

第一章 回归模型

§1.1. 线性回归——最基本

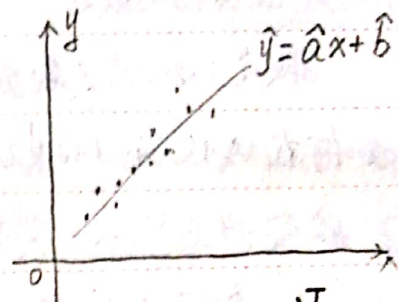
1. 模型背景

如果说给定样本 $D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$, 其中

$\vec{x}_i = (x_{i1}, x_{i2}, x_{i3}, \dots, x_{in})^T$ 是一个属性向量, 维数为 n .

若 y 与 \vec{x} 之间有线性关系, 即 $f(\vec{x}) = \omega^T \vec{x} + b$. $\vec{\omega} = (\omega_1, \omega_2, \omega_3, \dots, \omega_n)^T$.

那么我们就可以使用线性回归模型去对它进行拟合.



属性向量的分量 x_i 可能是连续的取值. 如果是离散的, 如高矮, 则人为数值化.

2. 评价指标: 均方误差 $(\omega^*, b^*) = \arg\min_{(\omega, b)} \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2$.

3. 模型求解:

11). 一对一的特殊情况: \vec{x} 只有 1 个分量, 目标函数 $E = \frac{1}{2} \sum_{i=1}^n (\omega x_i + b - y_i)^2$.

对参数 ω, b 求偏导, 解 $\begin{cases} \frac{\partial E}{\partial \omega} = \sum_{i=1}^n x_i (\omega x_i + b - y_i) = 0 \\ \frac{\partial E}{\partial b} = \sum_{i=1}^n (\omega x_i + b - y_i) = 0 \end{cases}$ 得极值点 $\begin{cases} \omega = \frac{\text{cov}(x, y)}{\text{var}(x)} \\ b = \bar{y} - \bar{x} \cdot \omega \end{cases}$

12). 一对多的一般情况: \vec{x} 有 n 个分量, 目标函数 $E = \frac{1}{2} \sum_{i=1}^n (\omega^T \vec{x}_i + b - y_i)^2$.

这是多元线性回归问题. 我们可以考虑用 $\hat{\omega} = (\omega^T; b)^T$ 代替

而 $X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} & 1 \\ x_{21} & x_{22} & \dots & x_{2n} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} & 1 \end{pmatrix}$ m 个数据点, 一个数据 n 个属性.

那么目标函数 $E(\hat{\omega}) = (y - X \hat{\omega})^T (y - X \hat{\omega})$. 求解方法有两种, 这里只介绍一下矩阵求导:

$$\frac{\partial E}{\partial \hat{\omega}} = 2X^T(X\hat{\omega} - y) = 0. \quad \text{那么解矩阵方程: } (X^T X) \hat{\omega} = X^T y.$$

讨论: ① 若 $X^T X$ 可逆: $\hat{\omega}^* = (X^T X)^{-1} X^T y$. 没有争议.

② 若 $X^T X$ 不可逆: 方程难解. 此时有多个目标解向量 $\hat{\omega}^*$.

选择哪一个最好, 还是要引入正则化项 (regularization).

4. 模型扩展:

如果发现 y 与 x 有指数关系, 那么 $\ln y = \omega x + b$, $y = e^{\omega x + b}$ (指数线性回归)

.....



§1.2. 岭回归 — 改进

1. 对线性回归的改进

假设：对于给定数据集，有一小部分数据有明显偏差，这显然会干扰 prediction。
如何在线性回归的基础上让机器自己剔除异常，就是 Ridge Regression 的优势。

2. 岭回归与 Lasso 回归的改进

岭回归 Ridge Regression 以及 Lasso Regression 都在线性回归的基础上加了罚项进行正则化。

Ridge: L_2 正则化 $J(w) = \frac{1}{2m} \sum_{i=1}^m (f_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$ θ 为系数向量

Lasso: L_1 正则化 $J(w) = \frac{1}{2m} \sum_{i=1}^m (f_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j|$

Lasso 回归使损失函数中许多 θ 都可以为 0。计算量相对更小。
这里， λ 叫正则化系数。过大易欠，过小则过。

3. 岭回归求解：

$$\frac{\partial J(w)}{\partial \theta} = 0 \Rightarrow 0 - X^T y - X^T y + 2X^T X \theta + 2\lambda \theta = 0$$

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

Lasso 回归求解：比较复杂

$$\theta = \begin{cases} (w_j - \frac{\lambda}{n_j}) / n_j, & w_j \geq \frac{\lambda}{n_j} \\ 0, & \text{else} \\ (w_j + \frac{\lambda}{n_j}) / n_j, & w_j < -\frac{\lambda}{n_j} \end{cases}$$

$$\begin{cases} w_j = \sum_{i=1}^n X_{ij}(y_i - \sum_k \theta_k X_{ik}) \\ n_j = \sum_{i=1}^n X_{ij}^2 \end{cases}$$

4. λ 值选取

(1). 岭迹法。由 $\theta = (X^T X + \lambda I)^{-1} X^T y$ ， θ 是 λ 的函数

θ 趋于稳定的 λ 就是岭迹

(2). 交叉验证 (建议)

