10-301/601: Introduction to Machine Learning Lecture 13 – Backpropagation

Henry Chai 2/24/25

Front Matter

- Announcements
 - Exam 1 viewings this week, Tuesday Thursday
 - See Piazza for complete details
 - Homework 4 released 2/17, due 2/26 at 11:59 PM

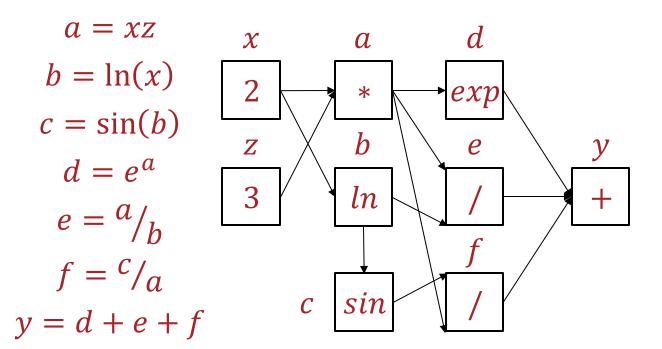
Given

Recall: Automatic Differentiation (reverse mode)

$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at x = 2, z = 3?

 First define some intermediate quantities, draw the computation graph and run the "forward" computation



Recall: **Automatic** Differentiation (reverse mode) Given

$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at x = 2, z = 3?

•
$$g_y = \frac{\partial y}{\partial y} = 1$$

• $g_d = g_e = g_f = 1$

 Then compute partial derivatives, starting from y and working back • $g_c = \frac{\partial y}{\partial c} = \frac{\partial y}{\partial f} \frac{\partial f}{\partial c} = g_f \left(\frac{1}{a}\right)$

•
$$g_c = \frac{\partial y}{\partial c} = \frac{\partial y}{\partial f} \frac{\partial f}{\partial c} = g_f \left(\frac{1}{a}\right)$$

$$g_{b} = \frac{\partial y}{\partial b} = \frac{\partial y}{\partial e} \frac{\partial e}{\partial b} + \frac{\partial y}{\partial c} \frac{\partial c}{\partial b}$$

$$= g_{e} \left(-\frac{a}{b^{2}} \right) + g_{c} (\cos(b))$$

$$g_{a} = \frac{\partial y}{\partial a} = \frac{\partial y}{\partial f} \frac{\partial f}{\partial a} + \frac{\partial y}{\partial e} \frac{\partial e}{\partial a} + \frac{\partial y}{\partial d} \frac{\partial d}{\partial a}$$

$$= g_{f} \left(\frac{-c}{a^{2}} \right) + g_{e} \left(\frac{1}{b} \right) + g_{d} (e^{a})$$

$$f$$

$$g_{x} = \frac{\partial y}{\partial x} = \frac{\partial y}{\partial b} \frac{\partial b}{\partial x} + \frac{\partial y}{\partial a} \frac{\partial a}{\partial x} = g_{b} \left(\frac{1}{x} \right) + g_{a} (z)$$

$$g_{z} = \frac{\partial y}{\partial z} = \frac{\partial y}{\partial a} \frac{\partial a}{\partial z} = g_{a} (x)$$

Computation Graph 10-301/601 Conventions

- The diagram represents an algorithm
- Nodes are rectangles with one node per intermediate variable in the algorithm
- Each node is labeled with the function that it computes (inside the box) and the variable name (outside the box)
- Edges are directed and do not have labels
- For neural networks:
 - Each weight, feature value, label and bias term appears as a node
 - We can include the loss function

Neural Network Diagram Conventions

- The diagram represents a *neural network*
- Nodes are circles with one node per hidden unit
- Each node is labeled with the variable corresponding to the hidden unit
- Edges are directed and each edge is labeled with its weight
- The diagram typically does not include any nodes related to the loss computation

Matrix Calculus

	Types of Derivatives	scalar	vector	matrix
Denominator	scalar	$\frac{\partial y}{\partial x}$	$\frac{\partial \mathbf{y}}{\partial x}$	$\frac{\partial \mathbf{Y}}{\partial x}$
	vector	$\frac{\partial y}{\partial \mathbf{x}}$	$rac{\partial \mathbf{y}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$
	matrix	$\frac{\partial y}{\partial \mathbf{X}}$	$rac{\partial \mathbf{y}}{\partial \mathbf{X}}$	$rac{\partial \mathbf{Y}}{\partial \mathbf{X}}$

Donominator

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.

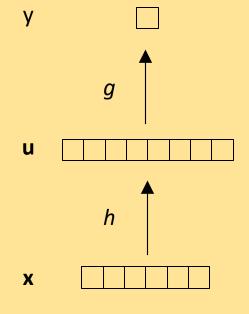
Types of Derivatives	scalar		
scalar	$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x}\right]$		
vector	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$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

	Types of Derivatives	scalar	vector
Matrix Calculus:	scalar	$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x}\right]$	$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} & \frac{\partial y_2}{\partial x} & \cdots & \frac{\partial y_N}{\partial x} \end{bmatrix}$
Denominator Layout	vector	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_N}{\partial x_2} \\ \vdots & & & & \\ \frac{\partial y_1}{\partial x_P} & \frac{\partial y_2}{\partial x_P} & \cdots & \frac{\partial y_N}{\partial x_P} \end{bmatrix}$

Matrix Calculus

Poll Question 1:

Suppose y = g(u) and u = h(x)



Which of the following is the correct definition of the chain rule?

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_N}{\partial x_2} \end{bmatrix}$$

Answers: $\frac{\partial y}{\partial \mathbf{x}} = \dots$

A.
$$\frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$$

$$\mathsf{B.} \ \frac{\partial \boldsymbol{y}}{\partial \mathbf{u}}^T \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$$

$$\mathsf{C.} \ \frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}^T$$

D.
$$\frac{\partial y}{\partial \mathbf{u}}^T \frac{\partial \mathbf{u}}{\partial \mathbf{x}}^T$$

(TOXIC) E.
$$(\frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}})$$

F. None of the above

Gradient Descent for Neural Network Training

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^{N}, \eta^{(0)}$
- 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 l = 1, ..., L
 - Compute $G^{(l)} = \nabla_{W^{(l)}} \ell_{\mathcal{D}} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ (???)
 - Update $W^{(l)}$: $W^{(l)}_{(t+1)} = W^{(l)}_{(t)} \eta_0 G^{(l)}$
 - Increment t: t = t + 1
- Output: $W_{(t)}^{(1)}, ..., W_{(t)}^{(L)}$

Computing Gradients

$$\begin{split} \ell_{\mathcal{D}}\left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}\right) &= \sum_{n=1}^{N} \ell^{(n)}\left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}\right) \\ \nabla_{W^{(l)}}\ell_{\mathcal{D}}\left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}\right) &= \begin{bmatrix} \frac{\partial \ell_{\mathcal{D}}}{\partial w_{1,0}^{(l)}} & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{1,1}^{(l)}} & \dots & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{1,d}^{(l-1)}} \\ \frac{\partial \ell_{\mathcal{D}}}{\partial w_{2,0}^{(l)}} & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{2,1}^{(l)}} & \dots & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{2,d^{(l-1)}}^{(l)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \ell_{\mathcal{D}}}{\partial w_{d^{(l)},0}^{(l)}} & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{d^{(l)},1}^{(l)}} & \dots & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{d^{(l)},d^{(l-1)}}} \end{bmatrix} \\ \frac{\partial \ell_{\mathcal{D}}}{\partial u_{d^{(l)},0}^{(l)}} &= \sum_{i=1}^{N} \frac{\partial \ell^{(n)}\left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}\right)}{\partial w_{1,1}^{(l)}} \end{split}$$

$$\frac{\partial v_D}{\partial w_{h,a}^{(l)}} = \sum_{n=1}^{\infty} \frac{(v_D + v_D)^{(l)}}{\partial w_{h,a}^{(l)}}$$

Computing Gradients: Intuition

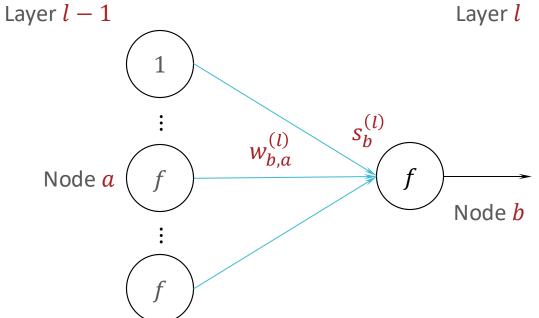
- A weight affects the prediction of the network (and therefore the error) through downstream signals/outputs
 - Use the chain rule!
- Any weight going into the same node will affect the prediction through the same downstream path
 - Compute derivatives starting from the last layer and move "backwards"
 - Store computed derivatives and reuse for efficiency (automatic differentiation)

2/24/25 **13**

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

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

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



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

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

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

Chain rule:
$$\frac{\partial \ell^{(n)}}{\partial w_{b,a}^{(l)}} = \frac{\partial \ell^{(n)}}{\partial s_b^{(l)}} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right)$$
$$s_b^{(l)} = \sum_{a=0}^{d^{(l-1)}} w_{b,a}^{(l)} o_a^{(l-1)} \rightarrow \frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} = o_a^{(l-1)}$$

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

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

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

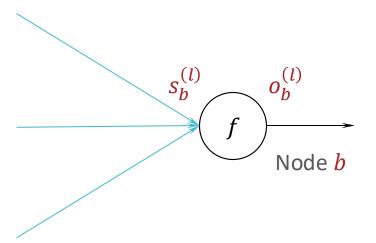
Insight: $w_{b,a}^{(l)}$ only affects $\ell^{(n)}$ via $s_b^{(l)}$

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

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

Layer *l*





Insight: $s_h^{(l)}$ only affects $\ell^{(n)}$ via $o_h^{(l)}$

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

$$o_b^{(l)} = f\left(s_b^{(l)}\right) \to \frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} = \frac{\partial f\left(s_b^{(l)}\right)}{\partial s_b^{(l)}}$$

"Vanishing gradients": repeatedly multiplying by something < 1 causes = $\frac{\exp\left(-s_b^{(l)}\right)}{\left(1 + \exp\left(-s_b^{(l)}\right)\right)^2} \le 1$ the gradient to approach zero!

$$= \frac{\exp\left(-s_b^{(l)}\right)}{\left(1 + \exp\left(-s_b^{(l)}\right)\right)^2} \le 1$$

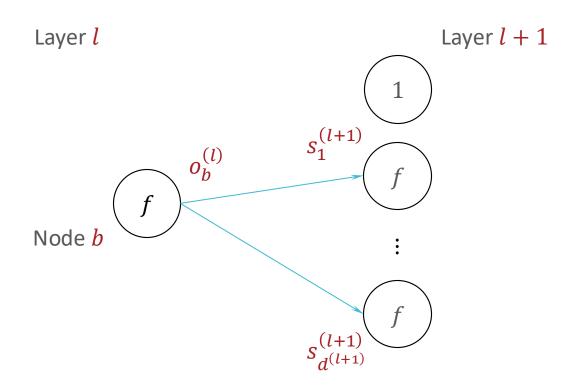
$$= \frac{1}{\left(1 + \exp\left(-s_b^{(l)}\right)\right)^2} \le 1$$

when
$$f(z) = \frac{1}{1 + \exp(-z)}$$

Recall: Other Activation Functions

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

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



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

Chain rule:
$$\frac{\partial \ell^{(n)}}{\partial o_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \frac{\partial \ell^{(n)}}{\partial s_c^{(l+1)}} \left(\frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} \right)$$
$$s_c^{(l+1)} = \sum_{b=0}^{d^{(l)}} w_{c,b}^{(l+1)} o_b^{(l)} \rightarrow \frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}} = w_{c,b}^{(l+1)}$$
$$\frac{\partial \ell^{(n)}}{\partial o_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} \left(w_{c,b}^{(l+1)} \right)$$

$$\begin{split} \boldsymbol{\delta}_b^{(l)} &= \frac{\partial \ell^{(n)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right) \\ &= \left(\sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} \left(w_{c,b}^{(l+1)} \right) \right) \left(1 - \left(o_b^{(l)} \right)^2 \right) \\ \boldsymbol{\delta}^{(l)} &\coloneqq \nabla_{\boldsymbol{s}^{(l)}} \ell^{(n)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) \\ &= W^{(l+1)^T} \boldsymbol{\delta}^{(l+1)} \odot \left(1 - \boldsymbol{o}^{(l)} \odot \boldsymbol{o}^{(l)} \right) \end{split}$$

where \odot is the element-wise product operation

Sanity check:
$$W^{(l+1)} \in \mathbb{R}^{d^{(l+1)} \times (d^{(l)}+1)}$$
 and

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

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

$$\frac{\partial \ell^{(n)}}{\partial w_{b,a}^{(l)}} = \delta_b^{(l)} \left(\frac{\partial s_b^{(l)}}{\partial w_{b,a}^{(l)}} \right) = \delta_b^{(l)} \left(o_a^{(l-1)} \right)$$

Computing Gradients

$$\nabla_{W^{(l)}}\ell^{(n)} = \boldsymbol{\delta}^{(l)}\boldsymbol{o}^{(l-1)^T}$$

Sanity check:
$$o^{(l-1)} \in \mathbb{R}^{\left(d^{(l-1)}+1\right) \times 1}$$
 and

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

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

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

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

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

$$\delta_{1}^{(L)} = \frac{\partial e\left(o_{1}^{(L)}, y^{(n)}\right)}{\partial s_{1}^{(L)}} = \frac{\partial}{\partial s_{1}^{(L)}} \left(o_{1}^{(L)} - y^{(n)}\right)^{2}$$
$$= 2\left(o_{1}^{(L)} - y^{(n)}\right) \frac{\partial o_{1}^{(L)}}{\partial s_{1}^{(L)}}$$

Backpropagation

- Input: $W^{(1)}, ..., W^{(L)}$ and $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$
- Initialize: $\ell_{\mathcal{D}}=0$ and $G^{(l)}=0\odot W^{(l)}$ \forall $l=1,\ldots,L$
- For n = 1, ..., N
 - Run forward propagation with $\boldsymbol{x}^{(n)}$ to get $\boldsymbol{o}^{(1)}$, ..., $\boldsymbol{o}^{(L)}$
 - (Optional) Increment $\ell_{\mathcal{D}}$: $\ell_{\mathcal{D}} = \ell_{\mathcal{D}} + \left(o^{(L)} y^{(n)}\right)^2$
 - Initialize: $\delta^{(L)} = 2 \left(o_1^{(L)} y^{(n)} \right) \frac{\partial o_1^{(L)}}{\partial s_1^{(L)}}$
 - For l = L 1, ..., 1
 - Compute $\boldsymbol{\delta}^{(l)} = W^{(l+1)^T} \boldsymbol{\delta}^{(l+1)} \odot (1 \boldsymbol{o}^{(l)} \odot \boldsymbol{o}^{(l)})$
 - Increment $G^{(l)}: G^{(l)} = G^{(l)} + \delta^{(l)} o^{(l-1)^T}$
- Output: $G^{(1)}$, ..., $G^{(L)}$, the gradients of $\ell_{\mathcal{D}}$ w.r.t $W^{(1)}$, ..., $W^{(L)}$

Backpropagation Learning Objectives

You should be able to...

- Differentiate between a neural network diagram and a computation graph
- Construct a computation graph for a function as specified by an algorithm
- Carry out the backpropagation on an arbitrary computation graph
- Construct a computation graph for a neural network, identifying all the given and intermediate quantities that are relevant
- Instantiate the backpropagation algorithm for a neural network
- Instantiate an optimization method (e.g. SGD) and a regularizer (e.g. L2)
 when the parameters of a model are comprised of several matrices
 corresponding to different layers of a neural network
- Apply the empirical risk minimization framework to learn a neural network
- Use the finite difference method to evaluate the gradient of a function
- Identify when the gradient of a function can be computed at all and when it can be computed efficiently
- Employ basic matrix calculus to compute vector/matrix/tensor derivatives.