

机器学习实战教程（六）：Logistic回归基础篇之梯度上升算法

🕒 2017年11月8日 18:42:19 🗨 93 👁 16,740 🌡 °C 📄 编辑



一、前言

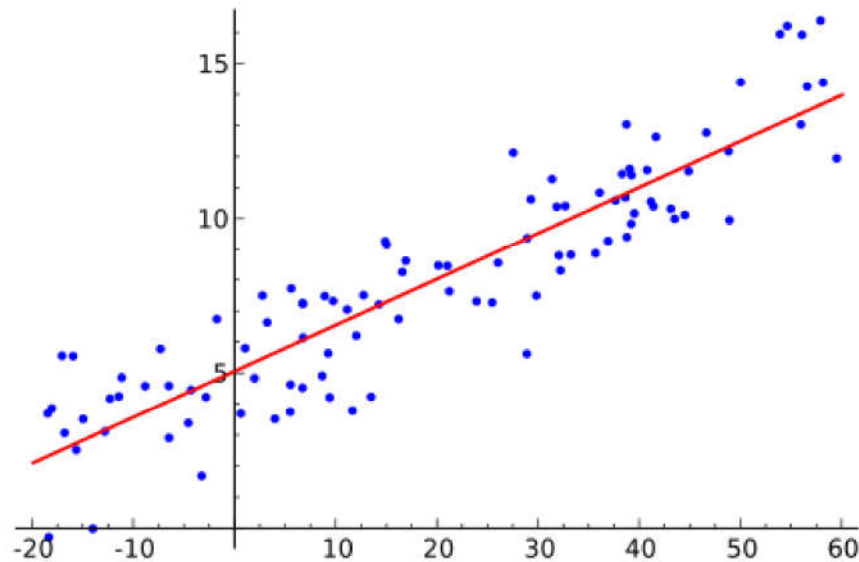
本文从Logistic回归的原理开始讲起，补充了书上**省略的数学推导**。本文可能会略显枯燥，理论居多，Sklearn实战内容会放在下一篇文章。**自己慢慢推导完公式，还是蛮开心的一件事。**

二、Logistic回归与梯度上升算法

Logistic回归是众多分类算法中的一员。通常，Logistic回归用于二分类问题，例如预测明天是否会下雨。当然它也可以用于多分类问题，不过为了简单起见，本文暂先讨论二分类问题。首先，让我们来了解一下，什么是Logistic回归。

1、Logistic回归

假设现在有一些数据点，我们利用一条直线对这些点进行拟合(该线称为最佳拟合直线)，这个拟合过程就称作为回归，如下图所示：



Logistic回归是分类方法，它利用的是Sigmoid函数阈值在[0,1]这个特性。**Logistic回归进行分类的主要思想是：根据现有数据对分类边界线建立回归公式，以此进行分类。**其实，Logistic本质上是一个基于条件概率的判别模型(Discriminative Model)。

所以要想了解Logistic回归，我们必须先看一看Sigmoid函数，我们也可以称它为Logistic函数。它的公式如下：

$$h_{\theta}(x) = g(\theta^T x)$$

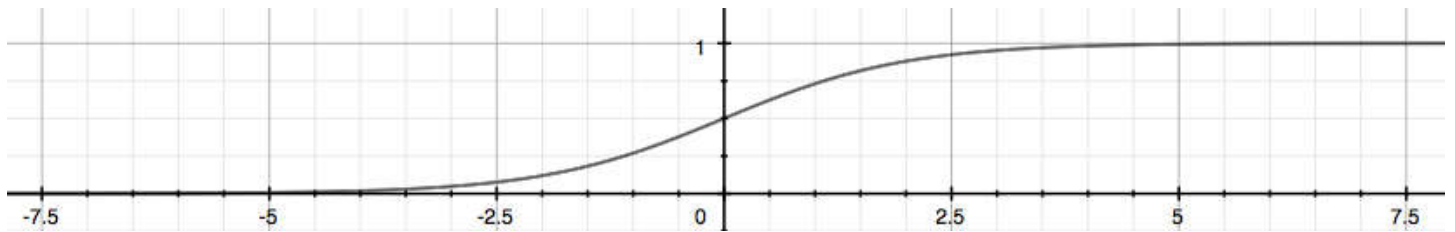
$$z = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

整合成一个公式，就变成了如下公式：

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}},$$

下面这张图片，为我们展示了Sigmoid函数的样子。



z 是一个矩阵， θ 是参数列向量(要求解的)， x 是样本列向量(给定的数据集)。 θ^T 表示 θ 的转置。 $g(z)$ 函数实现了任意实数到 $[0,1]$ 的映射，这样我们的数据集 $([x_0, x_1, \dots, x_n])$ ，不管是大于1或者小于0，都可以映射到 $[0,1]$ 区间进行分类。 $h_{\theta}(x)$ 给出了输出为1的概率。比如当 $h_{\theta}(x)=0.7$ ，那么说明有70%的概率输出为1。输出为0的概率是输出为1的补集，也就是30%。

如果我们有合适的参数列向量 $\theta([\theta_0, \theta_1, \dots, \theta_n]^T)$ ，以及样本列向量 $x([x_0, x_1, \dots, x_n])$ ，那么我们对样本 x 分类就可以通过上述公式计算出一个概率，如果这个概率大于0.5，我们就可以说样本是正样本，否则样本是负样本。

举个例子，对于"垃圾邮件判别问题"，对于给定的邮件(样本)，我们定义非垃圾邮件为正类，垃圾邮件为负类。我们通过计算出的概率值即可判定邮件是否是垃圾邮件。

那么问题来了！如何得到合适的参数向量 θ ？

根据sigmoid函数的特性，我们可以做出如下的假设：

$$\begin{aligned} P(y = 1 \mid x; \theta) &= h_{\theta}(x) \\ P(y = 0 \mid x; \theta) &= 1 - h_{\theta}(x) \end{aligned}$$

式即为在已知样本 x 和参数 θ 的情况下，样本 x 属性正样本($y=1$)和负样本($y=0$)的条件概率。理想状态下，根据上述公式，求出各个点的概率均为1，也就是完全分类都正确。但是考虑到实际情况，样本点的概率越接近于1，其分类效果越好。比如一个样本属于正样本的概率为0.51，那么我们就可以说明这个样本属于正样本。另一个样本属于正样本的概率为0.99，那么我们也可以说明这个样本属于正样本。但是显然，第二个样本概率更高，更具说服力。我们可以把上述两个概率公式合二为一：

$$\text{Loss}(h_{\theta}(x), y) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{(1-y)}$$

合并出来的Loss，我们称之为损失函数(Loss Function)。当 y 等于1时， $(1-y)$ 项(第二项)为0；当 y 等于0时， y 项(第一项)为0。为了简化问题，我们对整个表达式求对数，(将指数问题对数化是处理数学问题常见的方法)：

$$\text{Loss}(h_{\theta}(x), y) = y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x))$$

这个损失函数，是对于一个样本而言的。给定一个样本，我们就可以通过这个损失函数求出，样本所属类别的概率，而这个概率越大越好，所以也就是求解这个损失函数的最大值。既然概率出来

了，那么最大似然估计也该出场了。假定样本与样本之间相互独立，那么整个样本集生成的概率即为所有样本生成概率的乘积，再将公式对数化，便可得到如下公式：

$$J(\theta) = \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

其中， m 为样本的总数， $y^{(i)}$ 表示第 i 个样本的类别， $x^{(i)}$ 表示第 i 个样本，需要注意的是 θ 是多维向量， $x^{(i)}$ 也是多维向量。

综上所述，满足 $J(\theta)$ 的最大的 θ 值即是我们要求解的模型。

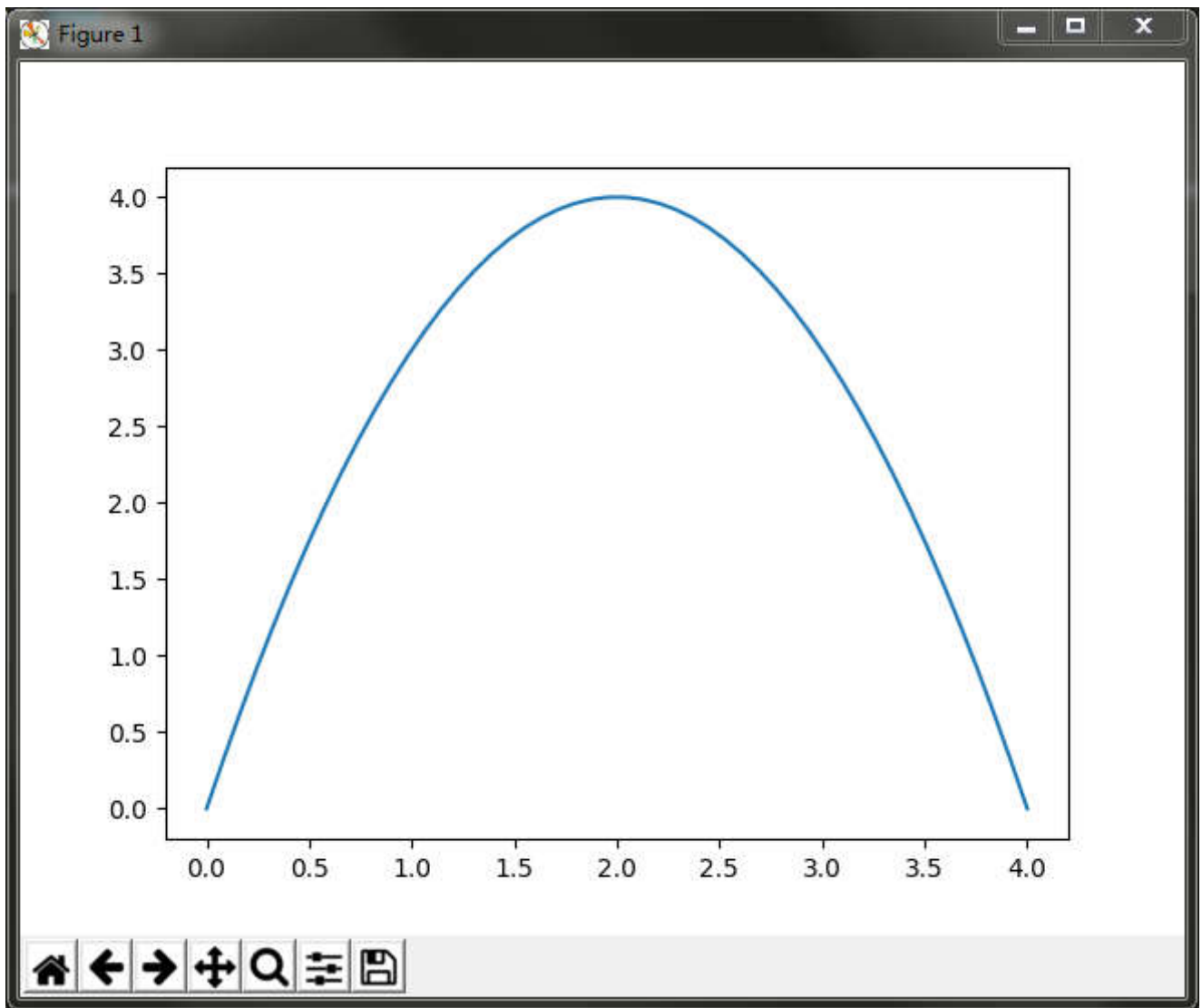
怎么求解使 $J(\theta)$ 最大的 θ 值呢？因为是求最大值，所以我们需要使用梯度上升算法。如果面对的问题是求解使 $J(\theta)$ 最小的 θ 值，那么我们就需要使用梯度下降算法。面对我们这个问题，如果使 $J(\theta) := -J(\theta)$ ，那么问题就从求极大值转换成求极小值了，使用的算法就从梯度上升算法变成了梯度下降算法，它们的思想都是相同的，学会其一，就也会了另一个。本文使用梯度上升算法进行求解。

2、梯度上升算法

说了半天，梯度上升算法又是啥？ $J(\theta)$ 太复杂，我们先看个简单的求极大值的例子。一个看了就会想到高中生活的函数：

$$f(x) = -x^2 + 4x$$

来吧，做高中题。这个函数的极值怎么求？显然这个函数开口向下，存在极大值，它的函数图像为：



求极值，先求函数的导数：

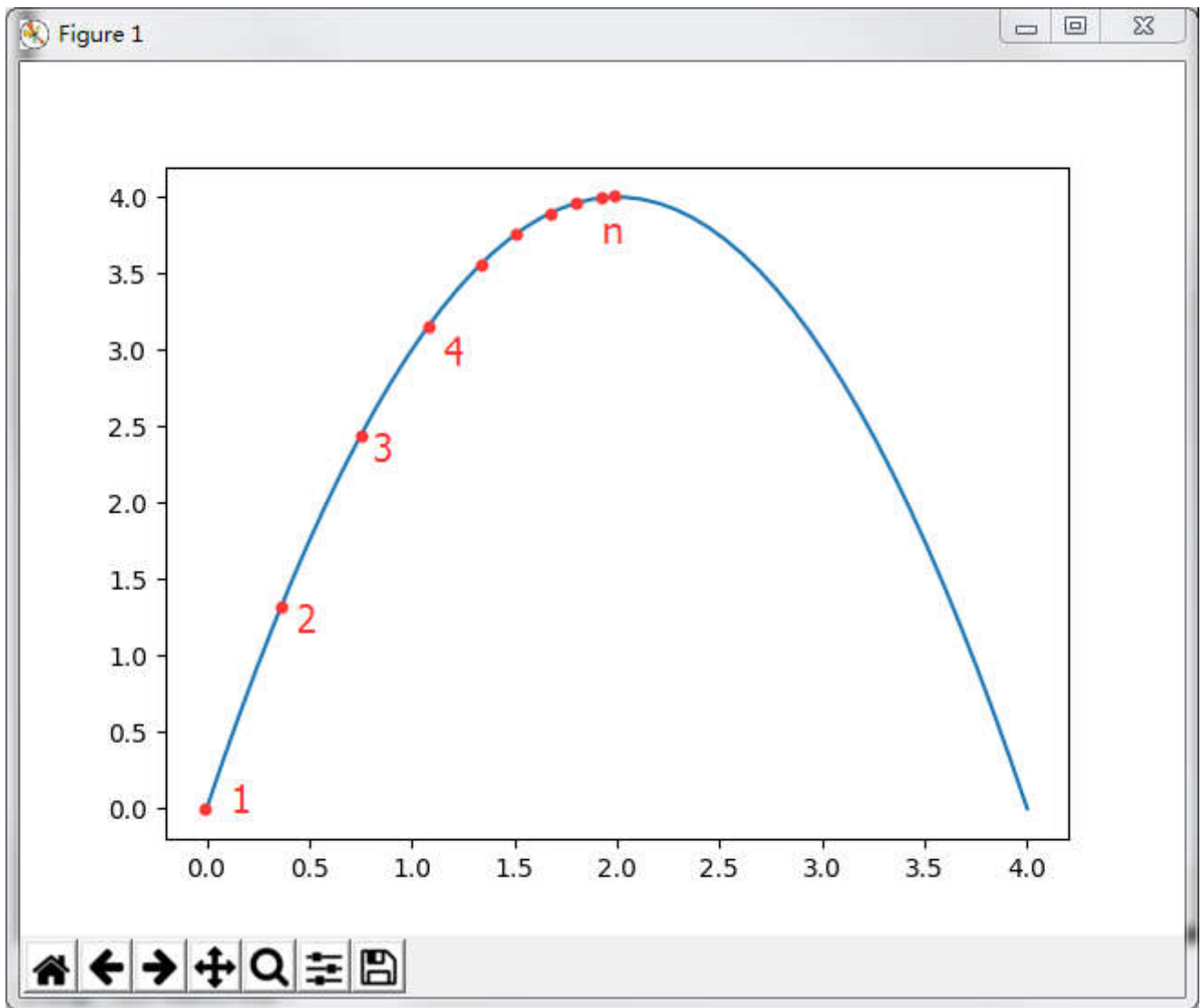
$$f'(x) = -2x + 4$$

令导数为0，可求出 $x=2$ 即取得函数 $f(x)$ 的极大值。极大值等于 $f(2)=4$

但是真实环境中的函数不会像上面这么简单，就算求出了函数的导数，也很难精确计算出函数的极值。此时我们就可以用迭代的方法来做。就像爬坡一样，一点一点逼近极值。这种寻找最佳拟合参数的方法，就是最优化算法。爬坡这个动作作用数学公式表达即为：

$$x_{i+1} = x_i + \alpha \frac{\partial f(x_i)}{\partial x_i}$$

其中， α 为步长，也就是学习速率，控制更新的幅度。效果如下图所示：



比如从(0,0)开始，迭代路径就是1->2->3->4->...->n，直到求出的x为函数极大值的近似值，停止迭代。我们可以编写Python3代码，来实现这一过程：

	Python
1	# -*- coding:UTF-8 -*-
2	"""
3	函数说明:梯度上升算法测试函数
4	
5	求函数 $f(x) = -x^2 + 4x$ 的极大值
6	
7	Parameters:
8	无
9	Returns:
10	无
11	Author:
12	Jack Cui
13	Blog:
14	http://blog.csdn.net/c406495762
15	Zhihu:
16	https://www.zhihu.com/people/Jack--Cui/
17	Modify:
18	2017-08-28
19	"""
20	def Gradient_Ascent_test():
21	def f_prime(x_old):
22	return -2 * x_old + 4
23	x_old = -1
24	x_new = 0
25	alpha = 0.01
26	presision = 0.00000001
	#f(x)的导数
	#初始值，给一个小于x_new的值
	#梯度上升算法初始值，即从(0,0)开始
	#步长，也就是学习速率，控制更新的幅度
	#精度，也就是更新阈值


```

27 while abs(x_new - x_old) > presision:
28     x_old = x_new
29     x_new = x_old + alpha * f_prime(x_old)           #上面提到的公式
30     print(x_new)                                     #打印最终求解的极值近似值
31
32 if __name__ == '__main__':
33     Gradient_Ascent_test()

```

代码运行结果如下：

```

94 if __name__ == '__main__':
95     Gradient_Ascent_test()

```

1.999999515279857
[Finished in 0.3s]

结果很显然，已经非常接近我们的真实极值2了。这一过程，就是梯度上升算法。那么同理， $J(\theta)$ 这个函数的极值，也可以这么求解。公式可以这么写：

$$\theta_j := \theta_j + \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

由上小节可知 $J(\theta)$ 为：

$$J(\theta) = \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

sigmoid函数为：

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}},$$

那么，现在我只要求出 $J(\theta)$ 的偏导，就可以利用梯度上升算法，求解 $J(\theta)$ 的极大值了。

那么现在开始求解 $J(\theta)$ 对 θ 的偏导，求解如下：

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial J(\theta)}{\partial g(\theta^T x)} * \frac{\partial g(\theta^T x)}{\partial \theta^T x} * \frac{\partial \theta^T x}{\partial \theta_j}$$

其中：

$$\frac{\partial J(\theta)}{\partial g(\theta^T x)} = y * \frac{1}{g(\theta^T x)} + (y - 1) * \frac{1}{1 - g(\theta^T x)}$$

再由：

$$\begin{aligned} g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} = \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\ &= \frac{1}{(1 + e^{-z})} \left(1 - \frac{1}{(1 + e^{-z})} \right) = g(z)(1 - g(z)) \end{aligned}$$

可得：

$$\frac{\partial g(\theta^T x)}{\partial \theta^T x} = g(\theta^T x)(1 - g(\theta^T x))$$

接下来，就剩下第三部分：

$$\frac{\partial \theta^T x}{\partial \theta_j} = \frac{\partial J(\theta_1 x_1 + \theta_2 x_2 + \cdots \theta_n x_n)}{\partial \theta_j} = x_j$$

综上所述：

$$\frac{\partial}{\partial \theta_j} J(\theta) = (y - h_{\theta}(x))x_j$$

因此，梯度上升迭代公式为：

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))x_j^{(i)}$$

知道了，梯度上升迭代公式，我们就可以自己编写代码，计算最佳拟合参数了。

三、Python3实战

1、数据准备

数据集已经为大家准备好，下载地址：[数据集下载](#)

这就是一个简单的数据集，没什么实际意义。让我们先从这个简单的数据集开始学习。先看下数据集有哪些数据：

1	-0.017612	14.053064	0
2	-1.395634	4.662541	1
3	-0.752157	6.538620	0
4	-1.322371	7.152853	0
5	0.423363	11.054677	0
6	0.406704	7.067335	1
7	0.667394	12.741452	0
8	-2.460150	6.866805	1
9	0.569411	9.548755	0
10	-0.026632	10.427743	0
11	0.850433	6.920334	1
12	1.347183	13.175500	0
13	1.176813	3.167020	1
14	-1.781871	9.097953	0
15	-0.566606	5.749003	1
16	0.931635	1.589505	1
17	-0.024205	6.151823	1
18	-0.036453	2.690988	1
19	-0.196949	0.444165	1
20	1.014459	5.754399	1
21	1.985298	3.230619	1
22	-1.693453	-0.557540	1

这个数据有两维特征，因此可以将数据在一个二维平面上展示出来。我们可以将第一列数据(X1)看作x轴上的值，第二列数据(X2)看作y轴上的值。而最后一列数据即为分类标签。根据标签的不同，对这些点进行分类。

那么，先让我们编写代码，看下数据集的分布情况：

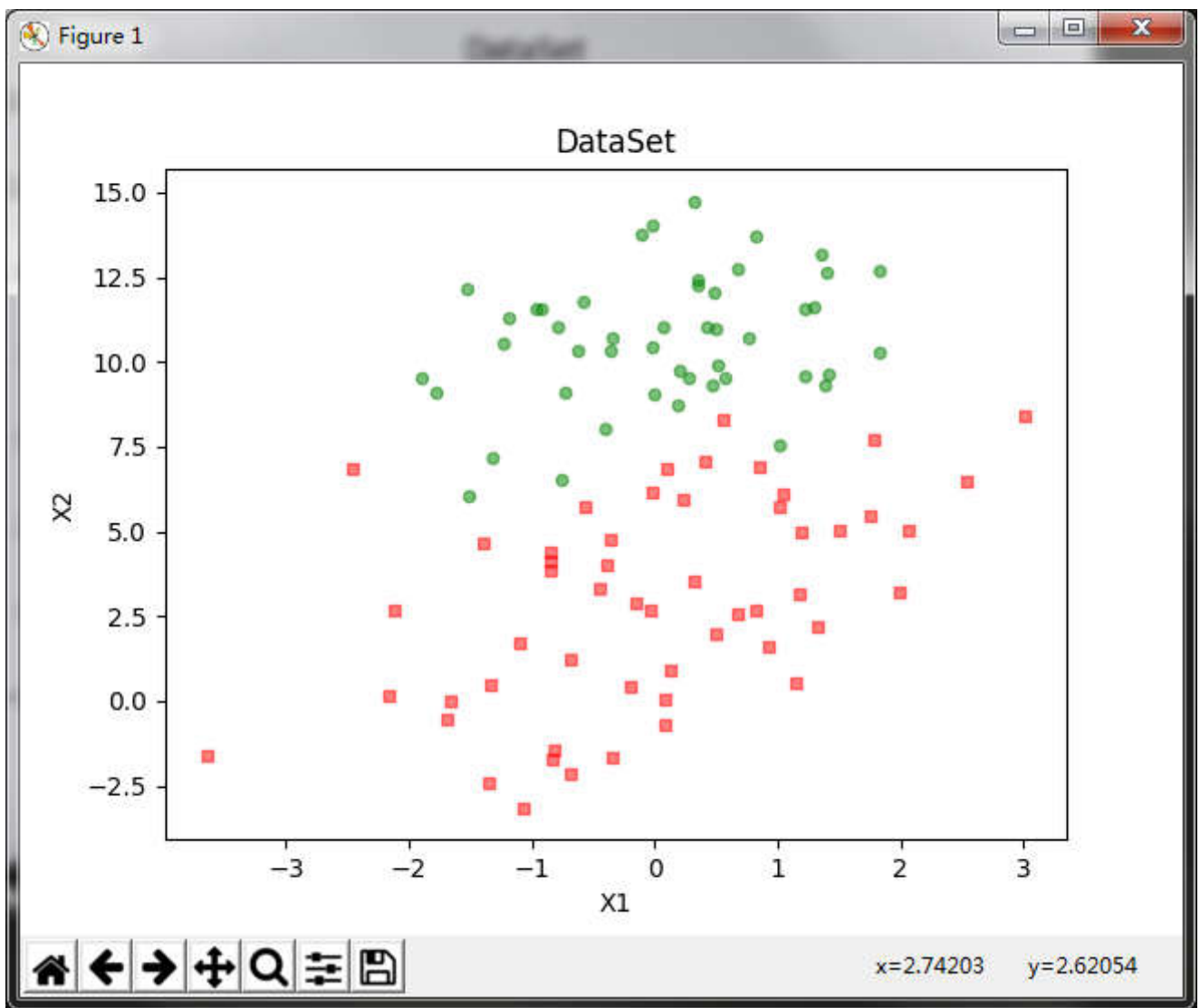
	Python
<pre> 1 # -*- coding:UTF-8 -*- 2 import matplotlib.pyplot as plt 3 import numpy as np 4 5 """ 6 函数说明:加载数据 7 8 Parameters: 9 无 10 Returns: 11 dataMat - 数据列表 12 labelMat - 标签列表 13 Author: 14 Jack Cui 15 Blog: 16 http://blog.csdn.net/c406495762 17 Zhihu: 18 https://www.zhihu.com/people/Jack--Cui/ 19 Modify: 20 2017-08-28 21 """ 22 def loadDataSet(): 23 dataMat = [] 24 labelMat = [] 25 fr = open('testSet.txt') 26 for line in fr.readlines(): 27 lineArr = line.strip().split() 28 dataMat.append([1.0, float(lineArr[0]), float(lineArr[1])]) </pre>	<pre> #创建数据列表 #创建标签列表 #打开文件 #逐行读取 #去回车，放入列表 #添加数据 </pre>

```

29         labelMat.append(int(LineArr[2]))          #添加标签
30     fr.close()                                     #关闭文件
31     return dataMat, labelMat                       #返回
32
33 """
34 函数说明:绘制数据集
35
36 Parameters:
37     无
38 Returns:
39     无
40 Author:
41     Jack Cui
42 Blog:
43     http://blog.csdn.net/c406495762
44 Zhihu:
45     https://www.zhihu.com/people/Jack--Cui/
46 Modify:
47     2017-08-28
48 """
49 def plotDataSet():
50     dataMat, labelMat = loadDataSet()              #加载数据集
51     dataArr = np.array(dataMat)                    #转换成numpy的array数
52     n = np.shape(dataMat)[0]                       #数据个数
53     xcord1 = []; ycord1 = []                       #正样本
54     xcord2 = []; ycord2 = []                       #负样本
55     for i in range(n):                             #根据数据集标签进行分类
56         if int(labelMat[i]) == 1:
57             xcord1.append(dataArr[i,1]); ycord1.append(dataArr[i,2]) #1为正样本
58         else:
59             xcord2.append(dataArr[i,1]); ycord2.append(dataArr[i,2]) #0为负样本
60     fig = plt.figure()
61     ax = fig.add_subplot(111)                       #添加subplot
62     ax.scatter(xcord1, ycord1, s = 20, c = 'red', marker = 's',alpha=.5) #绘制正样本
63     ax.scatter(xcord2, ycord2, s = 20, c = 'green',alpha=.5)          #绘制负样本
64     plt.title('DataSet')                            #绘制title
65     plt.xlabel('x'); plt.ylabel('y')                #绘制label
66     plt.show()                                       #显示

```

运行结果如下：



从上图可以看出数据的分布情况。假设Sigmoid函数的输入记为 z ，那么 $z = w_0x_0 + w_1x_1 + w_2x_2$ ，即可将数据分割开。其中， x_0 为全是1的向量， x_1 为数据集的第一列数据， x_2 为数据集的第二列数据。另 $z=0$ ，则 $0 = w_0 + w_1x_1 + w_2x_2$ 。横坐标为 x_1 ，纵坐标为 x_2 。这个方程未知的参数为 w_0 ， w_1 ， w_2 ，也就是我们需要的回归系数(最优参数)。

2、训练算法

在编写代码之前，让我们回顾下梯度上升迭代公式：

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

将上述公式矢量化：

$$\theta := \theta + \alpha X^T (\vec{y} - g(X\theta))$$

根据矢量化的公式，编写代码如下：

Python

```
1 # -*- coding:UTF-8 -*-
```

```

2  import numpy as np
3
4  """
5  函数说明:加载数据
6
7  Parameters:
8      无
9  Returns:
10     dataMat - 数据列表
11     labelMat - 标签列表
12 Author:
13     Jack Cui
14 Blog:
15     http://blog.csdn.net/c406495762
16 Zhihu:
17     https://www.zhihu.com/people/Jack--Cui/
18 Modify:
19     2017-08-28
20 """
21 def loadDataSet():
22     dataMat = []
23     labelMat = []
24     fr = open('testSet.txt')
25     for line in fr.readlines():
26         lineArr = line.strip().split()
27         dataMat.append([1.0, float(lineArr[0]), float(lineArr[1])])
28         labelMat.append(int(lineArr[2]))
29     fr.close()
30     return dataMat, labelMat
31
32 """
33 函数说明:sigmoid函数
34
35 Parameters:
36     inX - 数据
37 Returns:
38     sigmoid函数
39 Author:
40     Jack Cui
41 Blog:
42     http://blog.csdn.net/c406495762
43 Zhihu:
44     https://www.zhihu.com/people/Jack--Cui/
45 Modify:
46     2017-08-28
47 """
48 def sigmoid(inX):
49     return 1.0 / (1 + np.exp(-inX))
50
51
52 """
53 函数说明:梯度上升算法
54
55 Parameters:
56     dataMatIn - 数据集
57     classLabels - 数据标签
58 Returns:
59     weights.getA() - 求得的权重数组(最优参数)
60 Author:
61     Jack Cui
62 Blog:
63     http://blog.csdn.net/c406495762
64 Zhihu:
65     https://www.zhihu.com/people/Jack--Cui/
66 Modify:
67     -----

```

运行结果如图所示：

```

148 if __name__ == '__main__':
149     dataMat, labelMat = loadDataSet()
150     print(gradAscent(dataMat, labelMat))

[[ 4.12414349]
 [ 0.48007329]
 [-0.6168482 ]]
[Finished in 1.2s]

```

可以看出，我们已经求解出回归系数 $[w_0, w_1, w_2]$ 。

通过求解出的参数，我们就可以确定不同类别数据之间的分隔线，画出决策边界。

3、绘制决策边界

我们已经解出了一组回归系数，它确定了不同类别数据之间的分隔线。现在开始绘制这个分隔线，编写代码如下：

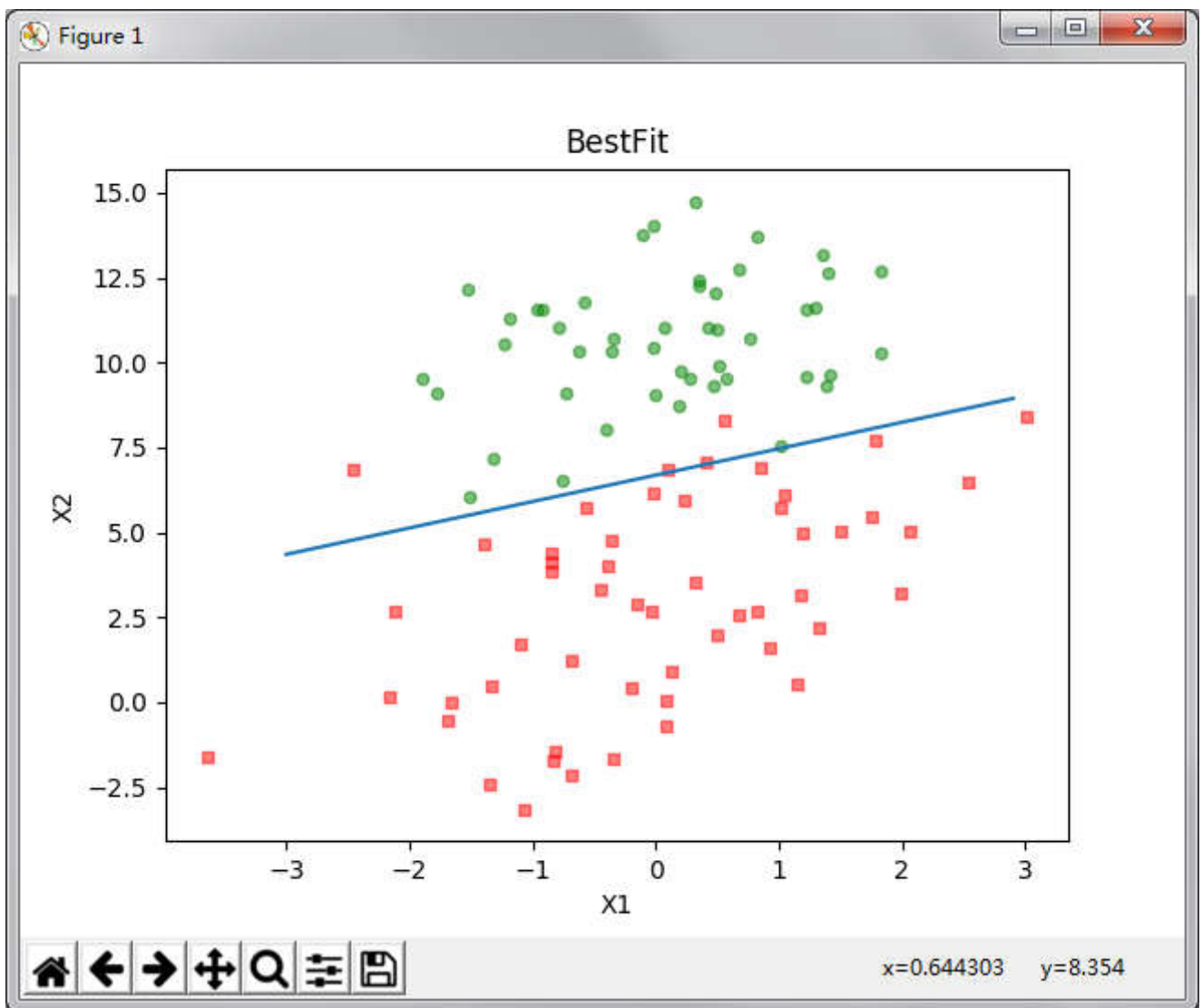
	Python
1	# -*- coding:UTF-8 -*-
2	import matplotlib.pyplot as plt
3	import numpy as np
4	
5	"""
6	函数说明:加载数据
7	
8	Parameters:
9	无
10	Returns:
11	dataMat - 数据列表
12	labelMat - 标签列表
13	Author:
14	Jack Cui
15	Blog:
16	http://blog.csdn.net/c406495762
17	Zhihu:
18	https://www.zhihu.com/people/Jack--Cui/
19	Modify:
20	2017-08-28
21	"""
22	def loadDataSet():
23	dataMat = []
24	labelMat = []
25	fr = open('testSet.txt')
26	for line in fr.readlines():
27	lineArr = line.strip().split()
28	dataMat.append([1.0, float(lineArr[0]), float(lineArr[1])])
29	labelMat.append(int(lineArr[2]))
30	fr.close()
31	return dataMat, labelMat
32	
33	
34	"""
35	函数说明:sigmoid函数
36	
37	Parameters:
38	inX - 数据
39	Returns:
40	sigmoid函数
41	Author:
42	Jack Cui
43	Blog:
44	http://blog.csdn.net/c406495762
45	Zhihu:
46	https://www.zhihu.com/people/Jack--Cui/

```

47 Modify:
48     2017-08-28
49 """
50 def sigmoid(inX):
51     return 1.0 / (1 + np.exp(-inX))
52
53 """
54 函数说明:梯度上升算法
55
56 Parameters:
57     dataMatIn - 数据集
58     classLabels - 数据标签
59 Returns:
60     weights.getA() - 求得的权重数组(最优参数)
61 Author:
62     Jack Cui
63 Blog:
64     http://blog.csdn.net/c406495762
65 Zhihu:
66     https://www.zhihu.com/people/Jack--Cui/

```

运行结果如下：



这个分类结果相当不错，从上图可以看出，只分错了几个点而已。但是，尽管例子简单数据集很小，但是这个方法却需要大量的计算(300次乘法)。因此下篇文章将对改算法稍作改进，从而减少计算量，使其可以应用于大数据集上。

四、总结

Logistic回归的一般过程：

- 收集数据：采用任意方法收集数据。
- 准备数据：由于需要进行距离计算，因此要求数据类型为数值型。另外，结构化数据格式则最佳。
- 分析数据：采用任意方法对数据进行分析。
- 训练算法：大部分时间将用于训练，训练的目的是为了找到最佳的分类回归系数。
- 测试算法：一旦训练步骤完成，分类将会很快。
- 使用算法：首先，我们需要输入一些数据，并将其转换成对应的结构化数值；接着，基于训练好的回归系数，就可以对这些数值进行简单的回归计算，判定它们属于哪个类别；在这之后，我们就可以在输出的类别上做一些其他分析工作。

其他：

- Logistic回归的目的是寻找一个非线性函数Sigmoid的最佳拟合参数，求解过程可以由最优化算法完成。
- 本文讲述了Logistic回归原理以及数学推导过程。
- 下篇文章将讲解Logistic回归的改进以及Sklearn实战内容。
- 如有问题，请留言。如有错误，还望指正，谢谢！

PS：如果觉得本篇本章对您有所帮助，欢迎关注、评论、赞！

本文出现的所有代码和数据集，均可在我的github上下载，欢迎Follow、Star：

<https://github.com/Jack-Cherish/Machine-Learning>

参考文献：

- 斯坦福大学的吴恩达《机器学习》：<https://www.coursera.org/learn/machine-learning>
- 《机器学习实战》第五章内容



微信公众号

分享技术，乐享生活：微信公众号搜索「JackCui-AI」关注一个在互联网摸爬滚打的潜行者。

不要为自己的努力道歉，这样太对不起自己的努力了。--- 凯《火影忍者》

