Cornell Bowers C·IS
College of Computing and Information Science

# Optimization

CS4782: Intro to Deep Learning
Varsha Kishore, Justin Lovelace, Gary Wei

Before we start:

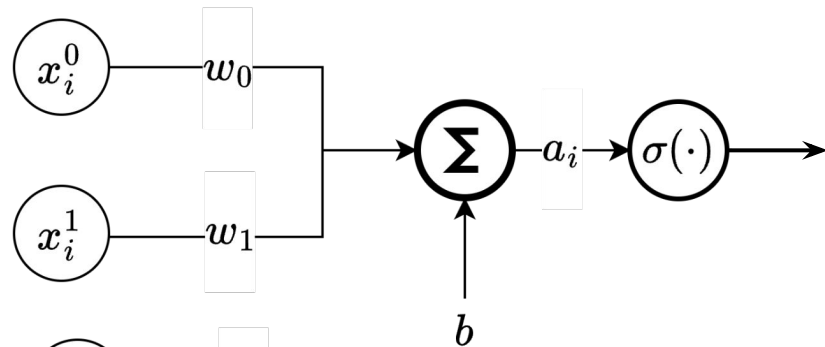Look at recap quizzes 1,2

Discuss with your neighbors

Choose option A, B, C

Which decision boundary CANNOT be learned by this network?

A.

B.

C.

$x_i^0$

$x_i^1$

$\hat{\mathbf{y}}_i$

## Course Announcement

- If you are in **5782**
  - Paper quizzes are mandatory (10%)

- If you are in **4782**
  - Paper quizzes are optional
  - If you do them, we will use the better grade with or without quizzes

# Agenda

- Backpropagation
- Optimizers
  - Gradient Descent
  - Stochastic Gradient Descent
  - SGD w. Momentum
  - AdaGrad
  - RMSProp
  - Adam
- Learning rate scheduling

# Calculus Review: The Chain Rule

Lagrange's Notation: $\text{If } h(x) = f(g(x)), \text{ then } h' = f'(g(x))g'(x)$

Leibniz's Notation: $\text{If } z = h(y), y = g(x), \text{ then } \frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$

# Calculus Review: The Chain Rule

Lagrange's Notation: If $h(x) = f(g(x))$, then $h' = f'(g(x))g'(x)$

Leibniz's Notation: If $z = h(y), y = g(x)$, then $\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$

Example: If $z = \ln(y), y = x^2$, then

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

$$= (\frac{1}{y})(2x) = (\frac{1}{x^2})(2x) = \frac{2}{x}$$

# Multivariate Chain Rule

# Multivariate Chain Rule



If $f(u)$ is $z = f(v(u), w(u))$, then

$$\frac{\partial f}{\partial u} = (\frac{\partial v}{\partial u}\frac{\partial z}{\partial v} + \frac{\partial w}{\partial u}\frac{\partial z}{\partial w})$$

# Backpropagation- An Example



$v = u^3$

$u$

$z = v \cdot w$

$w = u + u$

# Backpropagation- An Example

Forward



$$v = u^3$$

$$5$$
$$u$$

$$z = v \cdot w$$

$$w = u + u$$

# Backpropagation- An Example

Forward



$$125$$
$$v = u^3$$

$$5$$
$$u$$

$$1250$$
$$z = v \cdot w$$

$$10$$
$$w = u + u$$

# Backpropagation- An Example

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v}{\partial u} \cdot \frac{\partial z}{\partial v} + \frac{\partial w}{\partial u} \cdot \frac{\partial z}{\partial w} \right)$$

Forward
Backward



125
$v = u^3$

5
$u$

10
$w = u + u$

1250
$z = v \cdot w$

# Backpropagation- An Example

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v}{\partial u} \cdot \frac{\partial z}{\partial v} + \frac{\partial w}{\partial u} \cdot \frac{\partial z}{\partial w} \right)$$

Forward
Backward



125

$v = u^3$

$\frac{\partial z}{\partial v} = 10$

$\frac{\partial z}{\partial v} = w = 10$

5

$u$

1250

$z = v \cdot w$

10

$w = u + u$

# Backpropagation- An Example

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v}{\partial u} \cdot \frac{\partial z}{\partial v} + \frac{\partial w}{\partial u} \cdot \frac{\partial z}{\partial w} \right)$$

Forward
Backward

$\frac{\partial z}{\partial v} = w = 10$

125
$v = u^3$
$\frac{\partial z}{\partial v} = 10$

5
$u$

1250
$z = v \cdot w$

10
$w = u + u$
$\frac{\partial z}{\partial w} = 125$

$\frac{\partial z}{\partial w} = v = 125$

https://windowsontheory.org/2020/11/03/yet-another-backpropagation-tutorial/

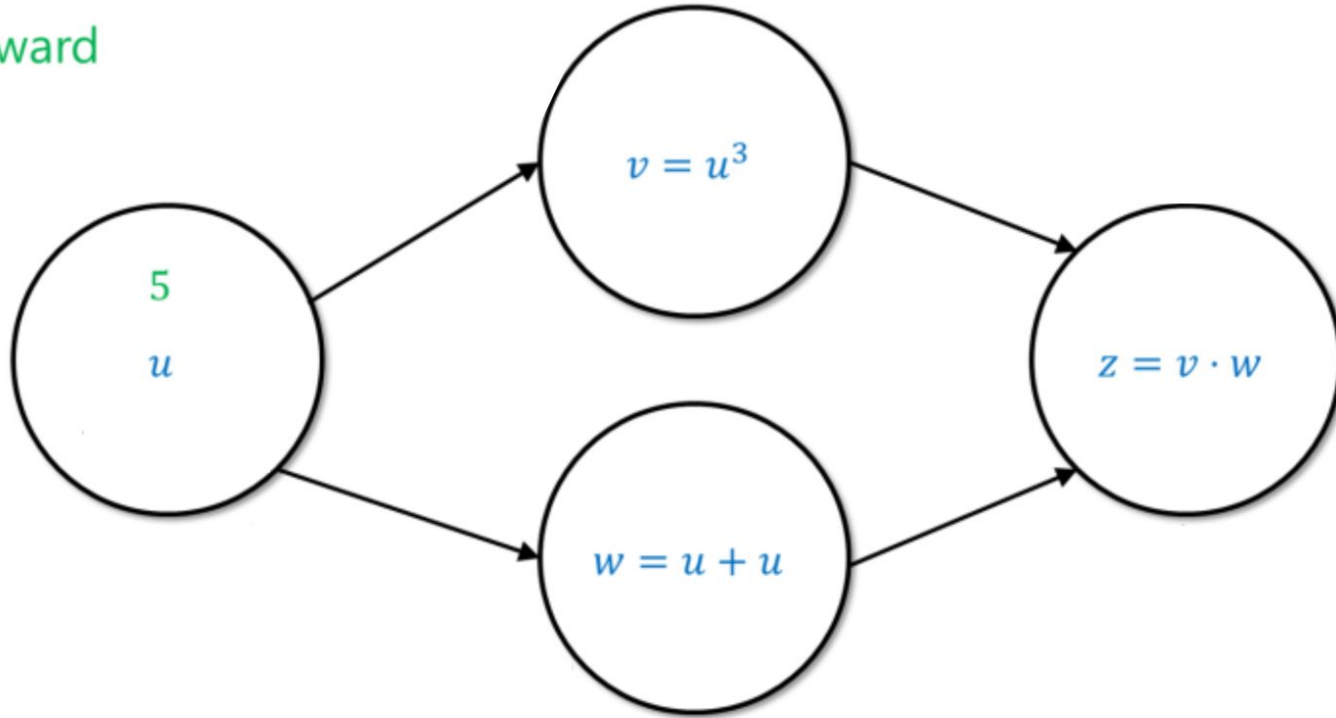# Backpropagation- An Example

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v}{\partial u} \cdot \frac{\partial z}{\partial v} + \frac{\partial w}{\partial u} \cdot \frac{\partial z}{\partial w} \right)$$

Forward
Backward

$$\frac{\partial v}{\partial u} \cdot \frac{\partial z}{\partial v} = 10 \cdot 3 \cdot u^2 = 750$$

$$\frac{\partial z}{\partial v} = w = 10$$

125
$v = u^3$
$$\frac{\partial z}{\partial v} = 10$$

5
$u$
$$\frac{\partial z}{\partial u} = 1000$$

1250
$z = v \cdot w$

10
$w = u + u$
$$\frac{\partial z}{\partial w} = 125$$

$$\frac{\partial w}{\partial u} \cdot \frac{\partial z}{\partial w} = 2 \cdot 125 = 250$$

$$\frac{\partial z}{\partial w} = v = 125$$

# Backpropagation- Key Idea



If you know
$$\frac{\partial z}{\partial v_1}, \frac{\partial z}{\partial v_2}, \frac{\partial z}{\partial v_3}$$

You can compute $\frac{\partial z}{\partial u}$

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v_1}{\partial u} \cdot \frac{\partial z}{\partial v_1} + \frac{\partial v_2}{\partial u} \cdot \frac{\partial z}{\partial v_2} + \frac{\partial v_3}{\partial u} \cdot \frac{\partial z}{\partial v_3} \right)$$
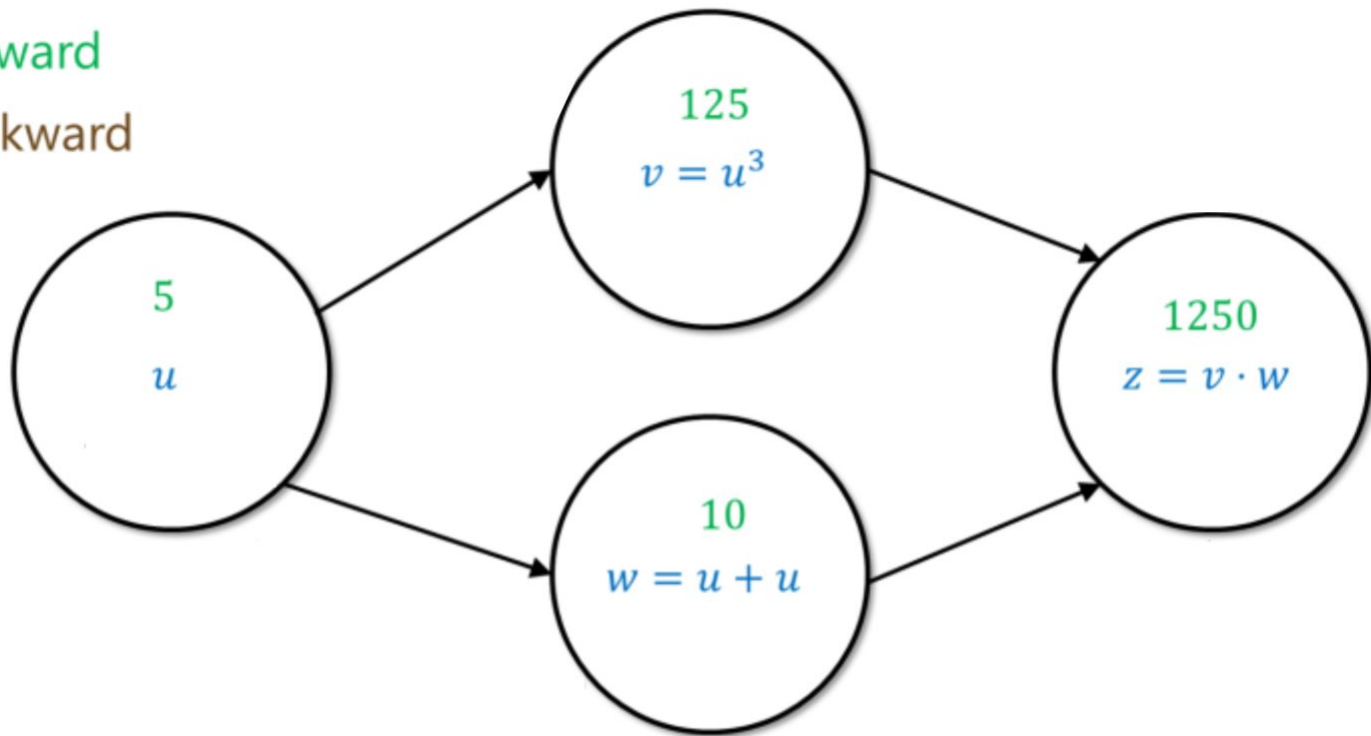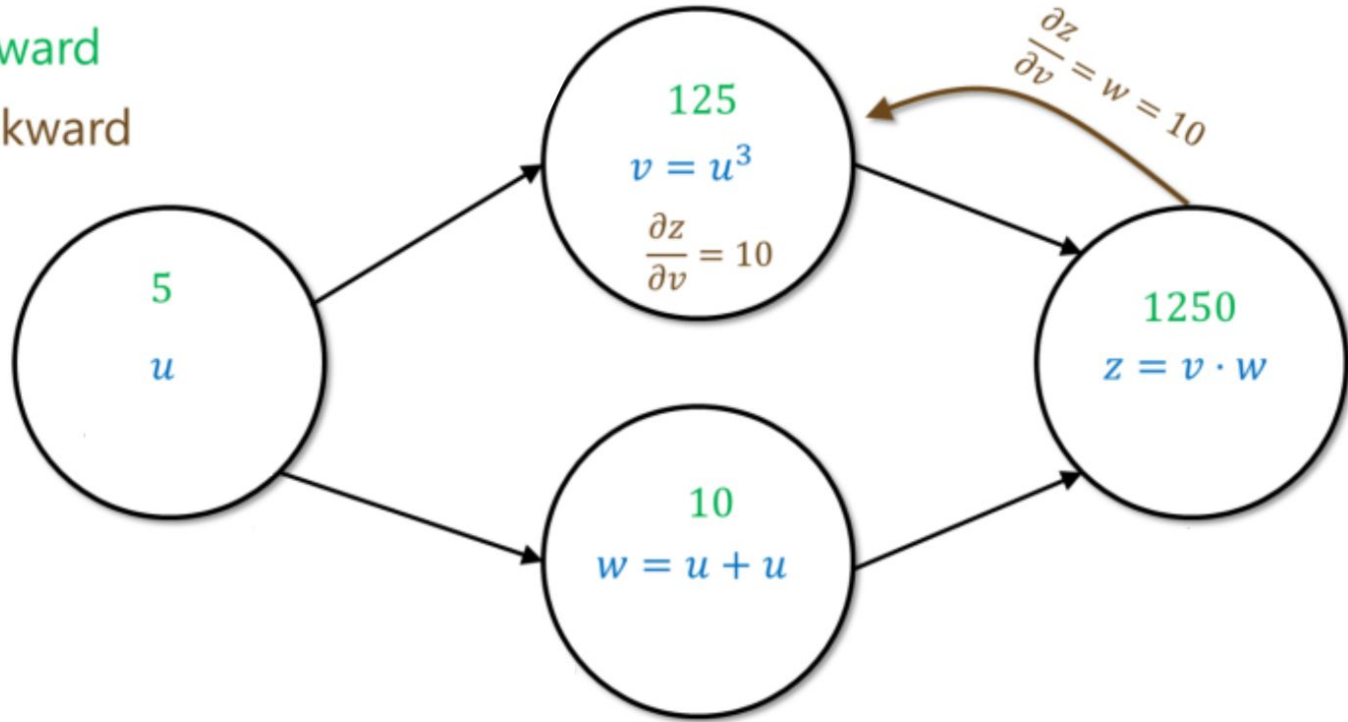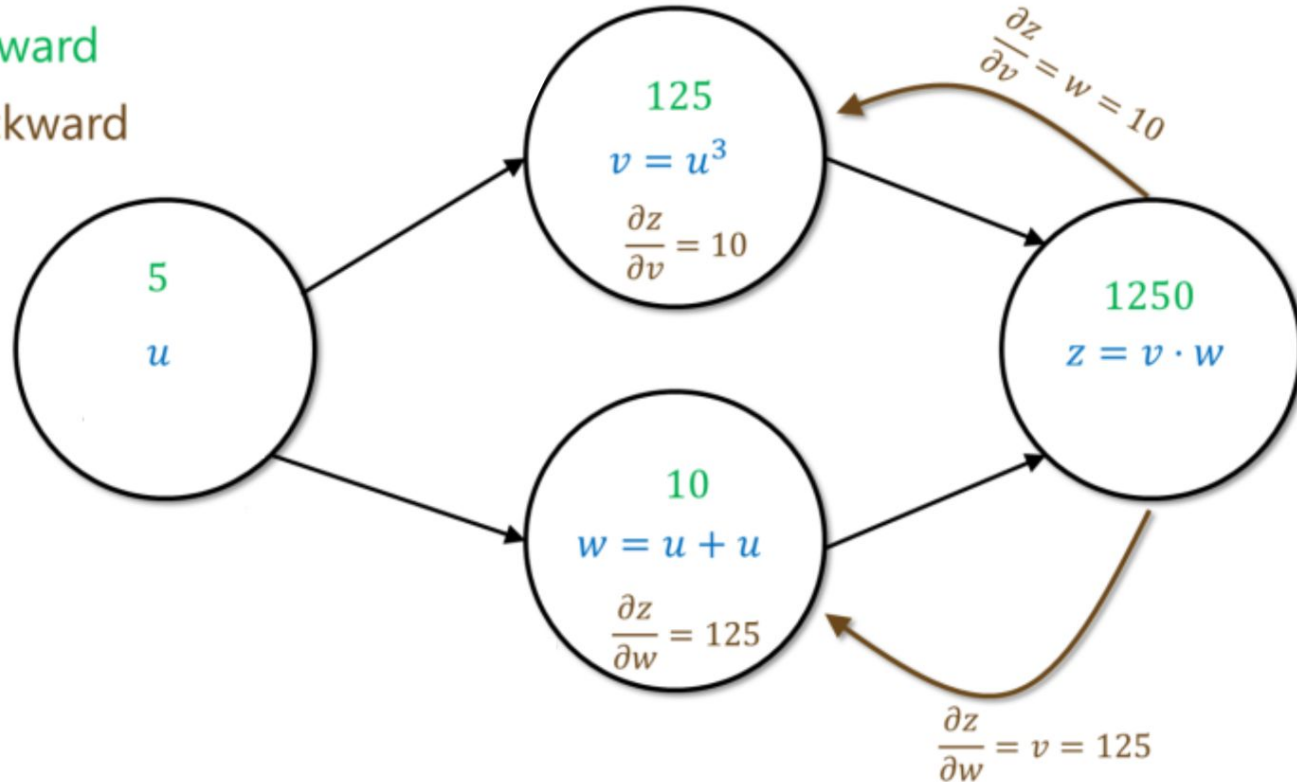
# Forward Pass - MLP

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

**Algorithm** Forward Pass through MLP

1: **Input:** input $\mathbf{x}$, weight matrices $\mathbf{W}^{[1]}, \ldots, \mathbf{W}^{[L]}$, bias vectors $\mathbf{b}^{[1]}, \ldots, \mathbf{b}^{[L]}$
2: $\mathbf{z}^{[0]} = \mathbf{x}$ ▷ Initialize input
3: **for** $l = 1$ **to** $L$ **do**
4: $\quad \mathbf{a}^{[l]} = \mathbf{W}^{[l]}\mathbf{z}^{[l-1]} + \mathbf{b}^{[l]}$ ▷ Linear transformation
5: $\quad \mathbf{z}^{[l]} = \sigma^{[l]}(\mathbf{a}^{[l]})$ ▷ Nonlinear activation
6: **end for**
7: **Output:** $\mathbf{z}^{[L]}$

# Forward Pass - MLP

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$



$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

# Forward Pass - MLP

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]} \mathbf{z}^{[0]}$$



$\sigma$

$\sigma$

$\sigma$

$x_1$

$x_2$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

# Forward Pass - MLP



$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\sigma$$

$$x_1$$

$$x_2$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

Cornell Bowers C·IS

# Forward Pass - MLP

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

# Forward Pass - MLP

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$



$x_1$

$x_2$

$\sigma$

$\mathbf{z}^{[0]} = \mathbf{x}$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

# Forward Pass - MLP

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

# Forward Pass - MLP

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$



$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

# Forward Pass - MLP



$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Forward Pass - MLP

**Algorithm** Forward Pass through MLP

1: **Input:** input $\mathbf{x}$, weight matrices $\mathbf{W}^{[1]}, \ldots, \mathbf{W}^{[L]}$, bias vectors $\mathbf{b}^{[1]}, \ldots, \mathbf{b}^{[L]}$
2: $\mathbf{z}^{[0]} = \mathbf{x}$        ▷ Initialize input
3: **for** $l = 1$ **to** $L$ **do**
4:      $\mathbf{a}^{[l]} = \mathbf{W}^{[l]}\mathbf{z}^{[l-1]} + \mathbf{b}^{[l]}$        ▷ Linear transformation
5:      $\mathbf{z}^{[l]} = \sigma^{[l]}(\mathbf{a}^{[l]})$        ▷ Nonlinear activation
6: **end for**
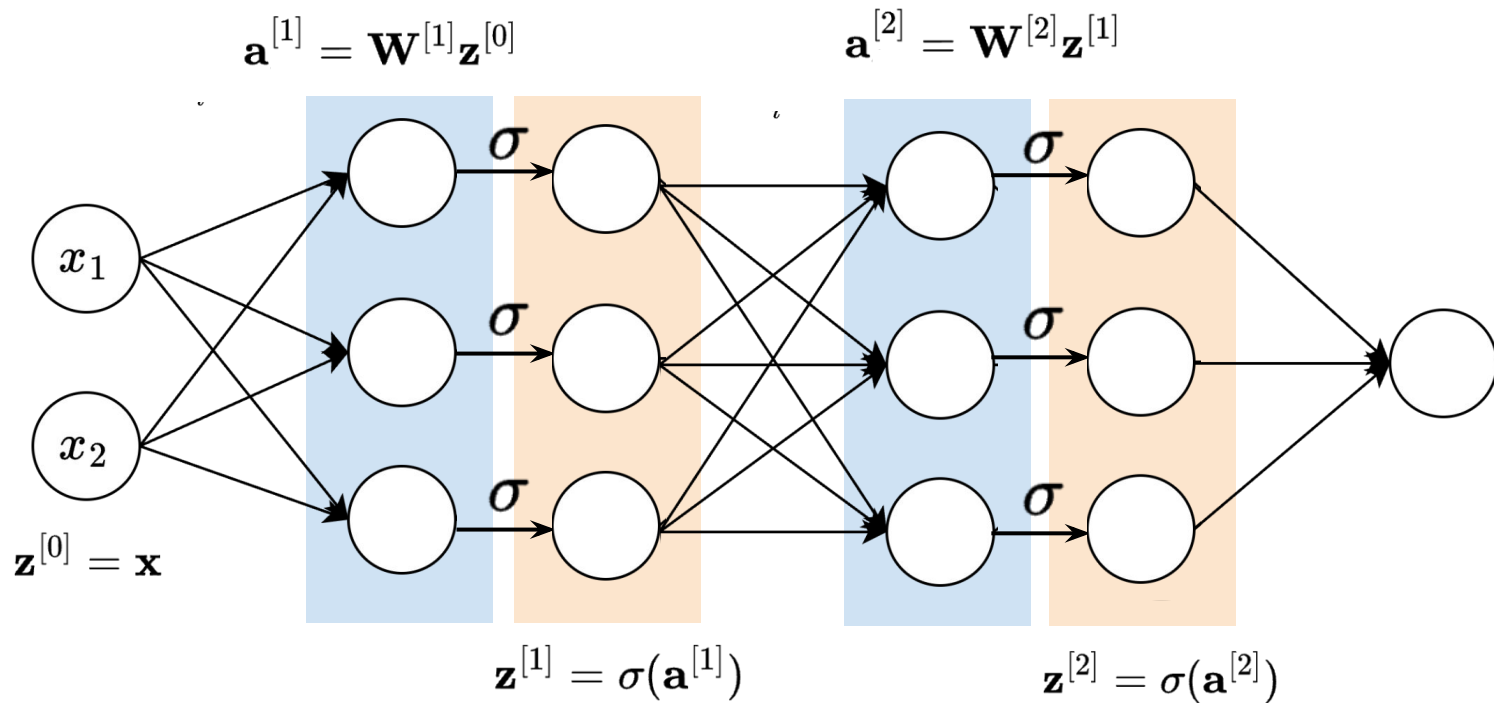7: **Output:** $\mathbf{z}^{[L]}$

# Forward Pass - MLP

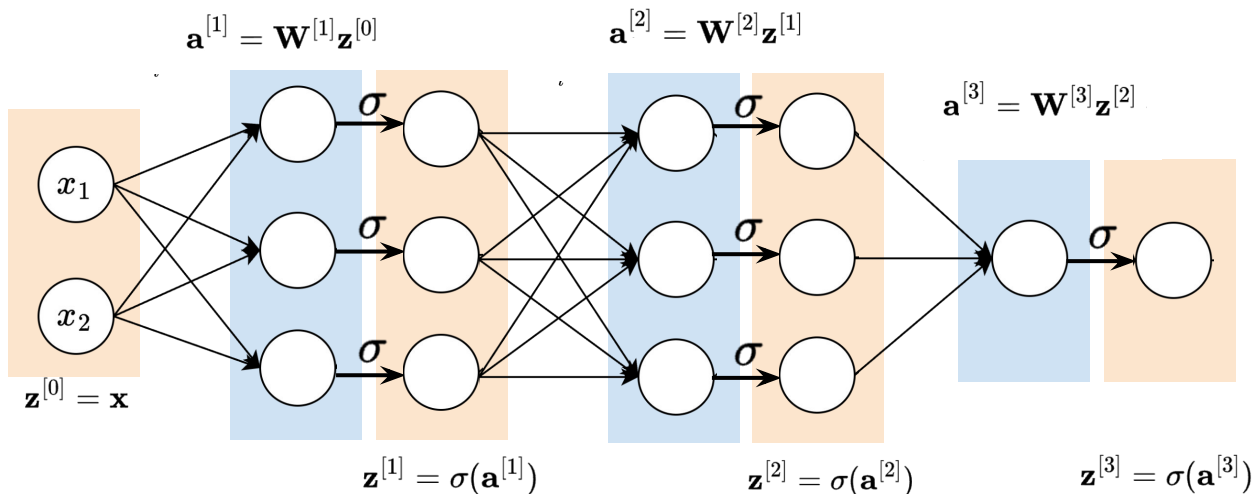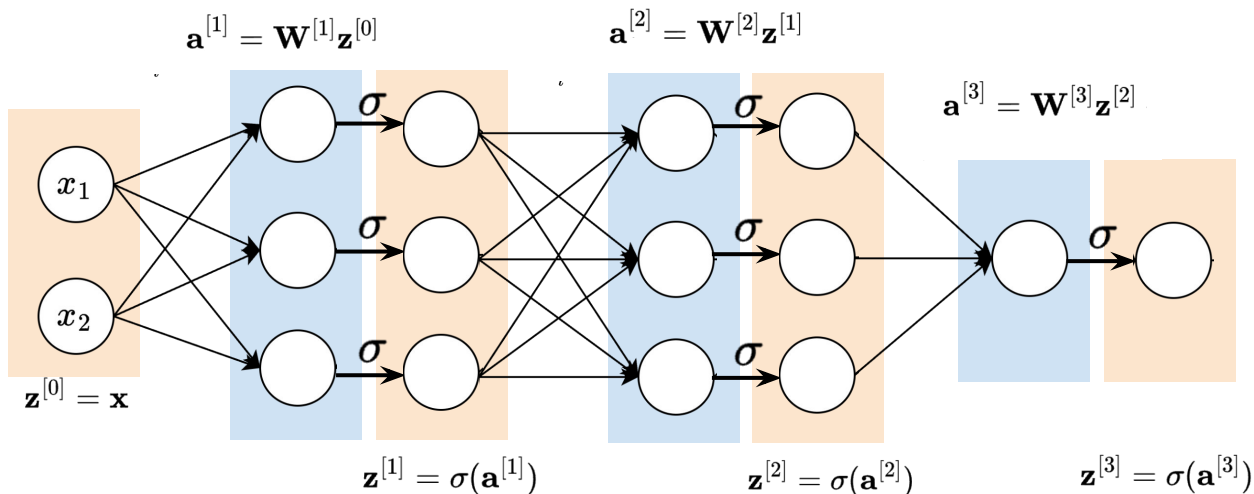**Algorithm** Forward Pass through MLP

1: **Input:** input $\mathbf{x}$, weight matrices $\mathbf{W}^{[1]}, \ldots, \mathbf{W}^{[L]}$, bias vectors $\mathbf{b}^{[1]}, \ldots, \mathbf{b}^{[L]}$

2: $\mathbf{z}^{[0]} = \mathbf{x}$            ▷ Initialize input

3: **for** $l = 1$ **to** $L$ **do**

4:     $\mathbf{a}^{[l]} = \mathbf{W}^{[l]}\mathbf{z}^{[l-1]} + \mathbf{b}^{[l]}$         ▷ Linear transformation

5:     $\mathbf{z}^{[l]} = \sigma^{[l]}(\mathbf{a}^{[l]})$         ▷ Nonlinear activation

6: **end for**

7: **Output:** $\mathbf{z}^{[L]}$



$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]}) \qquad \mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]}) \qquad \mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$



$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$!

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

$x_1$

$x_2$

$\sigma$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$ !

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$



Loss

slope of loss at $w^1$ is negative

one step of gradient descent

$w^1$    $w^{min}$

$0$    (goal)    $w$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$\sigma$    $\sigma$    $\sigma$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

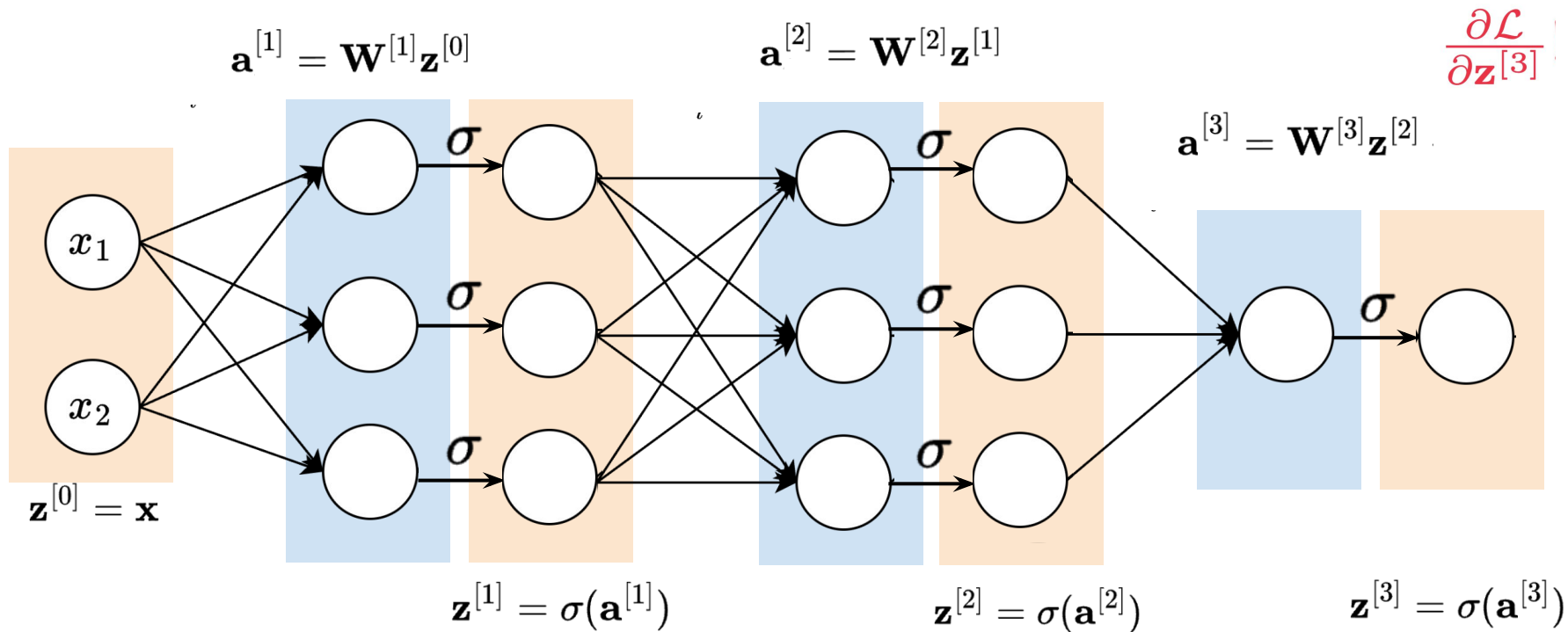$x_1$    $x_2$

# Backprop

Loss $= \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$

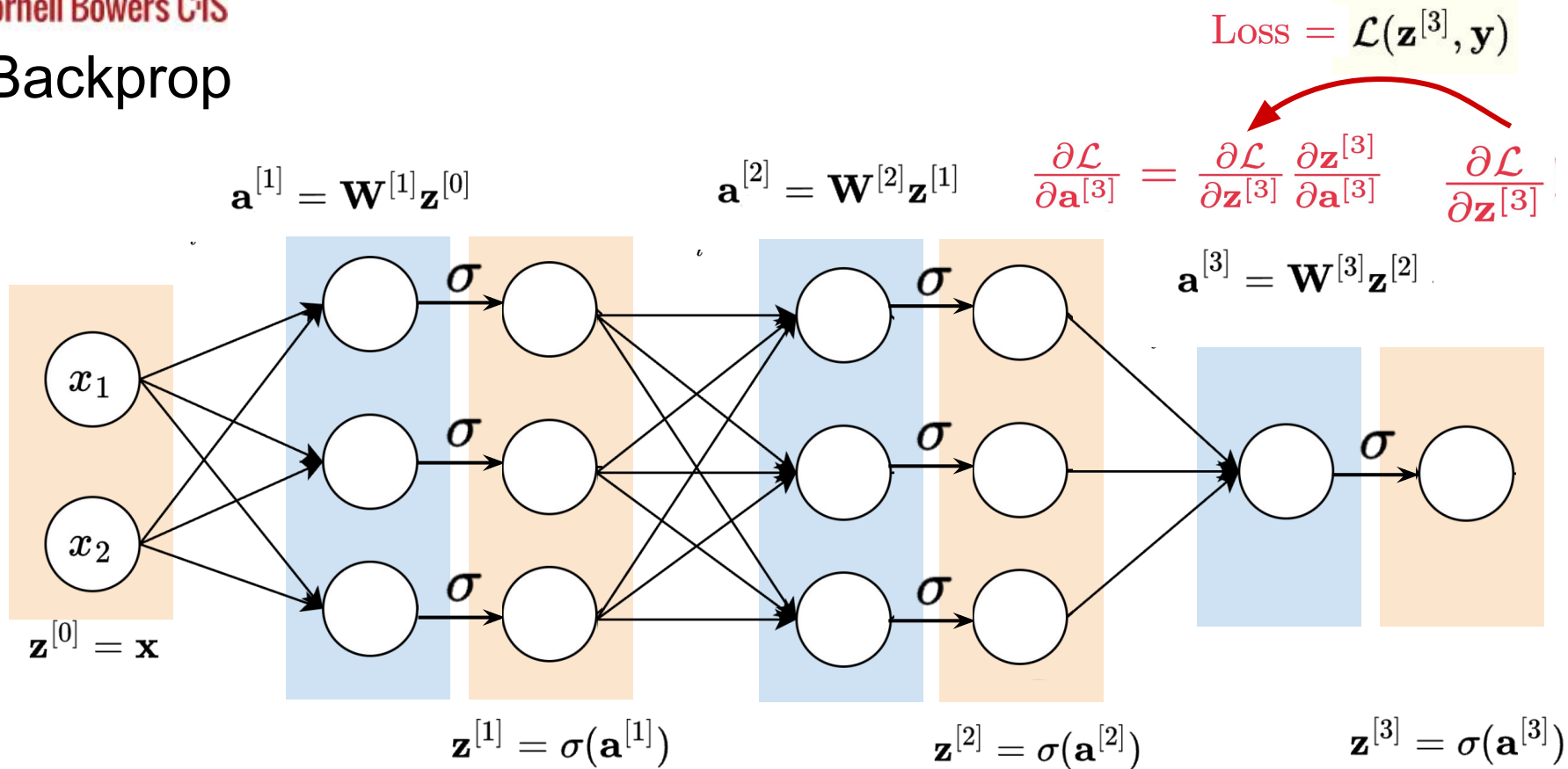We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$!

$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$

$$\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}; \mathcal{D}_{\mathrm{TR}}) = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w^{(0)}}(\mathbf{w}; \mathcal{D}_{\mathrm{TR}}) \\ \frac{\partial \mathcal{L}}{\partial w^{(1)}}(\mathbf{w}; \mathcal{D}_{\mathrm{TR}}) \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial w^{(m)}}(\mathbf{w}; \mathcal{D}_{\mathrm{TR}}) \end{bmatrix}, \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}; \mathcal{D}_{\mathrm{TR}}) \in \mathbb{R}^m$$

$\mathbf{z}^{[0]} = \mathbf{x}$

$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$  $\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$  $\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$
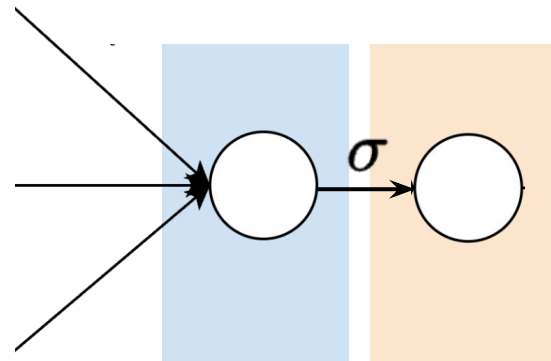
# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$



$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

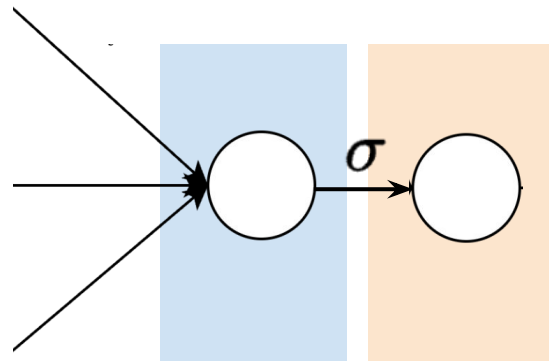$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]} \qquad \mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]} \qquad \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \frac{\partial \mathbf{z}^{[3]}}{\partial \mathbf{a}^{[3]}} \qquad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$



$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]}) \qquad \mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]}) \qquad \mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \frac{\partial \mathbf{z}^{[3]}}{\partial \mathbf{a}^{[3]}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \odot \sigma^{[3]'}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]} \mathbf{z}^{[2]}$$



$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\boxed{\delta^{[3]} =} \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \frac{\partial \mathbf{z}^{[3]}}{\partial \mathbf{a}^{[3]}}$$
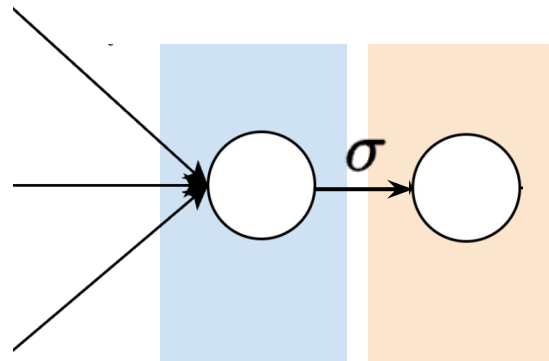
$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \odot \sigma^{[3]'}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]} \mathbf{z}^{[2]}$$



$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\boxed{\delta^{[3]} =} \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \frac{\partial \mathbf{z}^{[3]}}{\partial \mathbf{a}^{[3]}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \odot \sigma^{[3]'}$$

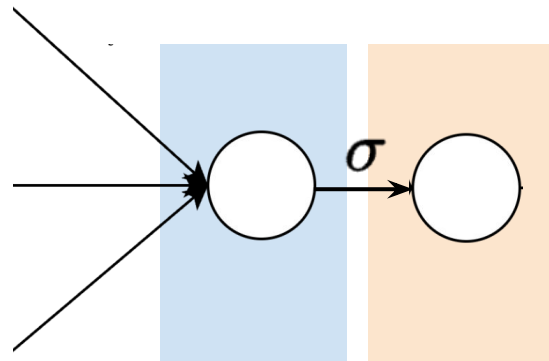$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]} \mathbf{z}^{[2]}$$



$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{W}^{[3]}}$$

$$= \delta^{[3]} (\mathbf{z}^{[2]})^T$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

For propagation to next layer:

$$\boxed{\delta^{[3]} = } \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \frac{\partial \mathbf{z}^{[3]}}{\partial \mathbf{a}^{[3]}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \odot \sigma^{[3]\prime}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]} \mathbf{z}^{[2]}$$

For weight updates:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{W}^{[3]}}$$

$$= \delta^{[3]} (\mathbf{z}^{[2]})^T$$

$$\sigma$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

For propagation to next layer:

$$\boxed{\delta^{[3]} =} \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \frac{\partial \mathbf{z}^{[3]}}{\partial \mathbf{a}^{[3]}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}} \odot \sigma^{[3]\prime}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$

For weight updates:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[3]}} =$$

$$= \delta$$

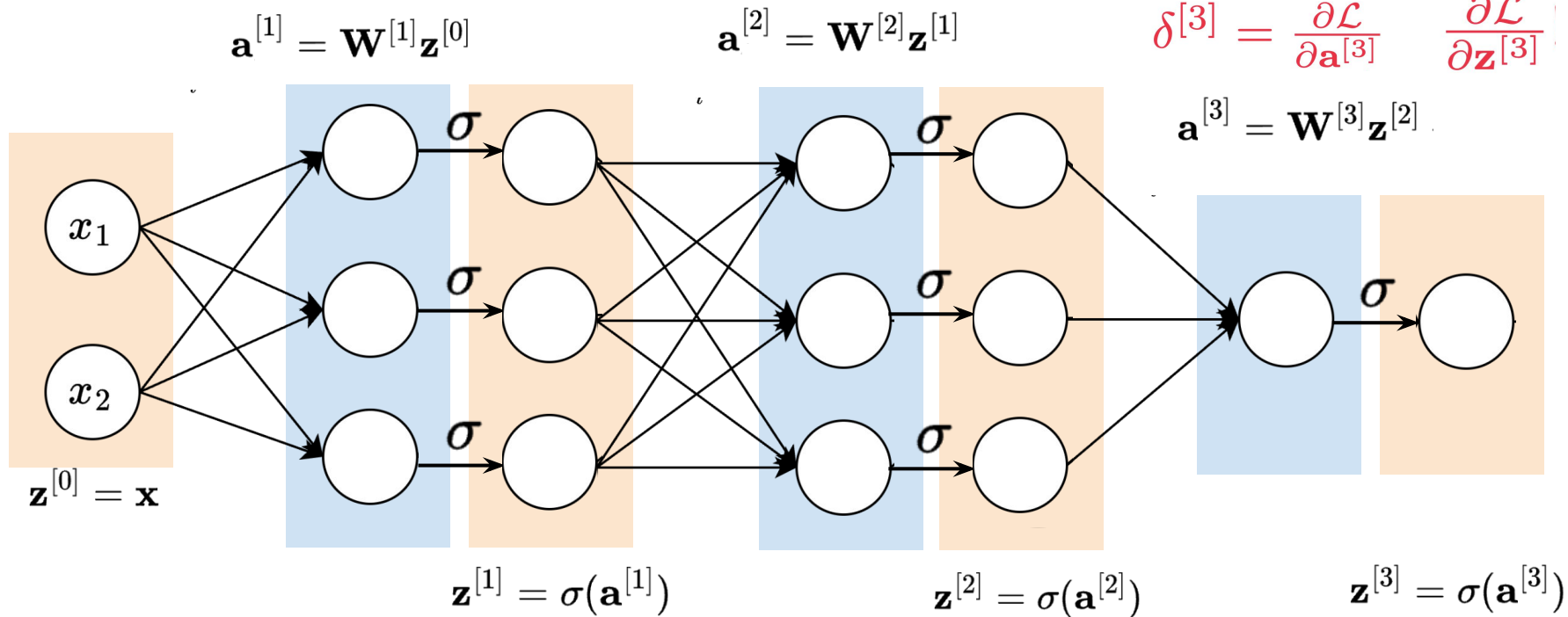### Backpropagation- Key Idea

If you know $\frac{\partial z}{\partial v_1}, \frac{\partial z}{\partial v_2}, \frac{\partial z}{\partial v_3}$

You can compute $\frac{\partial z}{\partial u}$

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v_1}{\partial u} \cdot \frac{\partial z}{\partial v_1} + \frac{\partial v_2}{\partial u} \cdot \frac{\partial z}{\partial v_2} + \frac{\partial v_3}{\partial u} \cdot \frac{\partial z}{\partial v_3} \right)$$

https://windowsontheory.org/2020/11/03/yet-another-backpropagation-tutorial

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]} \qquad \mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \qquad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$



$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]}) \qquad \mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]}) \qquad \mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

# Backprop

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]} \qquad \mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$$

$$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \qquad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$$



$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]}) \qquad \mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]}) \qquad \mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{z}^{[2]}} = (W^{[3]})^T \delta^{[3]}$$

**Cornell Bowers C·IS**

Backprop

$$\delta^{[2]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} \frac{\partial \mathbf{z}^{[2]}}{\partial \mathbf{a}^{[2]}}$$

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]} \mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]} \mathbf{z}^{[1]}$$

$$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \qquad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]} \mathbf{z}^{[2]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{z}^{[2]}} = (W^{[3]})^T \delta^{[3]}$$

Cornell Bowers C·IS

$$\delta^{[1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[1]}} \frac{\partial \mathbf{z}^{[1]}}{\partial \mathbf{a}^{[1]}}$$

$$\delta^{[2]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} \frac{\partial \mathbf{z}^{[2]}}{\partial \mathbf{a}^{[2]}}$$

$$\text{Loss} = \mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

$$\mathbf{a}^{[1]} = \mathbf{W}^{[1]} \mathbf{z}^{[0]}$$

$$\mathbf{a}^{[2]} = \mathbf{W}^{[2]} \mathbf{z}^{[1]}$$

$$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \qquad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$\mathbf{a}^{[3]} = \mathbf{W}^{[3]} \mathbf{z}^{[2]}$$

$$\mathbf{z}^{[0]} = \mathbf{x}$$

$$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$$

$$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$$

$$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[2]}} \frac{\partial \mathbf{a}^{[2]}}{\partial \mathbf{z}^{[1]}} = (W^{[2]})^T \delta^{[2]}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{z}^{[2]}} = (W^{[3]})^T \delta^{[3]}$$

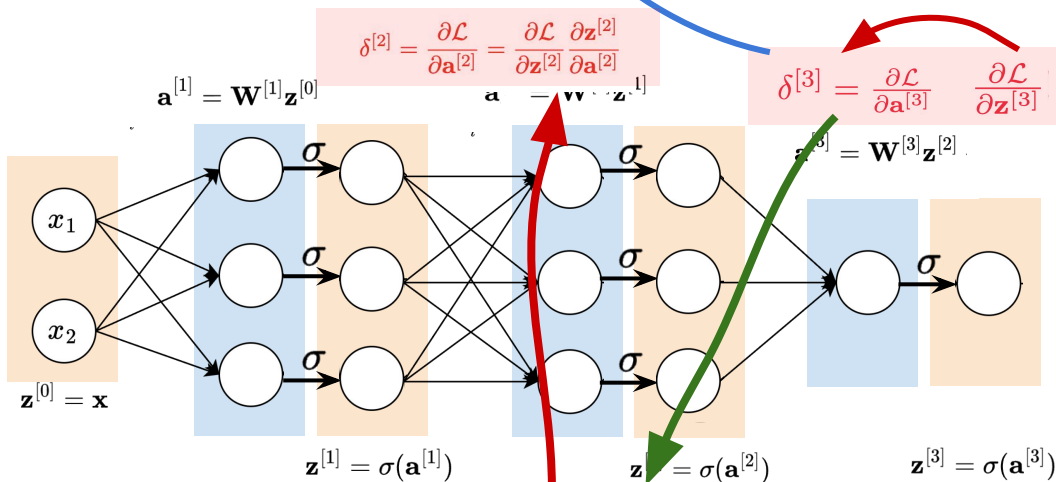# Backpropagation

**Algorithm** Backward Pass through MLP (Detailed)

1: **Input:** $\{\mathbf{z}^{[1]}, \ldots, \mathbf{z}^{[L]}\}$, $\{\mathbf{a}^{[1]}, \ldots, \mathbf{a}^{[L]}\}$, loss gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}}$

2: $\delta^{[L]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \frac{\partial \mathbf{z}^{[L]}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \odot \sigma^{[L]'}(\mathbf{a}^{[L]})$ ▷ Error term

3: **for** $l = L$ to 1 **do**

4:     $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]}(\mathbf{z}^{[l-1]})^T$ ▷ Gradient of weights

5:     $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]}$ ▷ Gradient of biases

6:     $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{z}^{[l-1]}} = (\mathbf{W}^{[l]})^T \delta^{[l]}$

7:     $\delta^{[l-1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} \frac{\partial \mathbf{z}^{[l-1]}}{\partial \mathbf{a}^{[l-1]}} = ((\mathbf{W}^{[l]})^T \delta^{[l]}) \odot \sigma^{[l-1]'}(\mathbf{a}^{[l-1]})$

8: **end for**

9: **Output:** $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1:L]}}, \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[1:L]}}$



$$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$$

$$\mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$$

We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$ !

$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$

$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$

$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$

$\mathbf{z}^{[0]} = \mathbf{x}$

$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$

$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$

$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$

# Backpropagation

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{W}^{[3]}}$$

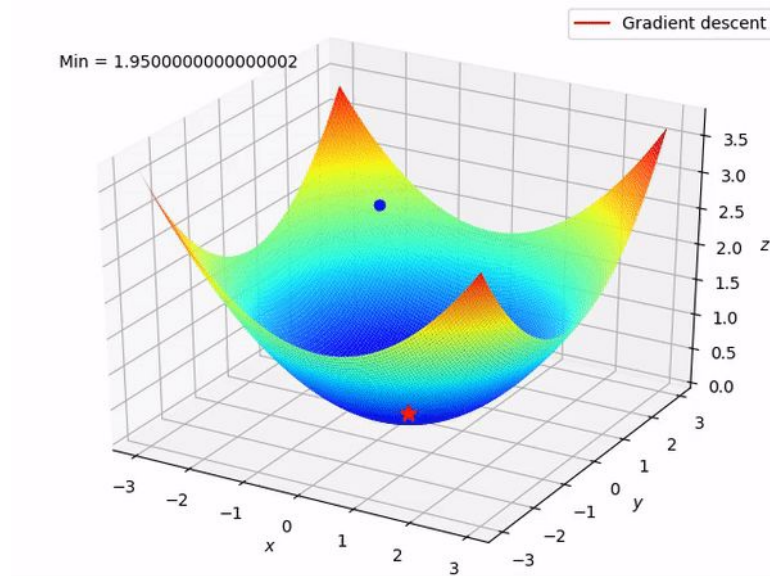$$= \delta^{[3]} (\mathbf{z}^{[2]})^T$$

**Algorithm** Backward Pass through MLP (Detailed)

1: **Input:** $\{\mathbf{z}^{[1]}, \ldots, \mathbf{z}^{[L]}\}$, $\{\mathbf{a}^{[1]}, \ldots, \mathbf{a}^{[L]}\}$, loss gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}}$

2: $\delta^{[L]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \frac{\partial \mathbf{z}^{[L]}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \odot \sigma^{[L]'}(\mathbf{a}^{[L]})$ ▷ Error term

3: **for** $l = L$ **to** 1 **do**

4: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]}(\mathbf{z}^{[l-1]})^T$ ▷ Gradient of weights

5: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]}$ ▷ Gradient of biases

6: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{z}^{[l-1]}} = (\mathbf{W}^{[l]})^T \delta^{[l]}$

7: $\quad \delta^{[l-1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} \frac{\partial \mathbf{z}^{[l-1]}}{\partial \mathbf{a}^{[l-1]}} = ((\mathbf{W}^{[l]})^T \delta^{[l]}) \odot \sigma^{[l-1]'}(\mathbf{a}^{[l-1]})$

8: **end for**

9: **Output:** $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1:L]}}, \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[1:L]}}$

$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$

$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$

$\mathbf{a}^{[2]} = \mathbf{W}^{[2]}\mathbf{z}^{[1]}$

$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$

$\mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$

We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$!



$\mathbf{z}^{[0]} = \mathbf{x}$

$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$

$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$

$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$

# Backpropagation

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{W}^{[3]}}$$

$$= \delta^{[3]} (\mathbf{z}^{[2]})^T$$

**Algorithm** Backward Pass through MLP (Detailed)

1: **Input:** $\{\mathbf{z}^{[1]}, \ldots, \mathbf{z}^{[L]}\}$, $\{\mathbf{a}^{[1]}, \ldots, \mathbf{a}^{[L]}\}$, loss gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}}$

2: $\delta^{[L]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \frac{\partial \mathbf{z}^{[L]}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \odot \sigma^{[L]'}(\mathbf{a}^{[L]})$ ▷ Error term

3: **for** $l = L$ to 1 **do**

4: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]}(\mathbf{z}^{[l-1]})^T$ ▷ Gradient of weights

5: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]}$ ▷ Gradient of biases

6: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{z}^{[l-1]}} = (\mathbf{W}^{[l]})^T \delta^{[l]}$

7: $\quad \delta^{[l-1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} \frac{\partial \mathbf{z}^{[l-1]}}{\partial \mathbf{a}^{[l-1]}} = ((\mathbf{W}^{[l]})^T \delta^{[l]}) \odot \sigma^{[l-1]'}(\mathbf{a}^{[l-1]})$

8: **end for**

9: **Output:** $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1:L]}}$, $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[1:L]}}$

$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$

$\mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$

We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$!



$\mathbf{a}^{[1]} = \mathbf{W}^{[1]} \mathbf{z}^{[0]}$

$\mathbf{a}^{[2]} = \mathbf{W}^{[2]} \mathbf{z}^{[1]}$

$\mathbf{a}^{[3]} = \mathbf{W}^{[3]} \mathbf{z}^{[2]}$

$\mathbf{z}^{[0]} = \mathbf{x}$

$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$

$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$

$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$

$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{z}^{[2]}} = (W^{[3]})^T \delta^{[3]}$

# Backpropagation

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{W}^{[3]}}$$

$$= \delta^{[3]} (\mathbf{z}^{[2]})^T$$

**Algorithm** Backward Pass through MLP (Detailed)

1: **Input:** $\{\mathbf{z}^{[1]}, \dots, \mathbf{z}^{[L]}\}$, $\{\mathbf{a}^{[1]}, \dots, \mathbf{a}^{[L]}\}$, loss gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}}$

2: $\delta^{[L]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \frac{\partial \mathbf{z}^{[L]}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \odot \sigma^{[L]'}(\mathbf{a}^{[L]})$ ▷ Error term

3: **for** $l = L$ **to** 1 **do**

4: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]} (\mathbf{z}^{[l-1]})^T$ ▷ Gradient of weights

5: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]}$ ▷ Gradient of biases

6: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{z}^{[l-1]}} = (\mathbf{W}^{[l]})^T \delta^{[l]}$

7: $\quad \delta^{[l-1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} \frac{\partial \mathbf{z}^{[l-1]}}{\partial \mathbf{a}^{[l-1]}} = ((\mathbf{W}^{[l]})^T \delta^{[l]}) \odot \sigma^{[l-1]'}(\mathbf{a}^{[l-1]})$

8: **end for**

9: **Output:** $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1:L]}}, \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[1:L]}}$

$\delta^{[2]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} \frac{\partial \mathbf{z}^{[2]}}{\partial \mathbf{a}^{[2]}}$

$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$

$\mathbf{a}^{[1]} = \mathbf{W}^{[1]} \mathbf{z}^{[0]}$

$\mathbf{z}^{[0]} = \mathbf{x}$

$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$

$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$

$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$

$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{z}^{[2]}} = (W^{[3]})^T \delta^{[3]}$

$\mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$

We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$!

# What is Optimization?



In deep learning, optimization methods attempt to find model weights that **minimize the loss function**.

# Loss function

Empirical Risk:

$$\mathcal{L}(\mathbf{w}_t) = \frac{1}{n} \sum_{i=1,\ldots,n} \ell(\mathbf{w}_t, \mathbf{x}_i)$$

$t$ : at time step $t$
$\mathbf{w}_t$: Model weights (parameters) at time $t$
$\mathbf{x}_i$: The i-th input training data

$\mathcal{L}$ : the Loss function (optimization target)
$\ell$ : per-sample loss



Local Minima

Saddle Point

Global Minima

# Gradient Descent (GD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

$\alpha$: the learning rate

$\nabla \mathcal{L}(\mathbf{w}_t)$: the gradient of Loss w.r.t. $\mathbf{w}_t$

# Demo

Gradient descent with global minimum

What are some potential problems with gradient descent?

# Convexity

- A function on a graph is **convex** if a line segment drawn through any two points on the line of the function, then it never lies below the curved line segment
- Convexity implies that every local minimum is **global minimum**.
- Neural networks are **not** convex!

Not convex

Convex

# Challenges in Non-Convex Optimization



Local Minima vs. Global Minima

Saddle Points

Vanishing gradient

# Demo

Gradient descent with local minimum

# Gradient Descent (GD)

$$\mathcal{L}(\mathbf{w}_t) = \boxed{\frac{1}{n} \sum_{i=1}^{n}} \ell(\mathbf{w}_t, \mathbf{x}_i)$$

$$\nabla\mathcal{L}(\mathbf{w}_t) = \boxed{\frac{1}{n} \sum_{i=1}^{n}} \nabla\ell(\mathbf{w}_t, \mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha\nabla\mathcal{L}(\mathbf{w}_t)$$

Full gradient: $\mathcal{O}(n)$ time => **Too expensive!**
- *Statistically, why don't we use 1 or a few samples from the training dataset to approximate the full gradient?*

# Gradient Descent (GD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

# Gradient Descent (GD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

Select **1** example randomly each time

# Gradient Descent (GD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

Select **1** example randomly each time

*Per-sample gradient is equivalent to full gradient in expectation!*

$$\mathbb{E}[\nabla \ell(\mathbf{w}_t, \mathbf{x}_i)] = \frac{1}{n} \sum_{i=1}^{n} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i) = \nabla \mathcal{L}(\mathbf{w}_t)$$

# Stochastic Gradient Descent (SGD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

Select **1** example randomly each time

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

*Per-sample gradient is equivalent to full gradient in expectation!*

$$\mathbb{E}[\nabla \ell(\mathbf{w}_t, \mathbf{x}_i)] = \frac{1}{n} \sum_{i=1}^{n} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i) = \nabla \mathcal{L}(\mathbf{w}_t)$$

# Cornell Bowers C·IS

## Stochastic Gradient Descent (SGD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

Select **1** example randomly each time

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

*Trade off convergence!*

*Per-sample gradients not necessarily points to the local minimum, introducing a **noise ball**...*

# Stochastic Gradient Descent (SGD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

Select **1** example randomly each time

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

Select a batch $\mathcal{B}_t$ of examples randomly each time, with *batch size $b$*

# Minibatch SGD

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

Select **1** example randomly each time

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

Select a batch $\mathcal{B}_t$ of examples randomly each time, with *batch size* $b$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{1}{b} \sum_{i \in \mathcal{B}_t} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

*Let's look at an example!*

Cornell Bowers C·IS

Local Minimum

Draw the gradients:
- Smaller learning rate
- Larger learning rate

Local Minimum

Local Minimum

Local Minimum

Cornell Bowers C·IS
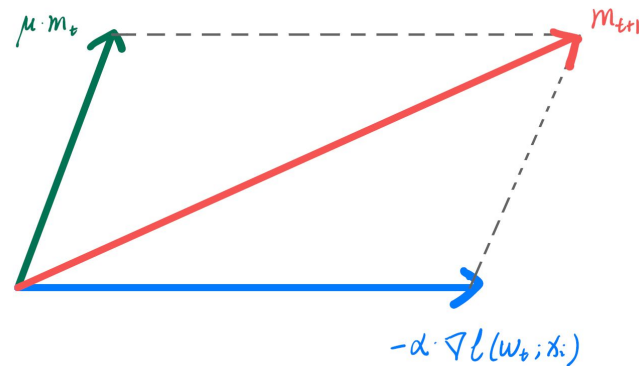
Local Minimum
Minibatch SGD
Momentum

# SGD with Momentum (Polyak, 1964)

*Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.*

# SGD with Momentum (Polyak, 1964)

*Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.*

# SGD with Momentum (Polyak, 1964)

*Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.*
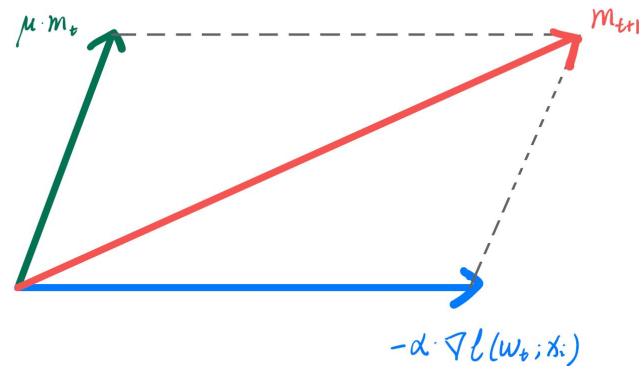
**SGD Update Rule**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

$$\longrightarrow$$

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

where $\mu \in [0, 1]$ is the momentum coefficient.

# SGD with Momentum (Polyak, 1964)

*Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.*



$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

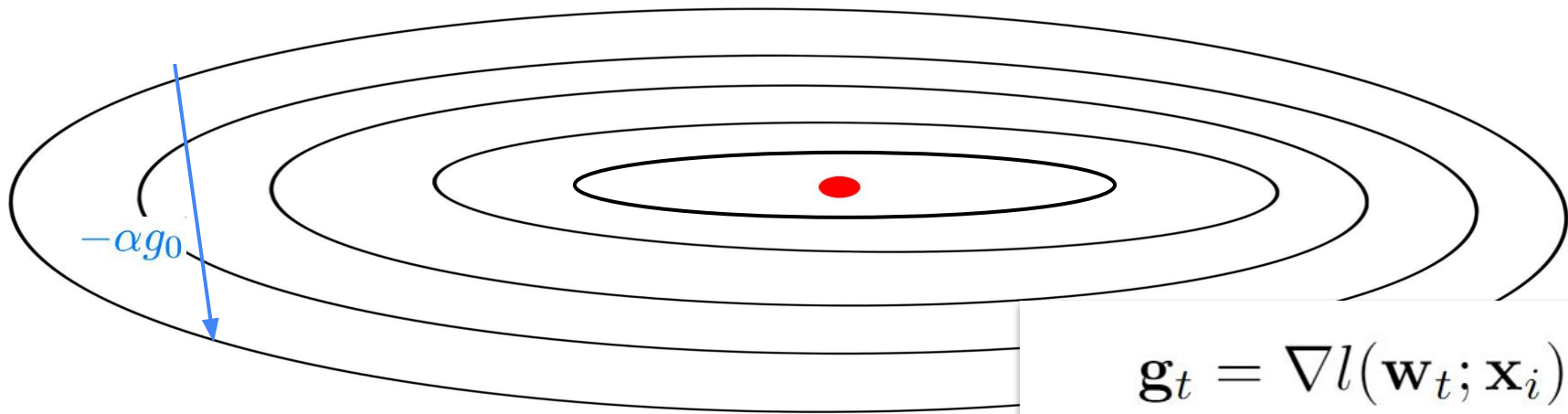where $\mu \in [0, 1]$ is the momentum coefficient.

# SGD with Momentum

*Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.*

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

# SGD with Momentum

Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.
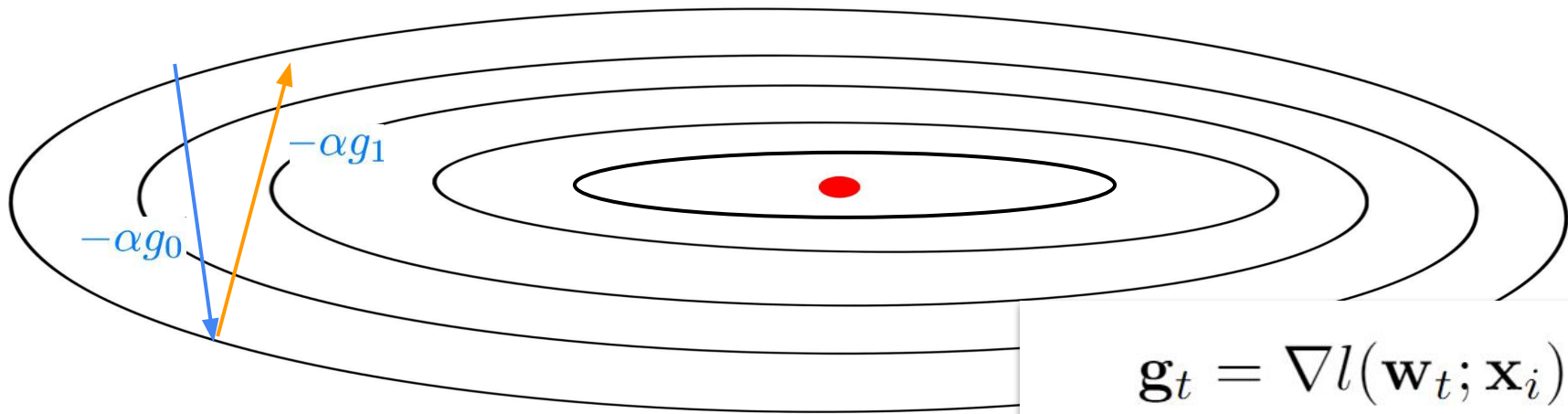
$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

# SGD with Momentum

*Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.*

$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$
$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$\mu \cdot m_t$

$m_{t+1}$

$-\alpha \cdot \nabla l(w_t; x_i)$

# SGD with Momentum

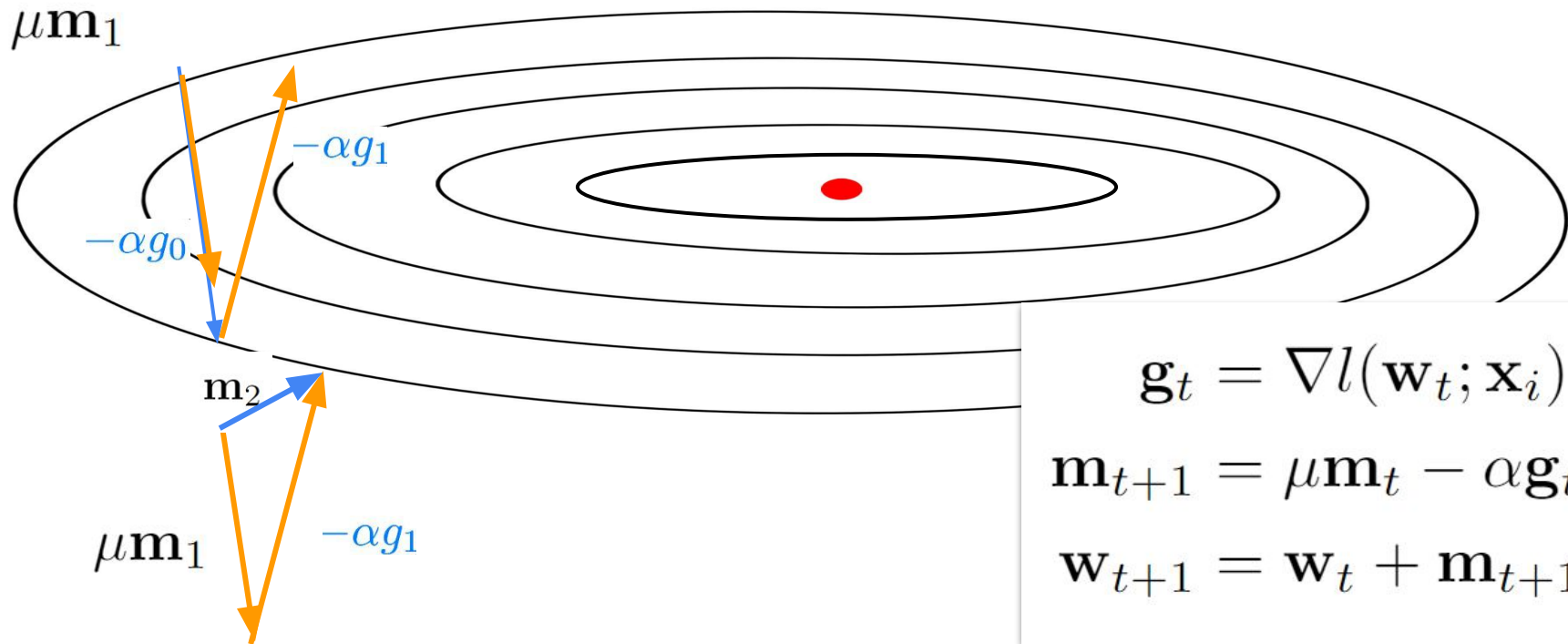*Compute an **Exponentially Weighted Moving Average (EWMA)** of the gradients as **momentum** and use that to update the weight instead.*

$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$= \mathbf{w}_t + \mu(\mu \mathbf{m}_{t-1} - \alpha \mathbf{g}_{t-1}) - \alpha \mathbf{g}_t$$

# SGD with Momentum

*Compute an **Exponentially Weighted Moving Average** (EWMA) of the gradients as **momentum** and use that to update the weight instead.*
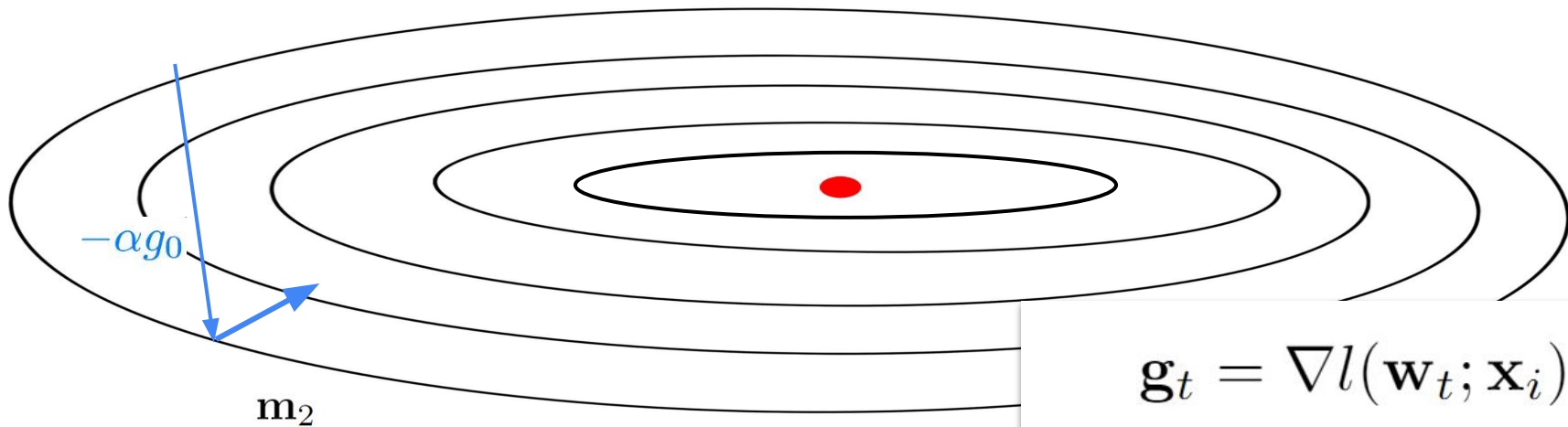
$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$= \mathbf{w}_t + \mu(\mu \mathbf{m}_{t-1} - \alpha \mathbf{g}_{t-1}) - \alpha \mathbf{g}_t$$

$$= \mathbf{w}_t + \mu(\mu(\mu \mathbf{m}_{t-2} - \alpha \mathbf{g}_{t-2}) - \alpha \mathbf{g}_{t-1}) - \alpha \mathbf{g}_t$$

$$= \mathbf{w}_t - \alpha \mathbf{g}_t - \mu \alpha \mathbf{g}_{t-1} - \mu^2 \alpha \mathbf{g}_{t-2} - \mu^3 \alpha \mathbf{g}_{t-3} - \dots$$

# SGD with Momentum

*Compute an **Exponentially Weighted Moving Average** (EWMA) of the gradients as **momentum** and use that to update the weight instead.*

$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$
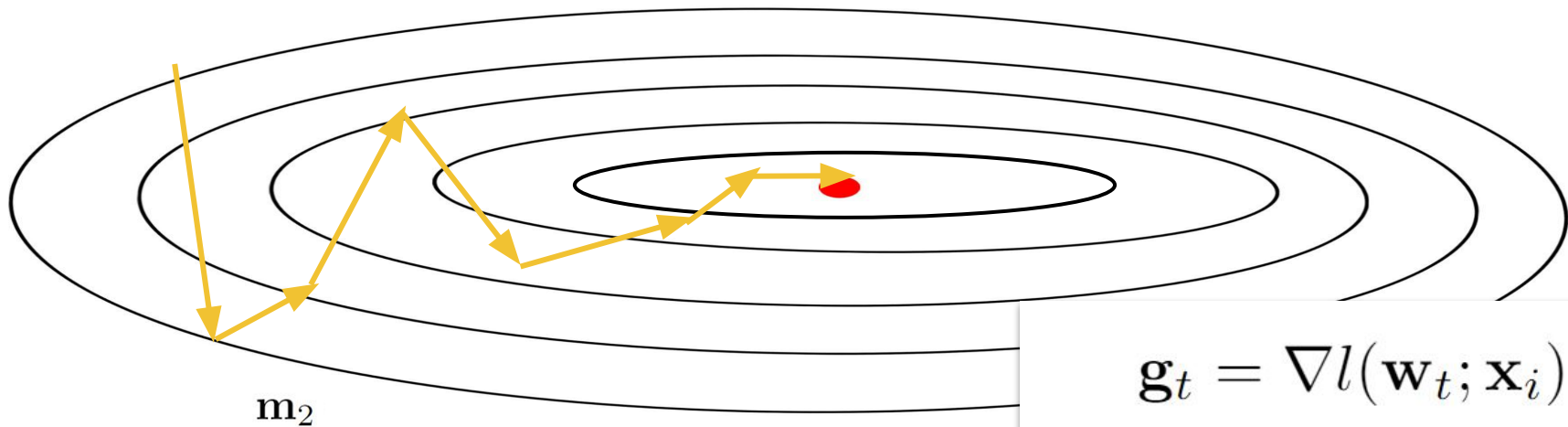
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$= \mathbf{w}_t + \mu(\mu \mathbf{m}_{t-1} - \alpha \mathbf{g}_{t-1}) - \alpha \mathbf{g}_t$$

$$= \mathbf{w}_t + \mu(\mu(\mu \mathbf{m}_{t-2} - \alpha \mathbf{g}_{t-2}) - \alpha \mathbf{g}_{t-1}) - \alpha \mathbf{g}_t$$

$$= \mathbf{w}_t - \alpha \mathbf{g}_t - \mu \alpha \mathbf{g}_{t-1} - \mu^2 \alpha \mathbf{g}_{t-2} - \mu^3 \alpha \mathbf{g}_{t-3} - \ldots$$

$$= \mathbf{w}_t - \alpha \sum_{i=0}^{t} \mu^i \mathbf{g}_{t-i}$$

Local Minimum
Minibatch SGD

$-\alpha g_0$

$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

$$\mathbf{m}_2 = \mu \mathbf{m}_1 - \alpha \mathbf{g}_1$$

$-\alpha g_0$

$\mathbf{m}_2$

$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

Local Minimum
Minibatch SGD
Momentum

$$\mathbf{m}_2 = \mu \mathbf{m}_1 - \alpha \mathbf{g}_1$$



$\mathbf{m}_2$

$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{m}_{t+1} = \mu \mathbf{m}_t - \alpha \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

$$\mathbf{m}_2 = \mu\mathbf{m}_1 - \alpha\mathbf{g}_1$$

Local Minimum
Minibatch SGD
Momentum

*Momentum converges almost always faster than standard SGD!*

$\mathbf{m}_2$

$$\mathbf{g}_t = \nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

# Quick Recap

**Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{1}{n} \sum_{i=1}^{n} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**Stochastic Gradient Descent**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$

**Minibatch SGD**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{1}{b} \sum_{i \in \mathcal{B}_t} \nabla \ell(\mathbf{w}_t, \mathbf{x}_i)$$
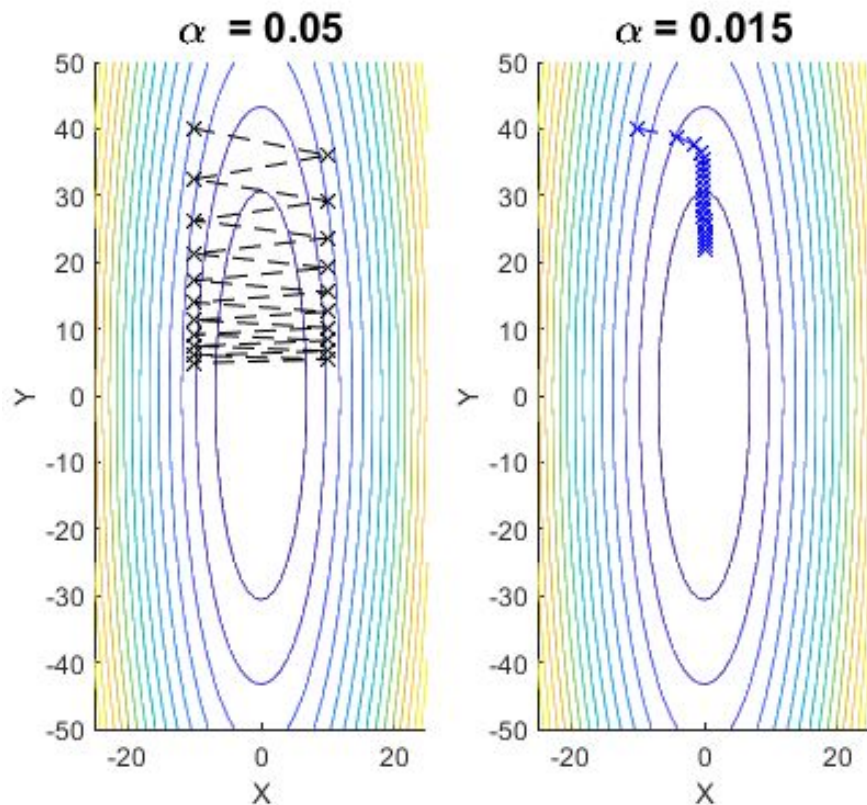
**SGD w. Momentum**

$$m_{t+1} = \mu m_t - \alpha \nabla l(w_t; x_i)$$
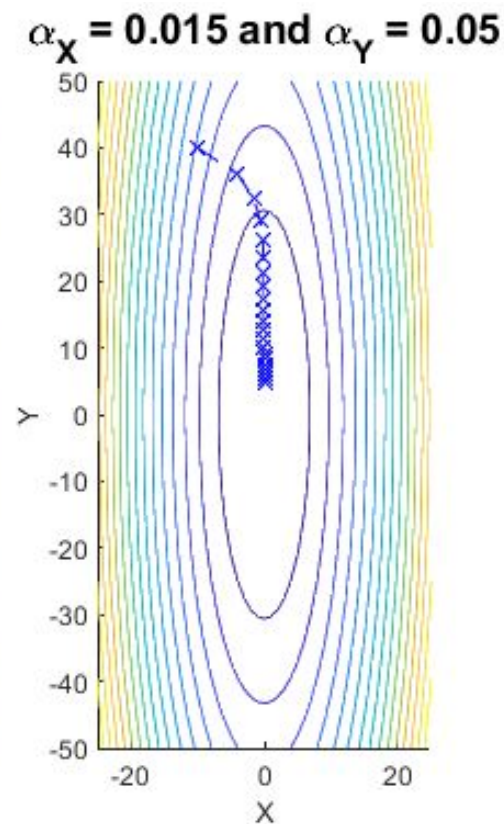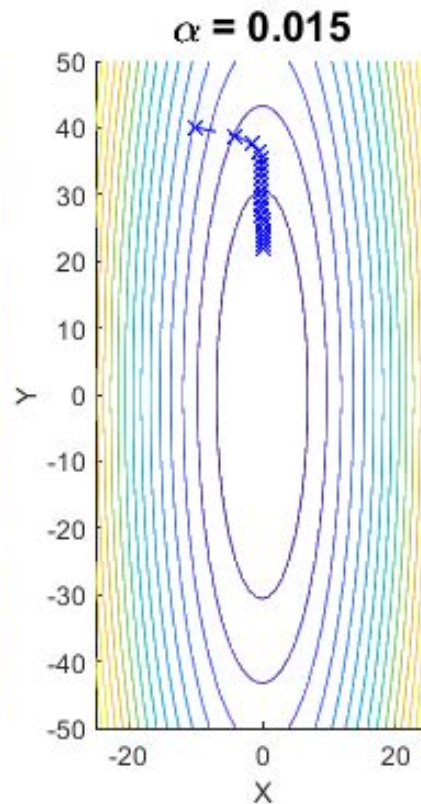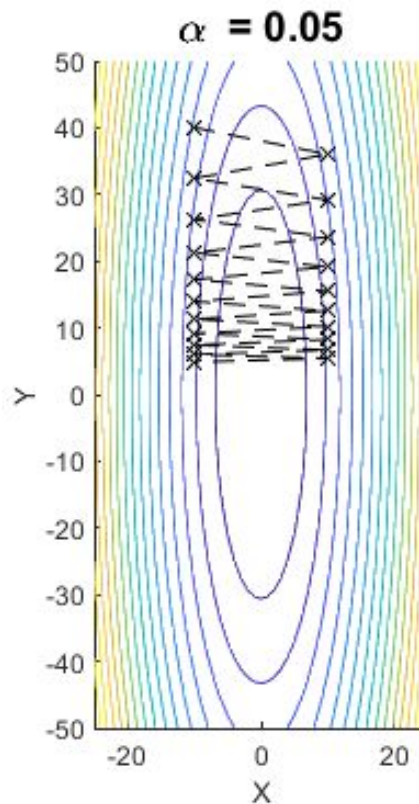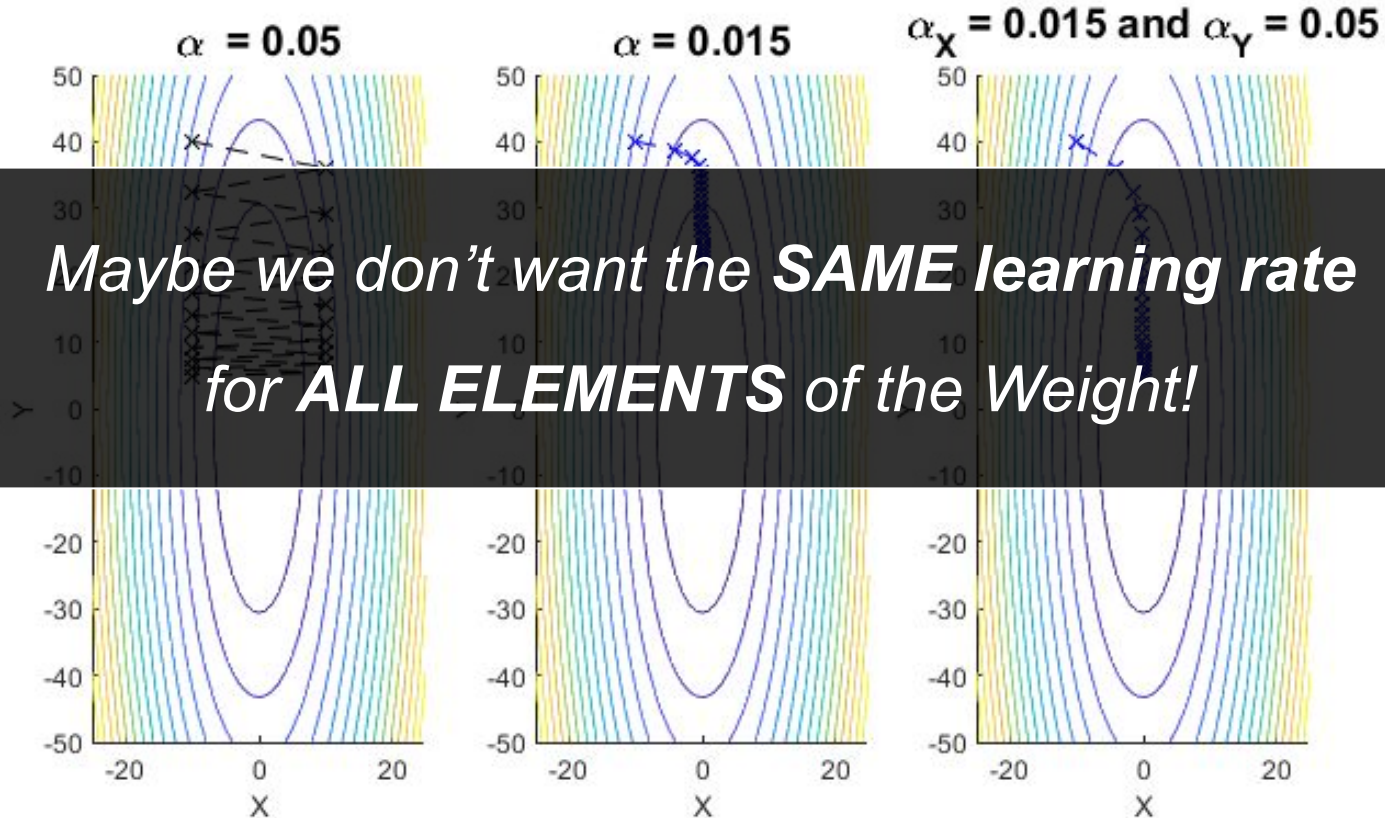
$$w_{t+1} = w_t + m_{t+1}$$

# Importance of Learning Rate



SGD (learning_rate=0.01)
step: 0: (-6.63, 43.953)

SGD (learning_rate=0.1)
step: 0: (-6.63, 43.953)

SGD(learning_rate=0.95)
step: 0: (-6.63, 43.953)

SGD(learning_rate=1.01)
step: 0: (-6.63, 43.953)

## Another example

# Another example

# Adaptive Learning Rate



*Maybe we don't want the **SAME learning rate** for **ALL ELEMENTS** of the Weight!*

## Adaptive Optimizers

*Different Learning Rate for each element of the Model Weights!*

# AdaGrad (Duchi et al. 2011)

More updates -> more decay

- Handle sparse gradients well
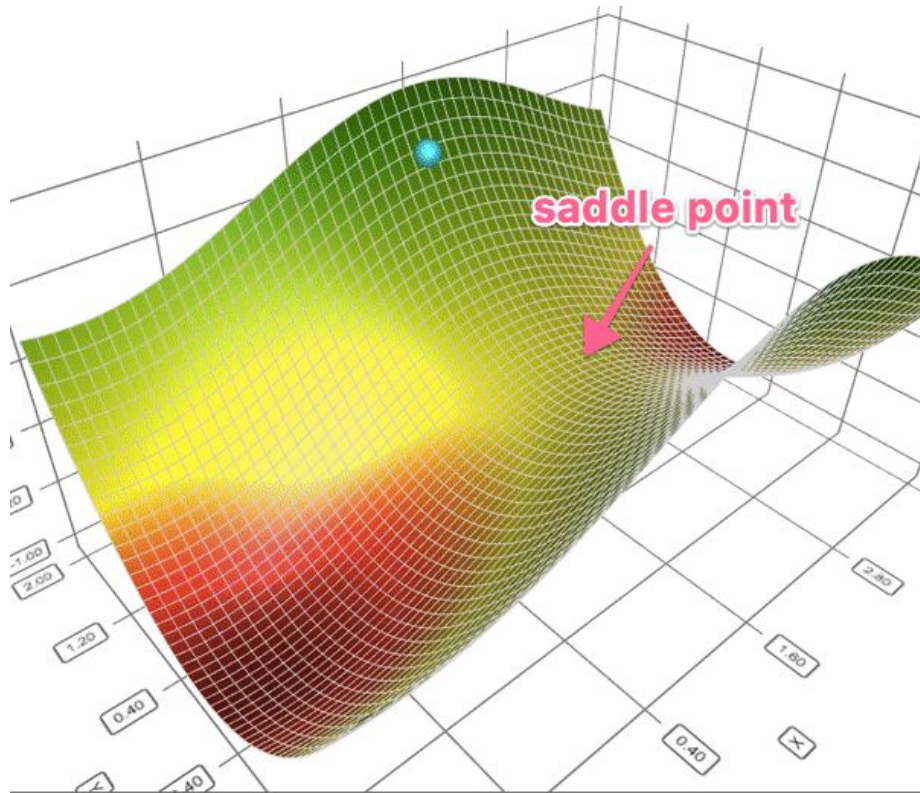  - *Sparse: The vector has 0 in most of the entries*

*Element-wise product*

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

SGD

Adagrad

saddle point

Gradient Descent
AdaGrad

# AdaGrad (Duchi et al. 2011)

More updates → more decay

- Handle sparse gradients well
  - *Sparse: The vector has 0 in most of the entries*

*Element-wise product*

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

SGD

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

Adagrad

***Exercise:***
*What's could be wrong with this optimizator?*
*(What would happen to the denominator.)*

# AdaGrad (Duchi et al. 2011)

More updates → more decay

- Handle sparse gradients well
  - *Sparse: The vector has 0 in most of the entries*

*Element-wise product*

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \mathcal{L}(\mathbf{w}_t)$$

SGD

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

Adagrad

**Issue**: *decays too aggressively!*

# RMSProp (Graves, 2013)

Keep an **exponential moving average** of the squared gradient for each element

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta \mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$

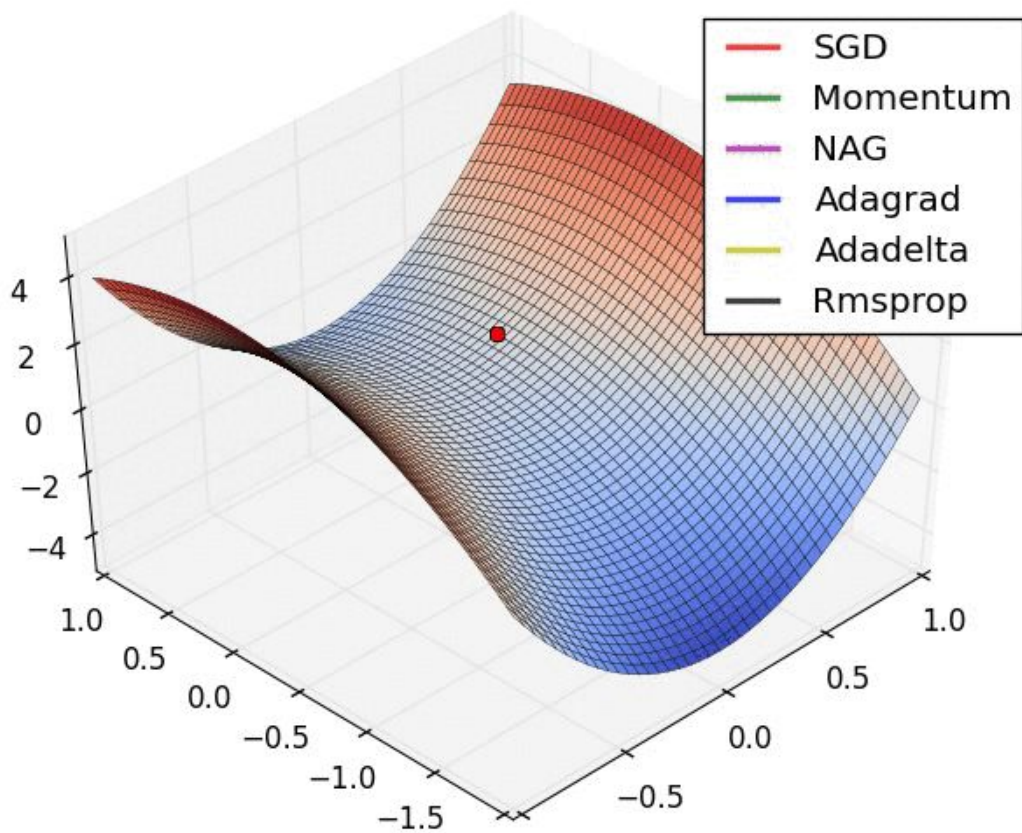$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

Adagrad

RmsProp

where $\beta \in [0, 1]$ the exponential moving average constant.

# Demo

Adagrad & RMSprop



Legend:
- SGD
- Momentum
- NAG
- Adagrad
- Adadelta
- Rmsprop

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

Momentum

$$\mathbf{v}_{t+1} = \beta\mathbf{v}_t + (1-\beta)\mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

RMSProp

$$\mathbf{v}_{t+1} = \beta \mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

RMSProp

**ADAM**
(**Ada**ptive **M**oment Estimate)

$$\mathbf{v}_{t+1} = \beta_2 \mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\boxed{\sqrt{\widehat{\mathbf{v}}_{t+1} + \epsilon}}} \odot \widehat{\mathbf{m}}_{t+1}$$

$$\mathbf{v}_{t+1} = \beta \mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\boxed{\sqrt{\mathbf{v}_{t+1} + \epsilon}}} \odot \mathbf{g}_t$$

RMSProp

**ADAM**
(**Ada**ptive **M**oment Estimate)

$$\mathbf{m}_{t+1} = \mu\mathbf{m}_t - \alpha\nabla l(\mathbf{w}_t; \mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{m}_{t+1}$$

Momentum

$$\mathbf{v}_{t+1} = \beta\mathbf{v}_t + (1 - \beta)\mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{v}_{t+1} + \epsilon}} \odot \mathbf{g}_t$$

RMSProp

$$\mathbf{m}_{t+1} = \beta_1\mathbf{m}_t + (1 - \beta_1)\mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta_2\mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\widehat{\mathbf{v}}_{t+1} + \epsilon}} \odot \widehat{\mathbf{m}}_{t+1}$$

**ADAM**
(**Ada**ptive **M**oment Estimate)

$$\mathbf{m}_{t+1} = \beta_1 \mathbf{m}_t + (1 - \beta_1)\mathbf{g}_t$$

$$\mathbf{v}_{t+1} = \beta_2 \mathbf{v}_t + (1 - \beta_2)\mathbf{g}_t^2$$

$$\widehat{\mathbf{m}}_{t+1} = \frac{\mathbf{m}_{t+1}}{1 - \beta_1^{t+1}}$$

$$\widehat{\mathbf{v}}_{t+1} = \frac{\mathbf{v}_{t+1}}{1 - \beta_2^{t+1}}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\alpha}{\sqrt{\widehat{\mathbf{v}}_{t+1} + \epsilon}} \odot \widehat{\mathbf{m}}_{t+1}$$

$$\mathbb{E}[\widehat{\mathbf{m}}_{t+1}]$$

$$\mathbb{E}[\widehat{\mathbf{v}}_{t+1}]$$

**ADAM**
(**Ada**ptive **M**oment Estimate)

# Optimizers Recap

- Gradient Descent
  - *Vanilla, costly, but for best convergence rate*
- Stochastic Gradient Descent
  - *Simple, lightweight*
- **Mini-batch SGD**
  - *balanced between SGD and GD*
  - ***1st choice for small, simple models***

- SGD w. Momentum
  - *Faster, capable to jump out local minimum*

- AdaGrad
- RMSProp
- **ADAM**
  - **JUST USE ADAM IF YOU DON'T KNOW WHAT TO USE IN DEEP LEARNING**