

1.  $\Theta_0$

$\mathcal{M}$

$L(\Theta, B)$

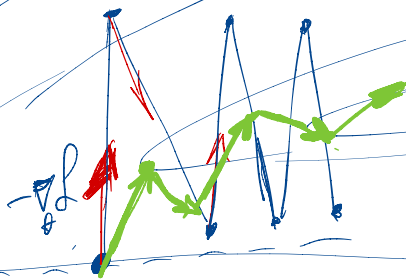
$$\bar{g} = \nabla_{\Theta} L \quad \mathcal{M}$$

$$\Theta_{t+1} = \Theta_t - \eta \bar{g}_t \quad \bar{g}_t = \nabla_{\Theta} L(\Theta_t, B_t)$$

Stochastic Gradient Descent

SGD

$\mathcal{LM}$



Моментум (Импульсный метод)

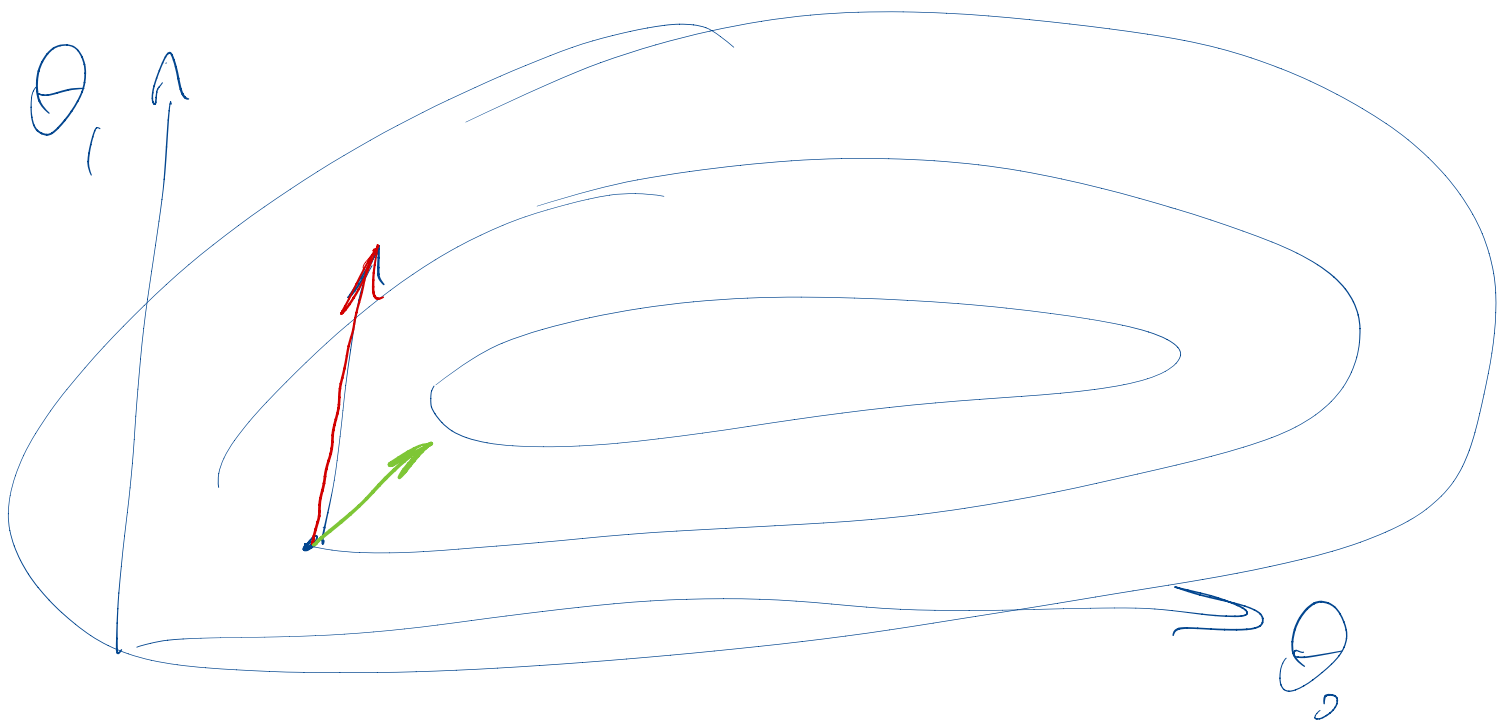
$$\bar{g}_t = \nabla_{\theta} L \quad \mathcal{M}$$

$$\theta_t: \mathcal{M}$$

$$\bar{m}_t = (1 - \beta) \bar{g}_t + \beta \bar{m}_{t-1} \quad \mathcal{M} \quad \beta = 0,9$$

$$\theta_{t+1} = \theta_t - \epsilon \bar{m}_t$$

$$3\mathcal{M}$$



~~$$S_t = (\bar{g}_t)^2$$~~

$$S_t = \beta S_{t-1} + (1 - \beta) (\bar{g}_t)^2$$

$$\theta_{t+1} = \theta_t - \eta \frac{\bar{g}}{\sqrt{S_t + \epsilon}}$$

Адаптивный градиент (АдаГрад)

# SGD + Ada Grad + Momentum Adam

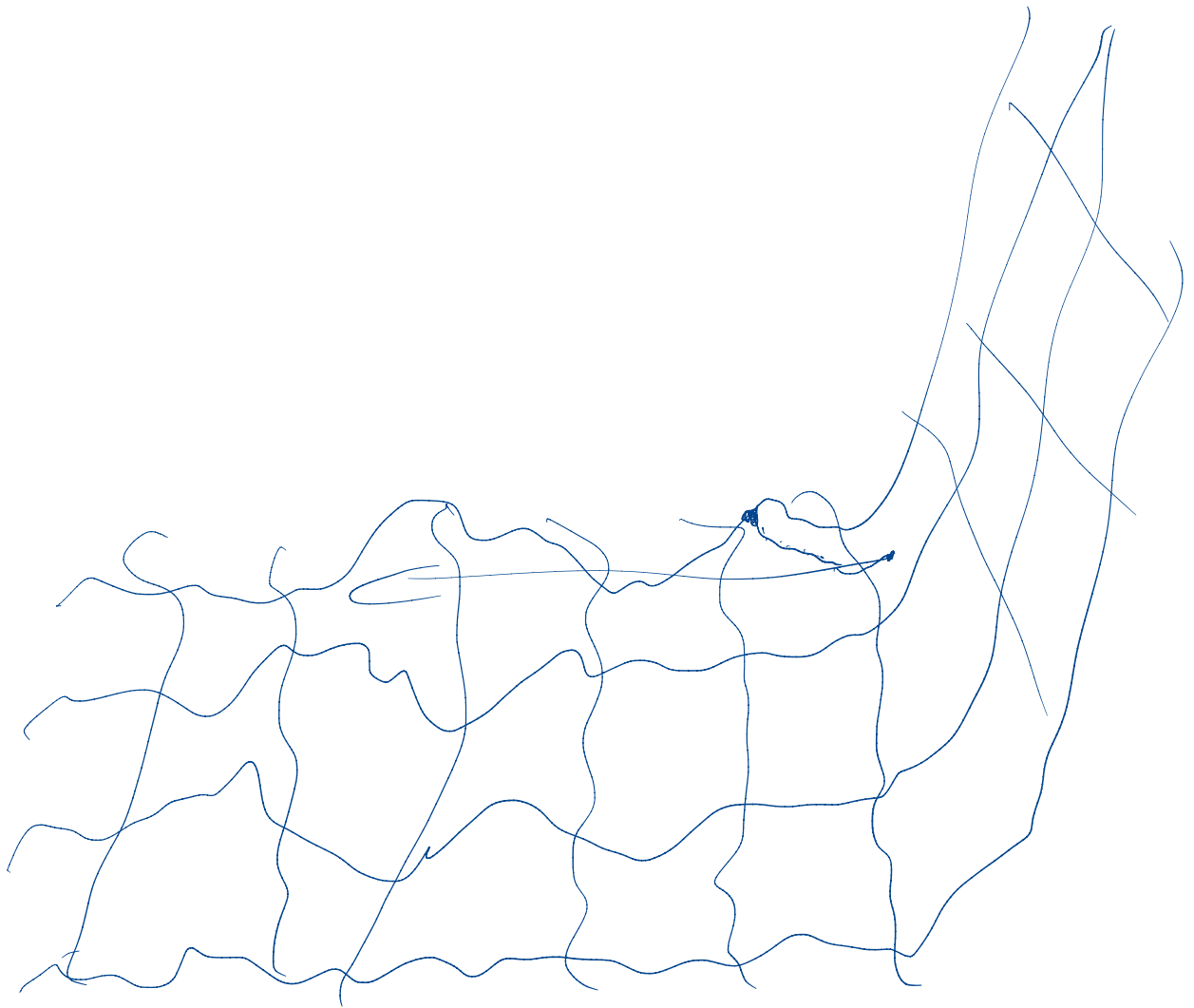
$$\theta_0 \quad \mu \quad \mu$$

$$\bar{m}_t = (1 - \beta_1) \bar{g}(\theta_t, b_t) + \beta_1 \bar{m}_{t-1}$$

$$S_t = \beta_2 S_{t-1} + (1 - \beta_2) \bar{g}_t^2$$

$$\mu$$

$$\theta_{t+1} = \theta_t - \eta \frac{m_t}{\sqrt{S_t} + \epsilon}$$



Gradient Clipping

$$\bar{g}^* = \frac{\bar{g}}{\|g\|} \cdot g$$

$$g_t^* = \begin{cases} g_t, & \text{if } \|g\| \leq L \\ Lg, & \text{if } \|g\| > L \end{cases}$$