

Machine Learning: Assignment #2

Logistic Regression

康瑞

id=1160300514

1160300514@stu.hit.edu.cn

marisukicandycandy@gmail.com

HIT — October 25, 2018

Contents

1	Details/Environments	2
2	设计思想	2
2.1	算法原理	2
2.1.1	Logistic Regression	2
2.1.2	梯度下降法求参数 W	3
2.1.3	牛顿法	3
2.2	算法实现	4
2.2.1	Gradient Descent	4
2.2.2	Newton Gradient Descent	4
3	实验结果与分析	4
3.1	Direct Gradient Descent	4
3.2	Newton's method of Gradient Descent	5
3.3	Test with Wine classify data from UCI	5
4	结论	6
5	参考文献	7
6	Codes	7
7	End	11

1 Details/Environments

实验目的：

理解逻辑回归模型，掌握逻辑回归模型的参数估计算法

实验要求：

实现两种损失函数的参数估计：梯度下降、牛顿法

- (1). 无惩罚项
- (2). 加入对参数的惩罚
- (3). X^l 各维度间的条件独立性假设验证

实验环境：python3.x + numpy

2 设计思想

2.1 算法原理

2.1.1 Logistic Regression

由条件概率公式及对条件分布的推导可以得出 Logistic Regression 的优化函数：

$$W_{MLE} = \operatorname{argmax}_W P(< X^1, Y^1 > \dots < X^L, Y^L > | W) = \operatorname{argmax}_W \prod_l P(< X^l, Y^l > | W) \quad (1)$$

$$W_{MLE} = \operatorname{argmax}_W \prod_l P(Y^l | X^l, W) \quad (2)$$

取二分类情况，带入对每个 label 的条件概率公式：

$$P(Y = 0 | X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i x_i)}, P(Y = 1 | X, W) = 1 - P(Y = 0 | X, W) \quad (3)$$

可以整理 Loss 函数，及 Loss 函数的矩阵表达式：

$$Loss = \ln \prod_l P(Y^l | X^l, W) = \sum_l Y^l \ln P(Y^l = 1 | X^l, W) + (1 - Y^l) \ln P(Y^l = 0 | X^l, W) \quad (4)$$

$$Loss = \sum_l Y^l \ln \frac{P(Y^l = 1 | X^l, W)}{P(Y^l = 0 | X^l, W)} + \ln P(Y^l = 0 | X^l, W) = \sum_l Y^l (W^T X^l) - \ln(1 + \exp(W^T X^l)) \quad (5)$$

$$Loss = Y * (W^T X) - \ln(1 + \exp(W^T X))^1 \quad (6)$$

因此，可以采用的 W 计算方式为梯度下降与牛顿法（由于共轭梯度法需要寻求 $Ax=b$ 形式的解析式，因此不易用共轭梯度法求解），下推导在牛顿法和梯度下降法中需要用到的计算表达式，和计算方法。

Loss 梯度及 Loss 的二阶梯度表达式：

$$\nabla Loss = Y * X - \frac{\exp(W^T X)}{1 + \exp(W^T X)} * X \quad (7)$$

¹Here, we suppose $X^l = [1, x^l]$ the l's sample of dataset x.

由式 (7) 得:

$$\nabla Loss = X * (Y - \frac{\exp(W^T X)}{1 + \exp(W^T X)}) \quad (8)$$

式 (8) 等价于 (9):

$$\nabla Loss = \sum_l X^l * (Y^l - \frac{\exp(\sum_i w_i X_i^l)}{1 + \exp(\sum_i w_i X_i^l)}) \quad (9)$$

考察矩阵的维度: W : (dim, 1), X : (dim, L), Y : (L, 1), $\nabla Loss$: (dim, 1), 重整式 (8), 可以写出可编程形式:

$$\nabla Loss = X * (Y^T - (\frac{\exp(W^T X)}{1 + \exp(W^T X)})^T) \quad (10)$$

对于 $\nabla Loss$ 为 (dim, 1) 阶矩阵, 因此, $\nabla Loss[i] = X[i] * (Y - (\frac{\exp(W^T X)}{1 + \exp(W^T X)})^T)$, 推导式 (9) 得二阶导数:

$$\nabla_{W_i W_j}^2 Loss = \frac{\partial \nabla_{W_i} Loss}{\partial W_j} \quad (11)$$

$$\nabla_{ij}^2 Loss = \sum_l -X_i^l X_j^l * \frac{\exp(W^T X^l)}{1 + \exp(W^T X^l)} \quad (12)$$

因此, $\nabla_{ij}^2 Loss$ 维度为 (dim, dim), 为 Hessian 矩阵, 计算时, 需要对每组数据计算生成矩阵: $(X_i^l X_j^l)$ 及参数计算式, 并将所有样本累加, 可得 Hessian 矩阵。

2.1.2 梯度下降法求参数 W

由式 (10) 及梯度下降方案 (无正则项): $W \leftarrow W + \lambda * \nabla_W Loss$ 。由 MAP 方法,

$$W = \operatorname{argmax}_W \ln P(W) \prod P(Y^l | X^l, W) \quad (13)$$

即增加 W 的假设先验, 可以假设 W 的先验分布为 Gaussian 分布, 即

$$W \sim N(0, \sigma) \quad (14)$$

, 可以得出 W 的含正则项的梯度下降法更新方程:

$$W \leftarrow W - \eta \lambda W + \eta \sum_l X (Y - \frac{\exp(W^T X)}{1 + \exp(W^T X)}) \quad (15)$$

2.1.3 牛顿法

牛顿法等价于求 $\min_w f(w)$ 的 W 点, 利用导数逼近的方式, 找寻最优点, 因此, Loss 为优化函数, 需要求得 W 的更新过程。

$$W \leftarrow W + \lambda \frac{\nabla_W Loss}{\nabla_{WW}^2 Loss} \quad (16)$$

, 正则化项为:

$$W \leftarrow W + \lambda \frac{\nabla_W Loss}{\nabla_{WW}^2 Loss} - \eta \lambda W \quad (17)$$

因而对 (17) 式做矩阵可以优化求解。

2.2 算法实现

2.2.1 Gradient Descent

Algorithm 1: GD

Input: (X, Y) , X is (dim, L) sample,
 $\text{dim} = \text{class} + 1$; Y is a row vector

Result: W , such that gradloss is
(almostly) minimized

```
while gradloss  $\geq 1e - 5$  do
     $\text{grad} \leftarrow X((Y - \frac{\exp(W^T X)}{1 + \exp(W^T X)})^T)$ 
     $W \leftarrow W - \alpha \text{grad}$ 
end
return  $W$ 
```

2.2.2 Newton Gradient Descent

Algorithm 2: Newton GD

Input: X, Y (X, Y are illustrated before)

Result: W

```
while gradLoss  $\geq 1e - 5$  do
     $H \leftarrow \sum_l -X_i^l X_j^l * \frac{\exp(W^T X^l)}{1 + \exp(W^T X^l)}$ 
     $G \leftarrow X * (Y^T - (\frac{\exp(W^T X)}{1 + \exp(W^T X)})^T)$ 
     $\text{step} \leftarrow H^{-1}G$ 
     $W \leftarrow W + \alpha \text{step}$ 
end
return  $W$ 
```

3 实验结果与分析

The dataset size ², using dimision, loss, α and learning rate are evaluated and illustrated in the result.

以下仅添加部分导出的结果。

3.1 Direct Gradient Descent

实验结果如下所示：图 1-4，所选对 gradLoss 的准出下界为 $1e-4$ ，学习率为 $5e-8$ 。实验中分别测得在有无正则化项（即认为 W 的分布满足以 0 为均值的正态分布）以及变量间的条件独立性假设的情况下的 4 个测试样例，由实验结果得出，正则化可以增强模型的泛化能力，而变量间的条件独立性虽在推导逻辑递归的过程中被使用作为假设，但在实际试验中的影响情况并不明显，变量间的条件独立性与否是由生成数据时确定， X 的样本生成为以 $[0,2]$ 、 $[2,0]$ 为中心的多维高斯分布，在 cov 矩阵为对角时变量之间相互独立，非对角对称阵时变量之间条件相关。

²size is the evaluate of dataset, set by initializer, usually 100-500(for visual consideration)

Logistic Regression: gussian:False ind:False testgrad test_res:0.979021

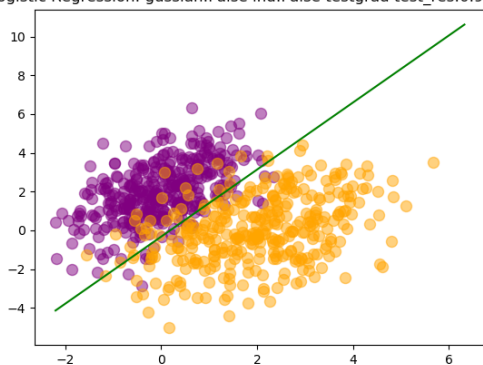


Figure 1: GD:no reg value not independent

Logistic Regression: gussian:False ind:True testgrad test_res:0.874126

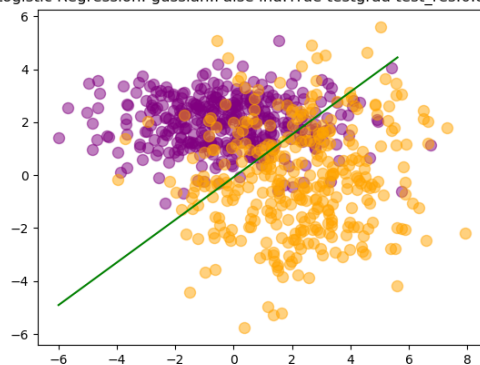


Figure 2: GD:no reg value independent

Logistic Regression: gussian:True ind:False testgrad test_res:0.982517

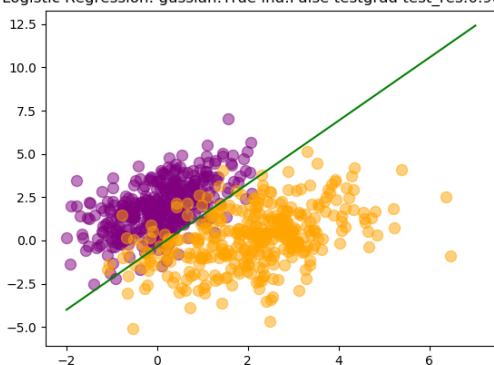


Figure 3: GD:with reg value not independent

Logistic Regression: gussian:True ind:True testgrad test_res:0.954545

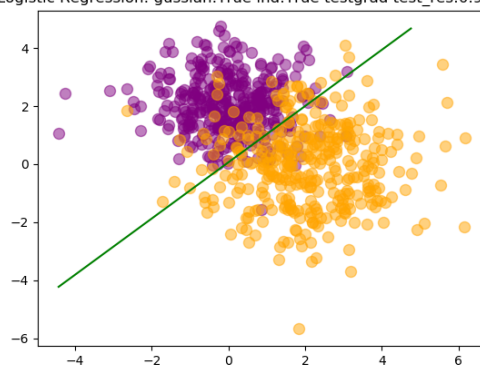


Figure 4: GD:with reg value independent

3.2 Newton's method of Gradient Descent

图 5-8 为 Newton 法用于对 Loss 的最小值求解，比较于直接的 Gradient Descent，Newton 法收敛更快。

3.3 Test with Wine classify data from UCI

测试数据来自 UCI 的 wine data, 该数据将红酒分为三类: 1, 2, 3, 并度量相应的指标, 相关指标一共 13 类, 因此, W 矩阵维数为 14, 13 个度量为:

- (1). Alcohol
- (2). Malic acid
- (3). Ash
- (4). Alcalinity of ash
- (5). Magnesium
- (6). Total phenols
- (7). Flavanoids
- (8). Nonflavanoid phenols
- (9). Proanthocyanins

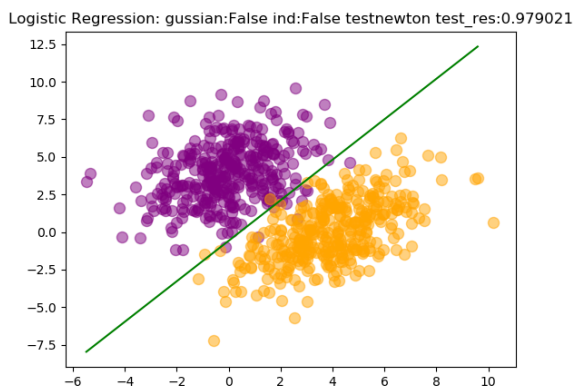


Figure 5: Newton: no reg val not independent

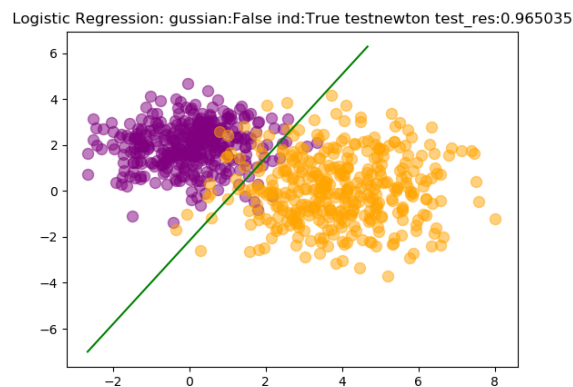


Figure 6: Newton: no reg val independent

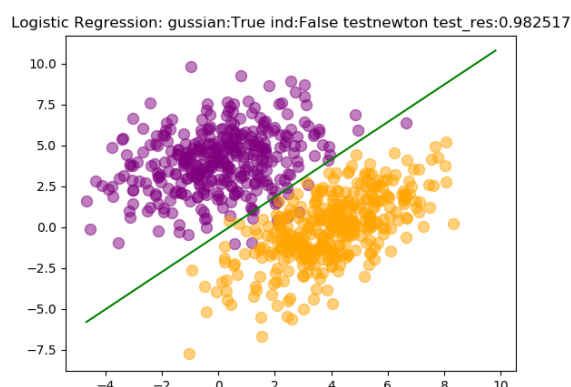


Figure 7: Newton: with reg val not independent

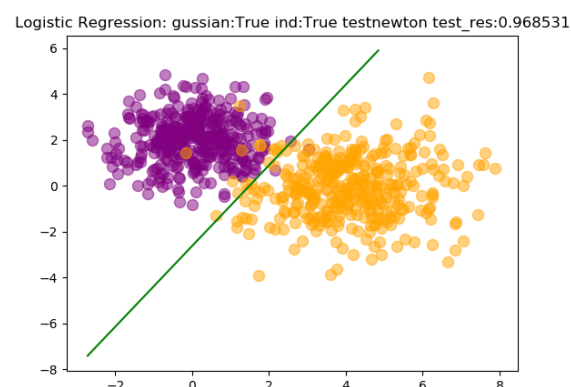


Figure 8: Newton:with reg val independent

(10). Color intensity

(11). Hue

(12). OD280/OD315 of diluted wines

(13). Proline

可以将此数据划分为三类，利用逻辑递归，首先度量分类类别 1，由梯度下降法得出相应的 W 矩阵，再计算类别 2 的矩阵，由这两个参数向量，可以采用决策树的结构来预测这三个分类（由于原代码实现的方式为 2 分类）得出结果预测准确率为 0.70.

4 结论

- (1). 利用多维的高斯分布模型，通过调整 COV 矩阵的元素可以生成变量间相关与无关的数据分布，在描点后的状态为椭圆（有无偏斜）与正圆。
- (2). Newton 法利用导数生成函数求最低点的方式逼近原函数的最值点，可以极大的加快模型拟合速度，但在维度过高时会出现由于 Hessian 阵求逆带来的浮点误差问题，导致模型发散，可以用拟牛顿法解决。
- (3). Logistic Regression 可以对分类问题做出很好的拟合效果，正则项添加后可增强模型的泛化能力，

而变量之间的条件独立性虽为 Logistic Regression 的推导条件之一，但对模型的拟合效果影响不大。

5 参考文献

References

[1] Christopher M.Bishop. *Pattern Recognition and Machine Learning*. Springer, 2007.

[2] 周志华. 机器学习. 清华大学出版社, 2016.

[2] [1]

6 Codes

hyperparams.py

```
1
2 import numpy as np
3 import random
4 import matplotlib.pyplot as plt
5
6
7 class logistic(object):
8     def __init__(self, xdim=2, Wgaussian_hypo=False, independent=True):
9         self.gaussian = Wgaussian_hypo
10        self.X_vdim = xdim # dim = 1+xdim: exp(WTX)
11        self.xvalue_independent = independent
12        self.X_data = None
13        self.X_test = None
14        self.X = None
15        self.test_label = None
16        self.Y = None # X:(dim, L)=X_data.T X_data:(L, dim) Y:(L, 1) X_test:(dim, l_test)
17        self.global_value_initializer()
18
19    def global_value_initializer(self):
20        """
21        Automatically generate train/test matrixes, if have the input from outside, apply val
22        directly.
23        """
24        if self.xvalue_independent and self.X_vdim == 2:
25            pos = np.random.multivariate_normal([0, 2], cov=[[1, 0], [0, 1]], size=500) # pos
26            neg = np.random.multivariate_normal([4, 0], cov=[[2, 0], [0, 2]], size=500) # neg
27        elif not self.xvalue_independent and self.X_vdim == 2:
28            pos = np.random.multivariate_normal([0, 4], [[3, 1], [1, 4]], size=500)
29            neg = np.random.multivariate_normal([4, 0], [[3, 2], [2, 4]], size=500)
30        else:
31            pos, neg = [], []
32            mapping = [(i+list(item1), 0) for item1 in pos]
33            for item in neg:
34                mapping.append((i+list(item), 1))
35            np.random.shuffle(mapping)
36            self.Y = []
37            self.X_data = []
38            for item in mapping: # shuffle: reduce
39                self.X_data.append(item[0])
40                self.Y.append(item[1])
41            self.X_test = np.matrix(self.X_data[int(len(self.X_data)*5/7):]).T
42            self.X_data = np.matrix(self.X_data[:int(len(self.X_data)*5/7)])
43            self.X = self.X_data.T
44            self.test_label = np.matrix(self.Y[int(len(self.Y)*5/7):]).T
45            self.Y = np.matrix(self.Y[:int(len(self.Y)*5/7)]).T
46
47    def generate_W(self):
48        return np.matrix([random.gauss(0, 0.1) for i in range(self.X_vdim+1)]).T
49
50    def crossMatrix(self, X): # x: (?, dim): self.X.T, matrix
51        """
52        Generate Matrix:(Xi*Xj) for each sample l, return val: matrix(dim, dim, L)
53        """
```

```

53     ans = []
54     for s in range(np.size(X, 0)):
55         tmpans = [] # X[s]: matrix(1, dim)
56         mask = np.eye(np.size(X, 1))
57         for i in range(np.size(X, 1)):
58             tmpans.append([np.sum(np.dot(X[s], mask[i].T)*np.dot(X[s], mask[j].T)) for j in
                             range(np.size(X, 1))])
59         ans.append(tmpans)
60     return ans
61
62     def Loss(self, W):
63         return np.array(self.Y)*np.array(np.dot(self.X.T, W)) - np.log(1 + np.exp(np.dot(W.T,
        self.X))).T
64
65     def gradLoss(self, W): # W:(dim,1) Xi:(dim, 1) Y:(L, 1), X:(dim, L)
66         return np.dot(self.X, (self.Y.T - (np.exp(np.dot(W.T, self.X))/(1 + np.exp(np.dot(W.T,
        self.X))))).T)
67
68     def gradgradLoss(self, W):
69         val = np.exp(np.dot(W.T, self.X))
70         factor = (val/(np.array(1 + val)*np.array(1 + val))).T
71         cross = self.crossMatrix(self.X.T)
72         hessian = np.matrix(cross[0])*np.sum(factor[0])
73         for i in range(1, np.size(factor, 0)):
74             hessian += np.matrix(cross[i])*np.sum(factor[i])
75         return hessian
76
77     def newton_log(self, threshold=1e-8, alpha=1, lambdas=1e-2):
78         W, val_hold = self.generate_W(), 0 # (dim, 1)
79         while True:
80             step = np.dot(self.gradgradLoss(W).I, self.gradLoss(W))
81             if self.gussian:
82                 W += alpha*(step - lambdas*W)
83             else:
84                 W += alpha*step
85             print(W, np.sum(self.gradLoss(W)))
86             if abs(np.sum(self.gradLoss(W))) <= threshold or val_hold - np.sum(self.gradLoss(W))
            < 1e-5:
87                 break
88             else: val_hold = np.sum(self.gradLoss(W))
89         print(self.gradLoss(W))
90         return W
91
92     def gradient_descent(self, threshold=1e-4, alpha=5e-8, lambdas=1e-2):
93         W = self.generate_W()
94         while True:
95             step = self.gradLoss(W)
96             if self.gussian:
97                 W += alpha*(step - lambdas*W)
98             else:
99                 W += alpha*step
100             print(W)
101             print(np.sum(self.gradLoss(W)))
102             if abs(np.sum(self.gradLoss(W))) <= threshold:
103                 break
104         return W
105
106     def test(self, W):
107         val = (1/(1+np.exp(np.dot(W.T, self.X_test)))).T
108         expect = [int(np.sum(val[i]) < 0.5) for i in range(np.size(val, 0))]
109         return 1 - abs(np.sum(abs(expect - self.test_label.T))/np.size(val, 0))
110
111     def pr_plot(self, W, flag=""):
112         if self.X_vdim == 2:
113             w0, w1, w2 = np.sum(W[0]), np.sum(W[1]), np.sum(W[2])
114             plt_X = np.linspace(np.min(self.X[1]), np.max(np.max(self.X[2])))
115             plt_Y = [-w0/w2 - w1/w2*x for x in plt_X]
116             mask = np.eye(3)
117             scatter = [(np.sum(np.dot(item, mask[1])), np.sum(np.dot(item, mask[2]))) for item in
                self.X_data]
118             print(scatter)
119             label_pt = self.Y.T
120             scat_pos = [[], []]
121             scat_neg = [[], []]
122             for i in range(len(scatter)):
123                 if np.sum(self.Y[i]) == 0:
124                     scat_neg[0].append(scatter[i][0])

```



```

125         scat_neg[1].append(scatter[i][1])
126     elif np.sum(self.Y[i]) == 1:
127         scat_pos[0].append(scatter[i][0])
128         scat_pos[1].append(scatter[i][1])
129     else:
130         print("Error while map/reduce the points.")
131     plt.scatter(scat_neg[0], scat_neg[1], s=75, alpha=.5, color='purple')
132     plt.scatter(scat_pos[0], scat_pos[1], s=75, alpha=.5, color='orange')
133     plt.plot(plt_X, plt_Y, color='green')
134     plt.title("Logistic Regression: gaussian: %s ind: %s %s test res: %f" % (str(self.gaussian)
135         ), str(self.xvalue_independent), flag, self.test(W)))
136     plt.savefig("LogisticRegression_gaussian=%s ind=%s %s.png" % (str(self.gaussian), str(
137         self.xvalue_independent), flag))
138     plt.show()
139 else:
140     pass # go to test
141
142 def wine_data_insert():
143     f = open("winedata.csv", "r")
144     classify, classify_test, X_data, X_test = [], [], [], []
145     for line in f:
146         line = line.split(",")
147         # if int(line[0]) == 3: continue
148         classify.append(int(line[0]) - 1)
149         tmp = [1]
150         for i in range(1, len(line)):
151             tmp.append(float(line[i]))
152         X_data.append((tmp, int(line[0]) - 1))
153     random.shuffle(X_data)
154     X_test = X_data[int(len(X_data)*5/7):]
155     classify_test = [item[1] for item in X_test]
156     X_test = [item[0] for item in X_test]
157     X_data = X_data[:int(len(X_data)*5/7)]
158     classify = [item[1] for item in X_data]
159     X_data = [item[0] for item in X_data]
160     f, fl = open("wine_test.txt", "w"), open("wine_test_label.txt", "w")
161     for (item, label) in zip(X_test, classify_test):
162         for it in item: f.write(str(it) + "\t")
163         f.write("\n")
164         fl.write(str(label) + "\n")
165     f.close(), fl.close()
166     fr, frl = open("wine_train.txt", "w"), open("wine_train_label.txt", "w")
167     for (item, label) in zip(X_data, classify):
168         for it in item: fr.write(str(it) + "\t")
169         fr.write("\n")
170         frl.write(str(label) + "\n")
171     fr.close(), frl.close()
172     return np.matrix(X_data), np.matrix(X_test).T, np.matrix(classify).T, np.matrix(
173         classify_test).T
174
175 def ReadFromStandardFile(mask=0):
176     X_train, X_test, Y_train, Y_test = [], [], [], []
177     f, fl = open("wine_test.txt", "r"), open("wine_test_label.txt", "r")
178     for line in f:
179         line = line.split()
180         X_test.append([float(line[i]) for i in range(0, len(line))])
181     for line in fl:
182         if mask != -1: Y_test.append([int(float(line) == mask)])
183         else: Y_test.append([int(line)])
184     fr, frl = open("wine_train.txt", "r"), open("wine_train_label.txt", "r")
185     for line in fr:
186         print(line)
187         line = line.split()
188         X_train.append([float(line[i]) for i in range(0, len(line))])
189     for line in frl:
190         Y_train.append([int(float(line) == mask)])
191     return np.matrix(X_train), np.matrix(Y_train), np.matrix(X_test).T, np.matrix(Y_test)
192
193 def decision_tree(): # precision: 0.705882
194     # classify label 1 first(original:2), then classify label 0 (ori:1), the other point is
195     # label 2 (ori:3)
196     W_cl1 = [[-0.19624184], [0.14048495], [-0.32750738], [0.04732983], [0.1123163],
197         [0.02439485], [0.0467235],

```

```

196         [0.09784779], [0.04087847], [-0.18535481], [-0.23524893], [-0.08212339],
197         [-0.05971632], [-0.00696094]]
198     W_cl0 = [[0.0187151], [0.00048866], [-0.10105207], [0.06552508], [-0.46800525],
199             [-0.04441048], [-0.02615863],
200             [0.13665488], [0.08469632], [-0.00764128], [-0.05588578], [-0.02218694],
201             [0.03718035], [0.01629198]]
202     W_cl1 = np.matrix(W_cl1)
203     W_cl0 = np.matrix(W_cl0)
204     X_test, Y_test = [], []
205     f, fl = open("wine_test.txt", "r"), open("wine_test_label.txt", "r")
206     for line in f:
207         line = line.split()
208         X_test.append([float(line[i]) for i in range(0, len(line))])
209     for line in fl:
210         Y_test.append(int(line))
211     X_test = np.matrix(X_test).T
212     val1 = (1/(1+np.exp(np.dot(W_cl1.T, X_test)))).T
213     expect1 = [int(np.sum(val1[i]) < 0.5) for i in range(np.size(val1, 0))]
214     val0 = (1/(1+np.exp(np.dot(W_cl0.T, X_test)))).T
215     expect0 = [int(np.sum(val0[i]) < 0.5) for i in range(np.size(val0, 0))]
216     print(expect0)
217     print(expect1)
218     print(Y_test)
219     ans = []
220     for i in range(0, len(expect0)):
221         if expect1[i] == 1:
222             ans.append(1)
223         elif expect0[i] == 1:
224             ans.append(0)
225         else:
226             ans.append(2)
227     cnt = 0
228     print(ans)
229     for i in range(0, len(ans)):
230         if ans[i] == Y_test[i]:
231             cnt += 1
232     print("Precision = %lf" % (cnt/len(ans)))
233     return cnt/len(ans)
234
235 if __name__ == "__main__":
236     """
237     Normal Train, test
238     """
239     log = logistic(Wgaussian_hypo=True, independent=True)
240     w = log.gradient_descent()
241     print("RES:")
242     print(w)
243     log.pr_plot(w, "testgrad")
244     """
245     Train, test from dataset: UCI: https://archive.ics.uci.edu/ml/machine-learning-databases
246     /wine
247     classify 0: W :0.90196
248     [[ 0.11116817]
249     [-0.18515697]
250     [ 0.05445636]
251     [-0.07500515]
252     [-0.41380713]
253     [-0.03442779]
254     [ 0.09983258]
255     [ 0.01055378]
256     [-0.01554494]
257     [ 0.02209726]
258     [-0.07036363]
259     [ 0.0266265 ]
260     [ 0.0548774 ]
261     [ 0.01650266]]
262     classify:1 W :0.882352941176
263     [[ 0.0088141 ]
264     [ 0.00994699]
265     [-0.1510287 ]
266     [ 0.0400127 ]
267     [ 0.15663027]
268     [ 0.02726139]
269     [ 0.10969113]
270     [ 0.27900054]

```

```

269     [ 0.06219452]
270     [ 0.15556723]
271     [-0.54909899]
272     [ 0.12975475]
273     [ 0.2285545 ]
274     [-0.0077714 ]]
275     class 2: 0.90196
276
277     class 1-> 0: 0.705882
278
279     Wine data/classification below:
280     """
281     log = logistic(xdim=13)
282     # log.X_data, log.X_test, log.Y, log.test_label = wine_data_insert()
283     log.X_data, log.Y, log.X_test, log.test_label = ReadFromStandardFile(mask=2) # 1
284     log.X = log.X_data.T
285     print(log.X)
286     print(log.Y)
287     print(log.test_label)
288     w = log.gradient_descent(threshold=1e-5)
289     print(log.test(w))
290     decision_tree()

```

7 End

Machine Learning: Logistic Regression

Copyright ©RuiKang(marisuki) 2018