño de la lesión, su composición,

posición y las características del escáner de TC con el que fueron generadas [96]. Las 275 lesiones analizadas fueron visibles en 1165 cortes, variando desde 2 hasta 21 por cada una. Primeramente, fue necesario obtener corte a corte una imagen binaria. Para esto se probaron diferentes métodos de segmentación basada en umbral. Se seleccionaron como métodos de binarización, el método de umbral de máxima verosimilitud [93], basado en modelos de mezcla de poblaciones, y el método de umbral de error mínimo [94], por ser los que presentaron un mejor resultado visual y por lo tanto se adaptaron mejor a las imágenes que se estudiaron.

El primer método fue propuesto por T. Kurita en 1992 [97]. Se utiliza para clasificar observaciones o datos en diferentes categorías o grupos, basándose en la probabilidad máxima (o máxima verosimilitud) de que pertenezcan a una población o clase específica. Se basa en la idea de que los datos observados se generan a partir de una mezcla de varias distribuciones de probabilidad. El objetivo fue estimar los parámetros de estas distribuciones para maximizar la verosimilitud de los datos observados. En este método se asumió que cada clase (objeto-fondo) de la que provienen los datos, está caracterizada por una distribución de probabilidad diferente y cada una de las clases tiene sus propios parámetros, (media y varianza).

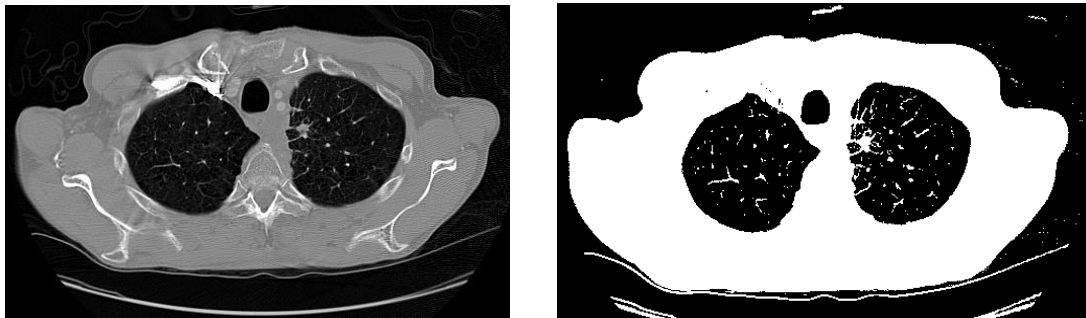
Una vez que los parámetros están estimados fue posible establecer el umbral de decisión para asignar nuevas observaciones a una clase específica. Esto se hizo considerando, para cada píxel, la probabilidad de pertenencia a cada clase y eligiendo la clase de máxima probabilidad. De esta forma, se asignó a cada píxel una clase según el umbral determinado, lo que divide la imagen completa en dos regiones. En resumen, se binarizó cada imagen 2D. Sin embargo, este método no fue satisfactorio para todas las imágenes. Por esta razón, en aquellas donde este método no mostró un resultado adecuado, entiéndase que presentó fallos, se aplicó un segundo método: umbral de error mínimo [94], con el mismo objetivo que el explicado anteriormente.

El método de umbral de error mínimo se desarrolló para determinar un umbral óptimo en la segmentación de imágenes. Se desarrolló por J. Kittler y J. Illingworth en 1986 [98] para abordar el problema de seleccionar un umbral de manera adaptativa y automática, teniendo en cuenta la distribución estadística de las intensidades de píxeles en la imagen.

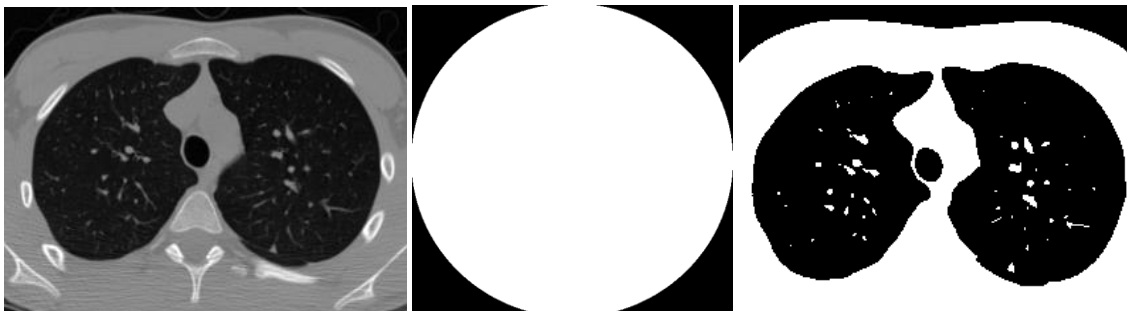
El método de Kittler-Illingworth [98] utiliza un criterio basado en la minimización del error de clasificación. En otras palabras, busca el umbral que minimiza la probabilidad de error de clasificación entre las clases objeto y fondo. La función del error tiene en cuenta la varianza intraclase y la diferencia de medias entre las dos clases.

Este método, según la literatura científica, es particularmente efectivo en situaciones donde las clases en la imagen están bien separadas y la distribución de intensidades puede modelarse de manera razonable como bimodal. Sin embargo, puede no ser tan robusto en casos de imágenes con distribuciones más complejas [98]. Por esta razón, para los objetivos del sistema CAD que se propone en este trabajo, se utilizó solo como segunda opción, ante un fallo del método de umbral de máxima verosimilitud.

La Figura 2.2 muestra un ejemplo de un corte de TC antes y después de aplicarle el método de máxima verosimilitud, siendo este exitoso. En la Figura 2.3, se muestra un ejemplo donde el método de máxima verosimilitud resultó fallido, pero el método de método de umbral de error mínimo resultó exitoso.

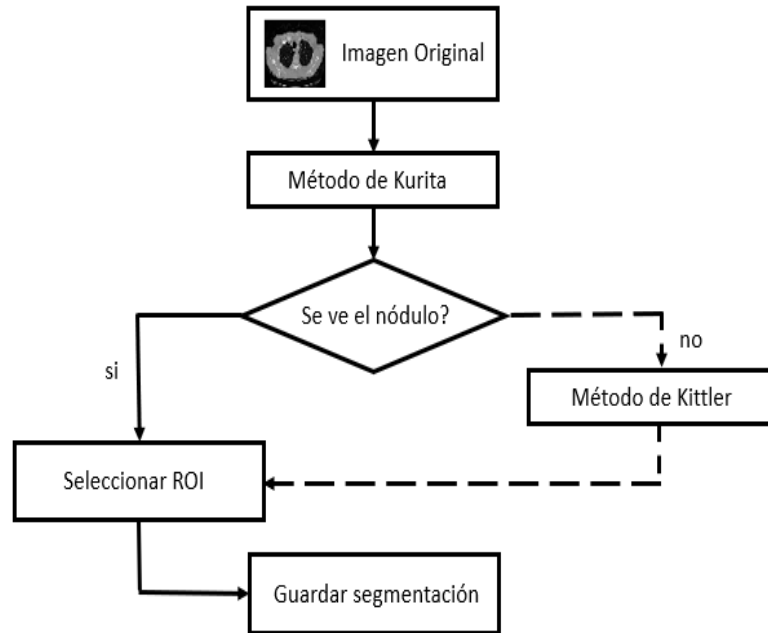


**Figura 2.2** Corte de un paciente antes y después de aplicar el método de máxima verosimilitud



**Figura 2.3** Fallo del método de máxima verosimilitud y éxito del método umbral de error mínimo.

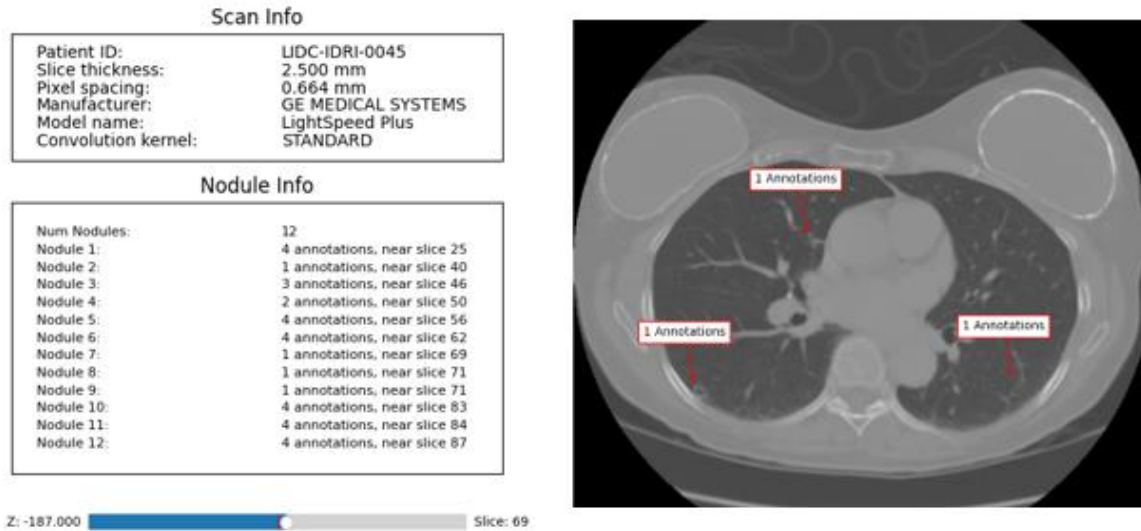
Luego de obtener la binarización de cada corte, se seleccionó manualmente la región de interés donde estaba el nódulo. Se comparó con las segmentaciones realizadas por 4 radiólogos expertos, que aparecen en la BD. Estas se visualizaron a partir de la herramienta *Pylidc*. Se guardó cada una de estas segmentaciones en formato DICOM. Este proceso se muestra en la Figura 2.4.



**Figura 2.4** Proceso de segmentación

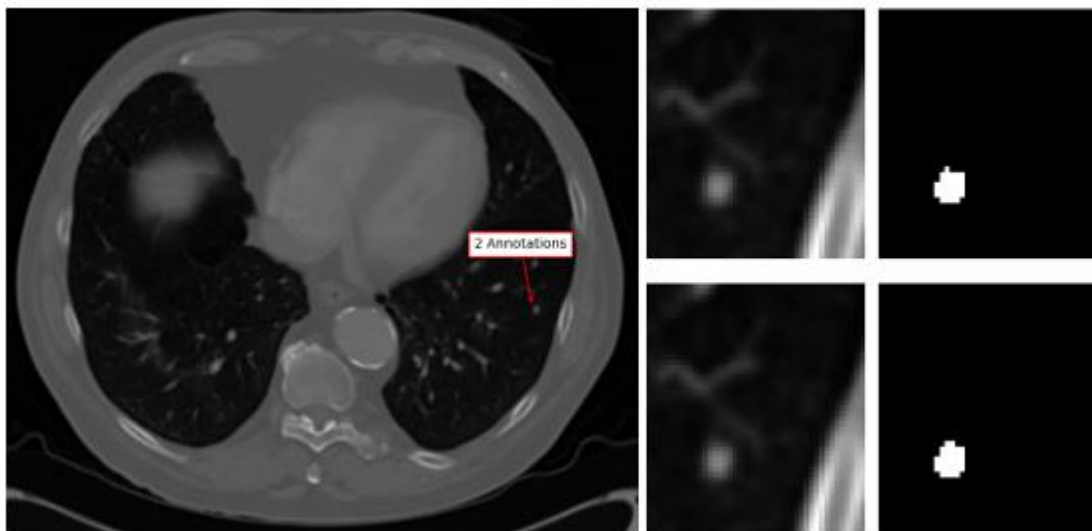
### 2.3.1 Evaluación de la calidad de la segmentación utilizando *Pylidc*

*Pylidc* posee una interfaz gráfica amigable que permite la navegación por los cortes de la TC y muestra las localizaciones de los nódulos, sus anotaciones, así como el corte de mejor visibilidad del nódulo, a criterio de los expertos [95]. La Figura 2.5 muestra el ambiente de esta herramienta. En el **Anexo 1** se puede encontrar el código para acceder a la interfaz presentada en dicha figura.



**Figura 2.5** Corte de TC con sus anotaciones cargado en *Pylidc*

A partir de *Pylidc* es posible acceder a la segmentación de un nódulo seleccionado de un paciente determinado, en el corte de mejor visualización. Se pueden observar tantas segmentaciones como anotaciones tenga el nódulo. Por esta razón, en este experimento sirvió como referencia para evaluar la calidad de la segmentación descrita en el paso anterior. La Figura 2.6 muestra visualmente las posibilidades que brinda esta herramienta. Se aprecia un ejemplo de la segmentación anotada de un nódulo en la BD.



**Figura 2.6** Ejemplo de segmentación de un nódulo anotado en la base de datos, obtenido con *Pylidc*

En el caso analizado, el nódulo tiene dos anotaciones en la base de datos, esto quiere decir que hubo dos radiólogos que vieron el nódulo. Por lo tanto, se podrá observar la segmentación que hizo cada uno de estos radiólogos.

## 2.4 Reconstrucción, etiquetado y preparación de los datos

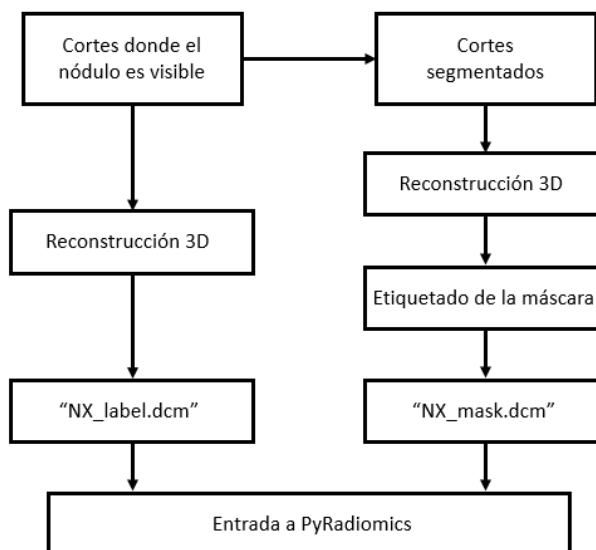
Entre los requisitos de PyRadiomics se encuentra que la entrada al extractor de características debe contener la imagen 3D y la máscara de segmentación, también en 3D, donde se aprecie el nódulo, para que sea posible extraer las características basadas en vóxeles. Por lo tanto, se hace necesario que, una vez segmentados los nódulos en 2D en todos los cortes donde es visible de cada tomografía, se reconstruya cada uno en 3D, y se guarde ese archivo en un formato compatible con la librería *SimpleITK* de Python, que es la librería con la que trabaja PyRadiomics para el manejo de las imágenes.

Se decidió mantener el formato DICOM para las reconstrucciones, porque es uno de los aceptados en la librería y el más completo, en el sentido de mantener, además de la imagen, sus metadatos [99]. La concatenación consistió en crear un nuevo archivo multiframe, que contuviera múltiples imágenes en un mismo conjunto de datos, donde los elementos de los píxeles de cada imagen DICOM original se extraen y se concatenan en una lista.

Una vez que se obtuvo la reconstrucción de las segmentaciones 3D a partir de los cortes 2D, se etiquetaron los píxeles en cada imagen. De esta forma se le asignó el valor de 1, a los píxeles que formaban parte de la región de interés, y que se observan en color blanco, y valor 0, para los correspondientes al color negro, que son los que pertenecen al fondo.

Los cortes sin segmentar de las imágenes de TC, que son aquellos donde se apreciaba el nódulo, se seleccionaron también, y se reconstruyeron de la misma forma que los segmentados. Sin embargo, a diferencia de los anteriores, estos quedaron sin etiquetar, como parte de la reconstrucción de la imagen completa.

Las reconstrucciones de las imágenes se nombraron de la forma “NX\_label.dcm”, mientras que las reconstrucciones de las máscaras se denominaron como “NXmask.dcm”. Ambas se guardaron en una misma carpeta, donde X es el número del nódulo que se analizaba en esas reconstrucciones. Estos pasos se muestran en el diagrama de la Figura 2.7.



**Figura 2.7** Explicación de los pasos de reconstrucción

### 2.4.1 Visualización de las anotaciones de malignidad mediante *Pylicd*

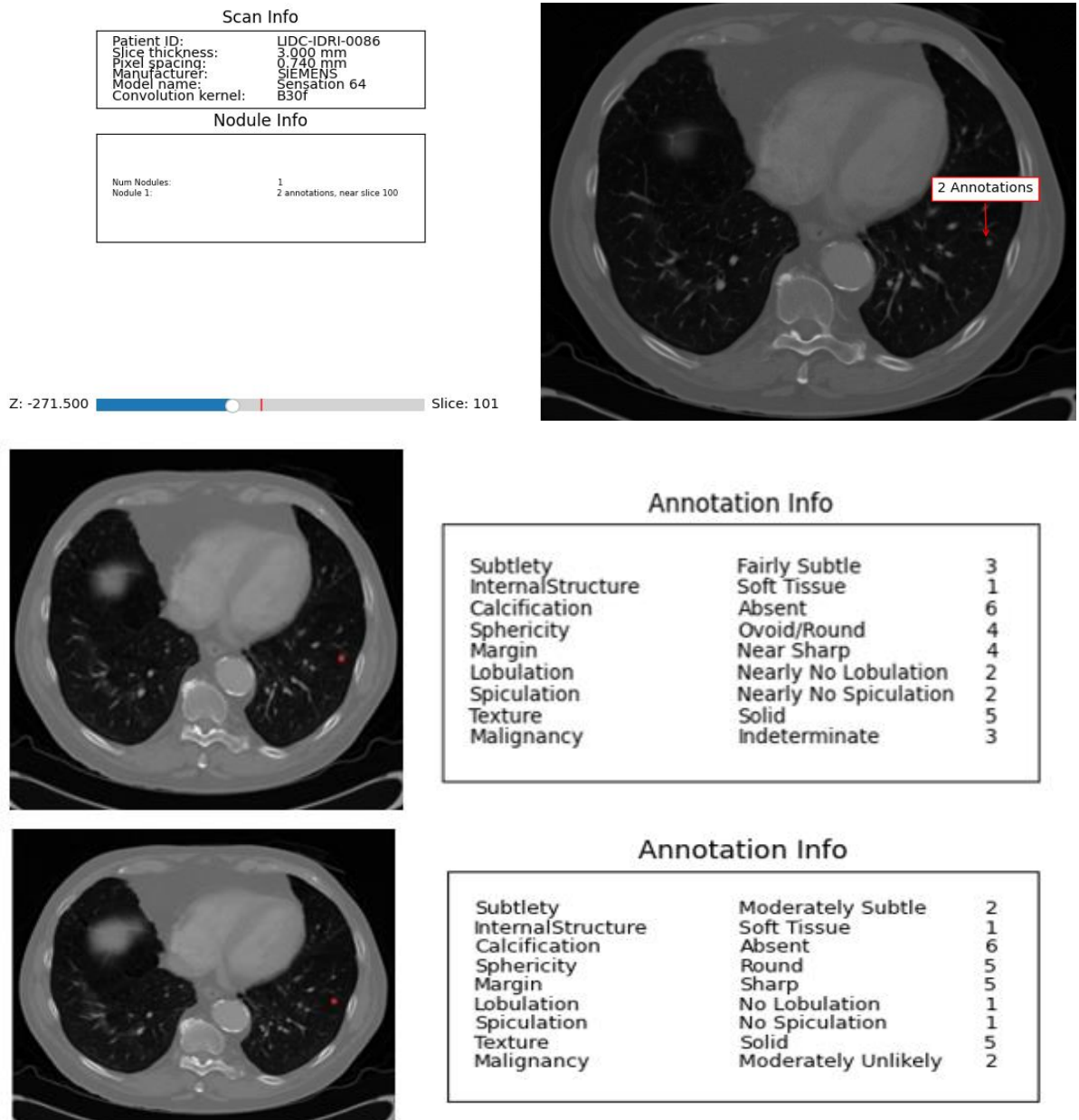
La herramienta *Pylicd* permite también la visualización de las anotaciones de malignidad de los nódulos realizadas por los radiólogos. El código para realizar esto se encuentra en el **Anexo 2**. Estas anotaciones se observan una a la vez, y están dadas en 5 grados de malignidad, donde:

1. Altamente improbable que sea maligno
2. Moderadamente improbable que sea maligno
3. Indeterminado
4. Moderadamente sospechoso de ser maligno
5. Altamente sospechoso de ser maligno

Sin embargo, a los efectos de este trabajo, cuando se utiliza un modelo de clasificación binaria, las clases solo son consideradas como maligno o benigno. En este caso, las categorías anteriores se agruparon de la siguiente forma:

- 0- Probablemente benigno: Agrupa los grados 1, 2 y 3.
- 1- Probablemente maligno: Agrupa la superposición con los grados 4 y 5.

Como hay nódulos que han sido anotados por más de un radiólogo, tienen más de una anotación y no siempre son coincidentes, ya que se utilizó la escala del 1 al 5. Es decir, cada nódulo en la BD tiene tantas anotaciones como radiólogos lo han visto. En este sentido, para los nódulos que presentaban más de una anotación diferente, se calculó el promedio de todas las anotaciones que aparecían en la BD y se tomó ese promedio como el grado de malignidad promedio del mismo. En la figura 2.8 se muestra un ejemplo.



**Figura 2.8** Ejemplo de visualización de las anotaciones de un nódulo

Este paciente tiene un nódulo ubicado en el pulmón izquierdo, que es anotado por dos radiólogos. El primero lo consideró como grado de malignidad 3, mientras que el segundo lo consideró como grado de malignidad 2. Como el valor promedio fue de 2.5, se redondeó y anotó la malignidad como 3. Si este resultado lo expresamos en notación binaria, sería grado 0, es decir benigno.

## 2.5 Extracción de características con PyRadiomics

En este paso se partió de realizar una nueva organización de los datos. Para esto se creó un archivo CSV, que almacenaba en cada fila los nombres de las imágenes y las máscaras de segmentación obtenidas en el paso anterior. Las imágenes originales y las máscaras de segmentación se encontraban en formato DICOM, ocupando las columnas “Subject\_Label” y “Subject\_Mask” del CSV respectivamente.

Cada CSV contó además con una columna denominada “Label”, que contiene las anotaciones de malignidad que se obtuvieron anteriormente. Este archivo se lee por el código en Python para confeccionar el *DataFrame* a partir de sus datos, donde se guardan las características extraídas.

Para la configuración de PyRadiomics se empleó un archivo de configuración YAML, donde se definieron las opciones que controlaban el proceso de extracción de características. Estas fueron:

- Tipo de imagen a utilizar (*imageType*),
- Ancho del *bin* utilizado en la discretización de las intensidades de píxeles (*binWidth*)
- Resolución espacial a la que se remuestrearon las imágenes (*resampledPixelSpacing*)

Además, se configuraron los tipos de características que se deseaban extraer. Los valores utilizados para configurar estos parámetros aparecen en el **Anexo 3**.

Para extraer las características en la plataforma PyRadiomics se utilizó parte del código del repositorio de GitHub [https://github.com/Astarakee/Radiomics\\_pipeline](https://github.com/Astarakee/Radiomics_pipeline) citado en [100]. A este código fue necesario especificarle las rutas a las imágenes y máscaras antes de realizar la extracción, así como la ruta donde se encontraba el CSV con los datos de trabajo. Una vez definido lo anterior, se establecieron las relaciones necesarias entre las columnas del CSV, así como las relaciones entre cada caso y el directorio donde están ubicadas las imágenes.

Estas tareas fueron divididas, definiendo funciones que utilizan librerías de Python como: os, re, csv, six, glob, SimpleITK. Previamente a la inicialización del extractor, la máscara etiquetada se remuestreó, para que coincidiera geométricamente en espacio, dirección y tamaño con la reconstrucción de los cortes de la TC, donde se veían los nódulos antes de segmentar. Posteriormente se cargó el extractor de características utilizando la configuración mostrada en el **Anexo 4** y se iteró a través de la lista de imágenes y máscaras.

Los resultados obtenidos se almacenaron en un diccionario, y utilizando la biblioteca ‘pandas’, se creó un *DataFrame*, que fue guardado en un nuevo archivo de formato CSV, para su posterior utilización, como se muestra en el **Anexo 5**.

## 2.6 Reducción de características

Una vez que son extraídas las características radiómicas en PyRadiomics, es necesario proceder a la reducción de la dimensionalidad del problema, a fin de facilitar la tarea de los clasificadores. Para esto se empleó el cálculo del coeficiente de correlación de Pearson, entre la anotación de la BD y cada característica radiómica extraída en todo el conjunto de imágenes.

El coeficiente de correlación de Pearson se utiliza para medir si dos conjuntos de datos se encuentran relacionados. En otras palabras, refleja el grado de correlación lineal entre dos variables [101]–[103]. Su ecuación de cálculo es la siguiente:

$$R = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (1)$$

En la ecuación,  $R$  describe el grado de correlación lineal entre dos variables,  $i$  es el número de la muestra,  $x_i$  e  $y_i$  son los valores de las variables  $i$ -ésima y  $j$ -ésima donde  $i \neq j$  [102]. El valor de  $R$  está comprendido entre -1 y +1. Cuanto mayor sea el valor absoluto de  $R$ , más fuerte será la correlación, aunque esto está sujeto a su significación estadística  $p$ , considerándose significativo  $p < 0.05$  para un 95 % de confianza [102].

Para el análisis de la correlación se utilizó la siguiente escala [103]:

Coefficiente de correlación:

- 0,8-1,0 correlación muy fuerte

- 0,6-0,8 correlación fuerte
- 0,4-0,6 correlación moderada
- 0,2-0,4 correlación débil
- 0,0-0,2 Correlación muy débil o ninguna correlación

En la implementación realizada, el archivo CSV que contenía las características obtenidas de PyRadiomics se cargó y se le eliminó la primera columna, de forma tal que el *DataFrame* leído solo contuviera la etiqueta de malignidad y las características de cada nódulo.

Para cada una de las características extraídas se calculó la correlación con respecto a la etiqueta de malignidad anotada, permitiendo conocer cuán fuerte o determinante resultaba cada característica para definir el grado de malignidad de la lesión. Como las anotaciones de la BD se daban en 5 grados de malignidad, se aplicó la correlación de Pearson para el conjunto de datos con estas 5 categorías de clasificación. Además, se realizó el cálculo de la correlación, pero para las anotaciones binarias de malignidad.

En ambos casos se seleccionaron solo aquellas características que presentaban una correlación mayor de 0.60 con la anotación de la BD. De esta forma, se obtuvo un vector de características de entrada para los clasificadores, más reducido que el inicial, lo cual redujo a su vez el costo computacional de todo el sistema.

## 2.7 Implementación de los clasificadores

El objetivo de cualquier clasificador es encontrar la frontera que permite separar clases. En este trabajo se realizó la clasificación a través de dos modelos: SVM y RF.

Ambos clasificadores recibieron el conjunto de características radiómicas reducidas para su entrenamiento, que son el resultado de la etapa de selección anterior. A partir de estas características se construyen los modelos de clasificación con las imágenes de prueba.

Para ambos clasificadores la implementación siguió los siguientes pasos:

1. Importar las clases *'RandomForestClassifier'* y *'SVC'* de *Scikit-Learn*.
2. Crear una instancia del clasificador.
3. Con la función *train\_test\_split* de *Scikit-Learn* se dividió el conjunto de datos en dos subconjuntos (el primero para entrenamiento del modelo predictivo y el segundo para la validación). A esta función se le fijó el parámetro *test\_size* en 0.2, para definir que

el 20% de los datos se escogiesen para validar. Esta división de los datos se realizó aleatoriamente, preservando la distribución de las etiquetas. Para garantizar que entre diferentes ejecuciones del código se obtuviera la misma división de los datos, se fijó el parámetro *random\_state* en un valor específico. De esta manera, se garantizó que el conjunto se dividiera de forma consistente y reproducible.

4. Con la función `fit()` se entrenó cada modelo y se ajustó a los datos de entrenamiento.
5. Con la función `predict()` se realizaron las predicciones en nuevos datos, utilizando el modelo previamente entrenado.

Luego de esto se obtuvieron las matrices de confusión y los resultados de las métricas.

Los clasificadores se entrenaron para clasificar en maligno-benigno, pero el vector de características de entrada se constituyó de dos formas diversas: en un primer caso las características seleccionadas de la correlación de Pearson para anotación binaria de malignidad y en el otro caso, la entrada consistió en las características seleccionadas a partir de la correlación de Pearson para la anotación de malignidad en 5 grados.

### 2.7.1 Implementación de máquinas de soporte vectorial (SVM)

Para implementar el SVM se utilizó la función de Python *SVC* de la biblioteca *Scikit-Learn* que permite obtener ese modelo. Los parámetros de esta función fueron:

- Tipo de kernel (lineal, radial y polinómico de grados 2,3 y 5)
- Parámetro de regularización - C (se utilizó C=1, que viene por defecto, lo que implicó una configuración equilibrada en la implementación del modelo).
- Parámetro de aleatoriedad en la construcción del modelo *random\_state* (Nulo, lo cual implicó que el modelo utilizara semillas aleatorias en cada ejecución).

El parámetro *kernel*, afecta cómo se mapean los datos a un espacio de características de mayor dimensión. El parámetro *C* controla la compensación entre obtener un límite de decisión suave y clasificar correctamente las muestras de entrenamiento. Un valor más grande de *C* puede llevar a un límite de decisión más complejo y, potencialmente, a sobreajuste. El valor de *C* escogido implicó un equilibrio en la clasificación correcta de ejemplos de entrenamiento, con la maximización del margen, permitiendo una cierta

flexibilidad en la clasificación errónea, con el objetivo de encontrar un compromiso que facilitara una buena generalización del modelo.

Para garantizar un comportamiento consistente se fijó el parámetro *random\_state* en *None*. Para obtener el mejor ajuste de este modelo se hicieron varias pruebas de ensayo y error, cambiando los tipos de kernel con los que se trabajaba y se escogió el mejor para los datos de prueba [104]. La implementación de este clasificador se encuentra en el **Anexo 6**.

### 2.7.2 Implementación de bosques aleatorios (RF)

En Python se utilizó la función *RandomForestClassifier* también de la biblioteca *Scikit-Learn* para obtener este modelo de clasificación. El modelo presenta varios hiperparámetros ajustables como son:

- *n\_estimators*
- *max\_depth*
- *min\_samples\_leaf*
- *random\_state*

El hiperparámetro *n\_estimators* controla la cantidad de árboles que se crean en el bosque, lo que ayuda a una mejora de rendimiento. Los valores por defecto de los parámetros que controlan el tamaño de los árboles (por ejemplo, *max\_depth*, *min\_samples\_leaf*) se ajustaron para reducir el consumo de memoria, la complejidad y el tamaño de los árboles. Se mantuvo *random\_state* en *None* para obtener un modelo consistente durante el ajuste. Se realizaron múltiples ensayos de prueba y error para encontrar el valor de los parámetros que se ajustaban mejor a los datos [105], [106]. En el **Anexo 7** se presenta la implementación de este clasificador.

## 2.8 Entrenamiento y validación

El conjunto de entrenamiento fue dividido en dos grupos, el 80% para el entrenamiento y el 20% para la validación. Para la validación del modelo se utilizaron como métricas la sensibilidad, precisión, exactitud y especificidad [107], definidas a continuación:

$$\text{Sensibilidad} = \frac{VP}{VP + FN} \quad (2)$$

$$\text{Precisión} = \frac{VP}{VP + FP} \quad (3)$$

$$\text{Especificidad} = \frac{VN}{VN + FP} \quad (4)$$

$$\text{Exactitud} = \frac{VP + VN}{VP + FP + VN + FN} \quad (5)$$

Los VP (verdaderos positivos) fueron nódulos identificados como positivos (malignos) por el sistema y consistentes con la anotación de la BD. Los VN (verdaderos negativos) son nódulos que el modelo identifica como negativos (benignos) y que también coinciden con el registro de la base de datos. Por otro lado, los FP (falsos positivos) son nódulos que el sistema detectó como positivos, al contrario de lo que muestra la BD. Los FN (falsos negativos) son nódulos que se detectan como negativos, cuando la anotación de la BD los muestra como positivos.

La sensibilidad mide la proporción de instancias positivas que fueron clasificadas correctamente entre todas las instancias reales positivas. También se conoce como tasa de verdaderos positivos. En contraposición, la especificidad es la proporción de instancias negativas, que fueron clasificadas correctamente entre todas las instancias reales negativas. La exactitud es la métrica que mide la proporción total de instancias clasificadas correctamente (tanto positivas como negativas) entre todas las instancias [108].

Los resultados de las métricas permitieron evaluar el rendimiento de los modelos.

## 2.9 Base de datos utilizada

El LIDC-IDRI ha creado una base de datos disponible pública y gratuita de imágenes de TC de tórax y sus anotaciones, realizadas por radiólogos expertos. Cuenta con 1018 tomografías computarizadas de tórax y sus anotaciones, de 1010 pacientes diferentes. Esta base de datos contiene 2669 lesiones etiquetadas como nódulos mayores o iguales de 3 mm, por al menos un radiólogo experto de 4 que participaron en el ejercicio de anotación de la BD. Contiene, además, 928 lesiones etiquetadas de esta forma por los cuatro radiólogos. Las notas de cada radiólogo sobre estas lesiones, incluyen contornos de nódulos y evaluaciones subjetivas de las características de los nódulos [80].

Los estudios que contiene son tomografías de dosis clínica baja, con un espesor de corte desde 0,7 mm hasta 5 mm. Los rangos de intervalos de reconstrucción de cortes están entre 0.45 mm y 5 mm. El tamaño del píxel se encuentra en un rango de 0,461 mm a 0,977 mm [80].

Los registros de diagnóstico LIDC-IDRI están en formato XML, uno para cada TC. Detallan la ubicación de los nódulos, el corte donde son más visibles y su grado de malignidad [80]. Esta BD que se utilizó para entrenar y validar el sistema CAD, puede encontrarse disponible en [Data from The Lung Image Database Consortium \(LIDC\) and Image Database Resource Initiative \(IDRI\): A completed reference database of lung nodules on CT scans \(LIDC-IDRI\) - The Cancer Imaging Archive \(TCIA\) Public Access - Cancer Imaging Archive Wiki](#) [79].

## 2.10 Hardware y software utilizado

Fue utilizado el lenguaje de programación *Python*, para el manejo de la anotación, la extracción de características de las lesiones y la corrida de los clasificadores. A través de los XML se empleó la asignación relacional de objetos (ORM del inglés *Object-relational mapping*) *Pylicd*, fue descargada libremente del sitio <https://github.com/notmatthancock/pylicd>, y se utilizó la plataforma integral de código abierto en Python *Pyradiomics* para la extracción de las características radiómicas.

Para ejecutar los programas y cálculos se utilizó una computadora con las siguientes prestaciones:

- Procesador (CPU): Intel® Core™ i5-4310M 2.70 GHz
- RAM: 8 GB DDR4
- Almacenamiento: HDD 500GB

## CAPÍTULO 3. RESULTADOS Y DISCUSIÓN

En el presente capítulo se analizan los resultados de la clasificación de nódulos pulmonares a partir de los clasificadores SVM y RF, utilizando características radiómicas. Se explican los dos experimentos realizados para dos vectores de características diferentes y se analizan varios modelos de clasificación con una salida en dos clases: maligno y benigno. Se evaluaron los modelos con la misma base de datos. Finalmente se discute el de mejor desempeño y se compara con resultados de la literatura científica.

### 3.1 Análisis de las características radiómicas extraídas

Con la aplicación de PyRadiomics a la segmentación en 3D de los nódulos, se obtuvieron vectores de 102 características radiómicas para cada nódulo. Estas características se dividen en:

- 14 características de forma
- 18 características estadísticas de primer orden
- 24 características de matriz de concurrencia de los niveles de gris (GLCM)
- 14 características de matriz de diferencia de nivel de gris (GLDM)
- 16 características de matriz de zona de tamaño de nivel de gris (GLSZM)
- 16 características de matriz de longitud de ejecución de nivel de gris (GLRLM)

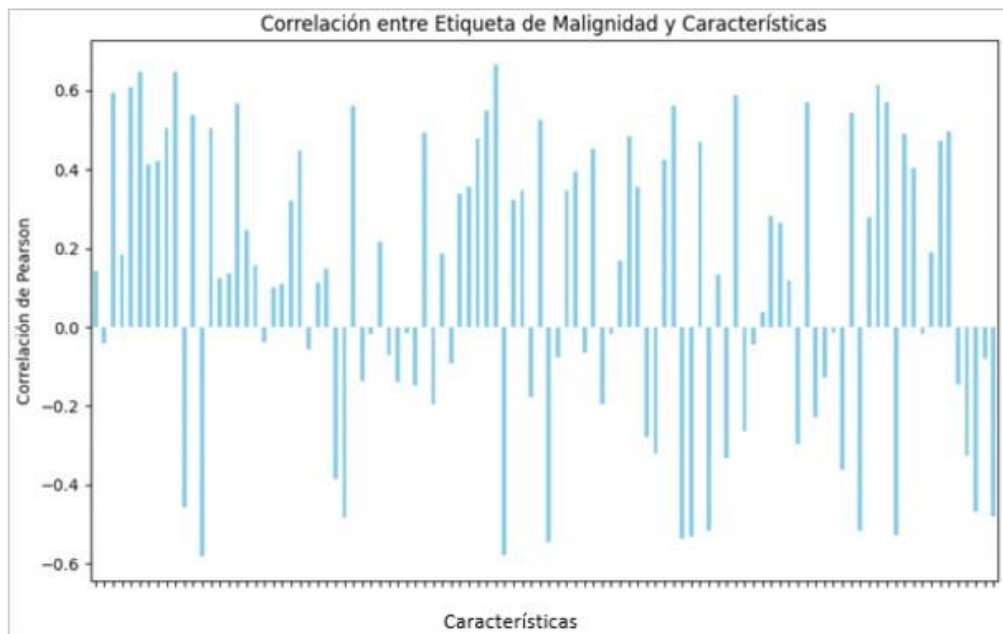
La tabla que contiene todas las características extraídas se encuentra en el **Anexo 8**.

### 3.2 Disminución de la dimensionalidad del problema

Entrenar un clasificador con una cantidad tan elevada de características puede causar un resultado desfavorable para la aplicación que se desea, no solo por la pérdida de eficiencia computacional, sino por la pérdida de eficacia para la tarea. Esto se debe a que no todas las características extraídas tienen igual poder de discriminación [92]. Dicho de otra manera, en la práctica es necesario reducir la dimensionalidad del problema de clasificación. Para esto se buscó discriminar entre las 102 características, cuáles fueron las estadísticamente

significativas en la determinación de malignidad. Para resolver este problema se calculó la correlación de Pearson entre cada característica y la anotación de la base de datos, en dos tipos de experimentos diferentes.

En el primer experimento se calculó la correlación de Pearson de cada característica con la anotación de la BD, para las etiquetas binarias: maligno-benigno. Los resultados se muestran en la Figura 3.1. Se puede apreciar que no todas las características fueron igual de sensibles para expresar las diferencias entre clases.



**Figura 3.1** Correlación entre las características y las etiquetas para 2 grados de malignidad

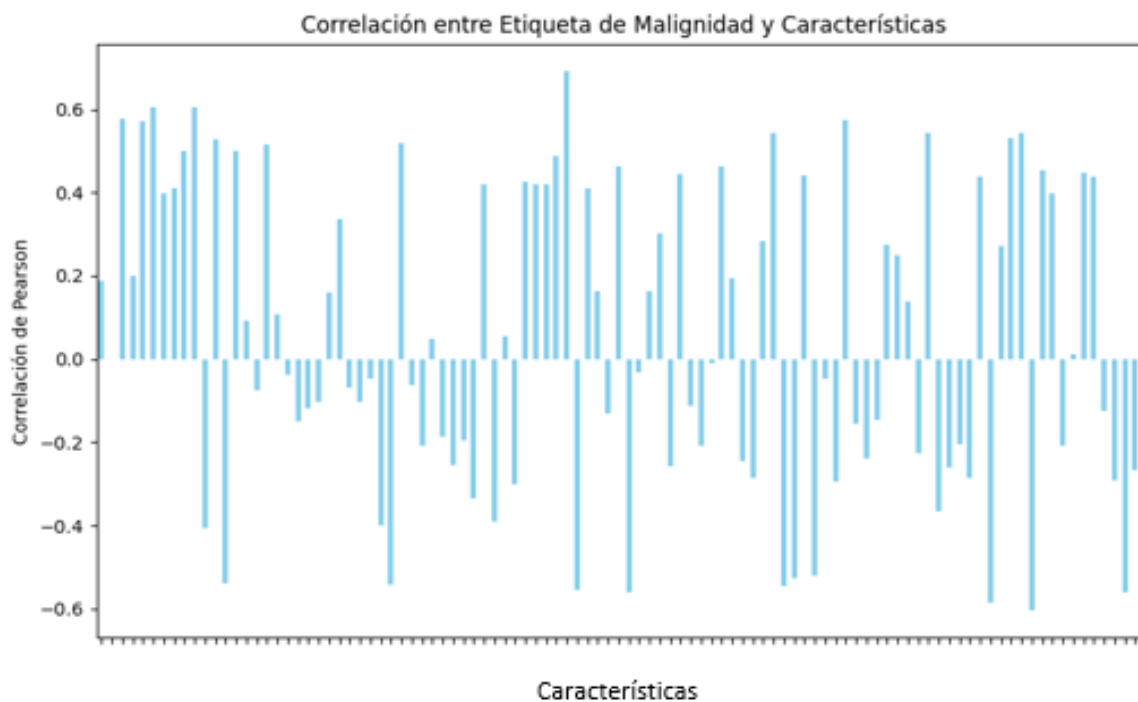
Debido a lo anterior se decidió seleccionar *a priori* un umbral de aceptación de 0.6, que se corresponde con una correlación de Pearson (  $R$  ) fuerte [101], y a partir de ahí, aceptar para la clasificación solo aquellas correlaciones donde  $R > 0,6$  y  $p < 0,05$ . En la Tabla 3.1, se aprecian los resultados de las correlaciones aceptadas como las más significativas. Como se puede apreciar en el primer experimento, el vector se redujo a solo cinco características, de las cuales tres son de forma, una es de GLCM y una de GLDM.

**Tabla 3.1** Características con una correlación de Pearson por encima de 0,6 y  $p < 0,05$  para las etiquetas de 2 grados de malignidad

Característica	Correlación
<i>original_shape_Maximum2DDiameterColumn</i>	0,608
<i>original_shape_Maximum2DDiameterRow</i>	0,647
<i>original_shape_MinorAxisLength</i>	0,646
<i>original_glcm_Imc1</i>	0,664
<i>original_gldm_DependenceEntropy</i>	0,613

Para el segundo experimento, se tuvo en cuenta el etiquetado en 5 grados de malignidad. La correlación mostró los resultados que se ilustran en la Figura 3.2 y las características fueron reducidas de 102 a solo 3, para las cuales  $R > 0,6$  y  $p < 0,05$ . Estas características seleccionadas coincidieron en parte con las que se seleccionaron en el primer experimento. Están divididas en dos características de forma y una de GLCM.

A pesar de que la selección fue parecida en ambos experimentos, los valores de correlación de Pearson no fueron exactamente iguales. Para el segundo análisis los valores de la correlación se pueden observar en la Tabla 3.2.



**Figura 3.2** Correlación entre las características y las etiquetas para 5 grados de malignidad

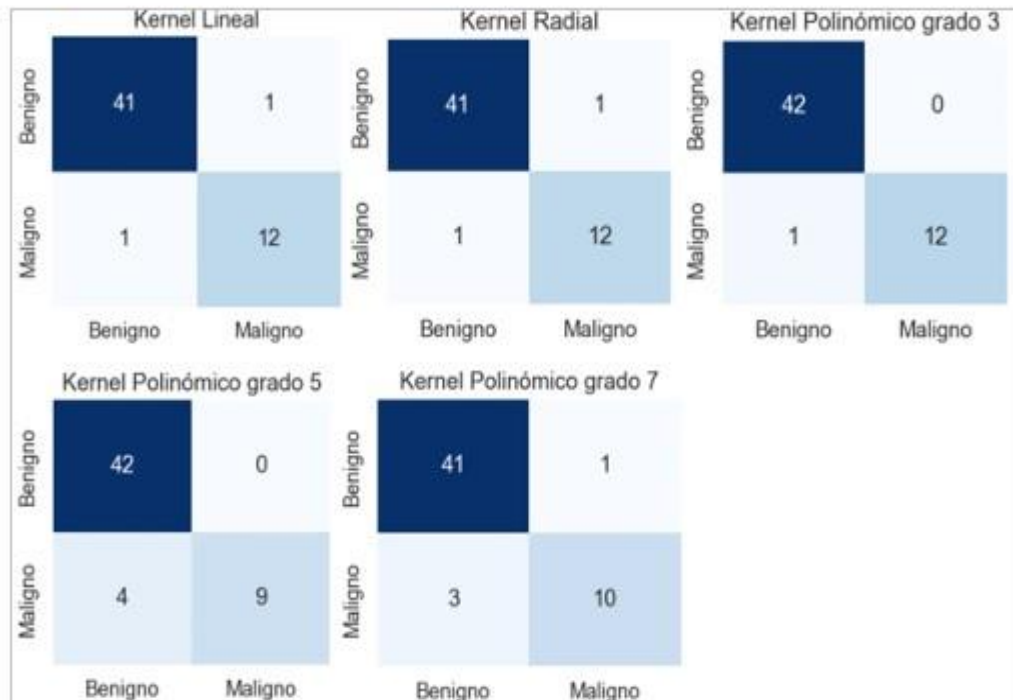
**Tabla 3.2** Características con una correlación por encima de 0.6 para las etiquetas de 5 grados de malignidad

Característica	Correlación
<i>original_shape_Maximum2DDiameterRow</i>	0,604
<i>original_shape_MinorAxisLength</i>	0,606
<i>original_glcm_Imc1</i>	0,692

Con el etiquetado en 5 grados de malignidad, hubo una mayor reducción de características determinantes de la clasificación. Sin embargo, los valores de correlación lineal de Pearson fueron ligeramente más bajos al analizarlos de forma individual, con respecto al primer experimento. En cambio, si se promedian los valores de correlación obtenidos en cada experimento, los resultados son similares y todos los valores entran dentro del rango de una correlación fuerte, por lo cual se espera que en el entrenamiento de los modelos se obtengan buenos resultados, tanto para aquellos cuyo vector de entrada contiene 5 características, como para cuando contiene solo 3.

### 3.3 Resultados de los clasificadores con etiquetado binario y no binario

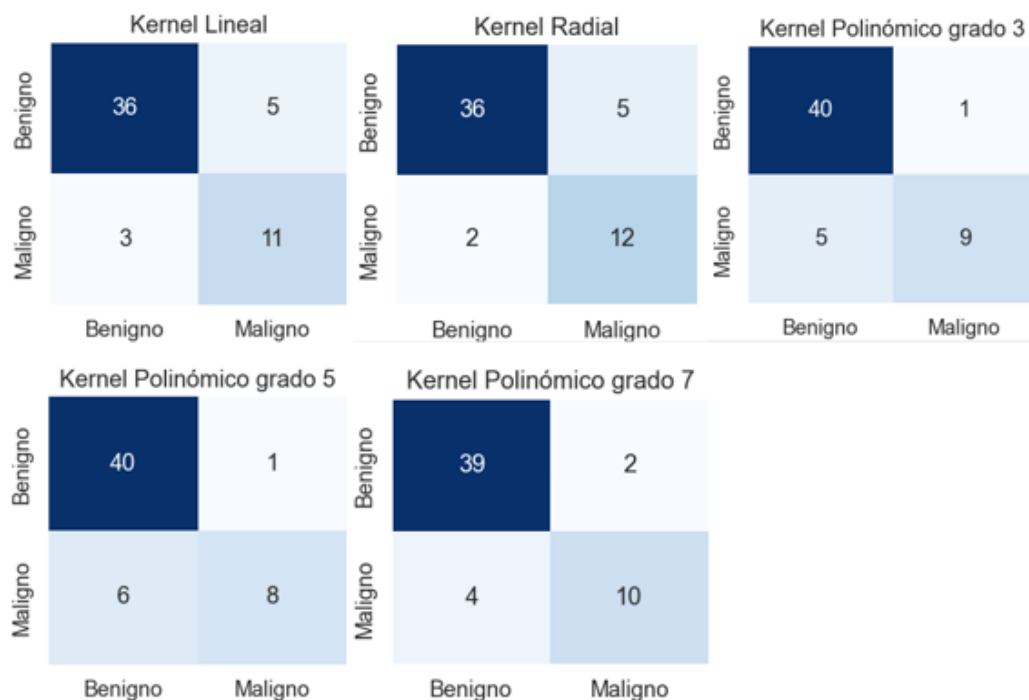
El clasificador SVM para varios tipos de kernel muestra los resultados que se aprecian en la Figura 3.3, para 5 características reducidas del etiquetado binario, y en la Figura 3.4, para 3 características reducidas del etiquetado no binario. Estas matrices de confusión permiten evaluar los modelos, proporcionando una visión más detallada de la clasificación. En el eje de las abscisas se expresa la predicción realizada por el modelo. El eje de las ordenadas muestra la anotación real de los datos. En la diagonal principal se muestran los aciertos de los mismos, mientras que, en la diagonal secundaria, en el extremo superior derecho se muestran los casos que clasificó incorrectamente como malignos (FP) y en el extremo inferior izquierdo, los que se clasificaron incorrectamente como benignos (FN). Las Tablas 3.3 y 3.4 resumen los resultados de las métricas, respectivamente.



**Figura 3.3** Matrices de confusión para los datos de validación de los modelos entrenados con etiquetado binario

**Tabla 3.3** Métricas de desempeño de los modelos entrenados con etiquetado binario

Modelo SVM	Precisión (%)	Sensibilidad (%)	Especificidad (%)	Exactitud (%)
Kernel lineal	92,31	92,31	97,62	96,36
Kernel radial	92,31	92,31	97,62	96,36
Kernel polinómico G3	100	92,31	100	98,18
Kernel polinómico G5	100	69,23	100	92,73
Kernel polinómico G7	90,91	76,92	97,61	92,73



**Figura 3.4** Matrices de confusión para los datos de validación de los modelos entrenados con etiquetado no binario

**Tabla 3.4** Métricas de desempeño de los modelos entrenados con etiquetado no binario

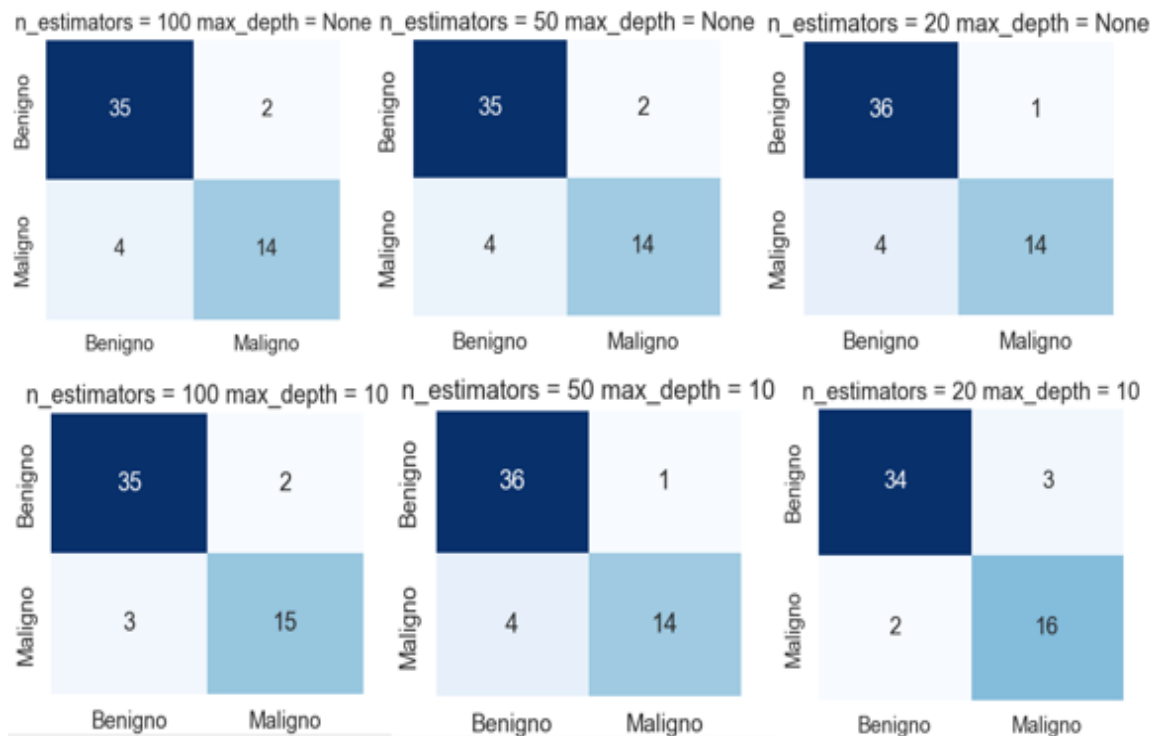
Modelo SVM	Precisión (%)	Sensibilidad (%)	Especificidad (%)	Exactitud (%)
Kernel lineal	68,75	78,57	87,80	85,45
Kernel radial	70,59	85,71	87,80	87,27
Kernel polinómico G3	90,00	64,29	97,56	89,09
Kernel polinómico G5	88,89	57,14	97,56	87,27
Kernel polinómico G7	83,33	71,43	95,12	89,09

En general, de los modelos de SVM entrenados, se aprecia que se comportaron mejor los del primer experimento. También se comprueba que cuando aumenta la especificidad, disminuye la sensibilidad para varios modelos.

Dentro de los modelos analizados de SVM, el clasificador que equilibra mejor los valores de sensibilidad y especificidad fue el modelo con kernel polinómico de grado 3, para el vector

de 5 características de entrada en etiquetado binario. Este fue capaz de clasificar correctamente todos los casos en los que el nódulo es benigno y solamente tuvo un fallo en la clasificación de un nódulo maligno, para los datos de validación.

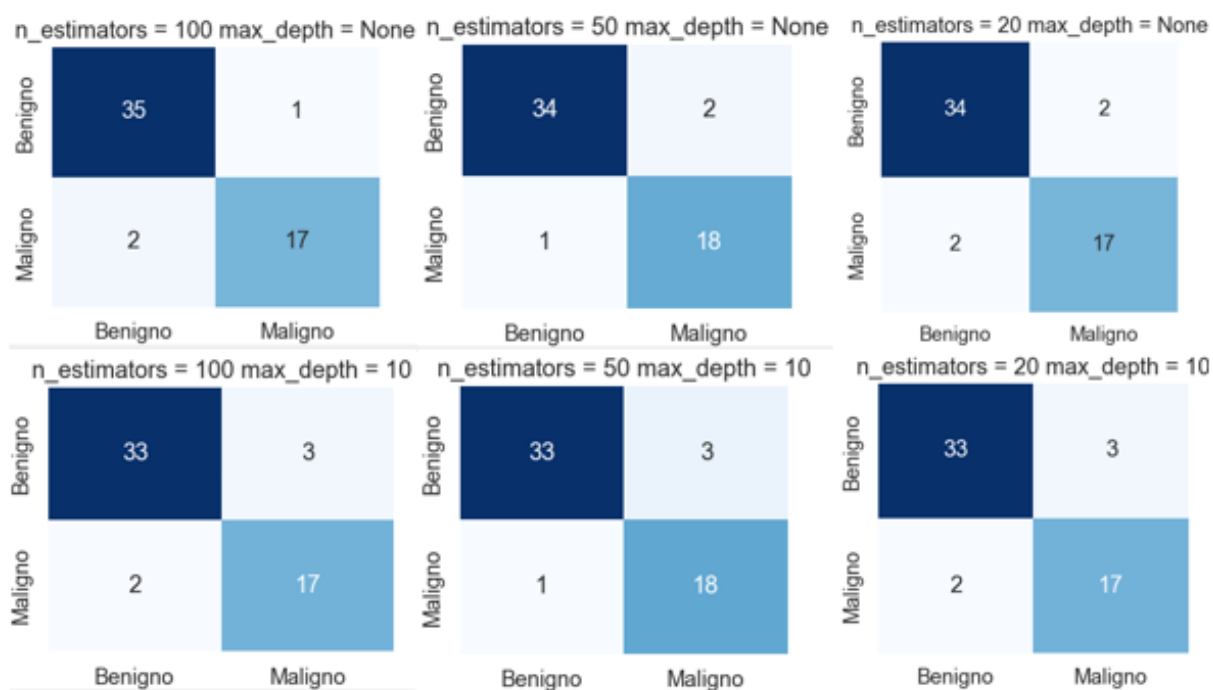
Para el entrenamiento del clasificador RF, se varió la cantidad de árboles de origen en 100, 50 y 20, para cada modelo, así como la profundidad de los mismos en *None* y 10, para regular su complejidad, según recomendaciones de [105], [106]. Los resultados obtenidos de estos modelos para las 5 características reducidas del etiquetado binario y 3 del no binario, se muestran en las matrices de confusión de las Figuras 3.5 y 3.6, respectivamente. En las Tablas 3.5 y 3.6 se resumen las métricas de los modelos entrenados.



**Figura 3.5** Matrices de confusión para datos de validación de los modelos RF entrenados con etiquetado binario

**Tabla 3.5** Métricas de desempeño de los modelos RF entrenados con etiquetado binario

Modelo RF	Precisión (%)	Sensibilidad (%)	Especificidad (%)	Exactitud (%)
<i>n_estimators</i> = 100 y <i>max_depth</i> = None	87,5	77,78	94,59	89,09
<i>n_estimators</i> = 50 y <i>max_depth</i> = None	87,5	77,78	94,59	89,09
<i>n_estimators</i> = 20 y <i>max_depth</i> = None	93,33	77,78	97,3	90,90
<i>n_estimators</i> = 100 y <i>max_depth</i> = 10	88,24	83,33	94,59	90,90
<i>n_estimators</i> = 50 y <i>max_depth</i> = 10	93,33	77,78	97,3	90,90
<i>n_estimators</i> = 20 y <i>max_depth</i> = 10	84,21	88,89	91,89	90,90

**Figura 3.6** Matrices de confusión para los datos de validación de los modelos RF entrenados con etiquetado binario

**Tabla 3.6** Métricas de desempeño de los modelos RF entrenados con etiquetado no binario

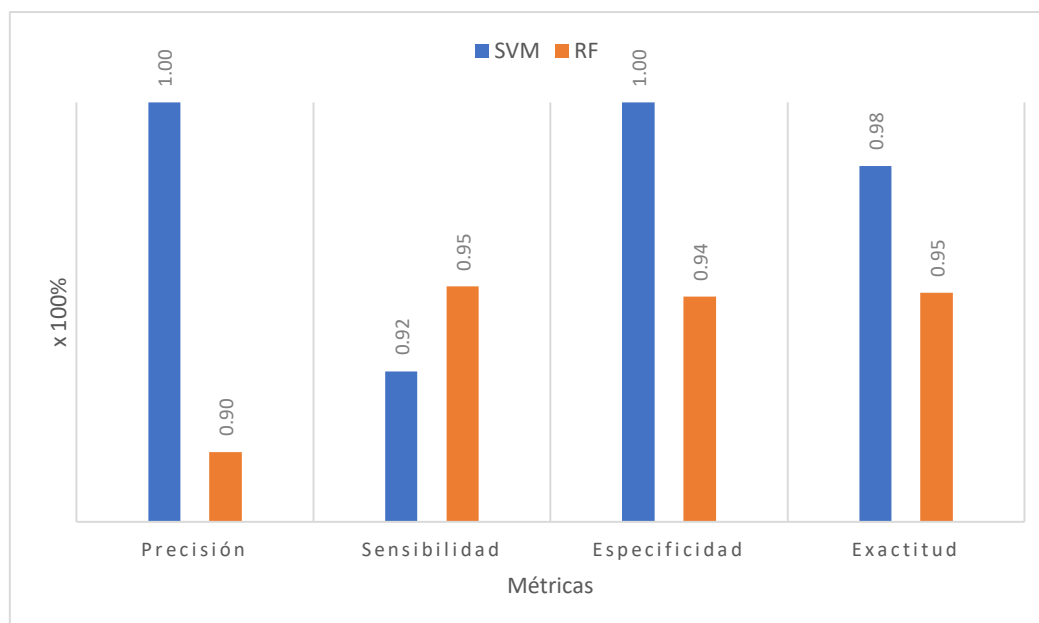
<b>Modelo RF</b>	<b>Precisión (%)</b>	<b>Sensibilidad (%)</b>	<b>Especificidad (%)</b>	<b>Exactitud (%)</b>
<i>n_estimators = 100 y max_depth = None</i>	94,44	89,47	97,22	94,55
<i>n_estimators = 50 y max_depth = None</i>	90	94,74	94,44	94,55
<i>n_estimators = 20 y max_depth = None</i>	89,47	89,47	94,44	92,72
<i>n_estimators = 100 y max_depth = 10</i>	85	89,47	91,67	90,91
<i>n_estimators = 50 y max_depth = 10</i>	85,71	94,74	91,67	92,72
<i>n_estimators = 20 y max_depth = 10</i>	85	89,47	91,67	90,91

A partir de los resultados de los modelos de RF obtenidos, se aprecia que fueron más bajos que los obtenidos con los modelos de SVM. Los mejores resultados se obtuvieron, en este caso, para vectores de entrada de 3 características reducidas con el etiquetado no binario.

De estos modelos, los de mejor resultado fueron los entrenados inicializando 100 y 50 árboles respectivamente, sin ninguna restricción en la profundidad de los mismos. Aunque entre estos dos modelos el primero presente valores más altos en las métricas de precisión y especificidad, la sensibilidad se vio más afectada con respecto al segundo modelo, y es precisamente esta métrica la que debe presentar un valor elevado a juicio de los radiólogos, para que el modelo pueda ser escogido. Por lo tanto, se evalúa el segundo modelo de la tabla 3.6, como el de mejor resultado para la tarea de clasificación requerida.

### 3.3.1 Valoración del mejor desempeño

En la figura 3.7 se grafican las métricas obtenidas de los modelos de SVM y RF entrenados y escogidos anteriormente como los de mejores resultados, para que sea más sencilla su comparación.



**Figura 3.7** Comparación entre los mejores modelos. En azul se grafica el modelo de SVM entrenado con kernel polinómico de grado 3, para las características reducidas del etiquetado binario. En rojo se grafica el modelo RF, entrenado con 50 árboles, para las características reducidas del etiquetado en cinco grados de malignidad.

Se puede apreciar que el modelo de RF escogido presenta una alta estabilidad en todas las métricas y posee además la mayor sensibilidad. El modelo de SVM por su parte, es el de mejor especificidad. Sin embargo, este resultado lo obtiene a expensas de una pérdida en sensibilidad.

Aunque se dispone de una data de diferente origen a la de entrenamiento-validación para realizar una prueba externa, que permita valorar de estos dos modelos, cuál es en realidad el mejor ajustado para tareas de clasificación de malignidad de nódulos pulmonares, por razones de tiempo esto no fue posible de ejecutar. Se trata de una data de PET-TC que requiere estandarizar sus datos de TC una vez separados de los de PET, y procesarlos para que entren a la plataforma PyRadiomics y al clasificador, bajo iguales parámetros que los utilizados en entrenamiento-validación, con la data LIDC-IDRI utilizada. En el presente trabajo fue precisamente la preparación de la data lo que consumió el mayor tiempo de experimentación. En estas circunstancias, la tarea de someter ambos sistemas a una prueba externa, queda como recomendación del Trabajo de Diploma y constituye la principal limitación del presente resultado.

### 3.4 Discusión

En la actualidad, para determinar la malignidad de los nódulos pulmonares se emplean biopsias, a partir de lo que subjetivamente resulta sospechoso para los radiólogos en las imágenes médicas [106], [109]. Con el avance de los métodos con inteligencia artificial, los expertos humanos están utilizando los sistemas CAD para reducir las tasas de falsos positivos y negativos y proporcionar una segunda opinión en las tareas de detección y clasificación [110].

Por otra parte, en los últimos años, las características radiómicas han demostrado su utilidad para diferenciar entre nódulos pulmonares benignos y malignos [27], [87], [92]. En comparación con los métodos cualitativos, basados puramente en la experiencia de los radiólogos, las características radiómicas extraídas de imágenes de TC están mostrando su potencial, al ser un método de diagnóstico no invasivo de alta precisión [93], a considerar en la Medicina personalizada moderna.

Respecto a la selección de características radiómicas buenas para separar las clases maligno-benigno, hay autores como [69] que afirman que la inclusión de un número elevado de características no ha mostrado mejorar el rendimiento de los clasificadores utilizados. En este estudio [69], se entrenaron modelos de aprendizaje supervisado con vectores de entrada que varían en su cantidad de características, desde 2 hasta 20. Entre ellas, *original\_shape\_Maximum2DDiameterRow* y *original\_shape\_MinorAxisLength* pertenecientes a los rasgos radiómicos de forma, fueron de las más fuertes para determinar la clasificación en ese trabajo. Ambas características coinciden con algunas de las seleccionadas en el presente estudio, para los dos tipos de etiquetado binario y no binario. En [69] obtuvieron los mejores resultados para un vector de características radiómicas entrada de 8 y 12 características, respectivamente. Aumentaron luego el número y comprobaron que no hubo mejoría en la clasificación.

En el presente estudio se aplicó un enfoque de cálculo de características radiómicas y se seleccionaron las más significativas, 5 en un caso y 3 en otro, para la clasificación, que son independientes, robustas y prominentes en los datos.

Si se decidiera escoger una mayor cantidad de características para entrenar los modelos no se apreciaría una mejora en los resultados, puesto que estas características no tienen un alto

nivel de importancia en los datos como las que ya se seleccionaron, existiendo la posibilidad de que incorporen redundancia a los mismos.

El modelo de SVM de base polinómica con grado 3, mostró grandes potencialidades para la clasificación. De la misma manera, el modelo de RF, basado en la creación de 50 árboles sin restricción de profundidad, también mostró buen potencial para la tarea. Sus resultados sugieren que ambos pueden ser sometidos a una prueba externa, que posibilite conocer cuál de los dos posee poder de generalización y cuál de los dos se desempeña mejor ante una data externa de diferente origen a la utilizada para entrenar / validar los modelos.

El costo computacional del entrenamiento de los modelos fue de 5 segundos por cada uno, tanto para los SVM como para los RF, para las prestaciones utilizadas (Procesador (CPU): Intel® Core™ i5-4310M 2.70 GHz, RAM: 8 GB DDR4, Almacenamiento: HDD 500GB) por lo que se puede considerar eficiente.

A partir de comparaciones de resultados con trabajos similares en la literatura científica, es posible ubicar los mejores modelos obtenidos en algún punto de referencia con relación a lo considerado como estándar a nivel mundial. Si comparamos los presentes resultados con los estudios [69], [89]–[94], presentados en la Tabla 1.1 del presente Trabajo de Diploma, que utilizan la misma base de datos y diversos clasificadores, se puede plantear que los sistemas propuestos en este trabajo, que funcionan sobre la base de un modelo de SVM y un RF respectivamente, presentan resultados comparables al estándar internacional.

Por otra parte, en el año 2022, en el grupo de investigación de la facultad que trabaja el procesamiento digital de imágenes y señales biomédicas, se desarrolló un primer intento de sistema para clasificar malignidad de nódulos pulmonares en imágenes de TC. Este sistema se implementó utilizando la misma BD, pero con un número más reducido de casos y no empleó métodos de inteligencia artificial. Para clasificar utilizó el criterio del cálculo de la dimensión fractal de los nódulos segmentados. Este enfoque fue muy utilizado a fines de los años 90 del pasado siglo y en la actualidad se está retomando dentro de sistemas con IA. Con ese trabajo se obtuvo una precisión, sensibilidad y especificidad del 96 % [95]. Este año se retomó el tema, pero esta vez con el enfoque más moderno que existe actualmente, que son las características radiómicas. En comparación con el trabajo precedente, la clasificación presenta resultados que superan a los anteriores, en cuanto a precisión y especificidad para el modelo de SVM.

El diagnóstico preciso de los nódulos pulmonares malignos mediante tomografía computarizada es fundamental y constituye uno de los mayores desafíos en la práctica clínica oncológica actual [111]. En este sentido, los modelos obtenidos brindan un método no invasivo de clasificación de malignidad, con eficacia para la tarea y con alta eficiencia computacional. De esta manera permitirían, una vez implementados en la práctica de rutina clínica, a través de un CAD que incluya todas las etapas descritas, ganar tiempo en el manejo de pacientes con tumores malignos agresivos. Se podrían detectar y clasificar los nódulos en su fase inicial, por lo que el trabajo realizado presenta impacto, no solo científico, sino social. Se debe de resaltar que, tanto el aprendizaje profundo como la radiómica, han demostrado un potencial para identificar tumores malignos en volúmenes de TC, y han dado como resultado un rendimiento comparable en diferentes conjuntos de datos [112]–[114]. Sin embargo, aún faltan pruebas concluyentes que respalden que un tipo de método sea superior al otro.

Aunque los descriptores radiómicos están diseñados, inspirados en características de imágenes radiológicas cuantitativas, los modelos basados en CNN se entrenan para detectar características abstractas con alta relevancia con respecto a las etiquetas de clase [100]. El rendimiento de las CNN pudiera en algunos casos superar al radiómico, especialmente cuando hay patrones complejos en los datos. Sin embargo, debido a la complejidad de sus arquitecturas y de sus resultados, por ejemplo en mapas de GradCam, a menudo se dificulta su interpretación clínica [115], además de que pueden requerirse recursos computacionales significativos para su entrenamiento y un número muy alto de imágenes de entrenamiento [116] para lograr resultados con poder de generalización.

La elección entre el enfoque radiómico y el de CNN depende de varios factores, como el tamaño del conjunto de datos, la interpretabilidad requerida y la complejidad de los patrones a aprender [113]. Una combinación de ambas técnicas también podría ser considerada para aprovechar las fortalezas de ambos métodos. En este sentido, una de las ventajas de los presentes resultados, con enfoque de radiómica y ML, es que han sido logrados a partir de un conjunto de datos relativamente pequeños y con prestaciones computacionales modestas.

### **3.5 Análisis económico y medioambiental**

El sistema propuesto, para ser utilizado en rutina clínica en el futuro, debe de ser sometido a revisión por la autoridad reguladora cubana (CECMED). Una vez certificado se podría introducir en el Sistema de salud cubano, e incluso comercializar. No obstante, para que sea posible su comercialización se requiere realizar una estimación de los requisitos mínimos de hardware para su funcionamiento en una entidad de salud y confeccionarle una interfaz de usuario.

En la actualidad existe un pequeño grupo de sistemas de detección y clasificación con IA, aprobados por entidades reguladoras internacionales, para su empleo en instalaciones médicas de manera profesional. La mayoría de ellos se localizan en Estados Unidos y Europa, donde se emplean como tercer lector. Proceden de varias compañías privadas que venden software biomédico, como son: Context Flow, Lunit, Sectra, entre otros. Todos poseen altos precios y pago de licencias a sus propietarios. El *Auto Lung Nodule Detection* de Samsung por ejemplo, tiene un costo de \$300 000 USD y sus licencias de software son de \$670.00 USD anuales. Estos costos dan una idea de lo que se ahorraría el Sistema de Salud Pública cubano al disponer de un CAD autóctono para clasificar nódulos pulmonares.

El sistema presentado en este estudio trabaja sobre la base de imágenes de TC. El paciente debe ser sometido a la exposición de radiaciones ionizantes (rayos X, en este caso) al realizarse una TC de tórax para conocer su condición. Esto implica un cierto riesgo de efecto biológico estocástico, al recibir una dosis de radiación de aproximadamente 8 mSv [117], [118]. En este sentido, para cuidar al hombre, como ente fundamental del medio ambiente, la adquisición de las tomografías deben estar regidas por los Principios de la Protección radiológica, de Justificación y Optimización de la práctica [119], [120]. Una vez obtenida la TC, el sistema CAD con IA no agrede al medio ambiente. No obstante, dada la posibilidad de obtención de un diagnóstico temprano de una enfermedad mortal, se justifica plenamente la exposición radiológica.

### **3.6 Conclusiones del capítulo**

Esta investigación permitió concluir que el empleo de características radiómicas, obtenidas a partir de imágenes de TC, es eficaz para la tarea de clasificar nódulos pulmonares en benigno-maligno, utilizando métodos de aprendizaje de máquinas. Estas características incluyen una amplia diversidad de rasgos, varios de los cuales son determinantes de la

---

clasificación. Además, se comprobó que clasificadores de ML como RF y SVM, alimentados con características radiómicas, clasifican bien las clases de nódulos pulmonares, realizando su entrenamiento con un número relativamente bajo de datos. No obstante, se requiere la realización de una prueba externa para evaluar el poder de generalización de los mejores modelos obtenidos, para conformar un sistema CAD. Los valores de sensibilidad y especificidad obtenidos por los modelos escogidos en validación, superan el 90 %, siendo del orden o superando a los obtenidos por sistemas CAD internacionales profesionales.

## CONCLUSIONES y RECOMENDACIONES

### Conclusiones

Se ha diseñado un sistema CAD para la clasificación de malignidad de nódulos pulmonares, basado en la extracción de características radiómicas a partir de PyRadiomics y clasificadores de ML.

Los modelos de clasificación incluidos en el CAD propuesto, han permitido la caracterización cuantitativa de las lesiones con eficacia y eficiencia computacional, cumpliéndose la hipótesis de investigación.

Se obtuvieron 102 características radiómicas, a partir del empleo de la plataforma PyRadiomics. Sin embargo, solo 5 fueron significativas para la clasificación, a partir de etiquetado binario y 3 utilizando el etiquetado en 5 clases.

Dos clasificadores de SVM y RF, cuyos modelos fueron programados en Python y entrenados para el 80 % de los datos de una BD internacional, demostraron ser eficaces para diferenciar, a partir de imágenes de TC, nódulos malignos y benignos.

Al evaluar los dos modelos de clasificación de mejor desempeño, se obtuvieron valores de precisión, sensibilidad, especificidad y exactitud por encima del 90%, los cuales igualan o superan el estándar internacional actual, así como resultados precedentes del grupo de investigación, donde se utilizaron métodos basados en fractales.

### Recomendaciones

Se recomienda evaluar los mejores modelos de clasificación con ML obtenidos, utilizando una base de datos externa, que permita demostrar el poder de generalización.

**REFERENCIAS BIBLIOGRÁFICAS**

- [1] B. S. Chhikara and K. Parang. Global Cancer Statistics 2022: the trends projection analysis. *Chemical Biology Letters Chem. Biol. Lett* 10(1): 1–16, 2023.
- [2] R. L. Siegel, K. D. Miller, H. E. Fuchs, and A. Jemal. Cancer statistics, 2022. *CA: A Cancer Journal for Clinicians* 72(1): 7–33, Jan. 2022. doi: <https://doi.org/10.3322/caac.21708>.
- [3] Ministerio de Salud Pública and Dirección de Registros Médicos y Estadísticas de Salud. Anuario estadístico de Salud 2021, 2022.
- [4] ONEI. Anuario Estadístico de Cuba 2021: 60, 2022.
- [5] D. M. Hansell, A. A. Bankier, H. MacMahon, T. C. McLoud, N. L. Müller, and J. Remy. Fleischner Society: Glossary of terms for thoracic imaging. *Radiology* 246(3): 697–722, 2008. doi: 10.1148/radiol.2462070712.
- [6] K. Loverdos, A. Fotiadis, C. Kontogianni, M. Iliopoulou, and M. Gaga. Lung nodules: A comprehensive review on current approach and management. *Annals of Thoracic Meedicine* 14(4): 226–238, 2019. doi: 10.4103/atm.ATM.
- [7] R. Li, C. Xiao, Y. Huang, H. Hassan, and B. Huang. Deep Learning Applications in Computed Tomography Images for Pulmonary Nodule Detection and Diagnosis: A Review, *Diagnostics*, 12(2). p.: 1–21, 2022. doi: 10.3390/diagnostics12020298.
- [8] T. M. Buzug. Milestones of Computed Tomography, in *Computed Tomography*, 2011, p.: 311–342.
- [9] D. Gu, G. Liu, and Z. Xue. On the performance of lung nodule detection, segmentation and classification. *Computerized Medical Imaging and Graphics* 89(February): 1–15, 2021. doi: 10.1016/j.compmedimag.2021.101886.
- [10] R. Hossain, C. C. Wu, P. M. de Groot, B. W. Carter, M. D. Gilman, and G. F. Abbott. Missed Lung Cancer. *Radiologic Clinics of North America* 56(3): 365–375, 2018. doi: 10.1016/j.rcl.2018.01.004.
- [11] K. Camacho-Sosa, M. C. Martí-Coruña, V. G. Ferreira Moreno, J. García Soto, L. Alonso Lemus, and I. E. Carreño Rolando. Eficacia de la tomografía en el seguimiento del cáncer de pulmón tratado con inmunoterapia cubana. *Revista Médica Electrónica* 44(2): 278–287, 2022.
- [12] A. Rivero Castro, Y. Rivera Suárez, Y. Borges González, and Y. Naranjo Gorrín. Algoritmo para la identificación de nódulos pulmonares solitarios en imágenes de tomografía de tórax. *Revista Cubana de Informática Médica* 7(1): 73–88, 2015.
- [13] M. Mazonakis and J. Damilakis. Computed tomography: what and how does it measure?. *European Journal of Radiology* 85(8): 1499–1504, 2016. doi: 10.1016/j.ejrad.2016.03.002.
- [14] P. J. Withers, C. Bouman, S. Carmignato, V. Cnudde, D. Grimaldi, C. K. Hagen, *et*

- al.* X-ray computed tomography. *Nature Reviews Methods Primers* 1(1): 18, Feb. 2021. doi: 10.1038/s43586-021-00015-4.
- [15] J. Wang, L. Lin, S. Zhao, X. Wu, and S. Wu. Research progress on computed tomography image detection and classification of pulmonary nodule based on deep learning. *Journal of biomedical engineering* 36(4): 670–676, 2019. doi: 10.7507/1001-5515.201806019.
- [16] J. Ma, Z. Zhou, Y. Ren, J. Xiong, L. Fu, Q. Wang, *et al.* Computerized Detection of Lung Nodules through Radiomics. *International Journal of Laboratory Hematology* 38(1): 42–49, 2016. doi: 10.1111/ijlh.12426.
- [17] N. Petrick, B. Sahiner, S. G. A. Iii, A. Bert, L. Correale, S. Delsanto, *et al.* Evaluation of computer-aided detection and diagnosis systems. *Medical Physics* 40(July): 1–17, 2013. doi: 10.1118/1.4816310.
- [18] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. Adiyoso Setio, F. Ciampi, M. Ghafoorian, *et al.* A survey on deep learning in medical image analysis. *Medical Image Analysis* 42(December 2012): 60–88, 2017. doi: 10.1016/j.media.2017.07.005.
- [19] J. D. López-Cabrera, L. A. López Rodríguez, and M. Pérez-Díaz. Classification of breast cancer from digital mammography using deep learning. *Inteligencia Artificial* 23(65): 56–66, 2020. doi: 10.4114/intartif.vol23iss65pp56-66.
- [20] S. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*, (2004). 2004.
- [21] A. Géron. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow*, O’Reilly Media, Inc., 2022.
- [22] J. Wang, H. Zhu, S. H. Wang, and Y. D. Zhang. A Review of Deep Learning on Medical Image Analysis. *Mobile Networks and Applications* 26(1): 351–380, 2020. doi: 10.1007/s11036-020-01672-7.
- [23] G. Chartrand, P. M. Cheng, E. Vorontsov, M. Drozdzal, S. Turcotte, C. J. Pal, *et al.* Deep learning: A primer for radiologists. *Radiographics* 37(7): 2113–2131, 2017. doi: 10.1148/rg.2017170077.
- [24] B. Chen, R. Zhang, Y. Gan, L. Yang, and W. Li. Development and clinical application of radiomics in lung cancer. *Radiation Oncology* 12(1): 154, Dec. 2017. doi: 10.1186/s13014-017-0885-x.
- [25] R. Wilson and A. Devaraj. Radiomics of pulmonary nodules and lung cancer. *Translational Lung Cancer Research* 6(1): 86–91, 2017. doi: 10.21037/tlcr.2017.01.04.
- [26] R. Paul, S. H. Hawkins, M. B. Schabath, R. J. Gillies, L. O. Hall, and D. B. Goldgof. Predicting malignant nodules by fusing deep features with classical radiomics features. *Journal of Medical Imaging* 5(1): 1–11, 2018. doi: 10.1117/1.jmi.5.1.011021.
- [27] P. Lambin, E. Rios-Velazquez, R. Leijenaar, S. Carvalho, R. G. P. M. Van Stiphout, P. Granton, *et al.* Radiomics: Extracting more information from medical images using advanced feature analysis. *European Journal of Cancer* 48(4): 441–446, 2012. doi:

- 10.1016/j.ejca.2011.11.036.
- [28] L. Deng and D. Yu. Deep learning: methods and applications. *Foundations and trends® in signal processing* 7(3–4): 197–387, 2014.
- [29] R. Gruetzemacher and A. Gupta. Using deep learning for pulmonary nodule detection & diagnosis. *AMCIS 2016: Surfing the IT Innovation Wave - 22nd Americas Conference on Information Systems*: 1–9, 2016.
- [30] B. Mahesh. Machine Learning Algorithms - A Review. *International Journal of Science and Research* 9(1): 381–386, 2020. doi: 10.21275/ART20203995.
- [31] P. Wang, E. Fan, and P. Wang. Comparative Analysis of Image Classification Algorithms Based on Traditional Machine Learning and Deep Learning. *Pattern Recognition Letters* 141(January): 61–67, 2020. doi: 10.1016/j.patrec.2020.07.042.
- [32] X. Zhang, Y. Zhang, G. Zhang, X. Qiu, and W. Tan. Deep Learning With Radiomics for Disease Diagnosis and Treatment: Challenges and Potential. *Frontiers in Oncology* 12(February): 1–25, 2022. doi: 10.3389/fonc.2022.773840.
- [33] C. N. de Innovación. Reunión del Consejo Nacional de Innovación (Acta 11 / 2022), 2022. <https://www.presidencia.gob.cu/es/documentos/reunion-del-consejo-nacional-de-innovacion-acta-11-2022/>
- [34] R. L. Xie, Y. Wang, Y. N. Zhao, J. Zhang, G. B. Chen, J. Fei, *et al.* Lung nodule pre-diagnosis and insertion path planning for chest CT images. *BMC Medical Imaging* 23(22): 1–16, 2023. doi: 10.1186/s12880-023-00973-z.
- [35] H. Mahersia and M. Zaroug. Lung Cancer Detection on CT Scan Images - A review techniques. *International Journal of Advanced Research in Artificial Intelligence* 4(4): 38–45, 2015.
- [36] J. John and M. G. Mini. Multilevel Thresholding Based Segmentation and Feature Extraction for Pulmonary Nodule Detection. *Procedia Technology* 24: 957–963, 2016. doi: 10.1016/j.protcy.2016.05.209.
- [37] J. T. Bushberg and J. M. Boone. *The essential physics of medical imaging*, Lippincott Williams & Wilkins, 2011.
- [38] G. E. T. Orozco. Tomografía Computarizada, *Revista Tame*, 2019. <https://es.scribd.com/doc/17184942/Principios-Fisicos-de-la-Tomografia-Computarizada>
- [39] V. Balac. Radiographic Imaging, in *Introduction to Radiologic and Imaging Sciences and Patient Care E-Book*, Elsevier Health Sciences, 2022, p.: 70.
- [40] S. Shin, M. W. Kim, K. H. Jin, K. M. Yi, Y. Kohmura, T. Ishikawa, *et al.* Deep 3D reconstruction of synchrotron X-ray computed tomography for intact lungs. *Scientific Reports* 13(1): 1–9, 2023. doi: 10.1038/s41598-023-27627-y.
- [41] E. Comission. *European Guidelines on Quality Criteria for Diagnostic Radiographic Images*, European Commission, EUR 16260 EN, giugno 1996, 1996. [Online]. Available: <https://publications.europa.eu/en/publication-detail/>

/publication/d229c9e1-a967-49de-b169-59ee68605f1a

- [42] F. Zarb, L. Rainford, and M. F. McEntee. Image quality assessment tools for optimization of CT images. *Radiography* 16(2): 147–153, 2010. doi: 10.1016/j.radi.2009.10.002.
- [43] E. Seeram. Computed tomography: Physical principles and recent technical advances. *Journal of Medical Imaging and Radiation Sciences* 41(2): 87–109, 2010. doi: 10.1016/j.jmir.2010.04.001.
- [44] F. Michael and M. McNitt-Gray. Tradeoffs in CT image quality and dose. *Med. Phys.* 33: 2154–2155, 2006.
- [45] J. T. Payne. CT radiation dose and image quality. *Radiologic Clinics of North America* 43(6): 953–962, 2005. doi: 10.1016/j.rcl.2005.07.002.
- [46] K. Gulliksrud, C. Stokke, and A. C. Trægde Martinsen. How to measure CT image quality: Variations in CT-numbers, uniformity and low contrast resolution for a CT quality assurance phantom. *Physica Medica* 30(4): 521–526, 2014. doi: 10.1016/j.ejmp.2014.01.006.
- [47] M. Firmino, A. H. Morais, R. M. Mendoca, M. R. Dantas, H. R. Hekis, and R. Valentim. Computer-aided detection system for lung cancer in computed tomography scans: Review and future prospects. *BioMedical Engineering Online* 13(41): 1–16, 2014.
- [48] H. H. Popper. Progression and metastasis of lung cancer. *Cancer and Metastasis Reviews* 35(1): 75–91, 2016. doi: 10.1007/s10555-016-9618-0.
- [49] Shallu, P. Nanglia, S. Kumar, and A. Kumar Luhach. Detection and Analysis of Lung Cancer using Radiomic Approach, in *Smart Computational Strategies: Theoretical and Practical Aspects*, Springer Singapore, 2019, p.: 1–286. doi: 10.1007/978-981-13-6295-8\_2.
- [50] D. Albano, R. Gatta, M. Marini, C. Rodella, L. Camoni, F. Dondi, *et al.* Role of f-fdg pet/ct radiomics features in the differential diagnosis of solitary pulmonary nodules: Diagnostic accuracy and comparison between two different pet/ct scanners. *Journal of Clinical Medicine* 10(21): 1–13, 2021. doi: 10.3390/jcm10215064.
- [51] G. Zheng, G. Han, and N. Q. Soomro. An inception module CNN classifiers fusion method on pulmonary nodule diagnosis by signs. *Tsinghua Science and Technology* 25(3): 368–383, 2020. doi: 10.26599/TST.2019.9010010.
- [52] Y. J. Jeong, C. A. Yi, and K. S. Lee. Nódulos pulmonares solitarios: Detección, caracterización y guías para su diagnóstico y tratamiento. *Radiologia* 50(3): 183–195, 2008. doi: 10.1016/S0033-8338(08)71964-7.
- [53] J. Gong, J. Liu, W. Hao, S. Nie, S. Wang, and W. Peng. Computer-aided diagnosis of ground-glass opacity pulmonary nodules using radiomic features analysis. *Physics in Medicine and Biology* 64(13): 1–11, 2019. doi: 10.1088/1361-6560/ab2757.
- [54] J. Z. Cheng, D. Ni, Y. H. Chou, J. Qin, C. M. Tiu, Y. C. Chang, *et al.* Computer-Aided

- Diagnosis with Deep Learning Architecture: Applications to Breast Lesions in US Images and Pulmonary Nodules in CT Scans. *Scientific Reports* 6(March): 1–13, 2016. doi: 10.1038/srep24454.
- [55] S. Singh, J. Maxwell, J. A. Baker, J. L. Nicholas, and J. Y. Lo. Computer-aided classification of breast masses: Performance and interobserver variability of expert radiologists versus residents. *Radiology* 258(1): 73–80, 2011. doi: 10.1148/radiol.10081308.
- [56] M. L. Giger, N. Karssemeijer, and J. A. Schnabel. Breast image analysis for risk assessment, detection, diagnosis, and treatment of cancer. *Annual Review of Biomedical Engineering* 15(April): 327–357, 2013. doi: 10.1146/annurev-bioeng-071812-152416.
- [57] S. Joo, Y. S. Yang, W. K. Moon, and H. C. Kim. Computer-aided diagnosis of solid breast nodules: use of an artificial neural network based on multiple sonographic features. *IEEE Transactions on Medical Imaging* 23(10): 1292–1300, 2004.
- [58] T. W. Way, B. Sahiner, H. P. Chan, L. Hadjiiski, P. N. Cascade, A. Chughtai, *et al.* Computer-aided diagnosis of pulmonary nodules on CT scans: Improvement of classification performance with nodule surface features. *Medical Physics* 36(7): 3086–3098, 2009. doi: 10.1118/1.3140589.
- [59] R. Agarwal, A. Shankhadhar, and R. K. Sagar. Detection of lung cancer using content based medical image retrieval. *International Conference on Advanced Computing and Communication Technologies, ACCT 2015-April(33)*: 48–52, 2015. doi: 10.1109/ACCT.2015.33.
- [60] S. Bathia, Y. Sinha, and L. Goel. Lung Cancer Detection: A deep learning approach. *Springer* 2(January): 699–705, 2019. doi: 10.1007/978-981-13-1595-4.
- [61] C. Dang, Y. Liu, H. Yue, J. X. Qian, and R. Zhu. Autumn Crop Yield Prediction using Data-Driven Approaches:- Support Vector Machines, Random Forest, and Deep Neural Network Methods. *Canadian Journal of Remote Sensing* 47(2): 162–181, 2021. doi: 10.1080/07038992.2020.1833186.
- [62] M. Schonlau and R. Y. Zou. The random forest algorithm for statistical learning. *Stata Journal* 20(1): 3–29, 2020. doi: 10.1177/1536867X20909688.
- [63] H. Chung and K. shik Shin. Genetic algorithm-optimized multi-channel convolutional neural network for stock market prediction. *Neural Computing and Applications* 32(12): 7897–7914, 2020. doi: 10.1007/s00521-019-04236-3.
- [64] R. Chauhan, K. K. Ghanshala, and R. C. Joshi. Convolutional Neural Network (CNN) for Image Detection and Recognition. *ICSCCC 2018 - First International Conference on Secure Cyber Computing and Communications*: 278–282, 2018. doi: 10.1109/ICSCCC.2018.8703316.
- [65] T. S. R. Mhathesh, J. Andrew, K. Martin Sagayam, and L. Henesey. *A 3d convolutional neural network for bacterial image classification*. 1167, Springer Singapore, 2021. doi: 10.1007/978-981-15-5285-4\_42.

- [66] S. Albawi, T. A. M. Mohammed, and S. Alzawi. Understanding of a Convolutional Neural Network. *Icet2017*: 1–6, 2017.
- [67] F. Chollet. *Deep Learning with Python*, 2021. doi: 10.1007/978-1-4842-5364-9.
- [68] N. Aloysius and M. Geetha. A review on deep convolutional neural networks. *Proceedings of the 2017 IEEE International Conference on Communication and Signal Processing, ICCSP 2017 2018-Janua*: 588–592, 2018. doi: 10.1109/ICCSP.2017.8286426.
- [69] H. Shakir, H. Rasheed, and T. M. Rasool Khan. Radiomic feature selection for lung cancer classifiers. *Journal of Intelligent and Fuzzy Systems* 38(5): 5847–5855, 2020. doi: 10.3233/JIFS-179672.
- [70] D. A. Pisner and D. M. Schnyer. Support vector machine. *Machine Learning: Methods and Applications to Brain Disorders*: 101–121, 2019. doi: 10.1016/B978-0-12-815739-8.00006-7.
- [71] M. R. Tomaszewski and R. J. Gillies. The biological meaning of radiomic features. *Radiology* 298(3): 505–516, 2021. doi: 10.1148/radiol.2021202553.
- [72] P. Lambin, R. T. H. Leijenaar, T. M. Deist, J. Peerlings, E. E. C. De Jong, J. Van Timmeren, *et al.* Radiomics: The bridge between medical imaging and personalized medicine. *Nature Reviews Clinical Oncology* 14(12): 749–762, 2017. doi: 10.1038/nrclinonc.2017.141.
- [73] C. P. Langlotz. Will artificial intelligence replace radiologists?. *Radiology: Artificial Intelligence* 1(3): 16–18, 2019. doi: 10.1148/ryai.2019190058.
- [74] A. D. de Leon, P. Kapur, and I. Pedrosa. Radiomics in Kidney Cancer: MR Imaging. *Magnetic Resonance Imaging Clinics of North America* 27(1): 1–13, 2019. doi: 10.1016/j.mric.2018.08.005.
- [75] V. Kumar, Y. Gu, S. Basu, A. Berglund, S. A. Eschrich, M. B. Schabath, *et al.* Radiomics: The process and the challenges. *Magnetic Resonance Imaging* 30(9): 1234–1248, 2012. doi: 10.1016/j.mri.2012.06.010.
- [76] H. Alsleem, R. Davidson, and M. Mi. Factors Affecting Contrast-Detail Performance in Computed Tomography : A Review. *Journal of Medical Imaging and Radiation Sciences* 44(2): 62–70, 2013. doi: 10.1016/j.jmir.2012.12.001.
- [77] Y. Li, Y. Jiang, H. Liu, X. Yu, S. Chen, D. Ma, *et al.* A phantom study comparing low-dose CT physical image quality from five different CT scanners. *Quantitative Imaging in Medicine and Surgery* 12(1): 766–780, 2022. doi: 10.21037/qims-21-245.
- [78] R. J. Gillies, P. E. Kinahan, and H. Hricak. Radiomics: Images are more than pictures, they are data. *Radiology* 000(0): 1–15, 2016.
- [79] K. Clark, B. Vendt, K. Smith, J. Freymann, J. Kirby, P. Koppel, *et al.* The Cancer Imaging Archive ( TCIA ): Maintaining and Operating a Public Information Repository2013. doi: 10.1007/s10278-013-9622-7.
- [80] S. G. Armato, G. McLennan, L. Bidaut, M. F. McNitt-Gray, C. R. Meyer, A. P.

- Reeves, *et al.* The Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI): A completed reference database of lung nodules on CT scans. *Medical Physics* 38(2): 915–931, 2011. doi: 10.1118/1.3528204.
- [81] G. Wu, A. Jochems, T. Refaee, A. Ibrahim, C. Yan, S. Sanduleanu, *et al.* Structural and functional radiomics for lung cancer. *European Journal of Nuclear Medicine and Molecular Imaging* 48(12): 3961–3974, 2021. doi: 10.1007/s00259-021-05242-1.
- [82] J. J. M. Van Griethuysen, A. Fedorov, C. Parmar, A. Hosny, N. Aucoin, V. Narayan, *et al.* Computational radiomics system to decode the radiographic phenotype. *Cancer Research* 77(21): e104–e107, 2017. doi: 10.1158/0008-5472.CAN-17-0339.
- [83] F. Binczyk, W. Prazuch, P. Bozek, and J. Polanska. Radiomics and artificial intelligence in lung cancer screening. *Translational Lung Cancer Research* 10(2): 1186–1199, 2021. doi: 10.21037/tlcr-20-708.
- [84] M. E. Mayerhoefer, A. Materka, G. Langs, I. Häggström, P. Szczypiński, P. Gibbs, *et al.* Introduction to radiomics. *Journal of Nuclear Medicine* 61(4): 488–495, 2020. doi: 10.2967/JNUMED.118.222893.
- [85] S. Ha, H. Choi, J. C. Paeng, and G. J. Cheon. Radiomics in Oncological PET / CT : a Methodological Overview. *Springer January*(53): 14–29, 2019.
- [86] S. Rizzo, F. Botta, S. Raimondi, D. Origgi, C. Fanciullo, A. G. Morganti, *et al.* Radiomics: the facts and the challenges of image analysis. *European Radiology Experimental* 2(36): 1–8, 2018. doi: 10.1186/s41747-018-0068-z.
- [87] K. Wei, S. Huifang, G. Zhou, R. Zhang, P. Cai, Y. Fan, *et al.* Potential Application of Radiomics for Differentiating Solitary Pulmonary Nodules. *OMICS J Radiol.* 5(2): 1–11, 2017. doi: 10.4172/2167-7964.1000218.Potential.
- [88] S. Hawkins, H. Wang, Y. Liu, A. Garcia, O. Stringfield, H. Krewer, *et al.* Predicting Malignant Nodules from Screening CT Scans. *Journal of Thoracic Oncology* 11(12): 2120–2128, 2016. doi: 10.1016/j.jtho.2016.07.002.
- [89] J. Ma, Q. Wang, Y. Ren, H. Hu, and J. Zhao. Automatic lung nodule classification with radiomics approach, in *Medical Imaging 2016: PACS and Imaging Informatics: Next Generation and Innovations*, 2016, 9789(9789906), p.: 978906-1-978906–6. doi: 10.1117/12.2220768.
- [90] J. Wang, X. Liu, D. Dong, J. Song, M. Xu, Y. Zang, *et al.* Prediction of malignant and benign of lung tumor using a quantitative radiomic method. *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBS 2016-October*: 1272–1275, 2016. doi: 10.1109/EMBC.2016.7590938.
- [91] D. Kumar, A. G. Chung, M. J. Shaifee, F. Khalvati, M. A. Haider, and A. Wong. Discovery radiomics for pathologically-proven computed tomography lung cancer prediction. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 10317 LNCS: 54–62, 2017. doi: 10.1007/978-3-319-59876-5\_7.
- [92] W. Choi, J. H. Oh, S. Riyahi, C. J. Liu, F. Jiang, W. Chen, *et al.* Radiomics analysis

- of pulmonary nodules in low-dose CT for early detection of lung cancer. *Medical Physics* 45(4): 1537–1549, 2018. doi: 10.1002/mp.12820.
- [93] C. H. Chen, C. K. Chang, C. Y. Tu, W. C. Liao, B. R. Wu, K. T. Chou, *et al.* Radiomic features analysis in computed tomography images of lung nodule classification. *PLOS ONE* 13(2): 1–13, 2018. doi: 10.1371/journal.pone.0192002.
- [94] L. Mao, H. Chen, M. Liang, K. Li, J. Gao, P. Qin, *et al.* Quantitative radiomic model for predicting malignancy of small solid pulmonary nodules detected by low-dose CT screening. *Quantitative Imaging in Medicine and Surgery* 9(2): 263–272, Feb. 2019. doi: 10.21037/qims.2019.02.02.
- [95] N. Amador Legón and M. Pérez Díaz. Use of fractals in determining the malignancy degree of lung nodules, *XX SIE, IV International Conference of UCLV*, p.: 1–68, 2023.
- [96] J. Tugwell-Allsup, B. W. Owen, and A. England. Low-dose chest CT and the impact on nodule visibility. *Radiography* 27(1): 24–30, 2021. doi: 10.1016/j.radi.2020.05.004.
- [97] T. Kurita, N. Otsu, and N. Abdelmalek. Maximum Likelihood Thresholding Based on Population Mixture Models. *Pattern Recognition* 25(10): 1231–1240, 1992.
- [98] J. Kittler and J. Illingworth. Minimum error thresholding. *Pattern Recognition* 19(1): 41–47, 1986. doi: 10.1016/0031-3203(86)90030-0.
- [99] Z. Yaniv, B. C. Lowekamp, H. J. Johnson, and R. Beare. SimpleITK Image-Analysis Notebooks: a Collaborative Environment for Education and Reproducible Research. *Journal of Digital Imaging* 31(3): 290–303, 2018. doi: 10.1007/s10278-017-0037-8.
- [100] M. Astaraki, G. Yang, Y. Zakko, I. Toma-Dasu, Ö. Smedby, and C. Wang. A Comparative Study of Radiomics and Deep-Learning Based Methods for Pulmonary Nodule Malignancy Prediction in Low Dose CT Images. *Frontiers in Oncology* 11(December): 1–12, 2021. doi: 10.3389/fonc.2021.737368.
- [101] P. Sedgwick. Pearson’s correlation coefficient. *BMJ (Online)* 345(7864): 1–2, 2012. doi: 10.1136/bmj.e4483.
- [102] K. Mei, M. Tan, Z. Yang, and S. Shi. Modeling of Feature Selection Based on Random Forest Algorithm and Pearson Correlation Coefficient. *Journal of Physics: Conference Series* 2219(1): 1–9, 2022. doi: 10.1088/1742-6596/2219/1/012046.
- [103] P. Chen, F. Li, and C. Wu. Research on Intrusion Detection Method Based on Pearson Correlation Coefficient Feature Selection Algorithm. *Journal of Physics: Conference Series* 1757(1): 1–10, 2021. doi: 10.1088/1742-6596/1757/1/012054.
- [104] F. Fernandez Castro, M. Perez Diaz, and R. Orozco Morales. Detection of lung nodules using support vector machine, *XX SIE, IV International Conference of UCLV*, p.: 1–45, 2022.
- [105] J. D. Kelleher and B. Tierney. Ciencia de datos, in *Ciencia de datos*, 2021, p.: 102.
- [106] W. Cao, R. Wu, G. Cao, and Z. He. *A Comprehensive Review of Computer-Aided Diagnosis of Pulmonary Nodules Based on Computed Tomography Scans*. 8, 2020.

- doi: 10.1109/ACCESS.2020.3018666.
- [107] G. Thibault, B. Fertil, and C. Navarro. Texture indexes and gray level size zone matrix: application to cell nuclei classification in Proceedings of the Pattern Recognition and Information Processing, in *International Conference on Pattern Recognition and Information Processing (PRIP)*, 2009, p.: 140–145.
- [108] M. Hossin and M. N. Sulaiman. A review on evaluation metrics for Data Classification Evaluations. *International Journal of Data Mining & Knowledge Management Process (IJDMP)* 5(2): 1–11, 2015.
- [109] Y. Balagurunathan, M. B. Schabath, H. Wang, Y. Liu, and R. J. Gillies. Quantitative Imaging features Improve Discrimination of Malignancy in Pulmonary nodules. *Scientific Reports* 9(May): 1–14, 2019. doi: 10.1038/s41598-019-44562-z.
- [110] K. Mehta, A. Jain, J. Mangalagiri, S. Menon, and P. Nguyen. Lung Nodule Classification Using Biomarkers, Volumetric Radiomics, and 3D CNNs. *Journal of Digital Imaging*(0123456789)2021. doi: 10.1007/s10278-020-00417-y.
- [111] Y. Xu, L. Lu, W. Lian, L. H. Schwartz, Z. Yang, Y. Xu, *et al.* Application of Radiomics in Predicting the Malignancy of Pulmonary Nodules. *American Journal of Roentgenology Application* 213(December): 1–8, 2019.
- [112] Y. Zhou, X. Xu, L. Song, C. Wang, J. Guo, Z. Yi, *et al.* The application of artificial intelligence and radiomics in lung cancer. *Precision Clinical Medicine* 3(May): 214–227, 2020. doi: 10.1093/pcmedi/pbaa028.
- [113] M. Avanzo, J. Stancanella, G. Pirrone, and G. Sartor. Radiomics and deep learning in lung cancer. *Springer* 196(May): 879–887, 2020. doi: 10.1007/s00066-020-01625-9.
- [114] D. Ardila, A. P. Kiraly, S. Bharadwaj, B. Choi, J. J. Reicher, L. Peng, *et al.* End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nature Medicine* 25(June): 954–961, 2019. doi: 10.1038/s41591-019-0447-x.
- [115] G. Varoquaux and V. Cheplygina. Machine learning for medical imaging: methodological failures and recommendations for the future. *Digital Medicine* 5(1): 1–8, 2022. doi: 10.1038/s41746-022-00592-y.
- [116] M. Fatan Serj, B. Lavi, G. Hoff, and D. Puig Valls. A Deep Convolutional Neural Network for Lung Cancer Diagnostic, *Computer Science*. p.: 1–10, 2018.
- [117] J. R. Mayo, J. Aldrich, and N. L. Müller. Radiation exposure at chest CT: A statement of the fleischner society. *Radiology* 228(1): 15–21, 2003. doi: 10.1148/radiol.2281020874.
- [118] C. H. McCollough, A. N. Primak, N. Braun, J. Kofler, L. Yu, and J. Christner. Strategies for Reducing Radiation Dose in CT. *Radiologic Clinics of North America* 47(1): 27–40, 2009. doi: 10.1016/j.rcl.2008.10.006.
- [119] L. J. Wilson and W. D. Newhauser. Justification and optimization of radiation exposures: a new framework to aggregate arbitrary detriments and benefits. *Radiation*

---

and Environmental Biophysics 59(3): 389–405, 2020. doi: 10.1007/s00411-020-00855-w.

- [120] J. R. Cooper and International Commission on Radiological Protection. Radiation protection principles.. Journal of radiological protection: official journal of the Society for Radiological Protection 32(1)2012. doi: 10.1088/0952-4746/32/1/N81.

**ANEXOS****Anexo 1 Código para visualizar el escaneo de las anotaciones de un paciente a partir de Pylide**

```
#Introducir el paciente que se está estudiando
pid = 'LIDC-IDRI-0086'
# Obtener el primer escaneo de tomografía computarizada (CT) del paciente seleccionado
scan = pl.query(pl.Scan).filter(pl.Scan.patient_id == pid).first()
# Obtener todas las anotaciones de nódulos del escaneo agrupadas
nods = scan.cluster_annotatations()
# Visualizar el escaneo con las anotaciones
scan.visualize(annotation_groups=nods)
```

**Anexo 2 Código para la visualización de las características de los nódulos (malignidad) a partir de Pylide**

```
#Introducir el paciente que se está estudiando
pid = 'LIDC-IDRI-0086'
# Obtener el primer escaneo de tomografía computarizada (CT) del paciente seleccionado
scan = pl.query(pl.Scan).filter(pl.Scan.patient_id == pid)
lista= scan.first()
scans = scan.first()
#Dentro del corchete se varía el número de la anotación que se desea visualizar
scans.annotations[1].visualize_in_scan()
```

**Anexo 3 Código para configurar el archivo YAML para PyRadiomics**

```
imageType: Original
binWidth: 25
resampledPixelSpacing:
featureClass:
shape: [ ]
firstorder: [ ]
```

```

glcm: [ ]
glrlm: [ ]
glszm: [ ]
gldm: [ ]

```

#### **Anexo 4 Código en Python para inicializar el extractor de características en PyRadiomics**

```

from radiomics import featureextractor

# Inicializando el extractor de características radiómicas
extractor = featureextractor.RadiomicsFeatureExtractor()
param_path = os.path.join(os.getcwd(), 'params.yaml')
extractor = featureextractor.RadiomicsFeatureExtractor(param_path)

```

#### **Anexo 5 Código en Python para salvar las características radiómicas en un archivo CSV**

```

features = extractor.execute(img_itk, mask_itk)
features_all = {}
for key, value in six.iteritems(features):
    if key.startswith('original') or key.startswith('wavelet') or \
    key.startswith('log'):
        features_all['Subject_ID'] = subject_name
        features_all['Subject_Label'] = subject_label
        features_all[key] = features[key]
        print("\t%s: %s" % (key, value))

df = pd.DataFrame(data=features_all, index=[ind])
if ind == 0:
    df.to_csv(radiomic_path_write, mode='a')
else:
    df.to_csv(radiomic_path_write, header=None, mode='a')

```

**Anexo 6 Código en Python para la implementación de los modelos de SVM**

```

def svm_kernel(kernel, poly_degree, c_val, class_weight):
    clf = svm.SVC(kernel=kernel, degree=poly_degree,
                  gamma='scale', C=c_val, tol=1e-1,
                  class_weight=class_weight,
                  probability=True,
                  random_state=None,
                  max_iter=-1)
    return clf

df = pd.read_csv('features_selection_01.csv')

X_train, X_test, y_train, y_test = train_test_split(
    df.drop('Subject_Label', axis = 1), # características
    df['Subject_Label'], # etiquetas
    test_size = 0.2, # porcentaje para prueba
    random_state = 8
)
## Dividir los datos en conjuntos de entrenamiento y validación
# X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42)

# Llamada a la función para crear el modelo
modelo_svm1 = svm_kernel(kernel='linear', poly_degree=3, c_val=1, class_weight=None)
modelo_svm2 = svm_kernel(kernel='rbf', poly_degree=3, c_val=1, class_weight=None)
modelo_svm3 = svm_kernel(kernel='poly', poly_degree=3, c_val=1, class_weight=None)
modelo_svm4 = svm_kernel(kernel='poly', poly_degree=5, c_val=1, class_weight=None)
modelo_svm5 = svm_kernel(kernel='poly', poly_degree=7, c_val=1, class_weight=None)

# Entrenar el modelo
modelo_svm1.fit(X_train, y_train)
modelo_svm2.fit(X_train, y_train)

```

## ANEXOS

```
modelo_svm3.fit(X_train, y_train)
modelo_svm4.fit(X_train, y_train)
modelo_svm5.fit(X_train, y_train)

# Realizar predicciones en el conjunto de validación de cada modelo variando el modelo a
calificar
predicciones1 = modelo_svm1.predict(X_test)
conf_matrix1 = confusion_matrix(y_test, predicciones1)

# Evaluar la precisión del modelo en el conjunto de validación
precision1 = accuracy_score(y_test, predicciones1)

VN1 = conf_matrix1[0][0]
FP1 = conf_matrix1[0][1]
FN1 = conf_matrix1[1][0]
VP1 = conf_matrix1[1][1]

Sensibilidad1 = VP1/(VP1+FN1)
Prec1 = VP1/(VP1+FP1)
Especificidad1 = VN1/(VN1+FP1)
Exactitud1 = (VP1+VN1)/(VP1+FP1+VN1+FN1)

print("\n Exactitud del modelo en conjunto de validación 1:", precision1)
print(f'Confusion Matrix:\n{conf_matrix1}')
print("Precisión:", Prec1)
print("Sensibilidad:", Sensibilidad1)
print("Especificidad", Especificidad1)
print("Exactitud", Exactitud1)

# Configurar el estilo para la matriz de confusión
sns.set(font_scale=1.2)
```

```
sns.heatmap(conf_matrix1, annot=True, fmt='g', cmap='Blues', cbar=False,
            xticklabels=['Benigno', 'Maligno'],
            yticklabels=['Benigno', 'Maligno'])
```

```
# Añadir etiquetas y título
plt.xlabel('Predicciones')
plt.ylabel('Valores reales')
plt.title('Kernel Lineal')
```

```
# Mostrar la imagen
plt.show()
```

### **Anexo 7 Código en Python para la implementación de los modelos de RF**

```
def random_forest(n_estimators, criterion, max_depth, class_weight):
```

```
    clf = RandomForestClassifier(criterion = criterion,
                               n_estimators = n_estimators,
                               max_depth = max_depth,
                               class_weight = class_weight,
                               n_jobs = -1,
                               random_state = None)
```

```
    return clf
```

```
df = pd.read_csv('features_selection_01.csv')
```

```
X_train, X_test, y_train, y_test = train_test_split(
    df.drop('Subject_Label', axis=1), # características
    df['Subject_Label'], # etiquetas
    test_size=0.2, # porcentaje para prueba
    random_state= 10 )
```

```
## Dividir los datos en conjuntos de entrenamiento y validación
```

```
# X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42)
```

## ANEXOS

```
# Llamada a la función para crear el modelo
modelo_rf1 = random_forest(n_estimators=100, criterion='gini', max_depth=None,
class_weight=None)
modelo_rf2 = random_forest(n_estimators=50, criterion='gini', max_depth=None,
class_weight=None)
modelo_rf3 = random_forest(n_estimators=20, criterion='gini', max_depth=None,
class_weight=None)
modelo_rf4 = random_forest(n_estimators=100, criterion='gini', max_depth=10,
class_weight=None)
modelo_rf5 = random_forest(n_estimators=50, criterion='gini', max_depth=10,
class_weight=None)
modelo_rf6 = random_forest(n_estimators=20, criterion='gini', max_depth=10,
class_weight=None)

# Entrenar el modelo
modelo_rf1.fit(X_train, y_train)
modelo_rf2.fit(X_train, y_train)
modelo_rf3.fit(X_train, y_train)
modelo_rf4.fit(X_train, y_train)
modelo_rf5.fit(X_train, y_train)
modelo_rf6.fit(X_train, y_train)

# Realizar predicciones en el conjunto de validación variando el nombre del modelo a
analizar
predicciones1 = modelo_rf1.predict(X_test)
conf_matrix1 = confusion_matrix(y_test, predicciones1)

# Evaluar la precisión del modelo en el conjunto de validación
precision1 = accuracy_score(y_test, predicciones1)
```

```
VN1 = conf_matrix1[0][0]
FP1 = conf_matrix1[0][1]
FN1 = conf_matrix1[1][0]
VP1 = conf_matrix1[1][1]

Sensibilidad1 = VP1/(VP1+FN1)
Prec1 = VP1/(VP1+FP1)
Especificidad1 = VN1/(VN1+FP1)
Exactitud1 = (VP1+VN1)/(VP1+FP1+VN1+FN1)

print("\nPrecisión del modelo en conjunto de validación 1:", precision1)
print(f'Confusion Matrix:\n{conf_matrix1}')
print("Precisión:", Prec1)
print("Sensibilidad:", Sensibilidad1)
print("Especificidad", Especificidad1)
print("Exactitud", Exactitud1)

# Configurar el estilo de la matriz de confusión
sns.set(font_scale=1.2)
sns.heatmap(conf_matrix1, annot=True, fmt='g', cmap='Blues', cbar=False,
            xticklabels=['Benigno', 'Maligno'],
            yticklabels=['Benigno', 'Maligno'])

# Añadir etiquetas y título
plt.xlabel('Predicciones')
plt.ylabel('Valores reales')
plt.title('n_estimators = 100 max_depth = None')

# Mostrar la imagen
plt.show()
```

**Anexo 8 Lista de características radiómicas extraídas****Características radiómicas**

---

original\_shape\_Elongation  
original\_shape\_Flatness  
original\_shape\_LeastAxisLength  
original\_shape\_MajorAxisLength  
original\_shape\_Maximum2DDiameterColumn  
original\_shape\_Maximum2DDiameterRow  
original\_shape\_Maximum2DDiameterSlice  
original\_shape\_Maximum3DDiameter  
original\_shape\_MeshVolume  
original\_shape\_MinorAxisLength  
original\_shape\_Sphericity  
original\_shape\_SurfaceArea  
original\_shape\_SurfaceVolumeRatio  
original\_shape\_VoxelVolume  
original\_firstorder\_10Percentile  
original\_firstorder\_90Percentile  
original\_firstorder\_Energy  
original\_firstorder\_Entropy  
original\_firstorder\_InterquartileRange  
original\_firstorder\_Kurtosis  
original\_firstorder\_Maximum  
original\_firstorder\_MeanAbsoluteDeviation  
original\_firstorder\_Mean  
original\_firstorder\_Median  
original\_firstorder\_Minimum  
original\_firstorder\_Range  
original\_firstorder\_RobustMeanAbsoluteDeviation  
original\_firstorder\_RootMeanSquared  
original\_firstorder\_Skewness  
original\_firstorder\_TotalEnergy  
original\_firstorder\_Uniformity  
original\_firstorder\_Variance  
original\_glcm\_Autocorrelation  
original\_glcm\_ClusterProminence  
original\_glcm\_ClusterShade  
original\_glcm\_ClusterTendency  
original\_glcm\_Contrast  
original\_glcm\_Correlation  
original\_glcm\_DifferenceAverage  
original\_glcm\_DifferenceEntropy  
original\_glcm\_DifferenceVariance  
original\_glcm\_Id  
original\_glcm\_Idm  
original\_glcm\_Idmn

ANEXOS

original\_glcm\_Idn  
original\_glcm\_Imc1  
original\_glcm\_Imc2  
original\_glcm\_InverseVariance  
original\_glcm\_JointAverage  
original\_glcm\_JointEnergy  
original\_glcm\_JointEntropy  
original\_glcm\_MCC  
original\_glcm\_MaximumProbability  
original\_glcm\_SumAverage  
original\_glcm\_SumEntropy  
original\_glcm\_SumSquares  
original\_glrlm\_GrayLevelNonUniformity  
original\_glrlm\_GrayLevelNonUniformityNormalized  
original\_glrlm\_GrayLevelVariance  
original\_glrlm\_HighGrayLevelRunEmphasis  
original\_glrlm\_LongRunEmphasis  
original\_glrlm\_LongRunHighGrayLevelEmphasis  
original\_glrlm\_LongRunLowGrayLevelEmphasis  
original\_glrlm\_LowGrayLevelRunEmphasis  
original\_glrlm\_RunEntropy  
original\_glrlm\_RunLengthNonUniformity  
original\_glrlm\_RunLengthNonUniformityNormalized  
original\_glrlm\_RunPercentage  
original\_glrlm\_RunVariance  
original\_glrlm\_ShortRunEmphasis  
original\_glrlm\_ShortRunHighGrayLevelEmphasis  
original\_glrlm\_ShortRunLowGrayLevelEmphasis  
original\_glszm\_GrayLevelNonUniformity  
original\_glszm\_GrayLevelNonUniformityNormalized  
original\_glszm\_GrayLevelVariance  
original\_glszm\_HighGrayLevelZoneEmphasis  
original\_glszm\_LargeAreaEmphasis  
original\_glszm\_LargeAreaHighGrayLevelEmphasis  
original\_glszm\_LargeAreaLowGrayLevelEmphasis  
original\_glszm\_LowGrayLevelZoneEmphasis  
original\_glszm\_SizeZoneNonUniformity  
original\_glszm\_SizeZoneNonUniformityNormalized  
original\_glszm\_SmallAreaEmphasis  
original\_glszm\_SmallAreaHighGrayLevelEmphasis  
original\_glszm\_SmallAreaLowGrayLevelEmphasis  
original\_glszm\_ZoneEntropy  
original\_glszm\_ZonePercentage  
original\_glszm\_ZoneVariance  
original\_gldm\_DependenceEntropy  
original\_gldm\_DependenceNonUniformity  
original\_gldm\_DependenceNonUniformityNormalized

ANEXOS

original\_gldm\_DependenceVariance  
original\_gldm\_GrayLevelNonUniformity  
original\_gldm\_GrayLevelVariance  
original\_gldm\_HighGrayLevelEmphasis  
original\_gldm\_LargeDependenceEmphasis  
original\_gldm\_LargeDependenceHighGrayLevelEmp-  
phasis  
original\_gldm\_LargeDependenceLowGrayLevelEmp-  
phasis  
original\_gldm\_LowGrayLevelEmphasis  
original\_gldm\_SmallDependenceEmphasis  
original\_gldm\_SmallDependenceHighGrayLevelEmp-  
phasis  
original\_gldm\_SmallDependenceLowGrayLevelEmp-  
phasis

---