

# Ridge Regression

~~Regularisation~~ Regularisation:

it is the technique in which we induce or add a information in the ML model. so that we can reduce overfitting

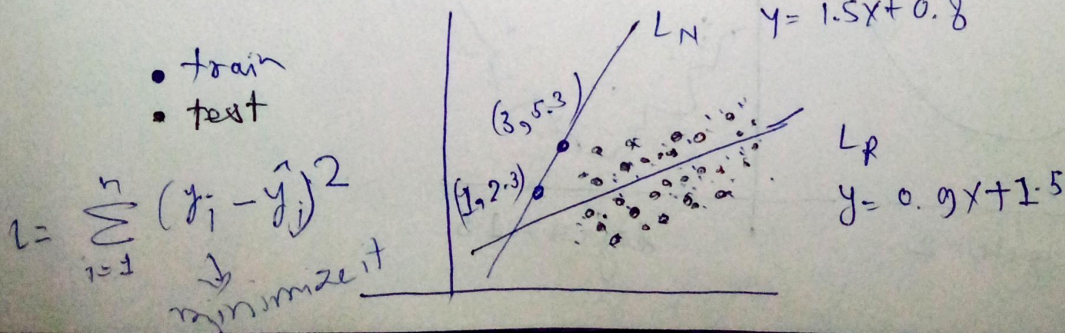
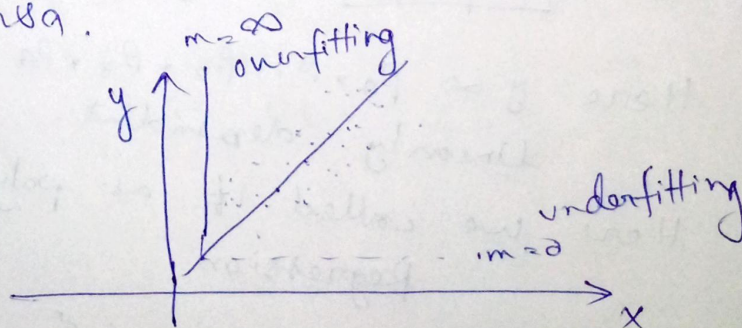
Regularisation techniques

1. Ridge ( $L_2$ ) -  $L_2$  Norm
2. LASSO ( $L_1$ ) -  $L_1$  Norm
3. ELastic Net (combination of  $L_1$  and  $L_2$ )

Overfitting:

Here the ML model exceptionally well on training data but not so well in testing data. it means the model variance is high.

in the linear regression a overfitted model. in which  $y = mx + b$ , here the  $m$  is very high and vice-versa.





in regularisation

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

$\lambda$  = hyperparameter (value from 0 to  $\infty$ )  
 $m$  = slope

loss $L_M$	loss $L_R$
$\lambda = 1$ $L = 0 + (1.5)^2$ $L = 2.25$	$\lambda = 1$ $(2.3 - 0.9 - 1.5)^2 + (5.3 - 2.7 - 1.5)^2$ $+ (0.9)^2$ $L = 2.03$

in 2D:

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

in 3D:

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

differentiate w.r.t  $b$  and  $m$  example

(i)  $\frac{\partial L}{\partial b} = 0$

↓  
get a  $b$

(ii)  $\frac{\partial L}{\partial m} = 0$

↓  
get a  $m$

(i) Here  $b = \bar{y} - m\bar{x}$   
 $\bar{y} = \frac{1}{n} \sum y_i$  mean  
 $m = \text{slope}$

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

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



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

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

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

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

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

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

in Ridge regression

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

in simple linear regression.

Here  $\lambda$  (alpha) is hyperparameter

if  $\lambda = 0$  (i) = (ii)



# Ridge Regression for nD Data:

Let us take an example.

$$x_1, x_2, \dots, x_n | y$$

with  $m$  rows. and weights are  $w_1, w_2, \dots, w_n$

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

$$= (XW - Y)^T (XW - Y)$$

$\hat{y}$  = output column. ( $m$  values)

$$W = [w_0, w_1, \dots, w_n]$$

$X$  = input column

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}_{n \times 1} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

\* Now the loss function:

$$L = (XW - Y)^T (XW - Y) \quad \text{--- for simple LR}$$

$$L = (XW - Y)^T (XW - Y) + \lambda \|W\|^2 \quad \text{--- for RR}$$

$$\because \lambda \|W\|^2 \equiv \lambda w_0^2 + \lambda w_1^2 + \lambda w_2^2 + \dots + \lambda w_n^2$$

or

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

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

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

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

$$L = W^T X^T X W - 2 W^T X^T Y + Y^T Y + \lambda W^T W$$

Squaring coefficients



$$\frac{\partial L}{\partial W} = 2X^T X W - 2X^T Y + 0 + 2\lambda W = 0$$

$$X^T X W - \cancel{X} W = X^T Y$$

$$(X^T X + \lambda I) W = X^T Y$$

$$W = (X^T X + \lambda I)^{-1} X^T Y$$

Ridge Regression using Gradient Descent:

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

in vector format

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

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

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$w_0, w_1, w_2, \dots, w_n$  (parameters)

$$w_0 = w_0 - \eta \frac{\partial L}{\partial w_0} ; \quad w_1 = w_1 - \eta \frac{\partial L}{\partial w_1} ; \quad w_n = w_n - \eta \frac{\partial L}{\partial w_n}$$

or



~~ΔL~~  

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\Delta L}{\Delta w}$$

here  $\frac{\Delta L}{\Delta w} \rightarrow$  gradient

$$\frac{\Delta L}{\Delta w} = \begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix}$$

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

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

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

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

$$\frac{\partial L}{\partial w} = \frac{1}{2} [2 X^T X W - \cancel{2 X^T Y} + 0] + \frac{1}{2} \times 2 \lambda W$$

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

$$W = \cancel{[w_0, w_1, \dots, w_n]} \rightarrow \text{starting}$$

in epochs

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