

$\theta$

$$\hat{y}_i = F(x_i, \theta)$$

$$\mathcal{T} = \{(x_i, y_i)\}_{i=1}^N$$

$$\mathcal{L} = \mathcal{L}(\theta, \{(x_i, y_i)\})$$

$$\vec{\nabla}_{\theta} \mathcal{L} \quad \vec{\theta}$$

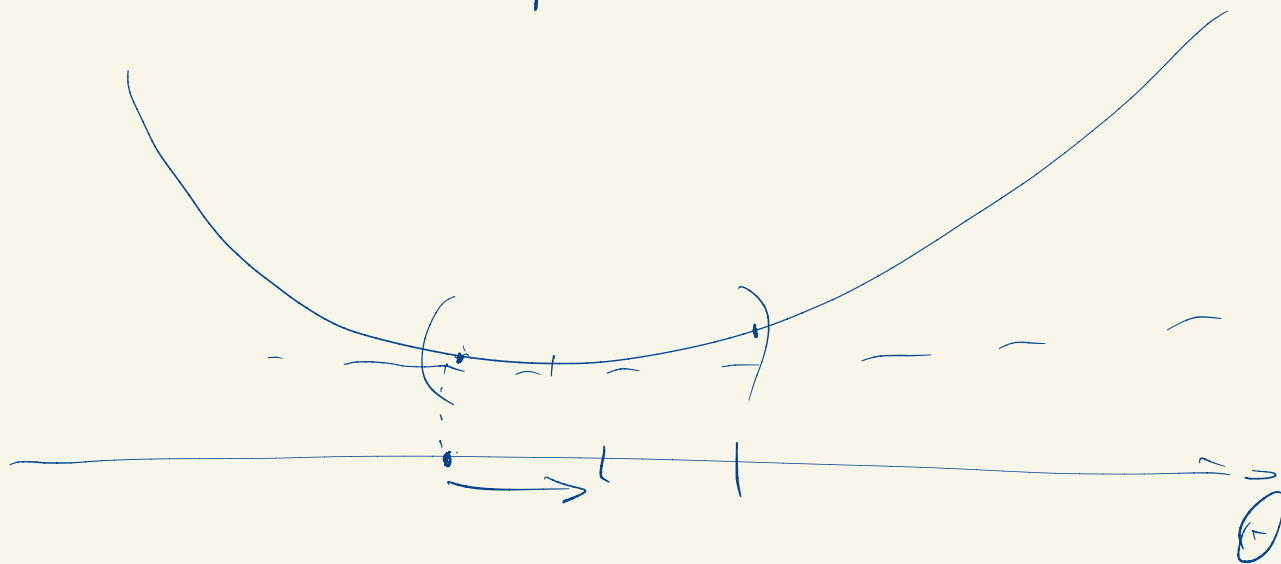
$$\textcircled{1} \theta_0 \sim \mathcal{N}(0, 1)$$

$\eta$  - max обьёмное  
learning rate

$$\textcircled{2} \vec{g} = \vec{\nabla}_{\theta} \mathcal{L}(\theta_t) = \vec{\nabla}_{\theta} \mathcal{L}_t = \vec{\nabla}_{\theta} \mathcal{L}(\theta_t, \mathcal{T})$$

$$\theta_{t+1} = \theta_t - \eta \vec{\nabla}_{\theta} \mathcal{L}_t$$

$$\textcircled{3} \mathcal{L} \leq \varepsilon \Rightarrow \text{stop} \Rightarrow \theta^*$$



SGD

Стохастический градиентный спуск  
Stochastic gradient descent

$$\textcircled{1} \theta_0 \sim \mathcal{N}(0, 1)$$

$n$  - batch size

$\eta$  - шаг обучения  
learning rate

$$\textcircled{2} g = \nabla_{\theta} L(\theta_t, D)$$

$D$  - выборка размера  $n$

$$\theta_{t+1} = \theta_t - \eta g$$

$$\textcircled{3} L \leq \varepsilon \Rightarrow \text{stop} \Rightarrow \theta^*$$

---

$$\textcircled{1} h^* = k n$$

$$\eta^* = k \eta$$

