

DETERMINACIÓN DEL GRADO SIMILITUD DE PATRONES  
PARENQUIMATOSOS EXTRAÍDOS DE IMÁGENES SINTÉTICAS CON  
IMÁGENES DE MAMOGRAFÍAS REALES

JUAN CARLOS PADILLA GARNICA

CIENCIA DE MATERIALES BIOLÓGICOS Y SEMICONDUCTORES (CIMBIOS)  
UNIVERSIDAD INDUSTRIAL DE SANTANDER  
FACULTAD DE CIENCIAS  
ESCUELA DE FISICA  
BUCARAMANGA  
2022

DETERMINACIÓN DEL GRADO SIMILITUD DE PATRONES  
PARENQUIMATOSOS EXTRAÍDOS DE IMÁGENES SINTÉTICAS CON  
IMÁGENES DE MAMOGRAFÍAS REALES

AUTOR:

JUAN CARLOS PADILLA GARNICA

TRABAJO DE GRADO PARA OPTAR POR EL TÍTULO DE FÍSICA

DIRECTORES:

DAVI ALEJANDRO MIRANDA MERCADO, PhD

SAID D. PERTUZ ARROYO, PhD

CIENCIA DE MATERIALES BIOLÓGICOS Y SEMICONDUCTORES (CIMBIOS)

UNIVERSIDAD INDUSTRIAL DE SANTANDER

FACULTAD DE CIENCIAS

ESCUELA DE FISICA

BUCARAMANGA

2022

El poder no se determina por tu tamaño, sino por el tamaño de tu corazón y tus sueños. *“Eiichiro Oda”*

# Agradecimientos

A mis abuelos Lucas Padilla, Teresa Garnica y a mi madre Patricia, todo lo que soy y todo lo que puedo llegar a ser, es gracias a ellos.

A mis amigos Johan y Sebastián por su amistad y apoyo incondicional durante toda mi vida.

A Angélica por su compañía y motivación durante mi carrera universitaria.

A la Universidad Industrial de Santander, por el desarrollo integral recibido y permitir mi formación como físico.

A cada uno de mis profesores de carrera, en especial a mi director David Alejandro Miranda gracias por sus enseñanzas y transmitir el entusiasmo y amor por la ciencia.

A mis compañeros de carrera, por su compañía y discusiones académicas interminables.

## CONTENIDO

	<b>pág.</b>
<b>1. INTRODUCCIÓN</b>	<b>14</b>
<b>2. MARCO TEÓRICO</b>	<b>17</b>
2.1. Anatomía de la mama	17
2.2. Mamografía	17
2.3. OpenVCT	18
2.4. OpenBreast	20
2.4.1. Estadística descriptiva e inferencial	24
2.4.2. Pruebas estadísticas	30
<b>3. Metodología</b>	<b>36</b>
<b>4. RESULTADOS Y DISCUSIÓN</b>	<b>42</b>
4.1. Prueba estadística: media muestral	45
4.2. Análisis del valor p para la media muestral	47
4.3. Prueba estadística: coeficiente de Correlación	51
4.4. Análisis del valor p para el coeficiente de correlación	52
4.5. Prueba estadística: distribución acumulativa	61
4.6. Análisis del valor p para distribución acumulativa	62
4.7. Determinación del grado de similitud	65
<b>5. CONCLUSIONES</b>	<b>69</b>
<b>BIBLIOGRAFÍA</b>	<b>71</b>

## LISTA DE FIGURAS

	<b>pág.</b>
Figura 1. Framework del programa OpenVct que simula los pasos para la generación y acumulación de mamografías sintéticas	19
Figura 2. Ejemplo de segmentación mamaria obtenida con OpenBreast. La región de la mamografía para analizar esta delimitada por el contorno de color rojo.	21
Figura 3. Cuatro posibles formas de seleccionar el área dentro de la mamá para realizar la extracción de características.	22
Figura 4. Zonas de rechazo o no rechazo de la hipótesis nula delimitada por el nivel de significancia	28
Figura 5. Zonas de rechazo o no rechazo de la hipótesis nula delimitada por el nivel de significancia	30
Figura 6.	33
Figura 7. Las gráficas en color azul y rojo representan la función de distribución acumulativa para cada una de las muestras, el segmento de línea roja punteada corresponde al estadístico de Kolmogorov-Smirnov el cual calcula la diferencia máxima entre las dos distribuciones acumulativas[70]	35
Figura 8. Densidades de las mamografías reales y mamografías sintéticas.	37
Figura 9. Representación de la información obtenida de manera general para cada una de las mamografías reales y sintética una vez realizada la extracción de características utilizando los 32 métodos.	38
Figura 10. División de la base de datos de mamografías reales y mamografías sintéticas en tres grupos de estudio: grupo de Casos, de Controles y Base completa. A cada una de las mamografías presentes en los tres grupos de estudio se realizó la extracción de características con cada uno de los 32 métodos.	39

Figura 11. Representación de los valores cuantitativos obtenidos con los métodos de extracción de características y el análisis realizado sobre el conjunto de mamografías reales y mamografías sintéticas.	40
Figura 12. Conjunto de mamografías reales pertenecientes al grupo Casos. Las mamografías presentan una densidad del 32 %, 47 %, 38 %, 25 % y 26 %, de manera respectiva.	43
Figura 13. Conjunto de mamografías reales pertenecientes al grupo Controles. Las mamografías presentan una densidad del 16 %, 29 %, 15 %, 8 % y 48 %, de manera respectiva.	44
Figura 15. Conjunto de mamografías sintéticas pertenecientes al grupo Controles Simulados. Las mamografías presentan una densidad del 15 %, 30 %, 15 %, 10 % y 35 %, de manera respectiva.	44
Figura 14. Conjunto de mamografías sintéticas pertenecientes al grupo Casos Simulados. Las mamografías presentan una densidad del 30 %, 35 %, 35 %, 25 % y 25 %, de manera respectiva.	45
Figura 16. Valores p obtenidos de la prueba t del estudiante para los tres grupos de estudio.	47
Figura 17. Valores p obtenidos de la prueba U Mann-Whitney para los tres grupos de estudio.	48
Figura 18. Gráficas de Dispersión de los valores cuantitativos del método de extracción de características 4 para el grupo Base completa y el método de extracción de características 22 del grupo Controles	58
Figura 19. Gráfica de dispersión de los valores cuantitativos del método de extracción de características 23 para el grupo Base completa	58
Figura 20. Gráficas de Dispersión de los valores cuantitativos del método de extracción de características 2 para el grupo Casos y el método de extracción de características 29 del grupo Controles	59
Figura 21. Gráfica de Dispersión de los valores cuantitativos del método de extracción de características 26 para el grupo Controles	59

Figura 22. Valores estadísticos de la prueba de Kolmogorov-Smirnov para dos muestras obtenidos en cada uno de los grupos analizados	63
Figura 23. Valores p obtenidos de la prueba estadística Kolmogorov-Smirnov para cada uno de los grupos analizados	64
Figura 24. Gráficas de distribución acumulativa para los grupos de Casos con método de extracción 4 y 15, y grupo de Control con método de extracción 4	65

## LISTA DE TABLAS

	<b>pág.</b>
Tabla 1. Métodos utilizados en la extracción de características de los patrones parenquimatosos de las mamografías sintéticas y mamografías reales.	23
Tabla 2. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Pearson para los métodos utilizados en el grupo de Base completa. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Pearson son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.	52
Tabla 3. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Pearson para los métodos utilizados en el grupo de Casos. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Pearson son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.	53
Tabla 4. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Pearson para los métodos utilizados en el grupo de Controles. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Pearson son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.	54

Tabla 5.	Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Spearman para los métodos utilizados en el grupo de Base completa. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Spearman son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.	55
Tabla 6.	Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Spearman para los métodos utilizados en el grupo de Casos. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Spearman son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa	56
Tabla 7.	Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Spearman para los métodos utilizados en el grupo de Controles. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Spearman son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa	57
Tabla 8.	Métodos para los cuales existe una correlación en los valores de las muestras para la prueba de correlación de Pearson	60
Tabla 9.	Métodos para los cuales existe una correlación en los valores de las muestras para la prueba de correlación de Spearman	60
Tabla 10.	Métodos para los cuales no existe una correlación en los valores de las muestras para la prueba de correlación de Pearson	61
Tabla 11.	Métodos para los cuales no existe una correlación en los valores de las muestras para la prueba de correlación de Spearman	61

Tabla 12. Métodos de extracción de características para los cuales existe una similitud en los patrones parenquimatosos de las mamografías reales con las mamografías sintéticas para el análisis estadístico con distribución normal de media muestral (16), coeficiente de correlación (4)(3) (2) y distribución acumulativa (22) (24)(23) .

67

Tabla 13. Métodos de extracción de características para los cuales existe una similitud en los patrones parenquimatosos de las mamografías reales con las mamografías sintéticas para el análisis estadístico con distribución no normal de media muestral (17) y coeficiente de correlación (5)(6)(7).

68

## RESUMEN

**TÍTULO:** Determinación del grado de similitud de patrones parenquimatosos extraídos de imágenes sintéticas con imágenes de mamografías reales.

**AUTOR:** Juan Carlos Padilla Garnica.

**PALABRAS CLAVES:** Mamografías sintéticas, virtual clinical trials (VCT), análisis parenquimatoso, prueba de hipótesis, media muestral, coeficiente de correlación, distribución acumulativa.

### **DESCRIPCIÓN:**

Los ensayos clínicos virtuales permiten obtener modelos sintéticos de la anatomía de las mamas. Generar un modelo sintético de una mama permite realizar todo tipo de pruebas y ensayos clínicos que dan la posibilidad probar nuevos sistemas de adquisición de imágenes, métodos y tratamientos para el cáncer de mama superando los inconvenientes y limitaciones de los ensayos clínicos reales.

De acuerdo con la similitud de las características anatómicas del modelo simulado con el modelo real, los resultados obtenidos del ensayo clínico virtual se pueden aplicar a las prácticas clínicas reales con seguridad. Por lo cual, es importante determinar la manera de lograr comparar el modelo real con el modelo sintético. En el presente trabajo se generó una base de mamografías sintéticas con el programa OpenVct con la mayor similitud posible a una base de mamografías reales. Con el programa OpenBreast se llevó a cabo un análisis parenquimatoso sobre las dos bases de mamografías y con los datos cuantitativos obtenidos se realizó un análisis estadístico inferencial utilizando las pruebas estadísticas de la media muestral, coeficiente de correlación y distribución acumulativa. De los resultados obtenidos se puede determinar que no existe un grado de similitud entre las mamografías reales y las mamografías sintéticas simuladas.

## **ABSTRACT**

**TITLE:** Determining the degree of pattern similarity extracted from synthetic images with images of real mammograms.

**AUTHOR:** Juan Carlos Padilla Garnica.

**KEY WORDS:** Synthetic mammograms, virtual clinical trials (VCT), parenchymal analysis, hypothesis test, the sample mean, correlation coefficient, cumulative distribution.

**DESCRIPTION:** Virtual clinical trials can generate synthetic models of human breasts anatomically similarly to natural breasts. Developing a synthetic model of a breast allows for performing all kinds of tests and clinical trials that give the possibility to test new systems of image acquisition, methods, and treatments for breast cancer, overcoming the drawbacks and limitations of current clinical trials.

According to the similarity of the anatomical characteristics of the simulated model with the real breast, the results obtained from the virtual clinical trial can be safely applied to real clinical practices. Therefore, it is crucial to compare the natural (real) and synthetic models. In the present work, a base of synthetic mammograms was generated with the OpenVct program with the most significant possible similarity to a base of real mammograms. With the OpenBreast program, a parenchymal analysis was carried out on the two mammogram bases. With the quantitative data obtained, a statistical inferential analysis was made using the statistical tests of the sample mean correlation coefficient and cumulative distribution. From the results obtained, it can be determined that there is no similarity between actual mammograms and simulated synthetic mammograms.

## 1. INTRODUCCIÓN

El cáncer de mama es una de las principales causas de muerte por cáncer en las mujeres [1] [2] [3]. El cáncer de mama de acuerdo con su evolución se ha clasificado en diferentes etapas. La detección del cáncer en una fase temprana está relacionada de forma directa con una mayor probabilidad de sobrevivir a la enfermedad [2][4][5]. Datos reportados en Estados Unidos indican una supervivencia de las pacientes del 89 % cuando son diagnosticadas en fases tempranas de la enfermedad ( I Y II ) y tan solo un 35 % de supervivencia en las mujeres con diagnóstico en las últimas fases ( III y IV ) [6]. Por lo anterior, en la actualidad se han desarrollado diferentes métodos que permiten realizar una detección temprana del cáncer de mama y analizar el cáncer en cualquiera de sus diferentes fases, entre estos métodos se encuentran: imagen de microondas, la detección mediante células tumorales circulantes, imágenes por resonancia magnética, tomografía por emisión de positrones, termografía, imágenes por impedancia eléctrica, mamografía digital de campo completo, mamografía con pantalla de película, entre otros[7] [8][9].

Con el avance de la tecnología y de las investigaciones sobre el cáncer de mama, se desarrollan y actualizan de forma constante nuevos métodos que permiten un estudio más profundo del cáncer de mama y a su vez lograr detectar el cáncer en cualquiera de sus fases y de esta manera lograr un aumento en el número de mujeres que sobreviven al cáncer de mama.

Sin embargo, todo nuevo método, sistema de imágenes médicas y en general toda práctica nueva relacionada con el estudio, detección y tratamiento del cáncer, debe pasar por varios ensayos clínicos antes de validar su utilización en la medicina. Estos ensayos clínicos son costosos, con un tiempo de ejecución extenso y según la técnica utilizada, a veces implican radiación ionizante de manera repetida sobre los voluntarios del ensayo clínico, lo cual puede generar afectaciones en la salud de

los pacientes [10][11]. Como una posible solución a la problemática anterior se han desarrollado los ensayos clínicos virtuales los cuales permiten modelar la anatomía humana, la adquisición, visualización y procesamiento de imágenes, análisis e interpretación de imágenes[12][13]. Los ensayos clínicos virtuales permiten realizar una evaluación cuantitativa y validar nuevos métodos de obtención de imágenes médicas antes de la ejecución de ensayos clínicos reales [12].

Los ensayos clínicos virtuales en el campo de las imágenes de mama, implican simulaciones de la anatomía de la mama, las cuales se utilizan para producir imágenes simuladas con o sin lesiones. Los ensayos clínicos y virtuales utilizan la mamografía como la principal modalidad de obtención de imágenes para la detección del cáncer. La mamografía se basa en la transmisión de rayos X de baja energía a través de la mama y aprovecha los diferentes coeficientes de atenuación existente en los diversos tejidos mamarios como tejido glandular, grasa, tejido fibroso y tumores para generar una imagen de la estructura interna de la mama[14].

El modelo sintético de mama es parte fundamental en los ensayos clínicos virtuales, de acuerdo a la calidad del modelo sintético y su similitud respecto a las características anatómicas un modelo de mama real, los resultados obtenidos en el ensayo clínico se pueden aplicar de manera satisfactoria en las prácticas clínicas reales, por lo tanto es importante tener un método o una estrategia que permita identificar un modelo sintético adecuado, es decir, con características anatómicas similares a un modelo real.

Una posible estrategia que permite comparar las similitudes en las características anatómicas de los modelos reales con los modelos sintéticos se basa en el diagnóstico asistido por computadora. De manera general, el diagnóstico asistido por computadora permite mejorar la interpretación de las imágenes médicas lo cual conlleva a una evaluación del riesgo del cáncer de mama en cualquiera de sus fases [15][16][17]. Los análisis mamográficos computarizados se dividen en cuatro etapas:

i) Segmentación de la mama, ii) detección de la región de interés (ROI), iii) extracción de características, y iv) puntaje de riesgo [18]. La realización de estos pasos sobre una imagen mamográfica se denomina análisis parenquimatoso. El análisis parenquimatoso se puede emplear tanto en imágenes mamográficas reales, como en imágenes mamográficas simuladas. Realizando la secuencia de los cuatro pasos es posible extraer un gran y variado conjunto de características y cantidades cuantitativas pertenecientes a la imagen mamográfica.

Lo anterior genera la posibilidad de utilizar el análisis parenquimatoso como una estrategia que permite comparar las mamografías sintéticas con las mamografías reales. Es de esperarse que la similitud entre los dos modelos mamarios real y sintético, de como resultado análisis parenquimatosos similares. En el presente trabajo se aborda la dificultad de comparar los modelos anatómicos reales con los modelos anatómicos sintéticos. Para lo cual se realizó un análisis parenquimatoso sobre los modelos anatómicos reales y los modelos anatómicos simulados, de esta manera se determinó un grado de similitud existente entre los dos modelos.

En el contenido del trabajo se presenta, primero, algunas de las mamografías reales con sus características anatómicas utilizadas para generar con el mayor grado de similitud posible las mamografías sintéticas. Segundo, los análisis estadísticos correspondientes a las cantidades cuantitativas extraídas con los análisis parenquimatosos realizados sobre cada una de las mamografías reales y mamografías sintéticas. Los análisis estadísticos presentados corresponden a: media muestral, coeficiente de correlación y distribución acumulativa. Seguido a esto se presenta una discusión de los resultados obtenidos a partir de las pruebas estadísticas realizadas sobre el conjunto de datos cuantitativos obtenidos. Para finalizar, se resume y se presenta toda la información importante obtenida con la realización del presente trabajo.

## **2. MARCO TEÓRICO**

### **2.1. Anatomía de la mama**

Los senos femeninos son un órgano glandular ubicado en el pecho, este órgano está representado por dos grandes eminencias hemisféricas que contienen la glándula mamaria la cual tiene la posibilidad de secretar leche [19] [20]. La superficie de la mama es convexa y tiene en el centro una pequeña área levantada llamada pezón[19].

La estructura anatómica del seno de una mujer adulta está compuesta por tejido fibroso, tejido graso, tejido glandular, vasos sanguíneos, nervios y conductos. El tejido glandular de la mama consiste en numerosos lóbulos, los cuales a su vez se dividen en pequeños lobulillos, estos lobulillos están conectados entre sí por un conducto lactífero subdividido que converge en el pezón. El conducto lactífero es el responsable de transportar la leche materna durante el tiempo de la lactancia [20] [21] El tejido graso cubre la mama y se encuentra ubicado entre el tejido glandular, este tejido es abundante y es el responsable de determinar la forma y el tamaño de la mama [19] [21]. El tejido fibroso da el soporte necesario a la mama para que esta se mantenga en su lugar[22]. El peso y la dimensión de los senos difieren entre individuos y en diferentes períodos de la vida[21] [23].

### **2.2. Mamografía**

La mamografía es un procedimiento que permite la adquisición de imágenes del interior de la mama mediante el uso de rayos X [24] [25]. La mamografía es la técnica más utilizada para la detección y diagnóstico de cáncer de mama [26]. El objetivo principal de la mamografía es detectar el cáncer en su fase más temprana, en una fase presintomática. Se ha demostrado que la supervivencia de la mujer aumenta de manera considerable si la enfermedad se detecta en su fase temprana[27][28].

Por el contrario, cuando se desarrollan los síntomas, el cáncer se vuelve invasivo y la supervivencia disminuye de manera considerable [29][30]

La mamografía se origina debido a la interacción de los diferentes tejidos de la mama con los rayos X. Los rayos X es un tipo de radiación electromagnética de alta energía, con longitudes de onda que oscilan entre  $10^{-8}$  [m] y  $10^{-12}$  [m] y con frecuencias entre  $10^{16}$  [Hz] y  $10^{21}$  [Hz] [31][32]. En la medida que los rayos X se propagan y penetran en la mama se generan diferentes procesos físicos, como lo son: procesos de absorción y procesos de dispersión [32][33]. Estas interacciones físicas provocan que los rayos X se atenúen a lo largo de los diferentes caminos ópticos que atraviesan la estructura anatómica interna de la mama, los rayos atenuados son registrados por un receptor el cual forma la imagen [25][24].

### **2.3. OpenVCT**

OpenVCT es un software desarrollado para generar modelos de mamas con características anatómicas personalizables. El software permite la adquisición, visualización y procesamiento de las imágenes adquiridas del modelo simulado de la mama. OpenVCT es un framework que consta de software gráfico para diseñar una secuencia de pasos de procesamiento para el proceso en cadena de VCT; software de gestión que coordina la ejecución del proceso en cadena, manipula y recupera fantasmas e imágenes utilizando una base de datos relacional; y un servidor que ejecuta los pasos individuales del proceso de acumulación de pacientes virtuales utilizando software optimizado para GPU [34]. La figura (??) resume el framework utilizado por OpenVCT.

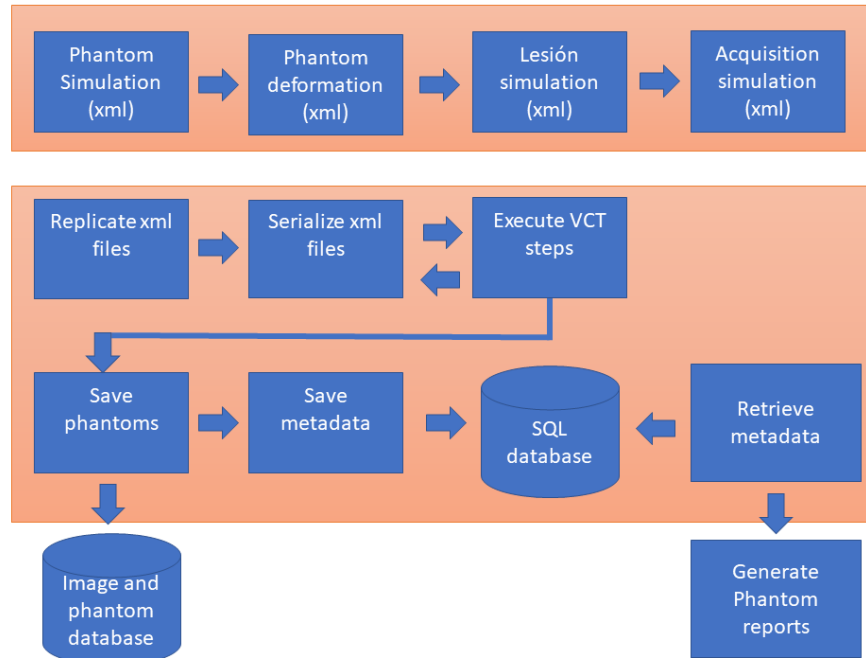


Figura 1. Framework del programa OpenVct que simula los pasos para la generación y acumulación de mamografías sintéticas

OpenVCT ejecuta cada paso del proceso de generar pacientes virtuales en el siguiente orden:

- Generador de modelos de mamaros: este software simula fantasmas mamaros antropomórficos que consisten en tejidos mamaros normales
- Deformador de modelos mamaros: aplicación que simula la compresión física de la mama mediante una técnica de mapeo 3D acelerada por GPU con mallas de compresión calculadas previamente.
- Inserción de lesiones: software opcional que inserta lesiones en el modelo deformado basándose en un archivo configurable que describe el tipo de lesión, la posición central, las dimensiones, la composición y los límites.
- Proyección de rayos X: una aplicación que simula imágenes mamográficas mediante la proyección de rayos X a través del fantasma de la mama. Supone un haz de rayos X poli energético sin dispersión y un modelo de detector ideal [34].

OpenVCT permite generar proyecciones de rayos X sobre la mama generando mamografías digitales de campo completo (FFDM) y tomosíntesis digital de la mama (DBT) utilizando el método de Siddon, el cual consiste en la representación de un conjunto de datos 3D como tres conjuntos de planos ortogonales. Para un rayo que viaja desde la fuente hasta un punto en el volumen, los vóxel a través de los cuales viaja el rayo se pueden determinar con la ayuda de estos planos. En la figura (??) se presenta el esquema en 2D del método Siddon. La longitud de la trayectoria radiológica RPL se calcula sumando la longitud recorrida por este rayo en cada vóxel, multiplicada por la densidad relativa de electrones del vóxel [35].

#### **2.4. OpenBreast**

La detección correcta de lesiones asintomáticas en mamografías permite detectar el cáncer de mama en las fases tempranas de su evolución permitiendo realizar un tratamiento oportuno de la enfermedad, aumentando la probabilidad del paciente de sobrevivir al cáncer de mama. La tarea del análisis de las imágenes mamografías recae en la gran mayoría de veces al radiólogo. El análisis y el consecuente dictamen médico de la mamografía puede ser acompañado con programas computacionales, los cuales están basados en algoritmos computarizados que realizan un análisis automático de imágenes [36],[37], la utilización de estos programas computacionales permite mejorar la detección de anomalías en la mama permitiendo al radiólogo generar una interpretación más acertada. Entre los diferentes programas computarizados para el análisis de imágenes mamográficas, se encuentra OpenBreast. OpenBreast es un software abierto basado en Matlab, su funcionamiento esta basado en cuatro pasos principales: segmentación mamaria, la detección de la región de interés (ROI), la extracción de características y la puntuación de riesgo[18]. A continuación, se realiza una breve descripción de cada función de OpenBreast.

- Segmentación mamaria: La segmentación mamaria tiene como objetivo separar y detectar de la mamografía la región mamaria. Para realizar esto, el programa debe seguir las siguientes tareas: detección del fondo, detección de

la pared torácica y detección del pezón[18]. En la imagen (2) se presenta el ejemplo de la segmentación mamaria para una mamografía digital de campo completo (FFDM).

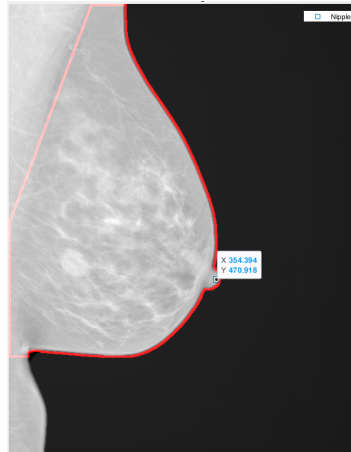


Figura 2. Ejemplo de segmentación mamaria obtenida con OpenBreast. La región de la mamografía para analizar esta delimitada por el contorno de color rojo.

- Detección de la región de interés (ROI): Una vez el programa selecciona y separa de la mamografía la región de la mama, se procede a detectar regiones en específico dentro de la mama para realizar la extracción de características. Existen diferentes maneras de seleccionar el área de interés dentro de la mama para realizar la extracción de características, entre las posibles regiones se encuentran: pecho completo, cuadrado más grande dentro del pecho, La región retroareolar (RA) y múltiples ROI siguiendo un muestreo basado en celosía[18].

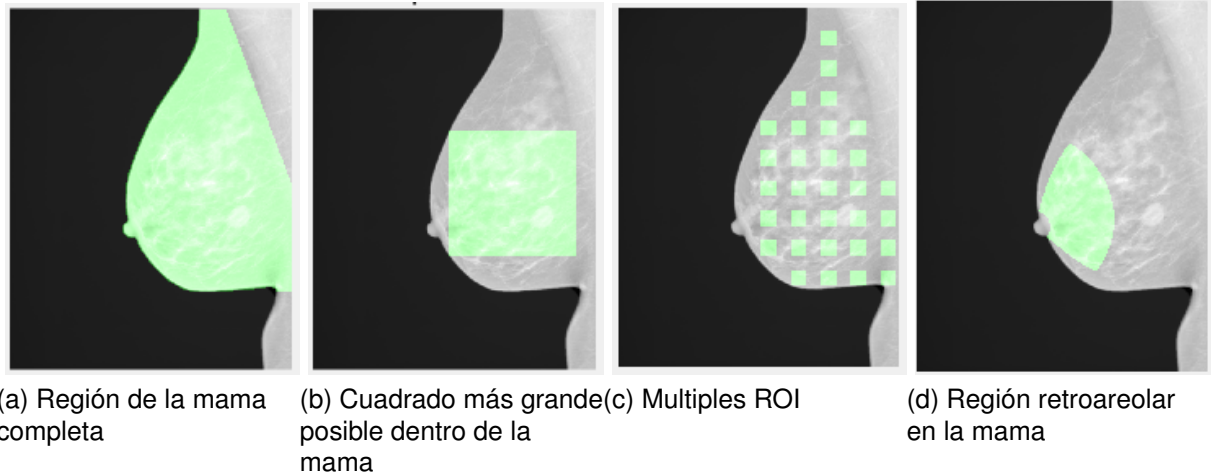


Figura 3. Cuatro posibles formas de seleccionar el área dentro de la mamá para realizar la extracción de características.

- Una vez seleccionada la región de interés (ROI) se procede a realizar la extracción de características, los diferentes métodos que realizan esta tarea permiten describir los diversos patrones de textura visual presentes en la mamografía en cantidades cuantitativas.

Puesto que la extracción de características es la etapa mas importante en el análisis computarizado de imágenes mamográficas, en el presente trabajo se utilizaron 32 métodos de extracción de características, los cuales se presentan en la tabla 1.

<b>Método</b>	<b>Nombre</b>	<b>Acrónimo</b>
1	Minimum gray-level value	STA
2	Maximum gray-level value	STA
3	Mean gray-level value	STA
4	Entropy	STA
5	Gray-level variance	STA
6	5-th percentile	STA
7	95-th percentile	STA
8	Balance 1	STA
9	Balance 2	STA
10	30-th percentile	STA
11	70-th percentile	STA
12	Skewness	STA
13	Kurtosis	STA
14	Gray-level range	STA
15	Energy	GLC
16	Correlation	GLC
17	Contrast	GLC
18	Homogeneity	GLC
19	Entropy	GLC
20	Short run emphasis	GLR
21	Long run emphasis	GLR
22	Gray-level non-uniformity	GLR
23	Run percentage	GLR
24	Run-length non-uniformity	GLR
25	Low gray-level run	GLR
26	High gray-level run	GLR
27	Gradient energy	SFA
28	Modified laplacian	SFA
29	Wavelet sum	GRA
30	Wavelet variance	GRA
31	Wavelet ratio	GRA
32	Gradient variance	SFA

Tabla 1. Métodos utilizados en la extracción de características de los patrones parenquimatosos de las mamografías sintéticas y mamografías reales.

Los métodos de extracción de características se pueden agrupar en diferentes categorías dependiendo de su estructura y manera de trabajar se clasifican en cinco grupos principales:

- **Características estadísticas (STA):** tienen como objetivo describir las

propiedades del histograma de intensidades de píxeles de nivel de gris[18][38][39].

- **Características de co-ocurrencia de nivel de gris (GLC):** tienen como objetivo describir la distribución estadística de las intensidades de imagen concurrentes en determinadas distancias y orientaciones de píxeles[18][40][41].
  - **Características de longitud de ejecución de nivel de gris (GLR):** tienen como objetivo describir la longitud y la distribución de valores de nivel de gris consecutivos en la imagen[18][42][43].
  - **Gradiente características basadas en características (GRA):** estiman el cambio de intensidad de píxeles en las direcciones horizontal y vertical de la imagen[18][40].
  - **Análisis de frecuencia espacial (SFA):** cálculo de descriptores basados en propiedades de textura en el dominio espacial, el dominio de frecuencia o ambos.[18][44][45].
- OpenBreast presenta la funcionalidad de utilizar un conjunto de imágenes de mamografías para entrenar un modelo que permite establecer una categoría de riesgo para cada mamografía de interés, estableciendo un alto o bajo riesgo de padecer cáncer de mama[18]. El presente trabajo de grado excluye de su estudio esta funcionalidad de OpenBreast.

**2.4.1. Estadística descriptiva e inferencial** En las investigaciones realizadas en los diferentes campos tanto científicos, sociales, económicos, etc. Es común obtener una gran cantidad de datos y variables a partir de una muestra pereciente a una población, los cuales con una sencilla observación es imposible extraer algún tipo de información relacionada con la muestra de estudio. Para ayudar a comprender mejor los estudios de investigación, se ha desarrollado a lo largo de los años la estadística. La estadística es una rama de las matemáticas la cual propicia el estudio, presentación y conclusiones de los datos obtenidos de una investigación. La estadística posee dos grandes ramas, la estadística descriptiva y la estadística

inferencial[46][47][48].

La estadística descriptiva se encarga de describir y resumir a través de diferentes técnicas la información contenida en el conjunto de datos extraídos de una muestra de estudio[49][47][50]. Por su parte la estadística inferencial cuenta con diferentes procedimientos que permiten: determinar parámetros de la muestra, extraer características de las distribuciones de los datos y realizar comparaciones formales entre dos o más grupos de muestras. La inferencia estadística permite extraer conclusiones sobre una población a través del análisis de la información de una muestra. Para realizar esto la inferencia estadística cuenta con dos métodos los cuales son: La prueba de hipótesis y los intervalos de confianza[49][51][52]. En el presente trabajo de grado se utiliza el método de prueba de hipótesis para realizar la comparación de manera formal y así determinar el grado de similitud entre las muestras de la base de mamografías real y la base de mamografías sintética.

**Prueba de hipótesis** La prueba de hipótesis es una manera formal de realizar y responder preguntas sobre las poblaciones involucradas en la investigación de acuerdo con las características obtenidas de cada una de las muestras extraídas de las poblaciones de estudio [53] [49][50]. La prueba de hipótesis permite comparar la tendencia central, la correlación, la distribución acumulativa y otros valores muestrales[50][54]. Una vez realizada la comparación y con la información obtenida de este proceso, es posible establecer conclusiones y hacer afirmaciones acerca de las muestras de estudio y sus poblaciones[50][46][47]. Para el desarrollo y ejecución de una prueba de hipótesis se deben considerar las características estadísticas (distribución normal o distribución no normal) de las muestras, así como clasificación del tipo de dato obtenido a partir del estudio de investigación. Los datos pueden ser categóricos o cuantitativos, a su vez pueden ser continuos, discretos, dependientes o independientes.

Una vez realizado el proceso anterior se establecen los métodos estadísticos ade-

cuados para desarrollar la prueba de hipótesis y establecer una conclusión. De acuerdo con el tipo de característica estadística obtenida de las muestras, se utiliza una prueba de hipótesis adecuada. Sin embargo, de manera general todas las pruebas de hipótesis comparten una misma estructura para su desarrollo. Los pasos a seguir para el desarrollo de una prueba de hipótesis son:

1. **Examinar los datos:** Una vez obtenida la información de la variable estadística de interés, se debe identificar la distribución de los datos. Los datos recolectados pueden presentar una distribución normal o una distribución que no es posible especificar, es decir una distribución no normal. Para el caso en el cual se presenta una distribución normal se utilizan pruebas de hipótesis paramétricas. Para el caso contrario, es decir una distribución no normal se utilizan pruebas de hipótesis no paramétricas. Además de identificar la distribución de los datos, es importante conocer el tipo de dato a trabajar. Entre los tipos mas comunes de datos están:
  - **Una muestra:** Este tipo de dato hace referencia a la obtención de varias medidas de una variable en una sola muestra, por ejemplo: la altura de los niños de tercer grado de primaria.
  - **Dos muestras emparejadas:** Se habla de muestras emparejadas cuando los datos obtenidos de una muestra están relacionados con los datos de la otra muestra, es decir, existe una correspondencia entre las mediciones de cada muestra, por ejemplo: Análisis de la presión arterial en un grupo de individuos antes y después de utilizar un medicamento.
  - **Dos muestras no emparejadas:** Los datos obtenidos de dos muestras independientes no están relacionados entre si, no existe ninguna correspondencia entre las mediciones de cada muestra, por ejemplo: Análisis del ritmo cardiaco de una población A frente a una población [53]
2. **Planteamiento de la hipótesis:** Para realizar y responder preguntas respecto a una variable de interés de una población se utilizan hipótesis. En una prueba

de hipótesis existen dos hipótesis las cuales deben ser mencionadas de manera explícita [54] [55]. La primera hipótesis y la cual se utiliza como punto de partida en una prueba de hipótesis, se define como hipótesis nula. La hipótesis nula afirma que no existe diferencia alguna en las variables poblacionales de las muestras de estudio. Por convención, se parte del hecho que la hipótesis nula es verdadera y a partir del desarrollo y análisis de la prueba de hipótesis es posible rechazar esta hipótesis nula [46][53]. La segunda hipótesis se define como hipótesis alternativa. La hipótesis alternativa afirma que existe una diferencia en las variables poblaciones de las muestras de estudio [53]. Con el desarrollo y análisis de la prueba de hipótesis es posible rechazar la hipótesis nula, por lo tanto, se acepta la hipótesis alternativa. Caso contrario, en el cual se acepte la hipótesis nula, lleva consigo el rechazo de la hipótesis alternativa. El rechazo de la hipótesis nula conlleva a que existe una diferencia estadísticamente significativa en el estudio realizado [46] [55] [53].

### **3. Elegir un grado de confianza:**

En toda investigación estadística existe la probabilidad de establecer conclusiones de los parámetros poblaciones por puro azar. Es posible elegir muestras de una determinada población que por puro azar no presenten ninguna diferencia entre ellas, lo cual conlleva a realizar una conclusión aparentemente verdadera, sin embargo, el resultado fue producto del azar y en realidad las muestras desde una perspectiva global exhiben una diferencia [51] [53]. De igual manera es posible obtener el resultado contrario, por puro azar se eligen muestras las cuales en sus parámetros poblaciones no se evidencia diferencia alguna, lo cual conlleva a realizar una conclusión errada acerca de las poblaciones.

En toda investigación es posible cometer alguno de los errores anteriores, puesto que es imposible analizar toda la población [53] y a su vez en los estudios estadísticos que involucren personas, se pueden presentar variaciones

más marcadas producto de la diversidad en los diferentes factores poblacionales.

Como existe la posibilidad de cometer algún tipo de los errores anteriores, los investigadores establecen un determinado valor de confianza, este valor puede ser del 90 %, 95 % ó 99 % de confianza en la respuesta. A su vez este valor de confianza se relaciona de manera directa con el nivel de significancia. El nivel de significancia permite al investigador aceptar o rechazar la hipótesis nula a través de una comparación con el valor p obtenido del estadístico de prueba. El nivel de significancia se representa con la letra griega  $\alpha$  y es igual a  $\alpha = 1 - (\text{valor de confianza expresado como probabilidad})$  [54] [56] , por ejemplo, un valor de confianza del 95 % expresado como probabilidad es representado como 0.95, por lo tanto el nivel de significancia para este valor de confianza es:  $\alpha = 1 - 0.95 = 0.05$ .

El nivel de significancia permite al investigador definir unas zonas de rechazo o no rechazo de la hipótesis nula y comparando estas zonas con el valor p obtenido del estadístico de prueba, aceptar o rechazar la hipótesis nula.

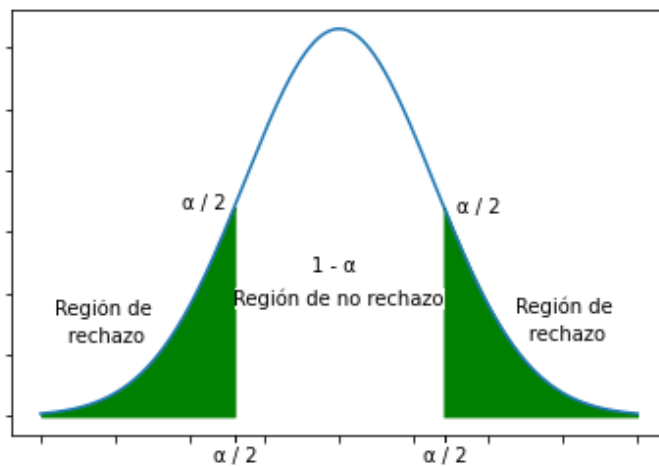


Figura 4. Zonas de rechazo o no rechazo de la hipótesis nula delimitada por el nivel de significancia

En la figura (4) se representan dos zonas, una zona de rechazo y una zona de no rechazo delimitadas por el nivel de significancia. Si el valor  $p$  del estadístico se ubica en la zona de no rechazo, el investigador acepta la hipótesis nula como verdadera, si el valor  $p$  se ubica en la zona de rechazo, el investigador rechaza la hipótesis nula y se acepta la hipótesis alternativa.

#### 4. **Calcular el estadístico de prueba y el valor $p$ :**

Una vez construida la hipótesis nula y la hipótesis alternativa de acuerdo con la variable de interés para el análisis estadístico de dos muestras, se procede a elegir el tipo de prueba de hipótesis de la cual se obtiene un estadístico de prueba que permite rechazar o aceptar la hipótesis nula. Existen diferentes tipos de prueba de hipótesis, cada uno de ellos diseñado para analizar las diferentes variables estadísticas presentes en las muestras, como lo son, su media muestral, la correlación entre las muestras, las distribuciones muestrales, proporción muestral, entre otros [50] [47]. Cada uno de estos tipos de prueba de hipótesis se elige de acuerdo con la finalidad del estudio estadístico de las dos muestras y con las características estadísticas de los datos provenientes de las dos muestras, es decir: Los datos son categóricos o cuantitativos, presentan una distribución normal o no normal, los datos provienen de muestras emparejadas o no emparejadas.

De manera general para todos los tipos de prueba de hipótesis, se calcula una estadística específica para cada una de las pruebas de hipótesis, esta estadística es obtenida a partir de los datos de las muestras. Esta estadística calculada puede tomar diferentes valores, sin embargo, va a tomar un valor específico de acuerdo con los datos extraídos de las muestras. Cada una de las pruebas de hipótesis siguen funciones de probabilidad tabuladas, por lo tanto, es posible obtener un estadístico de prueba el cual es una cantidad numérica y el denominado valor  $p$ . El valor  $p$  es un valor de probabilidad y está asociado de manera directa con el estadístico de prueba [51] [50] [53]. Para

rechazar o aceptar la hipótesis nula es posible utilizar el estadístico de prueba o el valor p. En el presente trabajo de grado se va a emplear el valor p en el desarrollo de la prueba de hipótesis [49] .

### 5. Aceptar o rechazar la hipótesis nula:

Con el uso de las zonas de rechazo o no rechazo construidas a partir del nivel de significancia elegido por el investigador y su comparación con el valor p obtenido de la prueba de hipótesis, es posible rechazar o aceptar la hipótesis nula. De manera tradicional se suele elegir un nivel de significancia de 0.05, si el valor de p es menor al valor del nivel de significancia, se considera que el análisis realizado sobre las muestras es estadísticamente significativo por lo cual se procede a rechazar la hipótesis nula y aceptar la hipótesis alternativa.

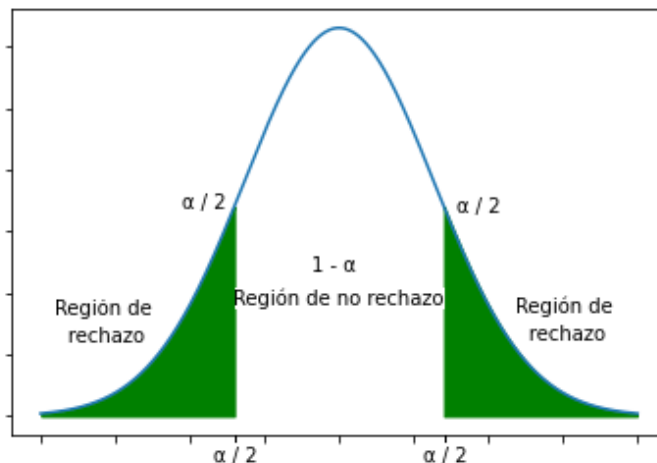


Figura 5. Zonas de rechazo o no rechazo de la hipótesis nula delimitada por el nivel de significancia

Si por el contrario el valor p obtenido es mayor al nivel de significancia, se afirma que existe suficiente evidencia para aceptar la hipótesis nula y rechazar la hipótesis alternativa.

**2.4.2. Pruebas estadísticas** La prueba de hipótesis permite analizar diferentes parámetros estadísticos de una población, como son: la media muestral, la corre-

lación, la distribución acumulada, entre otras. Para analizar cada uno de los parámetros estadísticos, existen diferentes pruebas estadísticas diseñadas para analizar cada uno de estos parámetros. Las pruebas estadísticas utilizadas en el presente trabajo se presentan a continuación:

**Prueba t del estudiante** El tipo más común de prueba de hipótesis paramétrica es la prueba t de Student. Las diferentes versiones de la prueba t de Student nos permiten probar si la media de la población de la que se extrajo nuestra muestra difiere de un valor esperado, o si las medias de dos poblaciones diferentes difieren [57] [58] [59]. El estadístico de prueba de la prueba t se conoce como valor t. Para utilizar cualquier tipo de prueba t de Student, se supone que las variables de población de interés tienen una distribución aproximadamente normal, puesto que la prueba utiliza la media, la desviación estándar, el tamaño de la muestra, y la media poblacional o el valor medio hipotético [60] [61]. En el presente trabajo, se utilizó una prueba t de Student para dos muestras de datos independientes. Para definir el estadístico de prueba t de dos muestras, es necesario hallar primero una desviación estándar combinada, es decir, una relación entre las desviaciones estándar de cada muestra y los tamaños de cada una, por lo tanto, la desviación estándar combinada es de la forma [57] [59]:

$$s = \sqrt{\frac{(n_1 - 1) s_1^2 + (n_2 - 1) s_2^2}{n_1 + n_2 - 2}} \quad (1)$$

Y el estadístico de prueba para la prueba t de Student es [57] [59]:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \quad (2)$$

Donde  $\bar{x}_1$  y  $\bar{x}_2$  corresponde al valor medio para el conjunto de datos de la muestra 1 y muestra dos de manera respectiva, s corresponde a la desviación estándar combinada y  $n_1$ ,  $n_2$  corresponde al tamaño de la muestra uno y muestra dos.

**Prueba de U Mann-Whitney** Para comparar dos muestras independientes y determinar si pertenecen a la misma población, primero se debe establecer su distribución, si al menos una muestra presenta en sus datos una distribución no normal, se aplica una prueba no paramétrica denominada prueba U de Mann-Whitney[57][58][62]. El siguiente paso es calcular el estadístico de prueba, se inicia agrupando las dos muestras de tamaño  $n_c$  y  $n_t$  para los datos de control y de prueba, en una sola muestra grande. Se ordenan los valores de menor a mayor de los datos en la muestra grande, cada valor ordenado va tener una posición denotada como  $(1, \dots, n_c + n_t)$ , luego de ello se calcula la suma de los rangos de cada muestra individual. Se utiliza  $R_t$  para representar la suma de los rangos de la muestra de prueba y  $R_c$  para la suma de los rangos de la muestra de control.  $R_c$  y  $R_t$  se utilizan para determinar los estadísticos de prueba [57][58][63], dados por:

$$U_t = n_t n_c + 0.5 n_t (n_t + 1) - R_t \quad (3)$$

$$U_c = n_t n_c + 0.5 n_c (n_c + 1) - R_c \quad (4)$$

Dependiendo del resultado obtenido para  $U_t$  y  $U_c$  se elige como estadístico de prueba aquel que presente el menor valor entre los dos, el valor arrojado por este estadístico está asociado con el valor p, el cual se puede determinar utilizando tablas estadísticas.

**Coefficiente de Correlación** El término correlación se usa en el lenguaje cotidiano para indicar asociación, este término se puede llevar al campo científico para proporcionar una medida de la fuerza y la dirección de la relación entre dos variables, es decir, la correlación se mide mediante una estadística llamada coeficiente de correlación el cual proporciona una medida de la fuerza y la dirección de la relación entre dos variables.

La correlación entre dos variables se puede visualizar utilizando un gráfico de dispersión. Las posibles correlaciones existentes entre dos variables son representadas en

los siguientes gráficos:

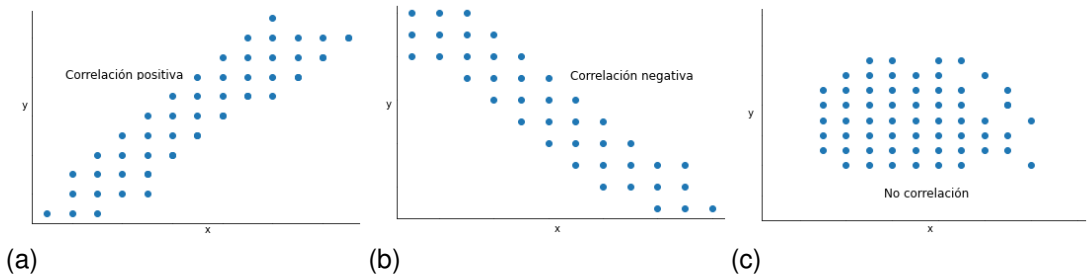


Figura 6

El gráfico A muestra una correlación positiva, lo que significa que a medida que aumenta la primera variable ( $x$ ), también aumenta la segunda variable ( $y$ ). El gráfico B muestra una correlación negativa (a veces llamada correlación inversa), lo que significa que a medida que  $x$  aumenta,  $y$  disminuye. El gráfico C muestra dos variables que no están correlacionadas entre sí de ninguna manera. El gráfico D muestra una relación curvilínea, sin embargo el coeficiente de correlación solo determina relaciones lineales, por lo cual una relación curvilínea u otro tipo de relación no se consideran correlaciones en términos estadísticos.

Para cuantificar la posible correlación entre las dos variables, se utiliza la letra  $r$  para denotar el coeficiente de correlación. El valor  $r$  puede estar en el rango de  $-1$  y  $+1$ . Cuando  $r$  es igual a  $-1$ , indica una relación negativa; cuando  $r$  es igual a  $0$ , indica que no hay relación; y cuando  $r$  es igual a  $+1$ , indica una relación positiva. Cuanto más cercano esté un coeficiente de correlación a  $0$ , más débil será la relación entre las dos variables.

De acuerdo con la distribución y tipo de cada variable es posible utilizar el coeficiente de correlación de Pearson o el coeficiente de correlación de Spearman.

### **Coeficiente de correlación de Pearson**

El coeficiente de correlación de Pearson se utiliza si las muestras analizadas pre-

sentan una distribución normal, el coeficiente de correlación de Pearson se define como[57][64][65][66]:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (5)$$

Donde  $x_i, y_i$  son las dos variables,  $i=1, \dots, n$  y  $n$  es el tamaño de la muestra. El coeficiente de correlación de Pearson proporciona un valor cuantitativo que establece la relación lineal entre las dos variables a comparar. El valor de  $r$  puede oscilar entre +1 (correlación positiva perfecta) y -1 (correlación negativa perfecta). Los valores de  $r$  cercanos a cero sugieren que no hay correlación lineal[64][65]. Con el uso de tablas es posible determinar el valor  $p$  asociado al estadístico calculado.

### **Coeficiente de correlación de Spearman**

El coeficiente de correlación de Spearman se utiliza si las muestras analizadas presentan una distribución no normal, el coeficiente de correlación de Spearman se define como:

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (6)$$

donde  $n$  es el número de muestras y  $\sum d_i^2$  es el suma de las diferencias al cuadrado en las puntuaciones de rango entre las dos muestras. El coeficiente de correlación de Spearman proporciona un valor cuantitativo que establece la relación lineal entre las dos variables a comparar. El valor de  $r$  puede oscilar entre +1 (correlación positiva perfecta) y -1 (correlación negativa perfecta). Los valores de  $r$  cercanos a cero sugieren que no hay correlación lineal[64][65].on el uso de tablas es posible determinar el valor  $p$  asociado al estadístico calculado.

**Distribución acumulativa:** Para determinar si dos distribuciones de datos provienen de una misma función de distribución acumulada, se utiliza la prueba de Kolmogorov-Smirnov para dos muestras. El estadístico de Kolmogorov-Smirnov de-

termina la diferencia máxima entre las funciones de distribución acumulativa de las dos muestras, de igual forma se calcula el valor p relacionado con el estadístico de la prueba utilizando tablas estadísticas. Una distribución acumulativa es de la forma[67][68]:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n 1_{[-\infty, x]}(X_i) \quad (7)$$

Donde  $1_{[-\infty, x]}(X_i)$  es la función indicadora, la cual es igual a 1 si  $X_i \leq x$  y 0 en caso contrario.

El estadístico Kolmogorov-Smirnov se define como[67][68][69]:

$$D_{n,m} = \sup_x |F_{1,n}(x) - F_{2,m}(x)| \quad (8)$$

Es posible visualizar de manera gráfica el estadístico de Kolmogorov-Smirnov como:

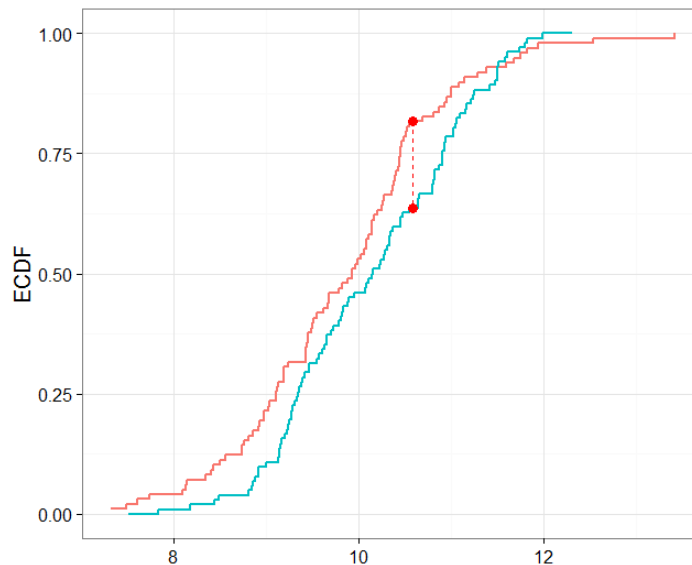


Figura 7. Las gráficas en color azul y rojo representan la función de distribución acumulativa para cada una de las muestras, el segmento de línea roja punteada corresponde al estadístico de Kolmogorov-Smirnov el cual calcula la diferencia máxima entre las dos distribuciones acumulativas[70]

### 3. Metodología

La base de datos de mamografías utilizada en el presente trabajo de grado corresponde a una base de datos obtenida del programa de cribado del Hospital Universitario de Tampere de Finlandia. Esta base de datos cuenta con 118 imágenes mamográficas divididas en dos grupos de 59 mamografías cada uno. Las mamografías están divididas en grupos puesto que un grupo presenta mamografías en las cuales las pacientes fueron detectadas con cáncer de mama. Este grupo fue denominado grupo de Casos. El otro grupo denominado como grupo de Controles, almacena 59 mamografías en las cuales las pacientes no presentaban signos de cáncer. Las imágenes mamográficas fueron obtenidas utilizando un sistema de imágenes Philips Healthcare. Este sistema de adquisición de imágenes mamográficas permite almacenar las mamografías en formato DICOM (.DCM). El formato DICOM es ampliamente utilizado en la medicina ya que permite almacenar un conjunto de metadatos los cuales brindan información mas detallada sobre la imagen mamográfica, como lo es: densidad de la mama, fuerza de compresión empleada, ancho y alto de la imagen, entre otros.

La información obtenida en cada una de las mamografías es importante ya que uno de los objetivos del presente trabajo de grado es generar una base de datos sintética con la mayor similitud posible a una base de mamografías reales, para lograr lo anterior es necesario conocer las características anatómicas de cada una de las mamografías presentes en la base de datos. Para acceder a la información contenida en cada imagen mamográfica se utilizó el programa Matlab y con ayuda de la función (.Dicom) fue posible realizar la extracción de la información más relevante de cada mamografía. La información obtenida de cada mamografía se registró en una tabla para facilitar la visualización de los datos. En la tabla se registró la información de las características anatómicas de cada una de las mamografías presentes en los grupos de Casos y Controles. En la gráfica 8a se presentan las densidades de cada

una de las mamografías reales.

Con el uso del programa OpenVct fue posible generar una base de mamografías sintéticas con la mayor similitud posible a la base de mamografías reales. Generando una mamografía sintética por cada una de las mamografías reales. De igual manera se conformó un grupo de mamografías sintéticas denominado Controles simulados y otro grupo de mamografías denominado Casos simulados. Se simularon en total 118 mamografías.

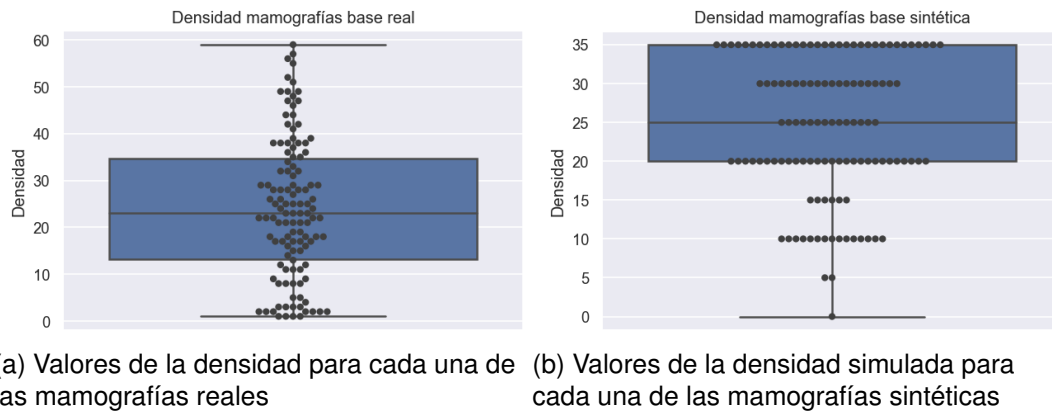


Figura 8. Densidades de las mamografías reales y mamografías sintéticas.

Obtenidas las bases de mamografías sintética y real, se utilizó el programa Open-Breast para realizar la extracción de características a cada una de las mamografías presentes en la base real y sintética. Para lograr la extracción de características sobre los patrones parenquimatosos en cada mamografía se realizaron los siguientes pasos: segmentación de la mama, elección de la región de interés (ROI) y la extracción de características.

Es importante recordar que la extracción de característica permite obtener una cantidad cuantitativa por cada uno de los métodos utilizados. Realizada la extracción de características a todas las imágenes mamográficas reales y sintéticas, se registró la información obtenida en una matriz de 118 filas por 32 columnas. Cada una de

las filas representa una mamografía presente en las bases de datos real y sintética. Las columnas representan la cantidad cuantitativa extraída de cada mamografía con cada uno de los métodos utilizados. En el presente trabajo de grado se utilizaron 32 métodos de extracción de características.

Para ilustrar el proceso anterior de obtención de los valores cuantitativos se presenta la siguiente figura:

	Método 1	Método 2	Método.....	Método 32
<b>Mamografía real 1</b>	xxxx	xxx	....	xxxx
....	....	....	....	.....
<b>Mamografía real 118</b>	xxxx	xxxx	....	xxxx
<b>Mamografía sintética 1</b>	xxxx	xxxx	....	xxxx
....	.....	.....	....	.....
<b>Mamografía sintética 118</b>	xxxx	xxxx	....	xxxx

Figura 9. Representación de la información obtenida de manera general para cada una de las mamografías reales y sintética una vez realizada la extracción de características utilizando los 32 métodos.

En cada una de las 118 mamografías reales y las 118 mamografías sintéticas se realizó la extracción de características utilizando 32 métodos. Cada uno de los métodos permitió obtener un valor cuantitativo.

Para analizar los datos obtenidos del proceso de extracción de características realizados a las mamografías presentes en las bases de datos real y sintética, se dividió el análisis en tres partes. La primera parte consistió en comparar el grupo de 59 mamografías reales denominado Controles con el grupo de 59 mamografías simuladas denominado Controles simulados. La segunda parte consistió en analizar el grupo de 59 mamografías reales denominados Casos con el grupo de 59 mamografías simuladas denominado Casos simulados. Para la tercera parte se analizó por completo la base de 118 mamografías real contra la base completa de 118 mamografías simuladas. La gráfica 10 permite visualizar los tres grupos de información analizados. A cada uno de los grupos se realizó un igual análisis estadístico.

Casos	Método 1	Método 2	Método.....	Método 32	Controles	Método 1	Método 2	Método.....	Método 32
Mamografía real 1	XXXXX	XXXX	....	XXXX	Mamografía real 1	XXXXX	XXXX	....	XXXX
Mamografía sintética 1	XXXXX	XXXX	....	XXXX	Mamografía sintética 1	XXXXX	XXXX	....	XXXX
.....	.....	.....	....	....	.....	.....	.....	....	....
.....	....	.....	....	....	.....	....	.....	....	....
Mamografía real 59	XXXXX	XXXX	....	XXXX	Mamografía real 59	XXXXX	XXXX	....	XXXX
Mamografía sintética 59	XXXXX	XXXX	....	XXXX	Mamografía sintética 59	XXXXX	XXXX	....	XXXX

(a) Grupo de mamografías reales y sintéticas pertenecientes al grupo de análisis Casos.

(b) Grupo de mamografías reales y sintéticas pertenecientes al grupo de análisis Controles.

Base completa	Método 1	Método 2	Método.....	Método 32
Mamografía real 1	XXXXX	XXXX	....	XXXX
Mamografía sintética 1	XXXXX	XXXX	....	XXXX
.....	.....	.....	....	....
.....	....	.....	....	....
Mamografía real 118	XXXXX	XXXX	....	XXXX
Mamografía sintética 118	XXXXX	XXXX	....	XXXX

(c) Grupo de mamografías reales y sintéticas pertenecientes al grupo de análisis Base completa.

Figura 10. División de la base de datos de mamografías reales y mamografías sintéticas en tres grupos de estudio: grupo de Casos, de Controles y Base completa. A cada una de las mamografías presentes en los tres grupos de estudio se realizó la extracción de características con cada uno de los 32 métodos.

Para realizar el análisis estadístico y lograr determinar el grado de similitud en los patrones parenquimatosos de las mamografías reales y mamografías sintéticas, se analizaron los valores cuantitativos arrojados por cada uno de los métodos de extracción de características utilizados sobre las mamografías reales en comparación con los valores cuantitativos arrojados por cada uno de los métodos de extracción de características utilizados sobre las mamografías sintéticas, es decir: Se analizó la información obtenida por el método de extracción número 1 aplicado sobre las mamografías reales contra la información obtenida por el método de extracción número 1 aplicado sobre las mamografías sintéticas. Es decir:

Casos	Método 1	Método 2	Método 3	Método....	Método 32
Mamografía real 1	xxxxx	xxxxx	xxxxx	....	xxxxx
Mamografía sintética 1	xxxxx	xxxxx	xxxxx	....	xxxxx
Mamografía real 2	xxxxx	xxxxx	xxxxx	....	xxxxx
Mamografía sintética 2	xxxxx	xxxxx	xxxxx	....	xxxxx
....	....	....	....	....	....
....	....	....	....	....	....
Mamografía real 59	xxxxx	xxxxx	xxxxx	....	xxxxx
Mamografía sintética 59	xxxxx	xxxxx	xxxxx	....	xxxxx

Figura 11. Representación de los valores cuantitativos obtenidos con los métodos de extracción de características y el análisis realizado sobre el conjunto de mamografías reales y mamografías sintéticas.

Los valores cuantitativos obtenidos por el método 1 sobre las mamografías reales (color azul) figura 11, se compararon con los valores cuantitativos obtenidos por el método 1 sobre las mamografías sintéticas (color amarillo) figura 11, lo anterior con cada uno de los 32 métodos y para cada grupo de Casos, Controles y Base completa.

Para lograr determinar el grado de similitud de los patrones parenquimatosos entre las mamografías reales y sintéticas se realizó un análisis estadístico descriptivo y un análisis estadístico inferencial. Para el análisis estadístico descriptivo se utilizó diagramas que representan cantidades como el valor medio, la desviación estándar, el valor central, entre otras características estadísticas. En el análisis estadístico inferencial se utilizó la prueba de hipótesis para determinar la media muestral, la correlación y la distribución acumulativa de los datos obtenidos de cada una de las mamografías reales y sintéticas.

Para aplicar una prueba de hipótesis a dos muestras, es necesario conocer la distribución de los datos de cada una de las muestras. Puesto que existen dentro de la prueba de hipótesis, pruebas paramétricas y no paramétricas, y la utilización de una u otra depende de la distribución de los datos que se analizan. Por lo anterior, se realizó una prueba de normalidad de Shapiro-Wilk utilizando la librería scipy de Python. Si los valores cuantitativos de las dos muestras a comparar (para efectos del presente trabajo de grado se hace referencia a muestra a la información obtenida por cada uno de los métodos de extracción de características, es decir las muestras

a comparar son datos de método 1 real contra datos de método 1 sintético, de esta manera con los 32 métodos.) presentan una distribución normal, se utilizan pruebas de hipótesis paramétricas. Si alguna de las dos muestras en sus datos presenta una distribución no normal o si las dos muestras presentan una distribución no normal en sus datos, se utilizan pruebas de hipótesis no paramétricas. Una vez seleccionado las muestras que presentan una distribución normal y una distribución no normal para cada uno de los grupos analizados (Casos, Controles y Base completa), se aplicó las pruebas de hipótesis paramétricas y no paramétricas para analizar las variables estadísticas de la media muestral, la correlación y la distribución acumulativa; En cada uno de los grupos, para el grupo de Casos, Controles y Base completa. Estas pruebas fueron:

**Distribución normal:**

Media muestral: Prueba t del estudiante.

Coefficiente de correlación: Coeficiente de correlación de Pearson.

**Distribución no normal:**

Media muestral: Prueba de U Mann Whitney.

Coefficiente de correlación: Coeficiente de correlación de Spearman.

**Distribución acumulativa:**

Se utilizó la prueba Z para dos muestras de Kolmogorov-Smirnov.

## 4. RESULTADOS Y DISCUSIÓN

Una primera aproximación para determinar el grado de similitud entre patrones parenquimatosos extraídos de mamografías reales y mamografías sintéticas, se basa en una inspección visual de las principales características que son posibles observar en las mamografías reales y sintéticas. Se presenta a continuación un conjunto de mamografías tomadas de manera aleatoria de la base de datos de mamografías reales, las mamografías pertenecen al grupo denominado como Casos 12 y al grupo denominado Controles 13. Utilizando el programa OpenVct fue posible generar un modelo sintético con la mayor similitud posible a las características anatómicas de las mamografías reales. Las mamografías sintéticas presentadas están divididas en un grupo denominado Casos simulados 14 y en un grupo denominado Controles simulados 15. Cada una de las mamografías sintéticas presentadas en cada uno de los dos grupos, corresponden a la simulación realizada correspondiente a las mamografías reales presentes en los dos grupos, es decir. La simulación realizada a la mamografía real del grupo Casos figura 12a corresponde a la mamografía sintética del grupo Casos Simulados figura 14a, así de manera sucesiva para cada una de las mamografías reales.

Es posible observar que las mamografías reales presentan de manera general una forma elipsoidal, aunque cada una tiene dimensiones diferentes lo cual genera diversos tamaños y formas entre las mamografías reales. Se observa de igual manera que existe una distribución sobre la mama del tejido parenquimatoso diferente en cada una de las mamografías, generando que el tejido denso no esté ubicado en un lugar específico en cada una de las mamografías.

La base de mamografías reales con la cual se trabajó, esta dividida en grupos de mamografías con cáncer y otro grupo de mamografías sin cáncer. Debido a la inexperiencia del autor del presente trabajo, no fue posible determinar aquellas caracte-



Figura 12. Conjunto de mamografías reales pertenecientes al grupo Casos. Las mamografías presentan una densidad del 32 %, 47 %, 38 %, 25 % y 26 %, de manera respectiva.

rísticas que presenta el tejido parenquimatoso y el tejido denso de una mama, las cuales a través de una inspección visual es posible determinar la sospecha de un posible cáncer de mama. Sin embargo, con el análisis cuantitativo presentado más adelante, fue posible determinar el grado de similitud existen entre las mamografías reales y las mamografías sintéticas, de igual manera fue posible establecer la capacidad del programa OpenVct para desarrollar modelos anatómicos simulados de la mama con alguna presencia de cáncer de mama.

Uno de los objetivos principales era generar modelos anatómicos simulados con la mayor similitud posible a modelos anatomicos reales, no fue posible modificar la forma elipsoidal estandar que trae por defecto el programa. Esto generó una limitación al momento de simular el volumen del modelo sintético, provocando que los volúmenes reales y sintéticos presenten una diferencia marcada entre ellos. Es posible observar en las figuras 14 y 15 que los modelos sintéticos generados presentan la misma forma elipsoidal y un volumen constante para todos los modelos.

Respecto a la característica anatomica de la densidad, los modelos simulados guardan una estrecha relación respecto a la densidad de las mamografías reales. Es posible observar que el tejido denso de los modelos simulados, se ubica alrededor de la región retroareolar. Este tejido denso es simulado de manera diferente aun

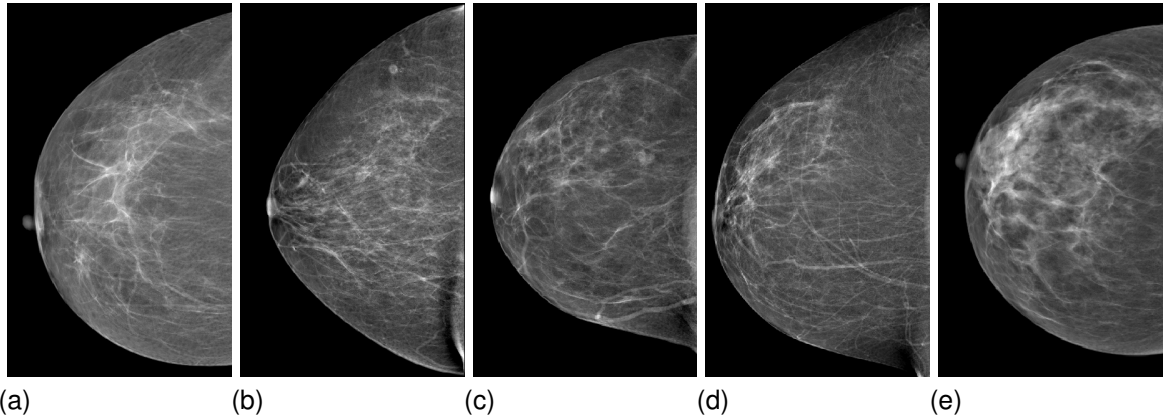


Figura 13. Conjunto de mamografías reales pertenecientes al grupo Controles. Las mamografías presentan una densidad del 16 %, 29 %, 15 %, 8 % y 48 %, de manera respectiva.

para los mismos valores de densidad, esto se evidencia en los modelos simulados con 35 % de densidad los cuales son 14b, 14c y 15e. El tejido denso en cada uno de los modelos anteriores presenta una distribución no uniforme y diferente para cada uno de los modelos. De igual forma es posible observar que el tejido denso aumenta de acuerdo al porcentaje de densidad establecido al momento de realizar la simulación, esto acorde a lo esperado entre el porcentaje de densidad de una mama y su relación con el tejido denso observado.

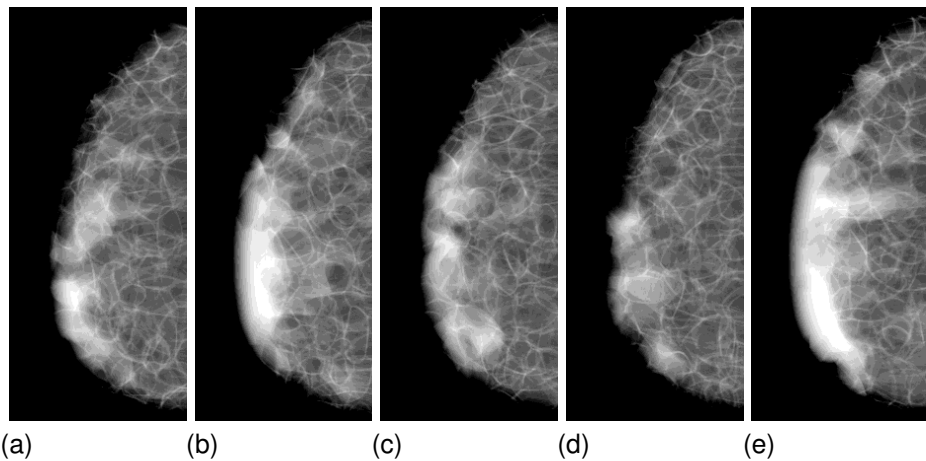


Figura 15. Conjunto de mamografías sintéticas pertenecientes al grupo Controles Simulados. Las mamografías presentan una densidad del 15 %, 30 %, 15 %, 10 % y 35 %, de manera respectiva.

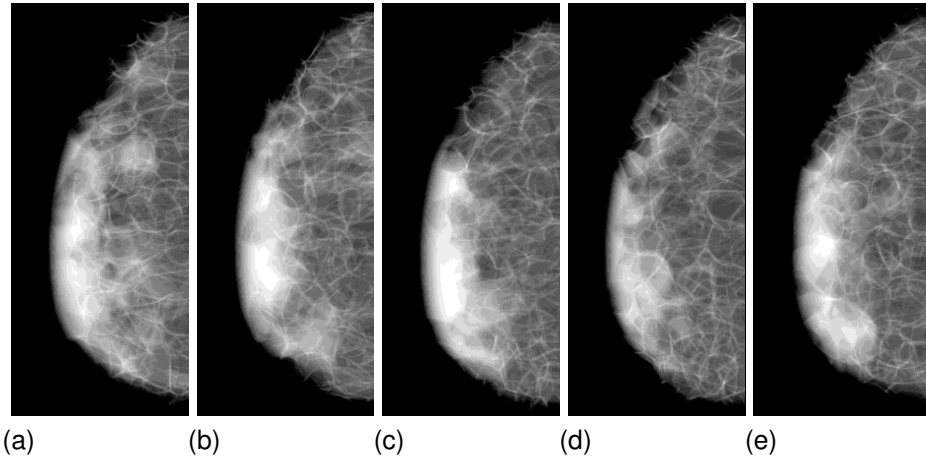


Figura 14. Conjunto de mamografías sintéticas pertenecientes al grupo Casos Simulados. Las mamografías presentan una densidad del 30 %, 35 %, 35 %, 25 % y 25 %, de manera respectiva.

Una vez realizados los anteriores análisis para los cuales se utilizó la estadística descriptiva, se da paso a interpretar los resultados obtenidos a través de la estadística inferencial; se clasificaron los valores cuantitativos obtenidos con los métodos de extracción de características entre distribuciones normales y distribuciones no normales tal y como se explica en la sección de Metodología (3). El siguiente paso fue determinar los parámetros poblacionales de las muestras pertenecientes al conjunto de mamografías reales y el conjunto de mamografías sintéticas. El primer parámetro poblacional a considerar para lograr determinar el grado de similitud entre las dos muestras, es la media muestral, para lo cual se utilizó la prueba de hipótesis. Teniendo en cuenta el parámetro de la media muestral, siguiendo los pasos expuestos en 25.

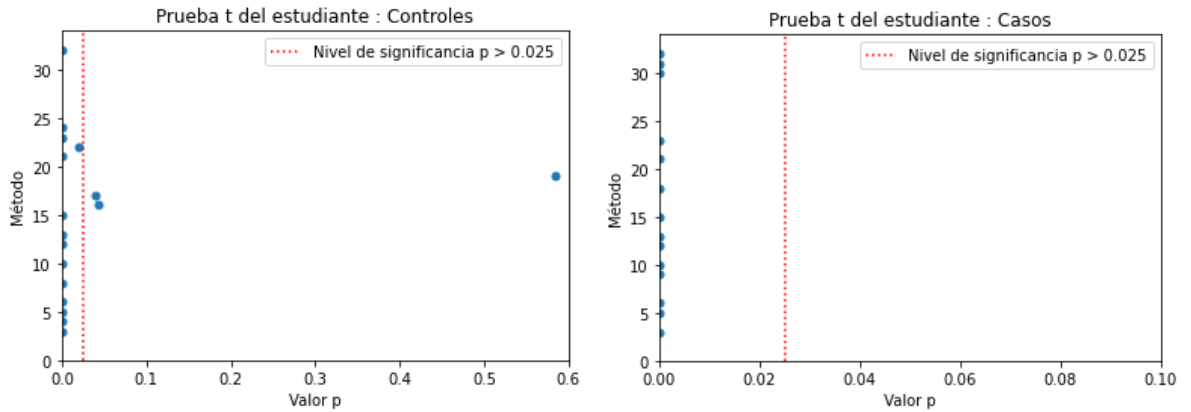
#### **4.1. Prueba estadística: media muestral**

Para realizar el análisis estadístico de la media muestral sobre las muestras de las mamografías reales y mamografías sintéticas. El primer paso corresponde a examinar los datos y determinar sus características estadísticas. De antemano se conoce que los datos son cuantitativos, continuos y provienen de dos muestras no empa-

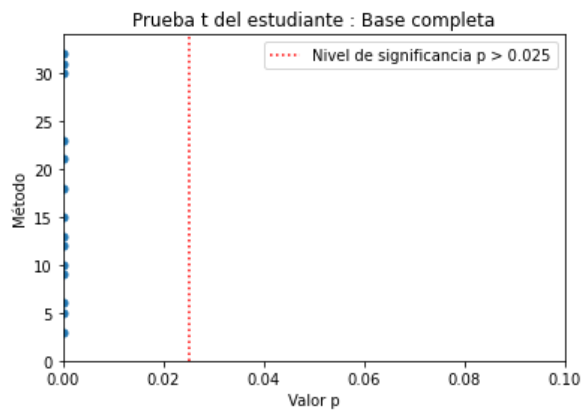
rejadas. Los datos obtenidos en su conjunto exhiben distribuciones normales y distribución no normales. El segundo paso corresponde a realizar el planteamiento de la hipótesis nula y la hipótesis alternativa para el conjunto de datos con distribución normal y con distribución no normal. La hipótesis por refutar es: No existe diferencia en la media muestral entre los valores cuantitativos obtenidos a través de la extracción de características realizado sobre la base de mamografías real y la base de mamografías sintéticas. Por lo cual la hipótesis alternativa corresponde a que si existe tal diferencia en la media muestral entre los valores cuantitativos obtenidos a través de la extracción de características realizado sobre la base de mamografías real y la base de mamografías sintéticas.

Para rechazar o aceptar la hipótesis nula se utilizó un grado de confianza del 95 % al cual corresponde un nivel de significancia  $\alpha = 0.05$ . Una vez establecida la hipótesis nula y el nivel de significancia, se procede a calcular el estadístico de prueba. Se analizó el parámetro poblacional de la media muestral para los datos que presentan una distribución normal y una distribución no normal. Por lo tanto, la prueba estadística elegida corresponde a la prueba t del estudiante para el conjunto de datos con distribución normal (2.4.2) y la prueba de U Mann-Whitney para el conjunto de datos con distribución no normal (2.4.2).

Utilizando los valores cuantitativos obtenidos a través de la extracción de características y realizando el computo a través de Python utilizando la ecuación (2) y (3), fue posible obtener los valores correspondientes al estadístico de la prueba del estudiante y la prueba de U Mann-Whitney, así como el valor p asociado a cada valor del estadístico obtenido. Los valores p obtenidos para la prueba t del estudiante se presentan en la figura 16 y los valores p obtenidos para la prueba de U Mann-Whitney se presentan en la figura 17.



(a) Valores p obtenidos para el grupo de Controles (b) Valores p obtenidos para el grupo de Casos

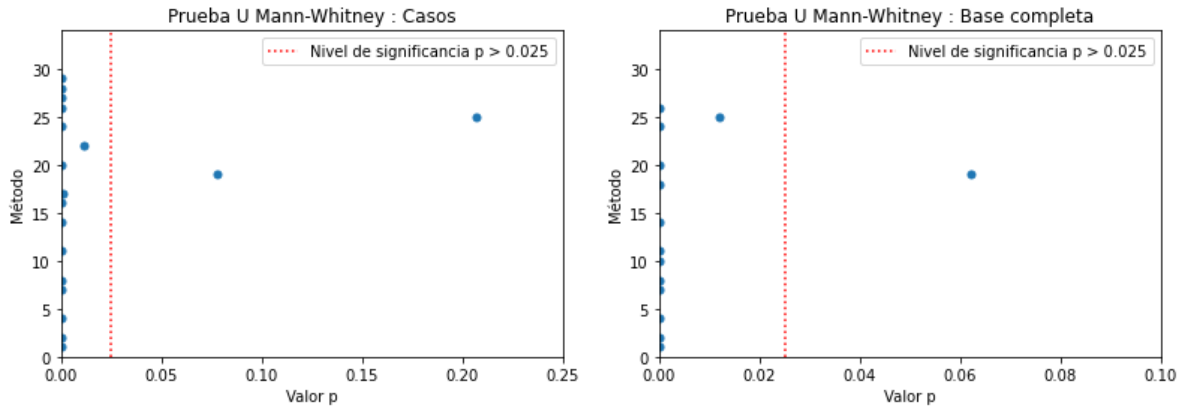


(c) Valores p obtenidos para el grupo de base completa

Figura 16. Valores p obtenidos de la prueba t del estudiante para los tres grupos de estudio.

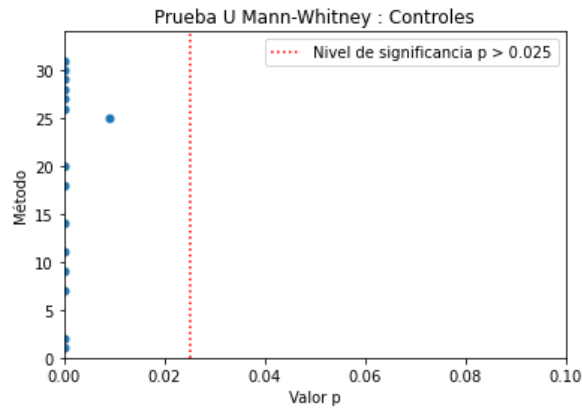
## 4.2. Análisis del valor p para la media muestral

En las gráficas presentes en las figuras (16) y (17) se ubicó en el eje horizontal el valor p correspondiente a cada método de extracción de características ubicado en el eje vertical. De igual manera se ubicó el nivel de significancia establecido en  $\alpha = 0.05$ . Para los grupos denominados Casos, Base completa y Controles (16) (17) se observa que la gran mayoría de métodos de extracción de características presentan un valor p de manera aproximada a cero. Se observa en las gráficas correspondiente al grupo de Casos 16b y al grupo Base completa 16c los diferentes valores p obtenidos de la prueba t del estudiante, sin embargo, todos los valores p



(a) Valores p obtenidos para el grupo de casos

(b) Valores p obtenidos para el grupo de base completa



(c) Valores p obtenidos para el grupo de controles

Figura 17. Valores p obtenidos de la prueba U Mann-Whitney para los tres grupos de estudio.

para los dos grupos mencionados son iguales a cero. Para el grupo de Controles es posible observar que, para tres métodos de extracción de características, su valor p es mayor a 0.025. Estos métodos de extracción de características son: El método 16 con un valor p de 0.043, el método 17 con un valor p de 0.04 y el método 19 con un valor p de 0.584. El valor p del método 22 presenta un valor de 0.019, por lo cual se ve próximo al nivel de significancia de 0.025, sin embargo, no logra pasar el umbral establecido.

Respecto a la gráfica de los valores p obtenidos de la prueba de U Mann-Whitney (17). Existen dos métodos para el grupo de Casos y un método para el grupo de

Base completa, para los cuales su valor p es mayor al nivel de significancia establecido. Estos métodos son: 25 y 19 para el grupo de Casos y el método 19 para el grupo Base completa.

Para el grupo de Casos el método 22 y el método 24 para el grupo de Base completa, sus valores p se acercan al nivel de significancia, con un valor p de 0.0111 y 0.012 de manera respectiva, sin embargo, estos valores no logran superar el nivel de significancia. En comparación con los valores estadísticos obtenidos con la prueba t del estudiante los cuales para un método presentaba un valor negativo, los valores estadísticos obtenidos con la prueba U Mann Whitney son todos positivos, por lo cual, existe una asociación que indica que la media muestral de los valores obtenidos de la extracción de características de las mamografías sintéticas es mayor a comparación de la media muestral de los valores obtenidos de la extracción de características de las mamografías reales. La importancia de elegir una prueba de dos colas para el presente trabajo se observa con el estadístico obtenido con el método número 17. Este estadístico tiene un valor de  $r = -2.0699$ . Debido al nulo conocimiento respecto a la diferencia de la media muestral de una muestra en comparación con la otra muestra analizada respecto a si era negativa o positiva, se eligió una prueba de dos colas para no perder información. Con los valores p obtenidos y el nivel de significancia establecido, es posible comparar estos dos valores para aceptar o rechazar la hipótesis nula de la prueba t del estudiante y la prueba de U Mann-Whitney.

Para los métodos de extracción de características que presentaron un valor p inferior al nivel de significancia, es posible rechazar la hipótesis nula establecida. Estos métodos para la prueba t del estudiante son: Todos los métodos del grupo de controles a excepción del método 16,17 y 19, todos los métodos del grupo Base completa y Casos. Para la prueba de U Mann-Whitney son: Todos los métodos del grupo Controles, los métodos de Base completa a excepción del método 19 y los métodos de casos a excepción del método 19,25. Los métodos para los cuales se obtuvieron

valores p de manera aproximada al valor cero y por lo tanto no superaron el nivel de significancia, es posible establecer que existe una diferencia entre las medias muestrales para los valores cuantitativos obtenidos con los métodos de extracción de características entre las mamografías sintéticas y las mamografías reales.

Los valores cuantitativos extraídos de las mamografías reales y las mamografías sintéticas no presentan una relación con la hipótesis nula establecida, por lo cual se considera el resultado como estadísticamente significativo y se rechaza la hipótesis nula. Por lo tanto, se acepta la hipótesis alternativa la cual considera que existe una diferencia en la media muestral en los valores cuantitativos obtenidos con los métodos de extracción de características en cada una de las mamografías reales y mamografías sintéticas. De igual manera esta diferencia encontrada entre las muestras pertenecientes al grupo de mamografías reales y al grupo de mamografías sintéticas es muy poco probable que se deba a un producto del azar, es decir, las muestras elegidas pertenecen a un pequeño grupo muestral que no representa las características de la población.

El método 19 y 25 del grupo Casos y el método 19 del grupo Base completa pertenecientes a la prueba de U Mann-Whitney se obtuvo un valor p mayor al nivel de significancia, es decir, se obtuvo una probabilidad la cual al compararla con las zonas de rechazo o no rechazo de la hipótesis nula, estos valores p se ubican en la zona de no rechazo de la hipótesis nula establecida. Por lo cual no es posible rechazar la hipótesis nula ya que los datos obtenidos a través del análisis de la prueba de hipótesis no proporcionan la evidencia suficiente para rechazar la hipótesis nula. Es posible afirmar que, para los métodos mencionados para las dos pruebas estadísticas, el parámetro estadístico de la media muestral de la población de mamografías real y de la población de mamografías sintéticas no presentan diferencia significativa.

Debido a los datos analizados y a la evidencia obtenida con el desarrollo de la

prueba de hipótesis, es muy poco probable que el resultado obtenido se deba a un producto del azar, es decir, que las muestras elegidas pertenecen a un pequeño grupo muestral el cual exhibe características que no representan las características estadísticas de la población en general.

#### **4.3. Prueba estadística: coeficiente de Correlación**

Para establecer el nivel de correlación existente entre las mamografías reales y las mamografías sintéticas. Se utilizó la prueba de hipótesis para analizar los valores cuantitativos obtenidos a través de la extracción de características de los patrones parenquimatosos presentes en las mamografías reales y sintéticas, y lograr determinar el tipo de correlación existente entre las mamografías mencionadas.

Para las muestras que presentaron en sus datos una distribución normal, se utilizó el coeficiente de correlación de Pearson (5). Para las muestras con una distribución no normal en sus datos, se utilizó el coeficiente de correlación de Spearman (6). Para lograr determinar la correlación entre las dos bases de mamografías se utilizó la prueba de hipótesis con el parámetro poblacional de correlación. El desarrollo de la prueba de hipótesis se llevó a cabo con los pasos mencionados en 25. Como primer paso se clasificaron los datos existentes. El segundo paso consiste en establecer la hipótesis nula, para el presente estudio se estableció la siguiente hipótesis nula, tanto para la prueba de hipótesis en la cual los datos presentan una distribución normal y una distribución no normal;  $H_0$  : No existe correlación entre los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y las mamografías sintéticas. Como hipótesis alternativa: Existe una correlación entre los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y las mamografías sintéticas.

Para esta prueba de hipótesis se utilizó un grado de confianza del 95% al cual corresponde un nivel de significancia  $\alpha = 0.05$ .

Con los valores cuantitativos obtenidos de la extracción de características, las ecuaciones de las pruebas de correlación descritas con anterioridad y el programa Python, fue posible obtener el valor estadístico y el valor p para cada uno de los métodos de extracción de características. A continuación, se presentan los valores estadísticos y los valores p obtenidos para cada una de las pruebas.

Método	Estadístico r	Valor p
3	-0,10	0,28
5	-0,08	0,38
6	0,36	6,24E-05
9	-0,07	0,42
12	0,36	7,82E-05
13	0,18	0,05
15	0,31	5,25E-04
16	0,33	2,09E-04
17	0,30	1,17E-03
21	0,01	0,92
22	-0,03	0,32
23	-0,43	1,40E-06
27	0,57	9,83E-12
28	0,43	1,06E-06
29	0,42	1,85E-06
30	0,16	0,08
31	0,38	2,49E-05
32	-0,03	0,74

Tabla 2. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Pearson para los métodos utilizados en el grupo de Base completa. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Pearson son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.

#### 4.4. Análisis del valor p para el coeficiente de correlación

En las tablas (2), (3) y (4) se presentan los tres grupos analizados con el coeficiente de correlación de Pearson, en las tablas (5), (6) y (7) se presentan los tres grupos analizados con el coeficiente de correlación de Spearman. En cada una de las tablas anteriores se registraron los valores estadísticos y los valores p obtenidos de la prueba de correlación para los diferentes métodos de extracción de características

Método	Estadístico r	Valor p
3	-0,15	0,25
5	0,42	0,00
6	-0,04	0,78
9	-0,07	0,58
10	-0,16	0,21
12	0,19	0,14
13	0,20	0,14
15	0,57	2,96E-06
18	-0,05	0,71
21	0,09	0,52
23	-0,44	5,50E-04
30	0,43	6,83E-04
31	0,63	8,12E-08
32	-0,06	0,68

Tabla 3. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Pearson para los métodos utilizados en el grupo de Casos. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Pearson son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.

utilizados. Los recuadros de los métodos en color verde indican un valor p menor al nivel de significancia de  $\alpha = 0.05$ , puesto que se trabajó con una prueba de dos colas, el nivel de significancia se debe dividir en dos. Los valores estadísticos (r) obtenidos con las pruebas de coeficiente de correlación representan la fuerza y el sentido de la correlación entre las dos variables.

En la tabla de base completa (2), se observa que existe para el método 23 una correlación fuerte en sentido negativo. Además, existe para este grupo un gran número de correlaciones en 10 de los 18 métodos utilizados sobre este grupo. Se evidencia que 9 de los 10 métodos correlacionados presentan una correlación positiva.

Para el grupo de Casos (3) se observa que el número de métodos correlacionados es bajo, representando tan solo 5 de 14 métodos utilizados sobre este grupo. Se observa para el método 15, 23 y 31 una fuerte correlación debido al valor estadístico r. Solo un método de los 5 correlacionados presenta un sentido negativo. En el grupo

Método	Estadístico r	Valor p
3	-0,20	0,13
4	0,14	0,28
5	-0,08	0,54
6	-0,01	0,93
8	0,23	0,07
10	0,01	0,92
12	0,06	0,66
13	0,05	0,72
15	0,09	0,50
16	-0,15	0,26
17	0,52	2,28E-05
19	0,12	0,35
21	-0,47	2,05E-04
22	0,32	0,01
23	-0,49	8,13E-05
24	-0,04	0,79
32	-0,03	0,82

Tabla 4. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Pearson para los métodos utilizados en el grupo de Controles. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Pearson son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.

de Controles (4) 4 métodos de un total de 17 las muestras analizadas presentan una correlación entre ellas. Los métodos 17, 21 y 23 presentan una fuerte correlación debido a que su valor r supera el valor de 0.40. Los métodos 21 y 23 presentan una correlación negativa.

En la tabla (5) los métodos 4 y 19 presentan una correlación fuerte positiva. Para el grupo Casos (6), los métodos 4,7,16,19,22 y 28 presentan una correlación en sus muestras, el método 4 y 19 presentan una correlación fuerte. El método 29 presenta una correlación negativa. Para el grupo de Controles (7) se observa que los métodos 18 y 20 presentan una correlación negativa.

Los valores estadísticos r obtenidos se pueden visualizar utilizando un diagrama de dispersión. En un diagrama de dispersión se gráfica en cada uno de los ejes una

Método	Estadístico r	Valor p
1	0,05	0,56
2	-0,01	0,89
4	0,51	2,94E-09
7	0,25	0,05
8	-0,02	0,81
10	-0,11	0,25
11	-0,13	0,16
14	0,09	0,33
18	0,09	0,35
19	0,43	1,33E-06
20	-0,25	0,05
24	-0,08	0,36
25	-0,07	0,45
26	0,03	0,72

Tabla 5. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Spearman para los métodos utilizados en el grupo de Base completa. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Spearman son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa.

de las dos muestras. La disposición de los puntos en el gráfico y la forma de la gráfica revela información acerca de la fuerza de correlación entre las dos muestras y el sentido de la correlación, si es un sentido positivo la gráfica se debe observar de manera ascendente, lo cual indica que a medida que una variable aumenta la otra también aumenta. En una correlación negativa la gráfica se observa de manera descendente, esto indica que a medida que una variable aumenta la otra variable disminuye. Se presentan varios diagramas de dispersión para visualizar el comportamiento entre las dos muestras analizadas.

En la figura 18a en la cual se representa la gráfica de dispersión de los valores cuantitativos obtenidos del método 4 de extracción para el grupo Base completa, es posible observar que los puntos están agrupados de forma compacta, esto se puede visualizar debido a los intervalos de los ejes X y Y en los cuales están distribuidos los puntos, con el valor  $r = 0.51$ . De igual manera, se observa en la figura 18a que la dispersión presenta una dirección ascendente, lo cual indica que los valores extraídos de los patrones parenquimatosos utilizando el método 4 aumentaban en las

Método	Estadístico r	Valor p
1	0,14	0,31
2	0,05	0,70
4	0,56	3,38E-06
7	0,30	0,02
8	0,03	0,81
11	-0,08	0,57
14	-0,06	0,64
16	0,37	0,004
17	0,08	0,55
19	0,52	2,59E-05
20	-0,15	0,25
22	0,34	0,01
24	-0,16	0,24
25	-0,10	0,45
26	-0,12	0,36
27	-0,07	0,58
28	-0,39	0,002
29	-0,08	0,53

Tabla 6. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Spearman para los métodos utilizados en el grupo de Casos. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Spearman son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa

mamografías reales a la par que aumentaban en las mamografías sintéticas pertenecientes al grupo Base completa. Es decir, se observa una correlación positiva.

Para la gráfica del método de extracción 22 para el grupo de Controles 18b, se observa un comportamiento similar al descrito con anterioridad, sin embargo, los puntos sobre la gráfica no se encuentran distribuidos de manera tan compacta, esto se ve reflejado en el valor estadístico calculado el cual corresponde a  $r=0.32$ . Aunque existe una correlación entre las dos muestras, la correlación no es igual de significativa como la observada con anterioridad. De igual manera la correlación observada presenta un comportamiento positivo, es decir, los valores extraídos de los patrones parenquimatosos utilizando el método 22 aumentaban en las mamografías reales a la par que aumentaban en las mamografías sintéticas pertenecientes al grupo Con-

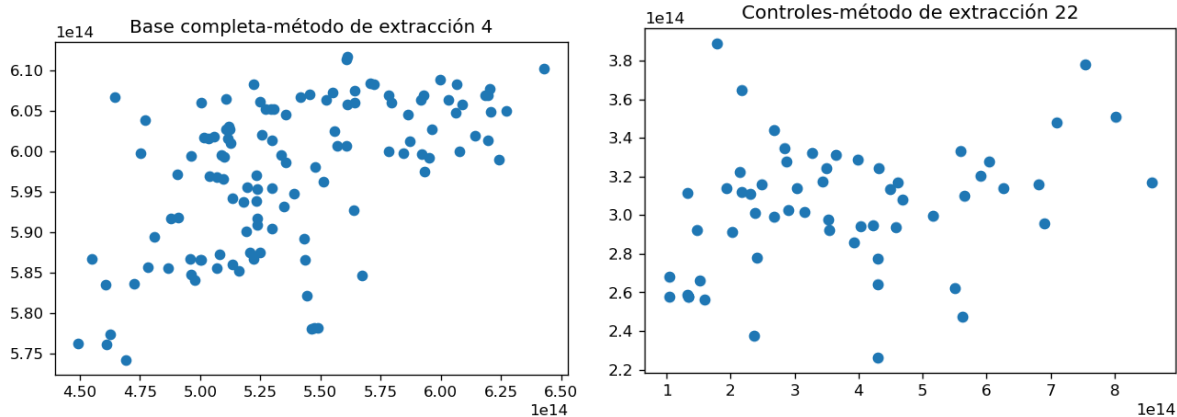
Método	Estadístico r	Valor p
1	-0,12	0,37
2	0,25	0,06
7	0,26	0,05
9	-0,12	0,36
11	-0,14	0,28
14	0,03	0,83
18	-0,31	0,02
20	-0,34	0,01
25	-0,28	0,03
26	0,18	0,17
27	-0,26	0,05
28	-0,02	0,91
29	-0,01	0,94
30	0,05	0,72
31	0,08	0,57

Tabla 7. Valor estadístico (r) y valor p obtenidos de la prueba de correlación de Spearman para los métodos utilizados en el grupo de Controles. En color verde se indican los métodos para los cuales los valores p de la prueba de correlación de Spearman son menores al nivel de significancia establecido, por lo cual, se rechaza la hipótesis nula y se acepta la hipótesis alternativa

troles.

La figura 19 en la cual se representa la dispersión de los valores cuantitativos obtenidos del método 23 para el grupo de Base completa, se observa que los puntos ubicados sobre la figura se encuentran distribuidos de manera compacta. A diferencia de las figuras anteriores, la figura de dispersión presenta una orientación descendente. Lo anterior se relaciona con el valor estadístico calculado el cual corresponde a un valor de  $r = -0.43$ . La magnitud representa la fuerza de la correlación entre las dos muestras y el signo menos indica la correlación negativa. Es decir, los valores extraídos de los patrones parenquimatosos utilizando el método 23 aumentaban en las mamografías reales a la par que disminuía en las mamografías sintéticas pertenecientes al grupo Base completa.

Los anteriores resultados se observaron para los métodos de extracción de carac-



(a) Método de extracción 4 grupo Base completa (b) Método de extracción 22 grupo Controles

Figura 18. Gráficas de Dispersión de los valores cuantitativos del método de extracción de características 4 para el grupo Base completa y el método de extracción de características 22 del grupo Controles

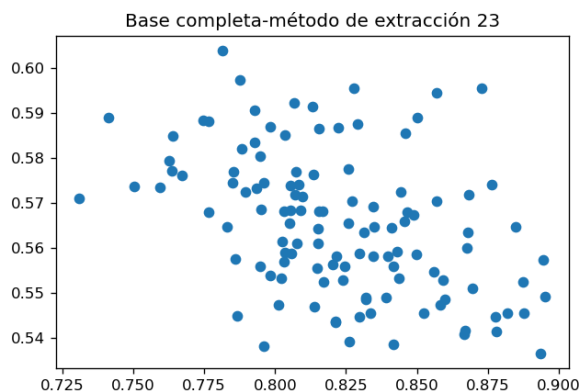
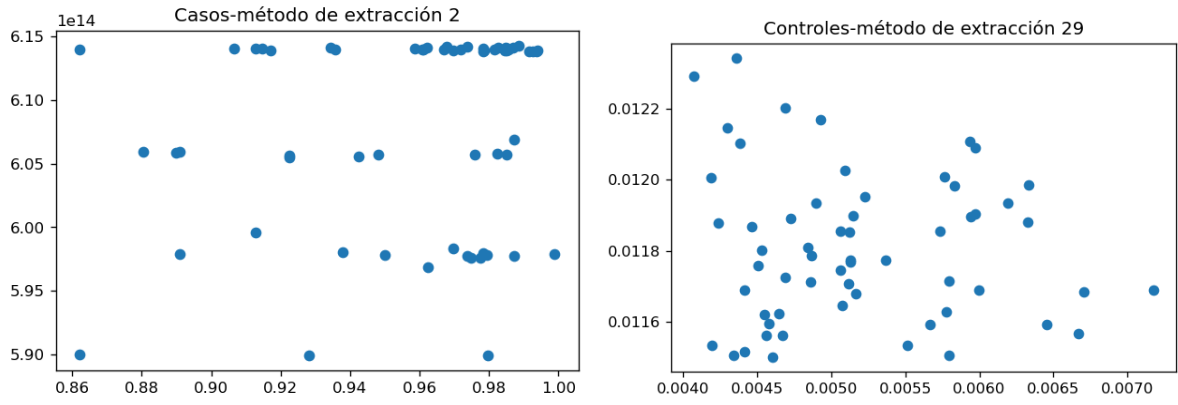


Figura 19. Gráfica de dispersión de los valores cuantitativos del método de extracción de características 23 para el grupo Base completa

terísticas que obtuvieron un valor p menor al nivel de significancia de  $\alpha = 0.025$ . Las gráficas de dispersión que permiten observar la correlación existente para los métodos los cuales presentaron un valor p mayor al nivel de significancia se presentan a continuación:

Se observa en las gráficas anteriores que la distribución de los puntos sobre el gráfico sucede de manera aleatoria, no presentan ningún orden y tampoco un sentido establecido. Este comportamiento era de esperarse debido al valor estadístico obtenido para cada uno de los métodos representados de manera gráfica. Los valores estadísticos obtenidos fueron:  $r=0.05$  para la gráfica 20a,  $r=0.01$  para la gráfica 20b



(a) Método de extracción 2 grupo Casos

(b) Método de extracción 29 grupo Controles

Figura 20. Gráficas de Dispersión de los valores cuantitativos del método de extracción de características 2 para el grupo Casos y el método de extracción de características 29 del grupo Controles

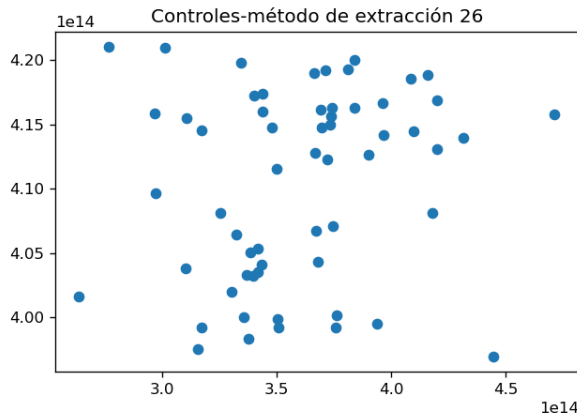


Figura 21. Gráfica de Dispersión de los valores cuantitativos del método de extracción de características 26 para el grupo Controles

y  $r=0.18$  para la gráfica 21. Los valores estadísticos representa un valor muy cercano al cero, lo cual indica que no existe una correlación entre las dos muestras, es decir, valores extraídos de los patrones parenquimatosos utilizando los métodos descritos en las mamografías reales y en las mamografías sintéticas no presentan una correlación.

De manera general, todas las gráficas de dispersión para todos los métodos presentan una correlación o una correlación nula de tal manera como se observó en las gráficas anteriores. El comportamiento esperado respecto a la correlación de

cada método depende del valor estadístico obtenido  $r$  y el valor  $p$ . En las gráficas presentadas se observa que existe una correlación para valores  $p$  menores al nivel de significancia y una no correlación para valores  $p$  mayores al nivel de significancia.

El último paso para el desarrollo de la prueba de hipótesis consiste en comparar el valor  $p$  obtenido de la prueba de correlación, tanto para la prueba de correlación de Pearson y la prueba de correlación de Spearman; Con el nivel de significancia establecido de  $\alpha = 0.05$ , de esta comparación surge la posibilidad de rechazar o aceptar la hipótesis nula. Los métodos que presentaron un valor  $p$  menor al nivel de significancia se presentan a continuación, para el coeficiente de correlación de Pearson y el coeficiente de correlación de Spearman.

<b>Coeficiente de correlación de Pearson</b>			
	<b>Casos</b>	<b>Controles</b>	<b>Base completa</b>
<b>Métodos</b>	5,15,23,30,31	17,21,22,23	6,12,15,16,17,23,27,28,29,31

Tabla 8. Métodos para los cuales existe una correlación en los valores de las muestras para la prueba de correlación de Pearson

<b>Coeficiente de correlación de Spearman</b>			
	<b>Casos</b>	<b>Controles</b>	<b>Base completa</b>
<b>Métodos</b>	4,7,16,19,22,28	18,20	4,19

Tabla 9. Métodos para los cuales existe una correlación en los valores de las muestras para la prueba de correlación de Spearman

El desarrollo y análisis de la prueba de hipótesis conlleva a establecer que existe una relación de correlación entre las muestras analizadas para los métodos indicados en las tablas (8) y (9), por lo cual se obtiene un resultado estadísticamente significativo y se rechaza la hipótesis nula establecida. Se determina que existe una correlación en los valores cuantitativos extraídos de los patrones parenquimatosos obtenidos con los métodos de extracción de características entre las mamografías reales y las mamografías sintéticas pertenecientes a los grupos de Casos, Controles y Base completa. Debido a los datos analizados y la evidencia obtenida con el desarrollo de la prueba de hipótesis, es muy poco probable que el resultado obteni-

do se deba a un producto del azar, es decir, que las muestras elegidas pertenezcan a un pequeño grupo el cual exhibe una correlación entre ellas, sin embargo, las poblaciones a comparar no presentan ningún tipo de correlación.

Los métodos para los cuales se obtuvo un valor p mayor al nivel de significancia se indican en las tablas 10 y 11, es decir, los valores p obtenidos de las pruebas de coeficiente de correlación de Pearson y Spearman se ubicaron en la zona de no rechazo de la hipótesis nula establecida. Por lo cual, no es posible rechazar la hipótesis nula ya que la información y análisis correspondiente de la prueba de hipótesis, no proporcionan la evidencia suficiente para rechazar la hipótesis nula. Es posible afirmar por lo tanto que no existe una correlación entre los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y las mamografías sintéticas para los métodos indicados en las tablas 10 y 11 .

<b>Coefficiente de correlación de Pearson</b>			
	<b>Casos</b>	<b>Controles</b>	<b>Base completa</b>
<b>Métodos</b>	3,6,9,10,12,13,18,21,32	3,4,5,6,8,10,12,13,15,16,19,24,32	3,5,9,13,21,22,30,32

Tabla 10. Métodos para los cuales no existe una correlación en los valores de las muestras para la prueba de correlación de Pearson

<b>Coefficiente de correlación de Spearman</b>			
	<b>Casos</b>	<b>Controles</b>	<b>Base completa</b>
<b>Métodos</b>	1,2,8,11,14,17,20,24,25,26,27,29	1,2,7,9,11,14,25,26,27,28,29,30,31	1,2,7,8,10,11,14,18,20,24,25,26

Tabla 11. Métodos para los cuales no existe una correlación en los valores de las muestras para la prueba de correlación de Spearman

#### **4.5. Prueba estadística: distribución acumulativa**

El último parámetro poblacional para analizar corresponde a la distribución acumulada. Para determinar si los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y las mamografías sintéticas presentan una misma distribución. Para tal fin se utilizó la prueba de hipótesis para analizar dicho

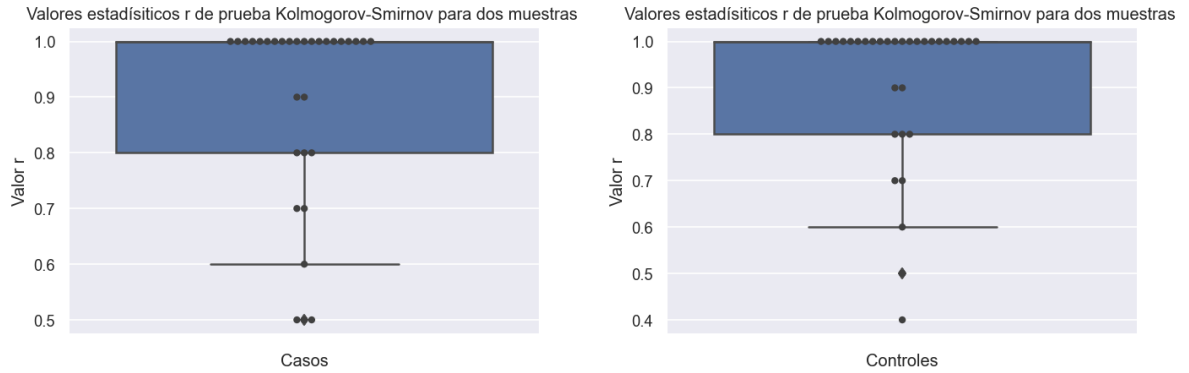
parámetro poblacional.

Para el desarrollo de la prueba de hipótesis es necesario establecer la hipótesis nula y la hipótesis alternativa. La hipótesis nula establecida fue: Los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y las mamografías sintéticas presentan una distribución acumulada igual. La hipótesis alternativa por lo tanto se definió como: Los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y las mamografías sintéticas presentan una distribución acumulada diferente. Para el desarrollo de la prueba se utilizó un grado de confianza del 95% y un nivel de significancia  $\alpha = 0.05$ . Con los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y sintéticas y utilizando la ecuación para el estadístico de Kolmogórov-Smirnov para dos muestras (7) y (8), se calcularon a través de Python los valores estadísticos y los valores p. Los estadísticos calculados para cada uno de los grupos de análisis se presentan a continuación:

#### **4.6. Análisis del valor p para distribución acumulativa**

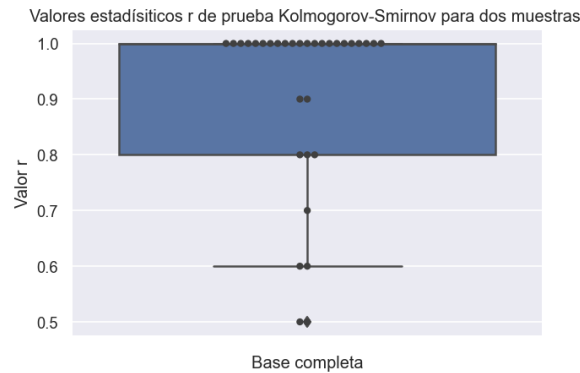
Se observa en los gráficos 22 que el valor estadístico  $r = 1$ , se presenta en la gran mayoría de los métodos de extracción de características utilizados en cada uno de los grupos analizados. De acuerdo con la teoría, un valor alto en el estadístico de Kolmogorov-Smirnov, representa un valor p bajo [67] [68] [69]. Los valores p obtenidos a través de Python se exponen en las siguientes figuras:

De manera general los valores p obtenidos en cada uno de los grupos analizados, presentan un valor de cero de manera aproximada. Ningún valor p obtenido de la prueba estadística logró pasar el umbral de significancia establecido en  $\alpha = 0.025$ . Los valores estadísticos y los valores p están relacionados en cada una de las pruebas de hipótesis, para este caso de la prueba de Kolmogorov-Smirnov, una cantidad alta en el valor estadístico va a representar un valor p bajo, esto se justifica con el hecho de que el estadístico de Kolmogorov-Smirnov calcula la diferencia máxima entre las funciones de distribución acumulada de las dos muestras. Un valor alto en



(a) Base completa

(b) Casos



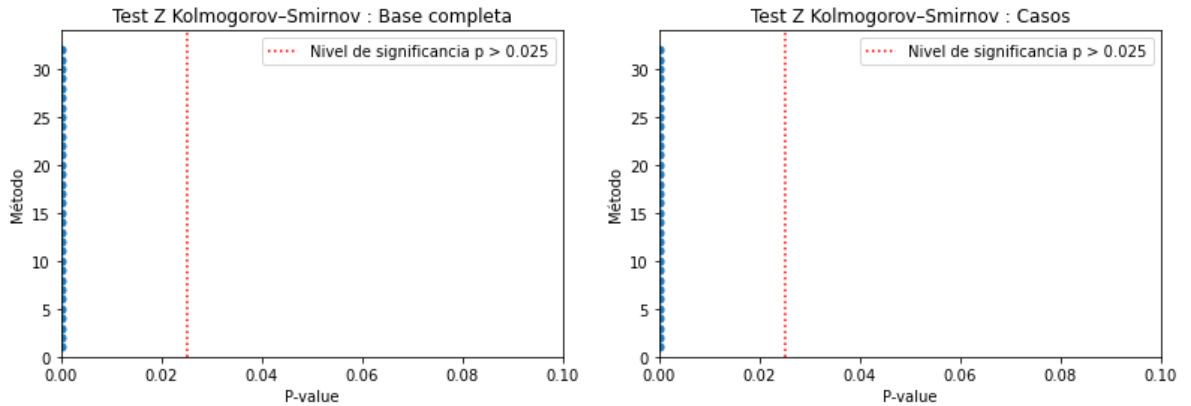
(c) Controles

Figura 22. Valores estadísticos de la prueba de Kolmogorov-Smirnov para dos muestras obtenidos en cada uno de los grupos analizados

el estadístico se ve reflejado en distribución acumuladas diferentes para cada una de las muestras. Este resultado es posible observarlo graficando las distribuciones acumuladas para un grupo de valores cuantitativos extraídos de los patrones parenquimatosos. Por ejemplo:

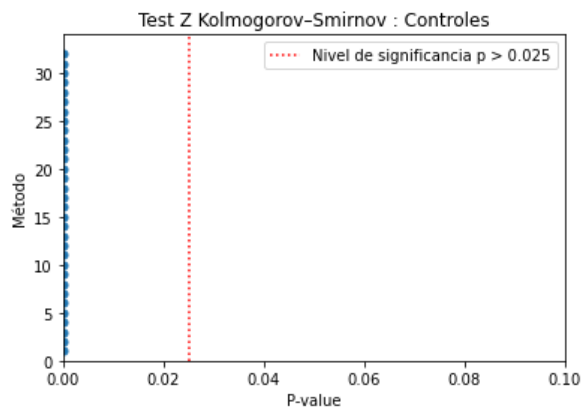
En las gráficas anteriores se observa que las distribuciones acumuladas de cada muestra difieren entre sí. Y puesto que el estadístico de Kolmogorov-Smirnov determina la diferencia máxima entre estas dos funciones, es de esperar obtener valores estadísticos altos y valores p bajos.

Con los valores p obtenidos se procede a compararlos con el nivel de significancia establecido, de tal manera que se acepte o rechace la hipótesis nula. Ningún



(a) Base completa

(b) Casos

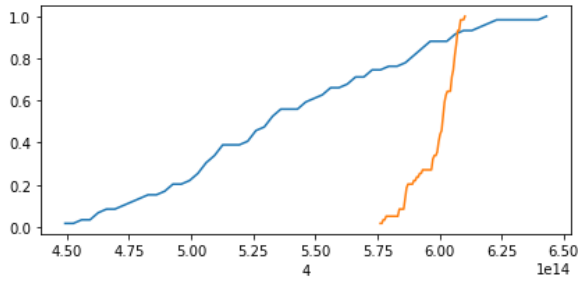


(c) Controles

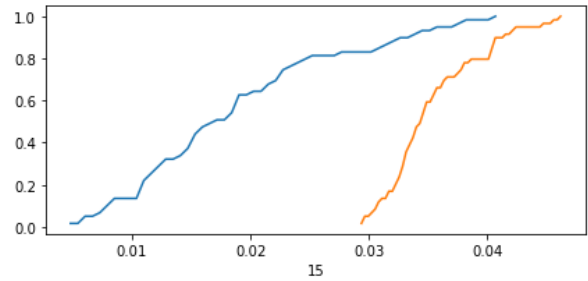
Figura 23. Valores p obtenidos de la prueba estadística Kolmogorov-Smirnov para cada uno de los grupos analizados

valor p obtenido superó el umbral del nivel de significancia establecido, por lo cual es posible establecer que las funciones de distribución acumulativas de los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales y mamografías sintéticas presentan distribuciones diferentes entre sí.

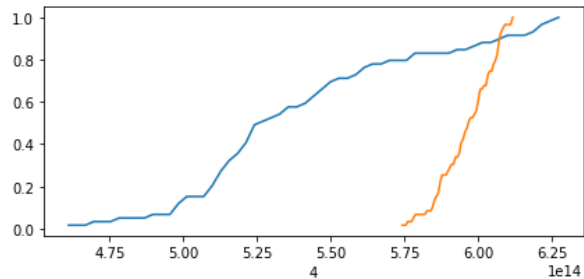
Los valores cuantitativos extraídos de las mamografías reales y las mamografías sintéticas no presentan una relación con la hipótesis nula establecida, por lo cual se considera el resultado como estadísticamente significativo y se rechaza la hipótesis nula. Por lo cual se acepta la hipótesis alternativa la cual considera que los valores cuantitativos extraídos de los patrones parenquimatosos de las mamografías reales



(a) Distribución acumulativa para el grupo de Casos método de extracción 4



(b) Distribución acumulativa para el grupo de Casos método de extracción 15



(c) Distribución acumulativa para el grupo de Control método de extracción 4

Figura 24. Gráficas de distribución acumulativa para los grupos de Casos con método de extracción 4 y 15, y grupo de Control con método de extracción 4

y las mamografías sintéticas presentan una distribución acumulada diferente entre sí.

#### 4.7. Determinación del grado de similitud

Realizados los análisis de media muestral, correlación y distribución acumulada sobre los valores cuantitativos extraídos de los patrones parenquimatosos presentes en las mamografías reales y mamografías sintéticas. Se determina el bajo nivel de similitud entre las bases de mamografías reales y mamografías sintéticas. El resultado más cercano a determinar una posible similitud fue el análisis de la correlación, en el análisis realizado un tercio de los métodos analizados presentó una fuerte correlación, es decir, entre las mamografías reales y sintéticas existe una correlación entre sus valores cuantitativos extraídos de los patrones parenquimatosos. Esto sig-

nifica que los valores cuantitativos presentaban una relación de proporcionalidad, a medida que un valor cuantitativo extraído de un patrón parenquimatoso de una mamografía real presenta un aumento, el valor cuantitativo extraído de un patrón parenquimatoso de una mamografía sintética aumenta en la misma medida. De igual forma en el análisis de la media muestral, se evidenció que no existe relación alguna entre los valores medios de los valores cuantitativos extraídos de las mamografías reales y las mamografías sintéticas. Por último se analizó la distribución muestral para los valores cuantitativos extraídos de las mamografías reales y mamografías sintéticas, se buscó determinar si la distribución de los datos presentaban una misma distribución acumulativa. Sin embargo, los resultados obtenidos indican que los valores extraídos de las mamografías reales presentan una distribución acumulativa diferente a la distribución acumulativa de los valores extraídos de las mamografías sintéticas. Por la evidencia recopilada es posible concluir que no existe un grado de similitud entre los patrones parenquimatosos extraídos de mamografías reales y mamografías sintéticas. Las mamografías reales y las mamografías sintéticas generadas difieren en gran medida en sus patrones parenquimatosos.

En las tablas (12) (13) se presenta de manera resumida los diferentes análisis estadísticos realizados sobre el conjunto de datos obtenidos de acuerdo con su distribución: distribución normal y distribución no normal. De igual manera, se presentan los métodos de extracción de características para los cuales existe una similitud en los patrones parenquimatosos de las mamografías reales con las mamografías sintéticas y los métodos para los cuales no existe una similitud en los patrones parenquimatosos.

Los resultados obtenidos pueden presentar varias razones. La primera de ellas corresponde a una razón estadística. Al momento en que se elige un grado de confianza del 95 % existe la posibilidad de no acertar en los resultados obtenidos por un factor del 5 % y que los resultados aparentemente verdaderos corresponden a un producto del azar.

Prueba estadística	Método de extracción de características					
	Existe similitud			No existe similitud		
	Casos	Controles	Base completa	Casos	Controles	Base completa
<b>Media muestral</b>	Ninguno	16,17,19	Ninguno	3, 5, 6, 9, 10, 12, 13, 15, 18, 21, 23, 30, 31, 32	3, 4, 5, 6, 8, 10, 12, 13, 15, 21, 22, 23, 24, 32	3, 5, 6, 9, 10, 12, 13, 15, 16, 21, 23, 30, 31, 32
<b>Coefficiente de correlación</b>	5, 15, 30, 31	17, 22	6, 12, 15, 16, 17, 27, 28, 29, 31	3, 6, 9, 10, 12, 13, 18, 21, 23, 32	3, 4, 5, 6, 8, 10, 12, 13, 15, 16, 19, 21, 23, 24, 32	3, 5, 9, 13, 21, 22, 23, 30, 32
<b>Distribución acumulativa</b>	Ninguno	Ninguno	Ninguno	-	-	-

Tabla 12. Métodos de extracción de características para los cuales existe una similitud en los patrones parenquimatosos de las mamografías reales con las mamografías sintéticas para el análisis estadístico con distribución normal de media muestral (16), coeficiente de correlación (4)(3) (2) y distribución acumulativa (22) (24)(23) .

Las mamografías sintéticas no simulaban de manera correcta el volumen correspondiente a las mamografías reales. Por limitaciones del programa no fue posible generar volúmenes similares a los volúmenes reales. Esta limitación al momento de generar los modelos anatómicos simulados genera la posibilidad de un posible fallo al momento de realizar la extracción de características sobre las mamografías sintéticas, obteniendo valores mayores o menores en comparación con los posibles valores obtenidos de mamografías sintéticas con volúmenes similares a mamografías reales. Aunque no fue posible generar volúmenes similares, la densidad simulada para cada mamografía sintética presenta una gran similitud a la densidad presentada en las mamografías reales.

Prueba estadística	Método de extracción de características					
	Existe similitud			No existe similitud		
	Casos	Controles	Base completa	Casos	Controles	Base completa
<b>Media muestral</b>	19,25	Ninguno	19	1, 2, 4, 7, 8, 11, 14, 16, 17, 20, 22, 24, 26, 27, 28, 29	1, 2, 7, 9, 11, 14, 18, 20, 25, 26, 27, 28, 29, 30, 31	1, 2, 4, 7, 8, 11, 14, 17, 18, 20, 22, 24, 25, 26, 27, 28, 29
<b>Coefficiente de correlación</b>	4, 7, 16, 19, 22	Ninguno	4, 19	1, 2, 8, 11, 14, 17, 20, 24, 25, 26, 27, 28, 29	1, 2, 7, 9, 11, 14, 18, 20, 25, 26, 27, 28, 29, 30, 31	1, 2, 7, 8, 10, 11, 14, 18, 20, 24, 25, 26

Tabla 13. Métodos de extracción de características para los cuales existe una similitud en los patrones parenquimatosos de las mamografías reales con las mamografías sintéticas para el análisis estadístico con distribución no normal de media muestral (17) y coeficiente de correlación (5)(6)(7).

## 5. CONCLUSIONES

Para determinar el grado de similitud entre patrones parenquimatosos extraídos de imágenes de mamografías reales con imágenes de mamografías sintéticas, se utilizó un estudio estadístico inferencial en el cual se analizaron los parámetros poblacionales de la media muestral, el coeficiente de correlación y la distribución acumulativa de las cantidades cuantitativas obtenidas con los diferentes métodos de extracción de características utilizados sobre los patrones parenquimatosos de las mamografías reales y mamografías sintéticas. El análisis realizado sugiere que las medias muestrales de las dos poblaciones difieren en sus valores medios, de igual manera, el análisis indica que no existe correlación entre los valores cuantitativos de las dos poblaciones y respecto a la distribución acumulativa, las poblaciones presentan distribuciones acumulativas diferentes entre ellas. El programa OpenVct fue desarrollado y probado en Estados Unidos, por lo cual los modelos anatómicos utilizados para desarrollar el programa fueron modelos de mamas de mujeres estadounidenses. Por lo tanto, el programa fue realizado bajo unos estándares específicos con el fin de desarrollar una población virtual emparejada con la población real utilizada. Los autores utilizan el parámetro de la densidad para comparar y encontrar semejanzas entre las mamas simuladas y las mamas reales, aunque los análisis realizados exponen que la base de datos simuladas y la base de datos reales son similares, respecto a la densidad [71]. En el presente trabajo se realizó el análisis parenquimatoso sobre la base de datos de mamografías reales y la base de mamografía sintéticas, con el consecuente análisis estadístico de las cantidades cuantitativas obtenidas, fue posible determinar que no existe un grado de similitud entre las mamografías reales y las mamografías sintéticas, respecto a los patrones parenquimatosos.

Con el análisis realizado sobre los grupos de Control y de Casos, se puede concluir que los modelos anatómicos simulados no reflejan las características de los patro-

nes parenquimatosos que permitan generar mamografías con incidencia de cáncer o mamografías sin presencia de cáncer. Es importante anotar que OpenVct tiene limitación para generar modelos sintéticos con volúmenes similares a los volúmenes de las mamografías reales incluidas en este estudio.

## BIBLIOGRAFÍA

- 1 Mario A González-Mariño. “Causas de muerte por cáncer de mama en Colombia”. En: *Revista de Salud Pública* 18 (2016), págs. 344-353 (vid. pág. 14).
- 2 Jennifer G Reeder y Victor G Vogel. “Breast cancer prevention”. En: *Advances in Breast Cancer Management, Second Edition* (2008), págs. 149-164 (vid. pág. 14).
- 3 Jack Cuzick y col. “Tamoxifen for prevention of breast cancer: extended long-term follow-up of the IBIS-I breast cancer prevention trial”. En: *The lancet oncology* 16.1 (2015), págs. 67-75 (vid. pág. 14).
- 4 Steven S Coughlin. “Social determinants of breast cancer risk, stage, and survival”. En: *Breast cancer research and treatment* 177.3 (2019), págs. 537-548 (vid. pág. 14).
- 5 Marie Sundquist, Lars Brudin y Göran Tejler. “Improved survival in metastatic breast cancer 1985–2016”. En: *The Breast* 31 (2017), págs. 46-50 (vid. pág. 14).
- 6 Asma Ghafoor y col. “Trends in breast cancer by race and ethnicity”. En: *CA: a cancer journal for clinicians* 53.6 (2003), págs. 342-355 (vid. pág. 14).
- 7 Bin Guo y col. “Microwave imaging via adaptive beamforming methods for breast cancer detection”. En: *Journal of Electromagnetic Waves and Applications* 20.1 (2006), págs. 53-63 (vid. pág. 14).
- 8 Bianca Mostert y col. “Circulating tumor cells (CTCs): detection methods and their clinical relevance in breast cancer”. En: *Cancer treatment reviews* 35.5 (2009), págs. 463-474 (vid. pág. 14).

- 9 Sharyl J Nass, I Craig Henderson, Joyce C Lashof y col. *Mammography and beyond: developing technologies for the early detection of breast cancer*. National Academy Press Washington, DC, 2001 (vid. pág. 14).
- 10 Ali Zamanian y CJHFE Hardiman. “Electromagnetic radiation and human health: A review of sources and effects”. En: *High Frequency Electronics* 4.3 (2005), págs. 16-26 (vid. pág. 15).
- 11 Francesca Pasi y col. “Effects of single or combined treatments with radiation and chemotherapy on survival and danger signals expression in glioblastoma cell lines”. En: *BioMed research international* 2014 (2014) (vid. pág. 15).
- 12 Andrew DA Maidment. “Virtual clinical trials for the assessment of novel breast screening modalities”. En: *International Workshop on Digital Mammography*. Springer. 2014, págs. 1-8 (vid. pág. 15).
- 13 Ehsan Abadi y col. “Virtual clinical trials in medical imaging: a review”. En: *Journal of Medical Imaging* 7.4 (2020), pág. 042805 (vid. pág. 15).
- 14 Maryellen L Giger, Nico Karssemeijer y Julia A Schnabel. “Breast image analysis for risk assessment, detection, diagnosis, and treatment of cancer”. En: *Annual review of biomedical engineering* 15 (2013), págs. 327-357 (vid. pág. 15).
- 15 Maryellen L Giger, Heang-Ping Chan y John Boone. “Anniversary paper: history and status of CAD and quantitative image analysis: the role of medical physics and AAPM”. En: *Medical physics* 35.12 (2008), págs. 5799-5820 (vid. pág. 15).
- 16 Matthias Elter y Alexander Horsch. “CADx of mammographic masses and clustered microcalcifications: a review”. En: *Medical physics* 36.6Part1 (2009), págs. 2052-2068 (vid. pág. 15).

- 17 Robert M Nishikawa. "Current status and future directions of computer-aided diagnosis in mammography". En: *Computerized Medical Imaging and Graphics* 31.4-5 (2007), págs. 224-235 (vid. pág. 15).
- 18 Said Pertuz y col. "Open framework for mammography-based breast cancer risk assessment". En: *2019 IEEE EMBS International Conference on Biomedical & Health Informatics (BHI)*. IEEE. 2019, págs. 1-4 (vid. págs. 16, 20, 21, 24).
- 19 Célia Freitas Da Cruz. "Automatic analysis of mammography images". En: *Engineering Faculty-Porto University, Oporto, Master in Bioengineering-Biomedical Engineering* (2011) (vid. pág. 17).
- 20 Harold Ellis y Vishy Mahadevan. "Anatomy and physiology of the breast". En: *Surgery (Oxford)* 31.1 (2013), págs. 11-14 (vid. pág. 17).
- 21 Rod R Seeley, Trent D Stephens y Philip Tate. *Anatomy and physiology*. McGraw-Hill, 2008, págs. 1056-1058 (vid. pág. 17).
- 22 Jéssica González Fernández y Carlos E Ugalde Ovaes. "La glándula mamaria, embriología, histología, anatomía y una de sus principales patologías, el cáncer de mama". En: *Revista médica de costa rica y centroamerica* 69.602 (2012), págs. 317-320 (vid. pág. 17).
- 23 Kandace P McGuire. "Breast anatomy and physiology". En: *Breast Disease*. Springer, 2016, págs. 1-14 (vid. pág. 17).
- 24 Joseph D Bronzino. *Biomedical Engineering Handbook 3*. Vol. 3. Taylor & Francis Group, 2006 (vid. págs. 17, 18).
- 25 Martin J Yaffe. "Digital mammography". En: *PACS*. Springer. 2006, págs. 363-371 (vid. págs. 17, 18).

- 26 Rangaraj M Rangayyan, Fabio J Ayres y JE Leo Desautels. "A review of computer-aided diagnosis of breast cancer: Toward the detection of subtle signs". En: *Journal of the Franklin Institute* 344.3-4 (2007), págs. 312-348 (vid. pág. 17).
- 27 Michael G Marmot y col. "The benefits and harms of breast cancer screening: an independent review". En: *British journal of cancer* 108.11 (2013), págs. 2205-2240 (vid. pág. 17).
- 28 Saleem Z Ramadan. "Methods used in computer-aided diagnosis for breast cancer detection using mammograms: a review". En: *Journal of healthcare engineering* 2020 (2020) (vid. pág. 17).
- 29 Olivia Jane Scully y col. "Breast cancer metastasis". En: *Cancer genomics & proteomics* 9.5 (2012), págs. 311-320 (vid. pág. 18).
- 30 Zehra Karapinar Senturk y Resul Kara. "Breast cancer diagnosis via data mining: performance analysis of seven different algorithms". En: *Computer Science & Engineering* 4.1 (2014), pág. 35 (vid. pág. 18).
- 31 Wen-Bin Jian. "Electromagnetic waves". En: (2000) (vid. pág. 18).
- 32 David J Griffiths. *Introduction to electrodynamics*. 2005 (vid. pág. 18).
- 33 Albert Shadowitz. *The electromagnetic field*. Courier Corporation, 2012 (vid. pág. 18).
- 34 Bruno Barufaldi y col. "OpenVCT: a GPU-accelerated virtual clinical trial pipeline for mammography and digital breast tomosynthesis". En: *Medical Imaging 2018: Physics of Medical Imaging*. Vol. 10573. International Society for Optics y Photonics. 2018, pág. 1057358 (vid. págs. 18, 19).
- 35 M De Greef y col. "Accelerated ray tracing for radiotherapy dose calculations on a GPU". En: *Medical physics* 36.9Part1 (2009), págs. 4095-4102 (vid. pág. 20).

- 36 Roberto Bellotti y col. "A completely automated CAD system for mass detection in a large mammographic database". En: *Medical physics* 33.8 (2006), págs. 3066-3075 (vid. pág. 20).
- 37 Guido Van Schie y col. "Mass detection in reconstructed digital breast tomosynthesis volumes with a computer-aided detection system trained on 2D mammograms". En: *Medical physics* 40.4 (2013), pág. 041902 (vid. pág. 20).
- 38 Zhimin Huo y col. "Computerized analysis of digitized mammograms of BRCA1 and BRCA2 gene mutation carriers". En: *Radiology* 225.2 (2002), págs. 519-526 (vid. pág. 24).
- 39 Zhimin Huo y col. "Computerized analysis of mammographic parenchymal patterns for breast cancer risk assessment: feature selection". En: *Medical Physics* 27.1 (2000), págs. 4-12 (vid. pág. 24).
- 40 Germán F Torres y S Pertuz. "Automatic detection of the retroareolar region in x-ray mammography images". En: *VII Latin American Congress on Biomedical Engineering CLAIB 2016, Bucaramanga, Santander, Colombia, October 26th-28th, 2016*. Springer. 2017, págs. 157-160 (vid. pág. 24).
- 41 Jun Wei y col. "Association of computerized mammographic parenchymal pattern measure with breast cancer risk: a pilot case-control study". En: *Radiology* 260.1 (2011), pág. 42 (vid. pág. 24).
- 42 Germán F Torres y S Pertuz. "Automatic detection of the retroareolar region in x-ray mammography images". En: *VII Latin American Congress on Biomedical Engineering CLAIB 2016, Bucaramanga, Santander, Colombia, October 26th-28th, 2016*. Springer. 2017, págs. 157-160 (vid. pág. 24).

- 43 Jun Wei y col. "Association of computerized mammographic parenchymal pattern measure with breast cancer risk: a pilot case-control study". En: *Radiology* 260.1 (2011), pág. 42 (vid. pág. 24).
- 44 Armando Manduca y col. "Texture features from mammographic images and risk of breast cancer". En: *Cancer Epidemiology and Prevention Biomarkers* 18.3 (2009), págs. 837-845 (vid. pág. 24).
- 45 Yuanjie Zheng y col. "Parenchymal texture analysis in digital mammography: a fully automated pipeline for breast cancer risk assessment". En: *Medical physics* 42.7 (2015), págs. 4149-4160 (vid. pág. 24).
- 46 Robert J Nordness. *Epidemiología y bioestadística secretos*. Elsevier, 2006 (vid. págs. 25, 27).
- 47 Gail F Dawson. *Interpretación fácil de la bioestadística: la conexión entre la evidencia y las decisiones médicas*. Elsevier Health Sciences, 2009 (vid. págs. 25, 29).
- 48 Santiago Fernández Fernández y col. *Estadística descriptiva*. Esic Editorial, 2002 (vid. pág. 25).
- 49 Clifford R Blair. *Bioestadística*. Pearson education, 2008 (vid. págs. 25, 30).
- 50 Ricardo Esper y RA Machado. *La investigación en medicina: bases teóricas y prácticas. Elementos de bioestadística*. 2008 (vid. págs. 25, 29).
- 51 Alfredo de Jesús Celis de la Rosa. *Bioestadística*. Inf. téc. 2004 (vid. págs. 25, 27, 29).
- 52 Erik Cobo, Pilar Muñoz y José Antonio González. *Bioestadística para no estadísticos*. Elsevier, 2007 (vid. pág. 25).

- 53 Andrew King y Robert Eckersley. *Statistics for biomedical engineers and scientists: How to visualize and analyze data*. Academic Press, 2019 (vid. págs. 25-27, 29).
- 54 Wayne W Daniel. *Bioestadística*. Limusa, 2003 (vid. págs. 25, 27, 28).
- 55 Francisca Rius Díaz y Francisco Javier Barón López. *Bioestadística*. Thomson, 2005 (vid. pág. 27).
- 56 Luis Prieto Valiente e Inmaculada Herranz Tejedor. *Bioestadística sin dificultades matemáticas*. Ediciones Díaz de Santos, 2010 (vid. pág. 28).
- 57 Andrew King y Robert Eckersley. *Statistics for biomedical engineers and scientists: How to visualize and analyze data*. Academic Press, 2019 (vid. págs. 31, 32, 34).
- 58 Barbara Hazard Munro. *Statistical methods for health care research*. Vol. 1. lippincott williams & wilkins, 2005 (vid. págs. 31, 32).
- 59 Winston Haynes. "Students t-test". En: *Encyclopedia of systems biology* (2013), págs. 2023-2025 (vid. pág. 31).
- 60 Prabhaker Mishra y col. "Application of student's t-test, analysis of variance, and covariance". En: *Annals of cardiac anaesthesia* 22.4 (2019), pág. 407 (vid. pág. 31).
- 61 Antoine Al-Achi. "The students t-test: a brief description". En: *Research & Reviews: Journal of Hospital and Clinical Pharmacy* 5.1 (2019), pág. 1 (vid. pág. 31).
- 62 Thomas W MacFarland y Jan M Yates. "Mann–whitney u test". En: *Introduction to nonparametric statistics for the biological sciences using R*. Springer, 2016, págs. 103-132 (vid. pág. 32).

- 63 S Yue y CY Wang. "Power of the Mann–Whitney test for detecting a shift in median or mean of hydro-meteorological data". En: *Stochastic Environmental Research and Risk Assessment* 16.4 (2002), págs. 307-323 (vid. pág. 32).
- 64 M Mukaka. "Statistics corner: A guide to appropriate use of correlation in medical research". En: *Malawi Med. J.* 24.3 (2012), págs. 69-71 (vid. pág. 34).
- 65 Joseph Lee Rodgers y W Alan Nicewander. "Thirteen ways to look at the correlation coefficient". En: *The American Statistician* 42.1 (1988), págs. 59-66 (vid. pág. 34).
- 66 Per Ahlgren, Bo Jarneving y Ronald Rousseau. "Requirements for a cocitation similarity measure, with special reference to Pearson's correlation coefficient". En: *Journal of the American Society for Information Science and Technology* 54.6 (2003), págs. 550-560 (vid. pág. 34).
- 67 Raul HC Lopes, ID Reid y Peter R Hobson. "The two-dimensional Kolmogorov-Smirnov test". En: (2007) (vid. págs. 35, 62).
- 68 Vance W Berger y YanYan Zhou. "Kolmogorov–smirnov test: Overview". En: *Wiley statsref: Statistics reference online* (2014) (vid. págs. 35, 62).
- 69 Sonja Engmann y Denis Cousineau. "Comparing distributions: the two-sample anderson-darling test as an alternative to the kolmogorov-smirnov test." En: *Journal of applied quantitative methods* 6.3 (2011) (vid. págs. 35, 62).
- 70 Simon (<https://stats.stackexchange.com/users/68268/simon>). *KolmogorovSmirnov test vs. t-test*. Cross Validated. URL:<https://stats.stackexchange.com/q/208517> (version: 2016-04-21). eprint: <https://stats.stackexchange.com/q/208517> (vid. pág. 35).

- 71 Bruno Barufaldi y col. "Developing populations of software breast phantoms for virtual clinical trials". En: *14th International Workshop on Breast Imaging (IWBI 2018)*. Vol. 10718. SPIE. 2018, págs. 481-489 (vid. pág. 69).