

$$\theta_t = \theta_{t-1} - \eta \nabla_{\theta} L(\{(x_i, y_i)\}_n, \theta_{t-1})$$

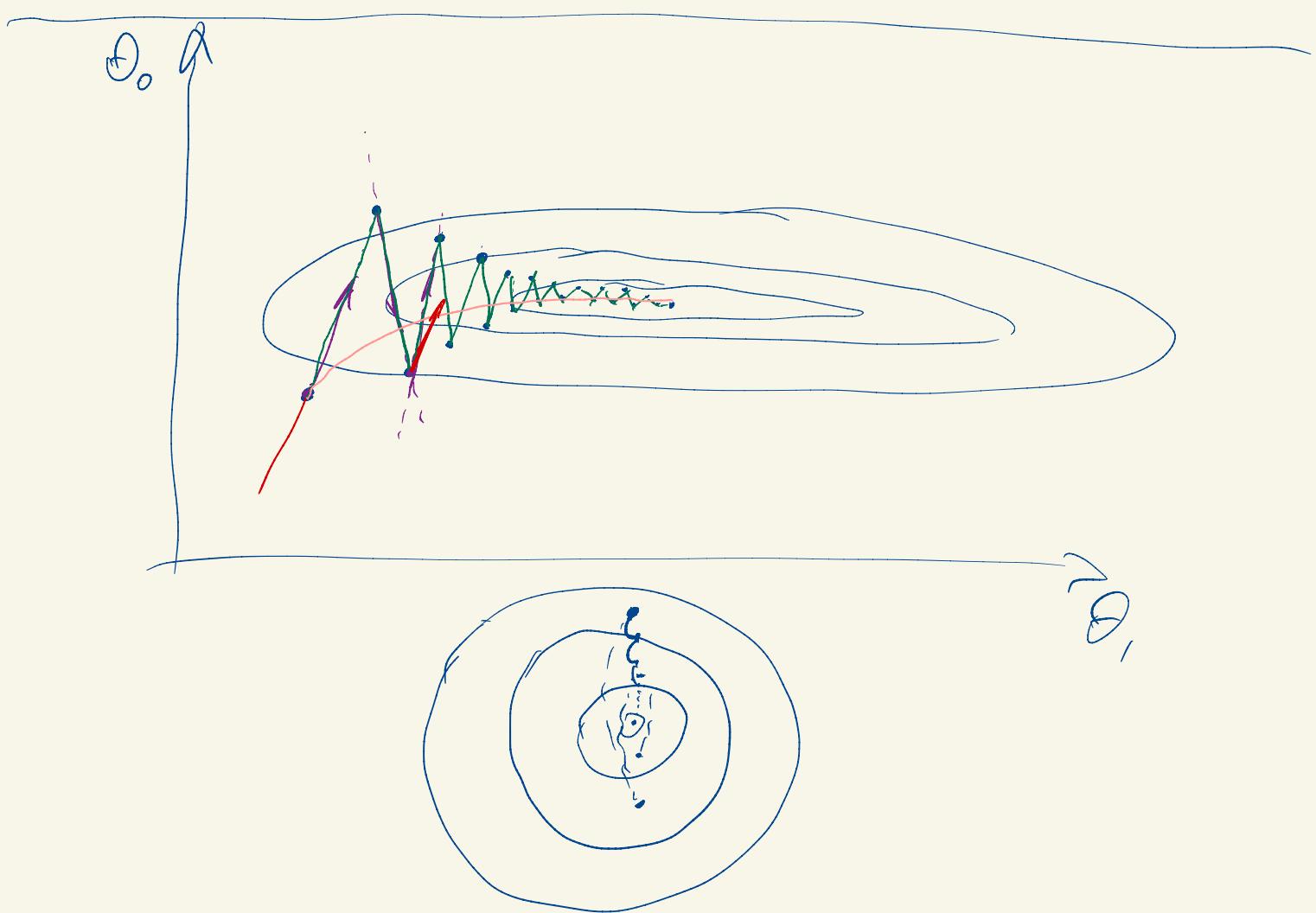
B_{t-1}

Итерации:

Финал: обучение на всем наборе данных
(1 итерация)

10'000

65:10



$$\theta_t = \theta_{t-1} - \eta \nabla_{\theta} \mathcal{L}(B_{t-1}, \theta_{t-1})$$

$$\left\{ \begin{array}{l} m_t = (1-\beta) \nabla_{\theta} \mathcal{L}(B_{t-1}, \theta_{t-1}) + \beta m_{t-1} \\ \theta_t = \theta_{t-1} - \eta m_t \end{array} \right.$$

SGD + momentum

умножающий вектор

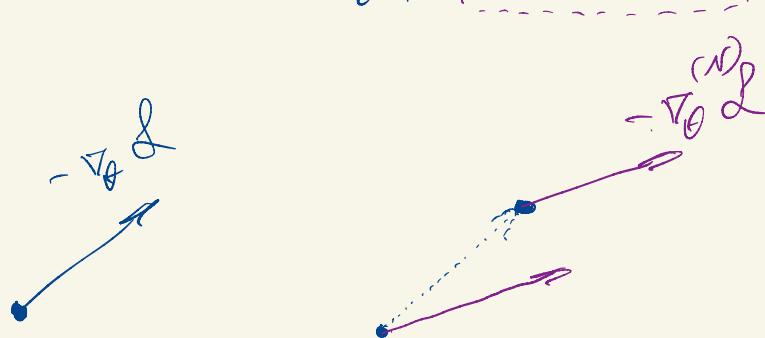
$$\beta = 0 \Rightarrow \text{SGD + momentum} = \text{SGD}$$

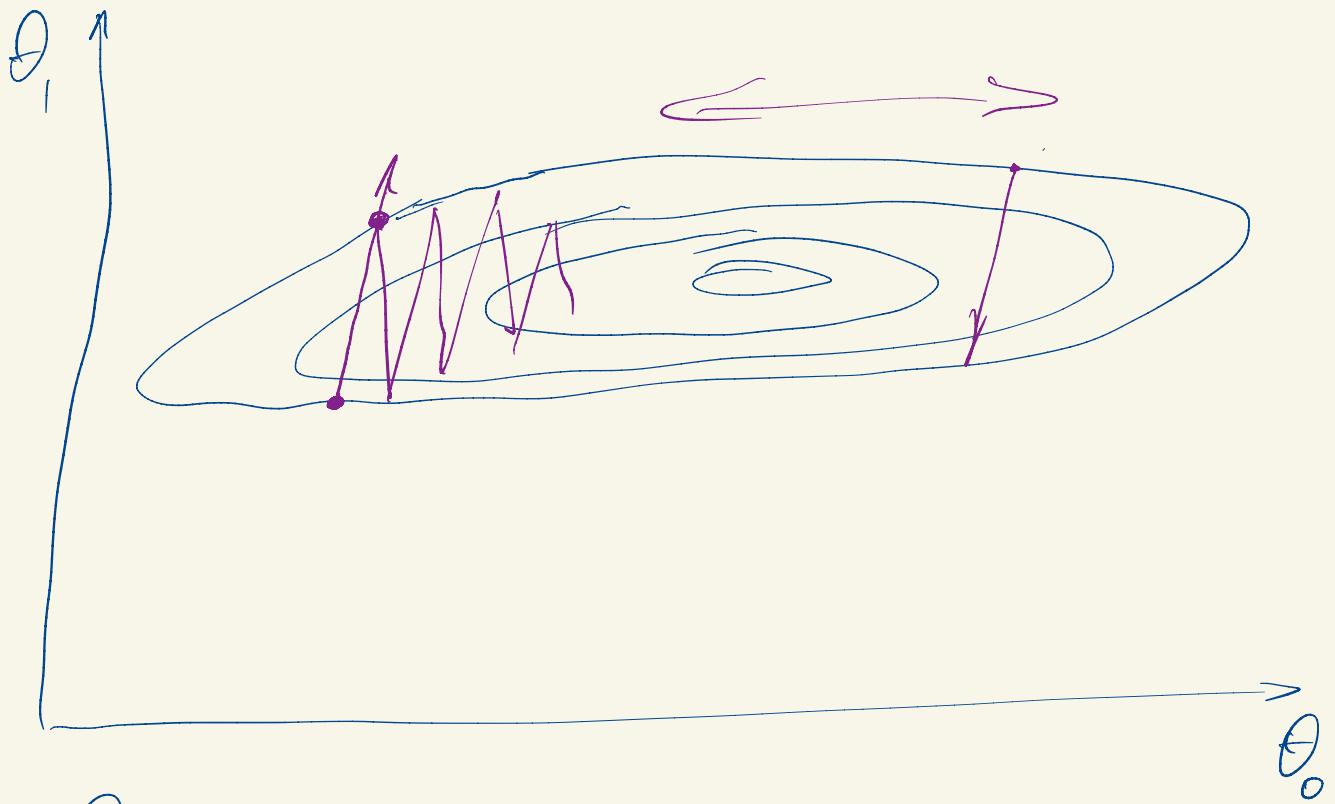
$$\beta = 1 \Rightarrow \text{коэффициент, определяющий неизменен}$$

Nesterov momentum

$$m_t = (1-\beta) \nabla_{\theta} \mathcal{L}(B_{t-1}, \theta_t) + \beta m_{t-1}$$

$$m_t^* = (1-\beta) \nabla_{\theta} \mathcal{L}(B_{t-1}, \theta_{t-1} - \frac{\beta m_{t-1}}{1-\beta}) + \beta m_{t-1}$$





$$\theta_t = \theta_{t-1} - \eta^m t$$

$$S_t = \beta S_{t-1} + (1-\beta) (\nabla_{\theta} L(\theta_{t-1}))^2$$

ноглаженное

$$\theta_t = \theta_{t-1} - \frac{\nabla L(\theta_{t-1})}{\sqrt{S + \epsilon}}$$

element-wise

Адаптивные шаги
 AdaGrad
 RMS Prop

$$\text{AdaGrad} + \text{Momentum} = \text{AdaM}$$

$$m_t = (1 - \beta_1) \nabla_{\theta} \mathcal{L}(\theta_{t-1}, \beta_{t-1}) + \beta_1 m_{t-1}$$

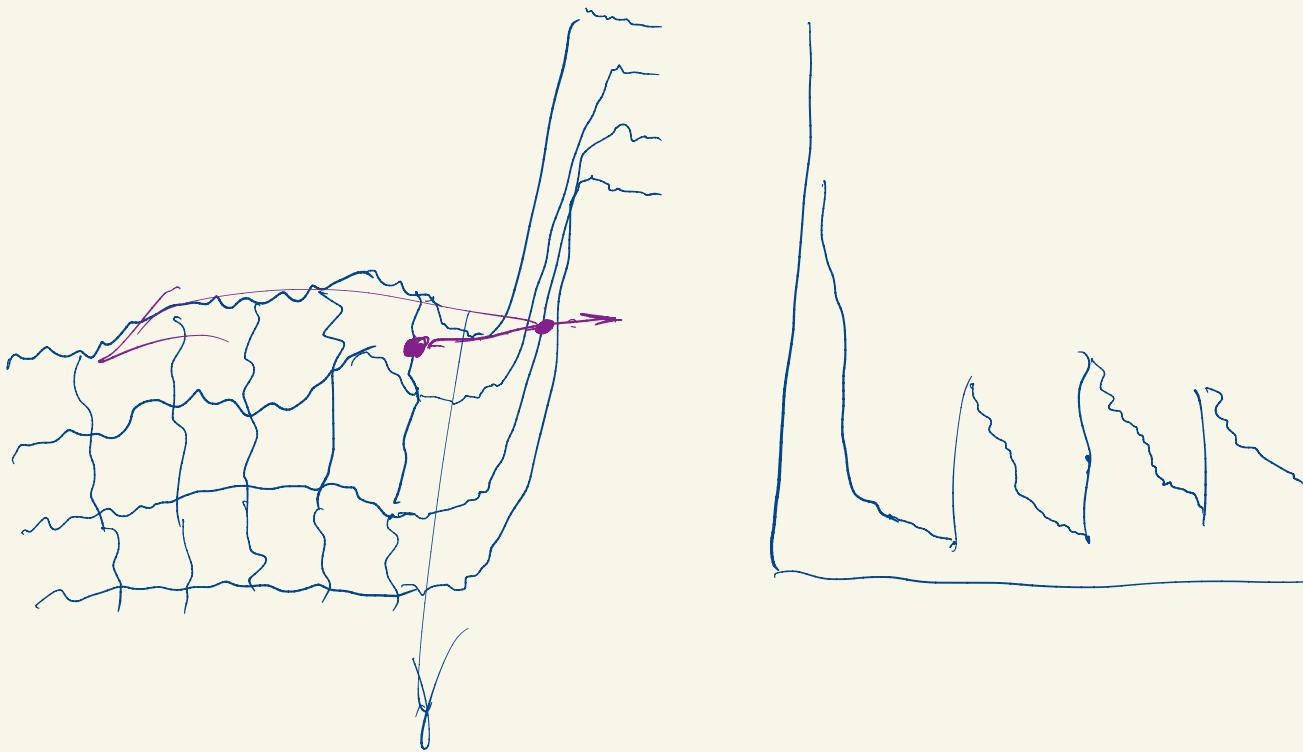
$$s_t = \beta_2 s_{t-1} + (1 - \beta_2) (\nabla_{\theta} \mathcal{L}(\theta_{t-1}, \beta_{t-1}))^2$$

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

Adagrad

$$\beta_1 = 0,9$$

$$\beta_2 = 0,99$$



Gradient clipping }

$$g_t = \nabla_{\theta} \mathcal{L}(B_t, \theta_t)$$

① global norm clipping

$$g_t^* = \frac{g_t}{\|g_t\|} \cdot S_g = 1$$

② component-wise clipping

$$g_{t,i}^* = \begin{cases} g_{t,i}, & \text{if } g_{t,i} \leq L_g \\ L_g, & \text{if } g_{t,i} > L_g \end{cases}$$