

Mean Squared Error Regression

(Also) Least-squares Regression



Use Mean-squared error (MSE)

$$\therefore \text{Criterion Function: } J(\underline{w}) = \frac{1}{N} \sum_{i=1}^N [g(\underline{x}_i) - y_i]^2 = \text{MSE}$$

$$\text{where } g(\underline{x}_i) = \underline{w}^T \underline{x}_i$$

Algebraic Method [Ordinary Least Squares (OLS)]

$$J(\underline{w}) = \frac{1}{N} \sum_{i=1}^N [\underline{w}^T \underline{x}_i - y_i]^2$$

$$\text{Sum} \left\{ \begin{bmatrix} (\underline{w}^T \underline{x}_1 - y_1)^2 \\ (\underline{w}^T \underline{x}_2 - y_2)^2 \\ \vdots \\ (\underline{w}^T \underline{x}_N - y_N)^2 \end{bmatrix} \right\} = \left\| \begin{bmatrix} \underline{w}^T \underline{x}_1 - y_1 \\ \vdots \\ \underline{w}^T \underline{x}_N - y_N \end{bmatrix} \right\|_2^2 = \left\| \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \cdot \underline{w} - \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \right\|_2^2 = \|\underline{X} \underline{w} - \underline{y}\|^2$$

$$\therefore J(\underline{w}) = \frac{1}{N} \|\underline{X} \underline{w} - \underline{y}\|^2 = \frac{1}{N} (\underline{X} \underline{w} - \underline{y})^T (\underline{X} \underline{w} - \underline{y})$$

$$\text{Minimize it: } \nabla_{\underline{w}} J(\underline{w}) = \frac{1}{N} [2 \underline{X}^T \underline{X} \underline{w} - 2 \underline{X}^T \underline{y}] = 0$$

$$\Rightarrow \underline{X}^T \underline{X} \hat{\underline{w}} = \underline{X}^T \underline{y}$$

$$\hat{\underline{w}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} \quad \text{if } \underline{X}^T \underline{X} \text{ nonsingular}$$

$$\underline{X}^- = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \leftarrow \text{Moore-Penrose (left) pseudoinverse of } \underline{X}$$

Gradient Descent Method [Least Mean Squares (LMS)]

Ex: Sequential GD

$$J(\underline{w}) = \frac{1}{N} \sum_{n=1}^N (\underline{w}^T \underline{x}_n - y_n)^2 = \sum_{n=1}^N J_n(\underline{w})$$

$$\therefore J_n(\underline{w}) = \frac{1}{N} [\underline{w}^T \underline{x}_n - y_n]^2$$

$$\nabla_{\underline{w}} J_n(\underline{w}) = \frac{2}{N} [\underline{w}^T \underline{x}_n - y_n] \cdot \underline{x}_n$$

Use sequential GD

weight update $\underline{w}(i+1) = \underline{w}(i) - \eta(i) \frac{2}{N} [\underline{w}(i)^T \underline{x}_n - y_n] \underline{x}_n$

Mean Squared Error Classification

Instead of $J(\underline{w}) = \frac{1}{N} \sum_{n=1}^N (\underline{w}^T \underline{x}_n - y_n)^2$

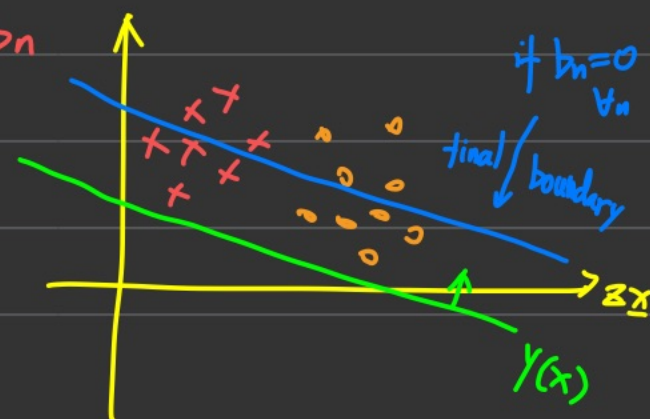
We use $J(\underline{w}) = \frac{1}{N} \sum_{n=1}^N [\underline{w}^T \underline{x}_n \underline{x}_n - b_n]^2$, $b_n > 0 \forall n$

\therefore only when all points are correctly classified ($\underline{w}^T \underline{x}_n \underline{x}_n > 0$)
and approaching to b_n

Then $J(\underline{w})$ approaches to zero.

Comment: b_n need to be properly chosen.

Algebraic: Pseudoinverse learning algorithm
Gradient Descent: Widrow-Hoff



Algebraic Method

because $\frac{1}{N} \sum_{i=1}^N [\underline{w}^T \underline{x}_i - y_i]^2 = \frac{1}{N} \|\underbrace{\underline{X}}_{N \times (f+1)} \underbrace{\underline{w}}_{(f+1) \times 1} - \underbrace{\underline{y}}_{N \times 1}\|^2 = \frac{1}{N} (\underline{X} \underline{w} - \underline{y})^T (\underline{X} \underline{w} - \underline{y})$ [in OLS]

f : features number

$$\begin{aligned} \therefore J(\underline{w}) &= \frac{1}{N} \sum_{n=1}^N [\underline{w}^T \underline{z}_n \underline{x}_n - b_n]^2, \quad b_n > 0 \quad \forall n \\ &= \frac{1}{N} \|\underbrace{\underline{Z} \underline{X}}_{N \times (f+1)} \underbrace{\underline{w}}_{(f+1) \times 1} - \underbrace{\underline{b}}_{N \times 1}\|^2 = \frac{1}{N} (\underline{Z} \underline{X} \underline{w} - \underline{b})^T (\underline{Z} \underline{X} \underline{w} - \underline{b}) \end{aligned}$$

Minimize it :

$$\underline{\hat{w}} = (\underbrace{\underline{Z} \underline{X}^T}_{(f+1) \times (f+1)} \underbrace{\underline{Z} \underline{X}}_{(f+1) \times (f+1)})^{-1} \underbrace{\underline{Z} \underline{X}^T \underline{b}}_{(f+1) \times N}$$

if $\underline{Z} \underline{X}^T \underline{Z} \underline{X}$ nonsingular

where $\underline{Z} \underline{X} = \begin{bmatrix} z_1 x_{10} & z_1 x_{11} & \dots & z_1 x_{1f} \\ z_2 x_{20} & z_2 x_{21} & \dots & z_2 x_{2f} \\ \vdots & \vdots & \ddots & \vdots \\ z_N x_{N0} & z_N x_{N1} & \dots & z_N x_{Nf} \end{bmatrix}$

Gradient Descent Method

Ex: Sequential GD

$$J(\underline{w}) = \sum_{n=1}^N J_n(\underline{w})$$

$$\therefore J_n(\underline{w}) = \frac{1}{N} [\underline{w}^T \underline{z}_n \underline{x}_n - b_n]^2$$

$$\nabla_{\underline{w}} J_n(\underline{w}) = \frac{2}{N} [\underline{w}^T \underline{z}_n \underline{x}_n - b_n] \cdot \underline{z}_n \underline{x}_n$$

weight update

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) \cdot \frac{2}{N} [\underline{w}(i)^T \underline{z}_n \underline{x}_n - b_n] \cdot \underline{z}_n \underline{x}_n$$