



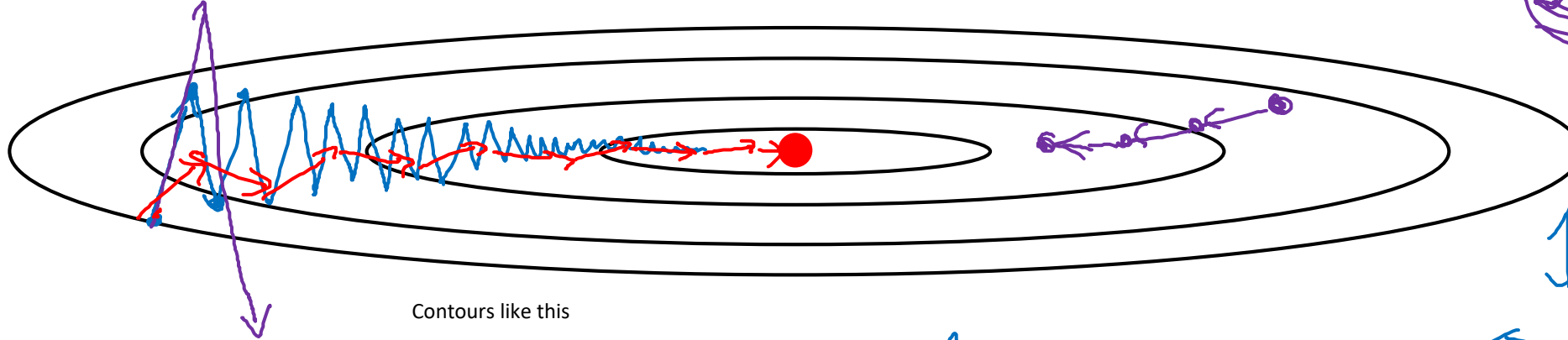
deeplearning.ai

Optimization Algorithms

Gradient descent with momentum

In one sentence, the basic idea is to compute an exponentially weighted average of your gradients, and then use that gradient to update your weights instead.

Gradient descent example



Contours like this

Momentum: And so the need to prevent the oscillations from getting too big forces you to use a learning rate that's not itself too large.

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 \rightarrow \uparrow velocity

\uparrow acceleration

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



\updownarrow slower learning
 \longleftrightarrow faster learning.

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

Implementation details

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

On iteration t :

Compute dW, db on the current mini-batch

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

it just affects what's the best value of the learning rate alpha.

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

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

Hyperparameters: α, β

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

average over last ≈ 10 gradients