

10-301/601: Introduction to Machine Learning

Lecture 14 – Backpropagation

Henry Chai

5/22/25

Front Matter

- Announcements:
 - HW3 released on 5/20, due 5/23 (tomorrow) at 11:59 PM
 - Quiz 2 on 5/23 (tomorrow) at 11:00 AM in BH A36 (here)
 - Study guide solutions partially released 5/21 (yesterday)
 - The remaining solutions to be released after recitation on 5/22 (today!)
 - Midterm on 5/30 at 9:30 AM in BH A36
 - Lectures 1 – 14 are in-scope; **next week's lectures will not be tested on the midterm**

Midterm Logistics

- Time and place:
 - Friday, 5/30 from 9:30 AM to 12:00 PM in BH A36 (here)
- Closed book/notes
 - 1-page cheatsheet allowed, both back and front; can be typeset or handwritten
 - No electronic devices allowed, **including calculators**

Midterm Coverage

- Lectures: 1 – 14 (through this week's lectures)
 - Foundations: probability, linear algebra, calculus
 - Important concepts: inductive bias, overfitting, model selection/hyperparameter optimization, regularization
 - Models: decision trees, kNN, Perceptron, linear regression, logistic regression, neural networks
 - Methods: (stochastic) gradient descent, closed-form optimization, backpropagation, MLE/MAP

Midterm Preparation

- Review midterm practice problems, to be posted on 5/26 to the course website (under [Schedule](#))
- Attend the exam review recitation on 5/29
- Review the homeworks and study guides
- Consider whether you understand the “Key Takeaways” for each lecture / section
- Write your cheat sheet

Recall: Loss Functions for Neural Networks

- Multi-class classification - cross-entropy loss
 - Express the label as a one-hot or one-of- C vector e.g.,

$$y = [0 \quad 0 \quad 1 \quad 0 \quad \dots \quad 0]$$

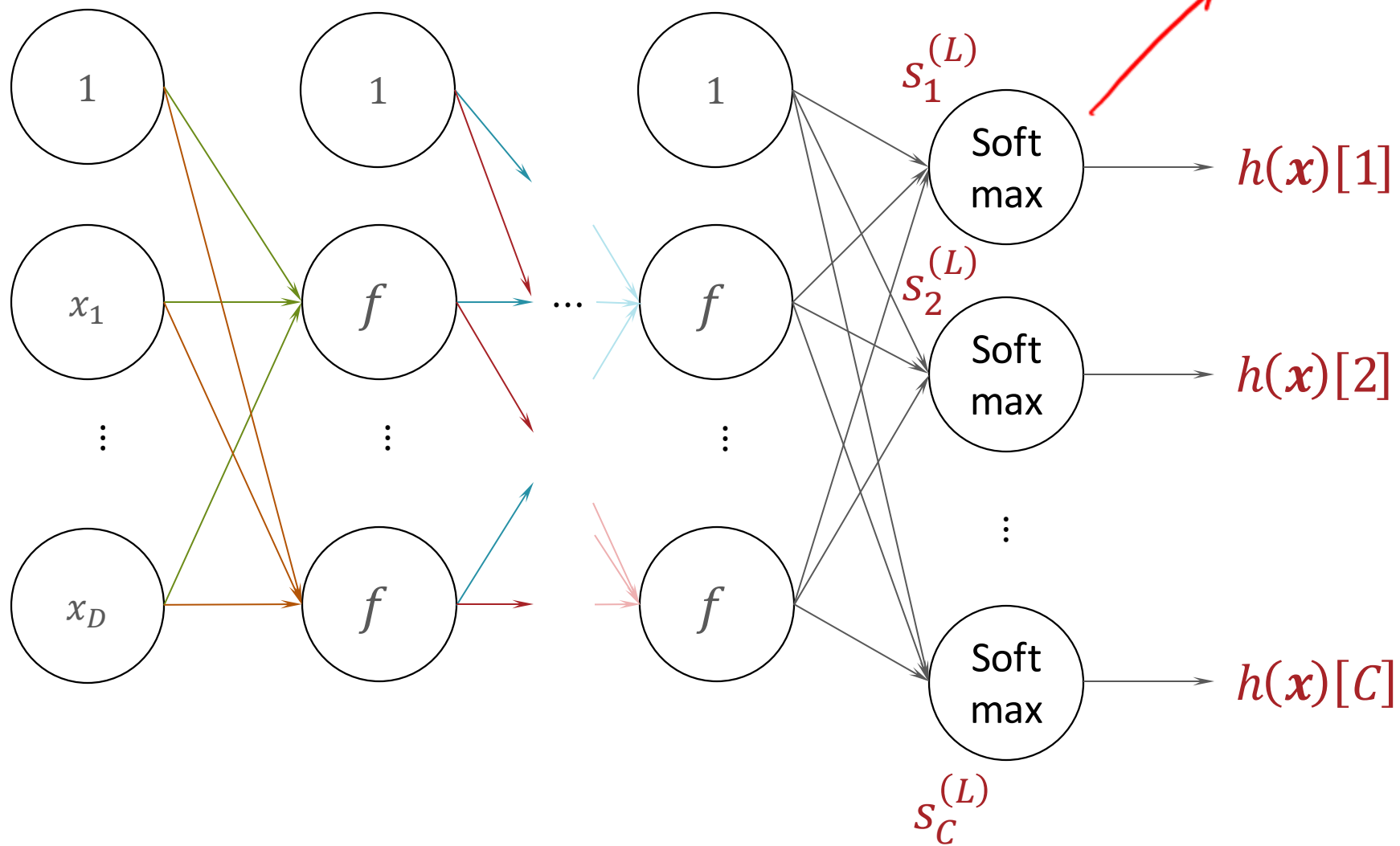
- Assume the neural network output is also a vector of length C

$$P(y[k] = 1 | \mathbf{x}, W^{(1)}, \dots, W^{(L)}) = h_{W^{(1)}, \dots, W^{(L)}}(\mathbf{x}^{(n)})[k]$$

- Then the cross-entropy loss is

$$\begin{aligned} \ell_{\mathcal{D}}(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}) &= - \sum_{n=1}^N \log P(y^{(n)} | \mathbf{x}^{(n)}, W^{(1)}, \dots, W^{(L)}) \\ &= - \sum_{n=1}^N \sum_{k=1}^C y[k] \log h_{W^{(1)}, \dots, W^{(L)}}(\mathbf{x}^{(n)})[k] \end{aligned}$$

Multi-dimensional Outputs



Recall: 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 $l = 1, \dots, L$
 - Compute $G^{(l)} = \nabla_{W^{(l)}} \ell_{\mathcal{D}}(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)})$ (???)
 - Update $W^{(l)}$: $W_{(t+1)}^{(l)} = W_{(t)}^{(l)} - \eta_0 G^{(l)}$
 - Increment t : $t = t + 1$
- Output: $W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}$

Matrix Calculus

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

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	$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial X_{11}} & \frac{\partial y}{\partial X_{12}} & \cdots & \frac{\partial y}{\partial X_{1Q}} \\ \frac{\partial y}{\partial X_{21}} & \frac{\partial y}{\partial X_{22}} & \cdots & \frac{\partial y}{\partial X_{2Q}} \\ \vdots & & & \vdots \\ \frac{\partial y}{\partial X_{P1}} & \frac{\partial y}{\partial X_{P2}} & \cdots & \frac{\partial y}{\partial X_{PQ}} \end{bmatrix}$

Matrix Calculus: Denominator Layout

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

The Chain Rule of Calculus

- If $y = f(z)$ and $z = g(x)$ then
computation graph is

$$\begin{array}{c} x \quad z \quad y \\ \square \longrightarrow \square \longrightarrow \square \end{array} \Rightarrow \frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \frac{\partial z}{\partial x}$$

- If $y = f(z_1, z_2)$ and $z_1 = g_1(x)$, $z_2 = g_2(x)$ then

$$\begin{array}{c} z_1 \quad y \\ x \quad \swarrow \quad \searrow \quad \nearrow \quad \nwarrow \quad \square \\ \square \quad \searrow \quad \swarrow \quad \nearrow \quad \nwarrow \quad \square \\ z_2 \end{array} \Rightarrow \frac{\partial y}{\partial x} = \frac{\partial y}{\partial z_1} \frac{\partial z_1}{\partial x} + \frac{\partial y}{\partial z_2} \frac{\partial z_2}{\partial x}$$

- If $y = f(\mathbf{z})$ and $\mathbf{z} = g(x)$ then

$$\begin{array}{c} z_1 \quad y \\ x \quad \swarrow \quad \searrow \quad \nearrow \quad \nwarrow \quad \square \\ \square \quad \searrow \quad \swarrow \quad \nearrow \quad \nwarrow \quad \square \\ \vdots \quad z_2 \\ \square \quad z_D \end{array} \Rightarrow \frac{\partial y}{\partial x} = \sum_{d=1}^D \frac{\partial y}{\partial z_d} \frac{\partial z_d}{\partial x}$$

Computing Gradients

$$\ell_{\mathcal{D}} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) = \sum_{n=1}^N \ell^{(n)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$$

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

$$= \begin{bmatrix} \frac{\partial \ell_{\mathcal{D}}}{\partial w_{1,0}^{(l)}} & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{1,1}^{(l)}} & \dots & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{1,d^{(l)-1}}^{(l)}} \\ \frac{\partial \ell_{\mathcal{D}}}{\partial w_{2,0}^{(l)}} & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{2,1}^{(l)}} & \dots & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{2,d^{(l)-1}}^{(l)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \ell_{\mathcal{D}}}{\partial w_{d^{(l)},0}^{(l)}} & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{d^{(l)},1}^{(l)}} & \dots & \frac{\partial \ell_{\mathcal{D}}}{\partial w_{d^{(l)},d^{(l)-1}}^{(l)}} \end{bmatrix}$$

$$\frac{\partial \ell_{\mathcal{D}}}{\partial w_{b,a}^{(l)}} = \sum_{n=1}^N \left[\frac{\partial \ell^{(n)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)}{\partial w_{b,a}^{(l)}} \right]$$

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 (dynamic programming)

Computing Partial Derivatives

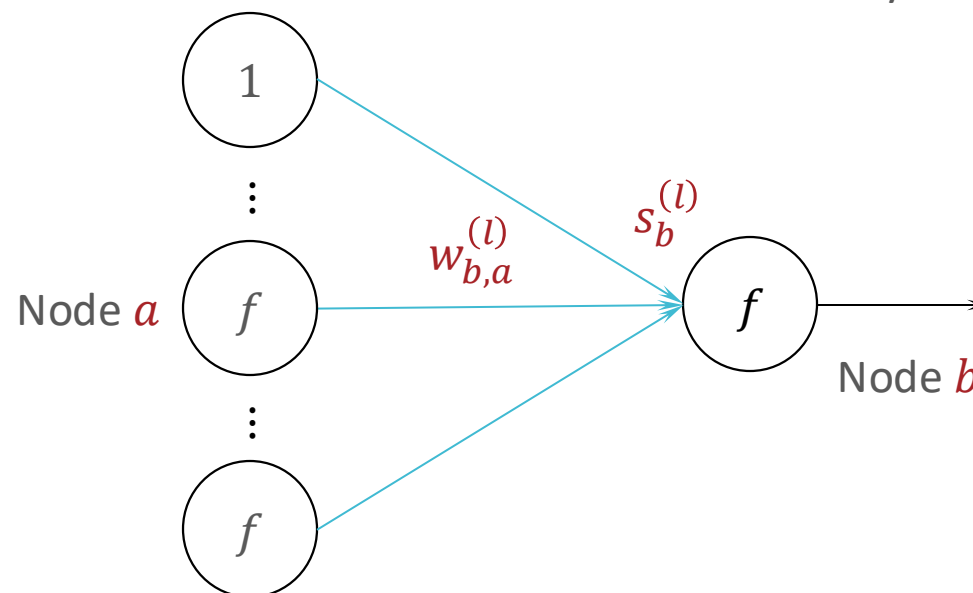
Computing $\nabla_{W^{(l)}} \ell_{\mathcal{D}} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$ reduces to computing

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

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

Layer $l - 1$

Layer l



Computing Partial Derivatives

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 $s_b^{(l)}$

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

$$\hookrightarrow s_b^{(l)} = \sum_{\alpha=0}^{d^{(l-1)}} o_{\alpha}^{(l-1)} w_{b,\alpha}^{(l)}$$

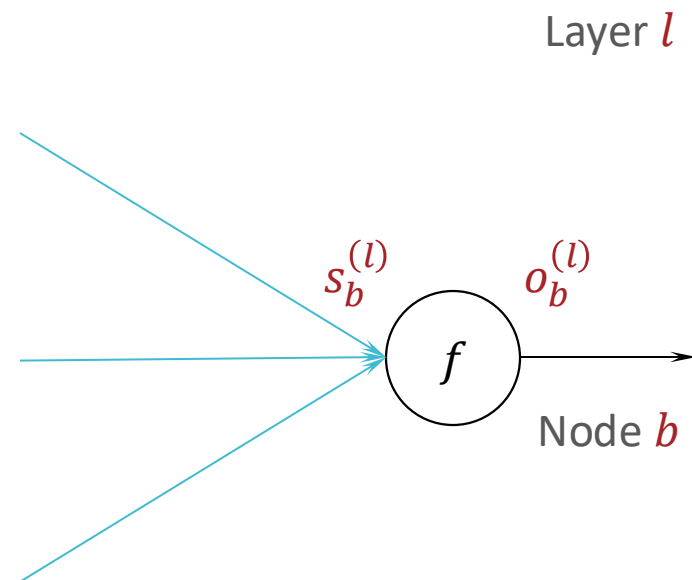
$$\delta_b^{(l)} := \frac{\partial \ell^{(n)}}{\partial s_b^{(l)}}$$

"sensitivity"

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

Computing Partial Derivatives

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



Computing Partial Derivatives

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

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

harder...

$$o_b^{(l)} = f(s_b^{(l)})$$

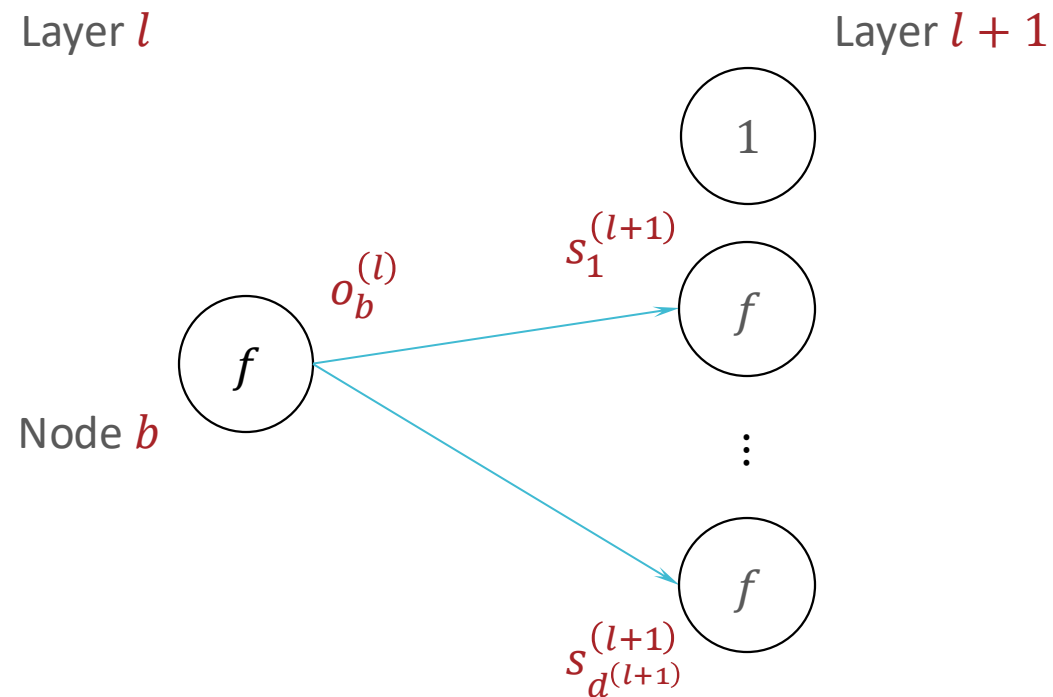
for example: $f(z) = \text{ReLU}(z)$
 $= \begin{cases} z & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$

$$\Rightarrow \frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} = \frac{\partial f(s_b^{(l)})}{\partial s_b^{(l)}}$$

for example $f(z) = \tanh(z)$
 $\frac{df}{dz} = 1 - (\tanh(z))^2$

Computing Partial Derivatives

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



Computing Partial Derivatives

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

$$\frac{\partial \ell^{(n)}}{\partial o_b^{(l)}} = \sum_{c=1}^{d^{(l+1)}} \frac{\partial \ell^{(n)}}{\partial s_c^{(l+1)}} \frac{\partial s_c^{(l+1)}}{\partial o_b^{(l)}}$$

$\delta_c^{(l+1)}$

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

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

Computing Partial Derivatives

$$\vec{s}^{(l)} = \begin{bmatrix} s_1^{(l)} \\ s_2^{(l)} \\ \vdots \\ s_{d^{(l)}}^{(l)} \end{bmatrix}$$

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

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

$$\delta^{(l)} := \nabla_{\vec{s}^{(l)}} \ell^{(n)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$$

if $f(z) = \tanh(z)$



Based solely on their shape alone, which of the following could be an expression for $\delta^{(l)} = \nabla_{\mathbf{s}^{(l)}} e(\mathbf{o}^{(l)}, \mathbf{y}^{(n)})$? Here \odot is the element-wise product operation.

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

0%

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

0%

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

0%

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

0%

Computing Partial Derivatives

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

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

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

Same size as

where \odot is the element-wise product operation

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

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

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

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

Computing Gradients

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

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

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

$\mathbf{o}^{(l-1)} \in \mathbb{R}^{(d^{(l-1)} + 1) \times 1}$

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

$\Rightarrow \boldsymbol{\delta}^{(l)} \mathbf{o}^{(l-1)T} \in \mathbb{R}^{d^{(l)} \times (d^{(l-1)} + 1)}$

Computing Partial Derivatives

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

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

$$\begin{aligned}\delta_1^{(L)} &= \frac{\partial e(o_1^{(L)}, y^{(n)})}{\partial s_1^{(L)}} = \frac{\partial}{\partial s_1^{(L)}} (o_1^{(L)} - y^{(n)})^2 \\ &= 2(o_1^{(L)} - y^{(n)}) \frac{\partial o_1^{(L)}}{\partial s_1^{(L)}} = 2(o_1^{(L)} - y^{(n)}) \left(1 - (o_1^{(L)})^2\right)\end{aligned}$$

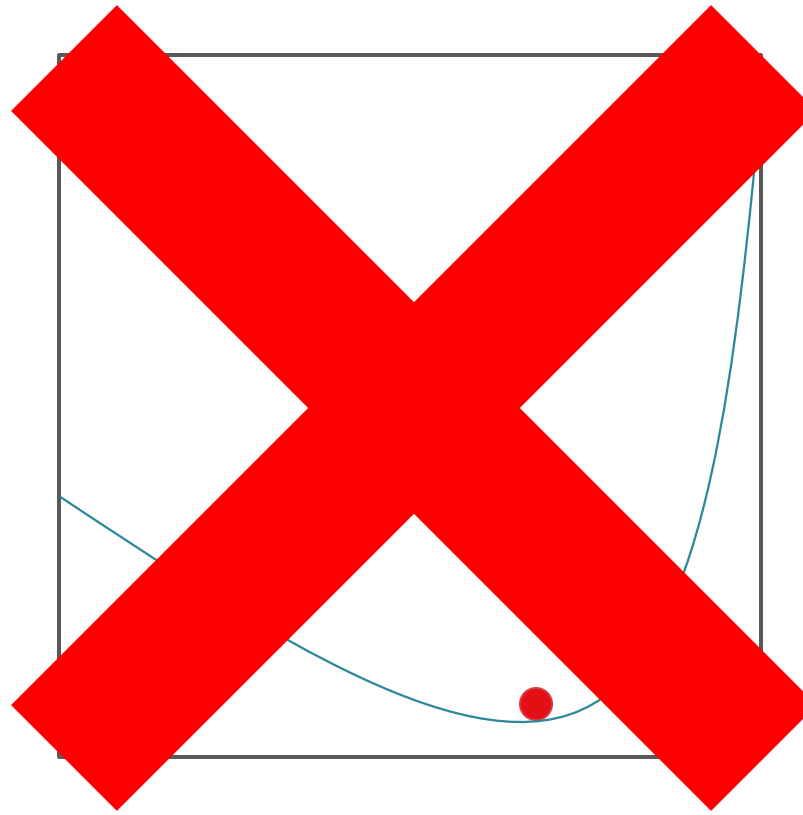
when $f(\cdot) = \tanh(\cdot)$

Back-propagation

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

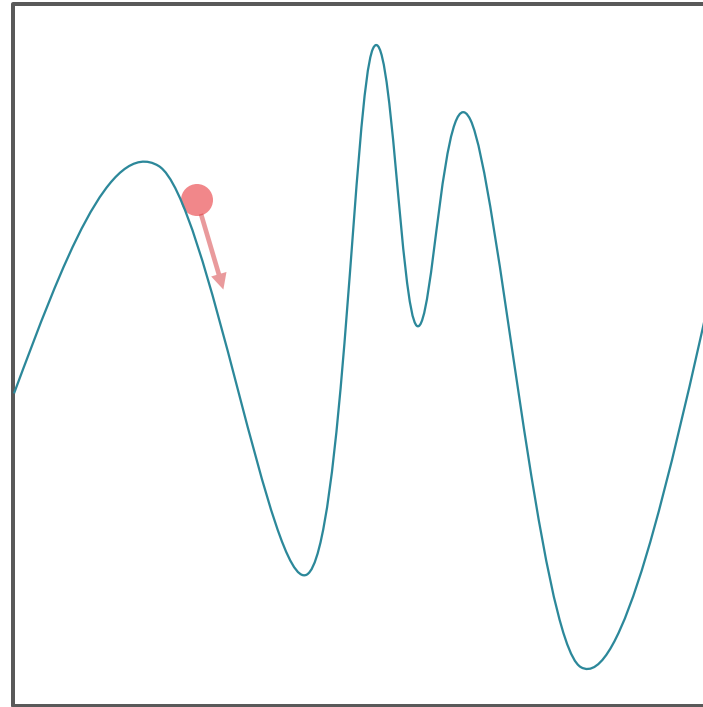
Recall: Gradient Descent

- 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 Neural Networks

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{SGD}^{(0)}$
- 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 a data point from $\mathcal{D}, (\mathbf{x}^{(n)}, y^{(n)})$
 - b. Compute the pointwise gradient using backpropagation

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

- c. Update $W^{(l)}: W_{t+1}^{(l)} \leftarrow W_t^{(l)} - \eta_{SGD}^{(0)} G^{(l)} \forall l$
 - d. Increment $t: t \leftarrow t + 1$
- Output: $W_t^{(1)}, \dots, W_t^{(L)}$

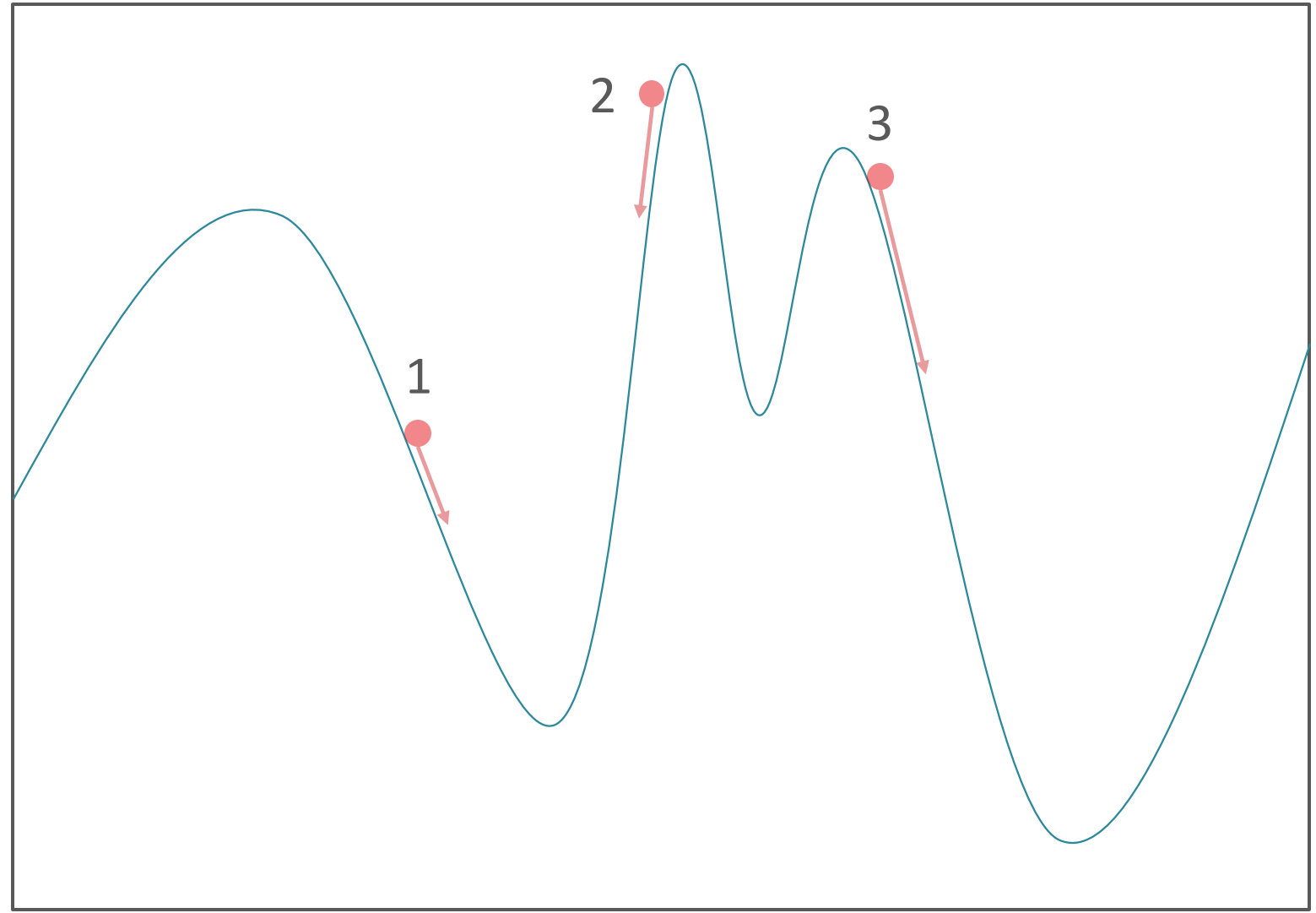
Mini-batch Stochastic Gradient Descent for Neural Networks

- 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*,
$$G^{(l)} = \frac{1}{B} \sum_{b=1}^B \nabla_{W^{(l)}} \ell^{(b)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) \quad \forall l$$
 - c. Update $W^{(l)}: W_{t+1}^{(l)} \leftarrow W_t^{(l)} - \eta_{MB}^{(0)} G^{(l)} \quad \forall l$
 - d. Increment $t: t \leftarrow t + 1$
- Output: $W_t^{(1)}, \dots, W_t^{(L)}$

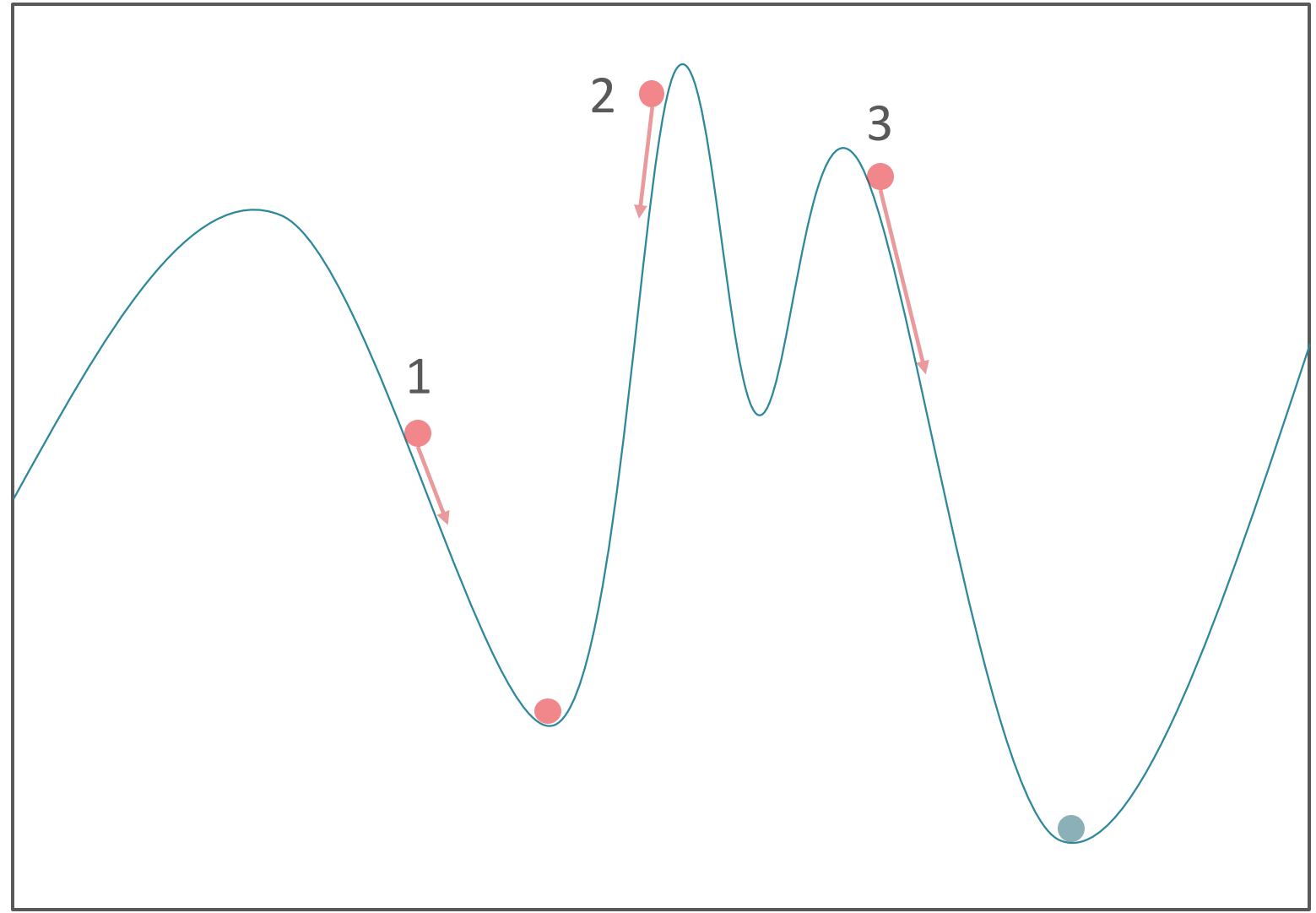
Random Restarts

- Run mini-batch gradient descent (~~with momentum & adaptive gradients~~) multiple times, each time starting with a ***different, random*** initialization for the weights.
- Compute the training error of each run at termination and return the set of weights that achieves the lowest training error.

Random Restarts



Random Restarts



Key Takeaways

- Backpropagation for efficient gradient computation
- Advanced optimization and regularization techniques for neural networks
 - SGD and Mini-batch gradient descent
 - Random restarts
 - Jitter & dropout act like regularization for neural networks by preventing them fitting the training dataset perfectly