

Spark逻辑回归(LR)模型使用介绍

```

In [2]: 1 #小案例学习spark LR模型的使用
2 from pyspark.ml.feature import VectorAssembler
3 import pandas as pd
4 sample_dataset = [
5     (0, "male", 37, 10, "no", 3, 18, 7, 4),
6     (0, "female", 27, 4, "no", 4, 14, 6, 4),
7     (0, "female", 32, 15, "yes", 1, 12, 1, 4),
8     (0, "male", 57, 15, "yes", 5, 18, 6, 5),
9     (0, "male", 22, 0.75, "no", 2, 17, 6, 3),
10    (0, "female", 32, 1.5, "no", 2, 17, 5, 5),
11    (0, "female", 22, 0.75, "no", 2, 12, 1, 3),
12    (0, "male", 57, 15, "yes", 2, 14, 4, 4),
13    (0, "female", 32, 15, "yes", 4, 16, 1, 2),
14    (0, "male", 22, 1.5, "no", 4, 14, 4, 5),
15    (0, "male", 37, 15, "yes", 2, 20, 7, 2),
16    (0, "male", 27, 4, "yes", 4, 18, 6, 4),
17    (0, "male", 47, 15, "yes", 5, 17, 6, 4),
18    (0, "female", 22, 1.5, "no", 2, 17, 5, 4),
19    (0, "female", 27, 4, "no", 4, 14, 5, 4),
20    (0, "female", 37, 15, "yes", 1, 17, 5, 5),
21    (0, "female", 37, 15, "yes", 2, 18, 4, 3),
22    (0, "female", 22, 0.75, "no", 3, 16, 5, 4),
23    (0, "female", 22, 1.5, "no", 2, 16, 5, 5),
24    (0, "female", 27, 10, "yes", 2, 14, 1, 5),
25    (1, "female", 32, 15, "yes", 3, 14, 3, 2),
26    (1, "female", 27, 7, "yes", 4, 16, 1, 2),
27    (1, "male", 42, 15, "yes", 3, 18, 6, 2),
28    (1, "female", 42, 15, "yes", 2, 14, 3, 2),
29    (1, "male", 27, 7, "yes", 2, 17, 5, 4),
30    (1, "male", 32, 10, "yes", 4, 14, 4, 3),
31    (1, "male", 47, 15, "yes", 3, 16, 4, 2),
32    (0, "male", 37, 4, "yes", 2, 20, 6, 4)
33 ]
34 columns = ["affairs", "gender", "age", "label", "children", "religiousness", "education", "occupation", "rating"]
35 # pandas构建dataframe, 方便
36 pdf = pd.DataFrame(sample_dataset, columns=columns)
37 df = spark.createDataFrame(pdf)
38 # 特征选取: affairs为目标值, 其余为特征值
39 df2 = df.select("affairs", "age", "religiousness", "education", "occupation", "rating")
40 colArray2 = ["age", "religiousness", "education", "occupation", "rating"]
41 df3 = VectorAssembler().setInputCols(colArray2).setOutputCol('features').transform(df2)
42 print('数据集: ')
43 df3.show(truncate=False)
44 # 随机切分 训练集和测试集
45 trainDF, testDF = df3.randomSplit([0.8, 0.2])
46 print('训练集: ')
47 trainDF.show(10, truncate=False)
48 print('测试集: ')
49 testDF.show(10, truncate=False)

```

数据集:

affairs	age	religiousness	education	occupation	rating	features
0	37	3	18	7	4	[37.0, 3.0, 18.0, 7.0, 4.0]
0	27	4	14	6	4	[27.0, 4.0, 14.0, 6.0, 4.0]
0	32	1	12	1	4	[32.0, 1.0, 12.0, 1.0, 4.0]

0	57	5	18	6	5	[57.0, 5.0, 18.0, 6.0, 5.0]
0	22	2	17	6	3	[22.0, 2.0, 17.0, 6.0, 3.0]
0	32	2	17	5	5	[32.0, 2.0, 17.0, 5.0, 5.0]
0	22	2	12	1	3	[22.0, 2.0, 12.0, 1.0, 3.0]
0	57	2	14	4	4	[57.0, 2.0, 14.0, 4.0, 4.0]
0	32	4	16	1	2	[32.0, 4.0, 16.0, 1.0, 2.0]
0	22	4	14	4	5	[22.0, 4.0, 14.0, 4.0, 5.0]
0	37	2	20	7	2	[37.0, 2.0, 20.0, 7.0, 2.0]
0	27	4	18	6	4	[27.0, 4.0, 18.0, 6.0, 4.0]
0	47	5	17	6	4	[47.0, 5.0, 17.0, 6.0, 4.0]
0	22	2	17	5	4	[22.0, 2.0, 17.0, 5.0, 4.0]
0	27	4	14	5	4	[27.0, 4.0, 14.0, 5.0, 4.0]
0	37	1	17	5	5	[37.0, 1.0, 17.0, 5.0, 5.0]
0	37	2	18	4	3	[37.0, 2.0, 18.0, 4.0, 3.0]
0	22	3	16	5	4	[22.0, 3.0, 16.0, 5.0, 4.0]
0	22	2	16	5	5	[22.0, 2.0, 16.0, 5.0, 5.0]
0	27	2	14	1	5	[27.0, 2.0, 14.0, 1.0, 5.0]

only showing top 20 rows

训练集:

affairs	age	religiousness	education	occupation	rating	features
0	32	1	12	1	4	[32.0, 1.0, 12.0, 1.0, 4.0]
0	37	3	18	7	4	[37.0, 3.0, 18.0, 7.0, 4.0]
0	57	5	18	6	5	[57.0, 5.0, 18.0, 6.0, 5.0]
0	32	2	17	5	5	[32.0, 2.0, 17.0, 5.0, 5.0]
0	57	2	14	4	4	[57.0, 2.0, 14.0, 4.0, 4.0]
0	22	4	14	4	5	[22.0, 4.0, 14.0, 4.0, 5.0]
0	37	2	20	7	2	[37.0, 2.0, 20.0, 7.0, 2.0]
0	22	2	17	5	4	[22.0, 2.0, 17.0, 5.0, 4.0]
0	27	4	14	5	4	[27.0, 4.0, 14.0, 5.0, 4.0]
0	37	1	17	5	5	[37.0, 1.0, 17.0, 5.0, 5.0]

only showing top 10 rows

测试集:

affairs	age	religiousness	education	occupation	rating	features
0	27	4	14	6	4	[27.0, 4.0, 14.0, 6.0, 4.0]
0	22	2	12	1	3	[22.0, 2.0, 12.0, 1.0, 3.0]
0	22	2	17	6	3	[22.0, 2.0, 17.0, 6.0, 3.0]
0	27	4	18	6	4	[27.0, 4.0, 18.0, 6.0, 4.0]
0	32	4	16	1	2	[32.0, 4.0, 16.0, 1.0, 2.0]
0	47	5	17	6	4	[47.0, 5.0, 17.0, 6.0, 4.0]
0	37	2	20	6	4	[37.0, 2.0, 20.0, 6.0, 4.0]
1	42	3	18	6	2	[42.0, 3.0, 18.0, 6.0, 2.0]

In [3]:

```
1 # 逻辑回归训练模型
2 from pyspark.ml.classification import LogisticRegression
3 lr = LogisticRegression()
4 model = lr.setLabelCol('affairs').setFeaturesCol('features').fit(trainDF)
5 model.transform(testDF).show(truncate=False)
```

affairs	age	religiousness	education	occupation	rating	features	rawPrediction
						probability	prediction
0	27	4	14	6	4	[27.0, 4.0, 14.0, 6.0, 4.0]	[1.4370128406197153, -1.4370128406197153]
0	22	2	12	1	3	[22.0, 2.0, 12.0, 1.0, 3.0]	[-2.6357535531749416, 2.6357535531749416]
0	22	2	17	6	3	[22.0, 2.0, 17.0, 6.0, 3.0]	[0.6549244841168171, -0.6549244841168171]
0	27	4	18	6	4	[27.0, 4.0, 18.0, 6.0, 4.0]	[2.3700715957681684, -2.3700715957681684]
0	32	4	16	1	2	[32.0, 4.0, 16.0, 1.0, 2.0]	[-4.628685738947233, 4.628685738947233]
0	47	5	17	6	4	[47.0, 5.0, 17.0, 6.0, 4.0]	[3.0503541549910054, -3.0503541549910054]
0	37	2	20	6	4	[37.0, 2.0, 20.0, 6.0, 4.0]	[4.9312737927811945, -4.9312737927811945]
1	42	3	18	6	2	[42.0, 3.0, 18.0, 6.0, 2.0]	[-0.598288626751545, 0.598288626751545]