

# Gradient and Optimization

Digression

# **DERIVATIVES AND GRADIENT**

# Derivatives

- Some derivation rules:

$(w^2)' = 2w$	$(f(w)^2)' = 2f(w) \cdot f'(w)$	$c' = 0$
$(w^a)' = aw^{a-1}$	$(f(w)^a)' = a[f(w)]^{a-1} \cdot f'(w)$	$(cw)' = c$
$(e^w)' = e^w$	$(e^{f(w)})' = e^{f(w)} \cdot f'(w)$	$(f(w) + g(w))' = f'(w) + g'(w)$
$(\ln(w))' = \frac{1}{w}$	$(\ln(f(w)))' = \frac{1}{f(w)} \cdot f'(w)$	

- If we are supplied a value for  $w$ , **say 5**, then the above become numbers.
  - We say, we obtain the derivative “**at point 5**”

# Partial Derivatives

- Suppose now we have a function of multiple variables, e.g.

$$f(w_1, w_2, w_3) = (w_1 w_2 w_3)^2$$

- It can also written as  $f(\mathbf{w})$ , where  $\mathbf{w}$  is  $[w_1, w_2, w_3]$
- This function has three partial derivatives:
  - $f'_{w_1}$  obtained by considering only  $w_1$  variable and  $w_2, w_3$  constant
  - $f'_{w_2}$  obtained by considering only  $w_2$  variable and  $w_1, w_3$  constant
  - $f'_{w_3}$  obtained by considering only  $w_3$  variable and  $w_1, w_2$  constant

$$f'_{w_1}(w_1, w_2, w_3) = 2(w_1 w_2 w_3)w_2 w_3 = 2w_1 w_2^2 w_3^2$$

$$f'_{w_2}(w_1, w_2, w_3) = 2(w_1 w_2 w_3)w_1 w_3 = 2w_1^2 w_2 w_3^2$$

$$f'_{w_3}(w_1, w_2, w_3) = 2(w_1 w_2 w_3)w_1 w_2 = 2w_1^2 w_2^2 w_3$$

Derivation rules are the same as those for a single variable.

# Other Notation

$$f'_{w_1} \quad \text{also denoted by} \quad \frac{\partial f}{\partial w_1}$$

$$f'_{w_2} \quad \text{also denoted by} \quad \frac{\partial f}{\partial w_2}$$

$$f'_{w_3} \quad \text{also denoted by} \quad \frac{\partial f}{\partial w_3}$$

# Gradient

- The gradient is the vector of partial derivatives.

$$\nabla_f(\mathbf{w}) = [2w_1w_2^2w_3^2, \quad 2w_1^2w_2w_3^2, \quad 2w_1^2w_2^2w_3]$$

- Now suppose we want to compute the gradient **at a point**, say  $\mathbf{w}=(1,2,3)$ , i.e.  $w_1=1$ ,  $w_2=2$ ,  $w_3=3$

$$\nabla_f(\mathbf{w}) = \nabla_f([1,2,3]) = [2 \cdot 1 \cdot 2^2 \cdot 3^2, \quad 2 \cdot 1^2 \cdot 2 \cdot 3^2, \quad 2 \cdot 1^2 \cdot 2^2 \cdot 3] = [81, \quad 36, \quad 24]$$

**OPTIMIZATION**

# Minimization Problem

$$\min_{\mathbf{w}} f(\mathbf{w})$$



# Iterative Method

- Start at some  $\mathbf{w}_0$ ; take a step down the **steepest slope**
- Fixed step size:

$$\mathbf{w} \leftarrow \mathbf{w} - \kappa \mathbf{v}$$

- $\mathbf{v}$  is a **vector** in the direction of the **steepest slope**.
  - **Steepest slope** at some point? – The gradient vector at that point.
  - E.g.  $f(w_1, w_2, w_3) = (w_1 w_2 w_3)^2$

$$\nabla_f(\mathbf{w}) = [2w_1 w_2^2 w_3^2, \quad 2w_1^2 w_2 w_3^2, \quad 2w_1^2 w_2^2 w_3]$$

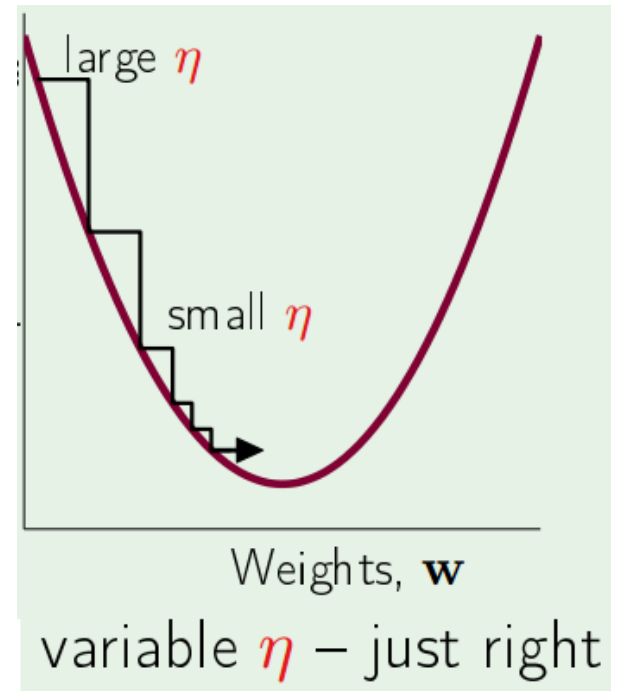
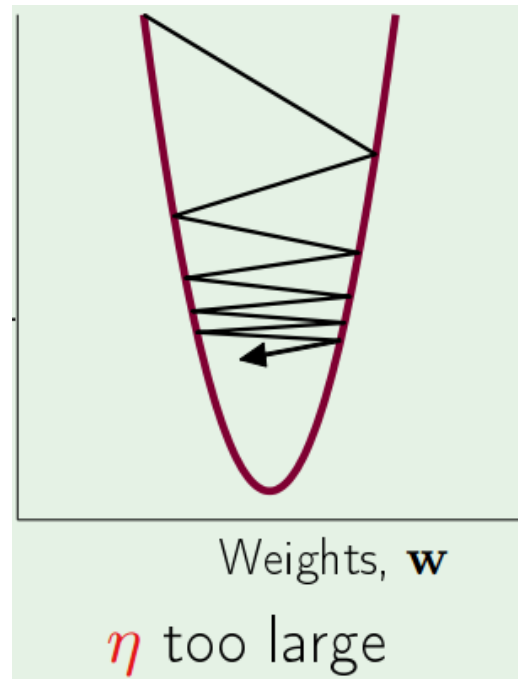
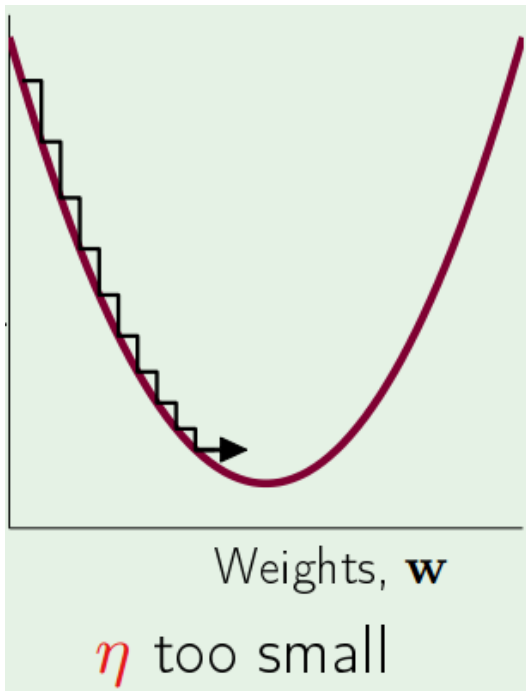
At point  $\mathbf{w}=(1,2,3)$ , i.e.  $w_1=1, w_2=2, w_3=3$ , we have

$$\begin{aligned}\nabla_f(\mathbf{w}) &= \nabla_f([1,2,3]) = [2 \cdot 1 \cdot 2^2 \cdot 3^2, \quad 2 \cdot 1^2 \cdot 2 \cdot 3^2, \quad 2 \cdot 1^2 \cdot 2^2 \cdot 3] \\ &= [81, \quad 36, \quad 24]\end{aligned}$$

$$\mathbf{w} \leftarrow [1,2,3] - 0.0001 * [81,36,24]$$

# Step Size (kappa, often denoted eta)

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{v}$$



# Gradient Descent Algorithm

Initialize  $\mathbf{w}=\mathbf{0}$

For  $t=0,1,2,\dots$  do

    Compute the gradient  $\nabla_f(\mathbf{w})$

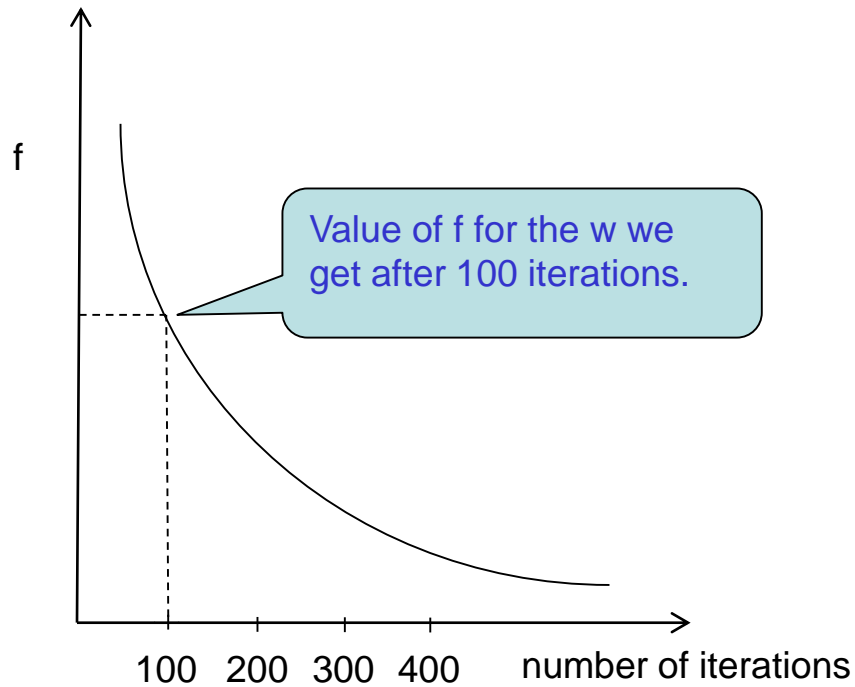
    Update the weights  $\mathbf{w} \leftarrow \mathbf{w} - \kappa \nabla_f(\mathbf{w})$

    Iterate with the next step until  $\mathbf{w}$  doesn't change too much  
    (or for a fixed number of iterations)

Return final  $\mathbf{w}$ .

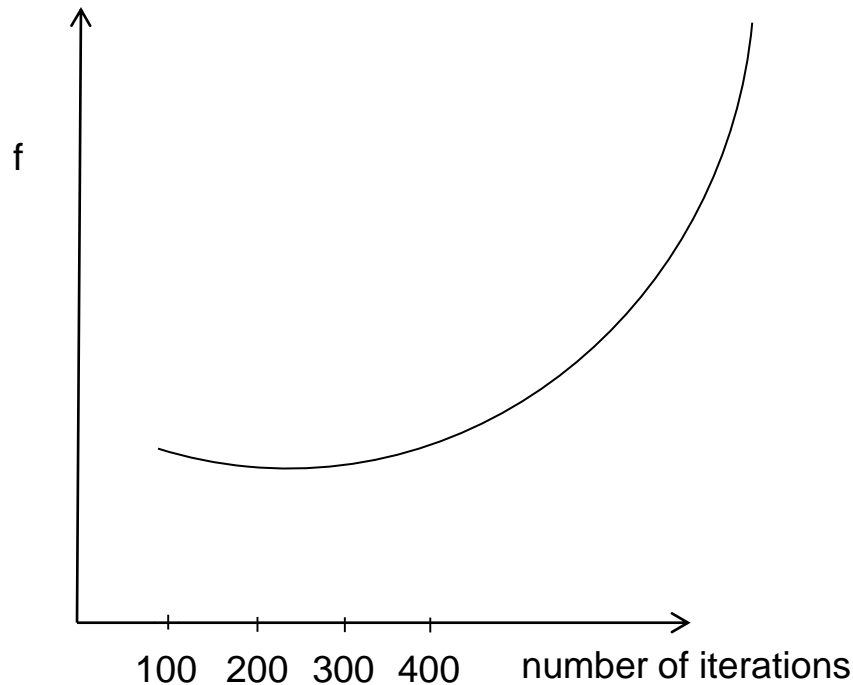
**HOW DO WE DETERMINE  $\kappa$ ?**

# $f(\mathbf{w})$ during iterations



A picture like this tells us that the gradient descent is working fine.

# $f(\mathbf{w})$ during iterations



A picture like this tells us that the gradient descent is NOT working fine.

We should use smaller  $\kappa$

If  $\kappa$  is small enough, a convex  $f(\mathbf{w})$  should decrease on every iteration. However, if  $\kappa$  is too small, it will take a long time to converge.

# Practically

Try

$\kappa=0.001$

$\kappa=0.01$

$\kappa=0.1$

$\kappa=1$

Plot or see  $f(\mathbf{w})$  for each one. If it decreasing with a reasonable speed, choose that value for  $\kappa$ .