

Linear Regression



1 LINEAR REGRESSION

1.1 Least Square

The dataset is $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$. The model is $f(w) = w^T x$. The loss function is

$$\begin{aligned}
 L(w) &= \|w^T x_i - y_i\|^2 \\
 &= (w^T x_1 - y_1, \dots, w^T x_n - y_n) \cdot (w^T x_1 - y_1, \dots, w^T x_n - y_n)^T \\
 &= (w^T X^T - Y^T) \cdot (Xw - Y) \\
 &= w^T X^T Xw - w^T X^T Y - Y^T Xw + Y^T Y
 \end{aligned} \tag{1}$$

Therefore,

$$\begin{aligned}
 \hat{w} &= \arg \min_w L(w) = 0 \\
 \therefore \frac{\partial a^T x}{\partial x} &= \frac{\partial x^T a}{\partial x} = a \\
 \therefore \frac{\partial L}{\partial w} &= 2X^T Xw - 2X^T Y = 0 \\
 \hat{w} &= \{X^T X\}^{-1} X^T Y
 \end{aligned} \tag{2}$$

1.2 MLE with Gaussian Noise

Given $y = w^T x + \eta$, $\eta \sim \mathcal{N}(0, \sigma^2)$, so $y \sim \mathcal{N}(w^T x, \sigma^2)$. Therefore,

$$\begin{aligned}
 L(w) &= \log p(Y|X, w) \\
 &= \log \prod_{i=1}^N p(y_i|x_i, w) \\
 &= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right) \right) \\
 \therefore \arg \max_w L(w) &= \arg \min_w \sum_{i=1}^N (y_i - w^T x_i)^2
 \end{aligned} \tag{3}$$

From Bayesian, having $w \sim \mathcal{N}(0, \sigma^2)$, therefore,

$$\begin{aligned}
 \hat{w} &= \arg \max_w \log p(w|Y) \\
 &= \arg \max_w \log(p(Y|w)p(w)) \\
 &= \arg \max_w (\log p(Y|w) + \log p(w)) \\
 \therefore p(Y|w) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right) \\
 \therefore p(w) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\|w\|^2}{2\sigma^2}\right), w \sim \mathcal{N}(0, \sigma^2) \\
 \therefore \hat{w} &= \arg \max_w \left(-\frac{(Y - w^T x)^2}{2\sigma^2} - \frac{\|w\|^2}{2\sigma^2} \right) \\
 &= \arg \min_w \underbrace{\left((Y - w^T x)^2 + \|w\|^2 \right)}_{\text{ridge regression}}
 \end{aligned} \tag{4}$$