# Стохастическая градиентная оптимизация

$$\Theta^* = \arg\min_{(\Theta)} \mathcal{L}(\mathcal{T}, \Theta)$$

(1) $\Theta_0$

(2) (Adam) $\beta_1, \beta_2, \ell, \quad C$ — предикат остановки

$$B, \quad t = 0$$

(3) 1) $\{X, Y\}^B - \mathcal{T}^{(B)}$

2) $\mathcal{L} = \mathcal{L}(\mathcal{T}^{(B)}, \Theta^{(t)})$

3) $g = \nabla_\Theta \mathcal{L} \qquad (\text{Adam}) \Rightarrow g^*$

4) $\Theta^{(t+1)} = \Theta^{(t)} - \ell g^*$

5) $C? \Rightarrow stop$

$$G(z) = \tanh(z)$$

tanh(z)

1

-1

$z$

$$\hat{y} = \varphi\left( \ldots \left( \theta^{(2)} G\left( \theta^{(1)} G\left( \theta^{(0)} X \right) \right) \right) \ldots \right)$$

$$z^{(0)}$$

$$\text{Var } z_i = \text{Var } \theta_i h_i = \text{Var } \theta_i \ \text{Var } h_i$$

$$z^{(1)} = \sum_{i=1}^{m} \theta_i h_i$$

$$\text{Var } z^{(1)} = \sum_{i=1}^{m} \text{Var} \theta_i \ \text{Var } h_i = m \ \text{Var } \theta_i \ \text{Var } h_i$$

$$\sigma = \tanh \qquad \text{Var } \overbrace{G(z^{(1)})}^{h} = \text{Var } z^{(1)}$$

$$\text{Var } h^{(\ell+1)} = m \ \text{Var } \underbrace{\theta^{(\ell+1)} h^{(\ell)}}_{c}$$

$$\text{Var } h^{(\ell+k)} = m^k c^k \ \text{Var } h^{(\ell)}$$

$$(m c) \approx 1$$

$$Var\,\theta^{(\ell)} \approx \frac{1}{m^{(\ell)}}$$

$$\theta \sim \mathcal{N}\left(0, \sigma^2 = \frac{1}{m}\right)$$

$$\theta \sim U\left(-\frac{\sqrt{m}}{2}; \frac{\sqrt{m}}{2}\right)$$

---

## Kaiming He 2015

ReLU: $\qquad Var\,\theta^{(\ell)} = \dfrac{2}{n^{(\ell-1)}}$

PReLU: $\qquad Var\,\theta^{(\ell)} = \dfrac{2}{(1+\alpha^2)\,n^{(\ell-1)}}$

$n$ — ширина слоя
$(\ell-1)$

PReLU



$y = z$

$y = \alpha z$