

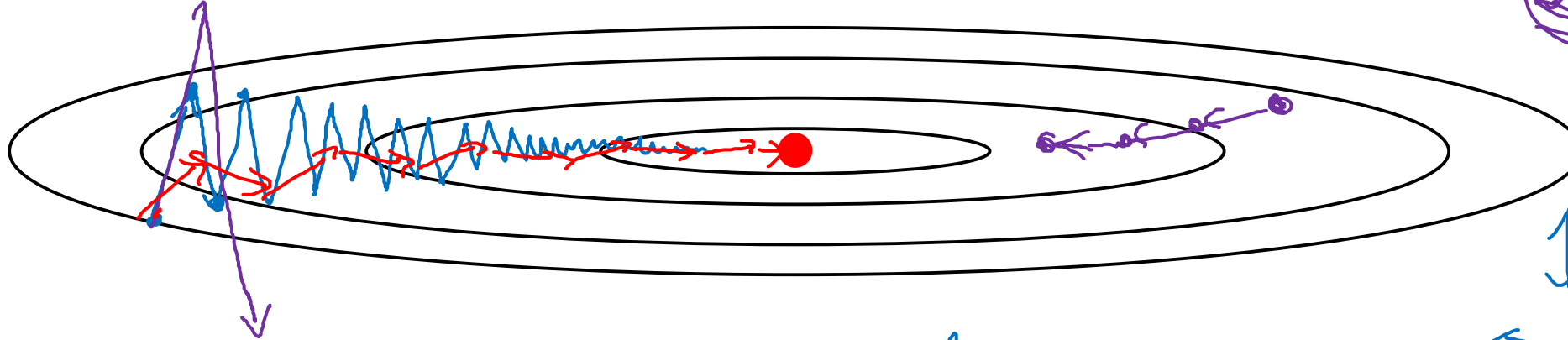


deeplearning.ai

Optimization Algorithms

Gradient descent
with momentum

Gradient descent example



↑ slower learning

↔ faster learning

Momentum:

On iteration t :

- Compute $\Delta W, \Delta b$ on current mini-batch.

- $V_{\Delta W} = \beta V_{\Delta W} + (1-\beta) \Delta W$

- $V_{\Delta b} = \beta V_{\Delta b} + (1-\beta) \Delta b$

friction → velocity → acceleration

- $W := W - \alpha V_{\Delta W}, \quad b := b - \alpha V_{\Delta b}$

To make the learning in the horizontal direction faster, to make the oscillation in the vertical direction smaller.

$$V_{\theta} = \beta V_{\theta} + (1-\beta) \theta_t$$

Implementation details

$$v_{dW} = 0, \quad v_{db} = 0$$

On iteration t :

Compute dW, db on the current mini-batch

$$\left. \begin{aligned} \rightarrow v_{dW} &= \beta v_{dW} + (1 - \beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1 - \beta) db \end{aligned} \right\} \quad \left| \quad \underbrace{v_{dW} = \beta v_{dW} + dW}_{\uparrow} \leftarrow$$

$$W = W - \underbrace{\alpha}_{\uparrow} \underbrace{v_{dW}}, \quad b = \underline{b} - \underbrace{\alpha}_{\uparrow} \underbrace{v_{db}}$$

$$\frac{v_{dW}}{1 - \beta^t}$$

Hyperparameters: α, β

$$\underline{\beta = 0.9}$$

average over last ≈ 10 gradients