# Ridge Regression

bias $\uparrow$  variance $\downarrow$

$\leftarrow$ under

$y = mx + b$     2D lines

m slope

$$L = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda m^2$$

$(m) <$

$\dfrac{\partial L}{\partial b} = 0 \rightarrow b$

$\dfrac{\partial L}{\partial m} = 0 \rightarrow m$

$$L = \sum_{i=1}^{n} (y_i - mx_i - b)^2 + \lambda m^2$$

$$b = \overline{y} - m\overline{x}$$

$\overline{y} \rightarrow$ y-mean

$\overline{x} \rightarrow$ x-mean

$m \rightarrow$ slope

$$L = \sum_{i=1}^{n} (y_i - mx_i - \overline{y} + m\overline{x})^2 + \lambda m^2$$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^{n} (y_i - mx_i - \overline{y} + m\overline{x})(-x_i + \overline{x}) + 2\lambda m = 0$$

$$= -2 \sum_{i=1}^{n} (y_i - \overline{y} - mx_i + m\overline{x})(x_i - \overline{x}) + 2\lambda m = 0$$

$$= \lambda m - \sum_{i=1}^{n} \left[(y_i - \overline{y}) - m(x_i - \overline{x})\right](x_i - \overline{x}) = 0$$

$\sum_{i=1}^{n} \quad \dots \quad + (x_i - \overline{x})^2 = 0$

$$= \lambda m - \sum_{i=1}^{n} {}^{I-1} (y_i - \bar{y})(x_i - \bar{x}) - m(x_1 - \bar{x})^2 = 0$$

$$= \lambda m - \sum_{i=1}^{n} (y_i - \bar{y})(x_i - \bar{x}) + m\sum_{i=1}^{n} (x_i - \bar{x})^2 = 0$$

$$= \lambda m + m\sum_{i=1}^{n} (x_i - \bar{x})^2 = \sum_{i=1}^{0} (y_i - \bar{y})(x_i - \bar{x})$$

$$= \boxed{m = \frac{\sum_{i=1}^{n} (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^{n} (x_i - \bar{x})^2 + \boxed{\lambda}}}$$

$\lambda \rightarrow$ hyperparalu. alpha

$\lambda = 0 \;\; = \lambda = 10, 100$

$$\boxed{m = \frac{\sum_{i=1}^{n} (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}}$$

$m \rightarrow$

$b = \bar{y} - m\bar{x}$

# Ridge Regression for 2D data

Thursday, June 3, 2021     6:38 AM

# Code

Thursday, June 3, 2021     6:39 AM

# Ridge Regression for nD data

$$\begin{array}{ccccc} w_0 & x_1 & x_2 & \cdots & x_n \end{array} \quad \boxed{y} \quad \rightarrow (n+1)$$

$m$ rows $\downarrow$

$w_1 \; w_2 \cdots w_n$

$\rightarrow$ Ridge

$$L = \sum_{i=1}^{m} (y_i - \hat{y}_i)^2$$

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

$m$ vals

$$y = \begin{bmatrix} - \\ - \\ - \\ - \end{bmatrix} \qquad W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \; (n \times 1) \qquad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1n} \\ 1 & x_{21} & x_{22} & \cdots & x_{2n} \\ 1 & \vdots & \vdots & \ddots & \vdots \\ \vdots & & & & \\ 1 & x_{m1} & x_{m2} & \cdots & x_{mm} \end{bmatrix}$$

Normal LR $\rightarrow$ Ridge

$$\boxed{L = (Xw - y)^T (Xw - y) + \lambda \|w\|^2}$$

$\lambda w_0^2 + \lambda w_1^2 + \lambda w_2^2 + \cdots + \lambda w_n^2$

$$\begin{bmatrix} w_0 & w_1 & w_2 & \cdots & w_n \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$W^T W \rightsquigarrow \lambda(w_0^2 + w_1^2 + w_2^2 \cdots w_n^2)$

$(a - b)^T = a^T - b^T$

$\dfrac{dL}{dw}$

$(ab)^T = b^T a^T$

$$L = \underline{(Xw - y)^T (Xw - y)} + \lambda W^T \underline{W}$$

$$L = \left[ (Xw)^T - (y)^T \right] (Xw - y) + \lambda W^T W$$

$$= (W^T X^T - y^T)(Xw - y) + \lambda W^T W$$

$$= W^T X^T X W - W^T X^T y - y^T X W + y^T y + \lambda W^T W$$

$$L = \underline{w^T X^T X \underline{w}} - 2 \underline{w^T} X^T y + \underline{y^T y} + \lambda w^T w$$

$$\frac{dL}{d\underline{w}} = \cancel{2} X^T X w - \cancel{2} X^T y + 0 + \cancel{2} \lambda w = 0$$

$$X^T X w + \textcircled{\lambda} \underline{w} = X^T y$$

$$(X^T X + \lambda \overset{\nearrow}{I}) w = X^T y$$

3 (4×4)

(n×1 , n×1)

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \qquad \boxed{w = (X^T X + \textcircled{\lambda I})^{-1} X^T y}$$

$$w = (X^T X)^{-1} X^T y \qquad \begin{bmatrix} \vert & \\ \vert & \end{bmatrix}$$

# Code

Vector form loss

$$L = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

$$L = (XW - y)^T (XW - y) + \lambda \|W\|^2$$

$$\boxed{L = (XW - y)^T (XW - y) + \lambda W^T W}$$

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} & X_{13} & \cdots & X_{1n} \\ 1 & X_{21} & X_{22} & X_{23} & \cdots & X_{2n} \\ 1 & X_{m1} & X_{m2} & X_{m3} & \cdots & X_{mn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$x_1 \, x_2 \cdots x_n \; \widehat{y}$

m rows

$\to (n+1)$

$W_0, W_1 \ldots W_n$ (parameters

$$W_0 = W_0 - \eta \frac{\partial L}{\partial W_0} \quad ; \quad W_1 = W_1 - \eta \frac{\partial L}{\partial W_1} \quad \cdots \quad W_n = W_n - \eta \frac{\partial L}{\partial W_n}$$

$$W_{new} = W_{old} - \eta \boxed{\frac{\Delta L}{\Delta W}} \to \text{gradient} \begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \partial L / \partial w_1 \\ \vdots \\ \partial L / \partial w_n \end{bmatrix}$$

$\textcircled{L}, \vec{w}$

$$L = \frac{1}{2}(XW - y)^T (XW - y) + \frac{1}{2}\lambda W^T W$$

$$= \frac{1}{2}(W^T X^T - y^T)(XW - y) + \frac{1}{2}\lambda W^T W$$

$$= \frac{1}{2}\left[ W^T X^T X W - W^T X^T y - y^T W X + y^T y \right] + \frac{1}{2}\lambda W^T W$$

$$= \frac{1}{2} \left[ w^T x^T x w \quad \cdots \right]$$

$$2 w^T x^T y$$

$$= \frac{1}{2} \left[ w^T x^T x w - \underbrace{2 y^T w x}_{} + y^T y \right] + \frac{1}{2} \lambda \, w^T w$$

$$x^T y$$

$$\frac{dL}{dw} = \frac{1}{2} \left[ 2 x^T x w - 2 \, y^T x \right] + \frac{1}{2} \cdot 2 \lambda w$$

$$x^T y$$

$$= \boxed{X^T X w - X^T X + \lambda w} = \frac{dL}{dw} \qquad \left( \frac{\Delta L}{x w} \right)$$

$$w_0 \quad w_1 \cdots \quad w_n$$

$$W = \left[ 0, 1 \cdots 1 \right]_{n \times m} \qquad starting$$

$$\underline{epochs}$$

$$W = W - \eta \frac{dL}{dw}$$

$$W \longrightarrow final \ answer$$

$$\underline{epoch \ times}$$

$$\hookrightarrow W = W - \eta \frac{dL}{dw}$$

$$\boxed{\frac{dL}{dw} = X^T X w - X^T y + \lambda w}$$

# Notes

Why is it called ridge

# 5 Key Understandings

$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \boxed{(\lambda \|w\|)^2}$$

Shrinkage coef

$$\lambda(w_1^2 + w_2^2 + \ldots + w_n^2)^2$$

Overfitting $\downarrow$
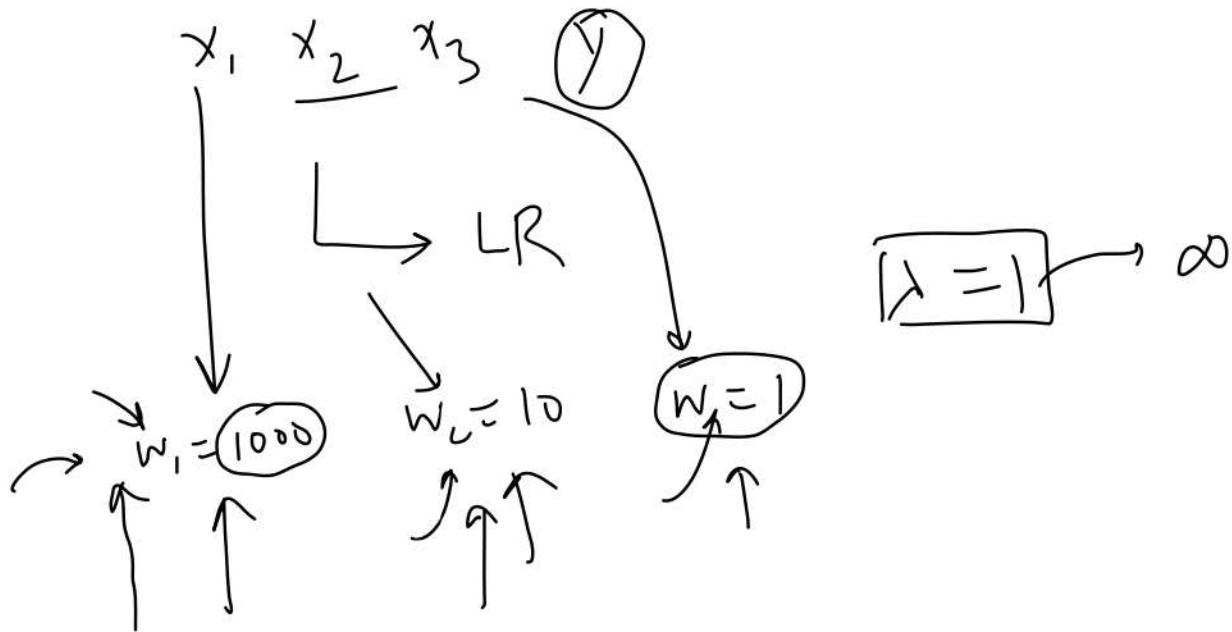
# 1. How the coefficients get affected?

$$\lambda \uparrow \qquad\qquad \lambda \rightarrow 0 \rightarrow \infty$$

$$\boxed{W_1, W_2 \ldots\ldots W_n} \leadsto \qquad \downarrow 0 \longrightarrow 0$$

# 2. Higher Values are impacted more

Never reaches 0

$$x_1 \quad x_2 \quad x_3 \qquad y$$

$$\rightarrow LR$$

$$w_1 = 1000$$

$$w_2 = 10$$

$$w_3 = 1$$

$$\lambda = 1 \rightarrow \infty$$

# 3. Bias Variance Tradeoff

Saturday, June 5, 2021    4:21 PM

Bias Variance          $(\lambda)$ ↓ ↑

$$\frac{\text{Bias} ↓ \quad \text{overfit Variance} ↑}{\text{Bias} ↑ \quad \text{underfitting Valaince} ↓}$$

← graph

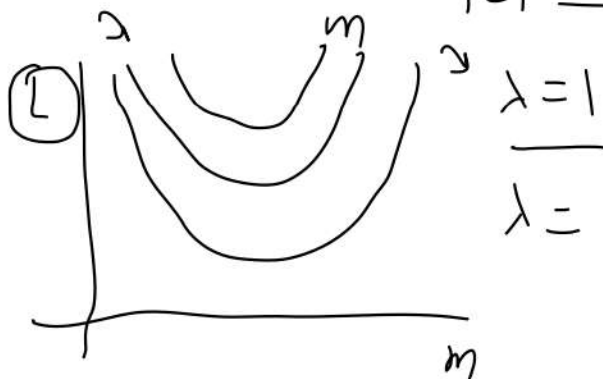0                    1000 $\lambda$

# 4. Impact on the Loss Function

Saturday, June 5, 2021    4:21 PM

$$\lambda \longrightarrow \qquad L = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \|w\|^2$$

$$X, \; y \rightarrow \boxed{m,b} \searrow \text{b is constant} \qquad \qquad b = -2.29$$

$$\downarrow \qquad \underline{b = 0}$$

$$m \qquad L = \sum_{i=1}^{n} (y_i - \underline{m} x_i)^2 + \lambda \underline{m}^2$$

$$\lambda = 1$$

$$\lambda =$$

$m$

# 5. Why called Ridge

Saturday, June 5, 2021     4:22 PM



Day 55 - Ridge Regression Page 279

# Practical Tip

Use ridge when there are more than 2 input cols

Ridge

$x_1 \; x_2 \; x_3 \cdots \quad \cdot \quad \cdot$

coefs

$>= 2$

$x_1 \quad y$