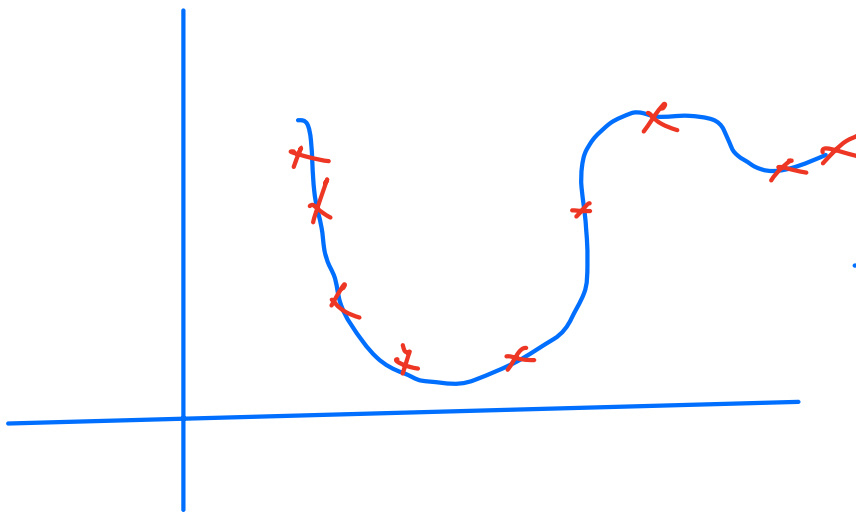


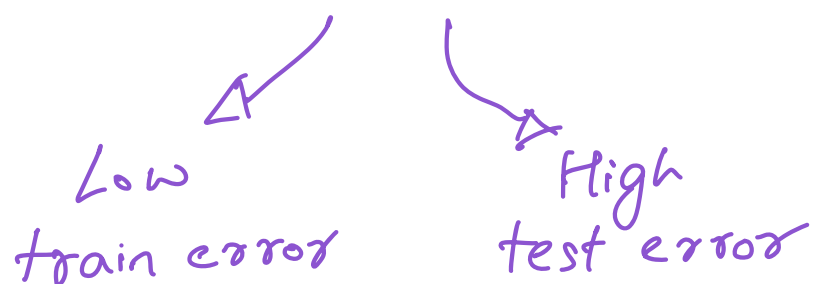
REGULARIZATION

$$y = w_1 x_1 + w_2 x_2^2 + w_3 x_1 x_2 + w_4 x_4^2$$



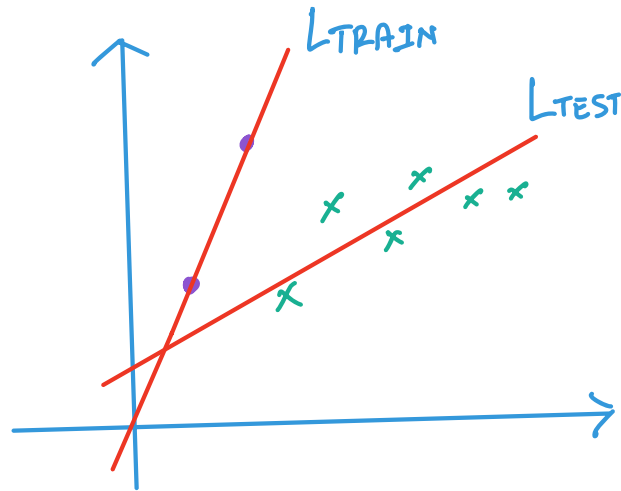
→ Can start overfitting!

- Regularization : Modification we make to the learning algorithm to prevent overfitting



① Ridge (L2)

Premise of overfitting : Too powerful model
fit on too little data



● : Train data
x : Test Data

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

With ridge (L2) regularization,

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

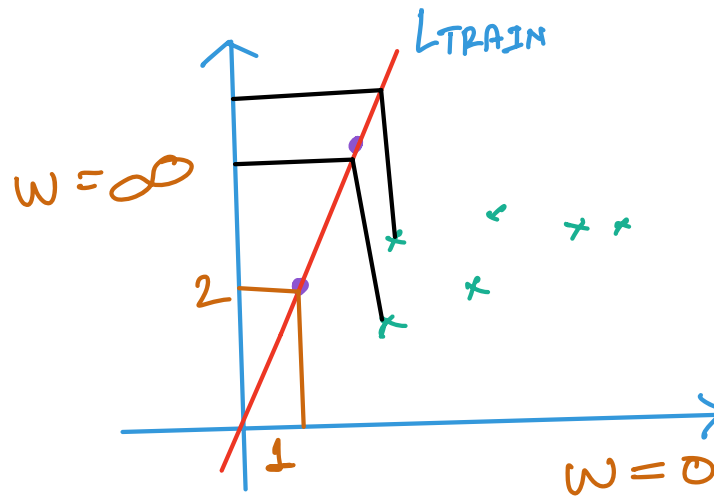
Weight Decay

Hyper param (1 to ∞)

Slope

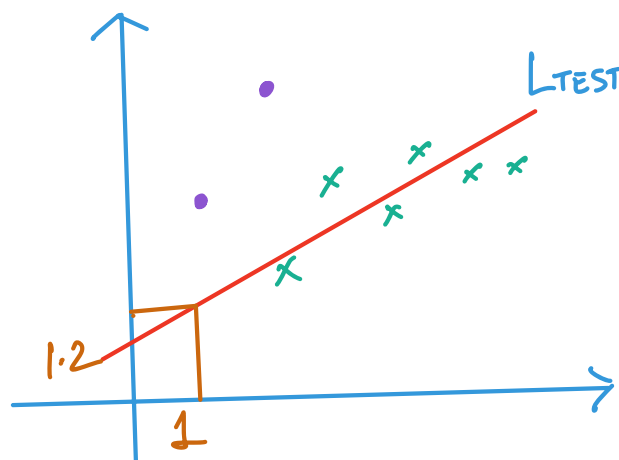
L_{TRAIN} : $y = \underbrace{2x}_{w=2} + 0.3$

$L_{\text{TRAIN}} = 0$



$L_{\text{TEST}} = \uparrow$

L_{TEST} : $y = \underbrace{1.2x}_{\text{slope has reduced!}} + 0.7$



Old line

$$y = 2x + 0.3$$

Say $\lambda = 1$,

$$L' = M^2 = 4$$

New line

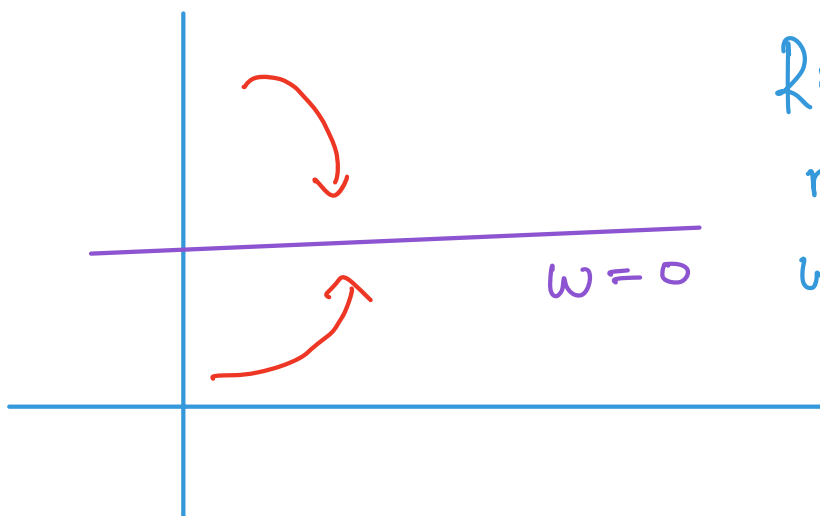
$$y = 1.2x + 0.7$$

Say $\lambda = 1$,

$L' = \text{small error}$

$$\begin{aligned} &+ \\ &(1.2)^2 \\ &= 1.5 + 1.4L \\ &= 2.94 \end{aligned}$$

\therefore New-line preferred!



Ridge tries to
make slope 0
while also
minimizing
 $\|\hat{y} - y\|_2^2$

Outcome : ① Reduce reliance on training data

② Slope tending towards 0
 \Rightarrow tending towards horizontal line

Multiple features

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

$$L^1 = \|\hat{y} - y\|_2^2 + \lambda(w_1^2 + w_2^2)$$

Lasso (L1) Regularization

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

Multi-Linear

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

$$L = \|\hat{y} - y\|_2^2 + \lambda |w_1 + w_2 + w_3|$$

→ While Ridge will tend towards 0,
Lasso will make those terms 0

Ridge

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

Lasso

$$m = \frac{\text{---} \pm \lambda}{\text{---}}$$

\Rightarrow If $w_2, w_3 = 0$

\Rightarrow We have eliminated
unimportant features

\therefore Feature Selection !

ELASTIC NET

$$L = \|\hat{y} - y\|_2^2 + \lambda_1 \|\omega\|^2 + \lambda_2 \|\omega\|$$

<u>Hyperparams</u>		<u>Default</u>
$\lambda = \lambda_1 + \lambda_2$		1

$l1_ratio = \frac{\lambda_1}{\lambda_1 + \lambda_2}$	0.5
---	-----

Used when :

- ① Unsure which one to use - L1 or L2?
- ② Large Datasets
- ③ Multi-collinearity