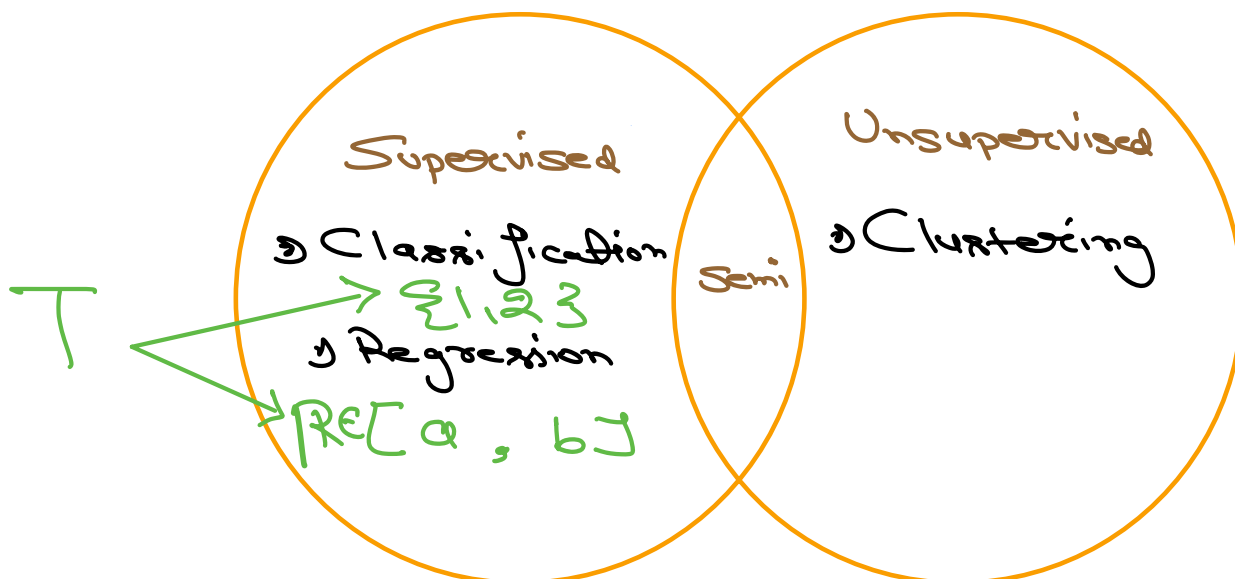


# Linear Regression 2

## Agenda

- ⇒ Recap
- ⇒ Model Interpretability
- ⇒ Feature importance
- ⇒ Scratch Implementation
  - Optimization
  - Calculate Gradient

## Types of ML



# Regression

⇒ Deals with prediction of continuous numeric values.

Ex: Stock price prediction

Car Value prediction

	$F_1$	$F_2$	$F_3$	$F_4$	$F_5 \dots F_d$	$y$
1						
2						
3						
...						
$i$	$x_i$					$y_i$
$j$						

$n \Rightarrow$  no. of rows sample

$d \Rightarrow$  no. of features

$i^{\text{th}}$  sample  $\rightarrow x_i$

$i^{\text{th}}$  label  $y_i$

$$y_i \in (0, \infty)$$

Regression

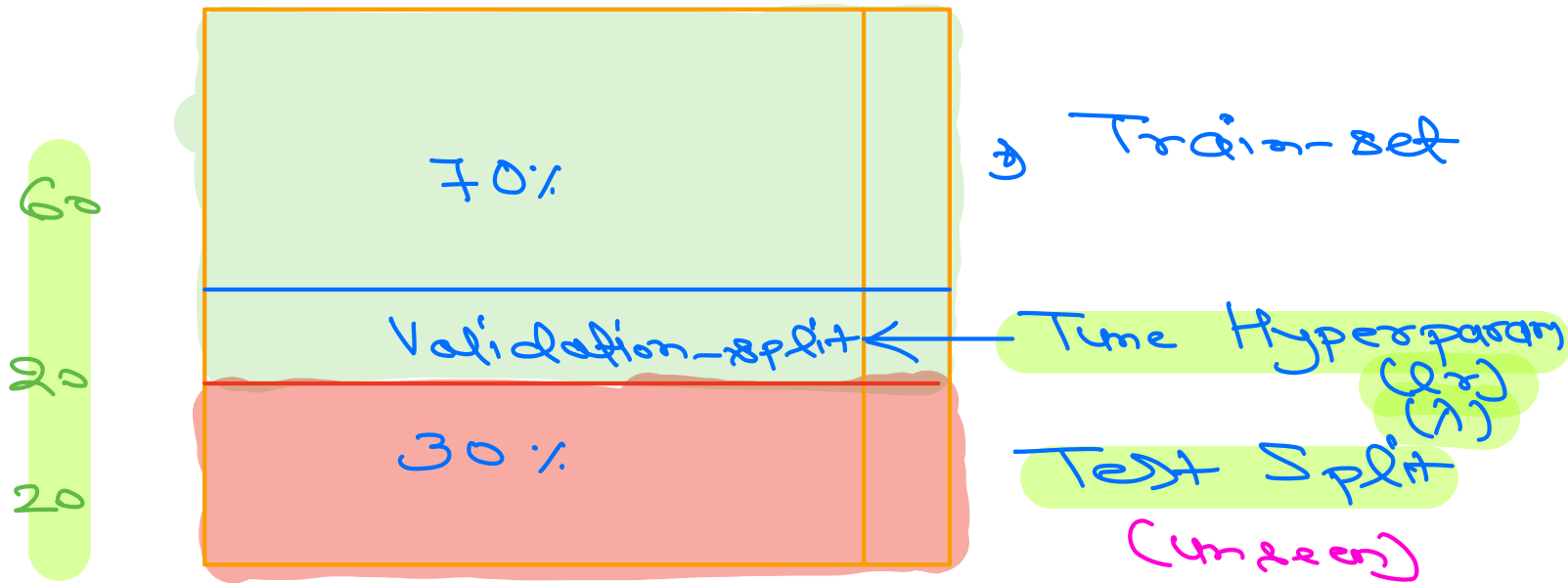
\* Linear Regression

Linear Hyperplane

$$w^T \cdot x + w_0 \Rightarrow y$$

- ① Missing value
- ② Convert Str to Numerical
- ② Normalizing

# Train Test Split



## Data Leakage

To avoid data leakage Evaluation must be done on Unseen.

Standardization  
Mean Value Imputation

$\bar{\mu}_g, \bar{\sigma}_g$   
 $\bar{\mu}_g$   
↓  
full dataset

$\mu_{train-set}$   
↓  
Imputing Test-set

$w^T x + w_0$   
↓  
Carried into

# Linear Regression

Single Variable L.R.

$$\hat{y} = w_0 + w_1(x)$$

Multi Variate L.R

$$\hat{y} = w_0 + w^T x$$

$\Rightarrow w$  is a Vector of 2 dim

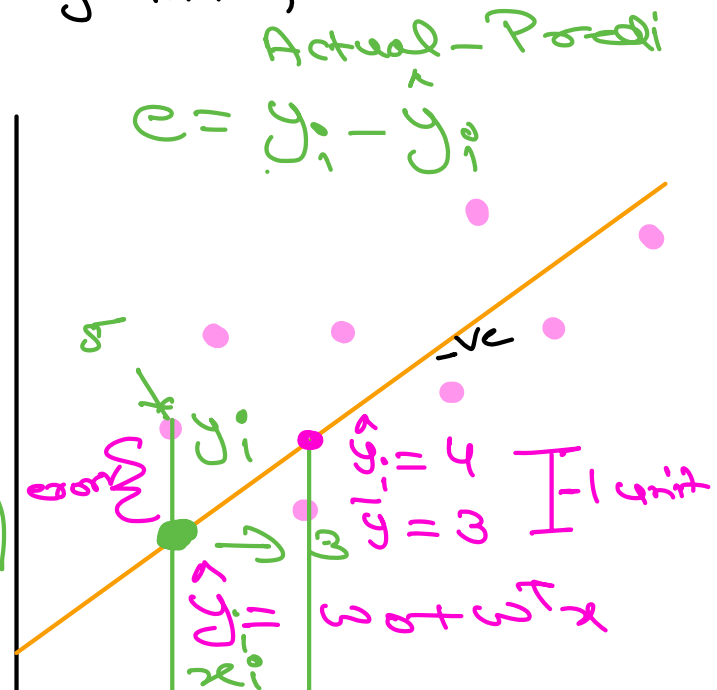
$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_0 \end{bmatrix}$$

$\Rightarrow$  num-feature  $x$

$\Rightarrow x$  is a Vector of 2 dim

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{bmatrix}$$

$\Rightarrow$  input Dataset



Goal: Find value of  $w_0$  and  $w_1$  for Best Fit line minimizing the Error

① Sum of all Error  $= \sum_{i=1}^n e_i$  ✗

② Sum of abs Error  $= \sum_{i=1}^n |e_i|$  ✓

③ Sum of Squ Error  $= \sum_{i=1}^n (e_i)^2$  ✓

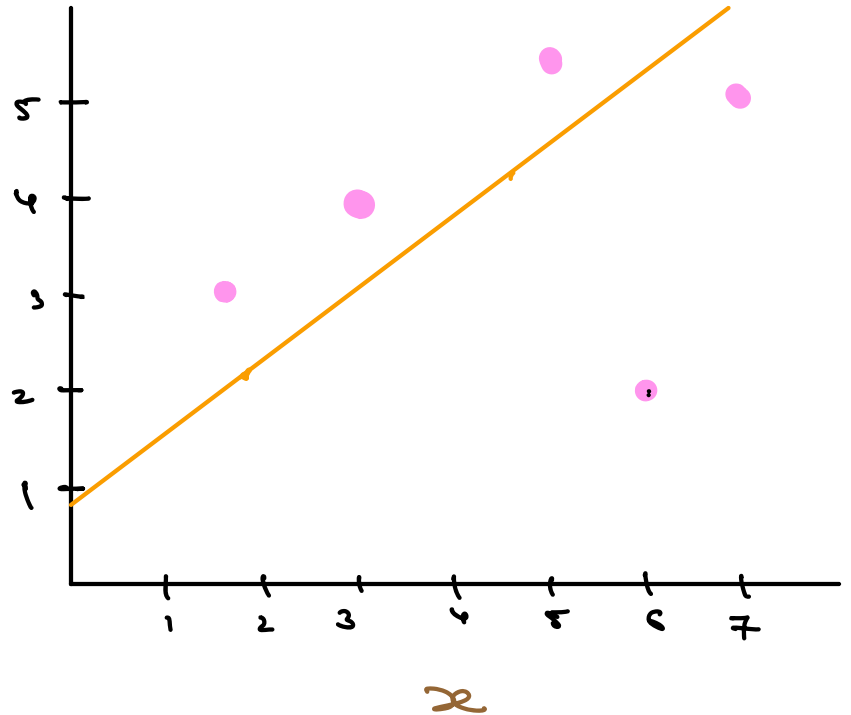
# Error or Residual

$i$	$y_i$	$\hat{y}_i$	$y_i - \hat{y}_i$
1	2	3	-1
2	4	2	2
3	5	1	4

④ 8

① 16

↓  
Prediction  
 $w^T x + w_0$



① Mean Absolute Error

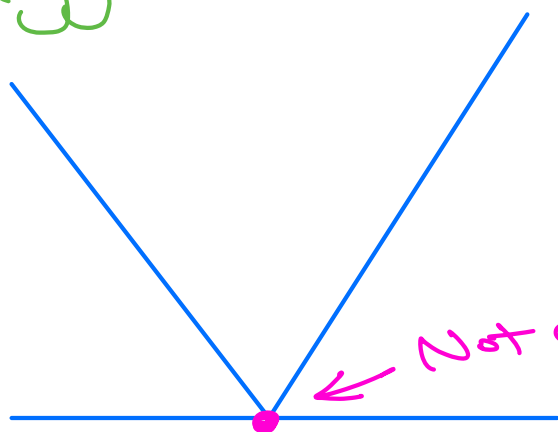
$$\frac{1}{3} \sum_{i=1}^3 |e_i| = \boxed{2.33} +$$

② Mean Squared Error

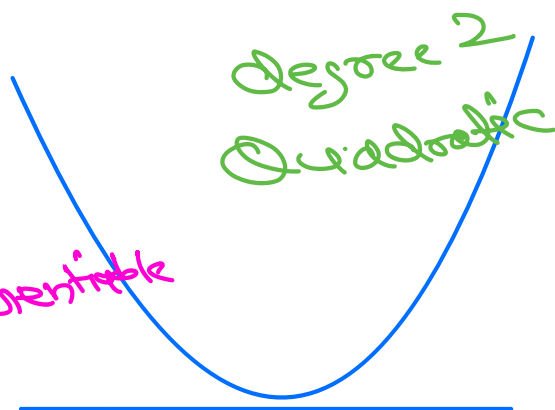
$$\frac{1}{3} \sum_{i=1}^3 (e_i)^2 = 7$$

$$21,000 \pm \sqrt{7}$$

Loss function must be  
differentiable



MAE



MSE

② Loss function as Eval Metric

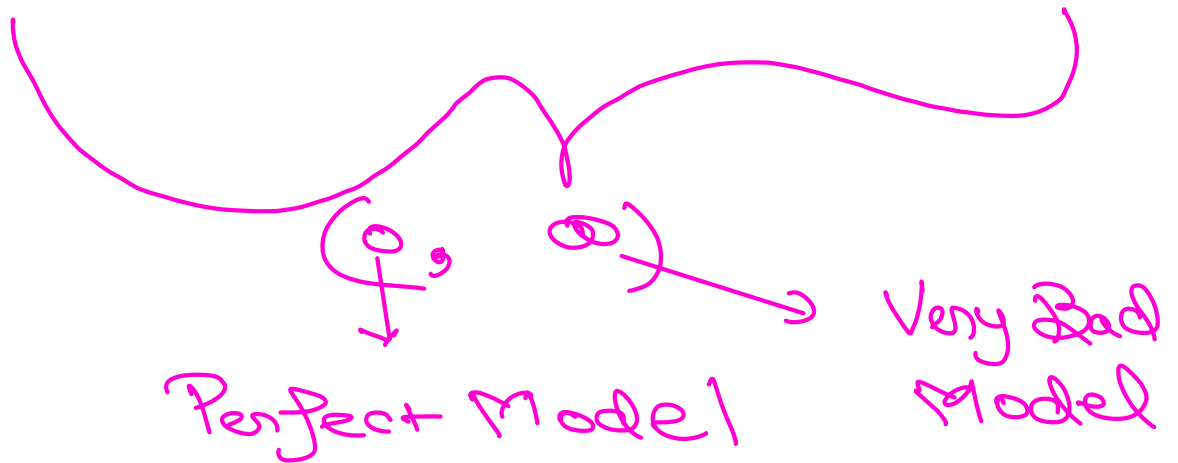
that helps  
find Best  
model

Metric  
that you  
report and

③ Regression  $\rightarrow$  MAE and MSE  
as Eval Metric

$$RMSE \Rightarrow \sqrt{MSE}$$

MSE or MAE or RMSE

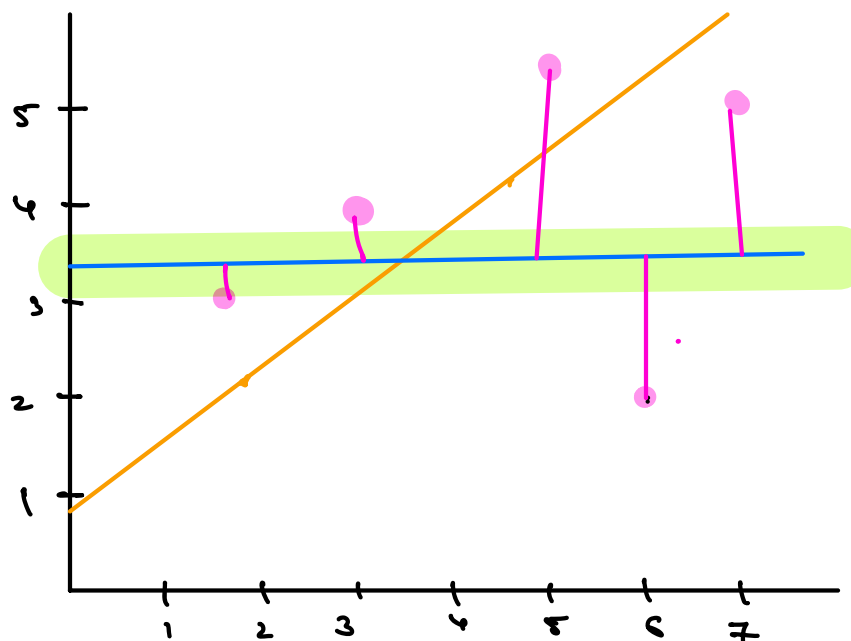


$M_1$   
 $M_2$  } Compare and decide which model

$M_1$  } **R<sup>2</sup> Score**  
(or squared Score)  
(Coefficient of Determination)

$$1 - \frac{SS_{\text{res}}}{SS_{\text{Total}}}$$

$$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$



Case-1:

$$1 - \frac{SS_{\text{res}}}{SS_{\text{Total}}}$$

$$e \perp R) \Rightarrow SS_{\text{res}} \Rightarrow 0$$

$$\Rightarrow 1 - \frac{0}{SS_{\text{Total}}} \Rightarrow 1$$

Good Model will have  
score close to 1

Case-2:

$$SS_{\text{res}} = SS_{\text{Total}}$$

$$\Rightarrow 1 - \frac{SS_{\text{res}}}{SS_{\text{Total}}}$$

$$\Rightarrow 0$$

Bad Model will have score  
close to Zero

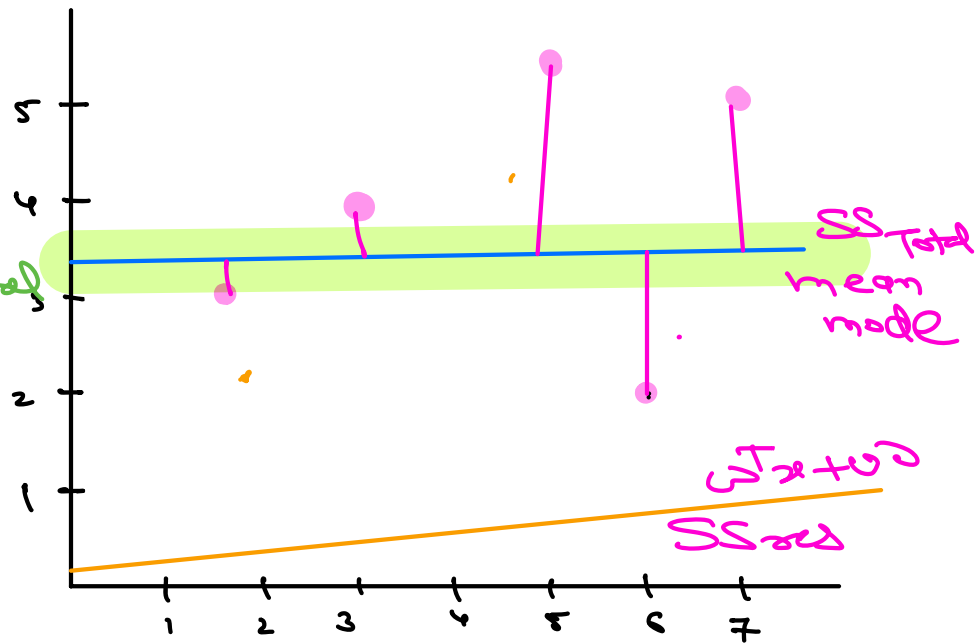


Case-3

$$SS_{Res} > SS_{Total}$$

$$1 - \frac{SS_{Res}}{SS_{Total}}$$

↓  
 $(-\infty, 0)$



$$R^2 \text{ score } \in (-\infty, 1)$$

↓  
 Practical Value  $\rightarrow (0, 1)$

# Model Interpretability and Weights

$$\hat{y} = w^T \cdot x + w_0$$

$$\hat{y} = w_1 x_1 + w_2 x_2 + \dots + w_0$$

Case-Study

$$\hat{y} = w_0 + w_1 x_1$$

$$+ \left[ \begin{array}{l} (-10000) * \text{Age} + 1 \\ (-10) * \text{odometer} + \end{array} \right] \quad \begin{array}{l} -10k \text{ (Age)} \\ -10 \text{ (odo)} \end{array}$$

»

»

$$w_d x_d$$

000 +1 @ -10 Dollar

Age +1 @ -10000 Dollar

Case  $\Rightarrow$  -ve  $\omega_+$

$x \uparrow$        $\hat{y} \downarrow$

$x \downarrow$        $\hat{y} \uparrow$

-ve  $\omega_T \Rightarrow x \propto \frac{1}{y}$

Case  $\Rightarrow$  +ve  $\omega_+$

$x \uparrow$        $\hat{y} \uparrow$

$x \downarrow$        $\hat{y} \downarrow$

-ve  $\omega_T \Rightarrow x \propto \hat{y}$

Case

$$w^T = 0$$

$$w_1 x_1 + \cancel{w_2 x_2} + w_0 = \hat{y}$$

$\downarrow$   
 $w_2 = 0$

feature has no impact  
on Target

Magnitude Cases

Age wt @ -10000

Odometer @ -10

engine\_capacity @ 1000

$$\begin{cases} \text{Age} + 1 @ \hat{y} = 10000 \\ \text{Odom} + 1 @ \hat{y} = 10 \end{cases}$$

Age (0, 20 years) +1 @ 1000

Qd @ (0, 100,000)

## Solution Scaling

Age (0, 20 years) @ (-1 1) 1.20

Qd @ (0, 100,000) @ (-1 1) 0.5

Now we can compare  
Magnitude weight to  
fairly Estimate importance

- ① Magnitude of Normalized feature weight can be used to determine feature importance

## Feature importance in linear regression is determined by :

4 options

Active Duration (Most preferred: 30 seconds)

Appears for 45 Secs

- A The magnitude of the regression coefficients. ✓
- B The number of observations in the dataset.
- C The correlation between the independent variables. ↗
- D The average squared difference between the predicted and actual values.

Assumption of LR

## Quiz time!

Time Left: 0s

When assessing model interpretability in Linear Regression, what is the impact of feature scaling?

51 users have participated

- A Feature scaling does not affect model interpretability 10%
- B Feature scaling improves model interpretability 27%
- ✓ C Feature scaling can help compare the magnitudes of different coefficients 63%

End Quiz Now

Consider the following Linear Regression model equation:  $y = 5.2x_1 - 3.8x_2 + 2.1x_3 + 0.01x_4 + 1.5$  if we were to drop one feature, which one would be the best choice ?

9 users have participated

- A  $x_1$  22%
- B  $x_2$  11%
- C  $x_3$  0%
- ✓ D  $x_4$  67%

End Quiz Now

min magnitude

# Optimization

## Gradient Descent

- ① Pick  $\bar{w}$  and  $w_0$  Randomly
- ② Calculate  $\frac{\partial L}{\partial \bar{w}}$  and  $\frac{\partial L}{\partial w_0}$
- ③ Update Step

$$\begin{aligned}\bar{w}_{t+1} &\rightarrow \bar{w}_t - \eta \frac{\partial L}{\partial \bar{w}} \\ w_{t+1} &\rightarrow w_t - \eta \frac{\partial L}{\partial w_0}\end{aligned}$$

Converging to minima  
Rench —

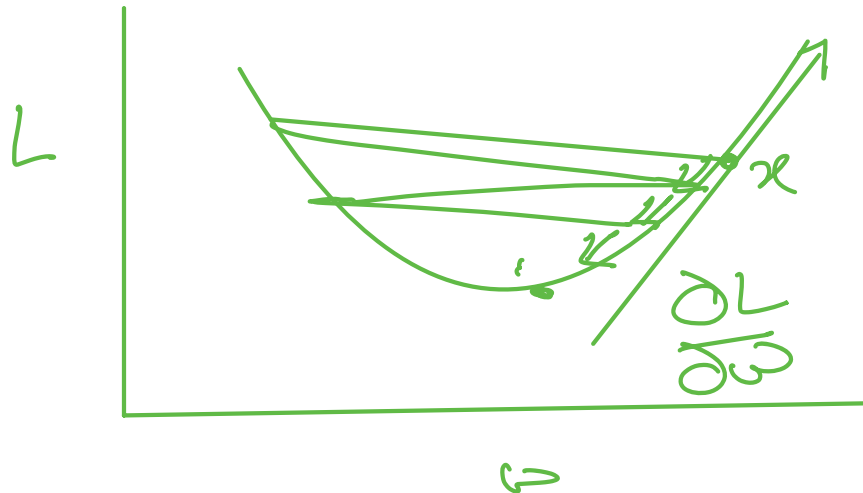
# How do we find Global Minima?

$$\text{Loss} \Rightarrow \text{MSE} \Rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$

In gradient descent, what does the gradient represent ?

43 users have participated

- ☒ A The direction of steepest increase of the cost function 35% ✓
- ☐ B The direction of steepest decrease of the cost function 58%
- ☐ C The number of training examples in the dataset 2%
- ☐ D The number of layers in the neural network 5%



$$\hat{y} = w^T \cdot x + w_0$$

Score ✓

Predict ✓



Loss of MSE of  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$

↓

$\frac{\partial L}{\partial \bar{y}}$  ?

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

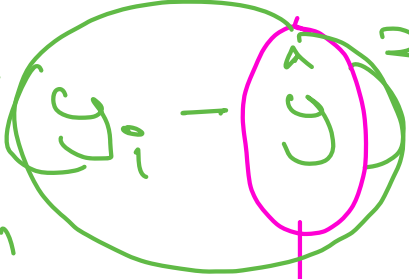
①  $\frac{\partial L}{\partial \bar{y}}$

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

↓

$\frac{1}{n} (y_1 - \omega_0 x_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_d x_d + \omega_0)$


$$\frac{\partial L}{\partial w_1} = \frac{1}{3} \sum_{i=1}^3 (y_i - \hat{y})^2 \quad (f(x))'^2$$



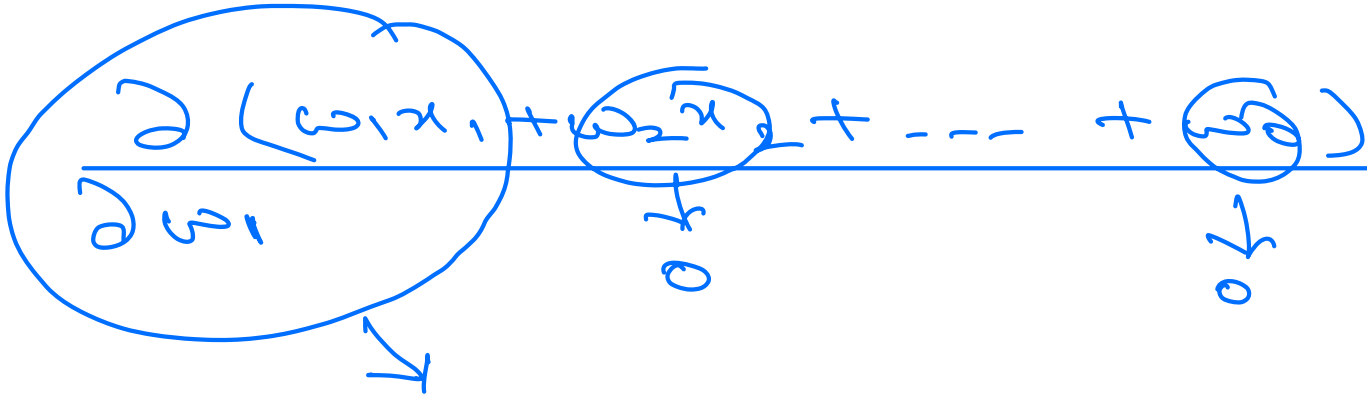
$w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$

$$f(g(x)) = f'(g(x)) \times g'(x)$$

$$\frac{\partial L}{\partial w_1} = 2 \times (y_i - \hat{y}) \times \left( \frac{\partial y_i}{\partial w_1} - \frac{\partial \hat{y}}{\partial w_1} \right)$$



$$\frac{\partial (w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b)}{\partial w_1}$$



$\frac{\partial \hat{y}}{\partial w_1}(x_1)$

$$\frac{\partial \mathcal{L}}{\partial w_2} \text{ ଓ } x_2$$

$$\frac{\partial \mathcal{L}}{\partial w_2} \text{ ଓ } x_2$$

$$\frac{\partial \mathcal{L}}{\partial w_0} \text{ ଓ } 1$$

ପ୍ରମୁଖ ଓ ପ୍ରମୁଖ

$$\frac{\partial \mathcal{L}}{\partial w_i} \text{ ଓ } \frac{\partial}{\partial w_i} \sum (y_i - \hat{y})^2 (-x_i)$$

$$0 - 2(y_i - \hat{y}) \cdot x_i$$

$$\frac{\partial \mathcal{L}}{\partial w_2}$$

$$-2(y_i - \hat{y}) \cdot x_i$$

$$\frac{\partial \mathcal{L}}{\partial w_0}$$

$$-2(y_i - \hat{y})$$

•

1



1-