

## A partir de las características celulares observadas, ¿el tumor es maligno o benigno?

El cáncer de mama es el cáncer más común entre las mujeres del mundo. Esta base de datos se centra en un conjunto de datos que contiene características clínicas y biométricas de pacientes con cáncer de mama. En este contexto, se utilizan variables como el radio promedio, la textura media y el área de las células, que son fundamentales para determinar si un tumor es benigno o maligno.

En la primera parte analizaremos la base de datos para limpiarla de valores atípicos y de valores faltantes en alguna columna de cada fila, y determinar la correlación de las variables con la salida esperada.

En la segunda parte amarramos la red neuronal con nuestros datos listos para trabajar y la entrenaremos, luego obtendremos su rendimiento en base a como performa con un subconjunto de la base de datos = "dado para la prueba", para después graficarla.

Por último, haremos la comparación de la red neuronal con la librería Scikit Learn.

## Análisis de los datos

Descripción de las columnas de la base de datos

1. **id (int)**: Identificador único del paciente.
2. **diagnosis (string)**: Diagnóstico del cáncer (M:maligno y B:benigno).
3. **radius\_mean (float)**: Promedio del radio de las células tumorales.
4. **texture\_mean (float)**: Promedio de la textura de las células tumorales.
5. **perimeter\_mean (float)**: Promedio del perímetro de las células tumorales.
6. **area\_mean (float)**: Promedio del área ocupada por las células tumorales.
7. **smoothness\_mean (float)**: Promedio de la suavidad de la superficie de las células tumorales.
8. **compactness\_mean (float)**: Promedio de la compacidad de las células tumorales.
9. **concavity\_mean (float)**: Promedio de la concavidad de las células tumorales.
10. **concave points\_mean (float)**: Promedio de los puntos cóncavos en las células tumorales.
11. **symmetry\_mean (float)**: Promedio de la simetría de las células tumorales.
12. **fractal\_dimension\_mean (float)**: Promedio de la dimensión fractal de las células tumorales.
13. **radius\_se (float)**: Error estándar del radio de las células tumorales.
14. **texture\_se (float)**: Error estándar de la textura de las células tumorales.
15. **perimeter\_se (float)**: Error estándar del perímetro de las células tumorales.
16. **area\_se (float)**: Error estándar del área ocupada por las células tumorales.
17. **smoothness\_se (float)**: Error estándar de la suavidad de las células tumorales.
18. **compactness\_se (float)**: Error estándar de la compacidad de las células tumorales.
19. **concavity\_se (float)**: Error estándar de la concavidad de las células tumorales.
20. **concave points\_se (float)**: Error estándar de los puntos cóncavos en las células tumorales.
21. **symmetry\_se (float)**: Error estándar de la simetría de las células tumorales.
22. **fractal\_dimension\_se (float)**: Error estándar de la dimensión fractal de las células tumorales.
23. **radius\_worst (float)**: Peor radio registrado de las células tumorales.
24. **texture\_worst (float)**: Peor textura registrada de las células tumorales.
25. **perimeter\_worst (float)**: Peor perímetro registrado de las células tumorales.
26. **area\_worst (float)**: Peor área ocupada por las células tumorales.
27. **smoothness\_worst (float)**: Peor suavidad registrada de las células tumorales.
28. **compactness\_worst (float)**: Peor compacidad registrada de las células tumorales.
29. **concavity\_worst (float)**: Peor concavidad registrada de las células tumorales.
30. **concave points\_worst (float)**: Peor cantidad de puntos cóncavos en las células tumorales.
31. **symmetry\_worst (float)**: Peor simetría registrada de las células tumorales.
32. **fractal\_dimension\_worst (float)**: Peor dimensión fractal de las células tumorales.

## ¿Que tipo de variable son?

1. Categórica: diagnosis.
2. Continuas: Todas las demás variables.

## Librerías a usar

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

Guardo en una variable el df e imprimo base de datos
```

```
In [3]: df = pd.read_csv("cancer_de_mama_limpio.csv", index_col=0)
df

Out [3]:
```

		diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	radius_worst	texture_worst
id	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2419	...	25.380	17.33
	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08960	0.07037	0.1812	...	24.990	23.41
	8430903	M	19.69	21.25	130.00	1203.0	0.10960	0.15980	0.19740	0.12790	0.2069	...	23.570	25.53
	8434301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	...	14.910	26.50
	8435402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1809	...	22.540	16.67
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
	826424	M	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1726	...	25.450	26.40
	926862	M	20.13	28.25	131.20	1261.0	0.09870	0.10340	0.14000	0.09701	0.1752	...	23.690	38.25
	926954	M	16.60	28.08	108.30	858.1	0.08485	0.10230	0.09251	0.03962	0.1590	...	18.890	34.12
	927241	M	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397	...	25.740	39.42
	87251	B	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587	...	9.456	30.37

569 rows x 15 columns

## imprimo las primeras 5 filas

```
In [54]: df.head()

Out [54]:
```

		id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	radius_worst	texture_worst
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2419	...	25.380	17.33	
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08960	0.07037	0.1812	...	24.990	23.41	
2	8430903	M	19.69	21.25	130.00	1203.0	0.10960	0.15980	0.19740	0.12790	0.2069	...	23.570	25.53	
3	8434301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	...	14.910	26.50	
4	8435402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1809	...	22.540	16.67	

## imprimo todas las columnas con las que vamos a trabajar

Cabe aclarar que en este análisis no se tuvo en cuenta la columna ID, ya que se considera irrelevante para la evaluación de las características predictivas.

```
In [64]: df.columns

Out [64]:
```

```
Index(['diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean',
       'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean',
       'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
       'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
       'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
       'fractal_dimension_se', 'radius_worst', 'texture_worst',
       'perimeter_worst', 'area_worst', 'smoothness_worst',
       'compactness_worst', 'concavity_worst', 'concave points_worst',
       'symmetry_worst', 'fractal_dimension_worst'],
      dtype='object')
```

## identificamos los tipos de datos de cada columna

```
In [4]: df.dtypes

Out [4]:
```

		diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	radius_se	texture_se	perimeter_se	area_se	smoothness_se	compactness_se	concavity_se	concave points_se	symmetry_se	fractal_dimension_se	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	concave points_worst	symmetry_worst	fractal_dimension_worst	
id	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2419	...	25.380	17.33																			
	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08960	0.07037	0.1812	...	24.990	23.41																			
	8430903	M	19.69	21.25	130.00	1203.0	0.10960	0.15980	0.19740	0.12790	0.2069	...	23.570	25.53																			
	8434301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	...	14.910	26.50																			
	8435402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1809	...	22.540	16.67																			
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...																			
	826424	M	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1726	...	25.450	26.40																			
	926862	M	20.13	28.25	131.20	1261.0	0.09870	0.10340	0.14000	0.09701	0.1752	...	23.690	38.25																			
	926954	M	16.60	28.08	108.30	858.1	0.08485	0.10230	0.09251	0.03962	0.1590	...	18.890	34.12																			
	927241	M	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397	...	25.740	39.42																			
	87251	B	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587	...	9.456	30.37																			

569 rows x 31 columns

## Como nuestra columna objetivo es diagnosis, reemplazamos la letras por valores numericos M:1 y B:0 e imprimimos las primeras 20 para chequear.

```
In [5]: # Maligno: Significa que es agresivo, potencialmente peligroso y puede diseminarse a otras partes del cuerpo.
# En el contexto del cáncer de mama, esto implica que las células cancerosas pueden crecer de manera descontrolada y causar daños en los tejidos cercanos.

# Benigno: Indica que no es agresivo y no representa un riesgo significativo para la salud.
# En el caso de los tumores benignos en el seno, aunque pueden causar molestias o preocupaciones, no se diseminan ni invaden otros tejidos.

df['diagnosis'] = df['diagnosis'].map({'M': 1, 'B': 0})

df.head(20)

Out [5]:
```

		id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	radius_worst	texture_worst
id	842302	M	1	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2419	...	25.38	17.33
	842517	M	1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08960	0.07037	0.1812	...	24.99	23.41
	8430903	M	1	19.69	21.25	130.00	1203.0	0.10960	0.15980	0.19740	0.12790	0.2069	...	23.57	25.53
	8434301	M	1	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	...	14.91	26.50
	8435402	M	1	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1809	...	22.54	16.67
	842796	M	1	12.45	15.70	82.57	471.0	0.12780	0.17000	0.15780	0.08089	0.2087	...	12.47	23.75
	9456202	M	1	28.25	19.96	119.60	2040.0	0.09363	0.10900	0.11270	0.07400	0.1794	...	25.88	27.66
	8458202	M	1	13.71	20.83	90.20	577.9	0.11990	0.16460	0.09366	0.05905	0.2196	...	17.06	28.14
	849881	M	1	33.00	21.82	87.50	519.8	0.12730	0.19320	0.16590	0.09353	0.2350	...	15.49	30.73
	84501001	M	1	12.46	24.04	83.97	475.9	0.11860	0.26660	0.22730	0.05843	0.2030	...	19.19	40.68
	845636	M	1	16.02	23.24	102.70	797.8	0.08206	0.03969	0.03299	0.03323	0.1558	...	18.09	33.88
	8461002	M	1	15.78	17.89	103.60	761.0	0.09710	0.12490	0.09654	0.06606	0.1184	...	20.42	27.66
	846226	M	1	19.17	24.80	132.40	1123.0	0.09740	0.24580	0.20850	0.11180	0.2397	...	20.96	29.94
	846381	M	1	15.85	23.95	103.70	782.7	0.08401	0.20200	0.09938	0.05364	0.1847	...	16.84	27.77
	84667401	M	1	13.73	22.61	93.60	578.3	0.11310	0.12390	0.21280	0.08025	0.2009	...	15.03	32.01
	8479006	M	1	14.54	27.54	96.73	658.8	0.11390	0.15950	0.16390	0.07394	0.2303	...	17.46	37.13
	849002	M	1	14.68	20.13	94.74	684.5	0.09867	0.07200	0.07395	0.05259	0.1586	...	19.07	30.88
	8482001	M	1	16.13	20.68	108.10	798.8	0.11700	0.20220	0.17220	0.10280	0.2164	...	20.96	31.48
	849014	M	1	19.81	22.15	130.00	1260.0	0.09831	0.10270	0.14790	0.09498	0.1582	...	27.32	30.88
	8510426	B	0	13.54	14.36	87.46	566.3	0.09779	0.08129	0.06664	0.04781	0.1885	...	15.11	19.26

20 rows x 31 columns

## Histograma para visualizar datos atípicos

```
In [3]: # Limpio el df
df = df.dropna()
df['diagnosis'] = df['diagnosis'].astype(int)

# Nuestro histogramas de las variables de entrada para visualizar valores atpicos
df.figure(figsize=(20, 17))

Out [3]:
```

```
array([['diagnosis': 'diagnosis'],
       ['diagnosis': 'radius_mean'],
       ['diagnosis': 'texture_mean'],
       ['diagnosis': 'perimeter_mean'],
       ['diagnosis': 'area_mean'],
       ['diagnosis': 'smoothness_mean'],
       ['diagnosis': 'compactness_mean'],
       ['diagnosis': 'concavity_mean'],
       ['diagnosis': 'concave points_mean'],
       ['diagnosis': 'symmetry_mean'],
       ['diagnosis': 'fractal_dimension_mean'],
       ['diagnosis': 'radius_se'],
       ['diagnosis': 'texture_se'],
       ['diagnosis': 'perimeter_se'],
       ['diagnosis': 'area_se'],
       ['diagnosis': 'smoothness_se'],
       ['diagnosis': 'compactness_se'],
       ['diagnosis': 'concavity_se'],
       ['diagnosis': 'concave points_se'],
       ['diagnosis': 'symmetry_se'],
       ['diagnosis': 'fractal_dimension_se'],
       ['diagnosis': 'radius_worst'],
       ['diagnosis': 'texture_worst'],
       ['diagnosis': 'perimeter_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
       ['diagnosis': 'area_worst'],
       ['diagnosis': 'smoothness_worst'],
       ['diagnosis': 'compactness_worst'],
       ['diagnosis': 'concavity_worst'],
       ['diagnosis': 'concave points_worst'],
       ['diagnosis': 'symmetry_worst'],
       ['diagnosis': 'fractal_dimension_worst'],
```



# Red neuronal

La red tendrá dos capas: una capa oculta y una capa de salida.

¿Cuántas neuronas habrá en cada capa?

- Capa de entrada: 9 neuronas
- Capa oculta: 3 neuronas
- Capa de salida: 1 neurona

La capa oculta utilizara la función ReLu y la capa de salida logistic.

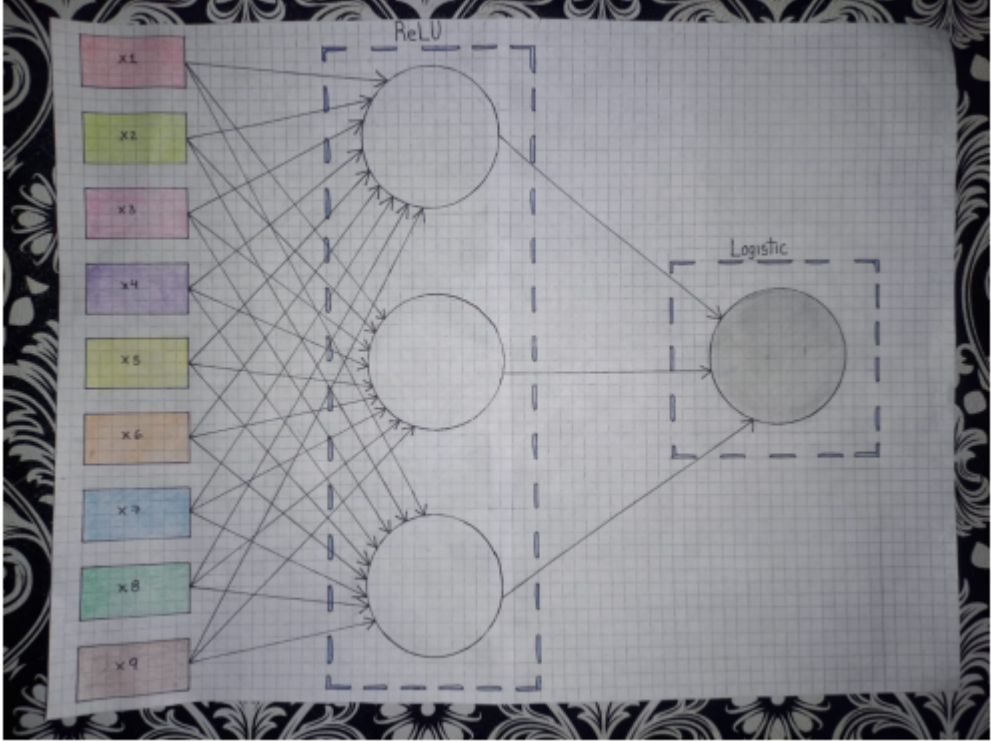
- ReLu: deja igual los valores de entrada positivos y = 0 y a los negativos = X si X > 0 y 0 si X < 0
- Logistic: transforma la salida en un valor entre 0 y 1, ideal para clasificación binaria.

Dibujo de la arquitectura de la red neuronal

```
In [2]: from PIL import Image
import matplotlib.pyplot as plt

# Cargar la imagen
imagen = Image.open('dibujo_red_neuronal.JPG')

# Mostrar la imagen
plt.imshow(imagen)
plt.axis('off')
plt.show()
```



## Red neuronal Libreria

```
In [116.. import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.neural_network import MLPClassifier
from sklearn.preprocessing import StandardScaler

# Cargar el archivo limpio de cáncer de mama
df = pd.read_csv('cancer_de_mama_limpio_top10.csv')

# Extraer variables de entrada (todas las filas, columnas de 2 a 10)
X = df.iloc[:, 2:11].values

# Extraer columna de salida (todas las filas, columna 'diagnosis')
Y = df.iloc[:, 1].values

# Normalizar las características
scaler = StandardScaler()
X = scaler.fit_transform(X)

# Separar los datos de entrenamiento y prueba
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=1/3, random_state=42)

# Crear la red neuronal
nn = MLPClassifier(solver='adam',
                  hidden_layer_sizes=(3, ), # Capas ocultas
                  activation='relu', # Función de activación
                  max_iter=150_000, # Máximo de iteraciones
                  learning_rate_init=0.01) # Tasa de aprendizaje inicial

# Entrenamiento
nn.fit(X_train, Y_train)

# Evaluación del modelo
print("Porcentaje de aciertos con test: ", (nn.score(X_test, Y_test) * 100))
print("Porcentaje de aciertos con train: ", (nn.score(X_train, Y_train) * 100))

Porcentaje de aciertos con test: 95.78947368421052
Porcentaje de aciertos con train: 96.04221635883906
```

## Red neuronal hecha a mano

```
In [4]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# Guarda el archivo en una variable
df2 = pd.read_csv('cancer_de_mama_limpio_top10.csv')

# Extraer variables de entrada (todas las filas, columnas de 2 a 10)
X = df2.iloc[:, 2:11].values

# Extraer columna de salida (todas las filas, columna 'diagnosis')
Y = df2.iloc[:, 1].values

# Normalizar las características con StandardScaler
scaler = StandardScaler()
X = scaler.fit_transform(X)

# Dividir en un conjunto de entrenamiento y uno de prueba
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=1/3)

# Número de registros de entrenamiento
n = X_train.shape[0]
print('Cantidad de filas de entrenamiento: ', n)

# Funciones de activación
relu = lambda x: np.maximum(x, 0)
logistic = lambda x: 1 / (1 + np.exp(-x))

# Semilla para reproducibilidad
np.random.seed(2)

# Construir red neuronal con pesos y sesgos
# inicializados aleatoriamente
w_hidden = (np.random.rand(3, 9) * 2) - 1
w_output = (np.random.rand(1, 3) * 2) - 1

b_hidden = (np.random.rand(3, 1) * 2) - 1
b_output = (np.random.rand(1, 1) * 2) - 1

# Funcion que corre la red neuronal con los datos de entrada para predecir la salida
def forward_prop(X):
    Z1 = w_hidden @ X + b_hidden
    A1 = relu(Z1)
    Z2 = w_output @ A1 + b_output
    A2 = logistic(Z2)
    return Z1, A1, Z2, A2

# Cálculo de precisión
def precision(X, Y):
    test_predictions = forward_prop(X.transpose())[3] # me interesa solo la capa de salida, A2
    test_comparisons = np.equal((test_predictions >= .5).flatten().astype(int), Y)
    accuracy = sum(test_comparisons.astype(int) / X.shape[0])
    print("Porcentaje de aciertos: ", (accuracy*100).round(2))

print('Pre entrenamiento: \n')
print('Test')
precision(X_test, Y_test)
print('Train')
precision(X_train, Y_train)

# Tasa de aprendizaje
L = 0.01

# Derivadas de las funciones de activación
#d_leaky_relu = lambda x: np.where(x > 0, x, 0.01)
d_relu = lambda x: x > 0
d_logistic = lambda x: np.exp(-x) / (1 + np.exp(-x)) ** 2

# Devuelve pendientes para pesos y sesgos
# usando la regla de la cadena
def backward_prop(Z1, A1, Z2, A2, X, Y):
    dC_dA2 = 2 * A2 - 2 * Y
    dA2_dZ2 = d_logistic(Z2)
    dZ2_dA1 = w_output
    dZ2_dW2 = A1
    dZ2_dB2 = 1
    dA1_dZ1 = d_relu(Z1)
    dZ1_dW1 = X
    dZ1_dB1 = 1

    dC_dW2 = dC_dA2 @ dA2_dZ2 @ dZ2_dW2.T

    dC_dB2 = dC_dA2 @ dA2_dZ2 * dZ2_dB2

    dC_dA1 = dC_dA2 @ dA2_dZ2 @ dZ2_dA1

    dC_dW1 = dC_dA1 @ dA1_dZ1 @ dZ1_dW1.T

    dC_dB1 = dC_dA1 @ dA1_dZ1 * dZ1_dB1

    return dC_dW1, dC_dB1, dC_dW2, dC_dB2

# Ejecutar descenso de gradiente
for i in range(150_000):
    # seleccionar aleatoriamente uno de los datos de entrenamiento
    idx = np.random.choice(n, 1, replace=False)
    X_sample = X_train[idx].transpose()
    Y_sample = Y_train[idx]

    # pasar datos seleccionados aleatoriamente a través de la red neuronal
    Z1, A1, Z2, A2 = forward_prop(X_sample)

    # distribuir error a través de la retropropagación
    # y devolver pendientes para pesos y sesgos
    dW1, dB1, dW2, dB2 = backward_prop(Z1, A1, Z2, A2, X_sample, Y_sample)

    # actualizar pesos y sesgos
    w_hidden -= L * dW1
    b_hidden -= L * dB1
    w_output -= L * dW2
    b_output -= L * dB2

# Cálculo de precisión post-entrenamiento
print('Post entrenamiento: \n')
print('Test')
precision(X_test, Y_test)
print('Train')
precision(X_train, Y_train)

Cantidad de filas de entrenamiento: 379
Pre entrenamiento:

Test
Porcentaje de aciertos: 65.26
Train
Porcentaje de aciertos: 61.48
Post entrenamiento:

Test
Porcentaje de aciertos: 96.84
Train
Porcentaje de aciertos: 96.57
```

## Comparación con scikit-learn

La red neuronal manual y la de scikit-learn mostraron un rendimiento muy similar, con una precisión del 95-96% en ambos conjuntos de Test y Train.

In [33]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from tqdm.auto import tqdm

# Cargar el archivo en una variable
df2 = pd.read_csv('cancer_de_mama_limpio_top10.csv')

# Extraer variables de entrada (todas las filas, columnas de 2 a 10)
X = df2.iloc[:, 2:11].values

# Extraer columna de salida (todas las filas, columna 'diagnosis')
Y = df2.iloc[:, 1].values

# Normalizar las características con StandardScaler
scaler = StandardScaler()
X = scaler.fit_transform(X)

# Dividir en un conjunto de entrenamiento y uno de prueba
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=1/3, random_state=42)

# Número de registros de entrenamiento
n = X_train.shape[0]
print('Cantidad de filas de entrenamiento: ', n)

# Funciones de activación
relu = lambda x: np.maximum(x, 0)
logistic = lambda x: 1 / (1 + np.exp(-x))

# Semilla para reproducibilidad
np.random.seed(2)

# Construir red neuronal con pesos y sesgos inicializados aleatoriamente
w_hidden = (np.random.rand(3, 9) * 2) - 1
w_output = (np.random.rand(1, 3) * 2) - 1
b_hidden = (np.random.rand(3, 1) * 2) - 1
b_output = (np.random.rand(1, 1) * 2) - 1

# Funcion que corre la red neuronal con los datos de entrada para predecir la salida
def forward_prop(X):
    Z1 = w_hidden @ X + b_hidden
    A1 = relu(Z1)
    Z2 = w_output @ A1 + b_output
    A2 = logistic(Z2)
    return Z1, A1, Z2, A2

# Tasa de aprendizaje
L = 0.01

# Derivadas de las funciones de activación
d_leaky_relu = lambda x: x > 0
d_logistic = lambda x: np.exp(-x) / (1 + np.exp(-x)) ** 2

# Devuelve pendientes para pesos y sesgos usando la regla de la cadena
def backward_prop(Z1, A1, Z2, A2, X, Y):
    dC_dA2 = 2 * A2 - 2 * Y
    dA2_dZ2 = d_logistic(Z2)
    dZ2_dA1 = w_output
    dZ2_dW2 = A1
    dZ2_dB2 = 1
    dA1_dZ1 = d_relu(Z1)
    dZ1_dW1 = X
    dZ1_dB1 = 1

    dC_dW2 = dC_dA2 @ dA2_dZ2 @ dZ2_dW2.T
    dC_dB2 = dC_dA2 @ dA2_dZ2 * dZ2_dB2
    dC_dA1 = dC_dA2 @ dA2_dZ2 @ dZ2_dA1
    dC_dW1 = dC_dA1 @ dA1_dZ1 @ dZ1_dW1.T
    dC_dB1 = dC_dA1 @ dA1_dZ1 * dZ1_dB1

    return dC_dW1, dC_dB1, dC_dW2, dC_dB2

# Listas para almacenar la precisión
accuracy_train_l = []
accuracy_test_l = []

# Ejecutar descenso de gradiente
for i in tqdm(range(150_000)):
    # Seleccionar aleatoriamente un registro de entrenamiento
    idx = np.random.choice(n, 1, replace=False)
    X_sample = X_train[idx].transpose()
    Y_sample = Y_train[idx]

    # Pasar datos seleccionados aleatoriamente a través de la red neuronal
    Z1, A1, Z2, A2 = forward_prop(X_sample)

    # Retropropagación y cálculo de pendientes para pesos y sesgos
    dW1, dB1, dW2, dB2 = backward_prop(Z1, A1, Z2, A2, X_sample, Y_sample)

    # Actualizar pesos y sesgos
    w_hidden -= L * dW1
    b_hidden -= L * dB1
    w_output -= L * dW2
    b_output -= L * dB2

    # Cálculo de precisión en el conjunto de prueba
    test_predictions = forward_prop(X_test.transpose())[3]
    test_comparisons = np.equal((test_predictions >= .5).flatten().astype(int), Y_test)
    accuracy_test = np.mean(test_comparisons) * 100
    accuracy_test_l.append(accuracy_test)

    # Cálculo de precisión en el conjunto de entrenamiento
    train_predictions = forward_prop(X_train.transpose())[3]
    train_comparisons = np.equal((train_predictions >= .5).flatten().astype(int), Y_train)
    accuracy_train = np.mean(train_comparisons) * 100
    accuracy_train_l.append(accuracy_train)

# Imprimir la precisión final
print("X_test ACCURACY post-training: ", accuracy_test_l[-1])
print("X_train ACCURACY post-training: ", accuracy_train_l[-1])

# Graficar la precisión
fmt_train = {
    'color': 'tab:blue',
    'ls': 'solid',
    'lw': 3,
}

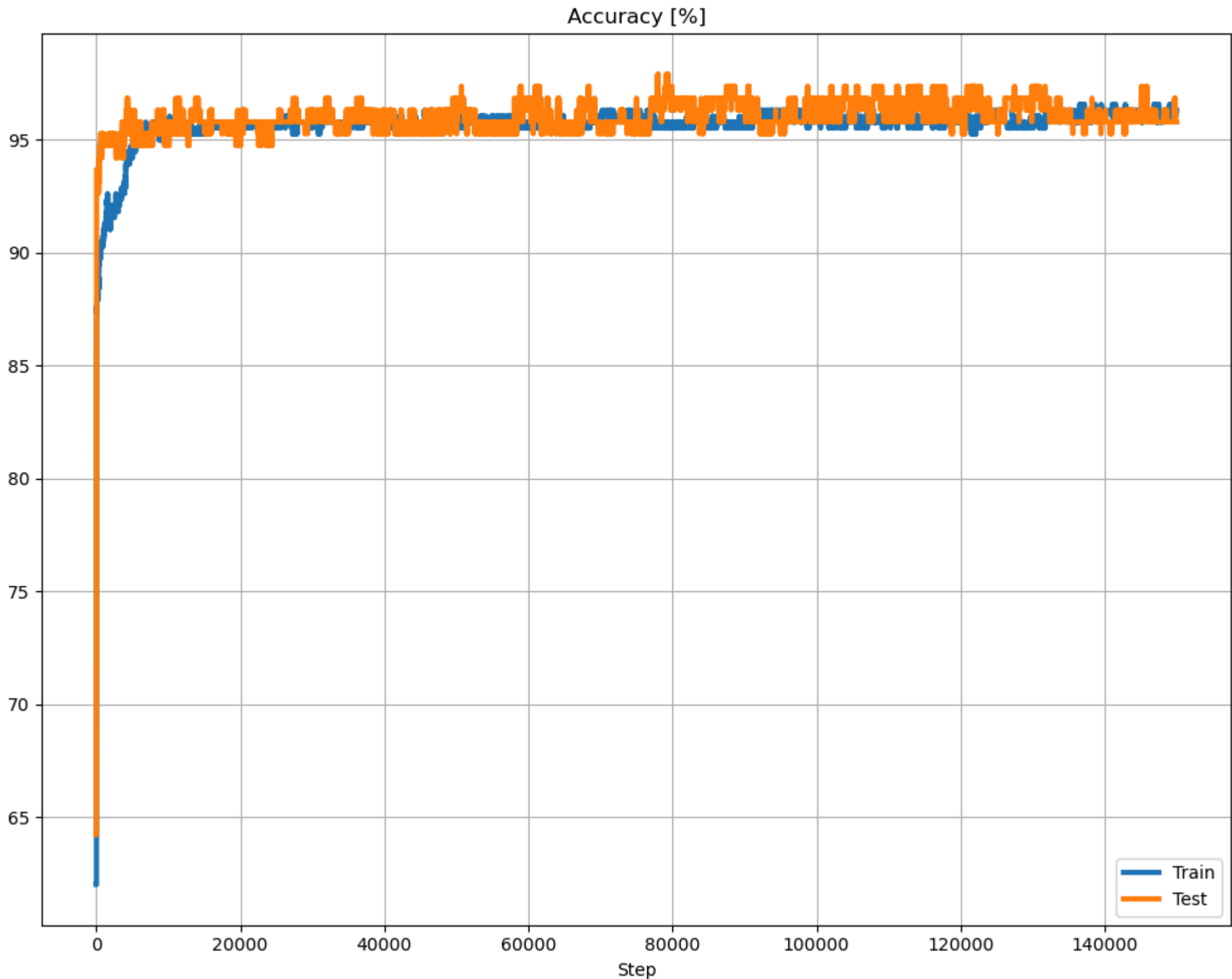
fmt_test = {
    'color': 'tab:orange',
    'ls': 'solid',
    'lw': 3,
}

fig, ax = plt.subplots(1, 1, figsize=(10, 8))

ax.plot(accuracy_train_l, label='Train', **fmt_train)
ax.plot(accuracy_test_l, label='Test', **fmt_test)

ax.grid(which='both')
ax.legend()
ax.set_title('Accuracy [%]')
ax.set_xlabel('Step')
```

Cantidad de filas de entrenamiento: 379  
0%| | 0/150000 [00:00<?, ?it/s]  
X\_test ACCURACY post-training: 95.78947368421052  
X\_train ACCURACY post-training: 96.30606860158312



Luego de desarrollar nuestra Red Neuronal, viendo el porcentaje de aciertos y el gráfico, ¿nuestro modelo presenta overfitting?.

El overfitting es cuando una red una de ajuste demasiado bien a los datos de entrenamiento. Una forma de detectar el sobreajuste es observando si el porcentaje de aciertos en el conjunto de entrenamiento es significativamente mayor que en el conjunto de prueba, especialmente si esta diferencia se amplía al aumentar las épocas de entrenamiento. También es indicativo de sobreajuste si el rendimiento en el conjunto de prueba se estabiliza mientras que el del entrenamiento continúa mejorando.

Existen varias estrategias para prevenir o detectar el sobreajuste en redes neuronales:

- Dropout: Esta técnica consiste en desactivar aleatoriamente algunas neuronas durante el entrenamiento, lo que ayuda a evitar la dependencia de neuronas específicas.
- Reducir la complejidad de la red neuronal: Una forma de lograrlo es disminuir el número de neuronas en las capas ocultas, lo que impide que la red aprenda en exceso sobre datos específicos.
- Ajustar el número de épocas: Limitar la cantidad de iteraciones de entrenamiento puede ayudar a detener el proceso antes de que la red continúe entrenándose sin mejoras en el conjunto de prueba.
- Aumentar la cantidad de datos: Contar con más muestras puede disminuir la probabilidad de sobreajuste, ya que se facilita la generalización.

## Conclusión Final

Desarrollar una red neuronal desde cero nos permitió entender en profundidad cómo funcionan las redes neuronales en cada paso del proceso, desde la inicialización de pesos y sesgos hasta el forward y el backward. Al construir manualmente cada función de activación y sus derivadas, adquirí una comprensión mucho más detallada de cómo los gradientes afectan el ajuste de los pesos y de cómo se optimizan los modelos de aprendizaje. Comparado con el uso la librería de scikit-learn, crear una red neuronal manualmente ofrece ventajas significativas en términos de aprendizaje. Trabajar desde cero ayuda a entender mejor las matemáticas y la lógica detrás del entrenamiento de los modelos. La principal desventaja de este enfoque es que resulta mucho más laborioso y es muy propenso a errores, ya que requiere gestionar manualmente cada cálculo y ajuste de parámetros.