



逻辑回归 & 决策树

Illusionna✉

19:58, Thursday 14th September, 2023

概述

Logistic Regression 和 Decision Tree 都可以用于样本分类任务，但两者原理不同，本次作业选用“乳腺癌”权威数据集，编写 Python

脚本，以准确率和查全率为重要指标，以 F_1 分数为综合指标，采用多种抽样方式划分数据集，较为全面地对比不同途径得到的结果。

相关代码：

https://gitee.com/Illusionna/OnlineSharing/tree/master/LogisticRegression_DecisionTree

I. 数据集

1.1 数据集来源

<https://archive.ics.uci.edu/dataset/15/breast+cancer+wisconsin+original>

Figure 1.1: 乳腺癌数据集视图

1.2 数据集描述

乳腺癌数据集一共收录 4024 条样本，包括 15 列属性，以及对应的状态标签 “Alive” “Dead”，同时注意到存在 “Undifferentiated” “anaplastic” “Grade IV” “Regional” 属性值的一类异常标签（非 “Alive” “Dead”）。

1.3 数据预处理

1.3.1 删除异常标签

一共筛选出 19 条异常样本，直接删除。

Table 1.1 异常样本

Age	Race	Marital Status	T Stage	N Stage	6th Stage	differentiate	Grade	A Stage	Tumor Size	Estrogen Status	Progesterone Status	Regional Node Examined	Reginol Node Positive	Survival Months	Status
52	Black	Single	T3	N3	IHC	Undifferentiated	anaplastic	Grade IV	Regional	100	Positive	Positive	23	17	16
38	White	Married	T3	N1	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	70	Positive	Positive	10	1	102
37	Black	Married	T3	N1	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	60	Positive	Positive	5	2	97
69	White	Married	T1	N1	IIA	Undifferentiated	anaplastic	Grade IV	Regional	20	Positive	Positive	17	3	85
59	White	Single	T2	N2	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	40	Negative	Negative	10	9	33
67	White	Married	T2	N3	IHC	Undifferentiated	anaplastic	Grade IV	Regional	21	Positive	Positive	25	17	75
39	White	Single	T3	N2	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	70	Positive	Positive	14	4	79
69	White	Widowed	T2	N3	IHC	Undifferentiated	anaplastic	Grade IV	Regional	28	Positive	Positive	14	13	13
58	White	Single	T1	N2	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	19	Negative	Negative	47	7	9
43	White	Married	T3	N2	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	75	Negative	Negative	10	5	59
62	Black	Separated	T4	N3	IHC	Undifferentiated	anaplastic	Grade IV	Regional	70	Positive	Positive	10	10	34
63	White	Married	T2	N1	IIB	Undifferentiated	anaplastic	Grade IV	Regional	25	Negative	Positive	10	2	27
50	White	Married	T3	N1	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	55	Positive	Positive	12	1	89
57	White	Married	T1	N1	IIA	Undifferentiated	anaplastic	Grade IV	Regional	18	Positive	Positive	5	1	98
40	White	Single	T3	N1	IIIA	Undifferentiated	anaplastic	Grade IV	Regional	52	Positive	Positive	11	9	88
43	White	Married	T4	N3	IHC	Undifferentiated	anaplastic	Grade IV	Regional	59	Negative	Positive	16	10	55
54	White	Divorced	T1	N1	IIA	Undifferentiated	anaplastic	Grade IV	Regional	16	Positive	Negative	19	2	101
45	White	Married	T1	N1	IIA	Undifferentiated	anaplastic	Grade IV	Regional	13	Positive	Positive	5	1	93
49	White	Married	T2	N1	IIB	Undifferentiated	anaplastic	Grade IV	Regional	28	Positive	Positive	12	3	71

1.3.2 字符标签编码

数据编码视图如 Figure 1.2 所示。

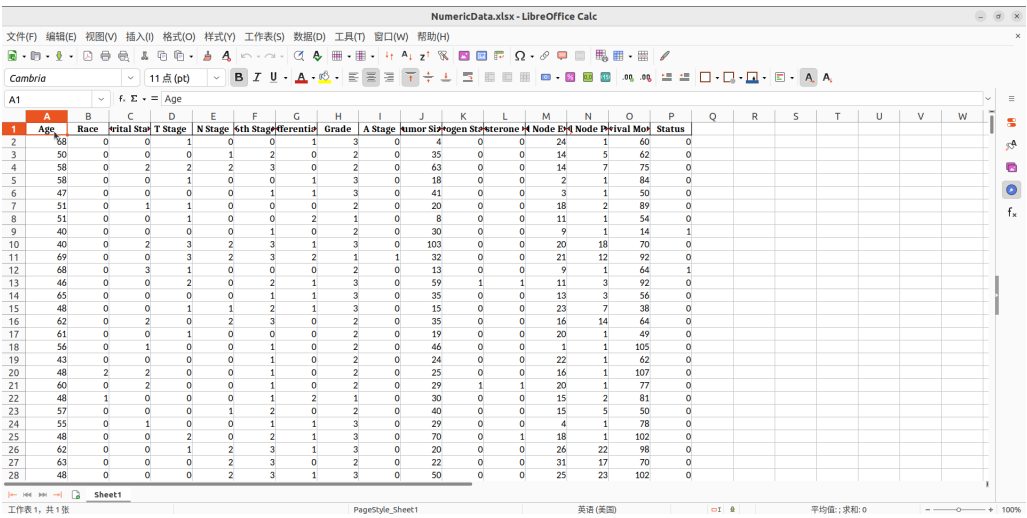


Figure 1.2: 字符属性值编码

给字符型分类变量按照 0 – 1 – 2 – 3 – ... 顺序编码，约定：

Table 1.2 编码规定

Encoding Attribute	0	1	2	3	4
Race	White	Other	Black		
Marital Status	Married	Single	Divorced	Widowed	Separated
T Stage	T2	T1	T3	T4	
N Stage	N1	N2	N3		
6th Stage	IIA	IIB	IIIA	IIIC	IIIB
differentiate	Moderately	Poorly	Well		
A Stage	Regional	Distant			
Estrogen Status	Positive	Negative			
Progesterone Status	Positive	Negative			
Status	Alive	Dead			

1.3.3 选择性标准化、归一化

标准化：

$$f(x) = \frac{x - \mu}{\sigma}$$

归一化：

$$f(x) = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

其中， x 为某个属性下的样本值， μ 为属性相应的样本均值， σ 为属性对应的样本标准差， x_{\min} 为属性相应的样本最小值， x_{\max} 为属性相应的样本最大值。

II. 抽样划分数据集

2.1 Random

乳腺癌数据集随机三七开，即 70% 的训练集，30% 的测试集。

2.2 Hierarchical

将乳腺癌数据集按照标签 “Alive” “Dead” 先划分开，再各从两者中取 70% 当训练集，30% 当测试集。

2.3 Leave-One-Out

乳腺癌数据集只留一个样本作测试集，剩余样本全部作训练集。

2.4 K-Fold Cross Validation

乳腺癌数据集留若干（假设 s 个）样本作测试集，剩余样本全部作训练集，使得 K 轮次循环后， $s \times K = N$ ，其中 $N = 4005$ 。

2.5 Bootstrap

一类非参数 Monte Carlo 方法概率密度抽样，从数据集中有放回的抽取样本 N 次，把每次抽到的样本构成训练集，剩下没被抽到的构成测试集。

$$\lim_{N \rightarrow +\infty} \frac{\text{Test}}{\text{Dataset}} \times 100\% = \lim_{N \rightarrow +\infty} (1 - \frac{1}{N})^N = \frac{1}{e}$$

III. 定义指标

由于精确率和召回率此消彼长，当两者出现矛盾时，此时无法确定以哪个评估指标为准，譬如，一个模型精确率高但召回率低，而另一个模型精确率低但召回率高。因此，需要定义诸多综合型指标。

$$\text{虚警率: FPrate} = \frac{FP}{N} \times 100\%$$

$$\text{准确度: Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \times 100\%$$

$$F_1 = (1 + \delta^2) \times \frac{\frac{TP}{TP + FP} \times \frac{TP}{TP + FN}}{\delta^2 \times \frac{TP}{TP + FP} + \frac{TP}{TP + FN}}$$

$$G - \text{Mean} = \sqrt{\frac{TP}{TP + FN} \times \frac{TN}{TN + FP}}$$

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + TN)(TP + FN)(TN + FP)(TN + FN)}}$$

IV. 逻辑回归

4.1 Logistic Regression Results

Table 4.1 逻辑回归比对结果

Processing	Shuffle	Epoch	Precision	Recall	Accuracy	F_1	MCC
None	Random	10	70.103%	41.975%	89.767%	52.510%	14.731%
		100	80.412%	40.415%	88.852%	53.793%	15.652%
		1000	75.424%	46.354%	89.018%	57.419%	17.780%
	Hierarchical	10	72.072%	43.716%	88.861%	54.422%	16.266%
		100	72%	39.344%	88.446%	50.883%	14.596%
		1000	81.818%	44.262%	90.025% ^{†‡}	57.447%	16.792%
	Bootstrap	10	77.5%	39.574%	88.882%	52.394%	14.973%
		100	74.265%	48.792%	90.195% ^{†‡}	58.892%	17.849%
		1000	79.464%	46.809%	91.964% ^{†‡}	57.329%	16.223%
Standardize	Random	10	56.757%	10.995%	84.526%	18.421%	3.802%
		100	76.190%	15.920%	85.108%	26.337%	6.144%
		1000	82.609%	20.321%	86.938%	32.618%	7.742%
	Hierarchical	10	76.667%	12.568%	86.118%	21.596%	4.657%
		100	75%	19.672%	86.783%	31.169%	7.277%
		1000	100% ^{†‡}	20.765%	87.947%	34.379%	8.081%
	Bootstrap	10	62.162%	11.005%	86.254%	18.699%	3.780%
		100	77.419%	23.645%	88.673%	36.226%	8.390%
		1000	84.746%	24.510%	88.828%	38.023%	8.937%
Normalize	Random	10	76.923%	35.928%	89.601%	48.980%	12.919%
		100	74.419%	51.892%	89.850%	61.146%	19.602%
		1000	82.022%	39.891%	89.517%	53.676%	15.121%
	Hierarchical	10	77.647%	36.066%	88.695%	49.254%	13.530%
		100	79.048%	45.355%	89.859%	57.639%	17.131%
		1000	82.828%	44.809%	90.191% ^{†‡}	58.156%	17.030%
	Bootstrap	10	80.531%	42.130%	89.890%	55.319%	15.762%
		100	74.603%	44.762%	89.772%	55.952%	16.409%
		1000	78.814%	44.286%	90.444% ^{†‡}	56.707%	16.159%
Both	Random	10	72.917%	37.838%	88.270%	49.822%	14.127%
		100	80.734%	45.833%	89.601%	58.472%	17.759%
		1000	80%	41.212%	90.183% ^{†‡}	53.543%	14.757%
	Hierarchical	10	77.528%	37.705%	88.861%	50.735%	14.151%
		100	75.258%	39.891%	88.861%	52.143%	14.915%
		1000	82.609%	41.530%	89.776%	55.273%	15.762%
	Bootstrap	10	76.364%	38.898%	89.215%	51.334%	14.366%
		100	80.357%	41.667%	89.821%	54.878%	15.581%
		1000	83.495%	40.376%	90.184% ^{†‡}	54.430%	15.012%

4.2 逻辑回归代码执行结果

```

35 pro.AbnormalProcessing(mode="int")
36 #
37 # 字符串编码
38 pro.EncodeData()
39 #
40 # 数据标准化、归一化处理
41 pro.Standardize()
42 pro.Normalize()
43 #
44 # 生成训练集、测试集
45 pro.Shuffle(rule="bootstrap")
46 #
47 # 加载训练集、测试集
48 train = LOADER('./tempData/Train.xlsx')
49 test = LOADER('./tempData/Test.xlsx')
50 #
51 lgs = LOGISTIC(
52     trainData = train.data,
53     trainLabels = train.labels,
54     testData = test.data,
55     testLabels = test.labels
56 )
57 #
58 # 逻辑回归
59 lgs.Logistic()
60 lgs.Score()

```

Annotation:

```

{'White': 0, 'Other': 1, 'Black': 2}
{'Married': 0, 'Single': 1, 'Divorced': 2, 'Widowed': 3, 'Separated': 4}
{'T2': 0, 'T1': 1, 'T3': 2, 'T4': 3}
{'N1': 0, 'N2': 1, 'N3': 2}
{'IIA': 0, 'IIB': 1, 'IIIA': 2, 'IIIC': 3, 'IIIB': 4}
{'Moderately differentiated': 0, 'Poorly differentiated': 1, 'Well differentiated': 2}
{'Regional': 0, 'Distant': 1}
{'Positive': 0, 'Negative': 1}
{'Positive': 0, 'Negative': 1}
{'Alive': 0, 'Dead': 1}

```

Processed result saved...

Shuffled result is saved to './tempData'...

Logistic Regression:

```

TP: 85
FN: 146
FP: 16
TN: 1223

Precision: 84.158 %
Recall: 36.797 %
Accuracy: 88.988 %
FPrate: 1.088 %
F1: 51.295 %
G-Mean: 6.768 %
MCC: 14.195 %

```

Figure 4.1: main.py 执行效果 (epoch=100)

V. 决策树

5.1 Decision Tree Results

Table 5.1 决策树比对结果

Shuffle	Precision	Recall	Accuracy	FPrate	F_1	G-Mean	MCC
Random	75.862%	49.162%	90.100%	2.329%	59.661%	7.690%	18.318%
Hierarchical	77.869%	51.913%	90.441%	2.244%	62.295%	7.635%	19.626%
Bootstrap	83.333%	53.488%	91.180%	1.534%	63.889%	7.046%	19.895%

决策树选取 Gini 系数作为评判指标，树的最大深度 $\text{maxDepth} = 4$ ，从决策树比对结果表格窥见：Bootstrap > Hierarchical > Random，即采用自助抽样得到的决策树不论是 Precision 还是 Recall 乃至 Accuracy 都较高。并且综合性指数 $F_1 = 63.889\%$ 胜过其余两个 Shuffle。

5.2 决策树代码执行结果

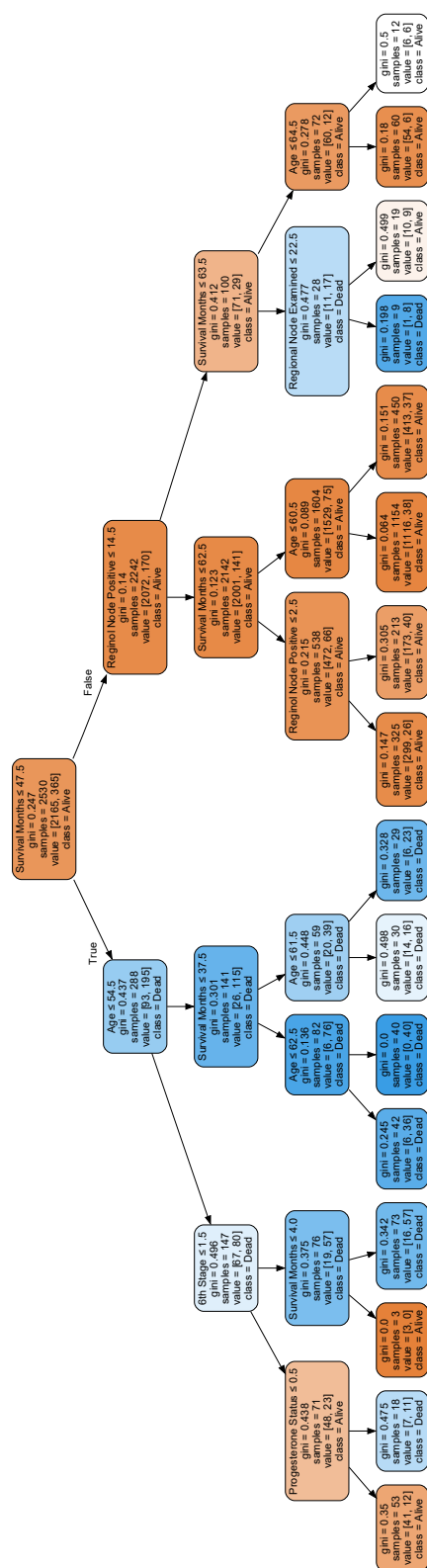


Figure 5.1: 决策树结果

VI. 结果分析

Logistic Regression 的比对结果 **Table 3.1** 反映，只对起源数据集进行编码处理，不进行标准化或归一化操作，采用 Bootstrap 抽样迭代 1000 次，逻辑回归结果 Accuracy 可达 91.964%，综合回归的结果：

None > Both > Normalize > Standardize

Bootstrap > Hierarchical > Random

一个有趣结果是，Standardize 处理采用 Hierarchical 抽样，训练迭代 1000 次，测试集的 Precision 居然满达 **100%**，测试集全部查准，笔者已将对应的训练集和测试集上传至文章概述提供的链接。

决策树祖先结点 SurvivalMonths ≤ 47.5 ，其 Gini 系数仅有 0.247，这意味分类的纯度相当高，决策树分类十分有效。

综上所述，逻辑回归和决策树对于分类问题（乳腺癌数据集）的预测结果效果都非常良好，最高可达 90% 左右的 Accuracy。

Code. 代码

Code1(Logistic): main.py

```

1  """
2  # System --> Linux & Python3.8.0
3  # File ----> main.py
4  # Author --> Illusionna
5  # Create --> 2023/09/12 19:22:39
6  """
7  # -*- encoding: utf-8 -*-
8
9
10 # -----
11 # -----
12 import os
13 from utils.Loader import LOADER

```



```
14 from utils.Processing import PROCESSING
15 from utils.LogisticRegression import LOGISTIC
16 # -----
17 # -----
18 def cls() -> None:
19     os.system('clear')
20     cls()
21 # -----
22 # -----
23 # 加载乳腺癌数据集.
24 io = './OriginalData/BreastCancer.xlsx'
25 # -----
26 lor = LOADER(io)
27 # -----
28 pro = PROCESSING(
29     data = lor.data,
30     labels = lor.labels,
31     attribute = lor.attribute
32 )
33 # -----
34 # 删除异常标签值.
35 pro.AbnormalProcessing(mode="int")
36 # -----
37 # 字符值编码.
38 pro.EncodeData()
39 # -----
40 # 数据标准化、归一化处理.
41 pro.Standardize()
42 pro.Normalize()
43 # -----
44 # 生成训练集、测试集.
45 pro.Shuffle(rule='bootstrap')
```

```
46 # -----
47 # 加载训练集、测试集.
48 train = LOADER('./tempData/Train.xlsx')
49 test = LOADER('./tempData/Test.xlsx')
50 # -----
51 lgs = LOGISTIC(
52     trainData = train.data,
53     trainLabels = train.labels,
54     testData = test.data,
55     testLabels = test.labels,
56     epoch = 1000
57 )
58 # -----
59 # 逻辑回归.
60 lgs.Logistic()
61 lgs.Score()
62 # -----
63 # -----
```

Code2(Decision): main.py

```
1 """
2 # System --> Windows & Python3.8.0
3 # File ----> main.py
4 # Author --> Illusionna
5 # Create --> 2023/09/14 03:41:34
6 """
7 # -*- encoding: UTF-8 -*-
8
9
10 import os
11 from utils.Loader import LOADER
```

```
12 from utils.Processing import PROCESSING
13 from utils.DecisionTree import DECISION
14 # -----
15 # -----
16 def cls() -> None:
17     os.system('cls ')
18     cls()
19 # -----
20 # -----
21 # 加载乳腺癌数据集.
22 io = './OriginalData/BreastCancer.xlsx'
23 # # -----
24 lor = LOADER(io)
25 # # -----
26 pro = PROCESSING(
27     data = lor.data,
28     labels = lor.labels,
29     attribute = lor.attribute
30 )
31 # -----
32 # 删除异常标签值.
33 pro.AbnormalProcessing(mode="int")
34 # -----
35 # 字符值编码.
36 pro.EncodeData()
37 # -----
38 # 数据标准化、归一化处理.
39 # pro.Standardize()
40 # pro.Normalize()
41 # -----
42 # 生成训练集、测试集.
43 pro.Shuffle(rule='bootstrap')
```

```
44 # -----
45 # -----
46 train = LOADER('./tempData/Train.xlsx')
47 test = LOADER('./tempData/Test.xlsx')
48 # -----
49 tree = DECISION(
50     trainData = train.data,
51     trainLabels = train.labels,
52     testData = test.data,
53     testLabels = test.labels
54 )
55 # -----
56 tree.Decision()
57 # -----
58 tree.Score()
59 tree.Chart(test.attribute)
60 # -----
61 # -----
```

Code3: Loader.py

```
1 """
2 # System --> Linux & Python3.8.0
3 # File ----> Loader.py
4 # Author --> Illusionna
5 # Create --> 2023/09/12 19:19:06
6 """
7 # -*- encoding: utf-8 -*-
8
9
10 import pandas as pd
11
```

```

12 class LOADER:
13     """
14     数据加载类.
15     """
16     def __init__(self, io: str) -> None:
17         self.io = io
18         LOADER.Loader(self)
19
20     def Loader(self) -> None:
21         """
22         Load the original data.
23         """
24         try:
25             df = pd.read_excel(self.io, header=None)
26             self.attribute = df.loc [0]. tolist ()
27             self.labels = df.loc [:( df.shape[1]-1)].drop(index=0).
                tolist ()
28             self.data = df.loc [:[ i for i in range(0, (df.shape
                [1]-1), 1) ]]. drop(index=0).values. tolist ()
29         except:
30             print("Failed to load data.")
31             print("数据加载失败.")
32             print("Check Excel table format and path.")
33             print("检查表格格式和路径.")
34             print("Reference link of picture format.")
35             print("参考格式图片链接.\n")
36             print(r"https:// gitee .com/Illusionna/OnlineSharing/raw
                /master/View_of_Excel.png")
37             print("\n")

```

Code4: Processing.py

```
1  """
2  # System --> Linux & Python3.8.0
3  # File -----> Processing.py
4  # Author --> Illusionna
5  # Create --> 2023/09/12 19:20:00
6  """
7  # -*- encoding: utf-8 -*-
8
9
10 import os
11 import random
12 import numpy as np
13 import pandas as pd
14 from typing import Literal
15
16 COLOURS = {
17     'default': '\033[39m',
18     'green': '\033[32m',
19     'pink': '\033[35m',
20     'red': '\033[31m',
21     'orange': '\033[33m'
22 }
23
24 class PROCESSING:
25     """
26     数据预处理类.
27     """
28     def __init__(self, data: list, labels: list, attribute: list) ->
29         None:
30         self.data = data
31         self.label = labels
32         self.attribute = attribute
```

```

32
33     RULE = Literal['int', 'float']
34     def AbnormalProcessing(self, mode:RULE) -> None:
35         """
36         异常值函数：处理标签异常的样本，视 mode 为正常。
37         """
38         remarkList = []
39         for i in range(0, len(self.label), 1):
40             if isinstance(self.label[i], eval(mode)) == True:
41                 remarkList.append(i)
42             print(COLOURS['red'], end='')
43             print('Abnormal Samples:')
44             print(COLOURS['orange'])
45             dictionary = {}
46             for i in range(0, len(remarkList), 1):
47                 temp = self.data[remarkList[i]]
48                 dictionary[f'Error{i+1}'] = temp
49             df = pd.DataFrame(dictionary).T
50             print(df)
51             transformerD = pd.DataFrame(self.data)
52             transformerD.drop(remarkList, inplace=True)
53             transformerL = pd.DataFrame(self.label)
54             transformerL.drop(remarkList, inplace=True)
55             print(COLOURS['green'])
56             print('Abnormal samples removed...\n')
57             print(COLOURS['default'])
58             self.data = transformerD.values.tolist()
59             self.label = transformerL.T.values.tolist()[0]
60
61     def EncodeData(self) -> None:
62         """
63         属性值编码函数：将字符型值按照 0,1,2,... 编码。

```

```

64         ?? ?? ??
65         matrix = []
66         for i in range(0, len(self.label), 1):
67             temp = list(self.data[i])
68             temp.append(self.label[i])
69             matrix.append(temp)
70         transformer = pd.DataFrame(matrix).T
71         statisticsList = []
72         write = []
73         for i in range(0, len(transformer), 1):
74             column = transformer.values.tolist()[i]
75             judge = (
76                 isinstance(column[0], int)
77                 |
78                 isinstance(column[0], float)
79             )
80             if judge:
81                 pass
82             else:
83                 temp = pd.value_counts(column)
84                 pos = 0
85                 for j in range(0, len(temp), 1):
86                     temp[j] = pos
87                     pos = ~pos
88                 write.append(dict(temp))
89                 statisticsList.append((i, dict(temp)))
90         print(COLOURS['red'], end='')
91         print('Annotation:')
92         print(COLOURS['orange'])
93         for i in range(0, len(write), 1):
94             print(write[i])
95         print(COLOURS['default'])

```



```

96     data = []
97     pos = 0
98     for i in range(0, len(matrix), 1):
99         sample = []
100         for j in range(0, len(matrix[i]), 1):
101             if j == statisticsList[pos][0]:
102                 sample.append(write[pos][matrix[i][j]])
103                 pos = ~pos
104             else:
105                 sample.append(matrix[i][j])
106         pos = 0
107         data.append(sample)
108     dictionary = {}
109     for i in range(0, len(self.attribute), 1):
110         dictionary[self.attribute[i]] = pd.DataFrame(data).T.
111             values.tolist()[i]
112         os.makedirs('./tempData', exist_ok=True)
113         pd.DataFrame(dictionary).to_excel('./tempData/
114             NumericData.xlsx', index=None)
115     print(COLOURS['green'])
116     print('Processed result saved...\n')
117     print(COLOURS['default'])
118     for i in range(0, len(data), 1):
119         self.label[i] = data[i][-1]
120         del data[i][-1]
121     self.data = data
122
123     def Standardize(self) -> None:
124         """
125         功能: 数据标准化.\n
126         函数: Z-score 法则.\n
127         参数: 要求 “纯” 数据矩阵, 不含表头属性和表格索引, 类

```

```

        型 list.
    """
126
127     self.data = pd.DataFrame(self.data).T.values.tolist()
128     temp = []
129     for i in range(0, len(self.data), 1):
130         tempList = []
131         avg = np.average(self.data[i])
132         var = np.var(self.data[i])
133         for j in range(0, len(self.data[i]), 1):
134             value = (self.data[i][j] - avg) / var
135             tempList.append(value)
136         temp.append(tempList)
137     self.data = temp
138     self.data = pd.DataFrame(self.data).T.values.tolist()
139
140     def Normalize(self) -> None:
141         """
142         功能: 数据归一化.\n
143         范围: 归一化后数据区间 [0, 1].\n
144         参数: 要求 “纯” 数据矩阵, 不含表头属性和表格索引, 类
            型 list.
145         """
146         self.data = pd.DataFrame(self.data).T.values.tolist()
147         temp = []
148         for i in range(0, len(self.data), 1):
149             tempList = []
150             for j in range(0, len(self.data[i]), 1):
151                 value = (self.data[i][j] - min(self.data[i])) / (
                    max(self.data[i]) - min(self.data[i]))
152                 tempList.append(value)
153             temp.append(tempList)
154         self.data = temp

```

```

155         self.data = pd.DataFrame(self.data).T.values.tolist()
156
157     MODE = Literal['random', 'hierarchical', 'LOO', 'KFCV', '
        bootsrap']
158     def Shuffle( self , rule:MODE="random", num:int=4) -> None
        or list:
159         """
160         数据洗牌：打乱数据集，划分训练集和测试集.\n
161         rule = {
162             "random": RandomShuffle( )\n
163             " hierarchical ": HierarchicalShuffle ( )\n
164             "LOO" --> LeaveOneOut: LOOShuffle( )\n
165             "KFCV" --> KFoldCrossValidation: KFCVShuffle( )\n
166             "bootsrap": BootsrapShuffle( )
167         }
168         """
169         os.makedirs("./tempData", exist_ok=True)
170         self.matrix = []
171         for i in range(0, len( self.data), 1):
172             temp = self.data[i]
173             temp.append(self.label[i])
174             self.matrix.append(temp)
175         if rule == "random":
176             PROCESSING.__RandomShuffle(self)
177         elif rule == "hierarchical":
178             PROCESSING.__HierarchicalShuffle(self)
179         elif rule == "LOO":
180             PROCESSING.__LOOShuffle(self, testSetIndex=num)
181         elif rule == "KFCV":
182             print("Shuffled ...\n")
183             return PROCESSING.__KFCVShuffle(self, K=num)
184         elif rule == "bootsrap":

```

```

185         PROCESSING.___BootsrapShuffle(self)
186     else :
187         PROCESSING.___RandomShuffle(self)
188     print(COLOURS['pink'])
189     print("Shuffled result is saved to './tempData'...\n")
190     print(COLOURS['default'])
191     self.data = pd.DataFrame(self.data).drop([len(self.data[0])
192                                             -1], axis=1)
193     self.data = self.data.values.tolist()
194
195     def ___SaveTempData(self, trainSet:list, testSet: list) -> None:
196         """
197         洗牌划分后的训练集和测试集寄存在 './tempData' 文件夹下.
198         """
199         trainSet = pd.DataFrame(trainSet).T.values.tolist()
200         testSet = pd.DataFrame(testSet).T.values.tolist()
201         trainDictionary = dict(zip(self.attribute, trainSet))
202         testDictionary = dict(zip(self.attribute, testSet))
203         dfTrain = pd.DataFrame(trainDictionary)
204         dfTest = pd.DataFrame(testDictionary)
205         dfTrain.to_excel('./tempData/Train.xlsx', index=False)
206         dfTest.to_excel('./tempData/Test.xlsx', index=False)
207
208     def ___RandomShuffle(self) -> None:
209         """
210         留出法抽样：数据集随即划分，默认 ratio = 0.7，即传统的前
211         后七三随机开.
212         """
213         random.shuffle(self.matrix)
214         sampleNumber = len(self.matrix)
215         ratio = 0.7
216         trainNumber = int(sampleNumber*ratio)

```

```

215         del ratio
216         trainSet = []
217         testSet = []
218         for i in range(0, trainNumber, 1):
219             trainSet.append(self.matrix[i])
220         for i in range(trainNumber, sampleNumber, 1):
221             testSet.append(self.matrix[i])
222         PROCESSING.___SaveTempData(self, trainSet, testSet)
223
224     def ___HierarchicalShuffle( self ) -> None:
225         """
226         分层抽样：默认 ratio = 0.7，从正样本中 70% 放到训练集，
                再从反样本中抽 70% 放到训练集，测试集是剩余 30% 的
                正反样本.
227         """
228         ratio = 0.7
229         trainSet = []
230         testSet = []
231         positiveSet = []
232         negativeSet = []
233         sampleNumber = len(self.matrix)
234         for i in range(0, sampleNumber, 1):
235             if self.matrix[i][-1] == 1:
236                 positiveSet.append(self.matrix[i])
237             else:
238                 negativeSet.append(self.matrix[i])
239         random.shuffle( positiveSet )
240         random.shuffle(negativeSet)
241         for i in range(0, int( len( positiveSet ) * ratio ), 1):
242             trainSet.append(positiveSet[i])
243         for i in range(0, int( len( negativeSet ) * ratio ), 1):
244             trainSet.append(negativeSet[i])

```

```

245         for i in range(int(len(positiveSet)*ratio), len(positiveSet)
246             , 1):
247             testSet.append(positiveSet[i])
248         for i in range(int(len(negativeSet)*ratio), len(negativeSet)
249             , 1):
250             testSet.append(negativeSet[i])
251         random.shuffle(trainSet)
252         random.shuffle(testSet)
253         PROCESSING.___SaveTempData(self, trainSet, testSet)
254
255     def ___LOOShuffle(self, testSetIndex:int=0) -> None:
256         """
257         留一法抽样：参数 num/testSetIndex 是选取一个测试样本的
258             索引.
259         """
260         sampleNumber = len(self.matrix)
261         trainSet = []
262         testSet = []
263         testSet.append(self.matrix[testSetIndex])
264         L = list(
265             set([i for i in range(0, sampleNumber, 1)])
266             -
267             set([testSetIndex])
268         )
269         for i in L:
270             trainSet.append(self.matrix[i])
271         PROCESSING.___SaveTempData(self, trainSet, testSet)
272
273     def ___KFCVShuffle(self, K:int=4) -> list:
274         """
275         K-折交叉验证抽样：参数 num/K 为数据集划分成子集的个
276             数.

```

```

273         """
274         random.shuffle( self.matrix)
275         sampleNumber = len(self.matrix)
276         tensor = []
277         trainSet = []
278         testSet = []
279         pos = 0
280         step = np.floor(sampleNumber/K)
281         for i in range(0, (K-1), 1):
282             matrix = []
283             while pos<(step*(i+1)):
284                 matrix.append(self.matrix[pos])
285                 pos = ~pos
286             tensor.append(matrix)
287         matrix = []
288         for t in range(pos, sampleNumber, 1):
289             matrix.append(self.matrix[pos])
290             pos = ~pos
291         tensor.append(matrix)
292         return tensor
293
294     def __BootsrapShuffle(self) -> None:
295         """
296         自助法概率抽样：数据集较小、难以有效划分训练集和测试集
                时有效，但会改变起源数据集的分布，引入估计偏差，数
                据量大时建议使用留出法和 K 折交叉验证法。
297         """
298         sampleNumber = len(self.matrix)
299         indexList = []
300         for i in range(0, sampleNumber, 1):
301             random.seed(random.random())
302             index = random.randint(0, (sampleNumber-1))

```

```
303         indexList.append(index)
304     indexListTrain = list(set(indexList))
305     indexListTest = list(
306         set([i for i in range(0, sampleNumber, 1)])
307         —
308         set(indexListTrain)
309     )
310     trainSet = []
311     testSet = []
312     for i in indexListTrain:
313         trainSet.append(self.matrix[i])
314     for i in indexListTest:
315         testSet.append(self.matrix[i])
316     PROCESSING.__SaveTempData(self, trainSet, testSet)
```

Code5: LogisticRegression.py

```
1     """
2     # System --> Linux & Python3.8.0
3     # File -----> LogisticRegression.py
4     # Author --> Illusionna
5     # Create --> 2023/09/12 19:19:32
6     """
7     # -*- encoding: utf-8 -*-
8
9
10    from numpy import sqrt
11    from sklearn.linear_model import LogisticRegression
12
13    class LOGISTIC:
14        def __init__(self, trainData: list, trainLabels: list, testData:
            list, testLabels: list, epoch:int=100) -> None:
```



```
15         self.trainX = trainData
16         self.trainY = trainLabels
17         self.testX = testData
18         self.testY = testLabels
19         self.epoch = epoch
20
21     def Logistic( self ) -> None:
22         model = LogisticRegression(
23             multi_class = "auto",
24             solver = "lbfgs",
25             max_iter = self.epoch
26         )
27         lgs = model.fit( self.trainX, self.trainY )
28         prediction = lgs.predict( self.testX )
29         self.TP = 0
30         self.FN = 0
31         self.FP = 0
32         self.TN = 0
33         for i in range(0, len(prediction), 1):
34             if ((prediction[i] == 1) and (self.testY[i] == 1)):
35                 self.TP = ~self.TP
36             elif ((prediction[i] == 0) and (self.testY[i] == 1)):
37                 self.FN = ~self.FN
38             elif ((prediction[i] == 1) and (self.testY[i] == 0)):
39                 self.FP = ~self.FP
40             elif ((prediction[i] == 0) and (self.testY[i] == 0)):
41                 self.TN = ~self.TN
42
43     def Score( self , delta: float=1) -> None:
44         print( ' Logistic Regression:\n' )
45         print( f' TP: {self.TP}' )
46         print( f' FN: {self.FN}' )
```

```

47     print(f'FP: {self.FP}')
48     print(f'TN: {self.TN}')
49     print('')
50     self.precision = self.TP/(self.TP+self.FP)
51     self.recall = self.TP/(self.TP+self.FN)
52     self.accuracy = (self.TP+self.TN)/ (self.TP+self.TN+self.
        FP+self.FN)
53     self.FPrate = self.FP / len(self.testY)
54     self.F1 = (1+delta**2) * (self.precision * self.recall) / ((
        delta**2)*self.precision + self.recall)
55     self.GMean = sqrt((self.TP+self.TN) / ((self.TP+self.FN) *
        (self.TN+self.FP)))
56     self.MCC = (self.TP*self.TN-self.FP*self.FN) / sqrt((self.
        TP+self.TN)*(self.TP+self.FN)*(self.TN+self.FP)*(self.
        TN+self.FN))
57     print('Precision: {:.3f} %'.format(100*self.precision))
58     print('Recall: {:.3f} %'.format(100*self.recall))
59     print('Accuracy: {:.3f} %'.format(100*self.accuracy))
60     print('FPrate: {:.3f} %'.format(100*self.FPrate))
61     print('F1: {:.3f} %'.format(100*self.F1))
62     print('G-Mean: {:.3f} %'.format(100*self.GMean))
63     print('MCC: {:.3f} %'.format(100*self.MCC))
64     print('')

```

Code6: DecisionTree.py

```

1     """
2     # System --> Windows & Python3.8.0
3     # File ----> DecisionTree.py
4     # Author --> Illusionna
5     # Create --> 2023/09/14 04:09:16
6     """

```

```
7  # -*- encoding: UTF-8 -*-
8
9
10 import os
11 import pydotplus
12 from numpy import sqrt
13 from sklearn import tree
14 from six import StringIO
15 from typing import Literal
16
17 def cls() -> None:
18     os.system('cls ')
19     os.environ["PATH"] += os.pathsep + './Graphviz/bin'
20     cls()
21
22 class DECISION:
23     def __init__(self, trainData: list , trainLabels: list , testData:
24         list , testLabels: list ) -> None:
25         self.trainX = trainData
26         self.trainY = trainLabels
27         self.testX = testData
28         self.testY = testLabels
29
30     CRITERION = Literal['gini', 'entropy', 'log_loss']
31     SPLITTER = Literal['best', 'random']
32
33     def Decision( self , criterion :CRITERION='gini', splitter:
34         SPLITTER='best', maxDepth:int=4) -> None:
35         clf = tree. DecisionTreeClassifier (
36             criterion = criterion ,
37             splitter = splitter ,
38             max_depth = maxDepth
```

```

37         )
38         self.clf = clf.fit(self.trainX, self.trainY)
39         prediction = self.clf.predict(self.testX)
40         self.TP = 0
41         self.FN = 0
42         self.FP = 0
43         self.TN = 0
44         for i in range(0, len(prediction), 1):
45             if ((prediction[i] == 1) and (self.testY[i] == 1)):
46                 self.TP = ~self.TP
47             elif ((prediction[i] == 0) and (self.testY[i] == 1)):
48                 self.FN = ~self.FN
49             elif ((prediction[i] == 1) and (self.testY[i] == 0)):
50                 self.FP = ~self.FP
51             elif ((prediction[i] == 0) and (self.testY[i] == 0)):
52                 self.TN = ~self.TN
53
54     def Score(self, delta: float=1) -> None:
55         print('Decision Tree:\n')
56         print(f'TP: {self.TP}')
57         print(f'FN: {self.FN}')
58         print(f'FP: {self.FP}')
59         print(f'TN: {self.TN}')
60         print('')
61         self.precision = self.TP/(self.TP+self.FP)
62         self.recall = self.TP/(self.TP+self.FN)
63         self.accuracy = (self.TP+self.TN)/(self.TP+self.TN+self.
64             FP+self.FN)
65         self.FPrate = self.FP / len(self.testY)
66         self.F1 = (1+delta**2) * (self.precision * self.recall) / ((
67             delta**2)*self.precision + self.recall)
68         self.GMean = sqrt((self.TP+self.TN) / ((self.TP+self.FN) *

```

```

        (self.TN+self.FP)))
67     self.MCC = (self.TP*self.TN-self.FP*self.FN) / sqrt((self.
        TP+self.TN)*(self.TP+self.FN)*(self.TN+self.FP)*(self.
        TN+self.FN))
68     print('Precision: {:.3 f} %'.format(100*self.precision))
69     print('Recall: {:.3 f} %'.format(100*self.recall))
70     print('Accuracy: {:.3 f} %'.format(100*self.accuracy))
71     print('FPrate: {:.3 f} %'.format(100*self.FPrate))
72     print('F1: {:.3 f} %'.format(100*self.F1))
73     print('G-Mean: {:.3 f} %'.format(100*self.GMean))
74     print('MCC: {:.3 f} %'.format(100*self.MCC))
75     print('')
76
77     def Chart(self, attribute: list) -> None:
78         image = StringIO()
79         del attribute[-1]
80         status = ['Alive', 'Dead']
81         tree.export_graphviz(
82             self.clf,
83             out_file = image,
84             feature_names = attribute,
85             class_names = status,
86             filled = True,
87             rounded = True,
88             special_characters = True
89         )
90         chart = pydotplus.graph_from_dot_data(image.getvalue())
91         chart.write_pdf('DecisionTree.pdf')

```