

Bias-Variance Trade-Off

Introduction

Let us talk about the weather. It rains only if it's a little humid and does not rain if it's windy, hot or freezing. In this case, how would you train a predictive model and ensure that there are no errors in forecasting the weather? You may say that there are many learning algorithms to choose from. They are distinct in many ways but there is a major difference in what we expect and what the model predicts. That's the concept of Bias and Variance Tradeoff.

Bias-Variance Trade-Off

A machine learning model's performance is evaluated based on how accurate is its prediction and how well it generalizes on another independent dataset it has not seen

The errors in a machine learning model can be broken down into 2 parts:

1. Reducible Error
2. Irreducible Error

Irreducible errors are errors that cannot be reduced even if you use any other machine learning model.

What exactly is Bias?

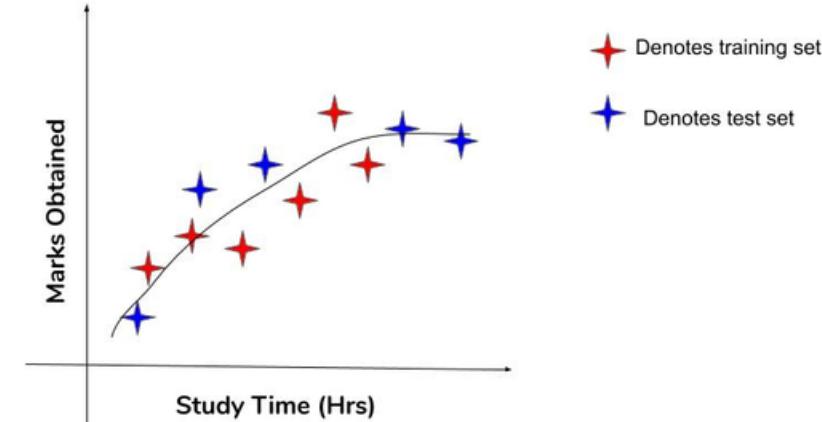
Bias is the inability of a machine learning model to capture the true relationship between the data variables. It is caused by the erroneous assumptions that are inherent to the learning algorithm. For example, in linear regression, the relationship between the X and the Y variable is assumed to be linear, when in reality the relationship may not be perfectly linear.

Bias is defined as the difference of the average value of prediction from the true relationship function $f(x)$

$$\text{bias}[\hat{f}(x)] = \mathbb{E}[\hat{f}(x)] - f(x)$$

What exactly is Bias?

This graph shows the original relationship between the variables. Notice, there is a limit to the marks you can get on the test. That is even if you study an extraordinary amount of time, there is always a certain ‘maximum mark’ you can score. You can see the line flattening beyond a certain value of the X-axis. So the relationship is only piecewise linear. This sort of error will not be captured by the vanilla linear regression model.



You can expect an algorithm like linear regression to have high bias error, whereas an algorithm like decision tree has lower bias. Why? because decision trees don't make such hard assumptions. So is the case with algorithms like k-Nearest Neighbours, Support Vector Machines, etc

- High Bias indicates more assumptions in the learning algorithm about the relationships between the variables.
- Less Bias indicates fewer assumptions in the learning algorithm.

What is the Variance Error?

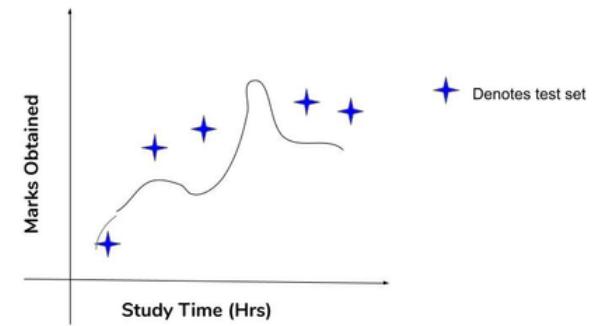
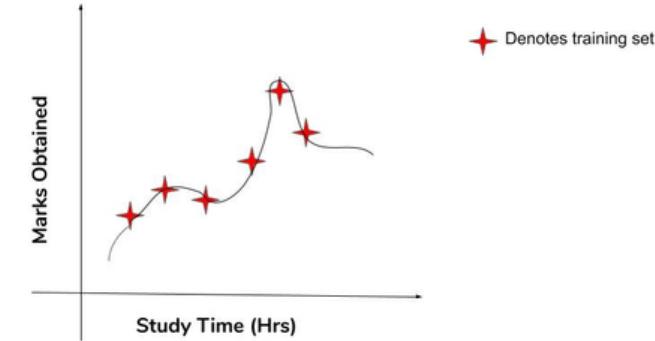
This is nothing but the concept of the model overfitting on a particular dataset. If the model learns to fit very closely to the points on a particular dataset, when it used to predict on another dataset it may not predict as accurately as it did in the first.

Variance is the difference in the fits between different datasets

Variance :- diff of fit on diff datasets

two set → train → fit error = 10
→ test → fit error = 100

$(100 - 10) = 90$ That is Variance



- Generally, nonlinear machine learning algorithms like decision trees have a high variance. It is even higher if the branches are not pruned during training.
- Low-variance ML algorithms: Linear Regression, Logistic Regression, Linear Discriminant Analysis.
- High-variance ML algorithms: Decision Trees, k-NN, and Support Vector Machines.

Mathematically

We have independent variables x that affect the value of a dependent variable y . Function f denotes the true relationship between x and y . In real life problems it is very hard to know this relationship. y is given by this formula along with some noise which is represented by the random variable ϵ with zero mean and variance $\sigma\epsilon^2$:

$$y = f(x) + \epsilon$$

$$\text{bias}[\hat{f}(x)] = \mathbb{E}[\hat{f}(x)] - f(x)$$

$$\text{var}(\hat{f}(x)) = \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$$

Bias is defined as the difference of the average value of prediction from the true relationship function $f(x)$.

Variance is defined as the expectation of the squared deviation of $\hat{f}(x)$ from its expected value $\mathbb{E}[\hat{f}(x)]$.

publication sharing

$$\begin{aligned}\mathbb{E}[(y - \hat{f}(x))^2] &= \mathbb{E}[(f(x) + \epsilon - \hat{f}(x))^2] \\ &= \mathbb{E}[(f(x) - \hat{f}(x))^2] + \mathbb{E}[\epsilon^2] + 2\mathbb{E}[(f(x) - \hat{f}(x))\epsilon] \\ &= \mathbb{E}[(f(x) - \hat{f}(x))^2] + \underbrace{\mathbb{E}[\epsilon^2]}_{=\sigma_\epsilon^2} + 2\mathbb{E}[(f(x) - \hat{f}(x))] \underbrace{\mathbb{E}[\epsilon]}_{=0} \\ &= \mathbb{E}[(f(x) - \hat{f}(x))^2] + \sigma_\epsilon^2\end{aligned}$$

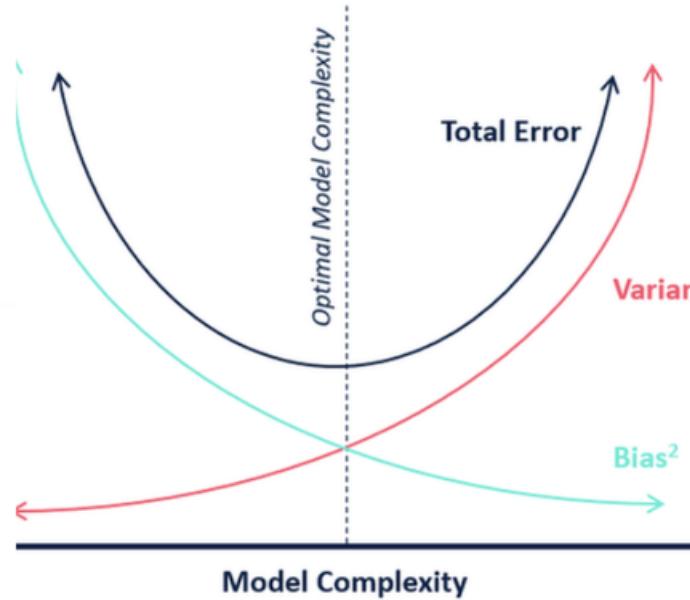
Mathematically

we proceed by expanding further using the linear property of expectation and independence of the random variables ϵ and f . Then using the properties of ϵ and the fact that when two random variables are independent, the expectation of their product is equal to the product of their expectations.

$$\begin{aligned}\mathbb{E}[(f(x) - \hat{f}(x))^2] &= \mathbb{E} \left[\left((f(x) - \mathbb{E}[\hat{f}(x)]) - (\hat{f}(x) - \mathbb{E}[\hat{f}(x)]) \right)^2 \right] \\ &= \mathbb{E} \left[(\mathbb{E}[\hat{f}(x)] - f(x))^2 \right] + \mathbb{E} \left[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2 \right] \\ &\quad - 2\mathbb{E} \left[(f(x) - \mathbb{E}[\hat{f}(x)]) (\hat{f}(x) - \mathbb{E}[\hat{f}(x)]) \right] \\ &= \underbrace{(\mathbb{E}[\hat{f}(x)] - f(x))^2}_{=\text{bias}[\hat{f}(x)]} + \underbrace{\mathbb{E} \left[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2 \right]}_{=\text{var}(\hat{f}(x))} \\ &\quad - 2(f(x) - \mathbb{E}[\hat{f}(x)]) \mathbb{E} \left[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)]) \right] \\ &= \text{bias}[\hat{f}(x)]^2 + \text{var}(\hat{f}(x)) \\ &\quad - 2(f(x) - \mathbb{E}[\hat{f}(x)]) (\mathbb{E}[\hat{f}(x)] - \mathbb{E}[\hat{f}(x)]) \\ &= \text{bias}[\hat{f}(x)]^2 + \text{var}(\hat{f}(x))\end{aligned}$$

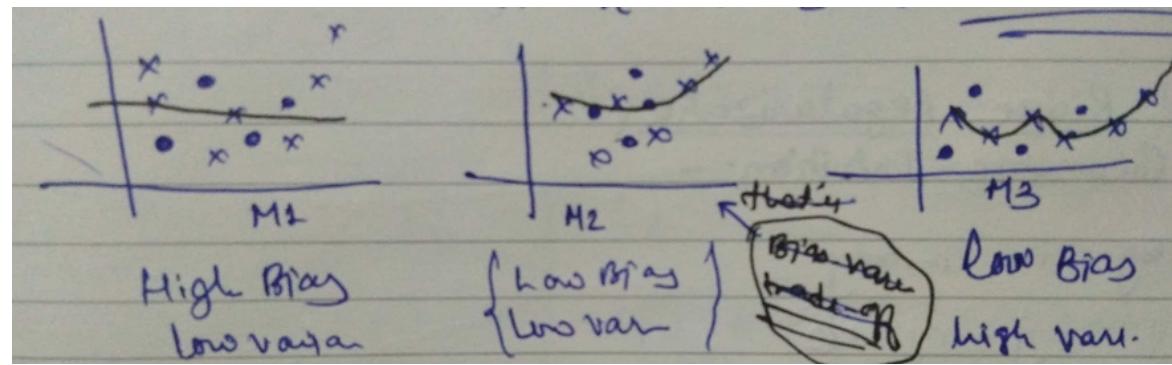
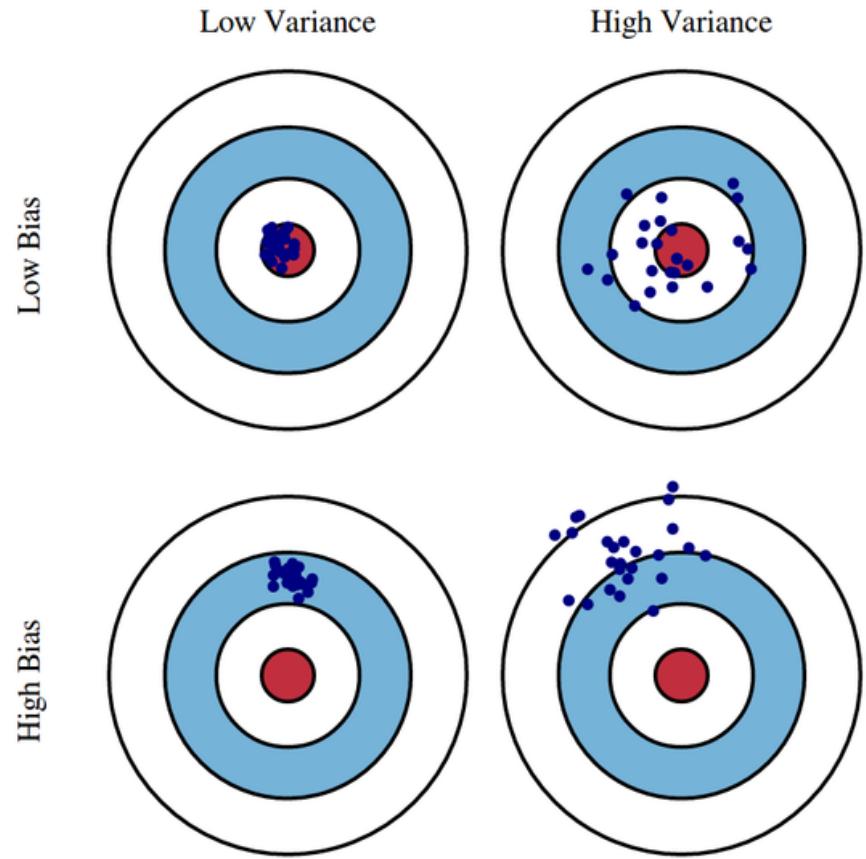
$$\mathbb{E}[(y - \hat{f}(x))^2] = \text{bias}[\hat{f}(x)]^2 + \text{var}(\hat{f}(x)) + \sigma_\epsilon^2$$

Bias – Variance Tradeoff



Example a Coin,
Student, Company...

How to Overcome:
Regularization,
Bagging, Boosting



Let's summarize

Overfitting → training data fit too well perform, but NOT on test data
Low Bias + high Variance

Underfitting → training data fit not good fit relationship
Capture only one type of data
High Bias + low variance

- 1.
2. If a model follows a complex machine learning model, then it will have high variance and low bias(overfitting the data).
3. You need to find a good balance between the bias and variance of the model we have used. This tradeoff in complexity is what is referred to as bias and variance tradeoff. An optimal balance of bias and variance should never overfit or underfit the model.
4. This tradeoff applies to all forms of supervised learning: classification, regression, and structured output learning.

Regularization

- This technique prevents the model from overfitting by adding extra information to it.
- It is a form of regression that shrinks the coefficient estimates towards zero. In other words, this technique forces us not to learn a more complex or flexible model, to avoid the problem of overfitting.
- In the Regularization technique, we reduce the magnitude of the independent variables by keeping the same number of variables”. It maintains accuracy as well as a generalization of the model

How does Regularization Work?

Regularization works by adding a penalty or complexity term or shrinkage term with Residual Sum of Squares (RSS) to the complex model.

Regularization

Techniques of Regularization

Mainly, there are two types of regularization techniques, which are given below:

- Ridge Regression
- Lasso Regression
- other is Elastic Net

Regularization

A simple relation for linear regression looks like this. Here Y represents the learned relation and β represents the coefficient estimates for different variables or predictors(X).

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

The fitting procedure involves a loss function, known as residual sum of squares or RSS. The coefficients are chosen, such that they minimize this loss function.

$$\text{RSS} = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2.$$

Now, this will adjust the coefficients based on your training data. If there is noise in the training data, then the estimated coefficients won't generalize well to the future data. This is where regularization comes in and shrinks or regularizes these learned estimates towards zero.

Ridge Regression

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

Above image shows ridge regression, where the RSS is modified by adding the shrinkage quantity. Now, the coefficients are estimated by minimizing this function. Here, **λ is the tuning parameter that decides how much we want to penalize the flexibility of our model.** The increase in flexibility of a model is represented by increase in its coefficients, and if we want to minimize the above function, then these coefficients need to be small. This is how the Ridge regression technique prevents coefficients from rising too high. Also, notice that we shrink the estimated association of each variable with the response, except the intercept β_0 . This intercept is a measure of the mean value of the response when $x_{i1} = x_{i2} = \dots = x_{ip} = 0$.

Ridge Regression

When $\lambda = 0$, the penalty term has no effect, and the estimates produced by ridge regression will be equal to least squares. However, as $\lambda \rightarrow \infty$, the impact of the shrinkage penalty grows, and the ridge regression coefficient estimates will approach zero. As can be seen, selecting a good value of λ is critical. Cross validation comes in handy for this purpose. The coefficient estimates produced by this method are also known as the L2 norm.

The coefficients that are produced by the standard least squares method are scale equivariant, i.e. if we multiply each input by c then the corresponding coefficients are scaled by a factor of $1/c$. Therefore, regardless of how the predictor is scaled, the multiplication of predictor and coefficient ($X_j\beta_j$) remains the same. However, this is not the case with ridge regression, and therefore, we need to standardize the predictors or bring the predictors to the same scale before performing ridge regression. The formula used to do this is given below.

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}},$$

Lasso Regression

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|.$$

Lasso is another variation, in which the above function is minimized. Its clear that this variation differs from ridge regression only in penalizing the high coefficients. It uses $|\beta_j|$ (modulus) instead of squares of β , as its penalty. In statistics, this is known as the L1 norm. Lets take a look at above methods with a different perspective. The ridge regression can be thought of as solving an equation, where summation of squares of coefficients is less than or equal to s. And the Lasso can be thought of as an equation where summation of modulus of coefficients is less than or equal to s. Here, s is a constant that exists for each value of shrinkage factor λ . These equations are also referred to as constraint functions.

Lasso Regression

Consider there are 2 parameters in a given problem. Then according to above formulation, the ridge regression is expressed by $\beta_1^2 + \beta_2^2 \leq s$. This implies that ridge regression coefficients have the smallest RSS(loss function) for all points that lie within the circle given by $\beta_1^2 + \beta_2^2 \leq s$.

Similarly, for lasso, the equation becomes, $|\beta_1| + |\beta_2| \leq s$. This implies that lasso coefficients have the smallest RSS(loss function) for all points that lie within the diamond given by $|\beta_1| + |\beta_2| \leq s$.

the lasso method also performs variable selection and is said to yield sparse models.

Lasso stands for least absolute shrinkage and selection operator

What does Regularization achieve?

A standard least squares model tends to have some variance in it, i.e. this model won't generalize well for a data set different than its training data. Regularization, significantly reduces the variance of the model, without substantial increase in its bias. So the tuning parameter λ , used in the regularization techniques described above, controls the impact on bias and variance. As the value of λ rises, it reduces the value of coefficients and thus reducing the variance. Till a point, this increase in λ is beneficial as it is only reducing the variance(hence avoiding overfitting), without loosing any important properties in the data. But after certain value, the model starts loosing important properties, giving rise to bias in the model and thus underfitting. Therefore, the value of λ should be carefully selected.

Ridge(Geometric Intuition)

I Ridge Regularization

Geometric - Intuition: -

e.g. in lin. reg

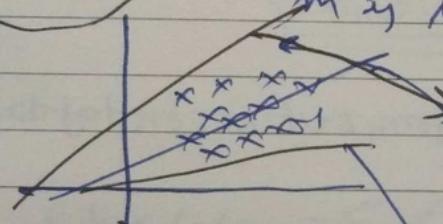
$$y = mx + b$$

वर्तमान में x, m, b को दिया गया है
इसका उपयोग करें।

- train & test

If m is high \rightarrow overfitting

m is low \rightarrow underfitting



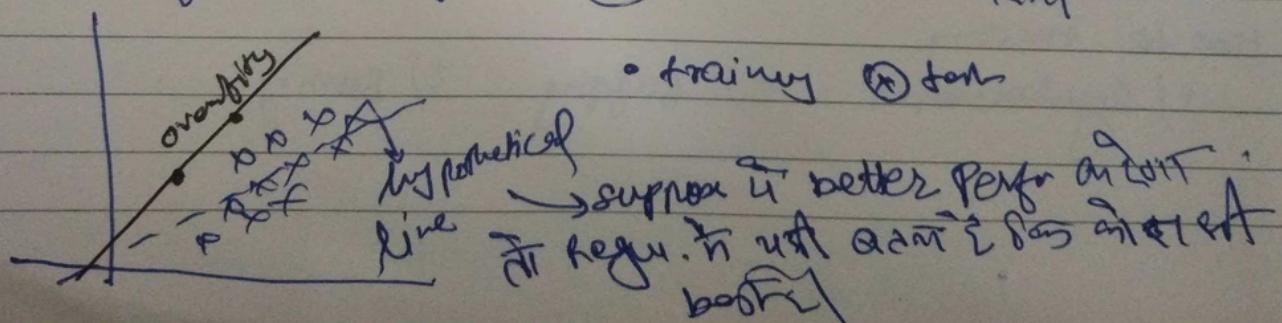
जब भी $m \uparrow$ हो $\rightarrow m \rightarrow \infty$

मूल वाले y को लाइन की ओर नहीं देख सकते।

वैसे भी $m \downarrow$ हो इसकी ओर नहीं देख सकते।

so for overfitting \rightarrow (m) को reduce करें।

- training & test

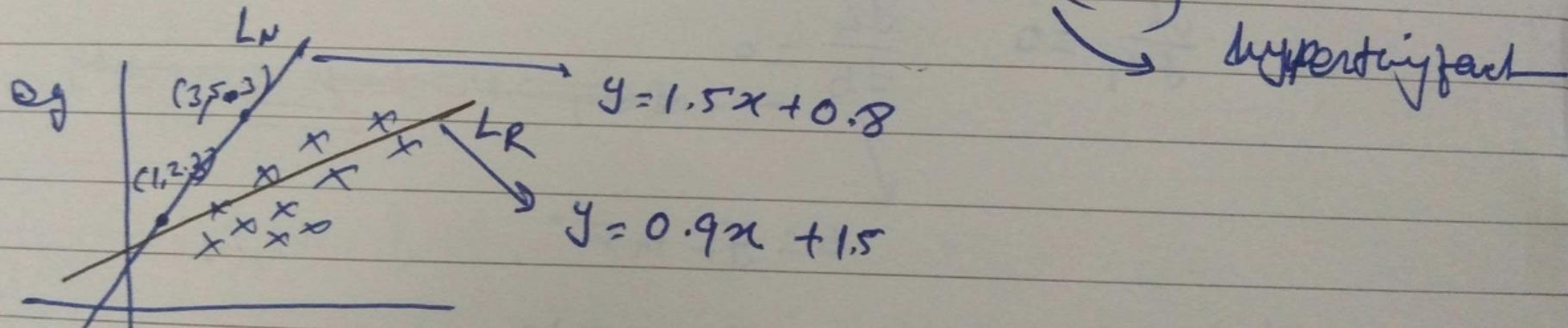


Ridge(Geometric Intuition)

$$L = \sum_{i=1}^n (\theta_i - \hat{y}_i)^2 \leftarrow \text{minimize}$$

In regularization

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



Let $\lambda = 1$

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

2.25

L_R

$$(2.3 - 0.9 - 1.5)^2 + (5.3 - 2.7 - 1.5)^2 + (0.9)^2$$

2.03

Ridge(Geometric Intuition)

$\lambda = 0.25$

$\lambda = 2.03$

$+ (0.9)^2$

$H_{\lambda=0.25}$ LR line better than L_2

$H_{\lambda=2.03}$ fit L_2 line to noisy data & it's more accurate
but $H_{\lambda=2.03}$ select most coefficients
 $L_{\lambda=2.03}$ is less due to regular

(Bias is greater but variance significantly
less than L_2)

It's called L_2 regular because

squared multiply $m_1^2 + m_2^2 + \dots$

L_2 norm

$$\sqrt{m_1^2 + m_2^2 + \dots}$$

Ridge(Mathematical Formulation)

mathematical formulation (Ridge-Reg)

① By OLS → Ordinary Least Squares Method

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

→ Least squares method → m must fit well

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

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

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

$$\hat{y}_i = mx_i + b$$

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

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

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

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

Ridge(Mathematical Formulation)

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

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

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

add $\frac{\lambda}{n}$ to \bar{y}

so b also chose

QH

n-p data

$$L = \sum (y_i - \hat{y}_i)^2$$

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

$$x \rightarrow mx(n+1)$$

$$y = mx_1$$

$$w = nx_1$$

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

Ridge(Mathematical Formulation)

Normal eqn.

$$L = (xw - y)^T (xw - y)$$

~~Model Eqn~~

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

$$(w_0^2 + w_1^2 + \dots + w_n^2)$$

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

minimise

$$\boxed{\frac{\partial L}{\partial w} = 0}$$

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

$$L = w^T x^T x w - 2 w^T x^T y + y^T y + \lambda w^T w$$

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

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

\uparrow
 \downarrow
 $\rightarrow w \rightarrow$

addecs *

vector form

Ridge(Using GD)

(b)

Ridge Reg using Gradient Descent

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

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

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

}

Ridge(Using GD)

Since $L = (x\omega - y)^T(x\omega - y) + d\omega^T\omega$

$$L' = \frac{1}{2} (\omega^T x^T - y^T) (x\omega - y) + \frac{1}{2} d\omega^T\omega$$

← Just multiply by 1/2

$$L' = \frac{1}{2} [\cancel{\omega^T x^T x\omega} - \cancel{2y^T \omega x} + y^T y] + \frac{1}{2} d\omega^T\omega$$

$$\frac{\partial L'}{\partial \omega} = \frac{1}{2} [2x^T x\omega - 2x^T y] + \frac{1}{2} \cdot 2d\omega$$

$$= x^T x\omega - x^T y + d\omega$$

$$= \frac{\partial L}{\partial \omega}$$

epoch }

$$\boxed{\omega = \omega - n \frac{\partial L}{\partial \omega}}$$

$$\omega = [\omega_0, \omega_1, \dots, \omega_m]$$

Initially

5 key Points (Ridge Regression)

5 key points : Ridge regression : —

1. How the coeffs get affected (by applying ridge reg.)

as $\lambda \uparrow$ $\lambda \rightarrow 0 \rightarrow \infty$

If $\lambda = 0 \rightarrow$ No regu.

$\lambda \uparrow$ then coeff.

shrink

toward $\rightarrow 0$
but $0 > \lambda \geq 0$

Can't $\lambda \uparrow$? single & losically ~~overfitting~~
~~multiple~~

5 key Points (Ridge Regression)

2. ^(coeff.) Higher Values are impacted more!

$$\begin{matrix} x_1 & x_2 & x_3 \\ \downarrow & \downarrow & \downarrow \\ w_1 = 1000 & 10 & 1 \end{matrix}$$

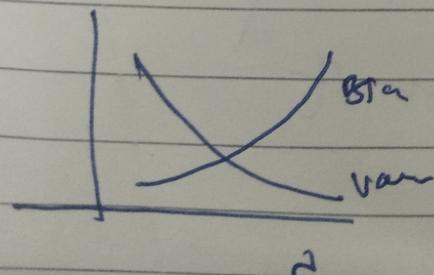
' ~~दूसरे~~ अंतर्रुपेक्षा की तरीके से coeff प्रभावित होता है और उनकी वाली तरफ से बढ़ता है।'

3. Effect on Bias-Variance Tradeoff

Bias \downarrow overfit var \uparrow
Bias \uparrow underfit var \downarrow

मानवानि यह एक चौथा है $\rightarrow L \downarrow \rightarrow$ overfit करने की कमी

overfit करने की कमी
(train data की लिये)



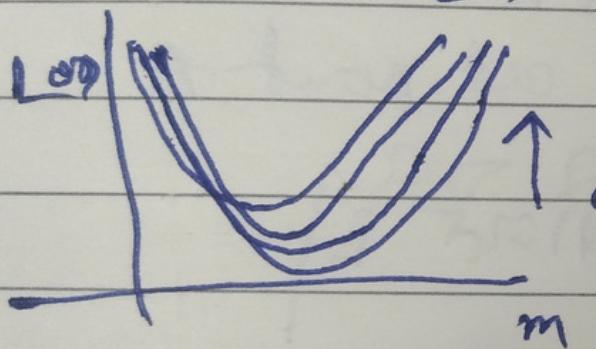
5 key Points (Ridge Regression)

4. Impact on Lossf²

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

let $b = \text{const}$ (say = 0)

$$L = \sum (y_i - m w_i)^2 + \lambda m^2$$



$\lambda \uparrow \Rightarrow \text{coeff } \downarrow \rightarrow \text{tends to 0}$

$\lambda \uparrow$ less dist acc ~~at~~ graph
close to sum (shrink)

←
towards centre at $(m=0)$

5 key Points (Ridge Regression)

5. Why called Ridge reg. → Ridge

soft vs hard constraint Ridge Reg.

one more Hard Constraint Ridge Reg.

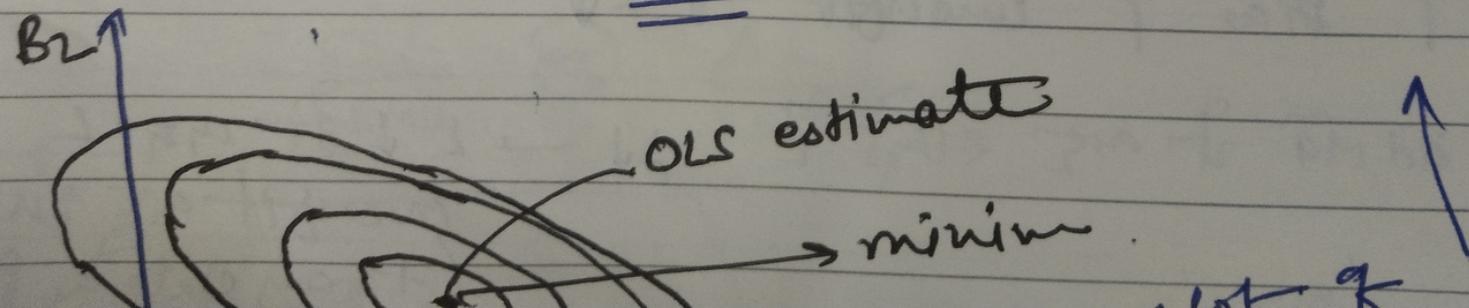
$$L_{\text{R}} = \text{MSE} + \lambda \|w\|^2 \quad \text{Let two coeff } B_1, B_2, \dots$$

$$\text{MSE} = \sum (y_i - (B_0 + B_1 x_{i1} + B_2 x_{i2}))^2$$

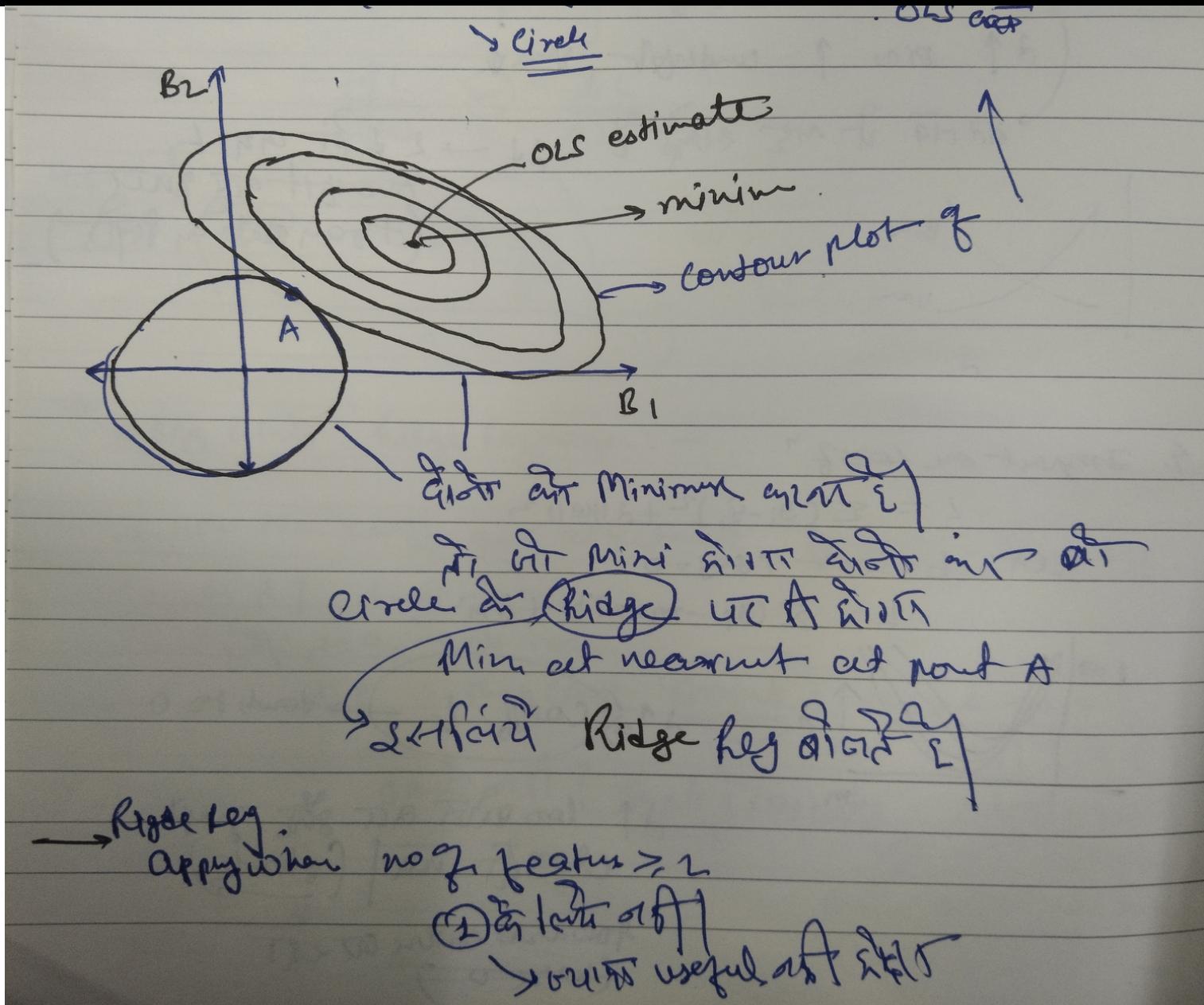
$$\lambda (B_1^2 + B_2^2)$$

→ circle

OLS est



5 key Points (Ridge Regression)



Lasso Regression (L1 Regularization)

Lasso Regression (L1 Regularization)

$$L = \text{MSE} + \gamma ||\omega||$$

$||\omega|| \rightarrow \text{L1 Norm}$

$$\gamma [|\omega_1| + |\omega_2| + \dots + |\omega_n|]$$

↑ soft coeff shrinking & soft outlier handling की तरह Ridge की तरह है।

do feature selection, dim.
automatically

do Lasso is prefer (mostly higher dim.)

key points

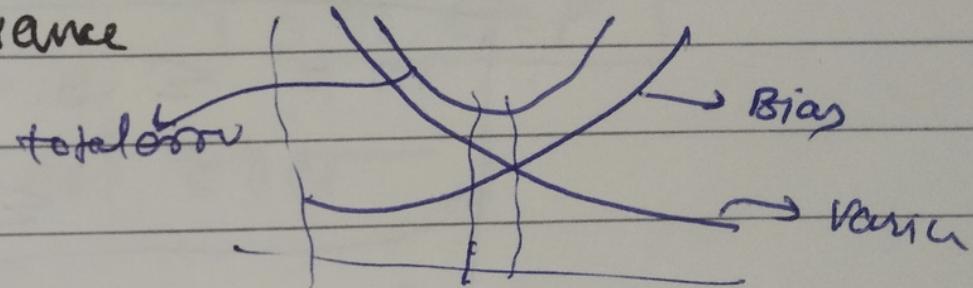
1. How are coeff affected? feature selection by λ
But very fit; may underfit

2. Higher coeff are affected more

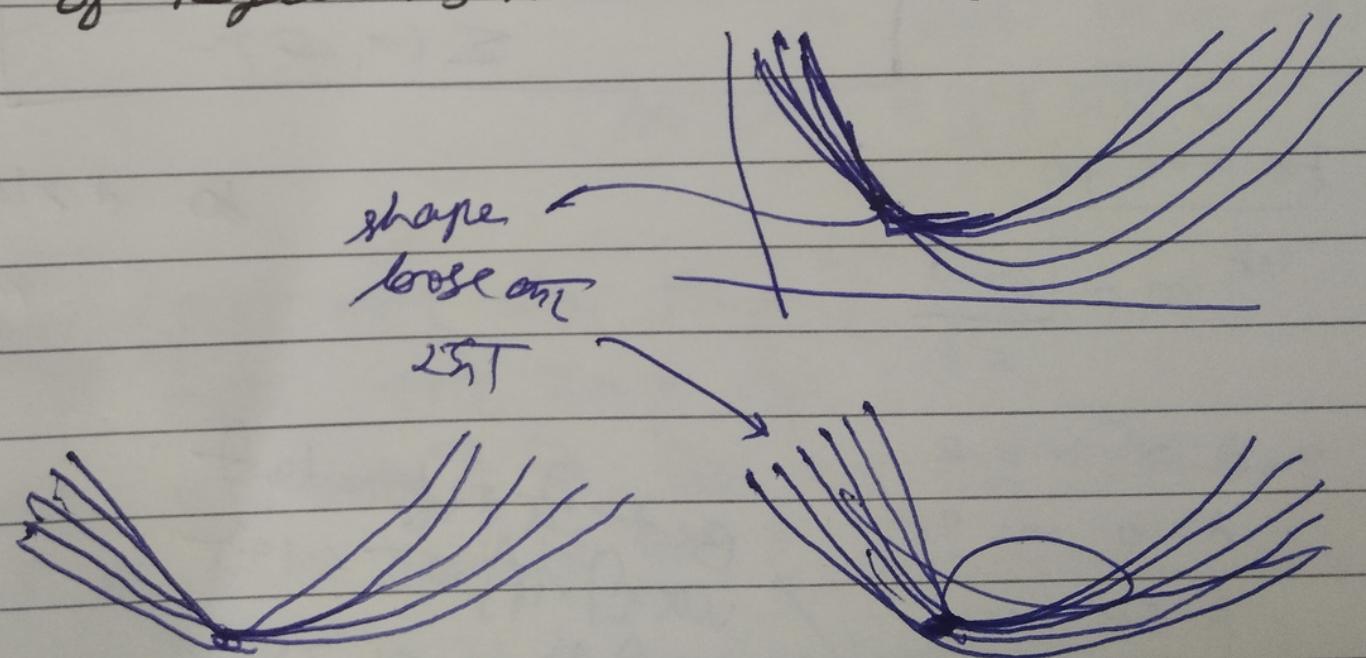
बड़े coeff पर जायदा effect -

Lasso Regression (L1 Regularization)

3. Impact on Bias & Variance



4. Effect of Regularization on Loss of f^*



Why Lasso Regression creates Sparsity ?

Q: why Lasso regression creates sparsity? ~~Hence $\alpha \uparrow$ $w \rightarrow 0$~~

& why Ridge not able to do?

Observe $w = 0$

for $\alpha \uparrow$ (high) all $w = 0$

Let consider case of 1-feature

$$y = mx + b$$

for Simple LR

for Ridge

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

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

for Lasso

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

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

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

Why Lasso Regression creates Sparsity ?

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

when
 $m > 0$

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

when
 $m < 0$

for +ve m

$$\text{let } m = \frac{100 - \lambda}{50}$$

so $n > 100$

$$\text{let } \lambda = 150$$
$$m = \frac{100 + \lambda}{50}$$

$$\lambda = 0 \quad m = 2$$

$$\lambda = 10 \quad m = 9$$

$$\lambda = 50 \quad m = 1$$

$$\lambda = 100 \quad m = 0$$

$$\lambda \geq 100 \quad m < 0$$

But λ in formula \neq
we will take $\lambda < 100$
and $\lambda > 0$

$$m = 5$$

Why Lasso Regression creates Sparsity ?

m → 2, 9₁₅, 11, 0₁ ← 5
 ← dec → Nao ↑ नाले पहुँची तो 2 आ
दी algo तो रक्खा
यही दीक्षा तो चले जा
दृश्यम् वर्तमान मिथुन

નો algo
અની એક સ્થાન
ને પરિસ્તિહાસ

o पर बढ़ो

for $m < 0$ should

$$m = \frac{-ve + \lambda}{+ve}$$

$$\text{Get } m = \frac{-100t + 2}{50}$$

| <u>m</u> | d | m |
|----------|----|---|
| 0 | -2 | |
| 50 | -1 | |
| 100 | 0 | |

so we often former

① +ve X

-5

same

o पर stop

Why not Ridge Regression?

Ridge Reg में sparsity को कैसे बढ़ावा दें ?

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

यहाँ परिस्थिति यह है कि
दो रेंज लेने के बाद $m=0$ हो जाएगा।
तो close नहीं पाहेंगे।

Ridge \leftarrow उपर्युक्त छोटे सारे दम फॉर्मूले
Lasso \leftarrow उपर्युक्त लार्ग दम फॉर्मूले
But इन्हीं वहाँ का डेटा है तो कैसे पता करी जीने फॉर्मूला
होता है ? So which should we use ?

That's why ElasticNet Reg

ElasticNet Regression

ElasticNet - ElasticNet Regression Comb. of Ridge + Lasso

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

By default $\alpha = \beta = 0.5$ दोनों की equal priority.

But $\alpha > \beta$ करा तो लसो

Where to use?

1) Test R-sq -

2) Input clm में Multi-collinearity होती है

In Sklearn two param

$$\begin{cases} \alpha = \alpha + \beta \\ l1_ratio = \frac{\alpha}{\alpha + \beta} \end{cases}$$

By default $\alpha = 1$, $l1_ratio = 0.5$

Why we need Elastic Reg?

Lasso regression it tries to shrink the coefficients to the absolute zero and if not possible to shrink to the absolute zero, then it eliminates the coefficient from the models. The ridge regression does not eliminate the coefficients from the model, which means it does not differentiate between important and less important predictive variables in the model and includes all of them by providing l2 penalty. It tries to shrink the unbiased coefficient by putting them with their squared magnitude into the model.

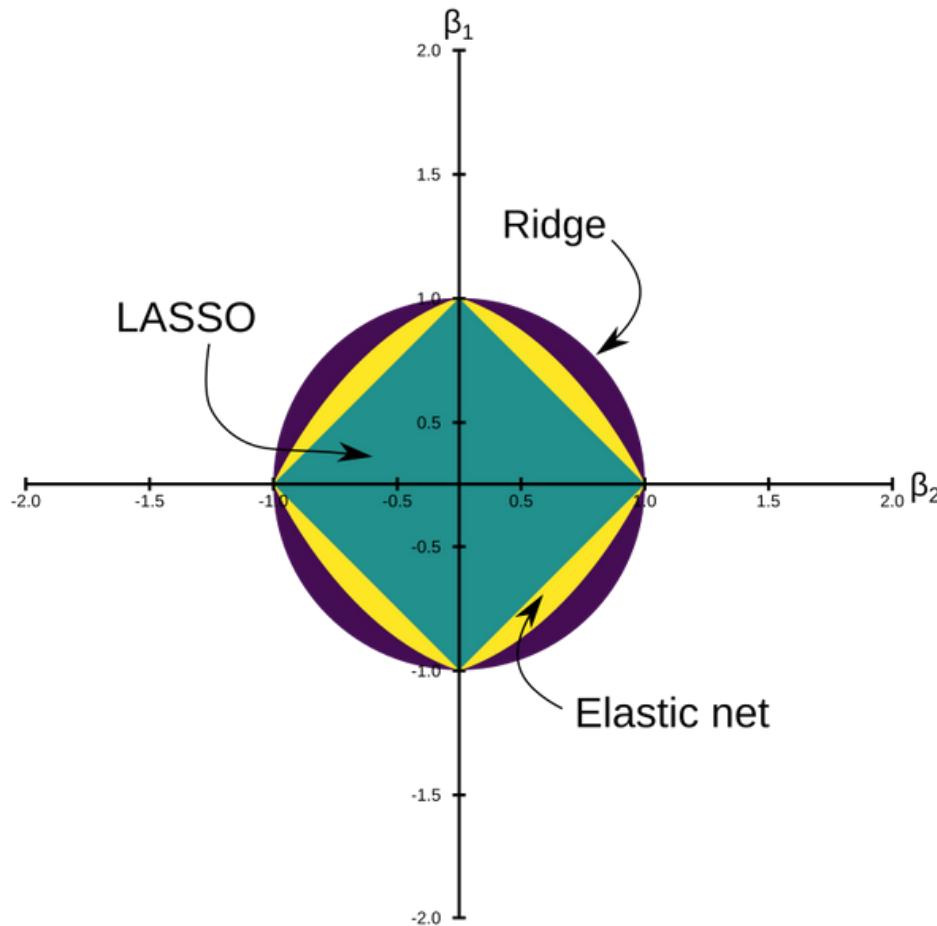
Why we need Elastic Reg?

But there are certain limitations of these models- ridge regression decreases the complexity of the model in performance but does not eliminate the unbiased variables hence we can increase the model's accuracy in a large dataset to a point. The new unbiased variable generated model can stop performing well. The lasso regression model picks the points according to the number of observations, not the predictor presented in the data. This kind of limitation can be handled and removed by the elastic net regression model where it includes both kinds of (l1 and l2) penalties in the model. (MULTICOLINIRITY) **incapability of lasso is choosing the number of predictors.**

How Elastic Reg Solve this Problem?

Elastic Net is a regularized regression model that combines l1 and l2 penalties, i.e., lasso and ridge regression. We have discussed the limitations of lasso regression, where we found the incapability of lasso is choosing the number of predictors. The elastic net includes the penalty of lasso regression, and when used in isolation, it becomes the ridge regression. In the procedure of regularization with an elastic net, first, we find the coefficient of ridge regression. After this, we perform a lasso algorithm on the ridge regression coefficient to shrink the coefficient.

How Elastic Reg Solve this Problem?



Here we can see that after performing the ridge regression, the lasso regression takes part in the procedure that considers all the variables from the dataset.

Mathematically we can represent the elastic net as follows.

$$\hat{\beta} \equiv \underset{\beta}{\operatorname{argmin}} (\|y - X\beta\|^2 + \lambda_2 \|\beta\|^2 + \lambda_1 \|\beta\|_1).$$