

31/10/23

Deep Learning

◦ Recap: RNN

◦ Back Propagation Through Time (BPTT)

◦ Gated Recurrent Units (GRUs)

References:

Chapters 9 & 10 in
Dive Into Deep Learning

$$\begin{aligned} \circ \text{RNN:} \quad & \circ \underline{z}_m^{(t)} = \sigma(\underline{\alpha}_m^T \underline{x}^{(t)} + \underline{\gamma}_m^T \underline{z}^{(t-1)}) \quad \text{--- (1)} \\ & \circ \hat{y}_k^{(t)} = g_k(\underline{\beta}_k^T \underline{z}^{(t)}) \quad \text{--- (2)} \\ & \circ \Theta = \{ \underline{\alpha}_1, \underline{\alpha}_2, \dots, \underline{\alpha}_M, \underline{\beta}_1, \underline{\beta}_2, \dots, \underline{\beta}_K, \\ & \quad \underline{\gamma}_1, \underline{\gamma}_2, \dots, \underline{\gamma}_M \} \end{aligned}$$

Aside:

$$p(y^{(t)} | x^{(1)}, \dots, x^{(t)})$$

$$\propto p(y^{(t)} | x^{(t)}, \underline{z}^{(t-1)})$$

◦ Just as with ANNs, we use an iterative approach to update the model parameters

$$\Theta^{(r)} = \Theta^{(r-1)} - \eta \underbrace{\nabla_{\Theta^{(r-1)}} \mathcal{L}(\Theta)}_{\text{where}}$$

$$\circ \mathcal{D} = \{ (\underline{x}^{(1)}, y^{(1)}), (\underline{x}^{(2)}, y^{(2)}), \dots, (\underline{x}^{(t)}, y^{(t)}), \dots, (\underline{x}^{(T)}, y^{(T)}) \}$$

$$\circ \mathcal{R}(\Theta) = \frac{1}{T} \sum_{t=1}^T d(y^{(t)}, \hat{y}^{(t)})$$

$$\circ \frac{\partial \mathcal{R}(\Theta)}{\partial \gamma_{mp}} = \frac{1}{T} \sum_{t=1}^T \frac{\partial d(y^{(t)}, \hat{y}^{(t)})}{\partial \hat{y}^{(t)}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial z_m^{(t)}} \cdot \frac{\partial z_m^{(t)}}{\partial \gamma_{mp}} \quad \text{--- (3)}$$

$$\circ \frac{\partial z_m^{(t)}}{\partial \gamma_{mp}} = \frac{\partial \sigma(\underline{\alpha}_m^T \underline{x}^{(t)} + \underline{\gamma}_m^T \underline{z}^{(t-1)})}{\partial \gamma_{mp}} + \frac{\partial [\sigma(\underline{\alpha}_m^T \underline{x}^{(t)} + \underline{\gamma}_m^T \underline{z}^{(t-1)})]}{\partial \underline{z}_m^{(t-1)}} \cdot \frac{\partial \underline{z}_m^{(t-1)}}{\partial \gamma_{mp}}$$

$$\circ \frac{\partial z_m^{(t)}}{\partial \gamma_{mp}} = \frac{\partial \sigma(\underline{\alpha}_m^T \underline{x}^{(t)} + \underline{\gamma}_m^T \underline{z}^{(t-1)})}{\partial \gamma_{mp}} +$$

$$\left[\frac{\partial f(\underline{x}^{(t)}, y^{(t)})}{\partial \underline{z}} \right]$$

$$\sum_{i=1}^{t-1} \frac{\partial z_m^{(i)}}{\partial \gamma_{mp}} \cdot \left[\prod_{j=i+1}^t \frac{\partial \sigma(\underline{\alpha}_m^T \underline{x}^{(j)} + \underline{\gamma}_m^T \underline{z}^{(j-1)})}{\partial \underline{z}_m^{(j-1)}} \right] \quad \text{--- (4)}$$

• Plugging (4) into (3) gives us $\frac{\partial R(t)}{\partial \gamma_{mp}}$

• Takeaway: * Computing $\frac{\partial R(t)}{\partial \gamma_{mp}}$ is expensive especially when t (or T) is large.

* We run into the issue of vanishing/exploding gradients

• A workaround: Truncated BPTT. Only a fixed number of time steps used in gradient computation.

• Gated Recurrent Units (GRUs)

$$z_m^{(t)} = \underbrace{s_m^{(t)}}_{\text{update gate}} z_m^{(t-1)} + (1 - s_m^{(t)}) \underbrace{\tilde{z}_m^{(t)}}_{\text{candidate } \tilde{z}_m}$$

$$\tilde{z}_m^{(t)} = \sigma(\alpha_m^T x^{(t)} + \underbrace{r_m^{(t)} \gamma_m^T}_{\text{reset gate}} z^{(t-1)})$$