

22/9/23

Deep Learning

◦ Recap: SGD, Momentum, Nesterov Momentum ✓

Reference: Chapter 8 in DL book by Goodfellow et al.

◦ Adaptive Learning Rate Algorithms

- Adaptive Gradient (Ada Grad) ✓

- RMS Prop ✓

- Adaptive Moment (Adam) ✓

◦ Recap - SGD: $R_{SGD}(\theta) = \sum_{i=1}^m d(y^i, \hat{y}^i)$; m samples out of n training samples

- Momentum: $v^{(r)} = \alpha v^{(r-1)} - \eta^{(r)} \nabla_{\theta^{(r-1)}} R_{SGD}(\theta)$ ($0 < \alpha < 1$)

$$\theta^{(r)} = \theta^{(r-1)} + v^{(r)}$$

- Nesterov momentum: $v^{(r)} = \alpha v^{(r-1)} - \eta^{(r)} \nabla_{\theta^{(r-1)}} R_{SGD}(\theta^{(r-1)} + \alpha v^{(r-1)})$

$$\theta^{(r)} = \theta^{(r-1)} + v^{(r)}$$

◦ Adaptive Learning Rate Algorithms:

- Adaptive Gradient (Ada Grad): Smaller steps along steeper directions & bigger steps along flatter directions

$$h^{(0)} = 0$$

→ cumulative gradient square.

$$h^{(r)} = h^{(r-1)} + \nabla_{\theta^{(r-1)}} \odot \nabla_{\theta^{(r-1)}} \quad \odot: \text{element-wise product}$$

$$\Delta \theta^{(r)} = \left(\frac{-\eta^{(r)}}{\delta + \sqrt{h^{(r)}}} \right) \odot \nabla_{\theta^{(r-1)}} \quad \rightarrow \text{adaptive learning rate}$$

$$\theta^{(r)} = \theta^{(r-1)} + \Delta \theta^{(r)} \quad \rightarrow \delta: \text{small stabilizing constant}$$

→ Element-wise division

Drawback: Since $h^{(r)}$ is cumulative, older noisy gradients could adversely impact the learning rate

- RMS Prop: To address this issue, RMS Prop proposes a weighting strategy

$$h^{(0)} = 0; \quad 0 \leq \rho < 1$$

$$h^{(r)} = \rho \cdot h^{(r-1)} + (1-\rho) \cdot \nabla_{\theta^{(r-1)}} \odot \nabla_{\theta^{(r-1)}} \quad \leftarrow \text{damps older gradients more.}$$

$$\Delta \theta^{(r)} = \frac{-\eta^{(r)}}{\sqrt{\delta + h^{(r)}}} \odot \nabla_{\theta^{(r-1)}}$$

$$\theta^{(r)} = \theta^{(r-1)} + \Delta \theta^{(r)}$$

- Adam: Combines momentum idea with RMSProp to further improve performance

$$\circ g^{(0)} = h^{(0)} = 0; \quad \beta_1, \beta_2 \in [0, 1)$$

$$\circ g^{(r)} = \beta_1 g^{(r-1)} + (1-\beta_1) \cdot \nabla_{\theta^{(r-1)}} \quad (\text{biased estimate of "first moment"})$$

$$\circ h^{(r)} = \beta_2 h^{(r-1)} + (1-\beta_2) \cdot \nabla_{\theta^{(r-1)}} \odot \nabla_{\theta^{(r-1)}} \quad (\text{biased "second moment"})$$

$$\circ \hat{g}^{(r)} = \frac{g^{(r)}}{(1-\beta_1^r)} \quad (\text{bias correction})$$

$\rightarrow \beta_1$ raised to the power r

$$\circ \hat{h}^{(r)} = \frac{h^{(r)}}{(1-\beta_2^r)}$$

$\rightarrow \beta_2$ raised to the power r

$$\circ \Delta \theta^{(r)} = \left(\frac{-\eta^{(r)}}{\delta + \sqrt{\hat{h}^{(r)}}} \right) \odot \hat{g}^{(r)} \quad \rightarrow \text{adaptive learning rate}$$

$$\circ \theta^{(r)} = \theta^{(r-1)} + \Delta \theta^{(r)}$$

◦ Second-order Taylor expansion of $R(\theta)$ about $(\theta + \Delta \theta)$ for univariate case

$$R(\theta + \Delta \theta) = R(\theta) + \Delta \theta \cdot R'(\theta) + \frac{(\Delta \theta)^2}{2} R''(\theta) + \dots$$