

10-701: Introduction to Machine Learning

Lecture 10 – Backpropagation

Nari Johnson, narij@andrew.cmu.edu

Slides adapted from Henry Chai

9/29/25

Recall: Stochastic Gradient Descent for Learning

- Input: $\mathcal{D} = \{\underline{(\mathbf{x}^{(n)}, y^{(n)})}\}_{n=1}^N, \eta^{(0)}$
- Initialize all weights $\underline{W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}}$ to small, random numbers and set $t = 0$
- While TERMINATION CRITERION is not satisfied
 - For $i \in \text{shuffle}(\{1, \dots, N\}) \leftarrow$
 - For $l = 1, \dots, L \leftarrow$
 - Compute $G^{(l)} = \nabla_{W^{(l)}} \ell^{(i)} \left(\underline{W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}} \right)$
 - Update $W^{(l)}: \underline{W_{(t+1)}^{(l)} = W_{(t)}^{(l)} - \eta_0 G^{(l)}}$
 - Increment $t: \underline{t = t + 1}$
- Output: $\underline{W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}}$

Matrix Calculus

| <i>Types of Derivatives</i> | scalar | vector | matrix |
|-----------------------------|--|---|---|
| scalar | $\frac{\partial y}{\partial x}$ | $\frac{\partial \mathbf{y}}{\partial x}$ | $\frac{\partial \mathbf{Y}}{\partial x}$ |
| vector | $\frac{\partial y}{\partial \mathbf{x}}$ | $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ | $\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$ |
| matrix | $\frac{\partial y}{\partial \mathbf{X}}$ | $\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$ | $\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$ |

Denominator

Numerator

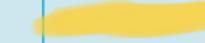
Table courtesy of Matt Gormley

Matrix Calculus: Denominator Layout

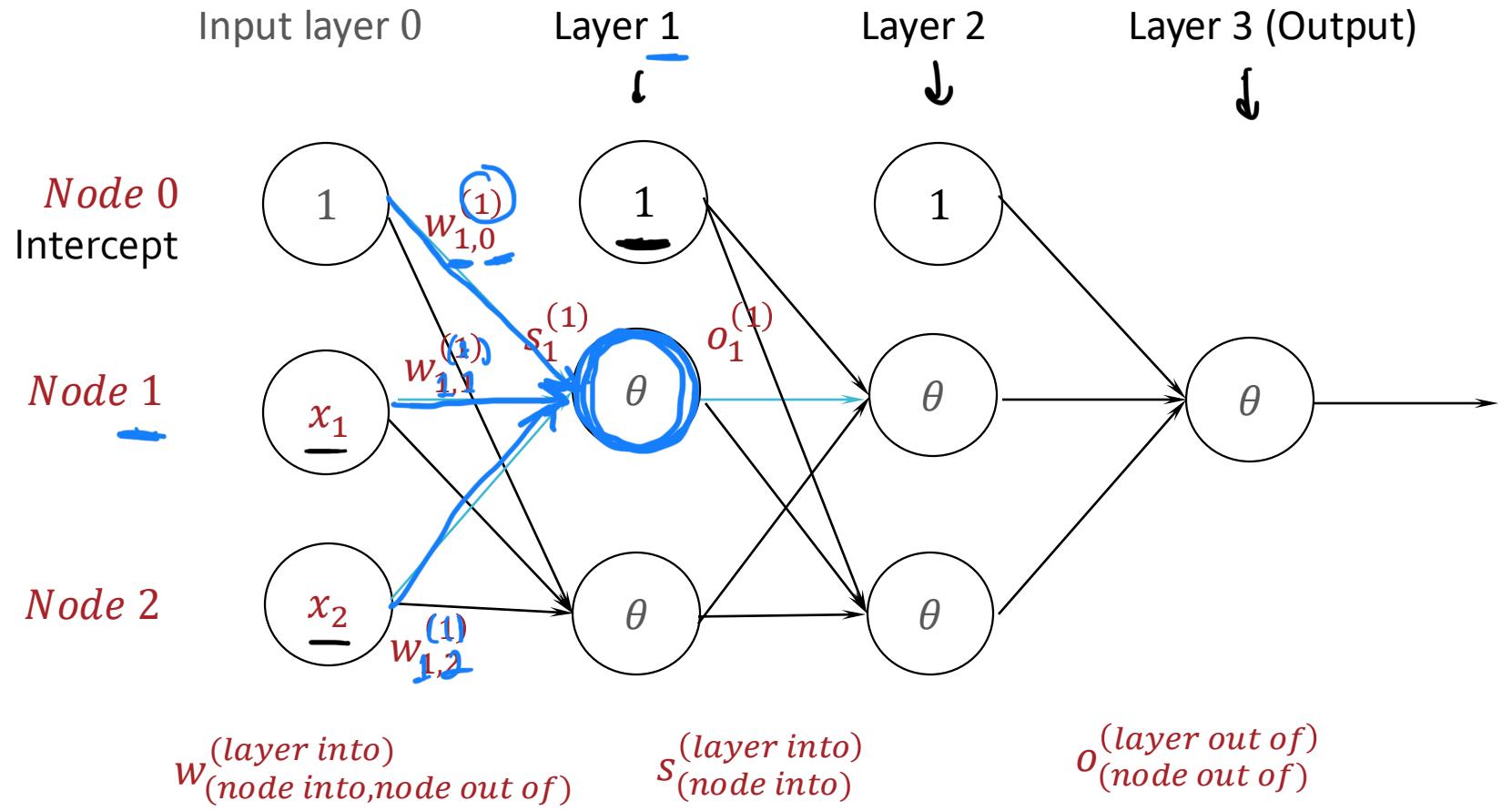
- Derivatives of a scalar always have the *same shape* as the entity that the derivative is being taken with respect to.

| <i>Types of Derivatives</i> | scalar |
|-----------------------------|---|
| scalar | $\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x} \right]$ |
| vector | $\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$ |
| matrix | $\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial X_{11}} & \frac{\partial y}{\partial X_{12}} & \cdots & \frac{\partial y}{\partial X_{1Q}} \\ \frac{\partial y}{\partial X_{21}} & \frac{\partial y}{\partial X_{22}} & \cdots & \frac{\partial y}{\partial X_{2Q}} \\ \vdots & & & \vdots \\ \frac{\partial y}{\partial X_{P1}} & \frac{\partial y}{\partial X_{P2}} & \cdots & \frac{\partial y}{\partial X_{PQ}} \end{bmatrix}$ |

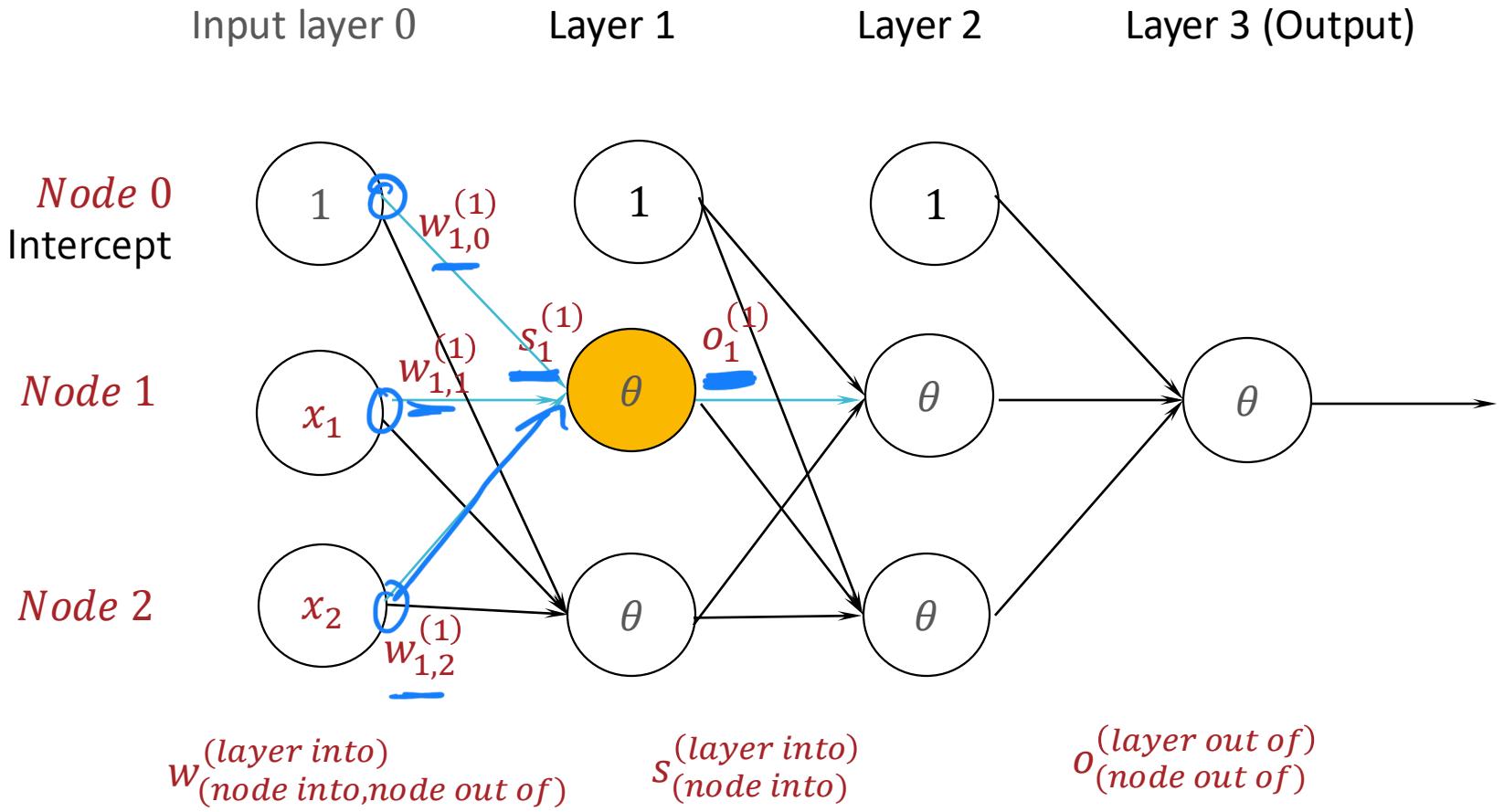
Matrix Calculus: Denominator Layout

| <i>Types of Derivatives</i> | scalar | vector |
|-----------------------------|--|---|
| scalar | $\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x} \right]$ | $\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x} & \frac{\partial y_2}{\partial x} & \cdots & \frac{\partial y_N}{\partial x} \end{bmatrix}$   |
| vector | $\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$  | $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_N}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_P} & \frac{\partial y_2}{\partial x_P} & \cdots & \frac{\partial y_N}{\partial x_P} \end{bmatrix}$ |

Recall our
neural
network
notation...



Recall our
neural
network
notation...

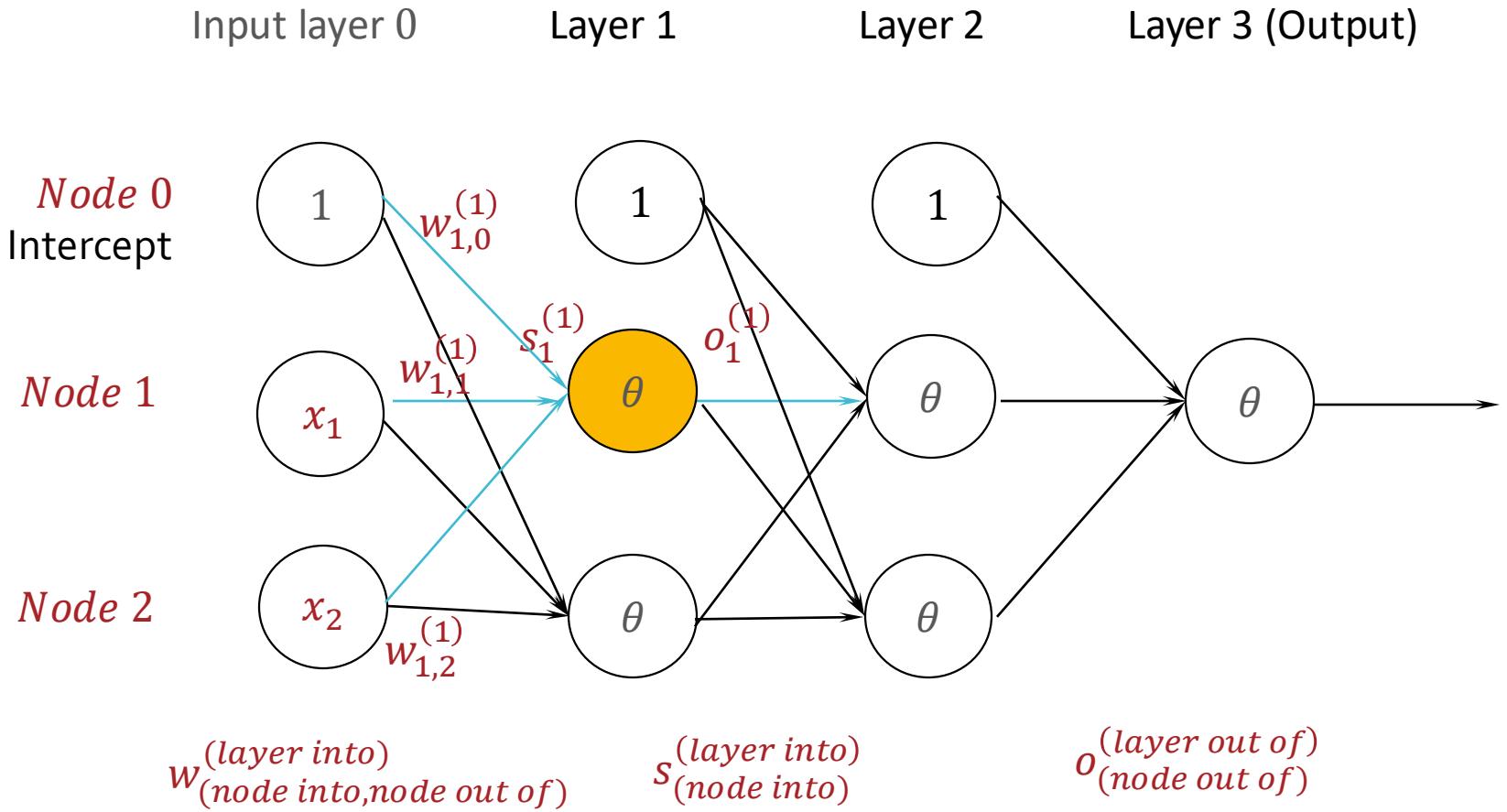


Write equations for the signal $s_1^{(1)}$ and output $o_1^{(1)}$ for the highlighted node:

$$s_1^{(1)} = w_{1,0}^{(1)}(1) + w_{1,1}^{(1)}x_1 + w_{1,2}^{(1)}x_2$$

$$o_1^{(1)} = \theta(s_1^{(1)}).$$

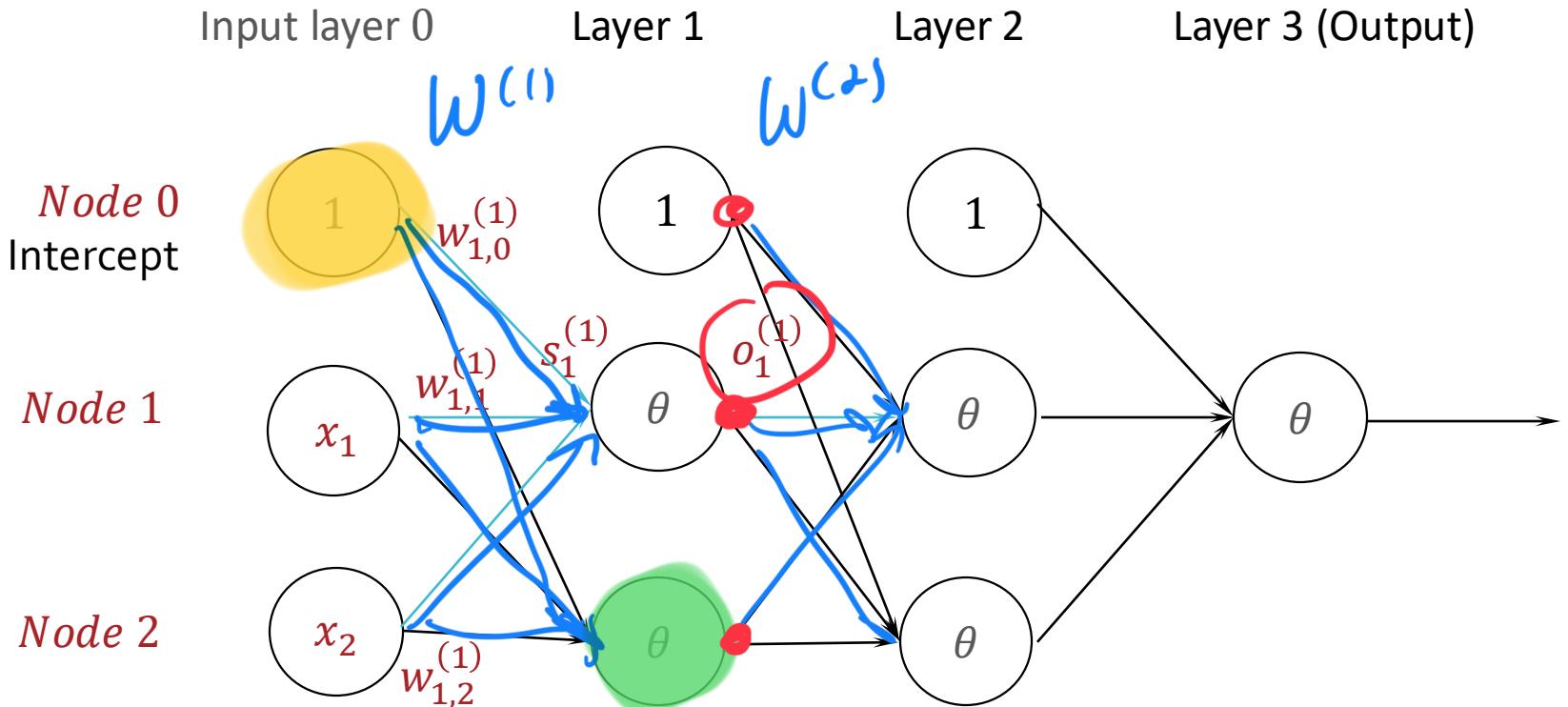
Recall our
neural
network
notation...



Write equations for the signal $s_1^{(1)}$ and output $o_1^{(1)}$ for the highlighted node:

$$s_1^{(1)} = w_{1,0}^{(1)} + w_{1,1}^{(1)}x_1 + w_{1,2}^{(1)}x_2 \quad o_1^{(1)} = \theta(s_1^{(1)})$$

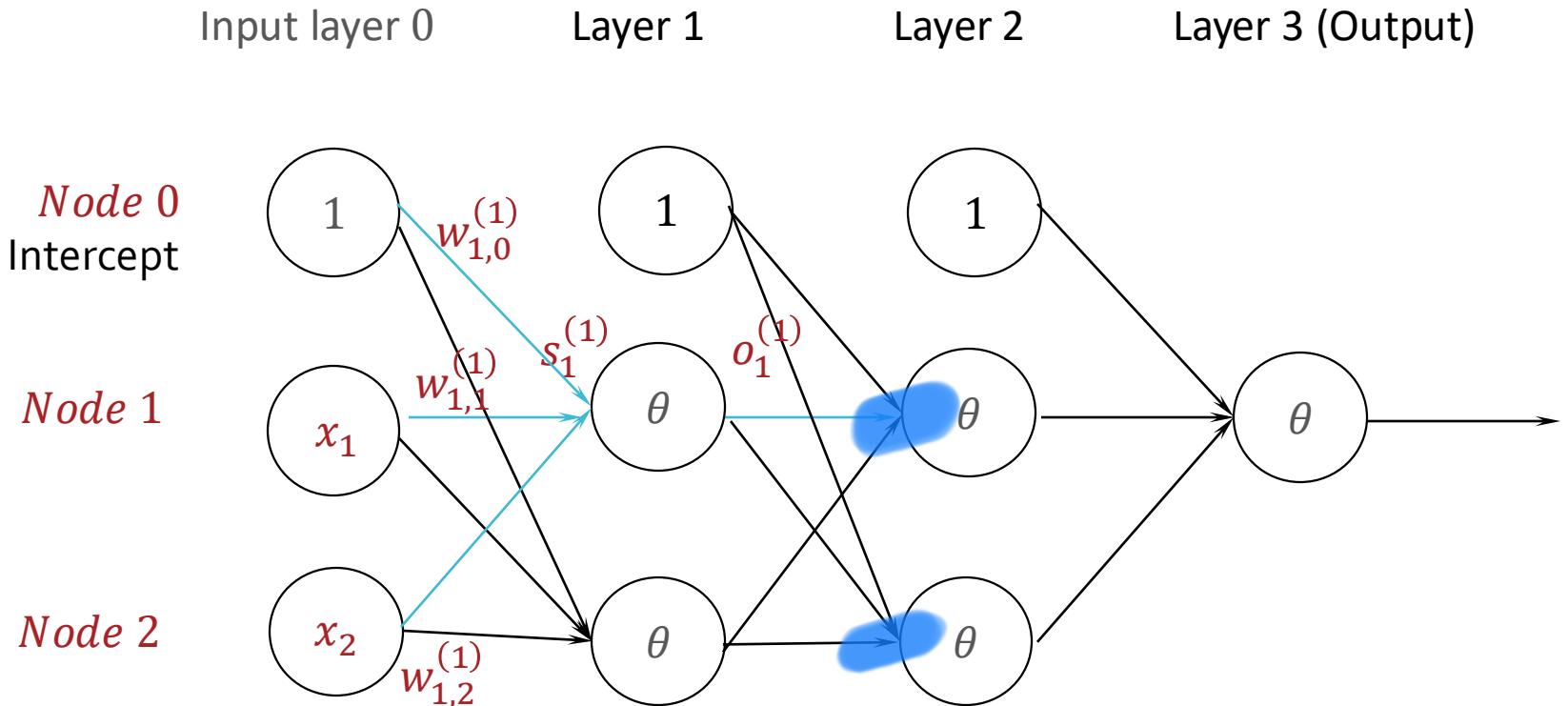
Recall our
neural
network
notation...



We can also write our weights, signals, and outputs using matrices:

$$w^{(1)} = \begin{matrix} \text{pointing into} \\ \left[\begin{matrix} w_{1,0}^{(1)} & w_{1,1}^{(1)} & w_{1,2}^{(1)} \\ w_{2,0}^{(1)} & w_{2,1}^{(1)} & w_{2,2}^{(1)} \end{matrix} \right] \end{matrix}$$
$$o_1^{(1)} = \begin{matrix} \text{pointing out of} \\ \left[\begin{matrix} 1 \\ o_1^{(1)} \\ o_2^{(1)} \end{matrix} \right] \end{matrix} .$$

Recall our
neural
network
notation...



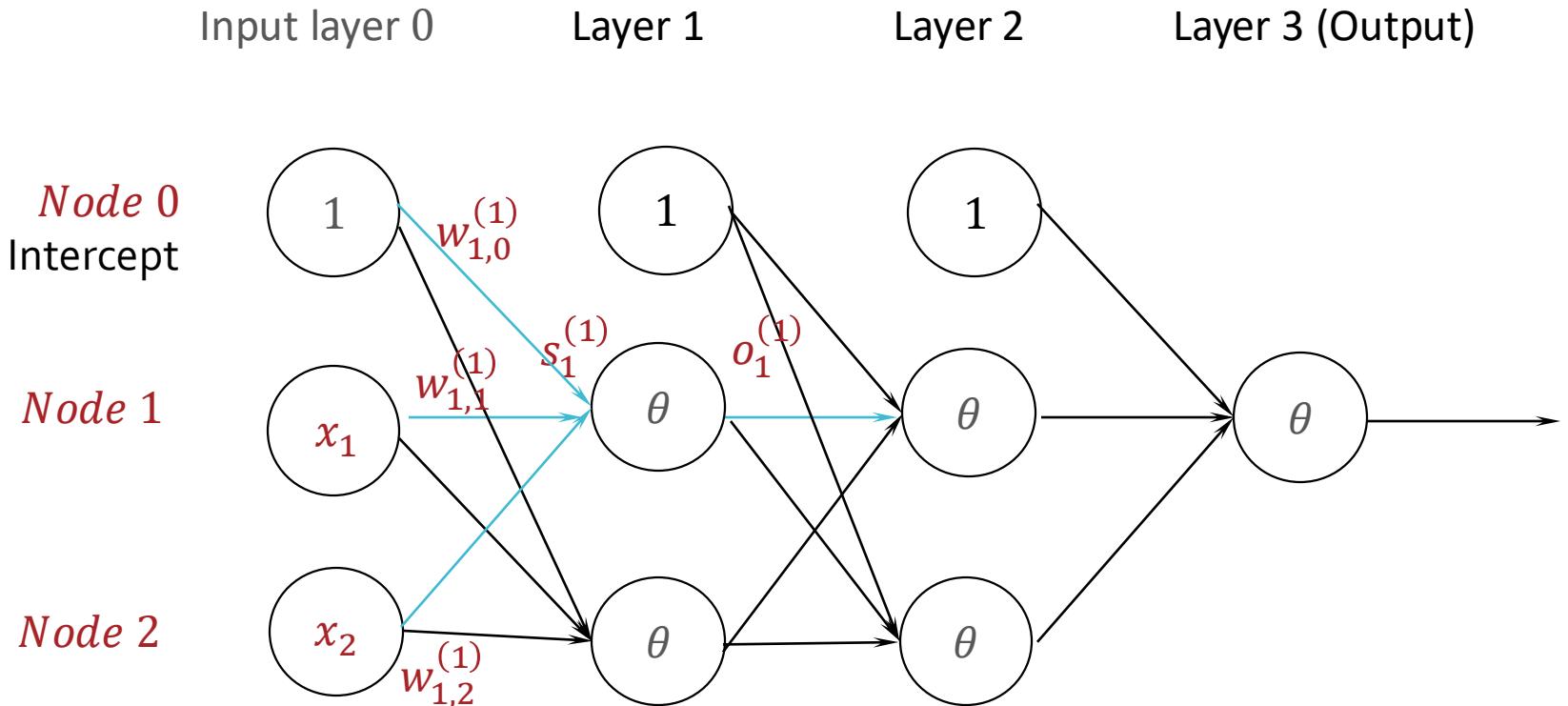
out of $W^{(1)} = \begin{bmatrix} w_{1,0}^{(1)} & w_{1,1}^{(1)} & w_{1,2}^{(1)} \\ w_{2,0}^{(1)} & w_{2,1}^{(1)} & w_{2,2}^{(1)} \end{bmatrix}$ into $\underline{o}^{(1)} = \begin{bmatrix} 1 \\ o_1^{(1)} \\ o_2^{(1)} \end{bmatrix}$

$\begin{bmatrix} 2,1 \end{bmatrix} \curvearrowleft \begin{bmatrix} 2,3 \end{bmatrix} \begin{bmatrix} 3,1 \end{bmatrix}$

Can you write $s^{(2)}$ as a matrix product of $W^{(2)}$ and $\underline{o}^{(1)}$?
What shape is each matrix?

$s^{(2)} = \underline{W}^{(2)} \underline{o}^{(1)}$.

Recall our
neural
network
notation...



Can you write $s^{(2)}$ as a matrix product of $W^{(2)}$ and $o^{(1)}$?
What shape is each matrix?

$$s^{(2)} = W^{(2)} o^{(1)}$$

$$s^{(2)}: [2, 1], W^{(2)}: [2, 3], o^{(1)}: [3, 1]$$

Recall: the chain rule

The chain rule tells us how to differentiate composite functions.

It is important for neural networks because the output of our neural network is a composite function of the network weights.

A simple example:

$$(x, a, b)$$

$$a = 2x^3 - 4x$$
$$b = \sin(a)$$
$$b = \sin(2x^3 - 4x)$$

$$\frac{\partial b}{\partial x} = \frac{\partial b}{\partial a} \cdot \frac{\partial a}{\partial x}$$

$$= (\cos(a)) \cdot (6x^2 - 4).$$

Recall: the chain rule

The chain rule tells us how to differentiate composite functions.

It is important for neural networks because the output of our neural network is a composite function of the network weights.

A simple example:

$$\begin{aligned} a &= 2x^3 - 4x \\ b &= \sin(a) \end{aligned}$$

$$\begin{aligned} \frac{\partial b}{\partial x} &= \frac{\partial b}{\partial a} \frac{\partial a}{\partial x} \\ &= \cos(a) (6x^2 - 4) \end{aligned}$$

Applying the chain rule to our neural network example

Recall from our toy neural network that:

$$\underline{s_1^{(1)}} = \underline{w_{1,0}^{(1)}} + \underline{w_{1,1}^{(1)}x_1} + \underline{w_{1,2}^{(1)}x_2} \quad \underline{o_1^{(1)}} = \theta(\underline{s_1^{(1)}})$$

Use the chain rule to calculate the following partial derivative.

Assume that the activation function $\theta(-) = \tanh(-)$

$$\frac{\partial o_1^{(1)}}{\partial w_{1,1}^{(1)}} =$$

Hint: $\frac{\partial}{\partial z} \tanh(z) = 1 - (\tanh(z))^2$

- A: $= \left(1 - (o_1^{(1)})^2 \right) (x_1)$ B: $= \left(1 - o_1^{(1)} \right) (w_{1,1}^{(1)})$
- C: $= \left(1 + o_1^{(1)} \right) (w_{1,1}^{(1)} x_1)$ D: $= o_1^{(1)} w_{1,1}^{(1)}$

$$\frac{\partial o_1^{(1)}}{\partial w_{111}^{(1)}} := \cancel{\frac{\partial o_1^{(1)}}{\partial s_1^{(1)}}} \cdot \underline{\frac{\partial s_1^{(1)}}{\partial w_{111}^{(1)}}}$$

-

w
*x*₁

$$o_1^{(1)} = \tanh(s_1^{(1)})$$

$$\begin{aligned} & \frac{\partial}{\partial s_1^{(1)}} (\tanh(s_1^{(1)})) \\ &= 1 - \underbrace{(\tanh(s_1^{(1)}))}_{}^2 \\ &= 1 - (o_1^{(1)})^2 \end{aligned}$$

Applying the chain rule to our neural network example

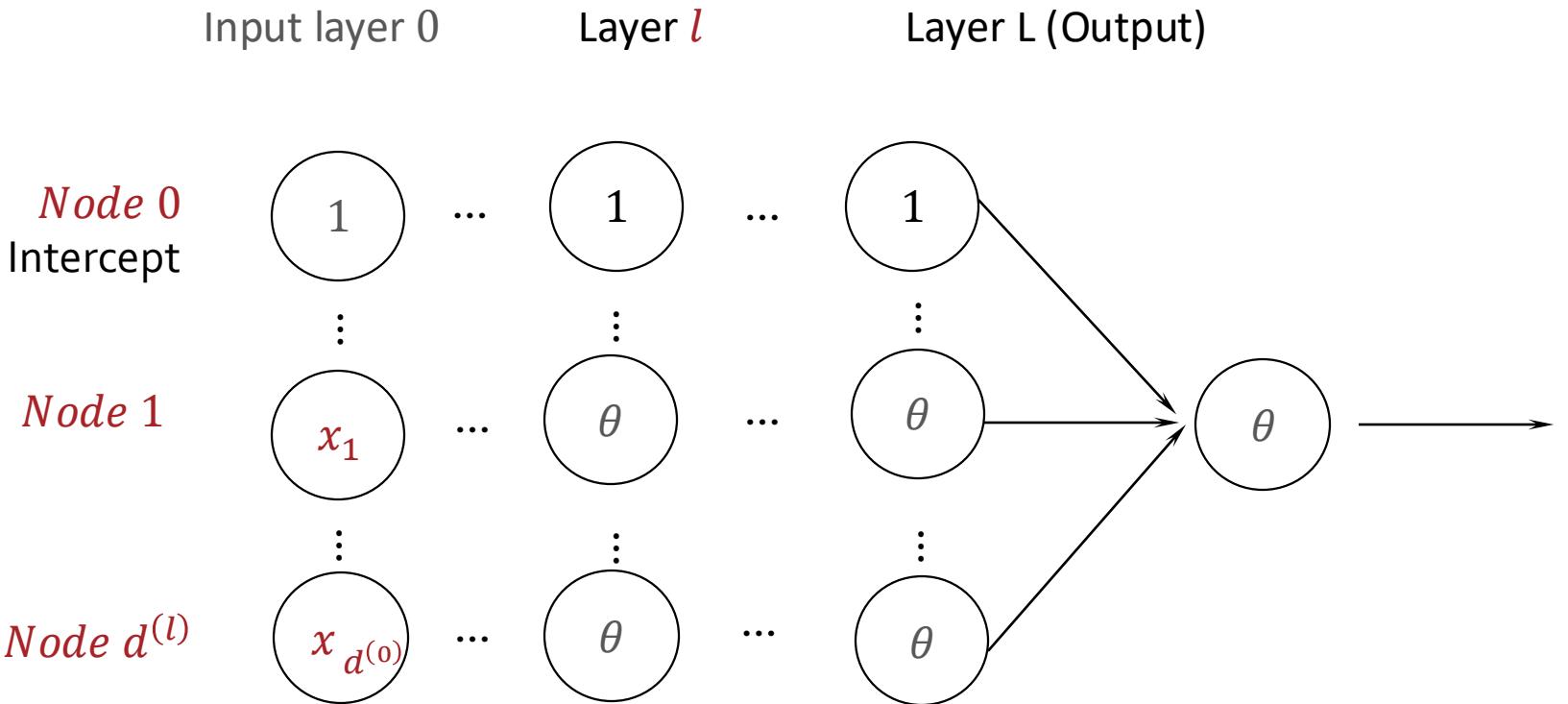
Recall from our toy neural network that:

$$s_1^{(1)} = w_{1,0}^{(1)} + w_{1,1}^{(1)}x_1 + w_{1,2}^{(1)}x_2 \quad o_1^{(1)} = \theta(s_1^{(1)})$$

Use the chain rule to calculate the following partial derivative.
Assume that the activation function $\theta(-) = \tanh(-)$

$$\begin{aligned}\frac{\partial o_1^{(1)}}{\partial w_{1,1}^{(1)}} &= \text{Hint: } \frac{\partial}{\partial z} \tanh(z) = 1 - (\tanh(z))^2 \\ &= \frac{\partial o_1^{(1)}}{\partial s_1^{(1)}} \frac{\partial s_1^{(1)}}{\partial w_{1,1}^{(1)}} \\ &= \left(1 - (\tanh(s_1^{(1)}))^2\right)(x_1) \\ &= \left(1 - (o_1^{(1)})^2\right)(x_1)\end{aligned}$$

Zooming out...



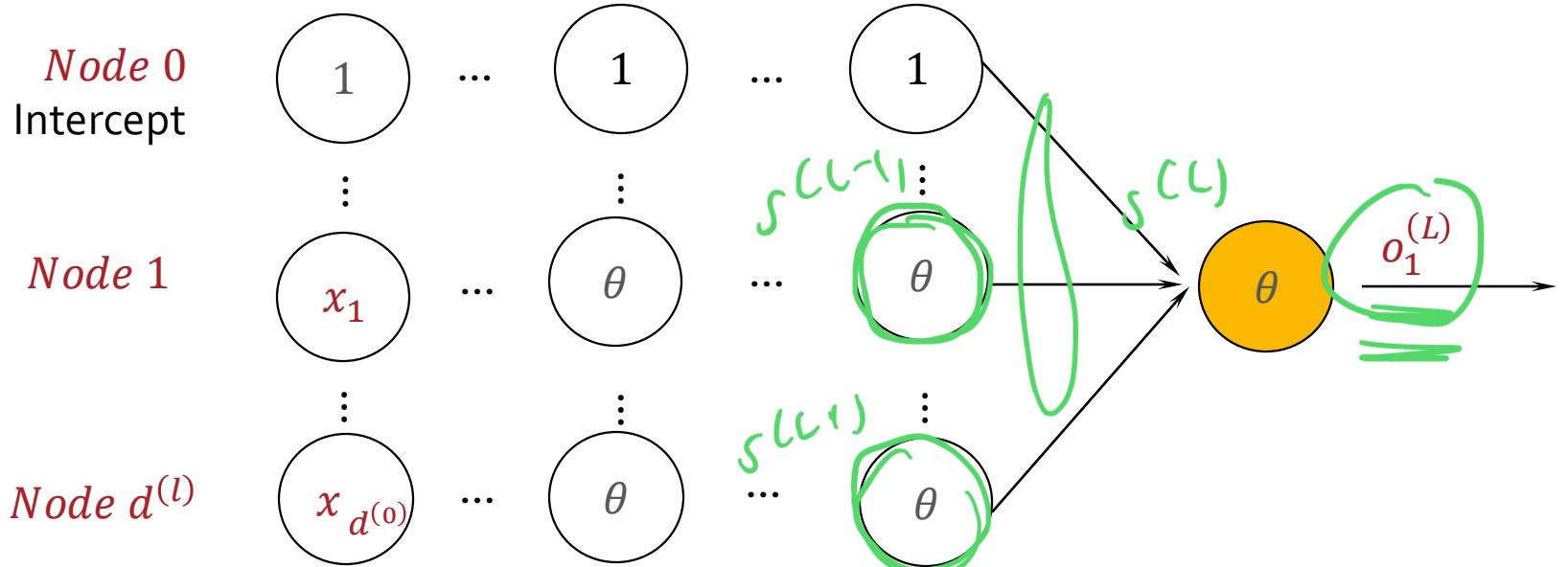
A neural network has L layers.

Each layer has $d^{(l)}$ nodes, and its own weight matrix $W^{(l)}$

The matrix $W^{(l)}$ connects layer $(l - 1)$ to layer l

Building intuition for backprop

Input layer 0 Layer l Layer L (Output)



The final loss is a function of the *output of the entire network*:

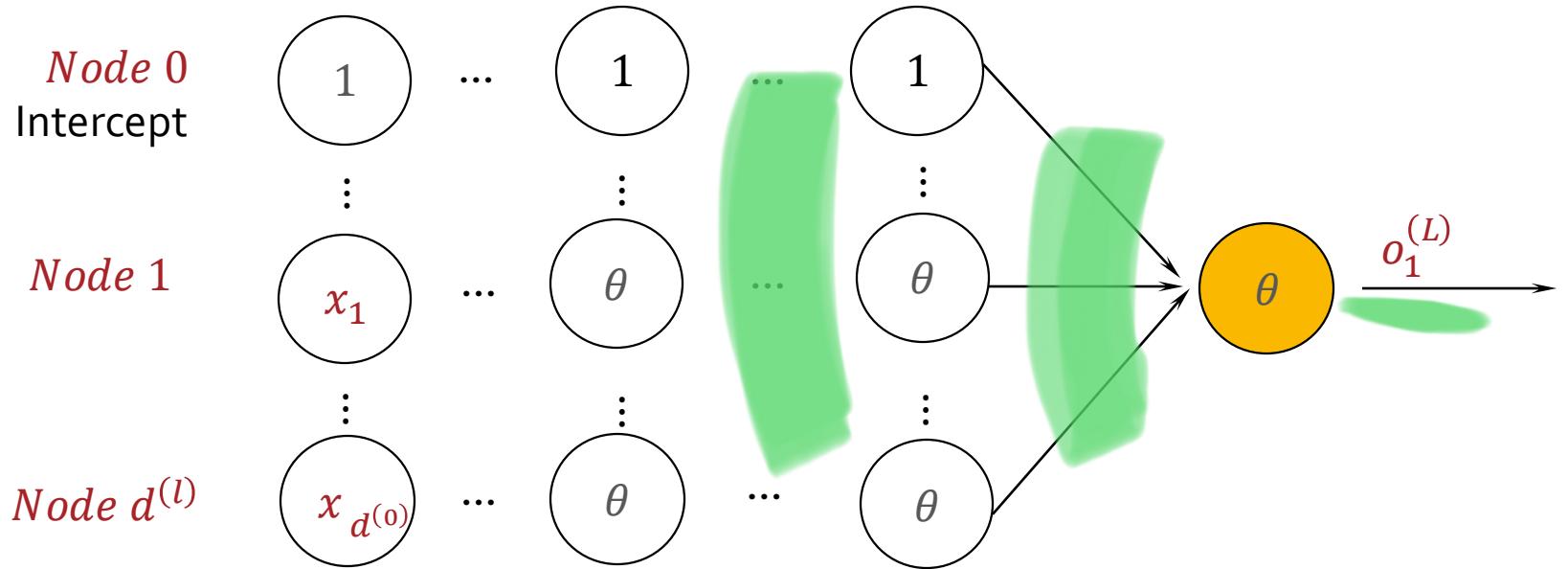
$$\ell^{(i)} = \left(\underline{o_1^{(L)}} - \underline{y^{(i)}} \right)^2$$

But the final output is a composite function of *all of the weights that came before...*

$$\begin{aligned} o^{(L)} &= \theta(s^{(L)}) \\ &= \theta(\overline{\theta}(W^{(L)}(\theta(s^{(L-1)})))) \\ &= \dots \end{aligned}$$

Building intuition for backprop

Input layer 0 Layer l Layer L (Output)



Idea: We can *apply the chain rule* to calculate the gradients at each layer.

$$\frac{\partial \ell^{(i)}}{\partial W^{(L)}} = \frac{\partial \ell^{(i)}}{\partial o^{(L)}} \left(\frac{\partial o^{(L)}}{\partial s^{(L)}} \right) \left(\frac{\partial s^{(L)}}{\partial W^{(L)}} \right)$$

$$\frac{\partial \ell^{(i)}}{\partial W^{(L-1)}} = \frac{\partial \ell^{(i)}}{\partial o^{(L)}} \left(\frac{\partial o^{(L)}}{\partial s^{(L)}} \right) \left(\frac{\partial s^{(L)}}{\partial o^{(L-1)}} \right) \left(\frac{\partial o^{(L-1)}}{\partial s^{(L-1)}} \right) \left(\frac{\partial s^{(L-1)}}{\partial W^{(L-1)}} \right)$$

Computing Gradients

$$\frac{\partial \ell^{(i)}}{\partial w^{(L-\delta)}} = \dots$$

$$\nabla_{W^{(l)}} \underline{\ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)} =$$

Computing Gradients

Recall: $W^{(l)}: [d^{(l)}, d^{(l-1)} + 1]$

$$\nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) =$$

$(l-1, l)$

$$\underline{W}^{(l)}: \left[\underline{d}^{(l)}, \underline{d}^{(l-1)} + 1 \right]$$

~~$d^{(l)}$~~

$$\begin{bmatrix} \downarrow & \downarrow & \dots & \downarrow \\ \frac{\partial \ell^{(i)}}{\partial w_{1,0}^{(l)}} & \frac{\partial \ell^{(i)}}{\partial w_{1,1}^{(l)}} & \dots & \frac{\partial \ell^{(i)}}{\partial w_{1,d^{(l-1)}}^{(l)}} \\ \frac{\partial \ell^{(i)}}{\partial w_{2,0}^{(l)}} & \frac{\partial \ell^{(i)}}{\partial w_{2,1}^{(l)}} & \dots & \frac{\partial \ell^{(i)}}{\partial w_{2,d^{(l-1)}}^{(l)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \ell^{(i)}}{\partial w_{d^{(l)},0}^{(l)}} & \frac{\partial \ell^{(i)}}{\partial w_{d^{(l)},1}^{(l)}} & \dots & \frac{\partial \ell^{(i)}}{\partial w_{d^{(l)},d^{(l-1)}}^{(l)}} \end{bmatrix}$$

Computing Partial Derivatives

Computing $\nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ reduces to computing

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}}$$

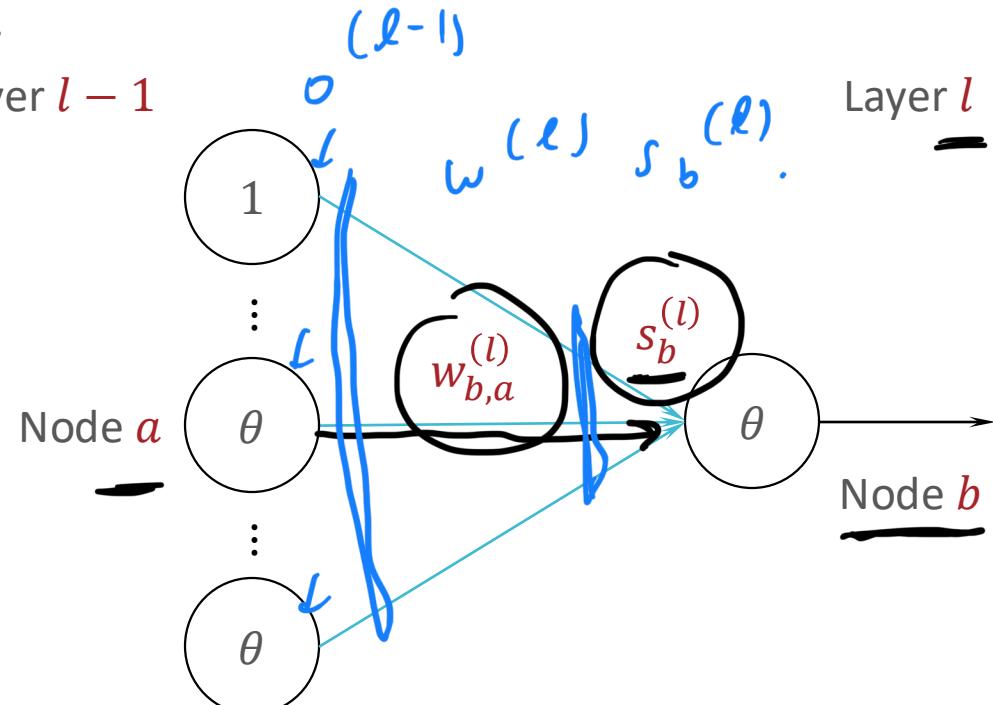


Computing Partial Derivatives

Computing $\nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ reduces to computing

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}}$$

Insight: $w_{b,a}^{(l)}$ only affects $\ell^{(i)}$ via $s_b^{(l)}$



Computing Partial Derivatives

Computing $\nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ reduces to computing

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}}$$

Insight: $w_{b,a}^{(l)}$ only affects $\ell^{(i)}$ via $s_b^{(l)}$

Chain rule: $\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}} = \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right)$

$$s_b^{(l)} = \sum_{c=0}^{\ell-1} w_{b,c}^{(l)} \cdot o_c^{(l)}$$

$$\frac{\partial}{\partial w_{b,a}^{(l)}} (\cdot) = o_a^{(l)}$$

Computing Partial Derivatives

Computing $\nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ reduces to computing

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}}$$

Insight: $w_{b,a}^{(l)}$ only affects $\ell^{(i)}$ via $s_b^{(l)}$

Chain rule: $\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}} = \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right)$

$$s_b^{(l)} = \sum_{a=0}^{d^{(l-1)}} w_{b,a}^{(l)} o_a^{(l-1)} \rightarrow \frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} = o_a^{(l-1)}$$

Compute outputs $\mathbf{o}^{(l)} \forall l \in \{0, \dots, L\}$ by forward propagation

Computing Partial Derivatives

Computing $\nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ reduces to computing

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}}$$

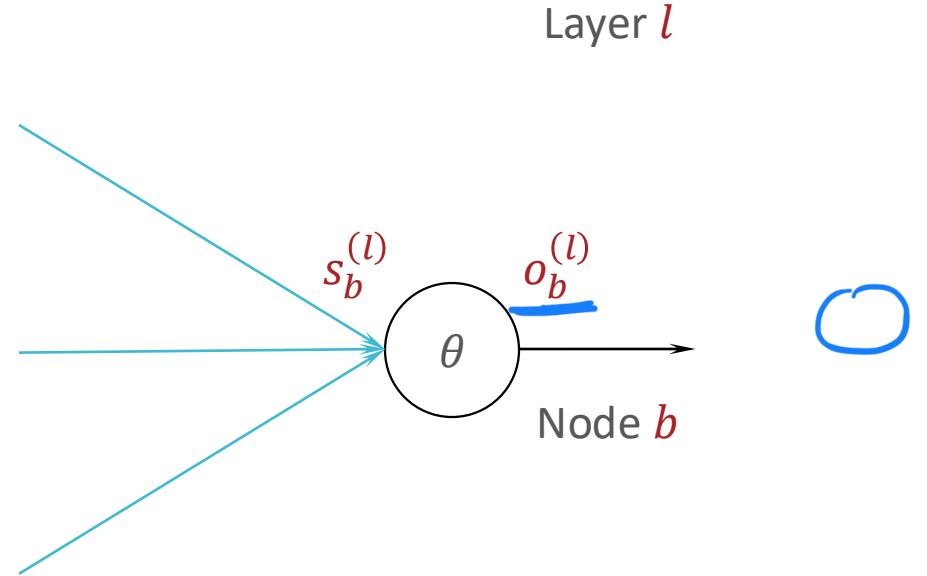
Insight: $w_{b,a}^{(l)}$ only affects $\ell^{(i)}$ via $s_b^{(l)}$

Chain rule:
$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}} = \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right)$$

$$\delta_b^{(l)} := \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}}$$

Computing Partial Derivatives

Insight: $s_b^{(l)}$ only affects $\ell^{(i)}$ via $o_b^{(l)}$



Computing Partial Derivatives

Insight: $s_b^{(l)}$ only affects $\ell^{(i)}$ via $o_b^{(l)}$

Chain rule: $\delta_b^{(l)} = \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right)$

$$o_b^{(l)} = \theta(s_b^{(l)}) \quad \theta = \tanh(z)$$

$$\begin{aligned} \frac{\partial}{\partial s_b^{(l)}} &: 1 - (\tanh(s_b^{(l)}))^2 \\ &= 1 - \underline{(o_b^{(l)})^2}. \end{aligned}$$

Computing Partial Derivatives

Insight: $s_b^{(l)}$ only affects $\ell^{(i)}$ via $o_b^{(l)}$

Chain rule: $\delta_b^{(l)} = \frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right)$

$$\begin{aligned} o_b^{(l)} &= \theta(s_b^{(l)}) \rightarrow \frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} = \frac{\partial \theta(s_b^{(l)})}{\partial s_b^{(l)}} \\ &= 1 - (\tanh(s_b^{(l)}))^2 \end{aligned}$$

when $\theta(\cdot) = \tanh(\cdot)$

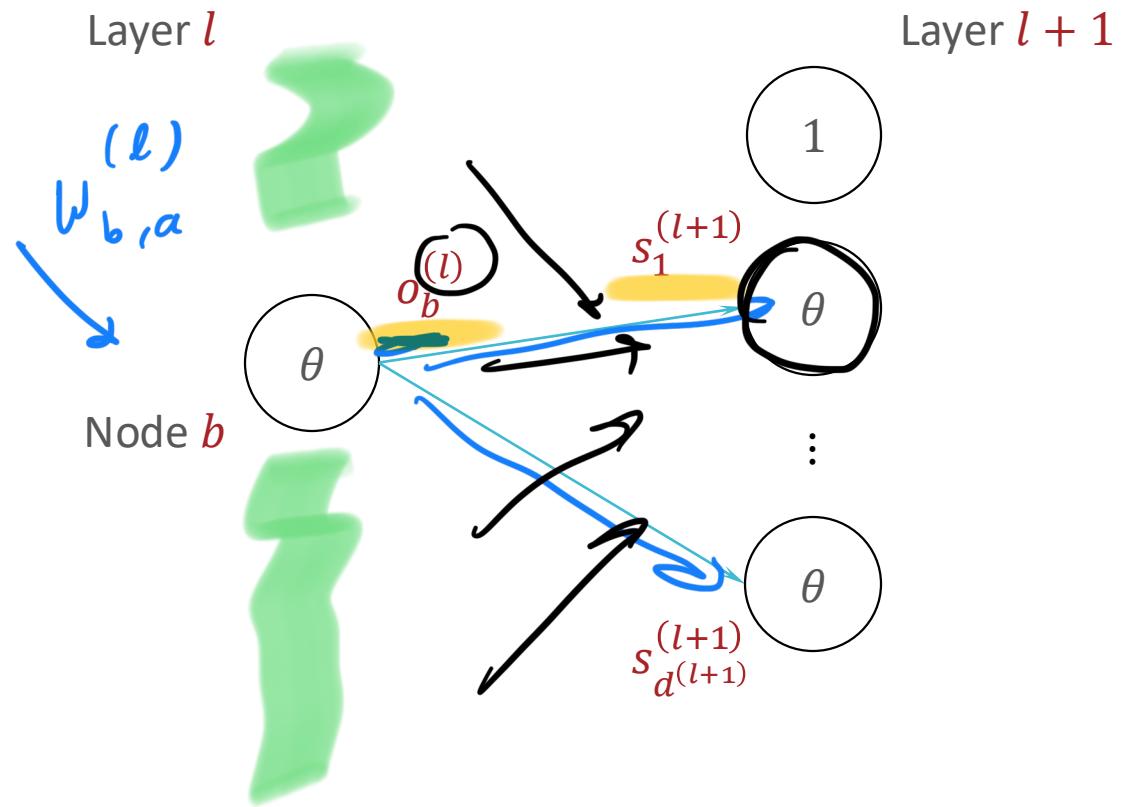
Computing Partial Derivatives

Recap:

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}} = \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right)$$
$$= \left(\frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \right) \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right) \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right)$$

Computing Partial Derivatives

Insight: $o_b^{(l)}$ affects $\ell^{(i)}$ via $s_1^{(l+1)}, \dots, s_{d^{(l+1)}}^{(l+1)}$



Computing Partial Derivatives

$\ell^{(i)}(s_1^{(l+1)}, \dots, s_d^{(l+1)})$

Insight: $o_b^{(l)}$ affects $\ell^{(i)}$ via $s_1^{(l+1)}, \dots, s_d^{(l+1)}$

Chain rule: $\frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \frac{\partial \ell^{(i)}}{\partial s_c^{(l+1)}} \left(\frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} \right)$

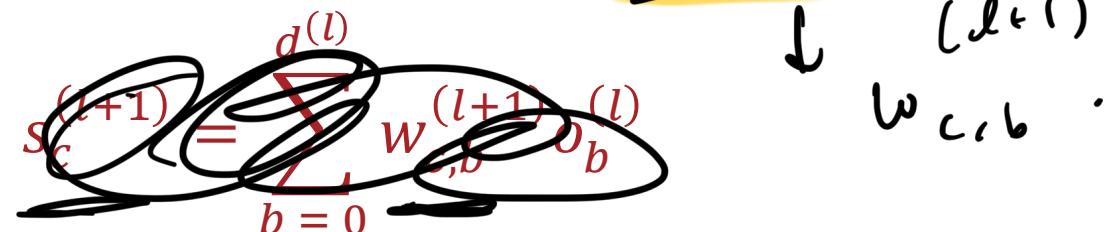
\downarrow

all nodes at layer $(l+1)$

Computing Partial Derivatives

Insight: $o_b^{(l)}$ affects $\ell^{(i)}$ via $s_1^{(l+1)}, \dots, s_{d^{(l+1)}}^{(l+1)}$

Chain rule: $\frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \frac{\partial \ell^{(i)}}{\partial s_c^{(l+1)}} \left(\frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} \right)$



$$s_c^{(l+1)} = \sum_{k=0}^{d^{(l)}} w_{c,k}^{(l+1)} o_k^{(l)}$$

$$\frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} = w_{c,b}^{(l+1)}$$

Computing Partial Derivatives

Insight: $o_b^{(l)}$ affects $\ell^{(i)}$ via $s_1^{(l+1)}, \dots, s_{d^{(l+1)}}^{(l+1)}$

Chain rule: $\frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \frac{\partial \ell^{(i)}}{\partial s_c^{(l+1)}} \left(\frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} \right)$

$$s_c^{(l+1)} = \sum_{b=0}^{d^{(l)}} w_{c,b}^{(l+1)} o_b^{(l)} \rightarrow \frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} = w_{c,b}^{(l+1)}$$

$$= \sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} \cdot w_{c,b}^{(l+1)}$$

$$\underline{\delta_b^{(l)}} := \overbrace{\frac{\partial \ell^{(i)}}{\partial s_b^{(l)}}}$$

Computing Partial Derivatives

Insight: $o_b^{(l)}$ affects $\ell^{(i)}$ via $s_1^{(l+1)}, \dots, s_{d^{(l+1)}}^{(l+1)}$

Chain rule: $\frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \frac{\partial \ell^{(i)}}{\partial s_c^{(l+1)}} \left(\frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} \right)$

$$s_c^{(l+1)} = \sum_{b=0}^{d^{(l)}} w_{c,b}^{(l+1)} o_b^{(l)} \rightarrow \frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} = w_{c,b}^{(l+1)}$$

$$= \sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} \left(w_{c,b}^{(l+1)} \right)$$

$$\delta_b^{(l)} := \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}}$$

Computing Partial Derivatives

Recap:

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}} = \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right)$$

$= \delta_b^{(l)}$

We just learned that we can write $\delta_b^{(l)}$ as a function of the terms $\delta_c^{(l+1)}$ at the next layer!

—

Computing Partial Derivatives

Pasting over what we calculated previously:

$$\underline{\delta_b^{(l)}} = \frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right)$$
$$= \left(\sum_{c=1}^{d^{(l+1)}} \underline{\delta_c^{(l+1)}} \left(\underline{w_{c,b}^{(l+1)}} \right) \right) \left(1 - \underline{\left(o_b^{(l)} \right)^2} \right)$$

Computing Partial Derivatives

$$\begin{aligned}
 \underline{\delta_b^{(l)}} &= \frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right) \quad \frac{\partial \ell^{(i)}}{\partial s_b^{(l)}} \\
 &= \left(\sum_{c=1}^{d^{(l+1)}} \underline{\delta_c^{(l+1)} (w_{c,b}^{(l+1)})} \right) \left(1 - (o_b^{(l)})^2 \right) \quad \text{↳ no de b-} \\
 \underline{\delta^{(l)}} &:= \nabla_{\underline{s^{(l)}}} \ell^{(i)} (W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}) \quad [\underline{\delta^{(l)}}, 1]
 \end{aligned}$$

How would you write the vector $\underline{\delta^{(l)}}$ as a function of the matrices $\underline{\delta^{(l+1)}}$, $\underline{W^{(l+1)}}$ and $\underline{o^{(l)}}$?

Hint: $\underline{s^{(l)}}: [d^{(l)}, 1]$, $\underline{o^{(l)}}: [d^{(l)} + 1, 1]$, $\underline{W^{(l)}}: [d^{(l)}, d^{(l-1)} + 1]$

Computing Partial Derivatives

$$\delta_b^{(l)} = \frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right)$$

$$= \left(\sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} (w_{c,b}^{(l+1)}) \right) \left(1 - (o_b^{(l)})^2 \right)$$

$$\underline{\delta^{(l)}} = \underline{W^{(l+1)^T} \delta^{(l+1)}} \odot (1 - \underline{o^{(l)}} \odot \underline{o^{(l)}})$$

where \odot is the element-wise product operation

Sanity check:

$$\delta^{(l+1)} \in \mathbb{R}^{d^{(l+1)} \times 1}$$

$$W^{(l+1)^T} \delta^{(l+1)} \in \mathbb{R}^{(d^{(l+1)} \times 1) \times 1}$$

$$\begin{bmatrix} o_1^{(l+1)} \\ o_2^{(l+1)} \\ \vdots \\ o_d^{(l+1)} \end{bmatrix} \odot \begin{bmatrix} o_1^{(l+1)} \\ o_2^{(l+1)} \\ \vdots \\ o_d^{(l+1)} \end{bmatrix} = \begin{bmatrix} (o_1^{(l+1)})^2 \\ (o_2^{(l+1)})^2 \\ \vdots \\ (o_d^{(l+1)})^2 \end{bmatrix}$$

Computing Partial Derivatives

Putting it all together:

$$\frac{\partial \ell^{(i)}}{\cancel{\partial w_{b,a}^{(l)}}} = \delta_b^{(l)} \left(\frac{\partial s_b^{(l)}}{\cancel{\partial w_{b,a}^{(l)}}} \right) = \underline{\delta_b^{(l)} \left(o_a^{(l-1)} \right)}$$
$$\frac{\partial \ell^{(i)}}{\partial \mathbf{W}^{(l)}} =$$

Computing Partial Derivatives

Putting it all together:

$$\frac{\partial \ell^{(i)}}{\partial w_{b,a}^{(l)}} = \delta_b^{(l)} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right) = \delta_b^{(l)} \left(\underline{\underline{o}_a^{(l-1)}} \right)$$

$$a \mapsto b$$

$$\frac{\partial \ell^{(i)}}{\partial W^{(l)}} = \underline{\delta^{(l)} o^{(l-1)^T}}$$

Sanity check:

$$\delta^{(l)} \in \mathbb{R}^{d^{(l)} \times 1} \text{ so}$$

$$\delta^{(l)} o^{(l-1)^T} \in \mathbb{R}^{d^{(l)} \times (d^{(l-1)} + 1)}, \text{ the same size as } W^{(l)}!$$

Computing Partial Derivatives

- Can recursively compute $\delta^{(l)}$ using $\delta^{(l+1)}$; need to compute the base case: $\underline{\delta^{(L)}}$

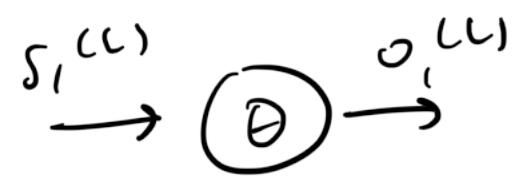
$$\frac{\partial \ell^{(i)}}{\partial s^{(l)}} \quad \frac{\partial \ell^{(i)}}{\partial s^{(l+1)}}.$$

$$\delta^{(l)} = W^{(l+1)^T} \delta^{(l+1)} \odot (1 - \underline{o^{(l)}} \odot o^{(l)})$$

- Assume the output layer is a single node and the error function is the squared error: $\delta^{(L)} = \delta_1^{(L)}$, $\underline{o^{(L)}} = \underline{o_1^{(L)}}$

$$\ell^{(i)}(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}) = (\underline{o_1^{(L)}} - y^{(i)})^2, \text{ and } \theta(\underline{z}) = z$$

Computing Partial Derivatives



- Assume the output layer is a single node and the error function is the squared error: $\delta^{(L)} = \delta_1^{(L)}$, $\mathbf{o}^{(L)} = o_1^{(L)}$

and $\theta(z) = z$, calculate $\delta_1^{(L)} = \frac{\partial \ell^{(i)}}{\partial s_1^{(L)}}$:

$$\ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) = \underbrace{\left(o_1^{(L)} - y^{(i)} \right)^2}_{\text{U}} \quad o_1^{(L)} = s_2^{(L)}$$

$$\frac{\partial \ell^{(i)}}{\partial s_1^{(L)}} = \underbrace{\frac{\partial \ell^{(i)}}{\partial o_2^{(L)}}}_{\text{U}} \cdot \underbrace{\frac{\partial o_2^{(L)}}{\partial s_2^{(L)}}}_{\text{U}}$$

$$\frac{\partial}{\partial o_2^{(L)}} \left(o_2^{(L)} - y^{(i)} \right)^2 = 2(o_2^{(L)} - y^{(i)}). = \delta_1^{(L)}.$$

Computing Partial Derivatives

- Assume the output layer is a single node and the error function is the squared error: $\delta^{(L)} = \delta_1^{(L)}$, $\mathbf{o}^{(L)} = o_1^{(L)}$

$$\ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) = \left(o_1^{(L)} - y^{(i)} \right)^2$$

$$\delta_1^{(L)} = \frac{\partial \ell^{(i)}}{\partial s_1^{(L)}} = \frac{\partial}{\partial s_1^{(L)}} \left(o_1^{(L)} - y^{(i)} \right)^2$$

$$= 2 \left(o_1^{(L)} - y^{(i)} \right) \frac{\partial o_1^{(L)}}{\partial s_1^{(L)}} = 2 \left(o_1^{(L)} - y^{(i)} \right)$$

when $\theta(z) = z$

Computing Partial Derivatives

We can compute $\underline{\delta^{(L)}}$!!! $\frac{\partial \ell^{(i)}}{\partial s^{(l)}}$

We can use these to compute $\underline{\delta^{(L-1)}}, \underline{\delta^{(L-2)}},$ all the way back to $\underline{\delta^{(1)}},$ using our equation from before....

$$\underline{\delta^{(l)}} = W^{(l+1)^T} \underline{\delta^{(l+1)}} \odot (1 - \underline{o^{(l)}} \odot \underline{o^{(l)}})$$

And we know how to use the $\underline{\delta^{(l)}}$ to compute the gradient:

$$\frac{\partial \ell^{(i)}}{\partial \underline{W^{(l)}}} = \underline{\delta^{(l)}} \underline{o^{(l-1)^T}}$$

Back-propagation

- Input: $\underline{W^{(1)}, \dots, W^{(L)}}$ and $(\underline{x^{(i)}}, \underline{y^{(i)}})$
- Run forward propagation with $\underline{x^{(i)}}$ to get $\underline{o^{(1)}, \dots, o^{(L)}}$
- (Optional) Compute $\ell^{(i)} = (\underline{o^{(L)} - y^{(i)}})^2$
- Initialize: $\underline{\delta^{(L)} = 2(o_1^{(L)} - y^{(i)})}$ $\frac{\partial \ell}{\sigma^{(L)}}$
- For $\underline{l = L - 1, \dots, 1}$
 - Compute $\underline{\delta^{(l)} = W^{(l+1)T} \delta^{(l+1)} \odot (1 - o^{(l)} \odot o^{(l)})}$
 - Compute $\underline{G^{(l)} = \delta^{(l)} o^{(l-1)T}}$
- Output: $\underline{G^{(1)}, \dots, G^{(L)}}$, the gradients of $\ell^{(i)}$ w.r.t $\underline{W^{(1)}, \dots, W^{(L)}}$

Key take-aways

- A weight affects the prediction of the network (and therefore the error) through downstream signals/outputs
 - Use the chain rule!
- Any weight going into the same node will affect the prediction through the same downstream path
 - Compute derivatives starting from the last layer and move “backwards”
 - Derive a recursive definition for the relevant partial derivatives
 - Automatic differentiation: store intermediate values and reuse for efficiency (dynamic programming)

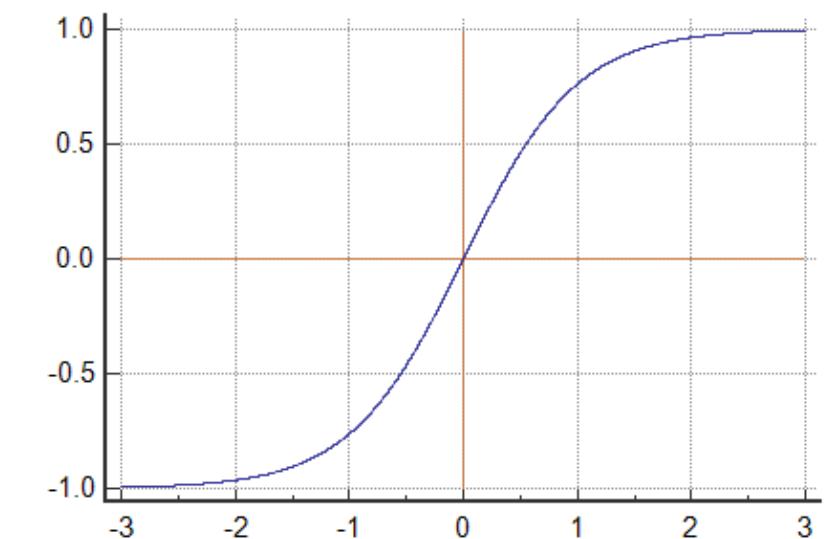
Aside: Vanishing Gradients

Insight: $s_b^{(l)}$ only affects $\ell^{(i)}$ via $o_b^{(l)}$

Chain rule: $\delta_b^{(l)} = \frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right)$

$$o_b^{(l)} = \theta(s_b^{(l)}) \rightarrow \frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} = \frac{\partial \theta(s_b^{(l)})}{\partial s_b^{(l)}} = 1 - (\tanh(s_b^{(l)}))^2 \leq 1$$

when $\theta(\cdot) = \tanh(\cdot)$

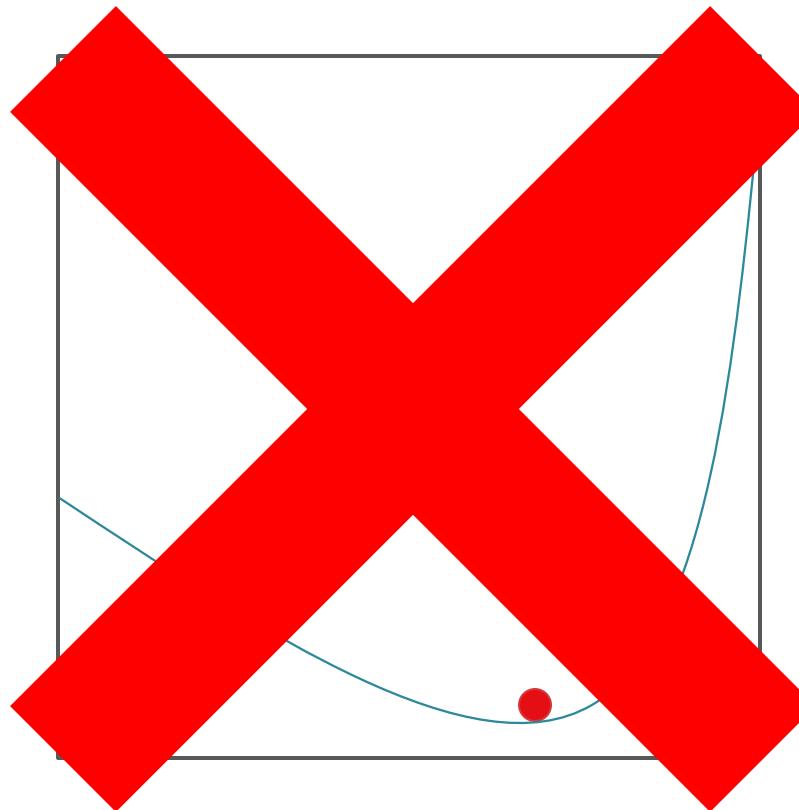


Recall: Other Activation Functions

| | | |
|--|--|---|
| Logistic, sigmoid, or soft step | | $\sigma(x) = \frac{1}{1 + e^{-x}}$ |
| Hyperbolic tangent (\tanh) | | $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ |
| Rectified linear unit (ReLU) ^[7] | | $\begin{cases} 0 & \text{if } x \le 0 \\ x & \text{if } x > 0 \end{cases} = \max\{0, x\} = x \mathbf{1}_{x>0}$ |
| Gaussian Error Linear Unit (GELU) ^[4] | | $\frac{1}{2}x \left(1 + \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) \right) = x\Phi(x)$ |
| Softplus ^[8] | | $\ln(1 + e^x)$ |
| Exponential linear unit (ELU) ^[9] | | $\begin{cases} \alpha(e^x - 1) & \text{if } x \le 0 \\ x & \text{if } x > 0 \end{cases}$ with parameter α |
| Leaky rectified linear unit (Leaky ReLU) ^[11] | | $\begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \ge 0 \end{cases}$ |
| Parametric rectified linear unit (PReLU) ^[12] | | $\begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \ge 0 \end{cases}$ with parameter α |

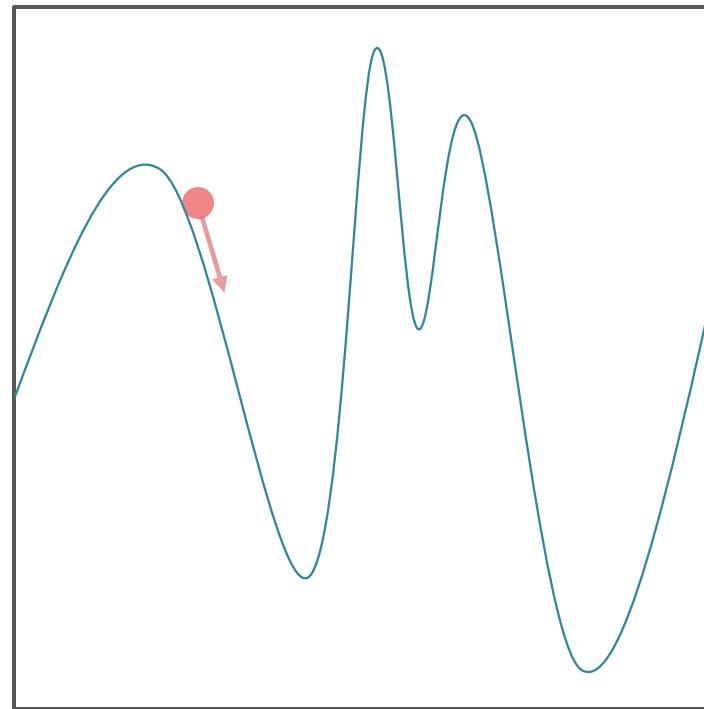
Recall: Gradient Descent

- Iterative method for minimizing functions
- Requires the gradient to exist everywhere



Non-convexity

- Gradient descent is not guaranteed to find a global minimum on non-convex surfaces



Stochastic Gradient Descent for Learning

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta^{(0)}$
- Initialize all weights $W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}$ to small, random numbers and set $t = 0$
- While TERMINATION CRITERION is not satisfied
 - For $i \in \text{shuffle}(\{1, \dots, N\})$
 - For $l = 1, \dots, L$
 - Compute $G^{(l)} = \nabla_{W^{(l)}} \ell^{(i)}(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)})$
 - Update $W^{(l)}$: $W_{(t+1)}^{(l)} = W_{(t)}^{(l)} - \eta_0 G^{(l)}$
 - Increment t : $t = t + 1$
 - Output: $W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}$

Mini-batch Stochastic Gradient Descent for Learning

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{MB}^{(0)}, B$
- 1. Initialize all weights $W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}$ to small, random numbers and set $t = 0$
- 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample B data points from $\mathcal{D}, \{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$
 - b. Compute the gradient w.r.t. the sampled batch,
$$G^{(l)} = \frac{1}{B} \sum_{b=1}^B \nabla_{W^{(l)}} \ell^{(b)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) \forall l$$
 - c. Update $W^{(l)}$: $W_{t+1}^{(l)} \leftarrow W_t^{(l)} - \eta_{MB}^{(0)} G^{(l)} \forall l$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $W_t^{(1)}, \dots, W_t^{(L)}$