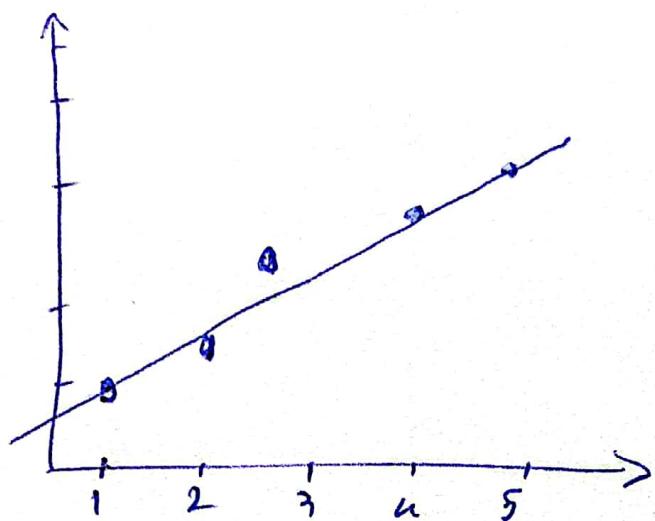


①

# "Calculus"

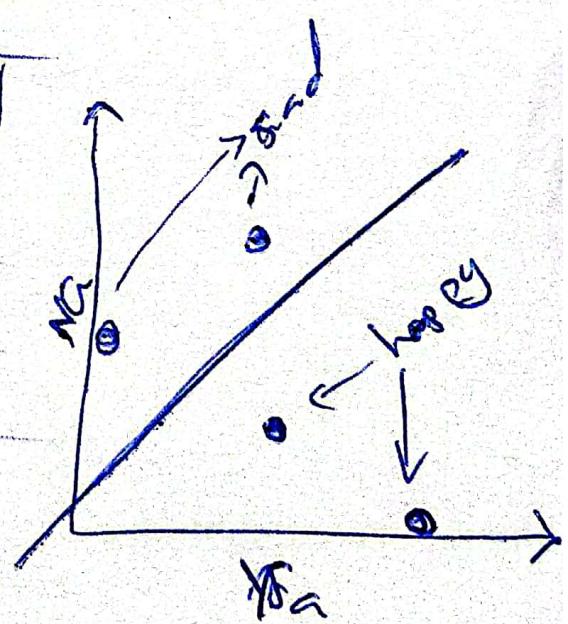
- A linear regression problem:



- A classification problem:

a sentiment analysis model

Sentence	Pa	Na	mood
Ya Ya Ya	3	0	happy
Na na	2	2	sad
Ya na Ya	2	1	happy
Na na m na	1	3	sad



math concepts used in training a model ②

- Gradients
- Derivatives
- Optimization
- Loss and cost functions
- Gradient Descent

models

Linear regression

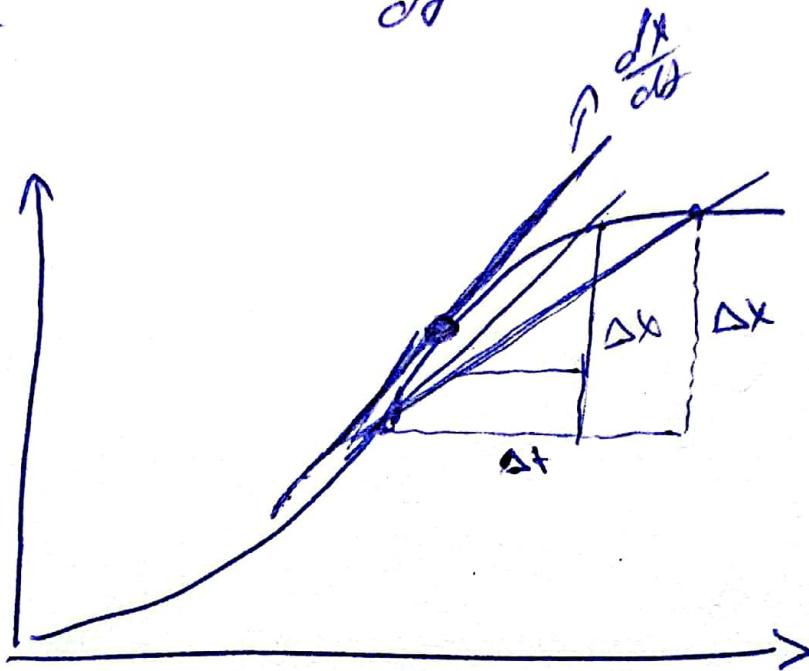
Classification

~~Line~~ "Derivative"

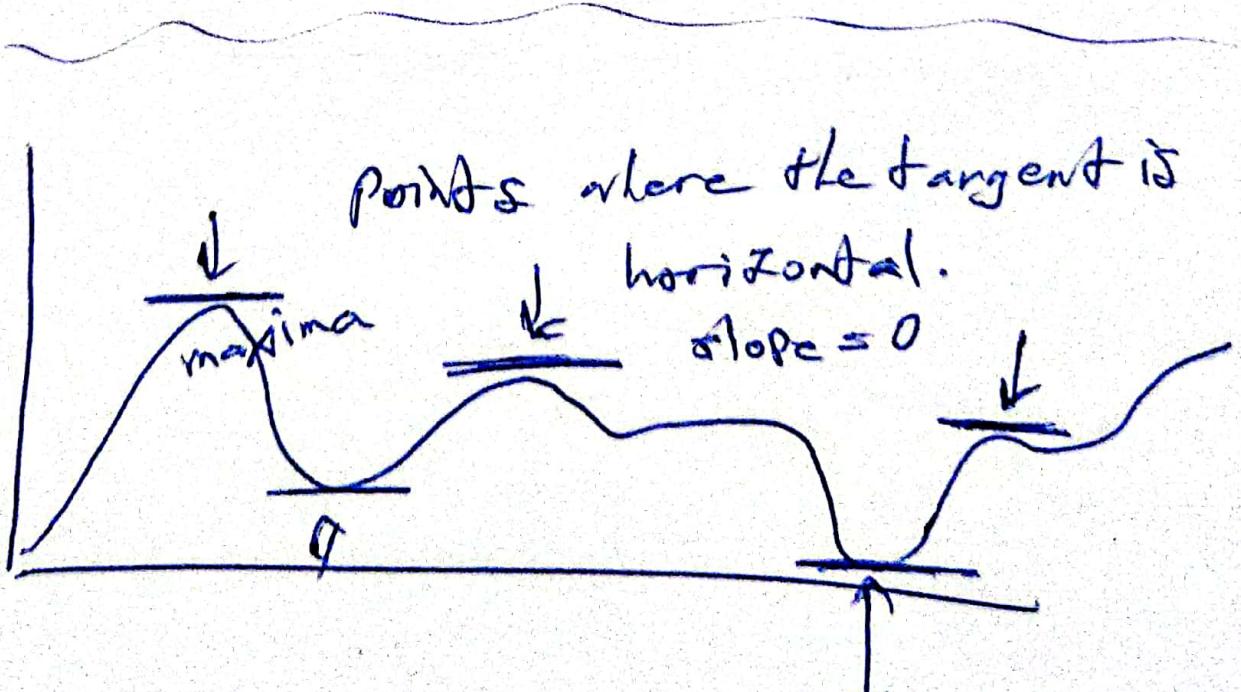
Derivatives: the instantaneous velocity.

$$\text{slope} = \frac{\text{rise}}{\text{run}} \Rightarrow \text{slope} = \frac{\text{change in distance } (\Delta x)}{\text{change in time } (\Delta t)}$$

$$\frac{\Delta x}{\Delta t} \rightarrow \rightarrow \rightarrow \frac{dx}{dt}$$

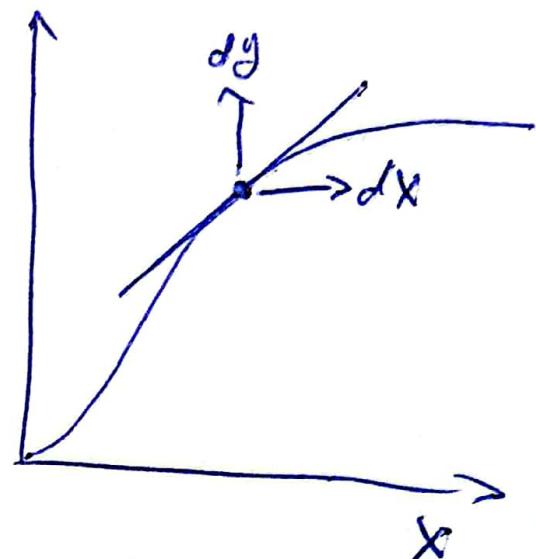


Derivative : the derivative of a function at a point is precisely the slope of the tangent at that particular point.



Q

slope at a point =  $\frac{dy}{dx}$



if:  $y = f(x)$

slope =  $f'(x)$

→ Lagrange's notation.

$$\frac{dy}{dx} = \frac{d}{dx} f(x)$$

↳ Leibniz's notation.

$$y = f(x) = c$$

$$f'(x) = 0$$

$$f(x) = ax + b$$

$$f'(x) = a$$

Quadratics  
 $y = f(x) = x^2$

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

cubics.  
 $y = f(x) = x^3$

$$f'(x) = 3x^2$$

(5)

$$y = f(x) = \frac{1}{x} = x^{-1}$$

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

so:

$$f(x) = x^n$$

$$f'(x) = \frac{d}{dx} f(x) = n x^{n-1}$$

Inverse function

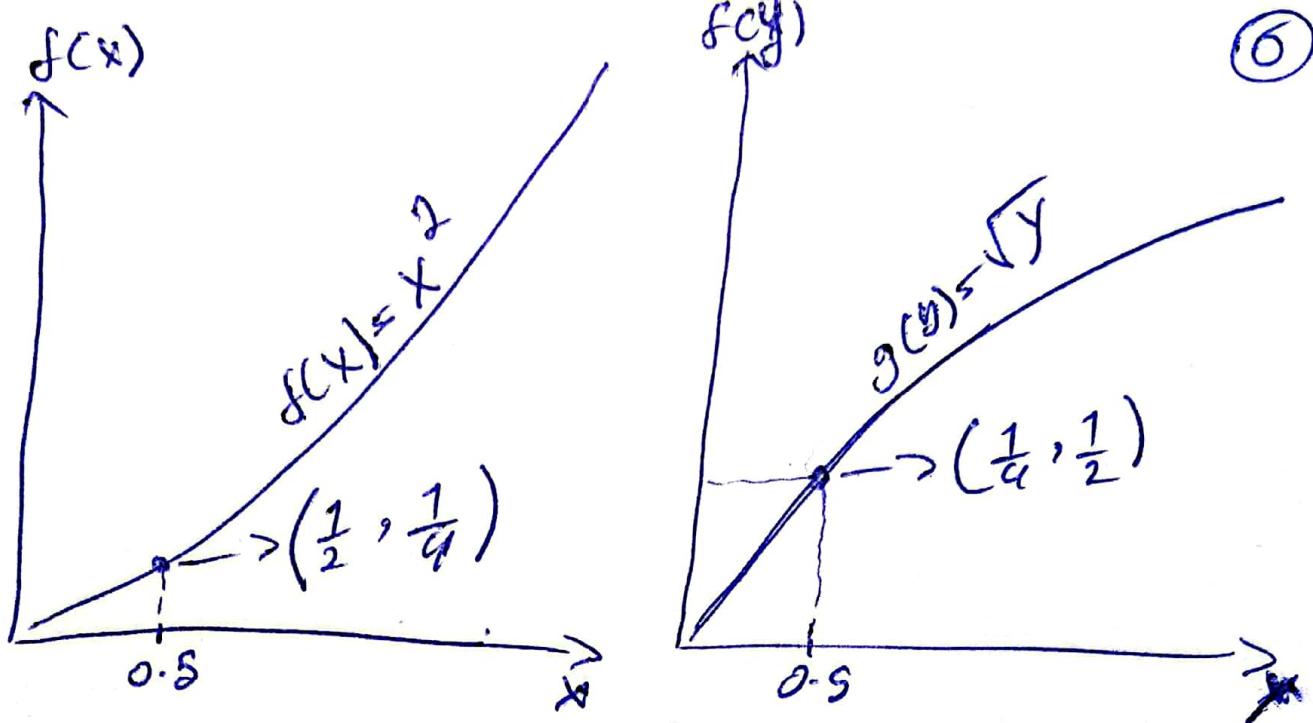
$$x \xrightarrow[f]{\downarrow x^2} x^2 \xrightarrow[g]{\downarrow \sqrt{x}} x \quad \text{for } x > 0$$

if  $g(x)$  and  $f(x)$  are inverses

$$\text{then: } g(x) = f^{-1}(x)$$

$$g(f(x)) = x$$

⑥



$$\text{if } g(y) = \sqrt{y} \quad g'(y) = \frac{1}{2\sqrt{y}}$$

- if  $f$  and  $g$  are inverse functions, then  
the derivative of  $g$  is ~~one~~ 1 over the  
derivative of  $f$ .

$$\text{if } f(x) = \sin(x)$$

$$\Rightarrow f'(x) = \cos(x)$$

$$\text{if } f(x) = \cos(x)$$

$$\Rightarrow f'(x) = -\sin(x)$$

"Euler's number"

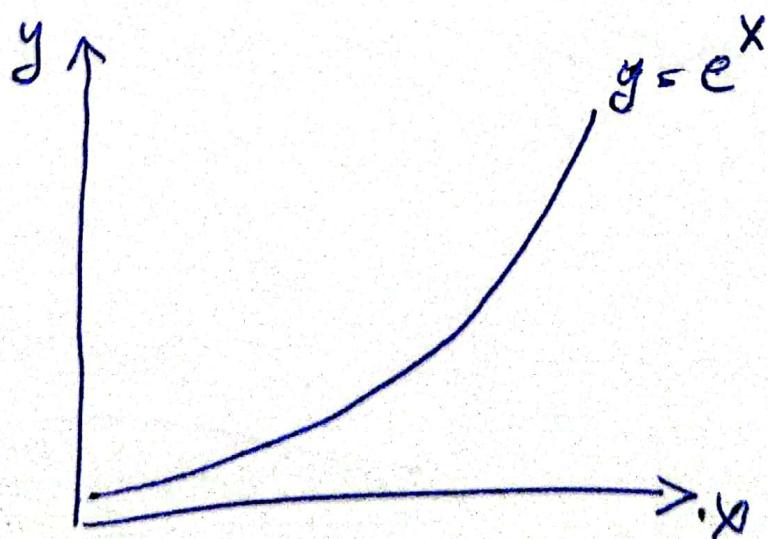
2.71828182...

$$\checkmark \frac{n}{\left(1 + \frac{1}{n}\right)^n} = \frac{1}{2} = \frac{10}{2592} = \frac{100}{2712} \dots \frac{\infty}{e}$$

if  $f(x) = e^x$

then:  $f'(x) = e^x$

"exponential function"



8

logarithm

$$e^? = 3 \Rightarrow \log_e(3)$$

$$e^x = X \Rightarrow \log(X) \quad \text{natural logarithm}$$

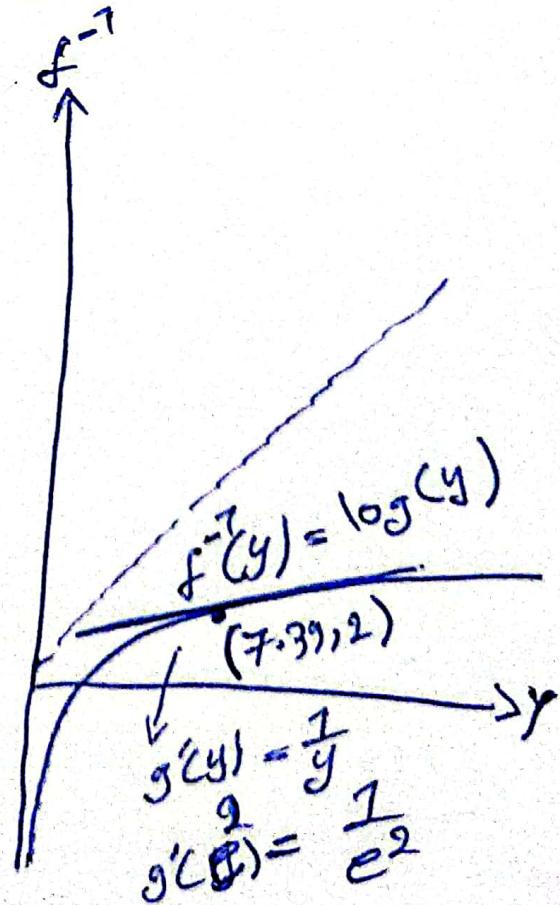
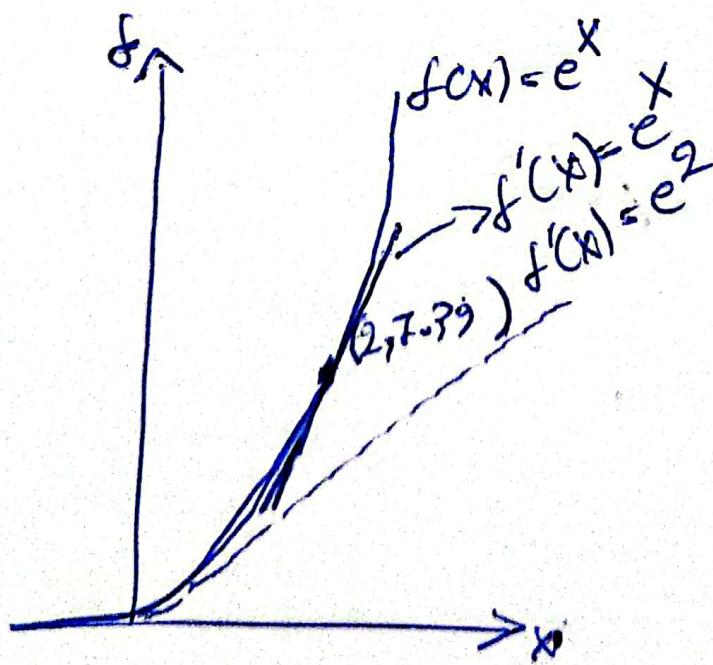
↓

$$\ln(X)$$

$$f(x) = e^x \quad f^{-1}(y) = \log(y)$$

$$e^{\log(x)} = x$$

$$\log(e^y) = y$$



$$\frac{d}{dy} f^{-1}(g) = \frac{1}{f'(f^{-1}(g))}$$

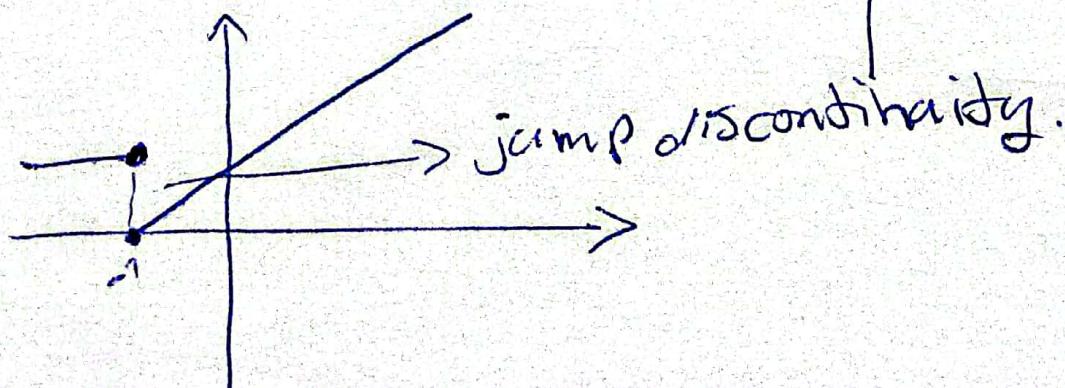
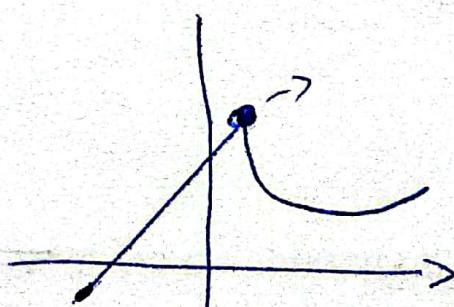
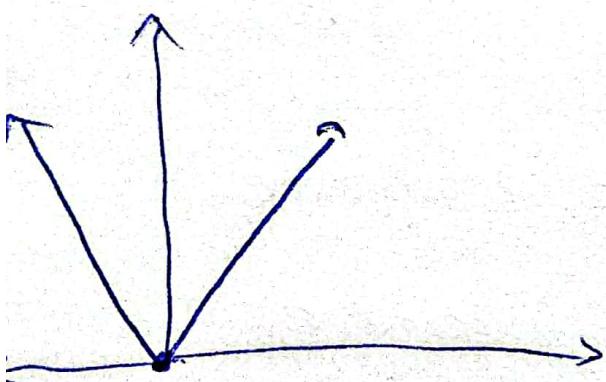
$$\frac{d}{dy} \log(g) = \frac{1}{e^{\log(g)}} = \frac{1}{g}$$

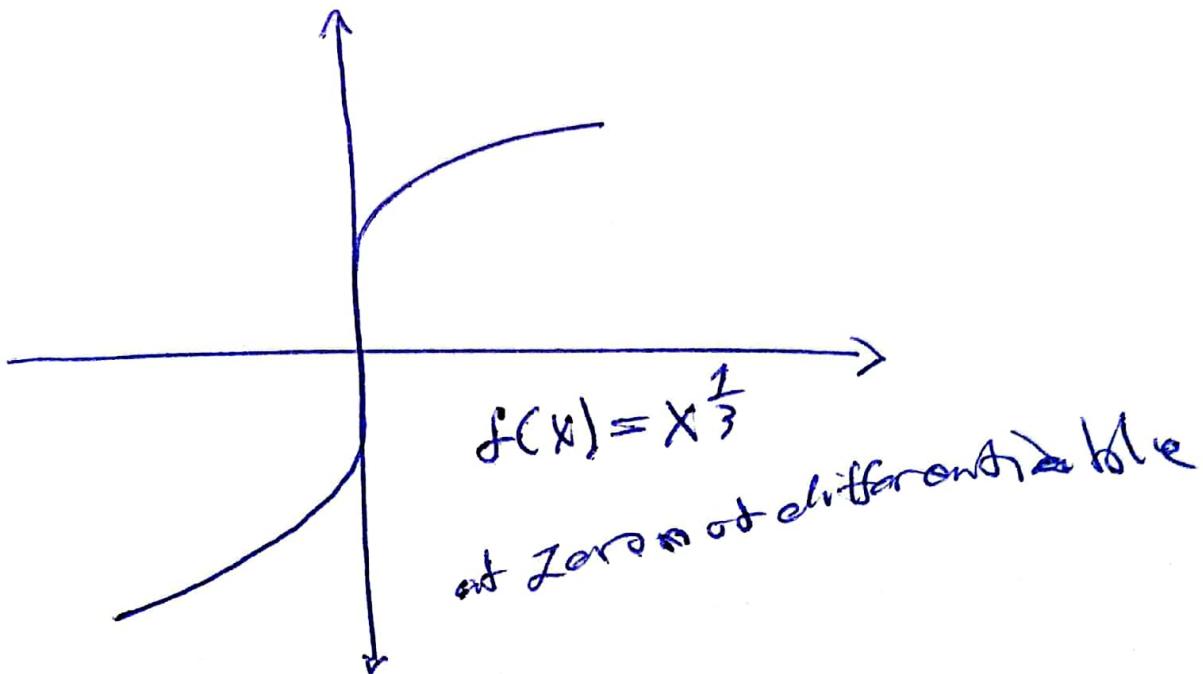
$$\frac{d}{dg} \log(g) = \frac{1}{g}$$

"non-differentiable functions"

$$f(x) = |x|$$

non-differentiable





### non-differentiable functions

- 1) any function with a cusp or a corner.
- 2) any function with a jump discontinuity
- 3) any function with a vertical tangent

### Properties of derivative

- 1) multiplication ~~of~~ by a scalar

~~if  $f = cg$~~  if  $f = cg'$

then  $f' = cg'$

2) sum rules

if  $f(x) = 2x$  and  
 $g(x) = x^2$

if  $f = g + h$

$$(f+g)'(x) = ?$$

then  $f' = g' + h'$   $\Rightarrow f'(x) + g'(x) = 2 + 2x$

3) the product rule

if  $f(x) = x e^x$

$$f'(x) = ?$$

if  $f = gh$

$$x e^x + (x \cdot e^x)$$

then  $f' = g'h + gh'$

~~$x e^x$~~

$$e^x + x e^x$$

4) the chain rule

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

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

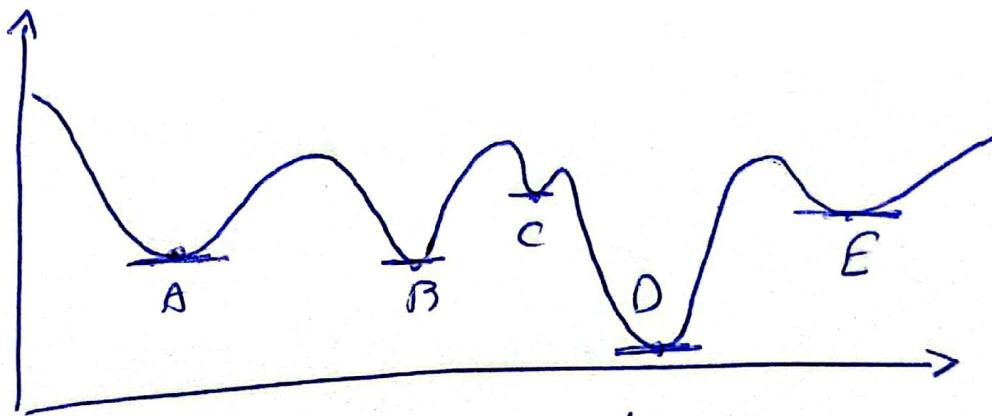
$$\frac{d}{dt} f(g(h(t))) = \frac{df}{dg} \cdot \frac{dg}{dh} \cdot \frac{dh}{dt}$$

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

if  $f(x) = e^{2x}$

$f'(x) = ?$

$$e^{2x} \cdot 2$$



slope = 0 in A, B, C & D ~~and E~~

D is global ~~minimum~~ minimum.

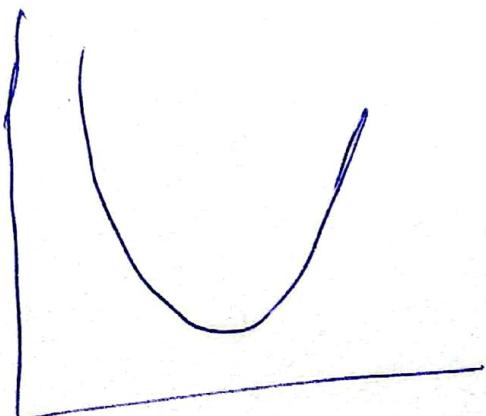
A, B, C and E are local minima.

"Two power line problem"

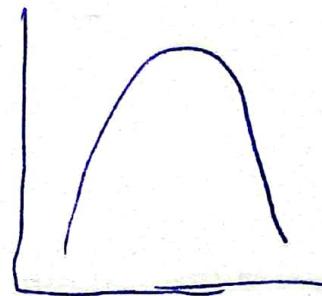
$$(x-a)^2 + (x-b)^2 = \text{total cost}$$

it  
is  
quadratic function

we should minimize  
the total cost.



or



$$\frac{d}{dx} [(x-a)^2 + (x-b)^2] = 0 \quad \begin{matrix} \text{minimum} \\ \text{or} \\ \text{maximum point} \end{matrix}$$

$$2(x-a) + 2(x-b) = 0$$

$$2x - 2a + 2x - 2b = 0$$

$$4x - 2a - 2b = 0$$

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

"Three point like problem"

$$\frac{d}{dx} \left[ (x-a)^2 + (x-b)^2 + (x-c)^2 \right] = 0$$

~~max~~  
minimum  
point

$$2(x-a) + 2(x-b) + 2(x-c) = 0$$

$$2x-a+2x-2b+2x-2c=0$$

$$3x - a - b - c = 0$$

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

"In general The square loss"

$$\text{minimize } (x-a_1)^2 + (x-a_2)^2 + \dots + (x-a_n)^2$$

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

"optimization of log loss"

$$P^7 \cdot (1-P)^3 = g(P)$$

coin toss  
 $T = 0.7, H = 0.3$

$$\frac{d}{dP} [P^7 \cdot (1-P)^3] = 0$$

$$7P^6(1-P)^3 + P^7 \cdot 3(1-P)(-1) = 0$$

$$= P^6 / (1-P)^2 / (7-10P) = 0$$

$\cancel{P \neq 0}$        $\cancel{1-P}$        $P = 0.7 \checkmark$

maximizing of  $g(P)$  is the same as

maximizing of  $\log(g(P))$  ~~so~~ so

$$\log(P^7 \cdot (1-P)^3) = 0$$

$$\Rightarrow \log P^7 + \log((1-P)^3) = 0$$

$$\Rightarrow \cancel{7} \log(P) + 3 \log(1-P) = G(P)$$

$$\frac{dG(P)}{dP} = \frac{d}{dP} \left[ 7 \log(P) + 3 \log(1-P) \right]$$

$$\Rightarrow 7 \frac{1}{P} + 3 \frac{1}{1-P} (-1)$$

$$\Rightarrow \frac{7(1-P) - 3P}{P(1-P)} = 0$$

$$7(1-P) - 3P = 0$$

$$7 - 7P - 3P = 0$$

$$-10P = -7$$

$$P = \frac{7}{10} = 0.7$$

Derivative of product is hard but  
the derivative of sum is much easier  
so we use logarithm to make a  
product problem to be a sum problem.

"Tangent Planes"

"Partial derivative"

$$f(x, y) = x^2 + y^2 \rightarrow \text{plot in 3-Dimensions}$$

$$\text{fix } y=4 \Rightarrow f(x, 4) = x^2 + 4^2$$

$$\frac{d}{dx} f(f(x, 4)) = 2x \rightarrow \begin{matrix} \text{partial derivative of } f \\ \text{with respect to } x \end{matrix}$$

$$\text{fix } x=2 \Rightarrow f(2, y) = 2^2 + y^2$$

$$\frac{d}{dy} (f(2, y)) = 2y \leftarrow \begin{matrix} \text{partial derivative of } f \\ \text{with respect to } y \end{matrix}$$

If we fix one variable, so the function  
will be a function with only one variable.  
and the slope of this function will be  
"Partial derivative"

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

$$\frac{\partial f}{\partial x} = \cancel{6} 6y^3 x \rightarrow \text{Partial derivative of } f \text{ with respect to } x.$$

$$\frac{\partial f}{\partial y} = 9x^2 y^2 \rightarrow \text{Partial derivative of } f \text{ with respect to } y.$$

Gradient

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

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

$$\text{Gradient} = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

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

a vector contains both  
partial derivatives.

in general Gradient  $\Rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$

this vector describes the slope of  
the two lines of the plane.  
 $\uparrow$   
tangent

if  $f(x, y) = x^2 + y^2$ , calculate the gradient  
 $f, \nabla f$ , at  $(2, 3)$  (19)

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

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

$$\nabla f = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

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

$$\text{so } \nabla f = \begin{bmatrix} 4 \\ 6 \end{bmatrix}$$

### Gradient and minima and maxima

The minimum of this function  ~~$f(x, y) = x^2 + y^2$~~  is where when both of the slopes of tangent lines are zero in other words the tangent plane in this case (minima) is parallel to the floor.

So minimum in this case is:

$$\frac{\partial f}{\partial x} = 0 \quad \text{and} \quad \frac{\partial f}{\partial y} = 0$$

50

$$2x=0 \quad \text{and} \quad 2y=0$$

$(x, y) = (0, 0)$  ← minimum point.

if  $f(x, y) = 85 - \frac{1}{90}x^2(3x-12)y^2(y-6)$

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

$$\frac{\partial f}{\partial x} = \cancel{-\frac{1}{90}x} - \frac{1}{90}x(3x-12)y^2(y-6)$$

$$\frac{\partial f}{\partial y} = -\frac{1}{90}x^2(x-6)y(3y-12)$$

find the minimum

$$-\frac{1}{90}x((3x-12)y^2(y-6)) = 0$$

?

$$-\frac{1}{90}x^2(x-6)y(3y-12) = 0$$

$$\begin{matrix} \swarrow x=0 & \downarrow x=6 & \downarrow y=0 & \searrow y=4 \end{matrix}$$

$$E(b, m) = 10m^2 + 3b^2 + 38 + 12mb - 92m - 2ab$$

$$\frac{dE}{dm} = 28m + 12b - 92$$

$$\frac{dE}{db} = 6b + 12m - 20$$

$$\begin{cases} 28m + 12b - 92 = 0 \\ (6b + 12m - 20 = 0) \times 2 \end{cases}$$

$$12b + 24m - 40 = 0$$

$$\begin{array}{r} - 12b + 28m - 92 = 0 \\ \hline - 4m + 2 = 0 \end{array}$$

$$m = \frac{1}{2}$$

$$b = \frac{7}{3}$$

$E(m = \frac{1}{2}, b = \frac{7}{3}) \approx 4.167$  the minimum value of the cost.

## "Gradient descent"

(22)

$$f(x) = e^x - \log(x)$$

$$f'(x) = e^x - \frac{1}{x}$$

$$e^x - \frac{1}{x} = 0$$

$$e^x = \frac{1}{x}$$

$$x \approx 0.5671\dots \rightarrow \text{omega constant.}$$

in gradient descent:

new point = old point -  $\underbrace{(\text{slope} \times \text{learning-rate})}_{\alpha}$

$$x_1 = x_0 - \underbrace{[f'(x_0) \times \text{learning-rate}]}_{\alpha}$$

$\alpha$

## Gradient Descent Algorithm

function:  $f(x)$ goal: find minimum of  $f(x)$ step 1: Define a learning rate ( $\alpha$ )choose a starting point  $(x_0)$ 

step 2:

$$\text{update: } x_k = x_{k-1} - \alpha f'(x_{k-1})$$

step 3:

repeat step 2 until you are close enough to the true minimum  $x^*$

$$\text{exp: } f(x) = e^{-\log(x)} \quad f'(x) = e^{-\frac{1}{x}}$$

$$\text{start: } x = 0.05 \quad \text{rate: } \alpha = 0.005$$

$$\text{find: } f'(0.05) = -18.9$$

$$\text{move by } -0.005f'(0.05) \quad x \rightarrow 0.1447$$

find:

$$f'(0.144J) = -5.7552$$

move by  $-0.005 f'(0.144J) \rightarrow 0.1735$

### Idea for gradient descent

Initial position:  $(x_0, y_0)$

direction of ~~gradient descent~~:  $\nabla f$

direction of ~~gradient descent~~:  $-\nabla f$

updated position:  $(x_0, y_0) - \alpha \nabla f$

learning-rate

better point

Expt:  $T = f(x, y) = 85 - \frac{1}{90} x^2(x-6)y^2(y-6)$

start:  $x = 0.5, y = 0.5$

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

$$\frac{\partial f}{\partial x} = -\frac{1}{90} x(3x-12)y^2(y-6)$$

$$\frac{\partial f}{\partial y} = -\frac{1}{90} x^2(x-6)y(3y-72)$$

$$\nabla f = \begin{bmatrix} -\frac{1}{90}x(3x-72)y^2(y-6) \\ -\frac{1}{90}x^2(x-6)y(8y-72) \end{bmatrix}$$

$$\nabla f(0.5, 0.6) = \begin{bmatrix} -0.2134 \\ -0.0935 \end{bmatrix}$$

move by:  $-0.05 \cdot \nabla f(0.5, 0.6)$



learning-rate

$$x \rightarrow 0.5057$$

$$y \rightarrow 0.6047$$

⋮  
⋮

repeat the process many times.

### Linear Regression: Gradient Descent

$$\nabla E = [28m + 72b - 92, 6b + 12m - 20]$$

initial

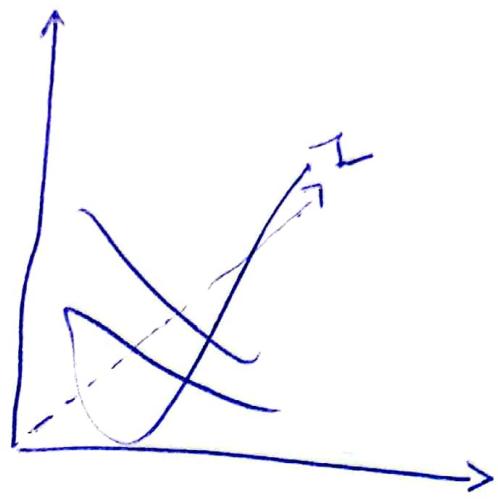
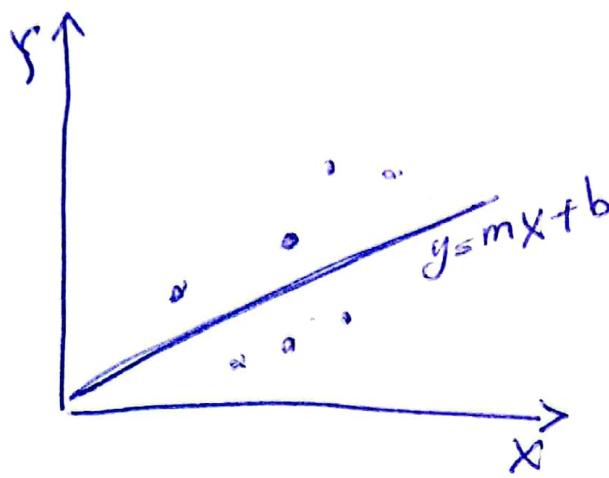
$m = ?$

$b = ?$

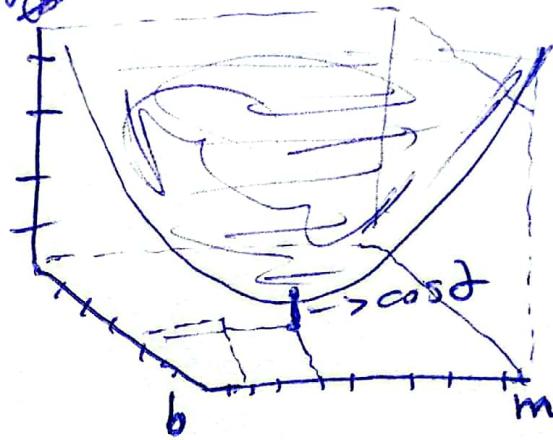
start with  $(m_0, b_0) \rightarrow$  some random points

iterate:

$$(m_{k+1}, b_{k+1}) = (m_k, b_k) - \alpha \nabla E(m_k, b_k)$$



square loss



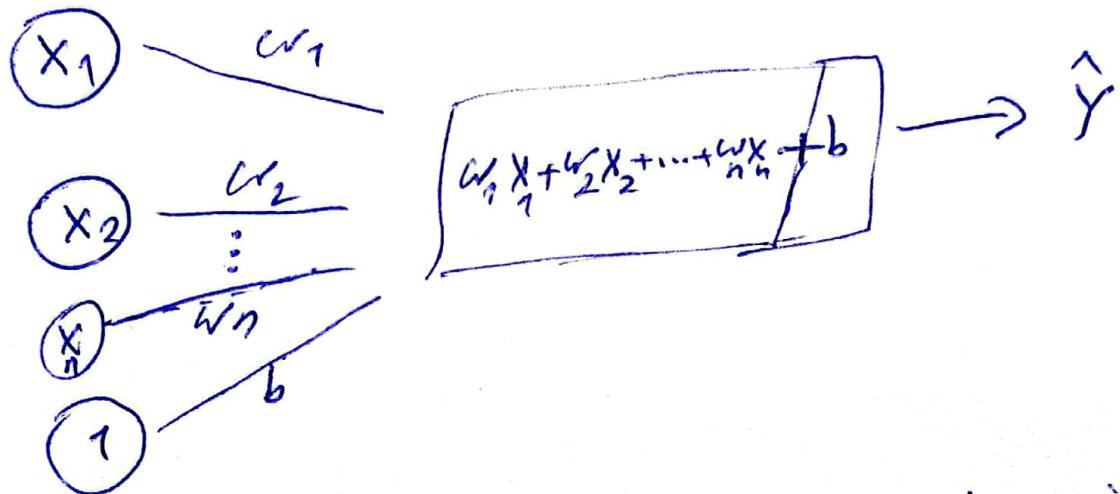
$$L(m, b) = \frac{1}{2n} [(mx_1 + b - y_1)^2 + \dots + (mx_n + b - y_n)^2]$$

cost function

$$\begin{bmatrix} m_n \\ b_n \end{bmatrix} = \begin{bmatrix} m_{n-1} \\ b_{n-1} \end{bmatrix} - \alpha \nabla L_i(m_{n-1}, b_{n-1})$$

# regression with a perceptron

## "Single Layer neural network Perceptron"



the goal is to find the optimal weights  
and bias ( $w_1, w_2, \dots, w_n$  and  $b$ )

↙  
how: minimize the error in the prediction  
with using the loss function

## regression with a perceptron

Prediction function:

$$\hat{y} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

Loss function

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

main goals

find ~~the~~  $w_1, w_2, \dots, w_n, b$  that give  $\hat{y}$  with the least error.

To find optimal values for:  $w_1, w_2, \dots, w_n, b$

we use gradient descent

$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1}$$

$$w_2 \rightarrow w_2 - \alpha \frac{\partial L}{\partial w_2}$$

$$\vdots \\ w_n \rightarrow w_n - \alpha \frac{\partial L}{\partial w_n}$$

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

# regression with a perceptron

Prediction function:

$$\hat{y} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

Loss function

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

main goals:

find ~~the~~  $w_1, w_2, \dots, w_n, b$  that give  $\hat{y}$  with the least error.

To find optimal values for:  $w_1, w_2, \dots, w_n, b$

we use gradient descent

$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1}$$

$$w_2 \rightarrow w_2 - \alpha \frac{\partial L}{\partial w_2}$$

$$\vdots \\ w_n \rightarrow w_n - \alpha \frac{\partial L}{\partial w_n}$$

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

$$\frac{\partial L}{\partial b}$$

~~$\hat{y} = w_1 x_1 + w_2 x_2 + b$~~

prediction function:  $\hat{y} = w_1 x_1 + w_2 x_2 + b$

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

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_1} = -(y - \hat{y}) \times x_1$$

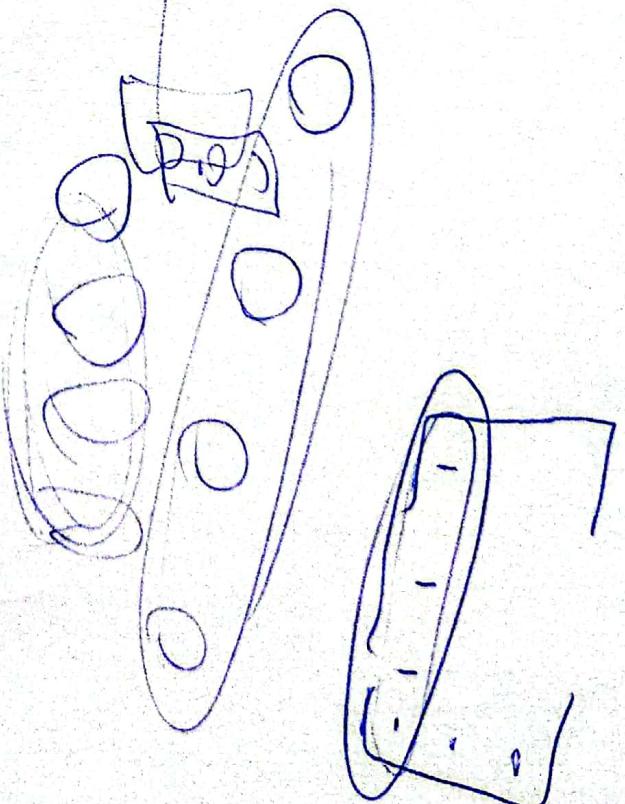
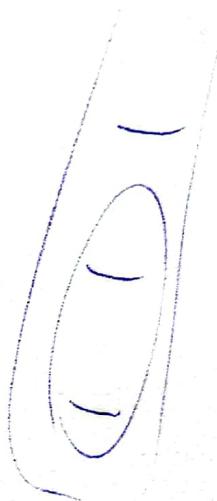
$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_2} = -(y - \hat{y}) \times x_2$$

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

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

$$\frac{\partial \hat{y}}{\partial w_1} = x_1$$

$$\frac{\partial \hat{y}}{\partial w_2} = x_2$$



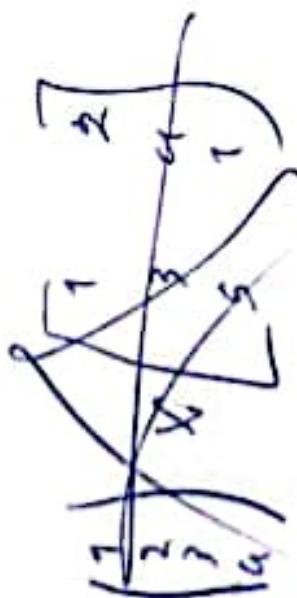
30

So: update the parameters

$$w_1 \rightarrow w_1 - \alpha (-x_1(y - \hat{y}))$$

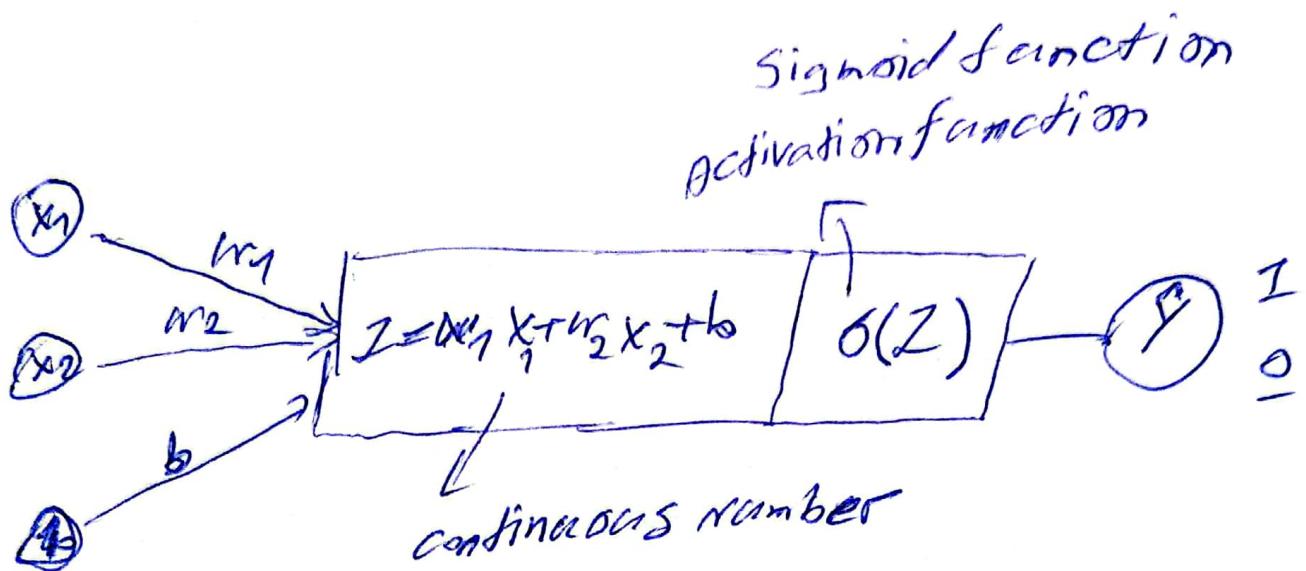
$$w_2 \rightarrow w_2 - \alpha (-x_2(y - \hat{y}))$$

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

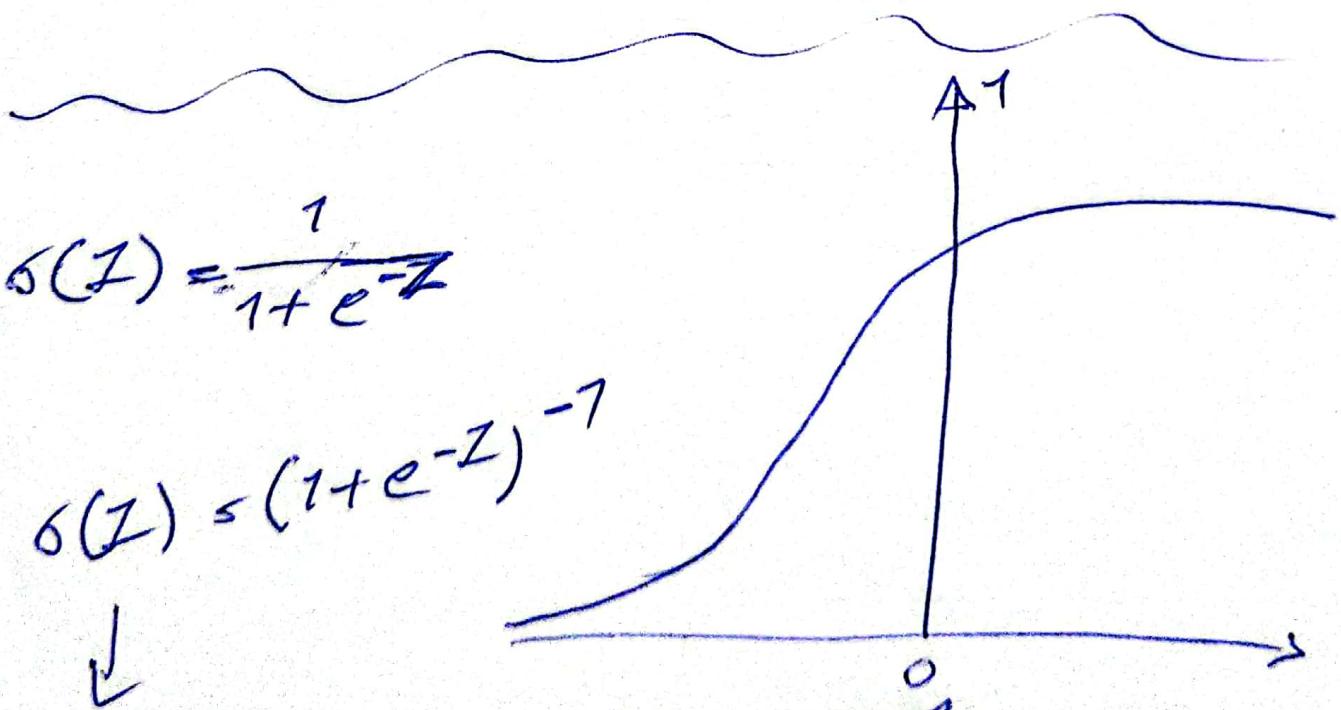


(31)

## Classification



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



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

$$\frac{d\sigma(Z)}{dZ} = \frac{d}{dZ} (1+e^{-Z})^{-1}$$

$$\begin{aligned}
 \frac{d}{dI} \delta(I) &= -1(1+e^{-I})^{-2-1} \cdot \left( \frac{d}{dI} (1+e^{-I}) \right) \quad (32) \\
 &= -1(1+e^{-I})^{-2} \cdot \left( \frac{d}{dI} (1) + \frac{d}{dI} (e^{-I}) \right) \\
 &= -1(1+e^{-I})^{-2} \cdot (0 + e^{-I} \left( \frac{d}{dI} (-I) \right)) \\
 &= -1(1+e^{-I})^{-2} (e^{-I}) (-1) \\
 &= (1+e^{-I})^{-2} (e^{-I}) \\
 &= \frac{1}{(1+e^{-I})^2} (e^{-I})
 \end{aligned}$$

$$\boxed{s = \frac{e^{-I}}{(1+e^{-I})^2}}$$

$$\begin{aligned}
 \frac{d}{dI} d(F) &= \frac{e^{-I} + 1 - 1}{(1+e^{-I})^2} \\
 &= \frac{1 + e^{-I} - 1}{(1+e^{-I})^2} \\
 &= \frac{1 + e^{-I}}{(1+e^{-I})^2} - \frac{1}{(1+e^{-I})^2}
 \end{aligned}$$

$$= \frac{1}{(1+e^{-z})} - \frac{1}{(1+e^{-z})^2}$$

$$= \frac{1}{(1+e^{-z})} - \left( \frac{1}{(1+e^{-z})} \right) \cdot \left( \frac{1}{(1+e^{-z})} \right)$$

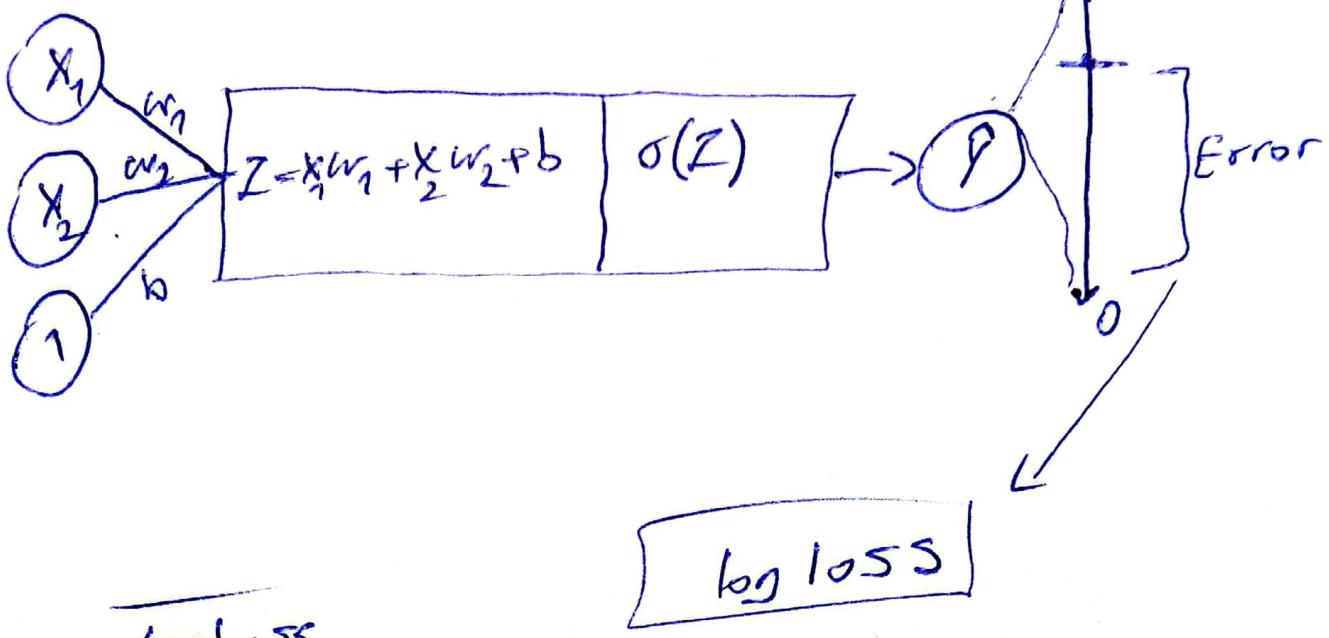
$$= \frac{1}{(1+e^{-z})} \left( 1 - \frac{1}{(1+e^{-z})} \right)$$

Recall that:  $e(z) = \frac{1}{1+e^{-z}}$

So:

$$\frac{d}{dz} \sigma(z) = \sigma(z)(1-\sigma(z))$$

$$\frac{d\hat{y}}{dy} = \hat{y}(1-\hat{y})$$



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

note:

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

To find optimal values for:

$$w_1, w_2, b$$

we need gradient descent:

$$w_1 = w_1 - \alpha \frac{dL}{dw_1}$$

$$w_2 = w_2 - \alpha \frac{dL}{dw_2}$$

$$b = b - \alpha \frac{dL}{db}$$

$\geq ?$

$$\frac{dL}{dw_1} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_1} \quad ?$$

$$\frac{dL}{dw_2} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_2} \quad ?$$

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{db} \quad ?$$

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

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

$$\frac{dL}{d\hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

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

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

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

$$\frac{d\hat{y}}{dw_1} = \hat{y}(1-\hat{y})x_1$$

$$\frac{d\hat{y}}{dw_2} = \hat{y}(1-\hat{y})x_2$$

$$\frac{dy}{db} = \hat{y}(1-\hat{y})$$

so:

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{db} = \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})} \cdot \hat{y}(1-\hat{y})$$

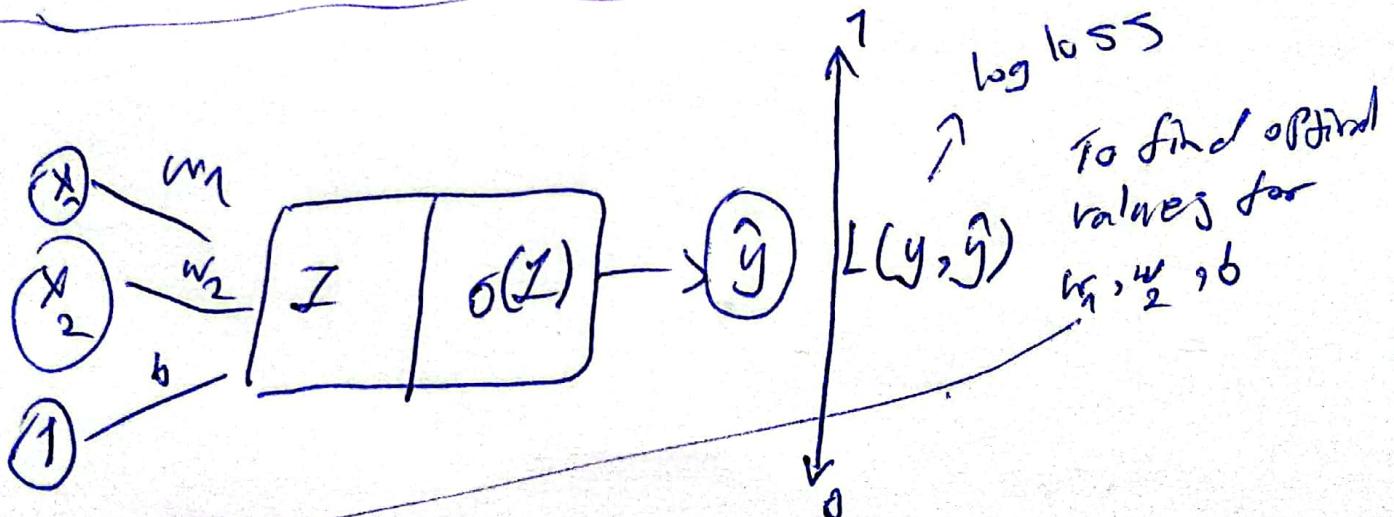
$$\frac{dL}{dw_1} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_1} = \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})} \cdot \hat{y}(1-\hat{y})x_1$$

$$\frac{dL}{dw_2} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_2} = \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})} \cdot \hat{y}(1-\hat{y})x_2$$

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

$$\frac{\partial L}{\partial w_1} = -(y - \hat{y}) X_1$$

$$\frac{\partial L}{\partial w_2} = -(y - \hat{y}) X_2$$



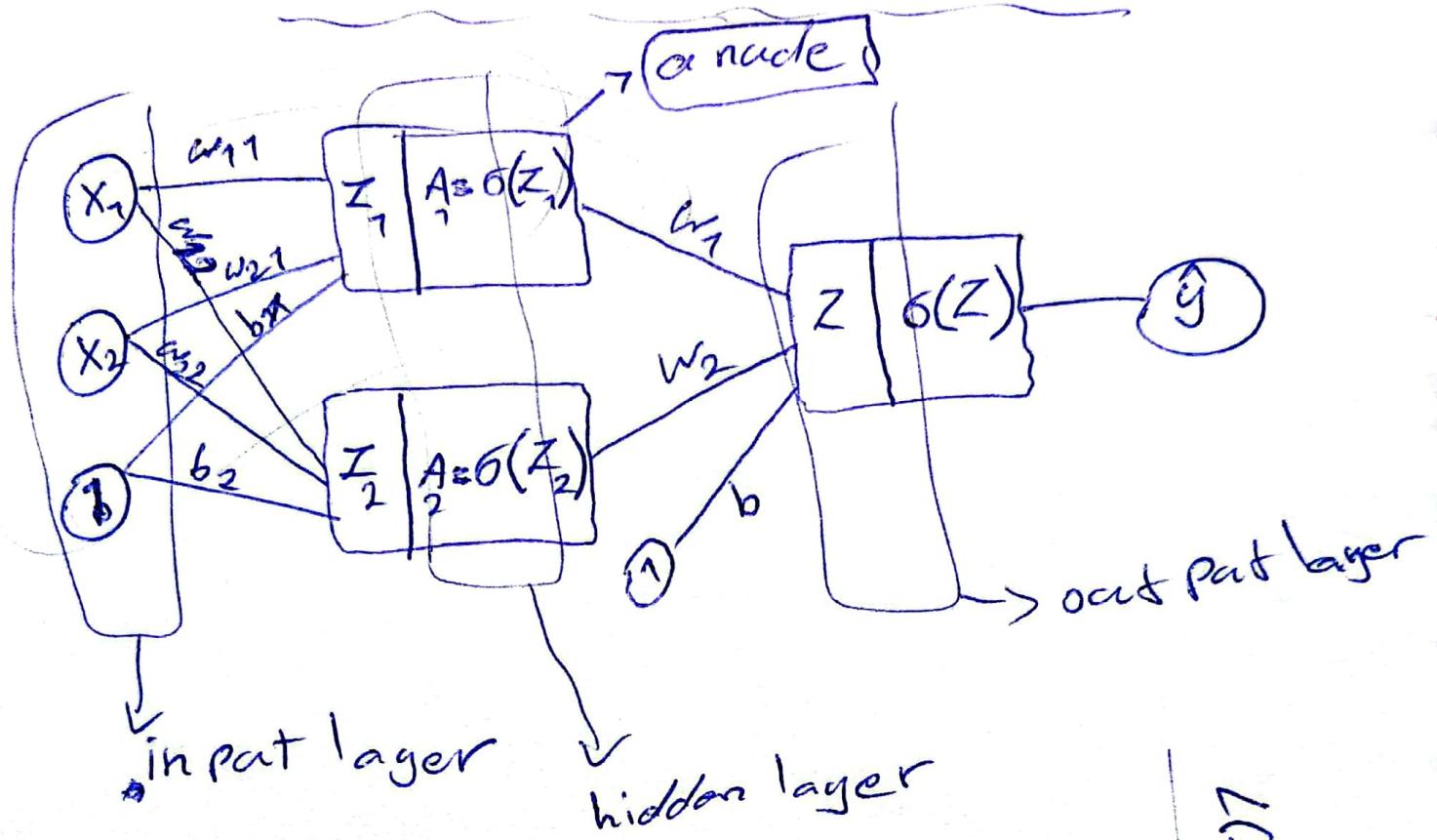
$$w_1 \rightarrow w_1 - \alpha (-x_1(y - \hat{y}))$$

$$w_2 \rightarrow w_2 - \alpha (-x_2(y - \hat{y}))$$

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

A neural network of depth 2

classification



$$\begin{cases} I_1 = x_1 w_{11} + x_2 w_{21} + b_1 \\ A_1 = \sigma(I_1) \end{cases}$$

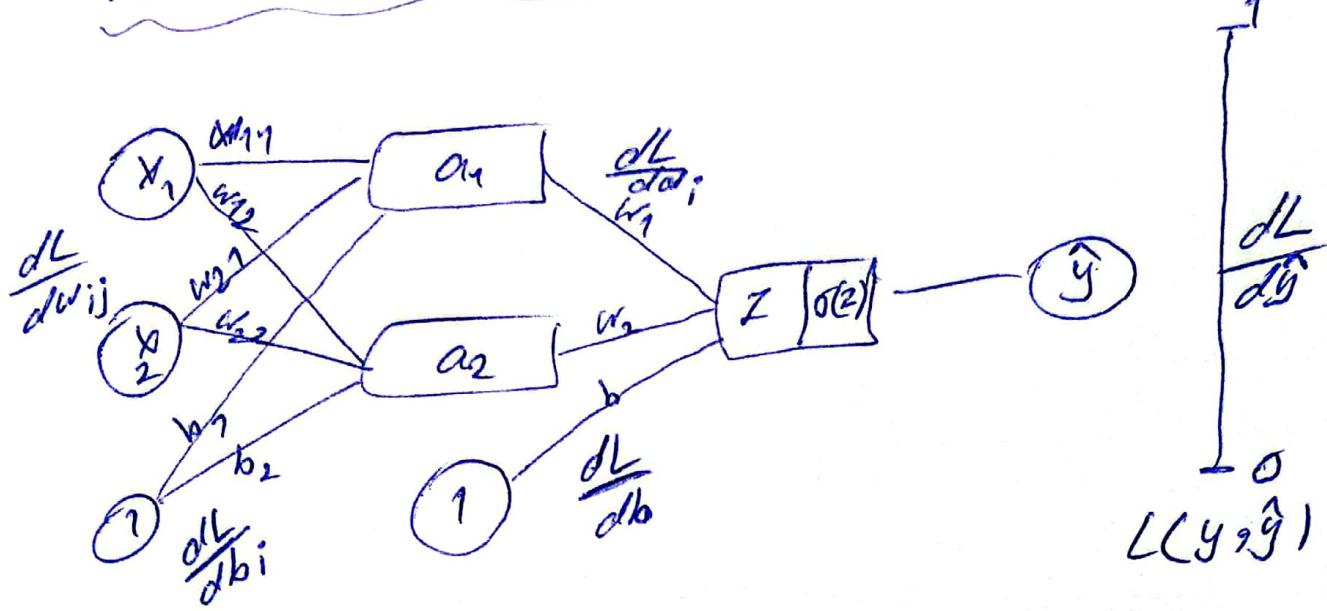
$$\begin{cases} I_2 = x_1 w_{12} + x_2 w_{22} + b_2 \\ A_2 = \sigma(I_2) \end{cases}$$

$$\begin{cases} Z = A_1 w_1 + A_2 w_2 + b \\ \hat{y} = \sigma(Z) \end{cases}$$

log loss

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

To reduce the loss function:  $L(y, \hat{y})$



$$\begin{aligned}
 \frac{\partial L}{\partial w_{11}} &= \frac{\partial I_1}{\partial w_{11}} \cdot \frac{\partial a_1}{\partial I_1} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial L}{\partial \hat{y}} \\
 &= x_1 \cdot a_1(1-a_1) \cdot w_1 \cdot \hat{y}(1-\hat{y}) \cdot \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})} \\
 &= -x_1 w_1 a_1 (1-a_1) (y-\hat{y})
 \end{aligned}$$

perform a gradient descent

$$w_{11} \rightarrow w_{11} - \alpha \frac{\partial L}{\partial w_{11}}$$

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

$$\frac{\partial L}{\partial b_1} = \frac{\partial Z_1}{\partial b_1} \cdot \frac{\partial a_1}{\partial Z_1} \cdot \frac{\partial Z}{\partial a_1} \cdot \frac{\partial g}{\partial Z} \cdot \frac{\partial L}{\partial g}$$

$$= 1 \cdot a_1(1-a_1) \cdot w_1 \cdot \cancel{g'(1-g)} \cdot \frac{-(y-\hat{y})}{\cancel{g(1-g)}}$$

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

Perform gradient descent

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

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

$$w_{12} \rightarrow w_{12} + \alpha x_2 w_1 a_1(1-a_1)(y-\hat{y})$$

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

$$w_{21} \rightarrow w_{21} + \alpha x_1 w_2 a_2(1-a_2)(y-\hat{y})$$

$$w_{22} \rightarrow w_{22} + \alpha x_2 w_2 a_2(1-a_2)(y-\hat{y})$$

$$b_2 \rightarrow b_2 + \alpha w_2 a_2(1-a_2)(y-\hat{y})$$

update  
the first  
layer of  
network

$$\frac{\partial L}{\partial w_i} = \frac{\partial Z}{\partial w_i} \cdot \frac{\partial \hat{y}}{\partial Z} \cdot \frac{\partial L}{\partial \hat{y}}$$

$$\alpha = \alpha_i \cdot \hat{y}(1-\hat{y}) \cdot \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})}$$

$$= -\alpha_i(y-\hat{y})$$

GD

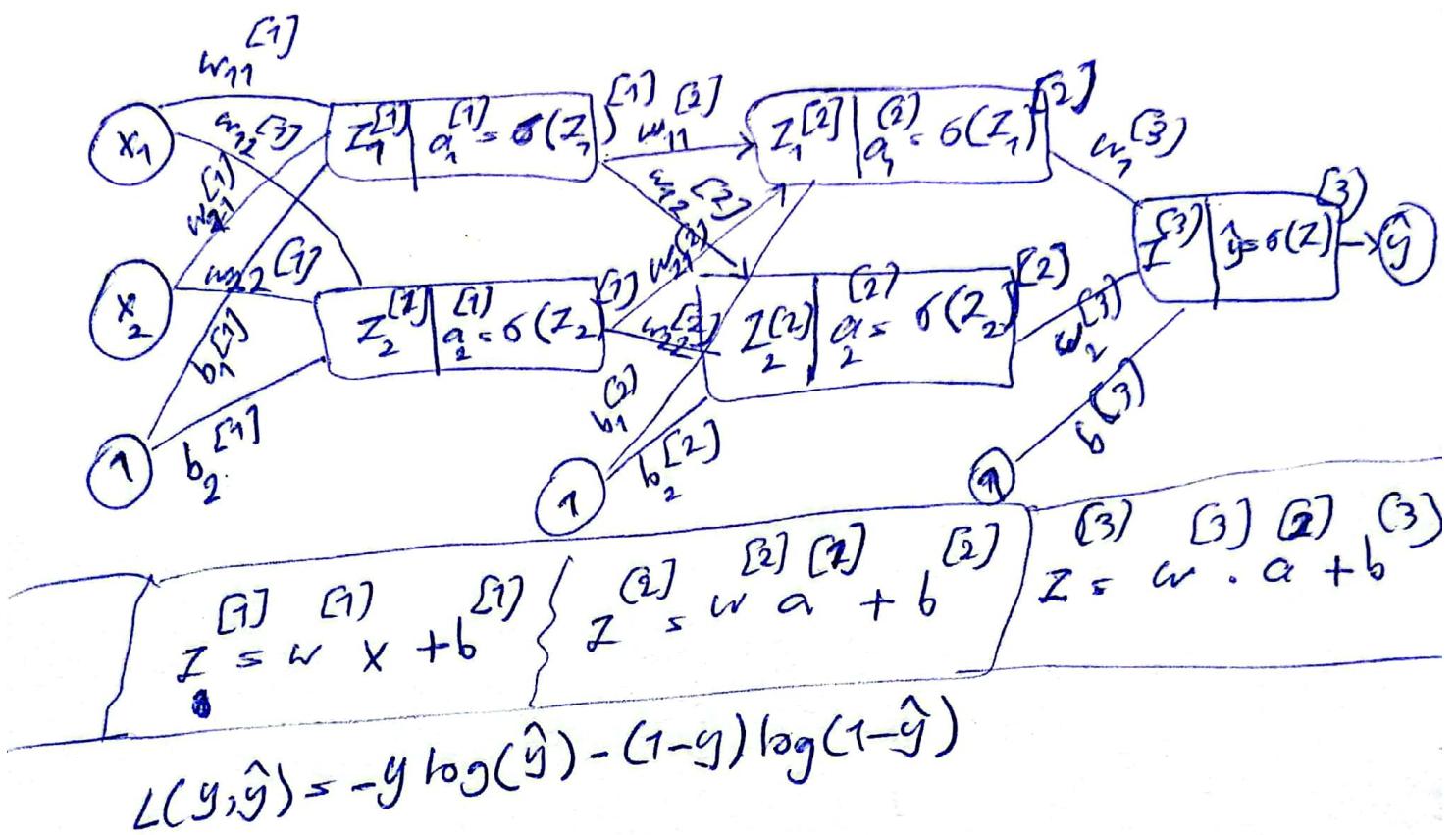
$$w_1 \rightarrow w_1 - \alpha(-\alpha_1(y-\hat{y}))$$

$$\rightarrow w_1 + \alpha \alpha_1 (y-\hat{y})$$

$$w_2 \rightarrow w_2 + \alpha \alpha_2 (y-\hat{y})$$

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

} update the  
second  
layer of  
NN

Back propagation in intro


$$\frac{dL}{dw^{(3)}} = \frac{dz^{(3)}}{da^{(3)}} \cdot \frac{da^{(3)}}{dI^{(3)}} \cdot \frac{dL}{da^{(3)}}$$

$$\frac{dL}{db^{(3)}} = \frac{dL}{db^{(3)}} \cdot \frac{da^{(3)}}{dI^{(3)}} \cdot \frac{dL}{da^{(3)}}$$

$$\frac{dL}{dw^{(2)}} = \frac{dL}{da^{(3)}} \cdot \frac{da^{(3)}}{dI^{(3)}} \cdot \frac{dI^{(3)}}{da^{(2)}} \cdot \frac{da^{(2)}}{dI^{(2)}} \cdot \frac{dI^{(2)}}{dw^{(2)}}$$

$$\frac{dL}{db^{(2)}} = \frac{dL}{da^{(3)}} \cdot \frac{da^{(3)}}{dI^{(3)}} \cdot \frac{dI^{(3)}}{da^{(2)}} \cdot \frac{da^{(2)}}{dI^{(2)}} \cdot \frac{dI^{(2)}}{dw^{(2)}}$$

$$\frac{dL}{da^{(1)}} = \frac{dL}{da^{(3)}} \cdot \frac{da^{(3)}}{dZ^{(3)}} \cdot \frac{dZ^{(2)}}{da^{(2)}} \cdot \frac{da^{(2)}}{dL^{(2)}} \cdot \frac{dL^{(2)}}{da^{(1)}} \cdot \frac{da^{(1)}}{dZ^{(1)}} \cdot \frac{dZ^{(1)}}{da^{(1)}} \quad (43)$$

$$\frac{dL}{da^{(1)}} = \frac{dL}{dZ^{(3)}} \cdot \frac{da^{(3)}}{dZ^{(3)}} \cdot \frac{dZ^{(2)}}{da^{(2)}} \cdot \frac{da^{(2)}}{dL^{(2)}} \cdot \frac{dL^{(2)}}{da^{(1)}} \cdot \frac{da^{(1)}}{dZ^{(1)}} \cdot \frac{dZ^{(1)}}{da^{(1)}}$$

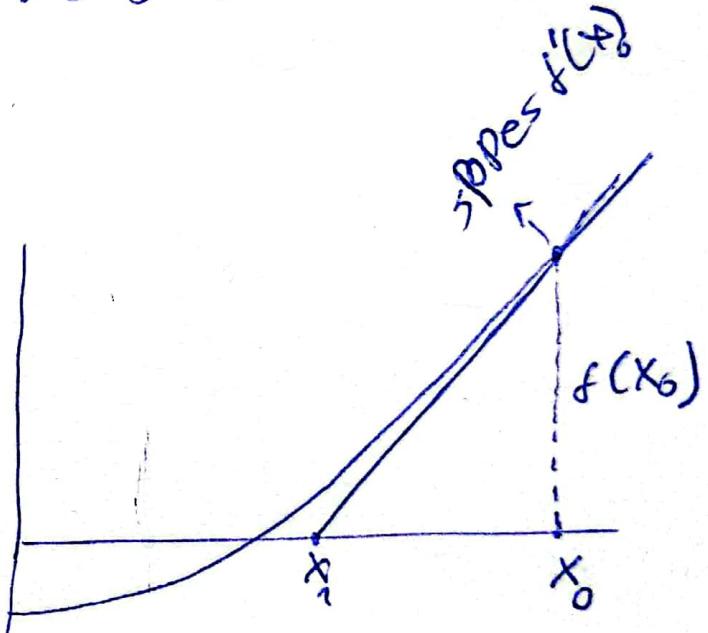
"newton's method"

- in principle it is used to find the zeros of a function.

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

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

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



Newton's method:

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

} NM for optimization:

Goal: minimize  $g(x)$

→ find zeros of

$g'(x)$

if:  $f(x) \rightarrow g'(x)$

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

NM

NMO

1) start with some  $x_0$

1) start with some  $x_0$

2) update:

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

2) updates

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

3) repeat 2) until you find  
the root

3) repeat 2) until you find  
the candidate for minimum.

Exps

$$g(x) = e^x - \log(x)$$

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

find the zero of derivative:

$$g'(x) = e^x - \frac{1}{x} \xrightarrow{\text{der}} e^x + \frac{1}{x^2}$$

$$x_{k+1} = x_k - \frac{e^x - \frac{1}{x}}{e^x + \frac{1}{x^2}}$$

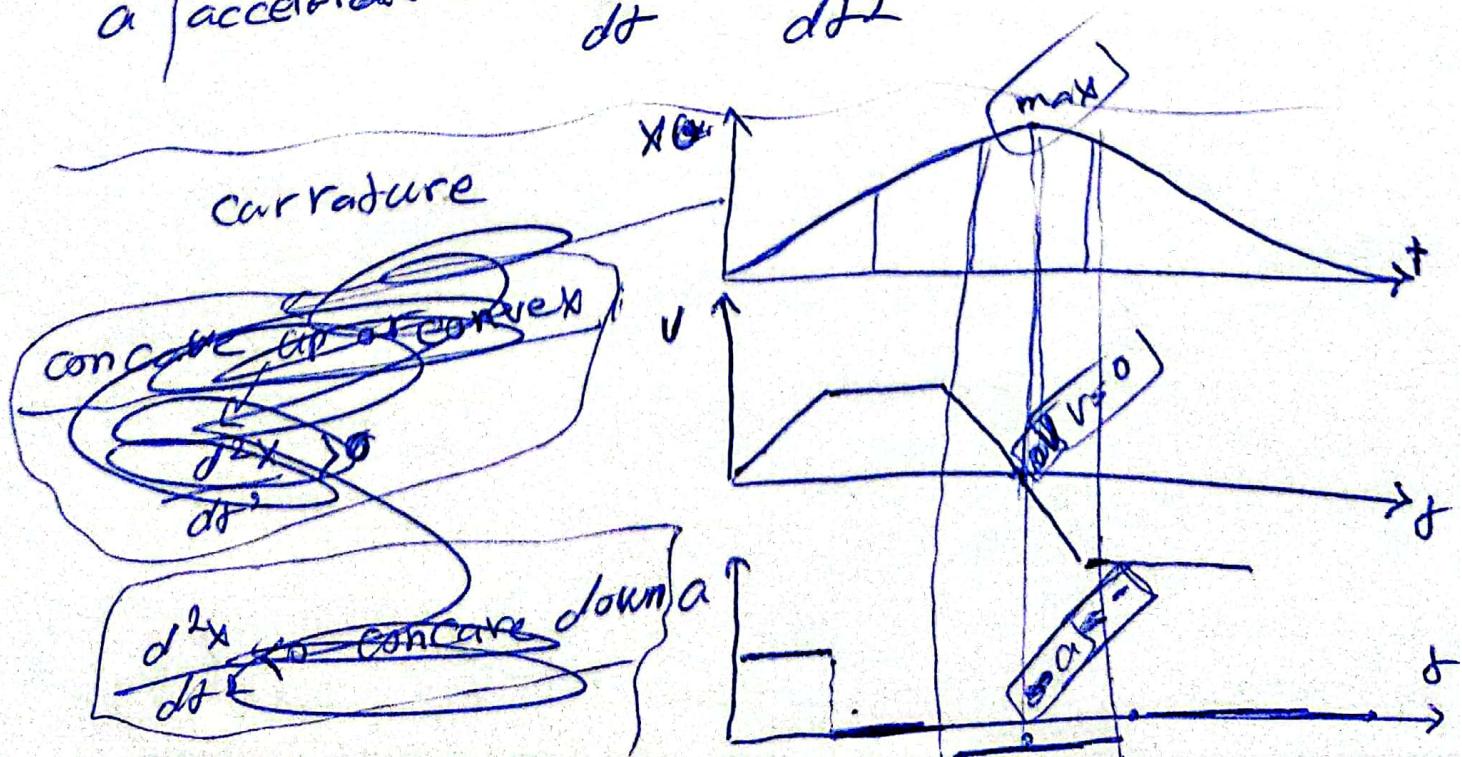
"The second derivative"

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

or

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

$x$	distance
$v$	velocity $\frac{dx}{dt}$
$a$	acceleration $\frac{dv}{dt} = \frac{d^2 x}{dt^2}$



$$\frac{d^2x}{dt^2} > 0 \quad \text{concave up or convex.}$$

$$\frac{d^2x}{dt^2} < 0 \quad \text{concave down.}$$

$\frac{d^2x}{dt^2} = 0$  need more information.

## The Hessian

$\mathbb{F}$ function	1 variable $f(x)$	2 variable $f(x, y)$
first derivative	$f'(x)$ rate of change of $f(x)$	$f_x(x, y)$ rate of change w.r.t $x$ $f_y(x, y)$ " " " w.r.t $y$
second derivative	$f''(x)$ rate of change of the rate of change of $f(x)$	$\nabla f = \begin{bmatrix} f_x(x, y) \\ f_y(x, y) \end{bmatrix}$
		$H = \begin{bmatrix} f_{xx}(x, y) & f_{xy}(x, y) \\ f_{yx}(x, y) & f_{yy}(x, y) \end{bmatrix}$

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

$$f_x(x, y) = 4x - y$$

$$f_y(x, y) = 6y - x$$

D

$$\frac{d^2 f}{dx^2} = f_{xx}(x, y) = 4$$

$$\frac{d^2 f}{dy^2} = f_{yy}(x, y) = 6$$

$$\text{Hessian} = \begin{bmatrix} 4 & -1 \\ -1 & 6 \end{bmatrix}$$

$$\frac{d^2 f}{dx dy} = f_{xy}(x, y) = -1$$

$$\frac{d^2 f}{dy dx} = f_{yx}(x, y) = -1$$

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

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

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

$$= (4-1)(6-1) - (-1)(-1)$$

$$= 12 - 10/1 + 23$$

$$\begin{cases} \lambda_1 = 6.447 \\ \lambda_2 = 3.55 \end{cases} > 0$$

if means the matrix  $\overset{(H)}{\text{is positive}}$

definite

so the function is concave up  
and the point  $(0,0)$  is minimum!

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

(49)

$$f_{xx}(x,y) = -4$$

$$f_x(x,y) = -4x - y \quad f_{xy}(x,y) = -1$$

$$f_y(x,y) = -6y - x \quad f_{yx}(x,y) = -1$$

$$f_{yy}(x,y) = -6$$

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

Hessian at (0,0)

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

gradient

$\sqrt{n}$  eigen values of Hessian  $M^n$

$$det(H(0,0) - I) = (-4-1)(-6-1) - (-1)(-1)$$

$$= 41^2 + 101 + 23 \quad \left\{ \begin{array}{l} l_1 = -3.59 \\ l_2 = -6.41 \end{array} \right\} \quad (0)$$

(0,0) is a maximum

## Saddle point

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

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

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

$$\det(H(0, 0) - I) = (4-1)(4-1) - 0$$

$\begin{cases} d_1 = -4 \\ d_2 = +4 \end{cases}$  we can not conclude anything

NN

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

or

$$x_k - f''(x_k)^{-1} \cdot f'(x_k)$$

2 variables:

$$\begin{bmatrix} x_{k+1} \\ y_{k+1} \end{bmatrix} = \begin{bmatrix} x_k \\ y_k \end{bmatrix} - H^{-1}(x_k, y_k) \cdot \nabla f(x_k, y_k)$$

exp:

(51)

$$f(x, y) = x^4 + 0.9y^4 + 4x^2 + 2y^2 - xy - 0.2x^2y$$

$$\begin{array}{c} x \\ \swarrow \\ f(x, y) \end{array} \quad 4x^3 + 8x - y - 0.4xy \quad [f'_x(x, y)]$$

$$\begin{array}{c} y \\ \searrow \\ f(x, y) \end{array} \quad \begin{array}{l} 3.2 \\ \cancel{0.9} \end{array} y^3 + 4y - x - 0.2x^2 \quad [f'_y(x, y)]$$

$$\begin{array}{c} x \\ \nearrow \\ f'_x(x, y) \end{array} \quad 12x^2 + 8 - 0.4y$$

$$\begin{array}{c} y \\ \searrow \\ f'_x(x, y) \end{array} \quad \cancel{10.8y^2 + 4} - 1 - 0.4x$$

$$\begin{array}{c} x \\ \nearrow \\ f'_y(x, y) \end{array} \quad \underline{-1 - 0.4x}$$

$$\begin{array}{c} y \\ \searrow \\ f'_y(x, y) \end{array} \quad 10.8y^2 + 4$$

$$H(0, 0) = \begin{bmatrix} 12x^2 + 8 - 0.4y & -1 - 0.4x \\ -1 - 0.4x & 10.8y^2 + 4 \end{bmatrix}$$

$$\nabla f = \begin{bmatrix} 4x^3 + 8x - y - 0.4xy \\ 3.2y^3 + 4y - x - 0.2x^2 \end{bmatrix}$$

(52)

start at some point  $\begin{bmatrix} x_0 \\ y_0 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$

$$\nabla f(4, 4) = \begin{bmatrix} 277.6 \\ 213.6 \end{bmatrix} \quad H(4, 4) = \begin{bmatrix} 198.4 & -2.6 \\ -2.6 & 157.6 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix} - \begin{bmatrix} 198.4 & -2.6 \\ -2.6 & 157.6 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 277.6 \\ 213.6 \end{bmatrix}$$

$$= \begin{bmatrix} 2.56 \\ 2.62 \end{bmatrix}$$

↓

$$\nabla f(2.56, 2.62) = \begin{bmatrix} 89.29 \\ 63.9 \end{bmatrix} \quad H(2.56, 2.62) = \begin{bmatrix} 86.83 & -2.032 \\ -2.032 & 69.39 \end{bmatrix}$$

$$\begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} 86.83 & -2.032 \\ -2.032 & 69.39 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 89.29 \\ 63.9 \end{bmatrix}$$

$$= \begin{bmatrix} 1.59 \\ 1.67 \end{bmatrix}$$

...  
needed 3 steps       $\rightarrow$  last zero

$$\begin{bmatrix} x_3 \\ y_3 \end{bmatrix} = \begin{bmatrix} 1.75 & x_{10}^{-1} \\ -2.05 & x_{10}^{-1} \end{bmatrix}$$