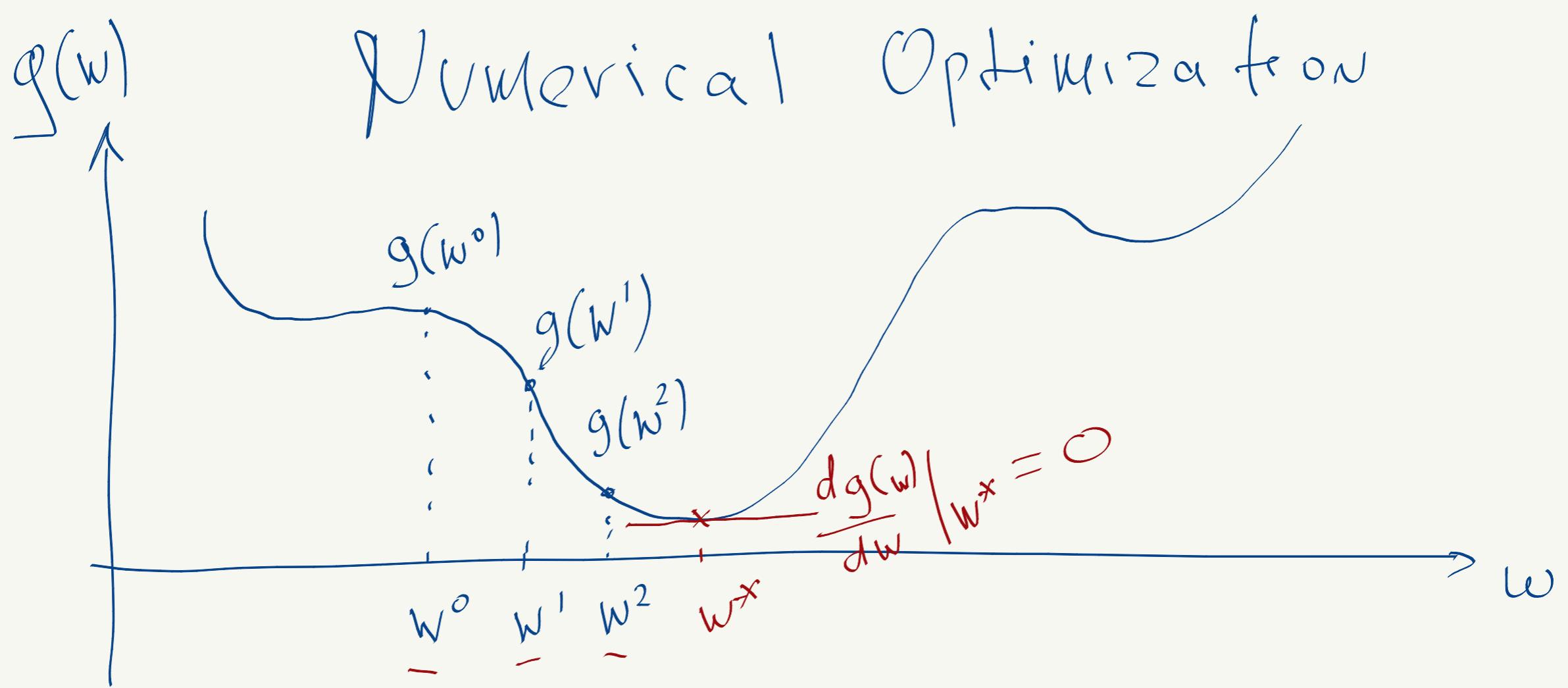
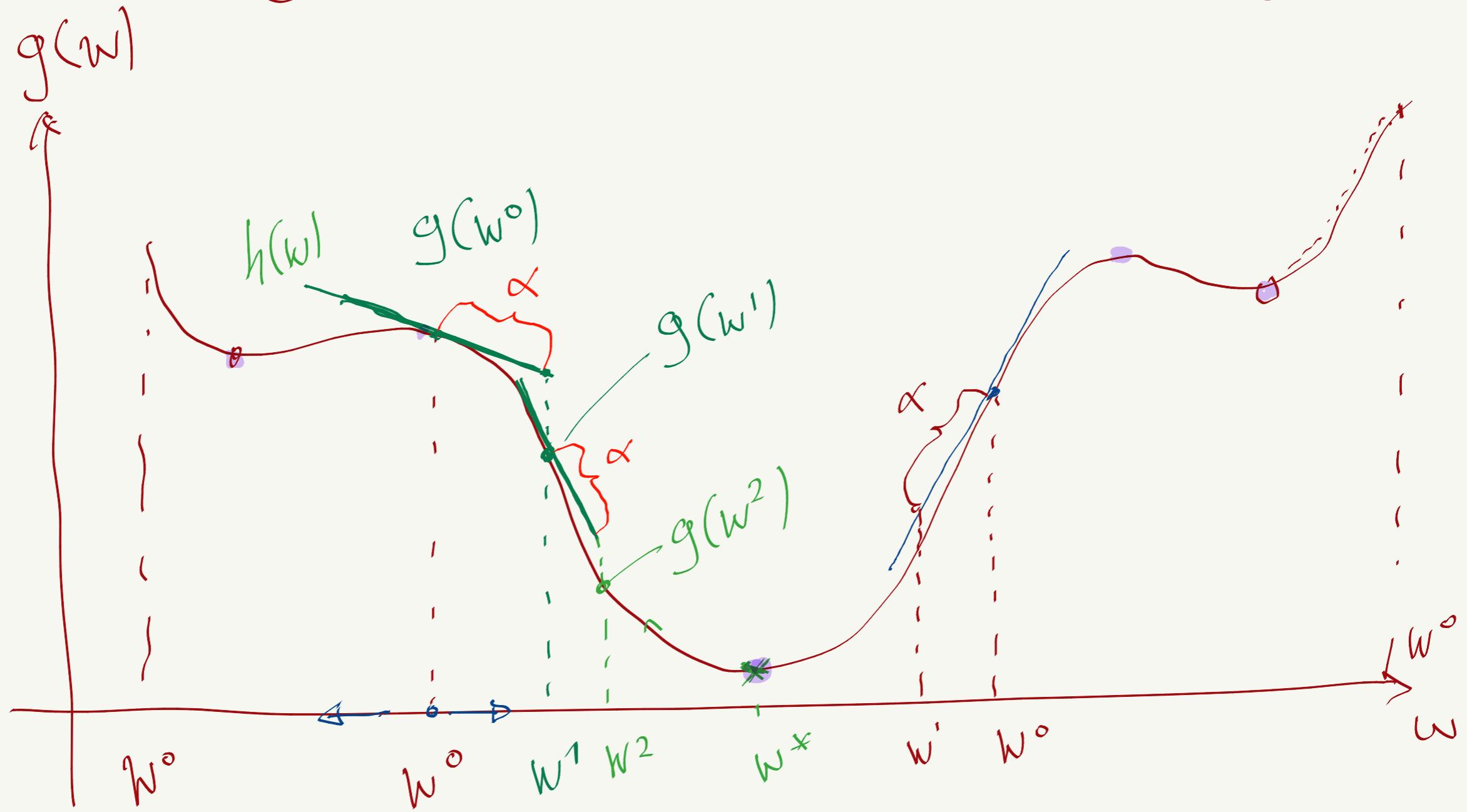


4/13/20  
Lecture 4



# 1st order optimization

## Gradient Descent (GD) algorithm



- Use a linear approximation to  $g$  to determine the next step.
- Linear approximation: 1st order Taylor approx.

$$\rightarrow h(\bar{w}) = \underline{g(\bar{w}^0)} + \nabla g(\bar{w}^0)^T (\bar{w} - \bar{w}^0)$$

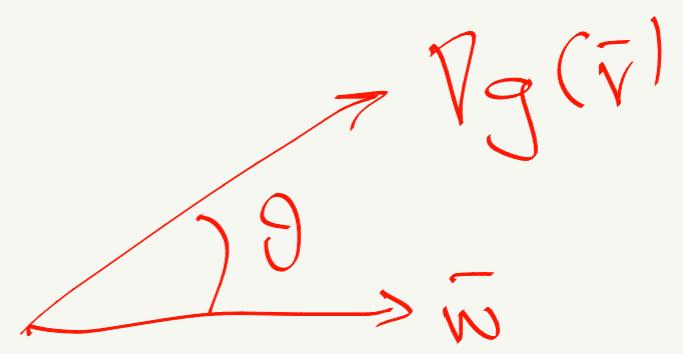
• Which direction is "downward" or "downhill"?

- △ accepts "trivial" answer for  $N=1$
- △  $N \geq 1$

$$h(\bar{w}) = \underline{\underline{g(\bar{w}^*)}} + \nabla g(\bar{w}^*)^T \bar{w} - \underline{\underline{\nabla g(\bar{w}^*)^T \bar{w}^*}}$$

do not  
depend  
on  $w$

$$\min_{\|w\|=1} \nabla g(\bar{v})^T \bar{w} = \min_{\|w\|=1} \|\nabla g(\bar{v})\| \cdot \|\bar{w}\| \cdot \cos \theta$$



clearly

$$-\nabla g(\bar{v}) \leftarrow \rightarrow w$$

when  $\theta = \pi$

$$\Rightarrow \bar{w} = \frac{-\nabla g(\bar{v})}{\|\nabla g(\bar{v})\|}$$

downhill direction

$$-\nabla g(\bar{w}^*)$$

# GD algorithm

Start  $\bar{w}^0$  step length or learning rate

$$\bar{w}^1 = \bar{w}^0 - \alpha \nabla g(\bar{w}^0)$$

$$\bar{w}^2 = \bar{w}^1 - \alpha \nabla g(\bar{w}^1)$$

.

$$\Rightarrow \boxed{\bar{w}^k = \bar{w}^{k-1} - \alpha \nabla g(\bar{w}^{k-1})}$$

$$\bar{w}^k - \bar{w}^{k-1} = -\alpha \nabla g(\bar{w}^{k-1})$$

for fixed  $\alpha$

more general form of GD

$$\boxed{\bar{w}^k = \bar{w}^{k-1} - \alpha^k \nabla g(\bar{w}^{k-1})}$$

$$\left. \begin{array}{l} \alpha^k = \frac{1}{k} \\ \alpha^k = \frac{k}{k+3} \end{array} \right\}$$

# When to stop iterating?

- ✓ ① When a stationary point is approximately reached

$$\|\nabla g(\bar{w})\| < \varepsilon \quad \leftarrow \text{small number}$$

✓ ②  $\frac{\|\bar{w}^{k+1} - \bar{w}^k\|}{\|\bar{w}^k\|} < \delta$

- 3. after a preset # of iterations

$$\bar{w}^{k+1} = \bar{w}^k - \alpha \nabla g(\bar{w}^k)$$

$$\left\{ \text{for } \nabla g(\bar{w}^k) = 0 \Rightarrow \bar{w}^{k+1} = \bar{w}^k \right.$$

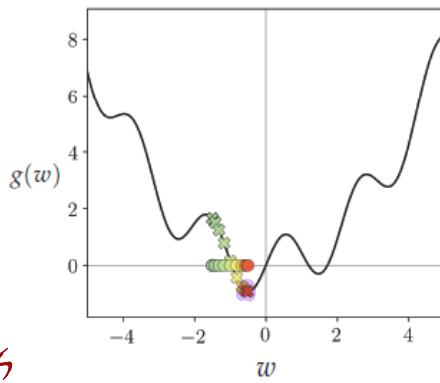
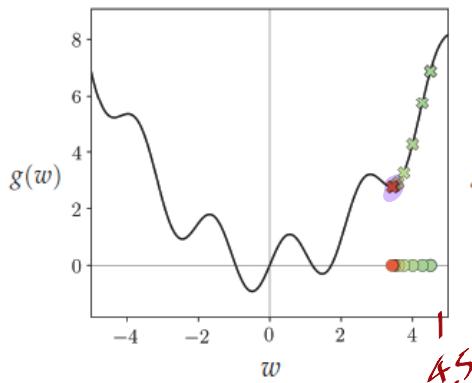
$$\bar{w}^{k+1} - \bar{w}^k = -\alpha \nabla g(\bar{w}^k)$$

$$\Rightarrow \underbrace{\|\bar{w}^{k+1} - \bar{w}^k\|}_{\substack{\text{small} \\ \text{large}}} = \alpha \underbrace{\|\nabla g(\bar{w}^k)\|}_{\substack{\text{small} \\ \text{large}}}$$

$$g(w) = \sin(3w) + 0.3w^2$$

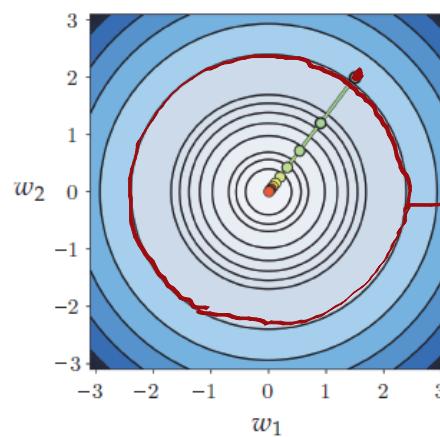
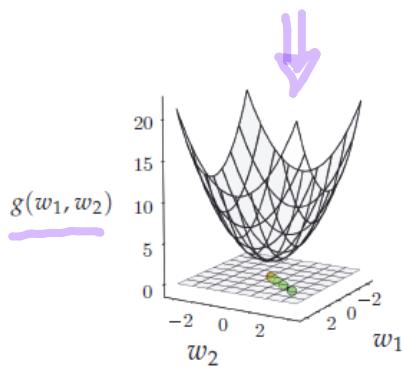
$w^* = 4.5$

$\alpha = 0.05$



3D plot

$\alpha = 0.1$



$w^* = -1.5$

$\alpha = 0.05$

"from above"

ISO-value curves

Figure 3.7 (top panels) Figure associated with Example 3.7. (bottom panels) Figure associated with Example 3.8. See text for details.

$$g(\bar{w}) = w_1^2 + w_2^2 + 2$$

$$\underline{g(w) = w^2}, \quad \underline{w^o = -2.5}$$

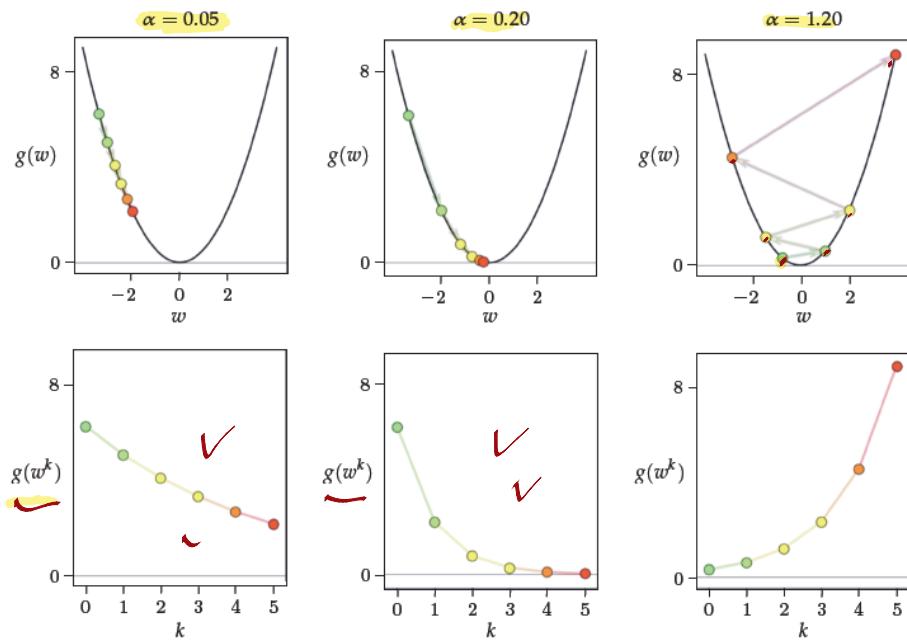


Figure 3.8 Figure associated with Example 3.9. See text for details.

Slow

ascend !

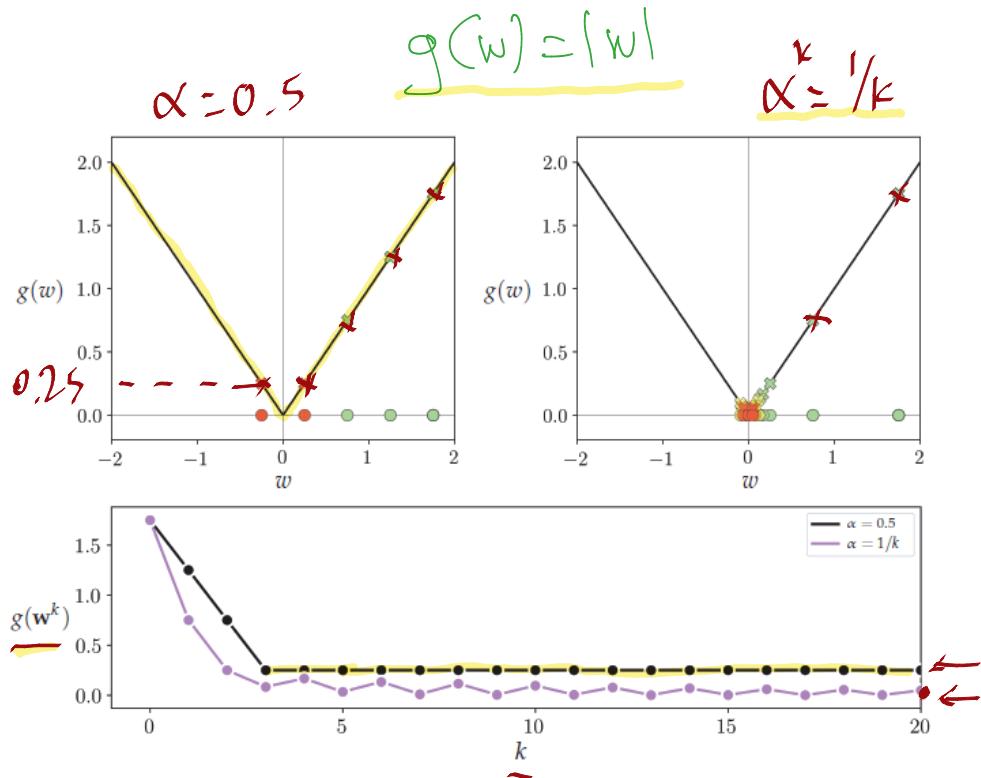


Figure 3.9 Figure associated with Example 3.10. See text for details.

20 iterations,  $w^0 = 2$

$$\frac{d}{dw} g(w) = \begin{cases} +1, & w > 0 \\ -1, & w < 0 \end{cases}$$

[Homework #3.6]

$$\Rightarrow g(\bar{w}) = w_1^2 + w_2^2 + 2 \sin(1.5(w_1 + w_2))^2 + 2$$

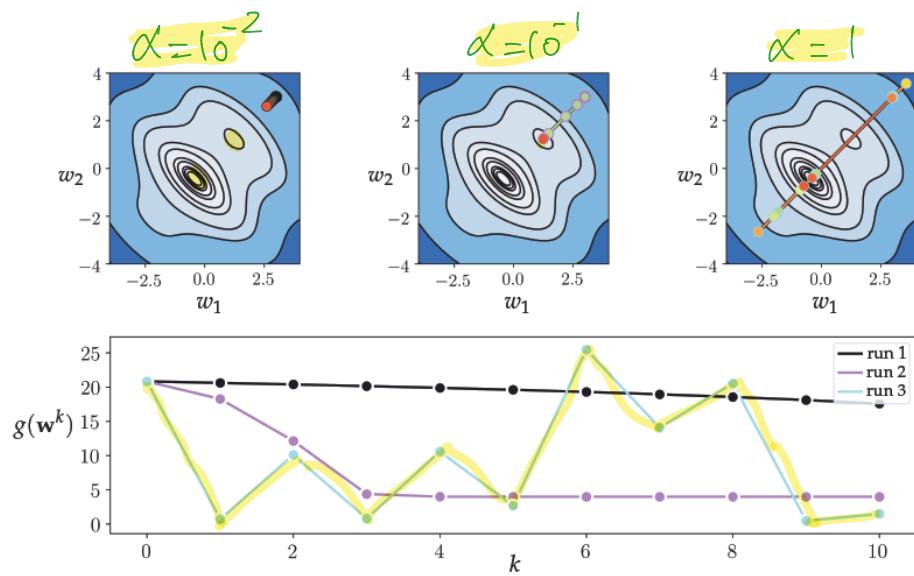


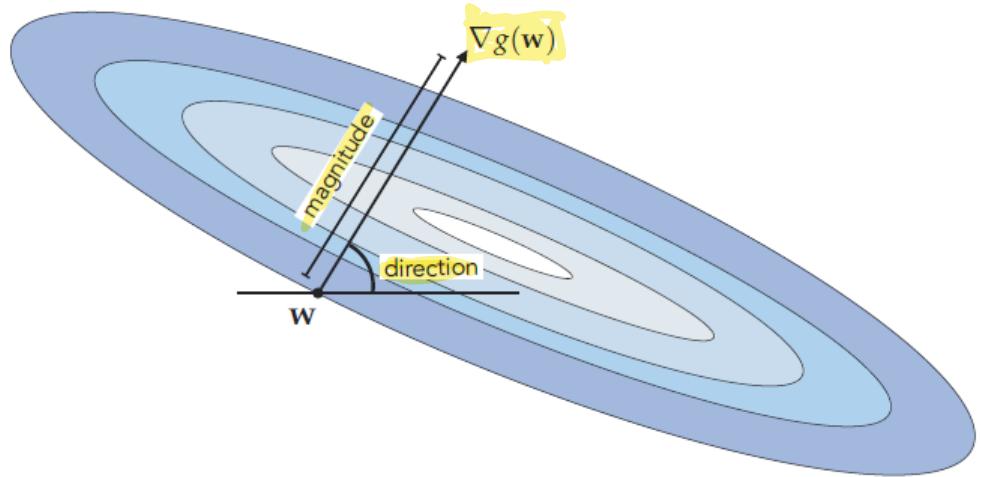
Figure 3.10 Figure associated with Example 3.11. See text for details.

$\rightarrow$  local min  $\sim [1.5 \ 1.5]^T$   
 $\rightarrow$  global "  $\sim [-0.5 \ -0.5]^T$   
 $\bar{w}^0 = [3 \ 3]^T$ , 10 steps

$\rightarrow \alpha = 10^{-2}$  : too small

$\rightarrow \alpha = 10^{-1}$  : local min

$\rightarrow \alpha = 1$  : ✓

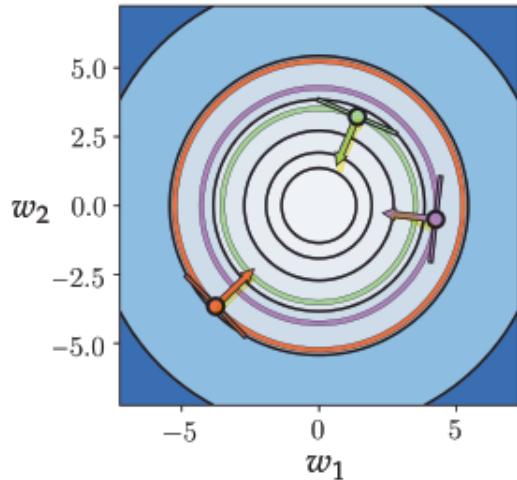


**Figure 3.11** The gradient vector of any arbitrary function at any point consists of a *magnitude* and a *direction*.

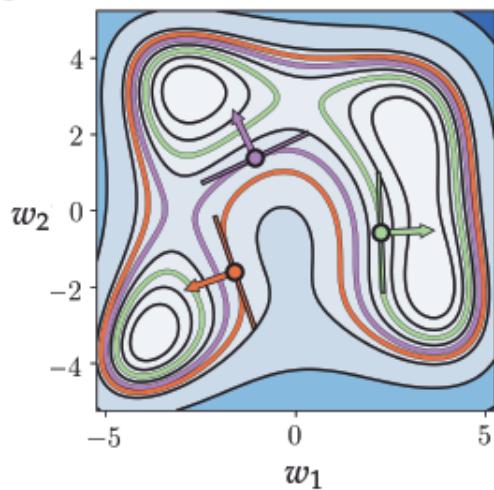
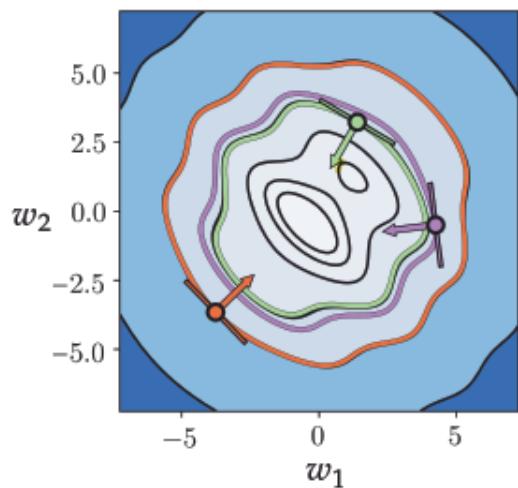
direction: zig-zagging

magnitude: slow-crawl

$$g(\bar{w}) = w_1^2 + w_2^2 + 2$$



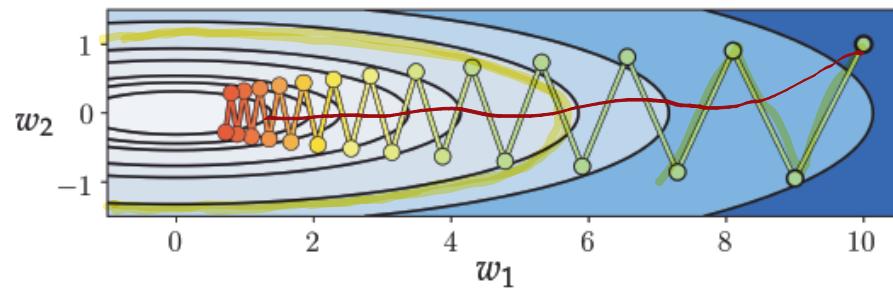
$$g(\bar{w}) = w_1^2 + w_2^2 + 2 \sin(1.5(w_1 + w_2)^2) + 2$$



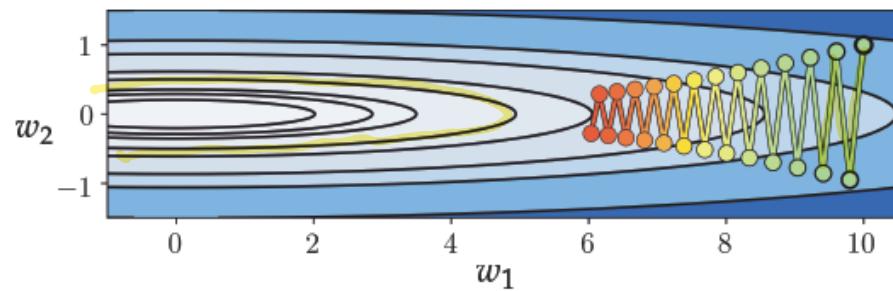
$$g(\bar{w}) = (w_1^2 + w_2^2 - 11)^2 + (w_1 + w_2^2 - 6)^2$$

**Figure 3.12** Figure associated with Example 3.12. Regardless of the function the negative gradient direction is always *perpendicular* to the function's contours. See text for further details.

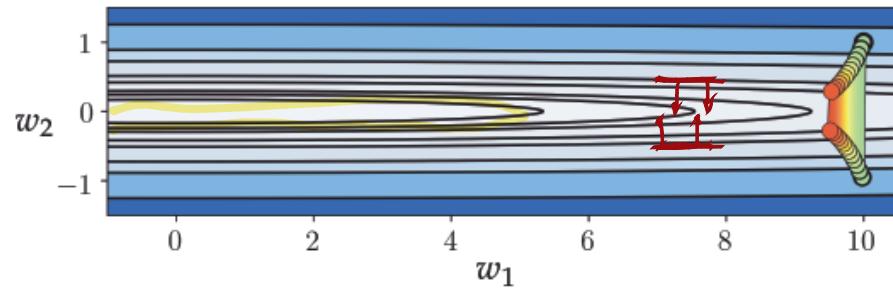
$$g(\bar{w}) = \bar{w}^\top C \bar{w}$$



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



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

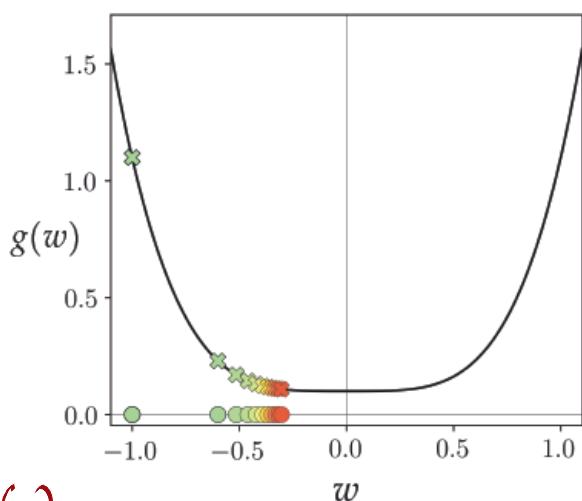


$$C = \begin{bmatrix} 0.01 & 0 \\ 0 & 12 \end{bmatrix}$$

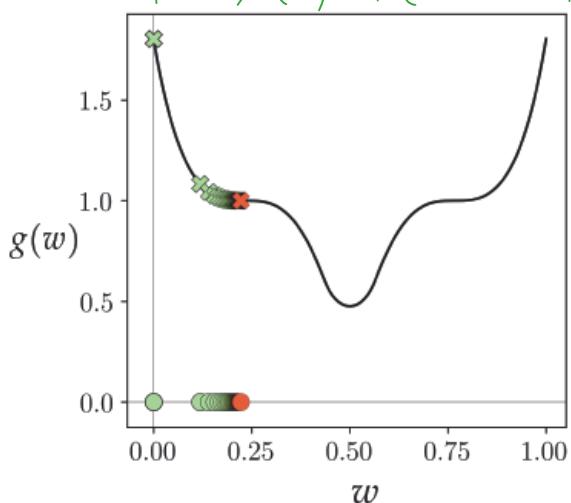
**Figure 3.13** Figure associated with Example 3.13, illustrating the zig-zagging behavior of gradient descent. See text for further details.

"Solution": momentum-accelerated GD (A.2)

$$\underline{g(w) = w^4 + 0.1}$$



$$\underline{g(w) = \max^2(0, 1 + (3w - 2.3)^3)} \\ + \underline{\max^2(0, 1 + (-3w + 0.7)^3)}$$



$\frac{-\alpha \nabla g(\cdot)}{\approx 0}$

left:  $\underline{\alpha = 10^{-1}}$   
10 steps

right:  $\underline{\alpha = 10^{-2}}$   
50 steps  
min at  $w = 1/2$   
saddle points at:  
 $w = 7/30$   
 $w = 23/30$

"solution": normalized GD (A3, A.4)

Adam optimizer (2 momentum approaches)