

ML  $\left\{ \begin{array}{l} \text{Classification} \rightarrow \text{Yes/No} \\ \text{Regression} \rightarrow \text{how much} \end{array} \right.$

## Linear Regression

eg. price of house based on size

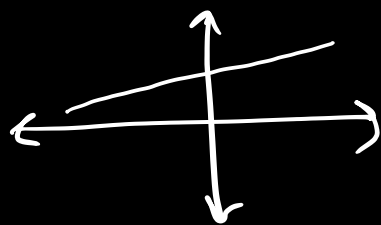
x-axis  $\rightarrow$  size      y-axis  $\rightarrow$  price

best fitting line  $\Rightarrow$  Linear Regression Model

## How to find line

select a random line and find error  
move the line

$$y = w_1 x + w_2$$



$w_1 \rightarrow$  slope

$w_2 \rightarrow$  y-intercept

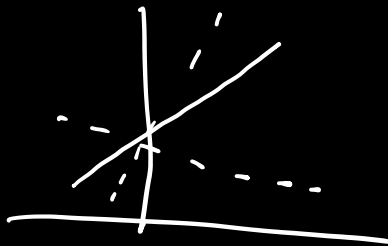
$$w_1 = \frac{\text{increase in vertical}}{\text{increase in horizontal}}$$

steepness  $\Rightarrow$  slope

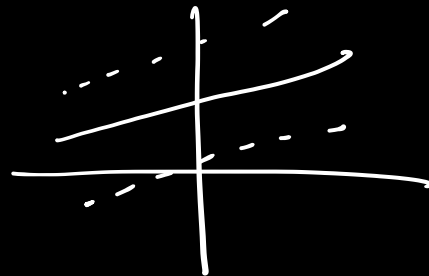
$$w_1 = \Delta y / \Delta x$$

We can move the line by change in  $w_1$  and/or  $w_2$

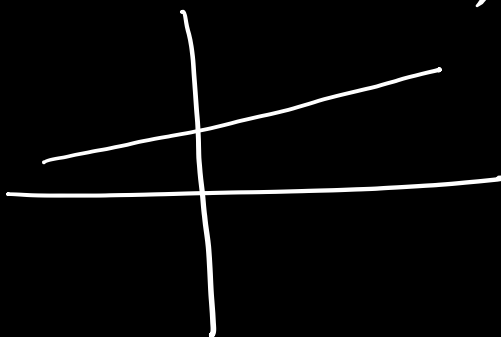
We can increase the slope or decrease the slope



if we can increase  $w_2$  or decrease  $w_1$



✓ Absolute Trick  
(p, q)



We want to move line towards (p, q)

$$y = w_1 x + w_2$$

$w_2$  to be increased by 1  
(y-intercept)  
 $w_1$  also to be increased by p

→ apply learning rate

$$y = (w_1 + p\alpha)x + (w_2 + 1\alpha)$$

If point is below the line

$$y = (w_1 - p\alpha)x + (w_2 - \alpha)$$

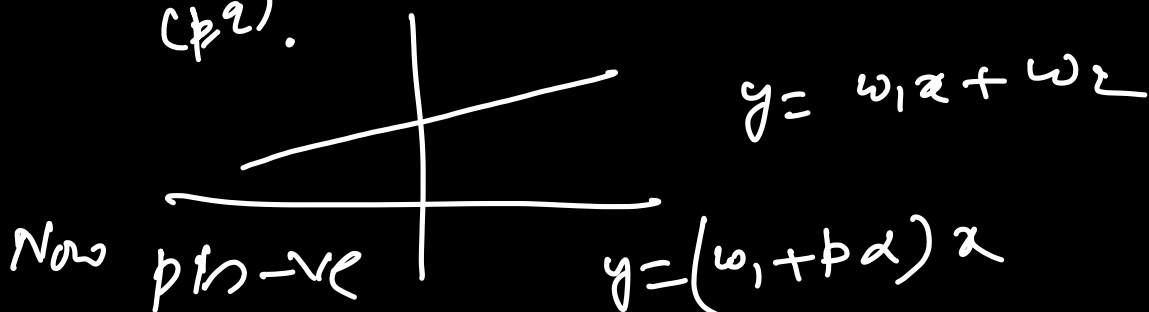
why subtract

If we reduce  $w_2$  line moves below

If we reduce  $w_1$ , line rotates in lower direction

If the point is left of y-axis

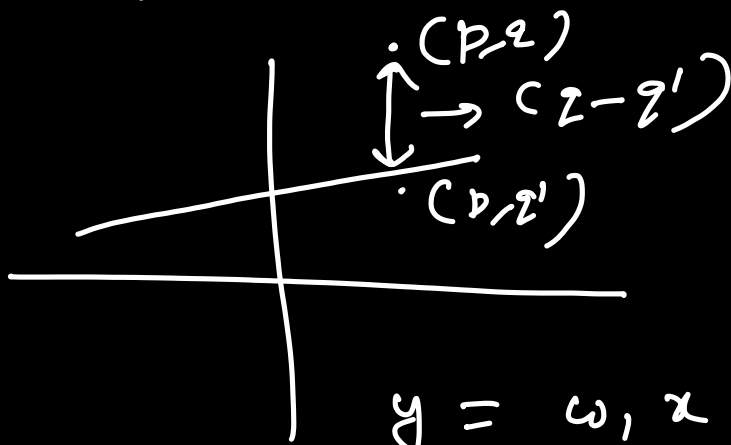
(p2).



but still add but since  $p < 0$   
slope is reduced

that's why we use  $p$

## ✓ Square Trick



$$y = w_1 x + w_2$$

$$p \alpha (z - z') + \alpha (z - z')$$

This helps to move line faster  
in the direction in comparison to  
constant move  $\alpha$

$$y = \left( w_1 + p(z - z')\alpha \right) x \\ + w_2 + (z - z')\alpha$$

If point is below line

$$z - z' < 0 \Rightarrow \text{subtraction}$$

No need of different rules

# Linear Regression

Draw a random line

Compute error

move line in one direction to  
see error

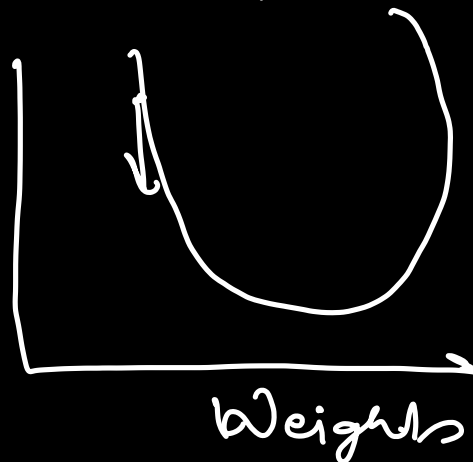
if error reduces then keep going  
else change direction

To minimize the error

We use Gradient Descent

Find the direction that can reduce  
the error

Error Function  
-ve gradient of  
error function



since gradient gives the direction  
increases most so we take  $-ve$  direction

$$w_i = w_i - \alpha \frac{\partial \text{Error}}{\partial w_i}$$

Now we have Errors

→ Mean Absolute Error

→ Mean Squared Error

Mean Absolute Error

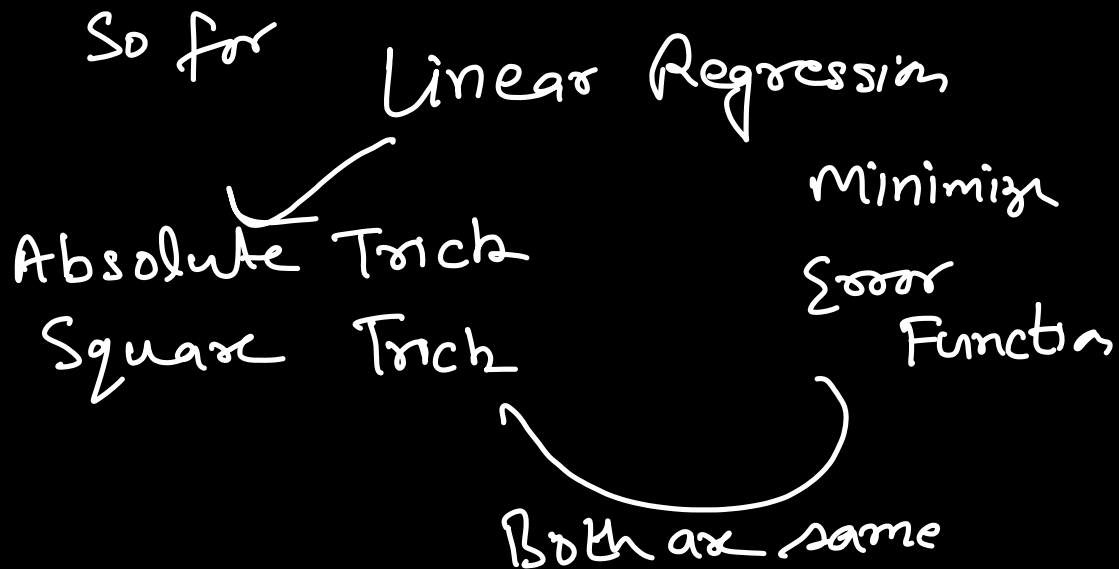
$$\text{Error} = \frac{1}{n} \sum_m |y - \hat{y}|$$

Mean Squared Error

$$\text{Error} = \frac{1}{2n} \sum_m (y - \hat{y})^2$$

$1/2$  is convenience

since no impact on minimization

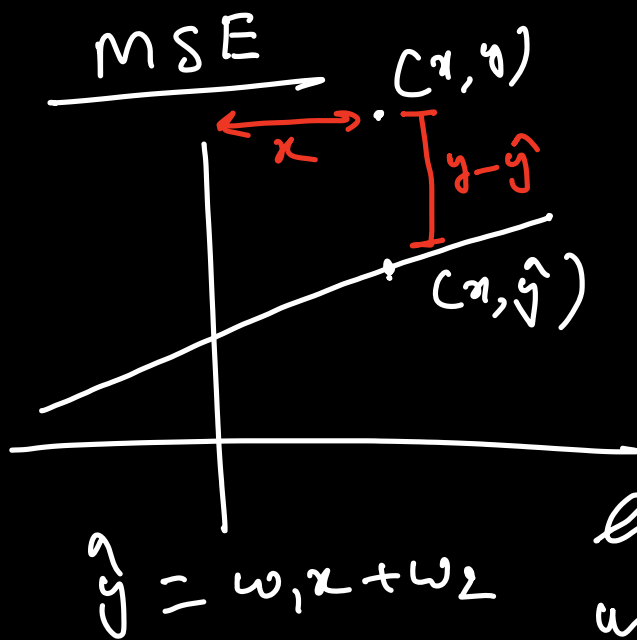


Gradient Descent on MAE

$\Rightarrow$  Absolute Trich (more constant)

Gradient Descent on MSE

$\Rightarrow$  Square Trich  
(move proportionally)



$$E = \frac{1}{2m} \sum (y - \hat{y})^2$$

$$= \frac{1}{2} (y - \hat{y})^2$$

Gradient Descent  
uses derivative

$$E = \frac{1}{2} (y - \hat{y})^2$$

and  $\hat{y} = w_1 x + w_2$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} = -(y - \hat{y}) \cdot x$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2} = -(y - \hat{y}) \cdot 1$$



## Gradient Step

$$w_i \rightarrow w_i - \frac{\partial}{\partial w} \text{Error}$$

$$w_1 = w_1 + (y - \hat{y}) \cdot x$$

$$w_2 = w_2 + (y - \hat{y})$$

earlier  
 $x = 1$   
 $y = 2$   
one  
learning  
rate

log MAE

$$E = |y - \hat{y}|$$

$$\hat{y} = w_1 x + w_2$$

$$\frac{\partial E}{\partial w_1} = \pm 1 \cdot x = \pm x$$

$$\frac{\partial E}{\partial w_2} = \pm 1 \cdot 1 = \pm 1$$

$$w_1 = w_1 \pm x$$

$$w_2 = w_2 \pm 1$$

\* learning  
rate

What's better MAE or MSE

MSE is a quadratic function  
and quadratic function has  
minimum in the middle

Now if we have many variables,

$$\begin{aligned} y &= w_1x + w_2x + w_3x + \dots + w_n \\ &= \sum_{i=1}^{n-1} w_i x_i + w_n \end{aligned}$$

$x_1, x_2, \dots, x_{n-1}$  are features

We can use MSE and use  
Gradient Descent  
of MAE

## Polynomial Regression

Consider a polynomial of higher degree

$$\hat{y} = w_1 x^3 + w_2 x^2 + w_3 x + w_4$$

Now compute MSE  
and Gradient Descent

## Regularization

works for both  
Regression and Classification

Useful technique to improve  
model

During training model takes error and minimize it

The problem is during minimization we can reduce so much that it overfits

$$m_1 \rightarrow 3x_1 + 4x_2 + 5 = 0$$

$$m_2 \rightarrow 2x_1^3 - 2x_1^2x_2 - 4x_2^3 + 3x_1^2 + 6x_1x_2 + 4x_2^2 + 5 = 0$$

$m_1$  is simple may have some error ( $E_1$ )

$m_2$  is complex may have less error ( $E_2$ )  
 $\rightarrow$

Now consider the complexity of the model to compute error

$$E_1 = E_1 + |3| + |4| \quad [\text{without constant}]$$

$$E_2 = E_2 + |2| + |-2| + |-4| + |3| + |6| + |4|$$

Adding absolute values of co-efficients  
to error term  
 $\Rightarrow$  L1 - Regularization

$\Rightarrow$  Error in  $E1$  is less than  $E2$   
So  $E1$  is less  $\Rightarrow M1$  is better

### L2 - Regularization

Add squares of the co-efficients  
instead of absolute values.

So, with L1, L2  
we punish the complex models

## Simple vs Complex Models

Medical  
Model

Sending  
Rocket on  
Moon

Video Recommendation  
or other simpler

Requires low error  
OK if Complex

Requires simple

Now how much to punish  
↳  $\lambda$  parameter

$$E1 = E1 + \lambda (\sum co - e_{ij})$$

$$E2 = E2 + \lambda (\sum co - e_{ij})$$

Now  $\lambda$  can decide punishment  
to complexity

## L1 / L2 Regularization

L1	L2
Computationally Inefficient (unless data is sparse)	Computationally Efficient (since LD can be easily applied due to derivative)
Sparse Outputs	Non-sparse Outputs
Feature Selection	No-Feature Selection