

日期:

线性模型

1. 线性回归

1. 回归

① 定义: 给定一组数据点 $(x_1, y_1), \dots, (x_n, y_n)$, 根据这些数据点研究 x 和 y 之间关系的分析方法就是回归。

② 线性回归:

(1) 利用数理统计中回归分析, 来确定 2 种或 2 种以上变量间相互依赖的定量关系的一种统计分析方法。

(2) 以一个线性函数 $f(x) = w^T x + b$ 描述两者关系。

③ 逻辑回归: 以逻辑函数描述

2. 线性模型表示形式

① 一般形式: $f(x) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$

$x = [x_1, x_2, \dots, x_d]^T$ 是由属性描述的样本, 其中 x_i 是 x 在第 i 个属性上的取值。

② 向量形式:

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

其中 $w = (w_1, w_2, \dots, w_d) = [w_1, w_2, \dots, w_d]^T$

3. 线性模型优点

日期:

/

① 线性简单, 易于建模.

② 可解释性

③ 非线性模型的基础: 引入特征映射和高维映射.

一个例子

- 综合考虑色泽、根蒂和敲声来判断西瓜好不好
- 其中根蒂的系数最大, 表明根蒂最要紧; 而敲声的系数比色泽大, 说明敲声比色泽更重要

$$f_{\text{好瓜}}(x) = 0.2 \cdot x_{\text{色泽}} + 0.5 \cdot x_{\text{根蒂}} + 0.3 \cdot x_{\text{敲声}} + 1$$

4. 数据.

① 给定 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$.

$$x_i = (x_{i1}, x_{i2}, \dots, x_{id}) \quad y_i \in \mathbb{R}.$$

② 离散属性处理.

{ 有序关系: 连续化为连续值

{ 无序关系: 有 k 个属性值, 则化为 k 维向量

5. 单一属性的线性回归流程

① 模型表示 $f(x_i) = w x_i + b$

② 代价函数 (loss function) 即为误差

③ 平方误差最小化模型求解方法: 最小二乘法

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

日期: /

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

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

最小化均方误差

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

分别对 w 和 b 求偏导可得

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

得:

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

$$\text{其中 } \bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

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

6. 多元线性回归.

① 数据集.

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$$

$$\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{id}) \quad (d \text{ 为属性维数})$$

$$= (x_{i1}, x_{i2}, \dots, x_{id})^T$$

$$= [x_{i1}, x_{i2}, \dots, x_{id}]^T$$

$$y_i \in \mathbb{R} \quad i = 1, 2, \dots, m$$

② 目标: $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b$ 使得 $f(\mathbf{x}_i)$ 逼近 y_i .

日期: /

其中: w, b 为向量形式: $\hat{w} = [w; b] = [w, b]^T$

$$X = \begin{pmatrix} \overbrace{x_{11} \ x_{12} \ \cdots \ x_{1d}}^{x_1} & 1 \\ x_{21} \ x_{22} \ \cdots \ x_{2d} & 1 \\ \vdots & \vdots \\ x_{m1} \ x_{m2} \ \cdots \ x_{md} & 1 \end{pmatrix} = \begin{pmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_m^T & 1 \end{pmatrix} = \begin{bmatrix} \hat{x}_1^T \\ \hat{x}_2^T \\ \vdots \\ \hat{x}_m^T \end{bmatrix} \in \mathbb{R}^{m \times (d+1)}$$

$$\hat{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}, 1]^T \quad y = (y_1; y_2; \dots; y_m)$$

③ 求解: 最小二乘模型.

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

$$\Rightarrow \hat{w}^* = \arg \min_{\hat{w}} \sum_{i=1}^m (\hat{w}^T \hat{x}_i - y_i)^2 = \arg \min_{\hat{w}} \sum_{i=1}^m (\hat{x}_i^T \hat{w} - y_i)^2$$

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

$$= \arg \min_{\hat{w}} \|X\hat{w} - y\|_2^2$$

$$= \arg \min \hat{w}^T X^T X \hat{w} - 2 \hat{w}^T X^T y + y^T y$$

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

$$\Rightarrow \frac{\partial E\hat{w}}{\partial w} = 2X^T (X\hat{w} - y) = 0$$

□ 若 $X^T X$ 是满秩矩阵或正定矩阵(矩阵可逆), 则

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

其中 $(X^T X)^{-1}$ 是 $X^T X$ 的逆矩阵, 线性回归模型为

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

日期: /

7. 对数线性回归

① 定义: 输出标记的对数为线性模型逼近的目标.

② 适用场合: 样本对 y 不是线性变化, 而是指数变化.

③ 表示形式: $\ln y = w^T x + b$

8. 广义线性模型

① 表示形式: $y = g^{-1}(w^T x + b)$

$g(\cdot)$ 称: 链接函数, 单调可微函数.

例 $g(\cdot) = \ln(\cdot)$ $\ln y = w^T x + b$

二. 二分类任务:

对于二分类任务, 输出标记为 $y \in \{0, 1\}$, 即我们更倾向于选择介于0和1之间的概率, 而线性回归的预测结果 $z = w^T x + b$ 一般是一个连续值, 因此, 我们需要将 $z = w^T x + b$ 做某种变换, 将其转换到输出介于0和1之间的值(类似于概率范围).

$\left\{ \begin{array}{l} \text{线性回归模型实际预测值 } z = w^T x + b \\ y \in \{0, 1\} \end{array} \right\} \xrightarrow{\text{链接函数}} z$

1. 最理想函数: 单位阶跃函数

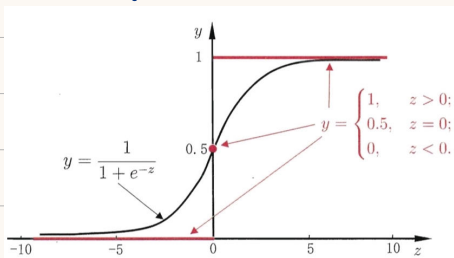
$$y = \begin{cases} 0, & z < 0 \\ 0.5, & z = 0 \\ 1, & z > 0 \end{cases}$$

/

缺点: 不连续, $x=0$ 不可导.

2. 对数几率函数 logistic function (sigmoid).

① 表示形式 $y = \frac{1}{1+e^{-z}} \Rightarrow$ 将 z 变为 $(0,1)$ 的 y



对数函数定义域 $(-\infty, +\infty)$, 值域 $(0, 1)$

② 几率(odds)定义: 指该事件发生的概率与
该事件不发生概率的比值.

$$Odds = \frac{p}{1-p} \quad (p \text{ 为事件发生概率})$$

$\left\{ \begin{array}{l} odds > 1, \text{事件发生概率增大} \\ odds < 1, \text{事件发生概率减小} \end{array} \right.$
 变化范围为 $(0, +\infty)$

变化范围为 $[0, +\infty)$

③ 对数 $\frac{1}{p}$ (log odds vs log-its) : $\log \frac{p}{1-p}$

$$\begin{cases} \log \frac{P}{1-P} > 0, \text{ 事件发生概率大于 } \\ \log \frac{P}{1-P} < 0, \text{ 事件发生概率小于 } \end{cases}$$

3. 对数归并

日期: /

$$y = g^{-1}(w^T x + b) \rightarrow g^{-1} = \frac{1}{1 + e^{-z}}$$

① 表现形式: $\ln \frac{y}{1-y} = w^T x + b$ 或 $y = \frac{1}{1 + e^{-(w^T x + b)}}$

其中 y 表示样本 x 作为正例的可能性。

② 本质: 是一种分类学习方法

③ 优点:

{ 无需事先假设数据分布
可得到“类别”的近似概率预测
可直接应用现有数值优化算法求得最优解

④. 概率解释:

y 表示样本 x 为正例可能性。(后验概率)

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

$$\Rightarrow \begin{cases} p(y=1|x) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}} = h_{\beta}(x) \\ p(y=0|x) = \frac{1}{1 + e^{w^T x + b}} = 1 - h_{\beta}(x) \end{cases} \quad \left. \begin{array}{l} \text{其中} \\ \beta = (w, b) \\ = [w, b]^T \end{array} \right\}$$

⑤ 最大似然法:

对于 $p(x|D)$

若 x 为变量, D 已知, $\therefore p(x|D)$ 为概率函数 [概率分布参数为 D 时, 不同样本 x 出现概率]。

日期: /

2. 当 x 已知, y 变量, $p(x|y)$ 为似然函数 [即概率分布的参数取不同值时, 某个样本 x 出现的概率], 在已有观测样本, 寻找数据分布的**最佳参数**

3. 求对数几率回归: 构造似然函数 + (牛顿法)

□ x_1, x_2, \dots, x_m 为独立同分布的采样, 定义**似然函数** L 为混合密度函数 (m 个样本同时出现):

$$\begin{aligned} L(\beta) &= p(y_1|x_1; \beta) \times p(y_2|x_2; \beta) \times \dots \times p(y_m|x_m; \beta) \\ &= \prod_{i=1}^m p(y_i|x_i; \beta) = \prod_{i=1}^m (h_{\beta}(x_i))^{y_i} (1 - h_{\beta}(x_i))^{1-y_i} \end{aligned}$$

思想: 找到一组参数使所有观测样本联合概率**最大化**.

$p(y_i|x_i; \beta)$: 在 β 条件下, x_i 预测正确的概率

$$\begin{aligned} \text{正确分类概率 } p(y_i|x_i; \beta) &= \begin{cases} p(y=1|x_i; \beta), & \text{if } y_i = 1, \\ p(y=0|x_i; \beta) = 1 - p(y=1|x_i; \beta) & \text{if } y_i = 0, \end{cases} \\ &= (p(y=1|x_i; \beta))^{y_i} \times (1 - p(y=1|x_i; \beta))^{1-y_i}, \\ \ln p(y_i|x_i; \beta) &= \begin{cases} \ln p(y=1|x_i; \beta), & \text{if } y_i = 1, \\ \ln p(y=0|x_i; \beta) = \ln[1 - p(y=1|x_i; \beta)] & \text{if } y_i = 0, \end{cases} \\ &= y_i \ln p(y=1|x_i; \beta) + (1 - y_i) \ln p(y=0|x_i; \beta) \end{aligned}$$

其中 $\beta = (w; b)$.

对数似然函数

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

概率中分号表示分号后的 w, b 是待估参数, 它们是确定的, 只是当前未知。它们不是随机变量。

日期: /

$$\ln L(\beta) = \sum_{i=1}^m \left[y_i \ln h_{\beta}(\mathbf{x}_i) + (1 - y_i) \ln(1 - h_{\beta}(\mathbf{x}_i)) \right]$$
$$\max_{\beta} \ln L(\beta) \Leftrightarrow \min_{\beta} -\ln L(\beta) \quad (\text{最大化转化为求最小化})$$

$$\begin{aligned} \ell(\beta) &= -\ln L(\beta) \\ &= \sum_{i=1}^m \left[y_i \ln h_{\beta}(\mathbf{x}_i) + (1 - y_i) \ln(1 - h_{\beta}(\mathbf{x}_i)) \right] \\ &= \sum_{i=1}^m \left[-y_i \beta^T \hat{\mathbf{x}}_i + \ln(1 + e^{\beta^T \hat{\mathbf{x}}_i}) \right] \quad \text{动手推导} \end{aligned}$$

□ 求解

$$\beta^* = \arg \min_{\beta} \ell(\beta)$$

□ 牛顿法第 $t+1$ 轮迭代解的更新公式

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

其中关于 β 的一阶、二阶导数分别为

$$\frac{\partial \ell(\beta)}{\partial \beta} = - \sum_{i=1}^m \hat{\mathbf{x}}_i (y_i - p_1(\hat{\mathbf{x}}_i; \beta)) \quad \text{动手推导}$$

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^m \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^T p_1(\hat{\mathbf{x}}_i; \beta) (1 - p_1(\hat{\mathbf{x}}_i; \beta)) \quad \text{动手推导}$$

对数几率回归-梯度下降法求解

□ 关于参数 β 的更新公式

$$\beta^{t+1} = \beta^t - \alpha \frac{\partial \ell(\beta)}{\partial \beta}$$

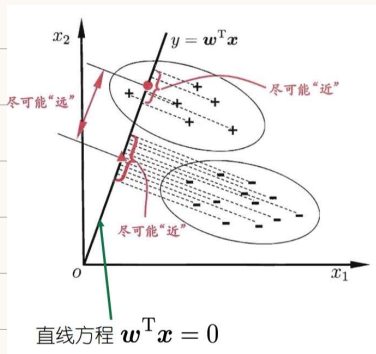
其中 α 是学习率. 关于 β 的一阶偏导数为

$$\frac{\partial \ell(\beta)}{\partial \beta} = - \sum_{i=1}^m \hat{\mathbf{x}}_i (y_i - p_1(\hat{\mathbf{x}}_i; \beta))$$

日期: /

4. 线性判别分析 LDA

① 基本思想: 给定的训练集, 设法将样例投影到一维适当选择的直线上, 使得同类样例的投影点尽可能接近, 异类样例投影点中心尽可能远离。[类内方差小, 类间方差大]



目标:

① 在保留尽可能多的类间信息的前提下

② 寻找最优投影向量 w , x 在 w 上投影为

$$y = w^T x$$

④ 目标函数

数据: $m \times d$ 维样本 x_1, \dots, x_m , 属于不同的类
 x_i 中两个不同类别 w_0, w_1 即

$$\begin{cases} X_0 = \{x_i \mid w(x_i) = w_0\} & |X_0| = n_1 \\ X_1 = \{x_i \mid w(x_i) = w_1\} & |X_1| = n_2 \end{cases}$$

日期: /

线性组合成 (x_1, x_2, \dots, x_m) 产生 n 个投影结果

$[y_1, \dots, y_m]$ 相应属于类 Y_0, Y_1 , 即 $y_i = w^T x_i$

投影前样本均值

$$\bar{x}_i = \frac{1}{n_i} \sum_{x \in X_i} x \quad i=0,1$$

投影后样本均值

$$\begin{aligned} \bar{y}_i &= \frac{1}{n_i} \sum_{y \in Y_i} y = \frac{1}{n_i} \sum_{x \in X_i} w^T x \\ &= w^T \bar{x}_i \end{aligned}$$

← 将原均值投影

定义: (2个指标) 投影后样本均值之差, 投影后第 i 类的类内散度 (类似于方差)

{ 投影后点的样本均值之差 $|\bar{y}_0 - \bar{y}_1| = |w^T(\bar{x}_0 - \bar{x}_1)|$
投影后第 i 类类内散度: $\tilde{s}_i^2 = \sum_{y \in Y_i} (y - \bar{y}_i)^2$

Fisher 线性判别准则要求

$$\max J(w) = \frac{|\bar{y}_0 - \bar{y}_1|^2}{\tilde{s}_0^2 + \tilde{s}_1^2} \quad \begin{array}{l} \rightarrow \text{分子可能大} \\ \rightarrow \text{分母可能小} \end{array}$$

定义类间散度矩阵 $S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T$, 可证明其为对称半正定的. 投影后两类样本均值之差展开为

$$|\bar{\mu}_0 - \bar{\mu}_1|^2 = w^T (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T w = w^T S_b w.$$

投影后第 i 内的类内散度为:

$$\begin{aligned} \tilde{s}_i^2 &= \sum_{y \in Y_i} (y - \bar{\mu}_i)^2 \\ &= \sum_{y \in Y_i} (w^T x - w^T \mu_i)^2 \\ &= \sum_{x \in X_i} w^T (x - \mu_i)(x - \mu_i)^T w = w^T \Sigma_i w. \end{aligned}$$

$$\tilde{s}_0^2 + \tilde{s}_1^2 = w^T S_w w.$$

$$J(w) = \frac{w^T S_b w}{w^T S_w w}$$