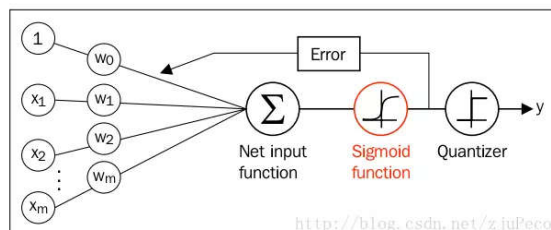
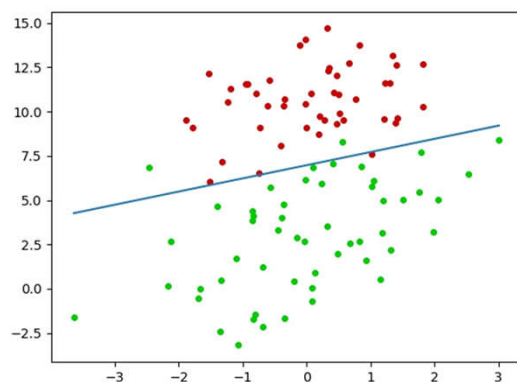


Python编程与人工智能实践



算法篇：逻辑回归 (logistic regression)



于泓
鲁东大学
信息与电气工程学院
2021.4.1

逻辑回归 (logistic regression)

- 回归：假设有一些数据点。利用一条直线对这些点进行拟合就叫做回归
- 逻辑回归：回归的目标是一个二值的结果 (0,1) ，是最常用的一种分类模型。

已知一组数据 \mathbf{X} 以及其相应的二值标签 \mathbf{Y}
利用这些数据构造一个函数，预测新的输入 X_i 的标签

$$X_1 \text{ ----} \rightarrow Y_1 \text{ (0)}$$

$$X_2 \text{ ----} \rightarrow Y_2 \text{ (1)}$$

$$X_3 \text{ ----} \rightarrow Y_3 \text{ (1)}$$

$$X_4 \text{ ----} \rightarrow Y_4 \text{ (0)}$$

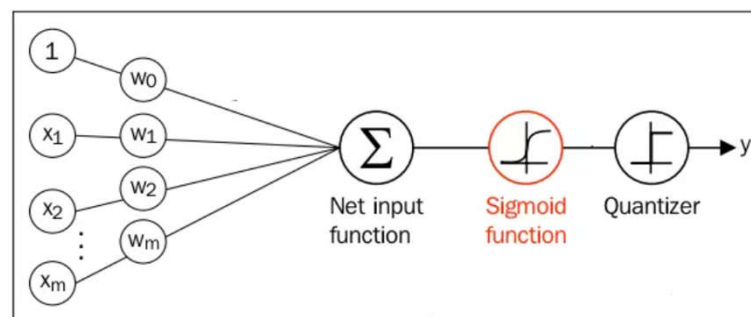
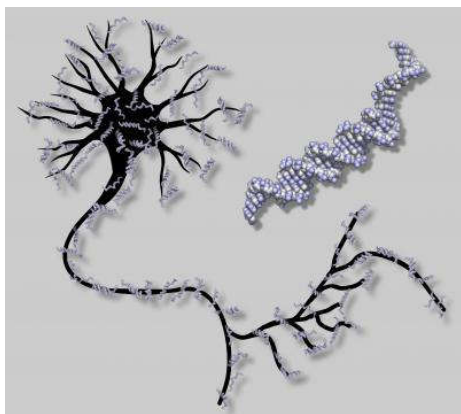
$$X_5 \text{ ----} \rightarrow Y_5 \text{ (1)}$$

$$X_6 \text{ ----} \rightarrow Y_6 \text{ (0)}$$

$$X_7 \text{ ----} \rightarrow Y_7 \text{ (1)}$$

.....

$$X_N \text{ ----} \rightarrow Y_N \text{ (1)}$$



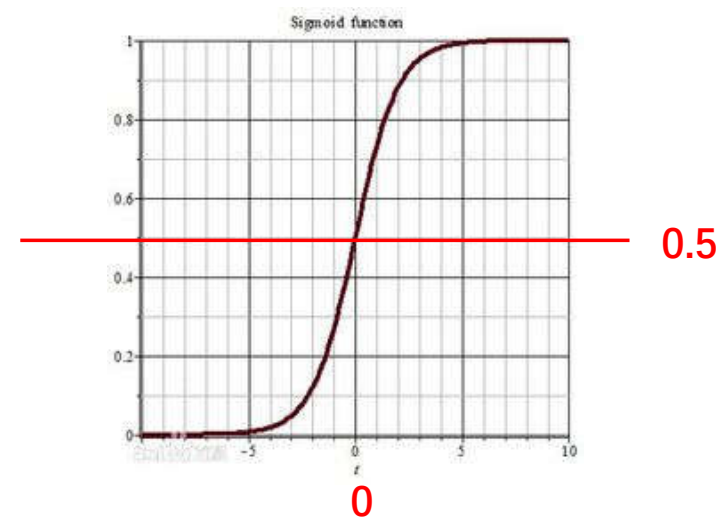
假设 X 是一个2维的特征

$$h(X; \mathbf{w}) = \text{sigmoid}(x_2 * w_2 + x_1 * w_1 + w_0) = \text{sigmoid}(X\mathbf{w})$$

其中

$$\sigma(z) = \text{sigmoid}(z) = \frac{1}{1 + e^{-z}}$$

输出范围 $[0,1]$, 可以把输出理解为 $Y=1$ 的概率
 概率 >0.5 输出标签1, 否则输出0



已知输入特征 X_i , 预测输出标签 $Y_i=1$ 的概率为:

$$P(Y_i=1|X_i;\mathbf{w})=h(X_i;\mathbf{w})=\sigma(X_i\mathbf{w})$$

已知输入特征 X_i , 预测输出标签 $Y_i=0$ 的概率为:

$$P(Y_i=0|X_i;\mathbf{w})=1-h(X_i;\mathbf{w})$$

已知输入特征 X_i 以及对应的标签 Y_i
预测准确的概率为:

$$P(Y_i, X_i;\mathbf{w})=[h(X_i;\mathbf{w})]^{Y_i} [1-h(X_i;\mathbf{w})]^{1-Y_i}$$

预测的越准 上述函数就越大

根据已知数据, 寻找最优参数 \mathbf{w} , 就是令

$$L(\mathbf{w})=\prod_{i=1}^N P(Y_i, X_i;\mathbf{w})$$

$$=\prod_{i=1}^N [h(X_i;\mathbf{w})]^{Y_i} [1-h(X_i;\mathbf{w})]^{1-Y_i} \quad \text{最大}$$

取log写成似然函数的形式

$$L(\mathbf{w})=\sum_{i=1}^N Y_i \log(h(X_i;\mathbf{w})) + (1-Y_i) \log(1-h(X_i;\mathbf{w}))$$

最大

再加一个“-”号上式就变成了熵的形式 (交叉熵)

$$L(\mathbf{w})=-\sum_{i=1}^N Y_i \log(h(X_i;\mathbf{w})) + (1-Y_i) \log(1-h(X_i;\mathbf{w}))$$

熵越小系统越确定,分类也就越准确

逻辑回归的最终目标变为：

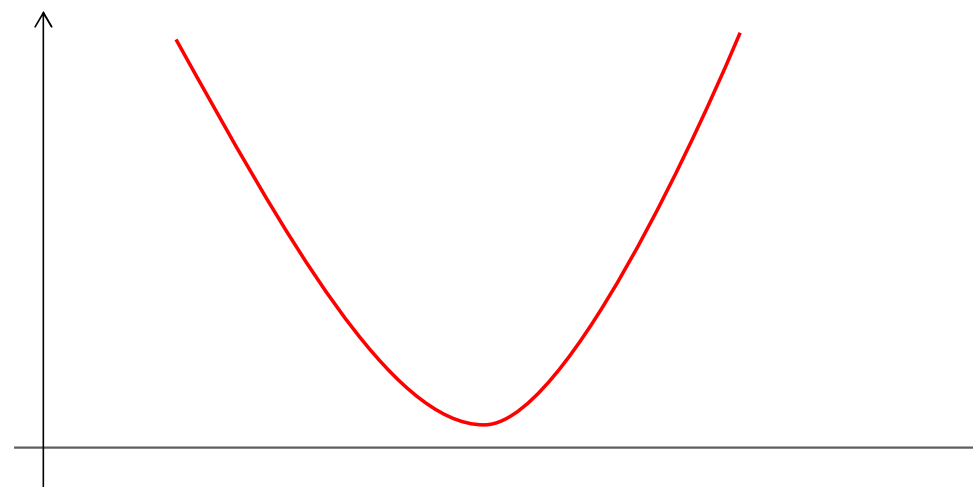
已知训练数据 X 以及相应标签 Y

求解参数 w 使交叉熵 $L(w)$ 最小

$$L(w) = -\sum_{i=1}^N Y_i \ln(h(X_i; w)) + (1 - Y_i) \ln(1 - h(X_i; w))$$

上述方程没有闭式解需要通过梯度下降的方法计算最优解。

$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial L(w)}{\partial w}$$



令 $z_i = X_i \mathbf{w}$

$$\begin{aligned}
 L(\mathbf{w}) &= -\sum_{i=1}^N Y_i \ln(h(X_i; \mathbf{w})) + (1 - Y_i) \ln(1 - h(X_i; \mathbf{w})) \\
 &= -\sum_{i=1}^N Y_i \ln \frac{h(z_i)}{1 - h(z_i)} + \ln(1 - h(z_i)) \\
 &= -\sum_{i=1}^N z_i Y_i - z_i - \ln(1 + e^{-z_i})
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial L(\mathbf{w})}{\partial z} \frac{\partial z}{\partial \mathbf{w}} = -\sum_{i=1}^N \left(Y_i - \frac{e^{-z_i}}{1 + e^{-z_i}} \right) \mathbf{x}_i^T \\
 &= -\sum_{i=1}^N \left(Y_i - \frac{1}{1 + e^{-z_i}} \right) \mathbf{x}_i = -\sum_{i=1}^N (Y_i - h(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i^T
 \end{aligned}$$

$$\begin{aligned}
 \ln \frac{h(z)}{1 - h(z)} &= \ln \frac{\frac{1}{1 + e^{-z}}}{1 - \frac{1}{1 + e^{-z}}} = \ln \frac{1}{e^{-z}} = z \\
 \ln(1 - h(z)) &= \ln \left(1 - \frac{1}{1 + e^{-z}} \right) = \ln \frac{e^{-z}}{1 + e^{-z}} = -z - \ln(1 + e^{-z})
 \end{aligned}$$

代入

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \alpha \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{w}_{\text{old}} + \alpha \sum_{i=1}^N (Y_i - h(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i^T$$

标签
实际输出

写成矩阵形式

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \alpha \mathbf{X}^T (\mathbf{Y} - h(\mathbf{X}; \mathbf{w}_{\text{old}}))$$

$$\mathbf{X} \in [N \times D]$$

$$\mathbf{Y} \in [N \times 1]$$

$$\mathbf{w} \in [D \times 1]$$

定义sigmoid 函数

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def sigmoid(z):
5     return 1.0/(1+np.exp(-z))
6

```

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \alpha \mathbf{X}^T (\mathbf{Y} - h(\mathbf{X}; \mathbf{w}_{\text{old}}))$$

$$\mathbf{X} \in [N \times D]$$

$$\mathbf{Y} \in [N \times 1]$$

$$\mathbf{w} \in [D \times 1]$$

```

def weight_update(datas, labs, w, alpha=0.01):
    z = np.dot(datas, w) # Nx1
    h = sigmoid(z)       # Nx1
    Error = labs - h     # Nx1
    w = w + alpha * np.dot(datas.T, Error)
    return w

```

进行训练, 返回w

```

def train_LR(datas, labs, n_epoch=2, alpha=0.005):
    N, D = np.shape(datas)
    w = np.ones([D, 1]) # Dx1
    # 进行n_epoch轮迭代
    for i in range(n_epoch):
        w = weight_update(datas, labs, w, alpha)
        error_rate = test_accuracy(datas, labs, w)
        print("epoch %d error %.3f%%" % (i, error_rate * 100))
    return w

```

```

def test_accuracy(datas, labs, w):
    N, D = np.shape(datas)
    z = np.dot(datas, w) # Nx1
    h = sigmoid(z)       # Nx1
    lab_det = (h > 0.5).astype(np.float)
    error_rate = np.sum(np.abs(labs - lab_det)) / N
    return error_rate

```

测试程序

```
def load_dataset(file):
    with open(file,"r") as f:
        lines = f.read().splitlines()

    # 取 lab 维度为 N x 1
    labs = [line.split("\t")[-1] for line in lines]
    labs = np.array(labs).astype(np.float)
    labs = np.expand_dims(labs,axis=-1) # Nx1

    # 取数据 增加 一维全是1的特征
    datas = [line.split("\t")[:-1] for line in lines]
    datas = np.array(datas).astype(np.float)
    N,D = np.shape(datas)
    # 增加一个维度
    datas = np.c_[np.ones([N,1]),datas]
    return datas,labs
```

数据增加
全是1的维度

testset.txt	logisticRegression.py	PeakDensity.
1	-0.017612	14.053064
2	-1.395634	4.662541
3	-0.752157	6.538620
4	-1.322371	7.152853
5	0.423363	11.054677
6	0.406704	7.067335
7	0.667394	12.741452
8	-2.460150	6.866805
9	0.569411	9.548755
10	-0.026632	10.427743
11	0.850433	6.920334
12	1.347183	13.175500
13	1.176813	3.167020
14	-1.781871	9.097953
15	-0.566606	5.749003
16	0.931635	1.589505
17	-0.024205	6.151823
18	-0.036453	2.690988
19	-0.196949	0.444165
20	1.014459	5.754399
21	1.985298	3.230619
22	-1.693453	-0.557540
23	-0.576525	11.778922
24	-0.346811	-1.678730
25	-2.124484	2.672471
26	1.217916	9.597015
27	0.722000	0.000000


```
def draw_design_line(datas, labs, w, name="0.jpg"):
    dic_colors={0:(.8,0,0),1:(0,.8,0)}

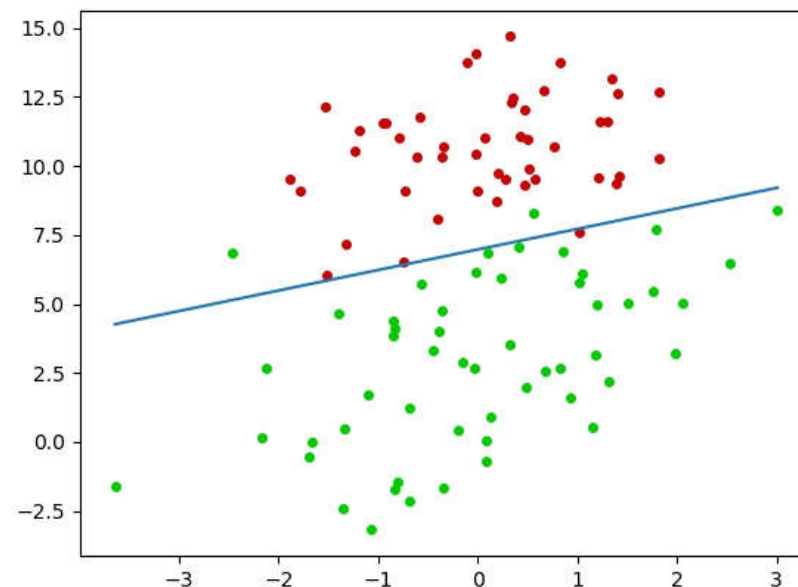
    # 画数据点
    for i in range(2):
        index = np.where(labs==i)[0]
        sub_datas = datas[index]
        plt.scatter(sub_datas[:,1],sub_datas[:,2],s=16.,color=dic_colors[i])

    # 画判决线
    min_x = np.min(datas[:,1])
    max_x = np.max(datas[:,1])
    w = w[:,0]
    x = np.arange(min_x,max_x,0.01)
    y = -(x*w[1]+w[0])/w[2]
    plt.plot(x,y)

    plt.savefig(name)

if __name__ == "__main__":
    ''' 实验1 基础测试数据 '''
    # 加载数据
    file = "testset.txt"
    datas, labs = load_dataset(file)

    weights = train_LR(datas, labs, alpha=0.002, n_epoch=550)
    print(weights)
    draw_design_line(datas, labs, weights)
```



随机梯度下降:

将数据分成小的Batch，每个Batch更新一次，可以增加训练的速度和稳定性

```
# 随机梯度下降
def train_LR_batch(datas, labs, batchsize, n_epoch=2, alpha=0.005):
    N, D = np.shape(datas)
    # weight 初始化
    w = np.ones([D, 1]) # Dx1
    N_batch = N // batchsize
    for i in range(n_epoch):
        # 数据打乱
        rand_index = np.random.permutation(N).tolist()
        # 每个batch 更新一下weight
        for j in range(N_batch):
            index = rand_index[j*batchsize:(j+1)*batchsize]
            batch_datas = datas[index]
            batch_labs = labs[index]
            w = weight_update(batch_datas, batch_labs, w, alpha)

        error = test_accuracy(datas, labs, w)
        print("epoch %d error %.2f%%" % (i, error*100))
    return w
```

```
'''实验2 马疝数据集'''  
# 加载训练数据  
train_file = "horse_train.txt"  
train_datas, train_labs = load_dataset(train_file)  
  
# 加载测试数据  
test_file = "horse_test.txt"  
test_datas, test_labs = load_dataset(test_file)  
  
# 梯度下降  
# weights = train_LR(train_datas, train_labs, alpha=0.001, n_epoch=90)  
# print(weights)  
  
# 随机梯度下降  
weights = train_LR_batch(train_datas, train_labs, batchsize=2, n_epoch=50, alpha=0.001)  
print(weights)  
  
acc = test_accuracy(test_datas, test_labs, weights)  
print(acc)  
  
# 截取几个维度画图  
index = [0, 4, 5]  
sub_datas = train_datas[:, index]  
sub_weights = weights[index]  
draw_desion_line(sub_datas, train_labs, sub_weights, name="horse.jpg")
```