

Redes Neuronales

Tarea 2: Dataset Breast Cancer Winsconsin (Original)

Realizado por: David Cevallos

Fecha: 2023-07-13

Enlace Google Colab: <https://colab.research.google.com/drive/1GYhFC8I4m7GvtLpAdiULtzYFD-TgK4tC>

En esta tarea analizaremos el dataset Breast Cancer Wisconsin (Original) mediante árboles de decisión y la estrategia k-folding cross validation.

Se debe tener cargado el fichero de nombre breast-cancer-wisconsin.data en la carpeta home previo a la ejecución del notebook.

El dataset se encuentra disponible en: <https://archive.ics.uci.edu/dataset/15/breast+cancer+wisconsin+original>

Consta de 9 descriptores y 699 observaciones. Comprende un problema de clasificación binaria.

```
1 # Importación de librerías
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import KFold
5 from sklearn import tree
6 import sklearn.metrics as mt
7 import imblearn.metrics as imt
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10
11 # Parte 1: Importación y limpieza de datos
12
13 # Función para redefinir valores de etiquetas
14 # 0 Benigno, 1 Maligno
15 # Tomaremos como referencia de clase positiva a 1 (Maligno)
16 def etiquetar(s):
17     if s == 2:
18         return 0
19     return 1
20
21 # Función para corregir valores inconsistentes encontrados en el dataset
22 def corregir(s):
23     if s == "?":
24         return 1
25     return int(s)
26
27 # Lectura de datos desde el archivo
28 data = pd.read_csv("/home/breast-cancer-wisconsin.data", header=None)
29
30 # Se presenta en pantalla inconsistencias encontradas respecto a los valores únicos
31 # que toman las columnas 6 y 10
32 print("--Antes de limpieza de datos")
33 print(data.dtypes)
34 print(np.unique(data[6]))
35 print(np.unique(data[10]))
36
37 # Resultados después de la limpieza de datos
38 print("--Después de limpieza de datos")
39 data[6] = data[6].map(lambda s: corregir(s))
40 y = data[10].map(lambda s: etiquetar(s)).to_numpy()
41 x = data.drop(columns=[0,10]).to_numpy()
42
43 # Se presenta en pantalla el resultado final de la limpieza de datos
44 print(np.unique(x[:,5].reshape(1,-1)))
45 print(np.unique(y))
46
```

```
--Antes de limpieza de datos
0      int64
1      int64
2      int64
3      int64
4      int64
5      int64
```

```

6     object
7     int64
8     int64
9     int64
10    int64
dtype: object
['1' '10' '2' '3' '4' '5' '6' '7' '8' '9' '?']
[2 4]
--Después de limpieza de datos
[ 1  2  3  4  5  6  7  8  9 10]
[0 1]

```

En el siguiente apartado de código definimos la función `calcularMetricas` para obtener los valores de métricas. Esta función es similar a la implementada en Matlab en clase.

```

1 # Función general para calcular métricas
2 # (Esta función es similar a la implementada en clase en Matlab para los ejemplos
3 # de PlayGolf, Tratamiento y Fisheriris)
4 # Calcula los valores de métricas y despliega la matriz de confusión y curva ROC
5 def calcularMetricas(ytrue, ypred):
6     print("Matriz de confusión")
7     cm = mt.confusion_matrix(ytrue, ypred)
8     labels = ["Benigno", "Maligno"]
9
10    # Si se desea se puede observar la matriz en forma no gráfica
11    pd_cm = pd.DataFrame(cm,
12                          #         columns=labels,
13                          #         index=labels)
14    # print(pd_cm)
15
16    # Desplegamos matriz de confusión de forma gráfica
17    ax= plt.subplot()
18    sns.heatmap(cm, annot=True, fmt='g', ax=ax,
19               square=True,
20               cmap="Blues",
21               linewidths=0.2,
22               annot_kws={"fontsize":12})
23    ax.set_xlabel("Categoría predecida", fontsize=8)
24    ax.set_ylabel("Categoría verdadera", fontsize=8)
25    ax.xaxis.set_ticklabels(labels, fontsize=8)
26    ax.yaxis.set_ticklabels(labels, fontsize=8)
27    cbar = ax.collections[0].colorbar
28    cbar.ax.tick_params(labelsize=8)
29    plt.yticks(rotation=0)
30    plt.xticks(rotation=90)
31    plt.show()
32
33    # Obtenemos las métricas correspondientes
34    precision = mt.precision_score(ytrue, ypred)
35    recall = mt.recall_score(ytrue, ypred)
36    f1score = mt.f1_score(ytrue, ypred)
37    accuracy = mt.accuracy_score(ytrue, ypred)
38    error = 1-accuracy
39    specificity = mt.specificty_score(ytrue, ypred)
40    fprm = 1-specificity
41    auc = mt.roc_auc_score(ytrue, ypred)
42
43    # Se presenta en pantalla las métricas obtenidas
44    print("Valores de métricas obtenidas con clase positiva Maligno:")
45    print("precision:", precision)
46    print("recall:", recall)
47    print("f1-score:", f1score)
48    print("accuracy:", accuracy)
49    print("error:", error)
50    print("specificity:", specificity)
51    print("false positive rate:", fprm)
52    print("auc:", auc)
53
54    # Desplegamos curva ROC (caso binario)
55    # Tomando como referencia de clase positiva a 1 (Maligno)
56    fig, ax = plt.subplots()
57    fpr, tpr, thresholds = mt.roc_curve(ytrue, ypred)
58    auc_value = round(auc,4)
59    plt.plot(fpr, tpr, color="blue", label="AUC="+str(auc_value))
60    plt.plot([0,1], [0,1], "--k")
61    plt.xlabel("False Positive Rate (FPR)")

```

```

62 plt.ylabel("True Positive Rate (TPR)")
63 plt.legend(loc="lower right", prop={'size': 6})
64 plt.grid()
65 plt.show()
66
67 # Otra alternativa para desplegar la curva ROC (caso binario)
68 # mt.RocCurveDisplay.from_predictions(ytrue, ypred)
69 # plt.show()
70
71 return np.array([[precision],[recall],[f1score],[accuracy],[error],[specificity],[fpr], [auc]])
72

```

Ahora aplicamos la estrategia k-fold cross validation y en cada iteración llamaremos a la función anteriormente creada. Al finalizar el bucle computaremos la media de cada métrica obtenida (suma de cada métrica dividido para el total de iteraciones k) para obtener por cada una la métrica media total (también conocida como macro metric).

```

1 # K-fold cross validation (10 iteraciones)
2 kfold = KFold(n_splits=10, shuffle=True, random_state=0)
3 num_metricas = 8
4 metricas = np.array([], dtype=np.int64).reshape(num_metricas,0)
5
6 # Obtenemos los índices de aquellas observaciones que corresponden a entrenamiento
7 # y test y realizamos las correspondientes iteraciones
8 for i, (train_index, test_index) in enumerate(kfold.split(x)):
9     print("-----")
10    print("Resultados para Fold Nro.",str(i+1))
11    print("-----")
12    # Definición y entrenamiento de árbol de decisión
13    modelo = tree.DecisionTreeClassifier(criterion="entropy", random_state=0)
14    modelo.fit(x[train_index],y[train_index])
15    # Obtención de etiquetas verdaderas y predecidas para la iteración
16    ytrue = y[test_index]
17    ypred = modelo.predict(x[test_index])
18    # Obtenemos las métricas para la iteración
19    resultado = calcularMetricas(ytrue, ypred)
20    metricas = np.hstack((metricas,resultado))
21
22 # Matriz con resultados de métricas
23 # Filas: Denota cada métrica (precision, recall, etc)
24 # Columnas: Número de k-fold (iteración)
25 print("")
26 print("-----")
27 print("Resultado final matriz de métricas obtenidas:")
28 print("Filas (Métrica) vs Columnas (Nro. de fold)")
29 print("-----")
30 print(metricas)
31 print("")
32 print("-----")
33 print("Se computa la media total para cada métrica:")
34 print("-----")
35 totales = np.mean(metricas,1).reshape(num_metricas,-1)
36 print("mean-precision:",totales[0,0])
37 print("mean-recall:",totales[1,0])
38 print("mean-f1-score:",totales[2,0])
39 print("mean-accuracy:",totales[3,0])
40 print("mean-error:",totales[4,0])
41 print("mean-specificity:",totales[5,0])
42 print("mean-false positive rate:",totales[6,0])
43 print("mean-auc:",totales[7,0])
44

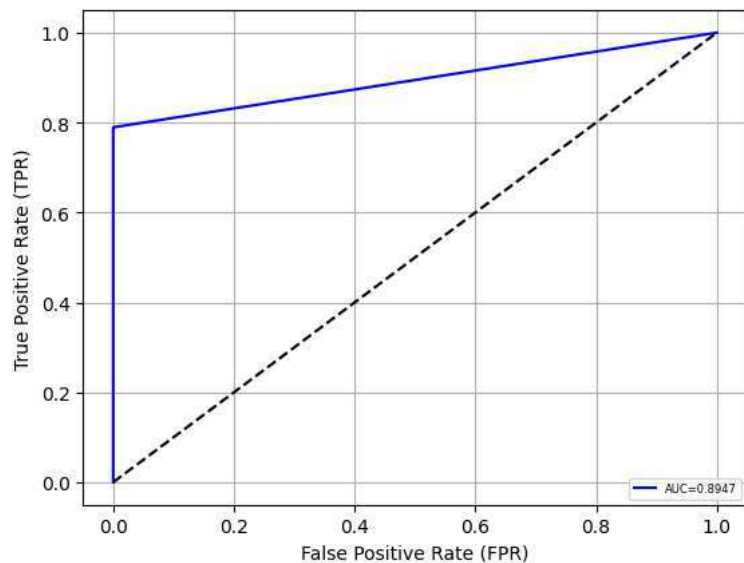
```

Benigno
Maligno

Categoría predecida

Valores de métricas obtenidas con clase positiva Maligno:

precision: 1.0
recall: 0.7894736842105263
f1-score: 0.8823529411764706
accuracy: 0.9420289855072463
error: 0.05797101449275366
specificity: 1.0
false positive rate: 0.0
auc: 0.8947368421052632



Resultado final matriz de métricas obtenidas:

Filas (Métrica) vs Columnas (Nro. de fold)

```
[[0.88888889 1. 0.80952381 0.95833333 0.85714286 1.
 1. 0.94444444 0.88 1.
[0.92307692 0.93103448 0.85 0.92 1. 0.96
 0.85 0.77272727 0.88 0.78947368]
[0.90566038 0.96428571 0.82926829 0.93877551 0.92307692 0.97959184
 0.91891892 0.85 0.88 0.88235294]
[0.92857143 0.97142857 0.9 0.95714286 0.92857143 0.98571429
 0.95714286 0.91428571 0.91428571 0.94202899]
[0.07142857 0.02857143 0.1 0.04285714 0.07142857 0.01428571
 0.04285714 0.08571429 0.08571429 0.05797101]
[0.93181818 1. 0.92 0.97777778 0.875 1.
 1. 0.97916667 0.93333333 1.
[0.06818182 0. 0.08 0.02222222 0.125 0.
 0. 0.02083333 0.06666667 0.
[0.92744755 0.96551724 0.885 0.94888889 0.9375 0.98
 0.925 0.87594697 0.90666667 0.89473684]]
```

Se computa la media total para cada métrica:

mean-precision: 0.9338333333333333
mean-recall: 0.8876312362773342
mean-f1-score: 0.907193051443822
mean-accuracy: 0.9399171842650104
mean-error: 0.060082815734989636
mean-specificity: 0.9617095959595959
mean-false positive rate: 0.03829040404040405
mean-auc: 0.9246704161184651