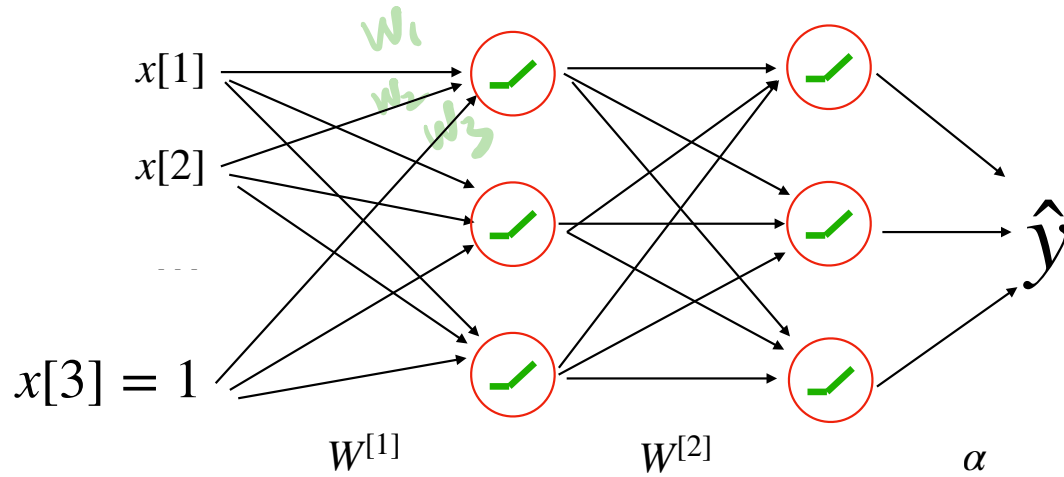


Neural Network: Training & Backpropagation

Announcements

Recap

A two layer fully connected feedforward NN:

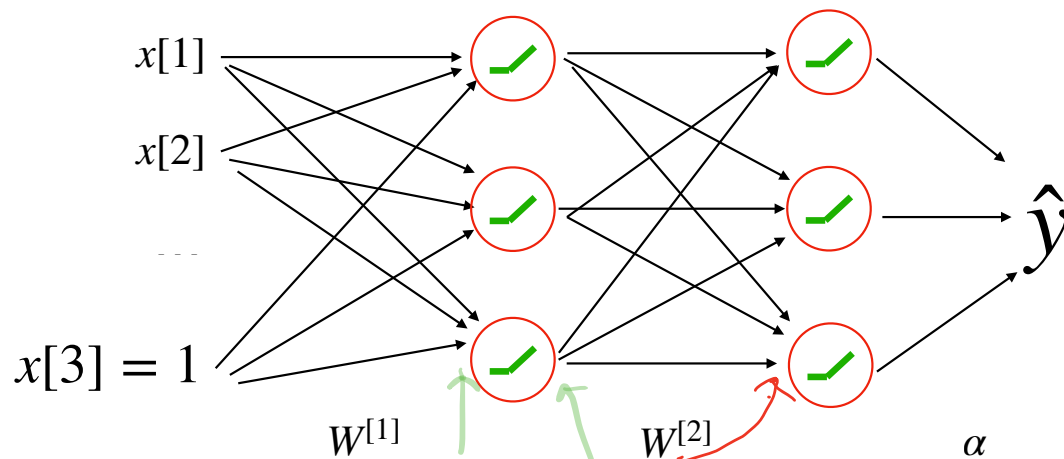


$$W^{[2]} = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix}$$

$$\text{ReLU}(x) = \max\{x, 0\}$$

Recap

A two layer fully connected feedforward NN:



$$h(x) := \alpha^T \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$

Outline of Today

1. Training NNs via SGD

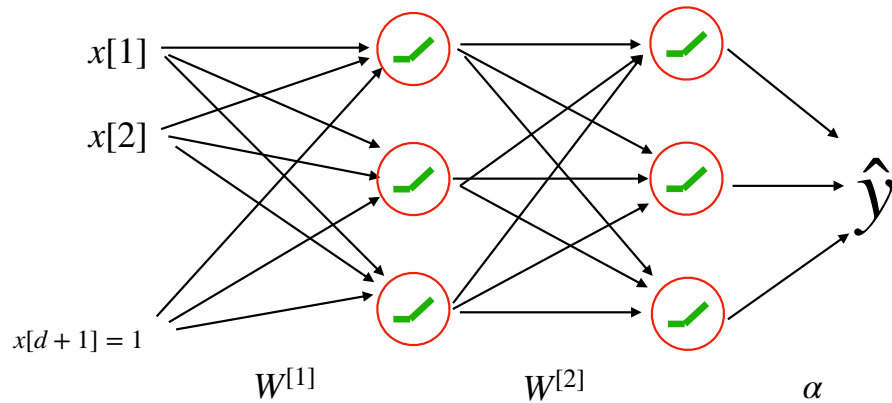
2. A Naive approach of computing gradients

3. Backpropagation: efficient way of computing gradients

Training neural network via SGD

Square loss on training example (x, y)

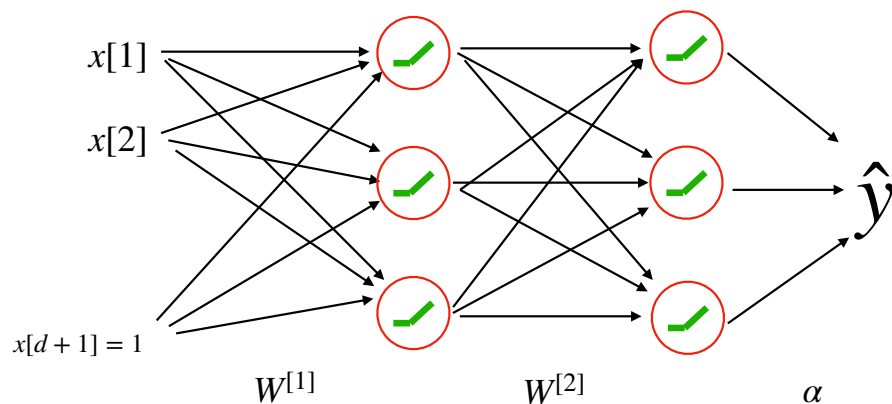
$$h(x) := \alpha^T \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$



Training neural network via SGD

Square loss on training example (x, y)

$$h(x) := \alpha^T \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$

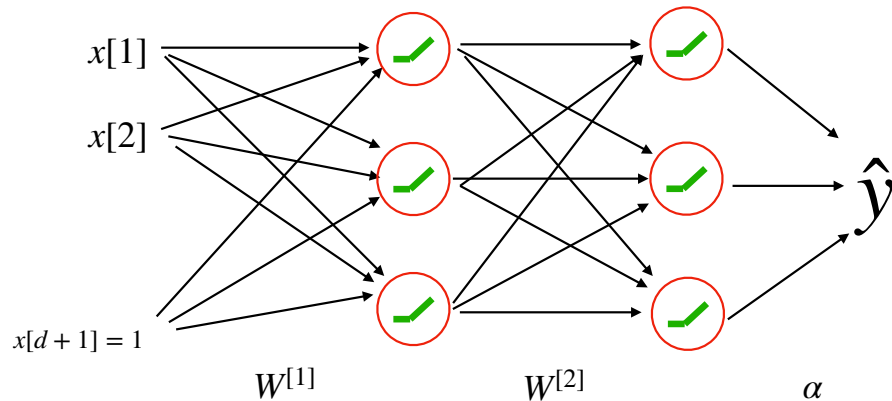


$$\ell(h(x), y) = (\hat{y} - y)^2, \text{ where } \hat{y} = h(x)$$

Training neural network via SGD

Square loss on training example (x, y)

$$h(x) := \alpha^T \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$



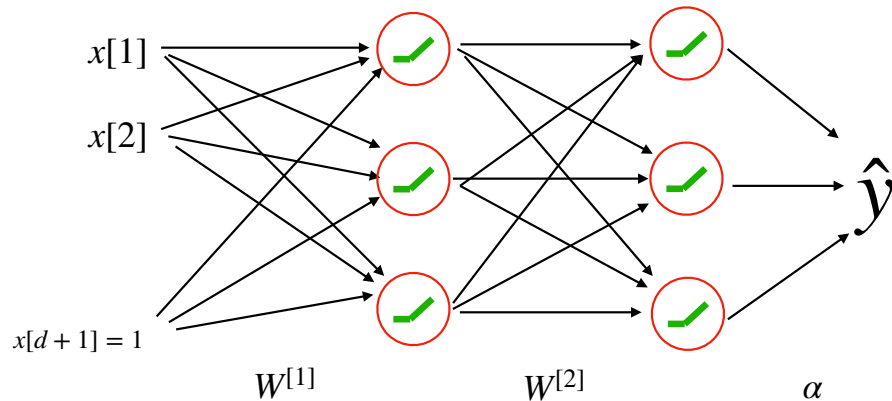
$$\ell(h(x), y) = (\hat{y} - y)^2, \text{ where } \hat{y} = h(x)$$

Trainable parameters $W^{[1]}, W^{[2]}, \alpha, b$

Training neural network via SGD

Square loss on training example (x, y)

$$h(x) := \alpha^T \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$



$$\ell(h(x), y) = (\hat{y} - y)^2, \text{ where } \hat{y} = h(x)$$

Trainable parameters $W^{[1]}, W^{[2]}, \alpha, b$

Compute gradients:

$$\frac{\partial \ell(h(x), y)}{\partial W^{[1]}} \quad \frac{\partial \ell(h(x), y)}{\partial W^{[2]}}$$

$$\frac{\partial \ell(h(x), y)}{\partial \alpha} \quad \frac{\partial \ell(h(x), y)}{\partial b}$$

Handwritten red notes and diagram illustrating the gradient calculation for a weight $W^{[1]}_{ij}$:

$$\frac{\partial \ell}{\partial W^{[1]}_{ij}} = \left[\frac{\partial \ell}{\partial c_{ij}} \right] \frac{\partial c_{ij}}{\partial W^{[1]}_{ij}}$$

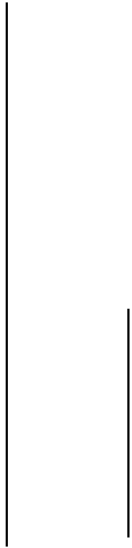
The diagram shows a node c_{ij} in a network, with an arrow pointing from the handwritten expression $\frac{\partial \ell}{\partial W^{[1]}_{ij}}$ to the node, and another arrow pointing from the node to the handwritten expression $\frac{\partial \ell}{\partial c_{ij}}$.

Training neural network via SGD

Mini-batch Stochastic gradient descent

$$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$$

For epoch $t = 1$ to T :



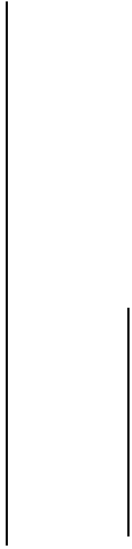
Training neural network via SGD

Mini-batch Stochastic gradient descent

$$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$$

// go through dataset multiple times

For epoch $t = 1$ to T :



Training neural network via SGD

Mini-batch Stochastic gradient descent

$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$ // go through dataset multiple times

For epoch $t = 1$ to T :

Randomly shuffle the data

Training neural network via SGD

Mini-batch Stochastic gradient descent

$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$ // go through dataset multiple times

For epoch $t = 1$ to T :

Randomly shuffle the data

// important (unbiased estimate of the true gradient)

Training neural network via SGD

Mini-batch Stochastic gradient descent

$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$ // go through dataset multiple times

For epoch $t = 1$ to T :

Randomly shuffle the data

// important (unbiased estimate of the true gradient)

Split the data into n/B many batches \mathcal{D}_i , each w/ size B

← 200

Training neural network via SGD

Mini-batch Stochastic gradient descent

$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$ // go through dataset multiple times

For epoch $t = 1$ to T :

Randomly shuffle the data

// important (unbiased estimate of the true gradient)

Split the data into n/B many batches \mathcal{D}_i , each w/ size B

For $i = 1$ to n/B

Training neural network via SGD

Mini-batch Stochastic gradient descent

$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$ // go through dataset multiple times

For epoch $t = 1$ to T :

Randomly shuffle the data

// important (unbiased estimate of the true gradient)

Split the data into n/B many batches \mathcal{D}_i , each w/ size B

For $i = 1$ to n/B

$$\text{Mini-batch gradient } g = \sum_{x,y \in \mathcal{D}_i} \nabla_{\theta} \ell(h_{\theta}(x), y) / B$$

Training neural network via SGD

Mini-batch Stochastic gradient descent

$\theta = [W^{[1]}, W^{[2]}, \alpha, b]$ // go through dataset multiple times

For epoch $t = 1$ to T :

Randomly shuffle the data

// important (unbiased estimate of the true gradient)

Split the data into n/B many batches \mathcal{D}_i , each w/ size B

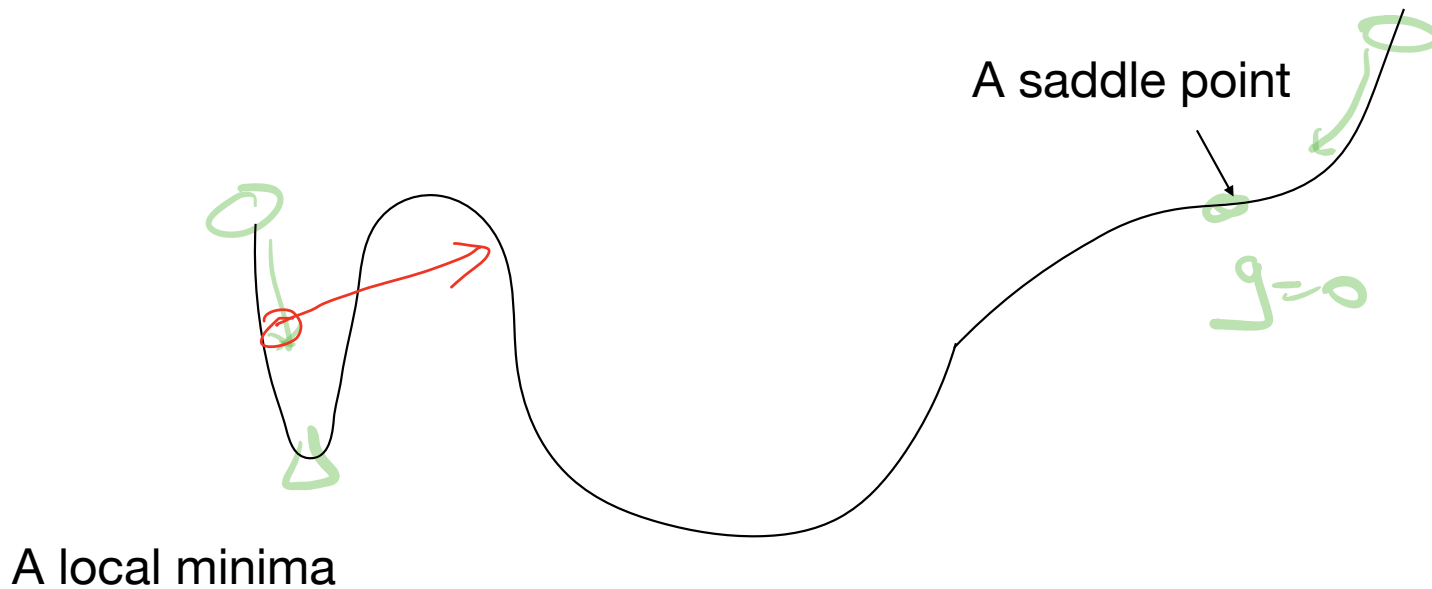
For $i = 1$ to n/B

$$\text{Mini-batch gradient } g = \sum_{x,y \in \mathcal{D}_i} \nabla_{\theta} \ell(h_{\theta}(x), y) / B$$

$$\theta = \theta - \eta g$$

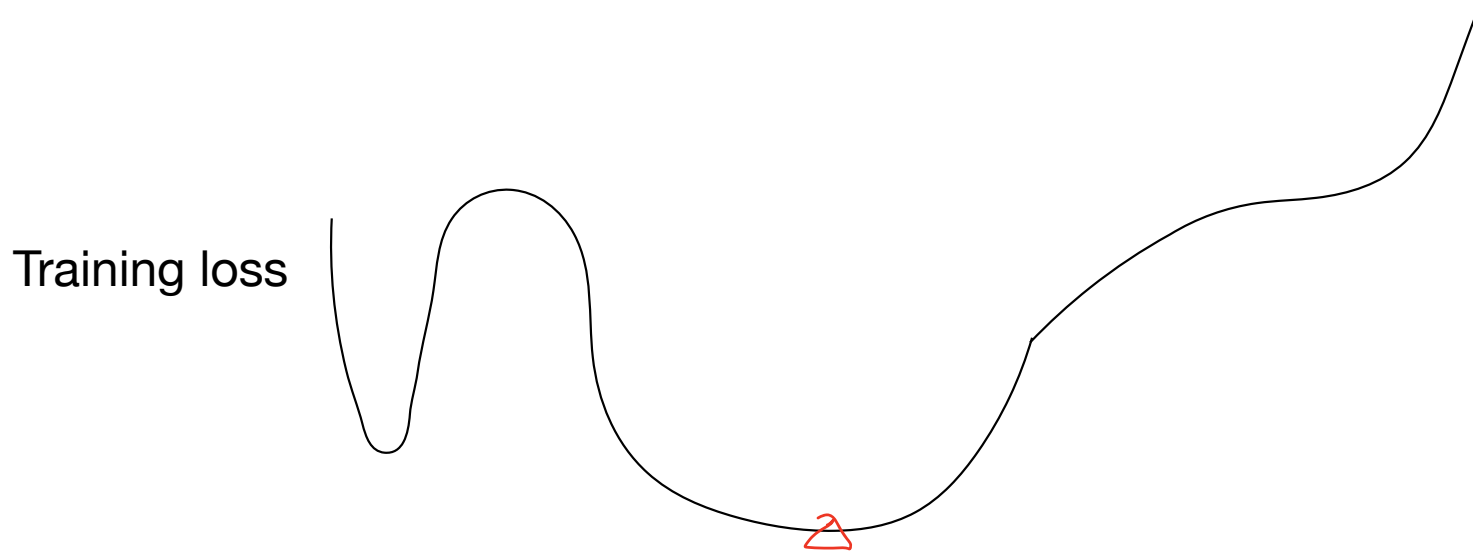
Training neural network via SGD

SGD helps avoiding local minima and saddle point



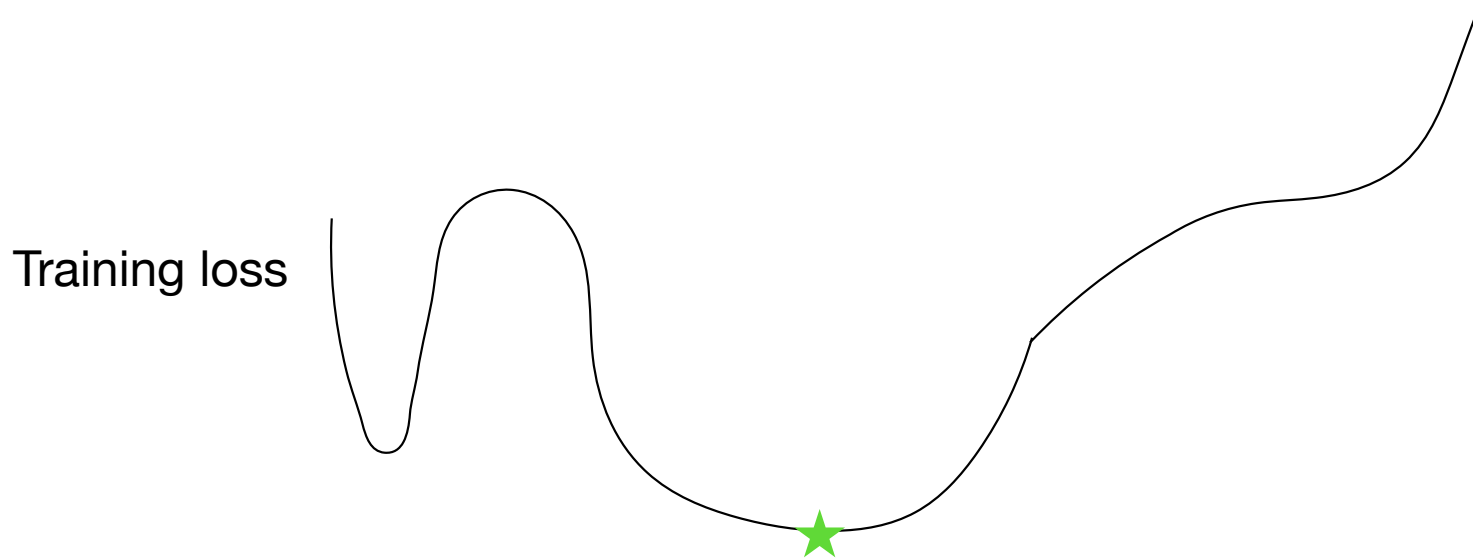
Training neural network via SGD

SGD tends to converge to a flat region



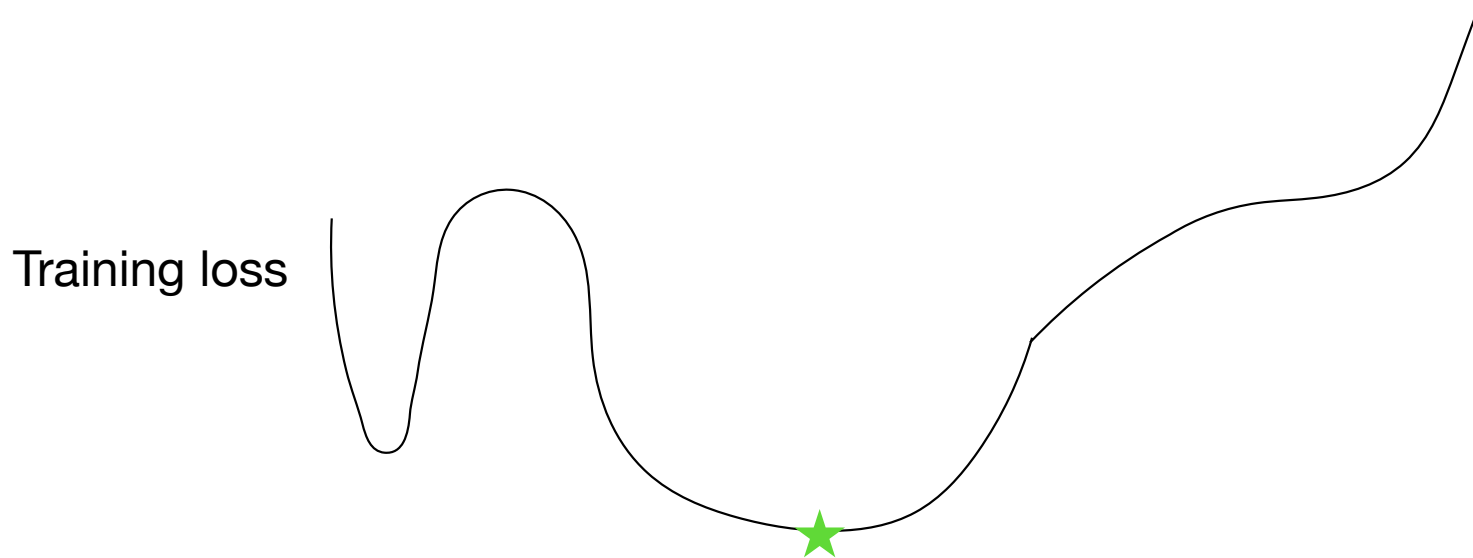
Training neural network via SGD

SGD tends to converge to a flat region



Training neural network via SGD

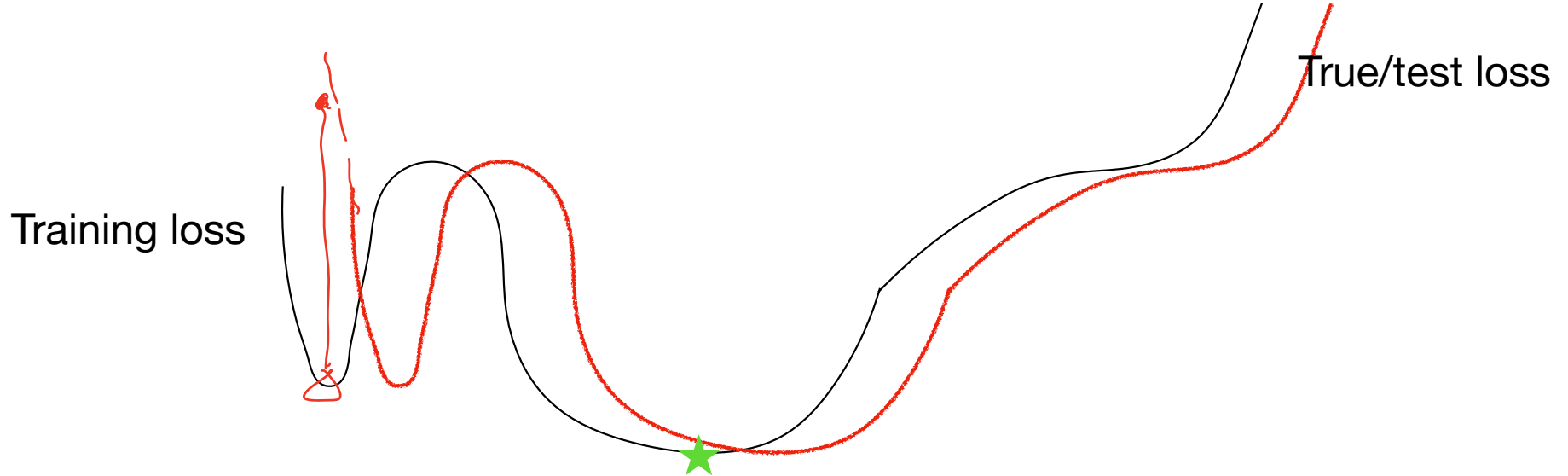
SGD tends to converge to a flat region



A flat local minima solution can help generalizes better to test data

Training neural network via SGD

SGD tends to converge to a flat region



A flat local minima solution can help generalizes better to test data

Outline of Today

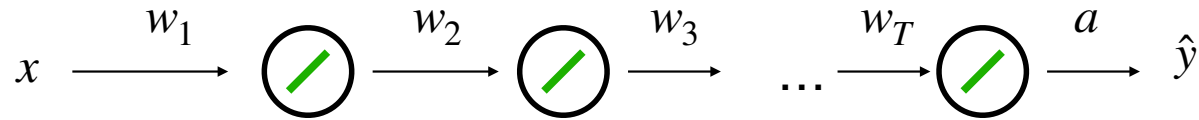
1. Training NNs via SGD

2. A Naive approach of computing gradients

3. Backpropagation: efficient way of computing gradients

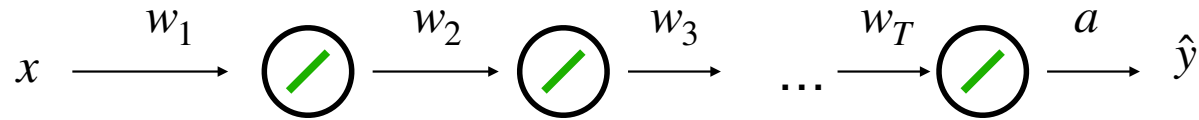
A naive algorithm

Consider the following one-dim case with identity transformation



A naive algorithm

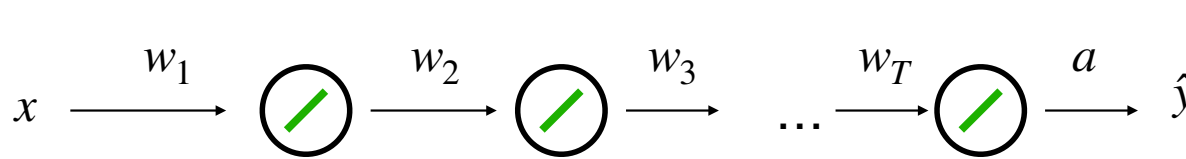
Consider the following one-dim case with identity transformation



$$\hat{y} = aw_T \dots w_2 w_1 x$$

A naive algorithm

Consider the following one-dim case with identity transformation



$$\hat{y} = aw_T \dots w_2 w_1 x$$

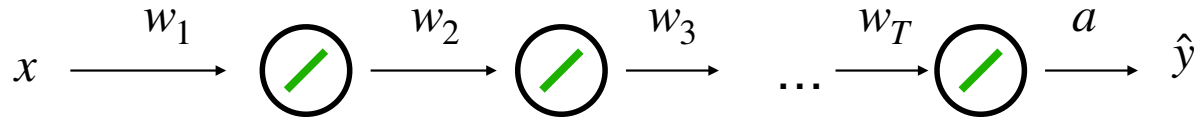
Let's compute derivatives $\partial \hat{y} / \partial w_i, \forall i = 1, \dots, T$

$$l(\hat{y}, y) = (\hat{y} - y)^2$$

$$\frac{\partial l}{\partial w_i} = \frac{\partial l}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_i}$$

A naive algorithm

Consider the following one-dim case with identity transformation

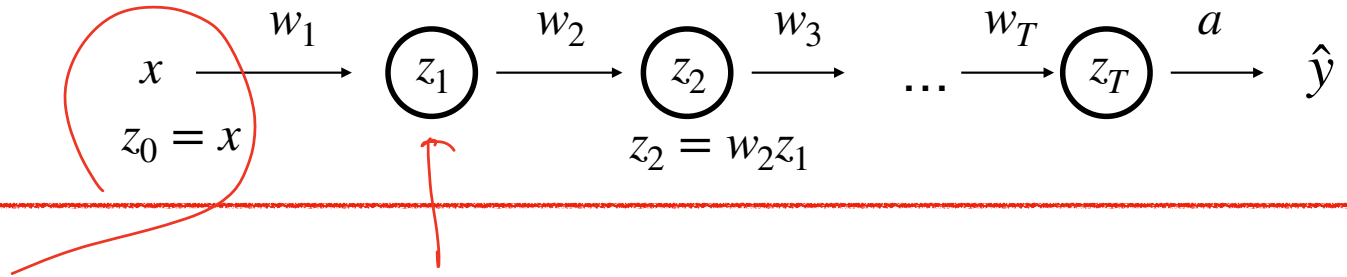


$$\hat{y} = aw_T \dots w_2 w_1 x$$

Let's compute derivatives $\partial \hat{y} / \partial w_i, \forall i = 1, \dots, T$

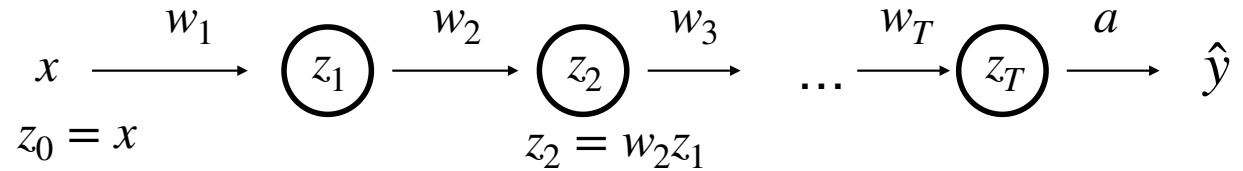
$$\text{Via chain rule: } \frac{\partial \ell}{\partial w_i} = \frac{\partial \ell}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_i}$$

A naive algorithm



$$z_1 = w_1 x = w_1 z_0$$

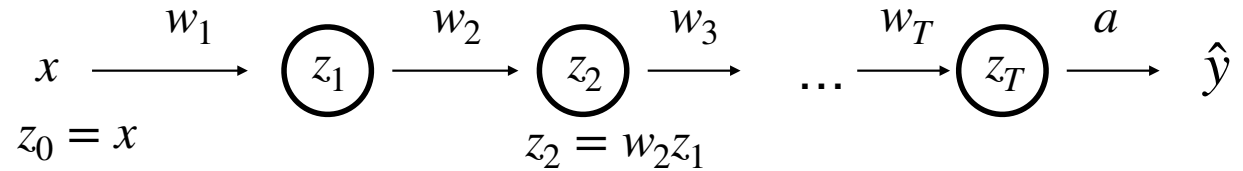
A naive algorithm



Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \cdot \frac{\partial z_T}{\partial z_{T-1}} \cdot \frac{\partial z_{T-1}}{\partial z_{T-2}} \cdot \dots \cdot \frac{\partial z_1}{\partial w_1}$$

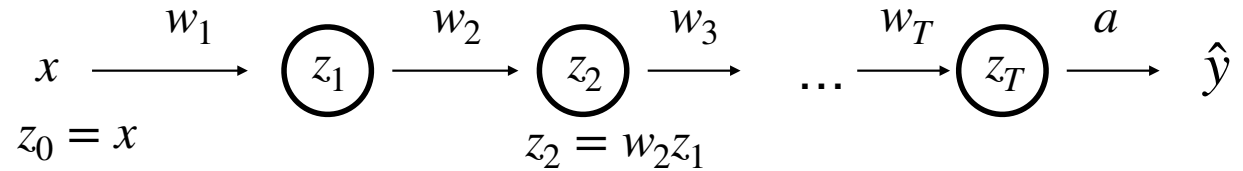
A naive algorithm



Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

A naive algorithm

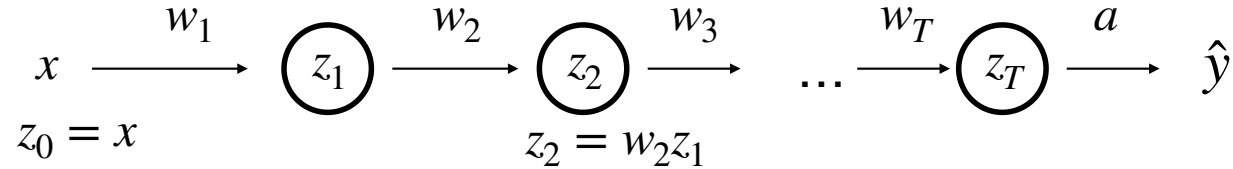


Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad // \text{ computation: T}$$

$$\frac{\partial \hat{y}}{\partial w_2}$$

A naive algorithm

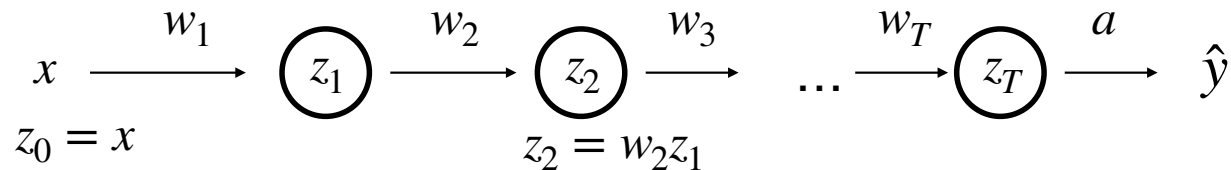


Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad // \text{ computation: } T$$

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

A naive algorithm



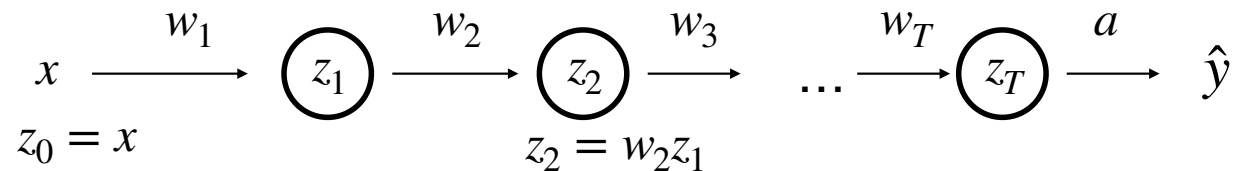
Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad // \text{ computation: } T$$

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2} \quad // \text{ computation: } T-1$$

$$\frac{\partial \hat{y}}{\partial w_1}$$

A naive algorithm



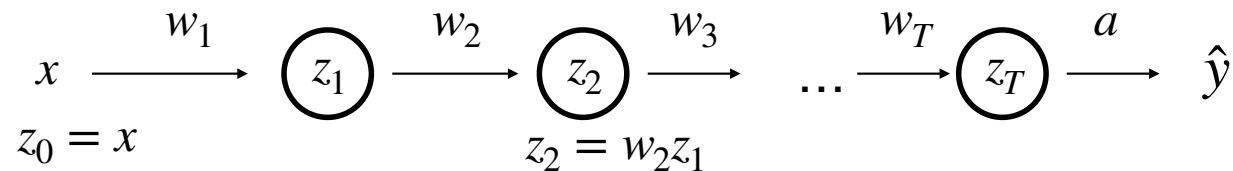
Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad // \text{ computation: T}$$

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2} \quad // \text{ computation: T-1}$$

$$\frac{\partial \hat{y}}{\partial w_T} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial w_T}$$

A naive algorithm



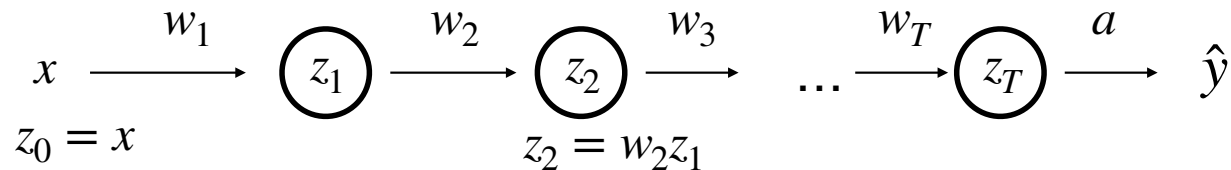
Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad // \text{ computation: } T$$

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2} \quad // \text{ computation: } T-1$$

$$\frac{\partial \hat{y}}{\partial w_T} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial w_T} \quad // \text{ computation: } 1$$

A naive algorithm



Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad // \text{ computation: } T$$

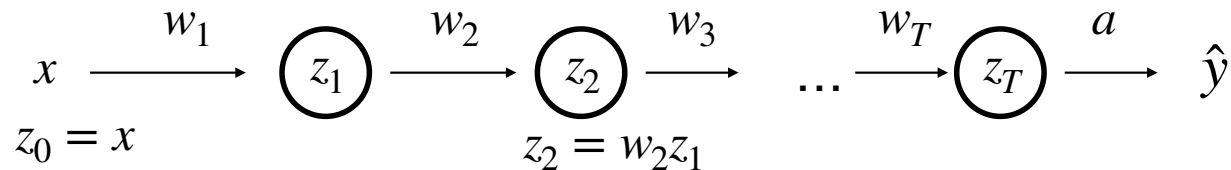
Total complexity:

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2} \quad // \text{ computation: } T-1$$

$$1 + 2 + \dots + T = O(T^2)$$

$$\frac{\partial \hat{y}}{\partial w_T} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial w_T} \quad // \text{ computation: } 1$$

A naive algorithm



Via chain rule:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad // \text{ computation: } T$$

Total complexity:

$$\frac{\partial \hat{y}}{\partial w_2} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial z_{T-1}} \dots \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2} \quad // \text{ computation: } T-1$$

$$1 + 2 + \dots + T = O(T^2)$$

$$\frac{\partial \hat{y}}{\partial w_T} = \frac{\partial \hat{y}}{\partial z_T} \frac{\partial z_T}{\partial w_T} \quad // \text{ computation: } 1$$

Quadratic in size of the graph!

Summary so far

What we did:

for each edge weight w_i , apply chain rule to calculate $\partial \hat{y} / \partial w_i$

Summary so far

What we did:

for each edge weight w_i , apply chain rule to calculate $\partial \hat{y} / \partial w_i$

What we got:

Able to compute gradient in running time $O((\text{size of graph})^2)$

Summary so far

What we did:

for each edge weight w_i , apply chain rule to calculate $\partial \hat{y} / \partial w_i$

What we got:

Able to compute gradient in running time $O((\text{size of graph})^2)$

Can we do better in running time?

Outline of Today

1. Training NNs via SGD

2. A Naive approach of computing gradients

3. Backpropagation: efficient way of computing gradients

...Hinton popularized what they termed a “backpropagation” algorithm ... in 1986.

...Hinton popularized what they termed a “backpropagation” algorithm ... in 1986.

...the algorithm propagated measures of the errors produced by the network’s guesses **backwards through its neurons, starting with those directly connected to the outputs.**

...Hinton popularized what they termed a “backpropagation” algorithm ... in 1986.

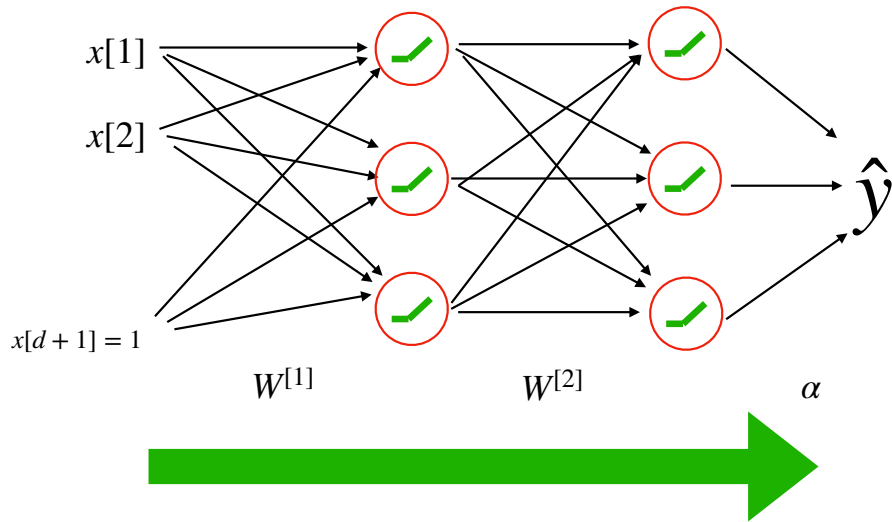
...the algorithm propagated measures of the errors produced by the network’s guesses **backwards through its neurons, starting with those directly connected to the outputs.**

This allowed networks with intermediate “hidden” neurons between input and output layers to **learn efficiently**, overcoming the limitations noted by Minsky and Papert.

Overview of backpropagation

Forward pass followed by a backward pass

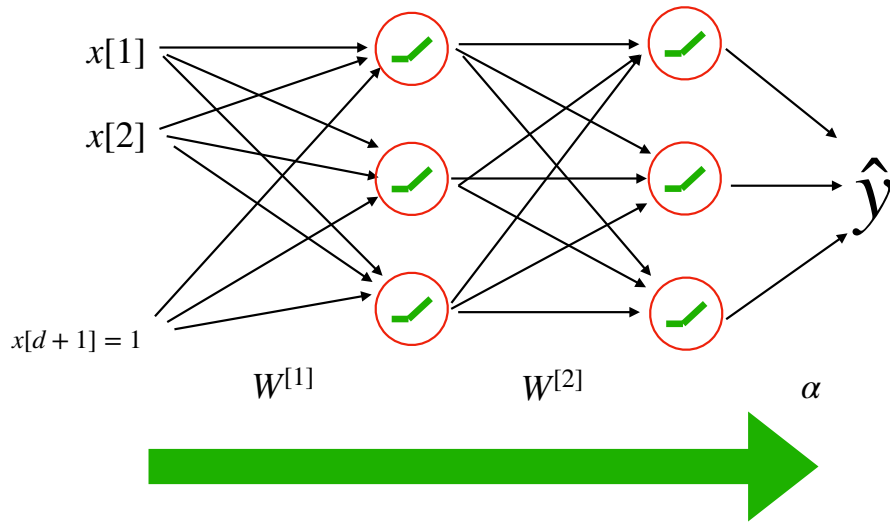
Forward pass:



Overview of backpropagation

Forward pass followed by a backward pass

Forward pass:

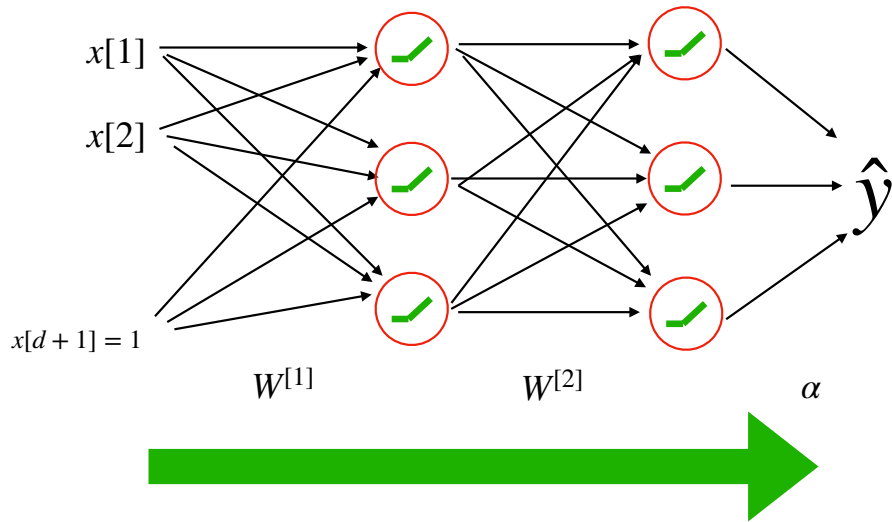


Store input & output of all neurons

Overview of backpropagation

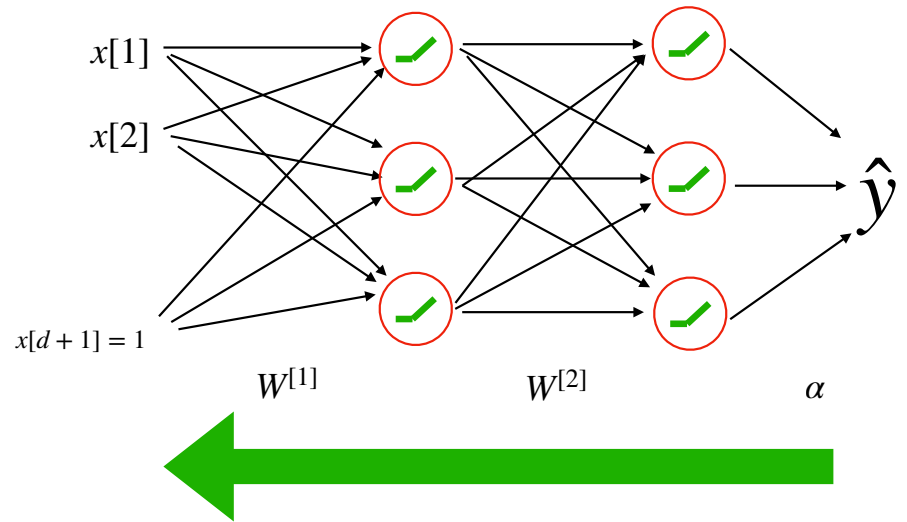
Forward pass followed by a backward pass

Forward pass:



Store input & output of all neurons

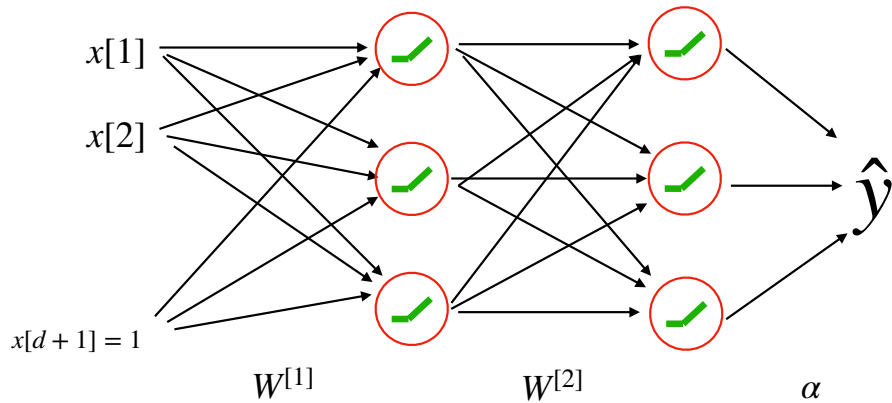
backward pass:



Overview of backpropagation

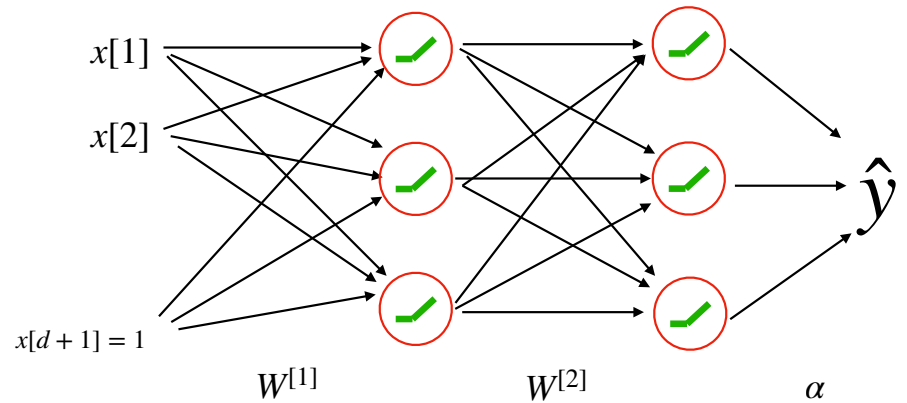
Forward pass followed by a backward pass

Forward pass:



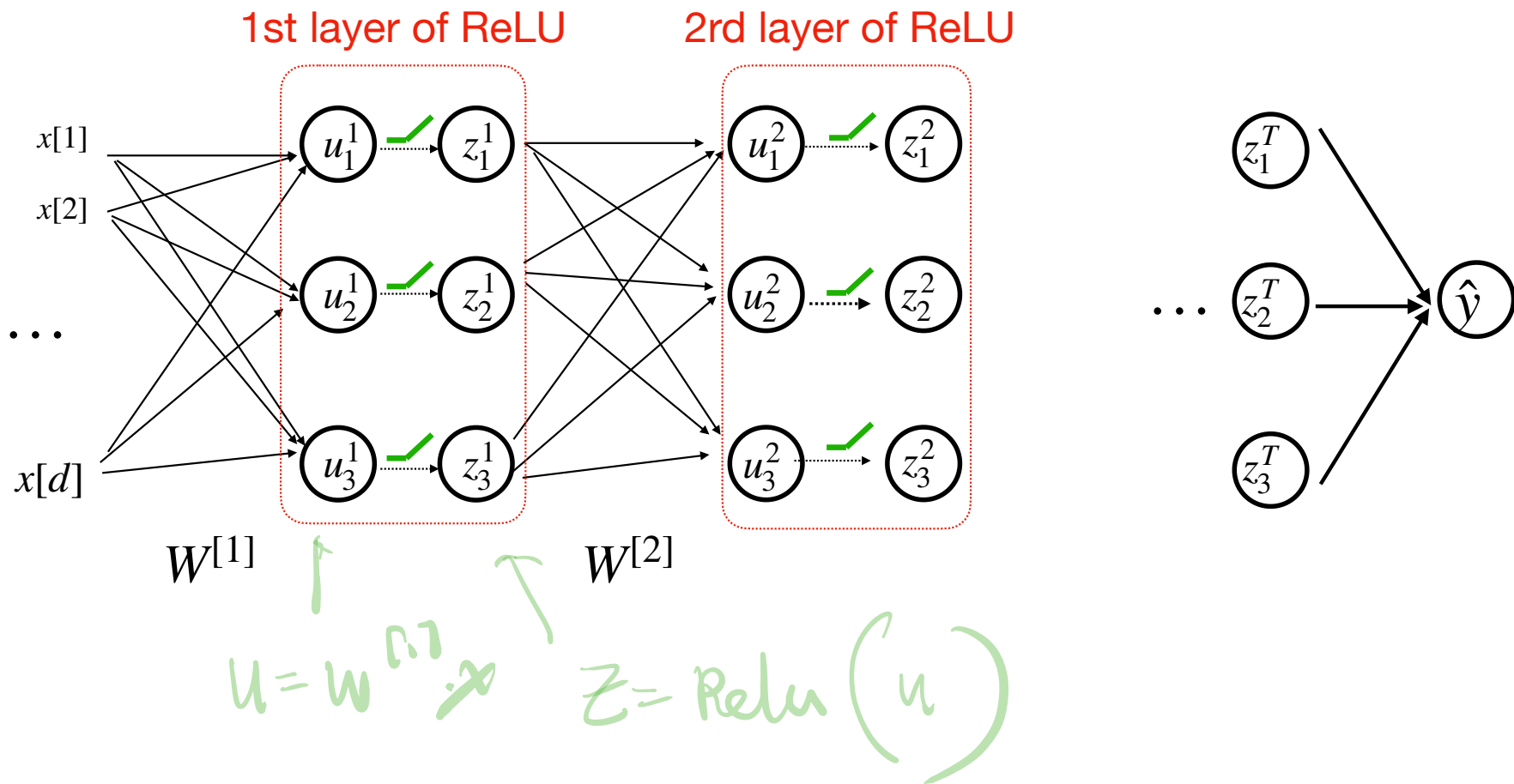
Store input & output of all neurons

backward pass:

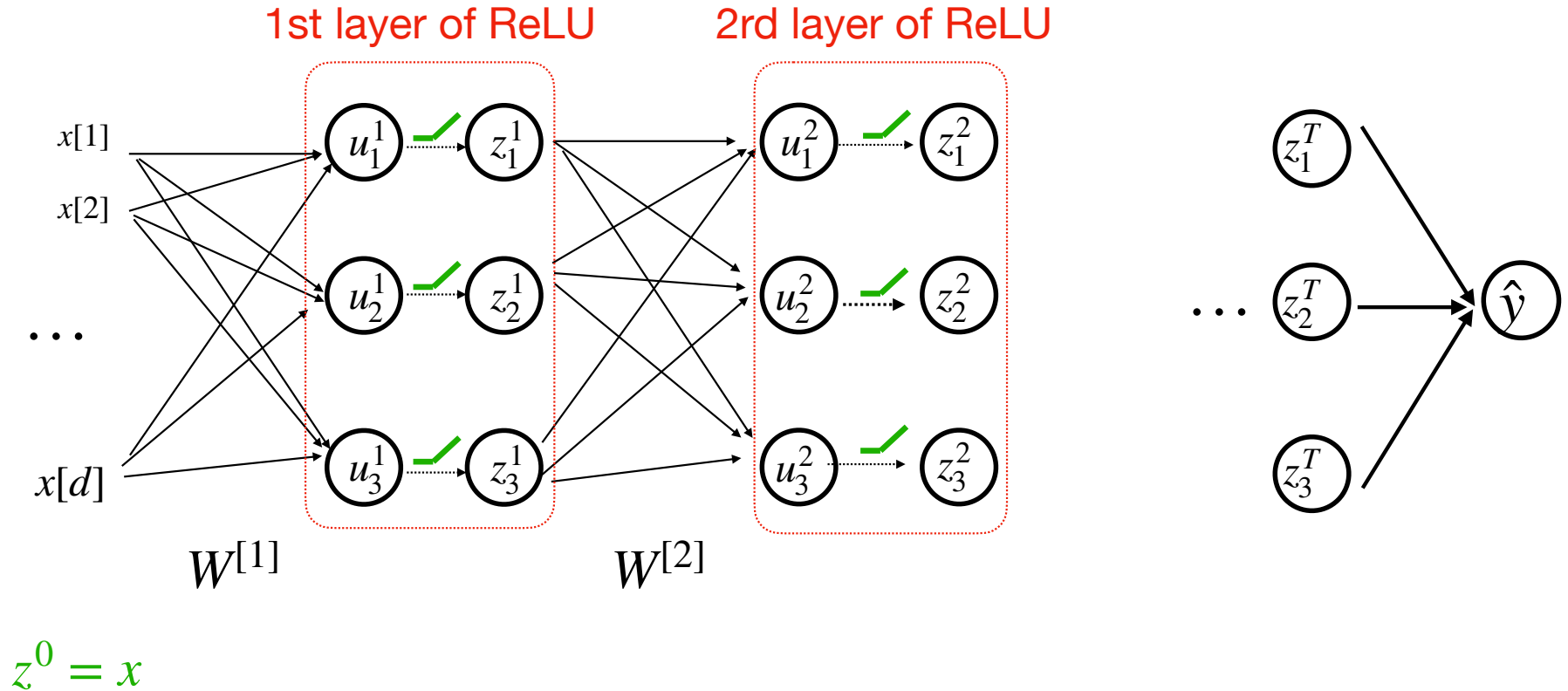


Compute derivatives

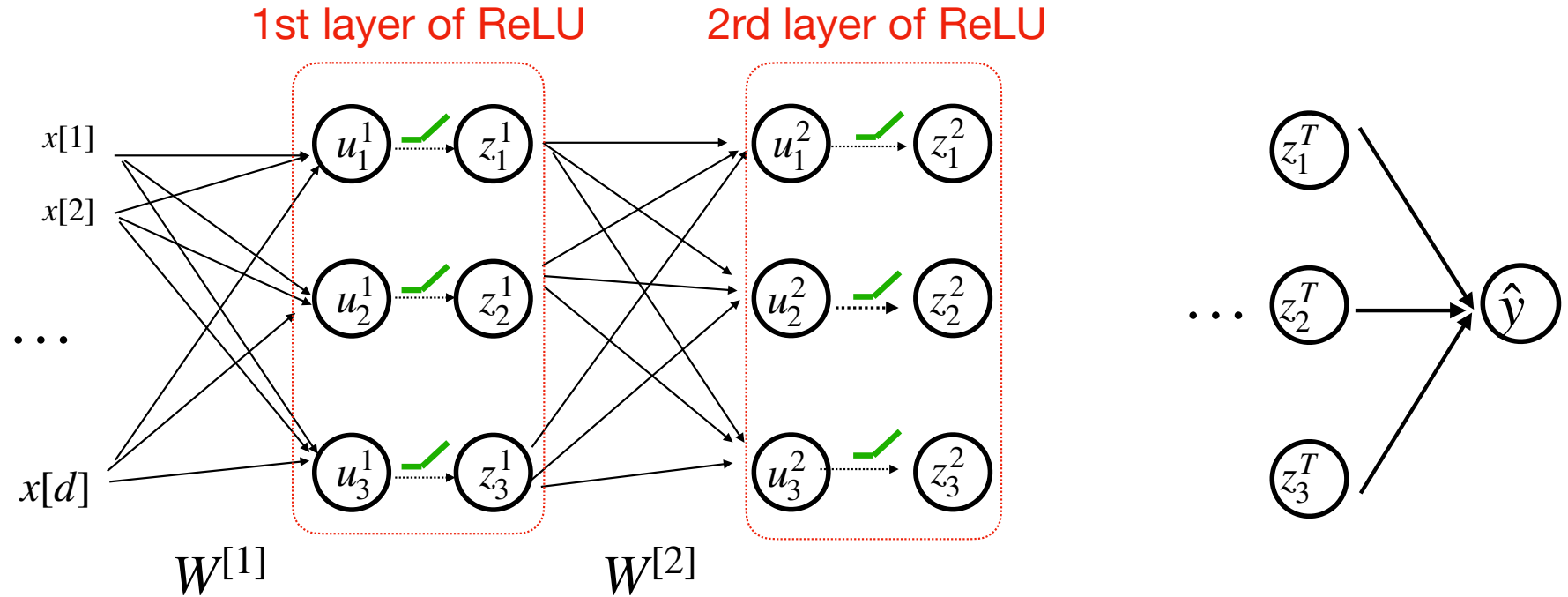
A Forward Pass: from $t = 0$ to T



A Forward Pass: from $t = 0$ to T

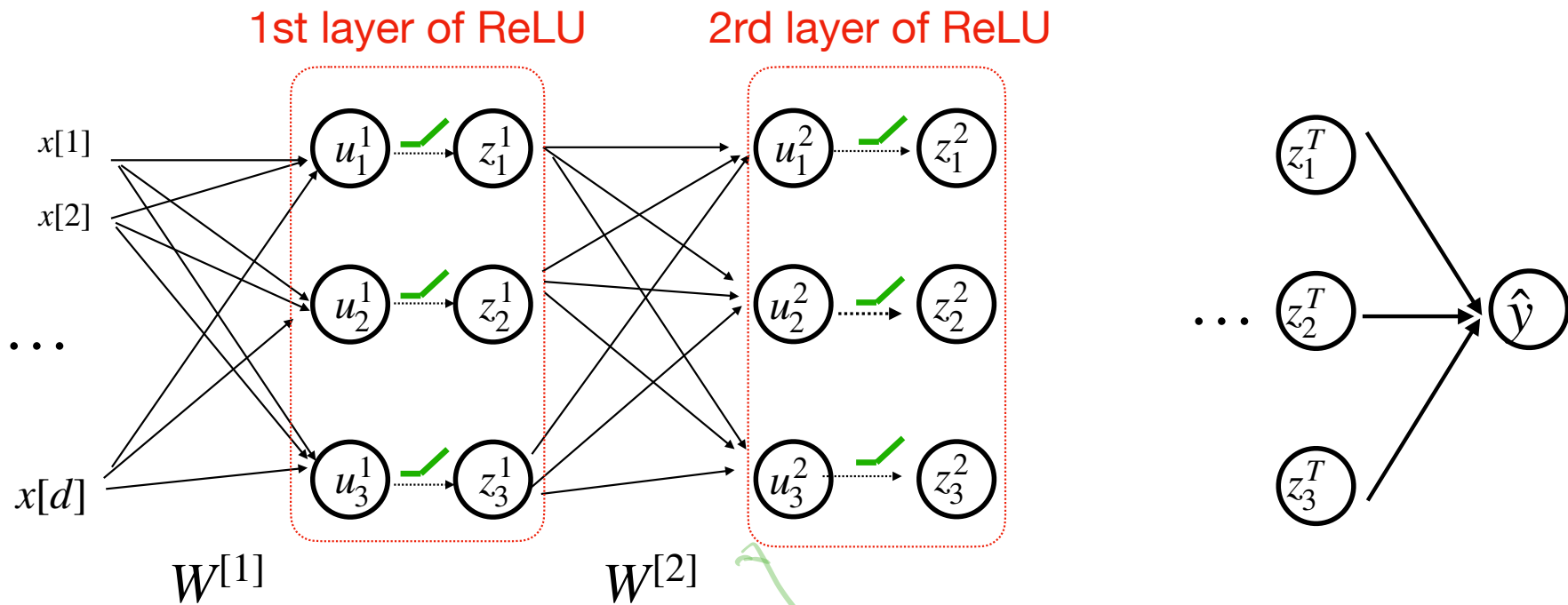


A Forward Pass: from $t = 0$ to T



$$z^0 = x \quad u^1 = W^{[1]}z^0$$

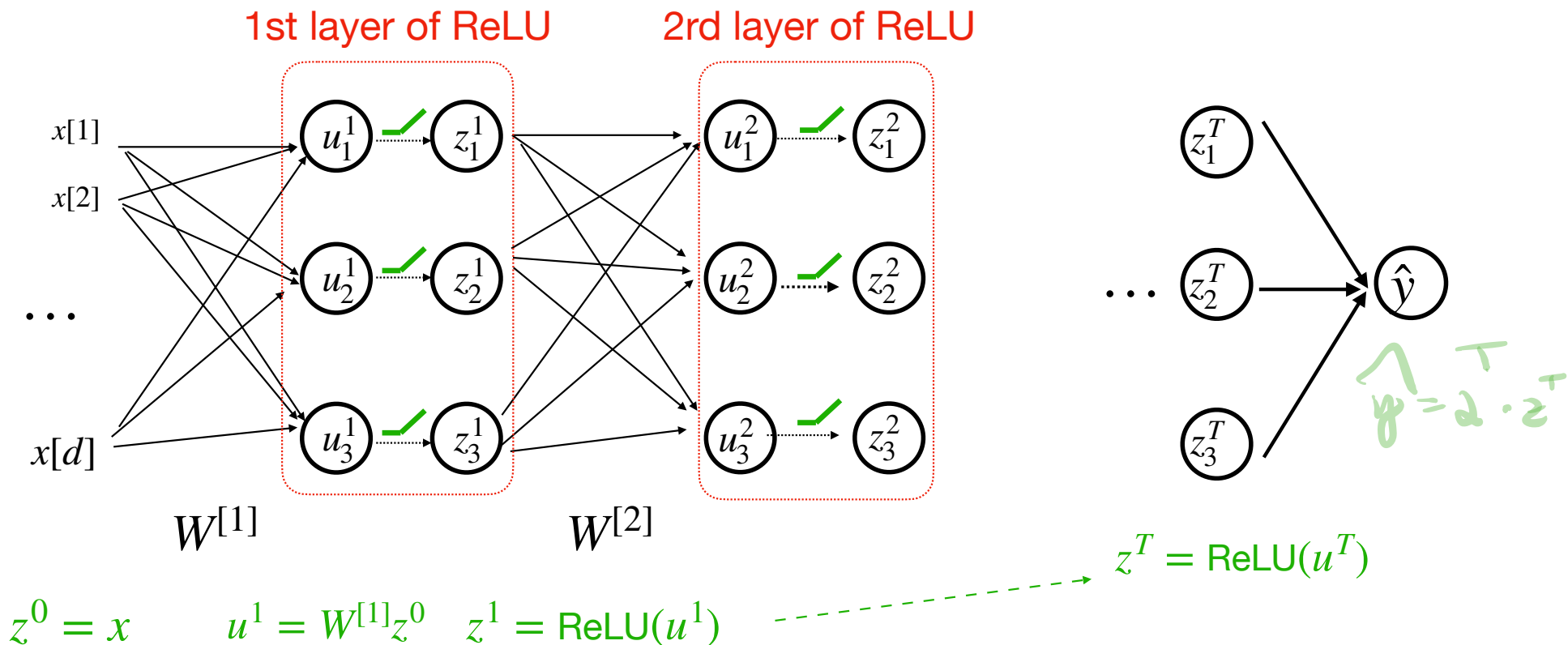
A Forward Pass: from $t = 0$ to T



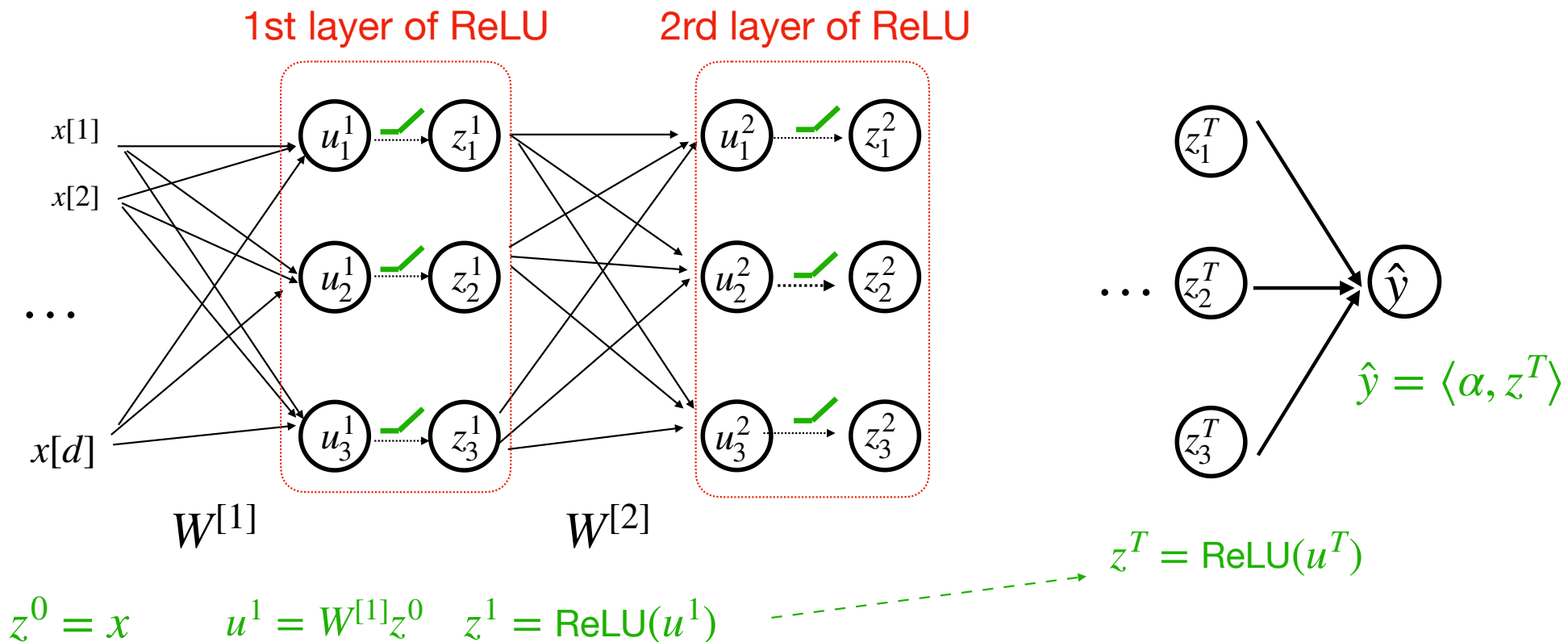
$$z^0 = x \quad u^1 = W^{[1]}z^0 \quad z^1 = \text{ReLU}(u^1)$$

$$u^2 = W^{[2]}z^1$$
$$z^2 = \text{ReLU}(u^2)$$

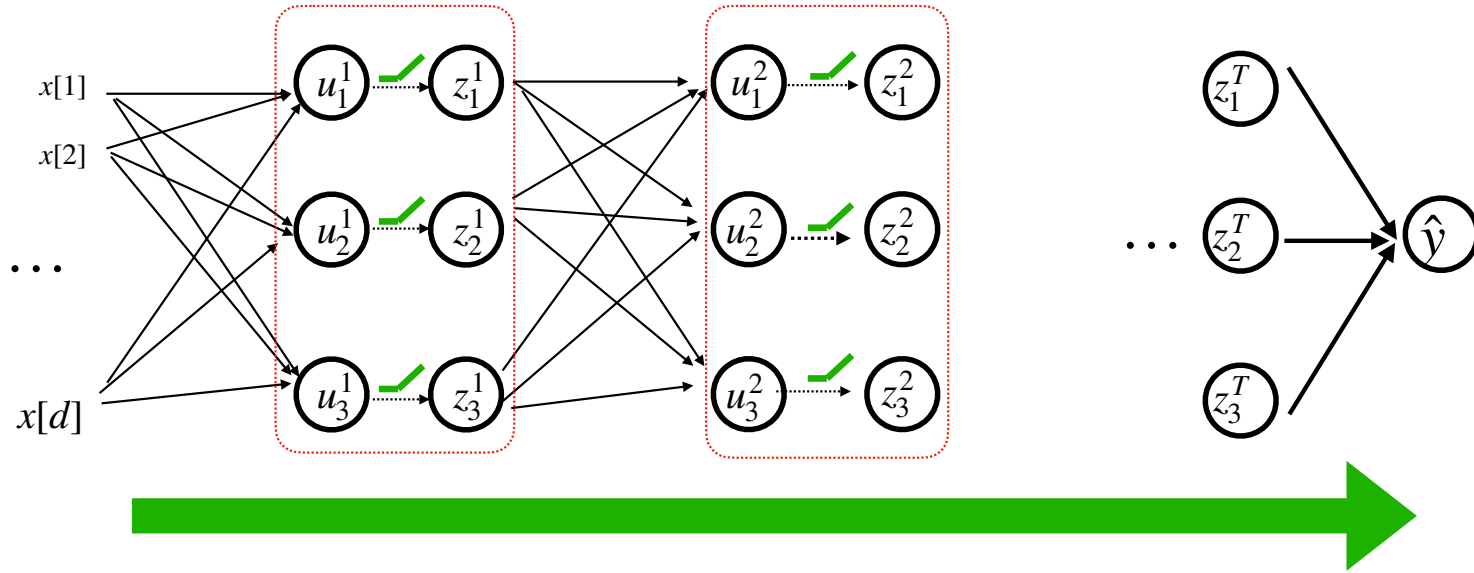
A Forward Pass: from $t = 0$ to T



A Forward Pass: from $t = 0$ to T

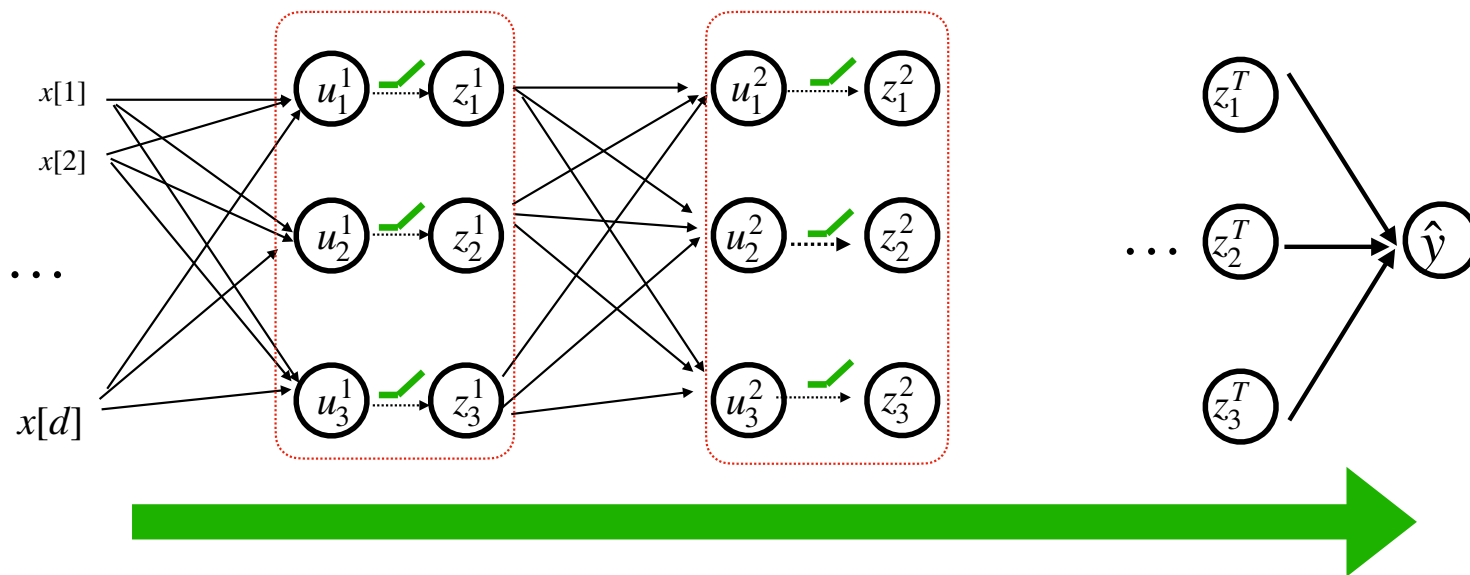


Summary of the forward pass



All nodes' values (i.e., z, u, \hat{y}) are computed and stored

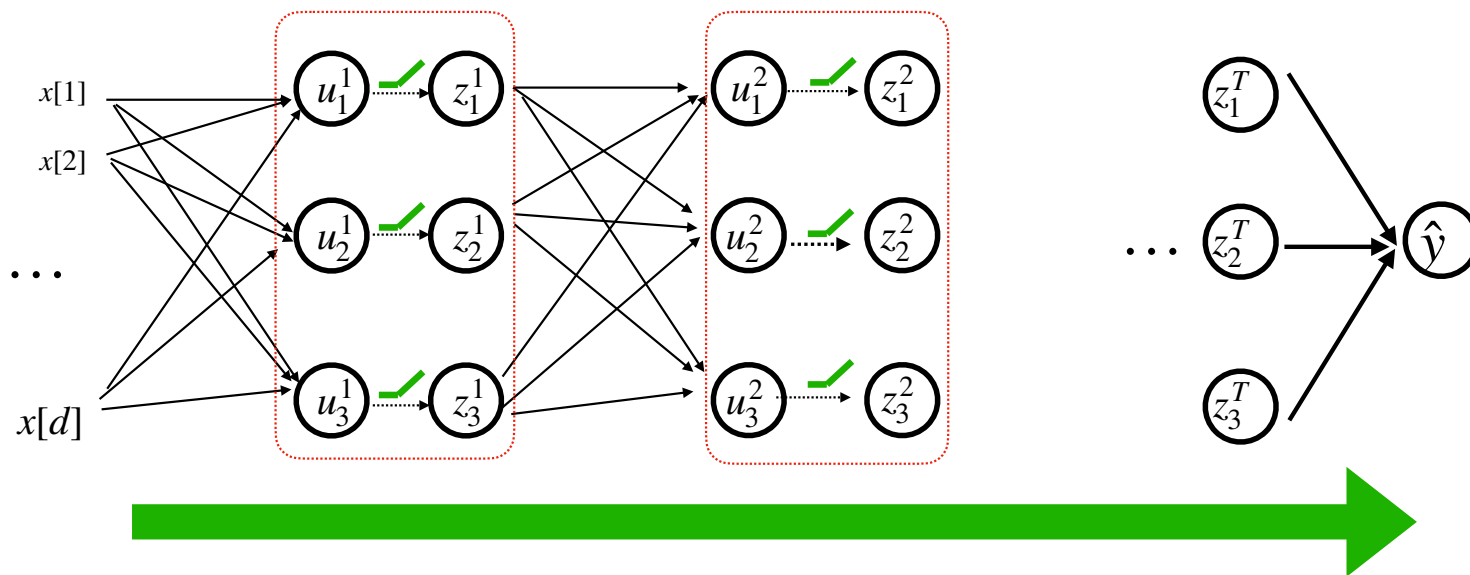
Summary of the forward pass



All nodes' values (i.e., z , u , \hat{y}) are computed and stored

Q: what is the computation complexity of the forward pass?

Summary of the forward pass



All nodes' values (i.e., z , u , \hat{y}) are computed and stored

Q: what is the computation complexity of the forward pass?

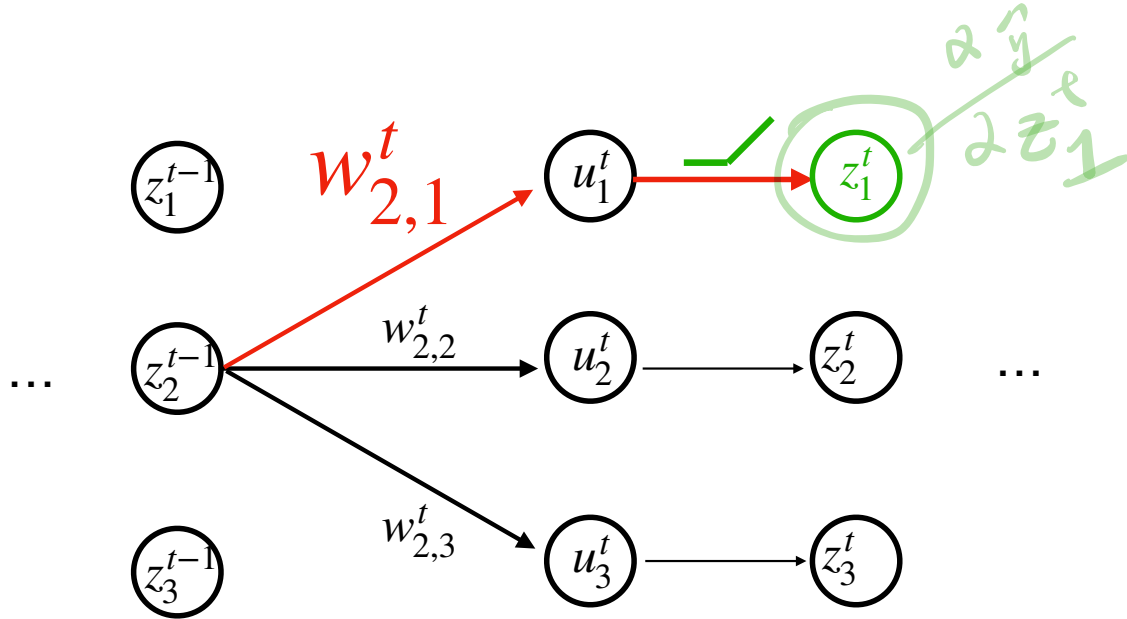
A: linear in # of Edges + # of nodes

The backward Pass

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .

The backward Pass

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .

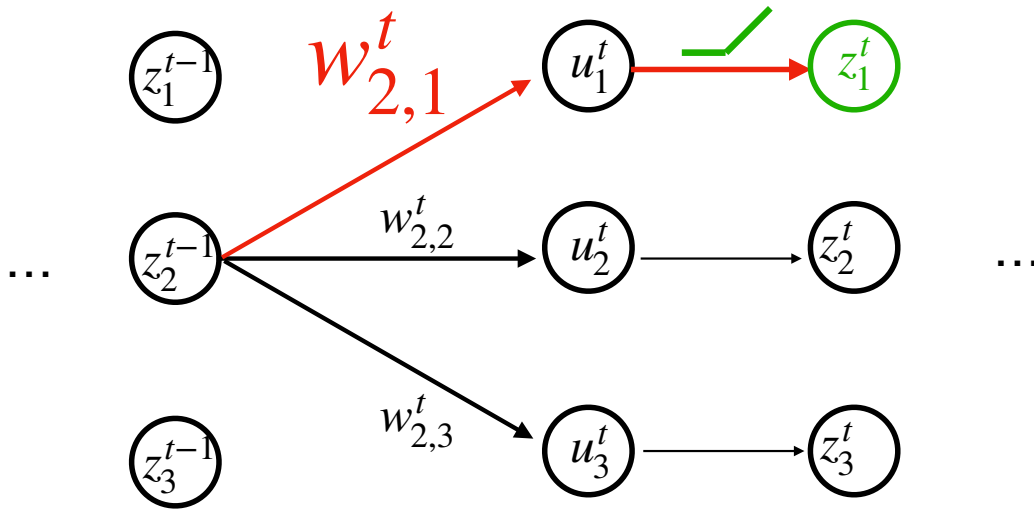


The backward Pass

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .

Proof:

WLOG consider $\partial \hat{y} / \partial w_{2,1}^t$



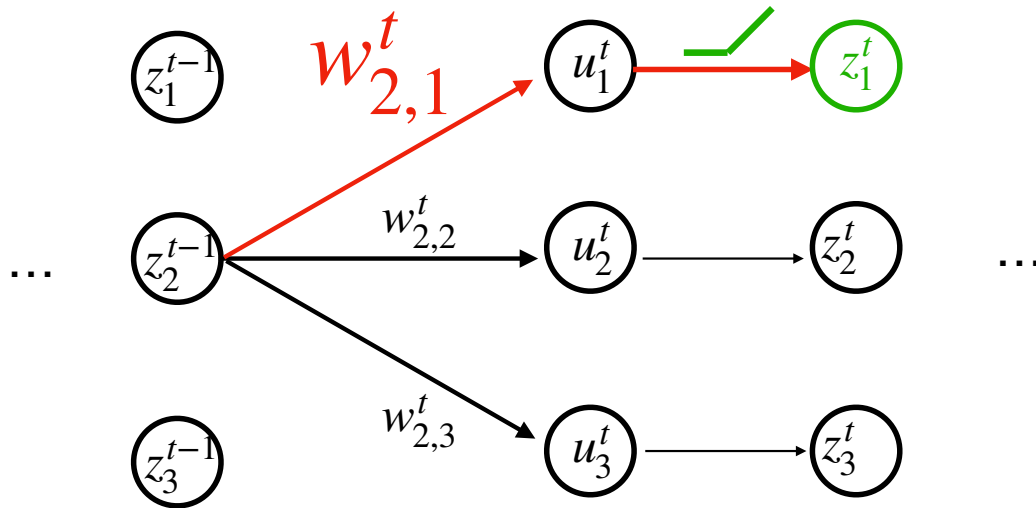
The backward Pass

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .

Proof:

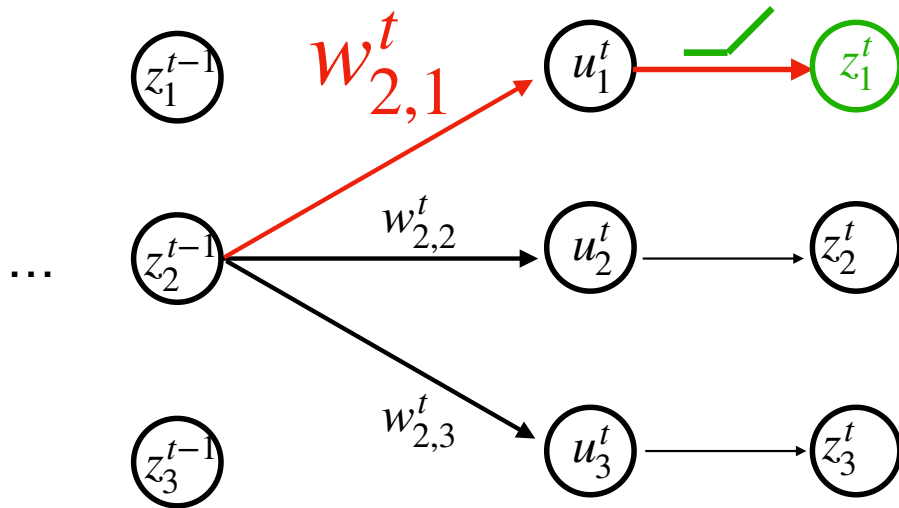
WLOG consider $\partial \hat{y} / \partial w_{2,1}^t$

$\partial \hat{y} / \partial w_{2,1}^t$



The backward Pass

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .



Proof:

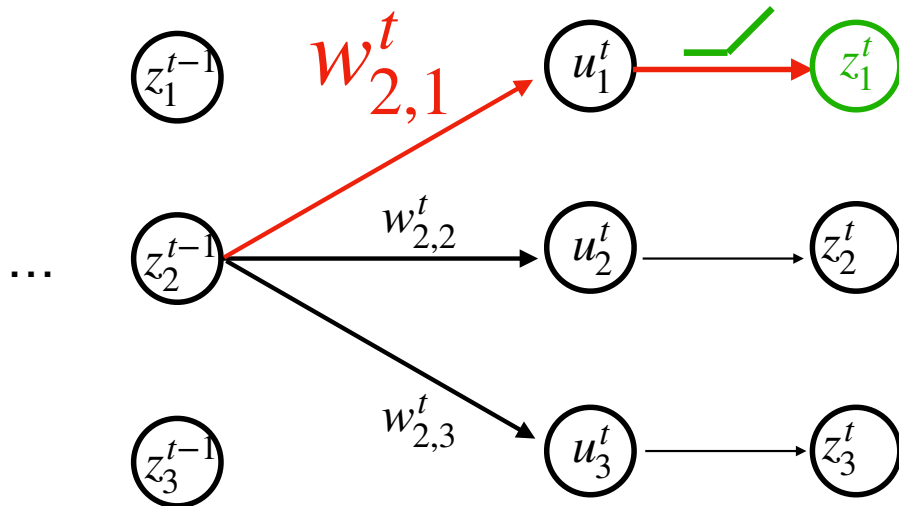
WLOG consider $\partial \hat{y} / \partial w_{2,1}^t$

$$\partial \hat{y} / \partial w_{2,1}^t = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$

The backward Pass

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .

Proof:

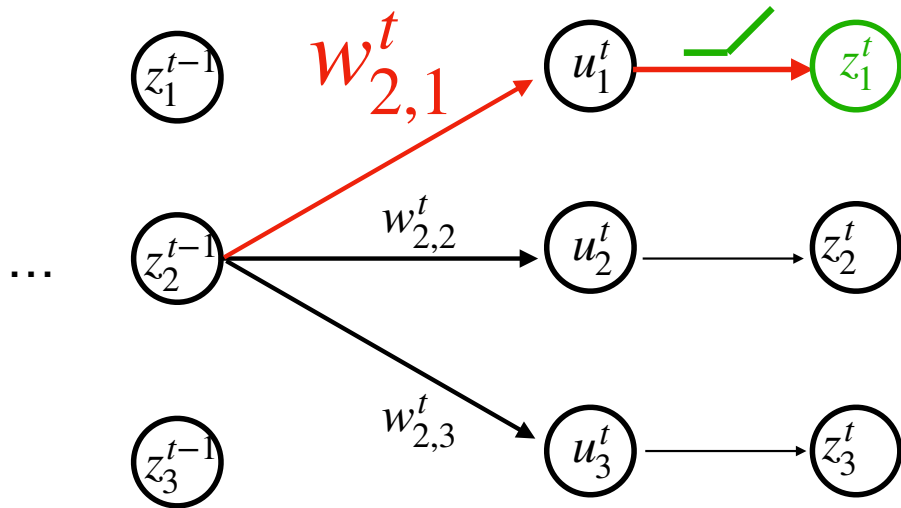


WLOG consider $\partial \hat{y} / \partial w_{2,1}^t$

$$\begin{aligned} \partial \hat{y} / \partial w_{2,1}^t &= \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t} \\ &= \frac{\partial \hat{y}}{\partial z_1^t} \cdot \sigma'(u_1^t) \cdot z_2^{t-1} \end{aligned}$$

The backward Pass

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .



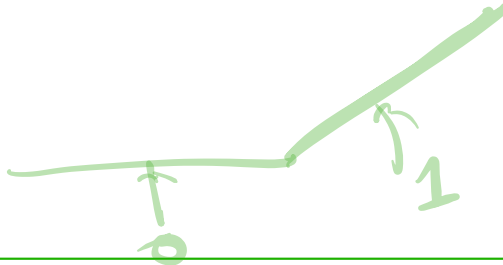
Proof:

WLOG consider $\partial \hat{y} / \partial w_{2,1}^t$

$$\partial \hat{y} / \partial w_{2,1}^t = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$

$$= \frac{\partial \hat{y}}{\partial z_1^t} \cdot \sigma'(u_1^t) \cdot z_2^{t-1}$$

Given by
assumption

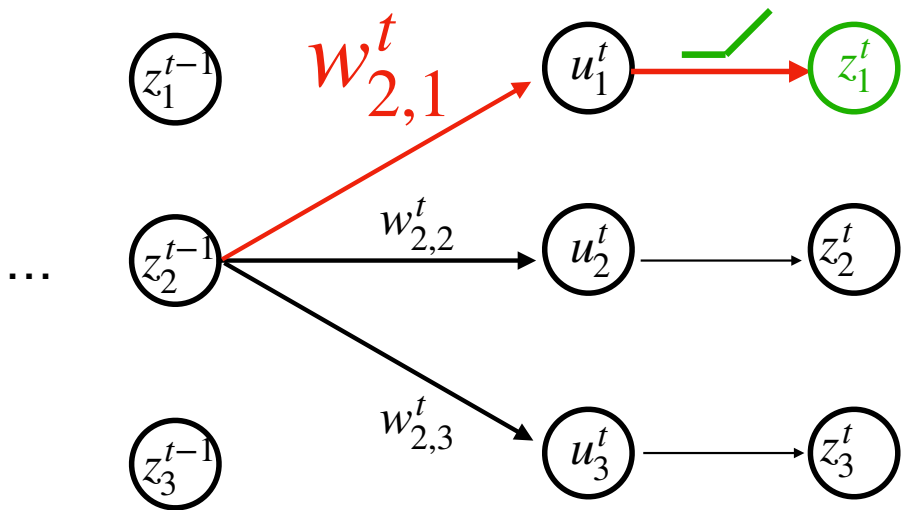


The backward Pass

$$\sigma(x) = \max\{x, 0\}$$

$$\sigma'(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ 0, & \text{else} \end{cases}$$

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .



Proof:

WLOG consider $\partial \hat{y} / \partial w_{2,1}^t$

$$\partial \hat{y} / \partial w_{2,1}^t = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$

...

$$= \left(\frac{\partial \hat{y}}{\partial z_1^t} \right) \cdot \left(\sigma'(u_1^t) \right) \cdot z_2^{t-1}$$

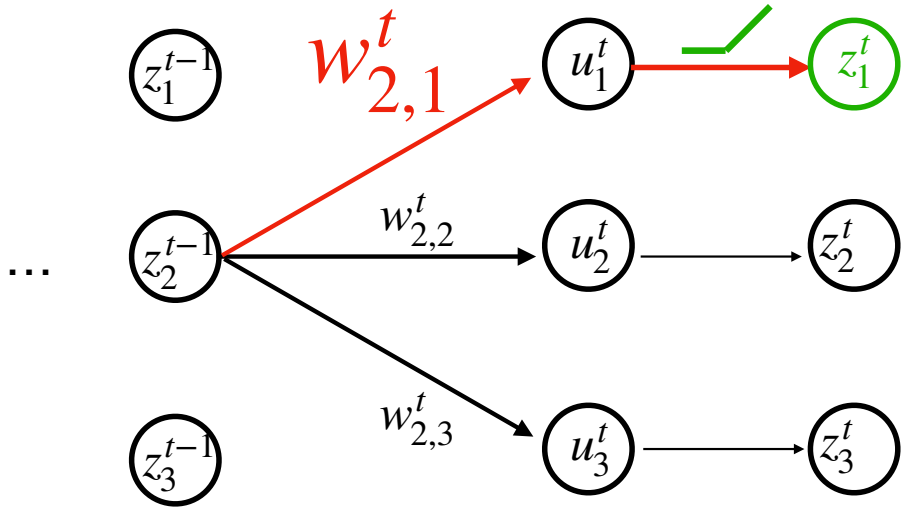
Given by
assumption

Derivative of
ReLU

The backward Pass

$$u_1^t = \dots z_2^{t-1} \cdot w_{2,1}^{t-1}$$

Claim: to compute $\partial \hat{y} / \partial w$, \forall edge w , it suffices to compute $\partial \hat{y} / \partial z$, \forall node z .



Proof:

WLOG consider $\partial \hat{y} / \partial w_{2,1}^t$

$$\partial \hat{y} / \partial w_{2,1}^t = \frac{\partial \hat{y}}{\partial z_1^t} \cdot \frac{\partial z_1^t}{\partial u_1^t} \cdot \frac{\partial u_1^t}{\partial w_{2,1}^t}$$

$$= \frac{\partial \hat{y}}{\partial z_1^t} \cdot \sigma'(u_1^t) \cdot z_2^{t-1}$$

Known from forward pass

Given by assumption

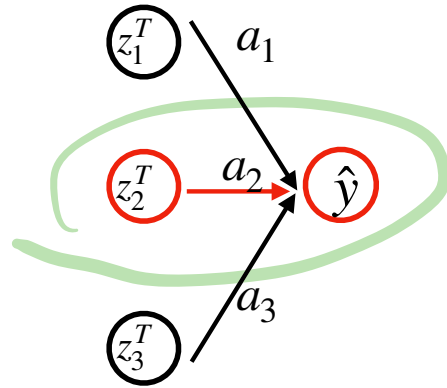
Derivative of ReLU

The backward Pass

We compute $\partial \hat{y} / \partial z^t$ backwards in time from $t = T$ to $t = 1$:

The backward Pass: base case

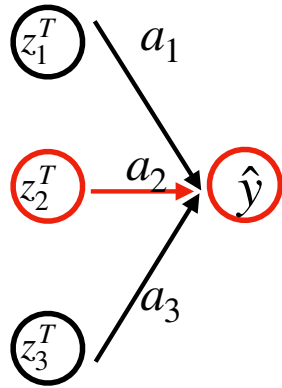
Base case: compute $\partial \hat{y} / \partial z^T$, for all node z at T -th Layer



$$\frac{\partial \hat{y}}{\partial z^T}$$

The backward Pass: base case

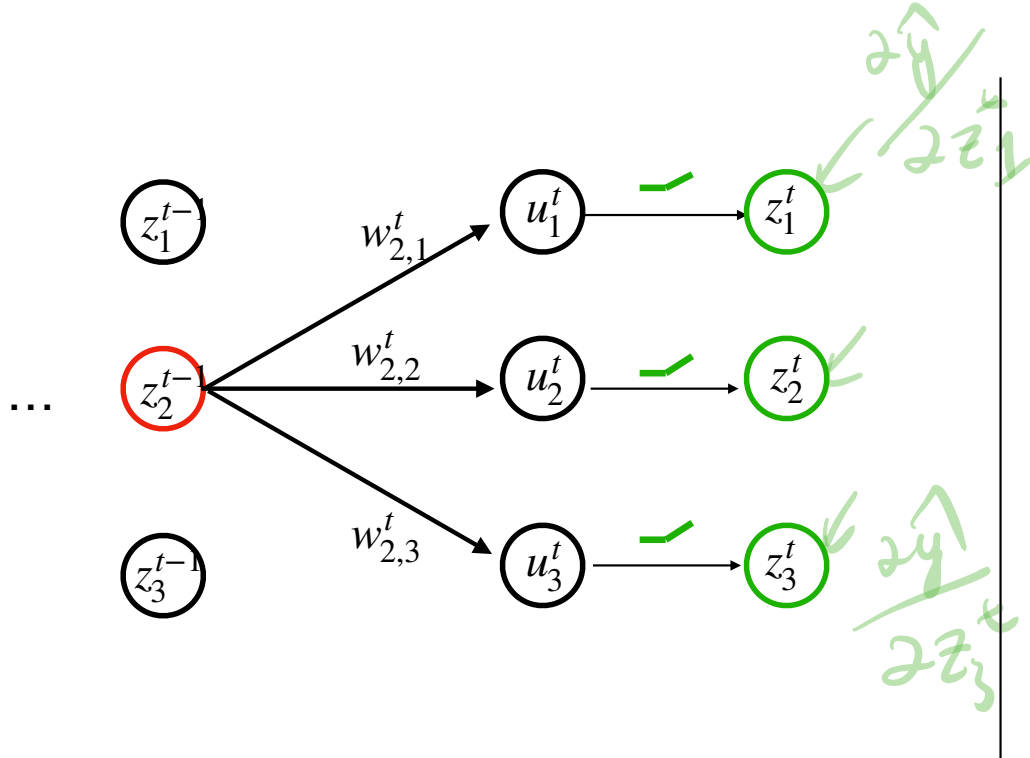
Base case: compute $\partial \hat{y} / \partial z^T$, for all node z at T -th Layer



$$\partial \hat{y} / \partial z_i^T = a_i$$

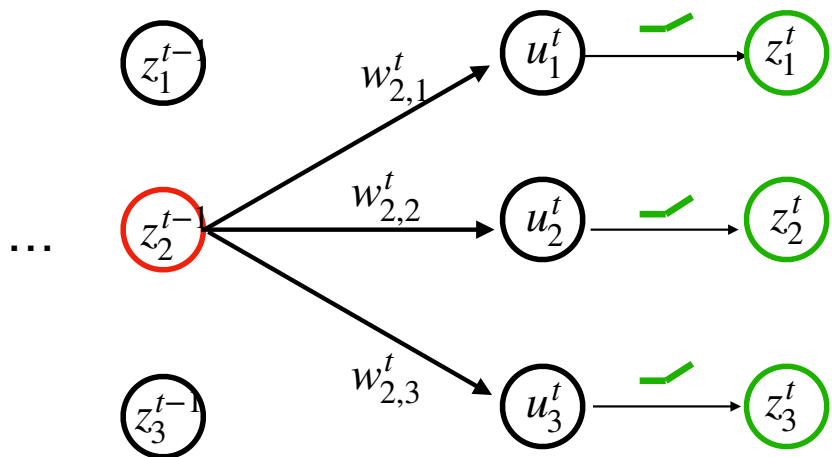
The backward Pass: induction step

Assume that we have computed $\partial \hat{y} / \partial z_i^t, \forall i$



The backward Pass: induction step

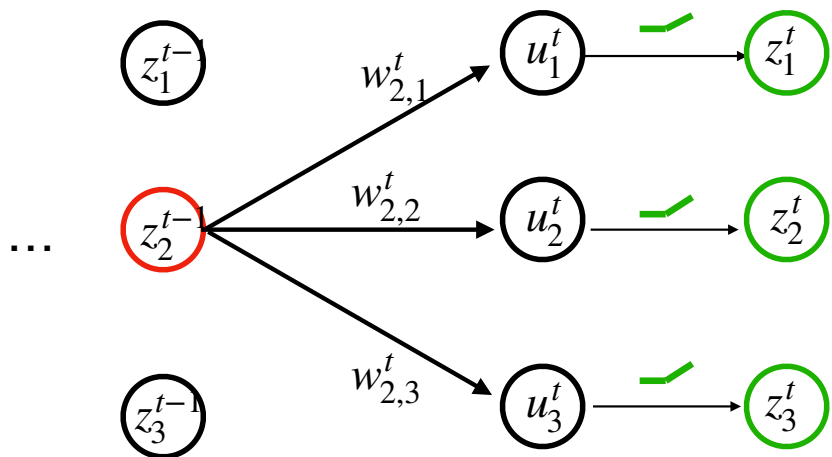
Assume that we have computed $\partial \hat{y} / \partial z_i^t, \forall i$



WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

The backward Pass: induction step

Assume that we have computed $\partial \hat{y} / \partial z_i^t, \forall i$

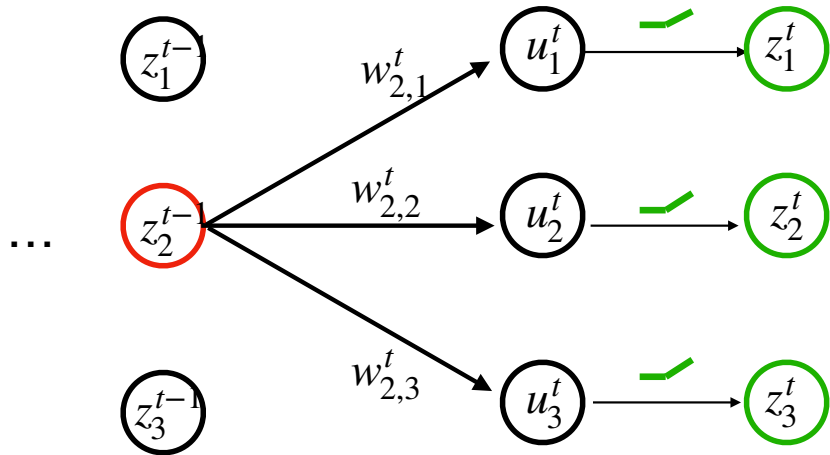


WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

Step 1: for all i , $\frac{\partial \hat{y}}{\partial u_i^t} = \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t}$

The backward Pass: induction step

Assume that we have computed $\partial \hat{y} / \partial z_i^t, \forall i$



WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

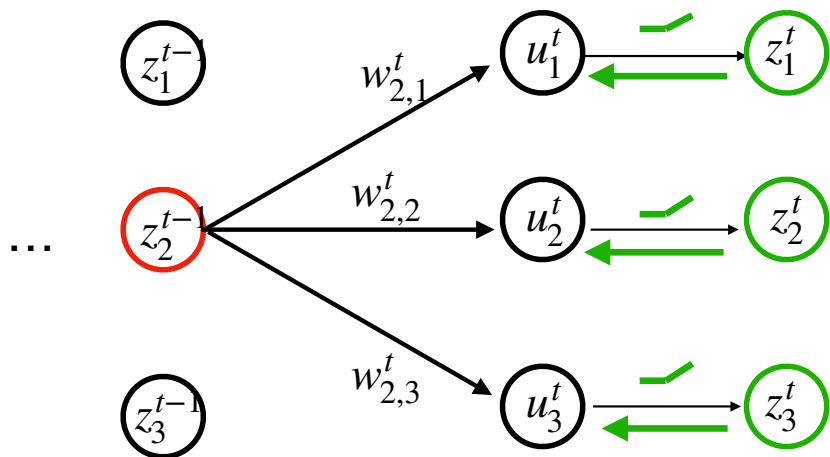
Step 1: for all i ,
$$\frac{\partial \hat{y}}{\partial u_i^t} = \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t}$$

$$= \frac{\partial \hat{y}}{\partial z_i^t} \cdot \sigma'(u_i^t)$$

Given by Assumption

The backward Pass: induction step

Assume that we have computed $\partial \hat{y} / \partial z_i^t, \forall i$



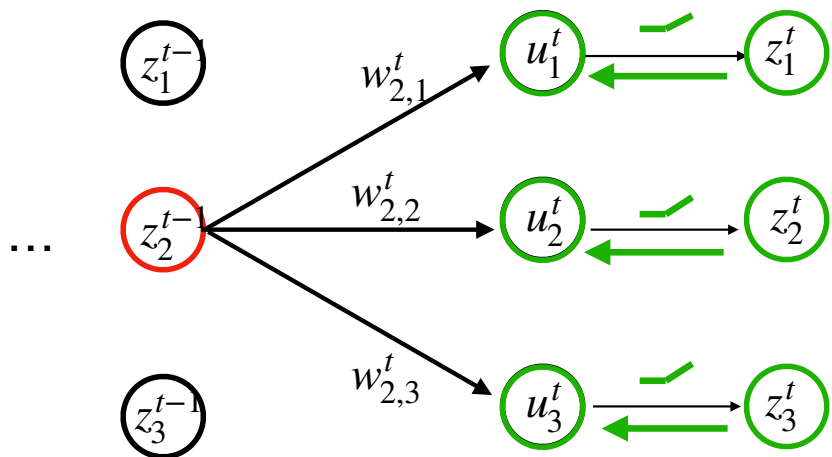
WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

Step 1: for all i ,
$$\frac{\partial \hat{y}}{\partial u_i^t} = \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t}$$

$$= \frac{\partial \hat{y}}{\partial z_i^t} \cdot \sigma'(u_i^t)$$

The backward Pass: induction step

Assume that we have computed $\partial \hat{y} / \partial z_i^t, \forall i$

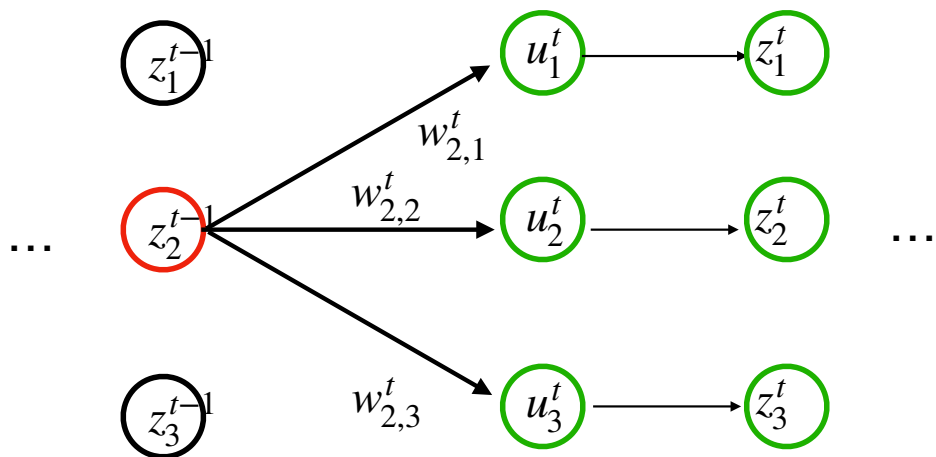


WLOG, consider $\partial \hat{y} / \partial z_2^{t-1}$

$$\begin{aligned} \text{Step 1: for all } i, \frac{\partial \hat{y}}{\partial u_i^t} &= \frac{\partial \hat{y}}{\partial z_i^t} \frac{\partial z_i^t}{\partial u_i^t} \\ &= \frac{\partial \hat{y}}{\partial z_i^t} \cdot \sigma'(u_i^t) \end{aligned}$$

The backward Pass: induction step

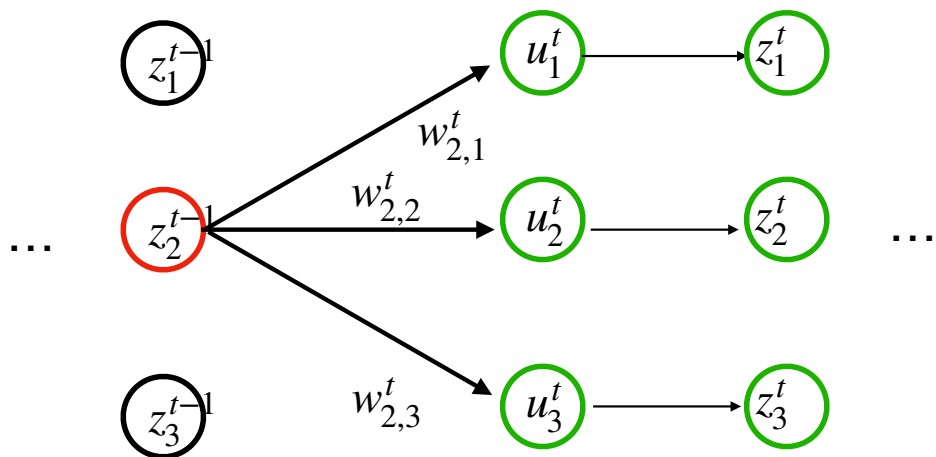
Assume that we have computed $\partial \mathcal{L} / \partial z_i^t, \forall i$



After step 1, we have $\partial \hat{y} / \partial u_i^t, \forall i$

The backward Pass: induction step

Assume that we have computed $\partial \mathcal{L} / \partial z_i^t, \forall i$



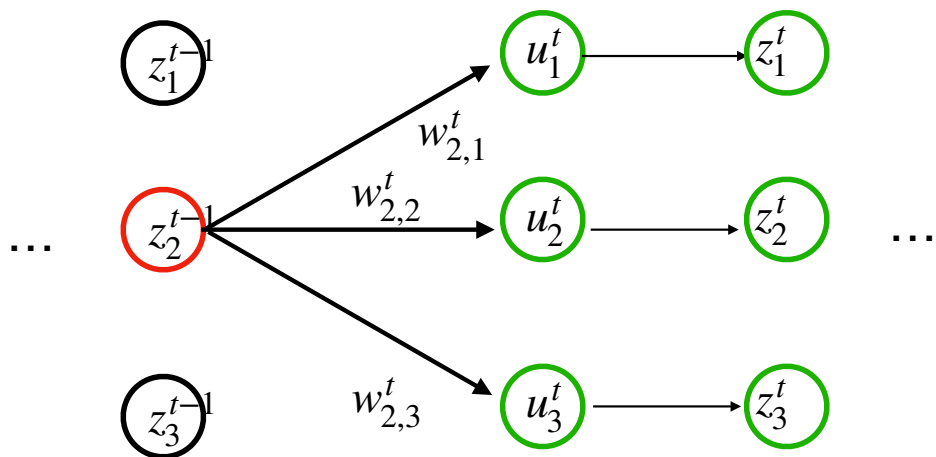
After step 1, we have $\partial \hat{y} / \partial u_i^t, \forall i$

Via multivariate chain rule:

$$\hat{y} = f(u_1(x), u_2(x), u_3(x))$$
$$\frac{\partial \hat{y}}{\partial x} = \frac{\partial \hat{y}}{\partial u_1} \frac{\partial u_1}{\partial x} + \frac{\partial \hat{y}}{\partial u_2} \frac{\partial u_2}{\partial x} + \frac{\partial \hat{y}}{\partial u_3} \frac{\partial u_3}{\partial x}$$

The backward Pass: induction step

Assume that we have computed $\partial \mathcal{L} / \partial z_i^t, \forall i$



After step 1, we have $\partial \hat{y} / \partial u_i^t, \forall i$

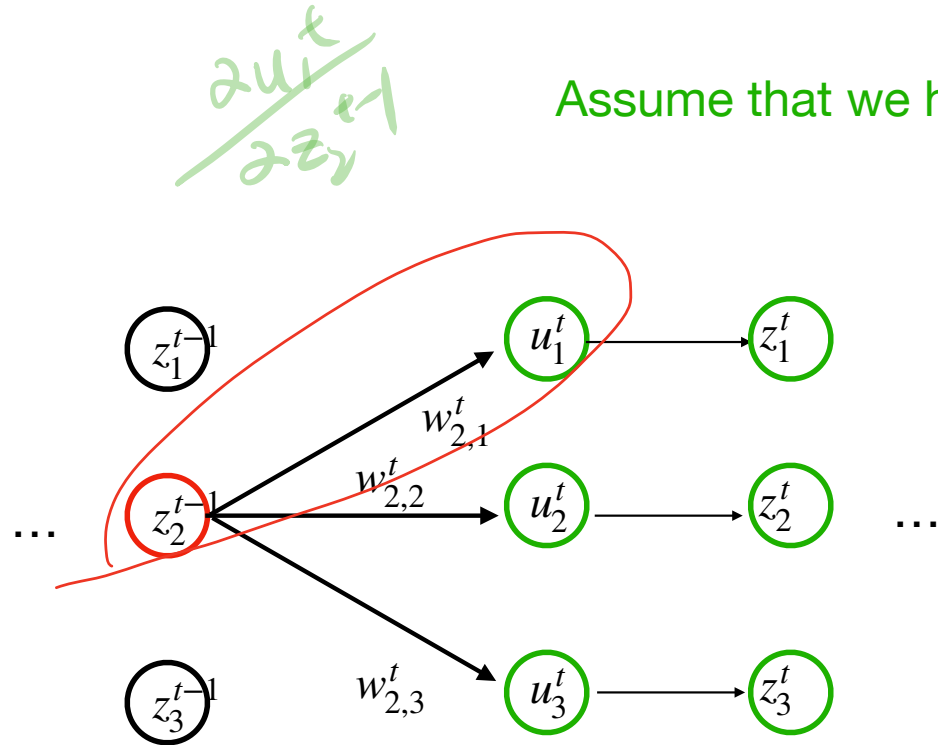
Via multivariate chain rule:

$$\text{Step 2: } \frac{\partial \hat{y}}{\partial z_2^{t-1}} = \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}}$$

Computed after step 1

The backward Pass: induction step

Assume that we have computed $\partial \mathcal{L} / \partial z_i^t, \forall i$



After step 1, we have $\partial \hat{y} / \partial u_i^t, \forall i$

Via multivariate chain rule:

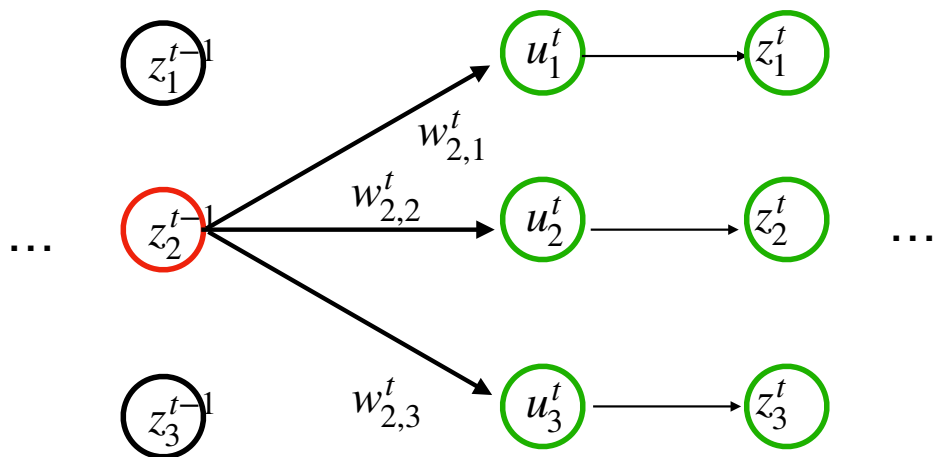
Step 2:

$$\frac{\partial \hat{y}}{\partial z_2^{t-1}} = \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}}$$

$$= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t$$

The backward Pass: induction step

Assume that we have computed $\partial \mathcal{L} / \partial z_i^t, \forall i$



After step 1, we have $\partial \hat{y} / \partial u_i^t, \forall i$

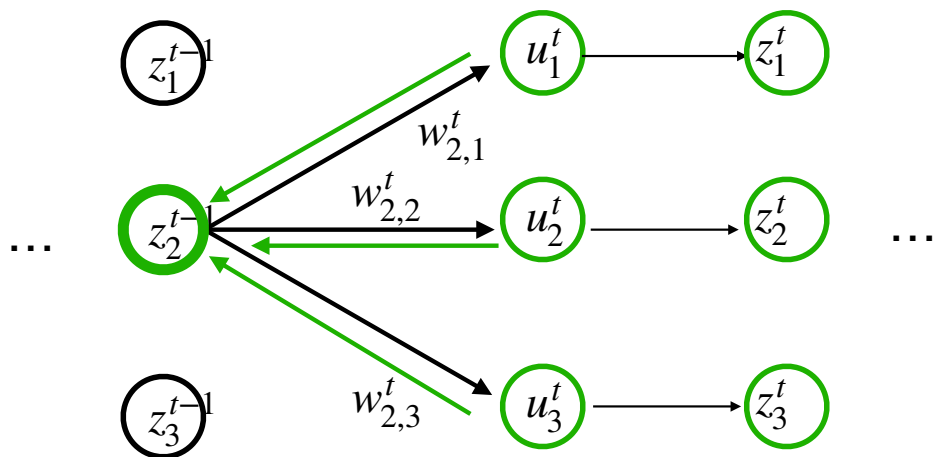
Via multivariate chain rule:

$$\begin{aligned} \text{Step 2: } \frac{\partial \hat{y}}{\partial z_2^{t-1}} &= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}} \\ &= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t \end{aligned}$$

We are done at node z_2^{t-1} !

The backward Pass: induction step

Assume that we have computed $\partial \mathcal{L} / \partial z_i^t, \forall i$



After step 1, we have $\partial \hat{y} / \partial u_i^t, \forall i$

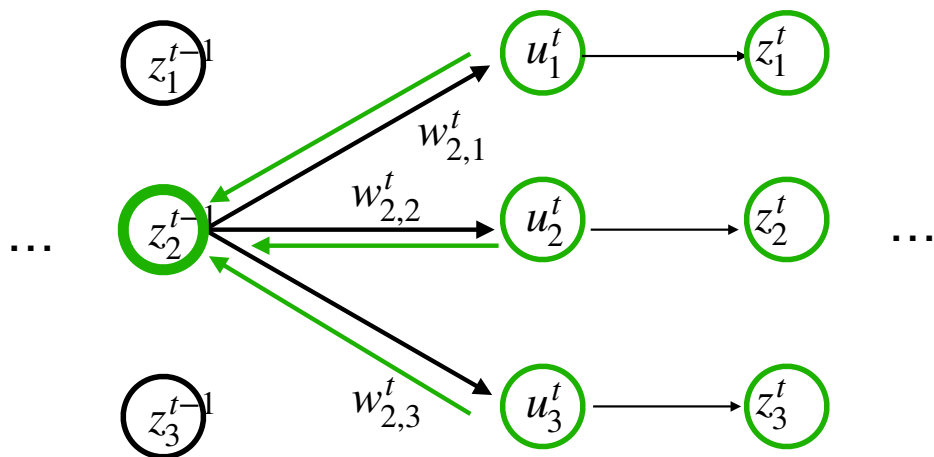
Via multivariate chain rule:

$$\begin{aligned} \text{Step 2: } \frac{\partial \hat{y}}{\partial z_2^{t-1}} &= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}} \\ &= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t \end{aligned}$$

We are done at node z_2^{t-1} !

The backward Pass: induction step

Assume that we have computed $\partial \mathcal{L} / \partial z_i^t, \forall i$



After step 1, we have $\partial \hat{y} / \partial u_i^t, \forall i$

Via multivariate chain rule:

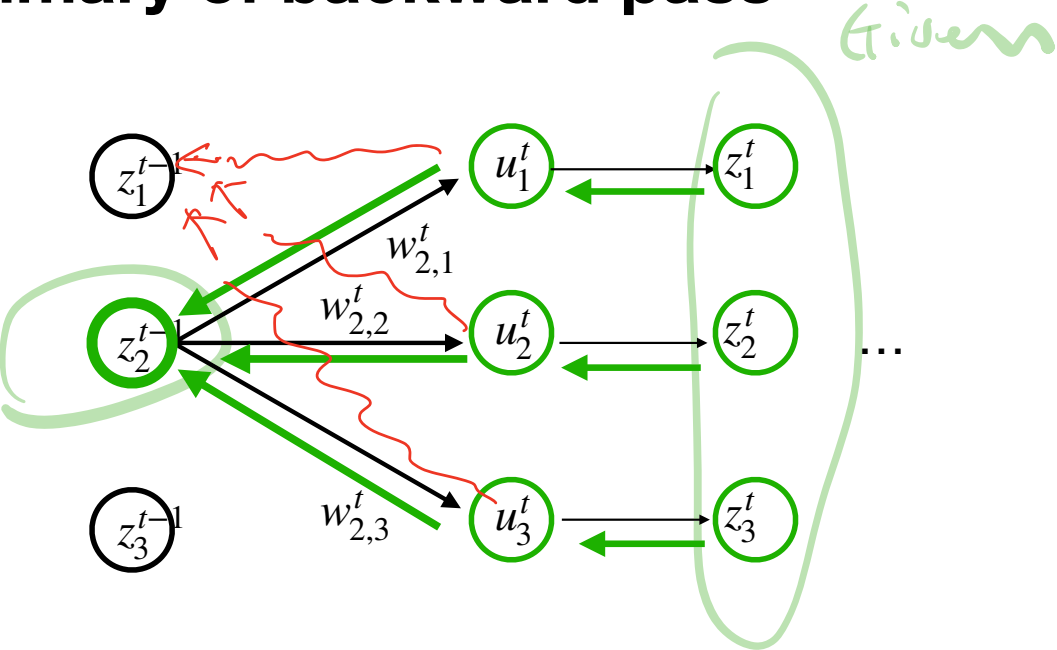
$$\begin{aligned} \text{Step 2: } \frac{\partial \hat{y}}{\partial z_2^{t-1}} &= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \frac{\partial u_i^t}{\partial z_2^{t-1}} \\ &= \sum_{i=1}^K \frac{\partial \hat{y}}{\partial u_i^t} \cdot w_{2,i}^t \end{aligned}$$

We are done at node z_2^{t-1} !

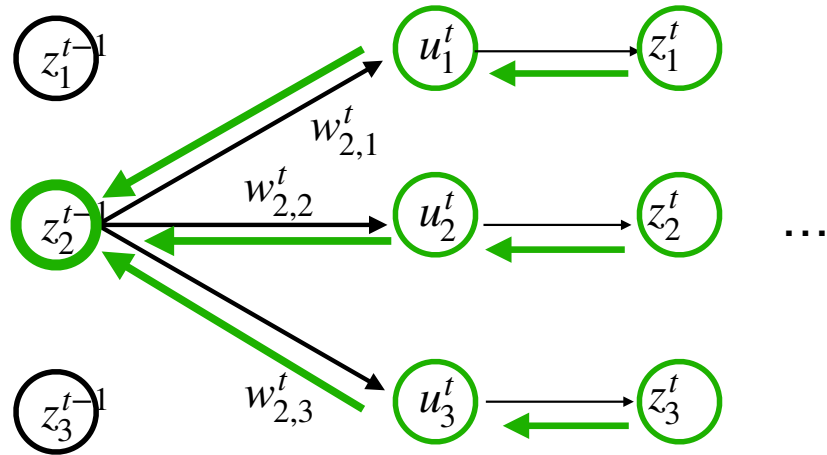
Repeat this for all $z_i^{t-1}, \forall i$

Summary of backward pass

$$\frac{\partial z_2}{\partial z_2^{t-1}}$$

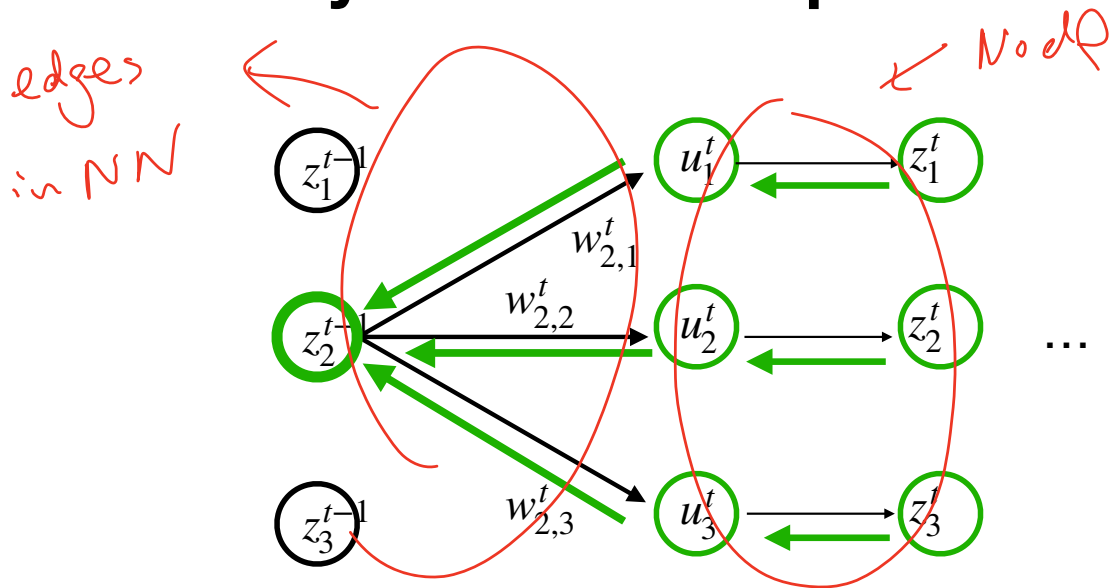


Summary of backward pass



The computation from $\partial \hat{y} / z^t$ to $\partial \hat{y} / z^{t-1}$ is the # of all edges in the sub-graph

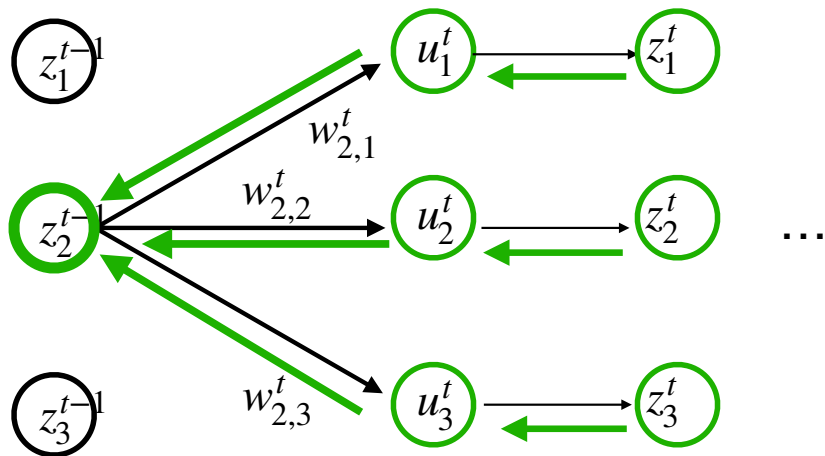
Summary of backward pass



The computation from $\partial \hat{y} / z^t$ to $\partial \hat{y} / z^{t-1}$ is the # of all edges in the sub-graph

Total computation: # of edges + # of nodes!

Summary of backward pass



The computation from $\partial \hat{y} / z^t$ to $\partial \hat{y} / z^{t-1}$ is the # of all edges in the sub-graph

Total computation: # of edges + # of nodes!

Exercise: can you express backward pass in matrix-vector format?

Summary for today

1. Naively compute all derivatives wrt edges using chain rule takes $(E + V)^2$ time

Summary for today

1. Naively compute all derivatives wrt edges using chain rule takes $(E + V)^2$ time
2. Backpropagation: forward pass & backward pass takes $O(E + V)$ time

Summary for today

1. Naively compute all derivatives wrt edges using chain rule takes $(E + V)^2$ time

2. Backpropagation: forward pass & backward pass takes $O(E + V)$ time

Forward pass: $x = z^0 \rightarrow u^1 \rightarrow z^1 \rightarrow \dots \rightarrow z^t \rightarrow u^{t+1} \rightarrow z^{t+1} \dots \rightarrow z^T \rightarrow \hat{y}$

Summary for today

1. Naively compute all derivatives wrt edges using chain rule takes $(E + V)^2$ time

2. Backpropagation: forward pass & backward pass takes $O(E + V)$ time

Forward pass: $x = z^0 \rightarrow u^1 \rightarrow z^1 \rightarrow \dots \rightarrow z^t \rightarrow u^{t+1} \rightarrow z^{t+1} \dots \rightarrow z^T \rightarrow \hat{y}$

Backward pass: $\frac{\partial \hat{y}}{\partial z^T} \rightarrow \frac{\partial \hat{y}}{\partial z^{T-1}} \rightarrow \dots \rightarrow \frac{\partial \hat{y}}{\partial z^1}$ + claim
 $\Rightarrow \frac{\partial \hat{y}}{\partial w}$