

2.线性模型

在求解参数模型参数时，涉及到模型选择问题，有两种策略，一种是经验风险最小化（如2.1），但容易出现过拟合，一种是结构风险最小化（如2.2），在经验风险最小化基础上加上一个正则化因子。

2.1 最小二乘法

解析解法

损失函数：

$$w = \arg \min \sum_{i=1}^n (y_i - f(x_i; w))^2 = \arg \min_w L$$
$$L = \sum_{i=1}^n (y_i - x_i^T w)^2 = \|y - Xw\|^2 = (y - Xw)^T (y - Xw)$$
$$\nabla_w L = 2X^T Xw - 2X^T y = 0 \Rightarrow w_{LS} = (X^T X)^{-1} X^T y$$

假设空间：

$$y_i = f(x_i; w) = w_0 + \sum_{j=1}^d x_{ij} w_j$$

如果原始变量由向量 x_i 组成，特征可以用基函数表示（类似于基本初等变化，对 X 进行空间变换），可以依旧使用上述公式求解。线性回归模型中的各个基函数如何选择，模型的函数就会有不同的结果。而且基函数的选择有时可以将线性模型扩展到一般的非线性形式，只要你将基函数定义为一个非线性函数就好，例如如果每一个基函数是三角函数，那么整个模型就是傅里叶变换的形式。

但是有个问题，一旦取逆不是满秩，那么它将不存在，不能使用解析解法，此时要用到梯度下降法。

梯度下降法

• 梯度下降法

考虑无约束优化问题 $\min_{\mathbf{x}} f(\mathbf{x})$

若能构造一个序列 $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$ 满足 $f(\mathbf{x}^{t+1}) < f(\mathbf{x}^t)$, $t = 0, 1, 2, \dots$

则不断执行该过程即可收敛到局部极小点

根据泰勒展式有 $f(\mathbf{x} + \Delta\mathbf{x}) \simeq f(\mathbf{x}) + \Delta\mathbf{x}^T \nabla f(\mathbf{x})$

于是, 欲满足 $f(\mathbf{x} + \Delta\mathbf{x}) < f(\mathbf{x})$, 可选择

$$\Delta\mathbf{x} = -\gamma \nabla f(\mathbf{x})$$

其中步长 γ 是一个小常数. 这就是梯度下降法

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla \mathcal{L}_n$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta (y_n - \mathbf{w}^{(\tau)T} \mathbf{x}_n) \mathbf{x}_n$$

其中 τ 表示迭代次数, η 是学习率参数

2.2 正则化

为了解决过拟合的问题, 引入正则化项, 加入罚函数, 第一项是经验误差, 第二项是惩罚项。

$$\mathbf{w}_{opt} = \arg \min_{\mathbf{w}} E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

惩罚项一般为范数惩罚,

$$\|\mathbf{w}\|_p = \left(\sum_{j=1}^d |w_j|^p \right)^{\frac{1}{p}} \quad 0 < p < \infty$$

$p=2$ 时, E_D 为线性函数时, 是岭回归。

$p=1$ 时, E_D 为线性函数时, 是LASSO回归。

岭回归

有时光采用训练样本进行无偏估计不好用, 要采取结构风险最小化, 岭回归的名称是因为求取参数的解析解的时候, 最后的表达式是在原来的基础上在求逆矩阵内部加上一个对角矩阵, 就好像一条“岭”一样。加上这条岭以后, 原来不可求逆的数据矩阵就可以求逆了。同时, 对角矩阵其实是由一个参数 λ 和单位对角矩阵相乘组成。 λ 越大, 说明偏差就越大, 原始数据对回归求取参数的作用就越小, 当 λ 取到一个合适的值, 就能在一定意义上解决过拟合的问题: 原先过拟合的特别大或者特别小的参数会被约束到正常甚至很小的值, 但不会为零。

- 岭回归 *Ridge regression*

$$w_{RR} = \arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

惩罚函数 $g(w) = \|w\|^2$ 将惩罚 w 中数值比较大的元素

$$\lambda \rightarrow 0 : w_{RR} \rightarrow w_{LS}$$

$$\lambda \rightarrow \infty : w_{RR} \rightarrow \vec{0}$$

$$\mathcal{L} = \|y - Xw\|^2 + \lambda \|w\|^2 = (y - Xw)^T (y - Xw) + \lambda w^T w$$

$$\nabla_w \mathcal{L} = -2X^T y + 2X^T Xw + 2\lambda w = 0$$

$$w_{RR} = (\lambda I + X^T X)^{-1} X^T y$$

λ 是超参数，不是越小越好/越大越好，需要提前设定。这里二范数的正则项偏向于把散点平滑地连在一起，减小（但不是变成0）不重要变量的 λ 值。**岭回归没有真正解决变量选择的问题。**

LASSO回归

不仅可以解决过拟合问题，而且可以在参数缩减过程中，将一些重复的没必要的参数直接缩减为零，也就是完全减掉了。这可以达到提取有用特征的作用。但是lasso回归的计算过程复杂，毕竟一范数不是连续可导的。

$$w_{LASSO} = \arg \min_w \|y - Ww\|_2^2 + \lambda \|w\|_1$$

关于稀疏解：

- 稀疏解

w 中的每个维度与数据 x 中的维度是对应的。

如果 $w_k = 0$, 则预测函数变为：

$$f(x, w) = x^T w = w_1 x_1 + \cdots + 0 \cdot x_k + \cdots + w_d x_d,$$

相当于预测将屏蔽了 k 个特征。

如果一个 w 满足大多数的维度 $= 0$, 则它被称为一个稀疏 *sparse* 解

缩减参数的目的：消除噪声特征，消除关联特征。

2.3 概率机器学习

线性回归的前提都是：

$$Y = Xw + \epsilon \quad \epsilon \sim N(0, \sigma^2 I) \quad Y \sim N(Xw, \sigma^2 I)$$

$$P(Y|Xw, \sigma^2 I) = \frac{1}{(\sqrt{2\pi}\sigma^2)^n} \exp\left\{-\frac{1}{2\sigma^2}(Y - Xw)^T(Y - Xw)\right\}$$

最小二乘回归：使用极大似然估计

$$w_{ML} = \arg \max_w \ln P(Y|\mu = Xw, \sigma^2 I)$$

$$= \arg \max_w -\frac{1}{2\sigma^2} \|Y - Xw\|^2 - \frac{n}{2} \ln(2\pi\sigma^2)$$

$$= \arg \min_w \|Y - Xw\|^2$$

会发现与之前的解析解相同。

岭回归：使用最大后验估计， w 服从正态分布

在最小二乘回归的基础上加入了参数的先验分布 $p(w)$ ：

$$w \sim N(0, \lambda^{-1} I)$$

$$p(w) = \left(\frac{\lambda}{2\pi}\right)^{\frac{d}{2}} e^{-\frac{\lambda}{2} w^T w}$$

• 岭回归的模型学习

$$w_{\text{MAP}} = \arg \max_w \ln p(y|w, X) + \ln p(w)$$

$$= \arg \max_w -\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) - \frac{\lambda}{2} w^T w + \text{const.}$$

$$\nabla_w \mathcal{L} = \frac{1}{\sigma^2} X^T y - \frac{1}{\sigma^2} X^T X w - \lambda w = 0$$

$w_{\text{MAP}} = (\lambda \sigma^2 I + X^T X)^{-1} X^T y$	参数 λ 反比于了参数先验分布的方差，即 λ 越大，我们假设模型应方差越小，所以岭回归可以鼓励学得一个光滑的模型。
--	---

LASSO回归：使用最大后验估计， w 服从拉普拉斯分布

$$w \sim Laplace(0, \lambda^{-1})$$

$$p(w) = \frac{\lambda}{2} e^{-\lambda|w|}$$

$$w_{MAP} = \arg \max P(Y | X, w) P(w)$$

$$= \arg \max \ln P(Y | X, w) \ln P(w)$$

$$= \arg \max \left(-\frac{1}{2\sigma^2} (Y - Xw)^T (Y - Xw) - \lambda |w| \right) + const$$

$$= \arg \min \left(\frac{1}{2\sigma^2} (Y - Xw)^T (Y - Xw) + \lambda |w| \right)$$

可以看出，与最小二乘法得出的结果一样