

## 正规方程

上一小节中，我们使用批量梯度下降算法，通过不断迭代以求得最佳参数 $\theta$ 的值。本小节将介绍另一种方法——正规方程（The Normal Equations）来计算出最佳参数 $\theta$ 的值。

在介绍正规方程法之前，我们先看看一些基本概念。

### Matrix Derivatives

对于一个 $m * n$ 的矩阵到实数的函数映射 $f: R^{m*n} \mapsto R$ ，其关于 $A$ 的导数为：

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

其中 $A$ 为 $m * n$ 的矩阵。

便于理解，我们不妨假设矩阵 $A$ 为：

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

函数映射 $f: R^{2*2} \mapsto R$ 为：

$$f(A) = \frac{3}{2} A_{11} + 5A_{12}^2 + A_{21} A_{22}$$

根据上述公式，我们可得：

$$\nabla_A f(A) = \begin{bmatrix} \frac{3}{2} & 10A_{12} \\ A_{22} & A_{21} \end{bmatrix}$$

对于 $n * n$ 矩阵 $A$ ，我们将矩阵 $A$ 对角线上元素的和定义为矩阵 $A$ 的迹：

$$tr A = \sum_{i=1}^n A_{ii}$$

其中若矩阵A为 $1 \times 1$ ，即为一实数，则其迹为本身， $tr A = A$ 。

一些常用性质如下：

$$\begin{aligned} tr AB &= tr BA \\ tr ABC &= tr CAB = tr BCA \\ tr A &= tr A^T \\ tr(A + B) &= tr(A) + tr(B) \\ tra A &= atr A \end{aligned}$$

结合矩阵导数的概念有如下性质：

$$\nabla_A tr AB = B^T \quad (1)$$

$$\nabla_{A^T} f(A) = (\nabla_A f(A))^T \quad (2)$$

$$\nabla_A tr ABA^T C = CAB + C^T AB^T \quad (3)$$

$$\nabla_A |A| = |A|(A^{-1})^T \quad (4)$$

其中等式（1）要求 $AB$ 为方阵；等式（3）要求 $ABA^T C$ 为方阵；等式（4）要求矩阵A为非奇异矩阵，即可逆； $|A|$ 表示矩阵A的行列式。

## Least Squares Revisited

好了，现在让我们开始介绍正规方程法，以找到最佳参数 $\theta$ 的值最小化代价函数 $J(\theta)$ 。

在给定训练集中，我们可构建一个维度为 $m \times n$ 的矩阵 $X$ ，其中 $m$ 为样本个数， $n$ 为每个样本的特征变量个数。

$$X = \begin{bmatrix} -(\mathbf{x}^{(1)})^T & - \\ -(\mathbf{x}^{(2)})^T & - \\ \vdots & \\ -(\mathbf{x}^{(m)})^T & - \end{bmatrix}$$

同样，向量 $Y$ 为：

$$Y = \begin{bmatrix} (\mathbf{y}^{(1)})^T \\ (\mathbf{y}^{(2)})^T \\ \vdots \\ (\mathbf{y}^{(m)})^T \end{bmatrix}$$

根据 $h_{\theta}(\mathbf{x}^{(i)}) = (\mathbf{x}^{(i)})^T \theta$ ，我们可得：

$$\begin{aligned}
X\theta - Y &= \begin{bmatrix} (x^{(1)})^T \theta \\ (x^{(2)})^T \theta \\ \vdots \\ (x^{(m)})^T \theta \end{bmatrix} - \begin{bmatrix} (y^{(1)})^T \\ (y^{(2)})^T \\ \vdots \\ (y^{(m)})^T \end{bmatrix} \\
&= \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ h_{\theta}(x^{(2)}) - y^{(2)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix}
\end{aligned}$$

又因为对于向量 $z$ ，有 $z^T z = \sum_i z_i^2$ 。故我们可得：

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

所以，我们对代价函数 $J(\theta)$ 求偏导，可得：

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

$$= \frac{1}{2} \nabla_{\theta} (\theta^T X^T X\theta - \theta^T X^T Y - Y^T X\theta + Y^T Y) \quad (2)$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr}(\theta^T X^T X\theta - \theta^T X^T Y - Y^T X\theta + Y^T Y) \quad (3)$$

$$= \frac{1}{2} \nabla_{\theta} (\text{tr} \theta^T X^T X\theta - 2\text{tr} Y^T X\theta) \quad (4)$$

$$= \frac{1}{2} (X^T X\theta + X^T X\theta - 2X^T Y) \quad (5)$$

$$= X^T X\theta - X^T Y \quad (6)$$

其中等式（1）类似于完全平方展开得到等式（2）；等式（2）应用 $\text{tr} A = A$ 得到等式（3）；等式（3）应用 $Y^T Y$ 为实数，且实数的转置为其本身，从而得到等式（4）；等式（4）应用 $\text{tr} ABA^T C = CAB + C^T AB^T$ 得到等式（5）。

最后，我们令该偏导为0可得：

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

从而，我们求出了参数 $\theta$ 的值。