

计算 LSTM 的梯度

16337341 朱志儒

设 LSTM 的损失函数为 $L(t)$, 有

$$L(t) = \sum_{t=1}^T L^{(t)}$$

对于两个隐藏状态 $h^{(t)}$ 和 $C^{(t)}$ 有

$$\delta_h^{(t)} = \frac{\partial L}{\partial h^{(t)}} = \left(\frac{\partial O^{(t)}}{\partial h^{(t)}} \right)^T \frac{\partial L^{(t)}}{\partial O^{(t)}} = V^T (\widehat{y^{(t)}} - y^{(t)})$$

$$\delta_c^{(t)} = \frac{\partial L}{\partial C^{(t)}} = \left(\frac{\partial h^{(t)}}{\partial C^{(t)}} \right)^T \frac{\partial L^{(t)}}{\partial h^{(t)}} = \delta_h^{(t)} \odot o^{(t)} \odot (1 - \tanh^2(C^{(t)}))$$

接着由 $\delta_c^{(t+1)}, \delta_h^{(t+1)}$ 反向推导 $\delta_c^{(t)}, \delta_h^{(t)}$ 。

$\delta_h^{(t)}$ 的梯度由本层 t 时刻的输出梯度误差和大于 t 时刻的误差两部分决定, 即:

$$\delta_h^{(t)} = \frac{\partial L}{\partial h^{(t)}} = \frac{\partial l(t)}{\partial h^{(t)}} + \left(\frac{\partial h^{(t+1)}}{\partial h^{(t)}} \right)^T \frac{\partial L^{(t+1)}}{\partial h^{(t+1)}} = V^T (\widehat{y^{(t)}} - y^{(t)}) + \left(\frac{\partial h^{(t+1)}}{\partial h^{(t)}} \right)^T \delta_h^{(t+1)}$$

$$\begin{aligned} \frac{\partial h^{(t+1)}}{\partial h^{(t)}} &= W_o^T [o^{(t+1)} \odot (1 - o^{(t+1)}) \odot \tanh(C^{(t+1)})] \\ &\quad + W_f^T [\Delta C \odot f^{(t+1)} \odot (1 - f^{(t+1)}) \odot C^{(t)}] + W_a^T \{\Delta C \odot i^{(t+1)} \\ &\quad \odot [1 - (a^{(t+1)})^2]\} + W_i^T [\Delta C \odot a^{(t+1)} \odot i^{(t+1)} \odot (1 - i^{(t+1)})] \end{aligned}$$

而 $\delta_c^{(t)}$ 的反向梯度误差由前一层 $\delta_c^{(t+1)}$ 的梯度误差和本层的从 $h^{(t)}$ 传回来的梯度误差两部分组成, 即:

$$\begin{aligned} \delta_c^{(t)} &= \left(\frac{\partial C^{(t+1)}}{\partial C^{(t)}} \right)^T \frac{\partial L}{\partial C^{(t+1)}} + \left(\frac{\partial h^{(t)}}{\partial C^{(t)}} \right)^T \frac{\partial L}{\partial h^{(t)}} \\ &= \left(\frac{\partial C^{(t+1)}}{\partial C^{(t)}} \right)^T \delta_c^{(t+1)} + \delta_h^{(t)} \odot o^{(t)} \odot (1 - \tanh^2(C^{(t)})) \\ &= \delta_c^{(t+1)} \odot f^{(t+1)} + \delta_h^{(t)} \odot o^{(t)} \odot (1 - \tanh^2(C^{(t)})) \end{aligned}$$

有了 $\delta_c^{(t)}$ 和 $\delta_h^{(t)}$, 计算 W_f 的梯度计算过程:

$$\frac{\partial L}{\partial W_f} = \sum_{t=1}^T [\delta_c^{(t)} \odot C^{(t-1)} \odot f^{(t)} \odot (1 - f^{(t)})] (h^{(t-1)})^T$$