

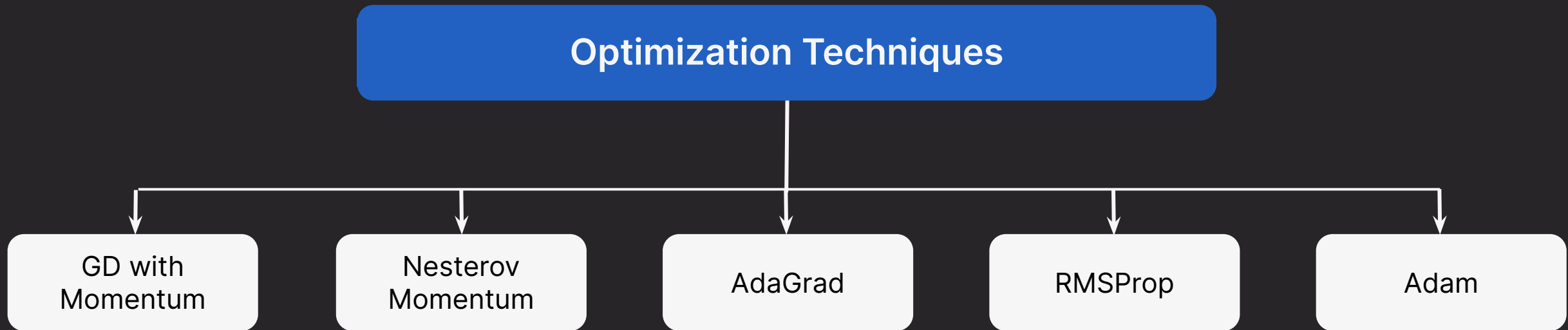
The background of the slide is a dark gray to black gradient, overlaid with a complex, white, abstract network pattern. This pattern consists of numerous small and large circular nodes connected by thin, light gray lines, creating a web-like structure that resembles a neural network or a data graph. The nodes are distributed across the entire frame, with some clusters being denser than others.

Understanding workings of Neural Networks

Video 6: Common Optimization Techniques – Part 1

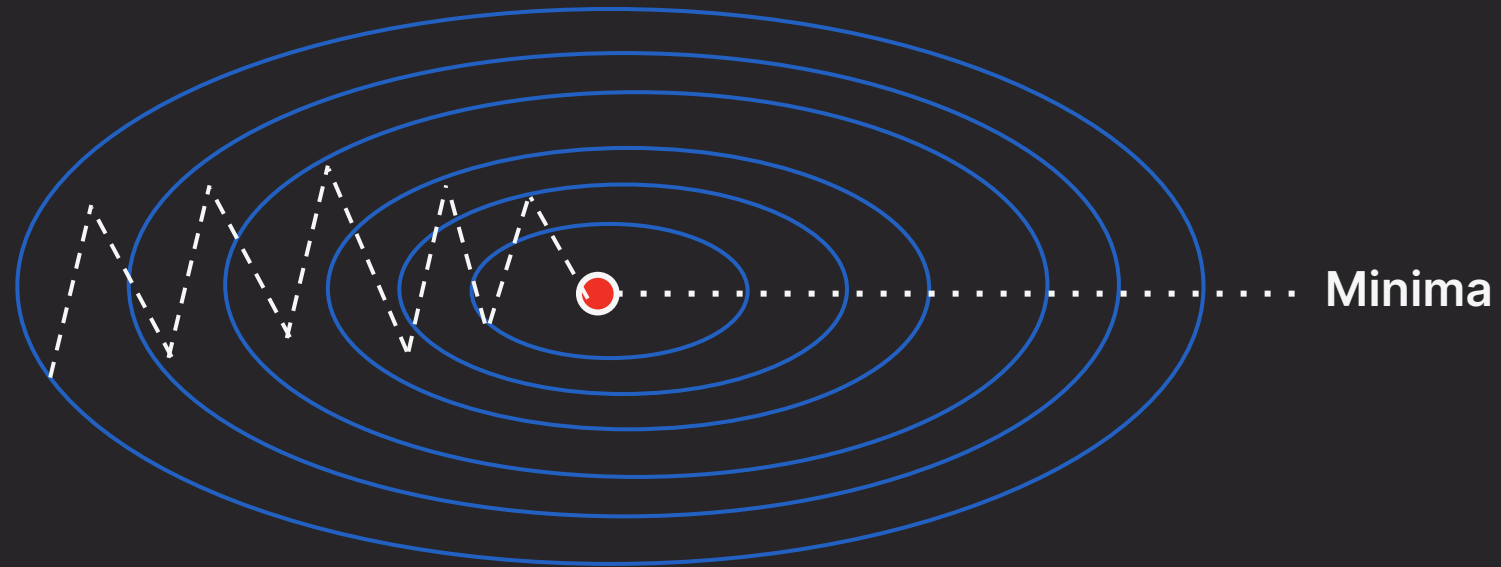
In Air

Common Optimization Techniques



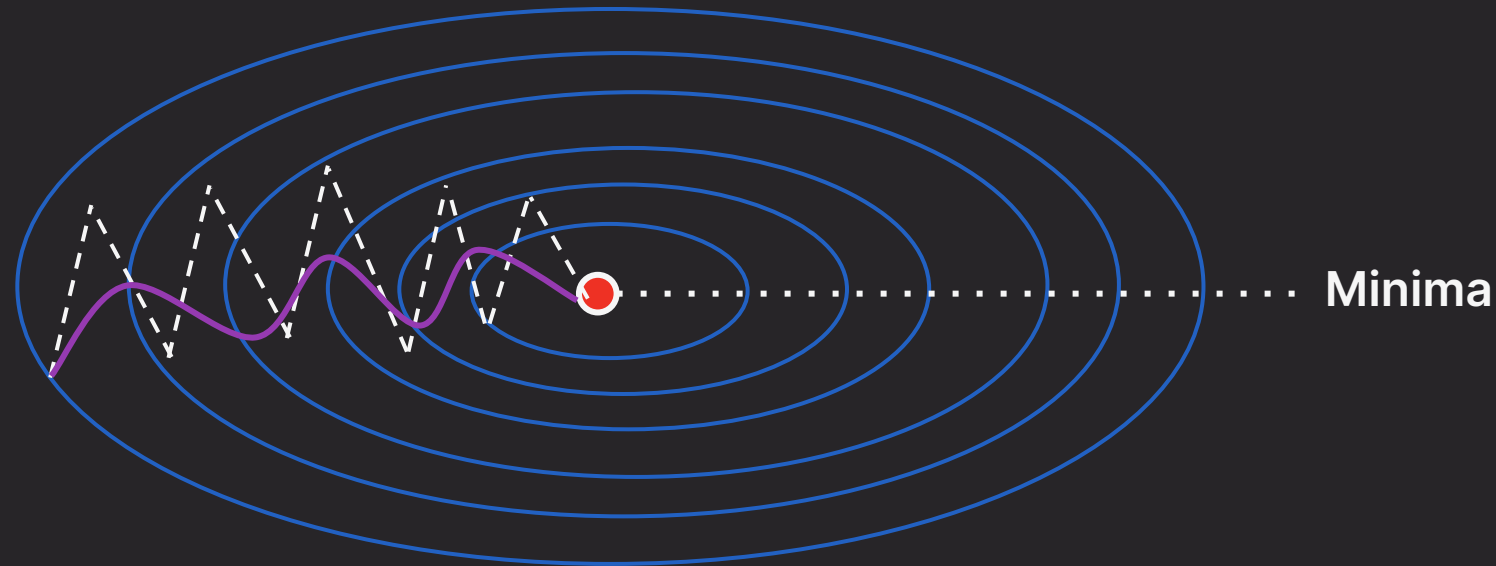
Gradient Descent with Momentum

Stochastic Gradient Descent with Momentum



--- GD

Stochastic Gradient Descent with Momentum



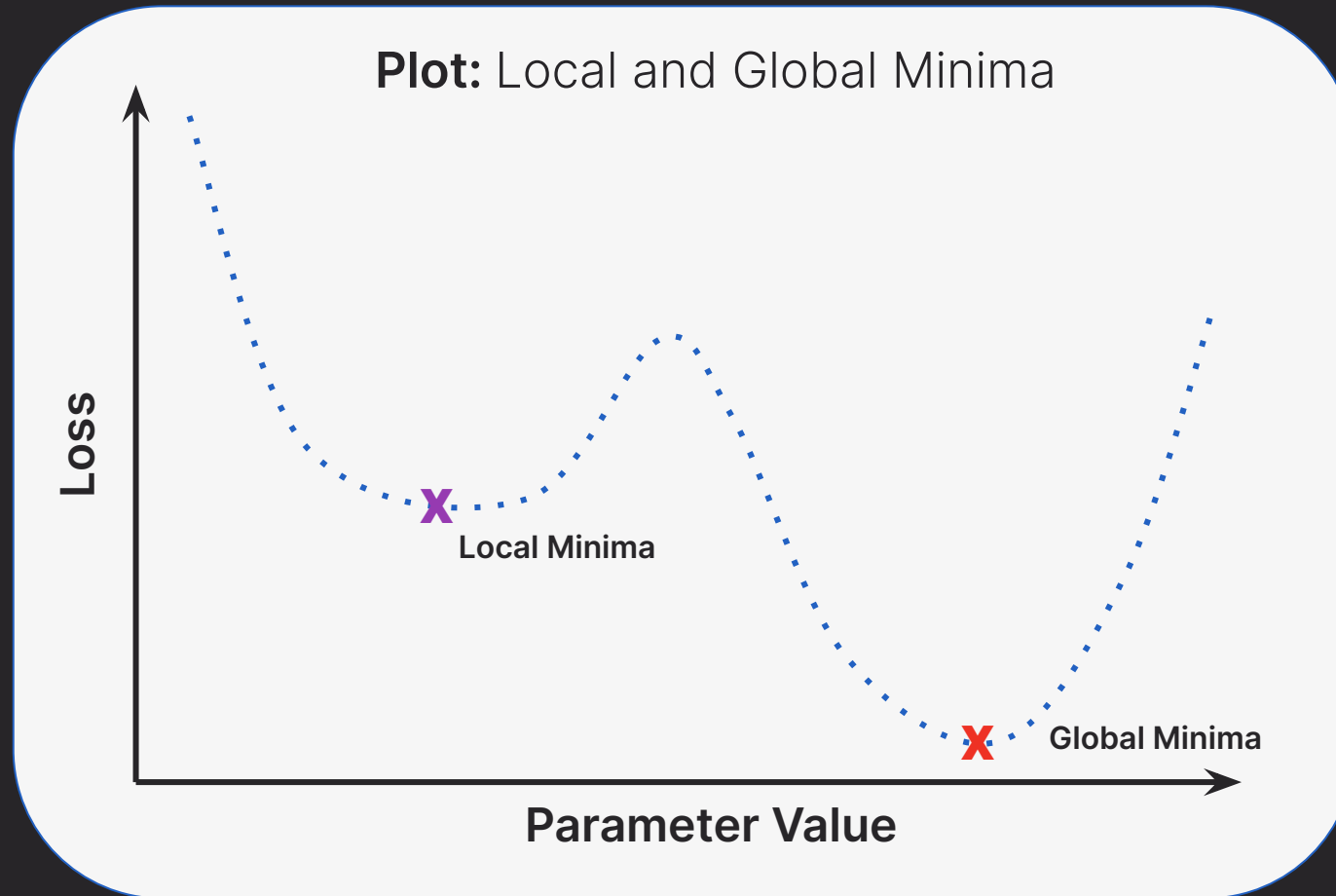
Approach: Finding the moving average of gradients

- Accelerate SGD
- Dampen the turbulence

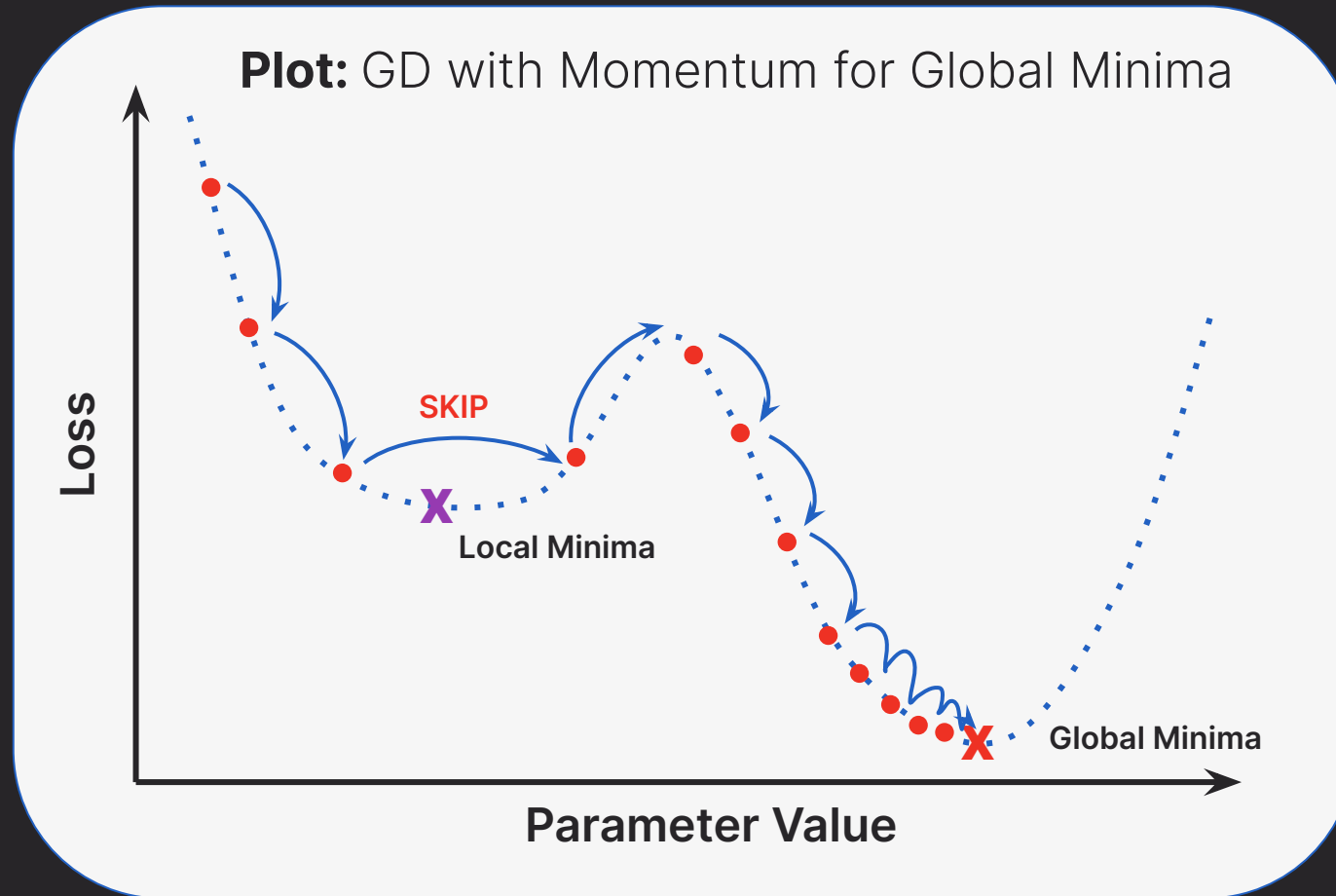
- - - - GD

— GD with Momentum

Stochastic Gradient Descent with Momentum



Stochastic Gradient Descent with Momentum



Stochastic Gradient Descent with Momentum

Moving average of gradients → Simple Moving Average



Computationally expensive procedure

Exponential Average

The Math Behind Exponential Average

1. Compute the gradient g_t at time step t .

$$g_t = \nabla_w L(w)$$

The Math Behind Exponential Average

1. Compute the gradient g_t at time step t .

$$g_t = \nabla_w L(w)$$

2. Calculate the moving average based on the previous timestamp value and the current gradient g_t .

$$v_t = \beta v_{t-1} + (1 - \beta)g_t$$

v_t = Moving Average

β = Momentum Term

$$\beta = 0.9 \sim \text{Avg of 10 time stamps} \sim \frac{1}{1 - \beta}$$

90% - weightage to the previous gradients

10% - weightage to current gradient

The Math Behind Exponential Average

1. Compute the gradient g_t at time step t .

$$g_t = \nabla_w L(w)$$

2. Calculate the moving average based on the previous timestamp value and the current gradient g_t .

$$v_t = \beta v_{t-1} + (1 - \beta)g_t$$

3. Update the weights w using the moving average and learning rate.

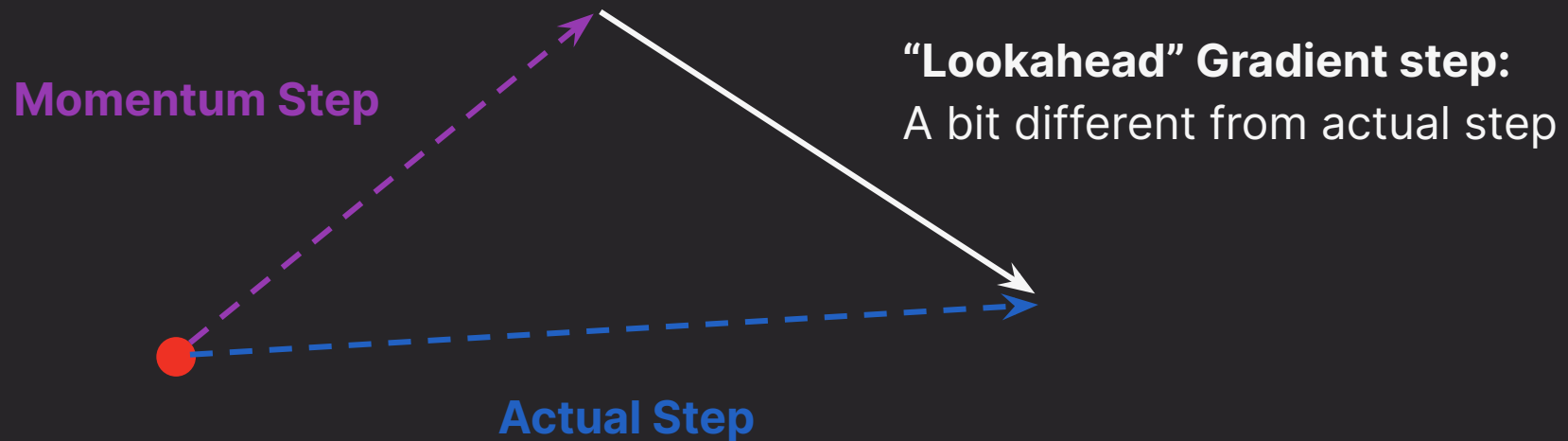
$$w = w - \eta v_t$$

$$\eta = \text{Learning Rate}$$

Nesterov Momentum

Nesterov Momentum

Nesterov refines momentum by looking at gradient or future steps.



The Math Behind Nesterov Momentum

1. Compute the gradient g_t at time step t .

$$w_{\text{lookahead}} = w - \beta v_{t-1}$$

The Math Behind Nesterov Momentum

1. Compute the gradient g_t at time step t .

$$w_{\text{lookahead}} = w - \beta v_{t-1}$$

2. Compute the gradient g_t at the lookahead weight $w_{\text{lookahead}}$

$$g_t = \nabla_w L(w_{\text{lookahead}})$$

The Math Behind Nesterov Momentum

1. Compute the gradient g_t at time step t .

$$w_{\text{lookahead}} = w - \beta v_{t-1}$$

2. Compute the gradient g_t at the lookahead weight $w_{\text{lookahead}}$

$$g_t = \nabla_w L(w_{\text{lookahead}})$$

3. Repeat steps 3 and 4 (same as SGD with Momentum).

$$v_t = \beta v_{t-1} + (1 - \beta)g_t$$

$$w = w - \eta v_t$$

Hands-on: GD with Momentum & Nesterov Momentum