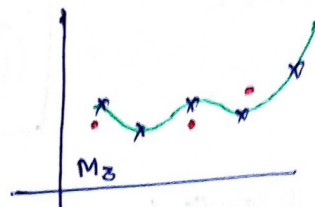
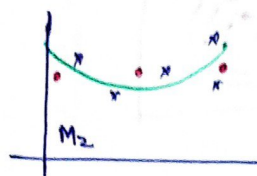
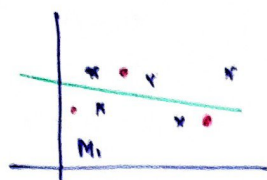


Bias variance Trade off

11:58 PM
Sat Apr



• test
x train

Bias

Inability of machine learning model to capture the relationship in the training data

	M_1	M_2	M_3
Bias	High	mid	low
Variance	low	Mid	High
	underfitting		overfitting

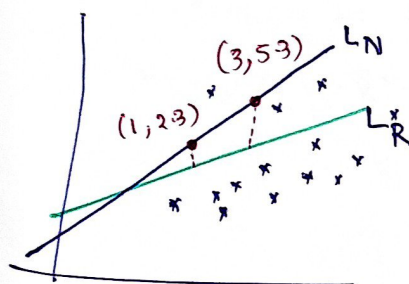
variance \rightarrow train data error - test data error

Sweet spot

1. Regularization
2. Bagging
3. Boosting

Ridge Regularization

P-1



$$L_N \quad y = 1.5x + 0.8$$

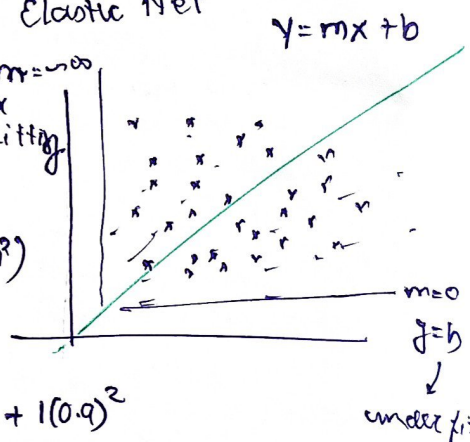
$$L_R \quad y = 0.9 + 1.5x$$

- ① Ridge (L2)
- ② Lasso (L1)
- ③ Elastic Net

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

$$= \sum (y_i - mx_i - b)^2 + \lambda(m^2)$$

hyperparameter $m = \infty$
yax overfitting



Loss L_N

$$\lambda = 1$$

line pass through both point so zero

$$0 + (1.5)^2 = 2.25$$

Loss L_R

$$\lambda = 1$$

$$= (2.3 - 0.9 - 1.5)^2 + (5.3 - 0.9 - 1.5)^2 + 1(0.9)^2$$

$$= (0.1)^2 + (1.4)^2 + (0.9)^2$$

$$= 2.03$$

Now ML can understand L_R how less loss so it choose L_R .

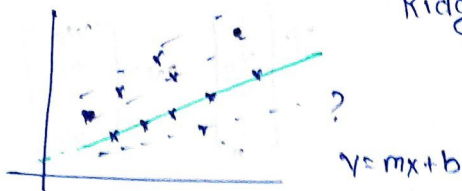
more axis.

$$\lambda(m_1^2 + m_2^2 + \dots + m_n^2)$$

we want to say our ML focus on other data bec. that is best fit line

② How can we do

- we add some extra data
- calculate loss



$$J = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda m^2 \quad \text{no change}$$

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

$$\begin{cases} \frac{\partial J}{\partial b} = 0 \\ \frac{\partial J}{\partial m} = 0 \end{cases}$$

affix $\frac{\partial J}{\partial b} = 0$

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

add $m \rightarrow 0$
here is changes.

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

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

$$= -2 \sum (y_i - \bar{y} - mx_i + m\bar{x})(x_i - \bar{x}) + 2\lambda m = 0$$

$$= \lambda m - \sum [(y_i - \bar{y}) - m(x_i - \bar{x})](x_i - \bar{x}) = 0$$

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

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

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

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

Enter

simple LR

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

Ridge Regression for n Data

$$\begin{matrix} & x_1 & x_2 & \dots & x_n \\ \text{rows} \downarrow & w_1 & w_2 & \dots & w_n \end{matrix}$$

$$J = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (XW - Y)^T (XW - Y)$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

$$W^T W = [w_0, w_1, w_2, \dots, w_n] \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$J = (XW - Y)^T (XW - Y) + \lambda \|W\|^2 \quad \{ \lambda w_0^2 + \lambda w_1^2 + \dots + \lambda w_n^2 \}$$

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

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

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

$$= W^T \underbrace{X^T X}_{\text{const}} W - 2W^T X^T Y + Y^T Y + \lambda W^T W$$

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

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

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

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

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

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

matrix derivative

$$I_{(n+1) \times (n+1)}$$

simple LR

$$W = (X^T X)^{-1} X^T Y$$

$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

solver = 'cholesky'

Ridge Regression using Gradient Descent

P3

$$J = \sum (y_i - \hat{y}_i)^2 + \lambda \|w\|^2$$

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

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

$$= \frac{1}{2} [W^T X^T X W - W^T X^T Y - Y^T X W + Y^T Y] + \frac{1}{2} \lambda W^T W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

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

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

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

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nn} \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$w_0 = w_0 - \eta \frac{\partial L}{\partial w_0}$$

$$w_1 = w_1 - \eta \frac{\partial L}{\partial w_1}$$

$$w_n = w_n - \eta \frac{\partial L}{\partial w_n}$$

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

$$\sum (y_i - \hat{y}_i)^2 + \lambda \|w\|^2$$

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

(shrinkage coef)

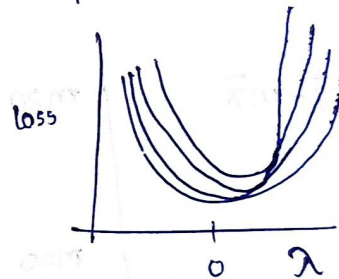
1. Affect of coefficient in shrink
tends to 0.

2. Higher values are impacted more

x_1	x_2	x_w	y
w_1	w_2	w_3	
1000	10	1	

This shrinks more

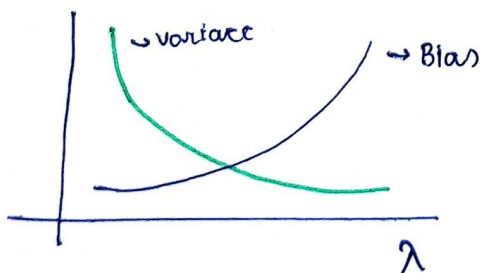
4. Impact on loss func $\lambda \uparrow$ coef \downarrow



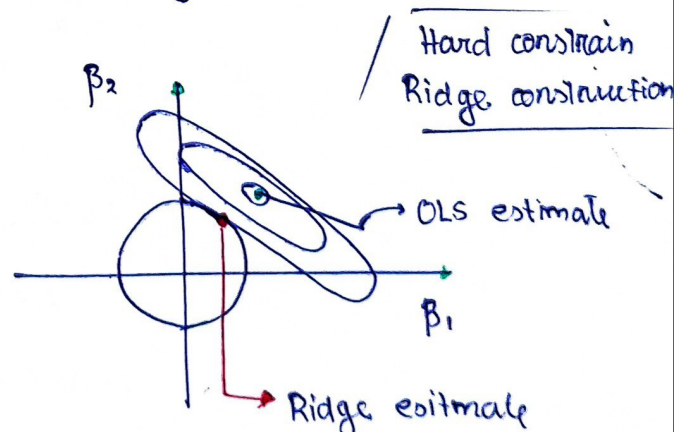
- it reduce & shrink
- then shift towards to zero

3. Bias variance trade off

$$\lambda \propto \frac{1}{\text{Bias}} \propto \text{variance}$$



5. why 'called Ridge



Lasso Regression

0056

6:25pm
Sat Apr 26

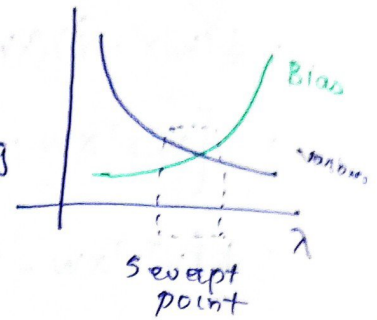
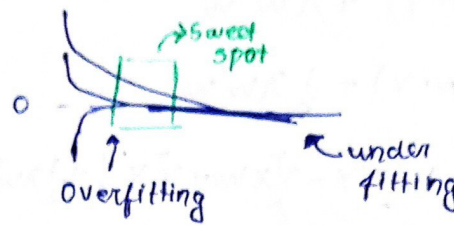
$$J = \text{MSE} + \lambda |w|$$

$$\lambda [|w_1| + |w_2| + |w_3| + \dots + |w_n|]$$

absolute

λ coefficient will zero in lasso

- Benefits less imp. features is zero.
- Apply the feature selection
- dimensionality reduction.



Sparsity

Simple

$$x/y \rightarrow y = mx + b$$

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

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

Ridge

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

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

Lasso

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

$$m \geq 0$$

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

$$m = 0$$

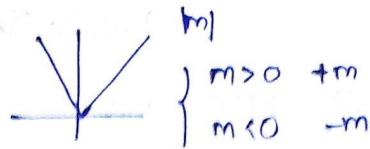
$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$m \leq 0$$

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

$$\frac{\partial J}{\partial m} = \sum (y_i - mx_i - \bar{y} + m\bar{x})^2 + 2\lambda|m|$$

to make simple



$$\begin{aligned} &= \sum (y_i - mx_i - \bar{y} + m\bar{x})^2 + 2\lambda m \\ &= 2\sum (y_i - mx_i - \bar{y} + m\bar{x})(-x_i + \bar{x}) + 2\lambda \\ &= -2\sum [(y_i - \bar{y}) - m(x_i - \bar{x})](x_i - \bar{x}) + 2\lambda \\ &= -\sum [(y_i - \bar{y})(x_i - \bar{x})] + m\sum (x_i - \bar{x})^2 + \lambda \end{aligned}$$

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

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

?

- why at the end all values are zero?
- why not values are -ve?
- why sparsity?

because Ridge λ is den so ~ 0 .
Lasso λ is neu so $= 0$.

let $m > 0$

$$\frac{K - \lambda}{K'}$$

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

$$K = \sum (y_i - \bar{y})(x_i - \bar{x})$$

$$K' = \sum (x_i - \bar{x})^2$$

$$K = 100$$

$$K' = 50$$

$\lambda = 0$	$\lambda = 10$	$\lambda = 50$	$\lambda = 100$	$\lambda > 100$
$m = 2$	$m = \frac{9}{5}$	$m = 1$	$m = 0$	$m = -1$

Then it use

$$\frac{K + \lambda}{K'}$$

$$\downarrow$$

$$\frac{100 + 100}{50} = 4$$

$m = 4$

$m \longrightarrow$

2 $\frac{9}{5}$ 1 0 (-1)

\downarrow

4

we want decrease the m
so why use
we stop after zero.

Elastic Net Regression

Ridge	$J = \text{MSE} + \lambda \ W\ ^2$	all IP in imp.
Lasso	$J = \text{MSE} + \lambda W $	use feature selection.
EN	$J = \text{MSE} + a \ W\ ^2 + b W $	multicollinearity

$$J = \sum (y_i - \bar{y}_i)^2 + a \|W\|^2 + b |W|$$

default

$$\lambda = 1$$

$$\alpha = 0.5$$

$$a = 0.5$$

$$b = 0.5$$

In sklearn

$$\lambda = a + b$$

$$\alpha\text{-ratio} = \frac{a}{a+b}$$

$$a = \alpha \lambda$$

$$b = \lambda - a$$

if $\alpha = 0.9$

90% Ridge 10% Lasso