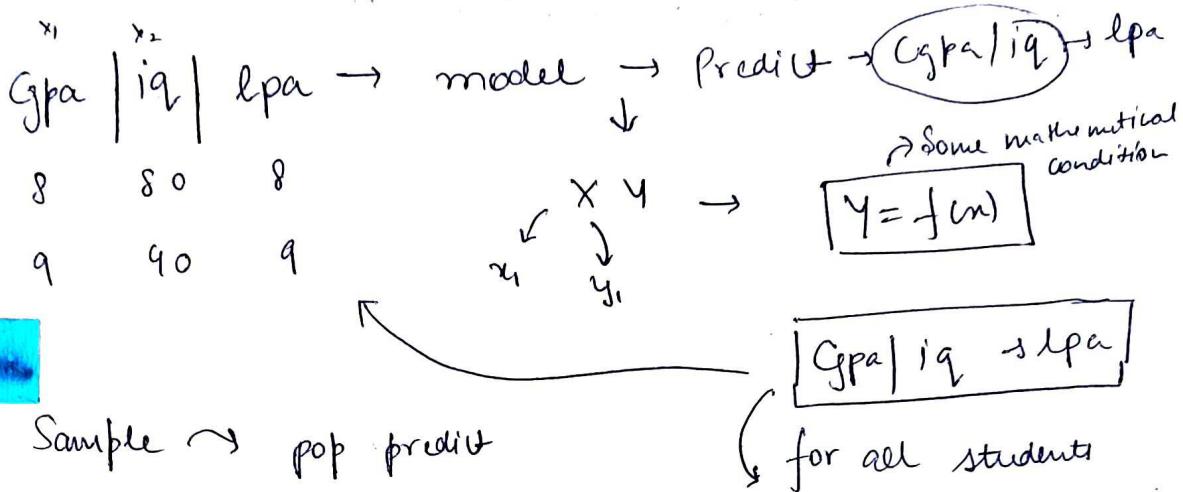


Regularization

The hidden truth (Trailer)

1000 std



$$y = f(x) + \text{error}$$

ml model

Irreducible error

$$y = f(x) \rightarrow \text{pop}$$

$$\hat{y} = f'(x)$$

$$f(x) - f'(x)$$

Reducible error

Bias Variance Trade off

reducible error

$$\text{Reducible error} = \text{Bias} + \text{Var}$$

Bias Variance Trade off

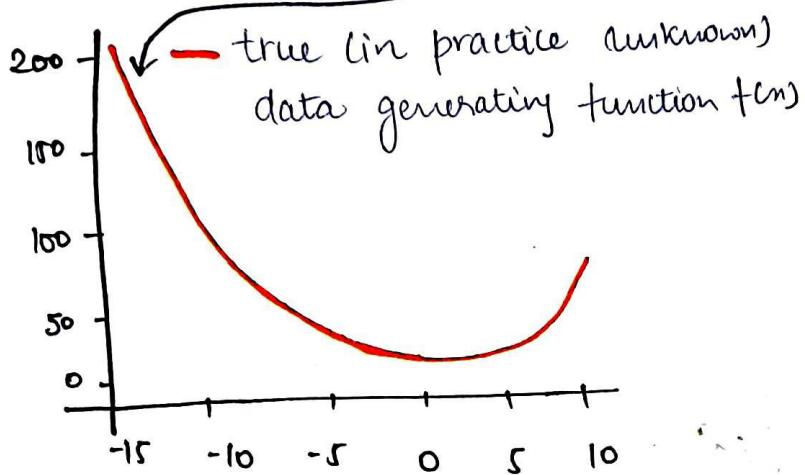
generally, Sample data $\xrightarrow{\text{ML Model}}$ Predict for population

but we do reverse

pop \rightarrow Sample

$$Y = f(x) = x^2 \quad [-15, 10] \quad \text{range}$$

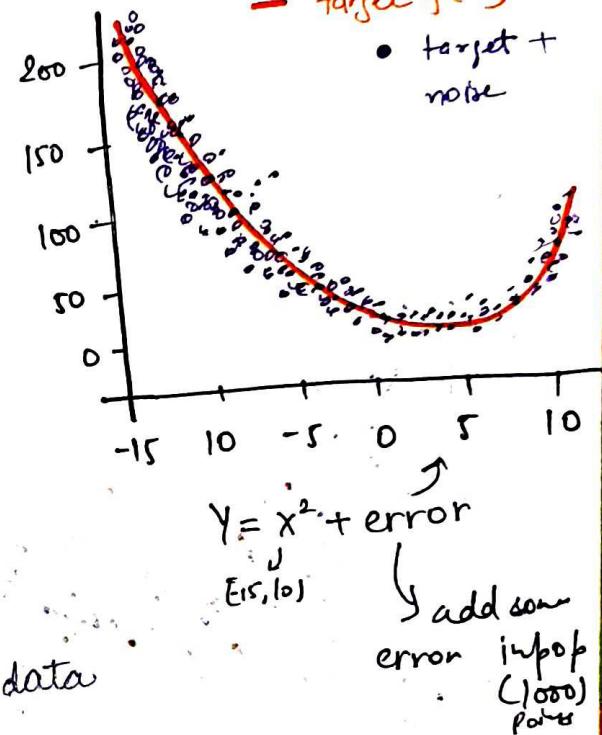
* In real world, we can't find this because we'll have sample data.



$$Y = x^2$$

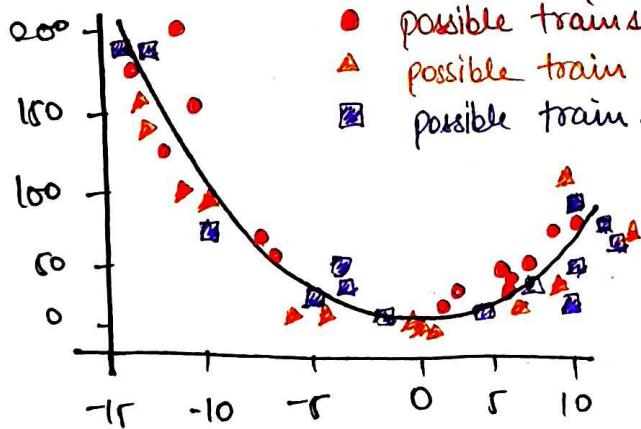
- target $f(x)$

- target + noise



Draw 3 random samples data

- true function $f(x)$
- possible train set 1 (sample 1)
- ▲ possible train set 2 (sample 2)
- possible train set 3 (sample 3)

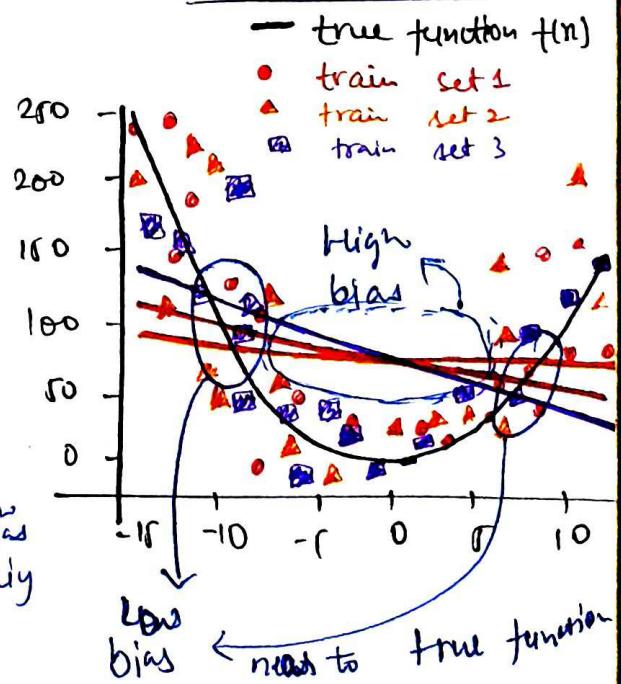


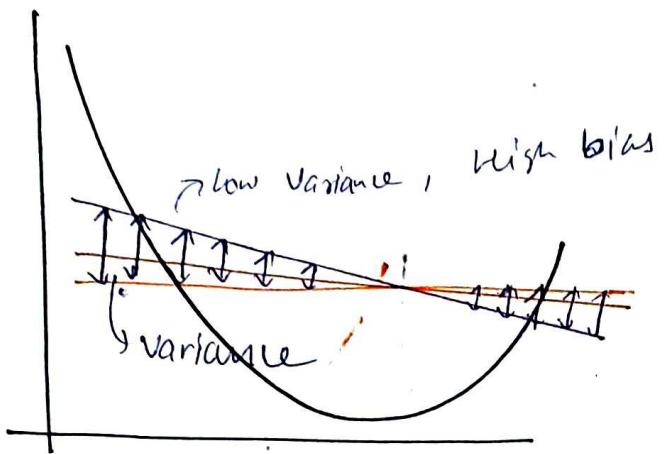
all three samples try to predict like true function.

Bias → The inability of a ML model to fit the training data (i) high bias (ii) low bias

Variance → ML model predict when training data is change.

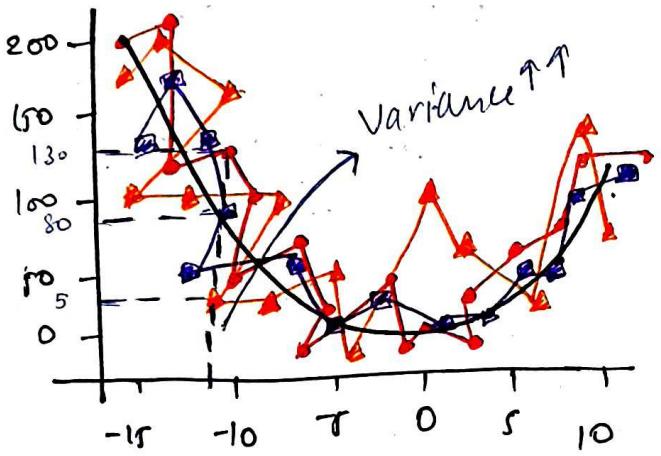
fitting Linear Regression





(1, 2, or 3)

Polynomial High degree



* High bias ✗

* Low bias ✓

* High Variance ✗

High Variance is closely related overfitting

High bias is closely related to underfitting

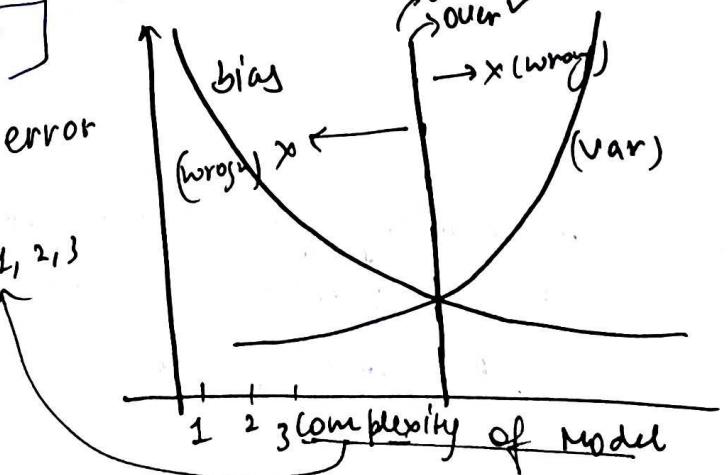
- * ideally output what we want is
 - * low bias → good result in training data
 - * low variance → change data but result is same

But not possible

[The "trade-off" in bias-variance trade-off refers to the fact that minimizing bias will usually increase variance versa.]

{ bias
variance }

Polynomial
degree → 1, 2, 3



Some Question

1. How would you define bias and variance mathematically?
2. How is bias and variance related to overfitting and underfitting mathematically?
3. Why is there a tradeoff between bias and variance mathematically?

Expected Value and Variance

Expected value represents the average outcome of a random variable over a large number of trials or experiments.

In a simple sense, the expected value of a random variable is the long term average value of repetitions of the experiment it represents. For example, the expected value of rolling a six-side die is 3.5 because over many rolls, we would expect to average about 3.5.

A die roll 1000 times

$X = \text{rolling die}$

$$E[X] = x_1 P(x_1) + x_2 P(x_2) + x_3 P(x_3) + x_4 P(x_4) + x_5 P(x_5) + x_6 P(x_6)$$

$$\begin{array}{c}
 X \xrightarrow{x_1 \rightarrow 1} \\
 \downarrow \quad \downarrow \quad \downarrow \\
 x_1 \quad x_2 \xrightarrow{x_2 \rightarrow 2} \\
 \downarrow \quad \downarrow \quad \downarrow \\
 x_3 \quad x_4 \xrightarrow{x_3 \rightarrow 3} \\
 \downarrow \quad \downarrow \quad \downarrow \\
 x_5 \quad x_6 \xrightarrow{x_4 \rightarrow 4} \\
 \downarrow \quad \downarrow \quad \downarrow \\
 \vdots \quad \vdots \quad \vdots
 \end{array}
 \quad E[X] = x_1 p(x_1) + x_2 p(x_2) + \dots + x_n p(x_n)$$

$$= 1 \times \frac{1}{6} + 2 \times \frac{1}{6} + 3 \times \frac{1}{6} + 4 \times \frac{1}{6} + 5 \times \frac{1}{6} + 6 \times \frac{1}{6}$$

$$= \frac{1+2+3+4+5+6}{6} = \frac{21}{6} = 3.5$$

Overall avg. mean

example

$$\boxed{5 \ 3 \ 4 \ 5 \ 3 \ 5} \rightarrow \text{mean} \downarrow$$

$$\hookrightarrow \frac{5+3+4+5+3+5}{6} = \frac{25}{6} = 4.166\overline{6}$$

* another method

$$\frac{3(5) + 2(3) + 1(4)}{6}$$

* ~~another~~ we can
also write like this

$$\frac{3}{6}(5) + \frac{2}{6}(3) + \frac{1}{6}(4)$$

* mean \rightarrow expected value

Absolute Random Variable $\rightarrow (X)$

$$E[X] = \sum_{i=1}^n x_i p(x_i)$$

Continuous Random Variable

$$E[X] = \int x_i f(x_i) dx$$

expected value $\xrightarrow{\text{estimate}}$ pop mean

$\text{Var}(X) \rightarrow \text{var of pop}$

$$\boxed{\text{Var}(X) = E[X^2] - (E[X])^2}$$

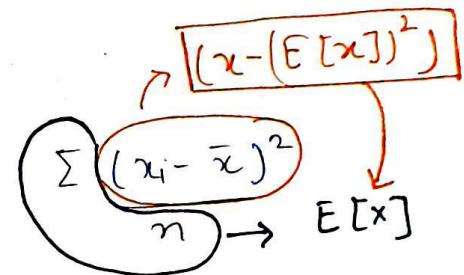
Derivation

$$\text{Var} = \frac{\sum (x_i - \bar{x})^2}{n}$$

\nwarrow sample mean

$$\boxed{E[(x - E[X])^2]} = \text{Var}(x)$$

\nwarrow pop mean



$$\begin{aligned}
 &= E[x^2 + (E[X])^2] - 2xE[X] \\
 &= E[x^2] + E[(E[X])^2] - E[2xE[X]] \\
 &= E[x^2] + E[(E[X])^2] - E[2]E[X]E[E[X]] \\
 &= E[x^2] + E[(E[X])^2] - 2E[X]E[X] \\
 &= E[x^2] + E[(E[X])^2] - 2(E[X])^2 \\
 &= E[x^2] + (E[X])^2 - 2(E[X])^2 \\
 &= E[x^2] - E[X]^2
 \end{aligned}$$

Random variable
 Constant
 Constant mean
 given X and Y independent
 Constant

$\therefore E[\text{constant}] = \text{constant}$

$$\boxed{\text{Var}(X) = E[X^2] - E[X]^2} \quad \text{hence proved}$$

True for both discrete Random Variable and continuous random variable.

* Another formula is

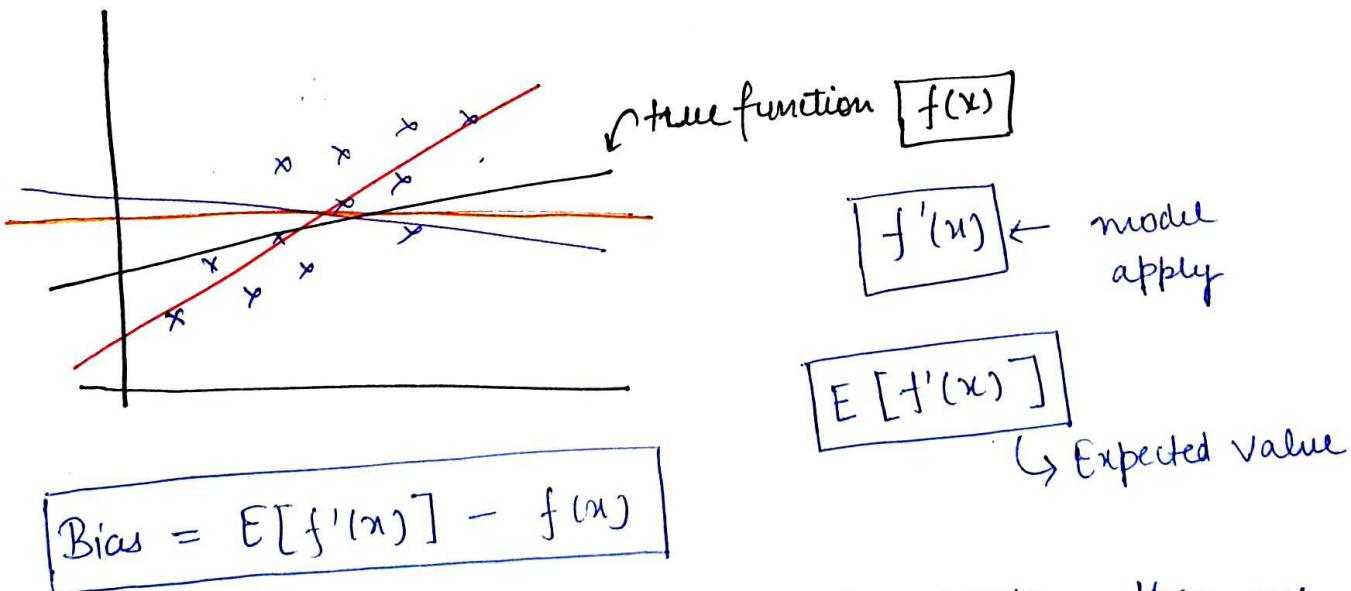
$$\text{Var}(x) = E[(x - E[x])^2]$$

we can use both formula.

What exactly are Bias and Variance Mathematically

Bias

In the context of Machine Learning and statistics, bias refers to the systematic error that a model introduce because it cannot capture the true relationship in data. It represents the difference between the expected prediction of our model and the correct value which we are trying to predict. More bias leads to underfitting, where the model does not fit the training data well.



- * If I am using Linear Regression Model, then we can calculate coefficient and intercept mean and compare these coefficient and intercept.

$$\boxed{E[f'(m)]} \xrightarrow[\text{way}]{\text{Linear Regression}} E[m] \quad \begin{matrix} \hookrightarrow \text{coefficient} \\ \hookrightarrow \text{generally formula} \end{matrix}$$

$$f(x) - f(m) = 0 \rightarrow \text{bias}$$

* If bias is zero then this is called unbiased predictor

$\boxed{\text{bias}=0}$

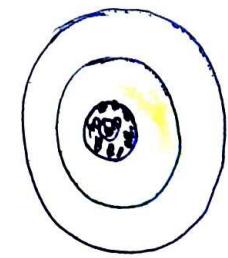
- \hookrightarrow unbiased predictor
- \hookrightarrow which means close to true function
(pop function)

Variance

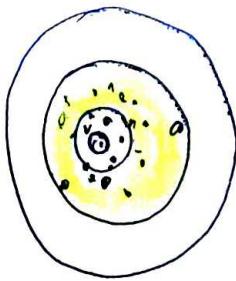
In the context of machine learning and statistics, variance refers to the amount by which the prediction of our model will change if we used a different training data set. In other words, it measures how much the predictions for a given point vary between different realizations of the models.

$$\boxed{\text{Var}(f'(x)) = E[(f'(m) - E[f'(x)])^2]}$$

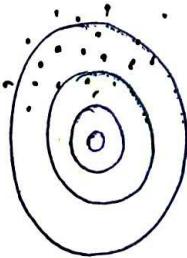
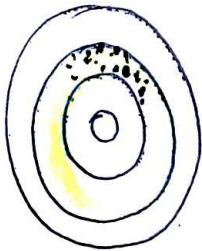
Low Bias



High Variance



High bias



Bias Variance Decomposition

Bias-variance decomposition is a way of analyzing a learning algorithm's expected generalization error with respect to a particular problem by expressing it as the sum of three very different quantities: bias, variance and irreducible error.

1. Bias: This is the error from erroneous assumption in the learning algorithms. High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting).

2. Variance: This is the error from sensitivity to small fluctuation in the training set. High variance can cause an algorithm to model the random noise in the training data, rather than the intended outputs (overfitting).

3. Irreducible Error: This is the noise term. This part of the error is due to the inherent noise in the problem itself, and can't be reduced by any model.

$$\text{Loss} = (\text{bias})^2 + \text{Variance} + \text{Var}(E)$$

reducible error

irreducible error
↳ epsilon

* irreducible error
 mean = 0
 var = σ^2
 constant assume

Derivation

$$mse = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n} = E[(y - \hat{y})^2]$$

$y = f(x) + \epsilon$ just change sign
 we can also write $y = \theta + \epsilon$

$$\hat{y} = f'(x) = \hat{\theta}$$

$$mse = E[(\theta + \epsilon - \hat{\theta})^2]$$

$$= E\left[\left(\underbrace{\theta - \hat{\theta}}_a + \underbrace{\varepsilon}_b\right)^2\right]$$

$$= E\left[\frac{(\theta - \hat{\theta})^2}{a^2} + \frac{\varepsilon^2}{b^2} + \frac{2\varepsilon(\theta - \hat{\theta})}{2ab}\right] \quad \therefore (a+b)^2 = a^2 + b^2 + 2ab$$

$$= E[(\theta - \hat{\theta})]^2 + E[\varepsilon^2] + \underbrace{E[2\varepsilon(\theta - \hat{\theta})]}_{E[x+y] = E[x] + E[y]}$$

$$= E[(\theta - \hat{\theta})]^2 + E[\varepsilon^2] + \underbrace{E[2]\ E[\varepsilon]\ E[\theta - \hat{\theta}]}_{E[x]\ E[y]} \quad \therefore E[xy] = E[x]E[y]$$

$$= E[(\theta - \hat{\theta})]^2 + E[\varepsilon^2] + 0 \quad \therefore E[\varepsilon] = 0 \text{ because we assume irreducible error mean is 0.}$$

$$\therefore \text{Var}(\varepsilon) = \sigma^2 = E[(\varepsilon - E[\varepsilon])^2]$$

$$\text{Var}(\varepsilon) = E[(\varepsilon - 0)^2] = E[\varepsilon^2]$$

$$mse = \underbrace{E[(\theta - \hat{\theta})^2]}_{\text{add}} + \text{Var}(\varepsilon)$$

$$E[(\theta - \hat{\theta})^2] = E\left[\underbrace{(\theta - E[\hat{\theta}])}_a + \underbrace{(E[\hat{\theta}] - \hat{\theta})}_b\right]^2 \quad \therefore \text{add and subtract with } E[\hat{\theta}]$$

Ques 2

$$E[(\theta - E[\hat{\theta}])^2 + (E[\hat{\theta}] - \hat{\theta})^2 + 2(\theta - E[\hat{\theta}])(E[\hat{\theta}] - \hat{\theta})]$$

$$\therefore (a+b)^2 = a^2 + b^2 + 2ab$$

$$= E[(\theta - E[\hat{\theta}])^2] + E[(E[\hat{\theta}] - \hat{\theta})^2] + E[2(\theta - E[\hat{\theta}])(E[\hat{\theta}] - \hat{\theta})]$$

$\therefore E[x+y] = E[x]+E[y]$

$$\therefore E[2(\theta - E[\hat{\theta}])(E[\hat{\theta}] - \hat{\theta})]$$

$$E[2] E[(\theta - E[\hat{\theta}])] E[(E[\hat{\theta}] - \hat{\theta})]$$

$$2 (\theta - E[\hat{\theta}]) \{ E[E[\hat{\theta}]] - E[\hat{\theta}] \}$$

$$2 (\theta - E[\hat{\theta}]) E[\hat{\theta}] - E[\hat{\theta}] = 0$$

$$= E[(\theta - E[\hat{\theta}])^2] + E[(E[\hat{\theta}] - \hat{\theta})^2] + \underset{\text{Variance}}{\downarrow} \text{Var}(\epsilon)$$

$$(\theta - E[\hat{\theta}])^2$$

$$(\theta - E[\hat{\theta}])^2$$

$$(\text{bias})^2$$

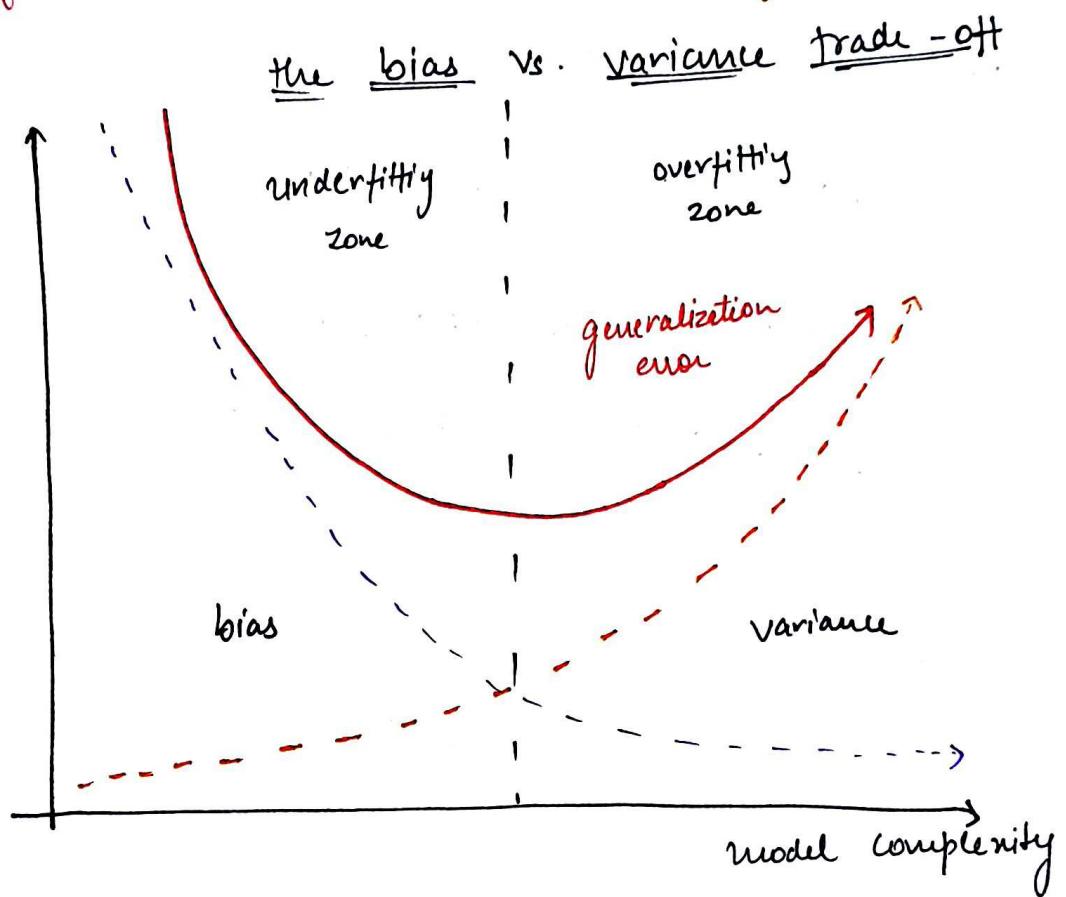
$$\text{mse} = (\text{bias})^2 + \text{Variance} + \text{irreducible error}$$

Hence proved

reducible
error

$$\text{mse} = \text{reducible error} + \text{irreducible error}$$

- * for improving bias we can use different ML algo.
or reduce
- * for reduce variance use Regularization.



When to Use Regularization?

1. Preventing Overfitting: Regularization is most commonly used as a tool to prevent overfitting. If your model performs well on the training data but poorly on the validation or test data, it might be overfitting and regularization could help.
2. High Dimensionality: Regularization is particularly useful when you have a high number of features compared to the number of data points. In such scenarios, models tend to overfit easily, and regularization can help by effectively reducing the complexity of the model.
3. Multicollinearity: When features are highly correlated (multicollinearity), it can destabilize your model and make the model's estimate sensitive to minor changes in the model. L₂ regularization (Ridge regression) can help in such cases by disturbing the coefficient estimate among correlated features.

: 4. Feature Selection: If you have a lot of features and you believe many of them might be irrelevant; L1 regularization (Lasso) can help. It tends to produce sparse solutions, driving the coefficients of irrelevant features to zero, thus performing feature selection.

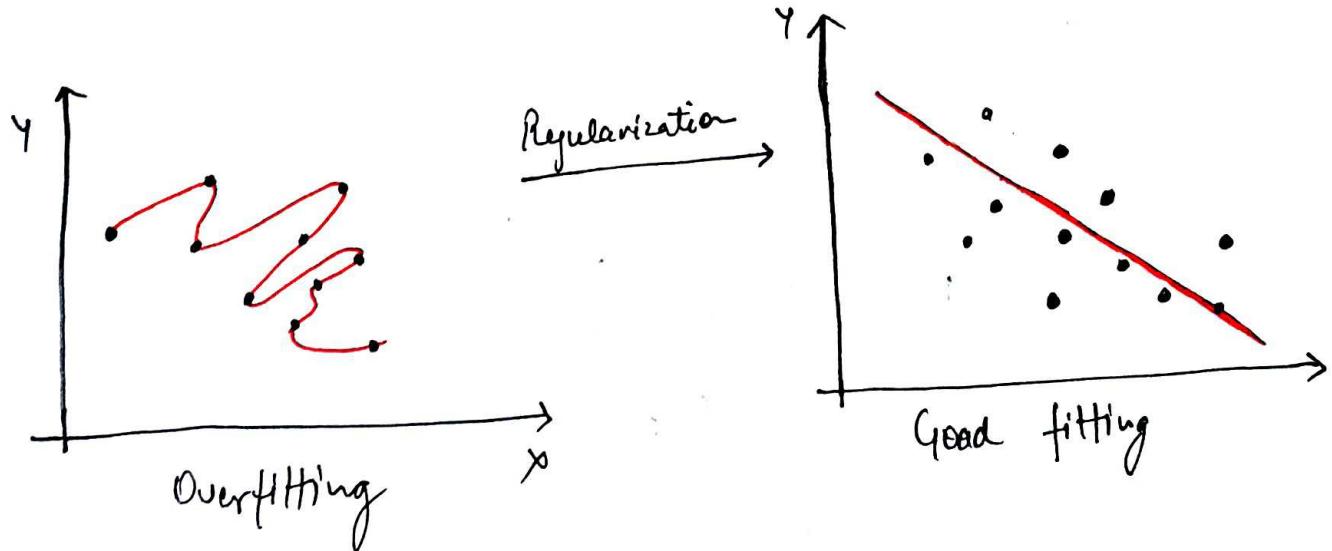
5. Interpretability: If model interpretability is important and you want a simpler model, regularization can help achieve this by constraining the model's complexity.

6. Model Performance: Even if you're not particularly worried about overfitting, regularization might still improve your model's out-of-sample prediction performance.

Regularization

- Regularization is one of the ways to improve our model to work on unseen data by ignoring the less important features.
- Regularization minimizes the validation loss and tries to improve the accuracy of the model.
- It avoids overfitting by adding a penalty to the model with high variance, thereby shrinking the beta coefficients to zero.

Regularization refers to techniques that are used to calibrate machine learning models in order to minimize the adjusted loss function and prevent overfitting or underfitting.



Ridge Regularization (L_2)

Also known as Ridge Regression, it modifies the over-fitted and under-fitted models by adding the penalty equivalent to the sum of the squares of the magnitude of coefficient.

Ridge This means that the mathematical function representing our machine learning model is minimized and coefficients are calculated. The magnitude of coefficient is squared and added. Ridge Regression performs regularization by shrinking the coefficients present. The function depicted below shows the cost function of ridge regression.

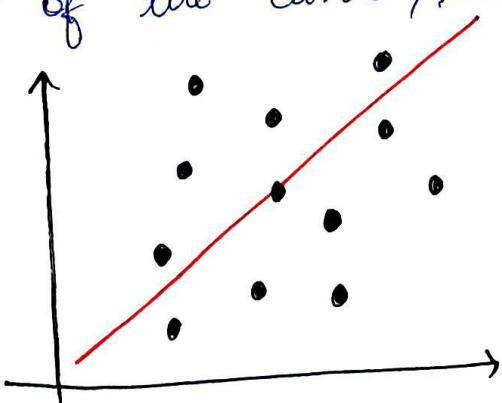
$$\text{cost function} = \text{Loss} + \lambda \times \sum \|w\|^2$$

Here:

Loss = sum of the squared residual

λ = Penalty for the errors

w = slope of the curve / line



In the cost function, the penalty term is represented by Lambda λ . By changing the values of the penalty function, we are controlling the penalty term. The higher the penalty, it reduce the magnitude of coefficients. It shrink the parameter. Therefore, it is used to prevent multicollinearity and it reduce the model complexity by coefficient shrinkage.

* Overfitting \rightarrow Linear Regression

$$y = mx + b$$

\hookrightarrow value of m is high in overfitting

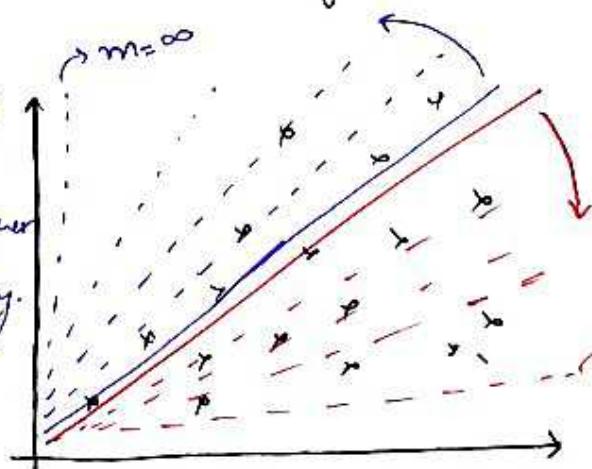
* Underfitting \rightarrow Linear Regression

$$y = mx + b$$

\hookrightarrow value of m is low in underfitting

* for calculate y

x is more important than other
so it is overfitting.

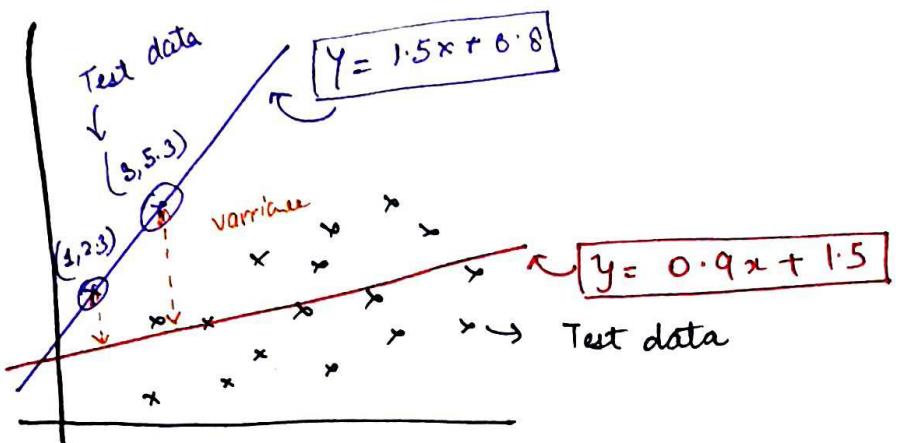


$$y = 0(x) + b$$

$$y = b$$

underfitting

- * less value less value
- x than y to calculate y



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

- * Train data \rightarrow 2 points and set purple best fit line
- * Test data \rightarrow more than train data and predict red best fit line
- * Red Variance ($\uparrow\uparrow$) between Train best fit line and Test best fit line.

Loss Ln (Train) $\lambda = 1$
 residual error = 0 (only two points)

$$0 + 1(1.5)^2$$

$$\boxed{\text{Loss} = 2.25}$$

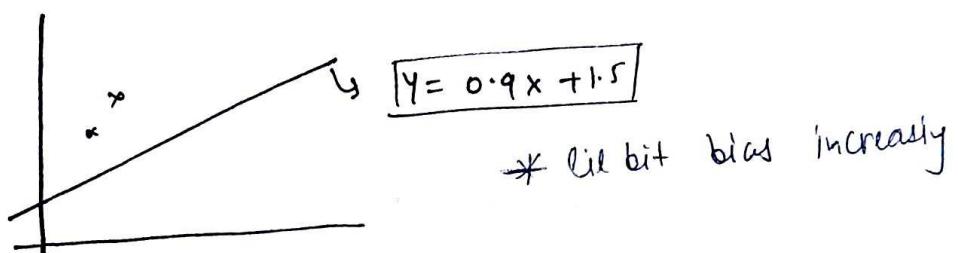
Loss Lf (Test)
 $\lambda = 1$

$$(y_i - \hat{y}_i)^2 + (y_i - \hat{y}_i)^2 + \lambda(m^2)$$

$$(2.3 - (0.9(1) + 1.5))^2 + (5.3 - (0.9(3) + 1.5))^2 + (0.9)^2$$

$$\boxed{\text{Loss} = 2.03}$$

* choose low loss value ($y = 0.9x + 1.5$) for Train and Test.



$L_2 \rightarrow$ because $\lambda(m_1 + m_2 + m_3)$
for multi feature

Derivation

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

$\frac{\partial L}{\partial b} = 0 \rightarrow b$

$\frac{\partial L}{\partial m} = 0 \rightarrow m$

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

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

where $\bar{y} \rightarrow y\text{-mean}$
 $\bar{x} \rightarrow x\text{-mean}$
 $m \rightarrow \text{slope}$

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

differentiate with respect to m

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

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

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

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

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

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

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

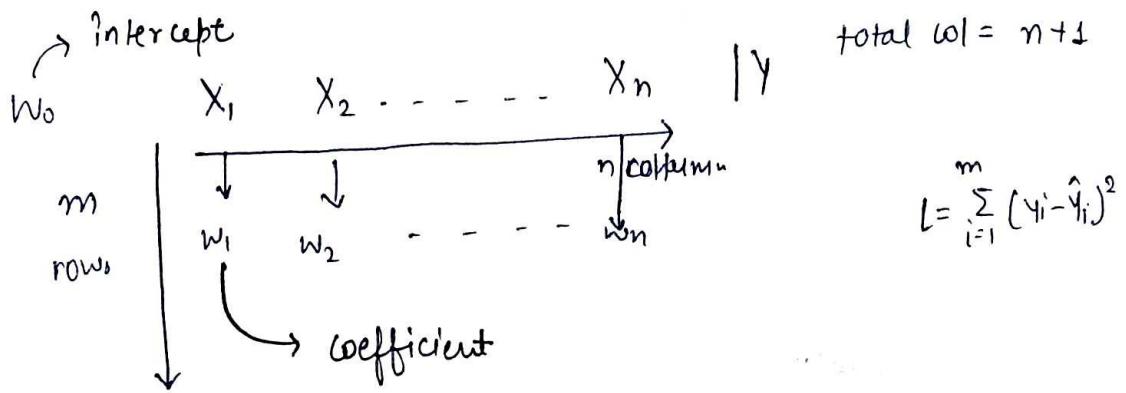
$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$

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

code

hyper parameter
(alpha)

Ridge Regression for nD data



Matrix form

$$= (Xw - y)^T (Xw - y)$$

$$y = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \\ \vdots \end{bmatrix}$$

m value

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$(n \times 1)$

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

$$L = (Xw - y)^T (Xw - y) + \lambda \underbrace{\|w\|^2}_{w^T w} \left\{ \begin{array}{l} \lambda w_0^2 + \lambda w_1^2 + \lambda w_2^2 + \dots \\ \dots + \lambda w_n^2 \\ \lambda (w_0^2 + w_1^2 + \dots + w_n^2) \end{array} \right.$$

$$\begin{bmatrix} w_0 & w_1 & w_2 & \dots & w_n \end{bmatrix} \rightarrow \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$L = (Xw - y)^T (Xw - y) + \lambda w^T w$$

$$L = \left[(x_w)^T - (y)^T \right] (x_w - y) + \lambda w^T w \quad \therefore (a-b)^T = a^T - b^T$$

$$= (N^T x^T - y^T) (x_w - y) + \lambda w^T w$$

$$= w^T x^T w x - w^T x + y - y^T x_w + y^T y + \lambda w^T w$$

↑
Same Proved in Multiple Linear Reg.

$$= w^T \overbrace{x^T x}^{\text{constant}} w - 2 w^T x^T y + y^T y + \lambda w^T w$$

differentiate with $\frac{dL}{dw}$

$$\frac{dL}{dw} = 2 x^T w - 2 x^T y + 0 + 2 \lambda w = 0$$

$$x^T x w + \lambda w = x^T y$$

$$(x^T x + \lambda I) w = x^T y$$

Identity matrix because we cannot extract w directly so we use identity matrix

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

if 3 col then (4×4)
 $n \times 1, n \times 1$

code

Vector derivatives

$$x^T B \rightarrow B$$

$$x^T b \rightarrow b$$

$$x^T x \rightarrow 2x$$

$$x^T B x \rightarrow 2Bx$$

Ridge Regression Using Gradient Descent

Vector form Loss

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

$$\therefore L = (xw - y)^T (xw - y) + \lambda \|w\|^2$$

$$L = (xw - y)^T (xw - y) + \lambda w^T w$$

$$\begin{matrix} & x_1 & x_2 & \dots & x_n & | & y \\ m \text{ rows} & \downarrow & & & & & \end{matrix}$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & x_{m3} & \dots & x_{mn} \end{bmatrix}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \rightarrow (n+1)$$

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

↓ learning rate

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

$w_{\text{new}} = w_{\text{old}} - \eta \boxed{\frac{\Delta L}{\Delta w}}$ → gradient $\begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix}$

$$L = (xw - y)^T (xw - y) + \lambda w^T w$$

$$L = (w^T x^T - y^T) (xw - y) + \lambda w^T w$$

Multiply $\frac{1}{2}$ both equations

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

$$L = \frac{1}{2} [w^T x^T x w - w^T x^T y - y^T w x + y^T y] + \frac{1}{2} \lambda w^T w$$

same proved in linear regression gradient descent

$$= \frac{1}{2} [w^T x^T x w - 2x^T w y + y^T y] + \frac{1}{2} \lambda w^T w$$

$$\frac{dL}{dw} = \frac{1}{2} [2x^T x w - 2x^T y] + \frac{1}{2} 2\lambda w$$

$$= x^T x w - x^T y + \lambda w = \frac{dL}{dw} \left(\frac{\Delta L}{\Delta w} \right)$$

$$w = [w_0, w_1, \dots, w_n] \quad \text{Starting}$$

Epochs

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

$$\boxed{\frac{dL}{dw} = X^T X w - X^T y + \lambda w}$$

5 key Understanding

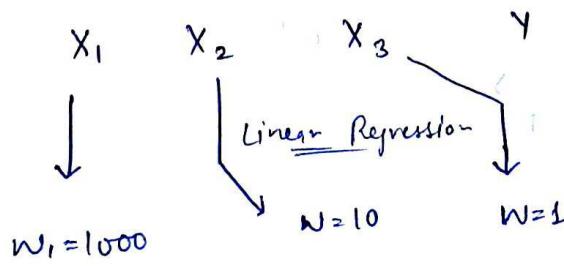
- How the coefficient get affected?

$$\lambda \uparrow \uparrow \quad \lambda \rightarrow 0 \rightarrow \infty$$

$$[w_0, w_1, \dots, w_n] \downarrow \downarrow \quad \text{Shrinks}$$

Code

- Higher Variance are impacted more



$$\lambda = 1 \rightarrow \infty$$

* w_1 play max role to find or predict y

* If λ value increase then w_1 value highly decrease because w_1 have large value.

Code

3. Bias - Variance trade off

Bias and Variance both depend on λ

$\lambda \downarrow = \text{Bias} \downarrow \text{ overfit Variance} \uparrow$

$\lambda \uparrow = \text{Bias} \uparrow \text{ Underfit Variance} \downarrow$

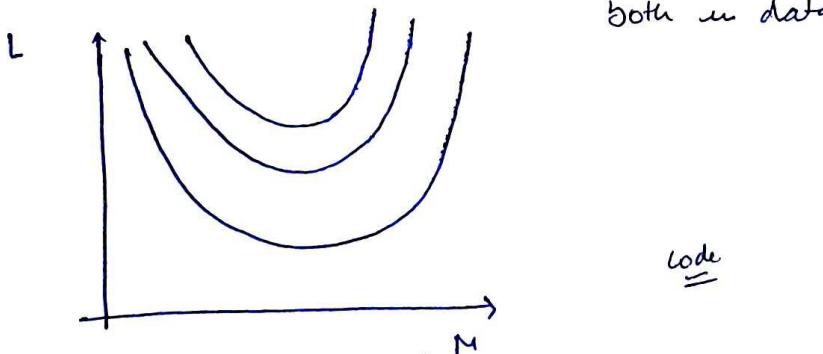
4. Impact on the Loss function

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

$x, y \rightarrow [m, b]$, b is constant ($b=6$)

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

↑ constant ↑ constant
both in data



5. Why called Ridge

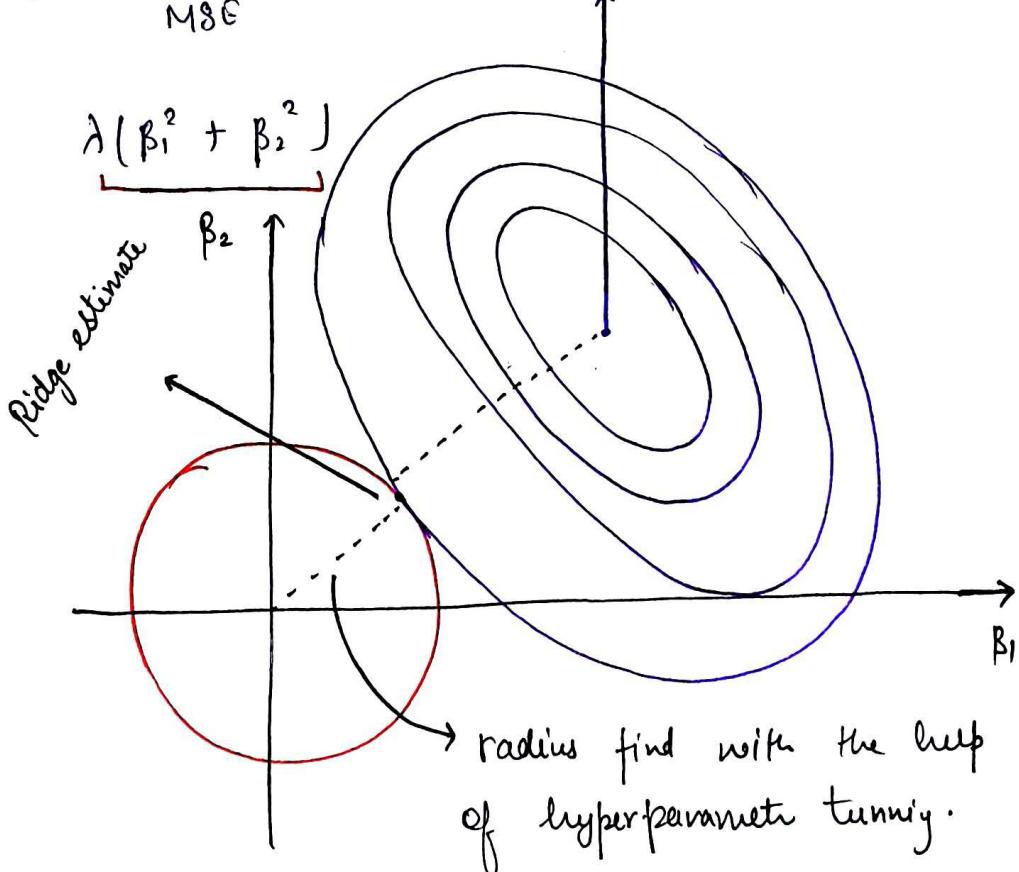
Hard Constraint Ridge Constraint
topic

$$2 \text{ cof } \beta_1 \quad \beta_2 \quad \beta_0$$

$$l = \underbrace{\text{MSE}}_{\lambda} + \underbrace{\lambda \|w\|^2}_{\lambda}$$

$$\sum_{i=1}^n \frac{(y_i - (\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}))^2}{\text{MSE}}$$

DLS estimate



Practice Tip

Use Ridge when there are 2 or more than 2 inputs.

$$x_1, x_2, \dots \geq 2$$

Lasso Regression (L1)

$$L = \text{MSE} + \lambda \|W\|_1$$

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

$\lambda \downarrow \rightarrow$ behaves like linear Regress \rightarrow overfitting

$\lambda \uparrow \rightarrow$ underfitting

$$x_1 \quad x_2 \quad x_3 \quad \dots \quad x_n \quad y$$

\hookrightarrow Lasso \rightarrow $\lambda \uparrow$ ~~less important column~~ \rightarrow correspondingly coeff. = 0
 Feature Selection

- * In Ridge Regression \rightarrow know less important column but still remain some value.
- * In Lasso \rightarrow less important column = 0
- # Use Lasso for large no. of columns.

Lasso Sparsity

$\lambda \uparrow \rightarrow w \rightarrow 0$
 feature \leftrightarrow corresponding coefficient

In Ridge Regression $\lambda \uparrow \rightarrow$ coefficient \rightarrow not equal 0.
 $\lambda \uparrow \rightarrow$ close to 0. but not equal to

In Lasso Regression $\lambda \uparrow \rightarrow$ coefficient becoming 0.
 \hookrightarrow feature Selection

Mathematical

Simple Linear Regression

$$X | Y \rightarrow y = mx + b$$

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

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

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

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

Ridge Regression

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

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

Lasso Regression

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

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

Let consider
 $m > 0$

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

just for
easy mathematics

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

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

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

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

eqn divide
by 2.

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

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

$$m = \frac{\sum (y_i - \bar{y}) (x_i - \bar{x}) - \lambda}{\sum (x_i - \bar{x})^2} \quad \text{if } m > 0$$

Conditions

for $m > 0$

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

for $m = 0$

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

[for $m < 0$]

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

Discuss Sparsity

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

Let

$$m = \frac{yx - \lambda}{x^2} \quad \begin{cases} yx = 100 \\ x^2 = 50 \\ \lambda = 0 \end{cases}$$

$$m = \frac{100}{50} = 2$$

$$\begin{array}{c|ccccc} \lambda = 0 & \lambda = 10 & \lambda = 50 & \lambda = 100 & \lambda = 150 \\ m = 2 & m = 9/5 & m = 1 & m = 0 & m = -1 \end{array}$$

$$m = \frac{yx + \lambda}{x^2} \quad \leftarrow \begin{matrix} \text{3rd condition} \\ m < 0 \end{matrix}$$

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

So, algorithm stop at 0 because after 0 value of m is increasing.
 2, 9/5, 1, 0, 5

$$m < 0$$

$$\lambda > 0$$

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

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

$$m = -2$$

$$\lambda = 0$$

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

$$m = -1$$

$$\lambda = 50$$

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

$$\lambda = 100$$

$$m = 1$$

$\hookrightarrow m$ is positive

so, formula use $m = \frac{-100 - \lambda}{50} = -5$

Why sparsity not in Ridge Regression?

$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2 + \lambda} \rightarrow$ In numerator, if value is 0, still not be 0 after calculation. Because lambda add in denominator
for 0, numerator should be 0. That's why Ridge Regression value near or close to 0.
and not exact 0!
Ridge

ElasticNet Regression

$$\text{Loss Function} = \sum (y_i - \hat{y}_i)^2 + \alpha \|w\|^2 + \beta \|w\|_1$$

$$\lambda = \alpha + \beta$$

$$L1\text{-ratio} = \frac{\alpha}{\alpha + \beta} \quad \text{or} \quad \frac{\alpha}{\lambda}$$

$\lambda, L1\text{-ratio}$
Hyper parameters

let

$$\begin{array}{l} \lambda = 1 \\ \downarrow \\ \alpha = 0.5 \\ b = 0.5 \end{array}$$

$$L1\text{-ratio} = 0.9$$

↳ 90% ridge and 10% lasso

when use?

↳ Input cols \rightarrow multicollinearity \rightarrow used ElasticNet

Code