

Vanishing gradients solved by LSTM.

$$\frac{\partial E_t}{\partial w} = \sum_{j=0}^t \frac{\partial E_t}{\partial h_j} \times \frac{\partial h_j}{\partial w} \times \frac{\partial h_t}{\partial h_j} \times \frac{\partial h_i}{\partial w}$$

$$\frac{\partial h_t}{\partial h_i} = \frac{\partial h_t}{\partial h_{t-1}} \times \frac{\partial h_{t-1}}{\partial h_{t-2}} \times \dots \times \frac{\partial h_{t+1}}{\partial h_i}$$

$$= \prod_{k=i}^{t-1} \frac{\partial h_{k+1}}{\partial h_k}$$

$$\frac{\partial E_t}{\partial w} = \frac{\partial E_t}{\partial h_t} \times \frac{\partial h_t}{\partial c_t} \times \frac{\partial c_t}{\partial c_{t-1}} \times \dots \times \frac{\partial c_2}{\partial c_1} \times \frac{\partial c_1}{\partial w}$$

$$= \frac{\partial E_t}{\partial h_t} \times \frac{\partial h_t}{\partial c_t} \times \left(\prod_{i=2}^t \frac{\partial c_i}{\partial c_{i-1}} \right) \times \frac{\partial c_1}{\partial w}$$

$$C_t = f_t \times C_{t-1} + i_t \times \tilde{C}_t$$

This is the key difference. The C_t is a function of 4 elements

$$f_t, C_{t-1}, i_t, \tilde{C}_t$$

Whereas in RNN case it is a function of previous states output function.

You can think that even though C_t is a function of n elements, but in those elements it contains the previous states output function.

Yes it contains but the diff is not all 4 elements are functions of previous outputs.

$$C_t = f_t \times C_{t-1} + i_t \times \tilde{C}_t$$

we want partial derivatives w.r.t C_{t-1}

$$\frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial}{\partial C_{t-1}} [f_t \times C_{t-1} + i_t \times \tilde{C}_t]$$

$$\frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial}{\partial C_{t-1}} [f_t \times C_{t-1}] + \frac{\partial}{\partial C_{t-1}} [\hat{y}_t \times \tilde{C}_t]$$

$$= \underbrace{\frac{\partial f_t}{\partial C_{t-1}} \times (C_{t-1})}_{\text{output } f_t} + \underbrace{\frac{\partial C_{t-1}}{\partial C_{t-1}} \times (f_t) + \frac{\partial \hat{y}_t}{\partial C_{t-1}} \times (\tilde{C}_t) + \frac{\partial \tilde{C}_t}{\partial C_{t-1}} \times (\hat{y}_t)}_{\text{still solve the problem of vanishing gradient}}$$

Even though the others can be 0 still these one f_t element which will solve the problem of vanishing gradient.