



实 验 报 告

(2017 / 2018 学年 第 2 学期)

课程名称	机器学习导论
实验名称	Logistic Regression Classifier
实验时间	2018 年 5 月 29 日
指导教师	王邦

姓名	游浩然	学号	U201515429
----	-----	----	------------

1 问题重述

- 钞票数据集 (Banknote Dataset) 涉及根据给定钞票的 4 个度量的特征。据此预测是真钞还是假钞。

2 Python Code for Logistic regression

2.1 logistic

```
1  #-*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  import numpy as np
7  import matplotlib.pyplot as plt
8  # from bigfloat import exp
9  import random
10 import csv
11 import os
12
13 class logistic():
14     def name(self):
15         return 'Logistic Model'
16
17     def sigmoid(self, x):
18         return 1.0 / (1 + np.exp(x))
19
20     def stoGradAscent(self, data, label, numIter):
21         m, n = np.shape(data)
22         weight = [float(1.0) for i in range(n)]
23         for j in range(numIter):
24             dataIndex = list(range(m))
25             for i in range(m):
26                 # alpha = 4 / (1.0+j+i) + 0.0001
27                 alpha = 0.00001
28                 randIndex = int(random.uniform(0, len(dataIndex)))
29                 h = self.sigmoid(np.sum(np.dot(data[randIndex], weight)))
30                 error = label[randIndex] - h
31                 weight += alpha * error * np.array(data[randIndex])
32                 del(dataIndex[randIndex])
33         return weight
34
35     def classifyVector(self, x, weight):
36         pred = self.sigmoid(np.sum(np.dot(x, weight)))
37         if pred > 0.5: return 1.0
38         else: return 0.0
39
40     def train(self, data, label, numIter):
41         self.train_data = data
42         self.train_label = label
43         self.train_weight = self.stoGradAscent(data, label, numIter)
44
45     def val(self, data, label):
46         self.val_data = data
47         self.val_label = label
48         error = 0
49         numTest = np.shape(data)[0]
50         for i in range(numTest):
51             pred = self.classifyVector(data[i], self.train_weight)
52             # print(pred, label[i])
53             if int(pred) != int(label[i]):
54                 error += 1
55         errorRate = (float(error) / numTest)
56         print('Error rate of val data : %f' % errorRate)
57
58     def test(self, data):
59         if os.path.exists('results.csv'):
60             os.remove('results.csv')
61         f = open('results.csv', 'a', newline='')
62         csv_write = csv.writer(f, dialect='excel')
63         i = 0
64         for vector in data:
65             result = []
66             i += 1
67             pred = self.classifyVector(vector, self.train_weight)
68             result.append(i)
69             for item in vector:
70                 result.append(item)
71             result.append(pred)
72             csv_write.writerow(result)
```

2.2 data

```
1  -*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  import random
7
8  def parse_file(filename, has_cls=True):
9      f = open(filename, 'r', encoding='gbk')
10     data, label = [], []
11     for line in f.readlines():
12         if has_cls == True:
13             data_str = line.strip().split(',')[::-1]
14             data_list = []
15             for data_item in data_str:
16                 data_list.append(float(data_item))
17             data.append(data_list)
18             label.append(float(line.strip().split(',')[0]))
19         else:
20             data_str = line.strip().split(',')
21             data_list = []
22             for data_item in data_str:
23                 data_list.append(float(data_item))
24             data.append(data_list)
25     return data, label
26
27 def divide(data, label):
28     num_train = int(0.8*len(data))
29     train_data, train_label = [], []
30     val_data, val_label = [], []
31     index = random.sample(range(len(data)), num_train)
32     for i in range(0, len(data)):
33         if i in index:
34             train_data.append(data[i])
35             train_label.append(label[i])
36         else:
37             val_data.append(data[i])
38             val_label.append(label[i])
39     return train_data, train_label, val_data, val_label
40
41 def dataset():
42     data, label = parse_file('train.txt')
43     train_data, train_label, val_data, val_label = divide(data, label)
44     test_data, test_label = parse_file('test.txt', has_cls=False)
45     print('number of train : ', len(train_data))
46     print('number of val : ', len(val_data))
47     print('number of test : ', len(test_data))
48     return train_data, train_label, val_data, val_label, test_data
```

2.3 Main

```
1  -*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  from data import dataset
7  from logistic import logistic
8
9  # load data
10 train_data, train_label, val_data, val_label, test_data = dataset()
11 # train
12 logistic = logistic()
13 logistic.train(train_data, train_label, 30)
14 # val
15 logistic.val(val_data, val_label)
16 # test
17 logistic.test(test_data)
```