

Linear Regression

Note Title

3/9/2010

Training Data with N observations

$$\{(\underline{x}_i, y_i), i=1, 2, \dots, N\}, \quad \underline{x}_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R}$$

data point \swarrow target variable associated with \underline{x}_i

$\begin{pmatrix} x_i(1) \\ x_i(2) \\ \vdots \\ x_i(d) \end{pmatrix}$ features/attributes

$$\hat{y} = f(\underline{x})$$

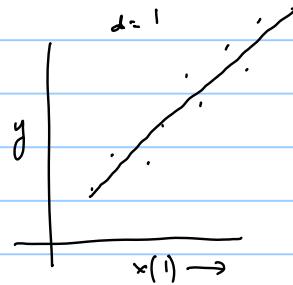
Learning problem: Given training data, come up with a "good f "

Linear Regression: $f(\underline{x})$ is a linear function of \underline{x}

$$\underline{x} = \begin{pmatrix} x(1) \\ x(2) \\ \vdots \\ x(d) \end{pmatrix}, \quad f(\underline{x}) = w_0 + w_1 x(1) + w_2 x(2) + \dots + w_d x(d)$$

$$f(\underline{x}_i) \approx y_i, \quad i=1, 2, \dots, N$$

$$\min_{w_0, w_1, \dots, w_d} \sum_{i=1}^N \left[(w_0 + w_1 x_i(1) + w_2 x_i(2) + \dots + w_d x_i(d)) - y_i \right]^2$$



Special case when $d=1$

$$f(x) = w_0 + w_1 x$$

$$\min_{w_0, w_1} \sum_{i=1}^N (w_0 + w_1 x_i - y_i)^2 \quad \text{— Least Squares formulation}$$

$E(w)$

$$\frac{\partial E}{\partial w_0} = \sum_{i=1}^N 2(w_0 + w_1 x_i - y_i) = 0 \Rightarrow N w_0 + w_1 \left(\sum_{i=1}^N x_i \right) = \sum_{i=1}^N y_i \quad \text{①}$$

$$\frac{\partial E}{\partial w_1} = \sum_{i=1}^N 2(w_0 + w_1 x_i - y_i) x_i = 0 \Rightarrow \left(\sum_{i=1}^N x_i \right) w_0 + w_1 \left(\sum_{i=1}^N x_i^2 \right) = \sum_{i=1}^N x_i y_i \quad \text{②}$$

$$\begin{bmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N x_i y_i \end{bmatrix} \quad \text{— Normal Equations}$$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

$$\begin{bmatrix} N & N\bar{x} \\ N\bar{x} & \sum_{i=1}^N x_i^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} N\bar{y} \\ \sum_{i=1}^N x_i y_i \end{bmatrix}$$

$$\begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \frac{1}{N} \sum_{i=1}^N x_i^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \bar{y} \\ \frac{1}{N} \sum_{i=1}^N x_i y_i \end{bmatrix}$$

$$Aw = b$$

$$w = A^{-1}b$$

$$\det(A) = \frac{1}{N} \sum_{i=1}^N x_i^2 - \bar{x}^2$$

$$= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \sigma_x^2$$

$$w = \frac{1}{\sigma_x^2} \begin{bmatrix} \frac{1}{N} \sum x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} \bar{y} \\ \frac{1}{N} \sum x_i y_i \end{bmatrix}$$

$$w_0 = \frac{1}{\sigma_x^2} \left[\frac{1}{N} \bar{y} \sum x_i^2 - \frac{1}{N} \bar{x} \sum x_i y_i \right]$$

$$w_1 = \frac{1}{\sigma_x^2} \left[-\bar{x} \bar{y} + \frac{1}{N} \sum x_i y_i \right] = \sigma_{xy}$$

$$w_1 = \frac{\sigma_{xy}}{\sigma_{xx}}$$

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$Xw = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \\ w_0 + w_1 x_3 \\ \vdots \\ w_0 + w_1 x_N \end{bmatrix}$$

Least Squares Objective: $\min_w \|Xw - y\|_2^2$ where $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

$$X^T X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} = \begin{bmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$$X^T y = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$(X^T X)w = X^T y \quad \text{--- Same as Normal Equations}$$

Training points x_i are d -dimensional, $x_i \in \mathbb{R}^d$, $x_i = \begin{bmatrix} x_i(1) \\ x_i(2) \\ \vdots \\ x_i(d) \end{bmatrix}$

$$X = \begin{bmatrix} 1 & x_1(1) & x_1(2) & \dots & x_1(d) \\ 1 & x_2(1) & x_2(2) & \dots & x_2(d) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n(1) & x_n(2) & \dots & x_n(d) \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$Xw = \begin{bmatrix} w_0 + w_1 x_1(1) + w_2 x_2(1) + \dots + w_d x_d(1) \\ w_0 + w_1 x_1(2) + w_2 x_2(2) + \dots + w_d x_d(2) \\ \vdots \\ w_0 + w_1 x_1(n) + w_2 x_2(n) + \dots + w_d x_d(n) \end{bmatrix}$$

Least Squares Objective: $\min_w \underbrace{\|Xw - y\|_2^2}_{g(w)}$

$$g(w) = \|Xw - y\|_2^2 = (Xw - y)^T (Xw - y) = w^T X^T X w - 2y^T X w + y^T y$$

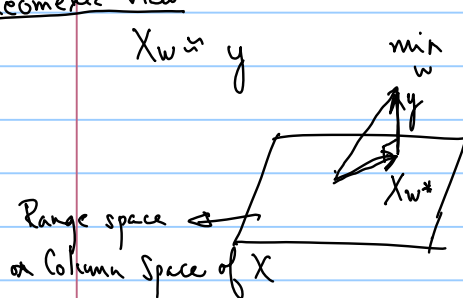
$$\nabla_w g = \frac{\partial g}{\partial w} = \begin{bmatrix} \frac{\partial g}{\partial w_0} \\ \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_d} \end{bmatrix} = 2X^T X w - 2X^T y = 0$$

$$\Rightarrow \boxed{X^T X w = X^T y} \text{ - Normal Equations}$$

BLUE

Best Linear Unbiased Estimation

Geometric View



$$Xw^* - y \perp Xw \text{ for all } w$$

$$(Xw)^T (Xw^* - y) = 0$$

$$\Rightarrow w^T X^T (Xw^* - y) = 0 \quad \forall w$$

$$\Rightarrow X^T X w^* - X^T y = 0$$

$$\Rightarrow \boxed{X^T X w^* = X^T y}$$

Instead of solving $\min_w \|Xw - y\|_2^2$

Use Regularization $\min_w \left(\|Xw - y\|_2^2 + \lambda \|w\|_2^2 \right), \lambda \geq 0$
↓
regularization parameter

Ridge Regression

$$\min_w \|Xw - y\|_2^2 + \lambda \|w\|_1 \text{ - Lasso}$$

Ridge Regression

Identical to solving

Solution: $(X^T X + \lambda I) w = X^T y$ - Ridge Regression

Lasso: $\min_w \|Xw - y\|_2^2$

st $\|w\|_1 \leq t$

st $\|w\|_1 \leq t$