

10-301/601: Introduction to Machine Learning

Lecture 13 – Backpropagation

Henry Chai & Zack Lipton

10/11/23

Front Matter

- Announcements
 - HW3 released 10/4, due 10/11 (today!) at 11:59 PM
 - HW4 released 10/11 (today!), due 10/25 (**after fall break**) at 11:59 PM
 - Project details will be released on 10/13 (Friday)
 - Midterm exam on 10/31 from 6:30 – 8:30 PM
 - If you have a conflict with this date/time fill out the conflict on Piazza ASAP
- Recommended Readings
 - Mitchell, [Chapters 4.1 – 4.6](#)

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

		Numerator		
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	$\frac{\partial y}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$
	matrix	$\frac{\partial y}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$

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	$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x} \right]$	$\frac{\partial \mathbf{y}}{\partial x} = \left[\frac{\partial y_1}{\partial x} \quad \frac{\partial y_2}{\partial x} \quad \dots \quad \frac{\partial y_N}{\partial x} \right]$
vector	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_N}{\partial x_2} \\ \vdots & \vdots & \dots & \vdots \\ \frac{\partial y_1}{\partial x_P} & \frac{\partial y_2}{\partial x_P} & \dots & \frac{\partial y_N}{\partial x_P} \end{bmatrix}$

Computing Gradients

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

Computing Gradients: Intuition

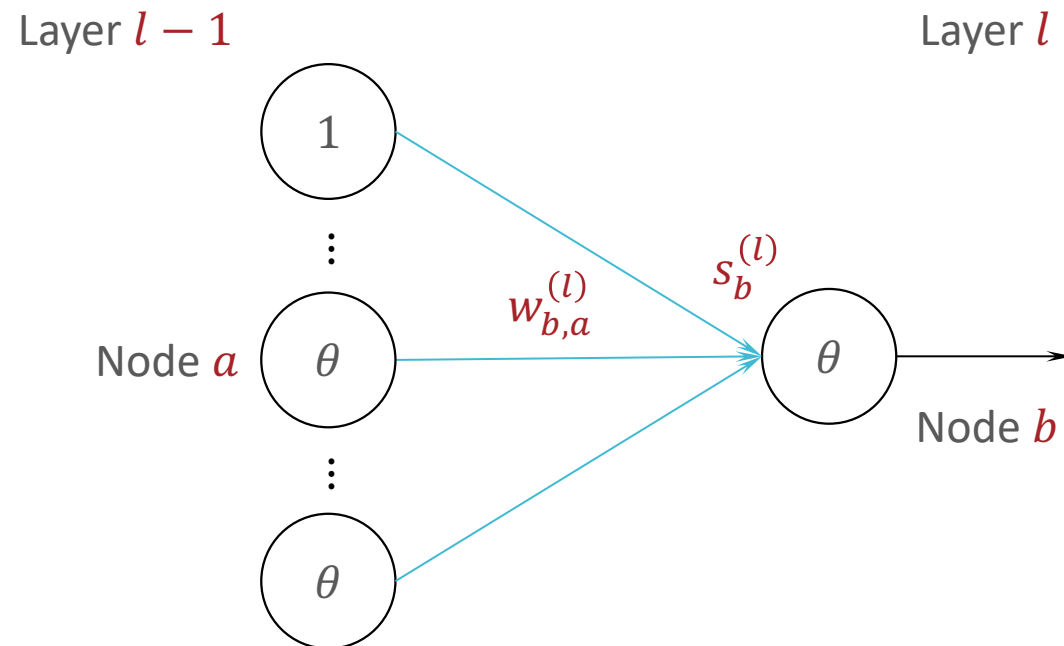
- A weight affects the prediction of the network (and therefore the error) through downstream signals/outputs
 - Use the chain rule!
- Any weight going into the same node will affect the prediction through the same downstream path
 - Compute derivatives starting from the last layer and move “backwards”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

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 $s_b^{(l)}$

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

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

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

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

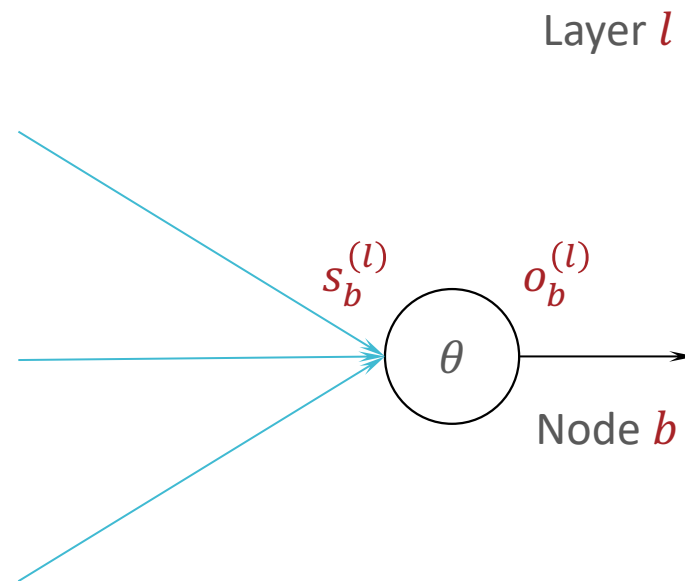
Insight: $w_{b,a}^{(l)}$ *only* affects $\ell^{(i)}$ via $s_b^{(l)}$

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

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

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 $o_b^{(l)}$

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

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

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

Aside: Vanishing Gradients


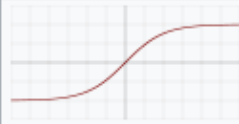

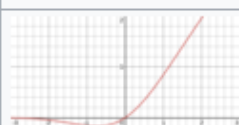
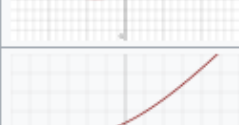



Insight: $s_b^{(l)}$ only affects $\ell^{(i)}$ via $o_b^{(l)}$

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

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

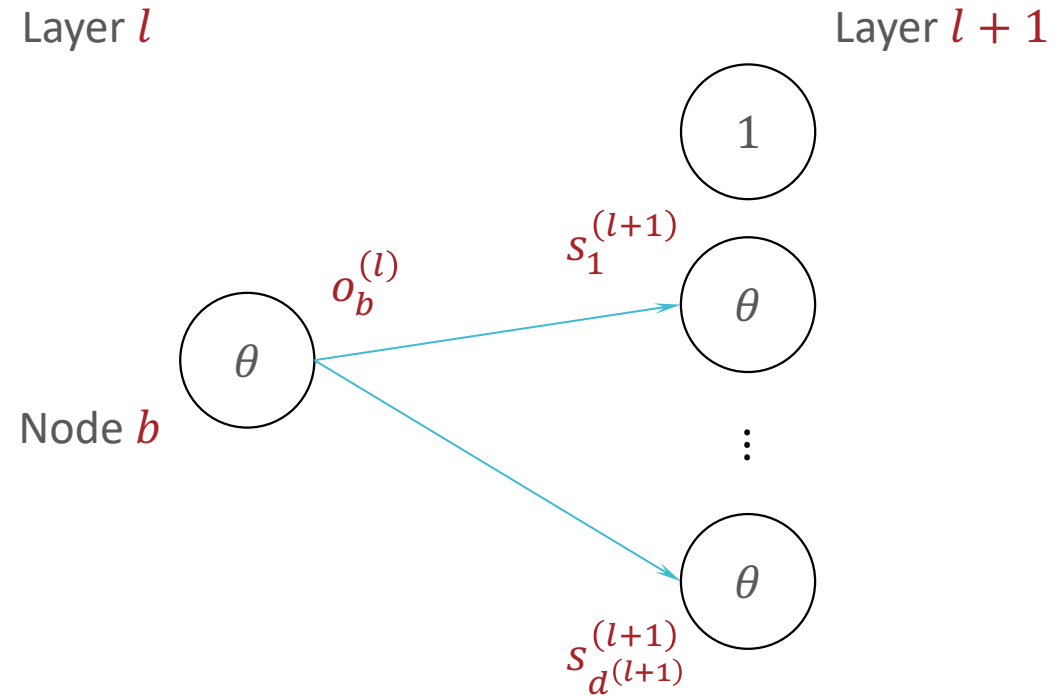
when $\theta(\cdot) = \tanh(\cdot)$

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]		$\begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max\{0, x\} = x \mathbf{1}_{x>0}$
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]		$\begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ with parameter α
Leaky rectified linear unit (Leaky ReLU) ^[11]		$\begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$
Parametric rectified linear unit (PReLU) ^[12]		$\begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$ with parameter α

Computing Partial Derivatives

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



Computing Partial Derivatives

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

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

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

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

Computing Partial Derivatives

$$\begin{aligned}\delta_b^{(l)} &= \frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right) \\ &= \left(\sum_{c=1}^{d^{(l+1)}} \delta_c^{(l+1)} \left(w_{c,b}^{(l+1)} \right) \right) \left(1 - \left(o_b^{(l)} \right)^2 \right) \\ \boldsymbol{\delta}^{(l)} &:= \nabla_{\mathbf{s}^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)\end{aligned}$$

Computing Partial Derivatives

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

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

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

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

Computing Gradients

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

$$\nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) = \boldsymbol{\delta}^{(l)} \mathbf{o}^{(l-1)T}$$

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

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

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

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: $\delta^{(L)} = \delta_1^{(L)}$, $\mathbf{o}^{(L)} = o_1^{(L)}$

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

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

$$= 2 \left(o_1^{(L)} - y^{(i)} \right) \frac{\partial o_1^{(L)}}{\partial s_1^{(L)}} = 2 \left(o_1^{(L)} - y^{(i)} \right) \left(1 - \left(o_1^{(L)} \right)^2 \right)$$

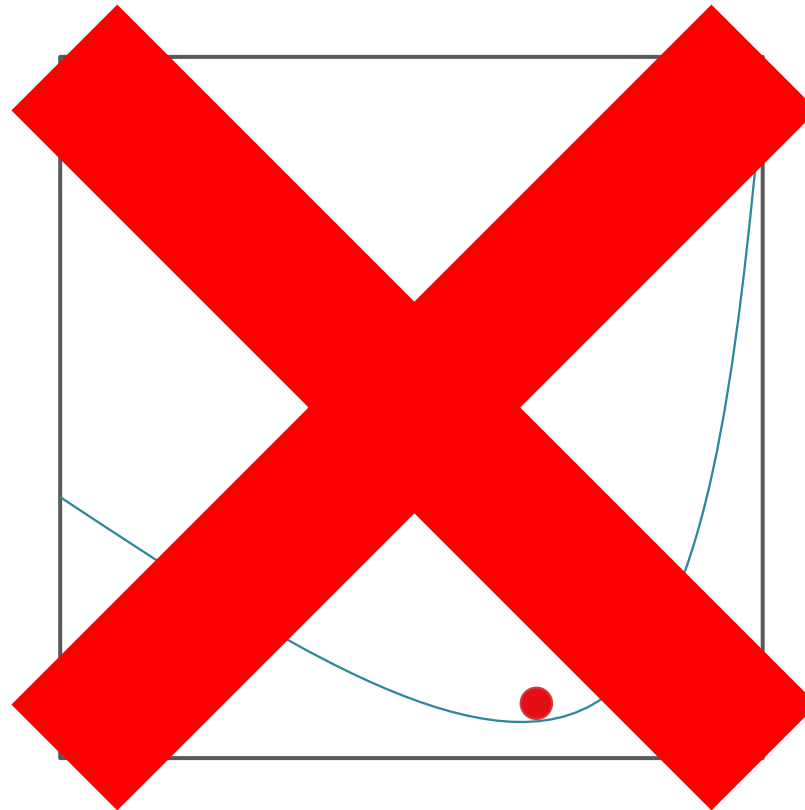
when $\theta(\cdot) = \tanh(\cdot)$

Back- propagation

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

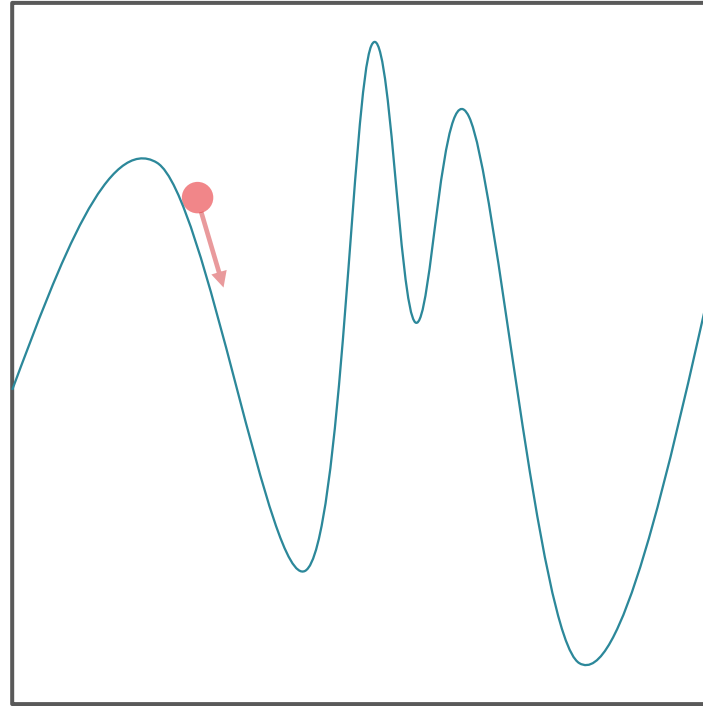
Recall: Gradient Descent

- 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 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 $G^{(l)} = \nabla_{W^{(l)}} \ell^{(i)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$
 - Update $W^{(l)}$: $W_{(t+1)}^{(l)} = W_{(t)}^{(l)} - \eta_0 G^{(l)}$
 - Increment t : $t = t + 1$
- Output: $W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}$

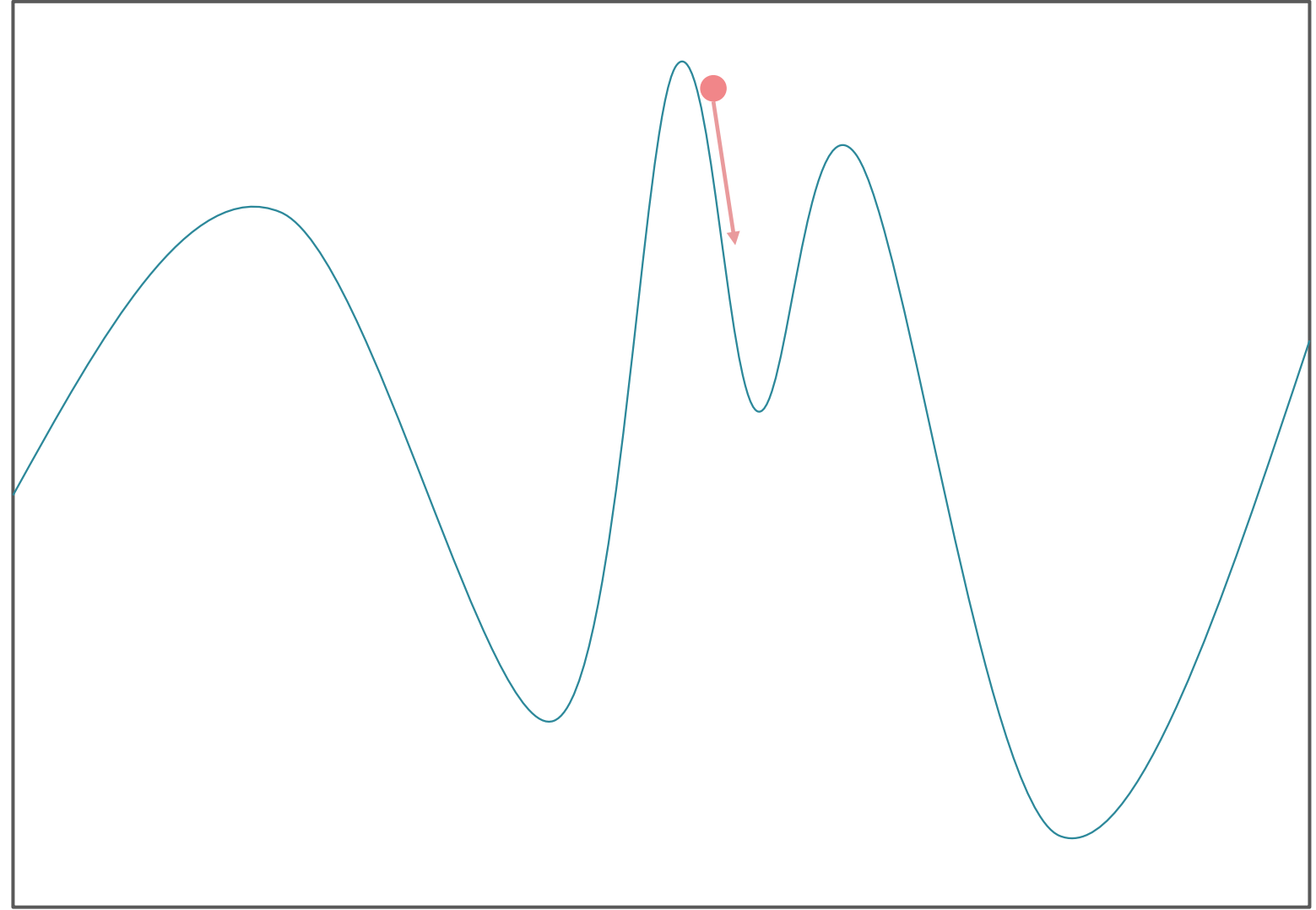
Mini-batch Stochastic Gradient Descent for Learning

- 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,
$$G^{(l)} = \frac{1}{B} \sum_{b=1}^B \nabla_{W^{(l)}} \ell^{(b)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) \quad \forall l$$
 - 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: $W_t^{(1)}, \dots, W_t^{(L)}$

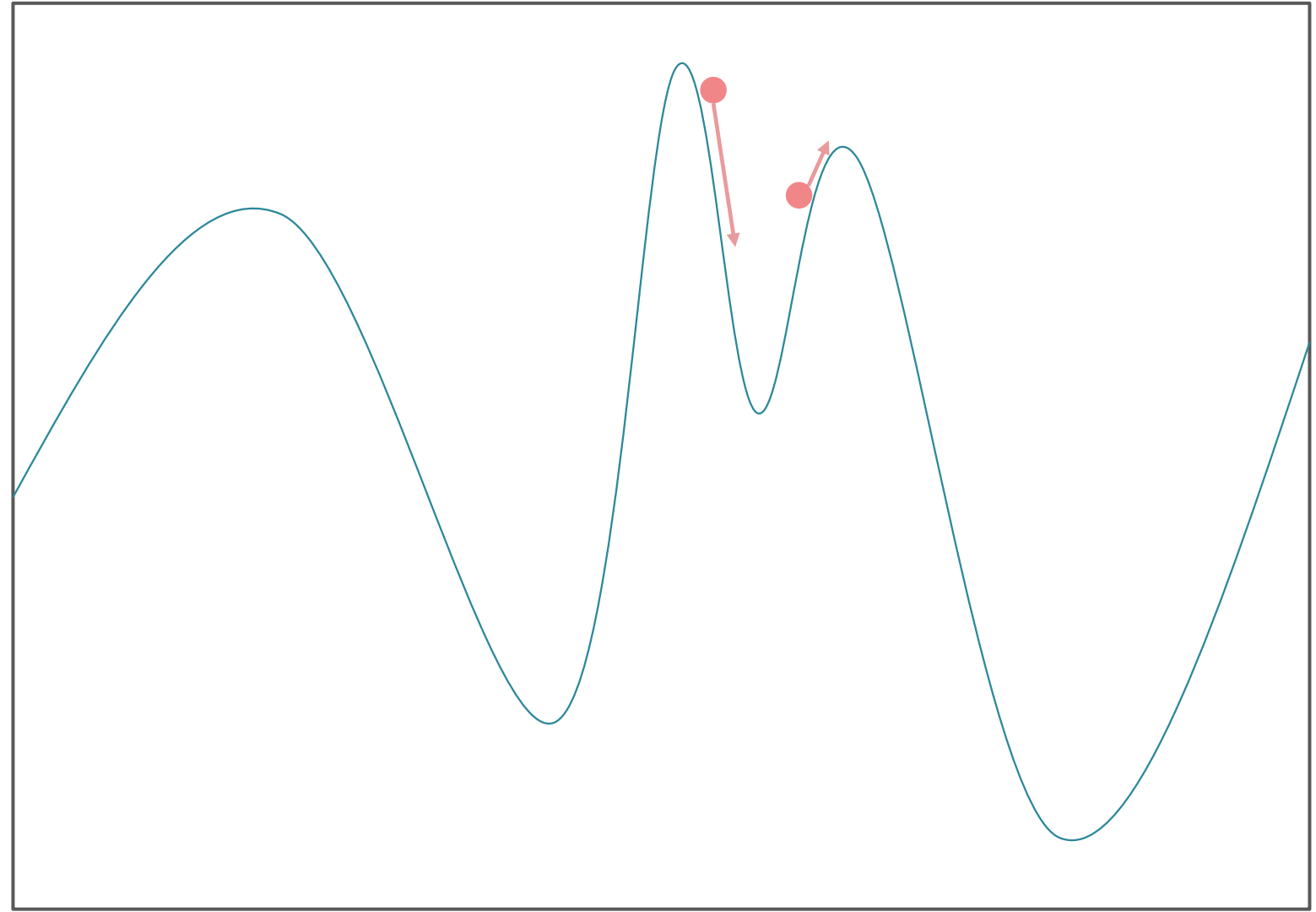
Mini-batch Stochastic Gradient Descent with Momentum for Learning

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{MB}^{(0)}, B, \beta$
- 1. Initialize all weights $W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}$ to small, random numbers and set $t = 0, G_{-1}^{(l)} = 0 \odot W^{(l)} \forall l = 1, \dots, L$
- 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*,
$$G_t^{(l)} = \frac{1}{B} \sum_{b=1}^B \nabla_{W^{(l)}} \ell^{(b)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) \forall l$$
 - c. Update $W^{(l)}$: $W_{t+1}^{(l)} \leftarrow W_t^{(l)} - \eta_{MB}^{(0)} \left(\beta G_{t-1}^{(l)} + G_t^{(l)} \right) \forall l$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $W_t^{(1)}, \dots, W_t^{(L)}$

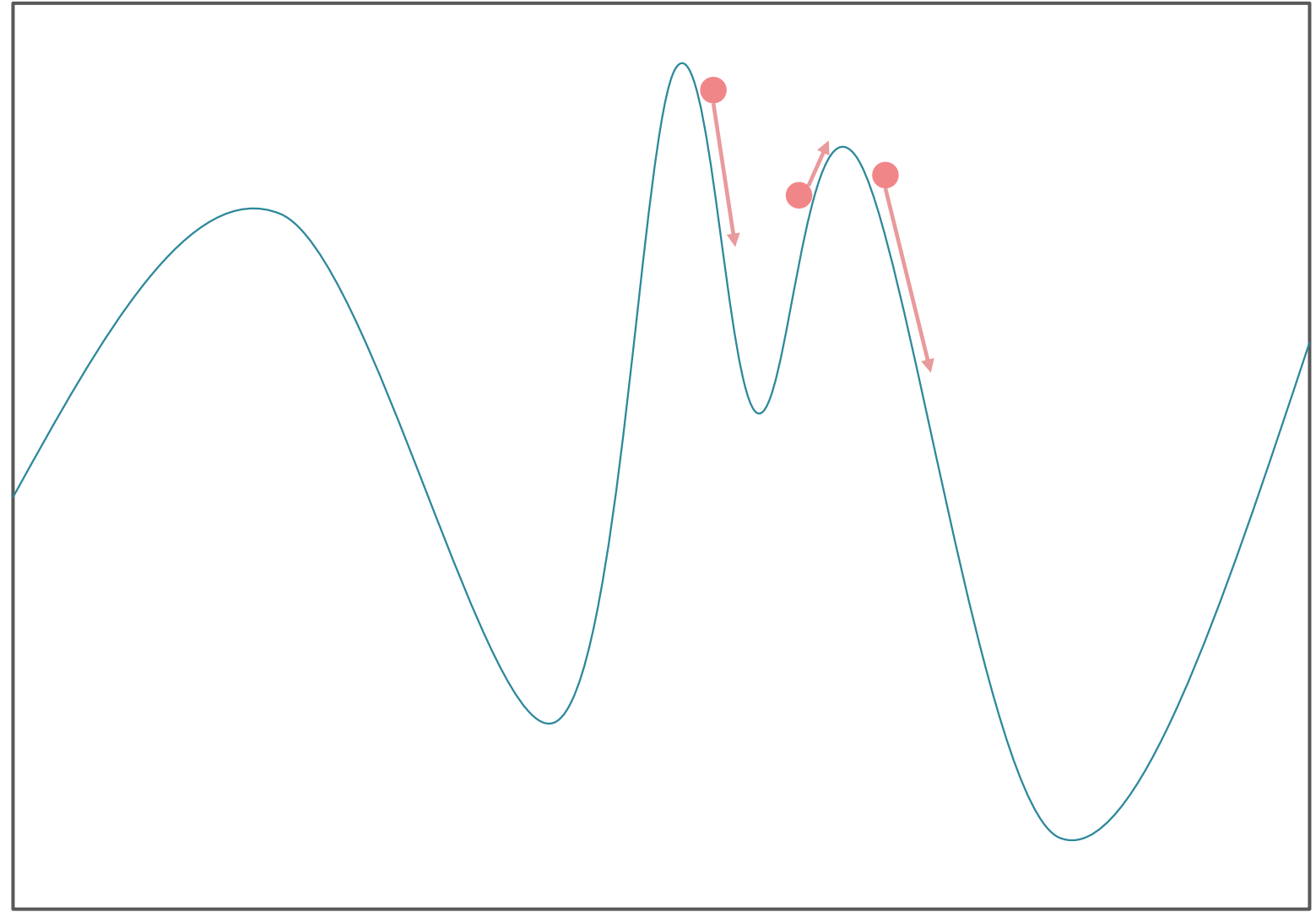
Mini-batch Stochastic Gradient Descent with Momentum for Learning



Mini-batch Stochastic Gradient Descent with Momentum for Learning



Mini-batch Stochastic Gradient Descent with Momentum for Learning



Mini-batch Stochastic Gradient Descent with Adaptive Gradients for Learning

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{MB}^{(0)}, B, \epsilon$
- 1. Initialize all weights $W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}$ to small, random numbers and set $t = 0, S_{-1}^{(l)} = 0 \odot W^{(l)} \forall l = 1, \dots, L$
- 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*,

$$G_t^{(l)} = \frac{1}{B} \sum_{b=1}^B \nabla_{W^{(l)}} \ell^{(b)} \left(W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right) \forall l$$

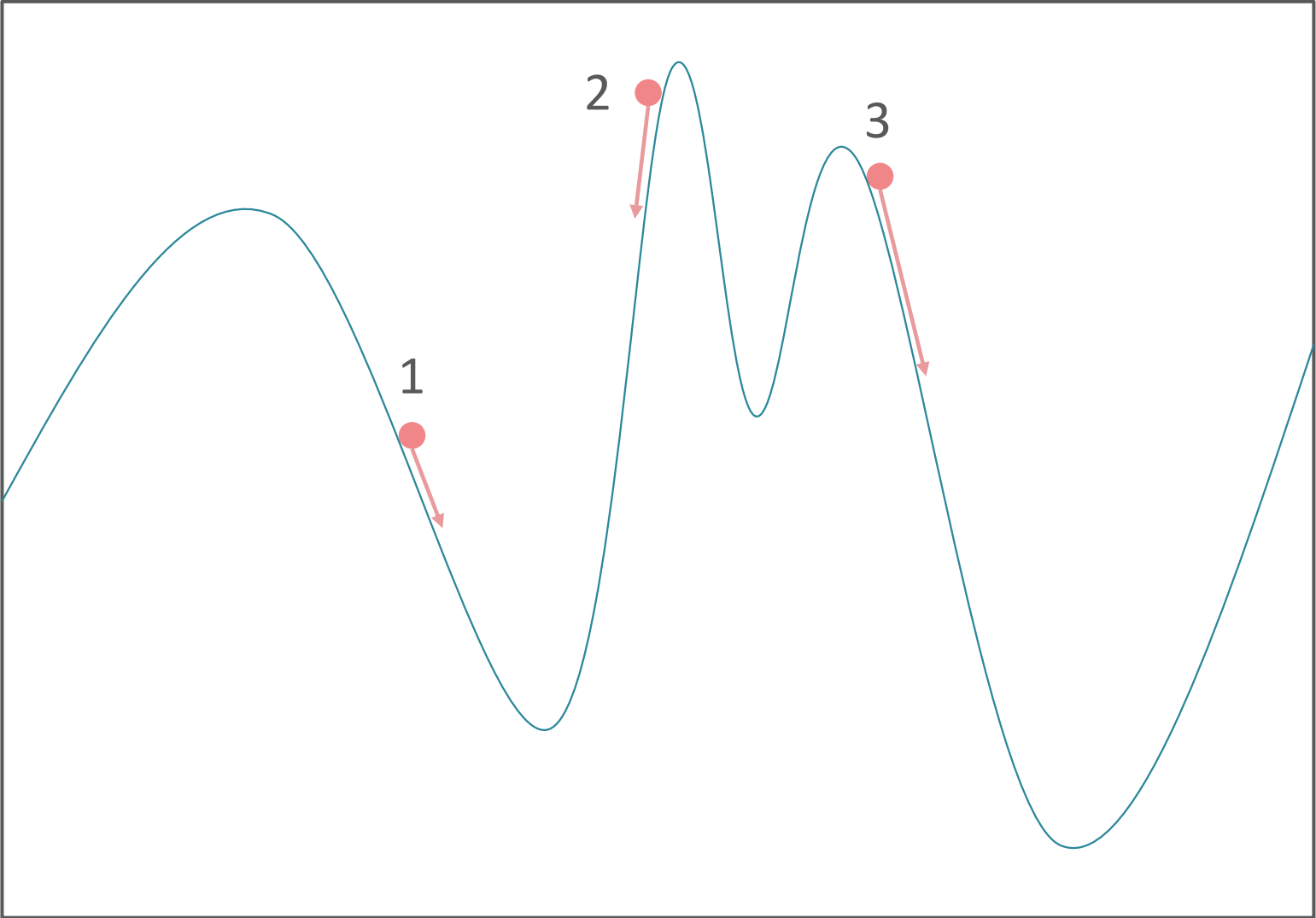
- c. Update $S^{(l)}: S_t^{(l)} = S_{t-1}^{(l)} + G_t^{(l)} \odot G_t^{(l)} \forall l$
- d. Update $W^{(l)}: W_{t+1}^{(l)} \leftarrow W_t^{(l)} - \frac{\eta_{MB}^{(0)}}{\sqrt{S_t^{(l)} + \epsilon}} \odot G_t^{(l)} \forall l$
- e. Increment $t: t \leftarrow t + 1$

- Output: $W_t^{(1)}, \dots, W_t^{(L)}$

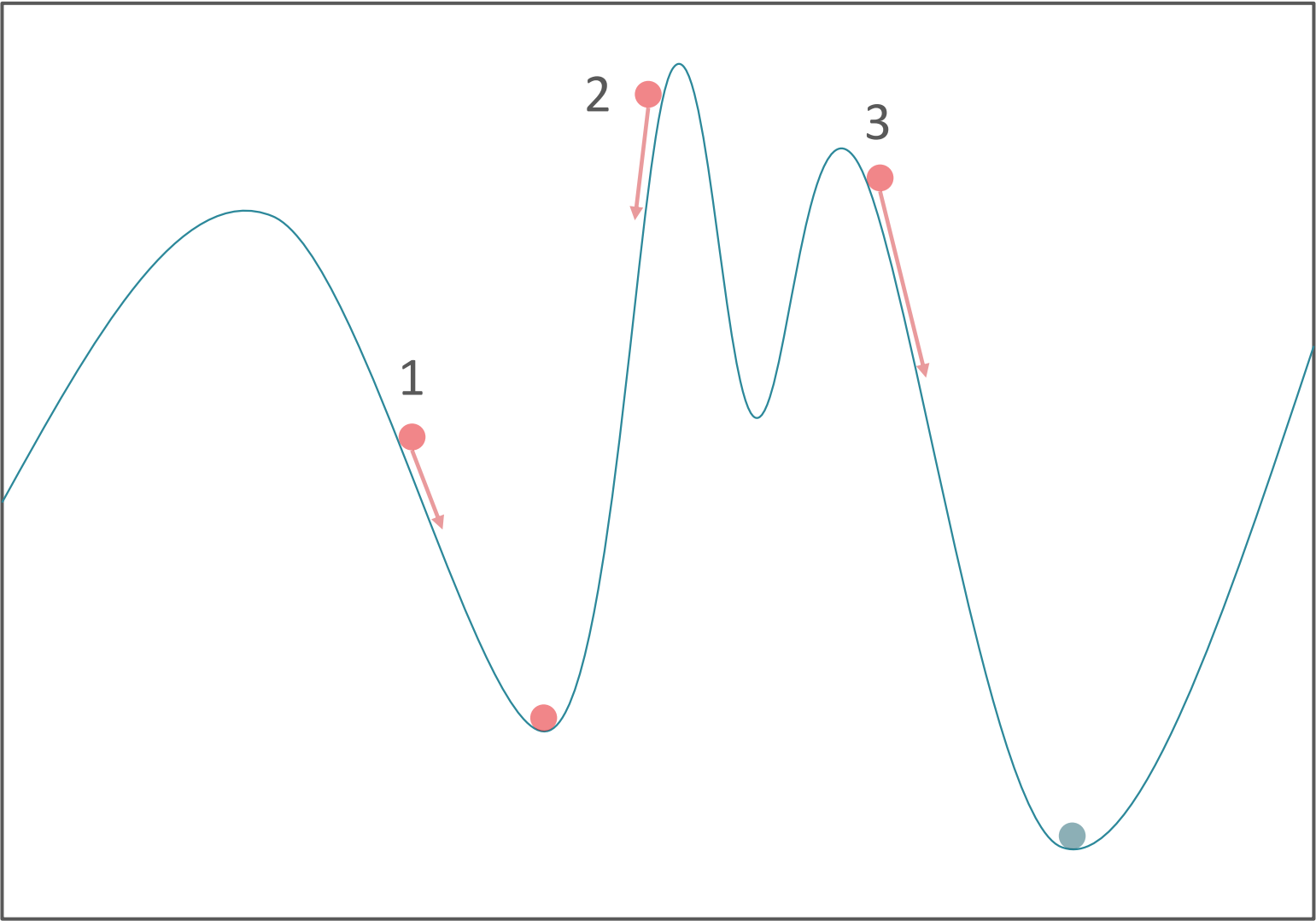
Random Restarts

- 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

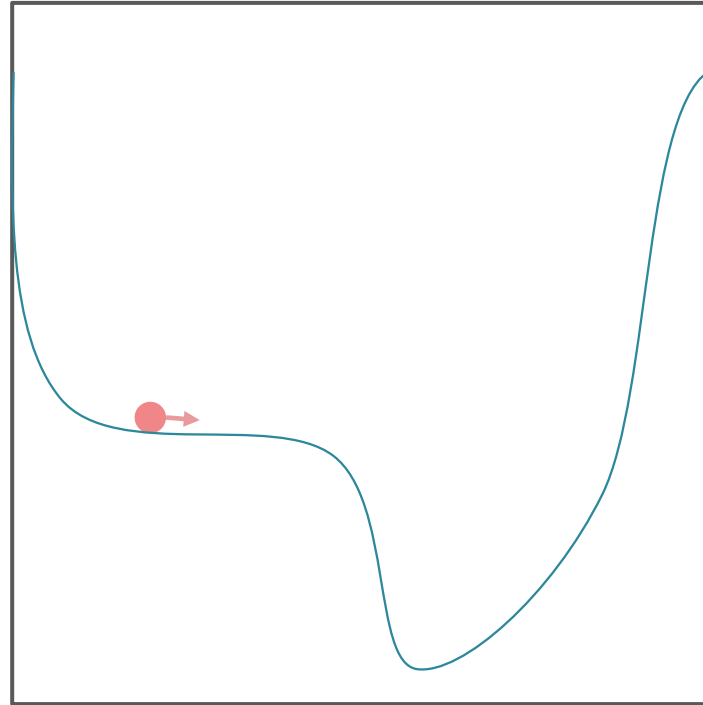


Random Restarts



Terminating Gradient Descent

- For non-convex surfaces, the gradient's magnitude is often not a good metric for proximity to a minimum



Terminating Gradient Descent “Early”

- For non-convex surfaces, the gradient’s magnitude is often not a good metric for proximity to a minimum
- Combine multiple termination criteria e.g. only stop if enough iterations have passed and the improvement in error is small
- Alternatively, terminate early by using a validation data set: if the validation error starts to increase, just stop!
 - Early stopping asks like regularization by **limiting how much of the hypothesis set** is explored

Neural Networks and Regularization

- Minimize $\ell_{AUG}^{(i)}(W^{(1)}, \dots, W^{(L)}, \lambda_C)$
 $= \ell^{(i)}(W^{(1)}, \dots, W^{(L)}) + \lambda_C r(W^{(1)}, \dots, W^{(L)})$

e.g. L2 regularization

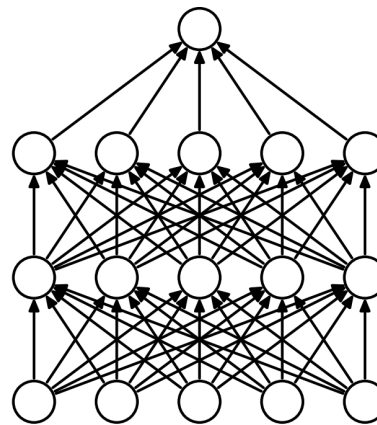
$$r(W^{(1)}, \dots, W^{(L)}) = \sum_{l=1}^L \sum_{i=0}^{d^{(l-1)}} \sum_{j=1}^{d^{(l)}} (w_{j,i}^{(l)})^2$$

Neural Networks and “Strange” Regularization (Bishop, 1995)

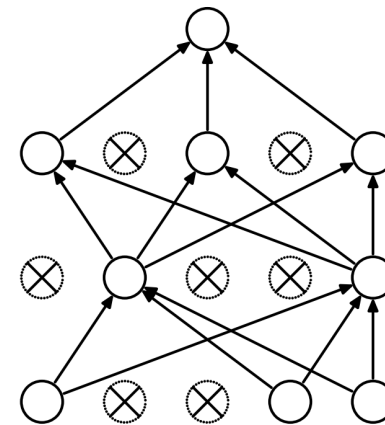
- Jitter
 - In each iteration of gradient descent, add some random noise or “jitter” to each training data point
 - Instead of computing the gradient w.r.t. $(\mathbf{x}^{(i)}, y^{(i)})$, use $(\mathbf{x}^{(i)} + \boldsymbol{\epsilon}, y^{(i)})$ where $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 I)$
 - Makes neural networks resilient to input noise
 - Has been proven to be equivalent to using a certain kind of regularizer r for some error metrics

Neural Networks and “Strange” Regularization (Srivastava et al., 2014)

- Dropout
 - In each iteration of gradient descent, randomly remove some of the nodes in the network
 - Compute the gradient using only the remaining nodes
 - The weights on edges going into and out of “dropped out” nodes are not updated



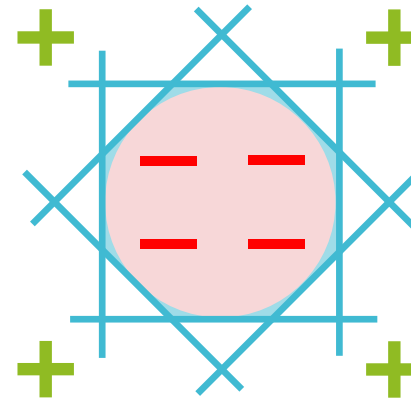
(a) Standard Neural Net



(b) After applying dropout.

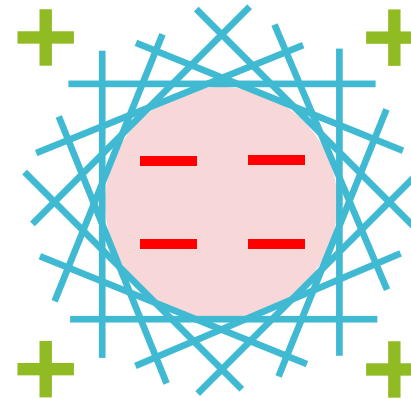
MLPs as Universal Approximators

- Theorem: any function that can be decomposed into perceptrons can be modelled exactly using a 3-layer MLP
- Any smooth decision boundary can be approximated to an arbitrary precision using a finite number of perceptrons



MLPs as Universal Approximators

- Theorem: any function that can be decomposed into perceptrons can be modelled exactly using a 3-layer MLP
- Any smooth decision boundary can be approximated to an arbitrary precision using a finite number of perceptrons



- Theorem: Any smooth decision boundary can be approximated to an arbitrary precision using a 3-layer MLP

NNs as Universal Approximators (Cybenko, 1989 & Hornik, 1991)

- Theorem: Any bounded, continuous function can be approximated to an arbitrary precision using a 2-layer (1 hidden layer) feed-forward NN if the activation function, θ , is continuous, bounded and non-constant.

NNs as Universal Approximators (Cybenko, 1988)

- Theorem: Any function can be approximated to an arbitrary precision using a 3-layer (2 hidden layers) feed-forward NN if the activation function, θ , is continuous, bounded and non-constant.

Deep Learning

- From Wikipedia's page on Deep Learning...

Definition [\[edit \]](#)

Deep learning is a class of **machine learning algorithms** that^{[11](pp199–200)} uses multiple layers to progressively extract higher level features from the raw input. For example, in image processing, lower layers may identify edges, while higher layers may identify the concepts relevant to a human such as digits or letters or faces.

- Deep learning = more than one layer

Key Takeaways

- Backpropagation for efficient gradient computation
- Advanced optimization and regularization techniques for neural networks
 - Momentum can be used to break out of local minima
 - Adagrad helps when parameters behave differently w.r.t. step sizes
 - Random restarts
 - Jitter & dropout act like regularization for neural networks by preventing them fitting the training dataset perfectly
- MLPs and neural networks of sufficient depth are universal approximators