



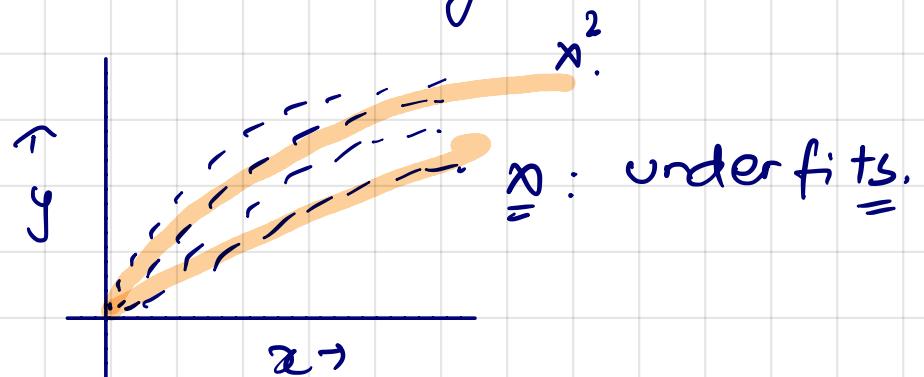
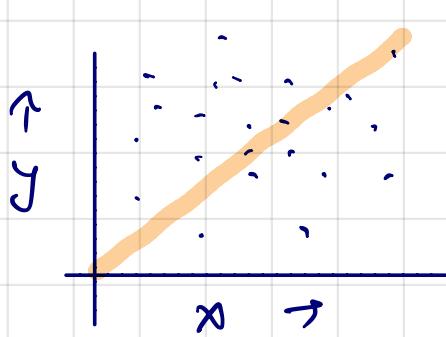
Regularization

## why Regularisation?

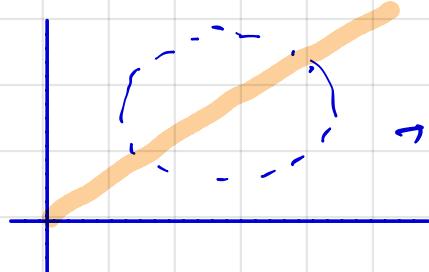
- Reduce Outliers
- Shrinking values.
- Lambda with beta.
- Scattered Points Train vs Test Error.
- Train Error ; Testing Error.

{

Test Error > Train Error : overfitting.



$x^2$  as power increases  
model overfits.



→ underfitting

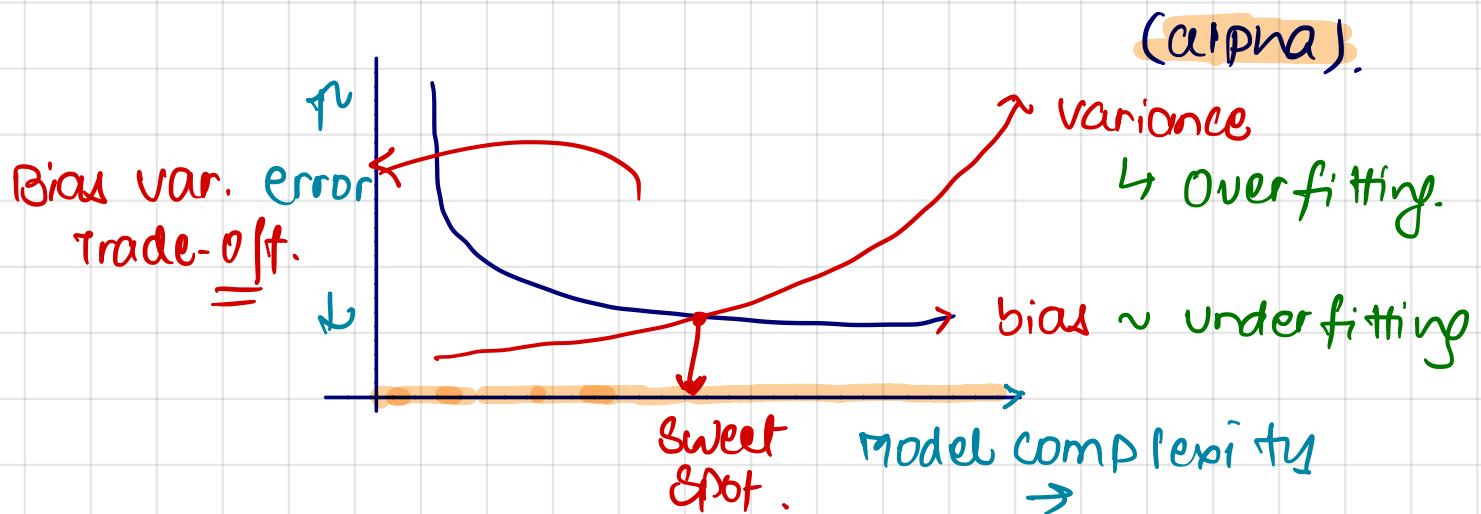
Most instances, LR appears to be an underfit.

It means it is not able to learn correct pattern from data.

Coefficients: Some coefficients get a lot of weightage and hence are higher or lower than the rest.

Regularization helps in making the model

- consistent
- avoid overfitting by introducing bias ( $\lambda$ ).



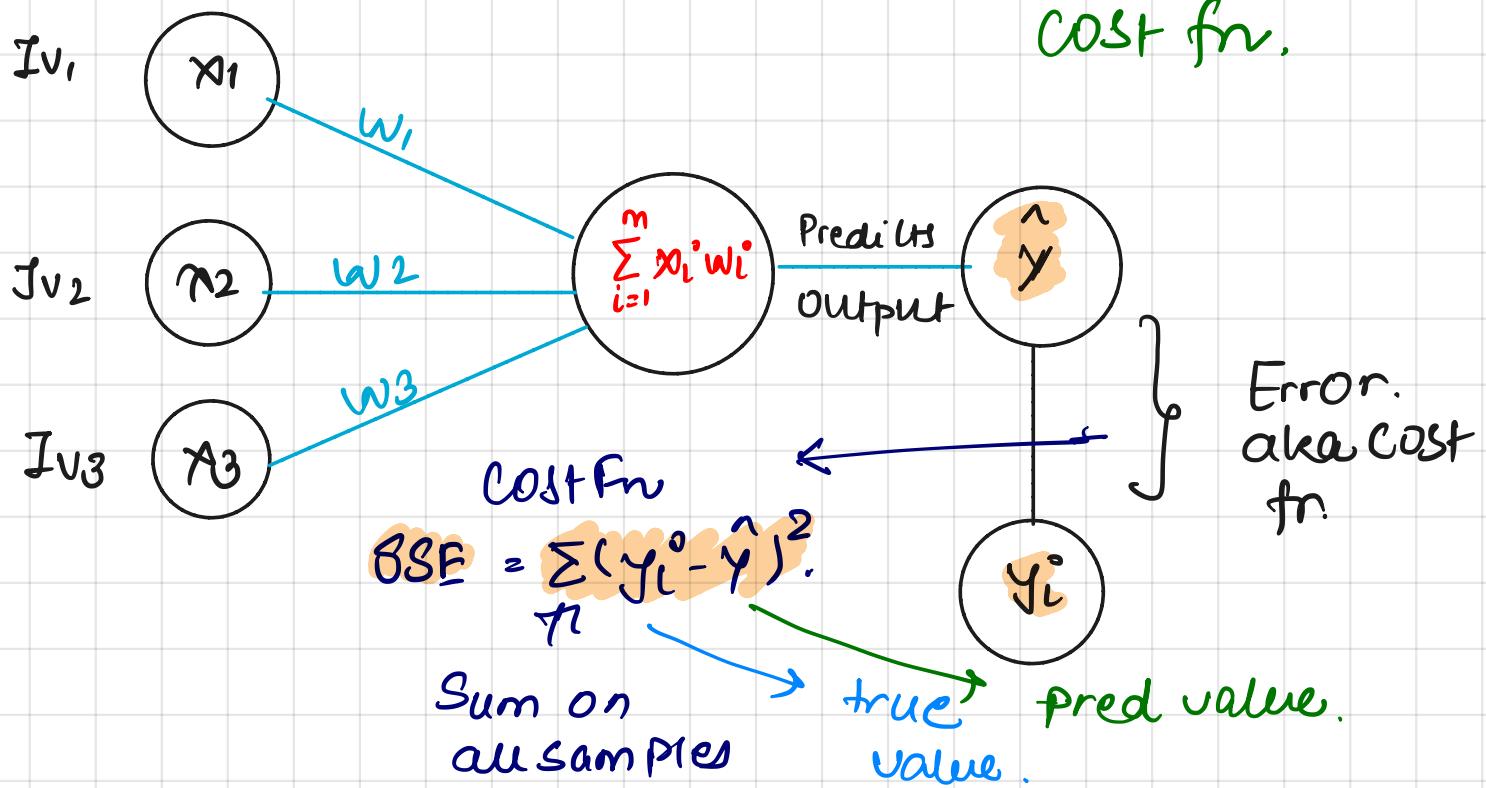
"Reg is a technique that reduces error by fitting a fn (Model) appropriately on the training set in order to minimize the overfit effect."

↳ Lasso Regression

Regularization → L2 ; Ridge Regression.

→ Dropout → NNets.  
Regularization

**COST Function**, : objective is to minimise the cost fn.



$$Cf = \min \left( \sum_{i=1}^m (y_i^0 - \hat{y})^2 \right)$$

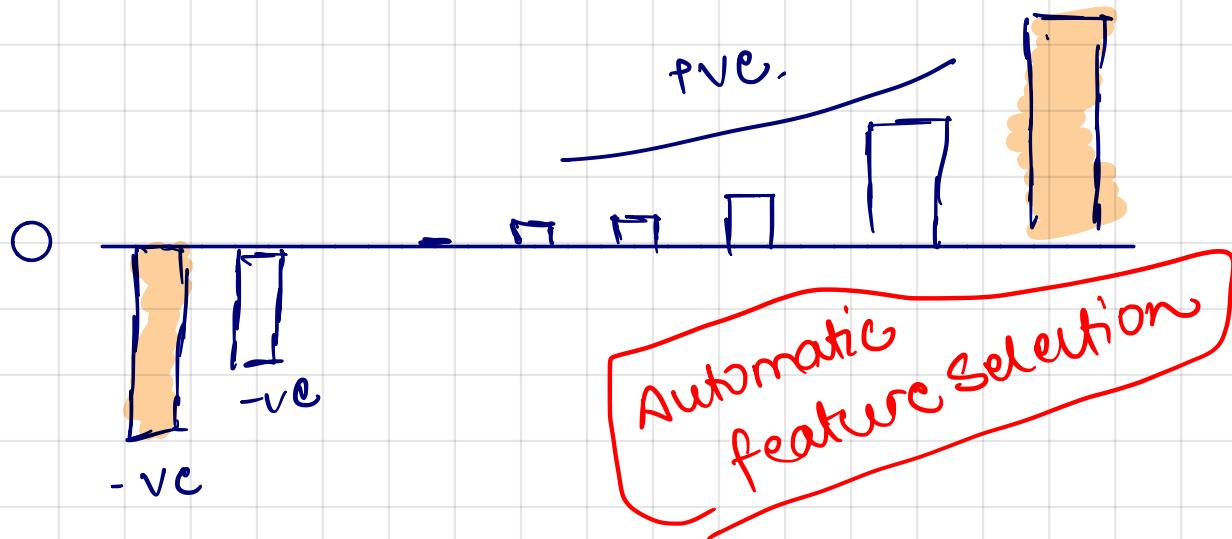
[ Cost Function =  $\frac{1}{m} \sum_{i=1}^m (y_i^0 - \hat{y})^2$  ]

Optimization problem.  $\approx \frac{1}{2m} \sum_{i=1}^m (h(\theta) - x_i^0)^2$

A cost fn is measure of error b/w what model predicts & what actual value is.

## Lasso Regression.

- Lasso adds **L1 Penalty**.
- Shrinkage parameter that reduces the **Coefficients to 0** of those predictors which don't contribute to model.



- Lasso will make some Coeff.  $\rightarrow 0$ . The Model will eventually become simple, overfitting is reduced
- At the same, underfitting is introduced  
The model will loose the predictive power.

Reducing  $\rightarrow 0$   $\Rightarrow$  less predictors  
Coefficient

loose prediction Power.

↓  
Model becomes simple

$$\text{Cost Fn} = \min \left( \sum_{i=1}^m (y_i^o - \hat{y})^2 \right)$$

$\Downarrow$

$$\text{Penalty.} \quad \Rightarrow \frac{\partial x}{\partial} \frac{\partial (y_i^o - \hat{y})}{\partial}$$

$$\text{Cost Fn (lasso)} = \min \left( \sum_{i=1}^m (y_i^o - \hat{y})^2 \right) + \lambda \|\theta\|_1$$

Here  $\theta$  = sum of absolute weights.

if  $\lambda = 0 \Rightarrow$  Cost fn  $\rightarrow$  LR.

$$\frac{d(C)}{dx} = 0, \quad \theta = \frac{3091}{\frac{d(3091)}{dx} = 0}$$

$\underline{\underline{\lambda}}$  : Ridge Regression.

$$\frac{d(\|\theta\|^2)}{dx} = 2\theta,$$

In case of Lasso:

$$\|w\| = |w_1| + |w_2| + \dots + |w_n|.$$

Ridge: SQ. Magnitude of coefficients.

$$\|\theta\| = (\|\theta_1\|^2 + \|\theta_2\|^2 + \dots + \|\theta_n\|^2)^{1/2}$$

$$CF = \sum_{i=1}^m (y_i^o - \hat{y})^2 + \lambda \|\theta\|^2$$

## when to use L<sub>1</sub> Ei L<sub>2</sub>

- L<sub>1</sub> is good when considering a cat. variable with many levels.  
eg ItemType in bigmart.
- L<sub>1</sub> is preferred when we are interested in fitting a linear model with fewer var.
- L<sub>1</sub> converges to 0, helps in feature selection and reduces overfitting where as L<sub>2</sub> does not encourage convergence towards zero but it gets close to zero and prevents overfitting
- Larger the features  $\rightarrow \alpha_2$  is better. Also, if predictors are highly collinear  $\Rightarrow$  Ridge.
- In lasso, it will drop one collinear var but in Ridge, it will take both the variables and jointly reduce the weight.

→  $\ell_2$  regularization is called "Weight Decay" as it forces the weights to decay towards 0.

Ridge

→ does not reduce the coeff. to 0.

$$\text{Ridge CF} = \sum (y_i - \hat{y})^2 + \lambda \cdot \sum \|\theta\|^2$$

Lasso

→ Reduces the coeff. to 0.

$$\text{Lasso CF} = \sum (y_i - \hat{y})^2 + \lambda \cdot \sum \|\underline{\theta}\|$$

Elastic Net

$$\text{Enet CF} = \left[ \sum (y_i - \hat{y})^2 + \lambda \cdot \sum \|\theta\|^2 + \lambda \cdot \sum \|\underline{\theta}\| \right]$$

Gives better prediction by combining

Penalty of Ridge & Lasso.

$$\Rightarrow (\alpha = 0) \text{ : Ridge, } (\alpha = 1) \text{ : Lasso}$$