

★ Gradient Descent for Ridge Regression

$$L = (XW - Y)^T (XW - Y) + \lambda W^T W$$

$$L = \frac{1}{2} [W^T X^T - Y^T] (XW - Y) + \frac{1}{2} \lambda W^T W$$

$$= \frac{1}{2} [W^T X^T X W - \underbrace{W^T X^T Y - Y^T X W}_{\text{}} + Y^T Y]$$

$$+ \frac{1}{2} \lambda W^T W$$

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

$$\frac{\partial L}{\partial W} = \frac{1}{2} [\cancel{2 X^T X W} - \cancel{2 X^T Y}] + \frac{1}{2} \cancel{2 \lambda W}$$

$$\boxed{\frac{\partial L}{\partial W} = X^T X W - X^T Y + \lambda W}$$

$$W_{old} = W$$

$$\boxed{W_{new} = W_{old} - \eta \frac{\partial L}{\partial W}}$$

key points.

Page No.

Date

- 1) when we will increase $\lambda \uparrow$ all the coef. will start to shrink but never be zero
- 2) Coef having larger value will shrink fast as we increase the λ as compared to Coef having lower value.

- 3) Bias Variance trade-off

$\lambda \downarrow$

Bias \downarrow overfit Variance \uparrow

$\lambda \uparrow$

Bias \uparrow underfitting Variance \downarrow

- 4) Impact on loss function.

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

As we increase the λ Loss function moves towards the origin

- 5) why called Ridge.