

## Segunda parte: validación

### ⚠ Nota ⚠

Las conclusiones que sacamos en este archivo se encuentran comentadas en el notebook principal, de modo que no afecta al entrenamiento y/o comparación de los modelos futuros.

Revisando la celda que entrena el modelo, vemos que hay un error:

```
12 # Entrenamos el modelo
13 history = network.fit(
14     X_train, train_labels,
15     epochs=100,
16     batch_size=5,
17     validation_data=(X_test, test_labels),
18     shuffle=True
19 )
```

En *validation\_data* deberían estar el segmento de datos de validación, no los de test. Esto se debe hacer así porque sino el modelo conoce los datos de test antes de hacer el propio test (data leakage) y puede afectar a los resultados del modelo.

Ahora, podemos aplicar, principalmente, dos tipos de validación:

- Validación simple.
- Validación cruzada (Cross-validation).

## Validación simple

Podemos utilizar varios tipos de validación simple. En esta práctica he optado por ver estos dos "tipos":

- El método *train\_valid\_test\_split*.
- Segmentar los datos a mano.

## Método *train\_valid\_test\_split*

Es una librería de *fast\_ml* que realiza la división de forma directa, solo le tenemos que indicar:

- X
- target

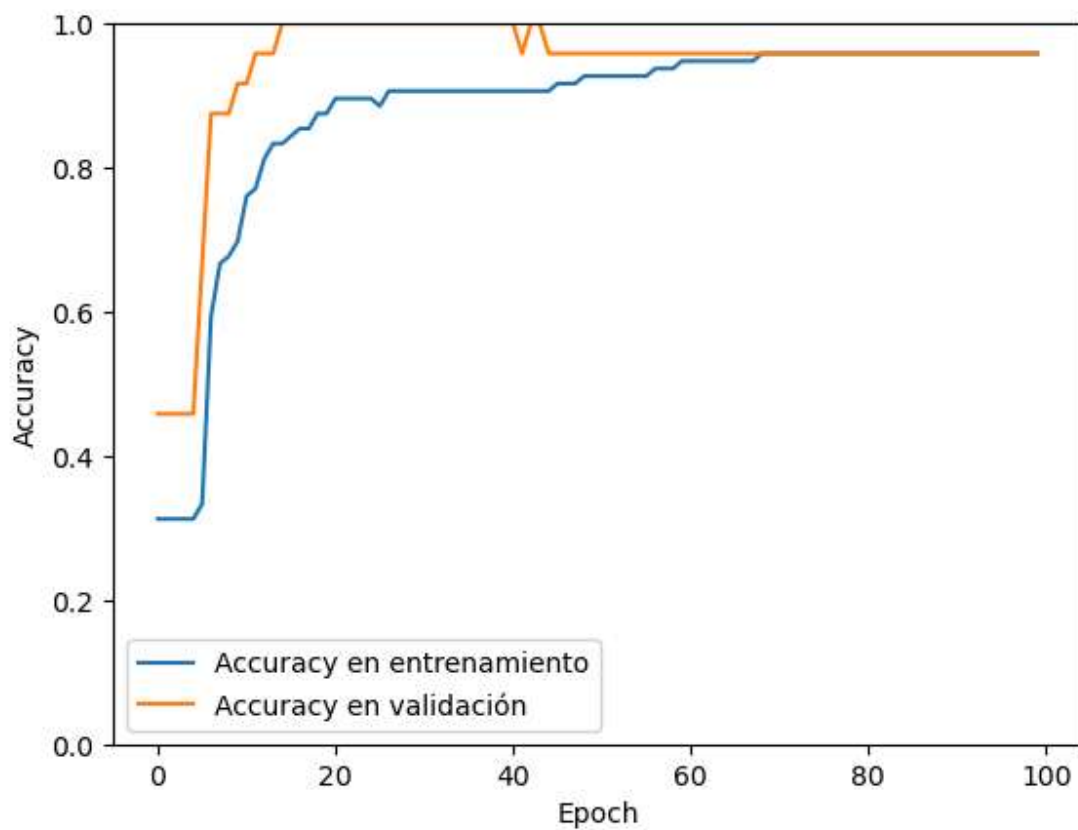
- Tamaño de train
- Tamaño de validation
- Tamaño de test

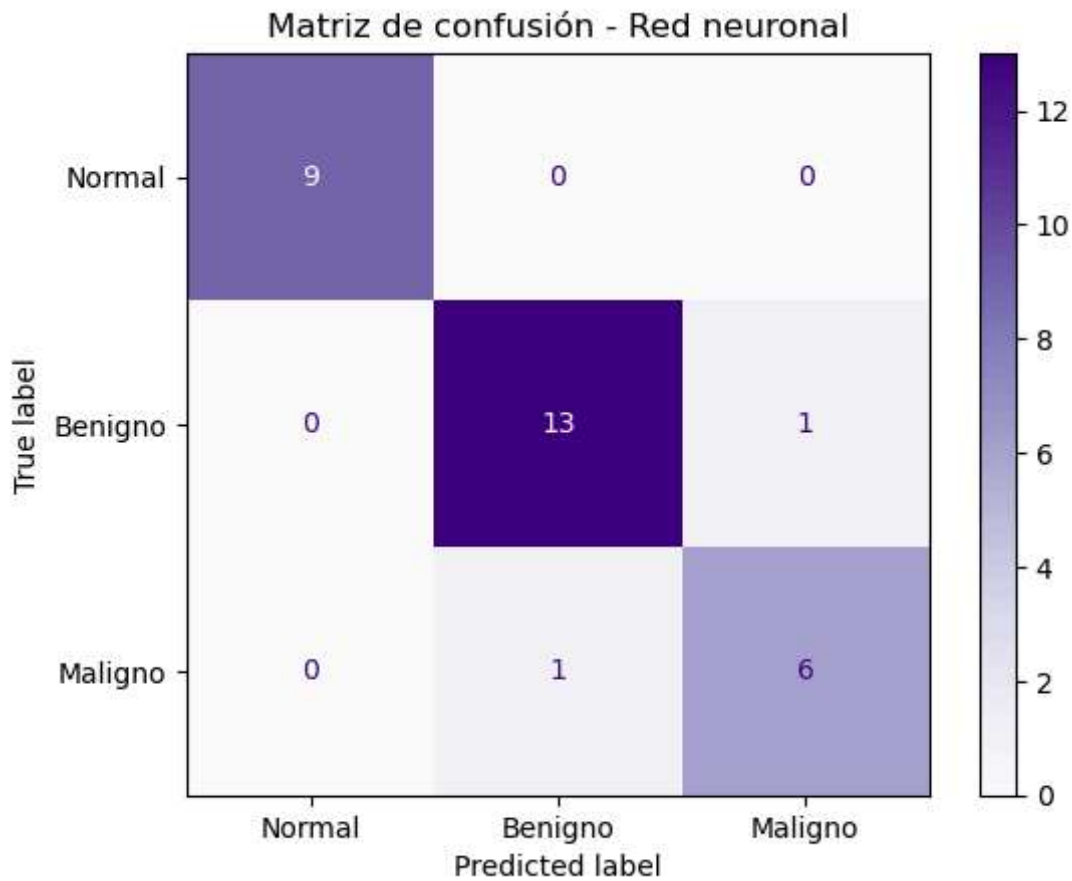
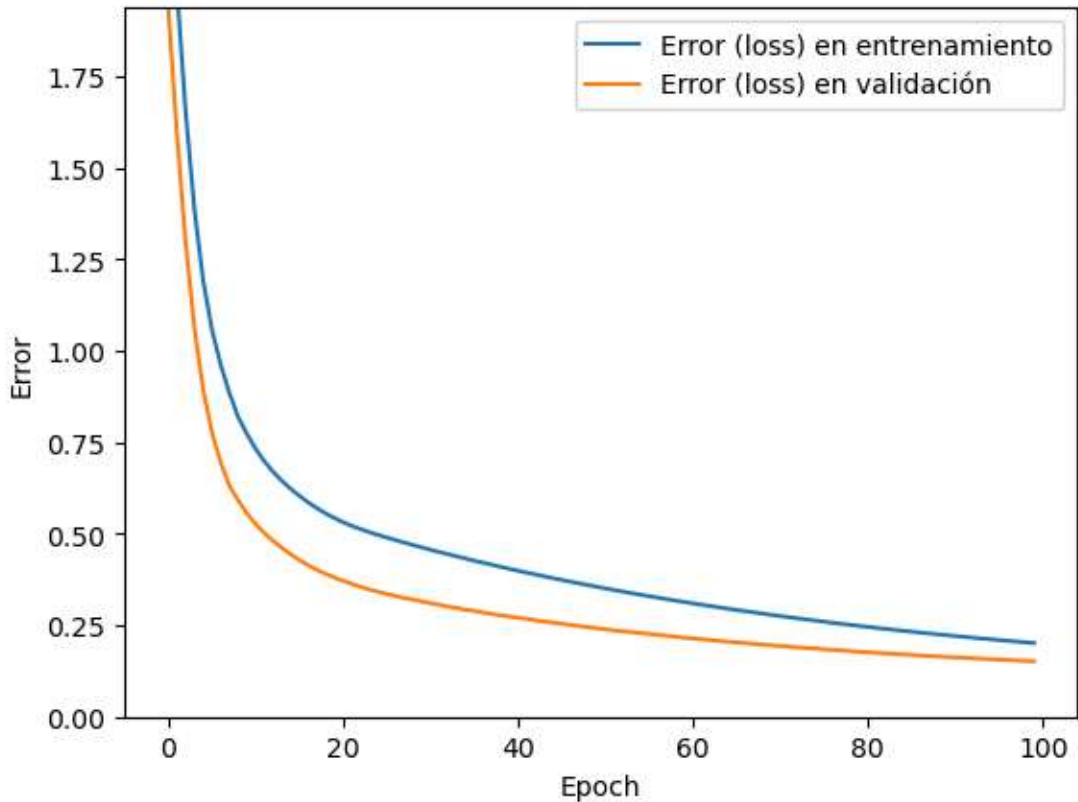
También podemos añadir "method" con el que podemos aleatorizar los datos o no, e indicar una "seed" para poder replicar los resultados.

Debemos usar el df que genera `pandas` y no los array de `numpy` que se usan en el programa original.

Lo malo de usar `train_valid_test_split` es que no existe el parámetro `stratify`, así existe la posibilidad de generar un bloque que solo tenga un tipo de tumor

Usando esta librería hemos obtenido los siguientes resultados con el modelo base:





	precision	recall	f1-score	support
Normal	1.00	1.00	1.00	9
Benigno	0.93	0.93	0.93	14
Maligno	0.86	0.86	0.86	7

accuracy			0.93	30
macro avg	0.93	0.93	0.93	30
weighted avg	0.93	0.93	0.93	30

## Segmentar los datos a mano

Esta opción es más recomendable en mi opinión porque no tienes que importar una librería para una sola función.

Para ello realizamos un primer *train\_test\_split* sobre **X** e **y** (suponiendo que ya están definidas las variables de objetivo) de modo que obtenemos:

- X\_train
- y\_train
- X\_test
- y\_test

Ahora realizamos otro *train\_test\_split* sobre **X\_train** e **y\_train** (preferiblemente con las mismas proporciones que hemos hecho la primera división), de modo que obtenemos:

- X\_train
- y\_train
- X\_val
- y\_val

Y una vez hechas estas dos divisiones nos quedamos con las siguientes variables:

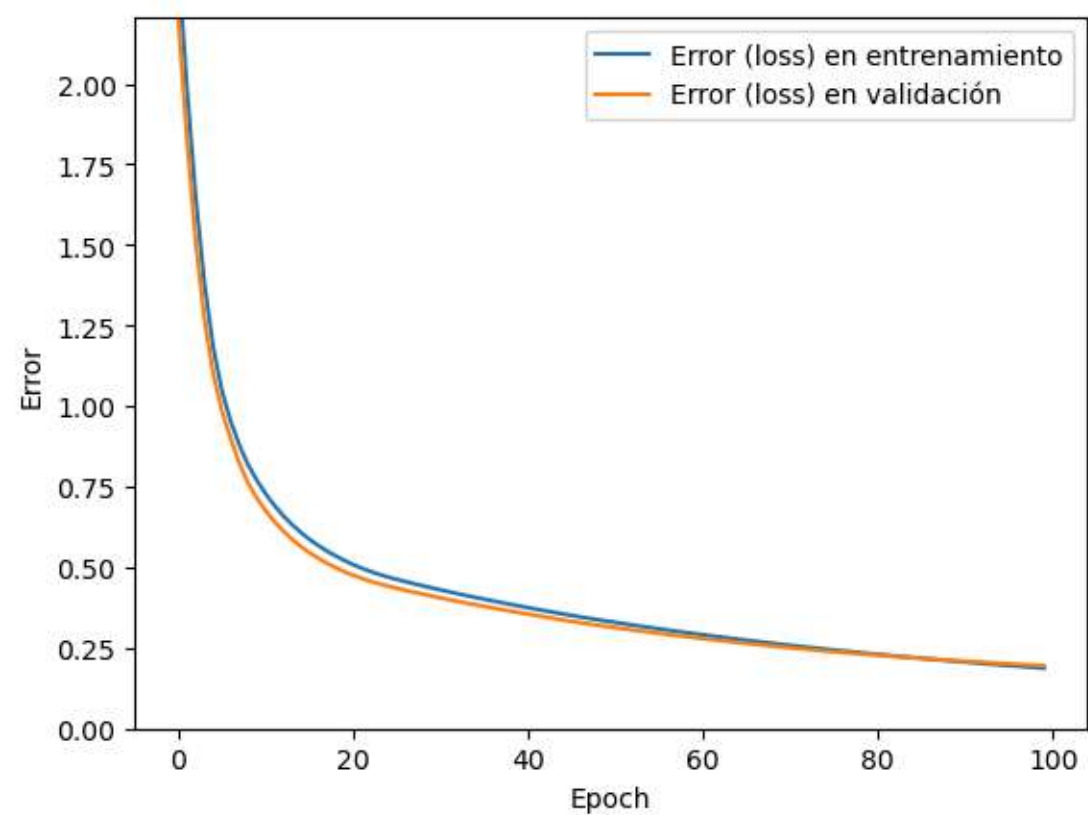
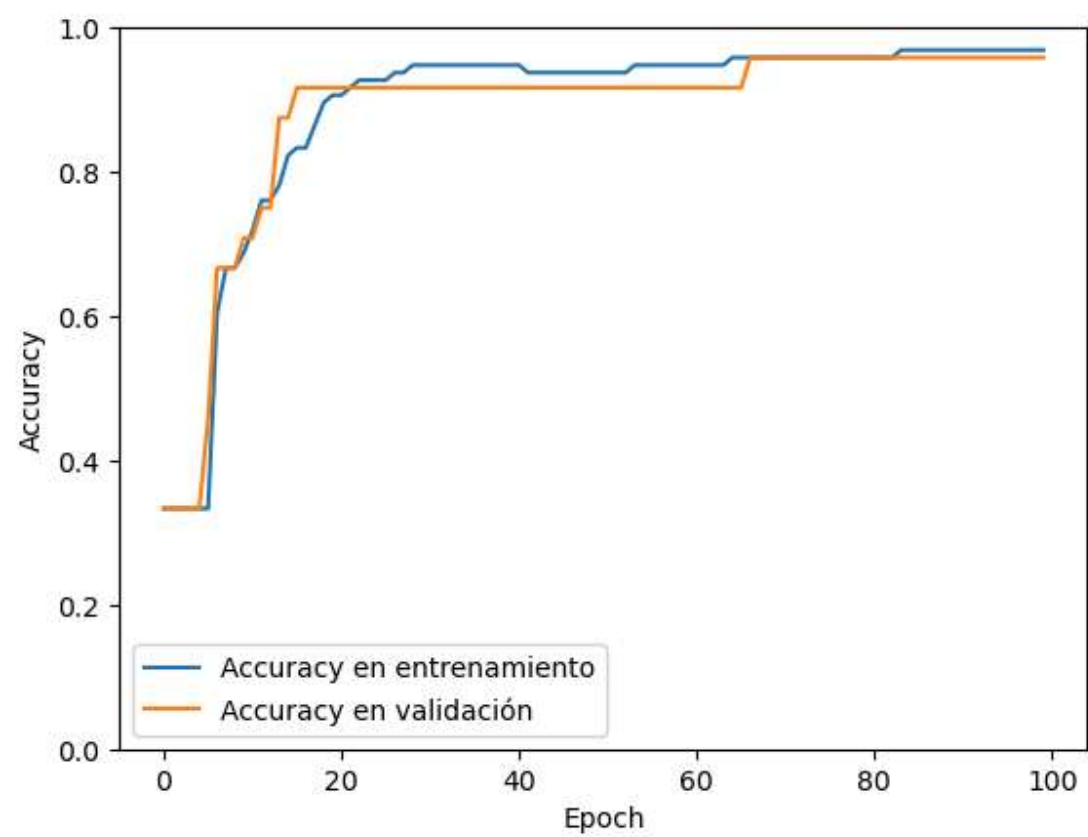
- X\_train
- y\_train
- X\_val
- y\_val
- X\_test
- y\_test

```

5 X_temp, X_test, y_temp, y_test = train_test_split(
6     X, Y,
7     test_size=0.2,
8     stratify=Y,
9     random_state=SEED,
10    shuffle = True
11 )
12
13 X_train, X_val, y_train, y_val = train_test_split(
14     X_temp, y_temp,
15     test_size=0.2,
16     stratify=y_temp,
17     random_state=SEED,
18     shuffle = True
19 )

```

Con la segmentación hecha a mano conseguimos los siguientes resultados:



	precision	recall	f1-score	support
Normal	1.00	1.00	1.00	10
Benigno	1.00	0.80	0.89	10

Maligno	0.83	1.00	0.91	10
accuracy			0.93	30
macro avg	0.94	0.93	0.93	30
weighted avg	0.94	0.93	0.93	30

## Validación cruzada (Cross-validation)

Haremos validación cruzada utilizando StratifiedKFold.

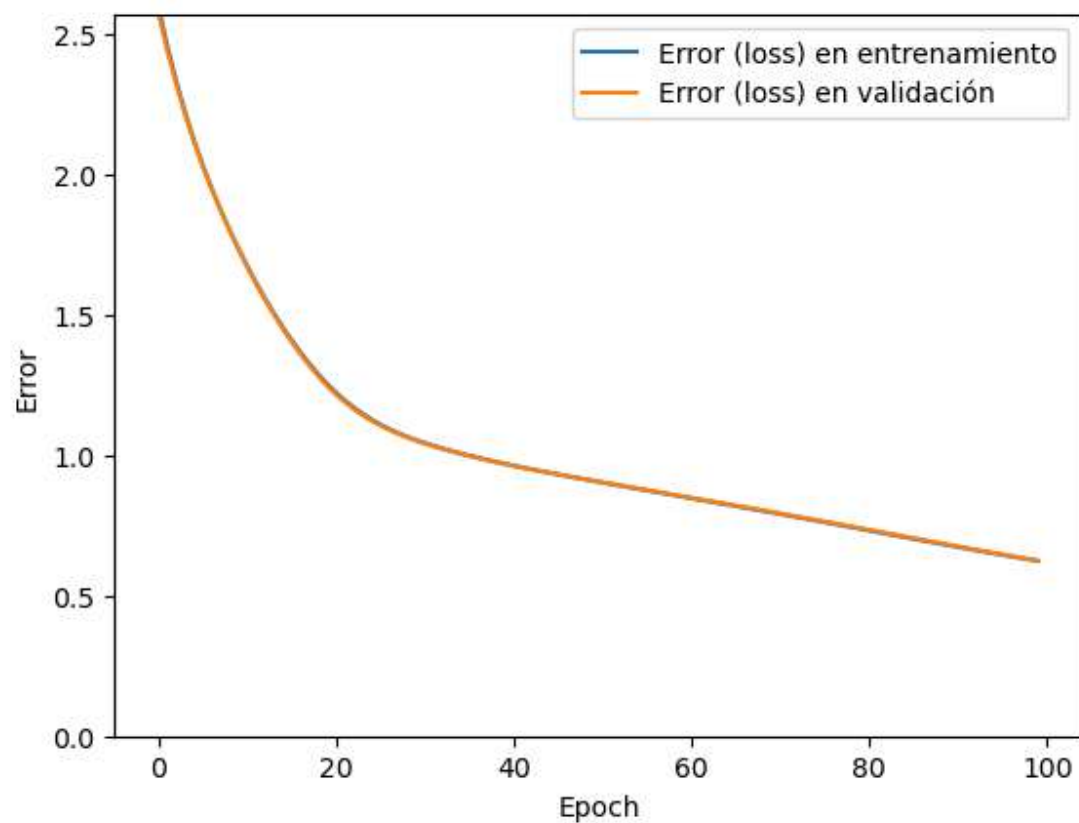
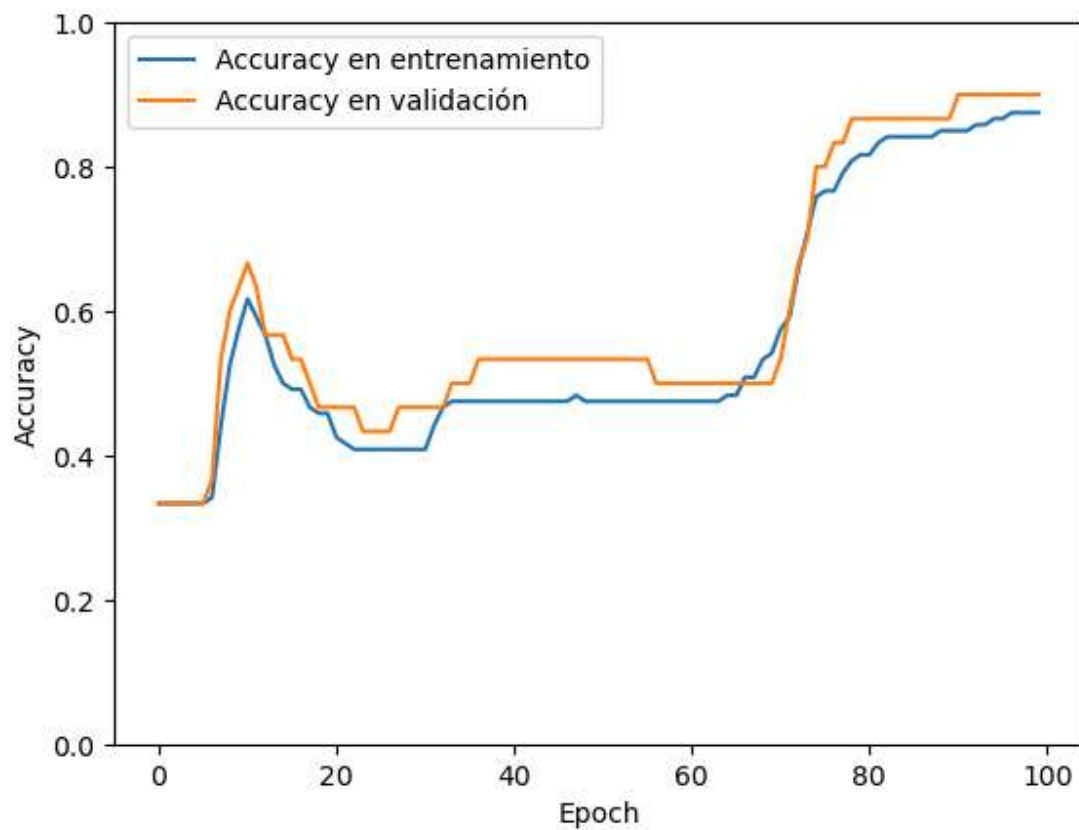
```

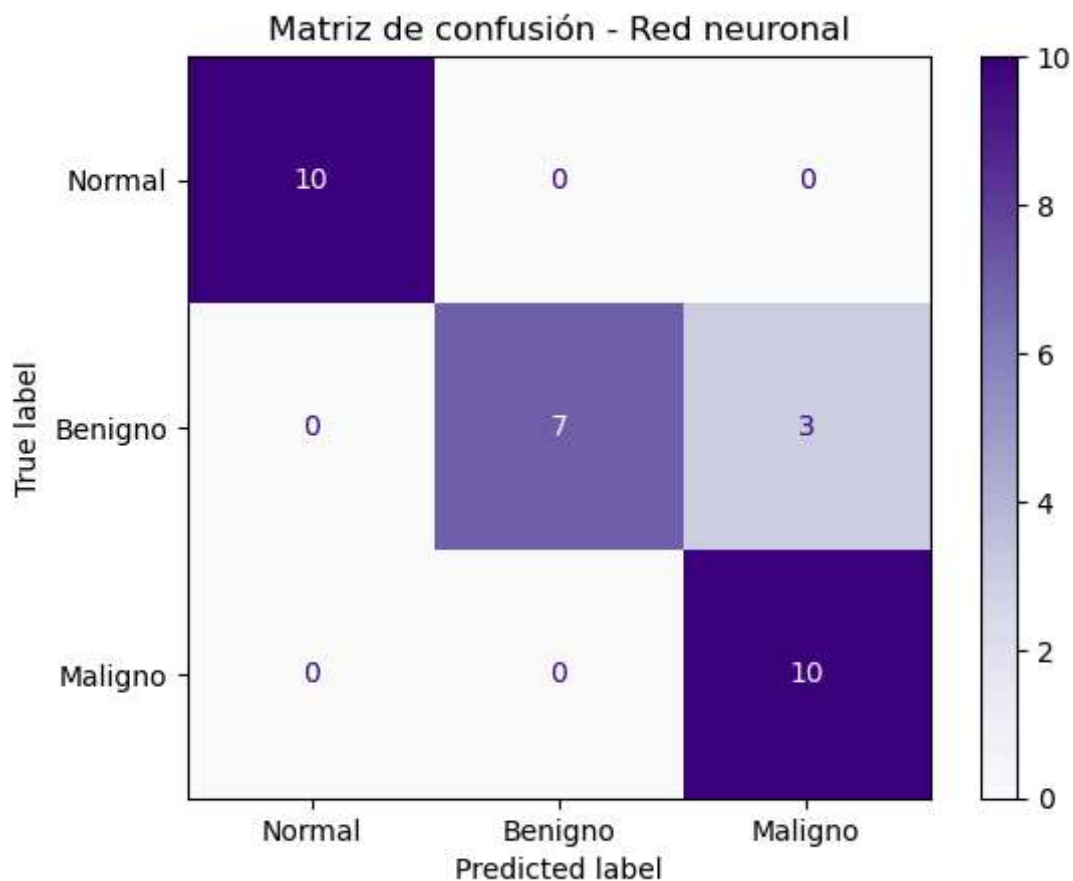
40 # Validación cruzada (StratifiedKFold)
41 skf = StratifiedKFold(
42     n_splits=5,
43     shuffle=True,
44     random_state=SEED
45 )
46
47 # entrenamos el modelo
48 for i, (train_index, test_index) in enumerate(skf.split(X, Y)):
49     # dividimos los datos
50     X_train, X_test = X[train_index], X[test_index]
51     y_train, y_test = Y[train_index], Y[test_index]
52
53     # conversion de tipos de dato
54     y_train = [classes[elem] for elem in y_train]
55     y_test = [classes[elem] for elem in y_test]
56     train_labels = to_categorical(y_train)
57     test_labels = to_categorical(y_test)
58
59     # creamos y entrenamos el modelo
60     model = baseline_model()
61     history = model.fit(
62         X_train, train_labels,
63         validation_data=(X_test, test_labels),
64         epochs=100,
65         batch_size=16,
66         verbose=0
67     )

```

Utilizamos 5 splits para imitar la proporción 80/20.

Con la validación cruzada hemos obtenido los siguientes resultados:





	precision	recall	f1-score	support
Normal	1.00	1.00	1.00	10
Benigno	1.00	0.70	0.82	10
Maligno	0.77	1.00	0.87	10
accuracy			0.90	30
macro avg	0.92	0.90	0.90	30
weighted avg	0.92	0.90	0.90	30

Vemos que los valores no son tan prometedores.

Sin embargo, simplemente aumentando el número de neuronas de la capa oculta del modelo y cambiando la función de activación de la capa de salida obtenemos un modelo casi perfecto.