

线性回归

① 例子: (a) 股价涨跌幅自回归 (b) 房价预测 $y = a_1x_1 + a_2x_2 + c$

② 原理: 用极大似然解释最小二乘

线性回归建模: $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$ (a) 样本 i 的似然函数

↑ 模型参数 ↑ 特征 ↑ 误差值

样本 i 的似然函数
 (b) 误差值 $\epsilon^{(i)}$ 的概率密度

$$P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\epsilon^{(i)})^2}{2\sigma^2}}$$

(a) 误差值 $\{\epsilon^{(1)}, \epsilon^{(2)}, \dots, \epsilon^{(m)}\}$ 服从某均值为 0 方差为 σ^2 的高斯分布。根据大数定理 (概率 (23)(24))

※ 假定 M 样本独立, 则 $\{\epsilon^{(i)}\}$ 独立分布 (不会相互影响)

※ 假定 m 样本同分布, 也合理因为某些特征未考虑

※ 均值 $\bar{\epsilon} = \sum \epsilon^{(i)} / m = 0$, 如果不为 0 可通过平移使其为 0

(c) 样本 i 似然函数值的概率密度

$$P(y^{(i)} | x^{(i)}, \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

误差值的高斯分布以 0 为中心, 似然函数以 $y^{(i)}$ (样本标签) 为中心, 方差都是 σ^2

(d) 所有样本 i 的似然函数

$$L(\theta) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

对数似然 (方便计算, 乘法变加法)

$$\ln L(\theta) = \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}} = m \ln \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

希望最大 (为什么希望 $\ln L(\theta)$ 最大?)

② 另种直观理解: 似然最大 \Leftrightarrow 相对错 (与特征) 最少

高斯分布中心点的概率密度最大, 而在这个中心点预估值与样本 y 值相同, 误差为 0

优化目标 $J(\theta) = \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$

③ 假设的内涵性, 简化性, 发散性

内涵性: 根据常理应该是正确的, 往往是正确的但不一定总是正确

简化性: 只是接近真实, 往往做了简化 (在词袋模型)

发散性: 在某个简化假设下推导出的结论, 假设不成立时有时也能用

④ 最小二乘求解:

$$\begin{aligned}\theta^* &= \arg \min_{\theta} J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2} (X\theta - Y)^T (X\theta - Y) \\ &= \frac{1}{2} (\theta^T X^T - Y^T) (X\theta - Y) = \frac{1}{2} (\theta^T X^T X \theta - \theta^T X^T Y - Y^T X \theta + Y^T Y)\end{aligned}$$

求偏导:

$$\begin{aligned}\frac{\partial J(\theta)}{\partial (\theta)} &= \frac{1}{2} \cdot \frac{\partial (\theta^T X^T X \theta - \theta^T X^T Y - Y^T X \theta + Y^T Y)}{\partial (\theta)} \\ &= \frac{1}{2} (\underbrace{2X^T X \theta}_{\text{线性(2)}} - \underbrace{X^T Y}_{\text{线性(2)}} - \underbrace{(Y^T X)^T}_{\text{线性(2)}} + 0) \\ &= X^T X \theta - X^T Y\end{aligned}$$

令 $X^T X \theta - X^T Y = 0$ 求驻点, 然而只有 $X^T X$ 可逆, 才有解.

为防止 $X^T X$ 不可逆, 同时防止过拟合, 在主对角线增加扰动因子 λI

$$(X^T X + \lambda I) \theta - X^T Y = 0$$

防止 θ_i 变得过大

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

一定有解 / $X^T X$ 半正定: 线性(6)
恒有 $\vec{v}^T (\lambda I) \vec{v} > 0$, $(X^T X + \lambda I)$ 正定

⑤ 正则项与防止过拟合: 与④使用的 λI 不同, 实际应用中, 使用下面三种之一

$$\begin{cases} \text{L1-norm (Lasso)} : J(\vec{\theta}) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j| \\ \text{L2-norm (Ridge)} : J(\vec{\theta}) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \\ \text{Elastic Net} : J(\vec{\theta}) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \left(\rho \cdot \sum_{j=1}^n |\theta_j| + (1-\rho) \sum_{j=1}^n \theta_j^2 \right) \end{cases}$$

|| 本质为假定 θ 服从高斯分布

Lasso 具有特征选择能力, 原理:

希望 $J(\theta)$ 极值, 同时希望 $\lambda \sum_{j=1}^n |\theta_j|$ 比较小, 即 $\sum_{j=1}^n |\theta_j| < r$, 转化为约束条件求极值 (拉格朗日) 问题, 这个约束条件使得 $|\theta_j|$ 不能太大, 当特征数 (阶数) 增多时, 一些与优化目标关系较弱特征对应的 $|\theta_j|$ 会非常接近零

⑥ 超参数 λ 的选择: 用 N-Fold 交叉验证做实验, 看哪个 λ 取值下训练出的模型 MSE (均方误差) 最小.

⑦ 伪逆矩阵：前面线性回归使用 λ 扰动因子保证公式有解，如果不用 λ 来解线性方程组 $A\vec{x} = \vec{B}$ ，另一个可选方法是使用伪逆矩阵。

| A 可逆时: $\vec{x} = A^{-1}\vec{B}$

| A 不可逆: $\vec{x} = A^+ \vec{B} = (A^T A)^{-1} A^T \vec{B}$ 其中 $A^+ = (A^T A)^{-1} A^T$ 叫做 A 的伪逆

A 可逆时: $A^+ = (A^T A)^{-1} A^T = A^{-1} (A^T)^{-1} A^T = A^{-1}$

如何得到 A^+ : ① 奇异值分解 $A_{m \times n} = U \Sigma V^T$ ② 直接令 $A^+ = V \Sigma^{-1} U^T$

⑧ L1-Norm 梯度下降

⑨ 批量梯度下降：把所有样本梯度加到一起然后下降。（梯度一定下降）

⑩ 随机梯度下降：拿到一个样本下降一次（梯度震荡下降）

优势：① 快 ② 不用等到所有样本，适合在线学习 ③ α 步长较小时，利用噪声样本跳出“洼地”

⑪ Mini-batch 梯度下降：若干样本下降一轮（折衷方法，实际中使用较多）

⑫ 下降步长 α ：

<a> 固定步长

 自适应 α ：参考斜率大小，参考一定步长时 $loss$ 是上升还是下降

<c> 回溯线性搜索

⑬ 样本与预估目标

预估目标与样本特征不是线性关系时也没关系（与是线性关系无关的）

事先对特征做变换，例如 $x \Rightarrow x^2$ ，不再仅使用一次函数，

用二次，三次函数来拟合

但次数过高的缺点是过拟合。

$$y = \theta_{11}x_1 + \theta_{12}x_1^2 + \theta_{13}x_1^3 + \theta_{21}x_2 + \theta_{22}x_2^2 + \dots$$

⑭ 过拟合解决：

<a> 正则项 删减无关特征，例如 PCA，one-hot 编码，事半功倍，基于一业务逻辑

⑮ 超参数判定系数

方法1: 均方误差 (MSE) $\sum_{i=1}^m (h_0(x^{(i)}), x^{(i)}) - y^{(i)})^2$ 其实就是残差平方和 RSS
 似然函数 预测的 $\hat{y}^{(i)}$ 样本真实值 $y^{(i)}$

方法2: 判定系数 (R^2)

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^m (\hat{y}_i - y_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2}$$

残差平方和 (Residual Sum of Square): 预测偏差.
 样本方差 (Total Sum of Square): 体现样本本身的发散程度

原理: 方法1 残差平方和太大有可能是因为样本本身震荡造成的

所以用 TSS ($m \times$ 样本方差) 来作为基准.

R² 值越大, 说明训练的越准确, 越小说明误差越大, 正常位于 [0, 1],

为负值时说明误差极大

方法3: 判定系数 (ESS), 又叫回归平方和.

$$ESS = \sum_{i=1}^m (\hat{y}_i - \bar{y})^2$$

预测值 样本均值

~~TSS = ESS + RSS~~

$$TSS \geq ESS + RSS$$

仅当无偏估计时才成立.

样本方差 (定值) 预测值 v.s 样本均值 残差平方和 (预测值 v.s 样本真实值)

n fold
va
bias

Logistic 回归 / Softmax 回归

① Logistic 回归用于 2 分类; Softmax 回归用于 N 分类: 两者都是从信息熵推导出来的分类器; 不建议用线性回归做分类, 如 PPT 中的例子

② LR 建模: 不像线性回归直接用 $\theta^T x$ 来对样本回归建模, 而是把 $\theta^T x$ 放在 sigmoid 函数中

Sigmoid 函数: $g(z) = \frac{1}{1+e^z}$

特性①



特性② $g'(z) = \dots = g(z) \cdot (1 - g(z))$

把 $z = \theta^T x$ 代入到 $g(z) = \frac{1}{1+e^z}$ 中, 得到建模用的函数.
模型参数 特征

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$$

③ LR 参数估计 (即 LR 参数求解)

假定 $P(y=1|x; \theta) = h_{\theta}(x)$ ① $P(y=0|x; \theta) = 1 - h_{\theta}(x)$ ② 推出
 则用一个公式统一表示 $P(y|x; \theta) = \frac{(h_{\theta}(x))^y (1-h_{\theta}(x))^{1-y}}$

扩展到多个样本, 得到似然函数

★ 也与二项分布一致

$$L(\vec{\theta}) = P(\vec{y} | X; \theta) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta) = \prod_{i=1}^m (h_{\theta}(x^{(i)})^{y^{(i)}} (1-h_{\theta}(x^{(i)})^{1-y^{(i)}})$$

求对数, 得到对数似然函数.

$$l(\vec{\theta}) = \ln(L(\vec{\theta})) = \sum_{i=1}^m (y^{(i)} \ln(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \ln(1-h_{\theta}(x^{(i)})))$$

关于 θ_j 求偏导,

$$\frac{\partial l(\vec{\theta})}{\partial \theta_j} = \dots = \sum_{i=1}^m (y^{(i)} - g(\theta^T x^{(i)})) \cdot x_j^{(i)}$$

偏导即为梯度方向, 沿梯度方向更新 θ_j (梯度上升, 让似然函数取最大值)

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

期望

学习

★ 似然函数取最大值 $\Rightarrow h_{\theta}(x^{(i)})^{y^{(i)}} (1-h_{\theta}(x^{(i)}))^{1-y^{(i)}}$ 是 0, 还是 1, h_{θ} 与 y 一致才能取最大值

证明

$y^{(i)}$ 取值 (0, 1) 决定 $h_{\theta}(x^{(i)})$ 该是 0, 还是 1, h_{θ} 与 y 一致才能取最大值

④ 为什么 LR 预估值可看作概率

$$\text{发生概率 } p = p(y=1 | \vec{x}; \theta) = h_{\theta}(\vec{x}) = g(\theta^T \vec{x}) = \frac{1}{1+e^{-\theta^T \vec{x}}}$$

$$\text{不发生概率 } 1-p = p(y=0 | \vec{x}; \theta) = 1 - h_{\theta}(\vec{x}) = 1 - g(\theta^T \vec{x}) = \frac{e^{-\theta^T \vec{x}}}{1+e^{-\theta^T \vec{x}}}$$

用对数概率公式来表示 (logit 函数)

$$\text{logit}(p) = \log \frac{p}{1-p} = \log \frac{h_{\theta}(x)}{1-h_{\theta}(x)} \xrightarrow{\text{代入}} \log \frac{1}{e^{-\theta^T x}} = \theta^T x$$

到这里可以看出，作者是想设计一个函数 $g(z)$ ，使得 $\log \frac{p}{1-p}$ 是线性的 线性函数

$$\text{即 } \log \frac{p}{1-p} = \theta^T x, \text{ 用这个等式可推出 } \begin{cases} g(z) = \frac{1}{1+e^z} \\ z = \theta^T x \end{cases}$$

（见本节③）

⑤ LR 损失函数：用负的对数似然函数作为损失函数

（似然函数取最大值时，损失函数取最小值）

对比线性回归，把 MSE (均方误差) 作为损失函数。

为了让损失函数好看些，用 $y=1$ 表示正例， $y=-1$ 表示负例，修改

$$\text{似然函数改写: } L(\theta) = \prod_{i=1}^m p_i^{y_i} (1-p_i)^{1-y_i} \text{ 改为 } L(\theta) = \prod_{i=1}^m p_i^{(y_i+1)/2} (1-p_i)^{-(y_i-1)/2}$$

$$\text{对数似然改写: } l(\theta) = \sum_{i=1}^m \ln(p_i^{y_i} (1-p_i)^{1-y_i}) \text{ 改为 } l(\theta) = \sum_{i=1}^m \ln(p_i^{(y_i+1)/2} (1-p_i)^{-(y_i-1)/2})$$

损失函数

$$\text{loss}(y, \hat{y}_i) = -l(\theta) = -\sum_{i=1}^m \left[\frac{1}{2}(y_i+1) \ln p_i - \frac{1}{2}(y_i-1) \ln(1-p_i) \right]$$

$$= \sum_{i=1}^m \left[\frac{1}{2}(y_i+1) \ln \frac{1}{p_i} - \frac{1}{2}(y_i-1) \ln \frac{1}{1-p_i} \right] \quad p_i = \frac{1}{1+e^{-\theta^T x}}$$

$$= \sum_{i=1}^m \left[\frac{1}{2}(y_i+1) \ln(1+e^{-\theta^T x}) - \frac{1}{2}(y_i-1) \ln(1+e^{\theta^T x}) \right] \quad 1-p_i = \frac{e^{\theta^T x}}{1+e^{-\theta^T x}}$$

$$= \begin{cases} \sum_{i=1}^m [\ln(1+e^{-\theta^T x})] & y_i = 1 \\ \sum_{i=1}^m [\ln(1+e^{\theta^T x})] & y_i = -1 \end{cases}$$

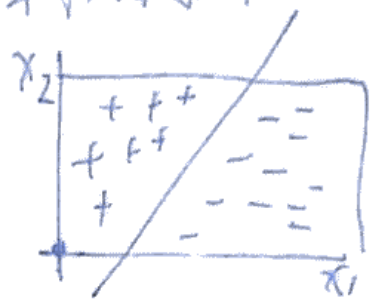
$$\Rightarrow \text{loss}(y, \hat{y}_i) = \sum_{i=1}^m [\ln(1+e^{-y_i \cdot \theta^T x})]$$

备注：也有用 MSE 做损失函数，不是不可以，只不过能用 LR 自带损失函数的，尽量用 LR 的

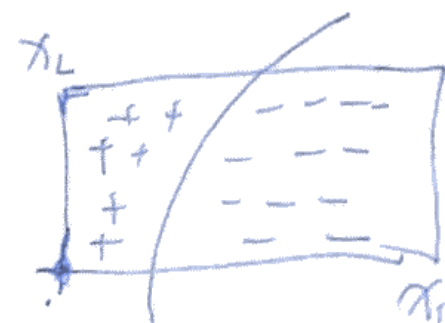
⑥ 梯度上升：沿似然函数正梯度上升，本质上与沿损失函数负梯度下降相同 → 见③：还等价于 $\hat{y}_{(i)}$ 与 $y_{(i)}$ 相对误差最小

⑦ 梯度提升：

提升前： $Z(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ 分类界面



提升后： $Z(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1 x_2 + \theta_5 x_2^2$ 分类界面



⑧ 线性回归，LR \iff 广义线性模型 \iff 指数族

⑨ Softmax 回归：K 分类，第 k 类的参数为 $\vec{\theta}_k$ ，组成二维矩阵 $\Theta_{K \times n}$

概率密度：
$$p(c=k|x;\theta) = \frac{e^{\theta_k^T x}}{\sum_{i=1}^K e^{\theta_i^T x}}$$

原因见⑩

* Soft-max 并不是直接用 $\frac{P_i}{P_1 + \dots + P_K}$ 来建模，而是用 $\frac{e^{P_i}}{e^{P_1} + \dots + e^{P_K}}$ 来近似求最大值。

假设样本间独立

概率就可以相乘，得到似然函数

似然函数：
$$L(\theta) = \prod_{i=1}^m \prod_{k=1}^K p(c=k|x^{(i)}, \theta)^{y_k^{(i)}} = \prod_{i=1}^m \prod_{k=1}^K \left[\frac{e^{\theta_k^T x}}{\sum_{l=1}^K e^{\theta_l^T x}} \right]^{y_k^{(i)}}$$

对数似然：
$$J(\theta) = \ln L(\theta) = \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} (\theta_k^T x - \ln \sum_{l=1}^K e^{\theta_l^T x})$$

$\theta_k = \begin{bmatrix} \theta_k^{(1)} \\ \theta_k^{(2)} \\ \vdots \\ \theta_k^{(n)} \end{bmatrix}$
 $y_k^{(i)} = k$ 时

$\theta_i = \begin{bmatrix} \theta_1^{(i)} \\ \theta_2^{(i)} \\ \vdots \\ \theta_n^{(i)} \end{bmatrix}$ $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

$y_0^{(i)} = 0, y_1^{(i)} = 0, \dots, y_k^{(i)} = 1, \dots, y_K^{(i)} = 0$
 one-hot 编码

特征用于类别 i 的模型参数

梯度：
$$\frac{\partial J(\theta)}{\partial \theta_k} = (y_k - p(y_k|x;\theta)) \cdot x$$

沿着这个梯度对似然概率做梯度上升

⑩ 用 $\frac{e^{p_i}}{e^{p_1} + e^{p_2} + \dots + e^{p_k}}$ 替代 $\frac{p_i}{p_1 + p_2 + \dots + p_k}$ 做为 Soft-Max 建模时概率密度函数

$$3 < 5 \Leftrightarrow e^3 < e^5$$

$$a = \ln \frac{e^3}{e^3 + e^5} \sim \frac{3}{3+5}$$

$$\frac{3}{3+5} < \frac{5}{3+5}$$

$$\frac{e^3}{e^3 + e^5} < \frac{e^5}{e^3 + e^5}$$

$$b = \ln \frac{e^5}{e^3 + e^5} \sim \frac{5}{3+5}$$

⑪ Soft-max 分类数 $k=2$ 时, 会退化成 logistic 回归.

$k=2$ 时, 令 $\theta = \theta_1 - \theta_2$.

$$p_1 = \frac{e^{\theta_1^T x}}{e^{\theta_1^T x} + e^{\theta_2^T x}} = \frac{1}{1 + e^{-\theta_1^T x}} = \frac{1}{1 + e^{-\theta^T x}}$$

这就是
LR 用的
sigmoid 函数, 记
 $z = \theta^T x$

⑫ 为条件 T 书极值问题 (见下页)

⑫ 约束条件下极值问题：拉格朗日乘子法

★ 题目例子：投骰子，6个点等概率 $P_1 = P_2 = \dots = P_6$ ，N次投置后均值为2.71828。

下一次投置出现点5的概率

建模：在给定的约束条件下，使信息熵取最大（~~这~~等价于在装箱问题中，不违反约束条件的前提下，装到几个箱子的概率尽可能趋于相同）

熵函数： $H(\vec{p}) = -\sum_{i=1}^6 p_i \ln p_i$ → 概率④ (信息熵，期望它能取到最大值)

约束条件： $\sum_{i=1}^6 p_i = 1$ 即 $(1 - \sum_{i=1}^6 p_i) = 0$

$\sum_{i=1}^6 p_i \cdot i = 2.71828$ 即 $(2.71828 - \sum_{i=1}^6 i \cdot p_i) = 0$

Lagrange函数： $L(\vec{p}, \lambda_1, \lambda_2) = -\sum_{i=1}^6 p_i \ln p_i + \lambda_1 (1 - \sum_{i=1}^6 p_i) + \lambda_2 (2.71828 - \sum_{i=1}^6 i \cdot p_i)$

对 p_i 求偏导，令偏导为0， $\frac{\partial L}{\partial p_i} = -\ln p_i - 1 + \lambda_1 - i \cdot \lambda_2 = 0$ ④

$$\Rightarrow \begin{cases} p_i = e^{-1-\lambda_1-i\lambda_2} & \text{⑤} \\ \lambda_1 = -0.0787 & \text{⑥} \\ \lambda_2 = 0.2808 & \text{⑦} \end{cases}$$

②：由公式④得到

⑥、⑦：把 $p_i = e^{-1-\lambda_1-i\lambda_2}$ 代入到 $J(\lambda) = (\sum_{i=1}^6 p_i - 1) + (\sum_{i=1}^6 i \cdot p_i - E)^2$

对 λ_1, λ_2 求偏导，得到梯度，

沿梯度下降，求出 λ_1, λ_2 代入⑤解出 p_i 。

如何理解最大熵模型？：见 →