# 线性回归 (Linear Regression)

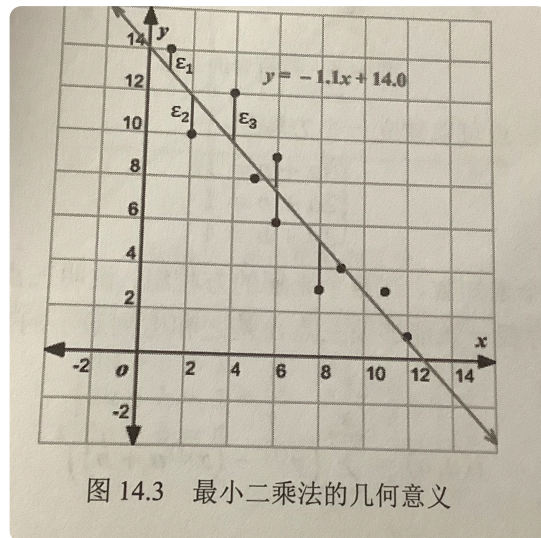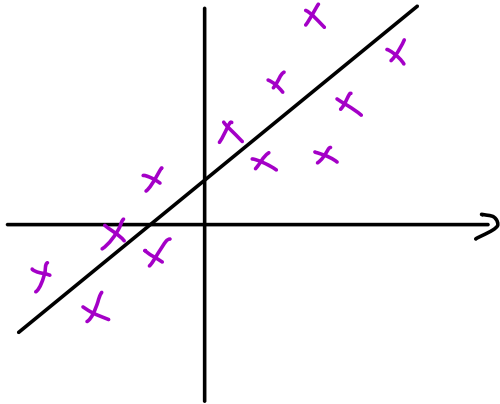## 1. Introduction.

- 2-d 线性拟合



图 14.3　最小二乘法的几何意义

- P维 线性拟合

前提:

样本 $D = \{(x_1, y_1), (x_2, y_2), \cdots, (x_n, y_n)\}$

其中, $x_i$ 是一个P维向量. 表达第 $i$ 个样本被观察的 P个特征.

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}_{P \times 1}$$

N个样本集合, 写作 $X = \begin{bmatrix} x_1 & x_2 & \cdots & x_N \end{bmatrix}^T_{P \times P}$

$y_i$ 表示第 $i$ 个样本的取值

$$Y = \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix}^T$$

假设:

回归到一条直线上. 即 $y = w^T x + b$

其中 $w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}$.

为了简化. 把偏置 $b$ 简化成 $w_0 x_{i0}$ 的形式. 则 $w$ 和 $x$ 向量又所成:

$$w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \qquad x_i = \begin{bmatrix} 1 \\ x_{i1} \end{bmatrix}$$

$$W = \begin{bmatrix} w_2 \\ \vdots \\ w_p \end{bmatrix}, \quad \begin{bmatrix} x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}$$

即原有的映射关系可以写成 $\boxed{y = w^T x}$

## 2. 最小二乘估计

- 如何估计上述的 $w$ ?

最小二乘估计给出以下公式

$$L(w) = \sum_{i=1}^{N} |\underbrace{w^T x_i}_{\text{估计值}} - \underbrace{y_i}_{\text{真实值}}|^2$$

展开上式:

$$L(w) = (w^T x_1 - y_1, w^T x_2 - y_2, \cdots w^T x_N - y_N) \cdot (w^T x_1 - y_1, w^T x_2 - y_2, \cdots w^T x_N - y_N)^T$$

$$= (w^T x_1, w^T x_2 \cdots w^T x_N) - (y_1, y_2, \cdots y_N)$$

$$= w^T (x_1, x_2, \cdots x_N) - (y_1, y_2, \cdots, y_N)$$

$$= (w^T x^T - y^T)(w^T x^T - y)^T$$

$$= (w^T x^T - y^T)(XW - Y) \qquad (w^T x^T y)^T$$

$$= w^T x^T X W - w^T x^T Y - Y^T X W + Y^T Y$$

$$\boxed{L(w) = w^T x^T X W - 2 w^T x^T Y + Y^T Y}$$

$$\frac{\partial}{\partial w} L(w) = 2 x^T X W - 2 x^T Y = 0$$

$$\Rightarrow \boxed{W = (x^T x)^{-1} \cdot x^T Y}$$

$$X = (x_1, x_2, \cdots, x_N)^T$$
$$Y = (y_1, y_2, \cdots y_N)^T$$

$(K)^T = K$. 实数转置为本身

$Y^T X W \to 1 \times P \times P \times P \times P \times 1$

$\because (x^T x)^T = x^T x \to$ 对称矩阵

对于对称矩阵 $\frac{\partial w^T A w}{\partial w} = 2Ax$

$$\frac{\partial x^T A}{\partial x} = A$$

- 从高斯噪声角度看最小二乘

  - 设噪声 $\varepsilon$ 服从高斯分布 $\varepsilon \sim N(0, \sigma^2)$

   则真实抽样值为 $y = w^T x + \underbrace{\varepsilon}_{\text{看成常数}}$

$$\Rightarrow y(x; w) \sim N(w^T x, \sigma^2)$$

MLE: $L(w) = \log P(y | x; w)$

$$= \log \prod_{i=1}^{N} P(y|x; w)$$

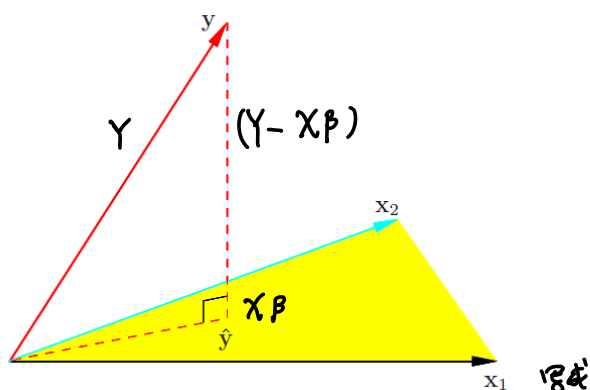$$P(y|x; w) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left[-\frac{(y - w^T x)^2}{2\sigma^2}\right]$$

$$= \sum_{i=1}^{N} \log P(y \mid x; w)$$

$$= \sum_{i=1}^{N} \log \left( \frac{1}{\sqrt{2\pi}\sigma} \exp\left[ -\frac{(y-w^T x_i)^2}{2\sigma^2} \right] \right)$$

$$= \sum_{i=1}^{N} \log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp\left[ -\frac{(y-w^T x_i)^2}{2\sigma^2} \right]$$

$$= \sum_{i=1}^{N} \left( \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2}(y-w^T x_i)^2 \right]$$

$$\Rightarrow \hat{w} = \underset{w}{argmax} \; L(w)$$

$$= \underset{w}{argmax} \; \sum_{i=1}^{N} \left( -\frac{1}{2\sigma^2}(y-w^T x_i)^2 \right]$$

$$= \underset{w}{argmin} \; \sum_{i=1}^{N} \left( y - w^T x_i \right)^2 \qquad \textcolor{red}{\rightarrow \text{同之前所差法式}}$$

· 最小二乘法几何意义



$$\chi^T(Y - \chi\beta) = 0$$
$$\chi^T Y - \chi^T \chi \beta = 0$$
$$\Rightarrow \beta = (\chi^T \chi)^{-1} \chi^T Y$$

$$X = span(X_1, \cdots, X_n) \Rightarrow f(w) = \chi\beta$$

3. 正则化 (Regulation)

· Background

Loss Function : $L(w) = \sum_{i=1}^{N} |w^T x_i - y_i|^2$

$$\Rightarrow \hat{w} = (\chi^T \chi)^{-1} \chi^T Y$$

$\chi \rightarrow n \times p$ 矩阵. 当 $n >> p$ 时.

(n个样本, $x_i \in \mathbb{R}^p$)

$$\chi = \begin{pmatrix} 0_{11} & 0_{21} \\ \underline{\hspace{1cm}} & \underline{\hspace{1cm}} \end{pmatrix}$$

<span style="color:blue">当样本数量不能远大于特征维度</span>

<span style="color:blue">数字解释:</span>

    hint, 不满秩是因为维特之间的阶 ⟺ $\chi^T \chi$ 存在线性相关组 ⟺ $\chi^T \chi$ 不可逆

(cond. 矩阵 $X$ 列 不满秩 ...)

(因变量过多了，出现引牌)
e.x. $x_1 + x_2 + x_3 + x_4 + x_5 = 1$
$x_2 + x_3 + x_4 = 10$

几何解释： **过拟合 (overfitting)**



→ 有无多同称

• 解决方法
① 土曾加数据
② 降维/特征提取（$P < A$）
③ 正则化

• 正则化框架

$$\arg\min_{w} \left[ L(w) + \lambda P(w) \right]$$

$L_1$: Lasso Regression    $P(w) = \|w\|$
$L_2$: Ridge Regression    $P(w) = \|w\|_2^2 = w^T w$

$L_2$:  $J(w) = \sum_{i=1}^{N} \|w^T x_i - y_i\|^2 + \lambda w^T w$

$= (w^T x_1, w^T x_2 \cdots w^T x_n)(y_1, y_2, \cdots, y_n)$

$= (w^T X^T - Y^T)(w^T X^T - Y^T)^T + \lambda w^T w$

$= (w^T X^T - Y^T)(Xw - Y) + \lambda w^T w$

$= w^T X^T X w - w^T X^T Y - Y^T X w + Y^T Y + \lambda w^T w$

$$= w^T X^T X w - 2 w^T X^T Y + \lambda w^T w + Y^T Y$$

$$= w^T (X^T X + \lambda I) w - 2 w^T X^T Y - Y^T Y$$

$$\hat{w} = \arg\min_{w} J(w)$$

$$\frac{d J(w)}{w} = 2(X^T X + \lambda I) w - 2 X^T Y = 0$$

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

<span style="color:red">避免了 $X^T X$ 不可逆</span>

• 从贝叶斯角看

取先验分布 $w \sim N(0, \sigma^2)$. 于是

$$p(w|y) = \frac{p(w) \cdot p(y|w)}{p(y)} \Rightarrow p(w|y) \propto p(w) \cdot p(y|w).$$

<span style="color:red">(最大后验概率)</span>

MAP: $\hat{w} = \arg\max_{w} \log p(w|y)$

<span style="color:red">$p(w) = \dfrac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\dfrac{w^2}{2\sigma^2}\right)$</span>

<span style="color:red">$p(w|y) = \dfrac{1}{\sqrt{2\pi}\sigma_0} \exp\left(-\dfrac{(y - w^T x)^2}{2\sigma_0^2}\right)$</span>

$$= \arg\max_{w} \log p(w) \cdot p(y|w)$$

$$= \arg\max_{w} \log\left(\frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{2\pi}\sigma_0}\right) + \log \exp\left[-\frac{w^2}{2\sigma^2} - \frac{(y - w^T x)^2}{2\sigma_0^2}\right]$$

$$= \arg\min_{w} \frac{(y - w^T x)^2}{2\sigma_0^2} + \frac{w^2}{2\sigma^2}$$

$$\boxed{= \arg\min_{w} (y - w^T x)^2 + \frac{2\sigma_0^2}{2\sigma^2} w^2}$$ <span style="color:red">$Lasso!$</span>

(LSE)最小二乘估计 $\Longleftrightarrow$ 极大似然估计 (noise 为 Gaussian Distribution)

Regularized LSE $\Longleftrightarrow$ MAP (先验分布为 Gaussian Distribution)