

# 线性回归基本原理

## 【参考资料】

吴恩达机器学习笔记 <http://www.ai-start.com/ml2014/html/week2.html>

## 1. 基本形式

假设模型特征为  $(x_1, x_1, \dots, x_n)$ ,  $n$  代表特征的数量;  $x^{(i)}$  代表第  $i$  个训练实例,  $y^{(i)}$  为对应的标签,  $m$  为样本数量;  $m \times n$  维矩阵  $X$  为特征矩阵。

线性回归模型为:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

引入  $x_0 = 1$ , 则有

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

此时模型中的参数是一个  $n + 1$  维的向量, 任何一个训练实例也是一个  $n + 1$  维的向量, 特征矩阵  $X$  的维度是  $m \times (n + 1)$ 。因此公式整体可以转化为矩阵形式

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

## 2. 梯度下降

线性回归的目标函数为

$$J(\theta_0, \theta_1 \dots \theta_n) = \frac{1}{2m} \sum_{i=1}^m \left( h_{\theta} \left( x^{(i)} \right) - y^{(i)} \right)^2$$

所以梯度更新的公式为:

$$\begin{aligned} \theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n) \\ &= \theta_j - \alpha \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m \left( h_{\theta} \left( x^{(i)} \right) - y^{(i)} \right)^2 \\ &= \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \left( \left( h_{\theta} \left( x^{(i)} \right) - y^{(i)} \right) \cdot x_j^{(i)} \right) \end{aligned}$$

## 3. 正规方程

除了梯度下降之外, 线性回归模型的参数还可以通过正规方程法求解, 即:

$$\theta = (X^T X)^{-1} X^T y$$

推导过程如下：

首先将目标函数

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left( h_{\theta} \left( x^{(i)} \right) - y^{(i)} \right)^2$$

改写为矩阵形式

$$J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

展开上式

$$\begin{aligned} J(\theta) &= \frac{1}{2} (X\theta - y)^T (X\theta - y) \\ &= \frac{1}{2} (\theta^T X^T - y^T) (X\theta - y) \\ &= \frac{1}{2} (\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta - y^T y) \end{aligned}$$

接下来，对 $J(\theta)$ 求偏导，需要用到以下几个矩阵的求导法则：

$$\begin{aligned} \frac{dAB}{dB} &= A^T \\ \frac{dX^T AX}{dX} &= 2AX \end{aligned}$$

所以有

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta} &= \frac{1}{2} \left( 2X^T X\theta - X^T y - (y^T X)^T - 0 \right) \\ &= \frac{1}{2} (2X^T X\theta - X^T y - X^T y - 0) \\ &= X^T X\theta - X^T y \end{aligned}$$

令 $\frac{\partial J(\theta)}{\partial \theta} = 0$ ，有

$$\theta = (X^T X)^{-1} X^T y$$

注意，对于 $(X^T X)^{-1}$ 不可逆的情况，比如特征之间相互不独立（即特征矩阵不满秩），正规方程法是无法使用的，不过在大多数情况下， $(X^T X)^{-1}$ 都是可逆的。

### 【梯度下降 v.s. 正规方程】

对比如下：

梯度下降	正规方程
需要选择学习率	不需要
需要多次迭代	一次运算得出
当特征数量 $n$ 大时也能较好适用	需要计算 $(X^T X)^{-1}$ ，如果特征数量 $n$ 较大则运算代价大，因为矩阵逆的计算时间复杂度为 $O(n^3)$ ，通常来说当 $n$ 小于10000 时还是可以接受的
适用于各种类型的模型	只适用于线性模型，不适合逻辑回归模型等其他模型

总结一下，只要特征变量的数目并不大，标准方程是一个很好的计算参数 $\theta$ 的替代方法。