

算法 8.5 RMSProp 算法

AdaGrad 的变种

Require: 全局学习速率 ϵ , 衰减速率 ρ

Require: 初始参数 θ

Require: 小常数 δ , 通常设为 10^{-6} (用于被小数除时的数值稳定)

初始化累积变量 $r = 0$

while 没有达到停止准则 **do**

从训练集中采包含 m 个样本 $\{x^{(1)}, \dots, x^{(m)}\}$ 的 minibatch, 对应目标为 $y^{(i)}$ 。

计算梯度: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

累积平方梯度: $r \leftarrow \rho r + (1 - \rho) g \odot g$ 只保留一定阶段的历史。

计算参数更新: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + r}} \odot g$ ($\frac{1}{\sqrt{\delta + r}}$ 逐元素应用)

应用更新: $\theta \leftarrow \theta + \Delta \theta$

end while

下, 它也许最好被看作结合 RMSProp 和具有一些重要区别的动量的变种。首先, 在 Adam 中, 动量直接并入了梯度一阶矩 (带指数加权) 的估计。将动量加入 RMSProp 最直观的方法是应用动量于缩放后的梯度。结合重放缩的动量使用没有明确的理论动机。其次, Adam 包括负责原点初始化的一阶矩 (动量项) 和 (非中心的) 二阶矩的估计修正偏置 (算法 8.7)。RMSProp 也采用了 (非中心的) 二阶矩估计, 然而缺



