# Computational Linguistics LT3233

Jixing Li

Lecture 9: Backpropagation and Computational Graph
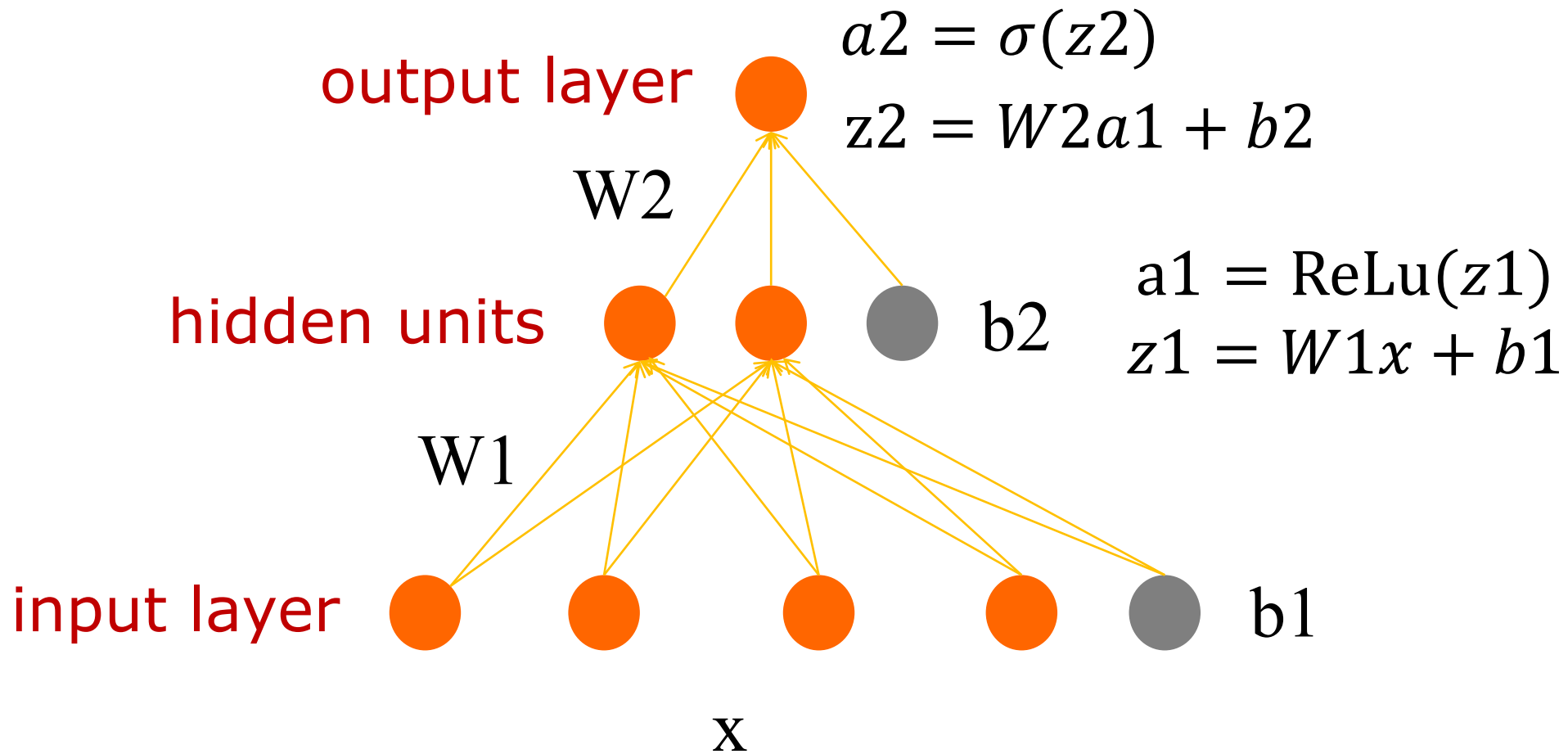
Slides adapted from Dan Jurafsky

# Lecture plan

- Overview of feedforward neural networks

- Backpropagation and computational graphs

- Implementing feedforward neural networks from scratch

- Short break (15 mins)

- Hands-on exercises

Ask for help if you need it:

- office hour: 3-5 pm Tuesdays at LI-5459

- Zoom meetings: by schedule

# Feedforward neural networks

Two-layer network with scalar output



output layer

$$a2 = \sigma(z2)$$
$$z2 = W2a1 + b2$$

W2

hidden units

$$a1 = \text{ReLu}(z1)$$
$$z1 = W1x + b1$$

b2

W1

input layer

b1

x

# Example

| chinese_name | major | gender | n1_male | n2_male | n1_uniqueness | n2_uniqueness |
|---|---|---|---|---|---|---|
| 林加敏 | LLA | F | 0.442 | -0.562 | 2.795 | 2.087 |

x = [0.442,-0.562,2.795,2.087]

W1 = [[1,3,2,4],
      [2,1,4,3]]

W2 = [-1,2]

b1 = [1,-1], b2 = 2

z1 = W1x+b1

a1 = ReLU(z1) = max(z1,0)

z2 = W2a1+b2

$a2 = \sigma(z2) = 1/(1+e^{-z2})$

$a1 = ReLu(z1)$

$z1 = W1x + b1$

$a2 = \sigma(z2)$

$z2 = W2a1 + b2$

W2

b2

W1

b1

x

# Example

$x = [0.442, -0.562, 2.795, 2.087]$

$W1 = [[1,3,2,4],$
$\qquad [2,1,4,3]]$

$W2 = [-1,2]$
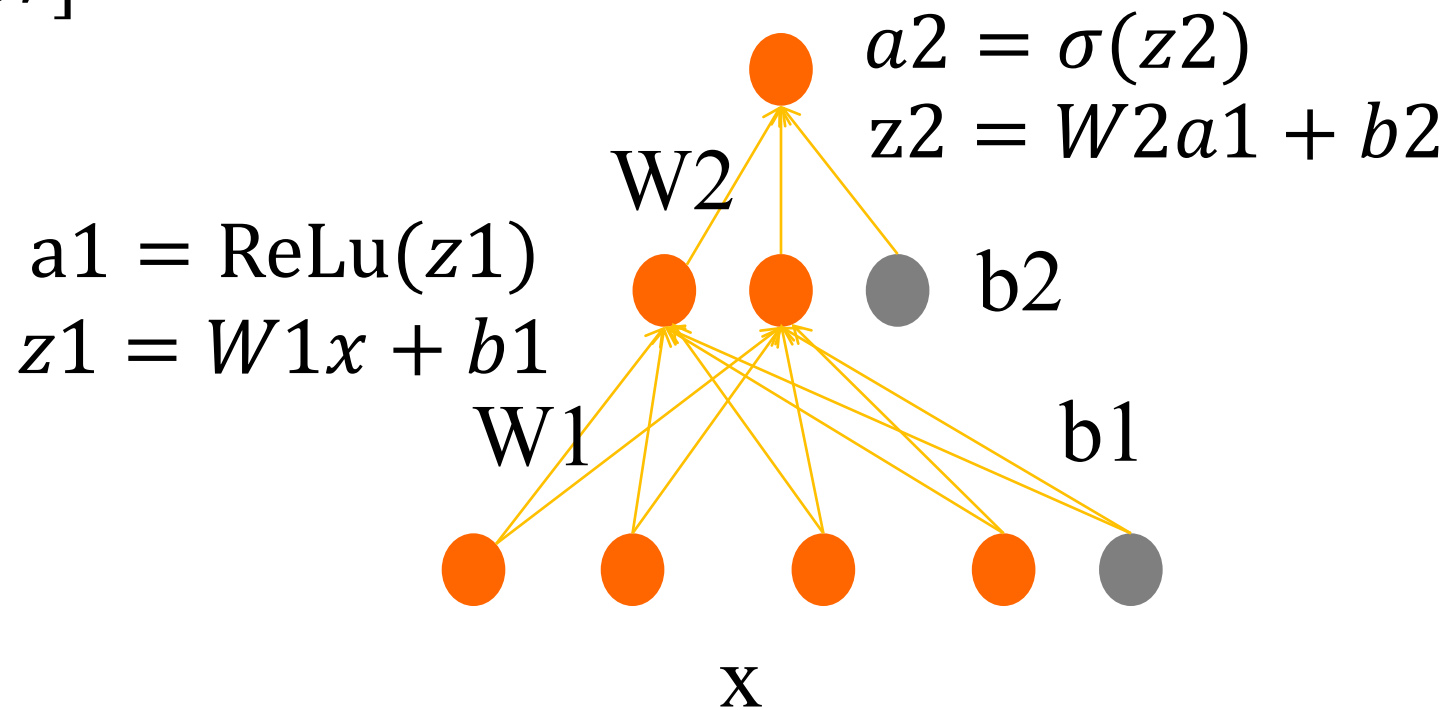
$b1 = [1,-1], b2 = 2$

$z1 = W1x + b1$

$a1 = ReLU(z1) = max(z1,0)$

$z2 = W2a1 + b2$

$a2 = \sigma(z2) = 1/(1+e^{-z2})$

$$\begin{bmatrix} 1,3,2,4 \\ 2,1,4,3 \end{bmatrix} \times \begin{bmatrix} 0.442 \\ -0.562 \\ 2.795 \\ 2.087 \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$z1 = \begin{bmatrix} 1\times0.442 + 3\times-0.562 + 2\times2.795 + 4\times2.087 \\ 2\times0.442 + 1\times-0.562 + 4\times2.795 + 3\times2.087 \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$z1 = \begin{bmatrix} 12.694 \\ 17.763 \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$
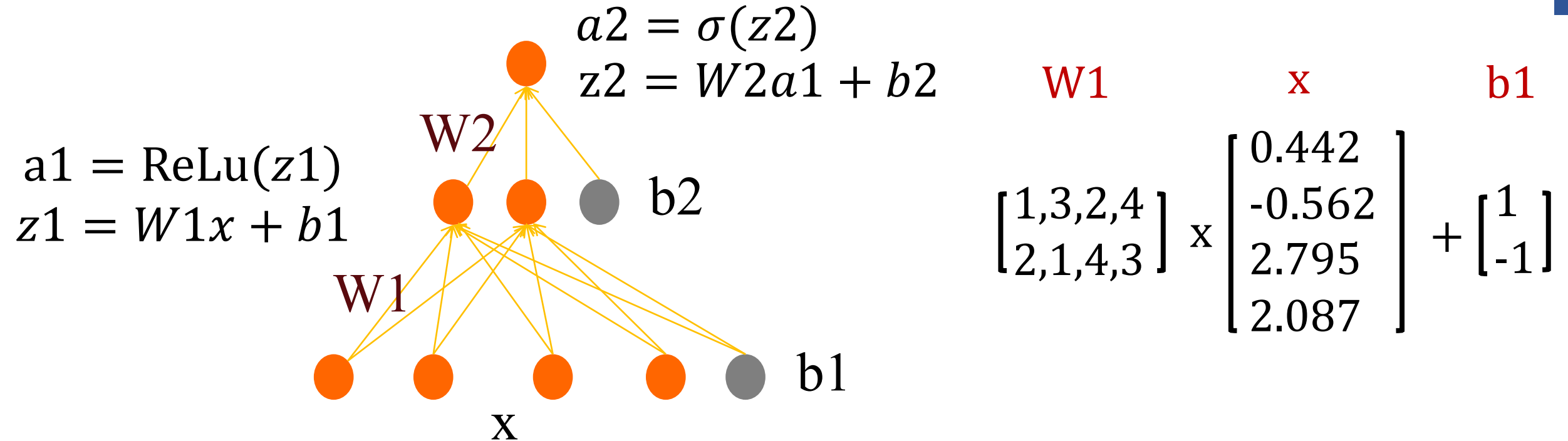
$$z1 = \begin{bmatrix} 13.694 \\ 16.763 \end{bmatrix}$$

$a1 = ReLU(z1) = max(z1,0) = z1$

$$z2 = W2a1 + b2 = [-1,2] \times \begin{bmatrix} 13.694 \\ 16.763 \end{bmatrix} + 2 = -1\times13.694 + 2\times16.763 + 2 = 21.832$$

$a2 = 1/(1+e^{-z2}) = 1/(1+e^{-21.832}) = 0.99$

# Compute the parameters

$$a2 = \sigma(z2)$$
$$z2 = W2a1 + b2$$

$$a1 = \text{ReLu}(z1)$$
$$z1 = W1x + b1$$

W2

b2

W1

x

b1

W1    x    b1

$$\begin{bmatrix} 1,3,2,4 \\ 2,1,4,3 \end{bmatrix} \times \begin{bmatrix} 0.442 \\ \text{-}0.562 \\ 2.795 \\ 2.087 \end{bmatrix} + \begin{bmatrix} 1 \\ \text{-}1 \end{bmatrix}$$

How to know the weights (W1,W2) and biases (b1,b2)?

→through error **backpropagation**
   which relies on **computation graphs**

© Jixing Li

# Gradient descent in logistic regression

[卓, 琳, Cheuk, Lam, LLA] → x = [0.5, 0.7, 0.5, 0.6, 0.8], y=1

**1. initialize** w **and** b**, set** η
w = [0, 0, 0, 0, 0], b = 0, η = 0.1

**2. compute** $\hat{y}$
$\hat{y}$ = σ(wx + b) = 0.5

**3. compute the gradients for** w **and** b
Gw = ($\hat{y}$ -y)x = (0.5– 1)[0.5, 0.7, 0.5, 0.6, 0.8] = [-0.25, -0.35, -0.25, -0.3, -0.4]
Gb = $\hat{y}$ -y = 0.5– 1 = -0.5

**4. update** w **and** b
$w_{t+1}$ = $w_t$ – ηGw = [0, 0, 0, 0, 0] – 0.1* [-0.25, -0.35, -0.25, -0.3, -0.4]
    = [0.025, 0.035, 0.025, 0.03, 0.04]
$b_{t+1}$ = $b_t$ – ηGb = 0- 0.1*(-0.5) = 0.05

# Backpropagation

$x = [0.442, -0.562, 2.795, 2.087]$, $y=1$

**1. initialize** $W1, W2$ **and** $b1, b2$**, set** $\eta$

$W1 = [[1,1,1,1],$
$\qquad [1,1,1,1]]$
$W2 = [1,1]$
$b1 = [1,1]$
$b2 = [1]$
$\eta = 0.1$

**2. forward propagation**

$z1 = W1x + b1$
$a1 = ReLU(z1) = \max(z1, 0)$
$z2 = W2a1 + b2$
$a2 = \sigma(z2) = 1/(1 + e\text{-}z2)$

**3. backpropagation**

$GW1$
$GW2$
$Gb1$
$Gb2$

**4. update** $W1, W2$ **and** $b1, b2$

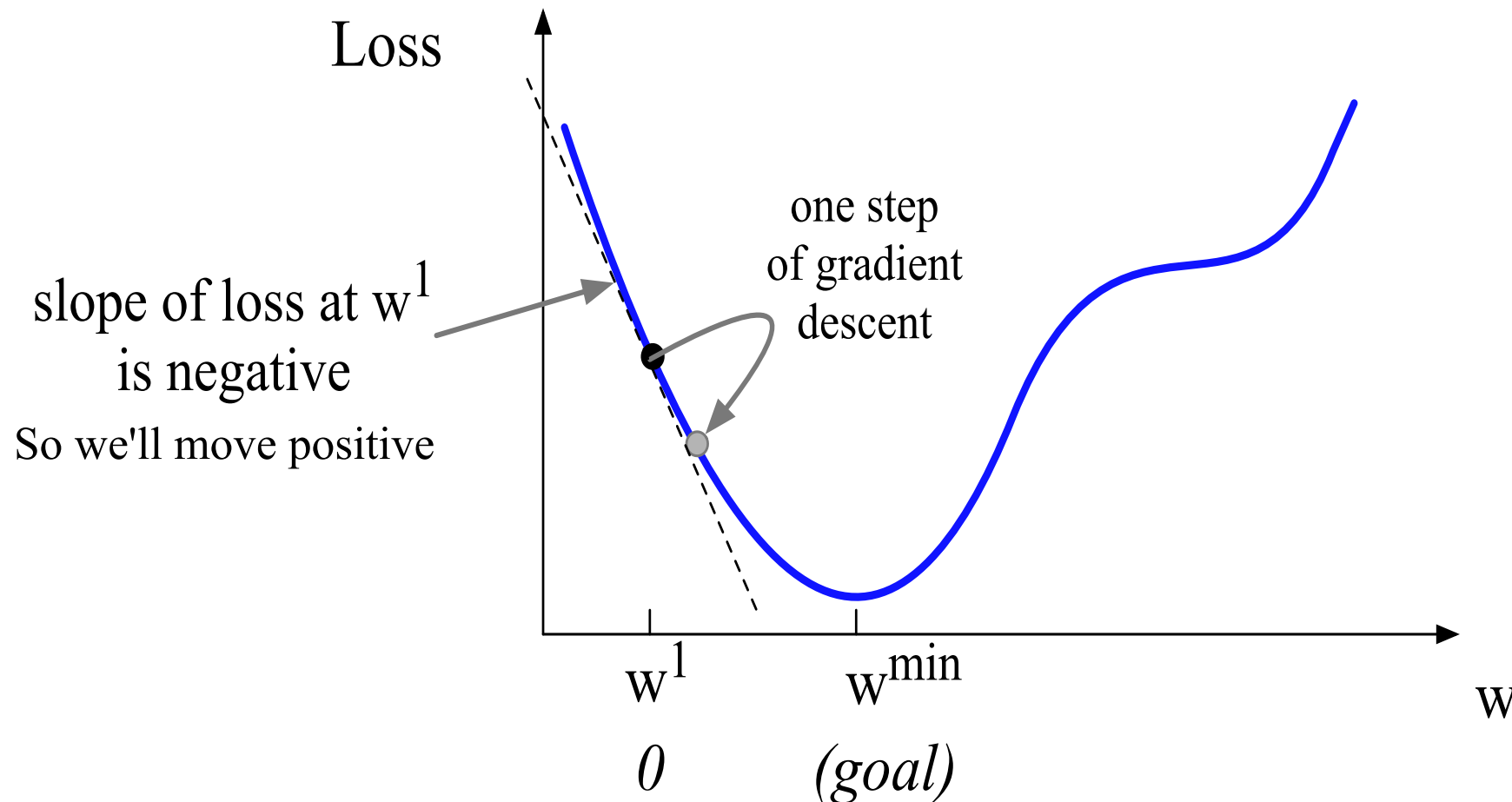$W1_{t+1} = W1_t - \eta GW1$
$W2_{t+1} = W2_t - \eta GW2$
$b1_{t+1} = b1_t - \eta Gb1$
$b2_{t+1} = b2_t - \eta Gb2$

© Jixing Li

# Gradient descent (again)

**Minimize loss:** Given the current $w$, move $w$ in the reverse direction from the slope of the function



concept of derivative

Loss

slope of loss at $w^1$ is negative

So we'll move positive

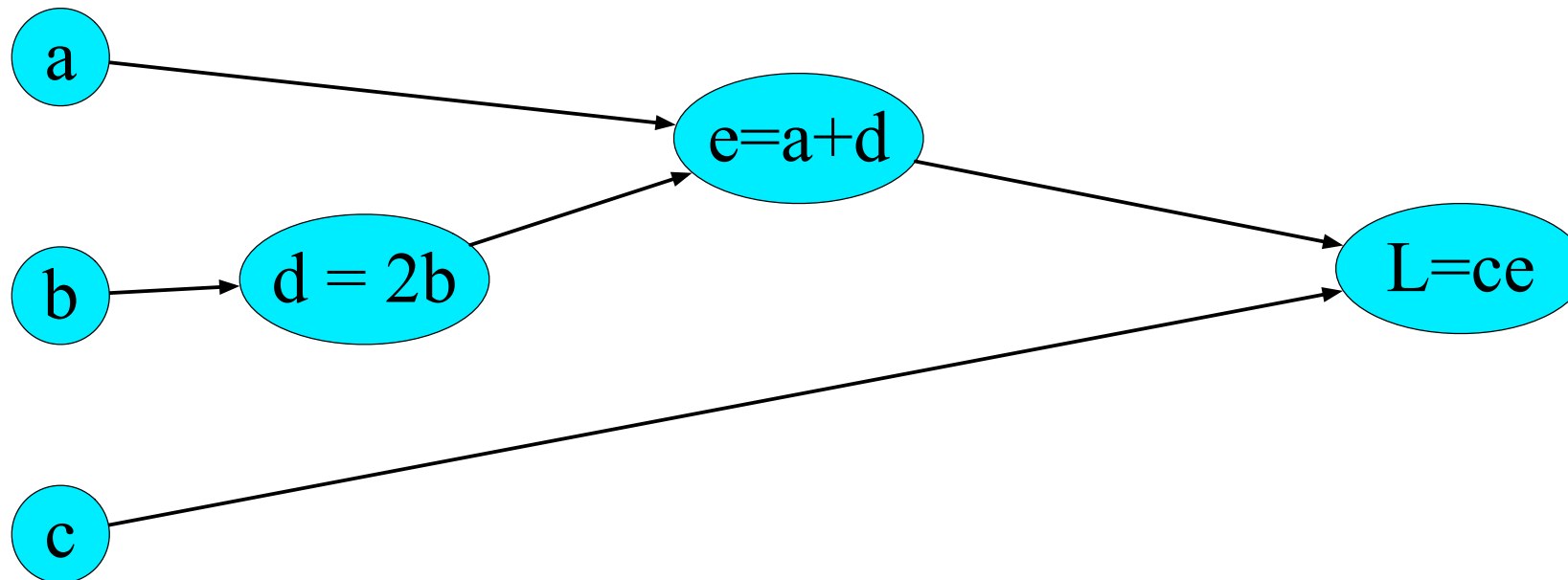one step of gradient descent

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

$0$     *(goal)*

# Computation graph

A computation graph represents the process of computing a mathematical expression

$$L(a,b,c) = c(a+2b)$$

$$
\begin{aligned}
d &= 2*b \\
e &= a+d \\
L &= c*e
\end{aligned}
$$

# Example

$$L(a,b,c) = c(a+2b)$$

$$
\begin{aligned}
d &= 2*b \\
e &= a+d \\
L &= c*e
\end{aligned}
$$

forward pass

3

a

e=5

d=2

1

b

e=d+a

d = 2b

L=-10

-2

c

L=ce

# Example

$$L(a, b, c) = c(a + 2b)$$

$$
\begin{aligned}
d &= 2 * b \\
e &= a + d \\
L &= c * e
\end{aligned}
$$

We want: $\frac{\partial L}{\partial a}$, $\frac{\partial L}{\partial b}$, and $\frac{\partial L}{\partial c}$

The derivative $\frac{\partial L}{\partial a}$ tells us how much a small change in $a$ affects $L$.

# The chain rule

Computing the derivative of a composite function:

$$f(x) = u(v(x))$$

$$\frac{df}{dx} = \frac{du}{dv} \cdot \frac{dv}{dx}$$

$$f(x) = u(v(w(x)))$$

$$\frac{df}{dx} = \frac{du}{dv} \cdot \frac{dv}{dw} \cdot \frac{dw}{dx}$$

© Jixing Li

# Example

$$f(x) = \frac{1}{x} \qquad \frac{df}{dx} = -1/x^2$$

$$L(a,b,c) = c(a+2b)$$

$$\begin{aligned} d &= 2*b \\ e &= a+d \\ L &= c*e \end{aligned}$$

$$f_c(x) = c + x \qquad \frac{df}{dx} = 1$$

$$f(x) = e^x \qquad \frac{df}{dx} = e^x$$

We want: $\frac{\partial L}{\partial a}$, $\frac{\partial L}{\partial b}$, and $\frac{\partial L}{\partial c}$

$$f_a(x) = ax \qquad \frac{df}{dx} = a$$

$$\frac{\partial L}{\partial c} = e$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial e}\frac{\partial e}{\partial a}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial e}\frac{\partial e}{\partial d}\frac{\partial d}{\partial b}$$

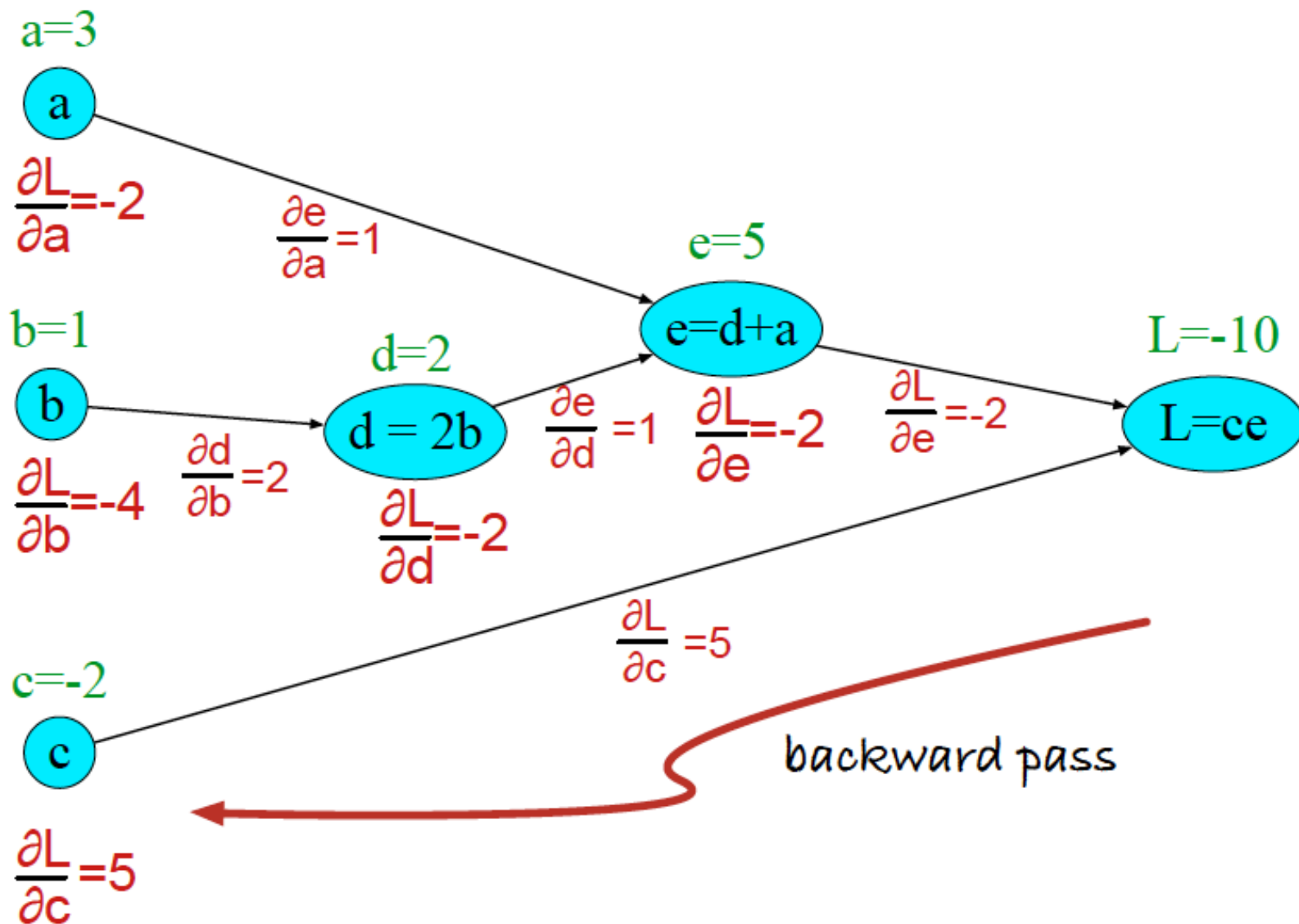$$L = ce \; : \quad \frac{\partial L}{\partial e} = c, \frac{\partial L}{\partial c} = e$$

$$e = a+d \; : \quad \frac{\partial e}{\partial a} = 1, \frac{\partial e}{\partial d} = 1$$

$$d = 2b \; : \quad \frac{\partial d}{\partial b} = 2$$

# Example

a=3

$$\frac{\partial L}{\partial a}=-2$$

$$\frac{\partial e}{\partial a}=1$$

$$L(a,b,c) = c(a+2b)$$

b=1

e=5

$$\frac{\partial L}{\partial b}=-4$$

d=2

$$\frac{\partial d}{\partial b}=2$$

e=d+a

L=-10

$$d = 2*b$$

$$\frac{\partial e}{\partial d}=1$$

$$\frac{\partial L}{\partial e}=-2$$

$$\frac{\partial L}{\partial e}=-2$$

d = 2b

$$e = a+d$$

L=ce

$$L = c*e$$

$$\frac{\partial L}{\partial d}=-2$$

$$\frac{\partial L}{\partial c}=5$$

c=-2

backward pass

c

$$\frac{\partial L}{\partial c}=5$$

# Backprop on a two-layer network



y

Sigmoid activation    σ

$W^{[2]}$    $b^{[2]}$

ReLU activation

$W^{[1]}$    $b^{[1]}$

$x_1$    $x_2$

$$z^{[1]} = W^{[1]}\mathbf{x} + b^{[1]}$$

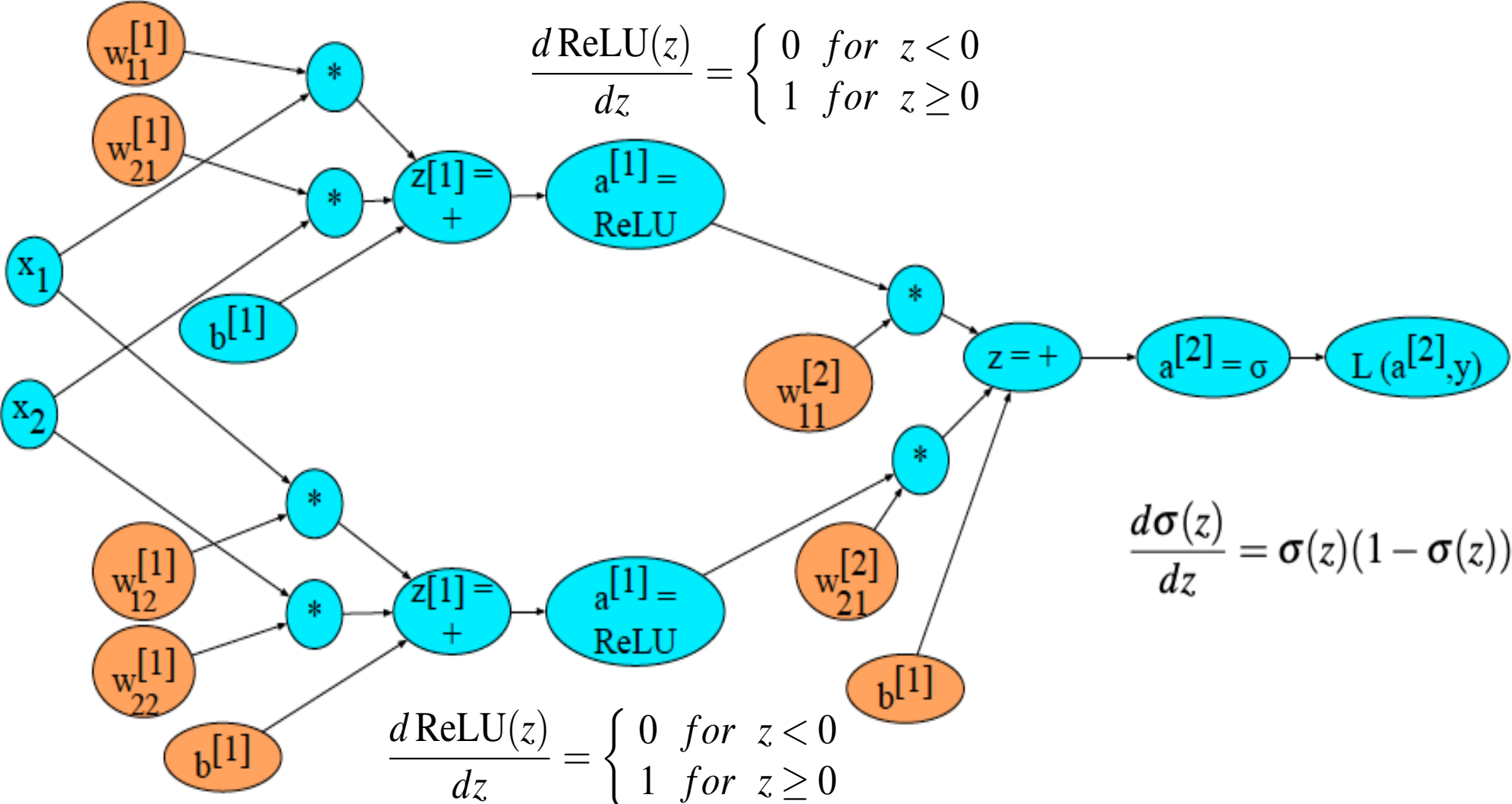$$a^{[1]} = \text{ReLU}(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

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

$$\hat{y} = a^{[2]}$$

$$\frac{d\,\text{ReLU}(z)}{dz} = \begin{cases} 0 & for \ z < 0 \\ 1 & for \ z \geq 0 \end{cases}$$

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

# Backprop on a two-layer network



$$\frac{d\,\text{ReLU}(z)}{dz} = \begin{cases} 0 & for \ \ z < 0 \\ 1 & for \ \ z \geq 0 \end{cases}$$

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

$$\frac{d\,\text{ReLU}(z)}{dz} = \begin{cases} 0 & for \ \ z < 0 \\ 1 & for \ \ z \geq 0 \end{cases}$$

© Jixing Li

# Starting off the backward pass

$$L(\hat{y}, y) = -(y \log(\hat{y}) + (1 - y)\log(1 - \hat{y}))$$

$$L(a, y) = -(y \log a + (1 - y)\log(1 - a))$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z}$$

$$
\begin{aligned}
z^{[1]} &= W^{[1]}\mathbf{x} + b^{[1]} \\
a^{[1]} &= \text{ReLU}(z^{[1]}) \\
z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\
a^{[2]} &= \sigma(z^{[2]}) \\
\hat{y} &= a^{[2]}
\end{aligned}
$$

$$\frac{\partial L}{\partial a} = -\left(\left(y \frac{\partial \log(a)}{\partial a}\right) + (1 - y)\frac{\partial \log(1 - a)}{\partial a}\right)$$

$$= -\left(\left(y\frac{1}{a}\right) + (1 - y)\frac{1}{1 - a}(-1)\right) = -\left(\frac{y}{a} + \frac{y - 1}{1 - a}\right)$$

$$\frac{\partial a}{\partial z} = a(1 - a)$$

$$\boxed{\frac{\partial L}{\partial z}} = -\left(\frac{y}{a} + \frac{y - 1}{1 - a}\right)a(1 - a) = \boxed{a - y}$$

# To do

- Optional reading: **SLP** Ch7.6
- Tutorial on backpropagation:

https://cs231n.github.io/optimization-2/

- Gentle introduction on derivatives:

https://www.khanacademy.org/math/ap-calculus-ab/ab-differentiation-1-new