

机器学习基础(1)-线性回归

给定由 d 个属性描述的示例 $x = (x_1, x_2, \dots, x_d)$, 其中 x_i 是 x 在第 i 个属性上的取值, 线性回归试图学得一个通过属性的线性组合来进行预测的函数, 即

$$f(x) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$$

一般用向量形式写成,

$$f(x) = w^T x + b$$

其中, $w = (w_1; w_2; \dots; w_d)$ 。 w 和 b 学得之后, 模型就得以确定。

给定数据集 $D = (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$, 其中, $x_i = (x_{i1}; x_{i2}; \dots; x_{id})$, $y_i \in R$ 。 “线性回归” (linear regression) 试图学得一个线性模型以尽可能准确的预测实际输出标记。

我们先考虑一种最简单的情况: 输入属性的数目只有一个。线性回归试图学得,

$$f(x_i) = wx_i + b, \text{ 使得 } f(x_i) \simeq y_i$$

如何确定 w 和 b 呢? 显然, 关键在于如何衡量 $f(x)$ 与 y 之间的差别。第二章中介绍过, 均方误差是回归任务中常用的性能度量, 因此我们可以试图让均方误差最小化, 即,

$$\begin{aligned} (w^*, b^*) &= \arg \min_{(w, b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \\ &= \arg \min_{(w, b)} \sum_{i=1}^m (y_i - wx_i - b)^2 \end{aligned}$$

均方误差有非常好的几何意义, 它对应了常用的欧几里得距离或简称“欧式距离” (Euclidean distance)。基于均方误差最小化进行模型求解的方法称为“最小二乘法” (least square method)。在线性回归中, 最小二乘法就是输入找到一条直线, 使所有样本到直线上的欧式距离之和最小。

求解 w 和 b 使,

$$E_{(w,b)} = \sum_{i=1}^m (y_i - wx_i - b)^2$$

最小化的过程, 称为线性回归模型的最小二乘“参数估计”(parameter estimation)。我们可以将 $E(w,b)$ 分别对 w 和 b 求导, 得到,

$$\begin{aligned}\frac{\partial E_{(w,b)}}{\partial w} &= 2 \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b) x_i \right) \\ \frac{\partial E_{(w,b)}}{\partial b} &= 2 \left(mb - \sum_{i=1}^m (y_i - wx_i) \right)\end{aligned}$$

然后, 另上面的式子为零, 从而求得 w 和 b 的最优解,

$$\begin{aligned}w &= \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} \left(\sum_{i=1}^m x_i \right)^2} \\ b &= \frac{1}{m} \sum_{i=1}^m (y_i - wx_i)\end{aligned}$$

更一般的情况是数据集 D , 样本由 d 个属性描述。此时我们试图学得,

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b, \text{ 使得 } f(\mathbf{x}_i) \simeq y_i$$

这称为“多元线性回归”(multivariate linear regression)。类似的, 可利用最小二乘法来对 w 和 b 进行估计。为了便于讨论, 我们把 w 和 b 吸入向量形式,

$$\hat{\boldsymbol{w}} = (\boldsymbol{w}; b)$$

相应的，把数据集 D 表示为一个 $m \times (d+1)$ 大小的矩阵 X ，其中，每行对应于一个示例，该行前 d 个元素对应于示例的 d 个属性值，最后一个元素恒置为1，即，

$$\boldsymbol{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} \boldsymbol{x}_1^T & 1 \\ \boldsymbol{x}_2^T & 1 \\ \vdots & \vdots \\ \boldsymbol{x}_m^T & 1 \end{pmatrix}$$

再把标记也写成向量形式 $\boldsymbol{y} = (y_1; y_2; \dots; y_m)$ ，则有，

$$\hat{\boldsymbol{w}}^* = \arg \min_{\hat{\boldsymbol{w}}} (\boldsymbol{y} - \boldsymbol{X}\hat{\boldsymbol{w}})^T (\boldsymbol{y} - \boldsymbol{X}\hat{\boldsymbol{w}})$$

当

$$|\boldsymbol{X}^T \boldsymbol{X}|$$

为满秩矩阵或正定矩阵时，可求得，

$$\hat{\boldsymbol{w}}^* = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$$

其中 $(\boldsymbol{X}^T \boldsymbol{X})^{-1}$ 是矩阵 $(\boldsymbol{X}^T \boldsymbol{X})$ 的逆矩阵。令 $\hat{\boldsymbol{x}}_i = (x_i, 1)$ 则最终学得的线性回归模型为

$$f(\hat{\boldsymbol{x}}_i) = \hat{\boldsymbol{x}}_i^T (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$$

然而，现实任务中 $\boldsymbol{X}^T \boldsymbol{X}$ 往往不是满秩的。例如我们在实际场景中，经常会遇到这样的数据集，数据的属性超过样例数，导致 \boldsymbol{X} 的列数大于行数， $\boldsymbol{X}^T \boldsymbol{X}$ 显然不满秩；也可能由于存在线性相关的属性，导致 $\boldsymbol{X}^T \boldsymbol{X}$ 不满秩。此时可解出多个 \boldsymbol{w} 的估计值，他们都能使均方误差最小化。这是，通常的做法是引入正则化项，求解出最优解 $\hat{\boldsymbol{w}}$ 。也可以这样理解，当数据属性远多于样例数时，更容易出现过拟合，通过引入正则化项，得到稀疏解，降低过拟合的风险。