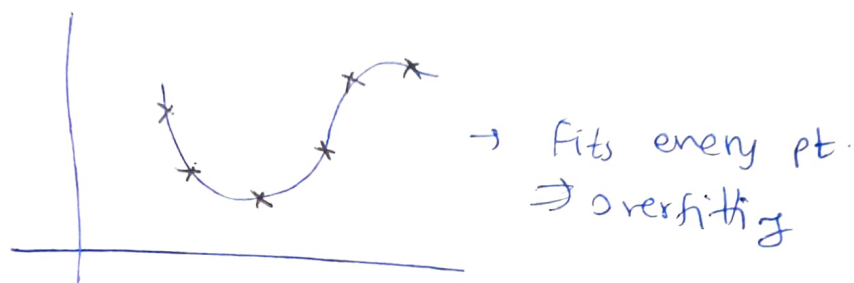


# Regularization

As no. of features  $\uparrow$ ,  $y = w_1 x_1 + w_2 x_2^2 + w_3 x_1 x_2 + w_4 x_4^2$   
 $\Rightarrow$  Polynomial

Overfitting can happen on data.



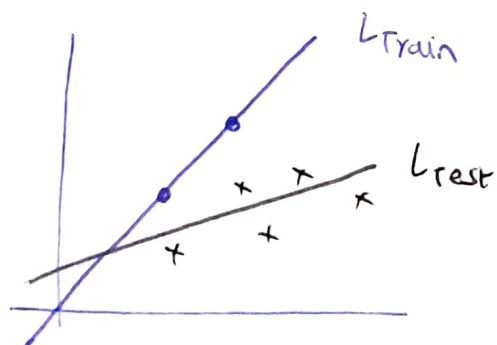
Regularization - Modification we make to the learning algo to prevent overfitting.

low train error

High test error.

## ① Ridge ( $L_2$ )

Premise of overfitting - Too powerful model fit on too little data.



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

with ridge ( $L_2$ ) regularization,

① Reduce power of model  $\rightarrow$  Here.

②  $\uparrow$  no. of samples.

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

wt. decay

slope

Hyperparam  
(1 to  $\infty$ )

preferring sol<sup>n</sup> where slope is 0.

7- how much imp. should be given to both parts of addit

1  $\rightarrow$  equal weightage

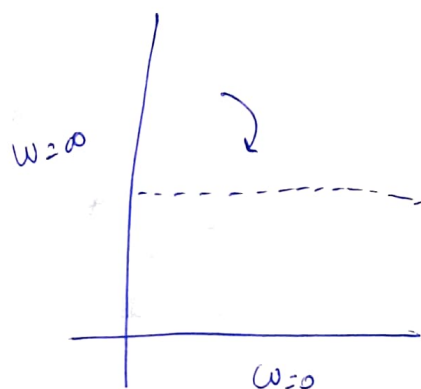
$\lambda \geq 2 \rightarrow$  more imp. to wt. decay.

Slope gets closer to 0.

$$y = 2x + 0.3$$

$$y = 1.2x + 0.7 \quad \left. \begin{array}{l} \end{array} \right\} \text{slope } \downarrow$$

slope  $\downarrow$ , bias  $\uparrow$



Ridge tries to make slope 0 & min.  $\|\hat{y} - y\|_2^2$

Outcome - ① Reduce reliance on training data

② slope tending  $\rightarrow 0$

Features (Multiple)

$$y = w_1 x_1 + w_2 x_2 + c$$

$$L' = \|\hat{y} - y\|_2^2 + \lambda (w_1^2 + w_2^2)$$

## ② Lasso ( $L_1$ ) Regularization

$$L = \|\hat{y} - y\|_2^2 + \lambda(w)$$

Multi-linear

$$y = w_1 x_1 + w_2 x_2 + w_3 x_3$$

$$L = \|\hat{y} - y\|^2 + \lambda |w_1 + w_2 + w_3|$$

⇒ while ridge will tend towards 0,  
lasso will make those terms 0.

Ridge

$$w = \frac{-}{- + \lambda}$$

as  $\lambda \uparrow$ ,  $w \downarrow$

For  $w = 0$ ,  $\lambda = \infty$   
↓  
not possible.

⇒ value of  $w$  can't be 0  
ever.

lasso

$$w = \frac{- \pm \lambda}{-}$$

as  $\lambda \uparrow$ , unimportant  
features will be  
eliminated.

Lasso

- Helps in feature selection.
- Prevents overfitting.