$J(w)$

$\eta$ small



$w^{(0)}$ $w^{(1)}$

$\omega$

$\overbrace{\phantom{\nabla J(w^{(t)})}}^{\Delta w^{(t)}}$

Weight
update
Rule
$:$
$$w^{(t+1)} \leftarrow w^{(t)} - \eta \cdot \nabla J(w^{(t)})$$

$\eta \equiv$ Learning Rate (global variable)

Delta
Correction
Rule
$:$
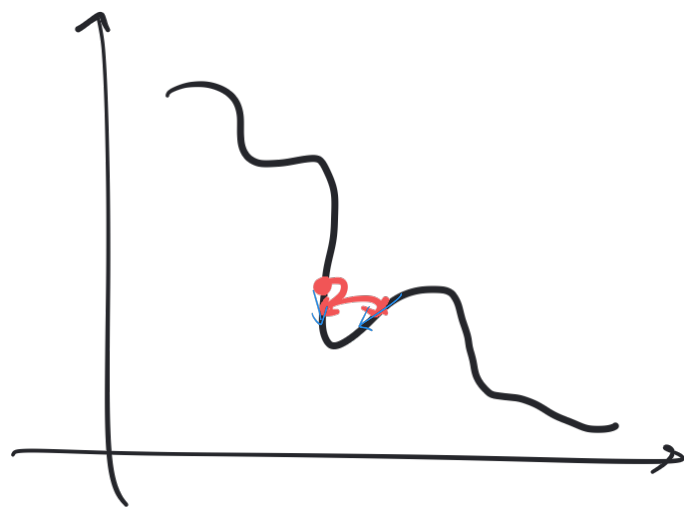$$\Delta w^{(t)} = -\eta \cdot \nabla J(w^{(t)})$$

$$\therefore \quad w^{(t+1)} \leftarrow w^{(t)} + \Delta w^{(t)}$$

# Momentum

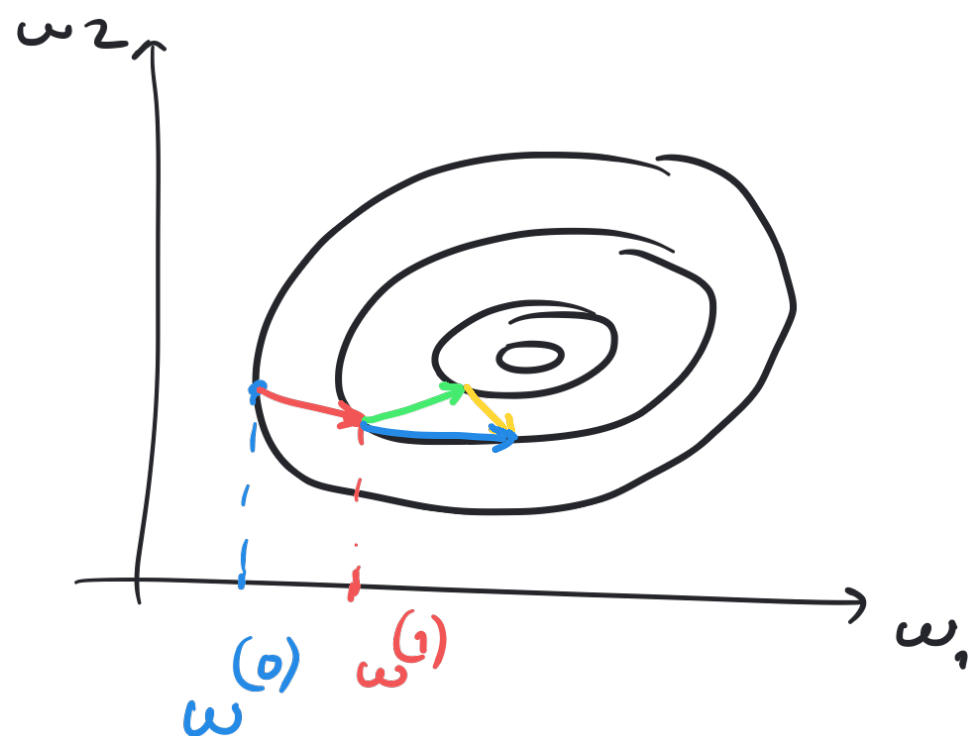Accelerate learning by adding information about previous gradient

$$\Delta w^{(t)} = -\gamma \cdot \nabla J(w^{(t)}) + \alpha \Delta w^{(t-1)}$$

$$(\alpha = 0.9)$$

→ Momentum may avoid local optima

→ it deaccelerates as the gradient sign changes.

Contours of objective function

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$w^{(1)} = w^{(0)} - \eta \cdot \nabla J(w^{(0)})$$

$$\longrightarrow -\eta \nabla J(w^{(1)})$$
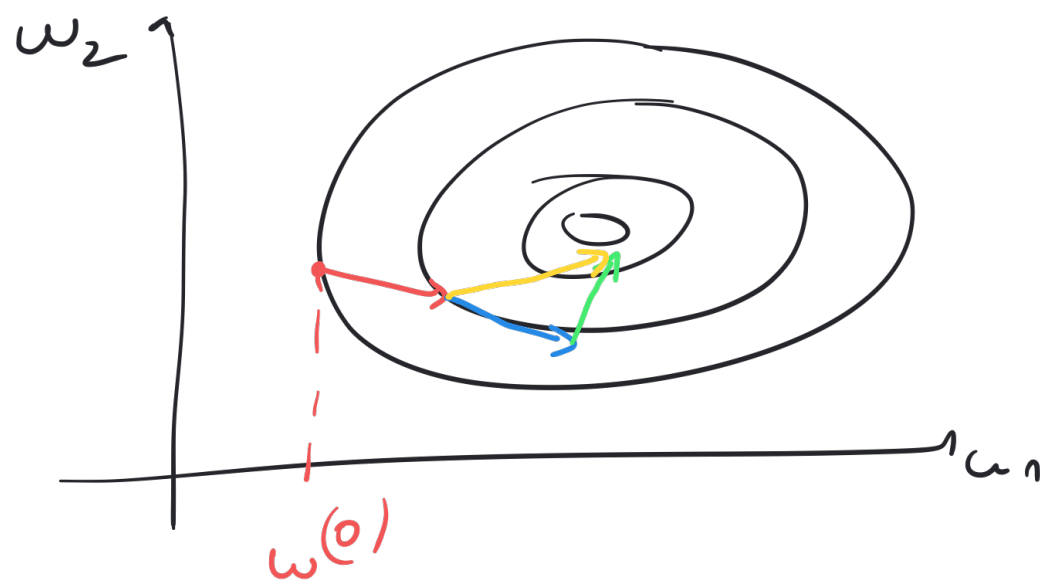
$$\longrightarrow \alpha \cdot \Delta w^{(1)}$$

$$\longrightarrow -\eta \nabla J(w^{(1)}) + \alpha \Delta \cdot w^{(1)} + w^{(1)} = w^{(2)}$$

# NESTEROV'S Momentum

We first add the correction rule
at $(t-1)$ and then compute the
gradient at that location.

$$\Delta w^{(t)} = m^{(t)} - \gamma \cdot \nabla J(m^{(t)})$$

$$m^{(t)} = w^{(t)} + \alpha \cdot \Delta w^{(t-1)}$$



$$\longrightarrow -\gamma \cdot \nabla J(w^{(0)}) = \Delta \cdot w^{(0)}$$

$$\longrightarrow \alpha \cdot \Delta w^{(0)}$$

# Adaptive Learning Rate

$$\Delta w^{(t)} = -\gamma^{(t)} \cdot \nabla J(w^{(t)})$$

① SCHEDULER — periodically schedule the learning rate decrease.

$$\delta_{ij}^{(t)} \equiv \text{gain for connection } w_{ij} \text{ at iteration } t$$

$$\Delta w_{ij}^{(t)} = -\eta \cdot \delta_{ij}^{(t)} \cdot \nabla J(w_{ij}^{(t)})$$

If $\nabla J(w_{ij}^{(t)}) \cdot \nabla J(w_{ij}^{(t-1)}) > 0$ :

<span style="color:blue">moving towards same minima</span>

<span style="color:blue">ADDITIVE INCREASE</span>

$$\delta_{ij}^{(t+1)} = \delta_{ij}^{(t)} + 0.01$$

ELSE :

<span style="color:blue">moving in different directions</span>

$$\delta_{ij}^{(t+1)} = \delta_{ij}^{(t)} \times 0.99$$

## RMS Prop : only DECREASES

the gain using a multiplicative
decrease strategy.

## Adam (2015)

① Adds momentum term ( to increase $\beta_1 = 0.9$
                                          speed)
② PERforms adaptive learning rate. $\beta_2 = 0.99$
                                      (multiplicative decrease

## Nadam : it uses Nesterov's momentum.