

《机器学习基础》实验报告

实验题目	对数几率回归算法实践		
实验时间	2024/04/07	实验地点	DS3401
<p>一、实验目的</p> <p>掌握线性模型、对率回归算法原理。</p>			
<p>二、实验项目内容</p> <p>1. 理解对率回归算法原理。</p> <p>2. 编程实现对数几率回归算法。</p> <p>3. 将算法应用于西瓜数据集、鸢尾花数据集或其它数据集分类问题。</p>			
<p>三、实验过程或算法（源程序）</p> <pre>1. import csv 2. import numpy as np 3. from sklearn.model_selection import train_test_split 4. 5. #读取西瓜数据集中的数据并返回 6. def loadDataset1(filename): 7. dataset = [] 8. labelset = [] 9. with open(filename, 'r', encoding='utf-8') as csvfile: 10. csv_reader = csv.reader(csvfile) 11. header = next(csv_reader) 12. for row in csv_reader: 13. data_row = [float(row[1]), float(row[2])] 14. dataset.append(data_row) 15. labelset.append(int(row[3])) 16. return dataset, labelset 17. 18. #读取鸢尾花数据集中的数据并返回（2 和 2p2 分别对应两种分类方式） 19. def loadDataset2(filename): 20. dataset = [] 21. labelset = [] 22. with open(filename, 'r', encoding='utf-8') as csvfile: 23. csv_reader = csv.reader(csvfile) 24. header = next(csv_reader) 25. for row in csv_reader: 26. data_row = [float(row[1]), float(row[2]),float(row[3]),float(row[4])] 27. dataset.append(data_row) 28. if row[5]=='Iris-setosa' or row[5]=='Iris-versicolor':</pre>			

```
29.         labelset.append(1)
30.     else:
31.         labelset.append(0)
32.     return dataset, labelset
33.
34. def loadDataset2p2(filename):
35.     dataset = []
36.     labelset = []
37.     with open(filename, 'r', encoding='utf-8') as csvfile:
38.         csv_reader = csv.reader(csvfile)
39.         header = next(csv_reader)
40.         for row in csv_reader:
41.             data_row = [float(row[1]), float(row[2]), float(row[3]), float(row[4])]
42.             if row[5]=='Iris-setosa':
43.                 labelset.append(1)
44.                 dataset.append(data_row)
45.             elif row[5]=='Iris-versicolor':
46.                 labelset.append(0)
47.                 dataset.append(data_row)
48.     return dataset, labelset
49.
50. #定义sigmoid 函数
51. def sigmoid(z):
52.     return 1.0 / (1 + np.exp(-z))
53.
54. #计算准确率
55. def test(dataset,labelset,w):
56.     data=np.mat(dataset)
57.
58.     y=sigmoid(np.dot(data,w))
59.     b,c=np.shape(y)
60.     #功能是查看矩阵或者数组的维数。
61.     rightcount=0
62.
63.     for i in range(b):
64.         flag=-1
65.         if y[i,0]>0.5:
66.             flag=1
67.         elif y[i,0]<0.5:
68.             flag=0
69.         if labelset[i] == flag:
70.             rightcount+=1
71.
72.     rightrate=rightcount/len(dataset)
73.     #print(rightrate)
74.     return rightrate
75.
```

```

76. #迭代求w
77. def training(dataset,labelset,highw):
78.     data=np.mat(dataset).astype(float)
79.     label=np.mat(labelset).transpose()
80.     w = np.ones((len(dataset[0]),1))
81.
82.     #步长
83.     n=0.0001
84.
85.     # 每次迭代计算一次正确率（在测试集上的正确率）
86.     # 达到highw 的正确率，停止迭代
87.     rightrate=0.0
88.     while rightrate<highw:
89.         c=sigmoid(np.dot(data,w))
90.         b=c-label
91.         change = np.dot(np.transpose(data),b)
92.         w=w-change*n
93.         #预测，更新准确率
94.         rightrate = test(dataset,labelset,w)
95.     return w
96.
97. #西瓜
98. dataset=[]
99. labelset=[]
100. filename = 'D:\zjw\demo\machine learning\watermelon_3a.csv'
101. dataset,labelset=loadDataset1(filename)
102. X_train, X_test, y_train, y_test = train_test_split(dataset, labelset, test_size=0.1, random_state=42)
103. w=training(X_train,y_train,0.90)
104. print("西瓜数据集: ")
105. print("若使得准确率大于 90%，则此时的 w 为: \n",w)
106. accuracy = test(X_test, y_test, w)
107. print("测试集上的准确率: %f" % (accuracy * 100) + "%")
108.
109. #鸢尾花
110. dataset=[]
111. labelset=[]
112. filename = 'D:\zjw\demo\machine learning\Iris-data.csv'
113. dataset,labelset=loadDataset2(filename)
114. X_train, X_test, y_train, y_test = train_test_split(dataset, labelset, test_size=0.2, random_state=42)
115. w1=training(X_train,y_train,0.95)
116. print("鸢尾花数据集: ")
117. print("对于第一次分类，准确率大于 0.95，则此时的 w 为: \n",w1)
118. r1= test(X_test, y_test, w1)
119. print("测试集上的准确率为:%f"%(r1*100)+"%")
120.
121. dataset,labelset=loadDataset2p2(filename)

```

```
122. X_train, X_test, y_train, y_test = train_test_split(dataset, labelset, test_size=0.2, random_state=42)
123. w2=training(X_train,y_train,0.95)
124. print("对于第二次分类, 准确率大于 0.95, 则此时的 w 为: \n",w2)
125. r2 = test(X_test, y_test, w2)
126. print("测试集上的准确率为:%f"%(r2*100)+"%")
```

四、实验结果及分析

1.控制台输出:

西瓜数据集:

若使得准确率大于90%, 则此时的w为:

`[[-0.27394248]`

`[0.90898622]]`

测试集上的准确率: 100.000000%

鸢尾花数据集:

对于第一次分类, 准确率大于0.95, 则此时的w为:

`[[0.22718722]`

`[0.93813229]`

`[-0.86834889]`

`[0.06346749]]`

测试集上的准确率为:86.666667%

对于第二次分类, 准确率大于0.95, 则此时的w为:

`[[-0.28867293]`

`[0.61324918]`

`[-0.35930641]`

`[0.5280983]]`

测试集上的准确率为:100.000000%

2.实验分析

(1) 对数几率算法原理

对数几率回归是一种广泛应用于二分类问题的统计学习方法。它通过假设数据服从伯努利分布, 并利用 Sigmoid 函数将线性模型的输出转化为样本属于某一类别的概率。具体来说, 对数几率回归计算给定输入特征下样本属于正类的概率, 并基于这个概率和真实标签构建损失函数。通过最小化损失函数, 模型可以学习到最优的权重和偏置项, 迭代得到最优解。

(2) 程序结果分析

鸢尾花数据集多分类问题处理: 将三个类进行多次分类, 以二叉树的形式对进行两次分类, 每次分类就是一次单独的分类;

准确率偏高: 可能是测试集样本数量较小, 仅有的几个测试用例全部预测成功, 同时也说明了模型训练效果好。