

Informe 1 - Análisis Exploratorio

Diagnóstico cáncer de mama

El cáncer de mama es una enfermedad que se ha convertido en una de las principales causas de muerte entre mujeres en todo el mundo, poniendo en relevancia la necesidad de un diagnóstico a tiempo y certero. Para tales efectos, se hace imperioso contar no sólo con el conocimiento, sino, además, con las tecnologías adecuadas que permitan prolongar las probabilidades de vida en aquellas mujeres que padezcan la enfermedad y disminuir el margen de error en su diagnóstico. Es este contexto que nos ha motivado a tomar el proyecto “Diagnóstico de cáncer de mama”.

Dicho lo anterior, el desglose del informe engloba tres apartados; el primero incluye los objetivos, alcances del proyecto y una breve contextualización de la data. Mientras que el apartado 2, contiene una descripción más detallada de los datos a utilizar (cantidad de datos, de dónde vienen, sus distintos atributos y etiquetas, etc.), con el objetivo de tener un mejor entendimiento de ellos. Y finalmente el último apartado contempla un análisis exploratorio del set de datos.

1. Objetivo, Alcance y Contextualización de los Datos

Objetivos

Construir un modelo predictivo para el diagnóstico de cáncer de mama tumores, para tales efectos, se propone desarrollar un modelo de clasificación binaria para predecir entre un tumor Maligno (M) o Benigno (B).

El presente proyecto cuenta con 569 muestras de las cuales 357 corresponden a imágenes benignas y 212 malignas. Cada muestra posee una identificación (ID number), el diagnóstico (diagnosis) y 30 atributos relacionados a 10 características.

Alcance

El alcance del presente informe es realizar proceso exploratorio del set de datos a fin de entenderlos y visualizar si estos nos permitirán construir un modelo lo suficientemente eficiente y con un margen de error muy bajo de la predicción del diagnóstico.

Contextualización de la Data

La base de datos utilizada en este trabajo corresponde al trabajo publicado en la revista científica *Biomedical Image Processing and Biomedical Visualization* durante Julio de 1993 por el Dr. William H. Wolberg del departamento de Cirugía general de la Universidad de Wisconsin, W. Nick Street y Olvi L. Mangasarian del departamento de Ciencias de computación de la Universidad de Wisconsin.

Los atributos contenidos en la base de datos fueron obtenidos en base a técnicas de procesamiento interactivo de imágenes en conjunto con un algoritmo de clasificación inductivo basado en programación lineal, basado en la definición de contornos tal de caracterizar tamaños, formas y texturas celular (**Figura 1**).

Loreto Mora - Oscar Hermosilla - Mauricio Narváez.

Figura 1: Límites Iniciales Aproximados de los Núcleos Celulares

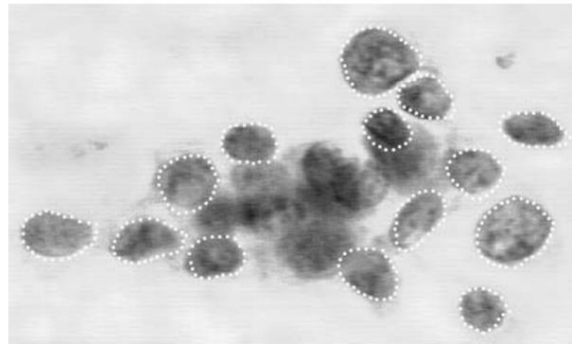


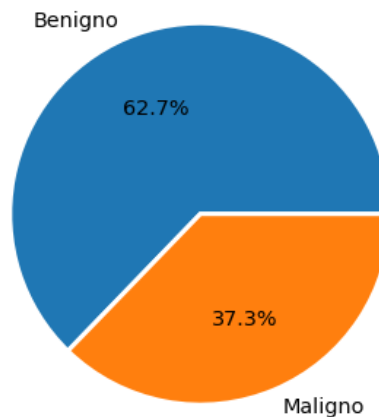
Figure 1: Initial Approximate Boundaries of Cell Nuclei
 The user first draws a rough initial outline of some cell nucleus boundaries. Each outline serves as the initial position for a deformable spline which converges to an accurate boundary of the nucleus.

Fuente: Biomedical Image Processing and Biomedical Visualization
 doi:10.1117/12.148698

2. Descripción de los Datos

La base de datos utilizada para este trabajo contiene 569 muestras, 357 de ellas benignas y 212 malignas (**Figura 2**).

Figura 2: Cantidad de Casos por Diagnóstico
 Cantidad de casos por diagnóstico



Fuente: Elaboración propia

Cada una de las muestras posee una identificación (ID number), el diagnóstico (diagnosis) y 30 atributos relacionados a 10 características, las cuales son:

1. **Radius:** Media de las distancias desde el centro hacia puntos del perímetro.
2. **Texture:** Desviación estándar de las intensidades en escala de grises.
3. **Perimeter:** Perímetro
4. **Area:** Área
5. **Smoothness:** Variación local de largo de los radios.
6. **Compactness:** $(\text{Perímetro})^2 / \text{Área} - 1.0$

Loreto Mora - Oscar Hermosilla - Mauricio Narváez.

7. **Concavity**: severidad de secciones cóncavas en la delineación.
8. **Concave points**: Número de secciones cóncavas en la delineación.
9. **Symmetry**.: Simetría
10. **Fractal dimension**: Coastline approximation" – 1

Atributos del Dataset

En la base de datos, y como producto del experimento científico del cual proviene la información (*Nuclear features extraction for breast tumor diagnosis*¹), figuran mediciones de 10 atributos base, registrando valores promedio de las mediciones (mean), medición del valor más alto (worst) y el error estándar de las mediciones de cada atributo (se), los que se describen a continuación:

Tabla 1: Descripción de 10 Atributos Base

	Promedio (mean)	Valor mas alto (worst)	Error estándar (se)
1	radius_mean	radius_worst	radius_se
2	texture_mean	texture_worst	texture_se
3	perimeter_mean	perimeter_worst	perimeter_se
4	area_mean	area_worst	area_se
5	smoothness_mean	smoothness_worst	smoothness_se
6	compactness_mean	compactness_worst	compactness_se
7	concavity_mean	concavity_worst	concavity_se
8	concave points_mean	concave points_worst	concave points_se
9	symmetry_mean	symmetry_worst	symmetry_se
10	fractal_dimension_mean	fractal_dimension_worst	fractal_dimension_se

Fuente: Elaboración propia

Al obtener una descripción rápida de los atributos disponibles, lo primero que podemos identificar es:

- Existencia de atributo “*Unnamed: 32*” que no contiene valores, por lo tanto, será eliminada de la base de dato.
- Atributo “*ID*” que no aporta información para nuestro modelo, por lo tanto, será eliminada de la base de dato.
- Los 31 atributos restantes no tienen valores nulos y son de tipo numéricos

¹ Street, W. N., Wolberg, W. H., & Mangasarian, O. L. (1993). .Biomedical Image Processing and Biomedical Visualization. doi:10.1117/12.148698

Loreto Mora - Oscar Hermosilla - Mauricio Narváez.

Luego de realizar la limpieza de datos correspondiente, nos quedamos con el siguiente Dataset:

Tabla 2: Dataset post Procesamiento de la Data (Limpieza)

#	Atributo	Cantidad Nulos	Tipo
0	diagnosis	569 non-null	object
1	radius_mean	569 non-null	float64
2	texture_mean	569 non-null	float64
3	perimeter_mean	569 non-null	float64
4	area_mean	569 non-null	float64
5	smoothness_mean	569 non-null	float64
6	compactness_mean	569 non-null	float64
7	concavity_mean	569 non-null	float64
8	concave points_mean	569 non-null	float64
9	symmetry_mean	569 non-null	float64
10	fractal_dimension_mean	569 non-null	float64
11	radius_se	569 non-null	float64
12	texture_se	569 non-null	float64
13	perimeter_se	569 non-null	float64
14	area_se	569 non-null	float64
15	smoothness_se	569 non-null	float64
16	compactness_se	569 non-null	float64
17	concavity_se	569 non-null	float64
18	concave points_se	569 non-null	float64
19	symmetry_se	569 non-null	float64
20	fractal_dimension_se	569 non-null	float64
21	radius_worst	569 non-null	float64
22	texture_worst	569 non-null	float64
23	perimeter_worst	569 non-null	float64
24	area_worst	569 non-null	float64
25	smoothness_worst	569 non-null	float64
26	compactness_worst	569 non-null	float64
27	concavity_worst	569 non-null	float64
28	concave points_worst	569 non-null	float64
29	symmetry_worst	569 non-null	float64
30	fractal_dimension_worst	569 non-null	float64

Fuente: Elaboración propia

Tabla 3: Resumen de Características del Dataset

Cantidad observaciones	569
Característica atributos	Numérico
Número de Atributos	30
Datos faltantes	No
Modelo	Clasificación

Fuente: Elaboración propia

Loreto Mora - Oscar Hermosilla - Mauricio Narváez.

3. Análisis Exploratorio de los Datos

Matriz de correlación

Mediante la siguiente matriz de correlación podremos identificar que atributos nos aportará más información, en función de que tan correlacionadas están entre ellas, por lo tanto, si el valor de correlación es alto indica que los atributos miden la misma característica de la imagen.

Figura 3: Matriz de correlación entre Atributos

radius_mean	100%	32%	100%	99%	17%	51%	68%	82%	15%	-31%	68%	10%	67%	74%	-22%	21%	19%	38%	10%	-4%	97%	30%	97%	94%	12%	41%	53%	74%	16%	1%
texture_mean	-32%	100%	33%	32%	-2%	24%	30%	29%	7%	-8%	28%	39%	28%	26%	1%	19%	14%	16%	1%	5%	35%	91%	36%	34%	8%	28%	30%	30%	11%	12%
perimeter_mean	100%	33%	100%	99%	21%	56%	72%	85%	18%	-26%	69%	-9%	69%	74%	-20%	25%	23%	41%	-8%	-1%	97%	30%	97%	94%	15%	46%	56%	77%	19%	5%
area_mean	-99%	32%	99%	100%	18%	50%	69%	82%	15%	-28%	73%	-7%	73%	80%	-17%	21%	21%	37%	-7%	-2%	96%	29%	96%	96%	12%	39%	51%	72%	14%	0%
smoothness_mean	-17%	-2%	21%	18%	100%	66%	52%	55%	56%	58%	30%	7%	30%	25%	33%	32%	25%	38%	20%	28%	21%	4%	24%	21%	81%	47%	43%	50%	39%	50%
compactness_mean	51%	24%	56%	50%	66%	100%	88%	83%	60%	57%	50%	5%	55%	46%	14%	74%	57%	64%	23%	51%	54%	25%	59%	51%	57%	87%	82%	82%	51%	69%
concavity_mean	-68%	30%	72%	69%	52%	88%	100%	92%	50%	34%	63%	8%	66%	62%	10%	67%	69%	68%	18%	45%	69%	30%	73%	68%	45%	75%	88%	86%	41%	51%
concave points_mean	-82%	29%	85%	82%	55%	83%	92%	100%	46%	17%	70%	2%	71%	69%	3%	49%	44%	62%	10%	26%	83%	29%	86%	81%	45%	67%	75%	91%	38%	37%
symmetry_mean	-15%	7%	18%	15%	56%	60%	50%	46%	100%	48%	30%	13%	31%	22%	19%	42%	34%	39%	45%	33%	19%	9%	22%	18%	43%	47%	43%	43%	70%	44%
fractal_dimension_mean	-31%	-8%	-26%	-28%	58%	57%	34%	17%	48%	100%	0%	16%	4%	-9%	40%	56%	45%	34%	35%	69%	-25%	-5%	-21%	-23%	50%	46%	35%	18%	33%	77%
radius_se	-68%	28%	69%	73%	30%	50%	63%	70%	30%	0%	100%	21%	97%	95%	16%	36%	33%	51%	24%	23%	72%	19%	72%	75%	14%	29%	38%	53%	9%	5%
texture_se	-10%	39%	-9%	-7%	7%	5%	8%	2%	13%	16%	21%	100%	22%	11%	40%	23%	19%	23%	41%	28%	-11%	41%	-10%	-8%	-7%	-9%	-7%	-12%	-13%	-5%
perimeter_se	-67%	28%	69%	73%	30%	55%	66%	71%	31%	4%	97%	22%	100%	94%	15%	42%	36%	56%	27%	24%	70%	20%	72%	73%	13%	34%	42%	55%	11%	9%
area_se	-74%	26%	74%	80%	25%	46%	62%	69%	22%	-9%	95%	11%	94%	100%	8%	28%	27%	42%	13%	13%	76%	20%	76%	81%	13%	28%	39%	54%	7%	2%
smoothness_se	-22%	1%	-20%	-17%	33%	14%	10%	3%	19%	40%	16%	40%	15%	8%	100%	34%	27%	33%	41%	43%	-23%	-7%	-22%	-18%	31%	-6%	-6%	-10%	-11%	10%
compactness_se	-21%	19%	25%	21%	32%	74%	67%	49%	42%	56%	36%	23%	42%	28%	34%	100%	80%	74%	39%	80%	20%	14%	26%	20%	23%	68%	64%	48%	28%	59%
concavity_se	-19%	14%	23%	21%	25%	57%	69%	44%	34%	45%	33%	19%	36%	27%	27%	80%	100%	77%	31%	73%	19%	10%	23%	19%	17%	48%	66%	44%	20%	44%
concave points_se	-38%	16%	41%	37%	38%	64%	68%	62%	39%	34%	51%	23%	56%	42%	33%	74%	77%	100%	31%	61%	36%	9%	39%	34%	22%	45%	55%	60%	14%	31%
symmetry_se	-10%	1%	-8%	-7%	20%	23%	18%	10%	45%	35%	24%	41%	27%	13%	41%	39%	31%	31%	100%	37%	-13%	-8%	-10%	-11%	-1%	6%	4%	-3%	39%	8%
fractal_dimension_se	-4%	5%	-1%	-2%	28%	51%	45%	26%	33%	69%	23%	28%	24%	13%	43%	80%	73%	61%	37%	100%	-4%	-0%	-0%	-2%	17%	39%	38%	22%	11%	59%
radius_worst	-97%	35%	97%	96%	21%	54%	69%	83%	19%	-25%	72%	-11%	70%	76%	-23%	20%	19%	36%	-13%	-4%	100%	36%	99%	98%	22%	48%	57%	79%	24%	9%
texture_worst	-30%	91%	30%	29%	4%	25%	30%	29%	9%	-5%	19%	41%	20%	20%	-7%	14%	10%	9%	-8%	-0%	36%	100%	37%	35%	23%	36%	37%	36%	23%	22%
perimeter_worst	-97%	36%	97%	96%	24%	59%	73%	86%	22%	-21%	72%	-10%	72%	76%	-22%	26%	23%	39%	-10%	-0%	99%	37%	100%	98%	24%	53%	62%	82%	27%	14%
area_worst	-94%	34%	94%	96%	21%	51%	68%	81%	18%	-23%	75%	-8%	73%	81%	-18%	20%	19%	34%	-11%	-2%	98%	35%	98%	100%	21%	44%	54%	75%	21%	8%
smoothness_worst	-12%	8%	15%	12%	81%	57%	45%	45%	43%	50%	14%	-7%	13%	13%	31%	23%	17%	22%	-1%	17%	22%	23%	24%	21%	100%	57%	52%	55%	49%	62%
compactness_worst	-41%	28%	46%	39%	47%	87%	75%	67%	47%	46%	29%	-9%	34%	28%	-6%	68%	48%	45%	6%	39%	48%	36%	53%	44%	57%	100%	89%	80%	61%	81%
concavity_worst	-53%	30%	56%	51%	43%	82%	88%	75%	43%	35%	38%	-7%	42%	39%	-6%	64%	66%	55%	4%	38%	57%	37%	62%	54%	52%	89%	100%	86%	53%	69%
concave points_worst	-74%	30%	77%	72%	50%	82%	86%	91%	43%	18%	53%	-12%	55%	54%	-10%	48%	44%	60%	-3%	22%	79%	36%	82%	75%	55%	80%	86%	100%	50%	51%
symmetry_worst	-16%	11%	19%	14%	39%	51%	41%	38%	70%	33%	9%	-13%	11%	7%	-11%	28%	20%	14%	39%	11%	24%	23%	27%	21%	49%	61%	53%	50%	100%	54%
fractal_dimension_worst	-1%	12%	5%	0%	50%	69%	51%	37%	44%	77%	5%	-5%	9%	2%	10%	59%	44%	31%	8%	59%	9%	22%	14%	8%	62%	81%	69%	51%	54%	100%

Fuente: Elaboración propia

De acuerdo con la figura 3, se puede identificar que algunas variables se encuentran con un alto porcentaje de correlación, por lo cual nos quedaremos solo con una de ellas.

- “Radius_mean” tiene una correlación de 100% con “perimeter_mean” y 99% con “area_mean”, 97% con “radius_worst”, 97% con “perimeter_worst”, 94% con “area_worst”, por lo tanto, nos quedaremos solo con “Radius_mean”.

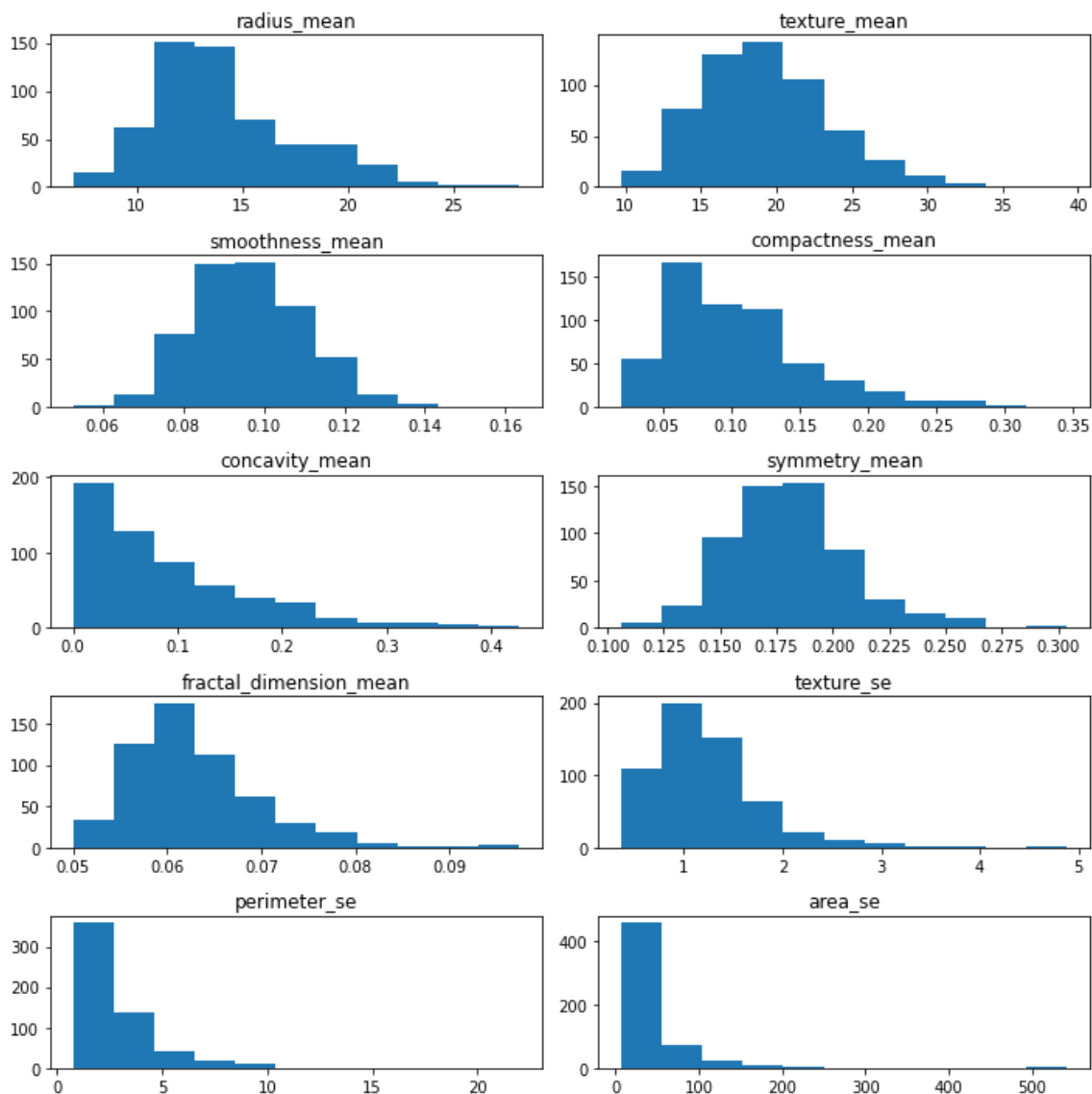
Loreto Mora - Oscar Hermosilla - Mauricio Narváez.

- “Texture_mean” tiene una correlación de 91% con “textura_worst”, por lo cual nos quedaremos solo con “texture_mean”.
- “perimeter_se”, tiene una correlación de 97% con “radius_se” y “áreas_se” tiene una correlación de 95% con “radius_se” por lo cual nos descartaremos “radius_se”.
- “concave points_mean”, tiene una correlación de 92% con “concave_mean”, por lo cual nos quedamos solo con “concave_mean”.

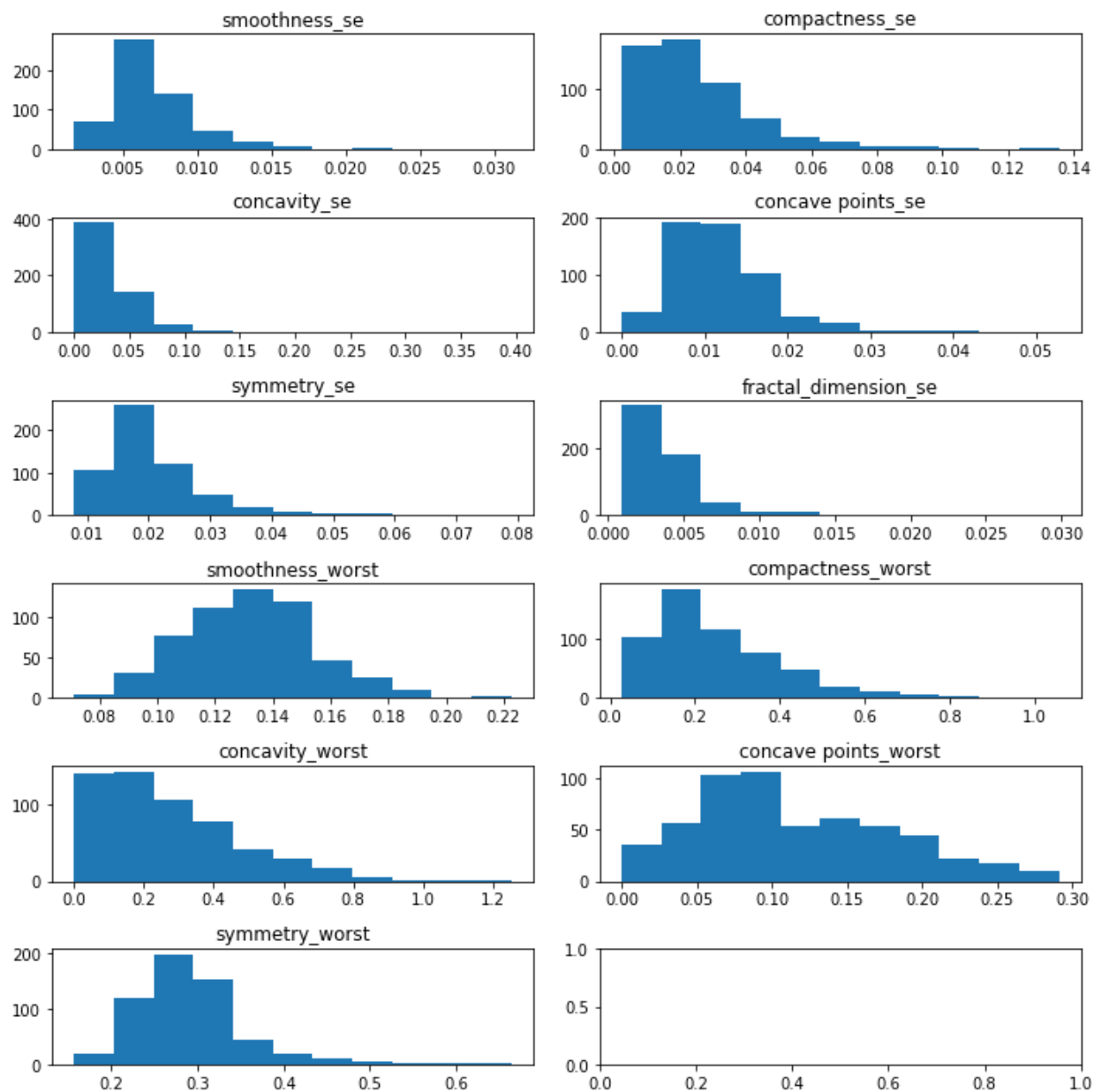
A continuación, presentaremos la distribución que posee cada una de las variables o atributos del Dataset.

Distribución Histogramas

Figura 4: Distribución de cada Atributo



Loreto Mora - Oscar Hermosilla - Mauricio Narváez.



Fuente: Elaboración propia

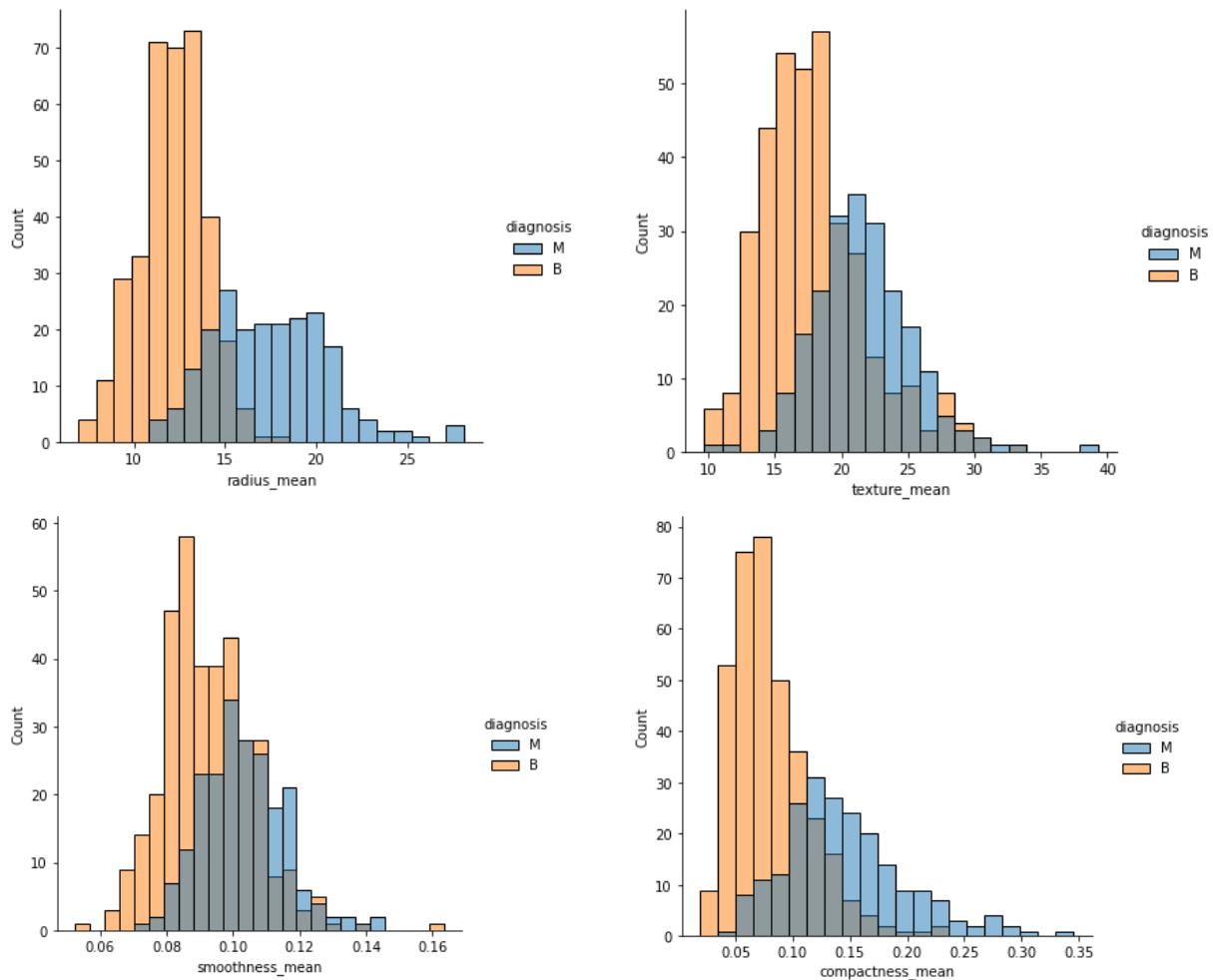
En función de la figura 4, al visualizar las distintas variables o atributos a través de los histogramas, podemos identificar que algunas tienen valores fuera de rango, por lo cual necesitaremos realizar una limpieza sobre ellas para que no generen ruidos al momento de clasificar. Por ejemplo: *concavity_mean*, *texture_se*, *area_se* o *perimeter_se*.

Loreto Mora - Oscar Hermosilla - Mauricio Narváez.

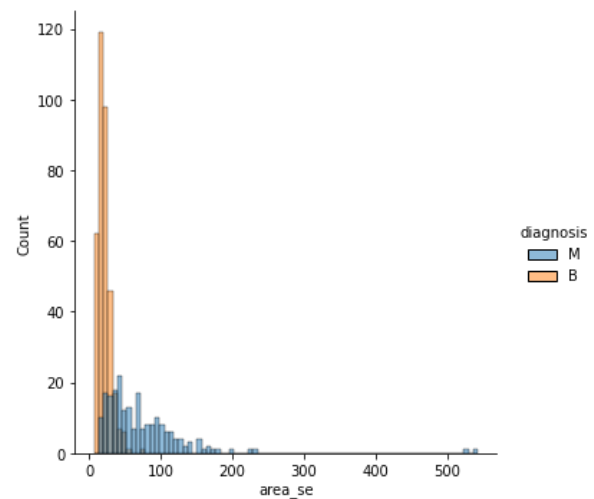
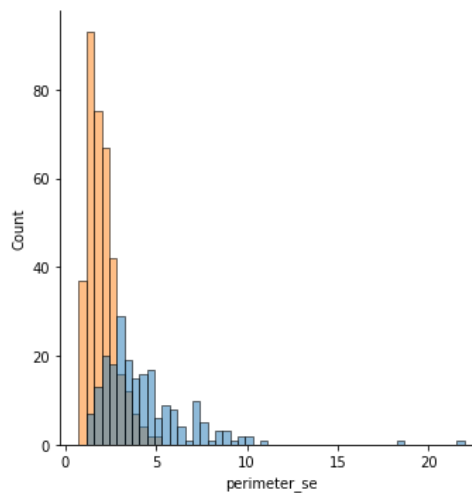
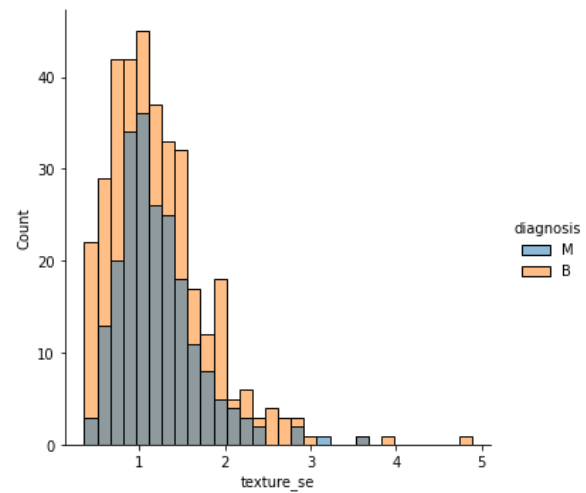
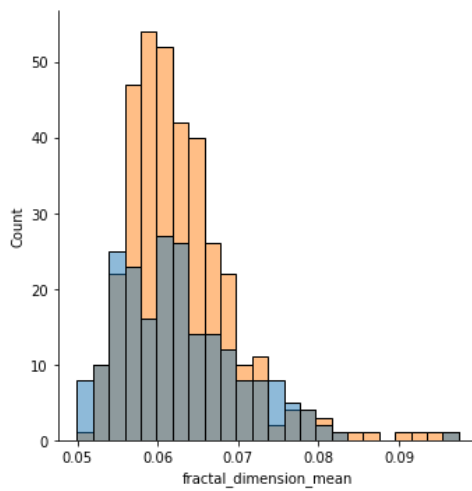
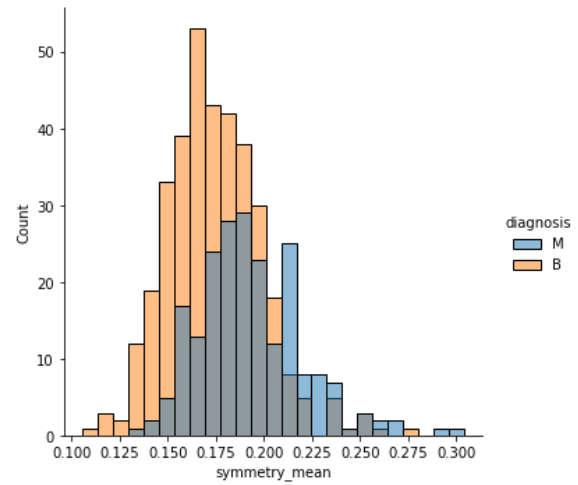
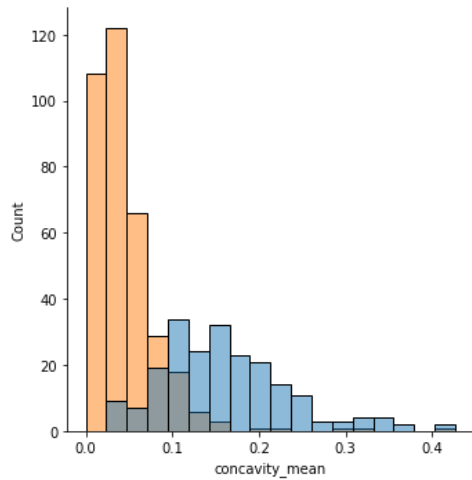
Comparación

A continuación, compararemos cada una de las variables versus nuestra variable dependiente, con la finalidad de identificar si la variable nos ayudará a clasificar de manera correcta si el diagnóstico es Maligno (M) o Benigno (B)

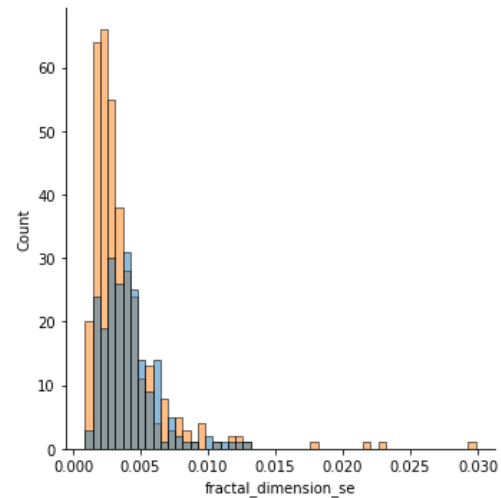
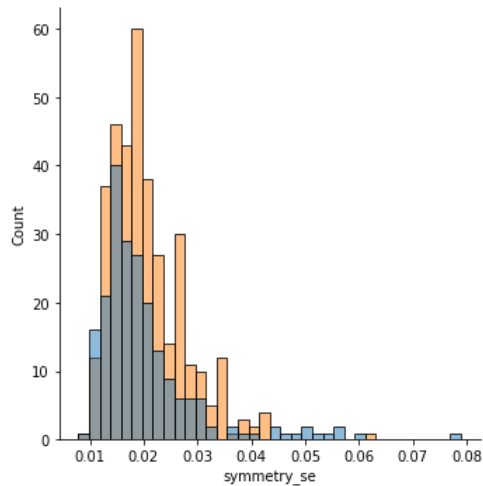
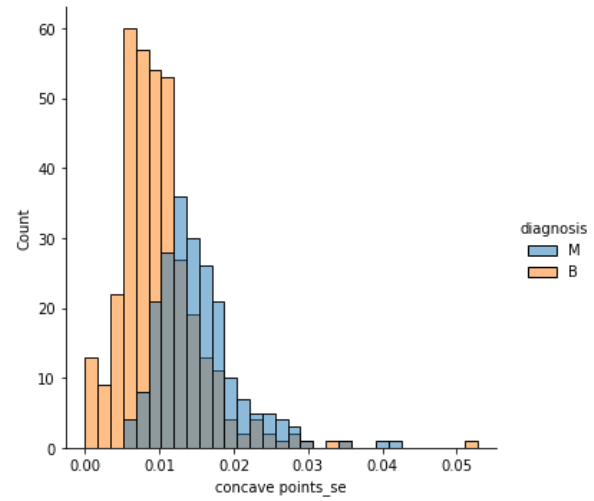
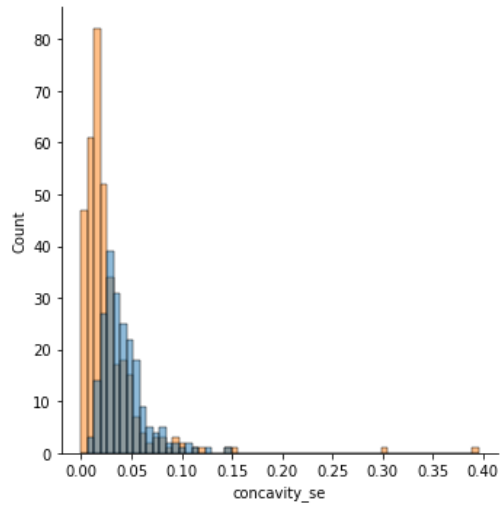
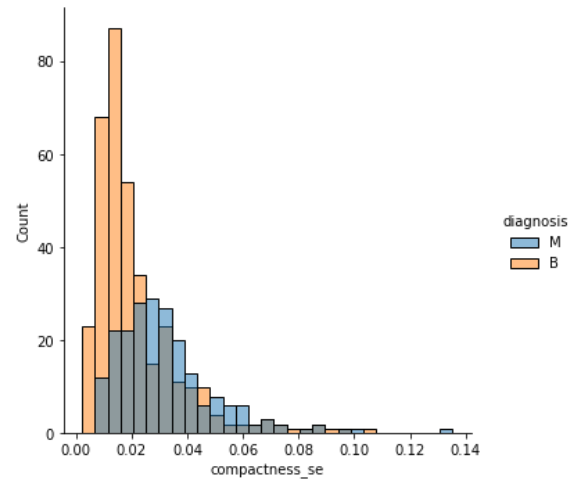
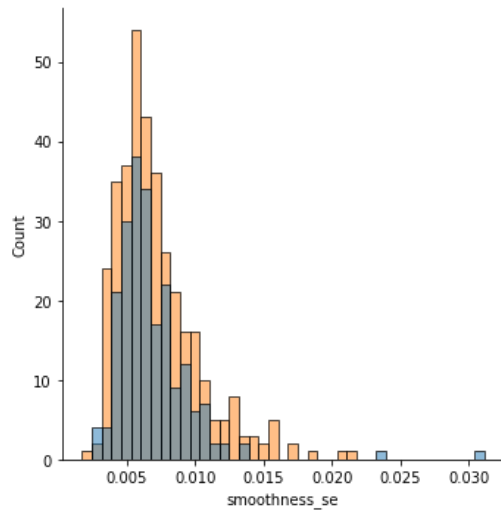
Figura 5: Análisis Comparativo entre los Atributos y la Variable Dependiente



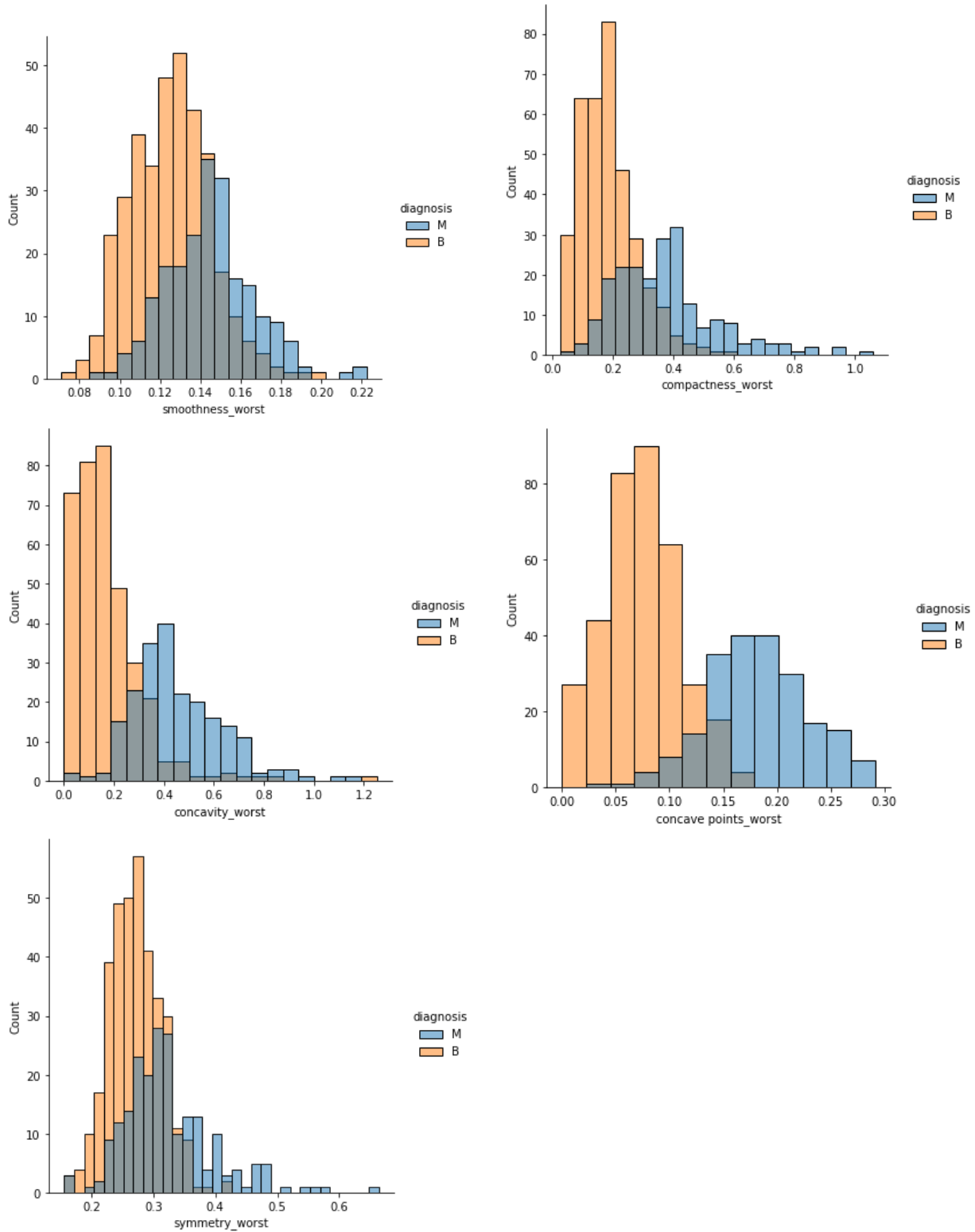
Loreto Mora - Oscar Hermosilla - Mauricio Narváez.



Loreto Mora - Oscar Hermosilla - Mauricio Narváz.



Loreto Mora - Oscar Hermosilla - Mauricio Narváez.



Fuente: Elaboración propia

Loreto Mora - Oscar Hermosilla - Mauricio Narváez.

Al visualizar cada uno de los gráficos presentados en la figura 5, podemos ver que existen variables que nos permitirán clasificar, tales como:

- radius_mean
- compactness_mean
- concavity_mean
- perimeter_se
- area_se
- compactness_worst
- concavity_worst
- concave points_worst
- fractal_dimension_worst

Mientras que otras variables no muestran una preferencia particular de un diagnóstico sobre otro, como son el caso de las siguientes variables:

- texture_mean
- smoothness_mean
- symmetry_mean
- fractal_dimension_mean
- texture_se smoothness_se
- compactness_se, concavity_se
- concave points_se, symmetry_se
- fractal_dimension_se, smoothness_worst
- symmetry_worst