# TOÁN CHO KHOA HỌC MÁY TÍNH

## MẠNG NƠ-RƠN NHÂN TẠO

**TS. Lương Ngọc Hoàng**

# Nội dung

1. Giới thiệu mạng nơ-rơn nhân tạo
2. Tính toán gradient với lan truyền ngược
3. Phân tích mạng nơ-rơn

# MẠNG NƠ-RƠN NHÂN TẠO

## ARTIFICIAL NEURAL NETWORK

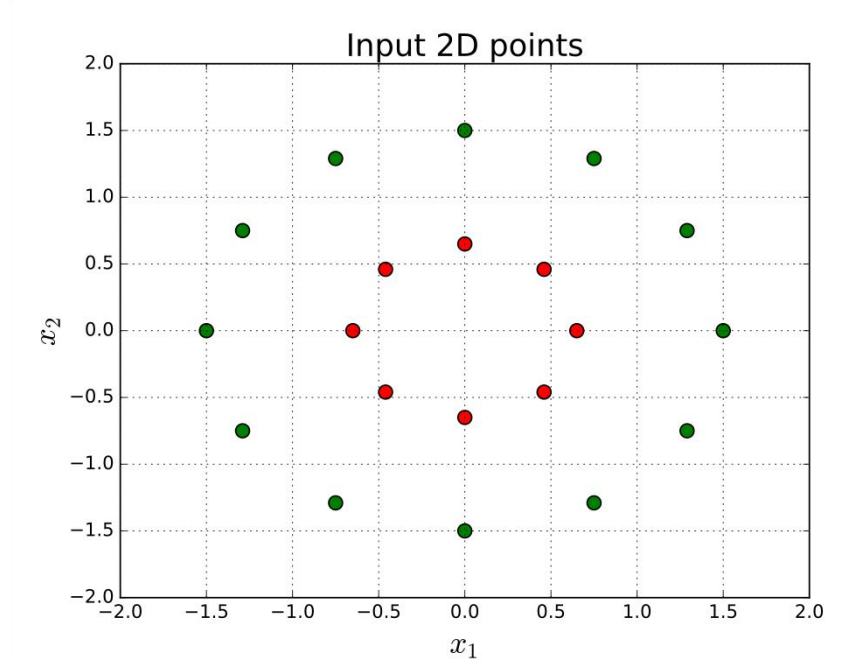Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Bài toán phân lớp

Bài toán phân lớp (classification):

Với mỗi điểm dữ liệu đầu vào $x = (x_1, x_2)$, ta muốn đầu ra của mô hình dự đoán nhãn (label) của $x$ là nhãn 0 (lớp đỏ), hay nhãn 1 (lớp xanh lá).



Ta muốn đầu ra của mô hình dự đoán xác suất điểm dữ liệu $x = (x_1, x_2)$ thuộc về lớp đỏ (nhãn 0), và lớp xanh lá (nhãn 1) là bao nhiêu.
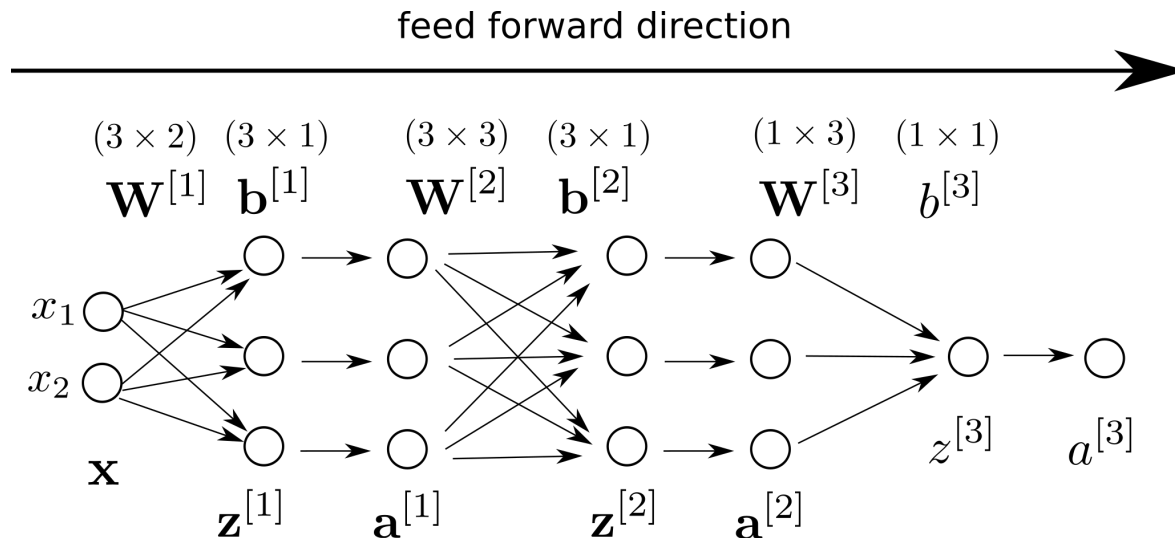
*Hồi quy logistic có phù hợp với tập dữ liệu này?*

- Biên quyết định của hồi quy logistic có dạng tuyến tính (đường thẳng).
- Tập dữ liệu trên không thể phân chia tuyến tính được.

# Mạng nơ-rơn (neural network)

feed forward direction →

$$(3 \times 2) \quad (3 \times 1) \qquad (3 \times 3) \quad (3 \times 1) \qquad\qquad (1 \times 3) \quad (1 \times 1)$$

$$\mathbf{W}^{[1]} \quad \mathbf{b}^{[1]} \qquad \mathbf{W}^{[2]} \quad \mathbf{b}^{[2]} \qquad\qquad \mathbf{W}^{[3]} \quad b^{[3]}$$

$x_1$
$x_2$

$\mathbf{x}$

$z^{[3]} \quad a^{[3]}$

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

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

$$a^{[1]} = g(z^{[1]})$$

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

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

$$z^{[3]} = W^{[3]}a^{[2]} + b^{[3]}$$

$$a^{[3]} = g(z^{[3]}) \in (0,1)$$

$g()$ là hàm kích hoạt (activation function) logistic sigmoid. $\quad g(z) = \dfrac{1}{1+e^{-z}}$

**Lớp ẩn** (hidden layer) thứ nhất:

$$W^{[1]} = \begin{bmatrix} w_{1,1}^{[1]} & w_{1,2}^{[1]} \\ w_{2,1}^{[1]} & w_{2,2}^{[1]} \\ w_{3,1}^{[1]} & w_{3,2}^{[1]} \end{bmatrix}, \; \boldsymbol{b}^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \end{bmatrix}$$
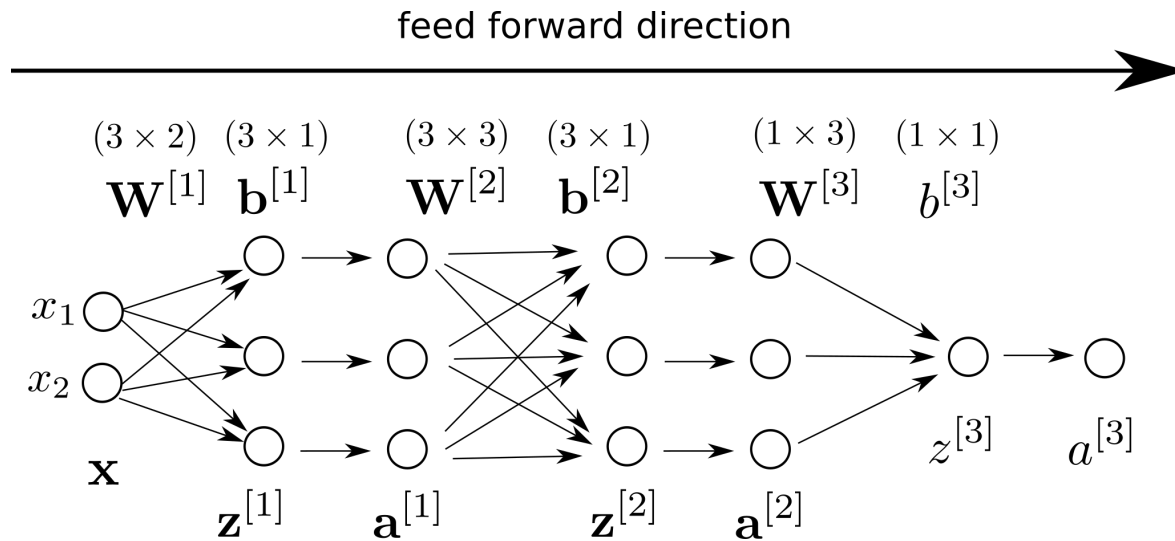
**Lớp ẩn** (hidden layer) :

$$W^{[2]} = \begin{bmatrix} w_{1,1}^{[2]} & w_{1,2}^{[2]} & w_{1,3}^{[2]} \\ w_{2,1}^{[2]} & w_{2,2}^{[2]} & w_{2,3}^{[2]} \\ w_{3,1}^{[2]} & w_{3,2}^{[2]} & w_{3,3}^{[2]} \end{bmatrix}, \; \boldsymbol{b}^{[2]} = \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \end{bmatrix}$$
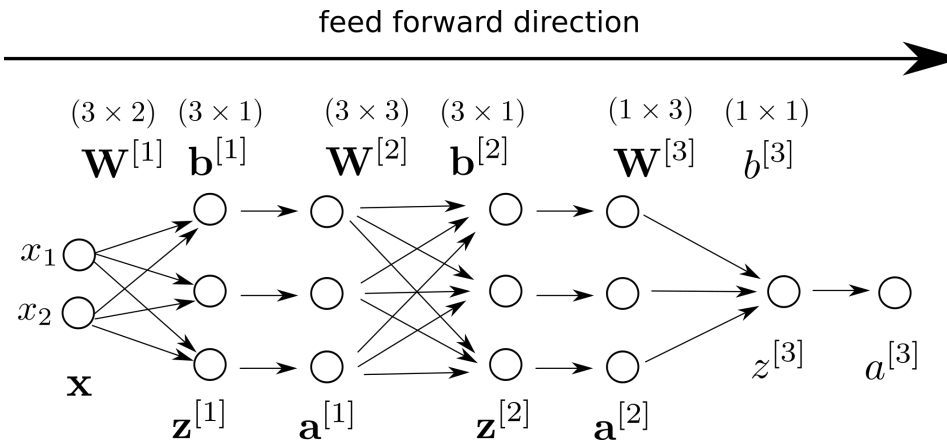
**Lớp đầu ra** (output layer):

$$W^{[3]} = \begin{bmatrix} w_{1,1}^{[3]} & w_{1,2}^{[3]} & w_{1,3}^{[3]} \end{bmatrix}, \; b^{[3]} = b_1^{[3]}$$

# Mạng nơ-rơn (neural network)

feed forward direction

$$W^{[1]} = \begin{bmatrix} w^{[1]}_{1,1} & w^{[1]}_{1,2} \\ w^{[1]}_{2,1} & w^{[1]}_{2,2} \\ w^{[1]}_{3,1} & w^{[1]}_{3,2} \end{bmatrix}, \boldsymbol{b}^{[1]} = \begin{bmatrix} b^{[1]}_1 \\ b^{[1]}_2 \\ b^{[1]}_3 \end{bmatrix}$$

$$(3 \times 2) \quad (3 \times 1) \quad (3 \times 3) \quad (3 \times 1) \quad (1 \times 3) \quad (1 \times 1)$$

$$\mathbf{W}^{[1]} \quad \mathbf{b}^{[1]} \quad \mathbf{W}^{[2]} \quad \mathbf{b}^{[2]} \quad \mathbf{W}^{[3]} \quad b^{[3]}$$

$x_1$

$x_2$

$z^{[3]} \quad a^{[3]}$

$\mathbf{x}$

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

**Lớp ẩn (hidden layer) :**

$$W^{[2]} = \begin{bmatrix} w^{[2]}_{1,1} & w^{[2]}_{1,2} & w^{[2]}_{1,3} \\ w^{[2]}_{2,1} & w^{[2]}_{2,2} & w^{[2]}_{2,3} \\ w^{[2]}_{3,1} & w^{[2]}_{3,2} & w^{[2]}_{3,3} \end{bmatrix}, \boldsymbol{b}^{[2]} = \begin{bmatrix} b^{[2]}_1 \\ b^{[2]}_2 \\ b^{[2]}_3 \end{bmatrix}$$

$$a^{[3]} = g(\mathbf{z}^{[3]}) = g(W^{[3]}\boldsymbol{a}^{[2]} + b^{[3]}) \qquad g(z) = \frac{1}{1+e^{-z}}$$

$$= g(W^{[3]}g(\mathbf{z}^{[2]}) + b^{[3]})$$

$$= g(W^{[3]}g(W^{[2]}\boldsymbol{a}^{[1]} + \boldsymbol{b}^{[2]}) + b^{[3]})$$

$$= g(W^{[3]}g(W^{[2]}g(\mathbf{z}^{[1]}) + \boldsymbol{b}^{[2]}) + b^{[3]})$$

$$= g(W^{[3]}g(W^{[2]}g(W^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}) + \boldsymbol{b}^{[2]}) + b^{[3]})$$

**Lớp đầu ra (output layer):**

$$W^{[3]} = \begin{bmatrix} w^{[3]}_{1,1} & w^{[3]}_{1,2} & w^{[3]}_{1,3} \end{bmatrix}, b^{[3]} = b^{[3]}_1$$

*Ta có cần các hàm kích hoạt tại các nơ-rơn của lớp ẩn?*

$$a^{[3]} = g(W^{[3]}(W^{[2]}(W^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}) + \boldsymbol{b}^{[2]}) + b^{[3]})$$

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Mạng nơ-rơn (neural network)

feed forward direction



$$g(z) = \frac{1}{1+e^{-z}}$$

**Lớp ẩn** (hidden layer) thứ nhất:

$$W^{[1]} = \begin{bmatrix} w_{1,1}^{[1]} & w_{1,2}^{[1]} \\ w_{2,1}^{[1]} & w_{2,2}^{[1]} \\ w_{3,1}^{[1]} & w_{3,2}^{[1]} \end{bmatrix}, \; \boldsymbol{b}^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \end{bmatrix}$$

**Lớp ẩn** (hidden layer) :

$$W^{[2]} = \begin{bmatrix} w_{1,1}^{[2]} & w_{1,2}^{[2]} & w_{1,3}^{[2]} \\ w_{2,1}^{[2]} & w_{2,2}^{[2]} & w_{2,3}^{[2]} \\ w_{3,1}^{[2]} & w_{3,2}^{[2]} & w_{3,3}^{[2]} \end{bmatrix}, \; \boldsymbol{b}^{[2]} = \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \end{bmatrix}$$
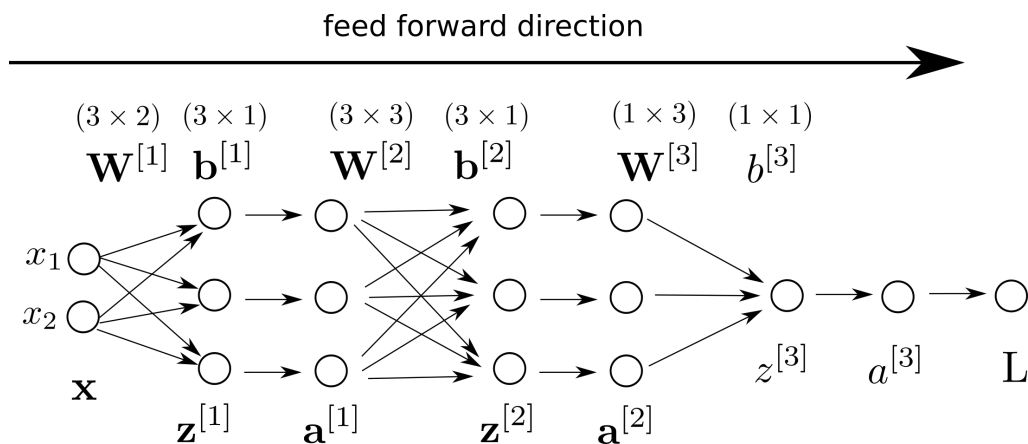
**Lớp đầu ra** (output layer):

$$W^{[3]} = \begin{bmatrix} w_{1,1}^{[3]} & w_{1,2}^{[3]} & w_{1,3}^{[3]} \end{bmatrix}, \; b^{[3]} = b_1^{[3]}$$

*Nếu bỏ các hàm kích hoạt tại các nơ-rơn của lớp ẩn:*

$$a^{[3]} = g\big(W^{[3]}\big(W^{[2]}\big(W^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}\big) + \boldsymbol{b}^{[2]}\big) + b^{[3]}\big)$$

$$= g\big(W^{[3]}\big(W^{[2]}W^{[1]}\boldsymbol{x} + W^{[2]}\boldsymbol{b}^{[1]} + \boldsymbol{b}^{[2]}\big) + b^{[3]}\big)$$

$$= g\big(\underbrace{W^{[3]}W^{[2]}W^{[1]}}_{(1\times3).(3\times3).(3\times2)}\underbrace{\boldsymbol{x}}_{2\times1} + \underbrace{W^{[3]}W^{[2]}}_{(1\times3).(3\times3)}\underbrace{\boldsymbol{b}^{[1]}}_{3\times1} + \underbrace{W^{[3]}}_{(1\times3)}\underbrace{\boldsymbol{b}^{[2]}}_{3\times1} + \underbrace{b^{[3]}}_{1\times1}\big)$$

$$= g\big(\underbrace{W}_{(1\times2)}\underbrace{\boldsymbol{x}}_{2\times1} + \underbrace{b}_{1\times1}\big)$$

$a^{[3]} = g(\boldsymbol{w}^T\boldsymbol{x} + b)$ là hồi quy logistic (chỉ có khả năng tạo các đường biên quyết định tuyến tính). → Hàm kích hoạt sử dụng phải là các hàm <span style="color:red">phi tuyến tính (non-linear)</span>.

# TÍNH TOÁN GRADIENT VỚI LAN TRUYỀN NGƯỢC

## GRADIENT COMPUTATION WITH BACKPROPAGATION

# Hàm mất mát

feed forward direction



Cho bài toán phân lớp, ta có thể sử dụng hàm Binary Cross Entropy (BCE) Loss như trong hồi quy logistic.
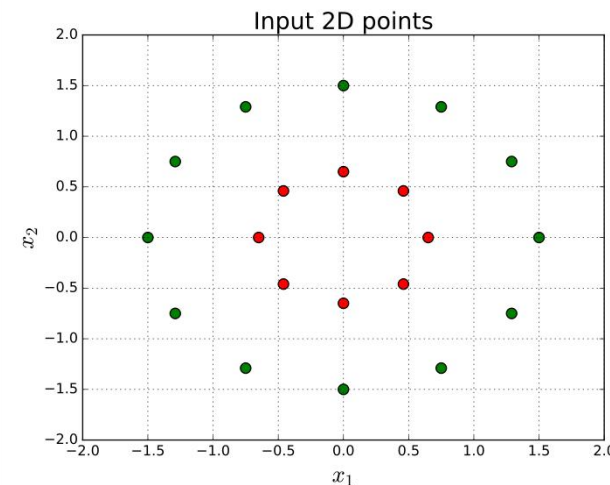
- Trên một điểm dữ liệu $i$:

$$L = -\left(y^{(i)} \log\left(a^{[3](i)}\right) + \left(1 - y^{(i)}\right) \log\left(1 - a^{[3](i)}\right)\right)$$

- Trên toàn bộ tập dữ liệu:

$$J = -\sum_{i=1}^{N}\left(y^{(i)} \log\left(a^{[3](i)}\right) + \left(1 - y^{(i)}\right) \log\left(1 - a^{[3](i)}\right)\right)$$

- Nếu ta có nhiều hơn 2 lớp, ta dùng hàm Cross Entropy (CE) Loss:

$$J = -\sum_{i=1}^{N}\sum_{j=1}^{C} y_j^{(i)} \log\left(a_j^{[3](i)}\right)$$



Ta muốn đầu ra của mô hình dự đoán xác suất điểm dữ liệu $x = (x_1, x_2)$ thuộc về lớp đỏ (nhãn 0), và lớp xanh lá (nhãn 1) là bao nhiêu.

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Lan truyền ngược (backpropagation)

Trong hồi quy logistic, ta có:

$$\frac{\partial L}{\partial z^{[3]}} = a^{[3]} - y$$

$$z^{[3]} = W^{[3]}\boldsymbol{a}^{[2]} + b^{[3]}$$

$$= \begin{bmatrix} w_{1,1}^{[3]} & w_{1,2}^{[3]} & w_{1,3}^{[3]} \end{bmatrix} \begin{bmatrix} a_1^{[2]} \\ a_2^{[2]} \\ a_3^{[2]} \end{bmatrix} + b^{[3]}$$
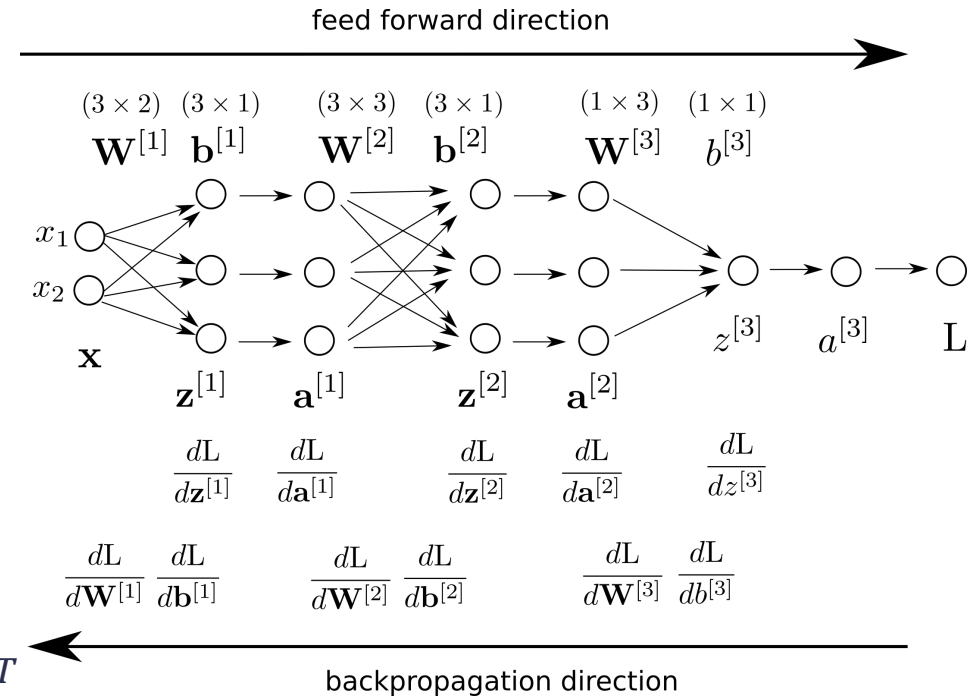
$$= w_{1,1}^{[3]}a_1^{[2]} + w_{1,2}^{[3]}a_2^{[2]} + w_{1,3}^{[3]}a_3^{[2]} + b^{[3]}$$

$$\frac{\partial z^{[3]}}{\partial W^{[3]}} = \begin{bmatrix} \frac{\partial z^{[3]}}{\partial w_{1,1}^{[3]}} & \frac{\partial z^{[3]}}{\partial w_{1,2}^{[3]}} & \frac{\partial z^{[3]}}{\partial w_{1,3}^{[3]}} \end{bmatrix} = \begin{bmatrix} a_1^{[2]} & a_2^{[2]} & a_3^{[2]} \end{bmatrix} = \left(\boldsymbol{a}^{[2]}\right)^T$$

$$\frac{\partial z^{[3]}}{\partial b^{[3]}} = 1$$

$$\textcolor{red}{\frac{\partial L}{\partial W^{[3]}}} = \frac{\partial L}{\partial z^{[3]}}\frac{\partial z^{[3]}}{\partial W^{[3]}} = (a^{[3]} - y)\left(\boldsymbol{a}^{[2]}\right)^T$$

$$\textcolor{red}{\frac{\partial L}{\partial b^{[3]}}} = \frac{\partial L}{\partial z^{[3]}}\frac{\partial z^{[3]}}{\partial b^{[3]}} = (a^{[3]} - y)$$

feed forward direction
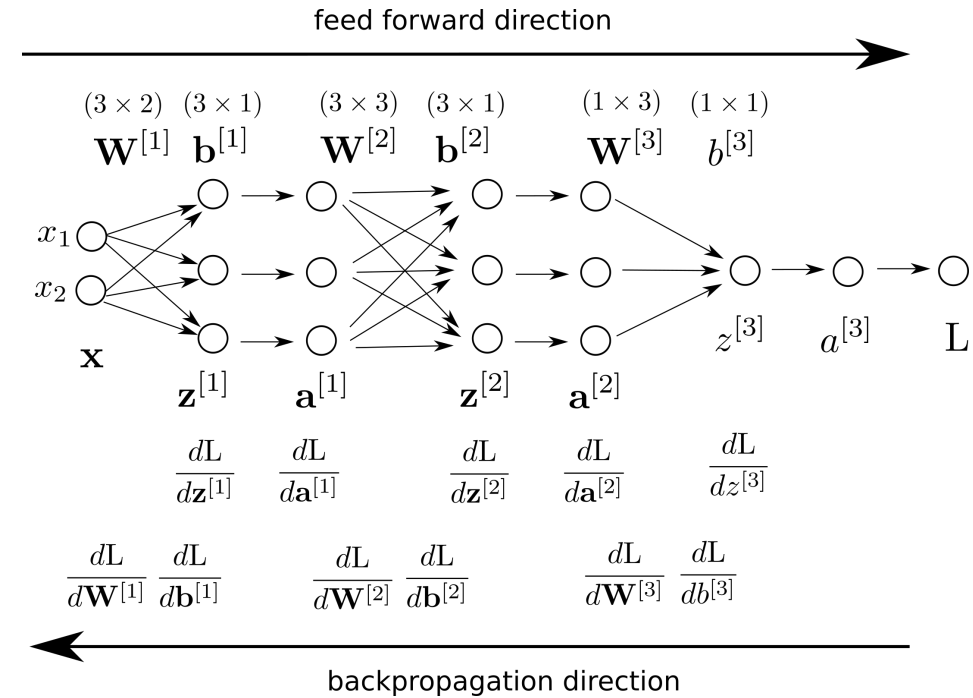
$(3 \times 2)$ $(3 \times 1)$ $(3 \times 3)$ $(3 \times 1)$ $(1 \times 3)$ $(1 \times 1)$

$\mathbf{W}^{[1]}$ $\mathbf{b}^{[1]}$ $\mathbf{W}^{[2]}$ $\mathbf{b}^{[2]}$ $\mathbf{W}^{[3]}$ $b^{[3]}$

$x_1$ $x_2$ $\mathbf{x}$

$\mathbf{z}^{[1]}$ $\mathbf{a}^{[1]}$ $\mathbf{z}^{[2]}$ $\mathbf{a}^{[2]}$ $z^{[3]}$ $a^{[3]}$ $\mathrm{L}$

$\frac{d\mathrm{L}}{d\mathbf{z}^{[1]}}$ $\frac{d\mathrm{L}}{d\mathbf{a}^{[1]}}$ $\frac{d\mathrm{L}}{d\mathbf{z}^{[2]}}$ $\frac{d\mathrm{L}}{d\mathbf{a}^{[2]}}$ $\frac{d\mathrm{L}}{dz^{[3]}}$

$\frac{d\mathrm{L}}{d\mathbf{W}^{[1]}}$ $\frac{d\mathrm{L}}{d\mathbf{b}^{[1]}}$ $\frac{d\mathrm{L}}{d\mathbf{W}^{[2]}}$ $\frac{d\mathrm{L}}{d\mathbf{b}^{[2]}}$ $\frac{d\mathrm{L}}{d\mathbf{W}^{[3]}}$ $\frac{d\mathrm{L}}{db^{[3]}}$

backpropagation direction

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Lan truyền ngược (backpropagation)

$$\frac{\partial L}{\partial \boldsymbol{a}^{[2]}} = \frac{\partial L}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial \boldsymbol{a}^{[2]}}$$

$$z^{[3]} = W^{[3]} \boldsymbol{a}^{[2]} + b^{[3]}$$

$$= \begin{bmatrix} w_{1,1}^{[3]} & w_{1,2}^{[3]} & w_{1,3}^{[3]} \end{bmatrix} \begin{bmatrix} a_1^{[2]} \\ a_2^{[2]} \\ a_3^{[2]} \end{bmatrix} + b^{[3]}$$

$$= w_{1,1}^{[3]} a_1^{[2]} + w_{1,2}^{[3]} a_2^{[2]} + w_{1,3}^{[3]} a_3^{[2]} + b^{[3]}$$

$$\frac{\partial z^{[3]}}{\partial \boldsymbol{a}^{[2]}} = \begin{bmatrix} \dfrac{\partial z^{[3]}}{\partial a_1^{[2]}} \\ \dfrac{\partial z^{[3]}}{\partial a_2^{[2]}} \\ \dfrac{\partial z^{[3]}}{\partial a_3^{[2]}} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[3]} \\ w_{1,2}^{[3]} \\ w_{1,3}^{[3]} \end{bmatrix} = \left( W^{[3]} \right)^T$$

$$\frac{\partial L}{\partial \boldsymbol{a}^{[2]}} = \frac{\partial L}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial \boldsymbol{a}^{[2]}} = \left( a^{[3]} - y \right) \left( W^{[3]} \right)^T$$

feed forward direction

$(3 \times 2)$ $(3 \times 1)$ $\quad$ $(3 \times 3)$ $(3 \times 1)$ $\quad$ $(1 \times 3)$ $(1 \times 1)$

$\mathbf{W}^{[1]}$ $\mathbf{b}^{[1]}$ $\quad$ $\mathbf{W}^{[2]}$ $\mathbf{b}^{[2]}$ $\quad$ $\mathbf{W}^{[3]}$ $b^{[3]}$



$x_1$ $\quad$ $x_2$ $\quad$ $\mathbf{x}$

$\mathbf{z}^{[1]}$ $\quad$ $\mathbf{a}^{[1]}$ $\quad$ $\mathbf{z}^{[2]}$ $\quad$ $\mathbf{a}^{[2]}$ $\quad$ $z^{[3]}$ $\quad$ $a^{[3]}$ $\quad$ L

$\dfrac{dL}{d\mathbf{z}^{[1]}}$ $\quad$ $\dfrac{dL}{d\mathbf{a}^{[1]}}$ $\quad$ $\dfrac{dL}{d\mathbf{z}^{[2]}}$ $\quad$ $\dfrac{dL}{d\mathbf{a}^{[2]}}$ $\quad$ $\dfrac{dL}{dz^{[3]}}$

$\dfrac{dL}{d\mathbf{W}^{[1]}}$ $\dfrac{dL}{d\mathbf{b}^{[1]}}$ $\quad$ $\dfrac{dL}{d\mathbf{W}^{[2]}}$ $\dfrac{dL}{d\mathbf{b}^{[2]}}$ $\quad$ $\dfrac{dL}{d\mathbf{W}^{[3]}}$ $\dfrac{dL}{db^{[3]}}$

backpropagation direction

$$\frac{\partial \boldsymbol{a}^{[2]}}{\partial \mathbf{z}^{[2]}} = \begin{bmatrix} \dfrac{\partial a^{[2]}}{\partial z_1^{[2]}} \\ \dfrac{\partial a^{[2]}}{\partial z_2^{[2]}} \\ \dfrac{\partial a^{[2]}}{\partial z_3^{[2]}} \end{bmatrix} = \begin{bmatrix} a_1^{[2]} \left( 1 - a_1^{[2]} \right) \\ a_2^{[2]} \left( 1 - a_2^{[2]} \right) \\ a_3^{[2]} \left( 1 - a_3^{[2]} \right) \end{bmatrix}$$

$$= \boldsymbol{a}^{[2]} \circ \left( 1 - \boldsymbol{a}^{[2]} \right)$$

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Lan truyền ngược

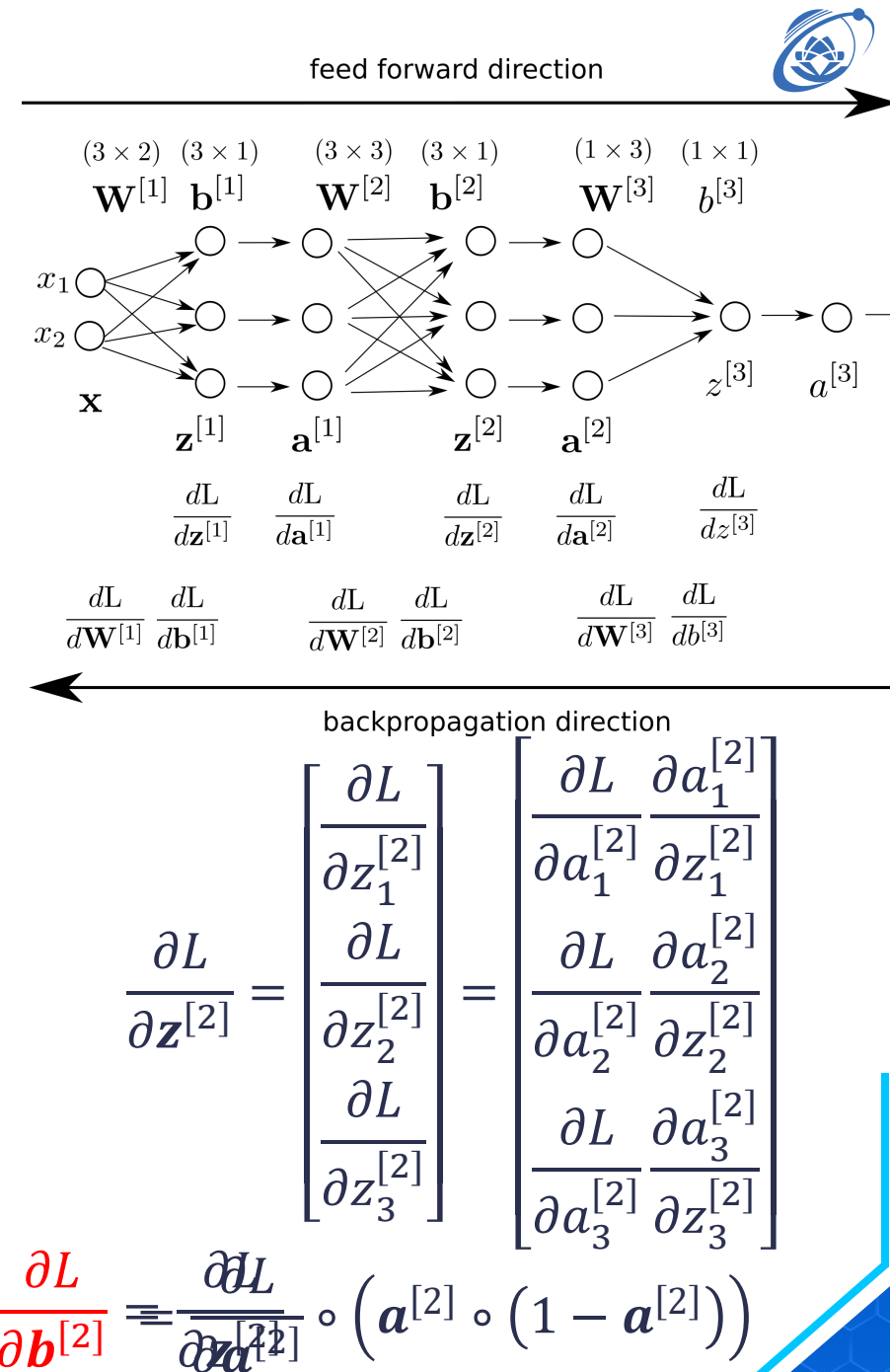$$\boldsymbol{z}^{[2]} = W^{[2]}\boldsymbol{a}^{[1]} + \boldsymbol{b}^{[2]}$$

$$\begin{bmatrix} z_1^{[2]} \\ z_2^{[2]} \\ z_3^{[2]} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[2]} & w_{1,2}^{[2]} & w_{1,3}^{[2]} \\ w_{2,1}^{[2]} & w_{2,2}^{[2]} & w_{2,3}^{[2]} \\ w_{3,1}^{[2]} & w_{3,2}^{[2]} & w_{3,3}^{[2]} \end{bmatrix} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix} + \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \end{bmatrix}$$



$$\frac{\partial L}{\partial W^{[2]}} = \begin{bmatrix} \frac{\partial L}{\partial w_{1,1}^{[2]}} & \frac{\partial L}{\partial w_{1,2}^{[2]}} & \frac{\partial L}{\partial w_{1,3}^{[2]}} \\ \frac{\partial L}{\partial w_{2,1}^{[2]}} & \frac{\partial L}{\partial w_{2,2}^{[2]}} & \frac{\partial L}{\partial w_{2,3}^{[2]}} \\ \frac{\partial L}{\partial w_{3,1}^{[2]}} & \frac{\partial L}{\partial w_{3,2}^{[2]}} & \frac{\partial L}{\partial w_{3,3}^{[2]}} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial z_1^{[2]}}\frac{\partial z_1^{[2]}}{\partial w_{1,1}^{[2]}} & \frac{\partial L}{\partial z_1^{[2]}}\frac{\partial z_1^{[2]}}{\partial w_{1,2}^{[2]}} & \frac{\partial L}{\partial z_1^{[2]}}\frac{\partial z_1^{[2]}}{\partial w_{1,3}^{[2]}} \\ \frac{\partial L}{\partial z_2^{[2]}}\frac{\partial z_2^{[2]}}{\partial w_{2,1}^{[2]}} & \frac{\partial L}{\partial z_2^{[2]}}\frac{\partial z_2^{[2]}}{\partial w_{2,2}^{[2]}} & \frac{\partial L}{\partial z_2^{[2]}}\frac{\partial z_2^{[2]}}{\partial w_{2,3}^{[2]}} \\ \frac{\partial L}{\partial z_3^{[2]}}\frac{\partial z_3^{[2]}}{\partial w_{3,1}^{[2]}} & \frac{\partial L}{\partial z_3^{[2]}}\frac{\partial z_3^{[2]}}{\partial w_{3,2}^{[2]}} & \frac{\partial L}{\partial z_3^{[2]}}\frac{\partial z_3^{[2]}}{\partial w_{3,3}^{[2]}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial L}{\partial z_1^{[2]}}a_1^{[1]} & \frac{\partial L}{\partial z_1^{[2]}}a_2^{[1]} & \frac{\partial L}{\partial z_1^{[2]}}a_3^{[1]} \\ \frac{\partial L}{\partial z_2^{[2]}}a_1^{[1]} & \frac{\partial L}{\partial z_2^{[2]}}a_2^{[1]} & \frac{\partial L}{\partial z_2^{[2]}}a_3^{[1]} \\ \frac{\partial L}{\partial z_3^{[2]}}a_1^{[1]} & \frac{\partial L}{\partial z_3^{[2]}}a_2^{[1]} & \frac{\partial L}{\partial z_3^{[2]}}a_3^{[1]} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial z_1^{[2]}} \\ \frac{\partial L}{\partial z_2^{[2]}} \\ \frac{\partial L}{\partial z_3^{[2]}} \end{bmatrix} \begin{bmatrix} a_1^{[1]} & a_2^{[1]} & a_3^{[1]} \end{bmatrix}$$

$$= \frac{\partial L}{\partial \boldsymbol{z}^{[2]}}\left(\boldsymbol{a}^{[1]}\right)^T$$

$$\frac{\partial L}{\partial \boldsymbol{z}^{[2]}} = \begin{bmatrix} \frac{\partial L}{\partial z_1^{[2]}} \\ \frac{\partial L}{\partial z_2^{[2]}} \\ \frac{\partial L}{\partial z_3^{[2]}} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial a_1^{[2]}}\frac{\partial a_1^{[2]}}{\partial z_1^{[2]}} \\ \frac{\partial L}{\partial a_2^{[2]}}\frac{\partial a_2^{[2]}}{\partial z_2^{[2]}} \\ \frac{\partial L}{\partial a_3^{[2]}}\frac{\partial a_3^{[2]}}{\partial z_3^{[2]}} \end{bmatrix}$$

$$\frac{\partial L}{\partial \boldsymbol{b}^{[2]}} = \frac{\partial L}{\partial \boldsymbol{z}^{[2]}} = \frac{\partial L}{\partial \boldsymbol{a}^{[2]}} \circ \left(\boldsymbol{a}^{[2]} \circ \left(1 - \boldsymbol{a}^{[2]}\right)\right)$$

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Lan truyền ngược

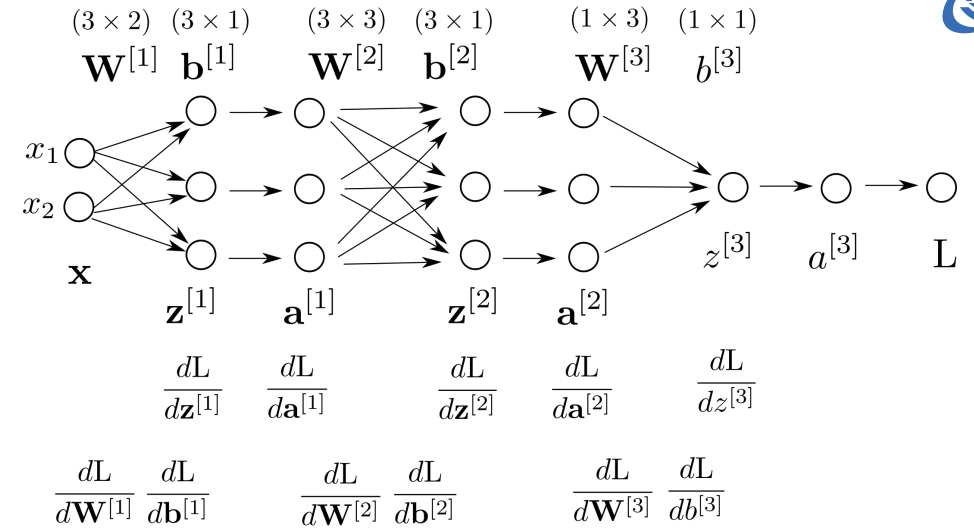$$\mathbf{z}^{[2]} = W^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

$$
\begin{bmatrix} z_1^{[2]} \\ z_2^{[2]} \\ z_3^{[2]} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[2]} & w_{1,2}^{[2]} & w_{1,3}^{[2]} \\ w_{2,1}^{[2]} & w_{2,2}^{[2]} & w_{2,3}^{[2]} \\ w_{3,1}^{[2]} & w_{3,2}^{[2]} & w_{3,3}^{[2]} \end{bmatrix} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix} + \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \end{bmatrix}
$$

$$
\frac{\partial \mathbf{z}^{[2]}}{\partial \mathbf{a}^{[1]}} = \begin{bmatrix} \dfrac{\partial \mathbf{z}^{[2]}}{\partial a_1^{[1]}} \\[8pt] \dfrac{\partial \mathbf{z}^{[2]}}{\partial a_2^{[1]}} \\[8pt] \dfrac{\partial \mathbf{z}^{[2]}}{\partial a_3^{[1]}} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial z_1^{[2]}}{\partial a_1^{[1]}} & \dfrac{\partial z_2^{[2]}}{\partial a_1^{[1]}} & \dfrac{\partial z_3^{[2]}}{\partial a_1^{[1]}} \\[8pt] \dfrac{\partial z_1^{[2]}}{\partial a_2^{[1]}} & \dfrac{\partial z_2^{[2]}}{\partial a_2^{[1]}} & \dfrac{\partial z_3^{[2]}}{\partial a_2^{[1]}} \\[8pt] \dfrac{\partial z_1^{[2]}}{\partial a_3^{[1]}} & \dfrac{\partial z_2^{[2]}}{\partial a_3^{[1]}} & \dfrac{\partial z_3^{[2]}}{\partial a_3^{[1]}} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[2]} & w_{2,1}^{[2]} & w_{3,1}^{[2]} \\ w_{1,2}^{[2]} & w_{2,2}^{[2]} & w_{3,2}^{[2]} \\ w_{1,3}^{[2]} & w_{2,3}^{[2]} & w_{3,3}^{[2]} \end{bmatrix} = \left(W^{[2]}\right)^T
$$

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

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

$$
\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[1]} & w_{1,2}^{[1]} \\ w_{2,1}^{[1]} & w_{2,2}^{[1]} \\ w_{3,1}^{[1]} & w_{3,2}^{[1]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \end{bmatrix}
$$



Tương tự, ta cũng tính được

$$\frac{\partial L}{\partial \mathbf{z}^{[1]}} = \frac{\partial L}{\partial \mathbf{a}^{[1]}} \circ \left(\mathbf{a}^{[1]} \circ \left(1 - \mathbf{a}^{[1]}\right)\right)$$

$$\frac{\partial L}{\partial W^{[1]}} = \frac{\partial L}{\partial \mathbf{z}^{[1]}} \left(\mathbf{a}^{[0]}\right)^T = \frac{\partial L}{\partial \mathbf{z}^{[1]}} \mathbf{x}^T$$

$$\frac{\partial L}{\partial \mathbf{b}^{[1]}} = \frac{\partial L}{\partial \mathbf{z}^{[1]}}$$

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# CÀI ĐẶT & PHÂN TÍCH

## IMPLEMENTATION & ANALYSIS

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Gradient Descent

**Hàm mất mát:**

$$L = -\left(y^{(i)}\log\left(a^{[3](i)}\right) + \left(1 - y^{(i)}\right)\log\left(1 - a^{[3](i)}\right)\right)$$

**Đạo hàm:**

$$\frac{\partial L}{\partial z^{[3]}} = a^{[3]} - y$$

$$\frac{\partial L}{\partial b^{[3]}} = \frac{\partial L}{\partial z^{[3]}}$$

$$\frac{\partial L}{\partial W^{[3]}} = \frac{\partial L}{\partial z^{[3]}}\left(\boldsymbol{a}^{[2]}\right)^T$$

$$\frac{\partial L}{\partial \boldsymbol{a}^{[2]}} = \frac{\partial L}{\partial z^{[3]}}\frac{\partial z^{[3]}}{\partial \boldsymbol{a}^{[2]}} = \left(a^{[3]} - y\right)\left(W^{[3]}\right)^T$$

$$\frac{\partial L}{\partial \boldsymbol{z}^{[2]}} = \frac{\partial L}{\partial \boldsymbol{a}^{[2]}} \circ \left(\boldsymbol{a}^{[2]} \circ \left(1 - \boldsymbol{a}^{[2]}\right)\right)$$

$$\frac{\partial L}{\partial \boldsymbol{b}^{[2]}} = \frac{\partial L}{\partial \boldsymbol{z}^{[2]}}$$

$$\frac{\partial L}{\partial W^{[2]}} = \frac{\partial L}{\partial \boldsymbol{z}^{[2]}}\left(\boldsymbol{a}^{[1]}\right)^T$$

$$\frac{\partial L}{\partial \boldsymbol{a}^{[1]}} = \frac{\partial \boldsymbol{z}^{[2]}}{\partial \boldsymbol{a}^{[1]}}\frac{\partial L}{\partial \boldsymbol{z}^{[2]}} = \left(W^{[2]}\right)^T \frac{\partial L}{\partial \boldsymbol{z}^{[2]}}$$

$$\frac{\partial L}{\partial \boldsymbol{z}^{[1]}} = \frac{\partial L}{\partial \boldsymbol{a}^{[1]}} \circ \left(\boldsymbol{a}^{[1]} \circ \left(1 - \boldsymbol{a}^{[1]}\right)\right)$$

$$\frac{\partial L}{\partial \boldsymbol{b}^{[1]}} = \frac{\partial L}{\partial \boldsymbol{z}^{[1]}}$$

$$\frac{\partial L}{\partial W^{[1]}} = \frac{\partial L}{\partial \boldsymbol{z}^{[1]}}\boldsymbol{x}^T$$

```python
def get_loss(y, a):
  return -1 * (y * np.log(a) + (1-y) * np.log(1-a))

def get_gradients(z1, a1, z2, a2, z3, a3, x, y, W1, b1, W2, b2, W3, b3):
  dz3 = a3 - y   # dL/dz_3
  db3 = dz3      # dL/db_3

  dW3 = dz3 * a2.T   # dL/dW_3
  da2 = dz3 * W3.T

  dz2 = da2 * (a2 * (1-a2))
  db2 = dz2

  dW2 = np.matmul(dz2, a1.T)
  da1 = np.matmul(W2.T, dz2)

  dz1 = da1 * (a1 * (1-a1))
  db1 = dz1
  dW1 = np.matmul(dz1, x.T)

  return dW1, db1, dW2, db2, dW3, db3
```
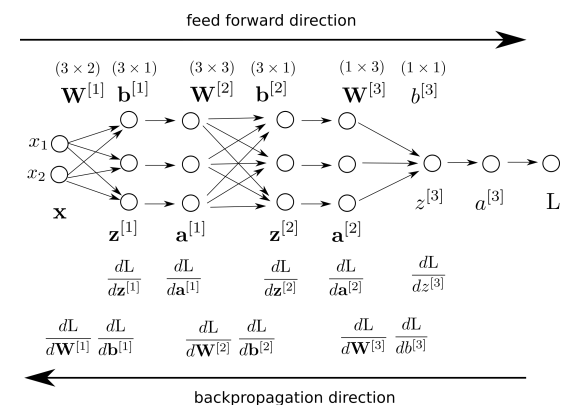
Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM

# Gradient Descent

```python
def gradient_descent(W1, b1, W2, b2, W3, b3, dW1, db1, dW2, db2, dW3, db3, alpha):
    W1 -= alpha
    b1 -= alpha
    W2 -= alpha
    b2 -= alpha
    W3 -= alpha
    b3 -= alpha

    return W1, b

def add_grad                    W1, db1,
    dW2, db2, dW3, db3):
    tdW1 += dW1
```
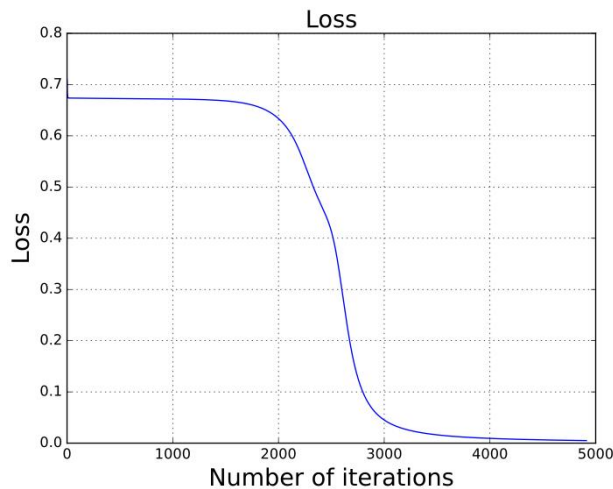


Loss



Decision boundary - Animated



Decision boundary - Contour plot

```python
alpha = 0.4
for i in range(20000):
    totalL = 0
    tdW1, tdb1, tdW2, tdb2, tdW3, tdb3 = get_zero_gradients(W1, b1, W2, b2, W3, b3)
    for j in range(X.shape[0]):
        x = X[j, :].reshape(2,1)
        z1, a1, z2, a2, z3, a3 = forward(x, W1, b1, W2, b2, W3, b3)
        L = (1.0 / 20) * get_loss(y[j], a3)
        totalL += L
    dW1, db1, dW2, db2, dW3, db3 = get_gradients (z1, a1, z2, a2, z3, a3, x, y[j], W1, b1, W2, b2, W3, b3)
    tdW1, tdb1, tdW2, tdb2, tdW3, tdb3 = add_gradients(tdW1, tdb1, tdW2, tdb2, tdW3, tdb3, dW1, db1, dW2, db2, dW3, db3)

    W1, b1, W2, b2, W3, b3 = gradient_descent(W1, b1, W2, b2, W3, b3, tdW1, tdb1, tdW2, tdb2, tdW3, tdb3, alpha)

    if totalL[0,0] < 0.005:
        break
```
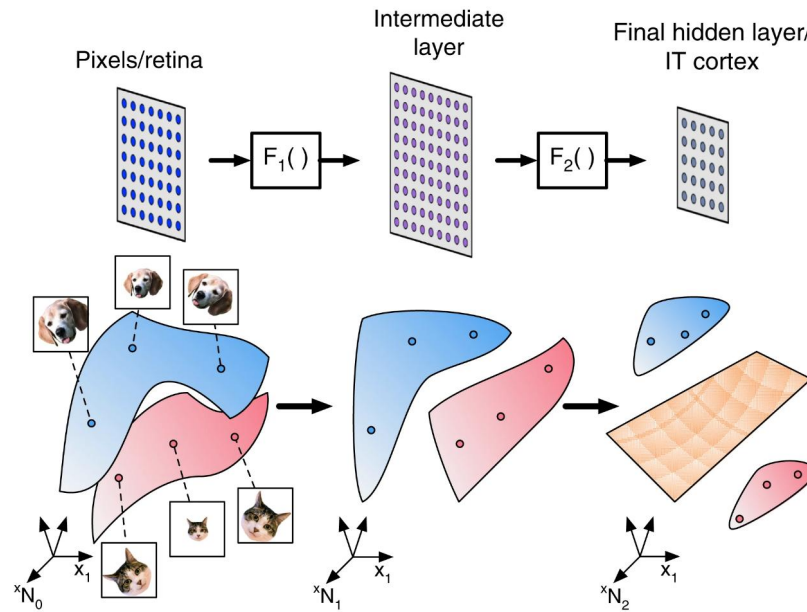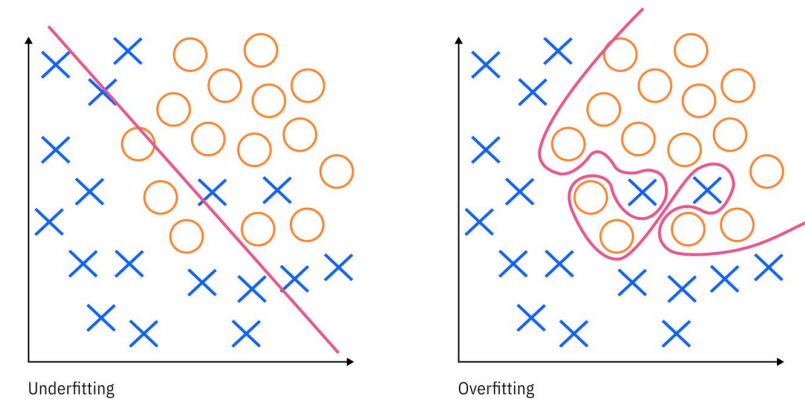
# Phân tích mạng nơ-rơn

- Mạng nơ-rơn nhân tạo (artificial neural network – ANN) với đủ độ sâu (số lớp) và đủ số lượng nơ-rơn ở mỗi lớp có thể tạo ra những biên quyết định (decision boundary) có độ phức tạp cao.



Cohen, U., Chung, S., Lee, D.D. *et al*. Separability and geometry of object manifolds in deep neural networks. *Nat Commun* **11**, 746 (2020). https://doi.org/10.1038/s41467-020-14578-5



Tim Mucci. https://www.ibm.com/think/topics/overfitting-vs-underfitting

- Mạng nơ-rơn có khả năng khớp rất tốt vào tập dữ liệu huấn luyện → rủi ro <span style="color:red">quá khớp (overfitting)</span>.
- Để khắc phục rủi ro quá khớp của mạng nơ-rơn, ta cần có tập dữ liệu huấn luyện đủ lớn.

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM