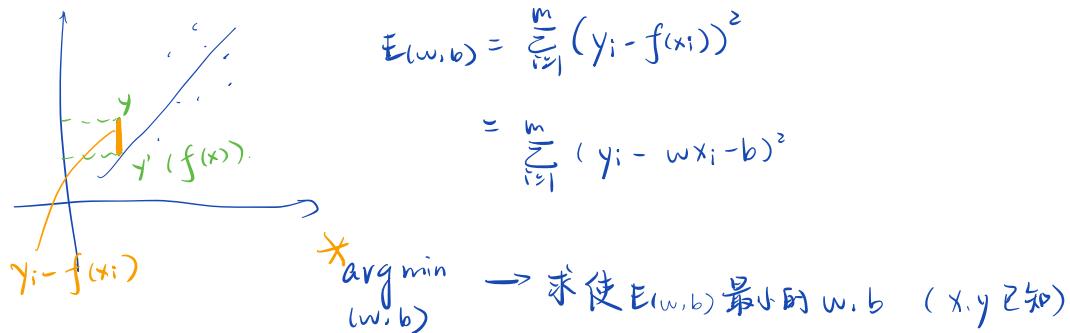


第三章 线性模型.

一. 线性回归

$$\left\{ \begin{array}{l} \text{连续值 } f(x) = w_1 x + b \\ \text{二值离散值 } (A: 1, B: 0) \\ \text{多值 } \begin{cases} \text{有序 } (A: 1, 2, \dots) & f(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b \\ \text{无序 } (1, 0, 0) \quad (0, 1, 0) \quad \dots & f(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + \\ & w_4 x_4 + w_5 x_5 + w_6 x_6 + b. \end{cases} \end{array} \right.$$

1. 最小二乘估计：基于均方误差最小化 模型求解。



2. 极大似然估计。先确定分布类型

估计概率分布的参数值。再求参数 → 极大似然。

连续型。

离散(连续型)随机变量 X . 概率质量函数 $P(x; \theta)$ (概率密度函数 $P(x; \theta)$).

θ — 待估计参数(可有多个)。

x_1, x_2, \dots, x_n 来自 X 的 n 个独立同分布 样本。

联合概率 $L(\theta) = \prod_{i=1}^n P(x_i; \theta)$ —— 样本的似然函数

已知量 不知量

使 $L(\theta)$ 最大的 θ^* 即 θ 的估计值

→ 逐项求对数函数化简 $\ln L(w, \sigma^2)$ 和 $L(w, \sigma^2)$ 最大值点相同

$$\ln \prod_i \Rightarrow \sum_i \ln$$

应用：

线性回归 : $y = wx + b + \epsilon$ → 不受控制的随机误差.

通常假设 $\epsilon \sim N(0, \sigma^2)$ 中心极限定理诠释.

$$\Rightarrow p(\epsilon) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right)$$

$$\text{代入 } \epsilon = y - wx - b$$

$$\Rightarrow p(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - (wx + b))^2}{2\sigma^2}\right)$$

利用极大似然估计. $L(w, b) = \prod_{i=1}^m p(y_i)$

$$\ln L(w, b) = \sum_{i=1}^m \ln \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - (wx + b))^2}{2\sigma^2}\right) \right]$$

$$= \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi}\sigma} + \sum_{i=1}^m \left(-\frac{(y - (wx + b))^2}{2\sigma^2} \right)$$

$$= m \ln \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y - (wx + b))^2$$

$$m. \alpha \text{ 为常数} \quad \text{故 } (\hat{w}, \hat{b}) = \underset{(w, b)}{\operatorname{argmax}} \ln L(w, b)$$

$$= \underset{(w, b)}{\operatorname{argmin}} \sum_{i=1}^m (y_i - w x_i - b)^2$$

公式 3.4.

求解 $w, b \Rightarrow$ 多元函数求最值点. $E = \sum_{i=1}^m (y_i - w x_i - b)^2$



凸集: 集合 $D \subset \mathbb{R}^n$ $\forall x, y \in D$, 与 $\alpha \in [0, 1]$ 有

$$\alpha x + (1-\alpha)y \in D \Rightarrow D \text{ 为凸集.}$$

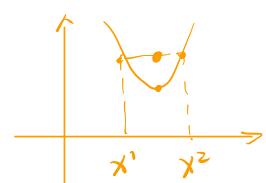
几何意义: 集合内任两点连线上任意点均 $\in D$

e.g. \emptyset , n 维欧式空间 \mathbb{R}^n .

凸函数: 设 D 非空凸集, f 定义在 D 上的函数. $\forall x^1, x^2 \in D$, $\alpha \in (0, 1)$.

$$\text{均有 } f(\alpha x^1 + (1-\alpha)x^2) \leq \alpha f(x^1) + (1-\alpha)f(x^2)$$

则 f 为 D 上的凸函数.



梯度(多元函数的一阶导数)

(列向量)

n 元函数 $f(x)$ 对自变量 $x = (x_1, \dots, x_n)^T$ 的各分量 i 的偏导数都存在.

则 $f(x)$ 在 x 处一阶可导

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \Rightarrow f(\mathbf{x}) \text{ 在 } \mathbf{x} \text{ 处的梯度}$$

Hessian (海塞) 矩阵 (多元 --- 二阶导数).

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \vdots & \ddots & & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \cdots & \ddots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix}$$

定理：设 $D \subset \mathbb{R}^n$ 非空开凸集， $f: D \subset \mathbb{R}^n \rightarrow \mathbb{R}$ ，且 $f(\mathbf{x})$ 在 D 上二阶 ...

若 $f(\mathbf{x})$ 的 Hessian 矩阵在 D 上半正定，则 $f(\mathbf{x})$ 是 D 上的凸函数.

(类似一元 --- 判断凹凸性).

证明 $E(w, b)$ 是关于 w 和 b 的凸函数 $\Rightarrow \nabla^2 E(w, b)$ 正定.

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

$$= \sum_{i=1}^m (-x_i) \cdot (y_i - w x_i - b)$$

$$= \sum_{i=1}^m \left[w \sum_{j=1}^m x_j^2 - \sum_{j=1}^m (y_j - b) x_j \right] \quad 3.5$$

$$\frac{\partial^2 E(w, b)}{\partial w^2} = 2 \frac{\partial}{\partial w} \left[w \sum_{i=1}^m x_i^2 - \underbrace{\sum_{i=1}^m (y_i - b) x_i}_{\text{常数项}} \right]$$

$$= \sum_{i=1}^m x_i^2$$

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

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

$$\frac{\partial^2 E(w, b)}{\partial b^2} = 2m$$

$$\frac{\partial^2 E(w, b)}{\partial w \partial b} = \frac{\partial}{\partial b} \left(\frac{\partial E(w, b)}{\partial w} \right)$$

$$= 2 \sum_{i=1}^m x_i$$

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

$$\therefore \nabla^2 E(w, b) = \begin{bmatrix} \sum_{i=1}^m x_i^2 & \sum_{i=1}^m x_i \\ \sum_{i=1}^m x_i & 2m \end{bmatrix}$$

半正定矩阵： 实对称矩阵 所有顺序主子式均为非负.

$$\left| \sum_{i=1}^m x_i^2 \right| > 0$$

$$\begin{vmatrix} \sum_{i=1}^m x_i^2 & \sum_{i=1}^m x_i \\ \sum_{i=1}^m x_i & 2m \end{vmatrix} = 4m \sum_{i=1}^m x_i^2 - 4 \left(\sum_{i=1}^m x_i \right)^2$$

$$= 4m \sum_{i=1}^m x_i^2 - 4 \cdot m \cdot \frac{1}{m} \left(\sum_{i=1}^m x_i \right)^2$$

$$= 4m \left(\sum_{i=1}^m x_i^2 - \sum_{i=1}^m x_i \bar{x} \right)$$

$$= 4m \sum_{i=1}^m (x_i^2 - 2x_i \bar{x} + x_i \bar{x})$$

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

$$= 4m \sum_{i=1}^m (x_i - \bar{x})^2 \geq 0.$$

$\therefore \nabla^2 E(w, b)$ 是半定矩阵. $\Rightarrow E(w, b)$ 凸函数.

凸优化定理(凸函数性质): $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 凸函数

x^* 全局极 $\Leftrightarrow \nabla f(x^*) = 0$

$$\nabla E(w, b) = \begin{bmatrix} \frac{\partial E(w, b)}{\partial w} \\ \frac{\partial E(w, b)}{\partial b} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\left\{ \begin{array}{l} \geq \left[w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b) x_i \right] = 0 \\ \geq \left(mb - \sum_{i=1}^m (y_i - w x_i) \right) = 0 \\ \hookrightarrow b = \frac{1}{m} \sum_{i=1}^m (y_i - w x_i) \end{array} \right. \quad 3.8$$

$$\Rightarrow b = \frac{1}{m} \sum_{i=1}^m y_i - w \bar{x} = \bar{y} - w \bar{x}$$

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

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

$$w \sum_{i=1}^m x_i^2 = \sum_{i=1}^m [x_i y_i - (\bar{y} - w \bar{x}) x_i]$$

$$w \sum_{i=1}^m x_i^2 - w \bar{x} \sum_{i=1}^m x_i = \sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i$$

$$\Rightarrow w = \frac{\sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i}{\sum_{i=1}^m x_i^2 - \bar{x} \sum_{i=1}^m x_i}$$

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

$$\bar{x} \sum y_i = \frac{1}{m} (\sum x_i)^2$$

$$\therefore w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} (\sum_{i=1}^m x_i)^2} \quad 3.7$$

逐项转化矩阵 numpy.
w 向量化.

$$\bar{y} \sum_{i=1}^m x_i = \bar{x} \sum_{i=1}^m y_i = m \bar{x} \bar{y} = \frac{m}{m} \bar{x} \bar{y}$$

$$\sum_{i=1}^m x_i \bar{x} = m \bar{x}^2 = \sum_{i=1}^m \bar{x}^2$$

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

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

令 $x = (x_1, \dots, x_n)^T$, $x_d = (x_1 - \bar{x}, \dots, x_n - \bar{x})^T$

$y = (y_1, \dots, y_n)^T$, $y_d = (y_1 - \bar{y}, \dots, y_n - \bar{y})^T$.

$$\text{B)} \quad w = \frac{x_d^T y_d}{x_d^T x_d}$$

机器学习三要素：

- { 模型：确定假设空间.
- 策略：损失函数 $L(w, b)$
- 算法：求解损失函数，确定最优模型.

多元线性回归.

$$\omega \in \mathbb{R}^d$$

$$f(x_i) = \omega^T x_i + b$$

$$= (\omega_1 \ \omega_2 \ \dots \ \omega_d) \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix} + b$$

$$= \omega_1 x_{i1} + \omega_2 x_{i2} + \dots + \omega_d x_{id} + b$$

将 ω 扩充到 $\hat{\omega} \in \mathbb{R}^{d+1}$

b 替换为 $\omega_{d+1} \cdot 1$

$$\Rightarrow f(x_i) = \omega_1 x_{i1} + \omega_2 x_{i2} + \dots + \omega_d x_{id} + \omega_{d+1} \cdot 1$$

$$f(x_i) = (\omega_1 \ \omega_2 \ \dots \ \omega_d \ \omega_{d+1}) \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \\ 1 \end{pmatrix}$$

$$\Rightarrow f(\hat{x}_i) = \hat{\omega}^T \hat{x}_i$$

$$\text{由最小二乘法, } E\hat{\omega} = \sum_{i=1}^m (y_i - f(\hat{x}_i))^2$$

$$= \sum_{i=1}^m (y_i - \hat{\omega}^T \hat{x}_i)^2$$

↓ 向量化

$$E_{\hat{w}} = (y_1 - \hat{w}^T \hat{x}_1, y_2 - \hat{w}^T \hat{x}_2, \dots, y_m - \hat{w}^T \hat{x}_m) \begin{pmatrix} y_1 - \hat{w}^T \hat{x}_1 \\ \vdots \\ y_m - \hat{w}^T \hat{x}_m \end{pmatrix}$$

$$\begin{pmatrix} y_1 - \hat{w}^T \hat{x}_1 \\ \vdots \\ y_m - \hat{w}^T \hat{x}_m \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} - \begin{pmatrix} \hat{x}_1^T \hat{w} \\ \hat{x}_2^T \hat{w} \\ \vdots \\ \hat{x}_m^T \hat{w} \end{pmatrix}$$

$$\text{令 } y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}, \quad \begin{pmatrix} \hat{x}_1^T \hat{w} \\ \hat{x}_2^T \hat{w} \\ \vdots \\ \hat{x}_m^T \hat{w} \end{pmatrix} = \begin{pmatrix} \hat{x}_1^T \\ \hat{x}_2^T \\ \vdots \\ \hat{x}_m^T \end{pmatrix} \hat{w} = X \cdot \hat{w}$$

$$y - X \hat{w}$$

$$\therefore E_{\hat{w}} = (y - X \hat{w})^T (y - X \hat{w}). \quad 3.9$$

① 证明 $E_{\hat{w}}$ 为凸函数

求解 \hat{w}^* . → 多元函数最值.

→ Hessian 矩阵.

② 凸函数求最值求解 \hat{w} .

$$\begin{aligned} \frac{\partial E_{\hat{w}}}{\partial \hat{w}} &= \frac{\partial}{\partial \hat{w}} [(y - X \hat{w})^T (y - X \hat{w})] \\ &= \frac{\partial}{\partial \hat{w}} [(\underbrace{y^T - \hat{w}^T X^T}_{y^T} y - \hat{w}^T X^T y + \hat{w}^T X^T X \hat{w})] \\ &= \frac{\partial}{\partial \hat{w}} (\underbrace{y^T y - y^T X \hat{w} - \hat{w}^T X^T y + \hat{w}^T X^T X \hat{w}}) \end{aligned}$$

$$= \frac{\partial}{\partial \hat{w}} (-\underbrace{y^T x \hat{w}}_{\text{与 } \hat{w} \text{ 无关}} - \hat{w}^T x^T y + \hat{w}^T x^T x \hat{w})$$

标量函数

\hat{w} : $d+1$ 元的 向量

标量函数对向量求导

矩阵微分公式:

$x \in \mathbb{R}^{n+1}$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 关于 x 的实值标量函数

$$\frac{\partial f(x)}{\partial x} = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \quad \frac{\partial f(x)}{\partial x^T} = ()$$

分子布局

分母布局

$$\frac{\partial E_{\hat{w}}}{\partial \hat{w}} = -y^T x - \frac{\partial \hat{w}^T x^T y}{\partial \hat{w}} + \hat{w}^T x^T x \quad \frac{\partial x^T A x}{\partial x} = (A + A^T)x$$

$$= -y^T x - y^T x + (x^T x + x^T x) \hat{w}$$

$$= 2x^T y + 2x^T x \hat{w}$$

$$= 2x^T (x \hat{w} - y)$$

3.10.

$$\frac{\partial^2 E_{\hat{w}}}{\partial \hat{w}^2} = \frac{\partial}{\partial \hat{w}} (2x^T (x \hat{w} - y))$$

$$= 2x^T X$$

$$\Rightarrow \nabla^2 E_{\hat{w}} = 2x^T X$$

假定 $X^T X$ 为正定矩阵, 则 $E_{\hat{w}}$ 为关于 \hat{w} 的凸函数.

$$\nabla E_{\hat{\omega}} = \frac{\partial E_{\hat{\omega}}}{\partial \hat{\omega}} = 2X^T(X\hat{\omega} - y) = 0$$

$$X^T X \hat{\omega} = X^T y$$

$$\Rightarrow \hat{\omega} = (X^T X)^{-1} X^T y \quad 3.11.$$

二. 对数几率回归.

在 线性模型 基础上 奠一个映射函数 进行分类.

$$\text{Sigmoid} : \quad f(x) = \frac{1}{1 + e^{-wx}} \quad \mathbb{R} \rightarrow (0, 1).$$

① 确定概率质量(密度)函数.

② 似然函数.

离散型 随变 $y \in \{0, 1\}$ 重建模.

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

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

令 $\beta = (w; b)$, $\hat{x} = (x; 1)$.

$$\text{则 } P(y=1 | \hat{x}; \beta) = \frac{e^{\beta^T \hat{x}}}{1+e^{\beta^T \hat{x}}} = p_1(\hat{x}; \beta)$$

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

$$\Rightarrow \text{随变} y \text{ 极限 } P(y | \hat{x}; \beta) = y \cdot p_1(\hat{x}; \beta) + (1-y) \cdot p_0(\hat{x}; \beta)$$

3.26

$$\text{或 } [p_1(\hat{x}; \beta)]^y [p_0(\hat{x}; \beta)]^{1-y}$$

$$\text{似然函数 } L(\beta) = \prod_{i=1}^m P(y_i | \hat{x}_i; \beta)$$

$$\text{对数似... } l(\beta) = \sum_{i=1}^m \ln P(\quad)$$

$$l(\beta) = \sum \ln (y_i p_1(\quad) + (1-y_i) p_0(\quad))$$

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

$$= \sum \ln (y_i e^{\beta^T \hat{x}_i} + (1-y_i)) - \ln (1+e^{\beta^T \hat{x}_i})$$

$$\because y_i \in \{0, 1\}$$

$$\therefore l(\beta) = \begin{cases} \sum (-\ln (1+e^{\beta^T \hat{x}_i})) & y_i=0 \\ \sum (\beta^T \hat{x}_i - \ln (1+e^{\beta^T \hat{x}_i})) & y_i=1 \end{cases}$$

$$\text{综合 } l(\beta) = \sum (y_i \beta^T \hat{x}_i - \ln(1+e^{-\dots}))$$

损失函数 \rightarrow 目标最小化

$$\text{最大化 } l(\beta) \Leftrightarrow \text{最小化 } -l(\beta)$$

信息论:

$$\text{自信息: } I(x) = -\log_b P(x).$$

$$b=2 \text{ 单位 bit} \quad b=e \text{ 单 nat}$$

信息熵 (自信息期望): 度量 \times 不确定性 \nearrow 不确定 \nearrow

$$H(x) = E[I(x)] = -\sum_x P(x) \log_b P(x). \quad (\text{离散型})$$

$$\text{约定: } P(x)=0 \Rightarrow P(x) \log_b P(x) = 0$$

相对熵 (KL 散度): 度量分布差异. 想象分布 $P(x)$

模拟 $q(x)$.

$$D_{KL}(P||Q) = \sum_x P(x) \log_b \left(\frac{P(x)}{Q(x)} \right)$$

$$= \sum_x P(x) (\log_b P(x) - \log_b Q(x))$$

$$= \sum_x P(x) \log_b P(x) - \sum_x P(x) \log_b Q(x).$$

↑
最小化
最小化交叉熵
想分布信息熵 反向熵.

summary:

对数几率回归 三要素

{ 模型：线性模型，输出范围 $[0, 1]$ 近似阶段的单调
策略：极大似然估计 / 信息论。
算法：梯度下降 / 牛顿法。

三、二分类线性判别分析.