

Lecture 3

UAT (Universal Approximation Theorem)

Cybenko 1989

$$\forall f \in C([0, 1]^d) \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

σ - activation f.

$$\sigma: \mathbb{R} \rightarrow \mathbb{R}$$

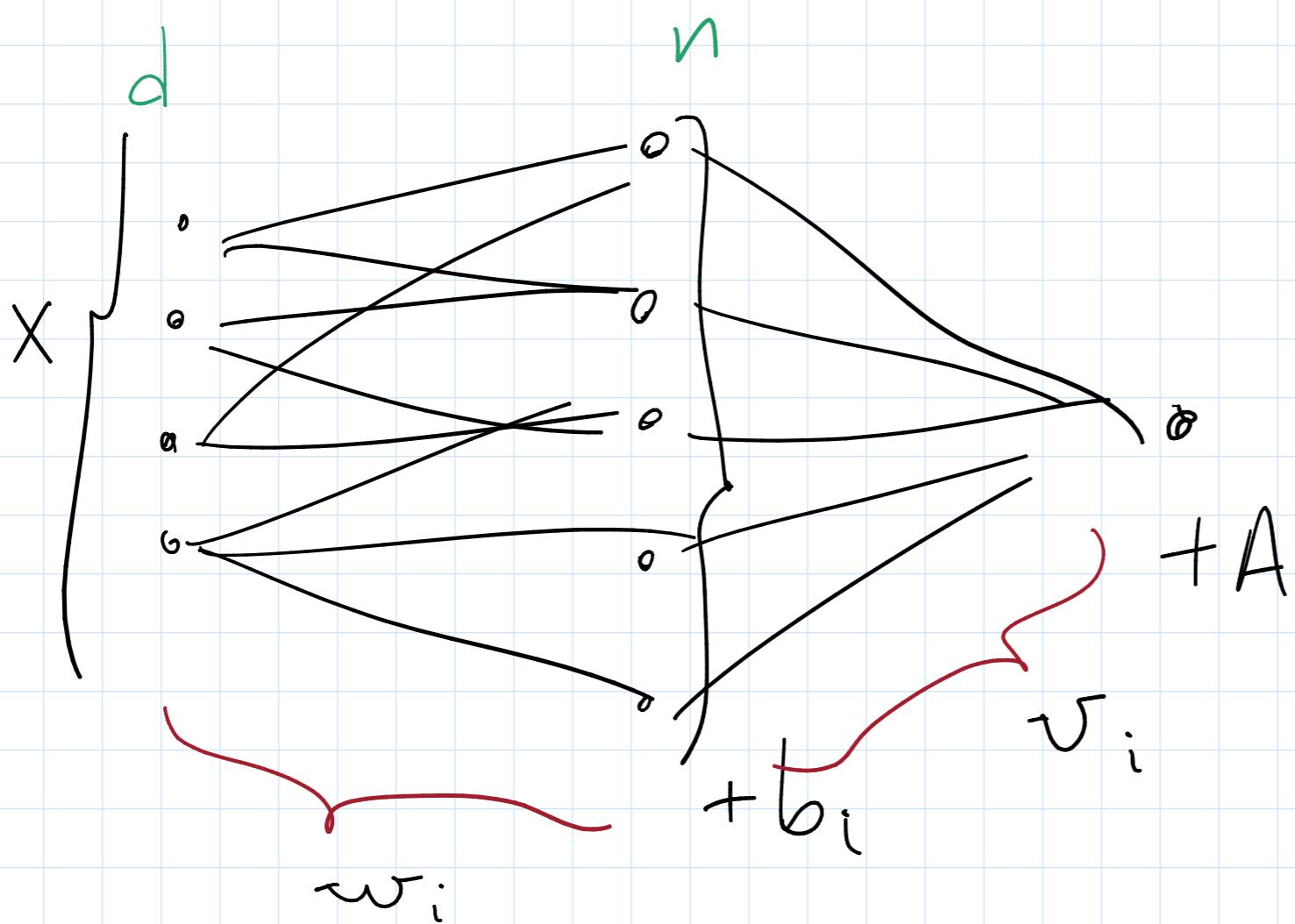
$$\begin{cases} \lim_{x \rightarrow -\infty} \sigma(x) = 0 \\ \lim_{x \rightarrow +\infty} \sigma(x) = 1 \end{cases} \Rightarrow \text{non-linear}$$

Fix error level $\forall \epsilon > 0$:

$$\exists n \in \mathbb{N}, \quad \left\{ w_i \right\}_{i=1}^n \quad w_i \in \mathbb{R}^d \quad \left\{ b_i \right\}_{i=1}^n \quad b_i \in \mathbb{R}$$

$$\left\{ v_i \right\}_{i=1}^n \quad v_i \in \mathbb{R} \quad A \in \mathbb{R}$$

$$g(x) = \sum_{i=1}^n v_i \sigma(w_i \cdot x + b_i) + A$$



$$\sup_{x \in [0,1]^d} |f(x) - g(x)| \leq \varepsilon$$

Optimization

$$Q^t(\omega) = \frac{1}{B} \sum_{i=1}^B L(f_\omega(x_{t_i}), y_{t_i})$$

t - opt. step

$$t_i \sim U(\{1, \dots, \ell\})$$

$$\{(x_i, y_i)\}_{i=1}^\ell$$

$$\begin{cases} g_t = \nabla_\omega Q^t(\omega_{t-1}) \\ m_t = \mu m_{t-1} + g_t \\ \omega_t = \omega_{t-1} - \eta_t \cancel{g_t} \cancel{m_t} \end{cases}$$

SGD + momentum

η_t — learning rate

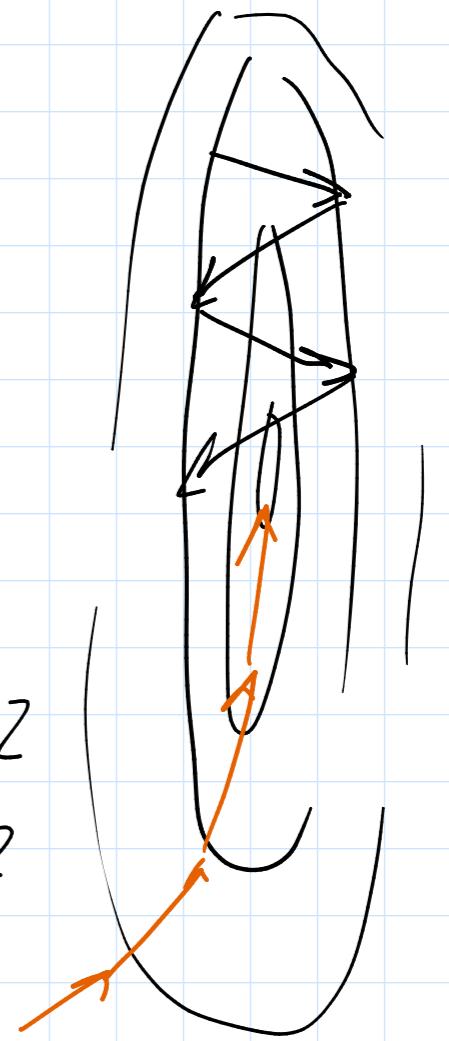
$$m_0 = 0, 0 < \mu \leq 1, \mu = 0, 9$$

$$\mu m_{t-1} + g_t (\mu-1)$$

$$Q^t(\omega) \rightarrow Q_\lambda^t(\omega) = Q^t(\omega) + \frac{\lambda}{2} \|\omega\|_2^2$$

$$\nabla_\omega \| \omega \|_2^2 = 2\omega$$

$$\begin{cases} g_t = \nabla_\omega Q^t(\omega_{t-1}) + \lambda \omega_{t-1} \\ m_t = \mu m_{t-1} + g_t \\ \omega_t = \omega_{t-1} - \eta_t m_t = \end{cases}$$



$$= w_{t-1} - \eta_t (\mu m_{t-1} + \nabla_w Q^t(w_{t-1}) + \lambda w_{t-1}) =$$

$$= w_{t-1} (1 - \eta_t \lambda) - \eta_t (\mu m_{t-1} + \nabla_w Q^t(w_{t-1}))$$

(1 - $\eta_t \lambda$) < 1
 m_t^{old}

Weights Decay

$$\lambda \sim 10^{-4}, 10^{-5}$$

Adam

1. Adaptive learning rate — AdaGrad, RMSProp
 2. Momentum

$$\left\{ \begin{array}{l} g_t = \nabla_{w^t} Q(w_{t-1}) + \lambda w_{t-1} \\ M_t = \beta_1 M_{t-1} + (1-\beta_1) g_t \\ \sqrt{v}_t = \beta_2 \sqrt{v}_{t-1} + (1-\beta_2) g_t \odot g_t \\ \hat{M}_t = \frac{M_t}{1-\beta_1^t} \quad \hat{\sqrt{v}}_t = \frac{\sqrt{v}_t}{1-\beta_2^t} \end{array} \right. \quad \text{--- element-wise prod.}$$

$M_0 = 0$

$\sqrt{v}_0 = 1$

$$w_t = w_{t-1} - \eta_t \frac{\hat{M}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

$\epsilon, \beta_1, \beta_2$ — hyperparameters $\epsilon \approx 10^{-8}$

$$\beta_1 = 0,9 \quad \beta_2 = 0,999$$

$$m_t = (1-\beta_1)g_t + \beta_1 m_{t-1} = (1-\beta_1)g_t + \beta_1(1-\beta_1)g_{t-1}$$

$$+ \beta_1^2 m_{t-2} = \sum_{i=0}^{t-1} (1-\beta_1)\beta_1^i g_{t-i}$$

$$\hat{m}_t = \frac{1}{1-\beta_1^t} m_t = \frac{1-\beta_1}{1-\beta_1^t} \sum_{i=0}^{t-1} \beta_1^i g_{t-i} =$$

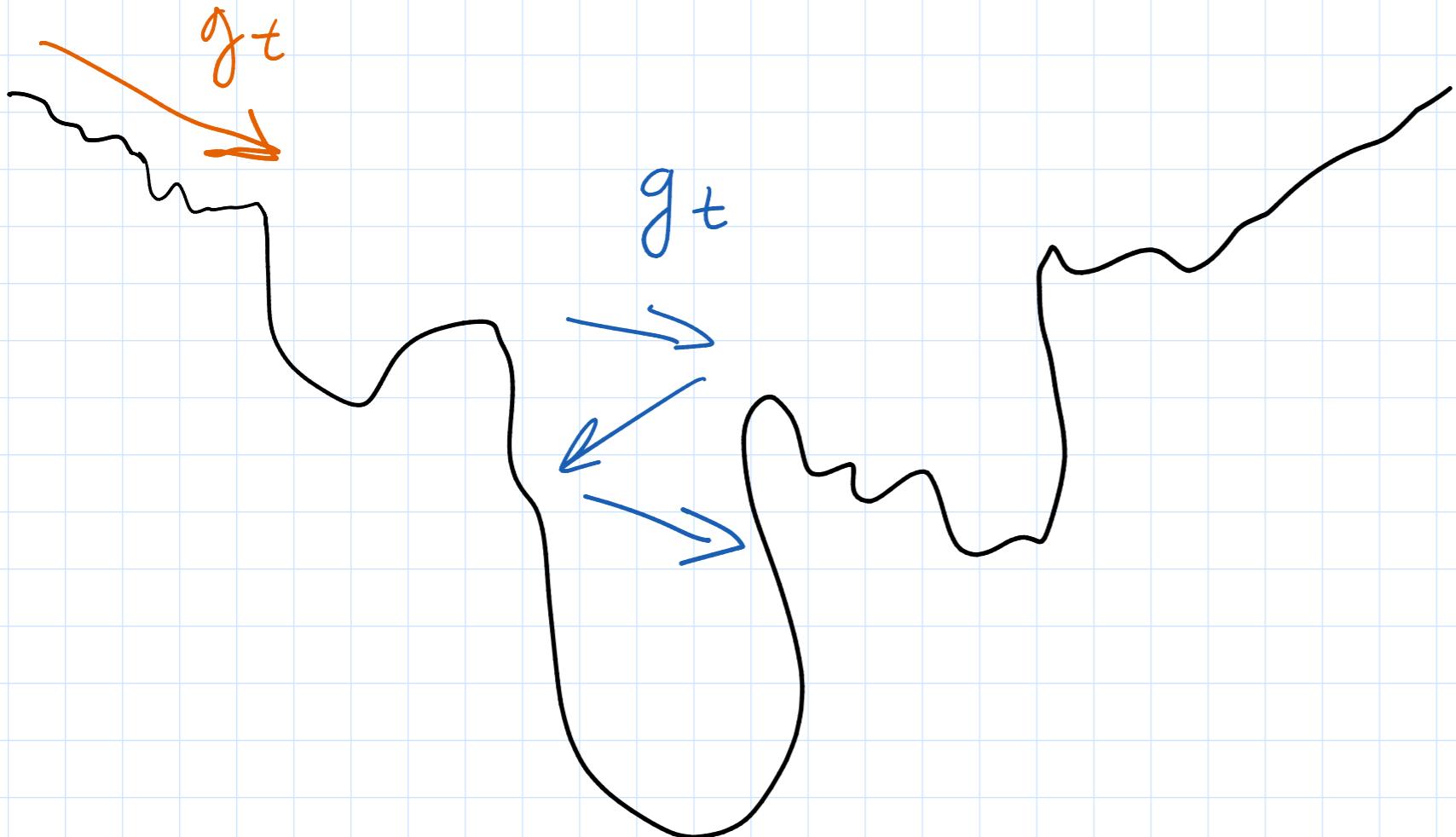
$$= \frac{\sum_{i=0}^{t-1} \beta_1^i g_{t-i}}{\sum_{j=0}^{t-1} \beta_1^j} \stackrel{||}{=} (1-\beta_1)(1 + \beta_1 + \dots + \beta_1^{t-1})$$

$$\alpha_i = \frac{\beta_1^i}{\sum_{j=0}^{t-1} \beta_1^j}$$

$$\sum_i \alpha_i = 1$$

$$\stackrel{||}{=} \sum_{i=0}^{t-1} \alpha_i g_{t-i} = \mathbb{E} g$$

$$\hat{\sigma}_t^2 = \mathbb{E}[g \circ g] = \text{Var}(g) + \mathbb{E}g \circ \mathbb{E}g$$



1. Far from optimum

$$\text{Var}(g) \ll \|g\|^2$$

$$w_t = w_{t-1} - \eta_t \frac{\hat{m}_t}{\sqrt{\hat{\sigma}_t + \epsilon}} \approx w_{t-1} - \eta_t$$

$$\frac{\mathbb{E} g}{\sqrt{\text{Var}(g) + \mathbb{E} g^2} + \epsilon} \approx 1$$

≈ 0

2. Close to optimum

$$\text{Var}(g) \gg 0 \quad \frac{\mathbb{E} g}{\sqrt{\text{Var}(g) + \mathbb{E} g^2} + \epsilon} < 1$$

SGD + momentum
CNN

Adam
RNN, Transformer

$$\begin{aligned}
 w_t &= w_{t-1} - \eta_t \frac{\hat{m}_t}{\sqrt{\hat{\sigma}_t + \epsilon}} = w_{t-1} - \eta_t \frac{m_t}{(1-\beta_1^t)\sqrt{\dots}} = \\
 &= w_{t-1} - \eta_t \frac{\beta_1 m_{t-1} + \nabla_w Q^t(w_{t-1}) + \lambda w_{t-1}}{(1-\beta_1^t)\sqrt{\dots}} = \\
 &= w_{t-1} \left(1 - \frac{\eta_t \lambda}{(1-\beta_1^t)\sqrt{\hat{\sigma}_t + \epsilon}} \right)
 \end{aligned}$$

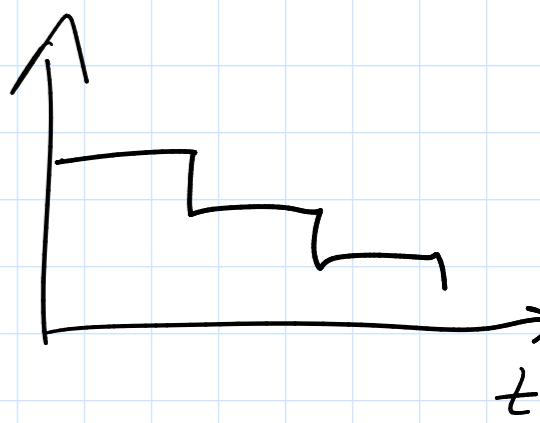
Adam

$$\left\{
 \begin{array}{l}
 g_t = \nabla_{\omega} Q^t(\omega_{t-1}) + \cancel{\lambda \omega_{t-1}} \\
 m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t \\
 v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t \odot g_t \\
 \hat{m}_t = \frac{m_t}{1-\beta_1^t} \quad \hat{v}_t = \frac{v_t}{1-\beta_2^t} \\
 \omega_t = \omega_{t-1} \left(1 - \lambda \eta_t\right) - \eta_t \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}
 \end{array}
 \right.$$

η_t - ? How to choose learning rate?

- $\eta_t = \eta_0$ - const scheduler Adam $\eta_t = 10^{-3}$

- Step LR

$$\eta_t = \begin{cases} \eta_0 & 0 < t \leq t_0 \\ \eta_1 & t_0 < t \leq t_1 \\ \vdots \\ \eta_0 > \eta_1 > \dots \end{cases}$$


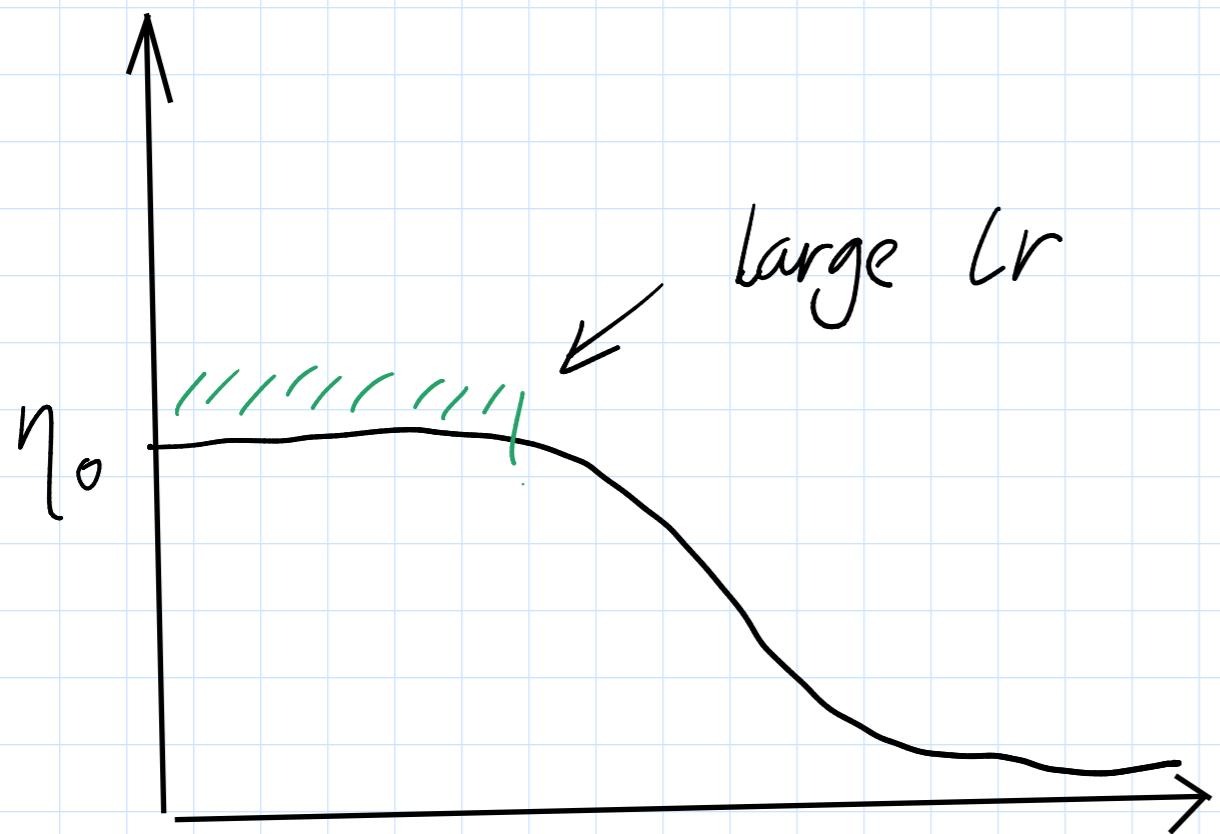
- Exponential LR

$$\eta_t = \gamma \cdot \eta_{t-1} \quad \gamma < 1$$


- Cosine LR

$$\eta_t = \eta_0 \frac{1}{2} \cos\left(\frac{\pi t}{T}\right)$$

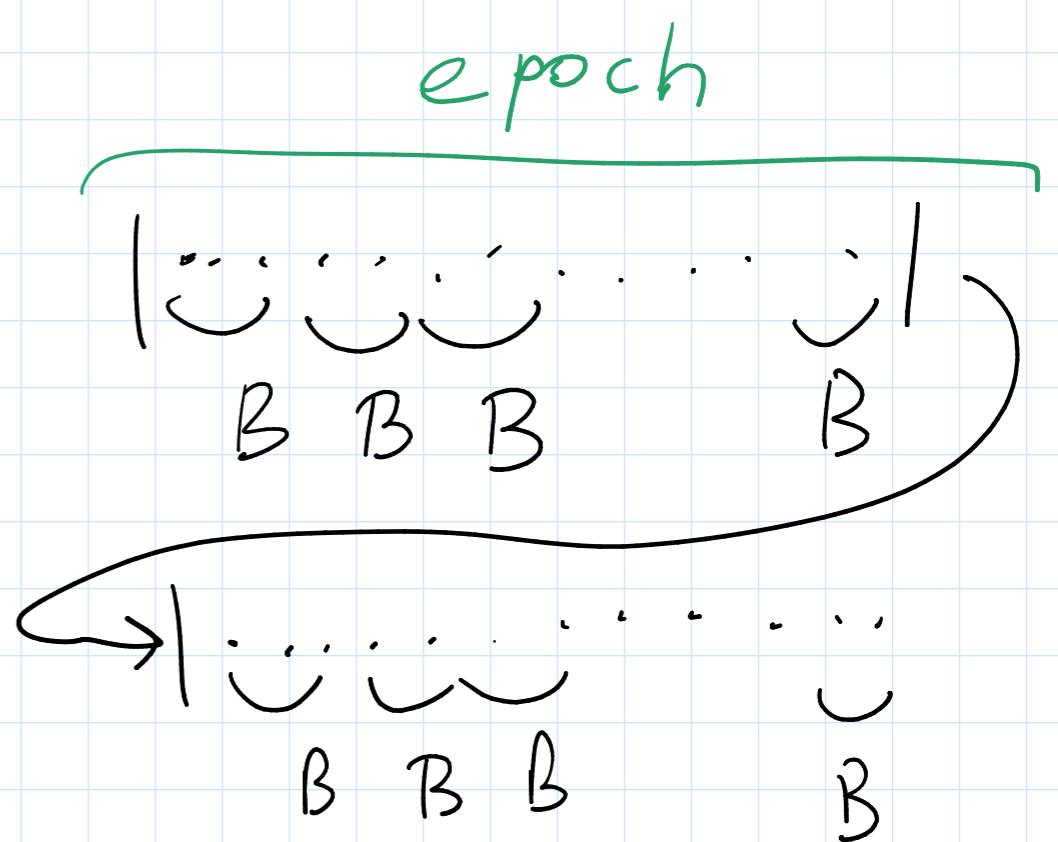
T - number of steps



- Warmup LR



first few epochs



- Reduce LR on plateau

