

Section 3 Linear Regression

最小二乘法 (Least Square Method)

数据集: $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

其中 $x_i \in \mathbb{R}^P$, $y_i \in \mathbb{R}$, $i = 1, 2, \dots, N$.

$$\text{令 } X = (x_1, x_2, \dots, x_N)^T, X \in \mathbb{R}^{N \times P}$$

$$Y = (y_1, y_2, \dots, y_N)^T, Y \in \mathbb{R}^{N \times 1}$$

最小二乘估计: $L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$

矩阵表达:

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

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

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

$$\stackrel{\Downarrow}{w^T X^T} - Y^T$$

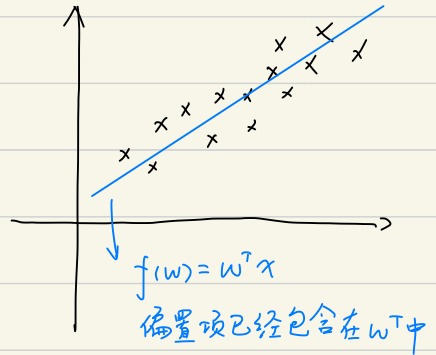
$$= (w^T X^T - Y^T) (Xw - Y)$$

$$= w^T X^T X w - \underbrace{w^T X^T Y}_{(w^T X^T)^T = Y^T X w} - Y^T X w + Y^T Y$$

每一项都为实数.

$$(w^T X^T)^T = Y^T X w$$

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



$$\begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} \Rightarrow Xw - Y$$

$$k) \quad W = \arg \min_W W^T X^T X W - 2W^T X^T Y + Y^T Y$$

$$\frac{\partial \mathcal{L}(W)}{\partial W} = \underbrace{2X^T X W}_{\downarrow \text{证明}} - 2X^T Y = 0 \Rightarrow \boxed{W = (X^T X)^{-1} X^T Y}$$

→ X^+ , 称为 X 的伪逆

简化问题, 令 $A = X^T X$. 求 $\frac{\partial}{\partial W} (W^T A W)$. 其中 $W \in \mathbb{R}^{p \times 1}$

假设 $W = (w_1, w_2, \dots, w_p)^T$.

$$A = \begin{pmatrix} a_{11} & \dots & a_{1p} \\ \vdots & & \vdots \\ a_{p1} & \dots & a_{pp} \end{pmatrix}$$

$$k) \quad W^T A W = \left(\sum_{i=1}^p w_i a_{i1} \quad \sum_{i=1}^p w_i a_{i2} \quad \dots \quad \sum_{i=1}^p w_i a_{ip} \right) (w_1 \ w_2 \ \dots \ w_p)^T$$

$$= \sum_{j=1}^p w_j \sum_{i=1}^p w_i a_{ij}$$

$$\frac{\partial (W^T A W)}{\partial w_k} = w_1 a_{k1} + w_2 a_{k2} + \dots + \left(\sum_{i=1}^p w_i a_{ik} + w_k a_{kk} \right) + \dots$$

$$= \sum_{i=1}^p w_i a_{ki} + \sum_{i=1}^p w_i a_{ik}$$

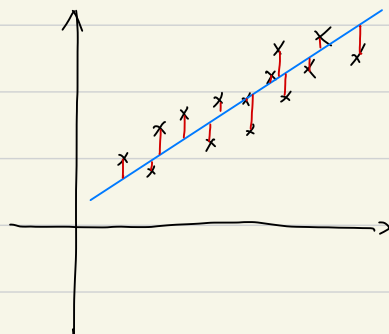
$$k) \quad \frac{\partial (W^T A W)}{\partial W} = \begin{pmatrix} \sum_{i=1}^p w_i a_{i1} \\ \sum_{i=1}^p w_i a_{i2} \\ \vdots \\ \sum_{i=1}^p w_i a_{ip} \end{pmatrix} + \begin{pmatrix} \sum_{i=1}^p w_i a_{1i} \\ \sum_{i=1}^p w_i a_{2i} \\ \vdots \\ \sum_{i=1}^p w_i a_{pi} \end{pmatrix} = A W + A^T W$$

由于 $A = X^T X$ 是对称矩阵, 即 $A^T = A$.

$$\text{所以} \quad \frac{\partial (W^T A W)}{\partial W} = 2AW \Rightarrow \frac{\partial (W^T X^T X W)}{\partial W} = 2X^T X W.$$

几何解释

①



红色的线段就是目标函数需要最小化的对象。即将误差均摊到每一个样本点。

② 我们注意到 $X = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^{n \times p}$ 。其中 $n > p$ 。

将 X 的列向量看作是一组基底，那么这组基底就在 n 维空间中构成了 p 维子空间。
则 XW 就是这 p 个列向量的线性组合，仍属于这个 p 维子空间。

线性回归问题就可以看作，寻找 p 个 n 维列向量的一个线性组合，使其与 Y 的距离最小。即该线性组合是 Y 在 p 个向量构成空间上的投影。
这个距离其实就是误差

由此可得：线性组合为 XW 。

由于线性组合是 Y 在子空间上的投影。

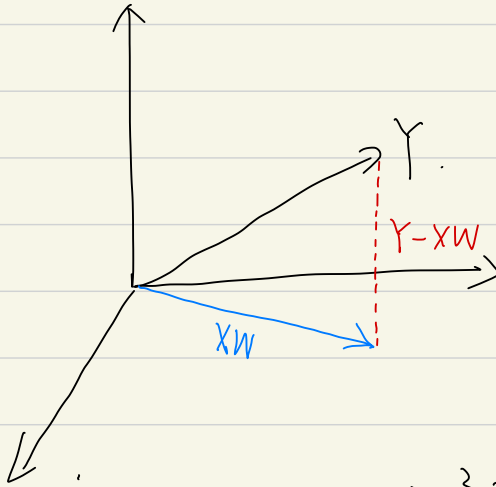
那么 $Y - XW$ 为子空间的一个法向量。

则 $X^T(Y - XW) = 0$ 法向量与基底内积为 0。

$$\Rightarrow W = (X^T X)^{-1} X^T Y$$

$n=3$, $p=2$ 为例。假设基为 $(1, 0, 0)^T$, $(0, 1, 0)^T$

构成的空间为 xoy 平面。



$$X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$Y = \begin{bmatrix} 3 \\ 2 \\ 2 \end{bmatrix}$$

$$W = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$W = (X^T X)^{-1} X^T Y = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

综上, ①的角度是将误差均摊到每个样本上
②的角度是将误差在向量空间中表示出来。

极大似然估计 = 最小二乘估计

已知: $X = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^{n \times p}$, $Y = (y_1, y_2, \dots, y_n) \in \mathbb{R}^{n \times 1}$
 $W = (w_1, w_2, \dots, w_p)^T \in \mathbb{R}^{p \times 1}$, $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

则可设 $y_i = f(w) + \varepsilon$, 其中 $f(w) = w^T x_i$.

则 $y_i = w^T x_i + \varepsilon$.

即 $y_i | x_i; w \sim \mathcal{N}(w^T x_i, \sigma^2)$.

定义 $L(w) = \log P(Y|X;w)$

$$= \sum_{i=1}^n \log P(y_i | x_i; w)$$

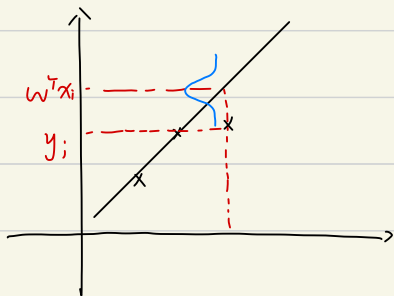
$$\hat{W}_{MLE} = \arg \max_w \sum_{i=1}^n \log P(y_i | x_i; w)$$

$$= \arg \max_w \sum_{i=1}^n \left(\frac{1}{2} \log \frac{1}{2\pi\sigma^2} - \frac{(y_i - w^T x_i)^2}{2\sigma^2} \right)$$

$$= \arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2$$

和最小二乘法的目标函数一致.

所以, 最小二乘估计其实等价于噪声为高斯分布的极大似然估计



如左图, 表示了两者等价之图形表达.

蓝色为高斯分布 $\mathcal{N}(0, \sigma^2)$, 误差就是其造成的.

$$y_i - w^T x_i$$

正则化 (Regularization)

岭回归 (Ridge Regression)

最小二乘估计的目标函数: $L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$

解析解: $\hat{w} = (X^T X)^{-1} X^T Y$.

引入正则化的原因:

这里需要: 注意 $(X^T X)^{-1}$.

由于 $X \in \mathbb{R}^{N \times P}$, 正常情况下 $N \gg P$.

如果 $N < P$.

由于 $r(X^T X) = r(X)$.

则 $r(X^T X)$ 必定小于 P . 但是 $X^T X \in \mathbb{R}^{P \times P}$

$\therefore X^T X$ 不可逆 数值上不可逆. 线性回归角度容易造成过拟合.

\therefore 解析解无法直接得到.

如何解决过拟合 \rightarrow

- ① 增加数据量
- ② 特征选择 / 特征提取 (降维).
- ③ 正则化 对参数空间的约束.

\rightarrow 正则化框架: $\arg \min_w [\underbrace{L(w)}_{\text{Loss Function}} + \lambda \underbrace{P(w)}_{\text{Penalty}}]$

正则化较为常用：为 L_1 正则化与 L_2 正则化

L_1 : Lasso. $P(W) = \|W\|_1$, 1-范数

L_2 : Ridge. 岭回归 $P(W) = \|W\|_2^2 = W^T W$. 2-范数 \approx 平方.
 \rightarrow 权值衰减.

应用 L_2 后.

$$\text{目标为, } \arg\min_W \left[\sum_{i=1}^N \|W^T x_i - y_i\|^2 + \lambda W^T W \right]$$

$$= \arg\min_W \left[(W^T X^T - Y^T)(XW - Y) + \lambda W^T W \right]$$

$$= \arg\min_W \left[W^T X^T X W - 2 W^T X^T Y + Y^T Y + \lambda W^T W \right].$$

$$= \arg\min_W \left[\underbrace{W^T (X^T X + \lambda I) W - 2 W^T X^T Y + Y^T Y}_{G(W)} \right].$$

$$\frac{\partial}{\partial W} G(W) = 2(X^T X + \lambda I)W - 2X^T Y = 0.$$

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

一定可逆 因为 $X^T X$ 半正定, 加上 I , 为正定, 一定可逆.

从贝叶斯角度看待岭回归

假设 w 服从一个先验分布: $w \sim \mathcal{N}(0, \Sigma)$, 且 Σ 为对角, 各向同性

则由贝叶斯定理:
$$p(w|Y) = \frac{p(Y|w) \cdot p(w)}{p(Y)}$$

MAP(最大后验估计):
$$\hat{w} = \underset{w}{\operatorname{argmax}} p(w|Y)$$
$$= \underset{w}{\operatorname{argmax}} p(Y|w) p(w)$$

由之前概率角度看待最小二乘估计可得:

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

$\Rightarrow p(Y|w) = \prod_{i=1}^N p(y_i|w)$ (y_i 之间相互独立, 联合密度是边缘密度的乘积)

则
$$\hat{w} = \underset{w}{\operatorname{argmax}} p(Y|w) p(w)$$

$$= \underset{w}{\operatorname{argmax}} [\log p(Y|w) p(w)]$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{(y_i - w^T x_i)^2}{2\sigma^2} \right] + \log \frac{1}{\sqrt{2\pi}|\Sigma|^{\frac{1}{2}}} - \frac{1}{2} w^T \Sigma^{-1} w$$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^N (y_i - w^T x_i)^2 + \sigma^2 w^T \Sigma^{-1} w$$

由于 Σ^{-1} 为对角且各向同性, 则 $\Sigma = \sigma_w^2 I$

$$\therefore \hat{w} = \underset{w}{\operatorname{argmin}} \sum_{i=1}^N (y_i - w^T x_i)^2 + \underbrace{\left(\frac{\sigma^2}{\sigma_w^2} \right)}_{L_2} w^T w$$

和上述 L_2 正则化后的目标函数相同。

最小二乘估计

综上 $\text{LSE} \iff \text{MLE} (\text{noise} \sim \text{Gaussian Distribution})$

Regularized LSE $\iff \text{MAP} (\text{noise} \sim \text{Gaussian Distribution}$
 $\text{prior} \sim \text{Gaussian Distribution})$

这个分布是各向同性。