

一、线性回归

线性回归的目的是试图学得一线性模型, 即

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

我们需要确定 w 和 b , 采用均方误差, 试图将均方误差最小化,

即

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

使用最小二乘法来求解均方误差最小化。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”。

我们将 $E(w, b)$ 分别对 w 和 b 求导, 得:

$$\frac{\partial E(w, b)}{\partial w} = 2 \sum_{i=1}^m (y_i - wx_i - b)(-x_i)$$
$$= 2 \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i \right)$$

$$\frac{\partial E(w, b)}{\partial b} = 2 \sum_{i=1}^m (y_i - wx_i - b)(-1)$$
$$= 2 \left(mb - \sum_{i=1}^m (y_i - wx_i) \right)$$

$$\text{令 } \frac{\partial E(w, b)}{\partial w} = 0, \frac{\partial E(w, b)}{\partial b} = 0, \text{得}$$

$$\frac{\partial E(w, b)}{\partial b} = 2 \left(mb - \sum_{i=1}^m (y_i - wx_i) \right) = 0$$

$$mb = \sum_{i=1}^m (y_i - wx_i)$$

$$b = \frac{1}{m} \sum_{i=1}^m (y_i - wx_i)$$

$$\frac{\partial E(w, b)}{\partial w} = 2 \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i \right) = 0$$

$$w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m y_i x_i + \sum_{i=1}^m b x_i = 0$$

$$w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m y_i x_i + \sum_{i=1}^m x_i \left(\frac{1}{m} \sum_{i=1}^m (y_i - wx_i) \right) = 0$$

$$W \sum_{i=1}^m X_i^2 - \sum_{i=1}^m y_i X_i + \sum_{i=1}^m X_i \left(\frac{1}{m} \sum_{i=1}^m y_i - \frac{1}{m} \sum_{i=1}^m W X_i \right) = 0$$

$$W \sum_{i=1}^m X_i^2 - \frac{W}{m} \left(\sum_{i=1}^m X_i \right)^2 = \sum_{i=1}^m y_i X_i - \sum_{i=1}^m y_i \bar{X}$$

$$W = \frac{\sum_{i=1}^m y_i (X_i - \bar{X})}{\sum_{i=1}^m X_i^2 - \frac{1}{m} \left(\sum_{i=1}^m X_i \right)^2}$$

多元线性
回归

更一般地, $f(X_i) = W^T X_i + b$, 使 $f(X_i) \simeq y_i$

同样的, 利用最小二乘法对 W, b 进行估计。便于讨论, 将 W 和 b 写成向量形式 $\hat{W} = (W; b)$, 将数据集 D 表示为一个 $m \times (d+1)$ 大小的矩阵 X , 即

$$X = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1n} & 1 \\ X_{21} & X_{22} & \dots & X_{2n} & 1 \\ \vdots & \vdots & & \vdots & \vdots \\ X_{m1} & X_{m2} & \dots & X_{md} & 1 \end{pmatrix} = \begin{pmatrix} X_1^T & 1 \\ X_2^T & 1 \\ \vdots & \vdots \\ X_m^T & 1 \end{pmatrix}$$

故此时, 有 $\hat{W}^* = \arg \min_{\hat{W}} (y - X\hat{W})^T (y - X\hat{W})$

$$\text{令 } E(\hat{W}) = (y - X\hat{W})^T (y - X\hat{W})$$

$$E(\hat{W}) = (y - X\hat{W})^T (y - X\hat{W})$$

$$= (y^T - \hat{W}^T X^T) (y - X\hat{W})$$

$$= y^T y - y^T X \hat{W} - \hat{W}^T X^T y + \hat{W}^T X^T X \hat{W}$$

对 \hat{W} 求导, 得

$$\frac{\partial E(\hat{W})}{\partial \hat{W}} = \frac{\partial (y^T y - y^T X \hat{W} - \hat{W}^T X^T y + \hat{W}^T X^T X \hat{W})}{\partial \hat{W}}$$

$$\frac{\partial \hat{W}^T X^T X \hat{W}}{\partial \hat{W}} = \frac{\partial (X^T X) (\hat{W}_1^2 + \hat{W}_2^2 + \dots + \hat{W}_{d+1}^2)}{\partial \hat{W}_i} = 2 X^T X \hat{W}$$

$$\frac{\partial y^T X \hat{W}}{\partial \hat{W}} = X^T y$$

$$\frac{\partial E(\hat{W})}{\partial \hat{W}} = 2 X^T X \hat{W} - X^T y - X^T y$$

如果 $X^T X$ 为满秩矩阵或正定矩阵, 令 $\frac{\partial E(\hat{w})}{\partial \hat{w}} = 0$, 有

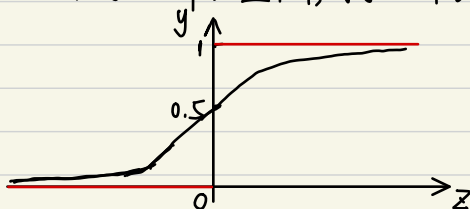
$$\hat{w} = (X^T X)^{-1} X^T y$$

二、对数几率回归

对数几率回归常用于分类任务。首先考虑二分类任务, 其输出标记 $y \in \{0, 1\}$, 而线性回归产生的预测值 $z = w^T x + b$ 是实值, 我们需要将 z 转换为 0/1 值。我们采用单位阶跃函数。

$$y = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$$

即若预测值 z 大于 0 则为正例, 小于 0 则为反例, 等于 0 则可任意判别



单位阶跃函数与对数几率函数

从上图可以得知单位阶跃函数不连续, 故我们采用一个常用的替代函数

$$y = \frac{1}{1 + e^{-z}}$$

$$z = w^T x + b$$

此时

$$\begin{aligned} y &= \frac{1}{1 + e^{-(w^T x + b)}} \\ \frac{1}{y} &= 1 + e^{-(w^T x + b)} \\ \frac{1-y}{y} &= e^{-(w^T x + b)} \\ \ln \frac{1-y}{y} &= -(w^T x + b) \\ \ln \frac{y}{1-y} &= w^T x + b \end{aligned}$$

若将 y 视为样本 x 作为正例的可能性, 则 $1-y$ 是其反例的可能, 两者的比值称为“几率”, 反映了 x 作为正例的可能性. 对几率取对数则得到“对数几率”:

$$\ln \frac{y}{1-y}$$

此时我们来确定前面的 w 和 b . 如果将 $y = \frac{1}{1+e^{-(w^T x + b)}}$ 中的 y 视为类后验概率估计 $p(y=1|x)$, 则

$$\ln \frac{p(y=1|x)}{p(y=0|x)} = w^T x + b$$

那么有

$$p(y=1|x) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}}$$

$$p(y=0|x) = \frac{1}{1 + e^{w^T x + b}}$$

于是, 我们可以通过极大似然法来估计 w 和 b . 给定数据集 $\{(x_i, y_i)\}_{i=1}^m$ 对率回归模型最大化“对数似然”

$$\ell(w, b) = \sum_{i=1}^m \ln p(y_i | x_i; w, b)$$

为了便于讨论, 令 $\beta = (w, b)$, $\hat{x} = (x; 1)$, 则 $w^T x + b$ 可简化为 $\beta^T \hat{x}$, 再令 $p_1(\hat{x}; \beta) = p(y=1 | \hat{x}; \beta)$, $p_0(\hat{x}; \beta) = p(y=0 | \hat{x}; \beta) = 1 - p_1(\hat{x}; \beta)$, 则

$$p(y_i | x_i; w, b) = y_i p_1(\hat{x}_i; \beta) + (1 - y_i) p_0(\hat{x}_i; \beta)$$

$$\begin{aligned} \ell(\beta) &= \sum_{i=1}^m \ln (y_i p_1(\hat{x}_i; \beta) + (1 - y_i) p_0(\hat{x}_i; \beta)) \\ &= \sum_{i=1}^m (-y_i \beta^T \hat{x}_i + \ln(1 + e^{\beta^T \hat{x}_i})) \end{aligned}$$

推导: 因为 $p_1(\hat{x}_i; \beta) = p(y=1 | \hat{x}_i; \beta) = \frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}}$

$$p_0(\hat{x}_i; \beta) = p(y=0 | \hat{x}_i; \beta) = \frac{1}{1 + e^{\beta^T \hat{x}_i}}$$

$$\ell(\beta) = \sum_{i=1}^m \ln \left(\frac{y_i e^{\beta^T \hat{x}_i} + 1 - y_i}{1 + e^{\beta^T \hat{x}_i}} \right)$$

$$\ell(\beta) = \begin{cases} \sum_{i=1}^m (-\ln(1 + e^{\beta^T \hat{x}_i})) & y_i = 0 \\ \sum_{i=1}^m (\beta^T \hat{x}_i - \ln(1 + e^{\beta^T \hat{x}_i})) & y_i = 1 \end{cases}$$

综合可得

$$\ell(\beta) = \sum_{i=1}^m (y_i \beta^T \hat{x}_i - \ln(1 + e^{\beta^T \hat{x}_i}))$$

所得出来的 $l(\beta)$ 是关于 β 的高阶可连续凸函数, 于是有

$$\beta^* = \arg \min l(\beta)$$

以牛顿法为例求解, 其第 $t+1$ 轮迭代解的更新公式为

$$\beta^{t+1} = \beta^t - \left(\frac{\partial^2 l(\beta^t)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial l(\beta^t)}{\partial \beta}$$

牛顿法原理:

牛顿法是一种线性化方法, 其基本思想是将非线性方程 $f(x)=0$ 逐步归结为某种线性方程来求解.

设已知方程 $f(x)=0$ 有近似根 x_k (假定 $f'(x_k) \neq 0$), 将函数 $f(x)$ 在点 x_k 展开, 有

$$f(x) \approx f(x_k) + f'(x_k)(x - x_k)$$

于是方程 $f(x)=0$ 可近似地表示为

$$f(x_k) - f'(x_k)(x - x_k) = 0$$

这是个线性方程, 记其根 x_{k+1} , 则 x_{k+1} 的计算公式为

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k=0, 1, \dots$$

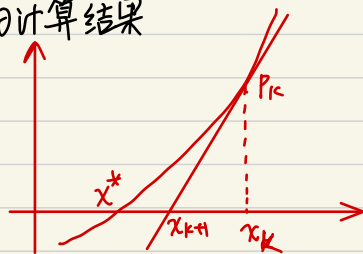
牛顿法的几何解释:

方程 $f(x)=0$ 的根 x^* 可解释为曲线 $y=f(x)$ 与 x 轴的交点的横坐标. 设 x_k 是根 x^* 的某个近似值, 过曲线 $y=f(x)$ 上横坐标为 x_k 的点 P_k 引切线, 并将该切线与 x 轴的交点的横坐标 x_{k+1} 作为 x^* 的新的近似值.

注意到切线方程为

$$y = f(x_k) + f'(x_k)(x - x_k)$$

此时求得值 x_{k+1} 就是牛顿公式的计算结果



设 $f(x)$ 二次连续可微, $x_k \in \mathbb{R}^n$, Hesse 矩阵 $\nabla^2 f(x_k)$ 正定. 我们在 x_k 附近用二次泰勒展开近似 $f(x)$.

$$f(x) = f(x_k) + \nabla f(x_k)^T (x - x_k) + \frac{1}{2} (x - x_k)^T \nabla^2 f(x_k) (x - x_k)$$

$$\text{令 } s = x - x_k \quad x = s + x_k$$

$$f(s + x_k) \approx f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T \nabla^2 f(x_k) s$$

令 $q^k(s) = f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T \nabla^2 f(x_k) s$ 为 $f(x)$ 的二次近似, 求解自变量迭代增量 s 使 $q^k(s)$ 最小化, 即对 $q^k(s)$ 求导并令导数等于 0:

$$\nabla q^k(s) = \nabla f(x_k) + \nabla^2 f(x_k) s = 0$$

$$\Rightarrow s = -[\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$

$$\text{代入 } x = s + x_k, \text{ 有 } x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$

对应到式中

$$\beta^* = x_k, [\nabla^2 f(x_k)]^{-1} = \left(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \nabla f(x_k) = \frac{\partial l(\beta)}{\partial \beta}$$

其关于 β 的一阶、二阶导数分别为

$$\frac{\partial l(\beta)}{\partial \beta} = - \sum_{i=1}^m \hat{x}_i (y_i - p_i(\hat{x}_i; \beta))$$

推导:

$$\begin{aligned} \frac{\partial l(\beta)}{\partial \beta} &= \frac{\partial \sum_{i=1}^m (-y_i \beta^T \hat{x}_i + \ln(1 + e^{\beta^T \hat{x}_i}))}{\partial \beta} \\ &= \sum_{i=1}^m \left(\frac{\partial (-y_i \beta^T \hat{x}_i)}{\partial \beta} + \frac{\partial (\ln(1 + e^{\beta^T \hat{x}_i}))}{\partial \beta} \right) \\ &= \sum_{i=1}^m \left(-y_i \hat{x}_i + \frac{1}{1 + e^{\beta^T \hat{x}_i}} \cdot e^{\beta^T \hat{x}_i} \cdot \hat{x}_i \right) \\ &= - \sum_{i=1}^m \hat{x}_i \left(y_i - \frac{1}{1 + e^{\beta^T \hat{x}_i}} \cdot e^{\beta^T \hat{x}_i} \right) \\ &= - \sum_{i=1}^m \hat{x}_i (y_i - p_i(\hat{x}_i; \beta)) \end{aligned}$$

$$\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^m \hat{x}_i \hat{x}_i^T p_i(\hat{x}_i; \beta) (1 - p_i(\hat{x}_i; \beta))$$

推导:

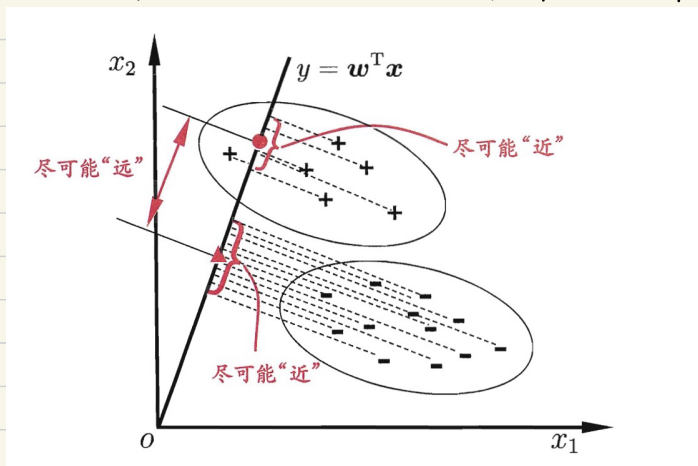
$$\begin{aligned} \frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} &= \frac{\partial \sum_{i=1}^m \hat{x}_i (y_i - \frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}})}{\partial \beta \partial \beta^T} \\ &= - \sum_{i=1}^m \hat{x}_i \frac{\partial (y_i - \frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}})}{\partial \beta^T} \\ &= - \sum_{i=1}^m \hat{x}_i \left(- \frac{\partial y_i}{\partial \beta^T} - \frac{\partial (\frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}})}{\partial \beta^T} \right) \end{aligned}$$

$$\begin{aligned}
 \text{第-项 } \frac{\partial y_i}{\partial \beta} &= 0, \text{ 第-项 } \frac{\partial \left(\frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}} \right)}{\partial \beta^T} = \frac{e^{\beta^T \hat{x}_i} \cdot \hat{x}_i (1 + e^{\beta^T \hat{x}_i}) - e^{\beta^T \hat{x}_i} \cdot e^{\beta^T \hat{x}_i} \cdot \hat{x}_i}{(1 + e^{\beta^T \hat{x}_i})^2} \\
 &= \hat{x}_i^T e^{\beta^T \hat{x}_i} \frac{(1 + e^{\beta^T \hat{x}_i}) - e^{\beta^T \hat{x}_i}}{(1 + e^{\beta^T \hat{x}_i})^2} \\
 &= \hat{x}_i^T \cdot \frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}} \cdot \frac{1}{1 + e^{\beta^T \hat{x}_i}} = \frac{\hat{x}_i^T}{1 + e^{\beta^T \hat{x}_i}} \\
 \text{则 } \frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} &= - \sum_{i=1}^m \hat{x}_i \left(- \hat{x}_i^T \cdot \frac{1}{1 + e^{\beta^T \hat{x}_i}} \cdot \frac{1}{1 + e^{\beta^T \hat{x}_i}} \right) \\
 &= \sum_{i=1}^m \hat{x}_i \hat{x}_i^T p_i(\hat{x}_i; \beta) (1 - p_i(\hat{x}_i; \beta))
 \end{aligned}$$

三、线性判别分析

LDA既是一种分类算法，同时也是一种降维方法。

线性判别分析(LDA)的思想是：给定训练集，设法将样例投影到一条直线上，使得同类样例的投影点尽可能接近，异类样例的投影点尽可能远离；在对新样本分类时，将其投影到同样的这条直线上，再根据投影点的位置来确定新样本的类别。



LDA示意图。

协方差矩阵的定义：对于两个变量的总体误差，期望值分别为E(X)和E(Y)的两个随机变量X与Y之间的协方差。

协方差Cov(X, Y)定义为：

$$\begin{aligned}
 \text{Cov}(X, Y) &= E(X - E(X))(Y - E(Y)) \\
 &= E(XY) - E(X)E(Y) \\
 &= E(XY) - EXEY
 \end{aligned}$$

设X = (X1, X2, ..., Xn)^T 为n维随机变量，协方差矩阵C = Cov(X) = Cov(X1, X2, ..., Xn)^T。

为n维随机变量X的协方差矩阵。

为记为D(X)，即Cij = Cov(Xi, Xj)。

为X的各分量Xi和Y的协方差。由于

Cij = Cji，所以协方差矩阵为对称阵。

非负定矩阵。

给定数据集 $D = \{(x_i, y_i) | i=1, \dots, m, y_i \in \{0, 1\}\}$ ，令 X_i, μ_i, Σ_i 分别表示第 $i \in \{0, 1\}$ 类样例的集合、均值向量、协方差矩阵。若将数据投影到直线 w 上，则两类样本的中心在直线上的投影分别为 $w^T \mu_0$ 和 $w^T \mu_1$ ，两类样本的协方差分别为 $w^T \Sigma_0 w$ 和 $w^T \Sigma_1 w$ 。由于直线是一维空间，因此 $w^T \mu_0, w^T \mu_1, w^T \Sigma_0 w$ 和 $w^T \Sigma_1 w$ 均为实数。

如果想要使同类样例的投影点尽可能接近,可以让同类样例投影点的协方差尽可能小,即 $W^T \Sigma_0 W + W^T \Sigma_1 W$ 尽可能小,而想要使异类样例的投影点尽可能远离,可以让类中心之间的距离尽可能大,即 $\|W^T \mu_0 - W^T \mu_1\|_2^2$ 尽可能大. 即我们同时考虑二者,则得到欲最大化的目标:

$$J = \frac{\|W^T \mu_0 - W^T \mu_1\|_2^2}{W^T \Sigma_0 W + W^T \Sigma_1 W}$$

$$= \frac{W^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T W}{W^T (\Sigma_0 + \Sigma_1) W}$$

推导:

$$J = \frac{\|W^T \mu_0 - W^T \mu_1\|_2^2}{W^T \Sigma_0 W + W^T \Sigma_1 W}$$

$$= \frac{\|(W^T \mu_0 - W^T \mu_1)^T\|_2^2}{W^T (\Sigma_0 + \Sigma_1) W}$$

$$= \frac{\|(\mu_0 - \mu_1)^T W\|_2^2}{W^T (\Sigma_0 + \Sigma_1) W}$$

$$= \frac{[(\mu_0 - \mu_1)^T W]^T (\mu_0 - \mu_1)^T W}{W^T (\Sigma_0 + \Sigma_1) W}$$

$$= \frac{W^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T W}{W^T (\Sigma_0 + \Sigma_1) W}$$

此时定义“类内散度矩阵”

$$S_w = \Sigma_0 + \Sigma_1$$

$$= \sum_{x \in X_0} (x - \mu_0)(x - \mu_0)^T + \sum_{x \in X_1} (x - \mu_1)(x - \mu_1)^T$$

以及“类间散度矩阵”

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T$$

则 $J = \frac{W^T S_b W}{W^T S_w W}$, 这就是 LDA 欲最大化的目标, 即 S_b 与 S_w 的“广义瑞利商”。

可以发现上式的分子和分母都是关于 W 的二次项, 因此它的解与 W 的长度无关, 只与其方向有关. 不失一般性, 令 $W^T S_w W = 1$, 则目标可改写为

$$\max_W -W^T S_b W \quad \text{s.t.} \quad W^T S_w W = 1$$

此时, 我们采用拉格朗日乘子法进行求解, 则

$$S_b W = \lambda S_w W$$

其中 λ 为拉格朗日乘子. 因为 $S_b W$ 的方向恒为 $\mu_0 - \mu_1$, 则令

$$S_b W = \lambda (\mu_0 - \mu_1)$$

$$\text{故 } W = S_w^{-1} (\mu_0 - \mu_1)$$

