

# Optimization

GD

SGD

mini batch SGD

$x, y$  - example

$B \gg 1$

$B \rightarrow$  batch size

$L$  - loss function

$w(\theta)$  - parameters

$\alpha_w(x) \rightarrow \hat{y}$

$L = \frac{1}{B} \sum_{i=1}^B (y_i - \hat{y}_i)^2 \rightarrow \min_w$

$\hat{y}_i = \alpha_w(x_i)$

$t$  - time

$L^t$

SGD:

$g_t = \nabla_w L^t(w_{t-1})$

$w_t = w_{t-1} - \eta_t g_t$

$\eta_t = \eta_r$

$m_t = \mu \cdot m_{t-1} + g_t$

$w_t = w_{t-1} - \eta_t m_t$

momentum

$\eta_r \sim 10^{-2} \rightarrow 10^{-7} = \eta$

$m_0 = 0$

$0 \leq \mu \leq 1$   $\mu = 0.8 - 0.9$

typ. value

$L \rightarrow L^* = L + \lambda \|w\|_2^2$

$\nabla_w \|w\|_2^2 = 2w$

$g_t \rightarrow \tilde{g}_t^* = \nabla_w L(w_{t-1}) + \lambda \underbrace{2w_{t-1}}_{\lambda}$

$\lambda$  PyTorch

weight-decay

Adam

(AdaGrad, RMSProp)

1) Momentum

2) Adaptive learning rate

$g_t = \nabla_w L(w_{t-1}) + \lambda w_{t-1}$

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

$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \rightarrow$  элемент  $m$  и  $v$

$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$   $\leftarrow$  сглаживаем

$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$

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

$\epsilon \rightarrow$  ЗАЩИТА ОТ  $\div 0$

$m_0 = v_0 = 0$

$\beta_1 = 0.9$

$\beta_2 = 0.999$

$\epsilon = 10^{-8.9}$

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

$= (1 - \beta_1) g_t + \beta_1 (1 - \beta_1) g_{t-1} + \beta_1^2$

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

$+ \dots + \beta_1^{t-1} (1 - \beta_1) g_1 + \beta_1^t m_0 =$

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

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

$= \frac{1 - \beta_1}{(1 - \beta_1) \sum_{k=0}^{t-1} \beta_1^k} \sum_{i=0}^{t-1} \beta_1^i g_{t-i} \quad \square$

$(1 - \beta_1)(1 + \beta_1 + \beta_1^2) \ominus$

$\ominus 1 + \beta_1 + \beta_1^2 - \beta_1 - \beta_1^2 - \beta_1^3$

$\square \frac{\sum_{i=0}^{t-1} \beta_1^i g_{t-i}}{\sum_{k=0}^{t-1} \beta_1^k} = \left\{ \alpha_i = \frac{\beta_1^i}{\sum_{k=0}^{t-1} \beta_1^k} \right\}$

$\alpha_i > 0, \sum_{i=0}^{t-1} \alpha_i = 1 \Rightarrow \sum_{i=0}^{t-1} \alpha_i g_{t-i} =$

$= E g$

$\hat{v}_t = E g^2 = Var(g) + (E g)^2$

$1: Var(g) \rightarrow 0$

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

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

2.  $Var(g) \neq 0$

$\frac{E g}{\sqrt{Var(g) + (E g)^2 + \epsilon}}$

$\uparrow$   $\Rightarrow \eta_t \cdot *$

быстро  $*$   $< 1$

AdamW

$g_t = \nabla_w L^t(w_{t-1})$  X

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

$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$

$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$

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

$(1 - \eta_t \lambda)$

weight decay

$\frac{\eta_t \lambda}{\sqrt{\hat{v}_t}}$

L BFGS - B

$\eta_t$

Constant

$\eta_t = \eta_0$

Step LR

$\eta_0: 1 \leq t \leq t_0$

$\eta_1: t_0 \leq t \leq t_1$

$\vdots$

Exponential LR

$\eta_t = \gamma \eta_{t-1}$

$\gamma < 1$

Cosine LR

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

$T = \text{epochs}$

1, 2, 3, 4

Linear warmup

450 500 600

warmup

Reduce LR on plateau

val loss

$LR \downarrow$

$LR \downarrow$

$t$

Cosine with Restarts

$t$