

機器與深度學習基礎知識初探 -導傳遞

黃志勝 Chih-Sheng (Tommy) Huang

義隆電子人工智慧研發部

國立陽明交通大學AI學院合聘助理教授



基礎機器學習與深度學習

- 1. 神經網路如何運作
- 2. 梯度下降法
- 3. 神經網路如何利用導傳遞找解
- 4. 神經網路太深層的問題
- 5. Batch Normalization在幹什麼



基礎機器學習與深度學習

- 神經網路如何利用導傳遞找解
- 神經網路太深層的問題:

梯度消失問題(Vanishing gradient)/梯度爆炸問題 (exploding gradient problem) 。

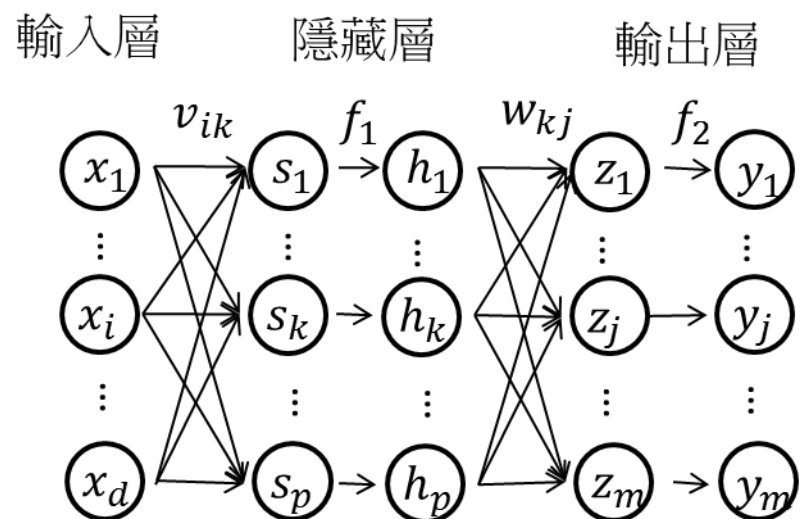
Activation Function為什麼要採用ReLU，而不是用Sigmoid。

Residual block克服神經網路不能太深層的問題。

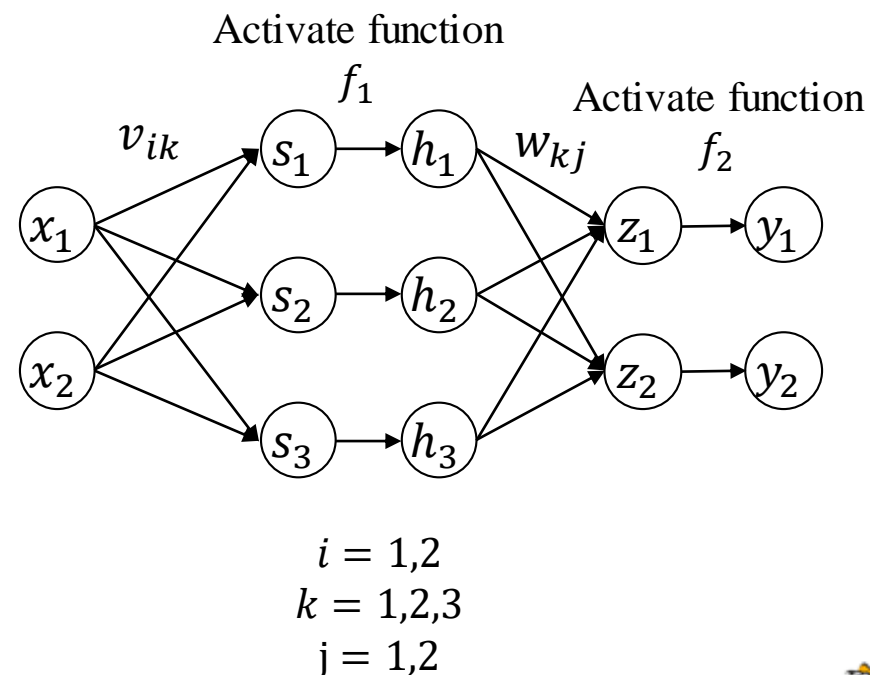
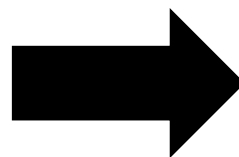


神經網路如何利用導傳遞找解

- 下圖為一個三層的MLP

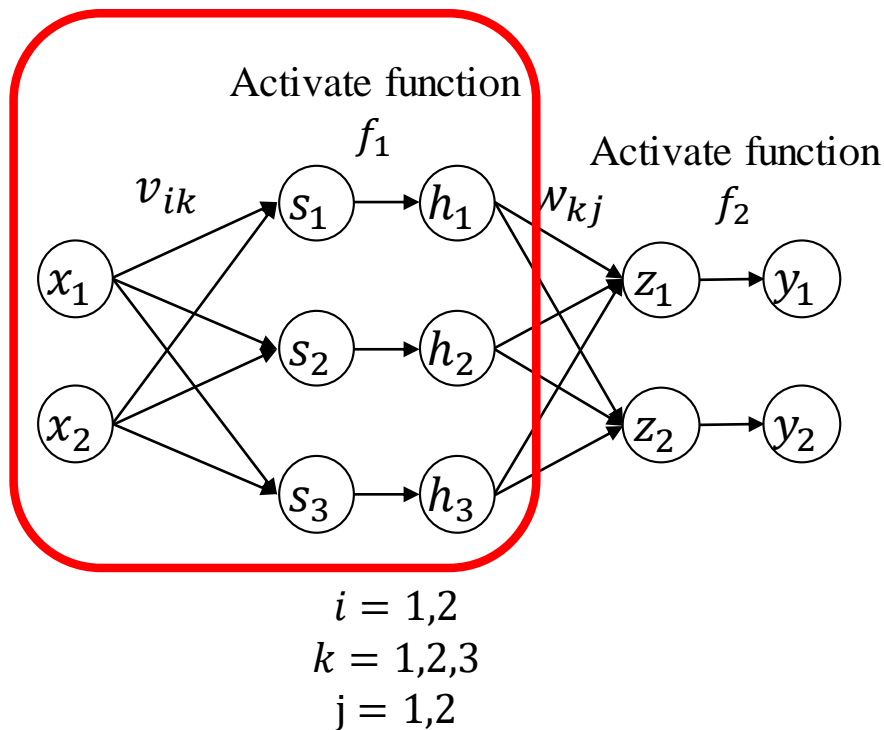


簡化



Forward propagation

輸入層→隱藏層



$$\begin{aligned} s_1 &= v_{11}x_1 + v_{21}x_2 = \mathbf{v}_1^T \mathbf{x} \\ s_2 &= v_{12}x_1 + v_{22}x_2 = \mathbf{v}_2^T \mathbf{x} \\ s_3 &= v_{13}x_1 + v_{23}x_2 = \mathbf{v}_3^T \mathbf{x} \end{aligned}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \mathbf{v}_k = \begin{bmatrix} v_{1k} \\ v_{2k} \end{bmatrix}, k = 1, 2, 3$$



$$s_k = \sum_{i=1}^2 v_{ik}x_i$$

$$\mathbf{s} = \mathbf{v}^T \mathbf{x}$$

$$\mathbf{s} = \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix}$$

$$\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \mathbf{v}_3 \end{bmatrix}$$

$$= \begin{bmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \end{bmatrix}$$

Activate function

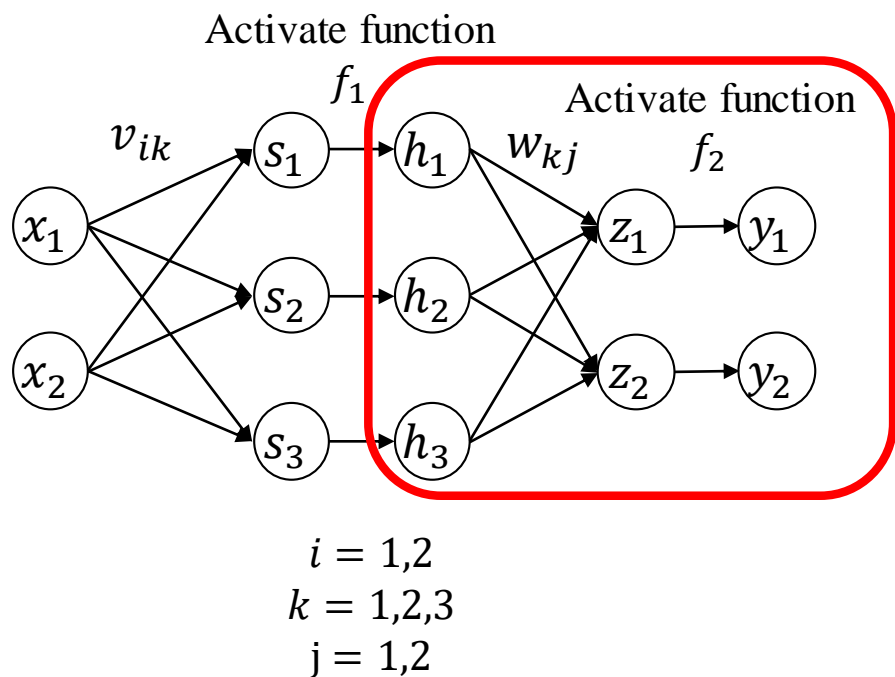
$$\mathbf{h} = f_1(\mathbf{s})$$

$$\mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix}$$



Forward propagation

隱藏層→輸出層



$$z_1 = w_{11}h_1 + w_{21}h_2 + w_{31}h_3 = \mathbf{w}_1^T \mathbf{h}$$

$$z_2 = w_{12}h_1 + w_{22}h_2 + w_{32}h_3 = \mathbf{w}_2^T \mathbf{h}$$

$$\mathbf{w}_j = \begin{bmatrix} w_{1j} \\ w_{2j} \\ w_{3j} \end{bmatrix}, j = 1, 2$$

Activate function

$$\mathbf{z} = \mathbf{w}^T \mathbf{h}$$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$\mathbf{w} = [\mathbf{w}_1 \quad \mathbf{w}_2]$$

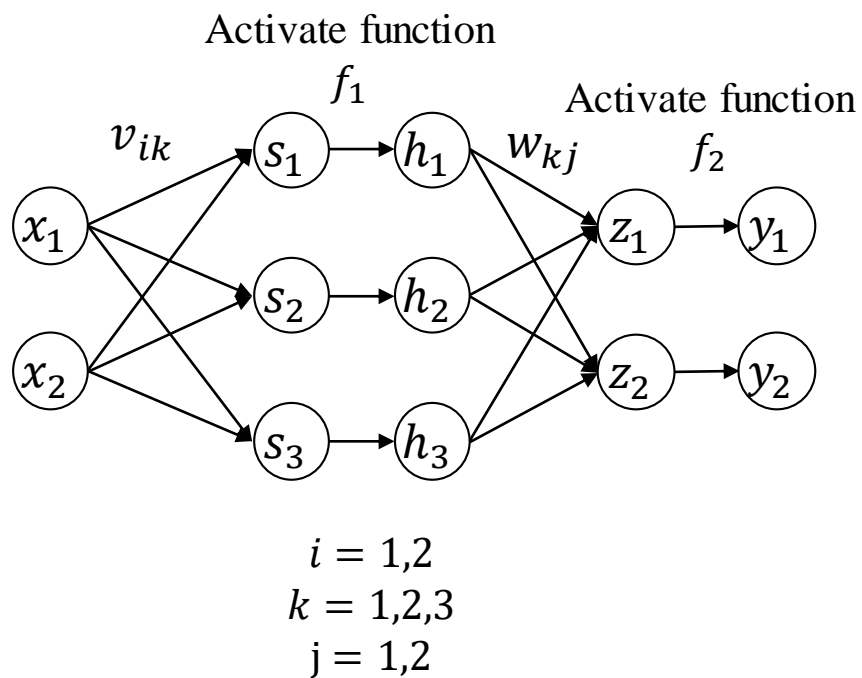
$$= \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

$$\mathbf{y} = f_2(\mathbf{z})$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$



Forward propagation



輸入層→隱藏層

$$\begin{aligned} \mathbf{s} &= \mathbf{v}^T \mathbf{x} \\ \mathbf{h} &= f_1(\mathbf{s}) \end{aligned}$$

隱藏層→輸出層

$$\begin{aligned} \mathbf{z} &= \mathbf{w}^T \mathbf{h} \\ \mathbf{y} &= f_2(\mathbf{z}) \end{aligned}$$

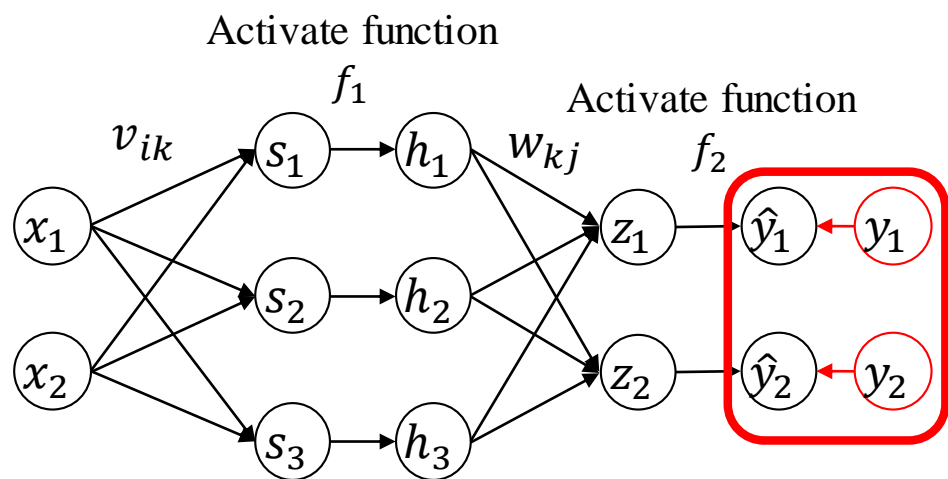
||

$$\mathbf{y} = f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x}))$$

參數為 \mathbf{v} 和 \mathbf{w} 。



Back-propagation



輸入層→隱藏層

$$\mathbf{s} = \mathbf{v}^T \mathbf{x}$$

$$\mathbf{h} = f_1(\mathbf{s})$$

隱藏層→輸出層

$$\mathbf{z} = \mathbf{w}^T \mathbf{h}$$

$$\mathbf{y} = f_2(\mathbf{z})$$

$$\hat{\mathbf{y}} = f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x})) \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

我們用SSE來當作loss function

$$\text{loss}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}})$$

利用gradient descent找最佳參數解(參數只有 \mathbf{w} 和 \mathbf{v})

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \alpha \nabla \mathbf{w}^{(t)}$$

$$\mathbf{v}^{(t+1)} = \mathbf{v}^{(t)} - \alpha \nabla \mathbf{v}^{(t)}$$

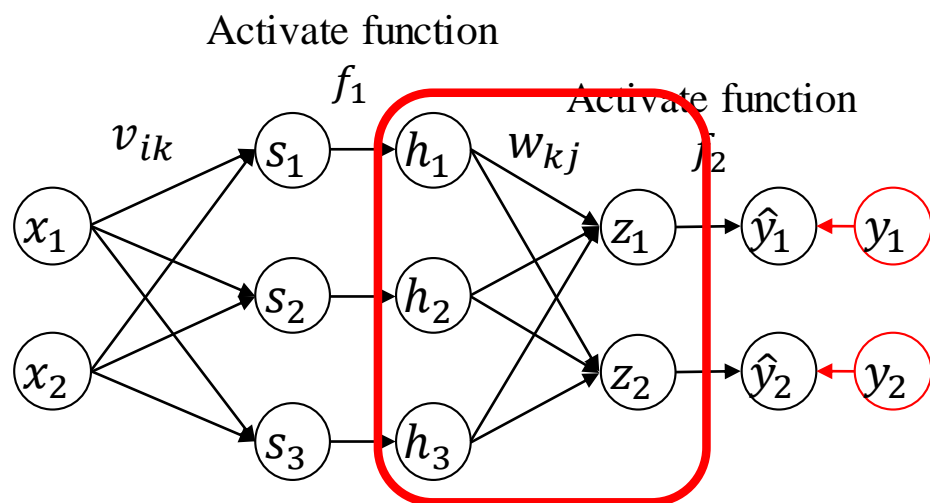
$\nabla \mathbf{w} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{w}}$ 輸出層→隱藏層參數為 \mathbf{w} ，對loss進行 \mathbf{w} 的偏微分

$\nabla \mathbf{v} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{v}}$ 隱藏層→輸入層參數為 \mathbf{v} ，對loss進行 \mathbf{v} 的偏微分



Back-propagation

輸出層→隱藏層



隱藏層→輸出層

$$\mathbf{z} = \mathbf{w}^T \mathbf{h}$$

$$\mathbf{y} = f_2(\mathbf{z})$$

參數為 \mathbf{w} ，因此我們對loss進行 \mathbf{w} 的偏微分

$$\begin{aligned} \text{loss}(\mathbf{y}, \hat{\mathbf{y}}) &= \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{z}))^T (\mathbf{y} - f_2(\mathbf{z})) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T \mathbf{h}))^T (\mathbf{y} - f_2(\mathbf{w}^T \mathbf{h})) \end{aligned}$$

Chain rule

$$\Delta \mathbf{w} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{w}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{w}}$$

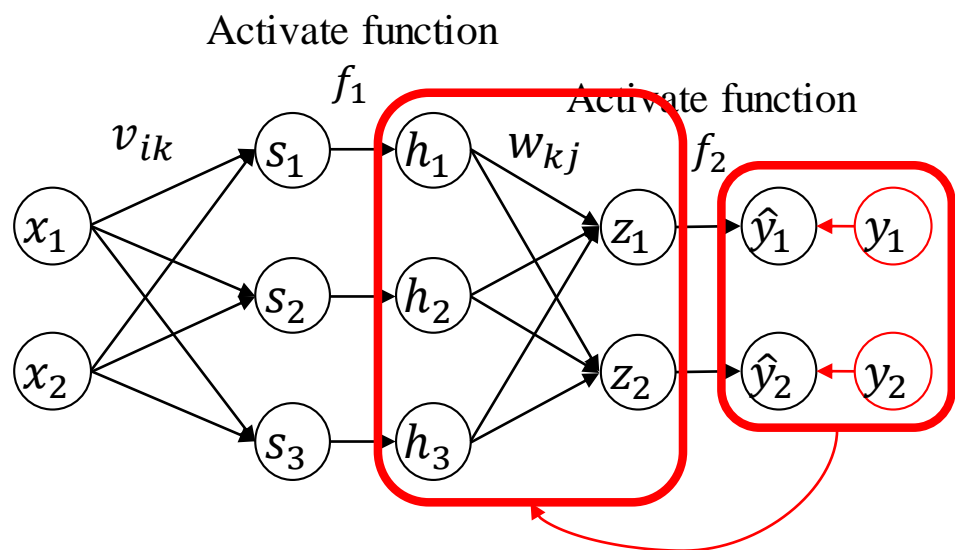
$$\begin{aligned} \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}} &= \frac{\frac{1}{2} (\mathbf{y} - f_2(\mathbf{z}))^T (\mathbf{y} - f_2(\mathbf{z}))}{\partial \mathbf{z}} \\ &= (\mathbf{y} - f_2(\mathbf{z})) \frac{\partial f_2(\mathbf{z})}{\partial \mathbf{z}} = (\mathbf{y} - \hat{\mathbf{y}}) f_2'(\mathbf{z}) \end{aligned}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}} = \frac{\partial \mathbf{w}^T \mathbf{h}}{\partial \mathbf{w}} = \mathbf{h}$$



Back-propagation

輸出層→隱藏層



隱藏層→輸出層

$$\mathbf{z} = \mathbf{w}^T \mathbf{h}$$

$$\mathbf{y} = f_2(\mathbf{z})$$

參數為 \mathbf{w} ，因此我們對loss進行 \mathbf{w} 的偏微分

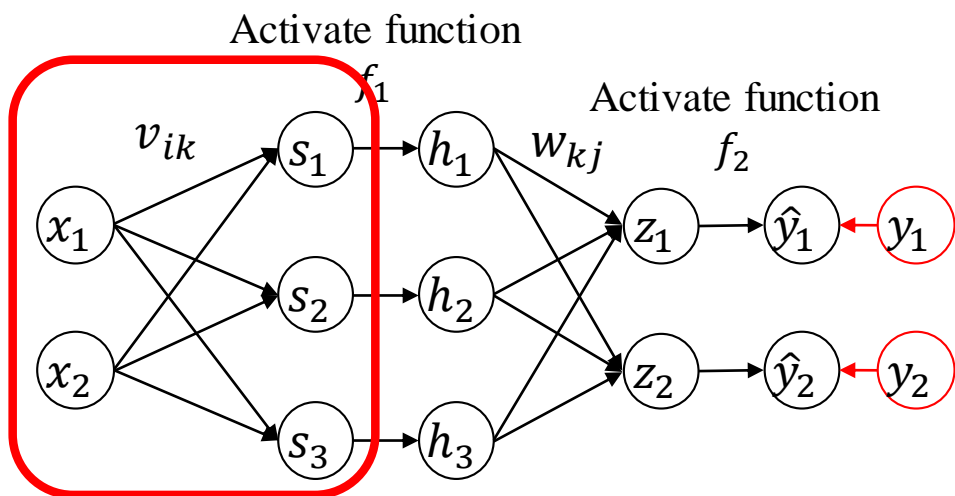
$$\begin{aligned} \text{loss}(\mathbf{y}, \hat{\mathbf{y}}) &= \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{z}))^T (\mathbf{y} - f_2(\mathbf{z})) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T \mathbf{h}))^T (\mathbf{y} - f_2(\mathbf{w}^T \mathbf{h})) \end{aligned}$$

$$\begin{aligned} \Delta \mathbf{w} &= \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{w}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{w}} \\ &= (\mathbf{y} - \hat{\mathbf{y}}) f_2'(\mathbf{z}) \mathbf{h} \end{aligned}$$



Back-propagation

隱藏層→輸入層



隱藏層→輸入層

$$\mathbf{s} = \mathbf{v}^T \mathbf{x}$$

$$\mathbf{h} = f_1(\mathbf{s})$$

參數為 \mathbf{v} ，因此我們對 loss 進行 \mathbf{v} 的偏微分

$$\begin{aligned} \text{loss}(\mathbf{y}, \hat{\mathbf{y}}) &= (\mathbf{y} - f_2(\mathbf{z}))^T (\mathbf{y} - f_2(\mathbf{z})) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{s})))^T (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{s}))) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x})))^T (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x}))) \end{aligned}$$

Chain rule

$$\Delta \mathbf{v} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{v}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{v}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{v}}$$

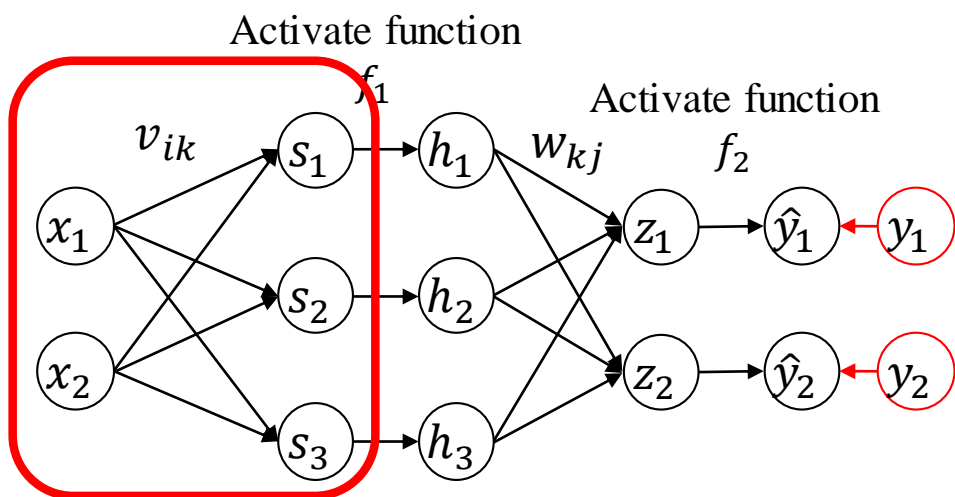
前面解過了

$$\frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}} = (\mathbf{y} - \hat{\mathbf{y}}) f_2'(\mathbf{z})$$



Back-propagation

隱藏層→輸入層



隱藏層→輸入層

$$\mathbf{s} = \mathbf{v}^T \mathbf{x}$$

$$\mathbf{h} = f_1(\mathbf{s})$$

參數為 \mathbf{v} ，因此我們對loss進行 \mathbf{v} 的偏微分

$$\begin{aligned} \text{loss}(\mathbf{y}, \hat{\mathbf{y}}) &= (\mathbf{y} - f_2(\mathbf{z}))^T (\mathbf{y} - f_2(\mathbf{z})) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{s})))^T (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{s}))) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x})))^T (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x}))) \end{aligned}$$

Chain rule

$$\Delta \mathbf{v} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{v}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{v}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{v}}$$

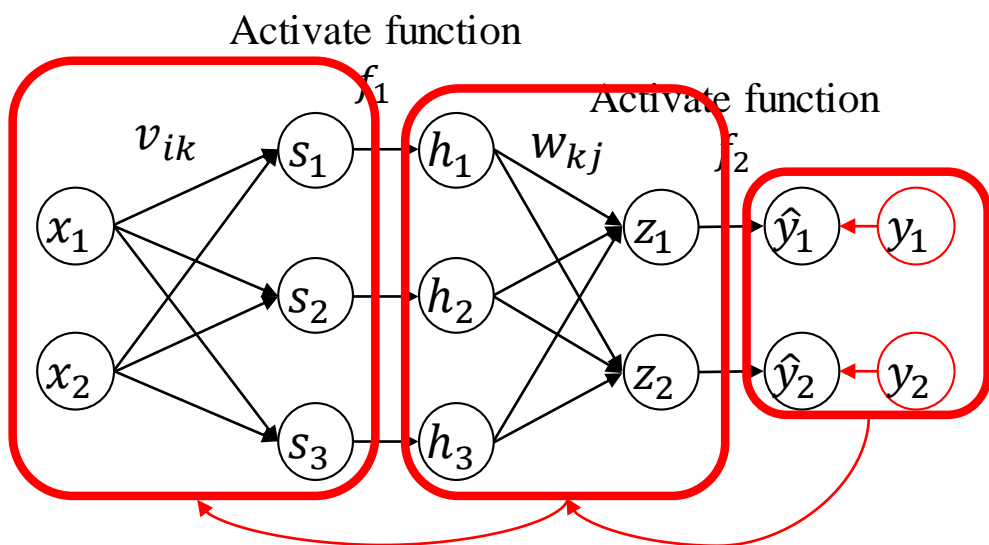
$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = \frac{\partial \mathbf{w}^T f_1(\mathbf{s})}{\partial \mathbf{s}} = \mathbf{w}^T f_1'(\mathbf{s})$$

$$\frac{\partial \mathbf{s}}{\partial \mathbf{v}} = \frac{\partial \mathbf{v}^T \mathbf{x}}{\partial \mathbf{v}} = \mathbf{x}$$



Back-propagation

隱藏層→輸入層



隱藏層→輸入層

$$\mathbf{s} = \mathbf{v}^T \mathbf{x}$$

$$\mathbf{h} = f_1(\mathbf{s})$$

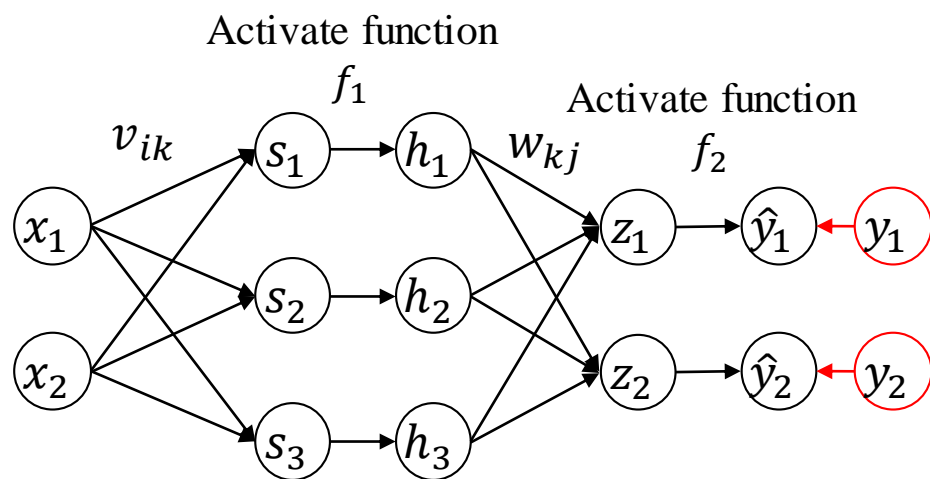
參數為 \mathbf{v} ，因此我們對 loss 進行 \mathbf{v} 的偏微分

$$\begin{aligned} \text{loss}(\mathbf{y}, \hat{\mathbf{y}}) &= (\mathbf{y} - f_2(\mathbf{z}))^T (\mathbf{y} - f_2(\mathbf{z})) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{s})))^T (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{s}))) \\ &= \frac{1}{2} (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x})))^T (\mathbf{y} - f_2(\mathbf{w}^T f_1(\mathbf{v}^T \mathbf{x}))) \end{aligned}$$

$$\begin{aligned} \Delta \mathbf{v} &= \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{v}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{v}} = \frac{\partial \text{loss}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \\ &= [(\mathbf{y} - \hat{\mathbf{y}}) f_2'(\mathbf{z})] [\mathbf{w}^T f_1'(\mathbf{s})] \mathbf{x} \end{aligned}$$



Back-propagation



輸入層→隱藏層

$$\begin{aligned} s &= v^T x \\ h &= f_1(s) \end{aligned}$$

隱藏層→輸出層

$$\begin{aligned} z &= w^T h \\ y &= f_2(z) \end{aligned}$$

$$w^{(t+1)} = w^{(t)} - \alpha \Delta w^{(t)}$$

$$v^{(t+1)} = v^{(t)} - \alpha \Delta v^{(t)}$$

輸出層→隱藏層參數為 w

$$\Delta w = (y - \hat{y}) f_2'(z) h$$

隱藏層→輸入層參數為 v

$$\Delta v = [(y - \hat{y}) f_2'(z)] [w^T f_1'(s)] x$$

前面層的gradient
會是由後面所有層的gradient和現在這層的疊乘。



神經網路太深層的問題

神經網路太深層的問題



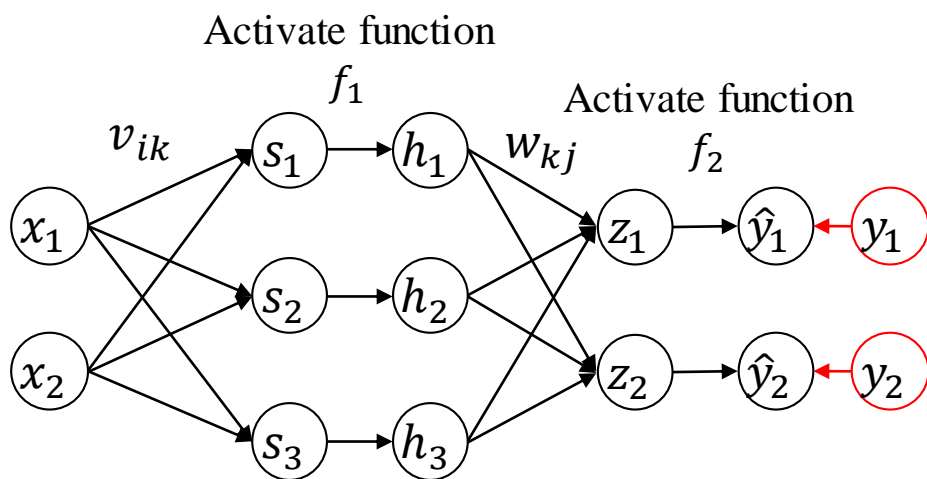
Sigmoid在神經網路太深層的問題

輸出層→隱藏層參數為 w

$$\Delta w = (y - \hat{y}) f_2'(z) h$$

隱藏層→輸入層參數為 v

$$\Delta v = [(y - \hat{y}) f_2'(z)] [w^T f_1'(s)] x$$



輸入層→隱藏層

$$s = v^T x$$

$$h = f_1(s)$$

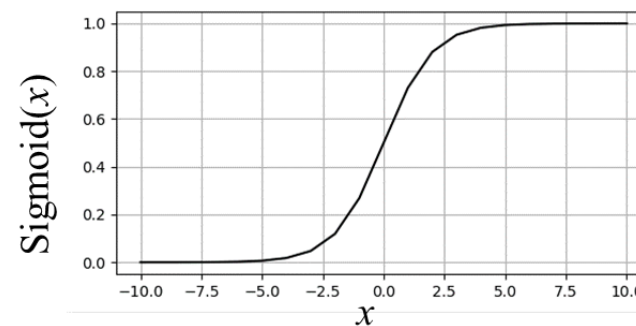
隱藏層→輸出層

$$z = w^T h$$

$$y = f_2(z)$$

$$\text{Sigmoid: } f(x) = \frac{1}{1+e^{-x}}$$

$$f' = f(x)(1 - f(x))$$

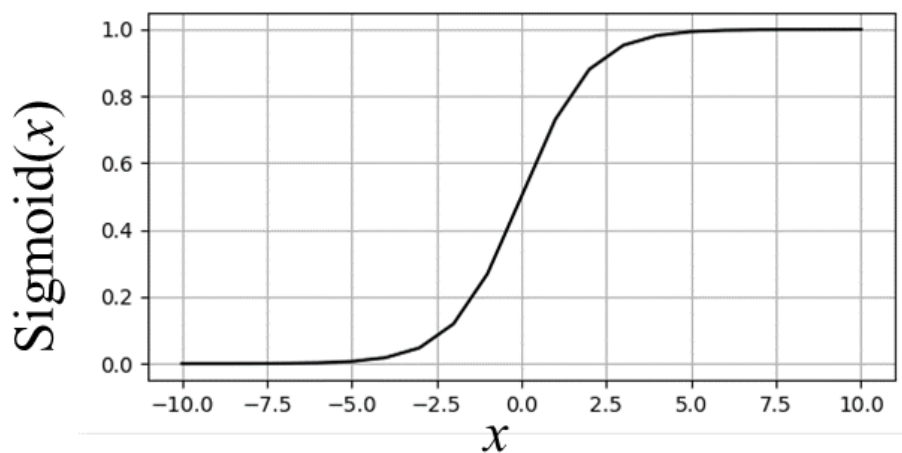


Sigmoid在神經網路太深層的問題

Sigmoid : $f(x) = \frac{1}{1+e^{-x}}$

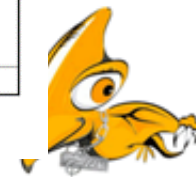
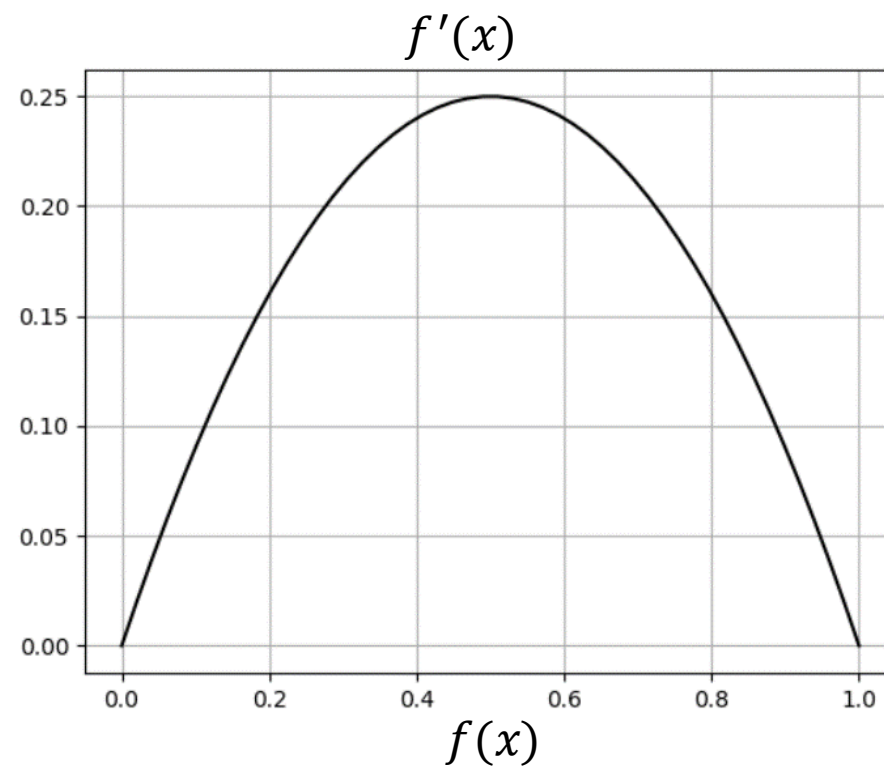
$$f'(x) = f(x)(1 - f(x))$$

$f(x)$ 輸出介於0~1



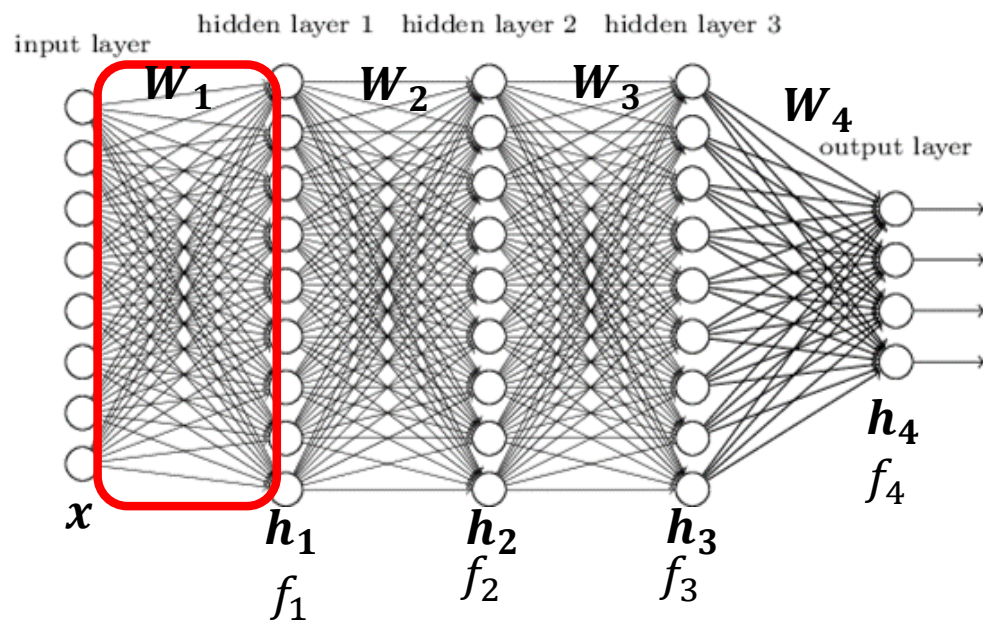
$f'(x)$ 會更小，最大值為0.25

$$f(x) = 0.5, f'(x) = 0.5 \times 0.5 = 0.25$$



Sigmoid在神經網路太深層的問題

MLP



$$\Delta W_1 = [(y - \hat{y}) f'_4(h_4)] [W_3^T f'_3(h_3)] [W_2^T f'_2(h_2)] [W_1^T f'_1(h_1)] x$$

所以當層數到100層時候，對於第一層的gradient會有100個activate function的導數相乘。

剛剛sigmoid已經說了，其導數最大為0.25。

所以 $0.25^{100} \approx 0$

這就是梯度消失問題(Vanishing gradient)

如果導數值單調大於1時，就會發生**梯度爆炸問題 (exploding gradient problem)**。



如何舒緩Gradient造成的問題

- 1. 重新設計網路架構: 更少的層。
- 2. Rectified Linear Activation (ReLU)
- 3. Gradient Clipping (Keras預設 clipnorm = 1.0和clipvalue = 0.5。)
- 4. Weight Regularization (L1 or L2 Regularizers)

實際解決神經網路太深層的方法: Residual block。

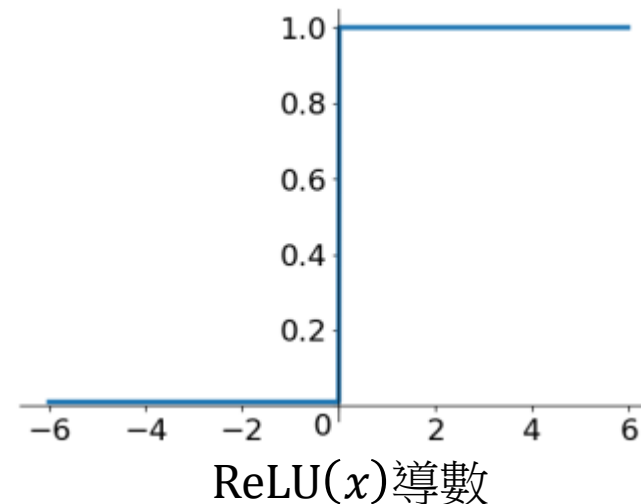
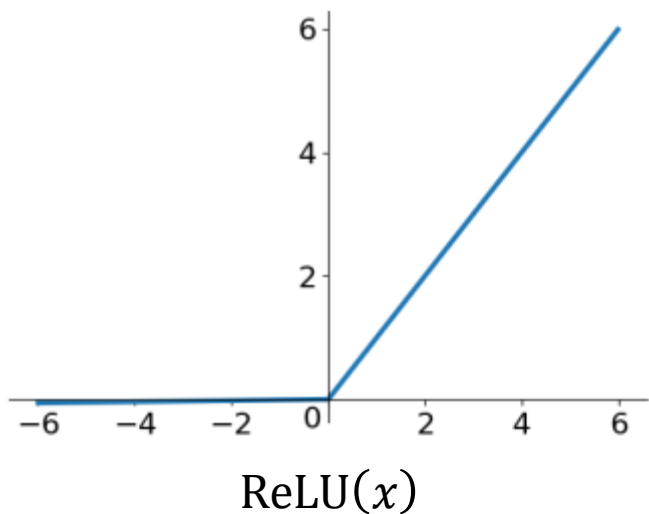
本次介紹不考慮RNN系列的設計。



ReLU在神經網路如何舒緩Gradient的問題

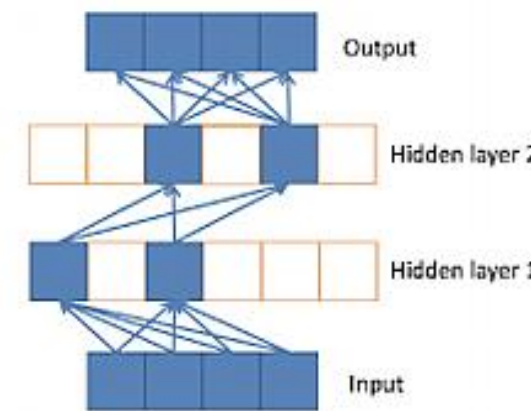
- $\text{ReLU}(x) = \max(0, x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{O.W.} \end{cases}$
- $\text{ReLU}(x)$ 的導數 = $\begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{O.W.} \end{cases}$

ReLU函數並不是全區間皆可微分，但是不可微分的部分可以使用Sub-gradient進行取代



ReLU在神經網路如何舒緩Gradient的問題

- ReLU激勵函數會使負數部分的神經元輸出為0，可以讓網路變得更加多樣性，如同Dropout的概念，可以緩解過擬合(Over fitting)之問題。
- **衍生Dead ReLU的問題**，當某個神經元輸出為0後，就難以再度輸出值，當遇到以下兩種情形時容易導致dead ReLU發生。
 - 初始化權重設定為不能被激活的數值。
 - 學習率設置過大，在剛開始進行誤差反向傳遞時，容易修正權重值過大，導致權重梯度為0，神經元即再也無法被激活。



ResNet

此Residual block有兩層，第一層權重為 w_1 ，第二層權重為 w_2
第一層輸出:

$$\text{relu}(w_1 x)$$

第二層輸出:

$$F(x) = w_2 \text{relu}(w_1 x)$$

最後element-add後結果:

$$H(x) = F(x) + x$$

假設神經網路太深對應用沒有幫助

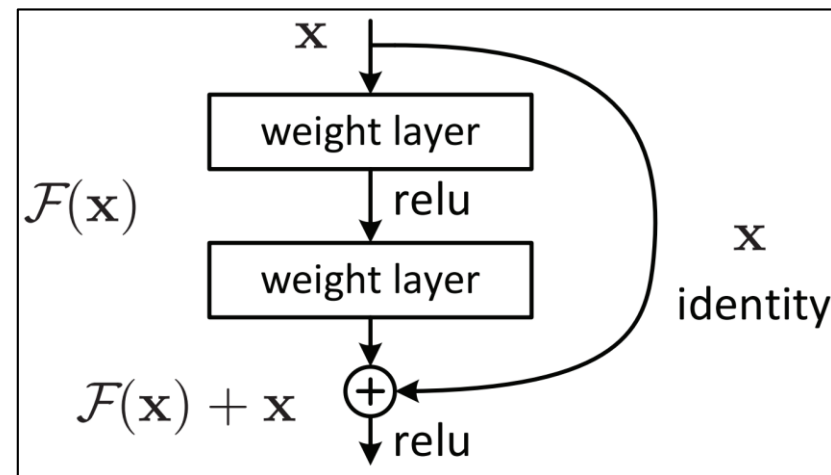
1. 沒有shortcut connection (residual) ，兩層的輸出結果叫做 $H(x)$

最理想狀況就是 $H(x) = \mathbf{F(x)} = \mathbf{x}$ (identity mapping)

2. 有shortcut connection , $H(x) = F(x) + x \rightarrow \mathbf{F(x) = H(x) - x}$

最理想狀況就是 $\mathbf{F(x) = 0}$

2的優化找weight解會比1容易。



ResNet

假設我們的*residual block*只有一層權重為 w

1. 沒有shortcut connection(residual) , 兩層的輸出結果叫做 $H(x)$

$$F(x) = wx = x$$

2. 有shortcut connection , $H(x) = F(x) + x \rightarrow F(x) = H(x) - x$

$$F(x) = wx = 0$$

Gradient:

$$1. \frac{\partial E}{\partial w} = \frac{(\hat{y}-y)^2}{\partial w} = \frac{(wx-y)^2}{\partial w} = 2(wx-y)x$$

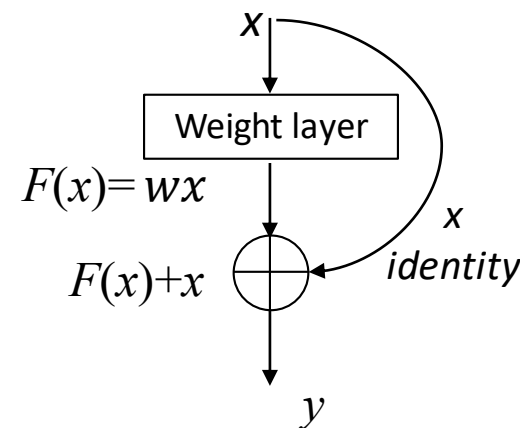
$$2. \frac{\partial E}{\partial w} = \frac{(\hat{y}-y)^2}{\partial w} = \frac{(wx+x-y)^2}{\partial w} = 2(wx-y)x + 2x^2$$

接近0

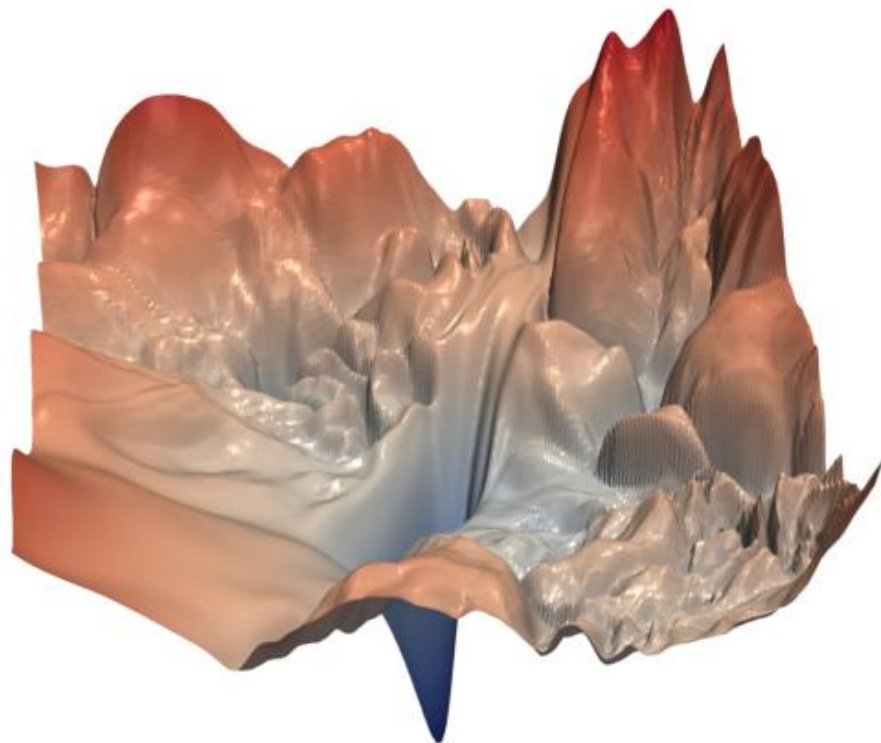
假設越後面層的Gradient很小，甚至Gradient vanish)。

$$\Delta W_1 = [(y - \hat{y})f_4'(h_4)][W_3^T f_3'(h_3)][W_2^T f_2'(h_2)][W_1^T f_1'(h_1)]x$$

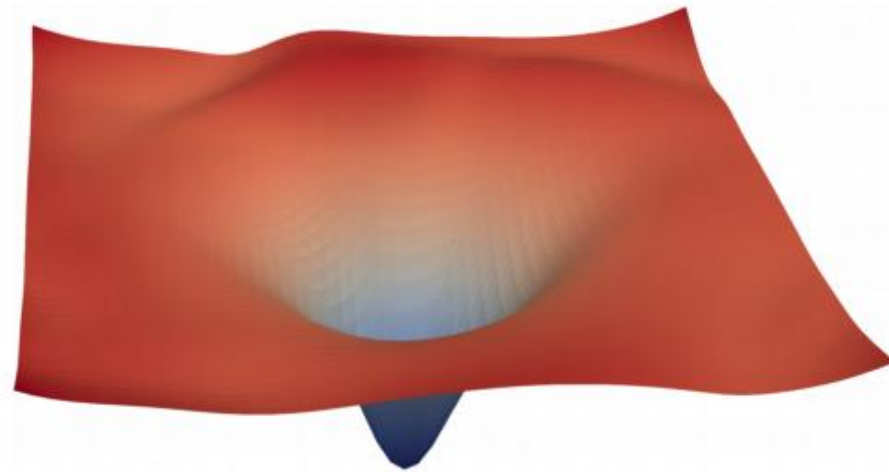
residual block



題外話: 有Residual block的loss space



(a) without skip connections



(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.



基礎機器學習與深度學習

- 神經網路如何利用導傳遞找解
- 神經網路太深層的問題:

梯度消失問題(Vanishing gradient)/梯度爆炸問題 (exploding gradient problem) 。

Activation Function為什麼要採用ReLU，而不是用Sigmoid。

Residual block克服神經網路不能太深層的問題。

