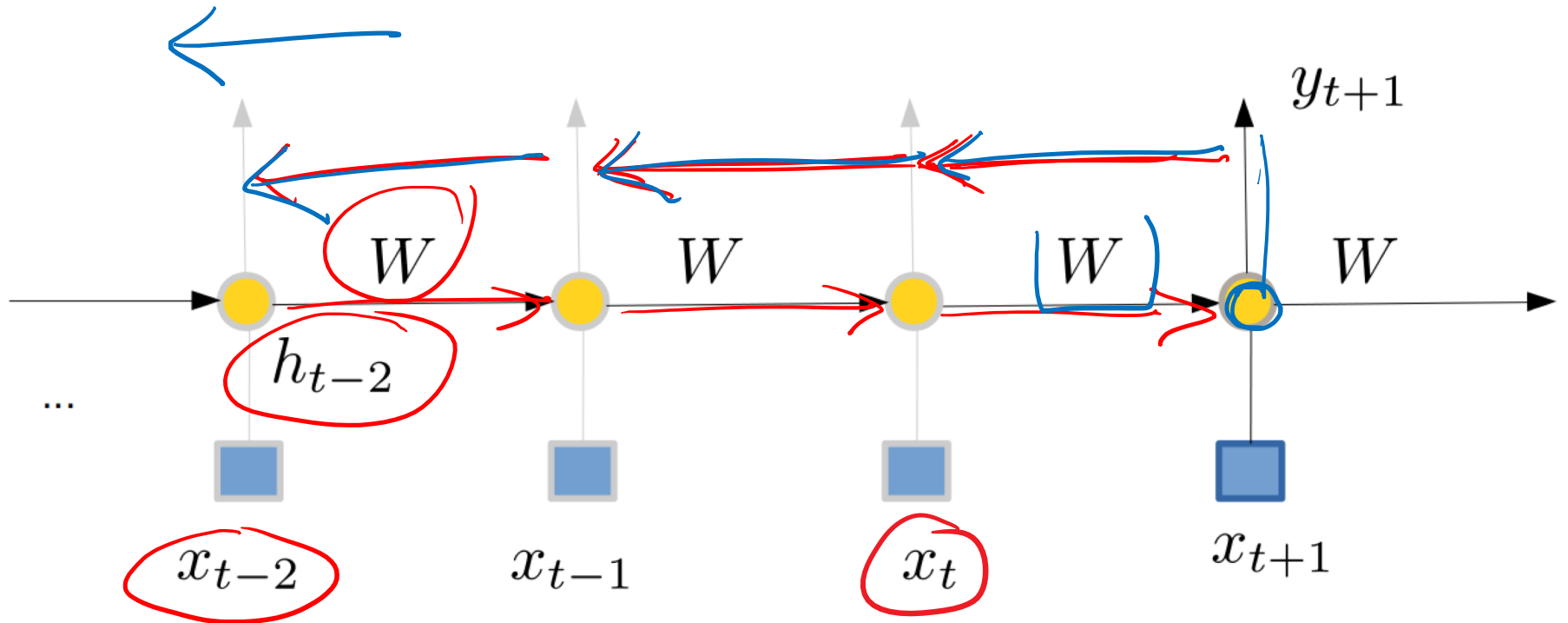

Understanding Vanishing Gradients

Heni Ben Amor, Ph.D.
Assistant Professor
Arizona State University

Deeper Look at Vanishing Gradients

We focus on the temporal connections



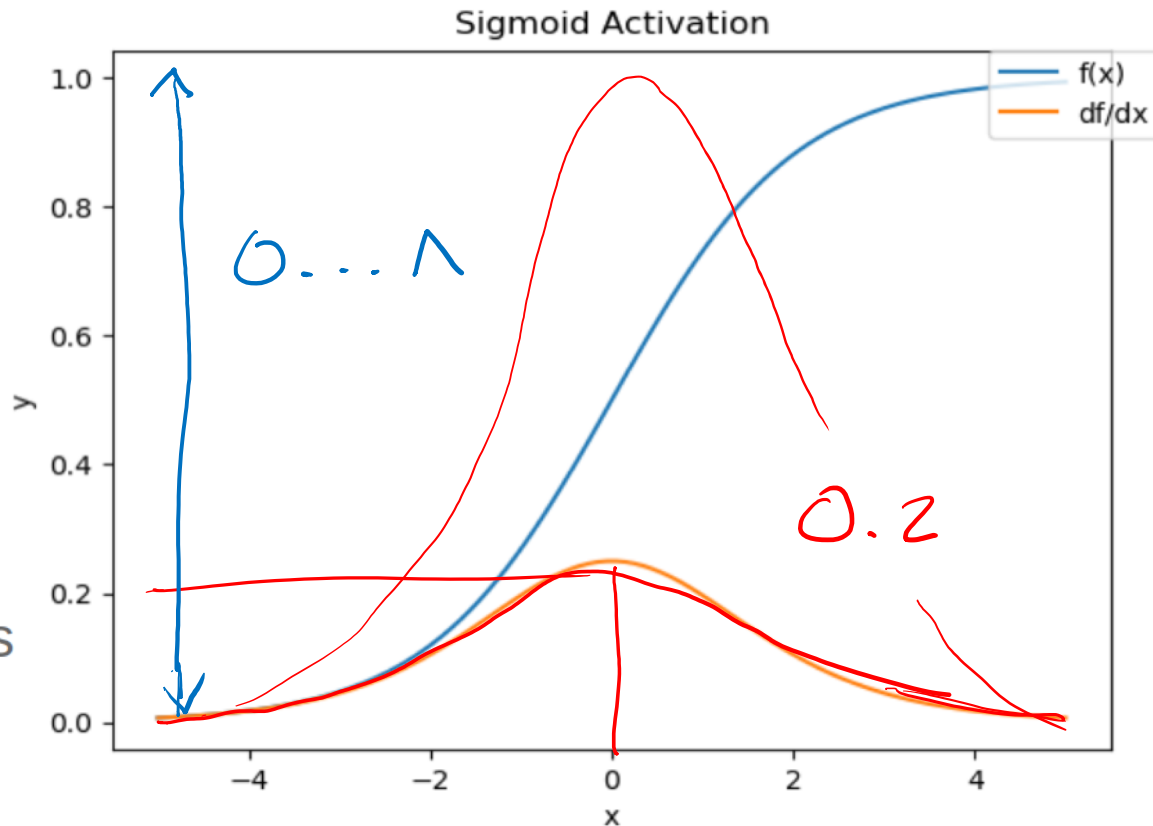
What happens when W is smaller than 1?

$$\sigma = \frac{1}{1 - e^{-x}}$$

Max value of df/dx is .25

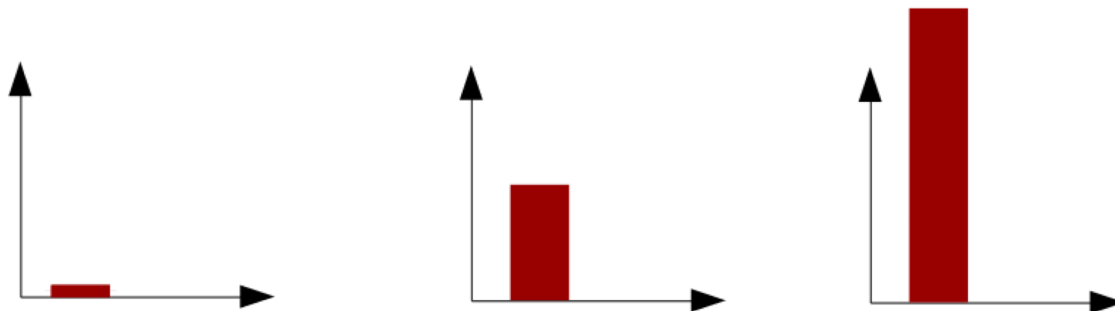
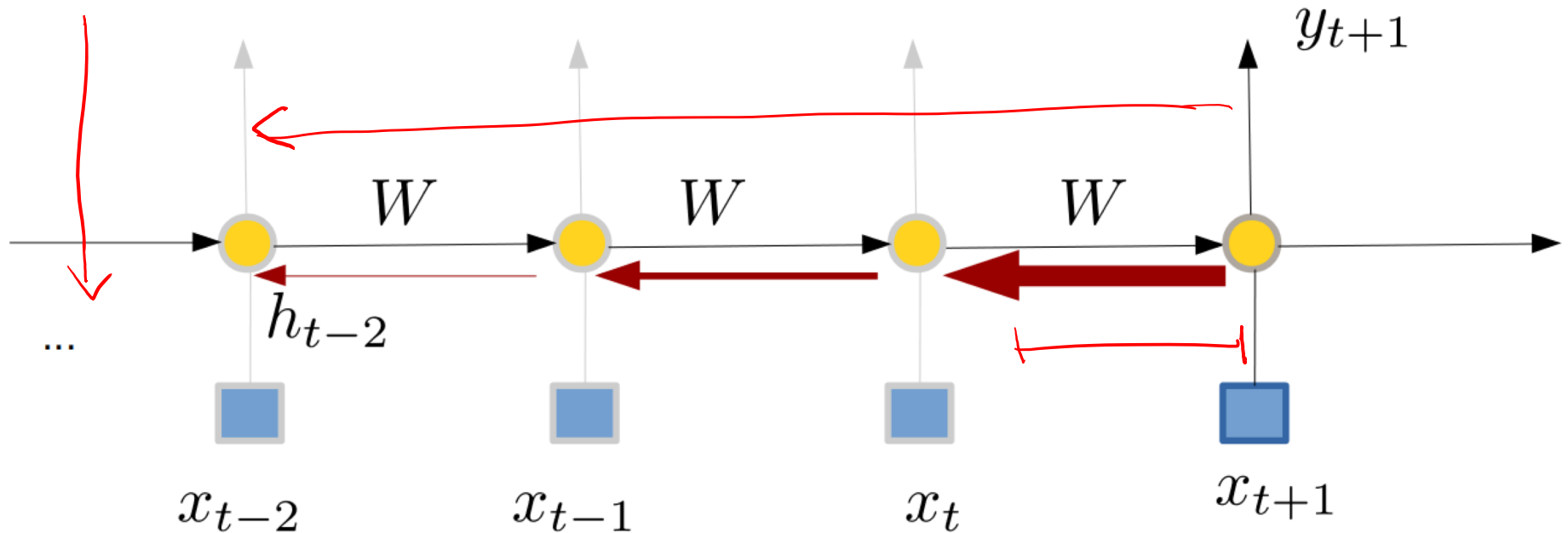
$$.25^2 = .0625$$

$$.25^5 = 0.00097$$



Vanishing Gradients Visualized

past timesteps have smaller influence



Gradient Magnitude

Theory Behind Vanishing Gradient

$$\nabla E = \left(\frac{\partial E}{\partial w_{\lambda 1}}, \frac{\partial E}{\partial w_{\lambda 2}}, \dots \right)$$

- BPTT: calculate gradient and propagate through time (TT)

$$\Theta = w_{ij} \quad \frac{\partial E}{\partial \Theta} = \frac{\partial \mathcal{L}}{\partial \theta} = \sum_{t=1}^T \frac{\partial \mathcal{L}_t}{\partial \theta}$$

\downarrow
variables / weights of RNN

- Apply chain rule to incorporate hidden state

$$\frac{\partial \mathcal{L}_t}{\partial \theta} = \sum_{k=1}^T \left(\frac{\partial \mathcal{L}_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial h_k} \cdot \frac{\partial h_k}{\partial \theta} \right)$$

- Temporal connections introduce dangerous product

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

Vanishing and Exploding Gradients

- | Challenge when training RNNs
- | Gradients quickly shrink to negligible values
- | Or, gradients may grow substantially and make learning unstable
- | An immediate result of the temporal connections
- | **Exponential growth in hidden state values**
- | Effect: learning is slow and yields poor results