

网络优化正则化

1. 小批量梯度下降

选 k 个训练样本, $t_k = \{x^k, y^k\}_{k=1}^k$

$$g_t(\theta) = \frac{1}{k} \sum_{(x^k, y^k) \in t_k} \frac{\partial L(y^k, f(x^k; \theta))}{\partial \theta}$$

$$g_t \triangleq g_t(\theta_{t+1})$$

$$\theta_{t+1} \leftarrow \theta_t - \alpha g_t \quad \alpha > 0 \text{ 学习率}$$

$$\Delta \theta_t \triangleq \theta_t - \theta_{t-1}$$

$\Delta \theta_t$ 和梯度 g_t 不需要完全一致

① 合适的 k (批量大小)

② 学习率 α

③ 参数更新方向

① 批量大小:

影响随机梯度的期望, 但影响方差:

批量小 \rightarrow 方差越大 \rightarrow 噪声越大 \rightarrow 训练越慢

\rightarrow 设置较小的学习率

学习率和批量满足线性缩放

② 学习率调整

学习率衰减 1) 分段常数衰减

2) 逆时衰减, $\alpha_t = \alpha_0 \frac{1}{1+\beta t}$

3) 指数衰减

$$\alpha_t = \alpha_0 \beta^t$$

4) 指数指数衰减

$$\alpha_t = \alpha_0 \exp(-\beta t)$$

5) 余弦

$$\alpha_t = \frac{1}{2} \alpha_0 (1 + \cos(\frac{t}{T}))$$

让学习率预热

参数随机

梯度较大

不收敛

所以逐渐减小学习率

周期性训练学习率, 避免陷入局部最小值

虽然短期损害优化过程

但长期收敛至全局最优值

$$m = \lfloor 1 + \frac{2}{\Delta T} \rfloor$$

循环学习率(后面)

超参数的随机梯度下降

$$\alpha_t = \alpha_{\min}^m + (\alpha_{\max}^m - \alpha_{\min}^m) (\max(0, 1 - \frac{t}{m}))$$

$$b = 1 - \frac{t}{\Delta T} - 2m + 1$$

1) Adagrad 算法

不同参数的收敛情况分别设置

$$G_t = \sum_{\tau=1}^t g_{\tau} \odot g_{\tau}$$

$$\Delta \theta_t = - \frac{\alpha}{\sqrt{G_t + \epsilon}} \odot g_t$$

$\left\{ \begin{array}{l} G_t \uparrow \quad \Delta \theta_t \downarrow \\ G_t \downarrow \quad \Delta \theta_t \uparrow \end{array} \right.$

整体都是随着迭代次数增加学习率减小

缺点: 一定次数迭代后若没有收敛, 由于此时很小, 将继续再收敛

即学习率下降

2) RMSprop

避免 Adagrad 中学习率单调下降过早衰减

$$G_t = \beta G_{t-1} + (1-\beta) g_t \odot g_t$$

$$= (1-\beta) \sum_{\tau=1}^t \beta^{t-\tau} g_{\tau} \odot g_{\tau}$$

指数移动平均

$$\Delta \theta_t = - \frac{\alpha}{\sqrt{G_t + \epsilon}} \odot g_t$$

3) Adadelta 算法

$$\Delta X_{t+1}^2 = \beta_1 \Delta X_t^2 + (1-\beta_1) \Delta \theta_{t+1} \cdot \Delta \theta_{t+1}$$

通过梯度平方而指数衰减移动平均
来调整学习率 G_t 和 RMSprop 一样

$$\Delta \theta_t = - \frac{\sqrt{\Delta X_{t+1}^2}}{\sqrt{G_t + \epsilon}} g_t$$

④ 更新优化方向

动量法:

$$\begin{aligned} \Delta \theta_t &= \rho \Delta \theta_{t+1} - \alpha g_t \\ &= -\alpha \sum_{\tau=1}^t \rho^{t-\tau} g_\tau \end{aligned}$$

负梯度的“加权移动平均”
 ρ 为动量因子。
 α 为学习率

2) Nesterov 加速梯度

动量法可拆分为两步

$$\hat{\theta} = \theta_{t+1} + \rho \Delta \theta_{t+1}$$

$$\theta_t = \hat{\theta} - \alpha g_t$$

不是在 $\hat{\theta}$ 点的梯度

$$\Delta \theta_t = \rho \Delta \theta_{t+1} - \alpha g_t(\theta_{t+1} + \rho \Delta \theta_{t+1})$$

3) Adam 算法

同时更新学习率和梯度。

RMSprop + 动量法

$$\Delta \theta_t = - \frac{\alpha}{\sqrt{G_t + \epsilon}} \hat{M}_t \quad (\alpha = 0.001, \alpha_t = \frac{\alpha}{\sqrt{t}})$$

$$\begin{cases} G_t = \beta_2 G_{t+1} + (1-\beta_2) g_t \cdot g_t & \beta_2 = 0.99 \\ M_t = \beta_1 M_{t+1} + (1-\beta_1) g_t \end{cases}$$

$$\hat{M}_t = \frac{M_t}{1-\beta_1} \quad \hat{G}_t = \frac{G_t}{1-\beta_2}$$

2. 参数初始化

全为 0 \rightarrow 对称权重

随机初始化 $\begin{cases} \text{太小: } \begin{cases} \text{多层网络容易消失} \\ \text{sigmoid 主非线性} \end{cases} \\ \text{太大: } \text{sigmoid 梯度为 0} \end{cases}$

保持每个神经元的输入输出方差一致

$\begin{cases} \text{高斯分布} \\ \text{均匀分布} \end{cases}$

类别

优化算法

$\begin{cases} \text{固定衰减} \alpha & \text{分段常数, 逐时, 指数, 余弦} \\ \text{周期} \alpha & \text{自适应 SGDR} \\ \text{自适应} \alpha & \text{Adagrad, RMSprop, Adadelta} \end{cases}$

梯度方向 (因此) 动量法, Nesterov, 梯度截断

综合

Adam

$$\Delta \theta_t = - \frac{\alpha_t}{\sqrt{G_t + \epsilon}} M_t$$

$$G_t = \psi(g_1, \dots, g_t)$$

$$M_t = \phi(g_1, \dots, g_t)$$

Xavier 初始化

神经元数量计算方差

$$\text{高斯: } N(0, \sqrt{\frac{2}{n^{in} + n^{out}}})$$

$$\text{均匀: } r = \sqrt{\frac{6}{n^{in} + n^{out}}}$$