

LSTM:

Forward Pass:

$$F_t = \sigma(w_{Fi} [A_{t-1}, x_t])$$

$$\tilde{C}_t = \tanh(w_{Ct} [h_{t-1}, x_t])$$

$$i_t = \sigma(w_{It} [h_{t-1}, x_t])$$

$$O_t = \sigma(w_{Ot} [h_{t-1}, x_t])$$

$$h_t = O_t \otimes \tanh(L_t)$$

$$\text{cell state } C_t = C_{t-1} \otimes F_t \oplus \tilde{C}_t \otimes i_t$$

Let ϵ_t be the error for time step $t \in T$

$$\frac{\partial \epsilon_T}{\partial w} = \frac{\partial \epsilon_T}{\partial C_T} \frac{\partial C_T}{\partial C_{T-1}} \frac{\partial C_{T-1}}{\partial C_{T-2}} \dots \frac{\partial C_1}{\partial w}$$

$$\frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial (C_t \otimes F_t \oplus \tilde{C}_t \otimes i_t)}{\partial C_{t-1}}$$

$$= \frac{\partial F_t}{\partial C_{t-1}} \cdot C_{t-1} \otimes F_t + \frac{\partial \tilde{C}_t}{\partial C_{t-1}} \cdot \tilde{C}_t + i_t \frac{\partial \tilde{C}_t}{\partial C_{t-1}}$$

$$\frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial (\sigma(w_{Fi} [h_{t-1}, x_t]))}{\partial C_{t-1}} \cdot C_{t-1} \otimes F_t + \frac{\partial (\sigma(w_{It} [h_{t-1}, x_t]))}{\partial C_{t-1}} \cdot \tilde{C}_t + \frac{\partial (\sigma(w_{Ot} [h_{t-1}, x_t]))}{\partial C_{t-1}} \cdot i_t$$

$$\frac{\partial \epsilon_T}{\partial w} = \frac{\partial \epsilon_T}{\partial C_T} \left(\prod_{t=2}^{T-1} \left(A_t + B_t + C_t + \epsilon_t \right) \right) \frac{\partial C_1}{\partial w}$$

Vanishing gradient is less likely ^{in LSTMs} due to there being ^{sum} ~~terms~~ of 4 terms.