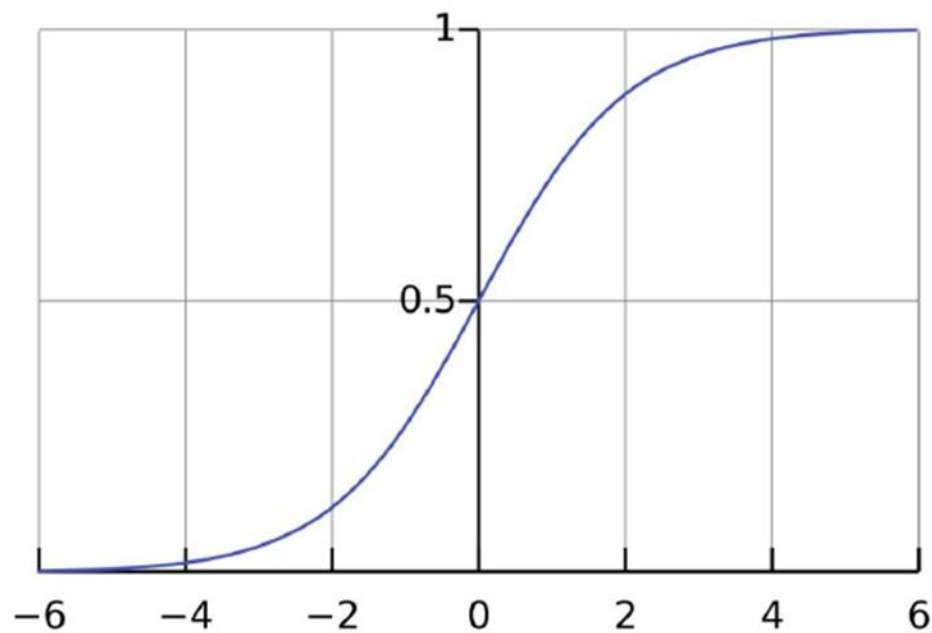
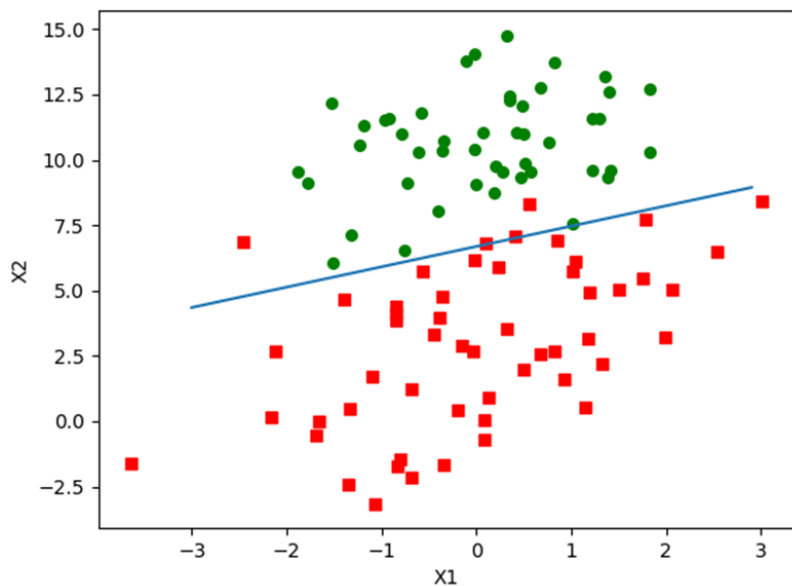


梯度下降

复习



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

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

思考

$$-\frac{1}{m} \sum_{i=1}^m \left(\left(y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)} \right)$$

- 这个式子是关于谁的偏导数？
- 这个式子什么时候等于0？
- 这个式子再求偏导怎么求？



再抢救一下

- 逻辑回归的代价函数

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log h_{\boldsymbol{\theta}}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(x^{(i)})) \right)$$

- 逻辑回归的代价函数需要求最小值，但是以求导函数何时为0的方式求最小值代价太大。

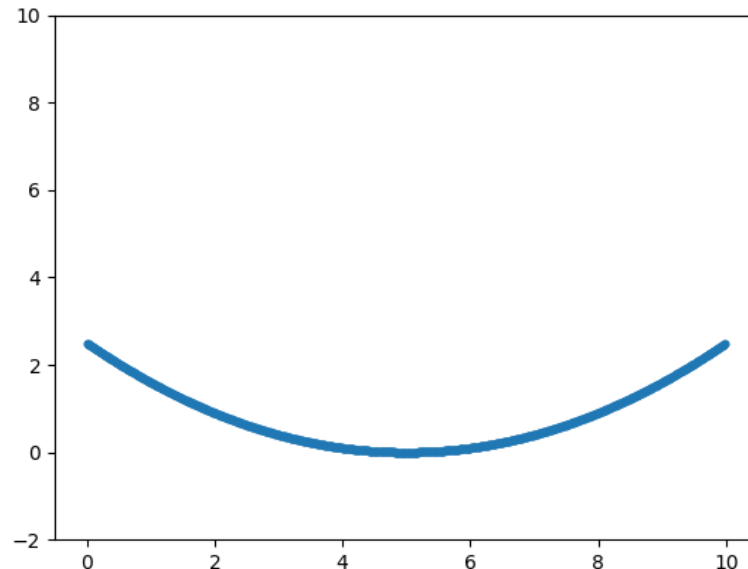
梯度下降

例

- 不使用导数为零处取得极值的结论计算 $f(x) = \frac{1}{10}(x-5)^2$ 什么时候取得最小值。

例

- 首先给定随机数 $x_0=7$ ，由于 $f'(7)=0.4>0$ ，说明在 $x=7$ 处函数处于递增状态，那么要找 x 为多少时可以取得最小值就应该找一个小于7值赋给 x ，即令 x_0 减去一个数作为 x_1 ，并考虑 x_1 处是否为最小值点。



例

- 可以令 $x_0 - f'(x_0) = x_1$, 那么 $x_1 = 6.6$, 由于 $f'(6.6) = 0.32$ 依然大于0, 说明应该继续减小。

第0次迭代

6.6

第1次迭代

6.2799999999999999

第2次迭代

6.0239999999999999

第3次迭代

5.8191999999999995

第4次迭代

5.65526

第20次迭代

5.018446744073709

第21次迭代

5.014757395258967

第22次迭代

5.011805916207174

第23次迭代

5.009444732965739

第24次迭代

5.007555786372591

第25次迭代

思考

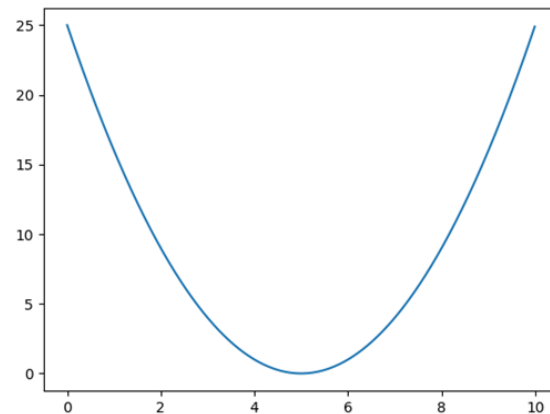
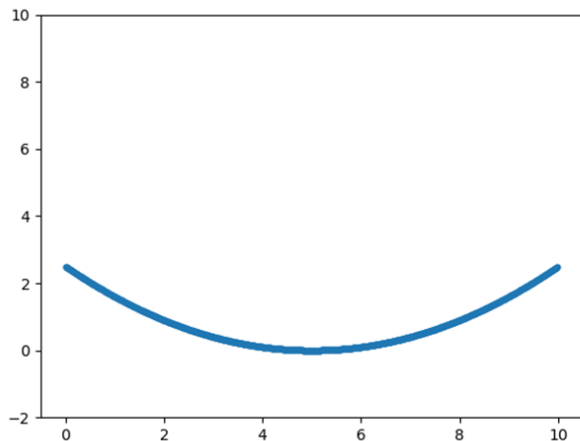
- 上面的例子中，一开始选择的 x_0 的导数大于0，所以我们可以令 x_0 减去它的导数。那么如果一开始选择的 x_0 的导数为负呢？

梯度下降

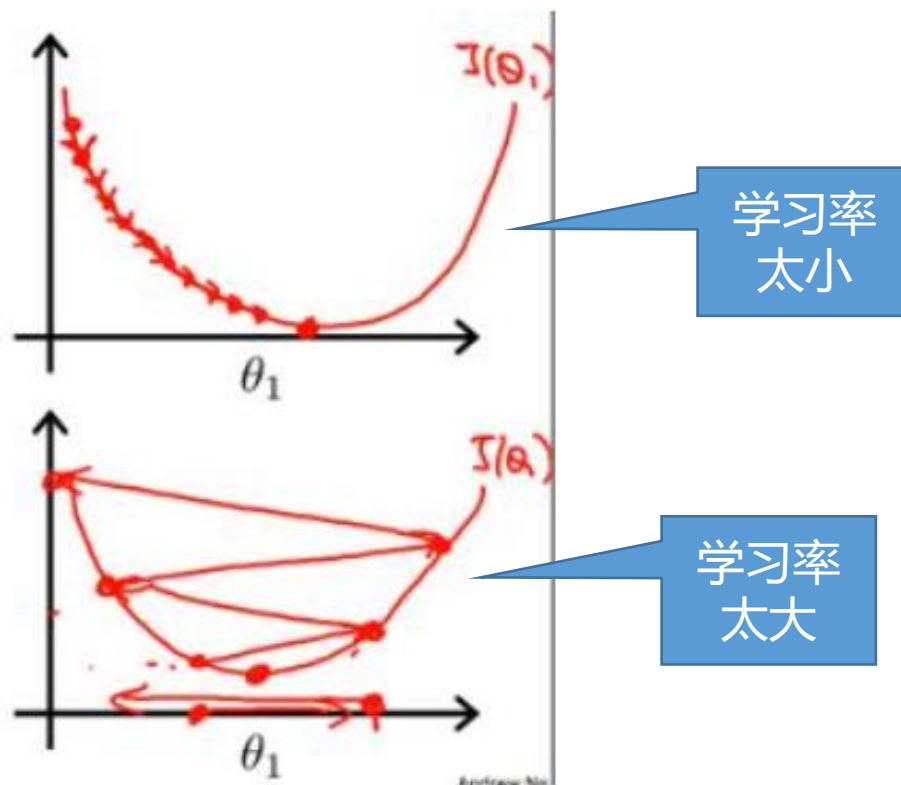
- 首先随机选择 x_0 ，计算函数在 x_0 处的导数，然后用 x_0 减去它自己的导数，并将得到的结果作为 x_1 ，并计算 x_1 处的导数，然后得到 $x_2 \cdots$ ，直到收敛。

学习率

- 如果去掉上面的例子中的 $1/10$ ，初始值还选7，后果会怎么样？
- $f'(7)=4$, 用 $7-f'(4)=3$ 取代7, $f'(3)=-4$, $3-f'(3) = 7$
- 为了控制 x 的变化速度，需要设定一个常数 α ，用 $x_0-\alpha f(x_0)$ 取代 x_0 。 α 称为学习率。



学习率的选择



学习率的选择

- 学习率的选择没有太好的方法，一般都是根据实际情况首先尝试一个常用的值，比如0.1，0.01，0.001等，根据代价的变化情况决定继续使用当前学习率还是对学习率进行调整。

例

```
learning_rate = 0.01
a = 3
for i in range(1000):
    a = a - learning_rate*(2 * a - 10)
    print(a)
```

4.9999999996423753

4.9999999996495278

4.9999999996565372

4.9999999996634065

问题

- 用梯度下降一定能找到最小值吗？
- 未必，找到的可能是局部最小值，即极小值。

例

- 利用上面的方法计算 $\sqrt{73}$ 。

例

- 根据题意 $x = \sqrt{73} \Rightarrow x^2 = 73 \Rightarrow x^2 - 73 = 0$
- 假设有 $f(x) = x^2 - 73, F'(x) = f(x)$, 此题目转换为求F(x)的极值(根据二阶导数可以知道是极小值)问题。
- 所以, 可以利用上一题的方法

```
learning_rate = 0.01
```

```
a = 3
```

```
for i in range(1000):
```

```
    a = a - learning_rate*(a*a - 73)
```

```
    print(a)
```

```
8.544003745317527
```

```
8.544003745317527
```

梯度下降

- 如果一个函数 $f(x)$ 至少有一个极小值，那么就可以利用下面的方法找到这个极小值：
 - 首先随机确定一个点 x_0 以及一个称为学习率的正数；
 - 计算 $f(x)$ 在 x_0 处的导数，并用 x_0 减去 $f(x)$ 在 x_0 处的倒数与学习率的乘积作为 x_1 ；
 - 重复第二步，直到 x 不再变化或者达到指定的迭代步数。
- 如果求极大值，则加上导数与学习率的乘积，称为梯度上升。

多元函数的梯度下降

- 用梯度下降的方法求函数 $z = (x-5)^2 + (y-3)^2$ 的极小值。

多元函数的梯度下降

首先，随机给定x和y的初始值，比如 $x=10$, $y=10$

计算偏导数, $\frac{\partial z}{\partial x} = 2x - 10 = 10$, $\frac{\partial z}{\partial y} = 2y - 6 = 14$

指定学习率 $\alpha=0.1$

$$x = 10 - 0.1 * 10 = 9, y = 10 - 0.1 * 14 = 8.6$$

$$x = 9 - 0.1 * 8 = 8.2, y = 8.6 - 0.1 * 11.2 = 7.48$$

...

多元函数的梯度下降

```
x = 10
y = 10
learning_rate = 0.1
for i in range(100):
    x = x - learning_rate * (2 * x - 10)
    y = y - learning_rate * (2 * y - 6)
    print("第" + str(i) + "次迭代")
    print("x = " + str(x))
    print("y = " + str(y))
```

多元函数的梯度下降

第0次迭代

$x = 9.0$

$y = 8.6$

第1次迭代

$x = 8.2$

$y = 7.4799999999999995$

第2次迭代

$x = 7.56$

$y = 6.584$

第3次迭代

$x = 7.048$

$y = 5.8671999999999995$

第4次迭代

$y = 3.013539969179684$

第28次迭代

$x = 5.007737125245534$

$y = 3.010831975343747$

第29次迭代

$x = 5.006189700196427$

$y = 3.0086655802749975$

第30次迭代

$x = 5.004951760157142$

$y = 3.006932464219998$

第31次迭代

$x = 5.003961408125713$

$y = 3.0055459713759984$

多元函数的梯度下降

16. 将长为2m的铁丝分成三段，依次围成圆、正方形与正三角形，三个图形的面积之和是否存在最小值？若存在，求出最小值。

【答案】设分成的三段分别为 x, y, z ，则有 $x + y + z = 2$ 及 $x, y, z > 0$ ，圆的面积为 $S_1 = \frac{1}{4\pi}x^2$ ，正方形的面积为 $S_2 = \frac{1}{16}y^2$ ，正三角形的面积为 $S_3 = \frac{\sqrt{3}}{36}z^2$ ，总面积 $S = \frac{1}{4\pi}x^2 + \frac{1}{16}y^2 + \frac{\sqrt{3}}{36}z^2$ ，则问题转化为在条件 $x + y + z = 2, x, y, z > 0$ 下，求函数 $\frac{1}{4\pi}x^2 + \frac{1}{16}y^2 + \frac{\sqrt{3}}{36}z^2$ 的最小值。令 $L = \frac{1}{4\pi}x^2 + \frac{1}{16}y^2 + \frac{\sqrt{3}}{36}z^2 + \lambda(x + y + z - 2)$ ，

$$\text{解得唯一条件极值点为} \begin{cases} x = \frac{2\sqrt{3}\pi}{\sqrt{3}\pi + 4\sqrt{3} + 9} \\ y = \frac{8\sqrt{3}}{\sqrt{3}\pi + 4\sqrt{3} + 9} \\ z = \frac{18}{\sqrt{3}\pi + 4\sqrt{3} + 9} \end{cases}$$

多元函数的梯度下降

- 由于 $x+y+z=2$ ，所以 $z=2-x-y$ ，所以总面积公式可以写为：

$$S = \frac{1}{4\pi} x^2 + \frac{1}{16} y^2 + \frac{\sqrt{3}}{36} (2-x-y)^2$$

- 求偏导：

$$\frac{\partial S}{\partial x} = \frac{1}{2\pi} x - \frac{\sqrt{3}}{18} (2-x-y)$$

$$\frac{\partial S}{\partial y} = \frac{1}{8} y - \frac{\sqrt{3}}{18} (2-x-y)$$

多元函数的梯度下降

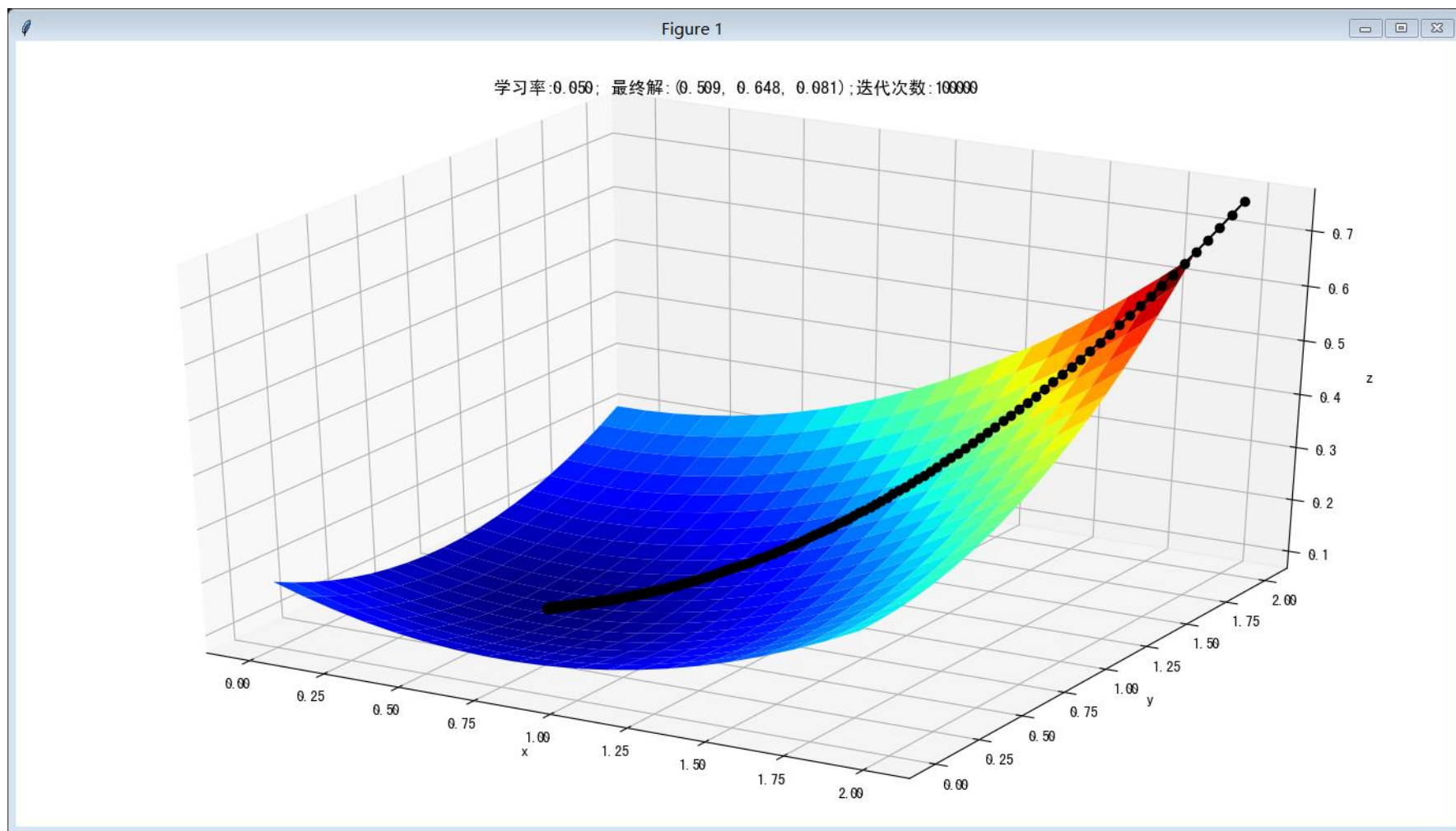
- 面积公式实现:

```
def f(x, y):  
    return (x * x) / (4 * math.pi) + (y * y) / 16 + \  
    (math.sqrt(3) * (2 - x - y) * (2 - x - y)) / 36
```

- 偏导数

```
def hx(x, y):  
    return x / (2 * math.pi) - (math.sqrt(3) * (2 - x - y)) / 18  
  
def hy(x, y):  
    return y / 8 - (math.sqrt(3) * (2 - x - y)) / 18
```

多元函数的梯度下降



多元函数梯度下降总结

- 首先生成一个 n 维随机向量 \mathbf{a}_0 ，其中 n 代表多元函数中有 n 个自变量。
- 分别求多元函数对这 n 个自变量在 \mathbf{a}_0 处的偏导数
这 n 个偏导数构成一个新的 n 维向量 \mathbf{b}_0 。
- 设定学习率 α 。
- 令 $\mathbf{a}_1 = \mathbf{a}_0 - \alpha \times \mathbf{b}_0$ ，并求多元函数在 \mathbf{a}_1 处的
 n 个偏导数，构成 \mathbf{b}_1 。
- 重复以上过程，直到收敛或者达到指定次数。

逻辑回归中的梯度下降

- 逻辑回归的代价函数为：

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right)$$

代价函数求偏导

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right)$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \cdot \frac{1}{h_{\theta}(\mathbf{x}^{(i)})} \cdot \frac{\partial}{\partial \theta_j} h_{\theta}(\mathbf{x}^{(i)}) - (1 - y^{(i)}) \cdot \frac{1}{1 - h_{\theta}(\mathbf{x}^{(i)})} \cdot \frac{\partial}{\partial \theta_j} h_{\theta}(\mathbf{x}^{(i)}) \right) \\ &= -\frac{1}{m} \sum_{i=1}^m \left(\left(y^{(i)} \cdot \frac{1}{h_{\theta}(\mathbf{x}^{(i)})} - (1 - y^{(i)}) \cdot \frac{1}{1 - h_{\theta}(\mathbf{x}^{(i)})} \right) \frac{\partial}{\partial \theta_j} h_{\theta}(\mathbf{x}^{(i)}) \right) \\ &= -\frac{1}{m} \sum_{i=1}^m \left(\left(y^{(i)} \cdot \frac{1}{g(\boldsymbol{\theta}^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \cdot \frac{1}{1 - g(\boldsymbol{\theta}^T \mathbf{x}^{(i)})} \right) \frac{\partial}{\partial \theta_j} g(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) \right) \end{aligned}$$

代价函数求偏导

$$\begin{aligned}\frac{\partial}{\partial \theta_j} g(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}} \right) \\&= -\frac{1}{\left(1 + e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}\right)^2} \cdot \frac{\partial}{\partial \theta_j} \left(1 + e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}\right) \\&= -\frac{e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}}{\left(1 + e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}\right)^2} \cdot \frac{\partial}{\partial \theta_j} \left(-\boldsymbol{\theta}^T \mathbf{x}^{(i)}\right) \\&= \frac{e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}}{\left(1 + e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}\right)^2} \cdot \mathbf{x}^{(i)}\end{aligned}$$

代价函数求偏导

$$\begin{aligned}g\left(\boldsymbol{\theta}^T \mathbf{x}^{(i)}\right) &= \frac{1}{1+e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}} \Rightarrow e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}} = \frac{1}{g\left(\boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)} - 1 \\&\Rightarrow -\frac{e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}}{\left(1+e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}\right)^2} = g^2\left(\boldsymbol{\theta}^T \mathbf{x}^{(i)}\right) \left(\frac{1}{g\left(\boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)} - 1\right) \\&= g\left(\boldsymbol{\theta}^T \mathbf{x}^{(i)}\right) \left(1 - g\left(\boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)\right)\end{aligned}$$

代价函数求偏导

$$\begin{aligned}& -\frac{1}{m} \sum_{i=1}^m \left(\left(y^{(i)} \cdot \frac{1}{g(\boldsymbol{\theta}^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \cdot \frac{1}{1 - g(\boldsymbol{\theta}^T \mathbf{x}^{(i)})} \right) g(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) (1 - g(\boldsymbol{\theta}^T \mathbf{x}^{(i)})) \mathbf{x}^{(i)} \right) \\&= -\frac{1}{m} \sum_{i=1}^m \left(\left(y^{(i)} (1 - g(\boldsymbol{\theta}^T \mathbf{x}^{(i)})) - (1 - y^{(i)}) g(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)} \right) \\&= -\frac{1}{m} \sum_{i=1}^m \left(\left(y^{(i)} - g(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)} \right) \\&= -\frac{1}{m} \sum_{i=1}^m \left(\left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)} \right) \\&= \frac{1}{m} \sum_{i=1}^m \left(\left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) \mathbf{x}^{(i)} \right)\end{aligned}$$

参数的迭代

- 根据梯度下降原理，要求代价函数的最小值，迭代过程如下

$$\begin{array}{l} \text{Repeat until convergence } \{ \\ \quad \theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \\ \} \end{array}$$

此处可以没有 $1/m$ ，
参见代码

逻辑回归实现

- 逻辑回归的实现主要分为三个部分：读取和处理数据，梯度下降，可视化。

```
def sigmoid(z):  
    return 1 / (1 + np.exp(-z))
```

```
def init_data():  
    data = np.loadtxt('data.csv')  
    dataMatIn = data[:, 0:-1]  
    classLabels = data[:, -1]  
    dataMatIn = np.insert(dataMatIn, 0, 1, axis=1)  
    return dataMatIn, classLabels
```

逻辑回归实现

```
def grad_descent(dataMatIn, classLabels):  
    dataMatrix = np.mat(dataMatIn)  
    labelMat = np.mat(classLabels).transpose()  
    m, n = np.shape(dataMatrix)  
    weights = np.ones((n, 1))  
    alpha = 0.001  
    maxCycle = 500  
  
    for i in range(maxCycle):  
        h = sigmoid(dataMatrix * weights)  
        weights = weights - alpha * dataMatrix.transpose() * (h - labelMat)  
    return weights
```

作业

- 保留鸢尾花数据集的前两个类别，用逻辑回归进行分类，只能调用numpy、pandas、matplotlib等基础库。
- 加注释。
- 下周二上课前通过qq发给我，文件以学号加姓名命名。