

HW2 (I promise tomorrow)

Requested Google Cloud Credits (Share instructions on how to retrieve them)

Review from last time

Long short-term Memory (LSTMs)

↓
main motivation → vanishing gradients issues in RNNs.

$$\vec{h}^{(t-1)}$$

$$\vec{c}^{(t-1)}$$

$$\vec{f}^{(t)} = \sigma(W_f \vec{h}^{(t-1)} + U_f \vec{e}^{(t)} + \vec{b}_f)$$

$$\vec{i}^{(t)} = \sigma(W_i \vec{h}^{(t-1)} + U_i \vec{e}^{(t)} + \vec{b}_i)$$

$$\vec{o}^{(t)} = \sigma(W_o \vec{h}^{(t-1)} + U_o \vec{e}^{(t)} + \vec{b}_o)$$

$$\vec{\tilde{c}}^{(t)} = \tanh(W_c \vec{h}^{(t-1)} + U_c \vec{e}^{(t)} + \vec{b}_c)$$

$$\vec{c}^{(t)} = \vec{f}^{(t)} * \vec{c}^{(t-1)} + \vec{i}^{(t)} * \vec{\tilde{c}}^{(t)}$$

$$\vec{h}^{(t)} = \vec{\sigma}^{(t)} * \tanh(\vec{c}^{(t)})$$

We will fix the issues of vanishing gradients

Let's focus on a LSTM with two time steps.

W_f	U_f	\vec{b}_f
W_i	U_i	\vec{b}_i
W_o	U_o	\vec{b}_o
W_c	U_c	\vec{b}_c

Let's figure out $\frac{\partial L}{\partial W_f}$

$$\frac{\partial L}{\partial W_f} = \frac{\partial L}{\partial \vec{f}^{(1)}} \underbrace{\frac{\partial \vec{f}^{(1)}}{\partial W_f}} + \frac{\partial L}{\partial \vec{f}^{(2)}} \cdot \underbrace{\frac{\partial \vec{f}^{(2)}}{\partial W_f}}_{\text{similar}}$$

we can compute

$$\vec{f}^{(1)} \cdot (1 - \vec{f}^{(1)}) \cdot (\vec{h}^{(0)})^T$$

$$\frac{\delta L}{\delta \vec{f}^{(1)}} = \frac{\delta L}{\delta \vec{c}^{(1)}} \frac{\delta \vec{c}^{(1)}}{\delta \vec{f}^{(1)}}$$

$$\frac{\delta L}{\delta \vec{f}^{(2)}} = \frac{\delta L}{\delta \vec{c}^{(2)}} \frac{\delta \vec{c}^{(2)}}{\delta \vec{f}^{(2)}}$$

$$\frac{\delta L}{\delta w_f} = \underbrace{\frac{\delta L}{\delta \vec{c}^{(1)}}}_{\checkmark} \underbrace{\frac{\delta \vec{c}^{(1)}}{\delta \vec{f}^{(1)}}}_{\checkmark} \underbrace{\frac{\delta \vec{f}^{(1)}}{\delta w_f}}_{\checkmark} + \underbrace{\frac{\delta L}{\delta \vec{c}^{(2)}}}_{\checkmark} \underbrace{\frac{\delta \vec{c}^{(2)}}{\delta \vec{f}^{(2)}}}_{\checkmark} \underbrace{\frac{\delta \vec{f}^{(2)}}{\delta w_f}}_{\checkmark}$$

\downarrow $\vec{c}^{(1)}$
 \downarrow $\vec{c}^{(1)}$

We now need to figure out $\frac{\delta L}{\delta \vec{c}^{(1)}}$ and

$$\frac{\delta L}{\delta \vec{c}^{(2)}}$$

$$\frac{\partial L}{\partial \vec{c}^{(1)}} = \underbrace{\left[\frac{\partial L}{\partial \vec{h}^{(1)}} \right]}_1 \underbrace{\frac{\partial \vec{h}^{(1)}}{\partial \vec{c}^{(1)}}}_{\vec{f}^{(2)}} + \left[\frac{\partial L}{\partial \vec{c}^{(2)}} \right] \underbrace{\frac{\partial \vec{c}^{(2)}}{\partial \vec{c}^{(1)}}}_{\vec{f}^{(2)}}$$

$$0^{(1)} * (1 - \tanh^2(c^{(1)}))$$

$$\begin{aligned} \frac{\partial L}{\partial \vec{c}^{(2)}} &= \frac{\partial L}{\partial \vec{h}^{(2)}} \frac{\partial \vec{h}^{(2)}}{\partial \vec{c}^{(2)}} \\ &= \frac{\partial L}{\partial \vec{y}^{(2)}} \cdot \frac{\partial \vec{y}^{(2)}}{\partial \vec{h}^{(2)}} \cdot \frac{\partial \vec{h}^{(2)}}{\partial \vec{c}^{(2)}} \end{aligned}$$

we're getting a recurrence relationship linking

$$\frac{\partial L}{\partial \vec{c}^{(t)}}$$

to

$$\frac{\partial L}{\partial \vec{c}^{(t+1)}}$$

$$\rightarrow \frac{\partial L}{\partial \vec{c}^{(\text{last } t)}}$$

→ will be
"easy" to
compute

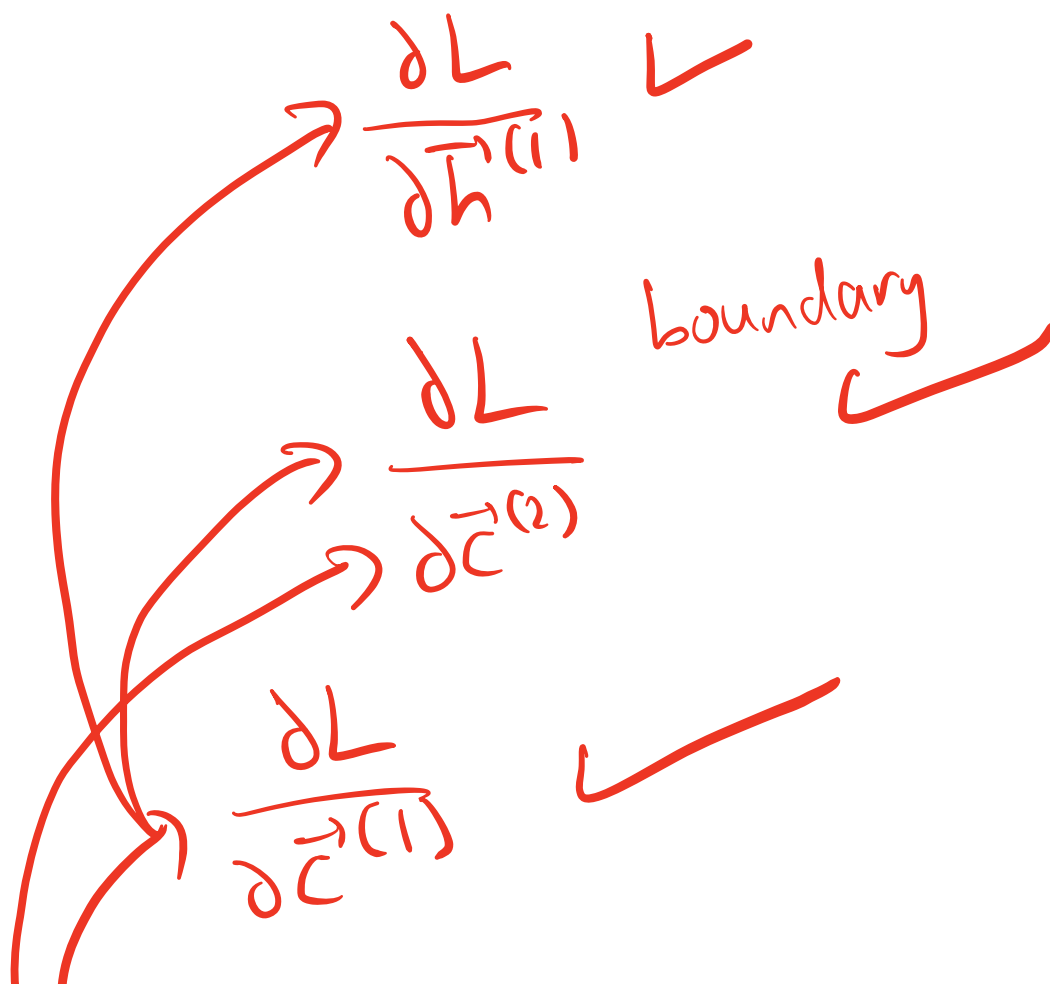
$$\frac{\partial L}{\partial \vec{h}^{(1)}} = \frac{\partial L}{\partial L^{(1)}} \frac{\partial L^{(1)}}{\partial \hat{y}^{(1)}} \frac{\partial \hat{y}^{(1)}}{\partial \vec{h}^{(1)}}$$


$$+ \frac{\partial L}{\partial \vec{h}^{(2)}} \left(\frac{\partial \vec{h}^{(2)}}{\partial \vec{o}^{(2)}} \frac{\partial \vec{o}^{(2)}}{\partial \vec{h}^{(1)}} + \frac{\partial \vec{h}^{(2)}}{\partial \vec{c}^{(2)}} \frac{\partial \vec{c}^{(2)}}{\partial \vec{f}^{(2)}} \frac{\partial \vec{f}^{(2)}}{\partial \vec{h}^{(1)}} \right. \\ \left. + \frac{\partial \vec{h}^{(2)}}{\partial \vec{x}^{(2)}} \frac{\partial \vec{c}^{(2)}}{\partial \vec{x}^{(2)}} \frac{\partial \vec{x}^{(2)}}{\partial \vec{h}^{(1)}} \right)$$

$$+ \frac{\partial \vec{h}^{(2)}}{\partial \vec{c}^{(2)}} \frac{\partial \vec{c}^{(2)}}{\partial \vec{c}^{(1)}} \frac{\partial \vec{c}^{(1)}}{\partial \vec{h}^{(1)}} \Bigg)$$

$$\frac{\partial L}{\partial \vec{h}^{(2)}} = \frac{\partial L}{\partial L^{(2)}} \frac{\partial L^{(2)}}{\partial \vec{y}^{(2)}} \frac{\partial \vec{y}^{(2)}}{\partial \vec{h}^{(2)}}$$

First compute $\frac{\partial L}{\partial \vec{h}^{(2)}}$ ✓ boundary




$$\frac{\partial L}{\partial W_h}$$

we can finally
compute

pytorch does all of that for you

↓
we love pytorch

Suppose we have 10 words
and word 2 is relevant for
word 10

with vanilla RNNs, we would
have gotten vanishing
gradients

$$\vec{c}^{(0)} = 0$$

$$\begin{aligned}\vec{c}^{(1)} &= \beta^{(1)} * c^{(0)} + \lambda^{(1)} * \hat{c}^{(1)} \\ &= 1 * 0 + 0 * \hat{c}^{(1)} \\ &= 0\end{aligned}$$

$$\lambda^{(1)} = 0$$

$$\beta^{(1)} = 1$$

$$\begin{aligned}\vec{c}^{(2)} &= \beta^{(2)} * c^{(1)} + \lambda^{(2)} * \hat{c}^{(2)} \\ &= 1 * 0 + 1 * \hat{c}^{(2)} \\ &= \hat{c}^{(2)}\end{aligned}$$

$$\lambda^{(2)} = 1$$

$$\beta^{(2)} = 1$$

$$\begin{aligned}\vec{c}^{(3)} &= \beta^{(3)} * c^{(2)} + \lambda^{(3)} * \hat{c}^{(3)} \\ &= \hat{c}^{(2)}\end{aligned}$$

$$\lambda^{(3)} = 0$$

$$f^{(3)} \approx 1$$

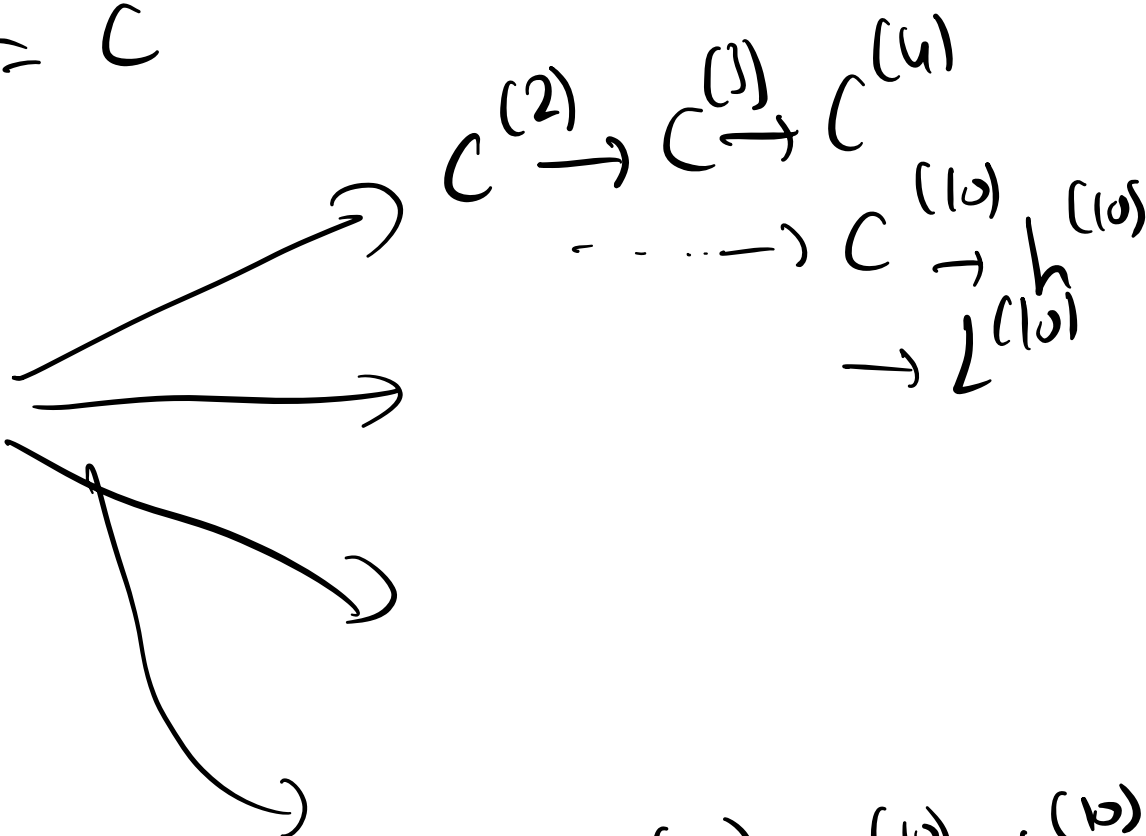
$$\forall t > 3$$

$$i^{(t)} = 0$$

$$f^{(t)} = 1$$

$$C^{(10)} = \tilde{C}^{(2)}$$

$$\frac{\partial L^{(10)}}{\partial C^{(2)}}$$



$$\text{path} = C^{(2)} \rightarrow C^{(3)} \rightarrow \dots \rightarrow C^{(10)} \rightarrow h^{(10)} \rightarrow L^{(10)}$$

$$\frac{\partial L^{(10)}}{\partial C^{(2)}} = \frac{\partial L^{(10)}}{\partial h^{(10)}} \frac{\partial h^{(10)}}{\partial C^{(10)}} \frac{\partial C^{(10)}}{\partial C^{(9)}} \dots \frac{\partial C^{(3)}}{\partial C^{(2)}}$$

$$= \frac{\partial L^{(10)}}{\partial h^{(10)}} * o^{(10)} * (1 - \tanh^2(c^{(10)}))$$

$$* f^{(10)} * f^{(9)} * f^{(8)} \dots * f^{(3)}$$

$$= \frac{\partial L^{(10)}}{\partial h^{(10)}} * o^{(10)} * (1 - \tanh^2(c^{(10)}))$$

$$* \prod_{i=3}^{10} f^{(i)}$$

$$(0.95)^8$$

$$\approx 0.663$$

Transformers

→ attention

"Attention is all you
need"

↓ paper that introduced
transformers

(by Google Folks)