# CS 6180 09/24

HW2 released tomorrow (will be a bit longer so you'll have 2 weeks instead of one)

## Review from last time

* Vanishing / Exploding gradients are a main limitation of RNNs.

* limit the # of time steps we are considering in backpropagation (reduce the number of sigmoids we are multiplying)

  losing some connections between words at the beginning of a sentence and later ones.

* Orthogonal Init

  mainly helps with exploding gradients still have the issue of the products of sigmoids.

## Long Short-term memory

main idea: store some info in a cell $\vec{c}^{(t)}$ (RAM in your computer)

to keep track of relevant info.

From one time step to another

* Forget gate: some info in the cell (RAM) will be forgotten

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

* Input gate: new info that will be stored in the cell.

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

* Output gate: new info that will be stored in the hidden state

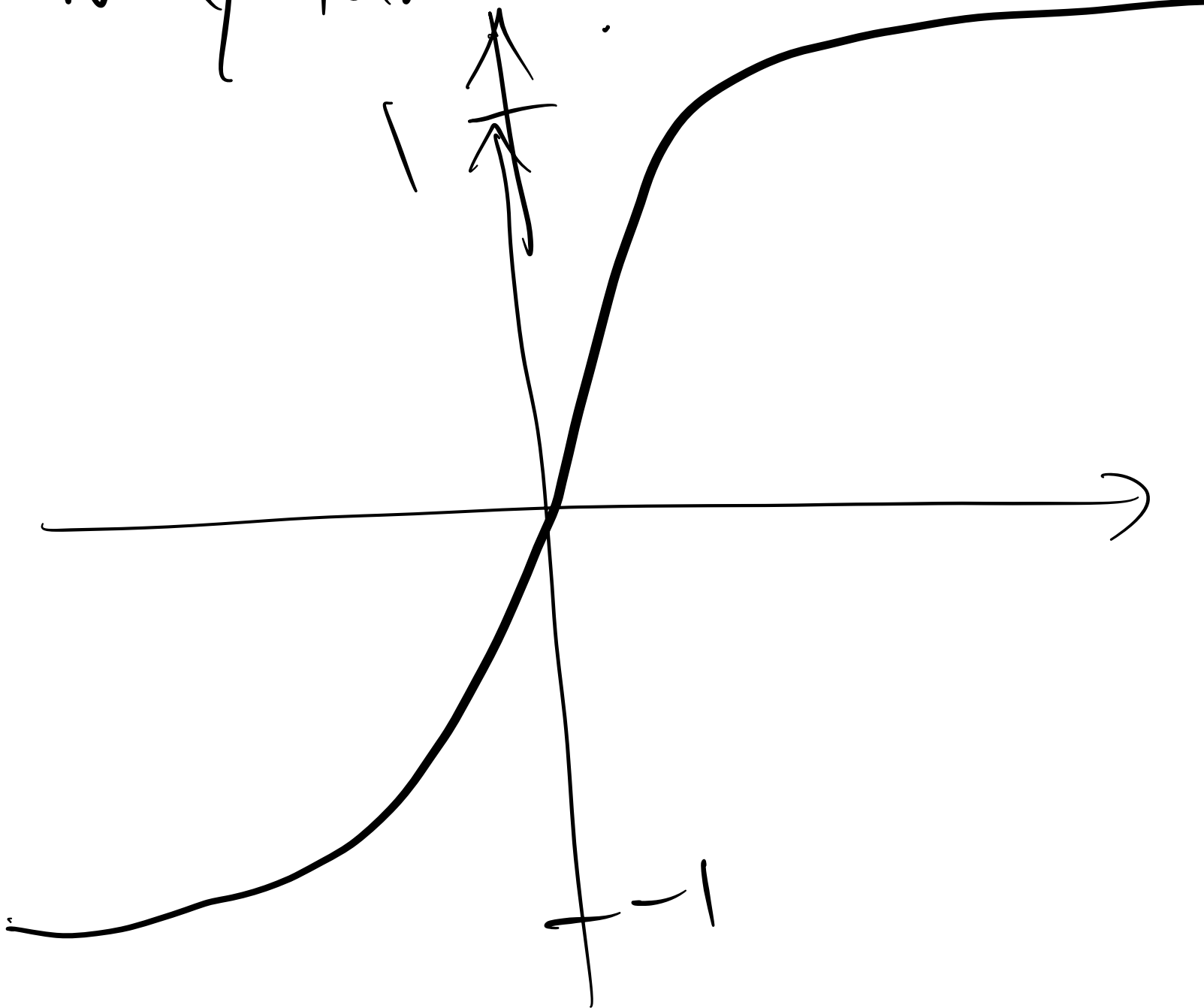$$\vec{o}^{(t)} = \sigma\left(W_o \vec{h}^{(t-1)} + U_o \vec{e}^{(t)} + \vec{b}_o\right)$$

Let's first start by updating the cell

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

element wise multiplication

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

# why tanh?

$1$

$t \sim 1$

$- - -$ "is fantastically bad"

$t-1$      $t$

$$C^{(t-1)} = 0.8$$

positive
sentiment

$$\tilde{C}^{(t)} = -0.6$$

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

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

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

$$= 1 * 0.8 + 1 * (-0.6)$$

$$= 0.2$$

## Remaining step:

Update the hidden state

$$\vec{h}^{(t)} = \vec{O}^{(t)} * \tanh(C_t)$$

$$\begin{cases} f^{(t)} & \text{forget gate} \\ i^{(t)} & \text{input gate} \\ O^{(t)} & \text{output gate} \\ \tilde{C}^{(t)} & \text{new info (tanh)} \end{cases}$$

$$c^{(t)} = f * c \quad + \lambda * c$$
$$h^{(t)} = o^{(t)} * \tanh(c^{(t)})$$

# Exercises:

## Show

$$\frac{\partial L}{\partial \vec{C}^{(t-1)}} = \frac{\partial L}{\partial \vec{C}^{(t)}} \cdot \vec{\beta}^{(t)}$$

then show

$$\frac{\partial L}{\partial \vec{C}^{(t-k)}} = \frac{\partial L}{\partial \vec{C}^{(t)}} \cdot \prod_{j=0}^{k-1} \vec{\beta}^{(t-j)}$$