

神经网络参数优化器

待优化参数 $w$ , 损失函数 $loss$ , 学习率 $lr$ , 每次迭代一个 $batch$ ,  $t$ 表示当前 $batch$ 迭代的总次数:

1. 计算 $t$ 时刻损失函数关于当前参数的梯度 $g_t = \nabla loss = \frac{\partial loss}{\partial (w_t)}$

2. 计算 $t$ 时刻一阶动量 $m_t$ 和二阶动量 $V_t$

3. 计算 $t$ 时刻下降梯度:  $\eta_t = lr \cdot m_t / \sqrt{V_t}$

4. 计算 $t+1$ 时刻参数:  $w_{t+1} = w_t - \eta_t = w_t - lr \cdot m_t / \sqrt{V_t}$

一阶动量: 与梯度相关的函数

二阶动量: 与梯度平方相关的函数

神经网络参数优化器:

待优化参数  $w$

损失函数  $loss$

学习率  $lr$

每次迭代一个  $batch$

$t$ 表示当前 $batch$ 迭代的总次数

注: 一阶动量: 梯度相关函数

二阶动量: 梯度平方相关函数

1. 计算 $t$ 时刻  $loss$  对  $w$  的梯度  $g_t = \nabla loss = \frac{\partial loss}{\partial (w_t)}$

2. 计算 $t$ 时刻一阶动量  $m_t$ , 二阶动量  $V_t$

3. 计算 $t$ 时刻下降梯度:  $\eta_t = lr * \frac{m_t}{\sqrt{V_t}}$

4. 计算 $t+1$ 时刻参数:  $w_{t+1} = w_t - \eta_t = w_t - lr * \frac{m_t}{\sqrt{V_t}}$

优化器:

SGD (无动量) base

$$m_t = g_t \quad V_t = 1 \quad \eta_t = lr * \frac{m_t}{\sqrt{V_t}} = lr * g_t$$

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

SGDM (SGD + 一阶动量)

$$m_t = \beta \cdot m_{t-1} + (1-\beta) \cdot g_t \quad V_t = 1 \quad \eta_t = lr * \frac{m_t}{\sqrt{V_t}} = lr * (\beta \cdot m_{t-1} + (1-\beta) \cdot g_t)$$

$$w_{t+1} = w_t - \eta_t = w_t - lr * (\beta \cdot m_{t-1} + (1-\beta) \cdot g_t)$$

Adagrad (SGD + 二阶动量)

$$m_t = g_t \quad V_t = \sum_{\tau=1}^t g_{\tau}^2 \text{ 开始到现在梯度的累积和} \quad \eta_t = lr * \frac{m_t}{\sqrt{V_t}} = lr * \frac{g_t}{\sqrt{\sum_{\tau=1}^t g_{\tau}^2}}$$

$$w_{t+1} = w_t - \eta_t = w_t - lr * \frac{g_t}{\sqrt{\sum_{\tau=1}^t g_{\tau}^2}}$$

RMSProp (SGD + 二阶动量)

$$m_t = g_t \quad V_t = \beta \cdot V_{t-1} + (1-\beta) \cdot g_t^2 \quad \eta_t = lr * \frac{m_t}{\sqrt{V_t}} = lr * \frac{g_t}{\sqrt{\beta \cdot V_{t-1} + (1-\beta) \cdot g_t^2}}$$

$$w_{t+1} = w_t - \eta_t = w_t - lr * \frac{g_t}{\sqrt{\beta \cdot V_{t-1} + (1-\beta) \cdot g_t^2}}$$

> Adam (同时结合 SGDM-一阶动量 和 RMSProp-二阶动量)

$$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t \quad \text{修正一阶变量的偏差: } \hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$V_t = \beta_2 \cdot V_{t-1} + (1 - \beta_2) \cdot g_t^2 \quad \text{修正二阶变量的偏差: } \hat{V}_t = \frac{V_t}{1 - \beta_2^t}$$

$$\begin{aligned} \eta_t &= \ln^* \frac{\hat{m}_t}{\sqrt{\hat{V}_t}} = \ln^* \frac{m_t}{1 - \beta_1^t} \bigg/ \sqrt{\frac{V_t}{1 - \beta_2^t}} \\ &= \ln^* \frac{\beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t}{1 - \beta_1^t} \bigg/ \sqrt{\frac{\beta_2 \cdot V_{t-1} + (1 - \beta_2) \cdot g_t^2}{1 - \beta_2^t}} \end{aligned}$$

$$W_{t+1} = W_t - \eta_t = W_t - \ln^* \frac{\frac{m_t}{1 - \beta_1^t}}{\sqrt{\frac{V_t}{1 - \beta_2^t}}}$$