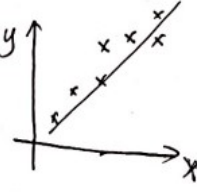


Linear Regression (线性回归)

1. Linear Regression



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

$$x_i \in \mathbb{R}^p, y_i \in \mathbb{R} \quad i=1, 2, \dots, N$$

$$X = (x_1, x_2, \dots, x_n)^T = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}_{N \times p}$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1}$$

$$f(w) = w^T X \quad w = \begin{pmatrix} w_1 \\ \vdots \\ w_p \end{pmatrix}$$

D: data(N samples)

(1) Algorithm: Least squares estimation (最小二乘估计)

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

$$\begin{aligned} L(w) &= \sum_{i=1}^N \|w^T x_i - y_i\|^2 \\ &= \sum_{i=1}^N (w^T x_i - y_i)(w^T x_i - y_i) \\ &= \underbrace{(w^T x_1 - y_1, \dots, w^T x_N - y_N)}_{\downarrow} \begin{pmatrix} w^T x_1 - y_1 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} \\ &= (w^T x_1, \dots, w^T x_N) - (y_1, \dots, y_N) (Xw - Y) \\ &= w^T (x_1, \dots, x_N) - (y_1, \dots, y_N) \cdot (Xw - Y) \\ &= (w^T X^T - y^T) \cdot (Xw - Y) \\ &= w^T X^T X w - w^T X^T Y - Y^T X w + Y^T Y = w^T X^T X w - 2w^T X^T Y + Y^T Y \end{aligned}$$

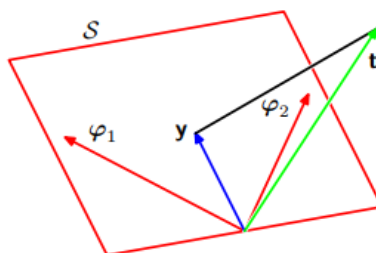
$$\hat{w} = \operatorname{argmin} L(w)$$

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

$$\hat{w} = (X^T X)^{-1} X^T Y \quad (X \text{ has to be a non-singular matrix, if } X \text{ is a singular matrix, use SVD/QR decomposition})$$

(2) Geometric interpretation of LSE (source: PRML) (几何角度解释最小二乘估计)

The idea is to separate errors in every dimensions, thus in an N -dimensional space whose axes are the values of t_1, \dots, t_N . The least-squares regression function is obtained by finding the orthogonal projection of the data vector \mathbf{t} onto the subspace spanned by the basis functions $\phi_j(\mathbf{x})$ in which each basis function is viewed as a vector ϕ_j of length N with elements $\phi_j(\mathbf{x}_n)$.
把误差分散到每个维度



(3) solving Loss function to obtain the minimum of loss function (求解损失函数的方法)

(Source: <http://cs229.stanford.edu/notes/cs229-notes1.pdf>)

- ① batch gradient descent: looks at every example in the entire training set on every step
(批量梯度下降法)

$$w_j := w_j - \alpha \frac{\partial L(w)}{\partial w_j} \quad \text{where } \alpha \text{ is the learning rate}$$

Starts with some “initial guess” for w , and that repeatedly changes w to make $L(w)$ smaller, until hopefully we converge to a value of w that minimizes $L(w)$

- ② Stochastic gradient descent: update the parameters according to the gradient of the error with respect to that single training example only.
(随机梯度下降法)

```
Repeat {  
  For  $l=1$  to  $N$ , {  
     $w_j := w_j - \alpha \frac{\partial L(w)}{\partial w_j}$  (for every  $j$ )  
  }  
}
```

2. Lasso/Ridge Regression (L1/L2 正则化)

In real life, $X \sim N \times p$ matrix, if the $N \ll p$, this will probably cause over fitting problem.

So there exist three solving method:

- Add more data (增加数据)
- Dimensionality reduction (降维)
- Regularization (正则化)

Regularization function's structure: $\text{argmin}(L(w) + \lambda P(w))$

L1: $P(w) = \|w\|_1$

L2: $P(w) = \|w\|_2^2 = w^T w$

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

For L2 (from frequency prospective):

$$\begin{aligned}
 \text{Goal function } J(w) &= \sum_{i=1}^N \|w^T x_i - y_i\|^2 + \lambda w^T w \\
 &= (w^T x_1 - y_1, \dots, w^T x_N - y_N) \begin{pmatrix} w^T x_1 - y_1 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} + \lambda w^T w \\
 &= (w^T X^T - Y) (Xw - Y) + \lambda w^T w \\
 &= w^T X^T X w - 2w^T X^T Y + Y^T Y + \lambda w^T w \\
 &= w^T (X^T X + \lambda I) w - 2w^T X^T Y + Y^T Y \\
 \hat{w} &= \text{argmin } J(w) \\
 \frac{dJ(w)}{dw} &= 2(X^T X + \lambda I)w - 2X^T Y = 0 \\
 \boxed{w} &= \boxed{(X^T X + \lambda I)^{-1} X^T Y}
 \end{aligned}$$

For L2 (from Bayesian perspective):

$$\begin{aligned}
 w &\sim N(0, \sigma^2) \\
 P(w|y) &= \frac{P(y|w) \cdot P(w)}{P(y)} \quad y = f(x) + \epsilon \quad \epsilon \sim N(0, \sigma^2) \\
 \hat{w} &= \text{argmax}_w P(w|y) \\
 &= \text{argmax}_w P(y|w) \cdot P(w) \quad \text{where } P(y|w) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y - w^T x)^2}{2\sigma^2}\right\} \\
 &\quad P(w) = \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left\{-\frac{\|w\|^2}{2\sigma_0^2}\right\} \\
 &= \text{argmax}_w \frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left\{-\frac{(y - w^T x)^2}{2\sigma^2} - \frac{\|w\|^2}{2\sigma_0^2}\right\}
 \end{aligned}$$

$$\begin{aligned}
&= \arg \max_w \log [P(y|w) P(w)] \\
&= \arg \max_w \log \left(\frac{1}{\sqrt{2\pi}\sigma_0} \cdot \frac{1}{\sqrt{2\pi}\sigma_0} \right) + \log \exp \left\{ -\frac{(y-w^T x)^2}{2\sigma_0^2} - \frac{\|w\|^2}{2\sigma_0^2} \right\} \\
&= \arg \min_w \left(\frac{(y-w^T x)^2}{2\sigma_0^2} + \frac{\|w\|^2}{2\sigma_0^2} \right) \\
&= \arg \min_w \left((y-w^T x)^2 + \frac{\sigma_0^2}{\sigma_0^2} \|w\|_2^2 \right)
\end{aligned}$$

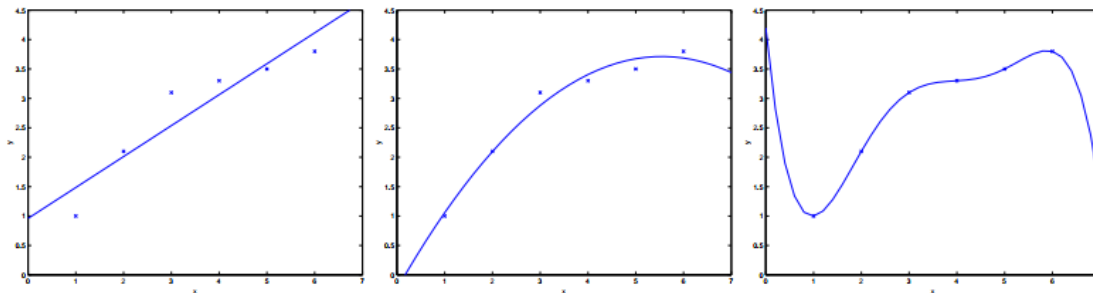
$$\begin{aligned}
w_{\text{MAP}} &= \arg \min_w \sum_{i=1}^N (y_i - w^T x_i)^2 + \frac{\sigma_0^2}{\sigma_0^2} \|w\|_2^2 \\
J(w) &= \sum_{i=1}^N \|w^T x_i - y_i\|^2 + \lambda w^T w
\end{aligned}$$

Regularized LSE \Leftrightarrow MAP (noise is Gaussian distribution)
 prior is Gaussian too)

3. Locally Weighted Linear Regression

(Source: <http://cs229.stanford.edu/notes/cs229-notes1.pdf>)

This method is used to solving under fitting problem and it is also a non-parametric method.



Steps:

1. Fit w to minimize $L(w) = \sum_{i=1}^N \theta \|W^T x_i - y_i\|^2$
2. output $W^T X$.