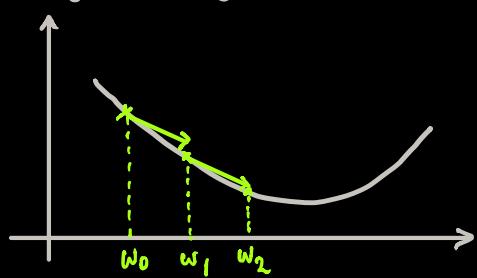


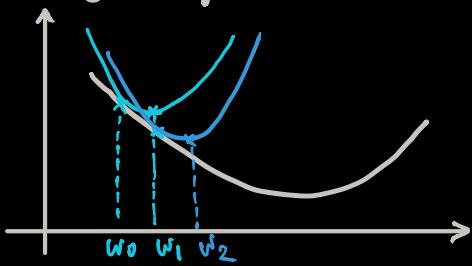
# MML 5

Gradient descent for loss function  $\ell(w)$  optimization  $[\min_w \ell(w)]$

Previously  
ML - Blink 4.4



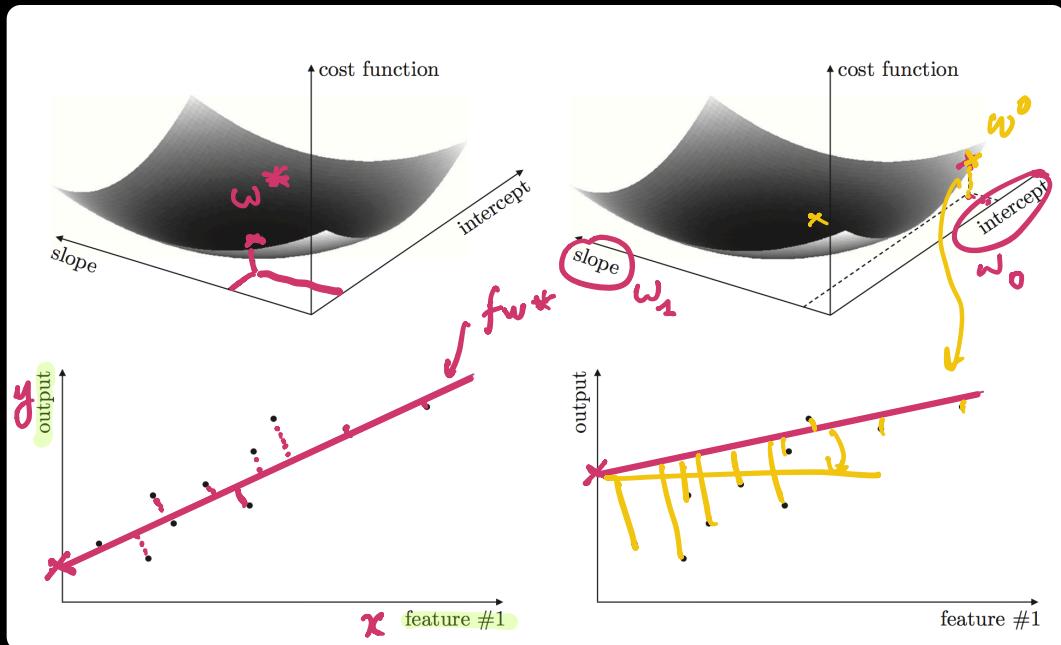
1 - Newton's method for optimization  
(second order)



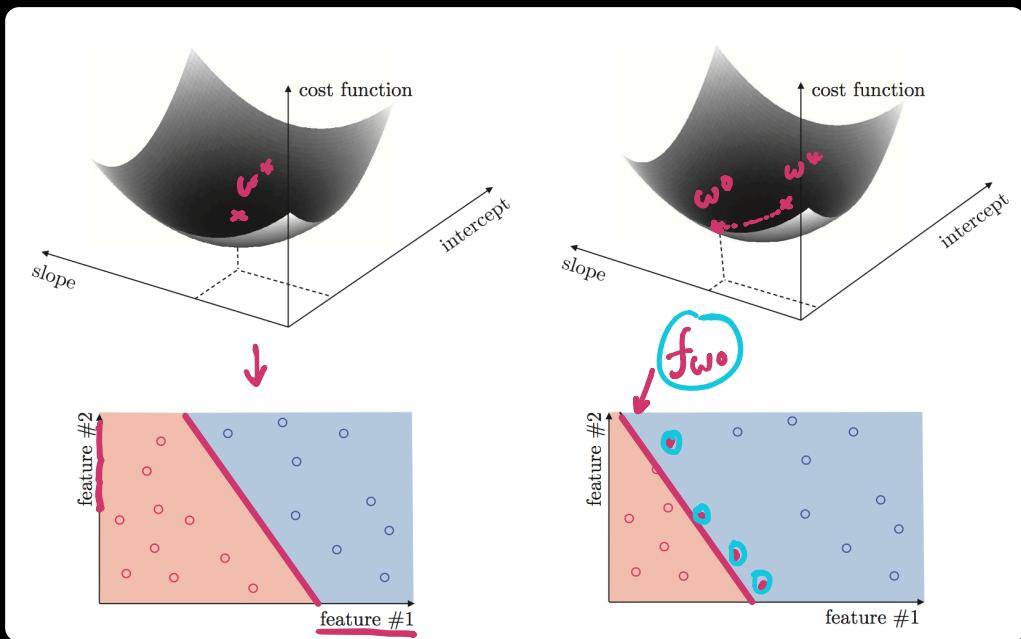
2 - Two major issues with gradient descent (zigzagging & vanishing  
direction ↑ magnitude)

# $\dot{\mathcal{M}} \otimes \dot{\mathcal{L}}_4 \rightarrow ?$

$$\arg \min_{\mathbf{w}} \overbrace{\mathcal{L}(\mathbf{w})}^{\mathcal{L}(\mathbf{w})} = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in \mathcal{D}} \underbrace{E(f_{\mathbf{w}}(\mathbf{x}^i), y^i)}_{E(f_{\mathbf{w}}(\mathbf{x}^i), y^i)} + R(\dots)$$

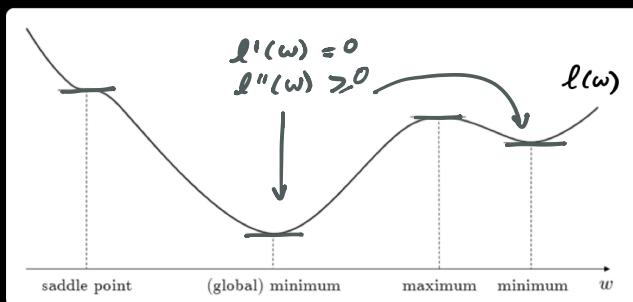


$$\arg \min_{\mathbf{w}} \overbrace{\mathcal{L}(\mathbf{w})}^{\mathcal{L}(\omega)} = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in \mathcal{D}} E(f_{\mathbf{w}}(\mathbf{x}^i), y^i) + R(\dots)$$

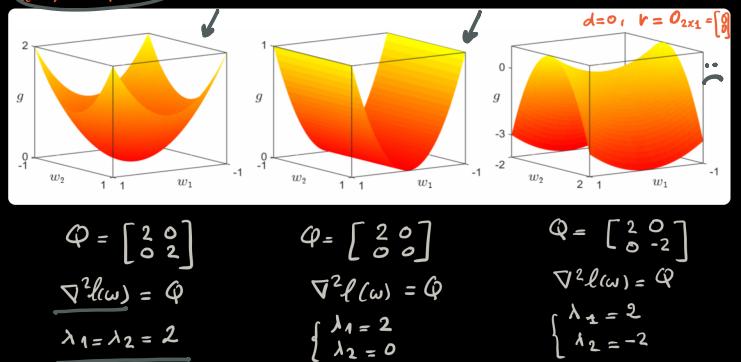


# Max $\mathcal{L}$ → ?

Univariate loss function  
 $w \in \mathbb{R}$



Multivariate loss function  
 $w \in \mathbb{R}^N$



**Algorithm 2.1** Gradient descent (with fixed step length)

**Input:** differentiable function  $\ell$ , fixed step length  $\alpha$ , and initial point  $\mathbf{w}^0$   
 $\varepsilon = 10^{-3}$   
 $k = 1$

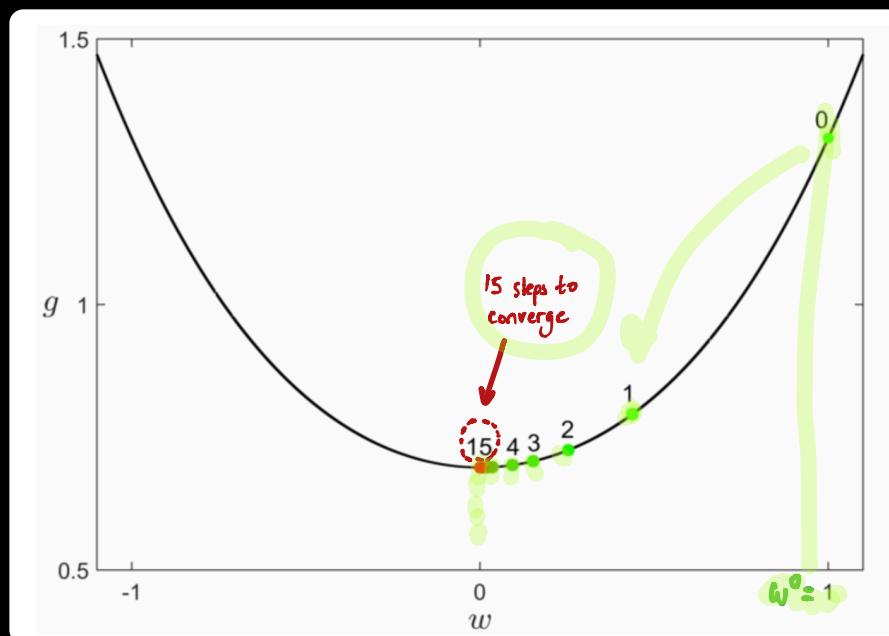
Repeat until stopping condition is met:  $\|\nabla \ell(\mathbf{w}^k)\|_2 \leq \varepsilon$

$$\underbrace{\mathbf{w}^k = \mathbf{w}^{k-1} - \alpha \nabla \ell(\mathbf{w}^{k-1})}_{k \leftarrow k + 1}$$

## Example 1

$w \in \mathbb{R}$  (scalar) univariate loss function

$$\{ l(w) = \log(1 + e^{w^2}) \quad [\text{derivative-calculator.net}]$$



- $l'(w) = \frac{2e^{w^2}w}{1+e^{w^2}}$
- $\alpha = 10^{-1} = 0.1$   
↑ (fixed by trial & error)
- $w^0 = 1$
- $\epsilon = 10^{-3} = 0.001$

[ derivative-calculator.net ]

YOUR INPUT:  
 $f(x) =$

$\ln(e^{x^2} + 1)$

Note: Your input has been rewritten/simplified.

Simplify Roots/zeros

FIRST DERIVATIVE:  
 $\frac{d}{dx}[f(x)] = f'(x) =$

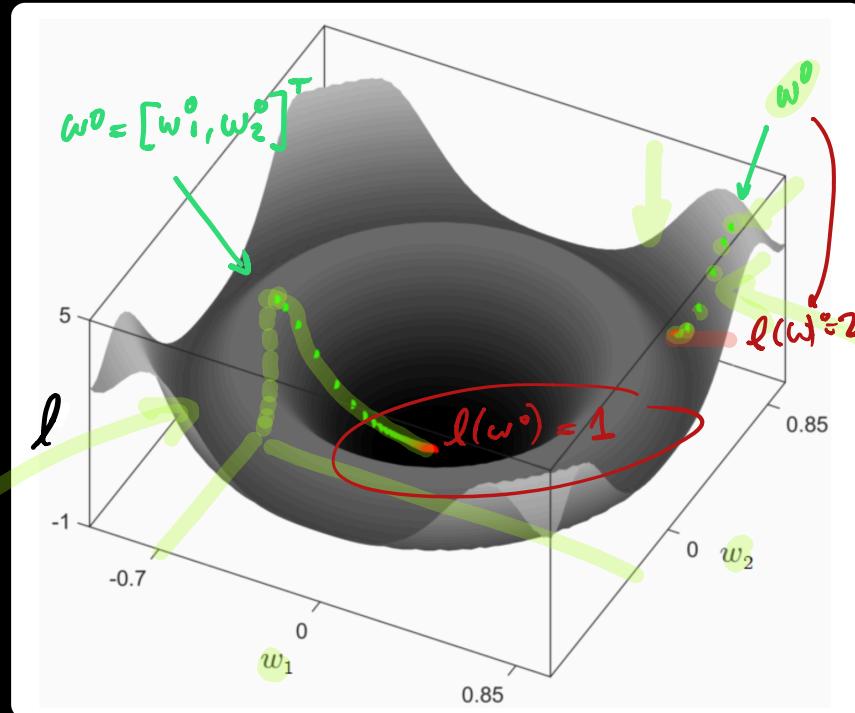
The steps of calculation are displayed.  
Click at any derivative  $\frac{d}{dx}[\dots]$  in order to show the rule that was applied.

$$\begin{aligned}
 & \frac{d}{dx} \left[ \ln(e^{x^2} + 1) \right] \\
 &= \frac{1}{e^{x^2} + 1} \cdot \frac{d}{dx} [e^{x^2} + 1] \\
 &= \frac{\frac{d}{dx} [e^{x^2}] + \frac{d}{dx} [1]}{e^{x^2} + 1} \\
 &= \frac{e^{x^2} \cdot \frac{d}{dx} [x^2] + 0}{e^{x^2} + 1} \\
 &= \frac{2xe^{x^2}}{e^{x^2} + 1}
 \end{aligned}$$

Simplify Roots/zeros

Example 2  $\{ w \in \mathbb{R}^2 \Rightarrow \text{multivariate loss function}$

$$\begin{cases} l(w) = -\cos(2\pi w^\top w) + 2w^\top w \\ \nabla l(w) = 4\pi \sin(2\pi w^\top w) w + 4w \end{cases}$$



- $\alpha = 10^{-3}$

- $\varepsilon = 10^{-3}$

- $w^0 = [0.85 \ 0.85]^\top \rightarrow \text{local minimum } \vdash$
- $w^0 = [-0.7 \ 0]^\top \rightarrow \text{global minimum } \ddot{\cup}$

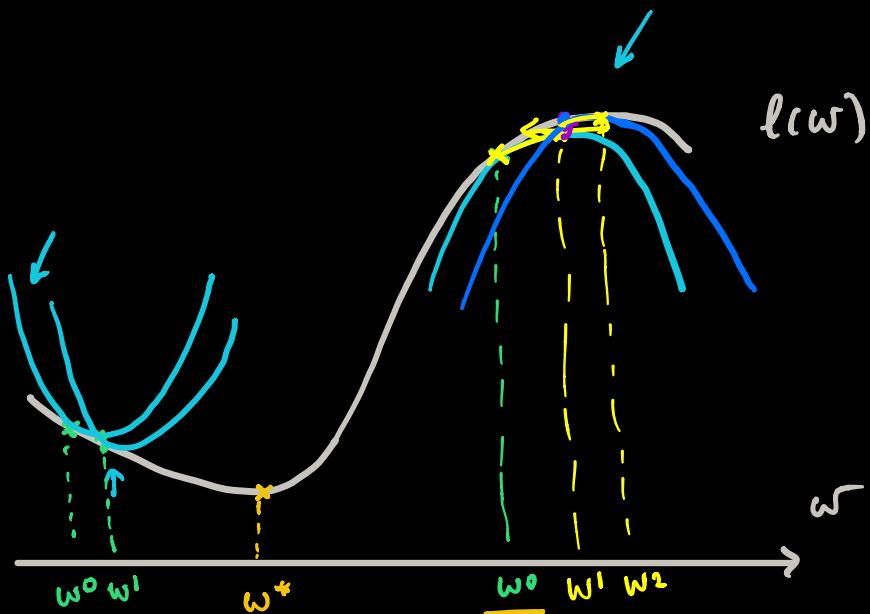
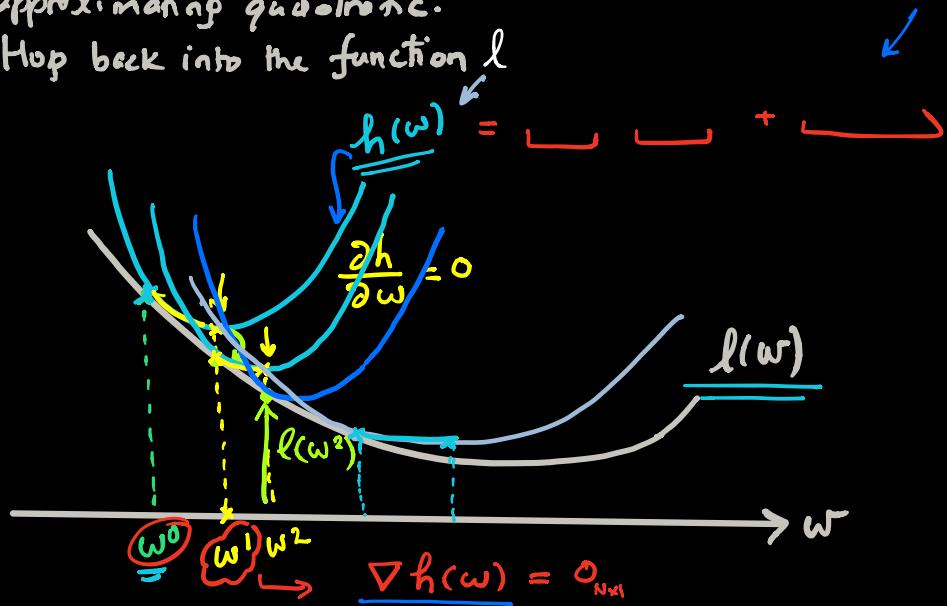
②

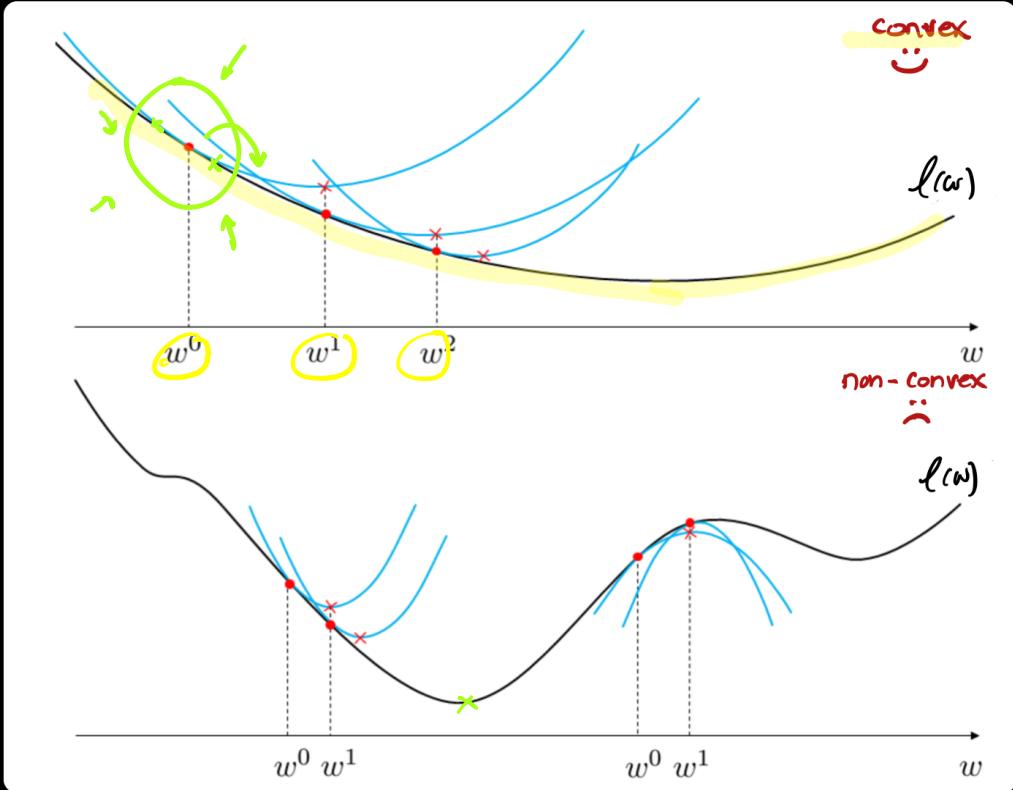
Basic technique 2 for finding  $\min_w \ell(w)$ :  
Newton's method

- Newton's method uses second-order (quadratic) approximation to locate the optimum  $w^*$  of  $\ell(w)$ .

### Steps :

- 1- Begin at an initial point  $w^0$
- 2- Travel to a stationary point (min/max) of the approximating quadratic.
- 3- Hop back into the function  $\ell$





\* Second-order Taylor approximation of  $\ell(\omega)$  around  $\omega^*$ :

$$h(\omega) = \ell(\omega^*) + \nabla \ell(\omega^*)^\top (\omega - \omega^*) + \frac{1}{2} (\omega - \omega^*)^\top \nabla^2 \ell(\omega^*) (\omega - \omega^*)$$

\* Compute approximating quadratic function:

① Set gradient to zero "0"

$$\nabla h(\omega) = \nabla \ell(\omega^*) + \nabla^2 \ell(\omega^*) (\omega - \omega^*) = \underline{\underline{0_{N \times 1}}}$$

$$\Rightarrow \nabla^2 \ell(\omega^*) \underline{\omega} = \nabla^2 \ell(\omega^*) \underline{\omega^*} - \nabla \ell(\omega^*) \in \mathbb{R}^N$$

$$\boxed{A \quad \underline{\omega} = b}$$

$\uparrow$   
 $\mathbb{R}^{N \times N}$

$\uparrow$   
 $\mathbb{R}^{N \times 1}$

$$\Rightarrow \text{Solve } N \text{ equations}$$

$$\omega^1 = \left\{ \begin{array}{c} \\ \vdots \end{array} \right.$$

② find next point  $\omega^1$

→  $\text{sol}^0$  is the point  $\omega^1$  traveled to by Newton's method.

③ repeat steps ① & ② until convergence

Solve for :

$$\nabla^2 \ell(\omega^1) \omega = \nabla^2 \ell(\omega^1) \omega^1 - \nabla \ell(\omega^1)$$

↓  
 $\text{sol}^0$  is  $\omega^2$

→ at  $k^{\text{th}}$  step (estimate  $\omega^k$ )

$$\nabla^2 \ell(\omega^{k-1}) \underbrace{\omega^k}_{\text{+}} = \nabla^2 \ell(\omega^{k-1}) \underbrace{\omega^{k-1}}_{\text{+}} - \underbrace{\nabla \ell(\omega^{k-1})}_{\text{O}}$$

- what happens when  $\nabla \ell(\omega^{k-1}) = \mathbf{0}_{N \times 1}$  ?

$$\omega^k = \omega^{k-1}$$

### Algorithm 2.2 Newton's method

**Input:** twice differentiable function  $\ell$ , and initial point  $w^0$

$k = 1$

Repeat until stopping condition is met:

Solve the system  $\nabla^2 \ell(w^{k-1}) w^k = \nabla \ell(w^{k-1}) w^{k-1} - \nabla \ell(w^{k-1})$  for  $w^k$ .

$k \leftarrow k + 1$

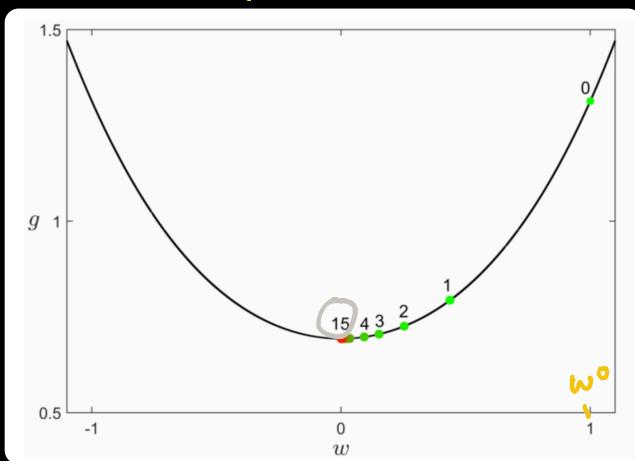
If  $\nabla^2 \ell(w^{k-1})$  is invertible?

$$\begin{cases} [\nabla^2 \ell(w^{k-1})]^{-1} \\ x x^{-1} = x \times \frac{1}{x} = 1 \end{cases}$$

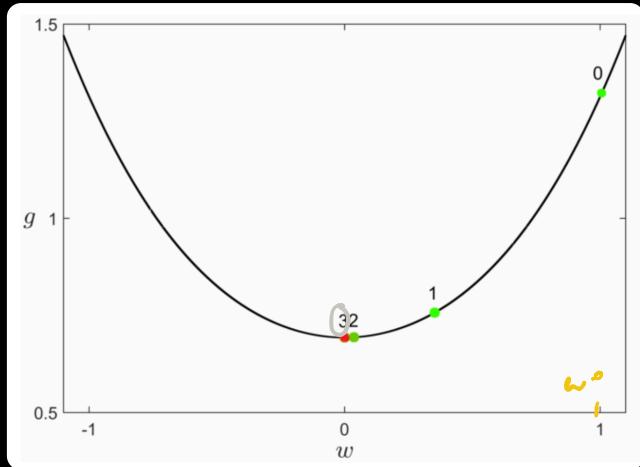
$$w^k = w^{k-1} - [\nabla^2 \ell(w^{k-1})]^{-1} \nabla \ell(w^{k-1})$$

$\alpha$

gradient descent  
(first order)



Newton method  
(second order)



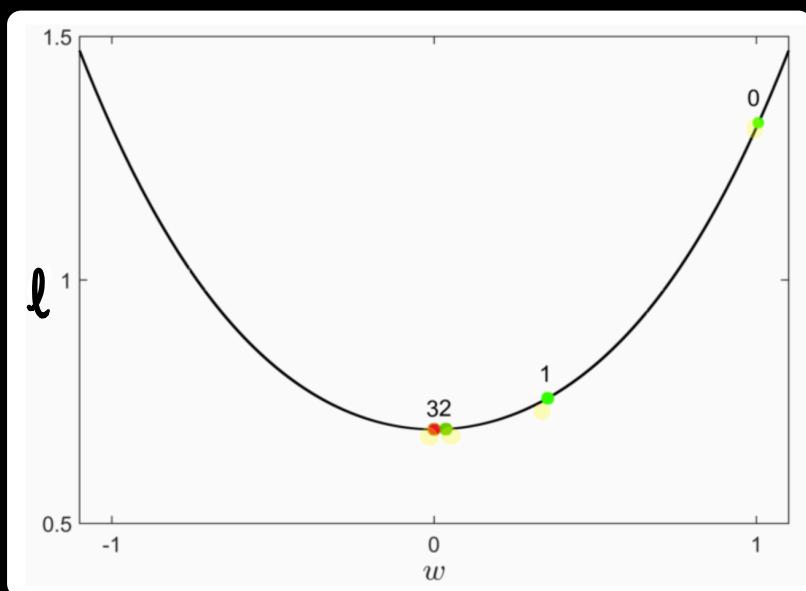
(+) if  $\ell$  is convex  $\Rightarrow$  faster  
 $\rightarrow \alpha$

$l_0^k \times l_0^K$

(-) non-convex  $\Rightarrow$  not converge  
\* inverting the Hessian matrix  
is computationally expensive

### Example 1

$$\left\{ \begin{array}{l} w \in \mathbb{R} \text{ (scalar) univariate loss function} \\ l(w) = \log(1 + e^{w^2}) \end{array} \right.$$



$$l''(w) = \frac{2e^{w^2}(2w^2 + e^{w^2} + 1)}{(1 + e^{w^2})^2}$$

$$\bullet w^0 = 1$$

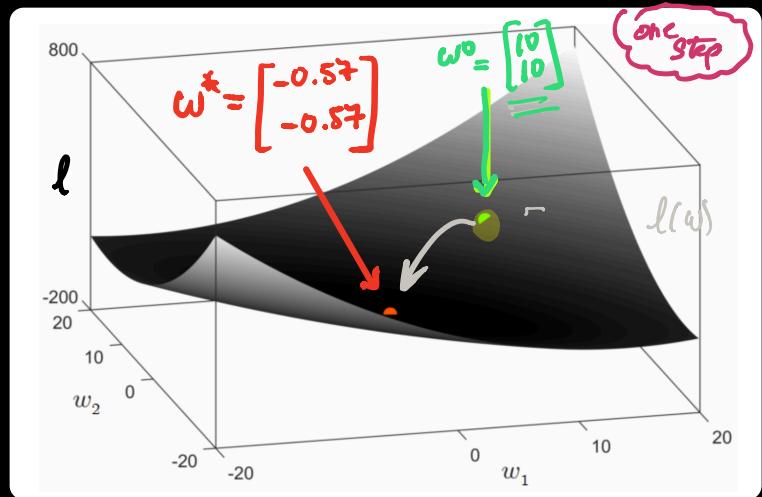
$$\bullet |l'(w)| < \varepsilon = 10^{-3} = 0.001$$

$$w^k = w^{k-1} - \underbrace{\frac{l'(w^{k-1})}{l''(w^{k-1})}}_{> 0} > 0$$

## Example 2

$$\ell(\omega) = \frac{1}{2} \omega^T Q \omega + r^T \omega + d$$

$\left\{ \begin{array}{l} Q = \begin{bmatrix} 1 & 0.75 \\ 0.75 & 1 \end{bmatrix}, r = [1 \ 1]^T, \text{ and } d = 0 \\ \omega^0 = [10 \ 10]^T \\ \cdot \nabla \ell(\omega) = Q\omega + r \\ \cdot \nabla^2 \ell(\omega) = Q \end{array} \right\}$



- Write the algebraic expression for  $\omega^1$  &  $\omega^2$ . What do you notice?

$$\Rightarrow \nabla^2 \ell(\omega^{k-1}) \omega^k = \nabla \ell(\omega^{k-1}) \omega^{k-1} - \nabla \ell(\omega^{k-1})$$

$\underbrace{\omega^k}_{Q \times \omega^k} = \underbrace{\omega^{k-1}}_{Q} - \underbrace{\nabla \ell(\omega^{k-1})}_{Q \omega^0 - r}$

$\underbrace{\omega^1}_{Q \omega^1 = -r} = \underbrace{\omega^0}_{Q^{-1}r}$

$\Rightarrow \boxed{\omega^1 = -Q^{-1}r}$        $\boxed{\omega^2 = -Q^{-1}r}$

I ❤️ M<sub>x</sub><sup>i</sup> [4]

### Mathematical notation

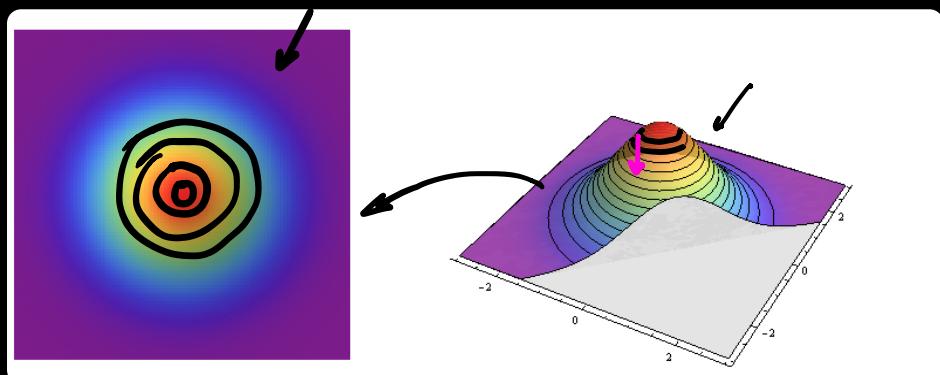
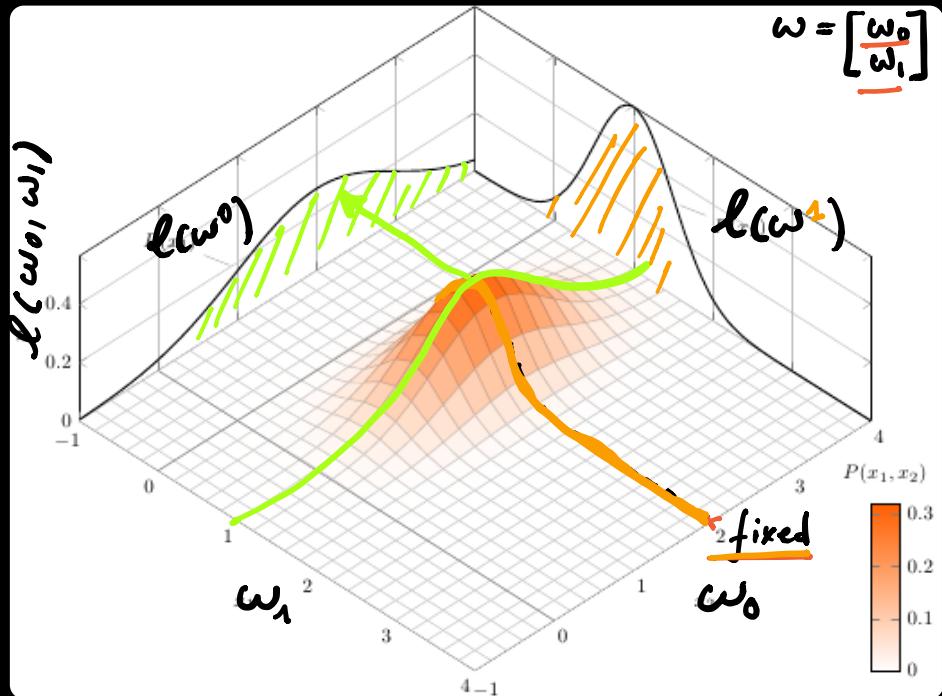
$$\begin{aligned} \mathbf{w}^k &= \mathbf{w}^{k-1} - \alpha_k \nabla l(\mathbf{w}^{k-1}) \\ \nabla^2 l(\mathbf{w}^{k-1}) \mathbf{w}^k &= \nabla^2 l(\mathbf{w}^{k-1}) \mathbf{w}^{k-1} - \nabla l(\mathbf{w}^{k-1}) \\ \mathbf{w}^k &= \mathbf{w}^{k-1} - [\nabla^2 l(\mathbf{w}^{k-1})]^{-1} \nabla l(\mathbf{w}^{k-1}) \end{aligned}$$

### Definition

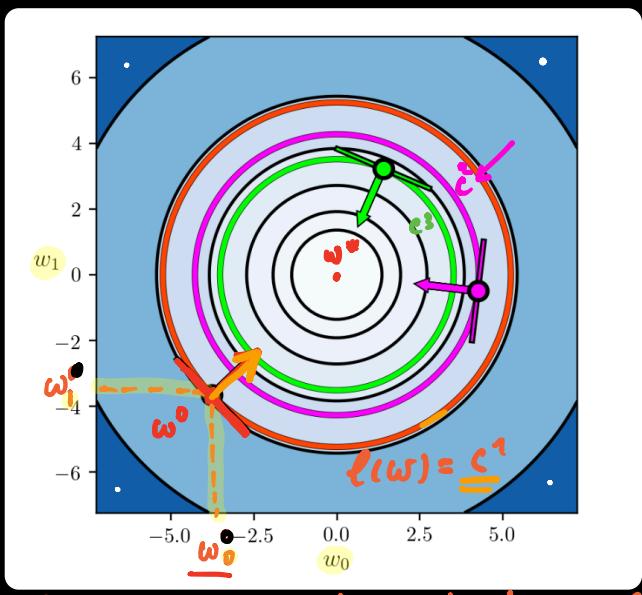
gradient descent for finding the optimal  $\mathbf{w}^*$   
 Newton's method update rule (second order optimization)  
 Newton's method update rule when the Hessian matrix  $\nabla^2 l(\mathbf{w}^{k-1})$  is invertible

[Matlab / Python demo]

MLR-exercise 2.14



$$l(\omega) = \omega_0^2 + \omega_1^2 + 2$$



**IMPORTANT PROPERTY  
OF THE NEGATIVE  
GRADIENT DIRECTION**

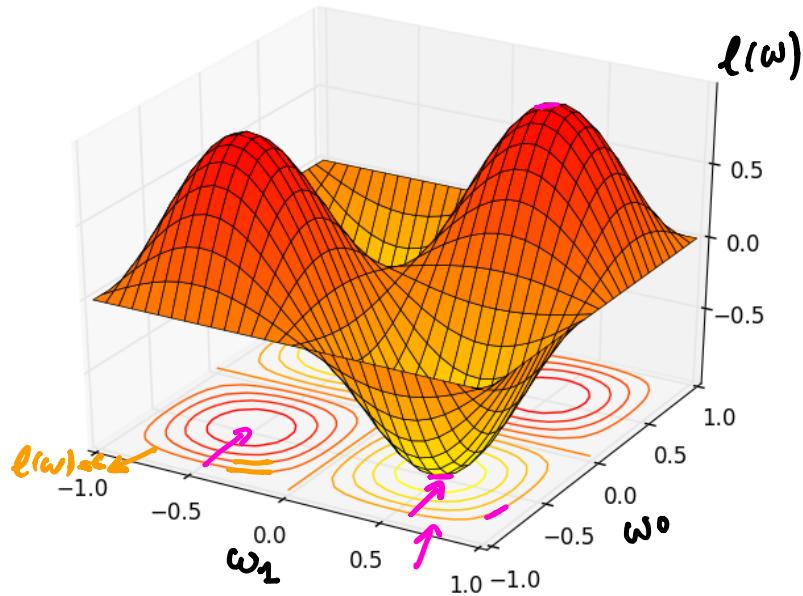
It is always perpendicular to the contours (level curves) of a function.

Universally true statement

↓  
It holds for any function and at all of its inputs.

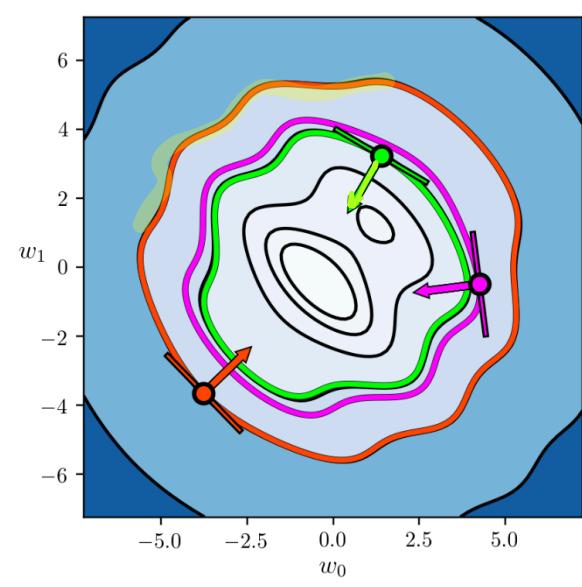
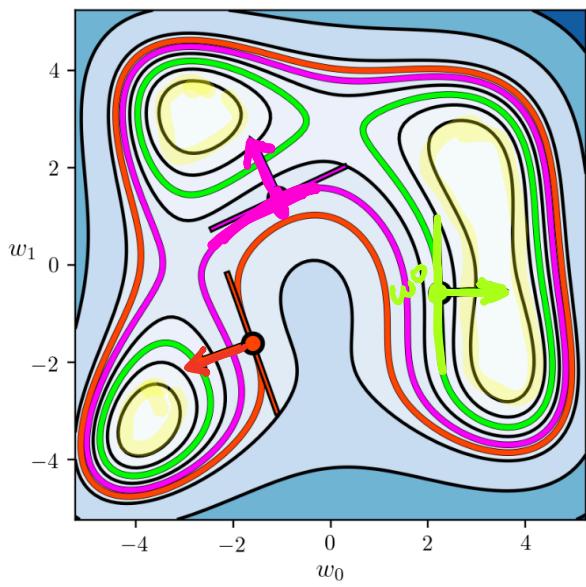
$\nabla l$  = normal vector to the level curve  $l(\omega) = C^1$

$$\omega = \begin{bmatrix} \omega_0 \\ \omega_1 \end{bmatrix}$$

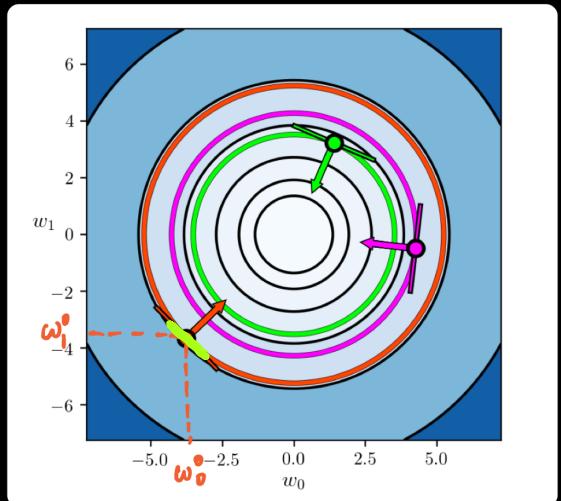


$$l(\omega) = (\omega_0^2 + \omega_1 - 1)^2 + (\omega_0 + \omega_1^2 - 6)^2$$

$$l(\omega) = \omega_0^2 + \omega_1^2 + 2 \sin(1.5(\omega_0 + \omega_1)^2) + 2$$



$$\ell(\omega) = \omega_0^2 + \omega_1^2 + 2$$



Let us confirm this fundamental property using rigorous mathematical reasoning!

- ①  $\ell(\omega)$  is differentiable at some input point  $a \Rightarrow a$  lies on the contour points where  $\ell(\omega) = \ell(a) = c$ .
- ② Let us take another point  $b$  close to  $a$  then essentially since:

$$\nabla \ell(a)^T (b - a) = 0$$

gradient is vector  
perpendicular to the  
( $a - b$ )

since the 1st order approximation  $h(b)$  around  $a$ :

$$h(b) = \ell(a) + \nabla \ell(a)^T (b - a)$$

same 0

$$\ell'(b) = \lim_{b \rightarrow a} \frac{\ell(b) - \ell(a)}{b - a}$$

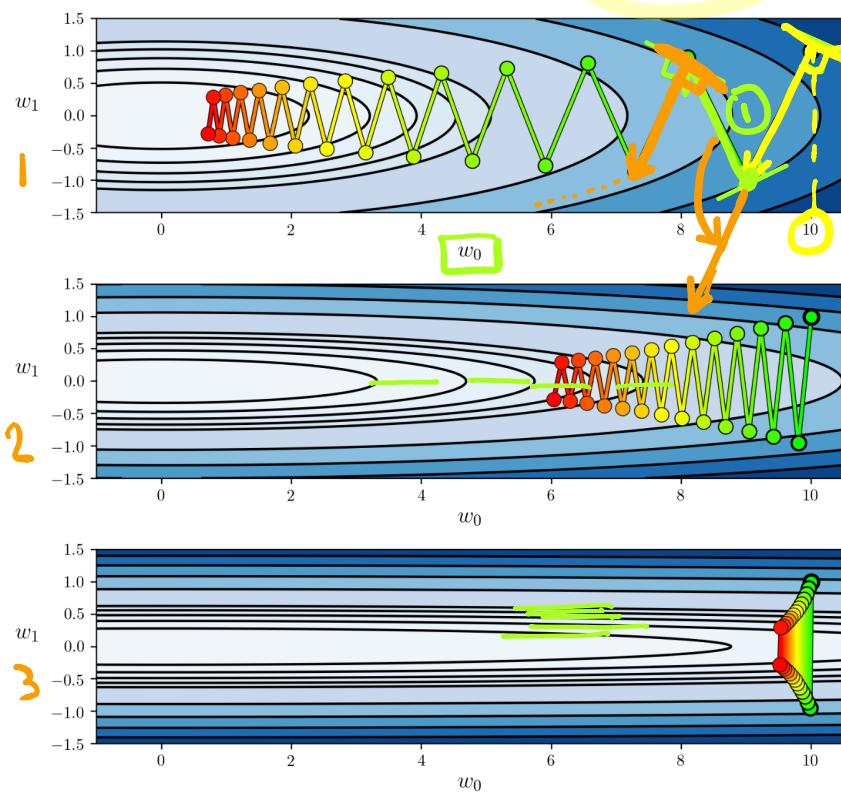
② The 'zig-zagging' behavior of gradient descent

:( It slows down convergence due to the oscillation of the negative gradient direction

Example :

$$\begin{cases} \ell(\omega) = a + b^T \omega + \omega^T C \omega ; \omega \in \mathbb{R}^2 \\ a = 0, b = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \end{cases}$$

$$\ell(\omega) = \omega^\top C \omega, \alpha = 0.1, \omega^0 = \begin{bmatrix} 10 \\ 1 \end{bmatrix}, 25 \text{ runs.}$$



$$C = \omega^0 \begin{bmatrix} 0.5 & 0 \\ 0 & 12 \end{bmatrix}$$

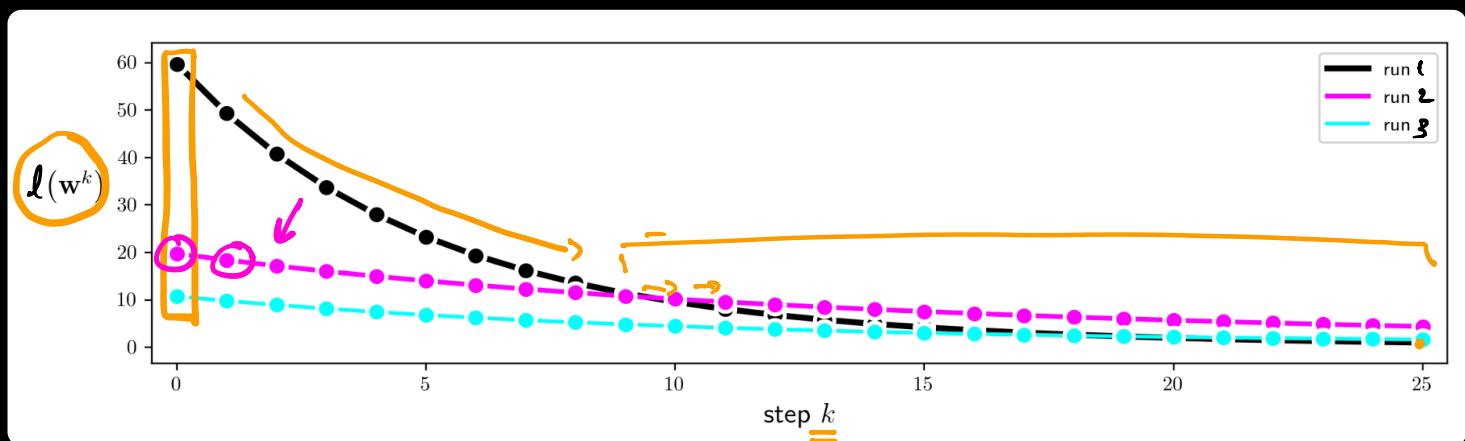
$$C = \begin{bmatrix} 0.1 & 0 \\ 0 & 12 \end{bmatrix}$$

$$C = \begin{bmatrix} 0.01 & 0 \\ 0 & 12 \end{bmatrix}$$

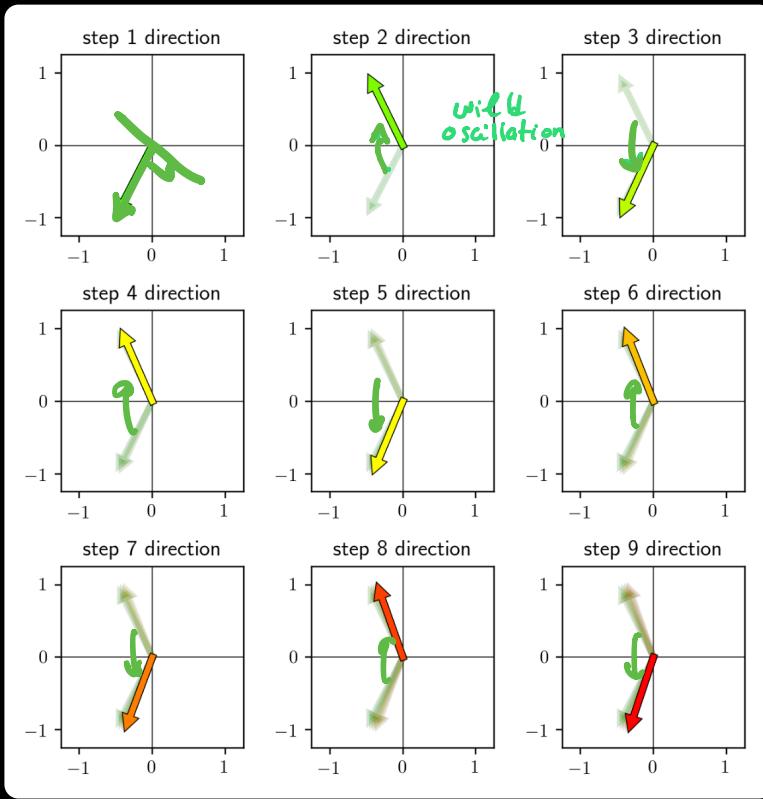
Variance along  $w_0$

- Shared global minimum at  $\omega^* = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$
- Zigzagging is caused by the negative gradient direction constantly pointing perpendicular to the contours of the loss function.

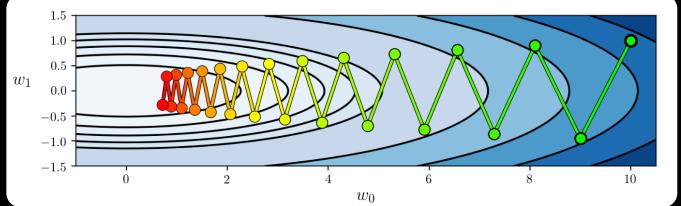
slow convergence caused by zigzagging



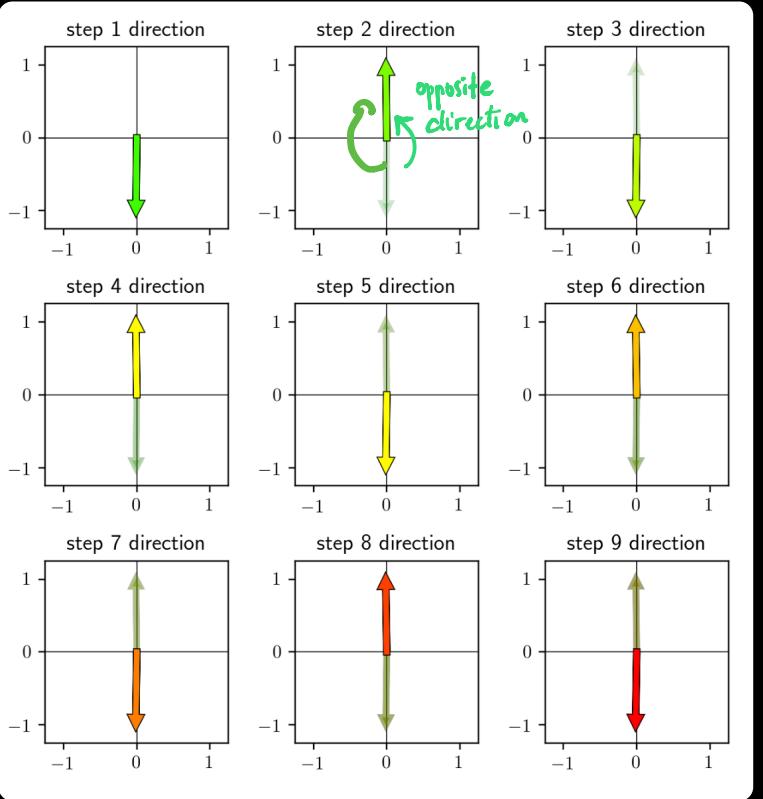
\* Plot the descent direction from the first 9 steps .



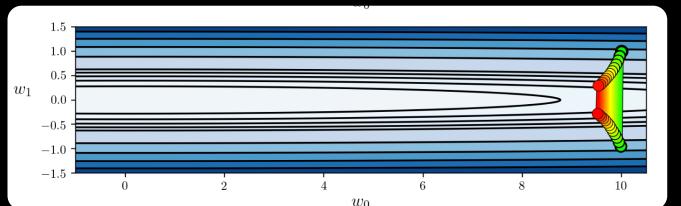
$$\ell(\omega) = \omega^T \begin{bmatrix} 0.5 & 0 \\ 0 & 12 \end{bmatrix} \omega$$



Massive zigzagging in the descent steps .

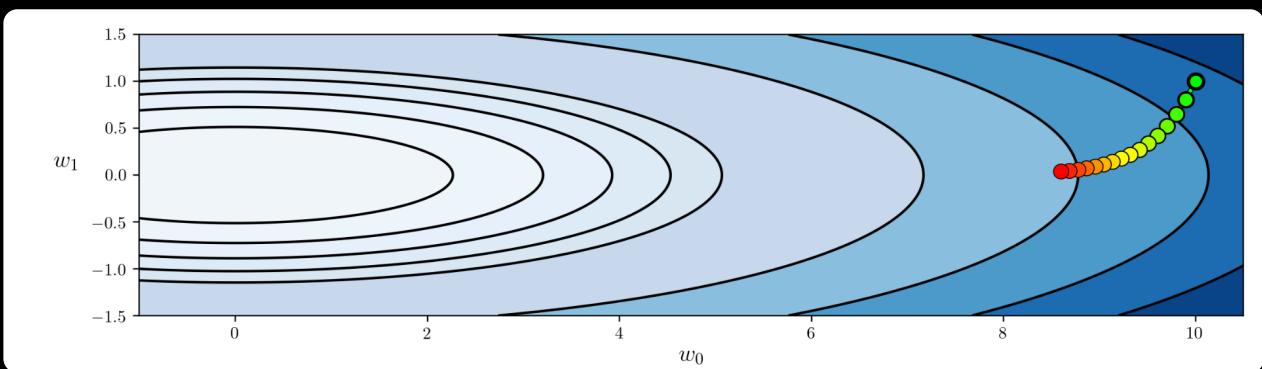


$$\ell(\omega) = \omega^T \begin{bmatrix} 0.01 & 0 \\ 0 & 12 \end{bmatrix} \omega$$

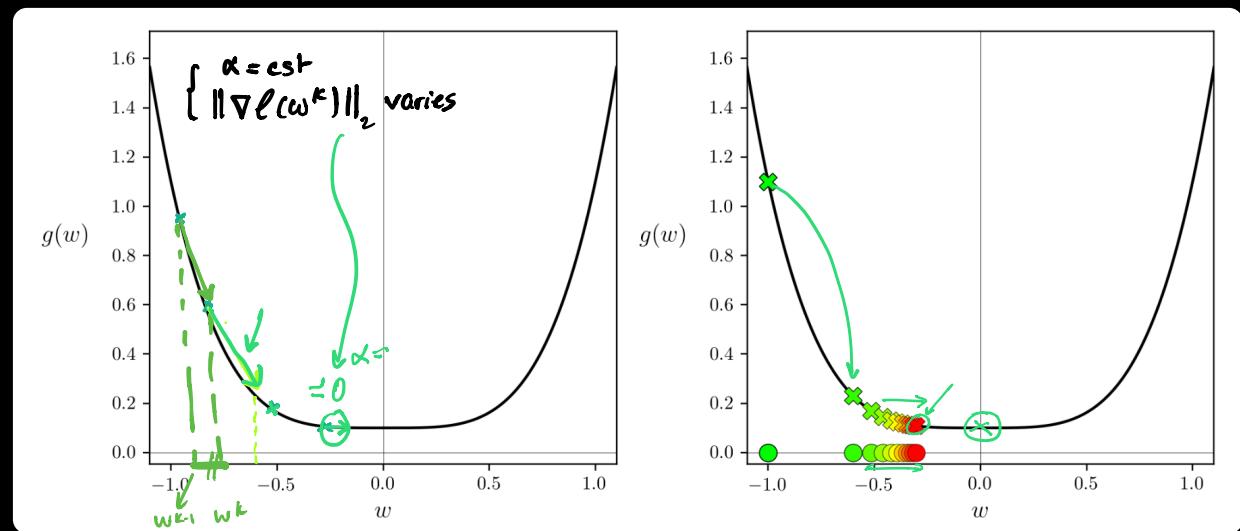


(+) Decreasing steplength value ( $\alpha$ ) can tone down the zigzagging behavior.

(-) Slow convergence



Issue 2 : the vanishing gradient  
“vanishing magnitude”

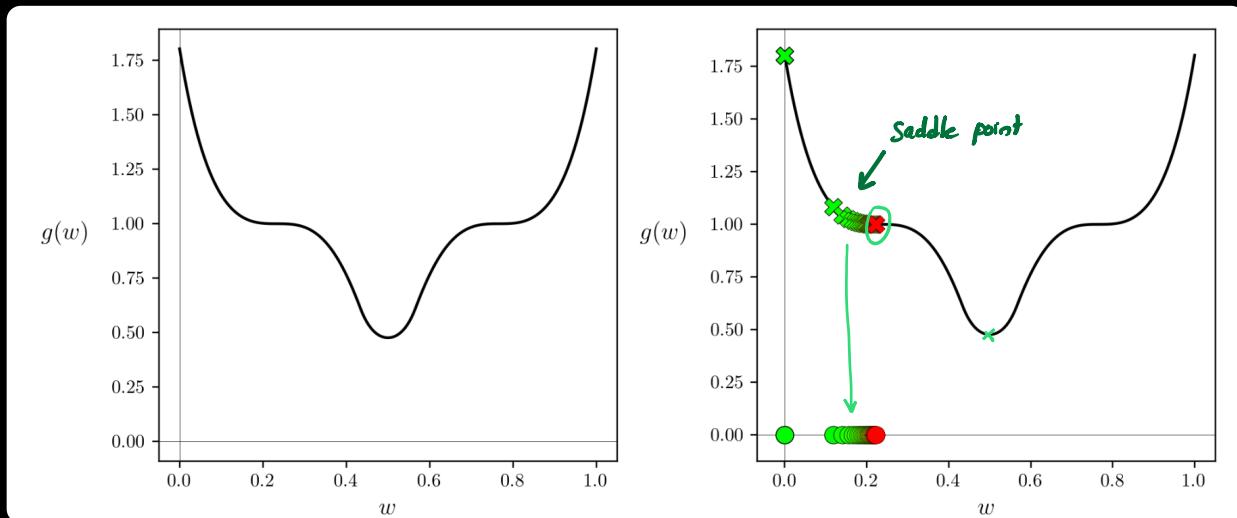


$$\omega^k = \omega^{k-1} - \alpha \nabla \ell(\omega^{k-1})$$

$$\Rightarrow \|\omega^{k-1} - \omega^k\|_2 = \underbrace{\alpha}_{\text{magnitude of the gradient}} \underbrace{\|\nabla \ell(\omega^{k-1})\|}_{} \quad \text{magnitude of the gradient}$$

The magnitude of the gradient vanishes at stationary points.

→ Slow-crawling behavior of gradient descent.



:( Initialization lies on a long narrow valley .  
No progress following 1000 steps of gradient descent!

