Creating a Logistic Regression Model to Predict Breast Cancer Recurrence

Read system parameters

Interact with the operating system

Work with multi-dimensional arrays and matrices

Import software libraries and load the dataset In [1]: # Import required libraries

import sys

import os

import numpy as np

View the first five records

<class 'pandas.core.frame.DataFrame'> RangeIndex: 286 entries, 0 to 285 Data columns (total 12 columns):

Column Non-Null Count Dtype

O recurrence 286 non-null int64

In [3]: | df.info()

2

3

4

0

Out[5]:

Out[6]:

0

0

201

df.describe()

count 286.000000

mean

std

min

25%

50%

75%

recurrence

0.297203

0.457828

0.000000

0.000000

0.000000

1.000000

mpl.rc('axes', labelsize=14) mpl.rc('xtick', labelsize=12) mpl.rc('ytick', labelsize=12)

recurrence

meno It 40

df.hist(figsize=(20,20))

plt.show()

200

175

150

50

250

200

150

100

50

40

60

40

Examine descriptive statistics

meno_pre

286.000000

0.524476

0.500276

0.000000

0.000000

1.000000

1.000000

Show histograms for each attribute in the dataset

df["recurrence"].value counts()

Name: recurrence, dtype: int64

Show descriptive statistics

age_decade

286.000000

46.643357

10.118183

20.000000

40.000000

50.000000

50.000000

import numpy as np
import pandas as pd # Manipulate and analyze data import matplotlib as mpl # Create 2D charts import scipy as sp
import sklearn # Perform scientific computing and advanced mathematics # Perform data mining and analysis import sklearn import seaborn as sns # Perform data visualization import matplotlib.pyplot as plt # Summarize software libraries used print('Libraries used in this project:') print('- NumPy {}'.format(np. version)) print('- Pandas {}'.format(pd.__version__)) print('- Matplotlib {}'.format(mpl.__version__)) print('- SciPy {}'.format(sp.__version__)) print('- Scikit-learn {}'.format(sklearn.__version__)) print('- Python {}\n'.format(sys.version)) # Read the raw dataset based on data from https://archive.ics.uci.edu/ml/datasets/breast+cancer print('----') print('Loading the dataset.') #PROJECT ROOT DIR = '.' #DATA PATH = os.path.join(PROJECT ROOT DIR, 'breast cancer data') #print('Data files in this project:', os.listdir(DATA PATH)) #data raw file = os.path.join(DATA PATH, 'breast-cancer.csv') df = pd.read csv("breast-cancer.csv") print('Loaded {} records\n'.format(len(df))) Libraries used in this project: - NumPy 1.19.2 - Pandas 1.1.3 - Matplotlib 3.3.2 - SciPy 1.5.2 - Scikit-learn 0.23.2 - Python 3.8.5 (default, Sep 3 2020, 21:29:08) [MSC v.1916 64 bit (AMD64)] Loading the dataset. Loaded 286 records Get acquainted with the data structure and preview the records In [2]: # Show the various features and their data types

1 age_decade 286 non-null int64 2 meno_pre 286 non-null int64 3 meno_lt_40 286 non-null int64 4 meno_ge_40 286 non-null int64 5 tumor_size 286 non-null int64 6 inv_nodes 286 non-null int64 7 node caps 286 non-null int64

node caps 286 non-null int64 286 non-null int64 deg_malig 286 non-null breast left int64 10 breast_right 286 non-null 286 non-null 11 irradiat dtypes: int64(12)

memory usage: 26.9 KB df.head() In [4]: Out[4]: deg_malig recurrence age_decade meno_pre meno_lt_40 meno_ge_40 tumor_size inv_nodes node_caps breast_left breast_right irra 0 30 0 32 0 40 42

0

0

meno_ge_40

286.000000

0.451049

0.498470

0.000000

0.000000

0.000000

1.000000

0

meno_lt_40

286.000000

0.024476

0.154791

0.000000

0.000000

0.000000

0.000000

100

80

20

160

140

120

100

80

60

40

42

62

42

tumor_size

286.000000

48.643357

10.118183

22.000000

42.000000

52.000000

52.000000

inv_nodes

286.000000

2.573427

3.451904

1.000000

1.000000

1.000000

4.000000

0

0

node_caps

286.000000

0.223776

0.417504

0.000000

0.000000

0.000000

0.000000

140

120

100

40

100

80

60

40

20

2

2

deg_malig

286.000000

2.048951

0.738217

1.000000

2.000000

2.000000

3.000000

0

0

breast_left breast_rig

286.00000

0.46853

0.49988

0.00000

0.00000

0.00000

1.00000

1.00000

286.000000

0.531469

0.499883

0.000000

0.000000

1.000000

1.000000

meno_pre

tumor size

1.000000 70.000000 1.000000 1.000000 1.000000 72.000000 25.000000 1.000000 3.000000 1.000000 max Use histograms to visualize the distribution of various features

125 60 80 100 40 60 75

age_decade

meno_ge_40

1.0 20 30 50 0.0 0.0 0.6 deg_malig node_caps inv_nodes 200 120 200 175 100 150 80 125 100 60 100 75 40 50 50 20 25 2.0 breast_left breast_right irradiat 140 140 200 120 120 150 100 100 100 60 60 40 40 50 20 20 0.8 Split the data into training and validation sets and labels # Import a function to split the dataset from sklearn.model selection import train test split # Specify the column to be included in the label set ('recurrence') X = df.iloc[:,1:]y = df.iloc[:,0]# Specify columns to be included in the training and validation sets (all other columns) # Split the training set, validation set, and labels for both X_train, X_test, y_train, y_test = train_test_split(X.values, y.values, random state = 0, stratify=y) # Compare the number of rows and columns in the original data to the training and validation sets X_train.shape, y_train.shape, X_test.shape, y_test.shape

0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0,

logistic_pred = logistic.predict(X_test)

0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,

Out[12]: LogisticRegressionCV(cv=5, max iter=1000, random state=0, scoring='f1')

In [15]: # Add columns to a copy of the test data to compare predictions against actual values. table = pd.DataFrame(X test.copy())

table["Predicted Recurrence"] = logistic pred table 0 1 2 3 4 5 6 7 8 9 10 Actual Recurrence Predicted Recurrence

X train array([[30, 0, ...,

((214, 11), (214,), (72, 11), (72,))

 $[50, 0, 0, \ldots, 1, 0, 1],$

 $[50, 0, 0, \ldots, 0, 1, 0],$ $[30, 1, 0, \ldots, 0, 1, 1],$

[50, 0, 1, ..., 1, 0, 0]], dtype=int64)

Preview the training data

Preview the labels

y_train

from sklearn.linear_model import LogisticRegressionCV logistic = LogisticRegressionCV(random_state=0, cv=5, scoring="f1", max_iter=1000) In [12]: logistic.fit(X_train,y_train)

Create a logistic regression model, and use the validation data and labels to score it.

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js

table["Actual Recurrence"] = y_test.copy()

0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1,

0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1], dtype=int64) **Build the model**

0, 0, 0, 1, 0, 0], dtype=int64) Test the model

View examples of the predictions compared to actual recurrence.

0 30 1 0 0 32 4 1 3 1 0 1 1 **1** 60 0 0 1 62 1 0 1 1 0 0 **2** 30 1 0 0 32 1 0 3 1 0 1 0 **3** 40 1 0 0 42 1 0 2 1 0 0 0 **4** 30 1 0 0 32 1 0 1 1 0 0 0

In [14]: logistic_pred 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0,

67 60 0 0 1 62 7 1 2 0 1 0 0 0

 60 0 0 1 62 1 0 1 0 1 1 0 40 1 0 0 42 7 1 3 0 1 40 1 0 0 42 1 0 2 1 0 0

71 30 1 0 0 32 1 0 3 1 0 0 1 0 72 rows × 13 columns from sklearn.metrics import roc auc score, classification report roc_auc_score(y_test,logistic_pred)

print(classification report(y test,logistic pred)) precision recall f1-score support

0 0.76 0.88 0.82 51 0.33 0.41 0.54 21 72

Out[20]: 0.6078431372549019

accuracy 0.72 0.65 0.61 0.61 72 macro avg weighted avg 0.70 0.72 0.70 72