# Deep into Neural Network from Gradient Descent

# Gradient 梯度



https://www.youtube.com/watch?v=npkI19rcpdY

在向量微積分中，純量場的梯度是一個向量場。 純量場中某一點的梯度指向在這點純量場增長最快的方向（當然要比較的話必須固定方向的長度），梯度的絕對值是長度為1的方向中函數最大的增加率，也就是說 ∇ f.
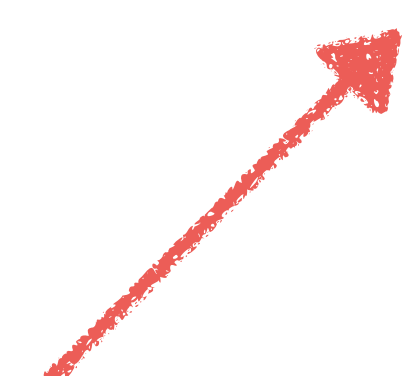
梯度- 維基百科，自由的百科全書 - Wikipedia
https://zh.wikipedia.org/zh-tw/梯度

$f = x^2 + y^2$

$grad = <2x, 2y>$

* 分別對 $x, y$ 偏微

* example:
when in point $<1, 2>$
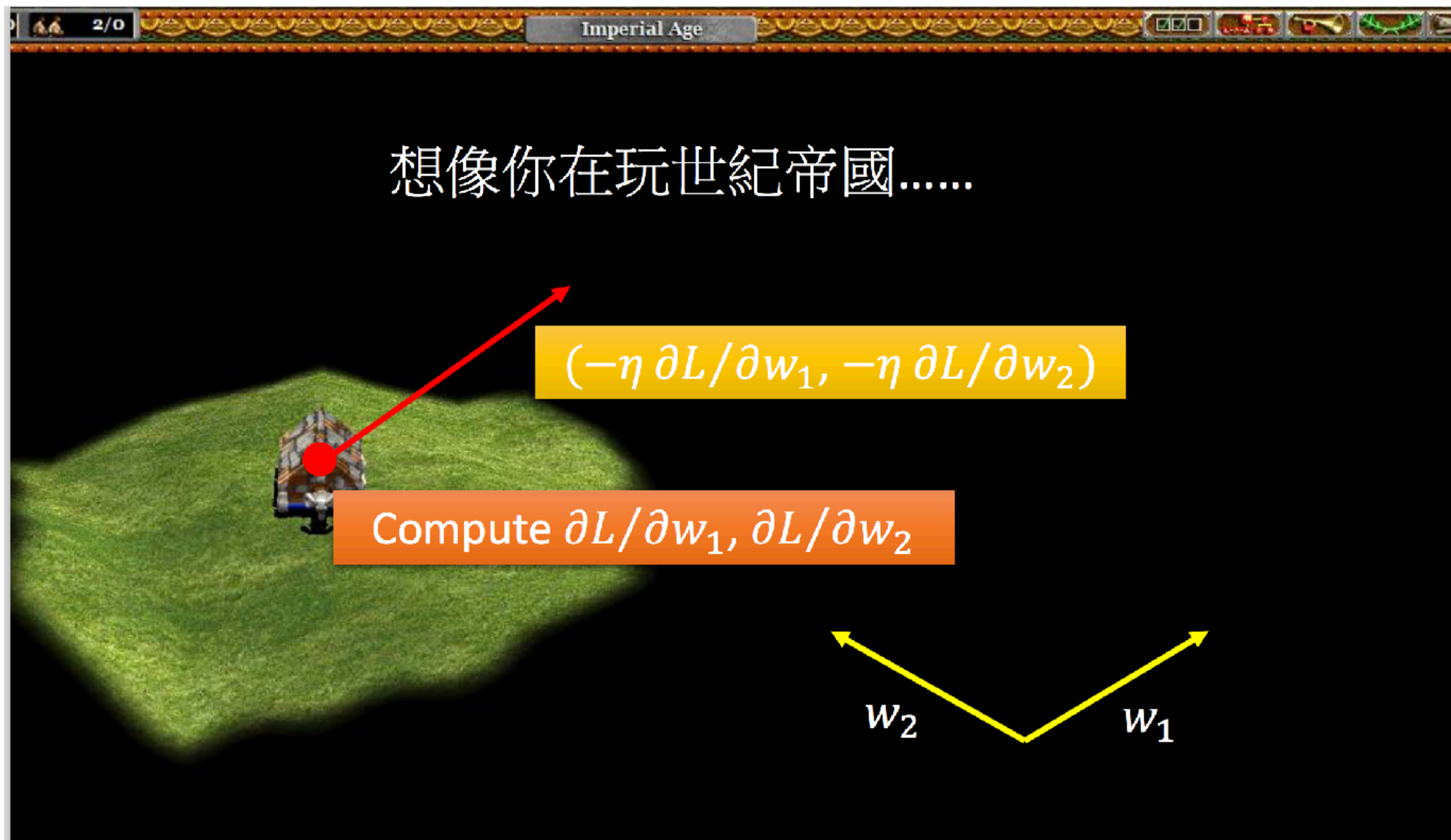
then grad = vector $<2, 4>$

就是在點(1, 2)時, 在<2, 4>的方向
坡度最陡 (最難爬的方向)

# Gradient Descent 梯度下降法

- 所以在X, Y點時, gradient的相反方向就是最好下降的方向

- 但跳太大步可能會跳到坑洞的另一方

- 跳太小步又很孬種(太慢了！)

- 跳的大小程度就是learning rate

- 每次都跳一樣大、越跳越小等等就是策略

- AdamOptimizer就是一種策略

- 每次前進前要用多少資料描繪當時可能的地形(1 vs mini-batch)



往當時下坡度最大的方向跳的程度

想像你在玩世紀帝國……

$(-\eta\, \partial L/\partial w_1, -\eta\, \partial L/\partial w_2)$

Compute $\partial L/\partial w_1, \partial L/\partial w_2$

$w_2$     $w_1$

一天搞懂深度學習

# 對方程式做Gradient Descent

$$\text{MSE}(X\ h_\theta) = \frac{1}{m}\sum_{i=1}^{m}(\theta^T \cdot x^{(i)} - y^{(i)})^2$$

w1*x(i)1 + w2*x(i)2 + ..wj*x(i)j + …

$$\frac{\partial}{\partial \theta_j}\text{MSE}(\theta) = \frac{2}{m}\sum_{i=1}^{m}(\theta^T \cdot x^{(i)} - y^{(i)})\, x_j^{(i)}$$

$$\nabla_\theta \text{MSE}(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_0}\text{MSE}(\theta) \\ \frac{\partial}{\partial \theta_1}\text{MSE}(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n}\text{MSE}(\theta) \end{pmatrix} = \frac{2}{m}X^T \cdot (X \cdot \theta - y)$$
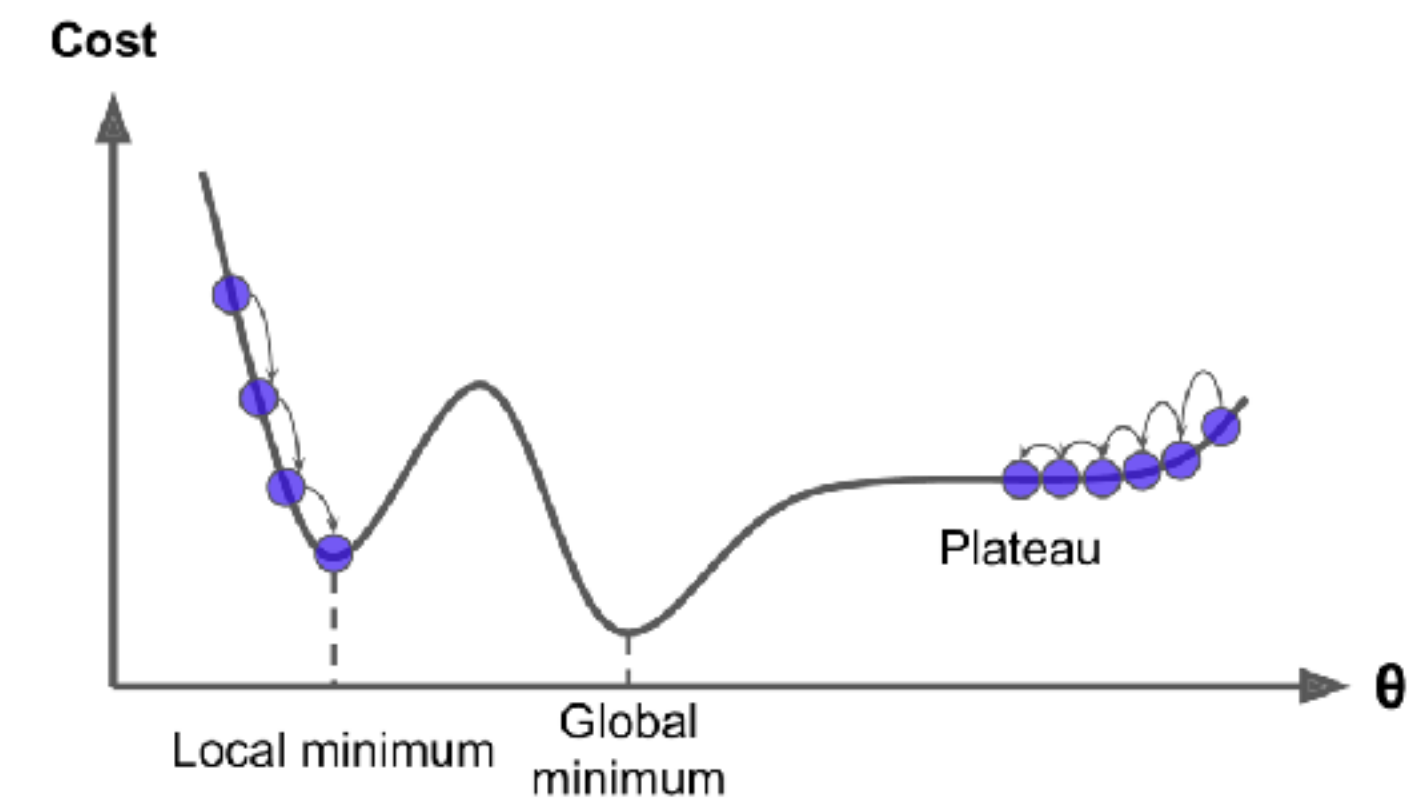
$$\theta^{(\text{next step})} = \theta - \eta\, \nabla_\theta \text{MSE}(\theta)$$

Cost

Plateau

Local minimum  Global minimum

θ

Figure 4-5. Gradient Descent pitfalls

$$X = \begin{bmatrix} [1, & 2, & 1] \\ [2, & 3, & 1] \end{bmatrix} \qquad \theta = \begin{bmatrix} [0.5] \\ [0.75] \\ [0.25] \end{bmatrix} \qquad y = \begin{bmatrix} [4] \\ [6] \end{bmatrix} \qquad X \cdot \theta = \begin{bmatrix} [2.25] \\ . \\ [3.5] \end{bmatrix}$$

$$(2, 3) \qquad\qquad (3, 1) \qquad\qquad (2, 1)$$

$$\text{MSE}(X\ h_\theta) = \frac{1}{m} \sum_{i=1}^{m} (\theta^T \cdot x^{(i)} - y^{(i)})^2$$

$$X \cdot \theta - y = error = \begin{bmatrix} [-1.75] \\ [-2.5] \end{bmatrix}$$

$$\nabla g = \left( \frac{2}{m} \sum_{i=1}^{m\ =\ \textcircled{2}} \left( \underbrace{x^{(i)} \cdot w}_{\downarrow \atop y\text{-pred}} - y \right) * x_1^{(i)} \right)$$

$$\underbrace{\qquad\qquad}_{error\ (都一樣)}$$

$$\cdots * x_2^{(i)}$$

$$\frac{\partial}{\partial \theta_j} \text{MSE}(\theta) = \frac{2}{m} \sum_{i=1}^{m} (\theta^T \cdot x^{(i)} - y^{(i)}) x_j^{(i)}$$

$$= \frac{2}{m} \begin{bmatrix} [-1.75 \times 1 + -2.5 \times 2] \\ [-1.75 \times 2 + -2.5 \times 3] \\ [-1.75 \times 1 + -2.5 \times 1] \end{bmatrix}$$

$$= \begin{bmatrix} [-6.75] \\ [-11] \\ [-4.25] \end{bmatrix}$$

$$\Rightarrow if\ \eta\ (learning\ rate) = 0.01$$

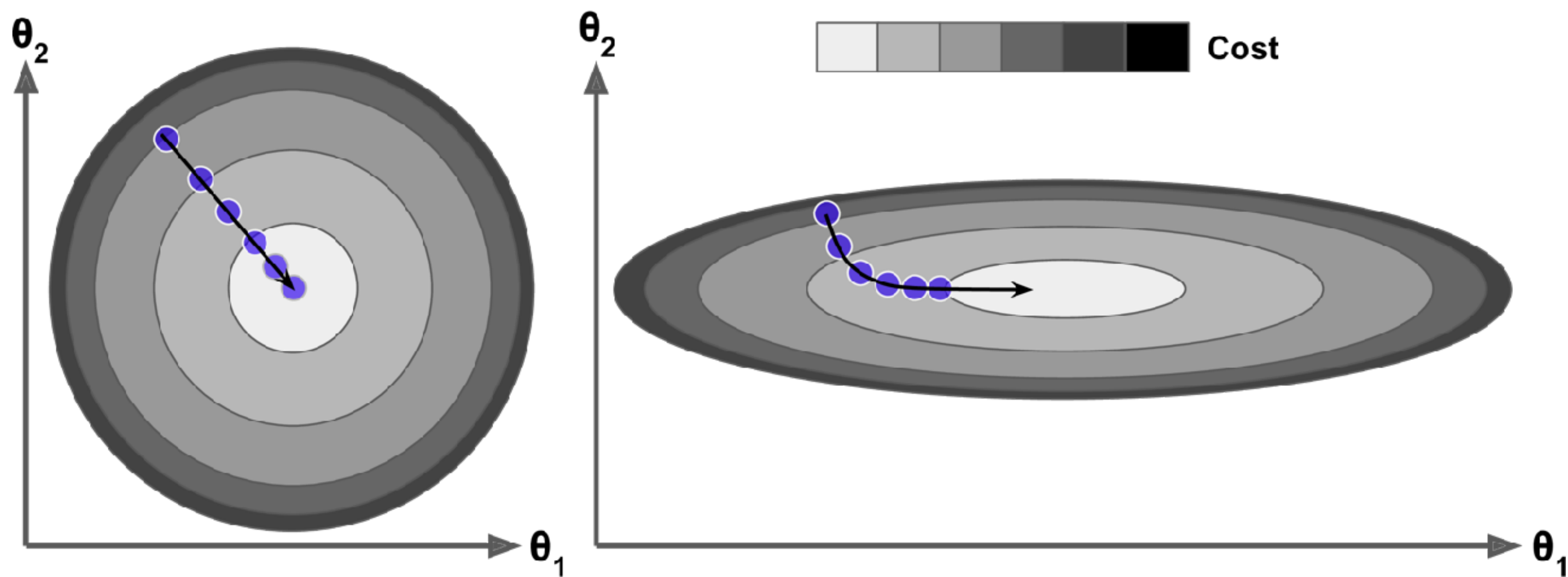$$\theta \leftarrow \theta - 0.01 \times \nabla g$$

$$= \begin{bmatrix} [0.5] \\ [0.75] \\ [0.25] \end{bmatrix} - 0.01 \begin{bmatrix} [-6.75] \\ [-11] \\ [-4.25] \end{bmatrix} = \begin{bmatrix} [0.5675] \\ [0.86] \\ [0.2925] \end{bmatrix}$$

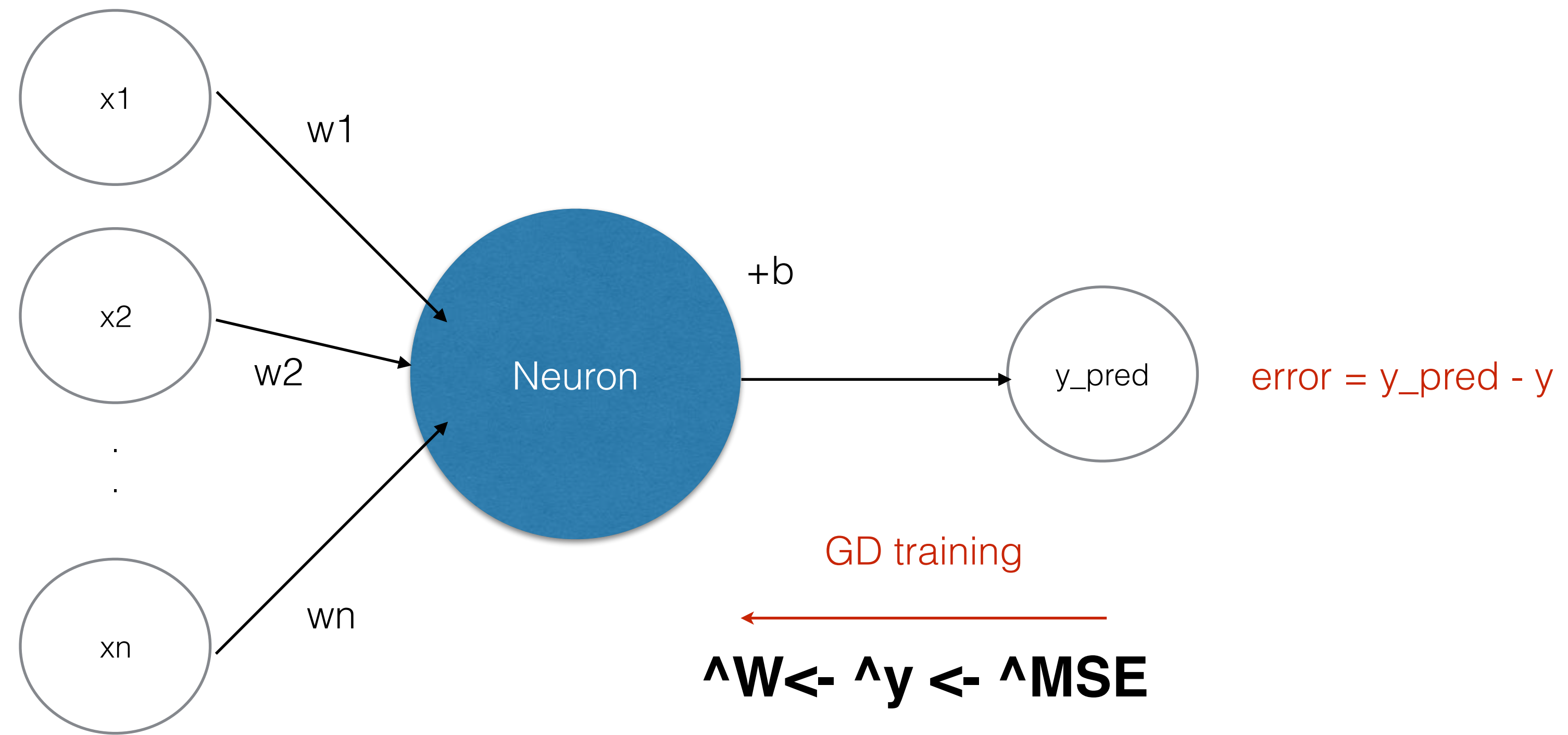$$\theta^{(next\ step)} = \theta - \eta \nabla_\theta \text{MSE}(\theta)$$

$$\theta\ next\ step \longrightarrow \cdots$$

# 做Gradient Descent前，要scaled features!

# 練習1: Linear Regression



最小化MSE訓練出W

# Logistic Regression



w1

x1

w2

x2

:
:

xn

wn

+b

Neuron

z

sigmoid

output

a

if >0.5
then 1
else 0

GD training

^W<-^z <- ^a(sigmoid 斜率) <- ^loss

$$H = -\sum_i (p_i . log_2(q_i))$$

-5*5 = -25
-4*6 = -24
-3*7 = -21

$$J(\theta) = -\frac{1}{m}\sum_{i=1}^{m}[y^{(i)}log(\hat{p}^{(i)}) + (1 - y^{(i)})log(1 - \hat{p}^{(i)})]$$

## 最小化Cross Entropy訓練出W

# Choosing Proper Loss

When using softmax output layer, choose cross entropy



Total Loss

Cross Entropy

Square Error

http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf

# 練習2: Softmax classification

$$P(y = j | \mathbf{x}) = \frac{e^{\mathbf{x}^\mathsf{T} \mathbf{w}_j}}{\sum_{k=1}^{K} e^{\mathbf{x}^\mathsf{T} \mathbf{w}_k}}$$

轉成機率



$$loss = -\sum_i \sum_c y_{c_i} \cdot log(y\_predicted_{c_i})$$

**最小化Cross Entropy訓練出W**

# Neural Network

**GD training**

**^z <- ^a <- ^w21 + ^w22 <- ^C**

+b

x1

x2

.
.

xn

$$P(y = j | \mathbf{x}) = \frac{e^{\mathbf{x}^\top \mathbf{w}_j}}{\sum_{k=1}^{K} e^{\mathbf{x}^\top \mathbf{w}_k}}$$

Softmax function
轉成機率

y1  0.3

y2  0.7

W1          W2

## 最小化Cross Entropy訓練出W2

## W1的改變會擴散到W1 -> Backpropagation

# 練習3: Build your neural network



input

4.3 sepal length

6.1 sepal width

4.2 petal width

5.3 petal length

**Dense Layer**

**relu**

w1

forward

**Dense Layer**

**relu**

w2

$$P(y = j|\mathbf{x}) = \frac{e^{\mathbf{x}^T\mathbf{w}_j}}{\sum_{k=1}^{K} e^{\mathbf{x}^T\mathbf{w}_k}}$$

**softmax**

**Dense Layer**

w3

**Dense Layer**

w4

output

0.2 Setosa

**0.7** Versicolour
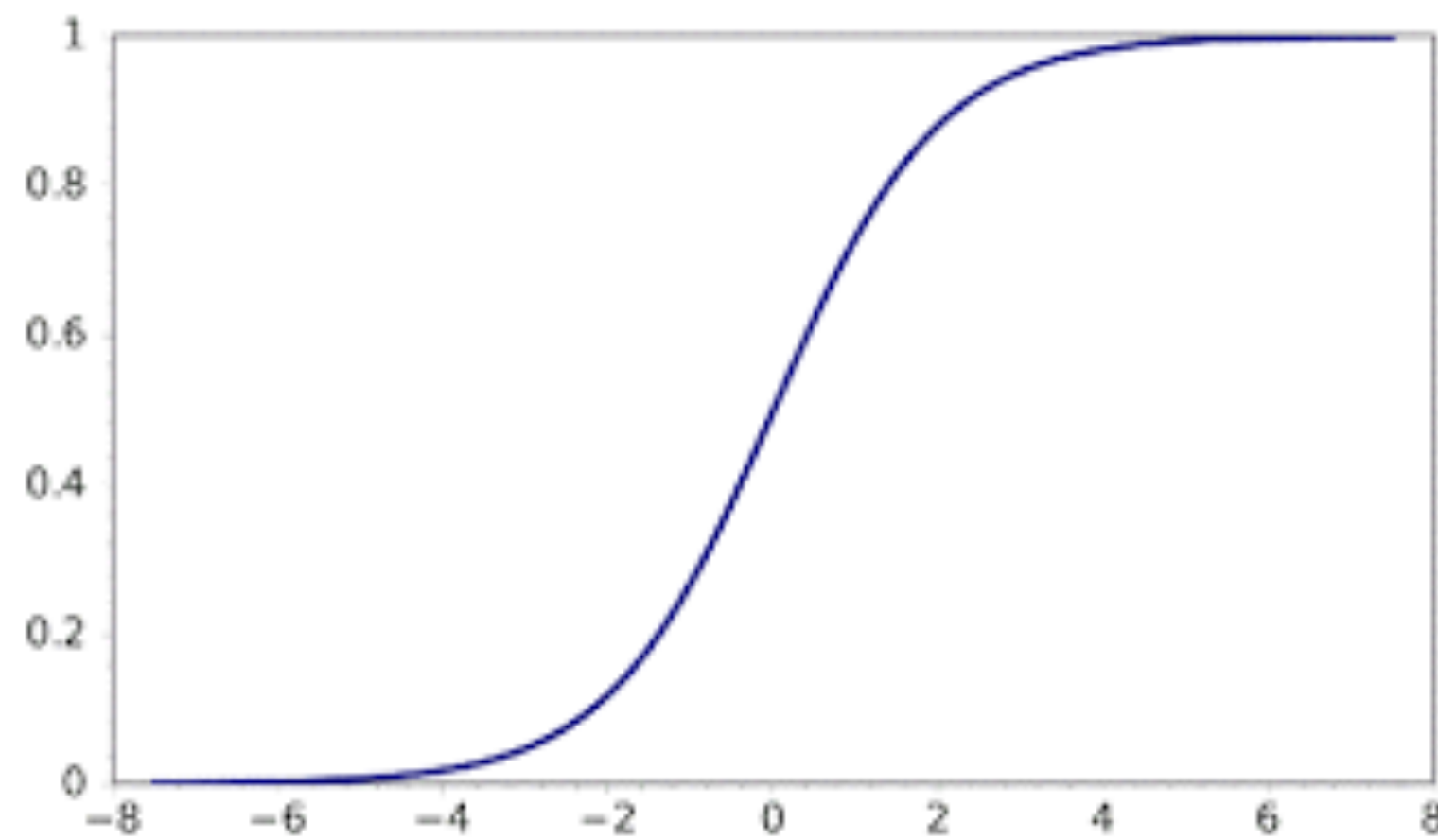
0.1 Virginica

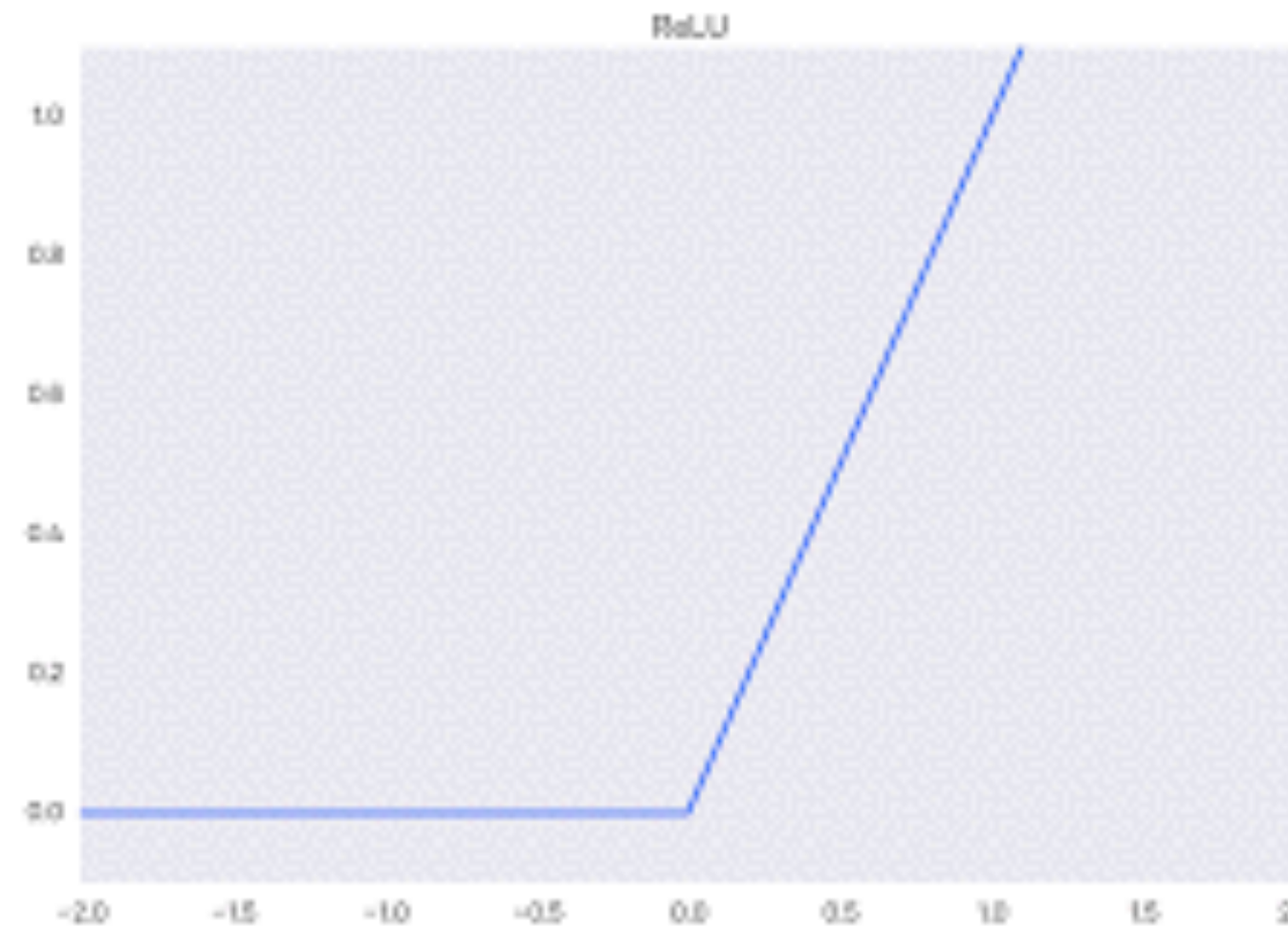SGD Backward Training

**loss function to minimize:**
**cross entropy**

## 深層網路的BP會層層往後更新

# Question

# Why from Sigmoid to Relu?

$$\text{ELU}_\alpha(z) = \begin{cases} \alpha(\exp(z) - 1) & \text{if } z < 0 \\ z & if z \geq 0 \end{cases}$$

ELU activation function ($\alpha = 1$)

**sigmoid**

**relu**

Leaky ReLU activation function

Leak

# Ways to improve DNN
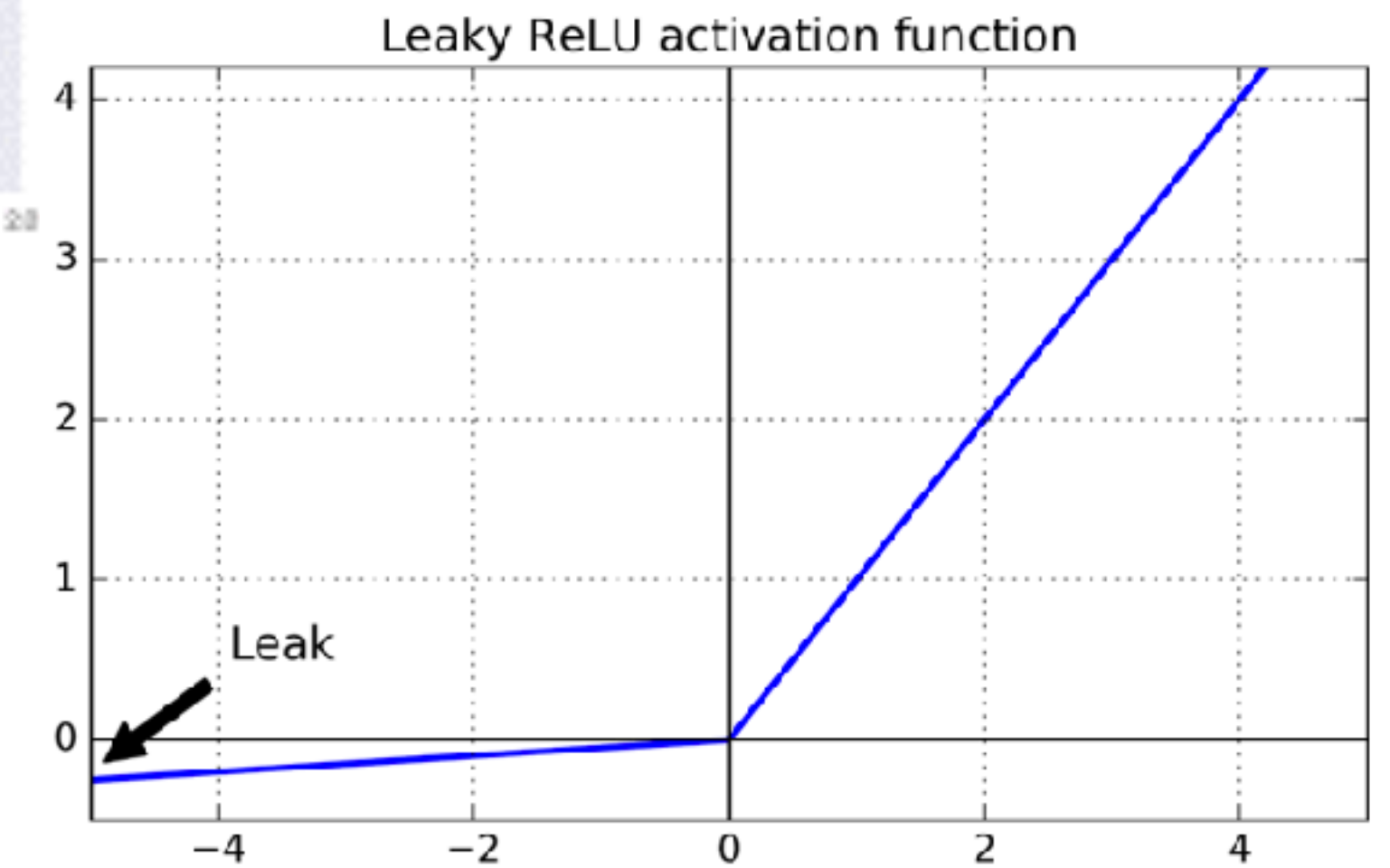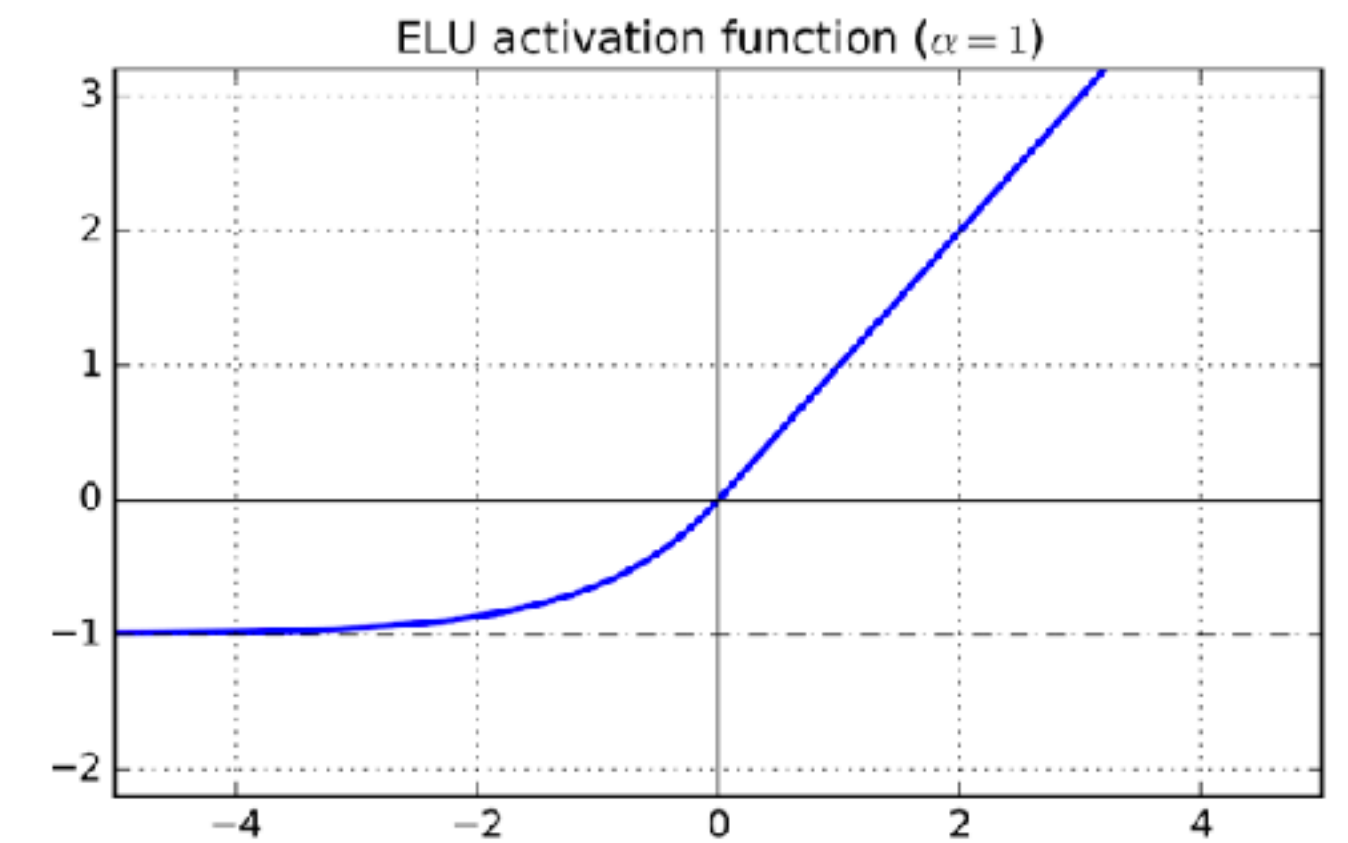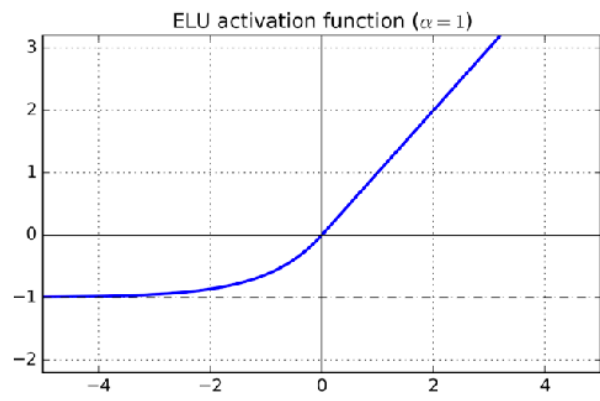
$$ELU_\alpha(z) = \begin{cases} \alpha(\exp(z) - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

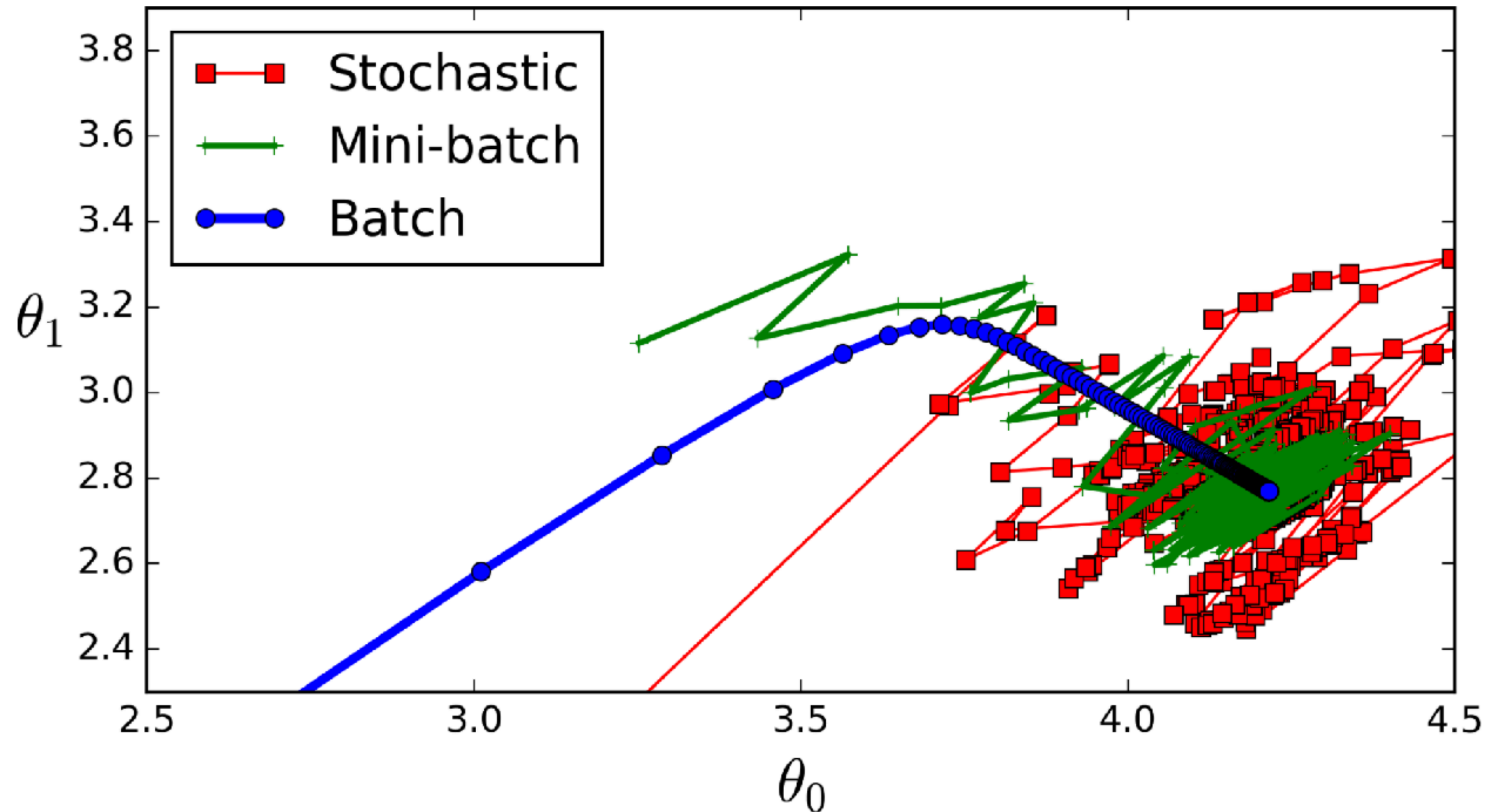ELU activation function ($\alpha = 1$)

- Weight Initialization - He initialization

- Activation function - ELU

- Normalization - Batch Normalization

- Regularization - Dropout

- Optimizer - Adam

- Check 一天搞懂深度學習

| Activation function | Uniform distribution [-r, r] | Normal distribution |
|---|---|---|
| Logistic | $r = \sqrt{\dfrac{6}{n_{inputs} + n_{outputs}}}$ | $\sigma = \sqrt{\dfrac{2}{n_{inputs} + n_{outputs}}}$ |
| Hyperbolic tangent | $r = 4\sqrt{\dfrac{6}{n_{inputs} + n_{outputs}}}$ | $\sigma = 4\sqrt{\dfrac{2}{n_{inputs} + n_{outputs}}}$ |
| ReLU (and its variants) | $r = \sqrt{2}\sqrt{\dfrac{6}{n_{inputs} + n_{outputs}}}$ | $\sigma = \sqrt{2}\sqrt{\dfrac{2}{n_{inputs} + n_{outputs}}}$ |

Equation 11-3. Batch Normalization algorithm

1. $\quad \mu_B = \dfrac{1}{m_B}\sum_{i=1}^{m_B} \mathbf{x}^{(i)}$

2. $\quad \sigma_B^2 = \dfrac{1}{m_B}\sum_{i=1}^{m_B}(\mathbf{x}^{(i)} - \mu_B)^2$

3. $\quad \mathbf{\hat{x}}^{(i)} = \dfrac{\mathbf{x}^{(i)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$

4. $\quad \mathbf{z}^{(i)} = \gamma\mathbf{\hat{x}}^{(i)} + \beta$

# Why using mini-batch?

# Too many matrix multiplication

$$\begin{bmatrix} y_1 \end{bmatrix} \qquad \left| \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \end{bmatrix} \begin{bmatrix} x_1 \end{bmatrix} \begin{bmatrix} b_1 \end{bmatrix} \right|$$

$$\begin{bmatrix} y_1 \end{bmatrix} \qquad \left| \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \end{bmatrix} \begin{bmatrix} x_1 \end{bmatrix} \begin{bmatrix} b_1 \end{bmatrix} \right|$$

$$\begin{bmatrix} y_1 \end{bmatrix} \qquad \left| \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \end{bmatrix} \begin{bmatrix} x_1 \end{bmatrix} \begin{bmatrix} b_1 \end{bmatrix} \right|$$

layer

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \mathrm{softmax} \left( \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \right)$$
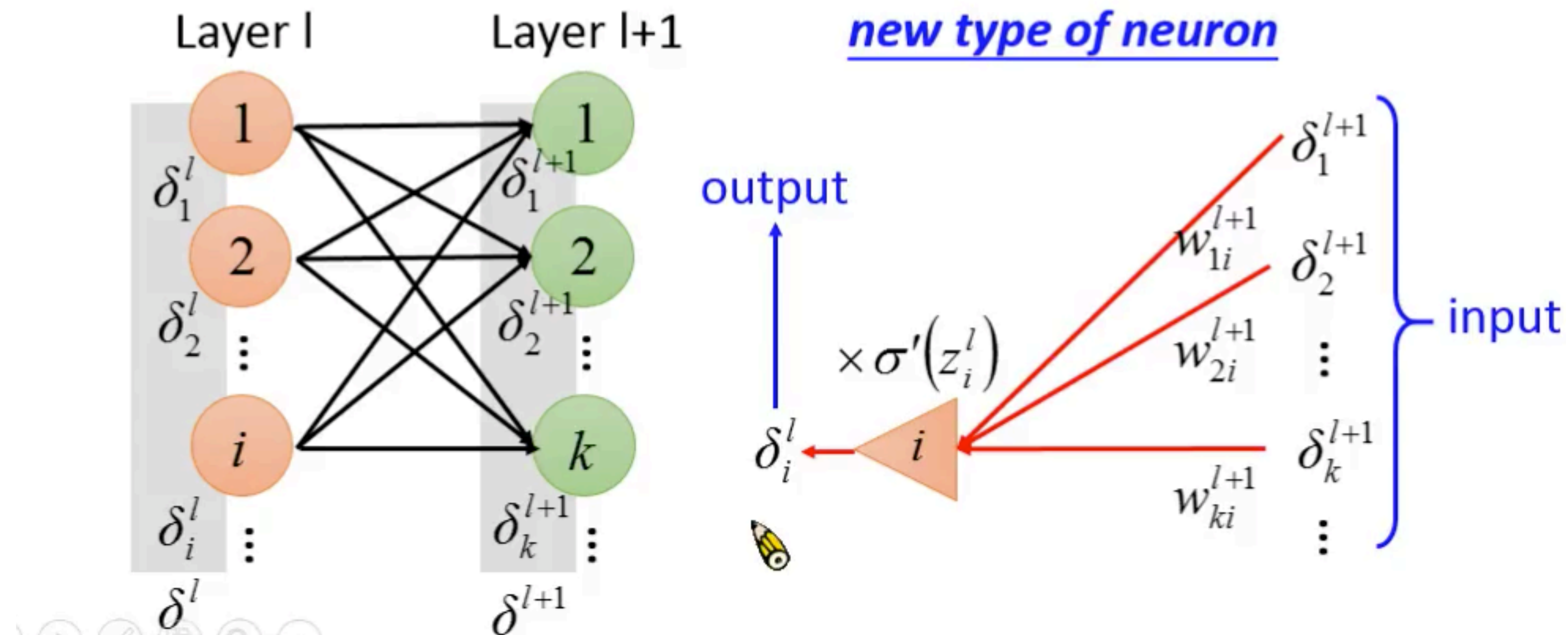
# What is your question?

# Others

# Backpropagation

$$\partial C^r / \partial w_{ij}^l \text{ - Second Term}$$

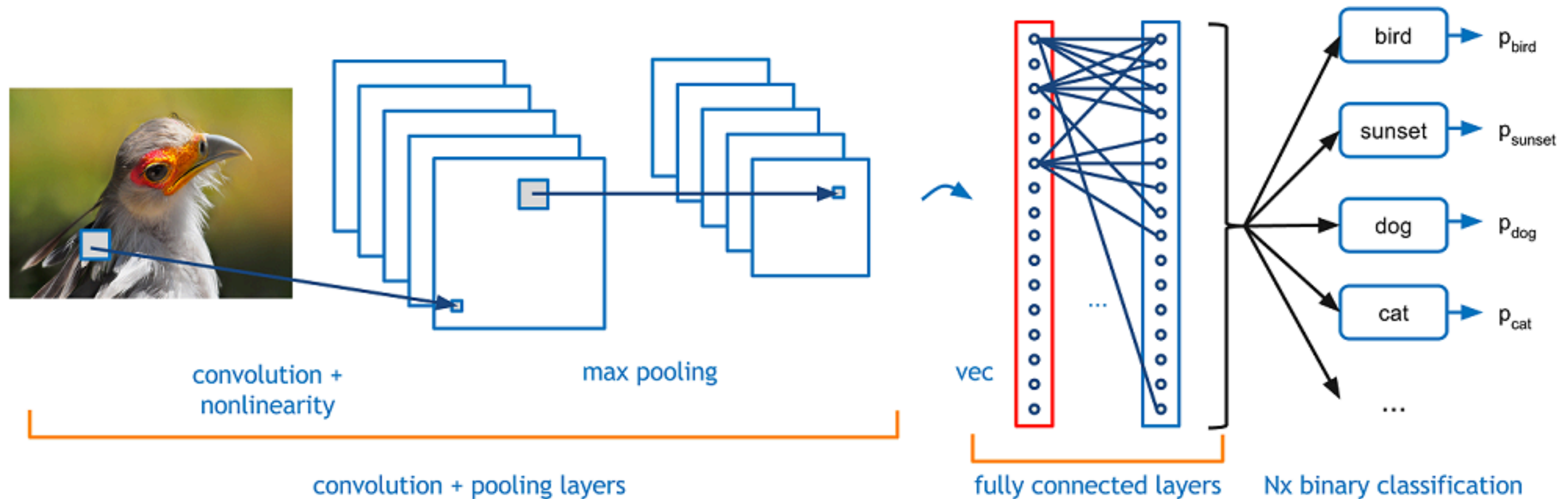$$\frac{\partial C^r}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \boxed{\frac{\partial C^r}{\partial z_i^l}} \rightarrow \delta_i^l \qquad\qquad \delta_i^l = \sigma'\left(z_i^l\right)\sum_k w_{ki}^{l+1}\delta_k^{l+1}$$

# Homework

- Try classify mnist with Keras

- Check the tensorflow tutorials - https://www.tensorflow.org/tutorials/mnist/beginners/

- Check the tensorflow template (in templates folder)

- Check CNN concept - https://www.youtube.com/watch?v=JiN9p5vWHDY

# Next time!



convolution +
nonlinearity

max pooling

vec

convolution + pooling layers

fully connected layers

Nx binary classification

bird → $p_{bird}$

sunset → $p_{sunset}$

dog → $p_{dog}$

cat → $p_{cat}$

...

**Convolution Neural Network**