

* Regularization *

Q. What is regularization and where should we use it?

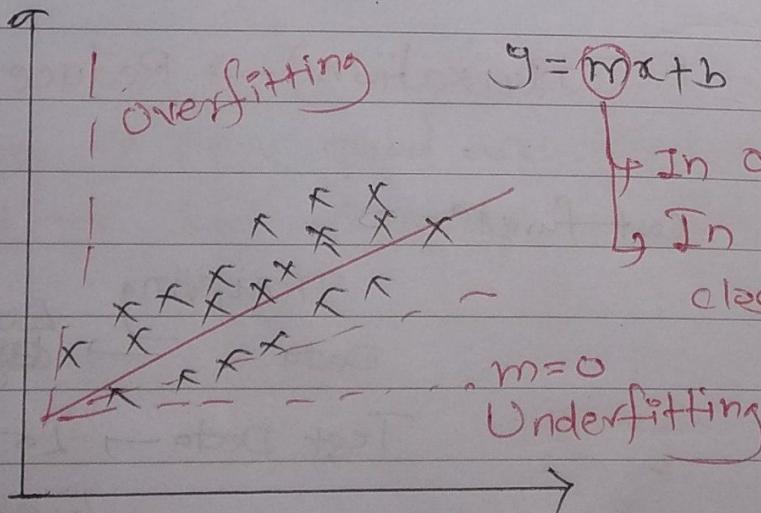
→ It is a technique that are used to calibrate machine Learning models in order to minimize the adjusted loss function and prevent overfitting or underfitting.

* Techniques to deal with overfitting

- ① Bagging
- ② Boosting
- ③ Regularization

* Overfitting:-

It is phenomenon or criteria where ML model performing exceptionally well on training data, but not so well on the testing data. (This means we have high variance in our ~~dataset~~ model.)



In overfitting m is high.
In Underfitting m is low or close to zero.

* Ridge - Regularization :-

Loss functn

$$L = \sum_{i=1}^n (\gamma_i - \hat{\gamma}_i)^2 + \lambda (\text{slope})^2 \rightarrow ①$$

Differentiate the above equⁿth w.r.t. b.

Let's rewrite the eqⁿ

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

$$\frac{\partial L}{\partial b} = \frac{\partial \sum (\gamma_i - mx_i - b)^2}{\partial b} + \lambda (m)^2 \underset{\text{const w.r.t } b}{=} 0$$

$$\Rightarrow 2\sum (\gamma_i - mx_i - b)(-1) + 0$$

$$\Rightarrow \sum_{i=1}^n \gamma_i - m \sum_{i=1}^n x_i - \sum_{i=1}^n b = 0$$

$$\Rightarrow n\bar{\gamma} - mn\bar{x} - nb = 0$$

$$\Rightarrow \bar{\gamma} - m\bar{x} - b = 0$$

$$\Rightarrow b = \bar{\gamma} - m\bar{x}$$

$\bar{\gamma}$ - mean of γ
 \bar{x} - mean of x
 m - slope.

Now, put this b in eqⁿ ①

$$\therefore L = \sum_{i=1}^n (\gamma_i - mx_i - \bar{\gamma} + m\bar{x})^2 + \lambda m^2$$

Differentiat w.r.t. m & equat it to zero

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

$$\Rightarrow 2 \sum_{i=1}^n (\gamma_i - mx_i - \bar{\gamma} + m\bar{x})(-x_i + \bar{x}) + 2m\lambda = 0$$

$$\Rightarrow -\lambda \sum_{i=1}^n (\gamma_i - \bar{\gamma} - m x_i + m \bar{x}) (x_i - \bar{x}) + \lambda m = 0$$

$$\Rightarrow \lambda m - \sum_{i=1}^n (\gamma_i - \bar{\gamma} - m x_i + m \bar{x}) (x_i - \bar{x}) = 0$$

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

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

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

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

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

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

~~if we ↑ λ then m will ↓~~

$$\text{W} \begin{array}{cccc|c} x_1 & x_2 & \cdots & x_n & | \gamma^{(n+1)} \\ \downarrow & \downarrow & & \downarrow & \\ \text{W}_1 & \text{W}_2 & \cdots & \text{W}_n & \end{array}$$

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

$$= (x \text{W} - \gamma)^T (x \text{W} - \gamma)$$

m values

$$\gamma = \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_m \end{bmatrix} \quad \text{W} = \begin{bmatrix} \text{W}_1 \\ \text{W}_2 \\ \vdots \\ \text{W}_n \end{bmatrix}_{n \times 1} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1n} \\ 1 & x_{21} & x_{22} & \cdots & x_{2n} \\ 1 & x_{31} & x_{32} & \cdots & x_{3n} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{m1} & x_{m2} & \cdots & x_{mn} \end{bmatrix}$$

$$L = (x\kappa - \gamma)^T (x\kappa - \gamma) + \lambda \|\kappa\|^2$$

$$\begin{aligned} & \lambda \kappa_0^2 + \lambda \kappa_1^2 + \lambda \kappa_2^2 + \dots + \lambda \kappa_n^2 \\ & \lambda [\kappa_0^2 + \kappa_1^2 + \kappa_2^2 + \dots + \kappa_n^2] \end{aligned}$$

We can rewrite as $\kappa^T \kappa$

$$[\kappa_0 \ \kappa_1 \ \kappa_2 \ \dots \ \kappa_n] \begin{bmatrix} \kappa_0 \\ \kappa_1 \\ \kappa_2 \\ \vdots \\ \kappa_n \end{bmatrix}$$

$$L = (x\kappa - \gamma)^T (x\kappa - \gamma) + \lambda \kappa^T \kappa$$

$$L = [(x\kappa)^T - (\gamma)^T] (x\kappa - \gamma) + \lambda \kappa^T \kappa \quad [\because (a-b)^T = a^T - b^T]$$

$$= (\kappa^T x^T - \gamma^T) (x\kappa - \gamma) + \lambda \kappa^T \kappa$$

$$= (\kappa^T x^T x\kappa - \kappa^T x^T \gamma - \gamma^T x\kappa + \gamma^T \gamma + \lambda \kappa^T \kappa)$$

both same

$$L = \frac{\kappa^T x^T x\kappa}{B} - 2 \frac{\kappa^T x^T \gamma}{B} + \frac{\gamma^T \gamma}{\text{const}} + \lambda \kappa^T \kappa$$

$$\frac{dL}{d\kappa} = f(x^T x\kappa) - f(x^T \gamma) + 0 + \cancel{f(\lambda \kappa)} = 0$$

$$\left\{ \begin{array}{l} \therefore x^T B x \rightarrow 2 B x \\ \therefore x^T B \rightarrow B \\ \therefore x^T x \rightarrow 2x \end{array} \right\} \rightarrow \text{vector derivative}$$

$$\Rightarrow x^T x\kappa + \lambda \kappa = x^T \gamma$$

$$\Rightarrow (x^T x + \lambda I)\kappa = x^T \gamma$$

$$\mathbf{w} = (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T \mathbf{y}$$

$$\text{In LR } \mathbf{w} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y} \rightarrow n \times 1, n \times 1$$

* Ridge with Gradient Descent :-

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

$$L = (\mathbf{x} \mathbf{w} - \mathbf{y})^T (\mathbf{x} \mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \vdots & & & & & \\ 1 & x_{m1} & x_{m2} & x_{m3} & \dots & x_{mn} \end{bmatrix}_{m \times n}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}_{n \times 1}$$

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

If we want to update

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

$$\dots, \quad w_n = w_n - \eta \frac{\partial L}{\partial w_n}$$

$$w_{\text{new}} = w_{\text{old}} - \eta \boxed{\frac{\partial L}{\Delta w}} \rightarrow \text{gradient}$$

$$\begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix}$$

Multiply $\frac{1}{2}$ through the loss and write it again

$$L = \frac{1}{2} (\mathbf{x} \mathbf{w} - \mathbf{y})^T (\mathbf{x} \mathbf{w} - \mathbf{y}) + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w}$$

$$= \frac{1}{2} (\mathbf{u}^T \mathbf{x}^T - \mathbf{y}^T) (\mathbf{x}^T \mathbf{u} - \mathbf{y}) + \frac{1}{2} \lambda \mathbf{u}^T \mathbf{u}$$

$\Leftrightarrow \frac{1}{2} (\mathbf{u}^T \mathbf{x}^T - \mathbf{y}^T)$

$$= \frac{1}{2} \left[\mathbf{u}^T \mathbf{x}^T \mathbf{x} \mathbf{u} - \mathbf{u}^T \mathbf{x}^T \mathbf{y} - \mathbf{y}^T \mathbf{x} \mathbf{u} + \mathbf{y}^T \mathbf{y} \right] + \frac{1}{2} \lambda \mathbf{u}^T \mathbf{u}$$

Same

$$= \frac{1}{2} \left[\mathbf{u}^T \mathbf{x}^T \mathbf{x} \mathbf{u} - 2 \cancel{\mathbf{u}^T \mathbf{x} \mathbf{x}} + \mathbf{y}^T \mathbf{y} \right] + \frac{1}{2} \lambda \mathbf{u}^T \mathbf{u}$$

$$\frac{dL}{du} = \frac{1}{2} \left[2 \mathbf{x}^T \mathbf{x} \mathbf{u} - \cancel{2 \mathbf{u}^T \mathbf{x} \mathbf{x}} \right] + \frac{1}{2} \lambda \mathbf{u}^T \mathbf{u}$$

$$\frac{dL}{du} = \mathbf{x}^T \mathbf{x} \mathbf{u} - \cancel{\mathbf{x}^T \mathbf{x}} + \lambda \mathbf{u} = \frac{dL}{du}$$

Starting pt $\mathbf{u}_0 = [0, 1, 1, \dots, 1]$

epochs $\mathbf{u}_t = \mathbf{u}_{t-1} - \eta \frac{dL}{du}$

$$\boxed{\frac{dL}{du} = \mathbf{x}^T \mathbf{x} \mathbf{u} - \mathbf{x}^T \mathbf{y} + \lambda \mathbf{u}}$$

⊗ 5 key points Regarding Ridge Reg.

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

$$\rightarrow \lambda (u_1^2 + u_2^2 + \dots + u_n^2)$$

{ shrinkage
of }

↳ Reduce overfitting.

(*) How the coef gets affected if $\lambda \uparrow$

$k_1, k_2, \dots, k_n \rightarrow$ coef get closer to 0, but not will be 0.

② Higher values are impacted more.

$x_1, x_2, x_3 \uparrow$

↳ Normal LR

$$k_1 = 1000 \quad k_2 = 10 \quad k_3 = 1$$

In Ridge $\lambda \rightarrow \infty$

$k_1 = 1000 \downarrow$ rapidly

$k_2 = 10 \downarrow$ slowly than k_1 ,

$k_3 = 1 \downarrow$ slowly than k_1 and k_2

③ Bias Variance Tradeoff

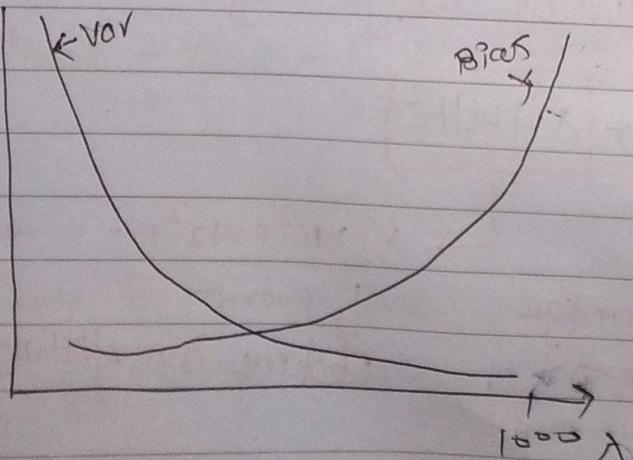
Bias & variance both depends on the value of ①

If $\lambda \rightarrow 0$ then,

Bias \downarrow Overfitting Variance \uparrow

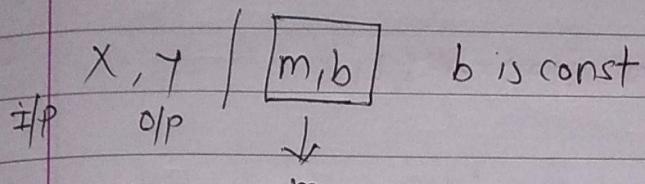
If $\lambda \rightarrow \infty$ then,

Bias \uparrow Underfitting Variance \downarrow



④ Impact on the Loss function of λ ?

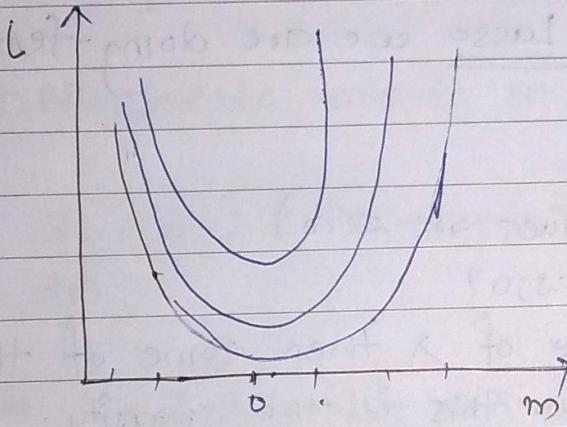
$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\|^2$$



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

where $b=0$ & slope = m

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



$\lambda = 1$
change λ

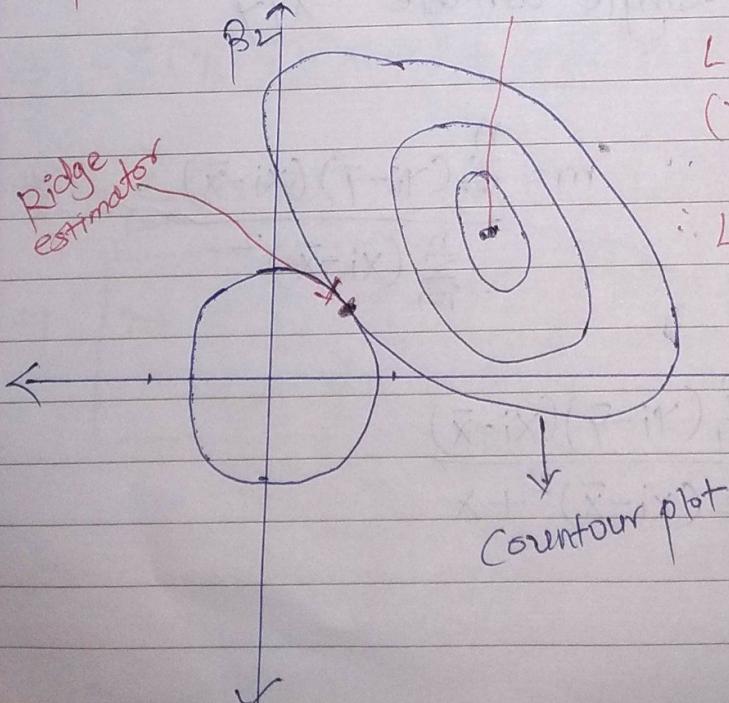
The graph will shrink
and minima shift towards
0.

⑤ Why called Ridge.

OLS estimator

we need to study

Hard constrain Ridge Reg.



$$L = RSS + \lambda \|w\|^2$$

$$(y_i - \hat{y}_i)^2 \quad \beta_0, \beta_1, \beta_2$$

$$\therefore L = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2})]^2 + \lambda (\beta_1 + \beta_2)^2$$

Contour plot

Practical Tip

Apply when more than 2 I/p variables.

* Lasso Regression $\hat{Y}(L_1)$

Loss function

$$L = \text{MSE} + \lambda ||\boldsymbol{\beta}||$$

$$\hookrightarrow [|\beta_1| + |\beta_2| + \dots + |\beta_n|]$$

- * Which one is most preferable among Lasso and Ridge.
 → Lasso → Because for high value of λ it will automatically bring the coef zero whose corresponding variable are less important. It will reduce the dimension of the data that with the help of Lasso we are doing feature selection

* Lasso sparsity (Some values are zero)

* What is sparsity in Lasso?

- If we increase the value of λ then some of the coefficients will be zero. This is the sparsity.

* Why sparsity in Lasso and why not in Ridge?

- Let's consider in single variable X/Y

$$Y = mx + b$$

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

$$\bar{y} = \text{mean}(Y)$$

$$\bar{x} = \text{mean}(X)$$

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

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

Lasso case $\Rightarrow b = \bar{y} - m\bar{x}$

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

$$L = \sum (\gamma_i - mx_i - b)^2 + \lambda |m| \quad [\because \hat{\gamma}_i = x_i m + b]$$

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

$$\left\{ \begin{array}{l} \therefore b = \bar{y} - m\bar{x} \end{array} \right.$$

Let's consider $m > 0$

$$\therefore |m| = m$$

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

Differentiate w.r.t. m & equat to 0.

$$\frac{dL}{dm} = 2 \sum (\gamma_i - mx_i - \bar{y} + m\bar{x}) (-x_i + \bar{x}) + 2\lambda = 0$$

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

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

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

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

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

Lasso

for $m=0$ for $m < 0$ for $m > 0$

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

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

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

for $m > 0 \wedge \lambda > 0$

$$m = \frac{yx - \lambda}{x^2} \quad \left\{ \begin{array}{l} yx = 100 \\ x^2 = 50 \end{array} \right.$$

$$m = \frac{100 - \lambda}{50} \quad \lambda = 0 \quad \left| \begin{array}{c} \lambda = 10 \\ m = 9 \end{array} \right. \quad \left| \begin{array}{c} \lambda = 50 \\ m = 1 \end{array} \right. \quad \left| \begin{array}{c} \lambda = 100 \\ m = 0 \end{array} \right.$$

~~now~~

$$\text{if } \lambda > 100 \Rightarrow m = -1$$

Now we need to use $m < 0$ formula. Let's take $\lambda = 150$

$$\therefore m = \frac{yx + \lambda}{x^2} = \frac{100 + 150}{50} = \frac{250}{50} = 5$$

2, 9, 15, 1, 0, -1, 5.

-ve

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

 $\lambda > 0$

$$m = \frac{-100 + \lambda}{50} \quad \lambda = 0 \quad \left| \begin{array}{c} \lambda = 50 \\ m = -1 \end{array} \right. \quad \left| \begin{array}{c} \lambda = 100 \\ m = 0 \end{array} \right. \quad \left| \begin{array}{c} \lambda = 150 \\ m = 1 \end{array} \right.$$

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

Ridge

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

\rightarrow This is denom term.

In order to become $m = 0$ num should be zero.

But here this is not possible. That is why

In Ridge coeff become close to zero, but not exactly 0. whereas, In Lasso the λ is numerators term that is why in Lasso coeff value will be zero.

* Elasticnet Regression *

* Its a combination of Ridge and Lasso.

$$L = \sum_{i=1}^n (\gamma_i - \hat{\gamma})^2 + \lambda \|k\|^2$$

Ridge

$$L = \sum_{i=1}^n (\gamma_i - \bar{\gamma})^2 + \lambda \|k\|$$

Lasso

$$\|k_1\| + \|k_2\| + \dots + \|k_n\|$$

$$k_1^2 + k_2^2 + \dots + k_n^2 \leftarrow$$

(Overfitting Reduce)

$\lambda \uparrow (\infty) \quad k \downarrow (0)$ but not 0

L₂ norm

Use where All col are imp.

$$L, \lambda \uparrow \quad k \rightarrow 0$$

feature selection

L₁ norm

Use where some of cols imp.

⊗ if we have big data and there are 1000 of cols in this case we cannot predict which cols are imp. so we use Elasticnet Reg here.

$$L = \sum_{i=1}^n (\gamma_i - \hat{\gamma}_i)^2 + a \|k\|^2 + b \|k\|$$

λ, a, b

$$\left\{ \begin{array}{l} \text{In sklearn } \lambda = a+b \quad \text{l1-ratio} = \frac{a}{a+b}, \\ \lambda = \frac{a}{a+b} \end{array} \right. \quad \left| \begin{array}{l} a = \lambda / \lambda \\ b = \lambda - a \end{array} \right.$$

By default $\lambda = 1, \text{l1-ratio} = 0.5$

⊗ if there is multicollinearity in input cols then we should use Elasticnet Reg

e.g. x_1 (height) | x_2 (weight)