

# Lasso Regularization

## L1 Penalty

As we know that in Ridge regg.

$$L = \text{MSE} + \text{Ridge term}$$

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

$\downarrow$   
 $\lambda (w_1^2 + w_2^2 + \dots + w_n^2)$

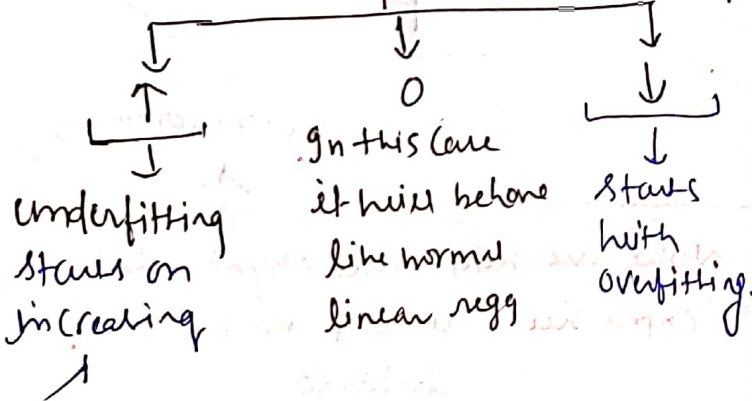
For Lasso we add L1 Penalty term as  $\rightarrow$   $\lambda \times$  Norm of (Coefficient Vector)

$\Downarrow$   
 $\lambda \|w\|$

then After applying Lasso loss function becomes.

$$L = \text{MSE} + \lambda \|w\|$$

$\downarrow$   
 $\lambda (|w_1| + |w_2| + |w_3| + \dots)$



$\Rightarrow$  what is benefit of Lasso regularization if it's making some of feature's co-eff. as zero?

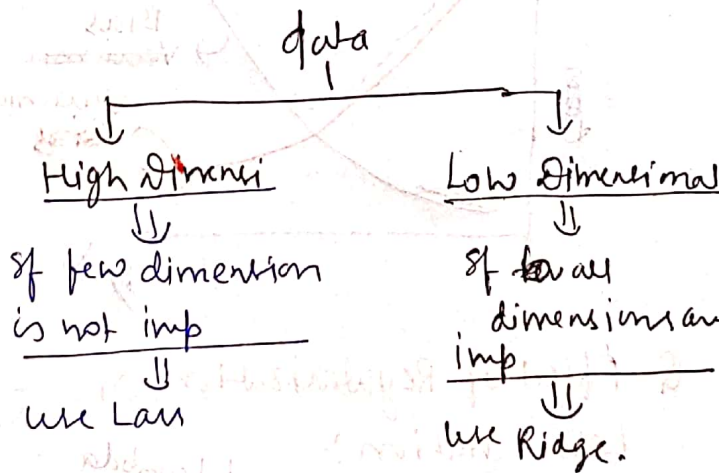
If we are working on Higher Dimensional data

$\rightarrow$  there would be chance of  $\times$  overfitting  $\rightarrow$  bcoz there would be some (value) of Co-eff. for each feature.

$\Downarrow$   
when we apply Lasso ( $L_1$ ) then it will make Co-eff.  $\Rightarrow 0$  for the features which are not impor.

$\Downarrow$   
Nothing but this process is called Feature Selection

$\Rightarrow$  which one should we prefer? ridge or Lasso?



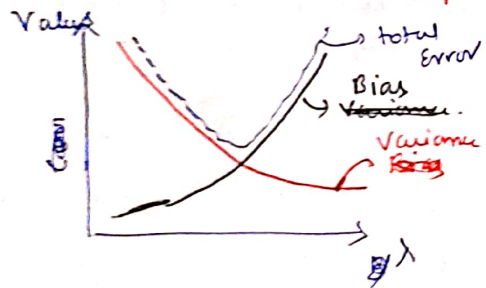
Q1. How Co-eff. are affected by increasing alpha value?

$\lambda \uparrow \rightarrow$  ~~Co-eff.~~ Co-eff.  $\downarrow$   
 $\downarrow$   
 for very high value Co-eff. becomes 0

for a certain value, it will make Co-eff. of non-important feature as 0.

$\downarrow$   
 that works as feature selection.

Q How  $\lambda$  in Lasso affect in Bias & Variance tradeoff



Q Effect of Regularization on Loss Function:-

of lambda

at certain value of  $\lambda$  loss becomes minimum but on increasing  $\lambda$  the minima of loss will shift towards higher value

Why Lasso creates sparsity?

Meaning of sparsity?

$\rightarrow$  on increasing  $\lambda$  some of the feature's Co-eff. becomes 0

If we remember in order to predict  $y$

$$Eq^n: y = mx + b$$

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

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

where  $\bar{y} = \text{mean}(y)$   
 $\bar{x} = \text{mean}(x)$

After applying ridge

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

Extra term  $\lambda$

Now, we will understand what extra we will add in Lasso?

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

$$m = ?$$

Loss func.:-  $\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda |m|$

in order to find minimum value of  $m \rightarrow$  we have to find  $\frac{\partial L}{\partial m}$   
 And we know that  $|m|$  mod function is not differentiable at 0.

let's take  $m > 0$

$$\text{Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda |m| \Rightarrow \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + \lambda m$$

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

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

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

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

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

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

$\therefore$  for Lasso

$\downarrow m > 0$

$\downarrow m = 0$

$\downarrow m < 0$

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

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

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

let take

$$\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) = yx$$

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

$$\text{then } m = \frac{yx - \lambda}{x^2}, \quad yx = 100, \quad x^2 = 50$$

$$m = 2, 9/5, 1, 0$$

$$\lambda = 0, 10, 50, 100, 150, \dots$$

but we cannot take  $m < 0$  in  $m > 0$  case, in this case automatically for  $\lambda = 150$ , it will go to this one return  $m = 0$



when  $m > 0$

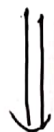
$$\lambda = 0, 10, 50, 100, 150$$

as we increase

$$\Rightarrow m = 2, 1/5, 1, 0, 5$$

on  $\lambda = 150$  only, it crosses 0 by using formula for  $m > 0$  it becomes ~~most~~ worst than when we started and it stopped to 0

Similarly



when

$$m < 0 \Rightarrow m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}$$

let's say to make  $m < 0 \Rightarrow$  we must have  $\sum (y_i - \bar{y})(x_i - \bar{x})$  to be negative.

let's

$$m = \frac{-100 + \lambda}{50}$$

$\lambda = 0$	$50$	$100$	$150$
$\downarrow$	$\downarrow$	$\downarrow$	
$m = -2$	$-1$	$0$	

After doing the  $m$  become 1 and we are calculating for  $m < 0$  So, it will not go beyond 0 So, it will stop.

why ridge regression do not create sparsity?

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2 + \lambda}$$

extra term.

$\therefore$  Here,  $\lambda$  is in denominator term so, it to being  $m = 0$ , depends on Numerator value, and Numerator value do not go to Zero (0) value. generally

That's why sparsity is in Lasso but not in Ridge.

where as

In Lasso,  $\lambda$  is in numerator term, So to being  $m = 0$ ,  $\lambda$  is in numerator plays important role.