$$\theta^* = \text{argmin } \mathcal{L}(\theta, \mathcal{T})$$

---

⓪ $\theta^{(0)}$:

$\eta$ – learning rate; $\eta \overset{def}{=} 10^{-4}$

$C$ – условие остановки.

① $g_t = \nabla_\theta \mathcal{L}(\theta^{(t)}, B^{(t)})$

$g^*$ (Adam)

② $\theta^{(t+1)} = \theta^{(t)} - g^* \eta$

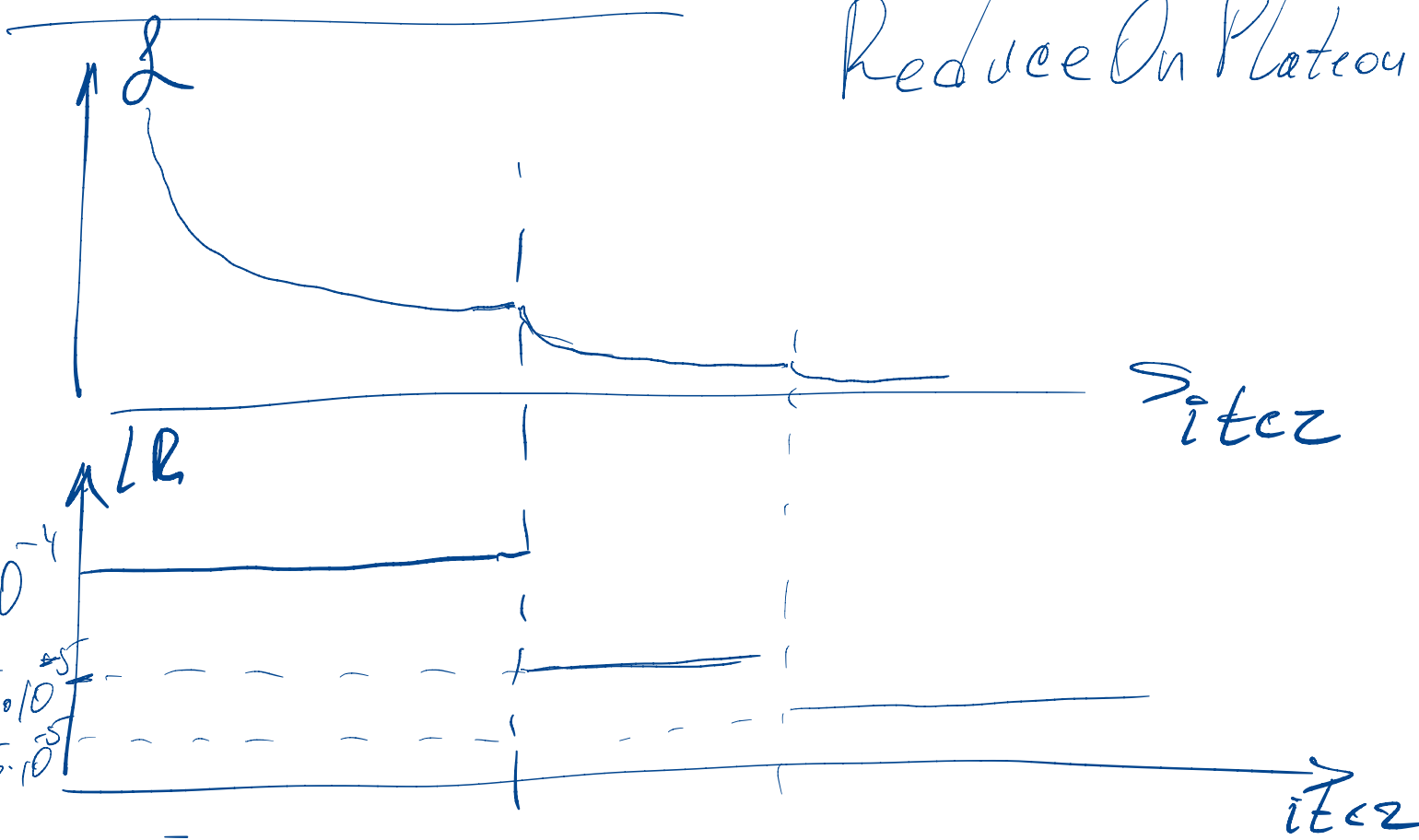③ $C? \Rightarrow$ stop

covariate shift

Стратегии изменения темпа обучения

LR schedulers

Reduce On Plateau

# Exponential LR



iter

# LR

$$l_2^{(t+1)} = \gamma \, l_2^{(t)}$$

iter

# Simulated annealing