

# LOGISTIC REGRESSION

田思成 工程科学学院 U201714461

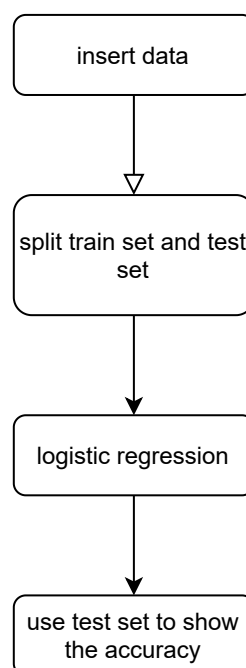
## 实验目的

1. 了解逻辑回归的概念

## 试验内容

使用逻辑回归的办法对乳腺癌数据进行处理，并且把数据分为训练集和测试集以便于测试。最后验证 logistic regression 的准确度。

流程如下：



代码如下：

```
1  #another way to insert data of breast cancer
2  from sklearn.datasets import load_breast_cancer
3  cancer = load_breast_cancer()
4  #show the data's format, so that we can deal with data
5  #print(cancer.data.shape)
6  #print(cancer.target.shape)
7  #show details of data
8  #print(cancer.DESCR)
9  #malignant = 0, benign = 1 if you wish, use print(cancer.target)
10 #split the train and test dataset
11 from sklearn.model_selection import train_test_split
12 X_train, X_test, y_train, y_test = train_test_split(cancer.data,
13 cancer.target, test_size=0.25, random_state=42)
13 # train_data: 所要划分的样本特征集
14 # train_target: 所要划分的样本结果 此处为Benign or malignant
15 # test_size: 样本占比，如果是整数的话就是样本的数量
16 # random_state: 是随机数的种子。随机抽取，random_state 保证每次数据可以重复。
```

```

17 # 随机数种子：其实就是该组随机数的编号，在需要重复试验的时候，保证得到一组一样的随机数。比
    如你每次都填1，其他参数一样的情况下你得到的随机数组是一样的。但填0或不填，每次都会不一样。
18 from sklearn.linear_model import LogisticRegression
19 log_reg = LogisticRegression()
20 log_reg.fit(X_train, y_train)
21 pred = log_reg.predict(X_test)
22 acc_score = log_reg.score(X_test, y_test)
23 print(acc_score)
24 list(cancer.target_names)
25 import pandas as pd
26 d = {'predictions': pred, 'real values': y_test}
27 data = pd.DataFrame(data=d)
28 print(data)
29 data.predictions == data['real values']
30 wrong_predictions= []
31 for i in range(0,143):
32     if data.predictions[i] != data['real values'][i]:
33         wrong_predictions.append(data.predictions[i])
34         print("wrongly diagnosed patient number:", i, 'as',
wrong_predictions[-1])
35     i=i+1
36

```

代码中注释已经写得较为清楚，便不对思路进行赘述。

## 实验结果

```

0.965034965034965
      predictions  real values
0              1             1
1              0             0
2              0             0
3              1             1
4              1             1
..          ...          ...
138            1             1
139            0             0
140            1             1
141            0             0
142            1             1

[143 rows x 2 columns]
wrongly diagnosed patient number: 20 as 1
wrongly diagnosed patient number: 58 as 1
wrongly diagnosed patient number: 77 as 1
wrongly diagnosed patient number: 112 as 0
wrongly diagnosed patient number: 120 as 0
PS C:\Users\Administrator\Desktop\大三下\机器学习\每周试验\2nd\Breast-Cancer-predictions-master>

```

可以注意到，在验证集占比为25%的情况下，准确度为96.5%左右。其中错误的个体序号分别为：20,58,77,112,120

## 试验总结和反思

本次试验主要采取了panda, sklearn这两个模块进行测试，实际上有投机取巧的意思，因为内置了logistic regression的函数，所以甚至不用大费周章的计算数据。但是我认为造轮子要比借轮子好，之后有时间会做一个自己提取数据进行计算的版本。

但总的来说还是比较好的完成了任务，可以较为完善的体现logistic regression的含义。