

Part 1:

1) 由 BPTT 相关性质, RNN 中参数的更新要以求出损失函数 L 关于这些参数的梯度为前提. 假设当前时刻为 t ($t > 1$).

$$\Rightarrow \nabla_{\mathbf{w}} L = \sum_{t=1}^T \nabla_{\mathbf{w}} L^{(t)} h^{(t)T}$$

其中 $\nabla_{\mathbf{w}} L = \frac{\partial L}{\partial \hat{y}^{(t)}} \otimes (\hat{y}^{(t)} - y)$, \otimes 是逐元素相乘符号.

$$\Rightarrow \nabla_{\mathbf{w}} L = \sum_{t=1}^T (\nabla_{\mathbf{h}} L^{(t)}) \frac{\partial h^{(t)}}{\partial \mathbf{h}^{(t-1)}}$$

若 $t = 1$, 则 $\nabla_{\mathbf{h}} L = \nabla_{\mathbf{h}^{(1)}} L = V^T \nabla_{\mathbf{w}} L$

若 $t \neq 1$, 则 $\nabla_{\mathbf{h}} L = \left(\frac{\partial h^{(t+1)}}{\partial h^{(t)}} \right)^T (\nabla_{\mathbf{h}^{(t+1)}} L) + \left(\frac{\partial O^{(t)}}{\partial h^{(t)}} \right)^T (\nabla_{\mathbf{w}} L)$
 $= W^T (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag}(1 - (h^{(t+1)})^2) + V^T (\nabla_{\mathbf{w}} L)$

其中 $\frac{\partial h^{(t)}}{\partial h^{(t-1)}} = W^T \text{diag}(1 - (h^{(t-1)})^2)$

$$\Rightarrow \nabla_{\mathbf{w}} L = \sum_{t=1}^T \text{diag}(1 - (h^{(t)})^2) (\nabla_{\mathbf{h}^{(t)}} L) h^{(t-1)T}$$

$$\Rightarrow \nabla_{\mathbf{w}} L = \sum_{t=1}^T \frac{\partial L}{\partial h} \nabla_{\mathbf{w}} h^{(t)}$$

$$= \sum_{t=1}^T \text{diag}(1 - (h^{(t)})^2) (\nabla_{\mathbf{h}^{(t)}} L) x^{(t)T}$$

2) 若使用 sigmoid 函数和 tanh 函数时可能会产生梯度消失问题.

因为这两个激活函数的导数均小于 1,

在沿时间反传传播误差的过程中, 反传值将会趋近于 0.

这会使前面的信息无法有效被利用.