

Course 2

Calculus

Week 1 - Derivatives and Optimization

Calculating the Slope

$$\text{Slope} = \frac{\text{rise} \uparrow}{\text{run} \rightarrow}, \quad \text{Slope} = \frac{\text{change in distance } (\Delta x)}{\text{change in time } (\Delta t)}$$

Derivatives

$$\text{Slope} = \frac{\text{change in distance}}{\text{change in time}} \rightarrow \text{Slope} = \frac{\Delta x}{\Delta t}$$

(vertical change)
(horizontal change)

$$\text{Slope at a point} = \frac{dx}{dt} \quad (\text{derivative})$$

Lagrange's notation

$$\text{Notation: Function: } y = f(x) \rightarrow \text{Derivative of } f: \frac{d f(x)}{dx} = f'(x)$$

Leibniz's notation

Some Common Derivatives

$$1) \quad y = f(x) = c \quad (\text{line}) \rightarrow f'(x) = 0$$

$$2) \quad f(x) = ax + b \quad (\text{line}) \rightarrow f'(x) = a$$

$$3) \quad f(x) = x^2 \quad (\text{quadratic}) \rightarrow f'(x) = 2x$$

$$4) \quad f(x) = x^3 \quad (\text{cubic}) \rightarrow f'(x) = 3x^2$$

$$5) \quad f(x) = \frac{1}{x} \rightarrow f'(x) = \frac{-1}{x^2} = -x^{-2}$$

6) Power Func : $f(x) = x^n \longrightarrow f'(x) = nx^{n-1}$

* Inverse Function: $g(x)$ and $f(x)$ are inverses $\rightarrow g(x) = f^{-1}(x)$
 $\rightarrow g(f(x)) = x$

7) Inverse : $f(x)$, $g(y) \longrightarrow g'(y) = \frac{1}{f'(x)}$
 are inverses

Ex: $f(x) = x^2$, $g(y) = \sqrt{y} = y^{\frac{1}{2}}$

$$\begin{array}{ccc} \downarrow & \downarrow & \Rightarrow g'(y) = \frac{1}{f'(x)} \\ f'(x) = 2x & g'(y) = \frac{1}{2} y^{-\frac{1}{2}} = \frac{1}{2\sqrt{y}} & \end{array}$$

\rightarrow at the point $(2, 4)$: $g'(4) = \frac{1}{4}$ $f'(2) = 4 \rightarrow g'(y) = \frac{1}{f'(x)}$

8) Trigonometric Functions:

• $f(x) = \sin(x) \longrightarrow f'(x) = \cos(x)$

• $f(x) = \cos(x) \longrightarrow f'(x) = -\sin(x)$

9) Exponential: $f(x) = e^x \longrightarrow f'(x) = e^x$

* Logarithm: $\log(x)$ is the inverse of e^x

10) Logarithm: $f(x) = \log(x) \longrightarrow f'(x) = \frac{d}{dx} \log x = \frac{1}{x}$

* Differentiable Functions:

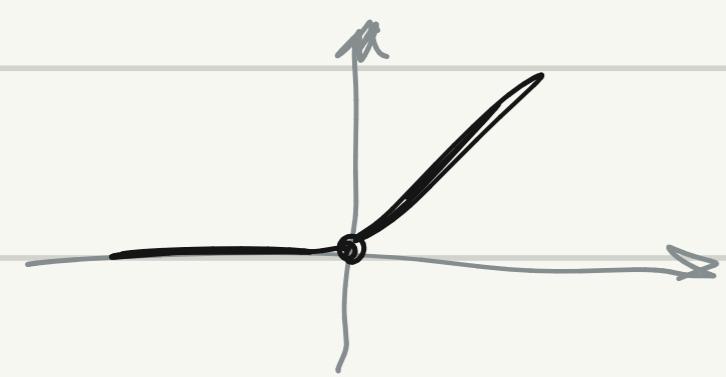
→ at a point: The derivative has to exist for that point.

→ at an interval: The derivative has to exist for every point in the interval.

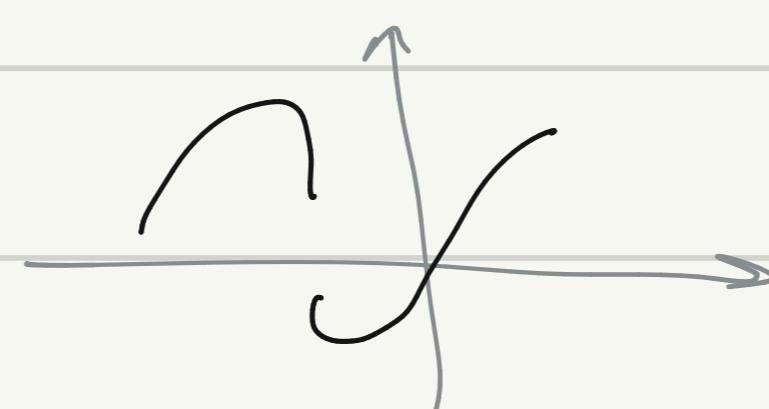
* Non Differentiable Functions:

Generally, when a Function has a corner or a cusp, the function is not differentiable at that point.

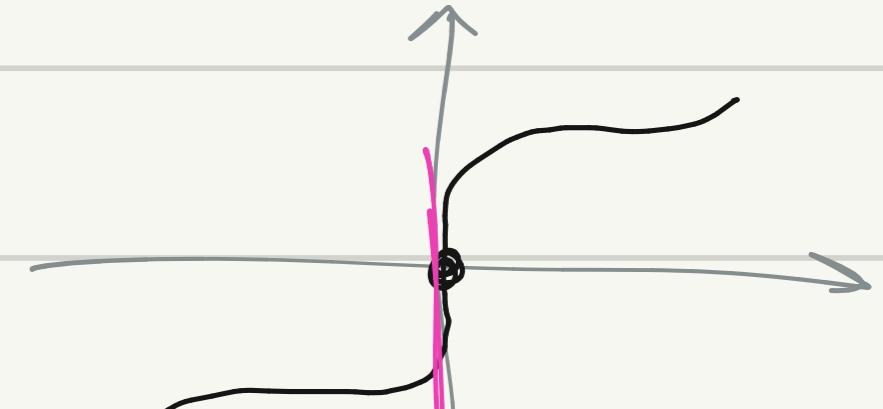
1) corners / cusps



2) jump discontinuity



3) Vertical tangents



11) The Sum Rule:

$$f = g + h \longrightarrow f' = g' + h'$$

12) The Product Rule:

$$f = gh \longrightarrow f' = g'h + gh'$$

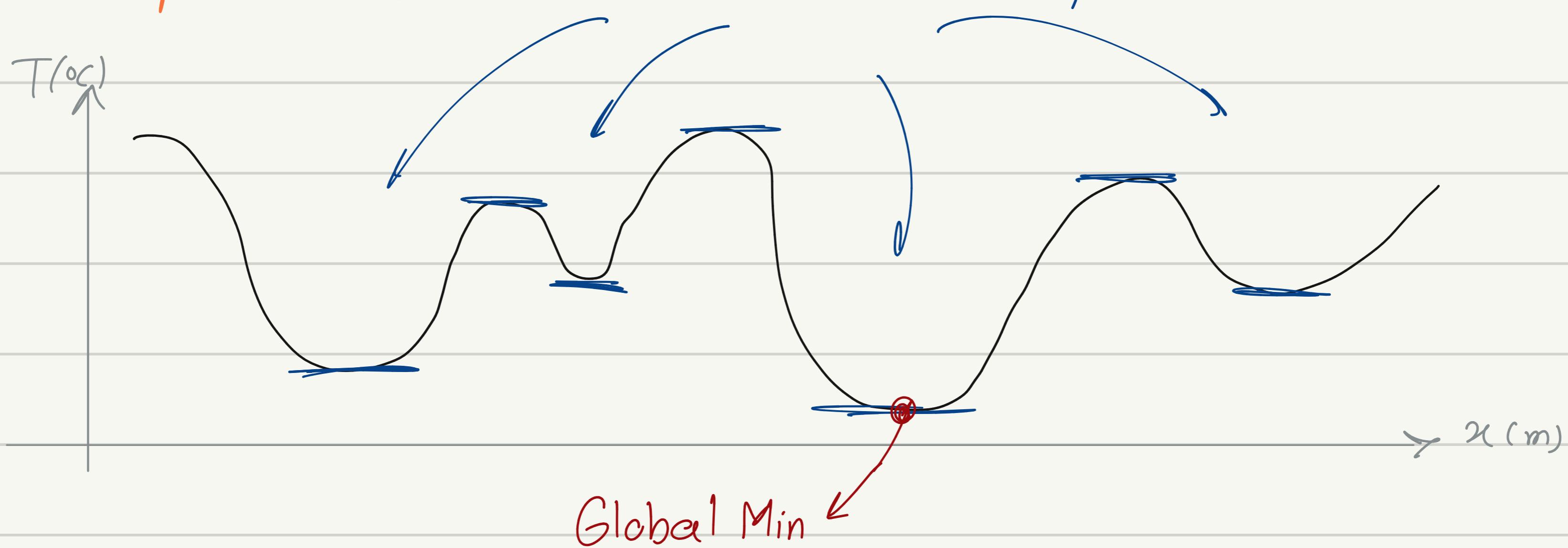
13) The chain Rule:

$$f(g(h(t))) \cdot g'(h(t)) \cdot h'(t)$$

$$f(g(h(t))) \longrightarrow \frac{d}{dt} f(g(h(t))) = \frac{df}{dg} \cdot \frac{dg}{dh} \cdot \frac{dh}{dt} =$$

Optimization

→ Multiple Minima:



→ Optimization of squared loss:

- Total cost: $(x-a)^2 + (x-b)^2$ Goal: Find \min_{so} $\frac{d}{dx} [(x-a)^2 + (x-b)^2] = 0$

Solution $\Rightarrow x = \frac{a+b}{2}$

- Cost Func: $(x-a)^2 + (x-b)^2 + (x-c)^2$ Goal: minimize func $\frac{d}{dx} \text{Cost Func} = 0$

Solution $\Rightarrow x = \frac{a+b+c}{3}$

So \Rightarrow Minimize $(x-a_1)^2 + (x-a_2)^2 + \dots + (x-a_n)^2$

Solution: $x = \frac{a_1 + a_2 + \dots + a_n}{n}$

→ optimization of log loss:

- Func: $g(p) = p^7(1-p)^3$ Goal: Minimize $g'(p) = 0$ * But it's hard to solve.

↓ Better way

$$G(p) = 7 \log(p) - 3 \log(1-p) \xrightarrow{\text{minimize}} G'(p) = 0 \quad * \text{Easier to solve.}$$

Week 2 : Gradients and Gradient Descent

Partial Derivatives

Functions with one variable : $f(x)$ $\xrightarrow{\text{derivative}}$ $f'(x) = \frac{df}{dx}$

* We calculate derivative for one variable functions with using the rules that we've learned in previous chapter.

Functions with two variables : $f(x, y)$ $\xrightarrow{\text{derivative}}$ $f'(x, y) = ?$

* We calculate derivative for two variables functions with Partial Derivatives.

What is partial derivatives ?

Ex. $f(x, y) = x^2 + y^2$ $\rightarrow f'(x, y) = ?$

step 1: Treat y as a constant $\Rightarrow f$ become one variable function

step 2: Differentiate the function using the normal rules.

step 3: Repeat step 1, 2 for the next variable (x)

$$f(x, y) = \underline{x^2} + \underline{y^2}$$

constant

and

$$f(x, y) = \underline{x^2} + \underline{y^2}$$

constant

$$\frac{\partial f}{\partial x} = 2x$$

$$\frac{\partial f}{\partial y} = 2y$$

⇒ Partial Derivatives Notation :

partial derivative of f with respect to y .

$$f(x, y)$$

partial derivative of f with respect to y .

$$f_x = \frac{\partial f}{\partial x}$$

$$f_y = \frac{\partial f}{\partial y}$$

pronounced
"die"

* The partial derivative is denoted by the symbol $\underline{\partial}$, which replaces the roman letter d used to denote a full derivative.

Ex. $f(x, y) = 3x^2y^3$

$$\frac{\partial f}{\partial x} = 3(2x)y^3 \quad / \quad \frac{\partial f}{\partial y} = 3x^2(3y^2)$$

② Gradient ②

$$f(x, y)$$

$$\longrightarrow$$

gradient :

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$$

Ex. Find the gradient of $f(x, y) = x^2 + y^2$ at $(2, 3)$:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x \\ 2y \end{bmatrix} \Big|_{(2, 3)} = \begin{bmatrix} 4 \\ 6 \end{bmatrix}$$

① Gradients and min / max ①

$$\underline{f'(x) = 0}$$

1) Functions of one variable : $f(x)$ \rightarrow min is when slope = 0
 slope = derivative \leftarrow

2) Functions of two variables : $f(x, y)$ \rightarrow min is when both slopes = 0

$$\frac{\partial f}{\partial y} = 0 \quad \text{and} \quad \frac{\partial f}{\partial x} = 0 \quad \begin{matrix} \leftarrow \\ \text{means} \end{matrix} \quad \underline{f'(x, y) = 0}$$

Ex. $f(x) = x^2$

$$f(x, y) = x^2 + y^2$$

$$f'(x) = 2x = 0$$

$$\frac{\partial f}{\partial x} = 2x = 0 \quad / \quad \frac{\partial f}{\partial y} = 2y = 0$$

$$x = 0 \quad \underbrace{\min}_{\text{min}}$$

$$(x, y) = (0, 0) \quad \underbrace{\min}_{\text{min}}$$

② Optimization with gradients ②

Week 3 : Optimization in Neural Networks and Newton Method

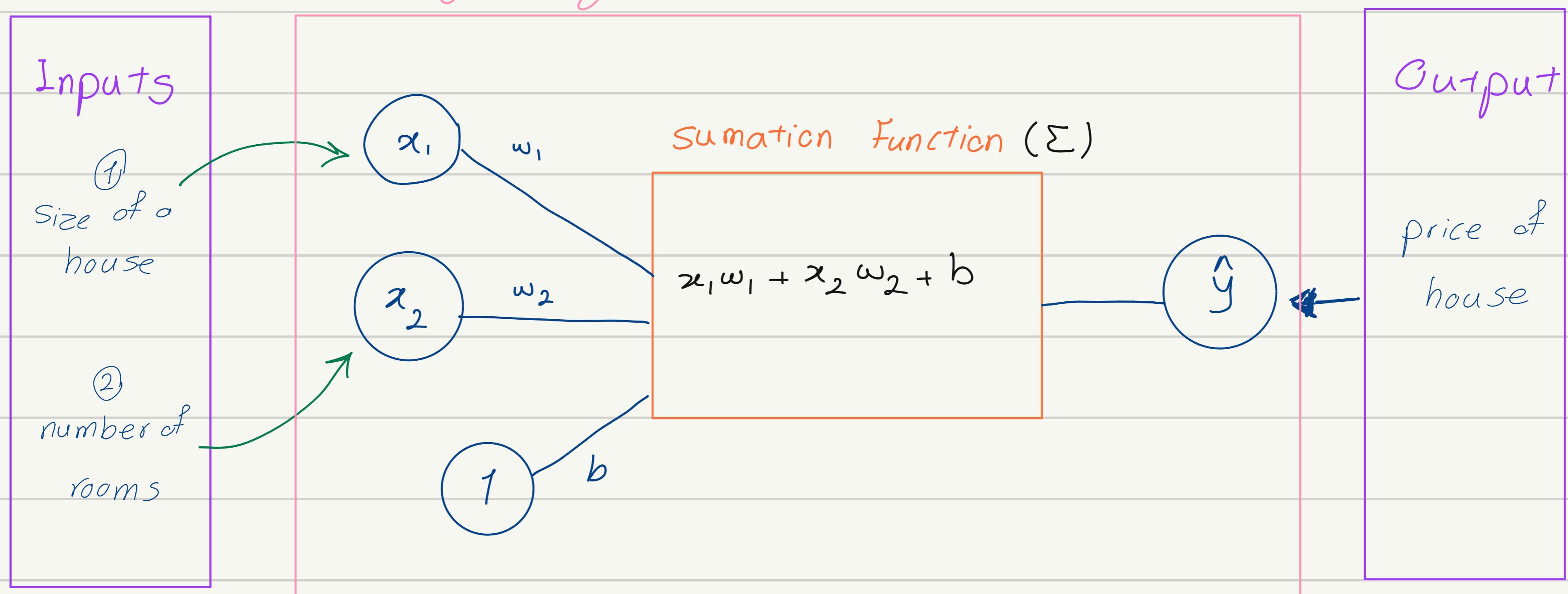
Regression with a Perceptron

Ex. We want to predict the price of a house from 2 inputs:

- 1) the size of the house
- 2) number of rooms

Features
=

Single layer neural network perceptron



* x_1, x_2 : inputs / features * \hat{y} : output

* Σ : Summation Function * $1, b$: A bias term

* w_1, w_2 : weights that determine how important the input is for the output.

Goal: Find weights and bias

So: $\hat{y} = x_1 w_1 + x_2 w_2 + b \rightarrow$ that will optimise the predictions.

or: Reduce the errors in the predictions. → what's error? The error is basically how far are you from the price of the house?

• Mean Squared Error •

Error : subtract the prediction from the actual value

$$\rightarrow \underset{\text{actual}}{y} - \underset{\text{prediction}}{\hat{y}}$$

\Rightarrow The lesser the error , the better the model.

Error's Problem ? Errors can be positive or negative so

adding them can get us zero or small numbers that can get

confusing . So for that we square the Error.

Squared Error : $\frac{1}{2} (y - \hat{y})^2$

What is $\frac{1}{2}$? When we take the derivative of $(y - \hat{y})^2$, we

get a lingering two . So we put $\frac{1}{2}$ there to cancel with that.
 $\cancel{2}$

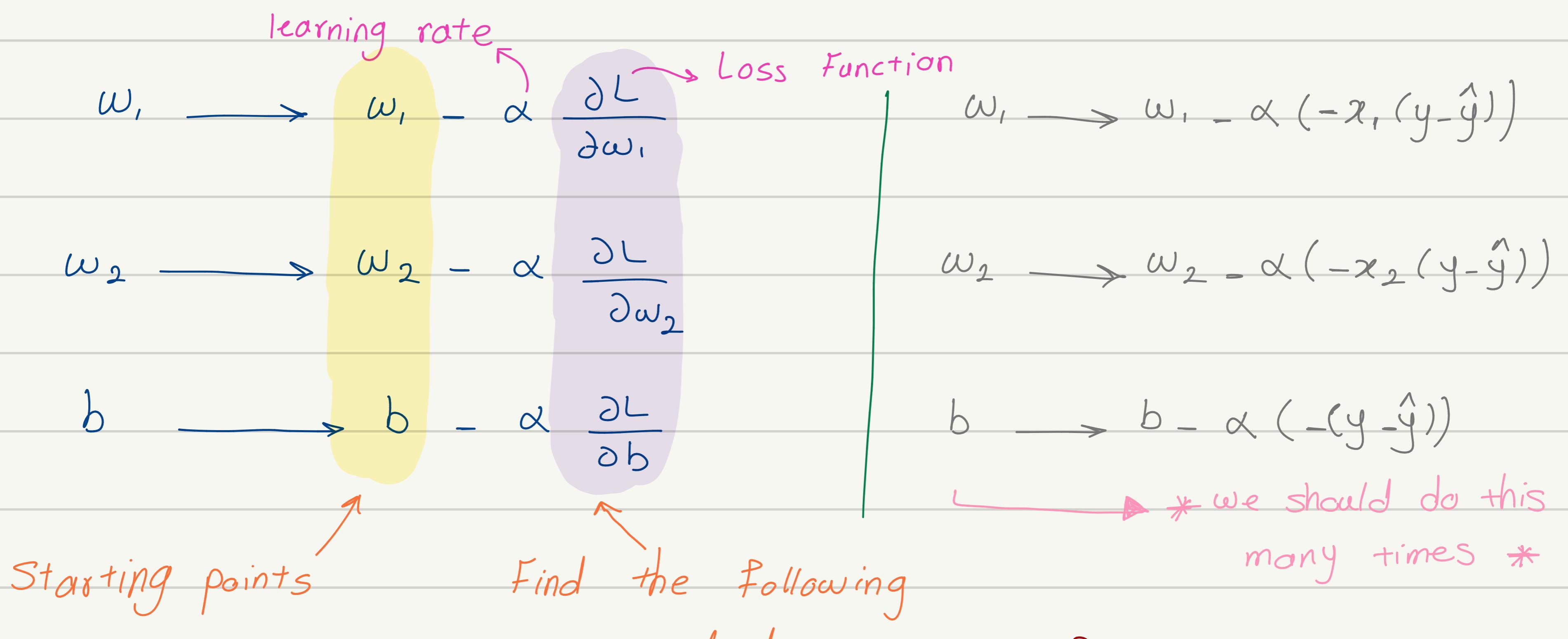
Loss Function : For every point is : $L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$

\Rightarrow Main Goal : Minimize the loss function

And as we said in last page: Find w_1, w_2, b that give \hat{y} with the least error.

So to find optimal values for: ω_1, ω_2, b

→ we must use Gradient Descent.



Finding Partial Derivatives : Using chain rule

■ Prediction Function:

$$\hat{y} = x_1 \omega_1 + x_2 \omega_2 + b$$

■ Partial Derivatives:

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} = -(y - \hat{y})$$

■ Loss Function :

$$\frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \omega_1} = -(y - \hat{y}) x_1$$

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

$$\frac{\partial L}{\partial \omega_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \omega_2} = -(y - \hat{y}) x_2$$

■ Calculation :

$$\frac{\partial L}{\partial \hat{y}} = -(y - \hat{y})$$

$$\frac{\partial \hat{y}}{\partial b} = 1$$

$$\frac{\partial \hat{y}}{\partial \omega_1} = x_1$$

$$\frac{\partial \hat{y}}{\partial \omega_2} = x_2$$

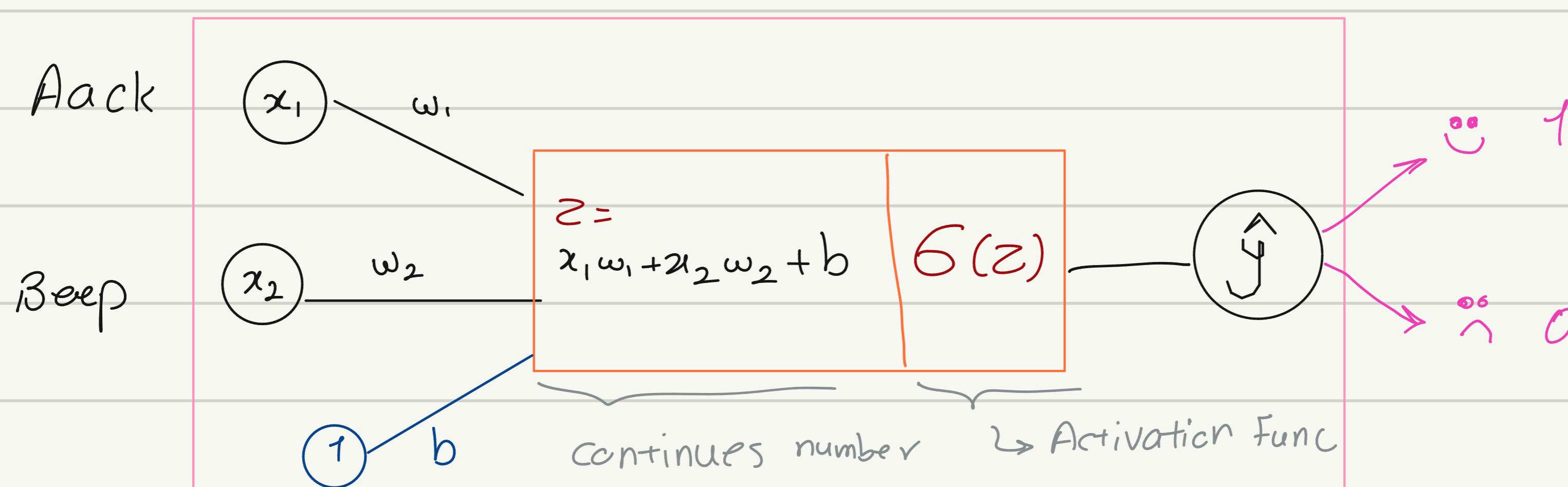
* continue

Classification with a Perceptron

Example: Alien Speaking (Emotion)

Sentence	Aack	Beep	Mood
Aack Aack Aack!	3	0	Happy 😊
Beep Beep !	0	2	Sad ☹
Aack Beep Beep Beep!	1	3	Sad ☹
Aack Beep Aack!	2	1	Happy 😊

* Models take numbers not words.

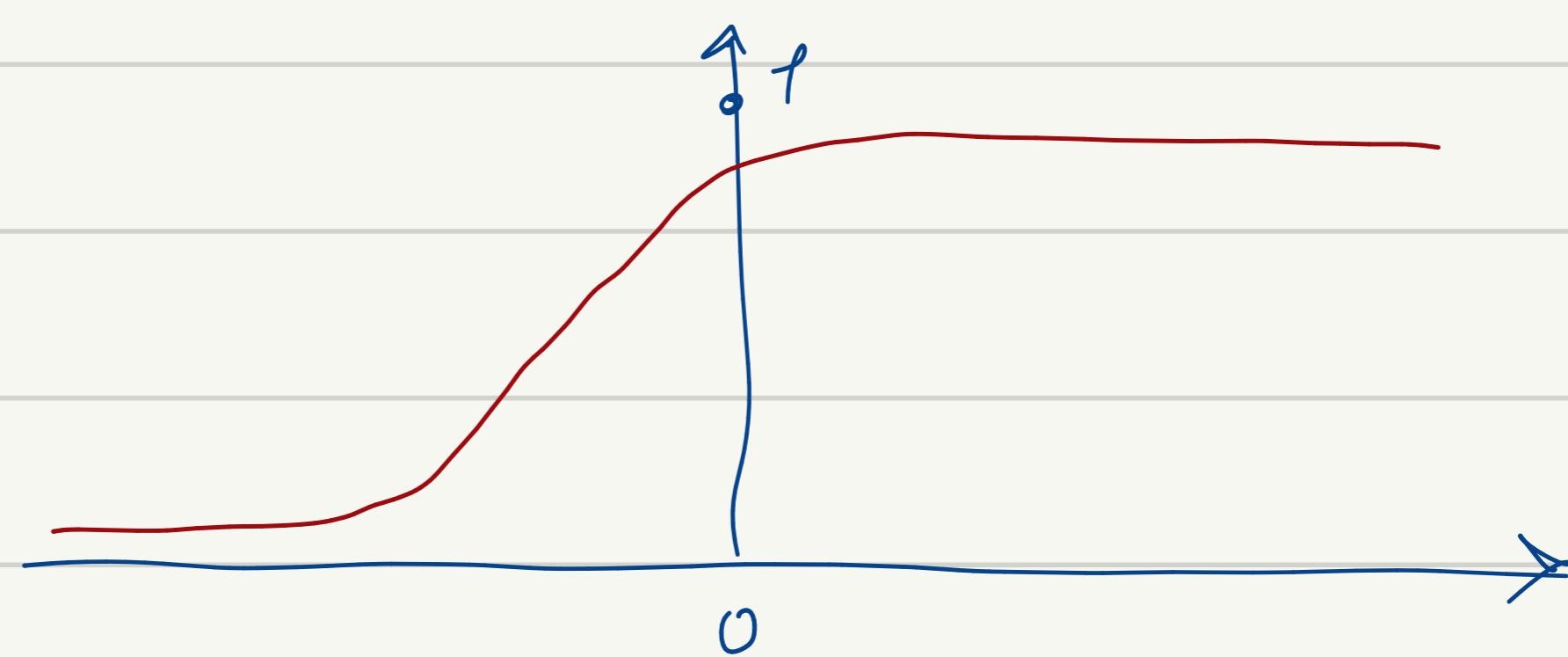


$\sigma(z)$: Sigmoid of z . This Function is going to take all the

numbers and crunch them into the interval 0, 1.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\Rightarrow 0 < \sigma(z) < 1$$



② Derivative of a Sigmoid Function ②

$$G(z) = \frac{1}{1+e^{-z}} \longrightarrow G(z) = (1+e^{-z})^{-1}$$

$$\begin{aligned} \text{So: } \frac{dG(z)}{dz} &= \frac{d}{dz} (1+e^{-z})^{-1} = -1(1+e^{-z})^{-2} (e^{-z})(-1) \\ &= \frac{1}{(1+e^{-z})^2} e^{-z} = \frac{e^{-z} + 1 - 1}{(1+e^{-z})^2} = \cancel{\frac{1+e^{-z}}{(1+e^{-z})^2}} - \frac{1}{(1+e^{-z})^2} \\ &= \frac{1}{1+e^{-z}} - \left(\frac{1}{1+e^{-z}}\right)\left(\frac{1}{1+e^{-z}}\right) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right) \\ &= G(z) (1 - G(z)) \end{aligned}$$

* $G(z) = \frac{1}{1+e^{-z}}$ $\xrightarrow{\text{derivative}}$ $\frac{d}{dz} G(z) = G(z) (1 - G(z))$ *

* The process and our goal is the same as Regression Except

in 1) Activation Function and 2) Loss Function. In classification

we use kind of Loss Function called Log Loss:

■ Prediction Function:

$$\hat{y} = G(\omega_1 x_1 + \omega_2 x_2 + b)$$

■ Main Goal:

Find ω_1, ω_2, b that give \hat{y} with the least error.

■ Loss Function (Log Loss):

$$L(y, \hat{y}) = -y \ln(\hat{y}) - (1-y) \ln(1-\hat{y})$$

* $L(y, \hat{y})$ is a large number if y and \hat{y} are far away from each other and small number if they are close to each other.

■ Finding optimal values for ω_1, ω_2, b :

* Using Gradient Descent *

$$\omega_1 \rightarrow \omega_1 - \alpha \frac{\partial L}{\partial \omega_1}$$

$$\omega_1 \rightarrow \omega_1 - \alpha (-x_1(y - \hat{y}))$$

$$\omega_2 \rightarrow \omega_2 - \alpha \frac{\partial L}{\partial \omega_2}$$

$$\omega_2 \rightarrow \omega_2 - \alpha (-x_2(y - \hat{y}))$$

$$b \rightarrow b - \alpha \frac{\partial L}{\partial b}$$

$$b \rightarrow b - \alpha (- (y - \hat{y}))$$

■ calculating partial derivatives:

→ Prediction Function:

→ Partial derivatives:

$$\hat{y} = G(\omega_1 x_1 + \omega_2 x_2 + b)$$

$$\frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \omega_1} = -(y - \hat{y}) x_1$$

→ Log Loss Function:

$$\frac{\partial L}{\partial \omega_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \omega_2} = -(y - \hat{y}) x_2$$

$$L(y, \hat{y}) = -y \ln(\hat{y}) - (1-y) \ln(1-\hat{y})$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} = -(y - \hat{y})$$

* Remember: $f(x) = \ln(x) \rightarrow f'(x) = \frac{1}{x}$ / $G(z) \rightarrow G'(z) = G(z)(1 - G(z))$

→ calculating:

$$\frac{\partial L}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}} = \frac{-y + y\hat{y} + \hat{y} - y\hat{y}}{\hat{y}(1-\hat{y})} = \frac{-(y - \hat{y})}{\hat{y}(1-\hat{y})}$$

$$\frac{\partial \hat{y}}{\partial \omega_1} = \hat{y} (1 - \hat{y}) x_1$$

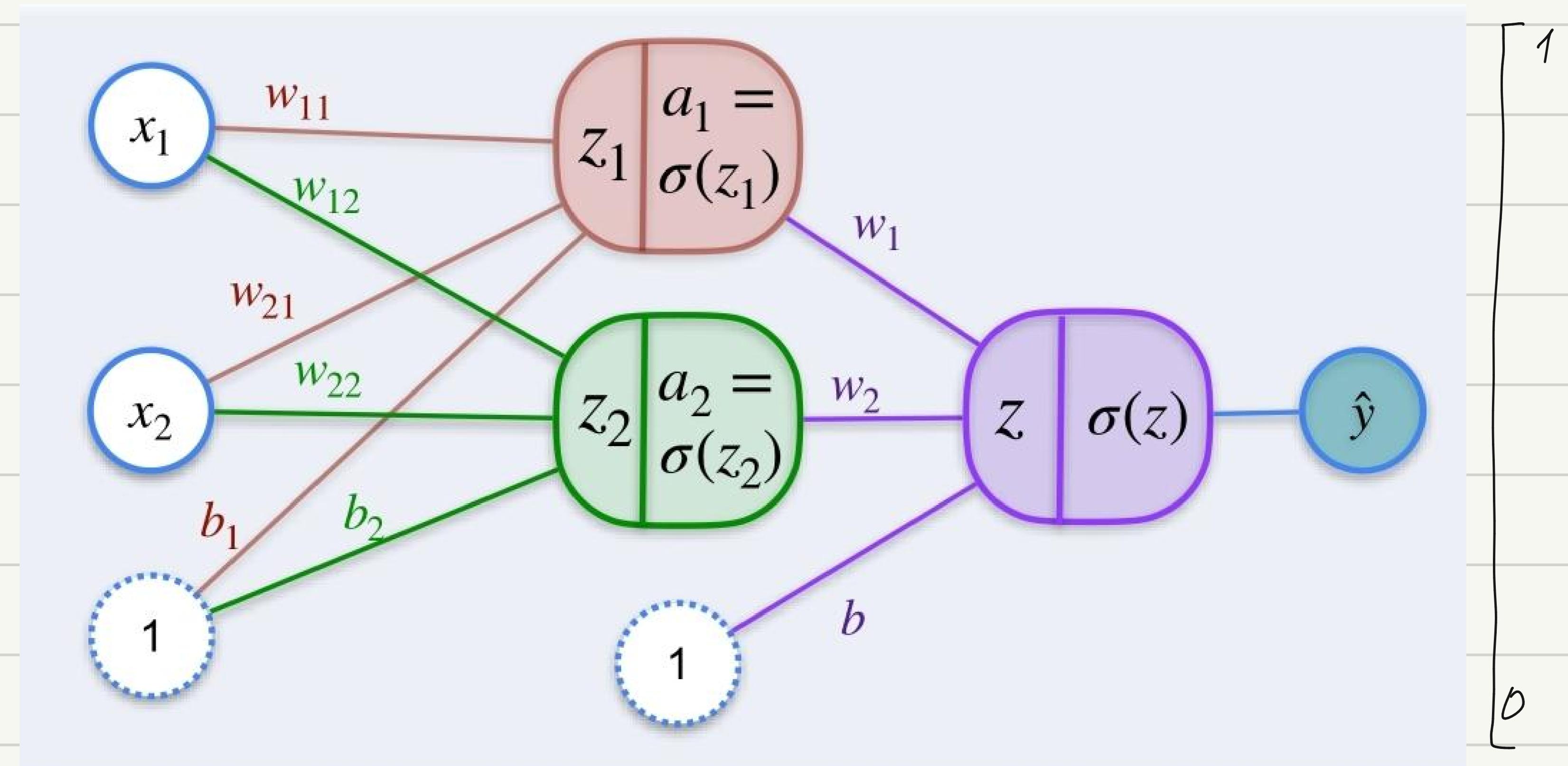
* Derivative of the inside
by ω_1

$$\frac{\partial \hat{y}}{\partial \omega_2} = \hat{y} (1 - \hat{y}) x_2$$

let's multiply
them.

$$\frac{\partial \hat{y}}{\partial b} = \hat{y} (1 - \hat{y})$$

③ 2, 2, 1 Neural Network



L(y, ŷ)

Neural network of depth 2 :

- One input layer
 - One hidden layer
 - One output layer

$$a_1 = g(z_1) \quad , \quad z_1 = x_1 \omega_{11} + x_2 \omega_{21} + b_1$$

$$a_2 = 6(z_2) \quad , \quad z_2 = x_1 \omega_{12} + x_2 \omega_{22} + b_2$$

$$\hat{y} = \theta(z) \quad , \quad z = a_1 w_1 + a_2 w_2 + b$$

$$\text{Log Loss} : L(y, \hat{y}) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

Goal : Adjust each of the weights and biases \rightarrow To reduce loss

How each one of the weights effects the Loss ?

for $w_{11}, w_{12}, w_{21}, w_{22}$ ← for w_1, w_2 ← function

for b_1, b_2 ← for b ←

? $\frac{\partial L}{\partial w_{ij}}$ / $\frac{\partial L}{\partial w_i}$

? $\frac{\partial L}{\partial b_i}$ / $\frac{\partial L}{\partial b}$

.. .. the biases ?

* using chain Rule .

$$L(y, \hat{y}) = -y\log(\hat{y}) - (1-y)\log(1-\hat{y})$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial z_1}{\partial w_{11}} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial L}{\partial \hat{y}}$$

$$\hat{y} = \sigma(z)$$

$$z = a_1 w_1 + a_2 w_2 + b$$

$$a_1 = \sigma(z_1)$$

$$z_1 = x_1 w_{11} + x_2 w_{21} + b_1$$

$$\frac{\partial L}{\partial w_{11}} = x_1 \cdot a_1 (1-a_1) \cdot w_1 \cdot \cancel{\hat{y}(1-\hat{y})} \cdot \frac{-(y - \hat{y})}{\cancel{\hat{y}(1-\hat{y})}}$$

$$= -x_1 w_1 a_1 (1-a_1) (y - \hat{y})$$

Perform gradient descent with

$$w_{11} \rightarrow w_{11} - \alpha \cdot x_1 w_1 a_1 (1-a_1) (y - \hat{y})$$

to find optimal value of w_{11} that gives the least error

$$L(y, \hat{y}) = -y\log(\hat{y}) - (1-y)\log(1-\hat{y})$$

$$\frac{\partial L}{\partial w_{21}} = \frac{\partial z_1}{\partial w_{21}} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial L}{\partial \hat{y}}$$

$$\hat{y} = \sigma(z)$$

$$z = a_1 w_1 + a_2 w_2 + b$$

$$a_1 = \sigma(z_1)$$

$$z_1 = x_1 w_{11} + x_2 w_{21} + b_1$$

$$\frac{\partial L}{\partial w_{21}} = x_2 \cdot a_1 (1-a_1) \cdot w_1 \cdot \cancel{\hat{y}(1-\hat{y})} \cdot \frac{-(y - \hat{y})}{\cancel{\hat{y}(1-\hat{y})}}$$

$$= -x_2 w_1 a_1 (1-a_1) (y - \hat{y})$$

Perform gradient descent with

$$w_{21} \rightarrow w_{21} - \alpha \cdot x_2 w_1 a_1 (1-a_1) (y - \hat{y})$$

to find optimal value of w_{21} that gives the least error

$$L(y, \hat{y}) = -y\log(\hat{y}) - (1-y)\log(1-\hat{y})$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial z_1}{\partial b_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial L}{\partial \hat{y}}$$

$$\hat{y} = \sigma(z)$$

$$z = a_1 w_1 + a_2 w_2 + b$$

$$a_1 = \sigma(z_1)$$

$$z_1 = x_1 w_{11} + x_2 w_{21} + b_1$$

$$\frac{\partial L}{\partial b_1} = 1 \cdot a_1 (1-a_1) \cdot w_1 \cdot \cancel{\hat{y}(1-\hat{y})} \cdot \frac{-(y - \hat{y})}{\cancel{\hat{y}(1-\hat{y})}}$$

$$= -w_1 a_1 (1-a_1) (y - \hat{y})$$

Perform gradient descent with

$$b_1 \rightarrow b_1 - \alpha (-w_1 a_1 (1-a_1) (y - \hat{y}))$$

to find optimal value of b_1 that gives the least error

In Summary (layer 1) :

$$\omega_{11} \longrightarrow \omega_{11} + \alpha \chi_1 \omega_1 a_1 (1 - a_1) (y - \hat{y})$$

$$\omega_{12} \longrightarrow \omega_{11} + \alpha \chi_2 \omega_1 a_1 (1 - a_1) (y - \hat{y})$$

$$\omega_{21} \longrightarrow \omega_{21} + \alpha \chi_1 \omega_2 a_2 (1 - a_2) (y - \hat{y})$$

$$\omega_{22} \longrightarrow \omega_{22} + \alpha \chi_2 \omega_2 a_2 (1 - a_2) (y - \hat{y})$$

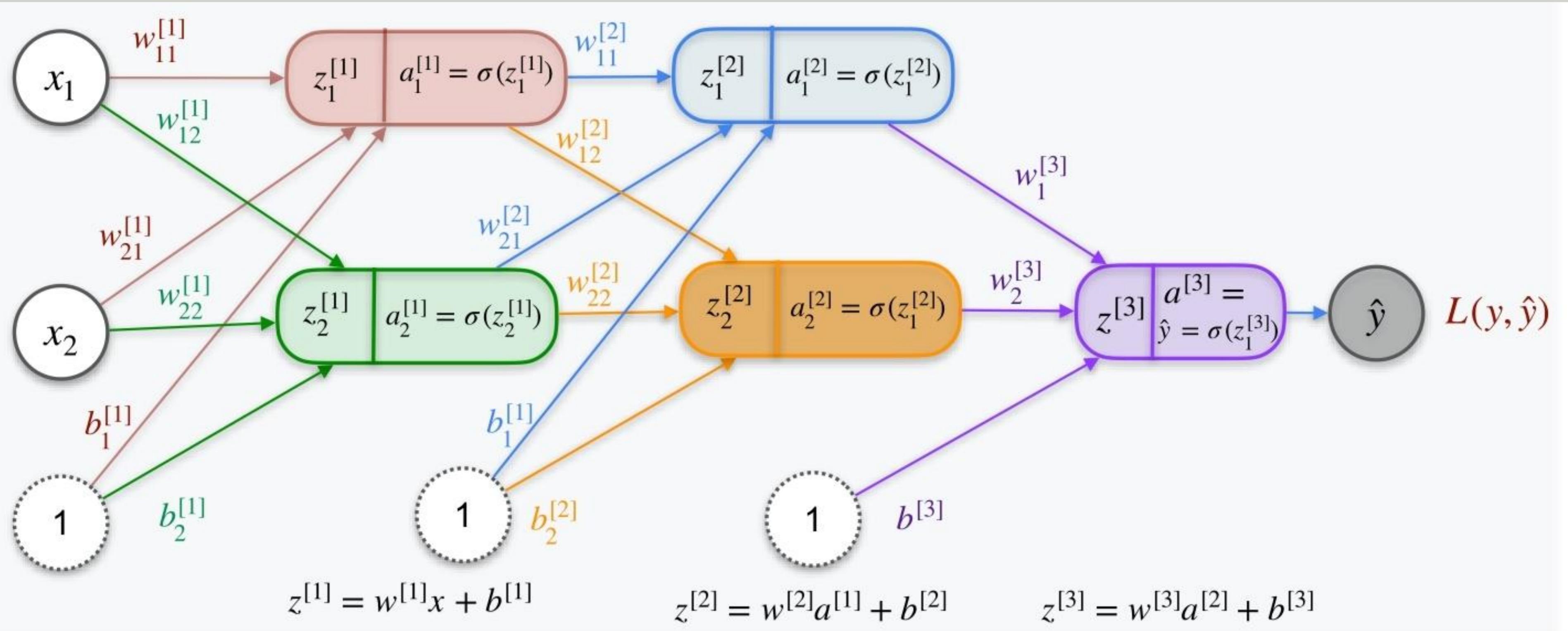
$$b_1 \longrightarrow b_1 + \alpha \omega_1 a_1 (1 - a_1) (y - \hat{y}) \quad , \quad b_2 \longrightarrow b_2 + \alpha \omega_2 a_2 (1 - a_2) (y - \hat{y})$$

Update Second Layer :

$$\omega_1 \longrightarrow \omega_1 + \alpha a_1 (y - \hat{y}) \quad , \quad \omega_2 \longrightarrow \omega_2 + \alpha a_2 (y - \hat{y})$$

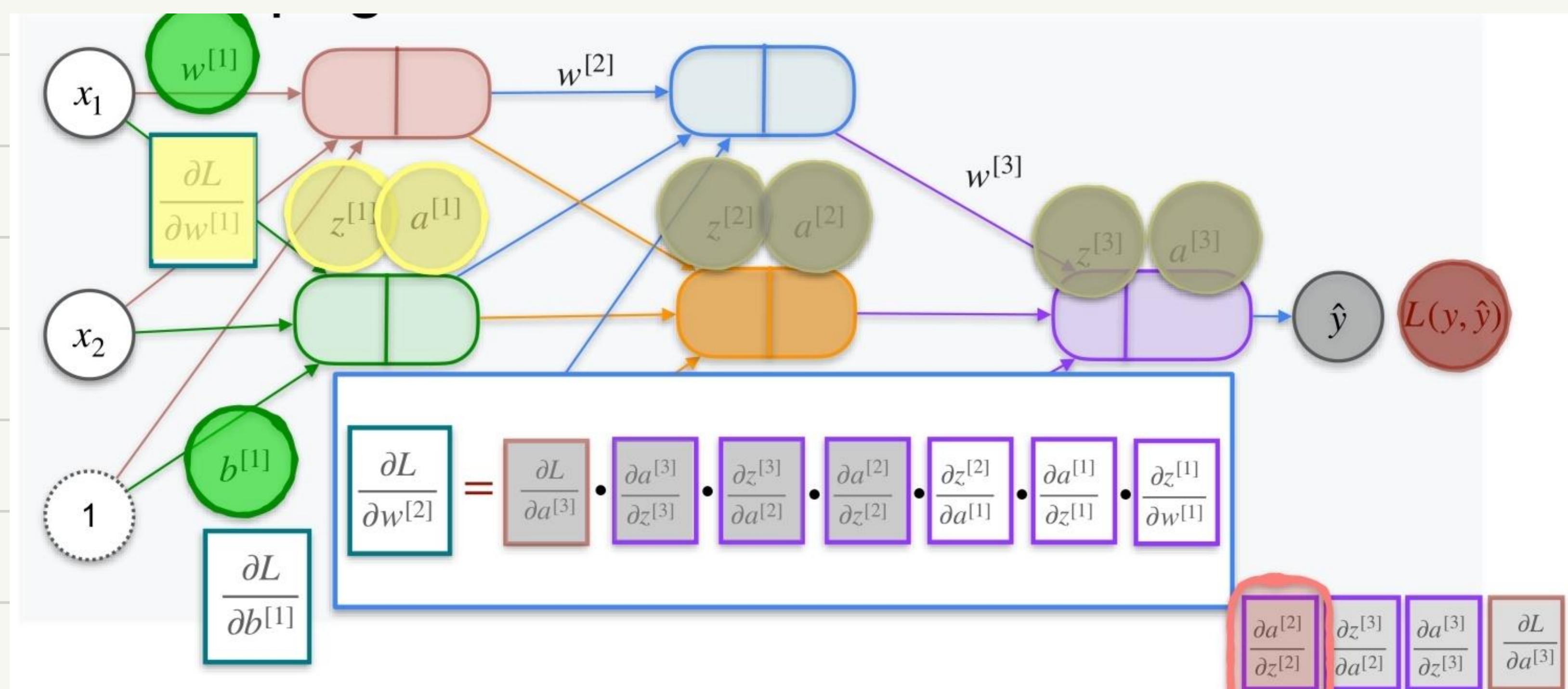
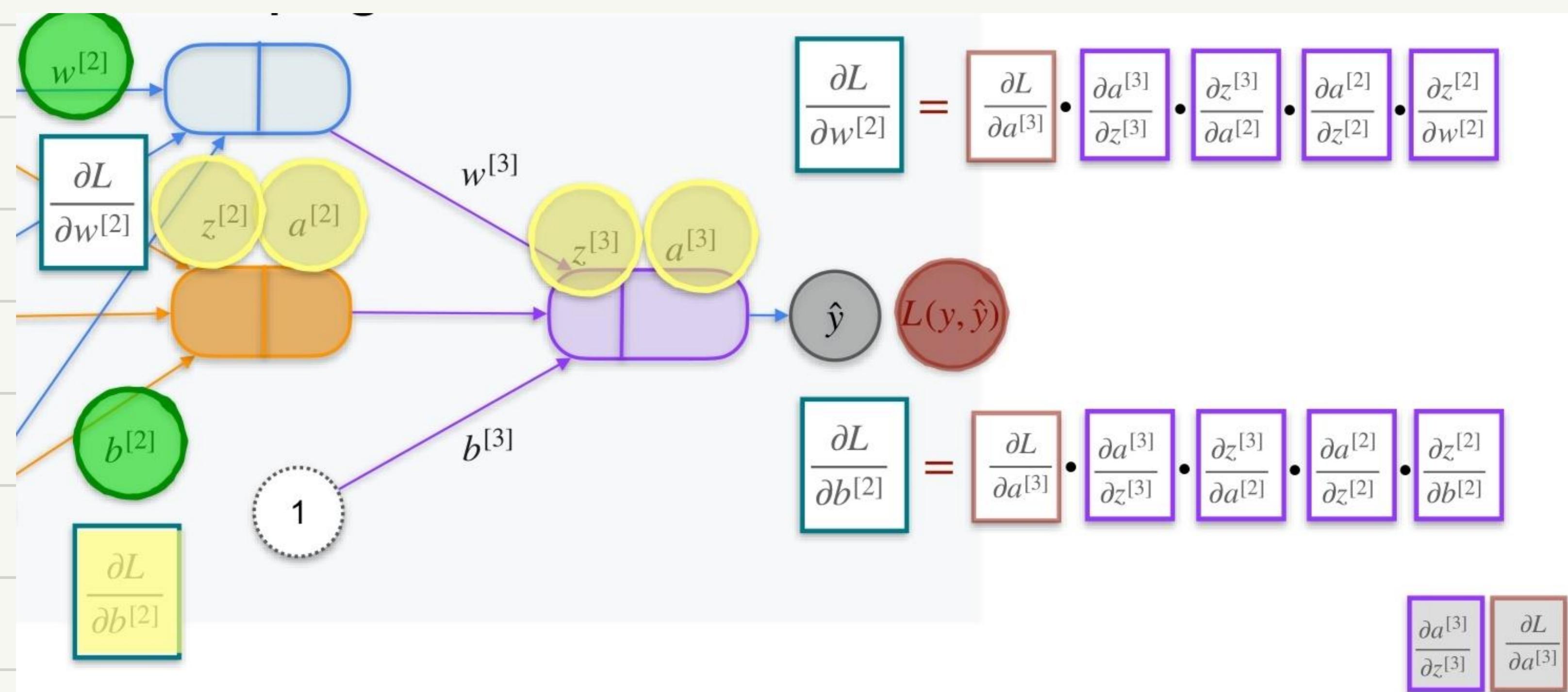
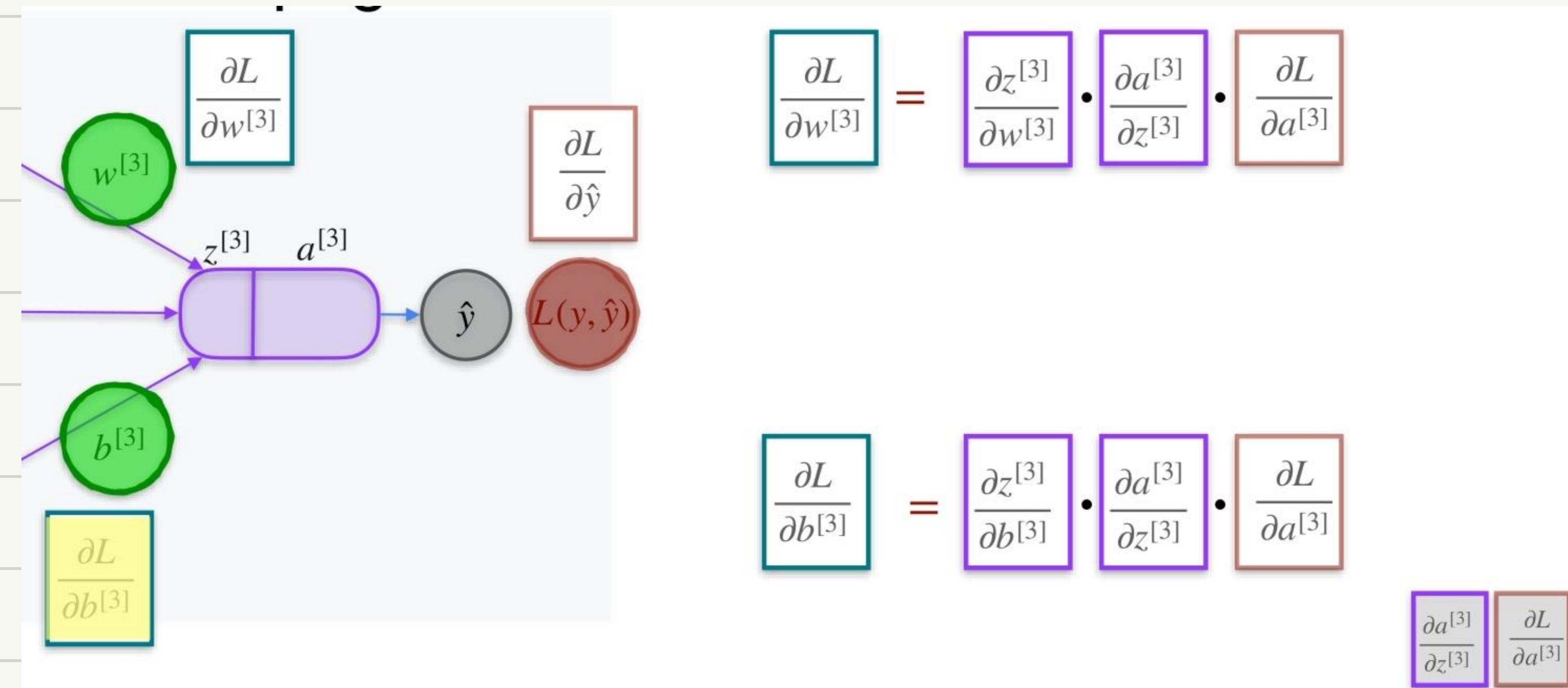
$$b \longrightarrow b + a (y - \hat{y})$$

• Training Model •

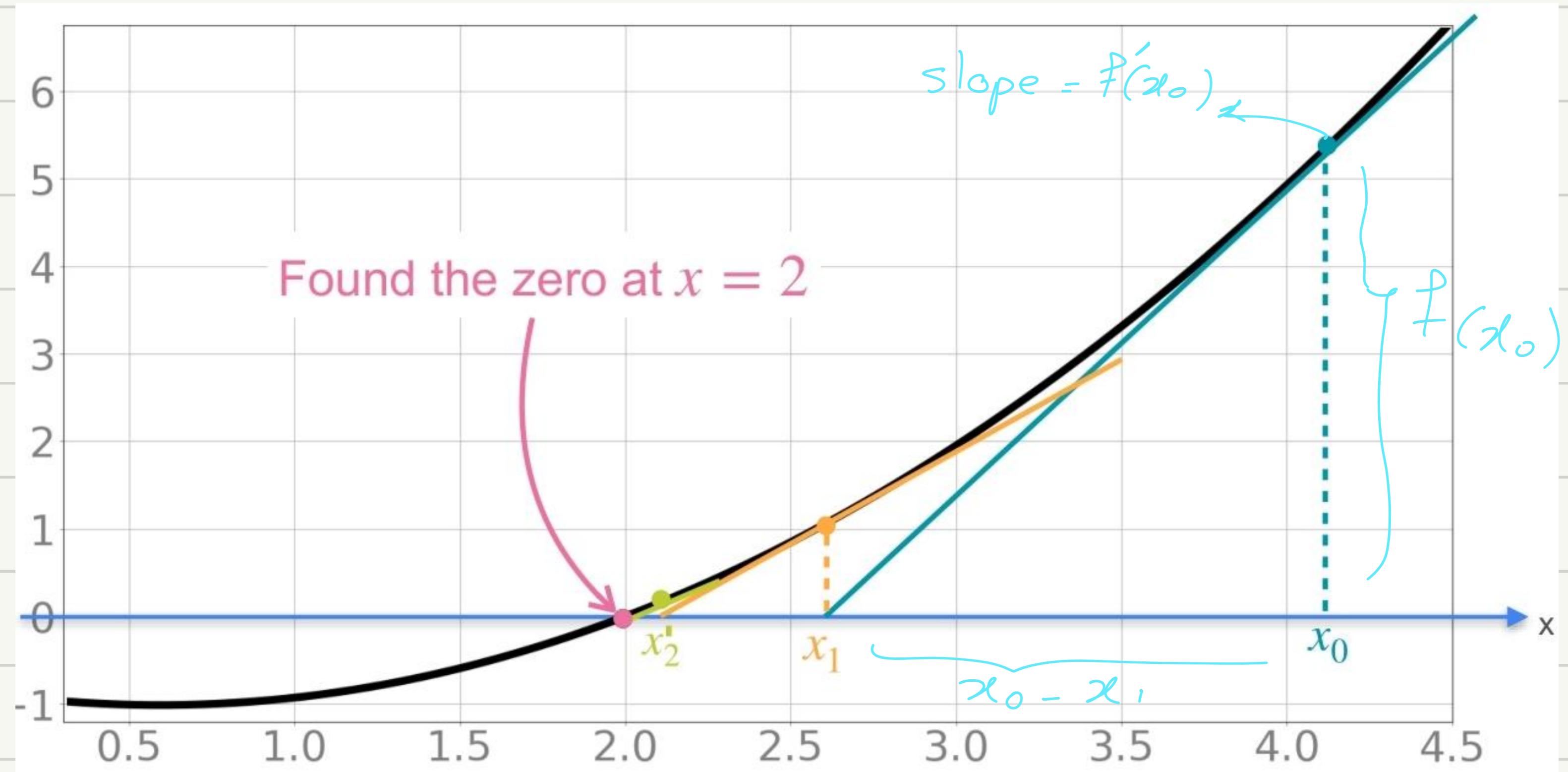


We want to minimize

• **Back Propagation** : is a gradient estimation method commonly used for training neural networks to compute the network parameter updates.



Newton's Method



$$\frac{f(x_0)}{x_0 - x_1} = f'(x_0) \rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

So in general : $x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$

Newton's Method for Optimization

Newton's Method

Goal: Find a zero of $f(x)$

1) Start with some x_0

2) Update :

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

3) Repeat 2 until you find the root

NM for Optimization

Goal: minimize $g(x) \rightarrow$ find zeros of $g'(x)$

$$f(x) \rightarrow g'(x), f'(x) \rightarrow (g'(x))'$$

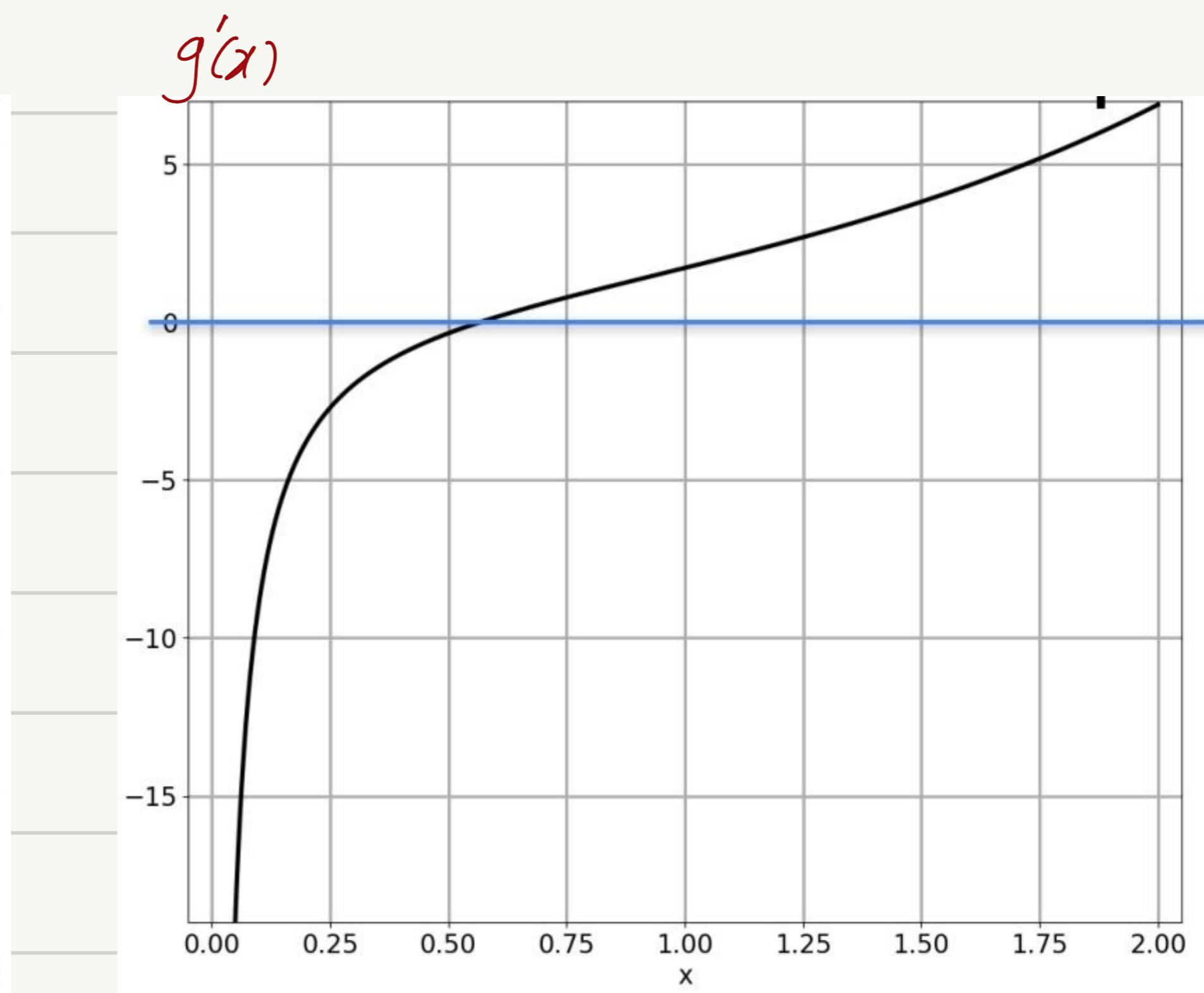
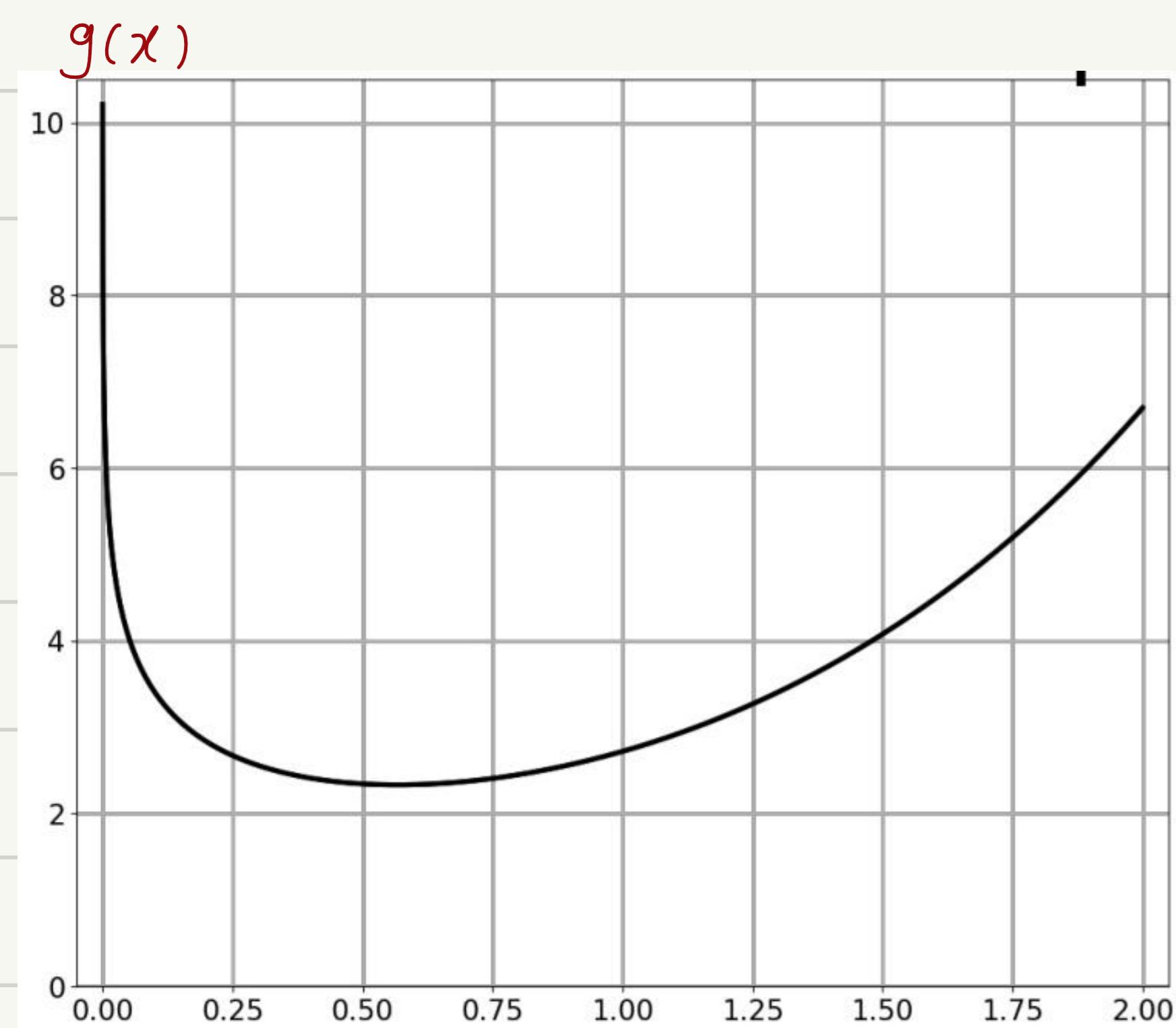
1) Start with some x_0

2) Update :

$$x_{k+1} = x_k - \frac{g'(x_k)}{(g'(x_k))'}$$

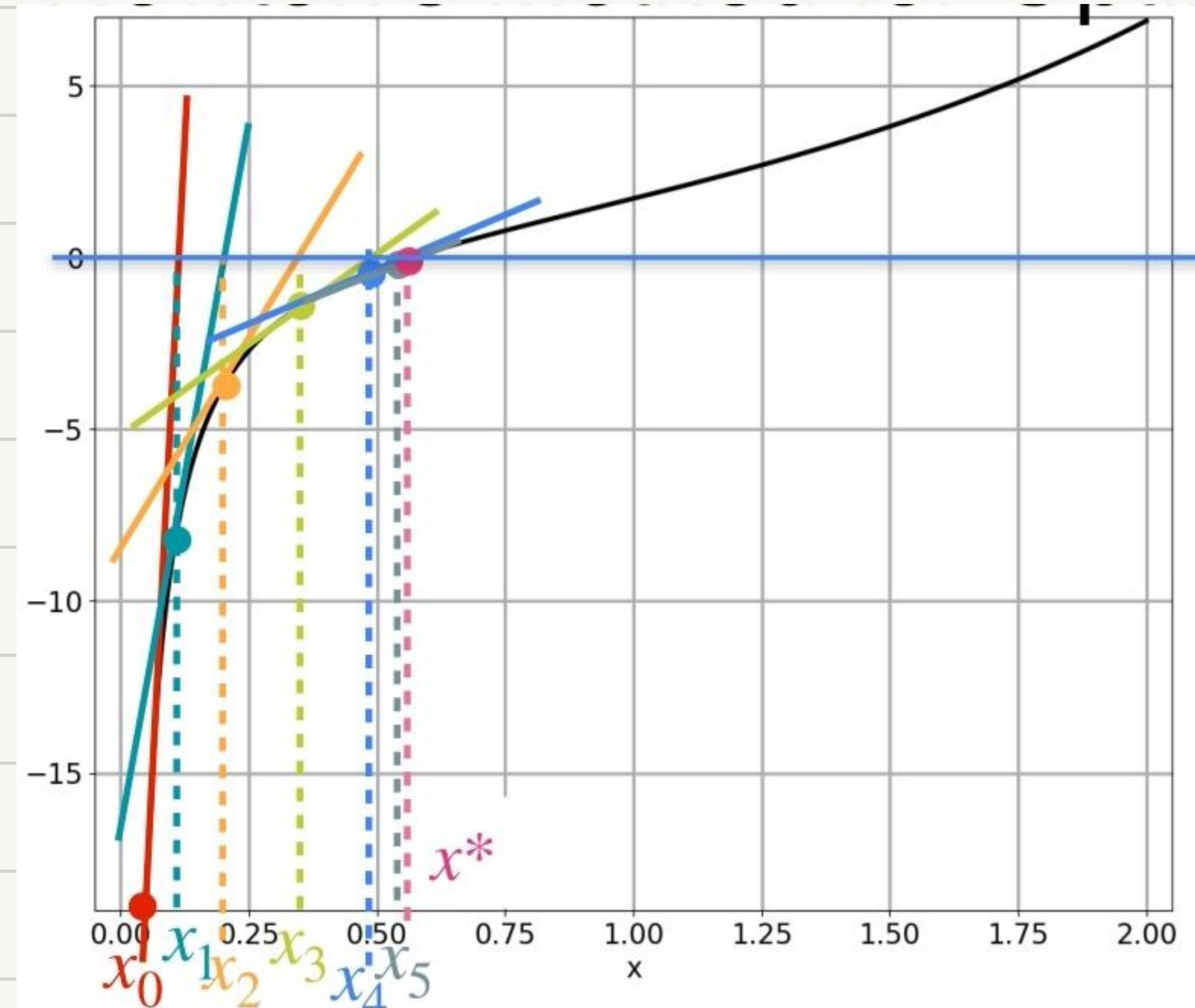
3) Repeat 2 until you find the root

Example : $g(x) = e^x - \log(x)$, $g'(x) = e^x - \frac{1}{x} = f(x)$



Min $\Rightarrow x^* = 0.5671$

Let's Start : $f'(x) = (g'(x))' = e^x + \frac{1}{x^2}$



$$x_0 = 0.05$$

$$x_1 = x_0 - \frac{g'(x_0)}{(g'(x_0))'} = 0.097$$

$$x_2 = x_1 - \frac{g'(x_1)}{(g'(x_1))'} = 0.183$$

$$x_3 = x_2 - \frac{g'(x_2)}{(g'(x_2))'} = 0.320$$

$$x_4 = x_3 - \frac{g'(x_3)}{(g'(x_3))'} = 0.477$$

$$x_5 = x_4 - \frac{g'(x_4)}{(g'(x_4))'} = 0.558$$

$$x_6 = x^* = x_5 - \frac{g'(x_5)}{(g'(x_5))'} = 0.567$$

min point of $g(x)$

④ Second Derivative → used for one variable

Newton's Method: $x_{k+1} = x_k - \frac{g'(x_k)}{(g'(x_k))'}$

Second derivative

$$\text{Leibniz notation: } \frac{d^2 f(x)}{dx^2} = \frac{d}{dx} \left(\frac{df(x)}{dx} \right)$$

$$\text{Lagrange notation: } f''(x)$$

⑤ Understanding Second Derivative

$x \rightarrow \text{Distance}$

$$v \rightarrow \text{Velocity} \rightarrow \frac{dx}{dt}$$

$$a \rightarrow \text{Acceleration} \rightarrow \frac{dv}{dt} = \frac{d^2 x}{dt^2} \quad (\text{second derivative})$$

⑥ Curvature

$$\frac{d^2 x}{dt^2} > 0$$

$$f''(0) > 0$$

Concave up or convex

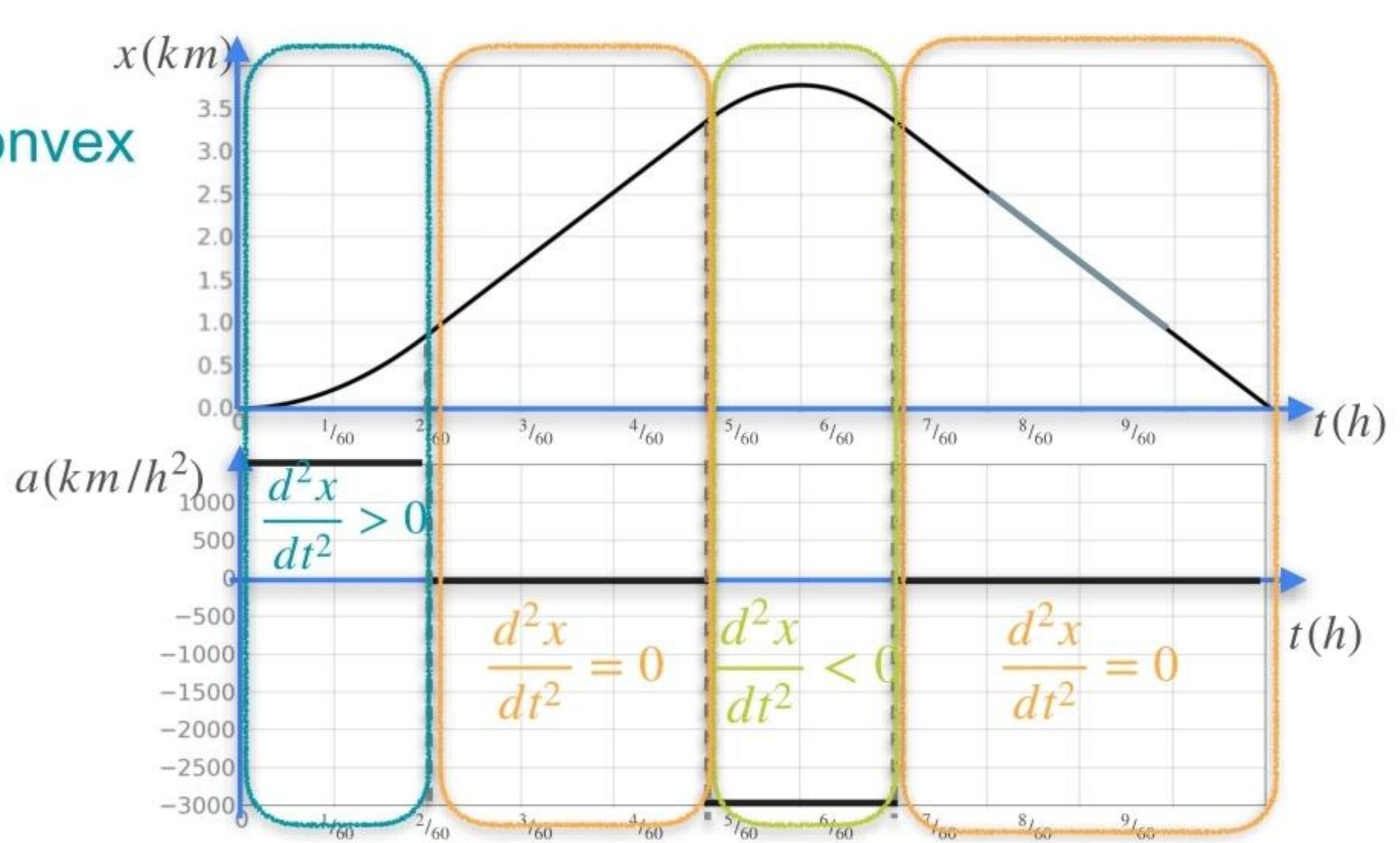
$$\frac{d^2 x}{dt^2} < 0$$

$$f''(0) < 0$$

Concave down

$$\frac{d^2 x}{dt^2} = 0$$

Need more information

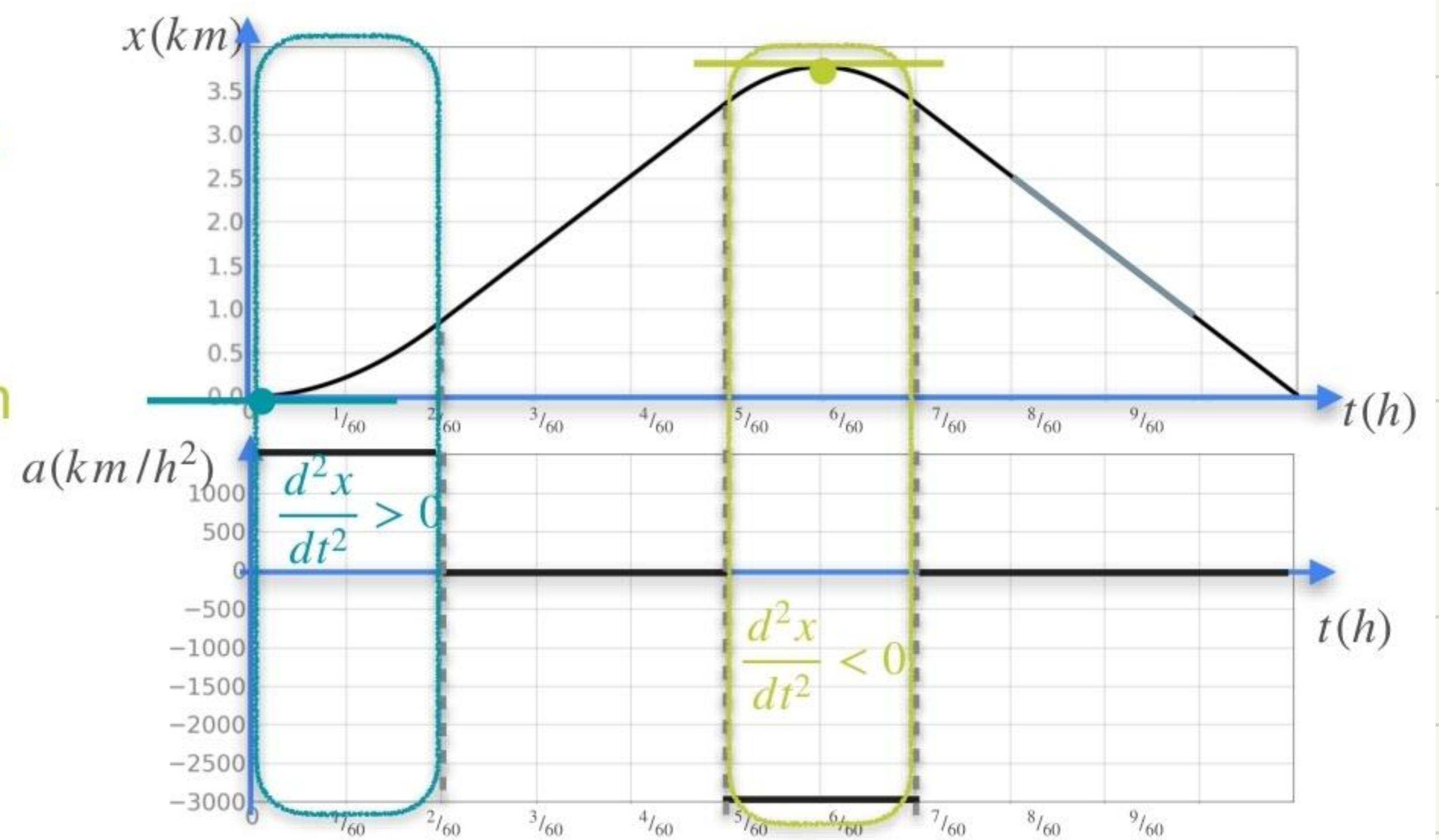


Second Derivative and Optimization

$\frac{d^2x}{dt^2} > 0$ (Local) Minimum

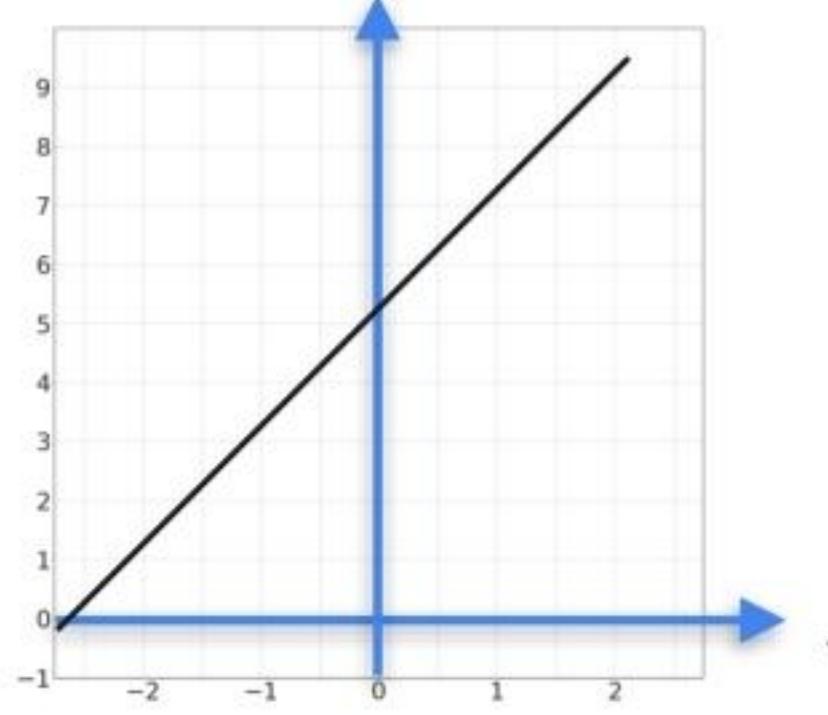
$\frac{d^2x}{dt^2} < 0$ (Local) maximum

$\frac{d^2x}{dt^2} = 0$ Inconclusive



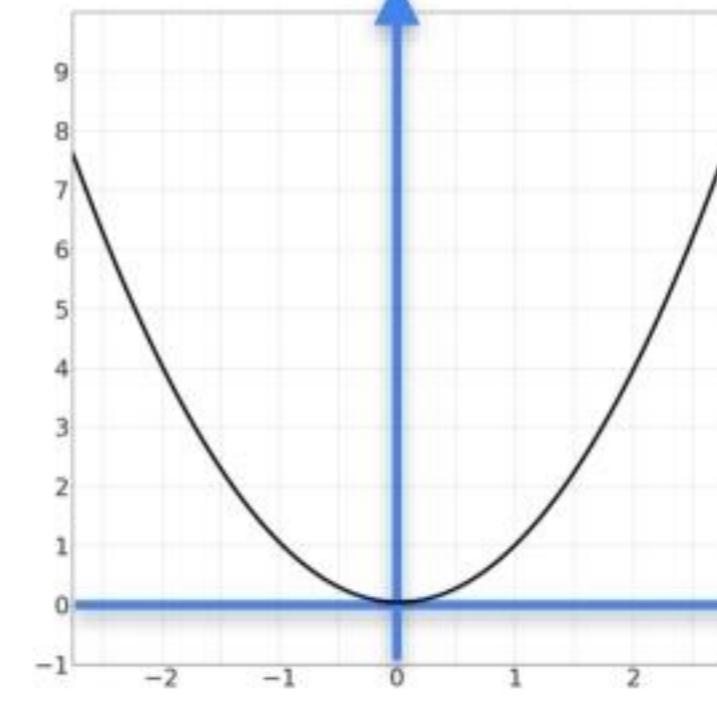
* Curvature *

First derivative



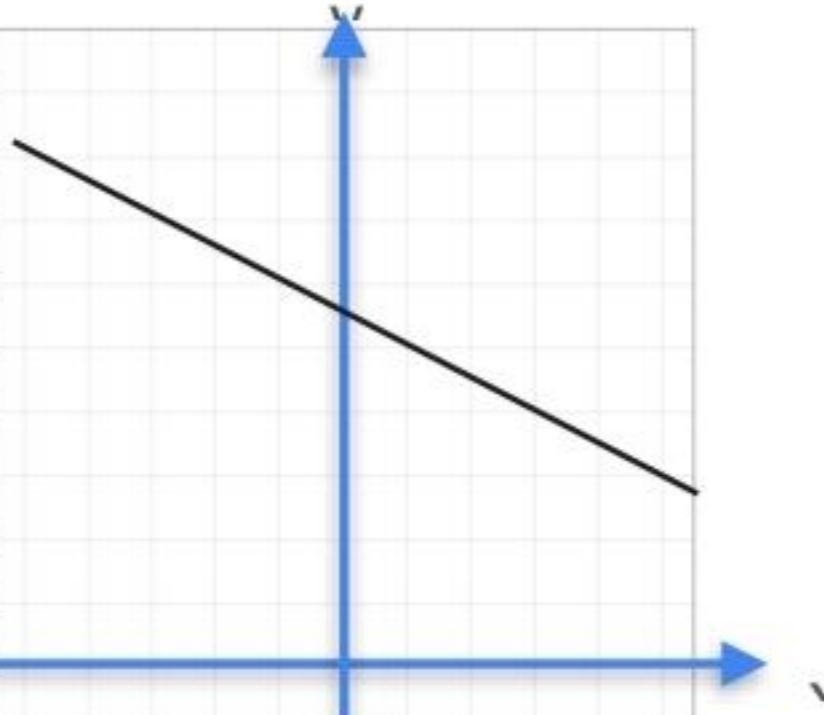
Increasing

Second derivative



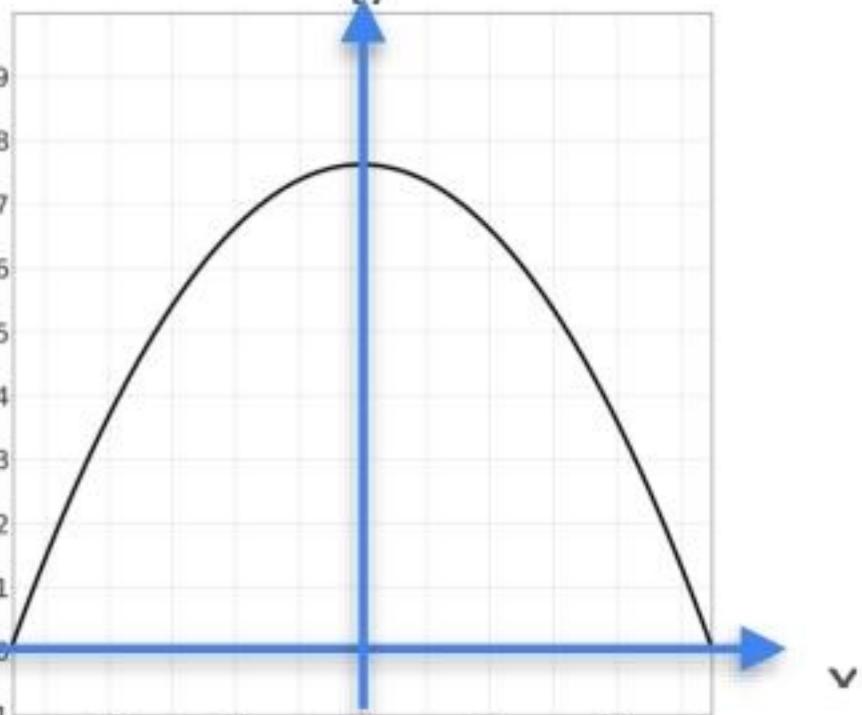
Concave up

$$f'(0) > 0$$



Decreasing

$$f'(0) < 0$$



Concave down

$$f''(0) > 0$$

$$f''(0) < 0$$

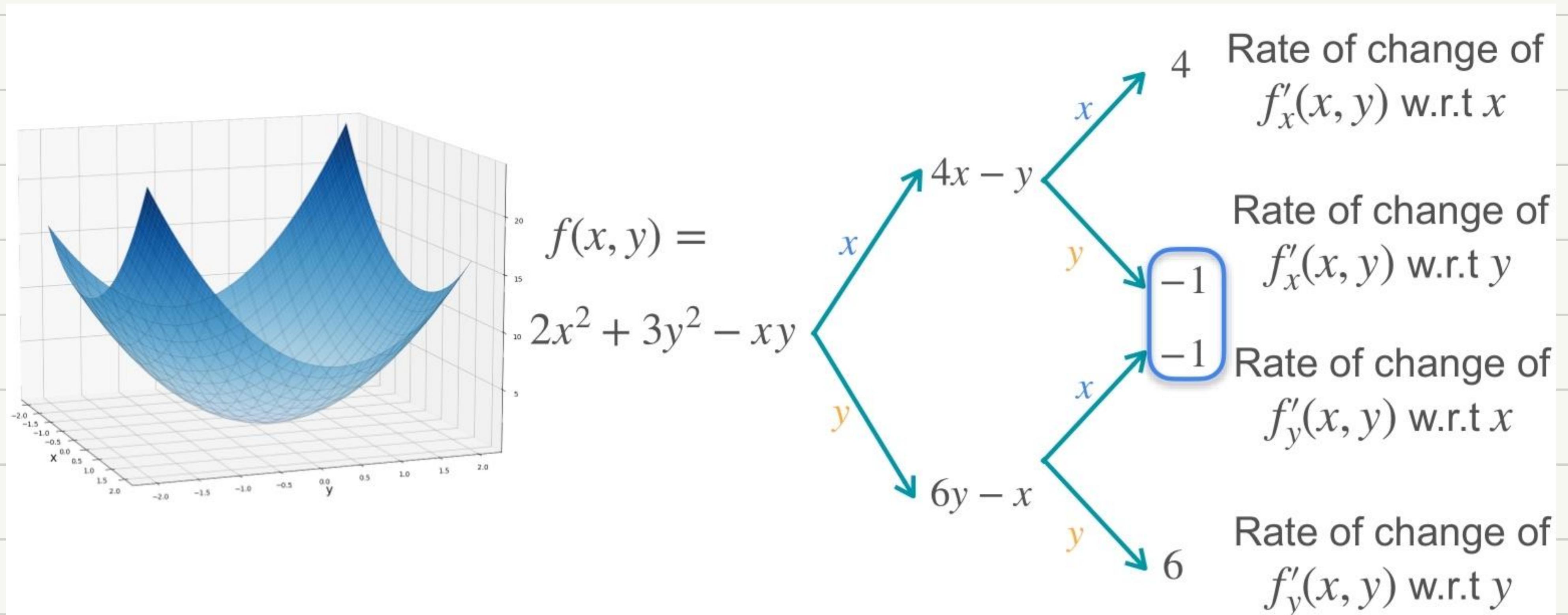
The Hessian (جیسن) ↪ used for multiple variables

Second Derivative

	1 variable	2 variables
Function	$f(x)$	$f(x, y)$
First derivative	$f'(x)$ Rate of change of $f(x)$	$f_x(x, y)$ Rate of change w.r.t x $f_y(x, y)$ Rate of change w.r.t y $\nabla f = \begin{bmatrix} f_x(x, y) \\ f_y(x, y) \end{bmatrix}$
Second derivative	$f''(x)$ Rate of change of the rate of change of $f(x)$??? $H = \begin{bmatrix} f_{xx}(x, y) & f_{xy}(x, y) \\ f_{yx}(x, y) & f_{yy}(x, y) \end{bmatrix}$

Example :

Hessian Matrix



What does this mean ?

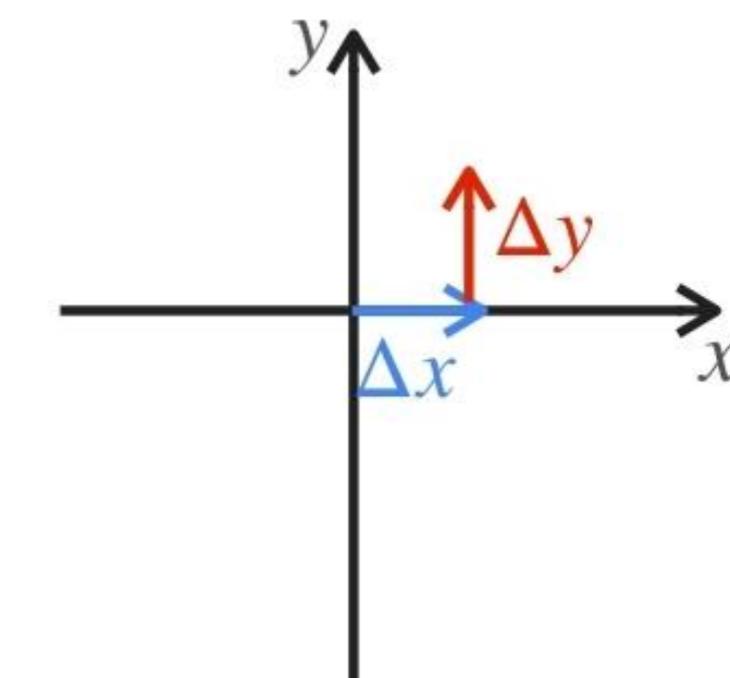
Rate of change of $f_x(x, y)$ w.r.t x
Rate of change of $f_y(x, y)$ w.r.t y

Change in the change in the function w.r.t tiny changes in x and y

Same idea as with one variable!

Rate of change of $f_x(x, y)$ w.r.t y
Rate of change of $f_y(x, y)$ w.r.t x

1. Change in the slope along one coordinate axis w.r.t tiny changes along an orthogonal coordinate axis
2. They are the same!
(In most cases)



* Notation *

Leibniz's notation

Lagrange's notation

Rate of change of
 $f'_x(x, y)$ w.r.t x

$$\frac{\partial^2 f}{\partial x^2}$$

$$f_{xx}(x, y)$$

Rate of change of
 $f'_y(x, y)$ w.r.t y

$$\frac{\partial^2 f}{\partial y^2}$$

$$f_{yy}(x, y)$$

Rate of change of
 $f'_x(x, y)$ w.r.t y

$$\frac{\partial^2 f}{\partial x \partial y}$$

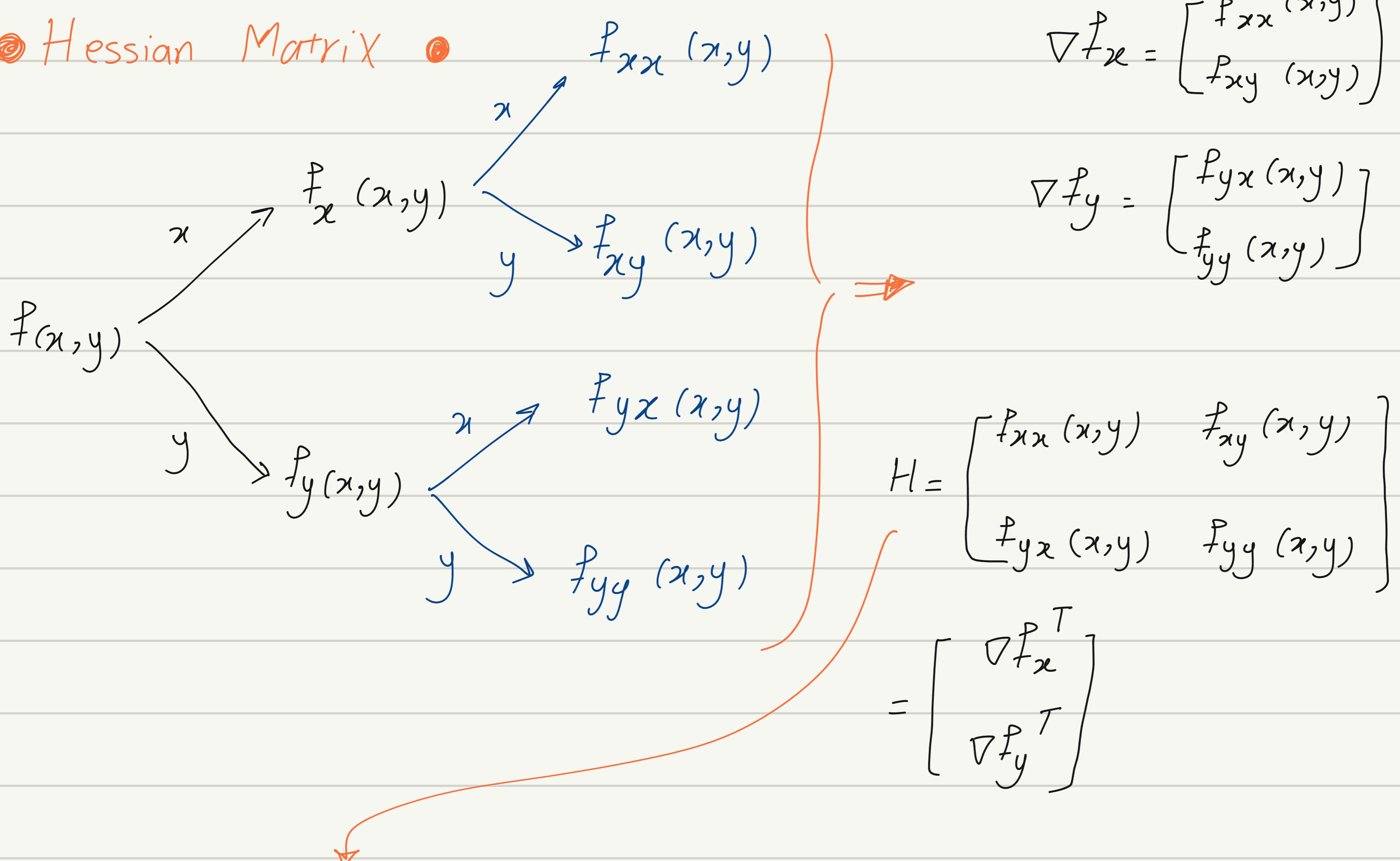
$$f_{xy}(x, y)$$

Rate of change of
 $f'_y(x, y)$ w.r.t x

$$\frac{\partial^2 f}{\partial y \partial x}$$

$$f_{yx}(x, y)$$

④ Hessian Matrix ④

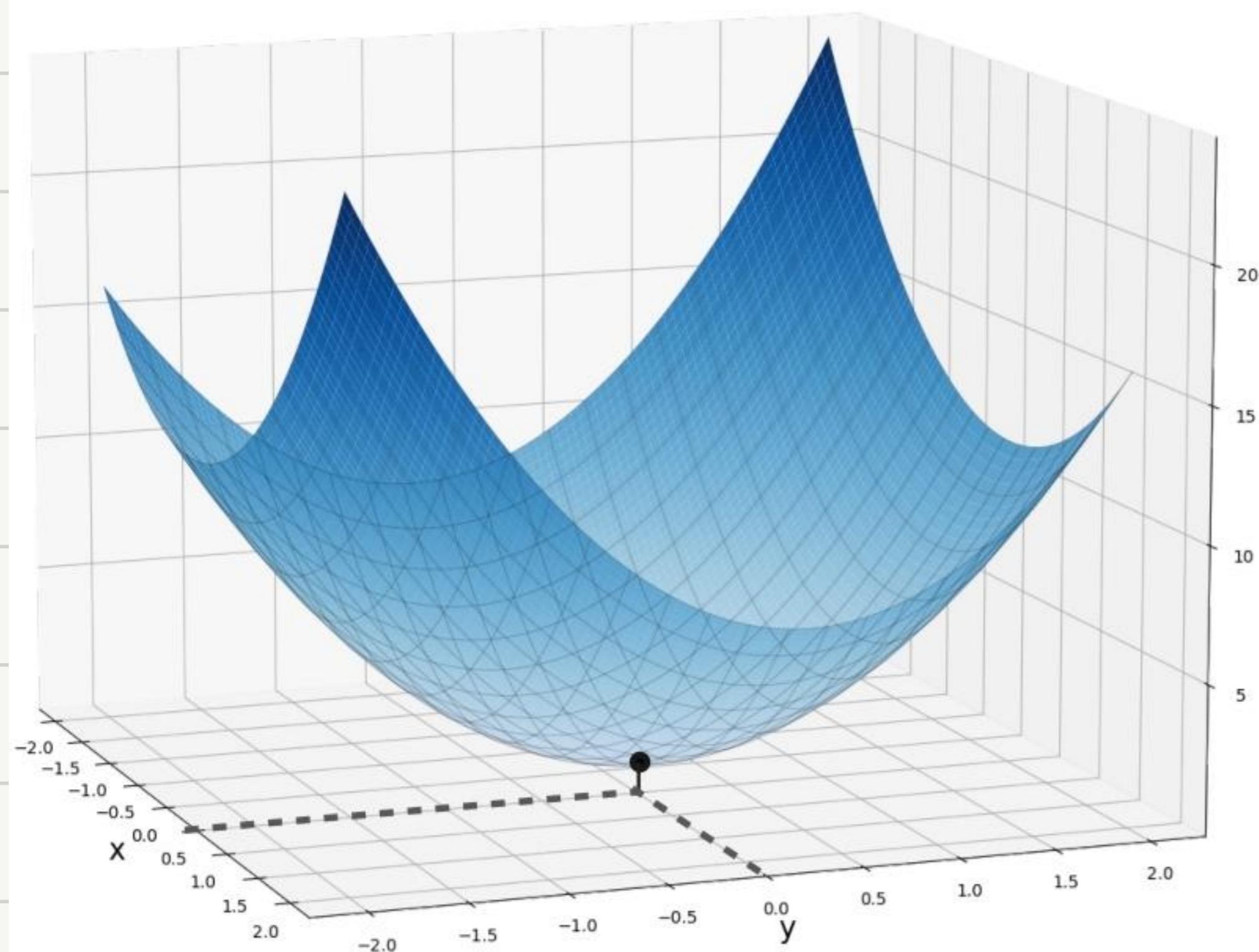


Hessian Matrix : All information about second derivatives

Back to previous page ↪

Hessians and Concavity

1) Concave up : both of eigen values are positive



$$f(x, y) = 2x^2 + 3y^2 - xy$$

$$H(0,0) = \begin{bmatrix} 4 & -1 \\ -1 & 6 \end{bmatrix}$$

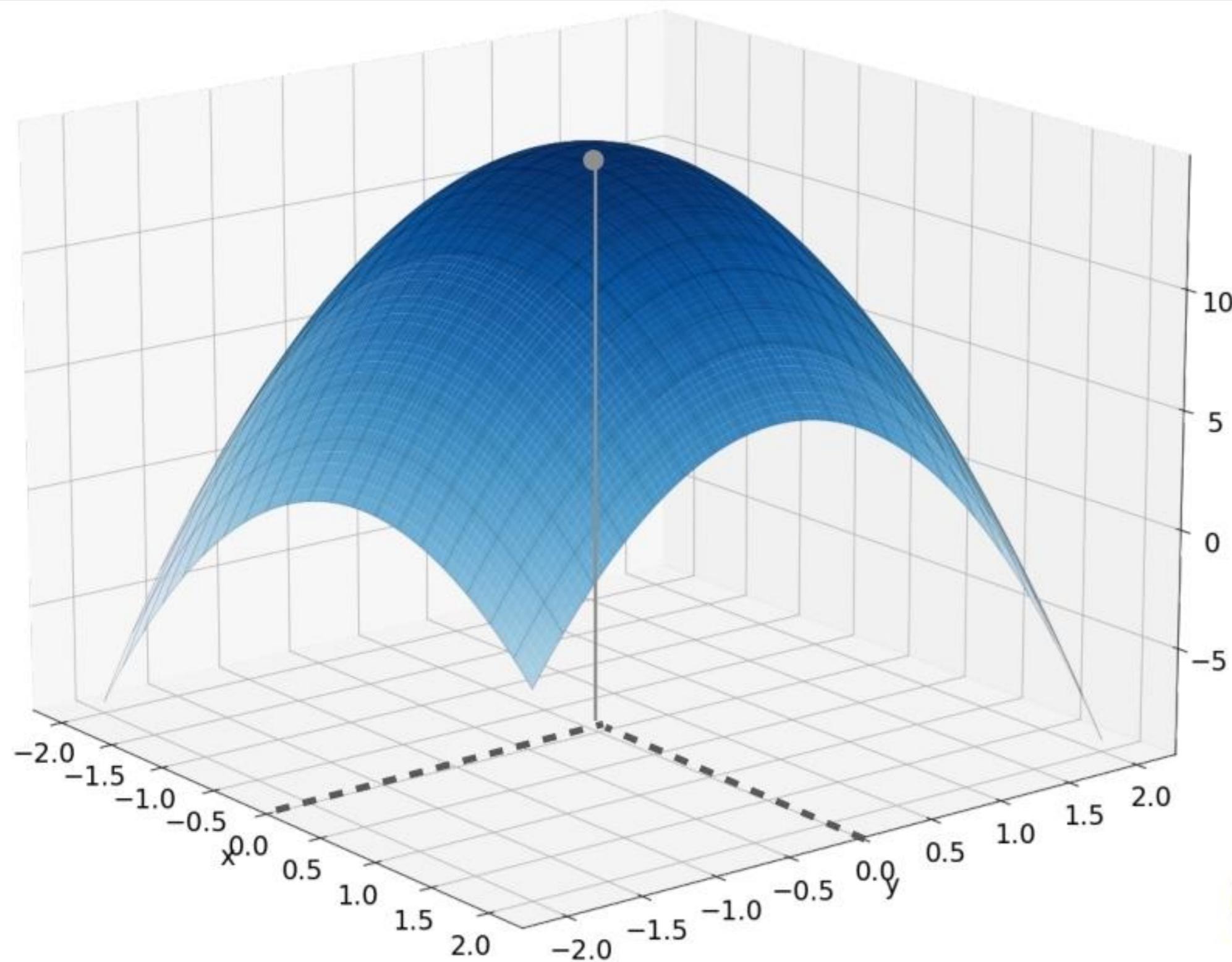
$$\begin{aligned} \det(H(0,0) - \lambda I) &= \det \left(\begin{bmatrix} 4-\lambda & -1 \\ -1 & 6-\lambda \end{bmatrix} \right) \\ &= (4-\lambda)(6-\lambda) - (-1)(-1) \\ &= \lambda^2 - 10\lambda + 23 \end{aligned}$$

$$\boxed{\lambda_1 = 6.41} \\ \boxed{\lambda_2 = 3.59}$$

(0,0) is a minimum!

$$> 0$$

2) Concave Down : both of eigen values are negative



$$f(x, y) = -2x^2 - 3y^2 - xy + 15$$

$$\nabla f(x, y) = \begin{bmatrix} -4x - y \\ -x - 6y \end{bmatrix}$$

$$H(0,0) = \begin{bmatrix} -4 & -1 \\ -1 & -6 \end{bmatrix}$$

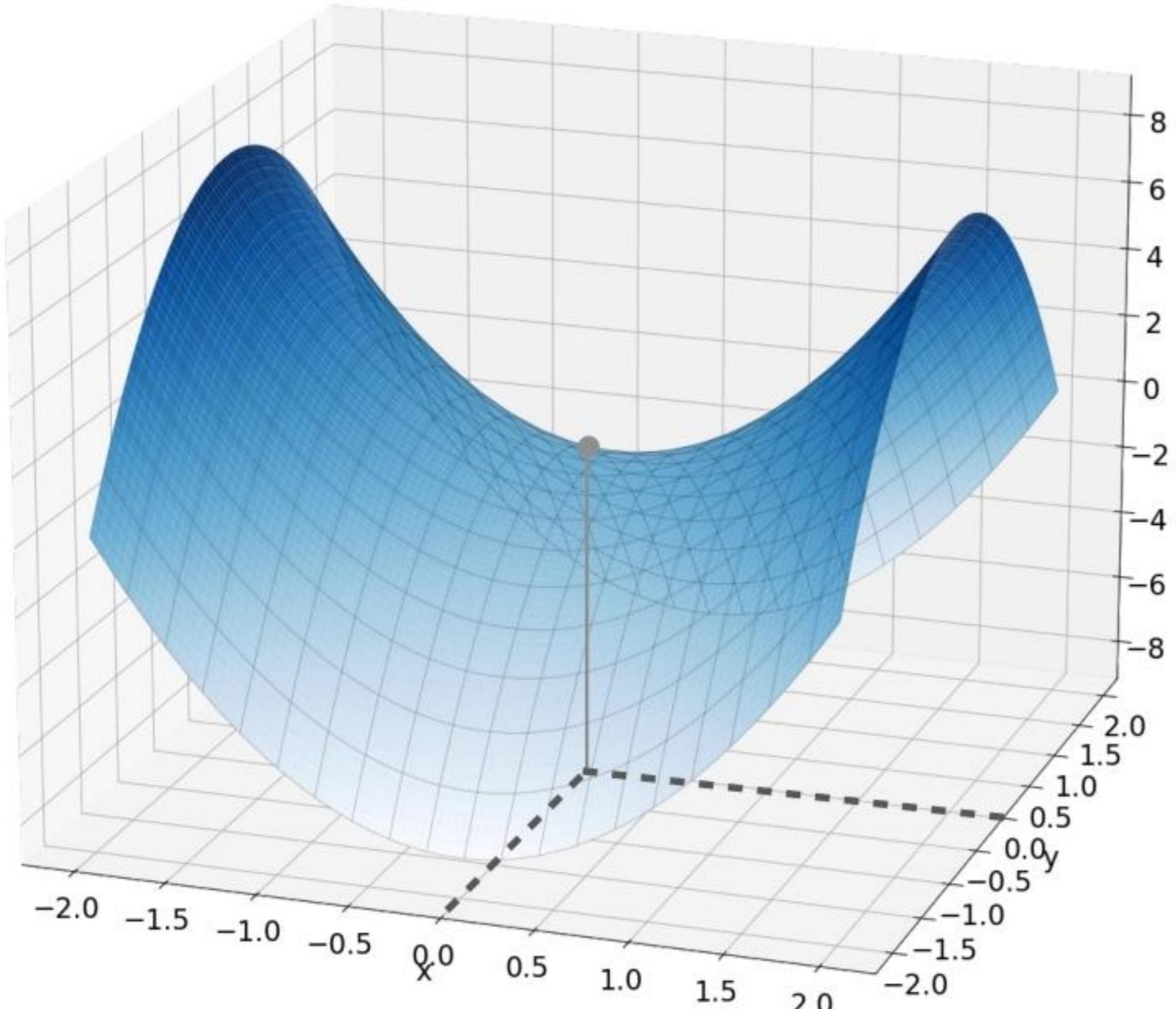
$$\begin{aligned} \det(H(0,0) - \lambda I) &= \\ &= (-4-\lambda)(-6-\lambda) - (-1)(-1) \\ &= \lambda^2 + 10\lambda + 23 \end{aligned}$$

$$\boxed{\lambda_1 = -3.49} \\ \boxed{\lambda_2 = -6.41}$$

(0,0) is a maximum!

$$< 0$$

3) Saddle point



$$f(x, y) = 2x^2 - 2y^2$$

$$\nabla f(x, y) = \begin{bmatrix} 4x \\ -4y \end{bmatrix}$$

$$H(0,0) = \begin{bmatrix} 4 & 0 \\ 0 & -4 \end{bmatrix}$$

$$\det(H(0,0) - \lambda I) =$$

$$(4-\lambda)(-4-\lambda) - 0 \quad \boxed{< 0}$$

$$\boxed{\lambda_1 = -4} \\ \boxed{\lambda_2 = 4}$$

(0,0) is saddle point

$$> 0$$

Summary

	1 variable $f(x)$	2 variables $f(x, y)$	More variables $f(x_1, x_2, \dots, x_n)$
(Local) minima	Happy face $f''(x) > 0$	Upper paraboloid $\lambda_1 > 0 \& \lambda_2 > 0$	All $\lambda_i > 0$
(Local) maxima	Sad face $f''(x) < 0$	Down paraboloid $\lambda_1 < 0 \& \lambda_2 < 0$	All $\lambda_i < 0$
Need more information	$f''(x) = 0$	Saddle point $\lambda_1 > 0 \& \lambda_2 < 0$ $\lambda_1 < 0 \& \lambda_2 > 0$ Or some $\lambda_i = 0$	Some $\lambda_i > 0$ and some $\lambda_j < 0$ OR At least one $\lambda_i = 0$

Newton's Method

1 variable : $x_{k+1} = x_k - \frac{f'(x_k)}{(f'(x_k))'}$

$$x_{k+1} = x_k - \underbrace{f''(x_k)^{-1}}_{\text{pink}} \cdot \underbrace{f'(x_k)}_{\text{green}}$$

2 variables : $\begin{bmatrix} x_{k+1} \\ y_{k+1} \end{bmatrix} = \begin{bmatrix} x_k \\ y_k \end{bmatrix} - \underbrace{H^{-1}(x_k, y_k)}_{\text{pink}} \cdot \underbrace{\nabla f(x_k, y_k)}_{\text{green}}$

So: changing the position is wrong : ~~$\nabla f(x, y) \cdot H^{-1}(x, y)$~~