

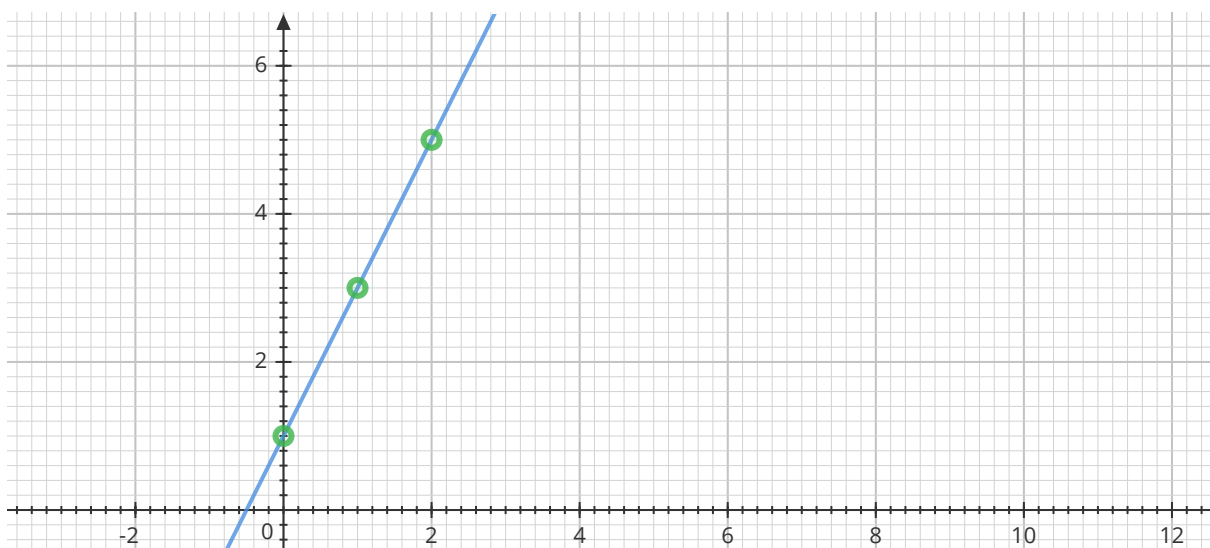
基本概念

假设我们已经测量到了一组数据：

序号	自变量 x	因变量 y
1	$x_1 = 0$	$y_1 = 1$
1	$x_2 = 1$	$y_2 = 3$
3	$x_3 = 2$	$y_3 = 5$
4	$x_4 = 3$	$y_4 = ?$

Table 1: 测量数据 (1)

如何根据前三次的记录预测以及自变量 x_4 来推算 y_4 的值呢？我们可以根据测量数据做出散点图，如下：



继而假设一条直线 $y = kx + b$ 表示点的变化趋势，要使直线尽可能地准确，就要求所有点到直线的距离之和——一个关于 k, b 的二元函数 $f(k, b)$ 最小，于是就能唯一地确定一组 (k, b) 用以表示最终的直线。这就是最简单的线性规划。

推广一：高维形式

对于二维平面中的点我们可以用一维的线去拟合；对于三维空间中的点我们可以用二维的平面去拟合；对于 n 维空间中的点，我们则可以用 $n - 1$ 维的超平面去拟合。形式如下：

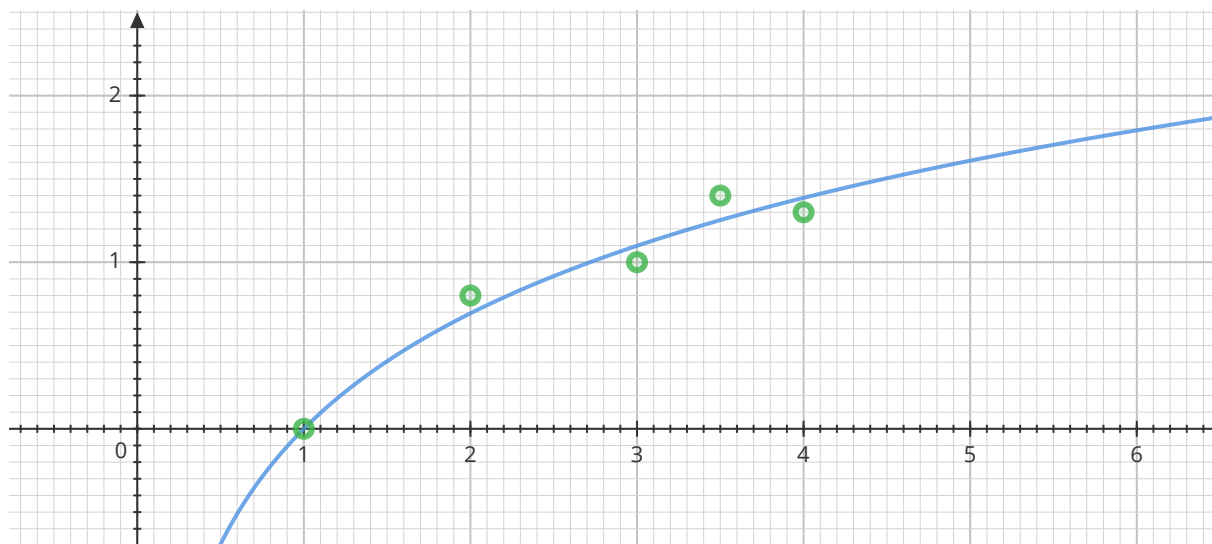
$$y = b + k_1x_1 + k_2x_2 + \dots + k_nx_n \quad (1)$$

b 可以看作 k_0x_0 ， $x_0 = 1$ ，所以上式也可写作：

$$y = k_0x_0 + k_1x_1 + k_2x_2 + \dots + k_nx_n \quad (2)$$

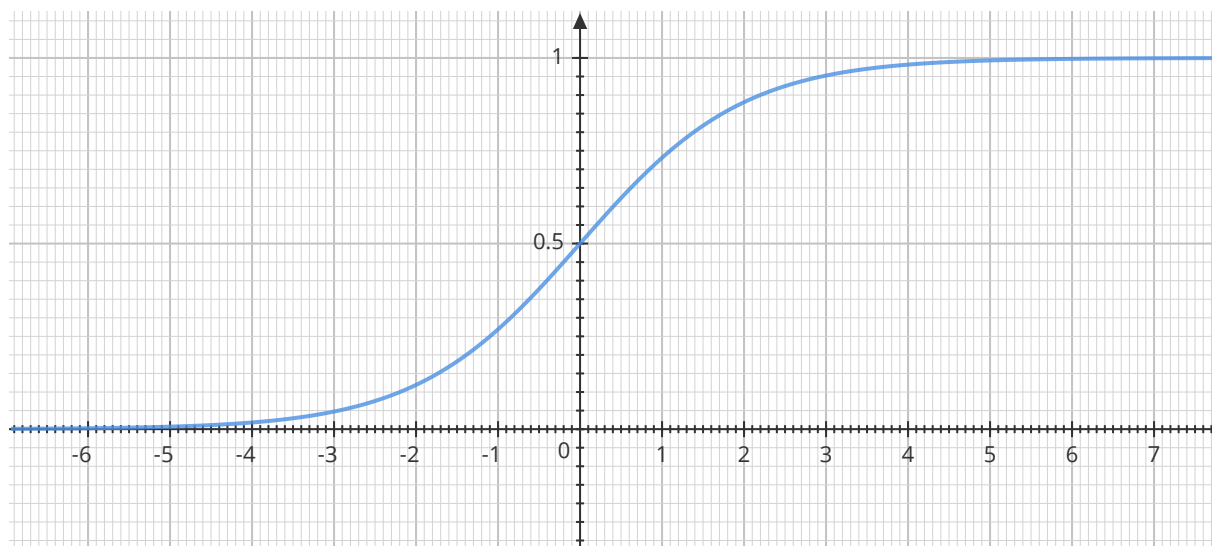
推广二：非线性数据

对于非线性数据，例如：



很明显最终的曲线应该是一个对数 $\ln(x)$ 。如何去用线性的函数表示非线性的图像呢？

我们假设（构造）另一个函数 $F = f(y) = e^y$ ，那么 $F - x$ 就是线性变化的了。这里我们注意到 $e^0 \Rightarrow \ln(1)$ 互为反函数。所以对于非线性的数据，我们可以根据数据变化的规律，来构造一个反函数，使其线性化。而对于实际情况来说，将线性的计算结果非线性化则更合理一些：逻辑回归中的sigmoid() 函数就是其中一种，它可以将线性的结果非线性化成类似于逻辑值。



到这里，我们得到了第一个关键步骤：**【一】寻找反函数，将线性结果非线性化。**

对于已经线性化的问题来说，如果想得到一个尽可能准确的直线，那我们可以采取简单线性回归中的待定系数法确定 $[k_0, k_1, \dots, k_n]$ 。也可以先假设一组系数 $[k_0, k_1, \dots, k_n]$ ，根据这组系数得到一个不那么准确的结果 \hat{y} ，比较误差 $\Delta(K) = L(y, \hat{y})$ ：

如果大于0，则将 $[k_0, k_1, \dots, k_n]$ ，适当调小一点；

如果小于0，则将 $[k_0, k_1, \dots, k_n]$ ，适当调大一点。

反复迭代之后，最终将误差控制在一个可以接受的范围之内，就算是找到了合适的拟合曲（直）线。

对于求误差函数最小值，我们可以采用梯度下降法。

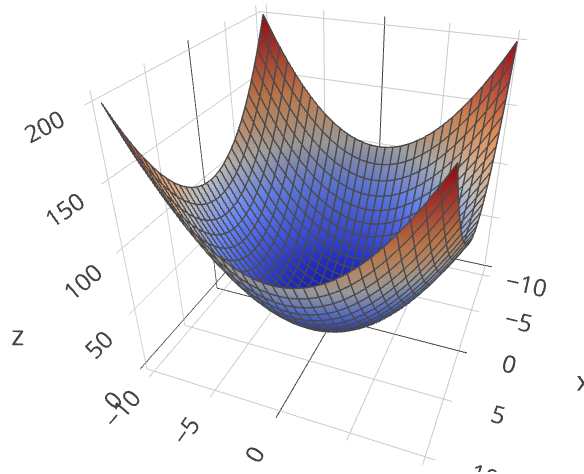
梯度下降

梯度是一个向量，而沿着梯度的方向就能最快达到误差函数的最小值：

$$\Delta(K) = f(k_1, k_2, k_3, \dots, k_n) \quad (3)$$

$$\nabla \Delta(K) = \left\langle \frac{\partial F}{\partial k_1}, \frac{\partial F}{\partial k_2}, \frac{\partial F}{\partial k_3}, \dots, \frac{\partial F}{\partial k_n} \right\rangle \quad (4)$$

而梯度下降一次就是： $K^{[i+1]} = K^{[i]} - \alpha \nabla \Delta(K^{[i]})$ ，其中 α 称作学习率或者步长，其值的选择需要慎重考虑。而梯度下降能够生效的前提是误差函数是凸函数，简单来说应该是应该是个碗形，其实是对于任一自变量 k_i ，误差函数的导数都应该是单调递增的。否则会出现多个极点，造成无法收敛。



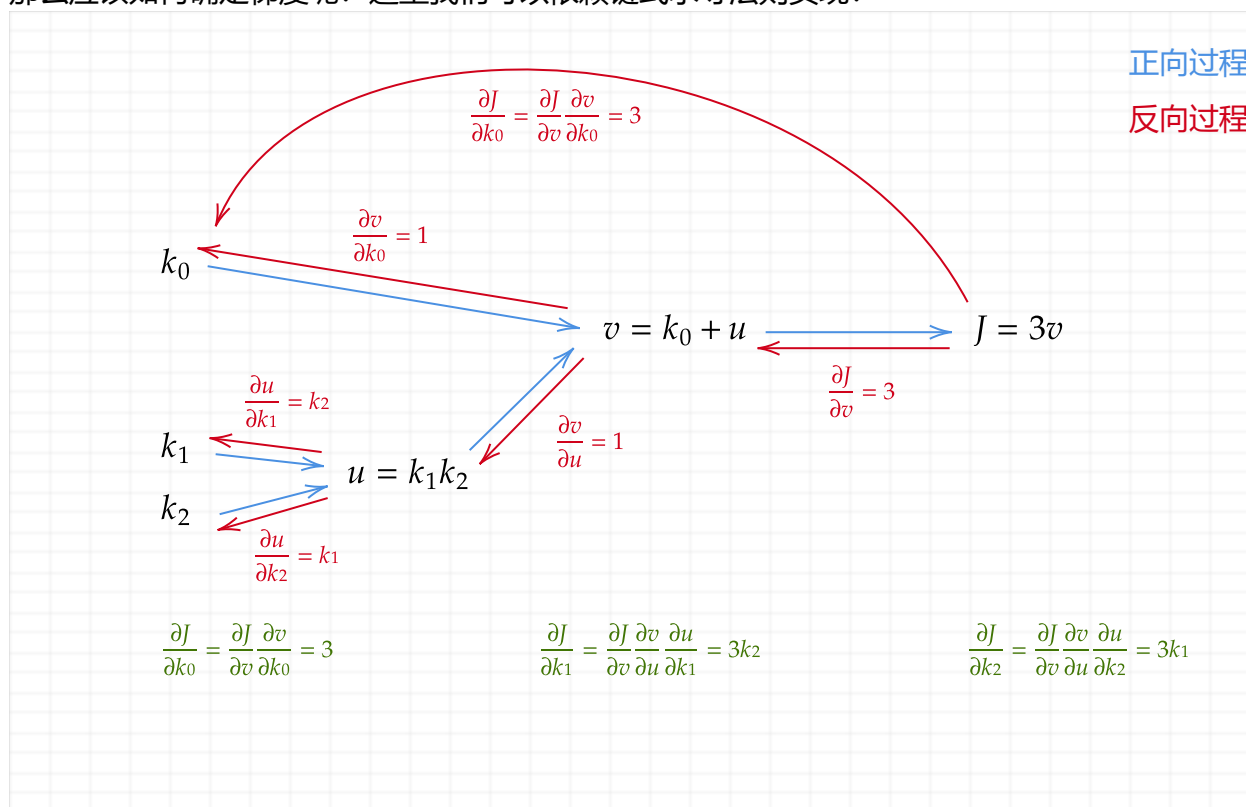
而对于有 m 组测量数据时，一般取其平均值： $J(K) = \frac{1}{m} \sum_{i=0}^m \Delta(K_i)$ ，也叫成本函数。

所以另外两个关键步骤：

【二】寻找合适的凸函数作为误差函数（成本函数）

【三】确定合适的学习率

那么应该如何确定梯度呢？这里我们可以依赖链式求导法则实现：



于是第四个关键步骤：**【四】根据链式求导法则，计算每个变量对应成本函数的偏微分**

接下来的事情就是让机器反复迭代了，以确定最合适的系数组合 K 。

逻辑回归的实现

例子来自于【[哔哩哔哩-大话神经网络，10行代码不调包，听不懂你打我！](#)】，讲的是去不去陪妹子看电影的故事。

假设我们有这么一组数据， A, B, C, D 四个人以前一起看电影的记录，来推测当 A 下次去看电影的时候， D 会不会一起去。

A	B	C	D
0	0	1	0
1	1	1	1
1	0	1	1
0	1	1	0
1	0	0	?

获取数据

一般传感器采集到的数据会被以文本的形式保存在SD 卡中，这里以.csv 文件为例：

```

1  '''data.csv
2  A,B,C,D # title
3  0,0,1,0
4  1,1,1,1
5  1,0,1,1
6  0,1,1,0
7  '''
8
9  import numpy as np
10
11  def loadData():
12      org_data = np.loadtxt(
13          './data.csv', # 文件名
14          delimiter=',', # 设置分隔符
15          usecols=(0, 1, 2, 3), # 需要用到的数据范围
16          skiprows=1 # 跳过标题行
17      )
18      exp_data = orgData[:, 0:3] # A,B,C 样本数据列
19      res_data = orgData[:, 3:4] # D 结果数据列

```

```

20 exp_data = np.insert(exp_data, 3, 1, axis=1)
21 # 在样本数据后附加一列常数项, 值为1
22 return exp_data, res_data
23
24 '''
25 exp_data=[
26     [0. 0. 1.]
27     [1. 1. 1.]
28     [1. 0. 1.]
29     [0. 1. 1.]
30 ]
31 res_data=[
32     [0.]
33     [1.]
34     [1.]
35     [0.]
36 ]
37 # 附加常数项后
38 exp_data=[
39     [0. 0. 1. 1.]
40     [1. 1. 1. 1.]
41     [1. 0. 1. 1.]
42     [0. 1. 1. 1.]
43 ]
44 '''

```

初始化权重

给权重 $w_0, w_1, w_2, \dots, w_n$ 设定一组合适的初始值, 可以有效地减少迭代的次数, 更快地找到答案, 这里我们设置一组随机数。

```

1 def initWeight(shape): # 给定输入数据的维度n
2     rows, cols = shape
3     np.random.seed(1)
4     weight = 2 * np.random.random((cols, 1))-1 # 生成一个nx1 的随机数矩
        阵, 元素数值介于[-1,1]
5     return weight

```

对于第 i 次测量, 有:

$$z_i = x_{0i}w_0 + x_{1i}w_1 + \dots + x_{ni}w_n = \text{np.dot}(\text{exp_data}[i], \text{weight})$$

关键点一：正则化函数

上文提到，逻辑回归中我们采用 *Sigmoid()* 函数

$$a = \frac{1}{1 + e^{-z}} \quad (5)$$

来将线性结果 z_i 非线性化，使其尽量贴近测量数据 res_data 。

```
1 def Sigmoid(z):
2     return 1/(1+np.exp(-z))
3 # numpy 中的“广播”特性可以很容易地将数值运算扩展到矩阵运算中
```

而 *Sigmoid()* 函数的导函数为

$$a' = a(1 - a) \quad (6)$$

关键点二：误差函数

依然如上文所述，我们需要构造一个 w_i 的凸函数来作为误差函数，这里我们选择：

$$L = loss(a, y) = -(y \log(a) + (1 - y) \log(1 - a)) \quad (7)$$

注意，在 *numpy* 中 *log()* 默认是以 e 为底的，也就是说：

$$L = loss(a, y) = -(y * \ln(a) + (1 - y) \ln(1 - a))$$

更改对数的底并不影响其凸函数的性质，其导函数为：

$$L' = -\left(\frac{y}{a} + \frac{-(1-y)}{(1-a)}\right) = -\frac{(1-a)y - a(1-y)}{a(1-a)} = \frac{a-y}{a(1-a)} \quad (8)$$

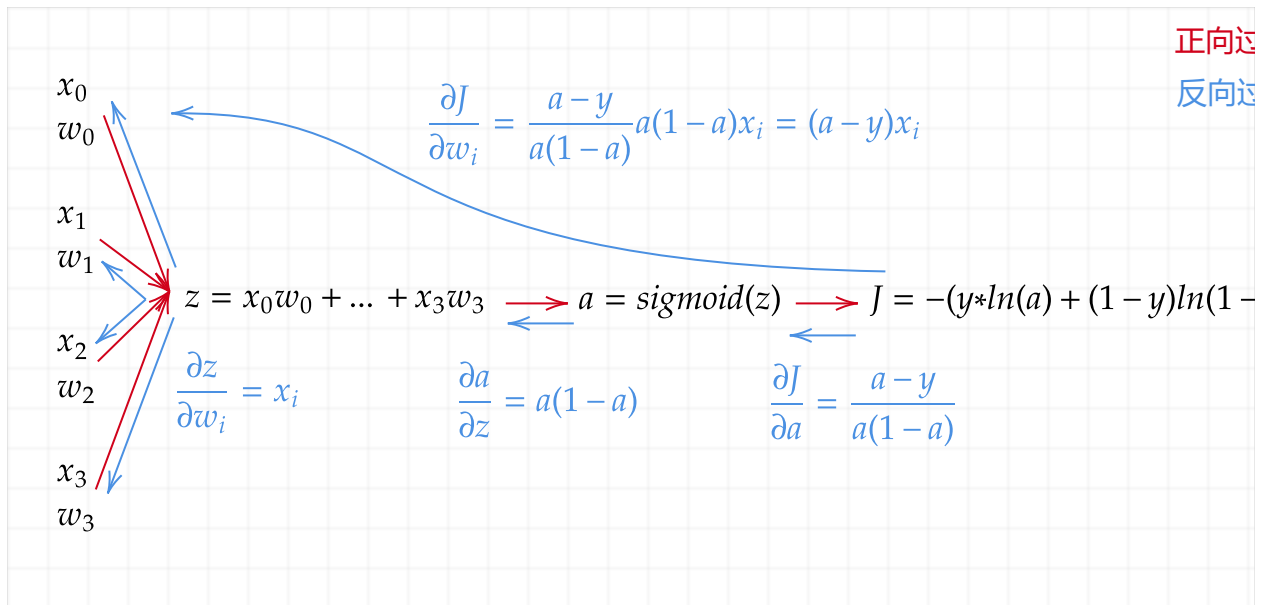
```
1 def Loss():
2     return -(y*np.log(a))-(1-y)*np.log(1-a)
3
4 def Cost(loss): # cost 去误差函数的平均值，函数在实际计算中的用处不大
5     return np.mean(loss, axis=0)
```

关键点三：确定学习率 α

α 一般取 0.1, 0.3, 0.01, 0.03, ...

关键点四：计算偏微分

对于每一次的采样数据，都有：



而对于全部的采样数据有：

$$\frac{\partial J}{\partial w_i} = \frac{1}{m} \sum_{j=0}^m (a_j - y_j) x_i \quad (9)$$

对应的代码实现为：

```

1  def Gradient(exp_data, res_data, a):
2      return np.dot(exp_data.T, a - res_data)/len(exp_data)
3      ...
4      # 这里需要注意一下dot 函数的作用
5      exp_data=[
6          [0. 0. 1. 1.]
7          [1. 1. 1. 1.]
8          [1. 0. 1. 1.]
9          [0. 1. 1. 1.]
10     ]
11
12     exp_data.T=[
13         [0. 1. 1. 0.]
14         [0. 1. 0. 1.]
15         [1. 1. 1. 1.]
16         [1. 1. 1. 1.]
17     ]
18

```



```

19 a=[
20     [0.19859385]
21     [0.24593466]
22     [0.1734943 ]
23     [0.27798934]
24 ]
25
26 y=[
27     [0.]
28     [1.]
29     [1.]
30     [0.]
31 ]
32
33 a-y=[
34     [ 0.19859385]
35     [-0.75406534]
36     [-0.8265057 ]
37     [ 0.27798934]
38 ]
39
40 np.dot(exp_data.T, a - res_data) = [ # 刚好对应着误差函数与x_i 的乘积之
    和
41     [-1.58057104]
42     [-0.47607599]
43     [-1.10398785]
44     [-1.10398785]
45 ]
46 '''

```

综上可得到具体的流程：

```

1 def main():
2     alpha = 0.3
3     exp_data, res_data = loadData()
4     weight = initWeight(exp_data.shape)
5
6     for i in range(1000):

```

```
7     a = Sigmoid(np.dot(exp_data, weight))
8     e = Loss(a, res_data) # 没啥用
9     grad = Gradient(exp_data, res_data, a)
10    weight -= grad*alpha
11
12    # 验证数据, 注意补上常数项1
13    print(1/(1+np.exp(-np.dot([[1, 0, 0, 1]], weight))))
14    pass
15
16
17    if __name__ == "__main__":
18        main()
19
20    '''
21    输出结果:
22    [[0.99872545]]
23    这里不是概率, 只表示去的可能性比较高
24    '''
```