

线性回归

2022年3月29日, 星期二 11:28

1、监督学习包括:

- 分类
- 线性回归 (找一个具体的值)

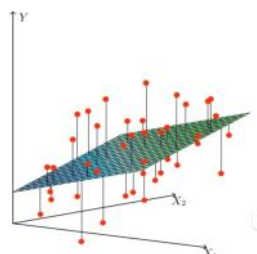
2、例子:

数据: 工资(x_1)、年龄 (x_2)

目标: 预测银行会贷多少款

考虑: 工资和年龄都会有影响,
那么他们各自的影响多大呢 (参数)

工资	年龄	额度
4000	25	2000
8000	30	70000
5000	28	35000
7500	33	50000
12000	40	85000



三、误差

• 误差的概率论意义:

误差 $\varepsilon^{(i)}$ 是独立并且具有相同分布, 并且服从均值为0, 方差为 σ^2 的高斯分布

$$\begin{aligned}h_{(\theta)}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\&= [1, x_1, x_2] \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \sum_0^2 \theta_i x_i = \theta^T x\end{aligned}$$

• 预测值与误差: $y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$, ε 为误差

1

误差服从高斯分布:

$$P(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\varepsilon^{(i)} - 0)^2}{2\sigma^2}}$$

2

将1式代入2式 (消去 $\varepsilon^{(i)}$), 得:

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \theta^T x^{(i)} - 0)^2}{2\sigma^2}}$$

3

似然函数:

$$L(\theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) = \prod_{i=1}^m \frac{1}{\theta\sqrt{2\pi}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\theta^2}}$$

4

解释: 什么样的参数跟我们的数据组合后恰好是真实值

对数函数:

$$\log L(\theta) = \log \prod_{i=1}^m \frac{1}{\theta\sqrt{2\pi}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\theta^2}}$$

5

解释: 乘法难解, 加法就容易了, 对数里面乘法可以转换成加法

- 展开化简5式:

$$\log \prod_{i=1}^m \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

$$= m \log \frac{1}{\sigma \sqrt{2\pi}} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

6

- 目标: 让似然函数越大越好, 也就是让蓝色部分越小越好 (概率论原理: 已经出现的概率越大越好)

- 目标函数 (损失函数、lossfunction):

$$\bar{J}(\theta) = \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \quad (\text{最小二乘法})$$

7

$$J \text{ 式} = \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 = \frac{1}{2} \sum_{i=1}^m (\overset{\text{预测值}}{h_{\theta}(x^{(i)})} - \overset{\text{真实值}}{y^{(i)}})^2 = \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

化成矩阵的形式, 等于矩阵的转置乘以自身

- 对目标函数求偏导:

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

标注 $\theta^T X^T X\theta$ 对 θ 求偏导得: $2X^T X\theta$

- 偏导等于0: $\theta = (X^T X)^{-1} X^T y$

偏导为0得极值点, 由实际意义得: $J(\theta)$ 函数取得最大值

8

问题来了: 根据公式8代入数据可直接求得 θ , 没有学习的过程, 而且不一定存在可逆矩阵

四: 真正机器学习的求解方法: 梯度下降

1、梯度下降

引入: 当我们得到了一个目标函数后, 如何进行求解? 直接求解? (并不一定可解, 线性回归可以, 但是其他问题不可以, 线性回归是个特例)

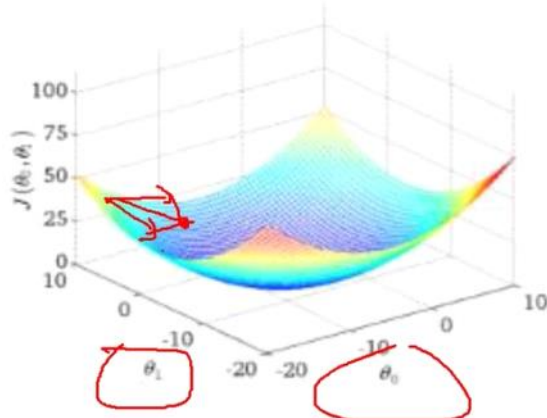
常规套路: 机器学习的套路就是交给机器一堆数据, 然后告诉他什么样的学习方式是对的 (目标函数), 然后让它

常规套路：机器学习的套路就是交给机器一堆数据，然后发告诉他什么样的学习方式是对的（目标函数），然后让它朝着这个方向去做

如何优化：一口吃不成一个胖子，我们要一步步的完成迭代（每次优化一点点，累计起来就是个大成绩了）

目标函数： $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

J是分别对两个参数求导更新， $\frac{\partial J}{\partial \theta_0}, \frac{\partial J}{\partial \theta_1}$ ，根据两个导数值（各自的方向）产生新的位置（新的方向）



寻找山谷最低点，也就是我们的目标函数终点（什么样的参数能使得目标函数达到极值点）

下山分几步走呢？（更新参数）

- 1、找到当前最合适的方向
- 2、走那么一小步，走快了容易“跌倒”
- 3、按照**方向与步伐**更新我们的参数

2、三种梯度下降方法（工程上的）

$$\text{目标函数: } J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2 \quad (\text{越小越好})$$
$$= \frac{1}{2m}$$

- 批量梯度下降（选取M个样本）： $\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^i$ 算出M个样本点的平均梯度（对一个参数 θ_j 进行求导）

更新 θ_j : $\theta_j = \theta_j + \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^i$
(容易得到最优解，但是由于每次考虑所有样本，速度很慢)

- 随机梯度下降: $\theta_j = \theta_j + (y^i - h_{\theta}(x^i)) x_j^i$ (算一个样本)

(每次找一个样本迭代，速度快，但不一定每次都朝着收敛方向)

- 小批量梯度下降法(miniBatch):

$$\theta_j = \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (y^{(k)} - h_{\theta}(x^{(k)})) x_j^k$$

选择10个样本，一般选64

(每次更新选择一小部分数据来算，实用!!!)

3、学习率 (步长)

学习率对结果会产生巨大影响，一般小一些。就是上面的 α

如何选择：从小的时候，不行再小

批量处理：32,64,128都可以，很多时候还得考虑内存和效率

正则化

2022年4月30日, 星期六 17:07

1、8式($X^T X$)不一定可逆, 表现为可能会出现过拟合

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

2、加入正则化项

$$J(\theta) = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \theta^T \theta \quad \text{正则项}$$

$$\begin{aligned} J(\theta): &= (\theta^T X^T - Y^T)(X\theta - Y) + \lambda \theta^T \theta \\ &= \theta^T X^T X \theta - 2\theta^T X^T Y + Y^T Y + \lambda \theta^T \theta \\ &= \theta^T (\theta^T \theta + \lambda I) \theta - 2\theta^T X^T Y + Y^T Y \end{aligned}$$

$$\begin{aligned} \text{对 } J(\theta) \text{ 求导: } \quad \frac{\partial J(\theta)}{\partial \theta} &= 2(\theta^T \theta + \lambda I) \theta - 2\theta^T Y = 0 \\ \theta &= (X^T X + \lambda I)^{-1} X^T Y \end{aligned}$$

这样的话($X^T X + \lambda I$)⁻¹就是可逆的了。就可以防止过拟合

加入正则项相当于什么呢?

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

相当于加了一个约束条件, 就是拉格朗日乘数法, 在约束条件下求最优解。

$$\begin{aligned} \text{约束条件: } \theta_1^2 + \theta_2^2 + \theta_3^2 + \dots \leq w \\ J(\theta) = \sum_{i=1}^m (\quad)^2 \end{aligned} \quad \downarrow$$

$$L(\theta) = J(\theta) + \lambda (\theta_1^2 + \theta_2^2 + \dots - w)$$

$$L(\theta) = J(\theta) + \lambda (\theta_1^2 + \theta_2^2 + \dots + w)$$

$$\begin{cases} \frac{\partial L}{\partial \theta_1} = 0 \\ \frac{\partial L}{\partial \theta_2} = 0 \\ \vdots \\ \frac{\partial L}{\partial \lambda} = 0 \end{cases}$$

在梯度下降中的表现

L2正则化是怎么避免过拟合的呢？先推导一下看看：

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial C0}{\partial \theta} + \frac{\lambda}{n} \theta^2$$

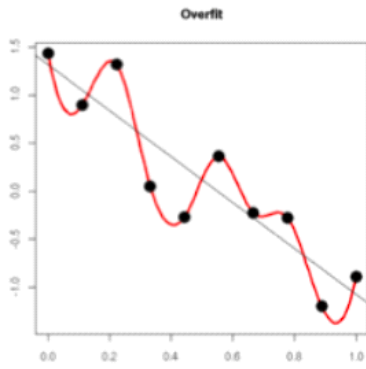
$$\frac{\partial J(\theta)}{\partial b} = \frac{\partial C0}{\partial b}$$

可以发现L2正则化项对b的更新没有影响，但是对θ的更新有影响：

$$\begin{aligned} \theta &\xrightarrow{\text{更新}} \theta - \frac{\partial C0}{\partial \theta} - \frac{\lambda}{n} \theta^2 \\ &= \left(1 - \frac{\lambda}{n}\right) \theta - \frac{\partial C0}{\partial \theta} \end{aligned}$$

在不使用L2正则化时，求导结果中w前系数为1，现在w前面系数为 $1 - \eta\lambda/n$ ，因为 η 、 λ 、 n 都是正的，所以 $1 - \eta\lambda/n$ 小于1，它的效果是减小w，这也就是**权重衰减 (weight decay)**的由来。当然考虑到后面的导数项，w最终的值可能增大也可能减小。

到目前为止，我们只是解释了L2正则化项有让w“变小”的效果，但是还没解释为什么w“变小”可以防止overfitting？一个所谓“显而易见”的解释就是：更小的权值w，从某种意义上说，表示网络的复杂度更低，对数据的拟合刚刚好（这个法则也叫做奥卡姆剃刀），而在实际应用中，也验证了这一点，L2正则化的效果往往好于未经正则化的效果。当然，对于很多人（包括我）来说，这个解释似乎不那么显而易见，所以这里添加一个稍微数学一点的解释（引自知乎）：



过拟合的时候，拟合函数的系数往往非常大，为什么？如下图所示，过拟合，就是拟合函数需要顾忌每一个点，最终形成的拟合函数波动很大。在某些很小的区间里，函数值的变化很剧烈。这就意味着函数在某些小区间里的导数值（绝对值）非常大，由于自变量值可大可小，所以只有系数足够大，才能保证导数值很大。