

# Convolutional Neural Network

# CNN介紹

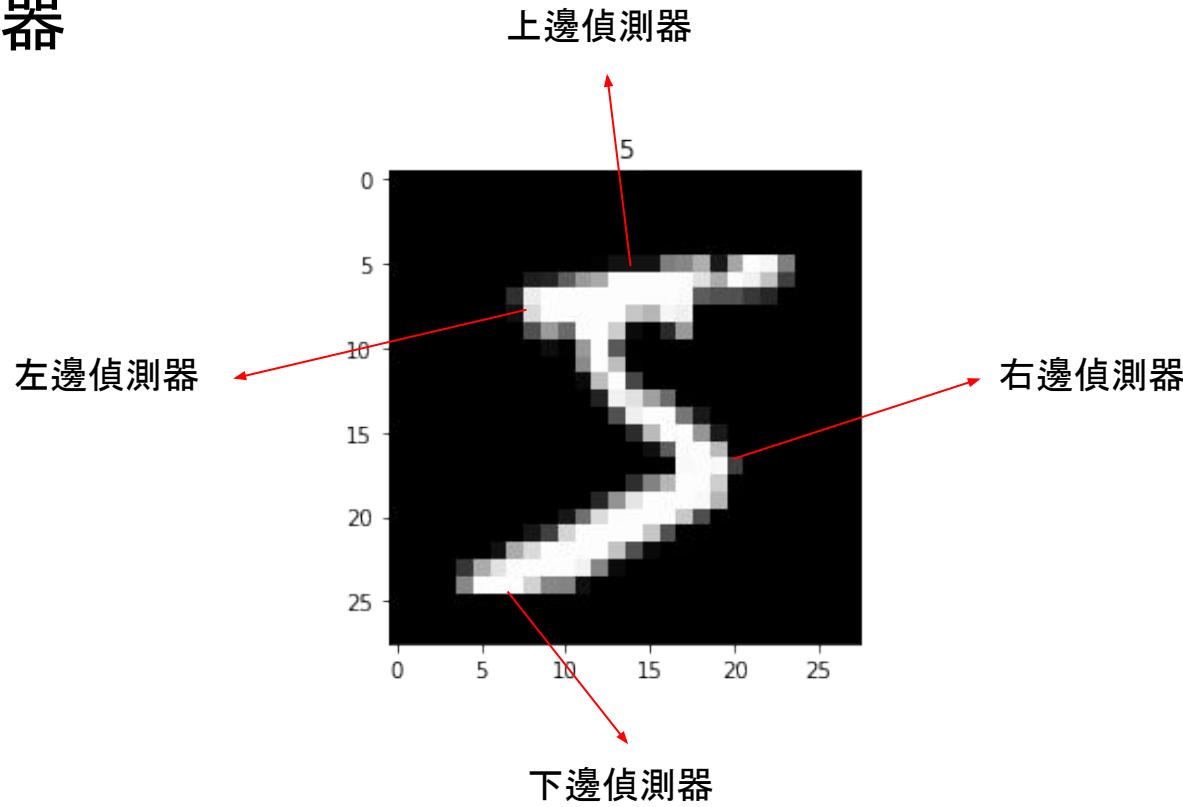


# 我們如何分辨

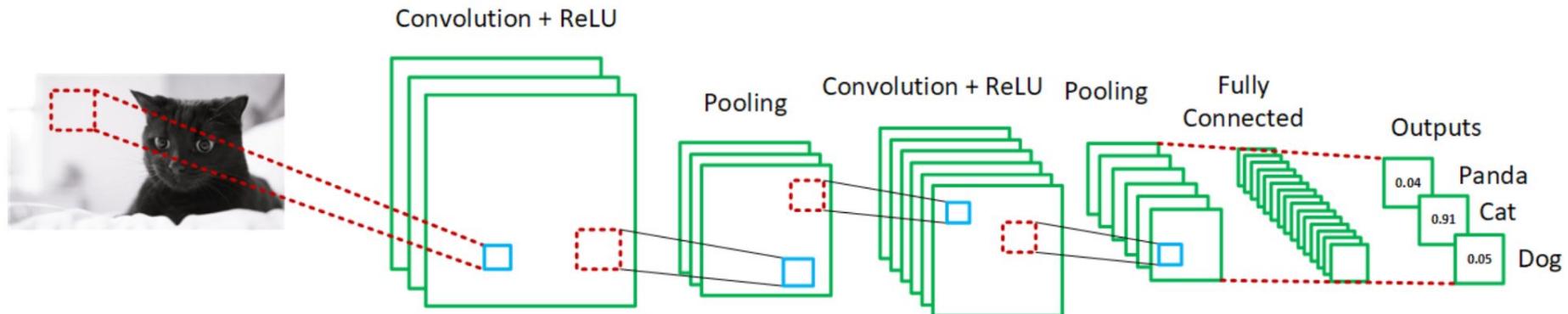
- 許多特徵能讓我們分辨狗和貓
- 鬍鬚、臉型、眼睛、毛、顏色等等
- 模型可以分辨這些特徵嗎？

CNN

# 特徵偵測器



# 用CNN預測



# Convolution (卷積)



# Convolution 運算

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	1	-1
-1	1	3
2	1	1

Output or feature map

# Convolution 運算

$$(1 \times 0) + (0 \times 1) + (1 \times 0) + (1 \times 1) + (0 \times 0) + (0 \times -1) + (0 \times 0) + (1 \times 1) + (1 \times 0) = 2$$

1x0	0x1	1x0	0	1
1x1	0x0	0x-1	1	1
0x0	1x1	1x0	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2		

Output or feature map

# Convolution 運算

$$(0 \times 0) + (1 \times 1) + (0 \times 0) + (0 \times 1) + (0 \times 0) + (1 \times -1) + (1 \times 0) + (1 \times 1) + (0 \times 0) = 1$$

1	0x0	1x1	0x0	1
1	0x1	0x0	1x-1	1
0	1x0	1x1	0x0	0
1	0	0	1	0
0	0	1	1	0

Input image

0	1	0
1	0	-1
0	1	0

\*

=

2	1	

Filter or Kernel

Output or feature map

# Convolution 運算

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	1	-1

Output or feature map

# Convolution 運算

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	1	-1
-1	1	3
2	1	1

Output or feature map

# 特徵偵測

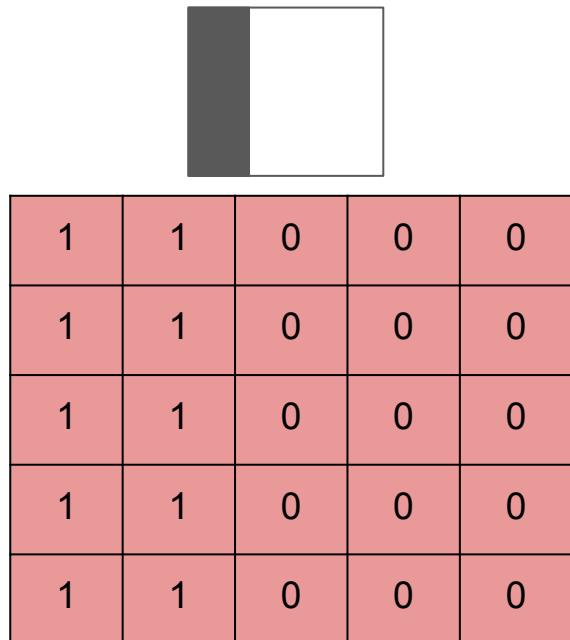
- Filter or kernel 特徵偵測器
- Feature map = 偵測結果

1	0	-1
1	0	-1
1	0	-1

=

Vertical edge detector

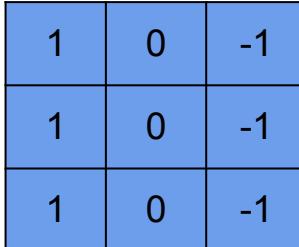
# Vertical edge detector



1	1	0	0	0
1	1	0	0	0
1	1	0	0	0
1	1	0	0	0
1	1	0	0	0

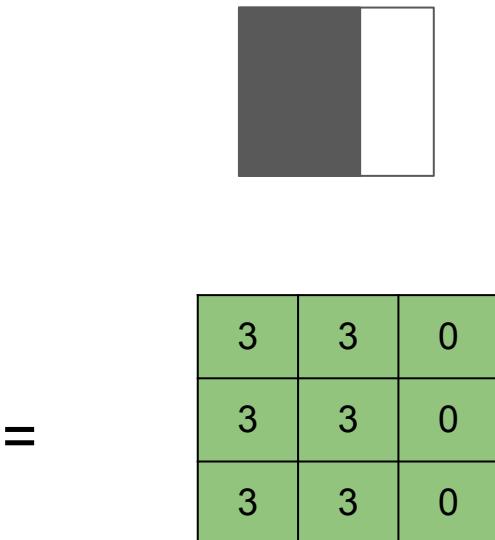
Input image

\*



1	0	-1
1	0	-1
1	0	-1

Filter or Kernel



3	3	0
3	3	0
3	3	0

=

3	3	0
3	3	0
3	3	0

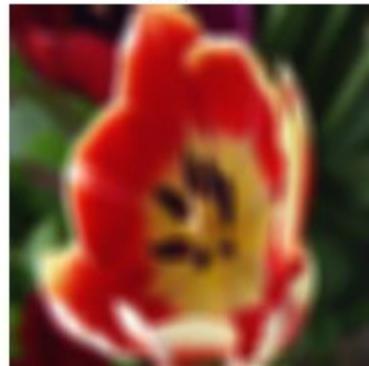
Output or feature map

# 各種 feature map

original



smooth



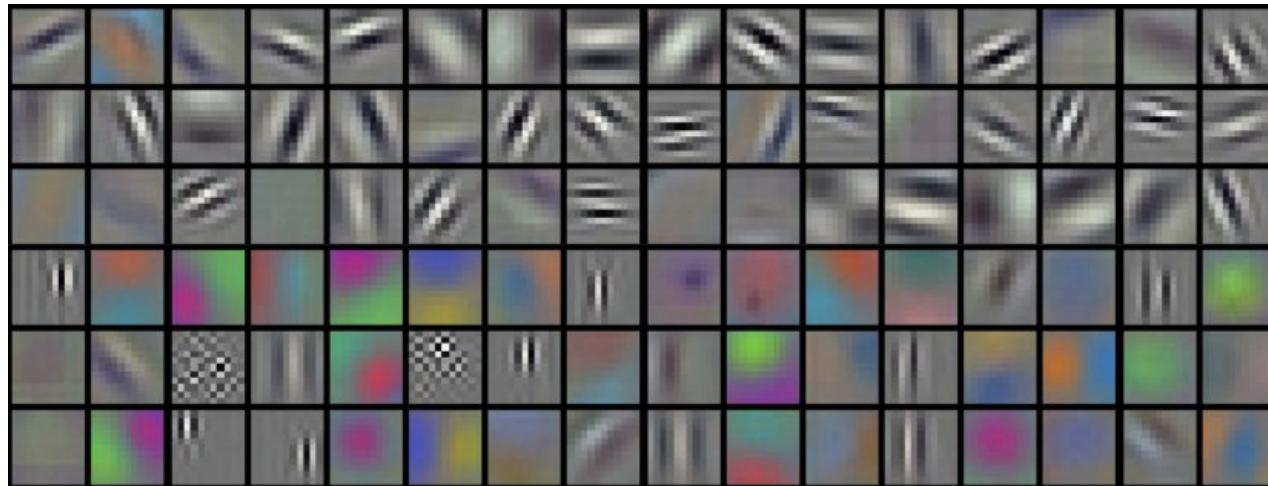
edge



intensity



# Filters 長怎樣？



# 計算 feature map 的大小

$$\text{Feature Map Size} = n - f + 1 = m$$
$$\text{Feature Map Size} = 5 - 3 + 1 = 3$$

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

\*

0	1	0
1	0	-1
0	1	0

=

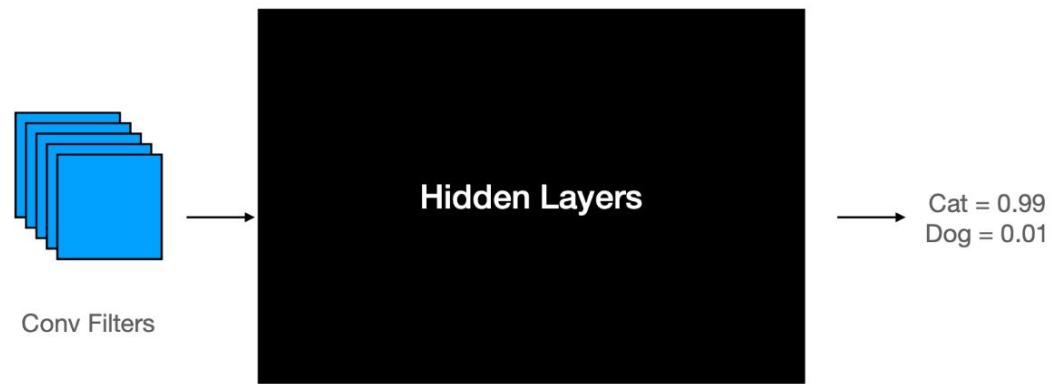
2	1	-1
-1	1	3
2	1	1

$5 \times 5$   
 $n \times n$

$3 \times 3$   
 $f \times f$

$3 \times 3$   
 $m \times m$

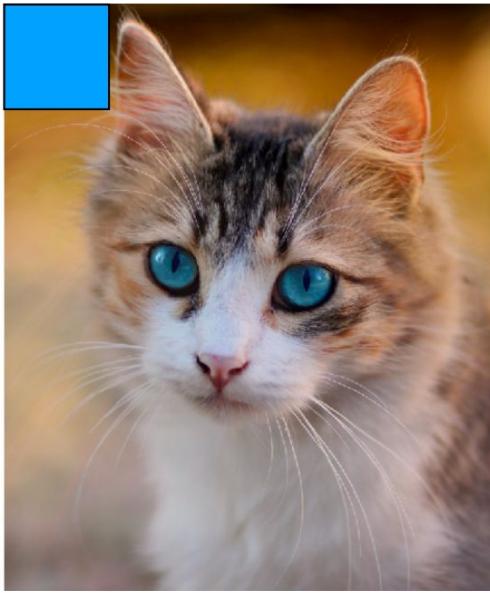
# 如何分類 image



Uses the Feature Maps we output to determine the image class

We perform Convolution operations with the input image and our filters

# 如何分類 image



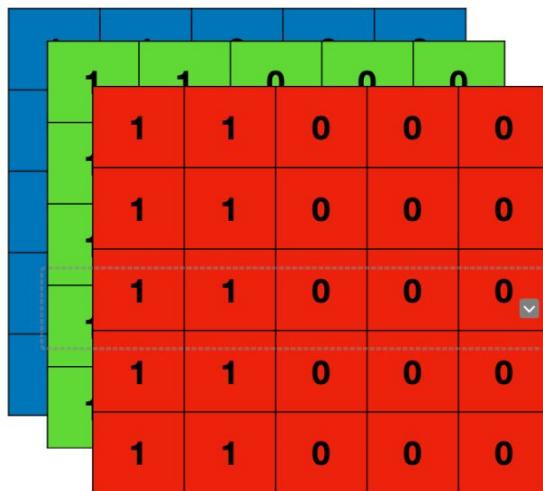
We perform Convolution operations with the input image and our filters

# Convolution demo

[https://deeplizard.com/learn/video/YRhxdVk\\_sls](https://deeplizard.com/learn/video/YRhxdVk_sls)

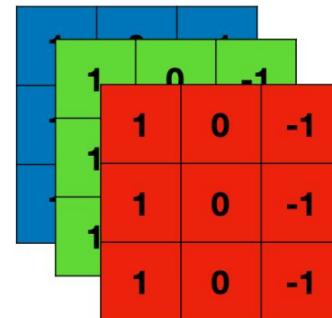
# Convolution on Color Images

# 如何計算彩色 image 的 convolution



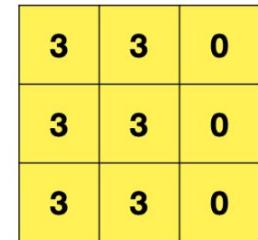
Input Image

\*



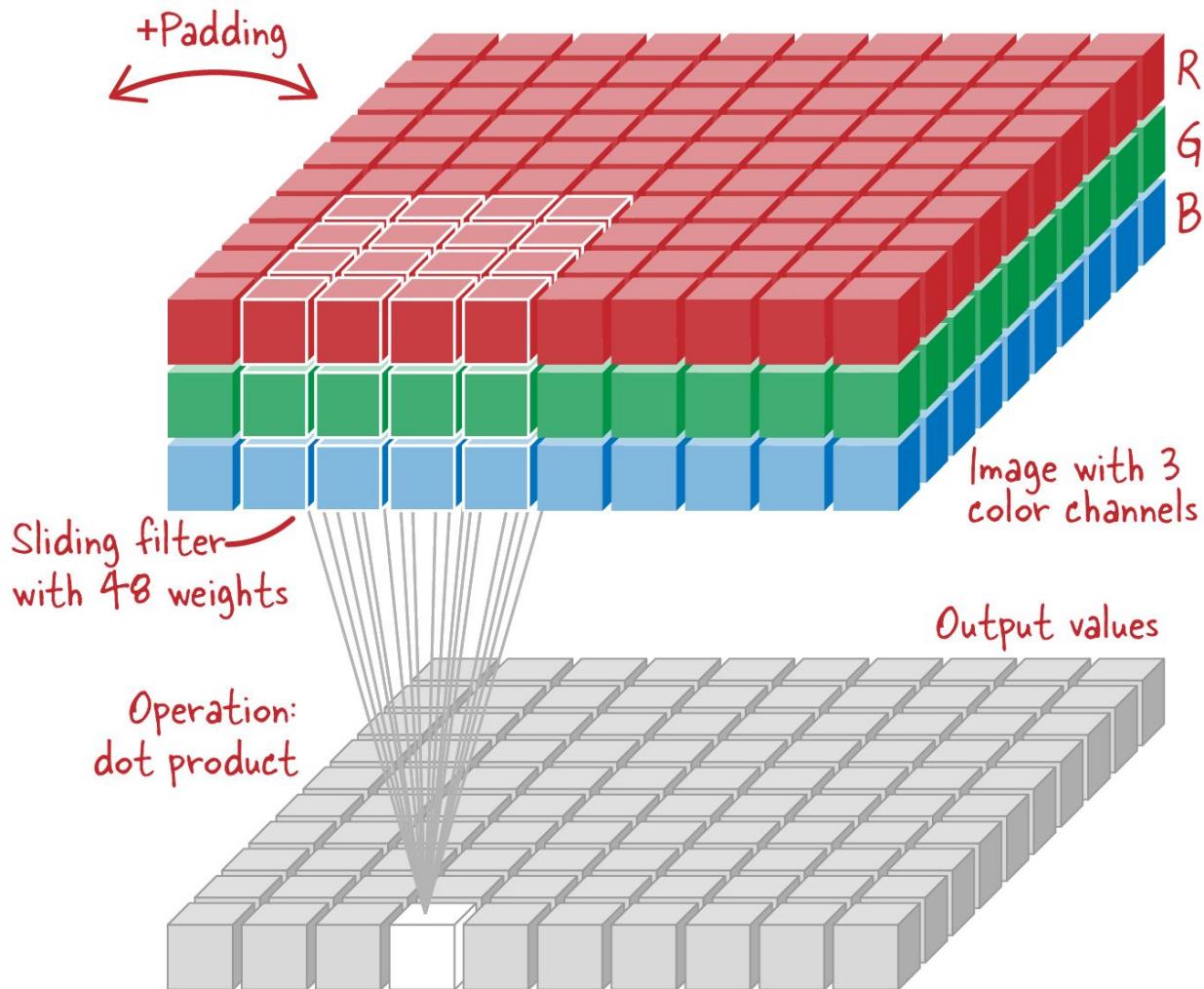
Filter or Kernel

=

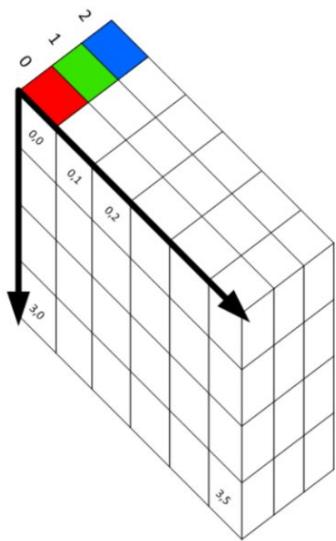


Output or Feature Map

我們可以調整不同頻道的 filter 來關注某個顏色的特徵

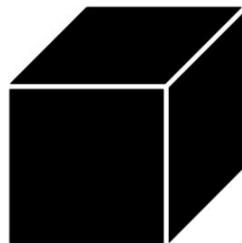


# 計算 feature map size



Input Image  
 $5 \times 5 \times 3$

\*



Filter or Kernel  
 $3 \times 3 \times 3$

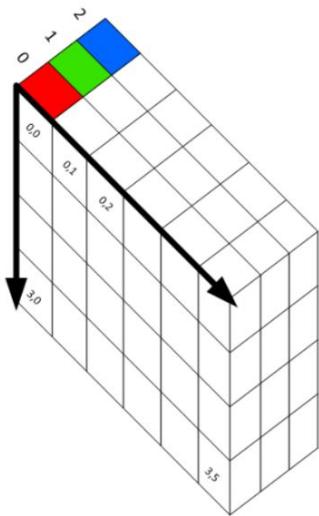
=

3	3	0
3	3	0
3	3	0

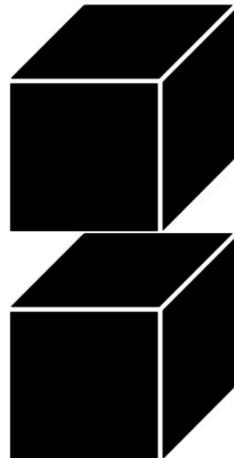
Output or Feature Map  
 $3 \times 3$

# 計算 feature map size

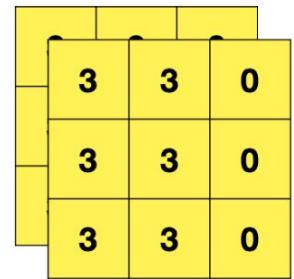
$$(n \times n \times nc) * (f \times f \times nc) = (n - f + 1) \times (n - f + 1) \times nf$$
$$(5 \times 5 \times 3) * (3 \times 3 \times 3) = 3 \times 3 \times 2$$



\*



=



Input Image  
5 x 5 x 3

2 Filters or Kernels  
3 x 3 x 3 x 2

Output or Feature Map  
3 x 3 x 2

# Parameters for convolution filter

# Convolution filter 有哪些超參數

- Kernel size ( $k \times k$ )
- Depth
- Padding
- Stride

# Kernel size

1	1	0	0	0	0
1	1	0	0	0	0
1	1	0	0	0	0
1	1	0	0	0	0
1	1	0	0	0	0
1	1	0	0	0	0

6 x 6

\*

1	0	-1	0	1
1	0	-1	0	1
1	0	-1	0	1
1	0	-1	0	1
1	0	-1	0	1

5 x 5

=

3	3
3	3

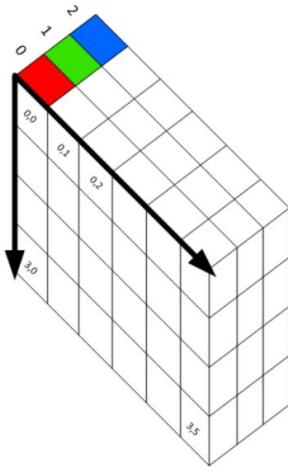
Output or Feature Map

$$\text{Feature Map Size} = n - f + 1 = m$$

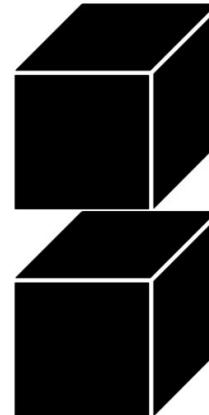
$$\text{Feature Map Size} = 6 - 5 + 1 = 2$$

# Depth

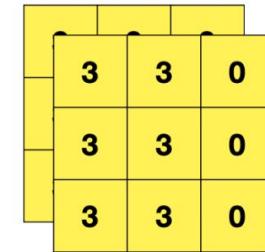
- Channels
- Filters / kernels



\*



=



Input Image  
 $5 \times 5 \times 3$

2 Filters or Kernels  
 $3 \times 3 \times 3 \times 2$

Output or Feature Map  
 $3 \times 3 \times 2$

# Padding

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

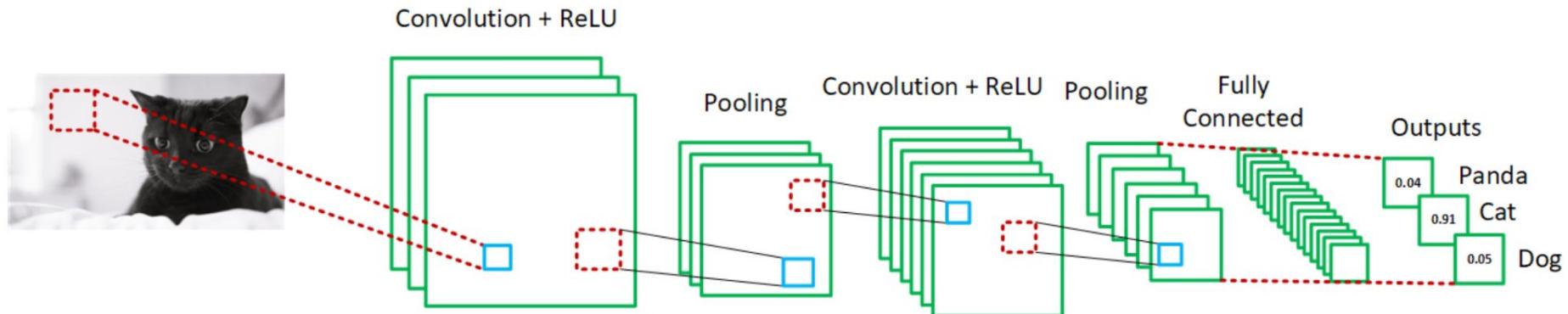
Filter or Kernel

=

2	1	-1
-1	1	3
2	1	1

Output or feature map

# 隨著convolution次數增加.....



# 如何避免image變太小

0	0	0	0	0	0	0
0	1	0	1	0	1	0
0	1	0	0	1	1	0
0	0	1	1	0	0	0
0	1	0	0	1	0	0
0	0	0	1	1	0	0
0	0	0	0	0	0	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

# 如何避免image變太小

0	0	0	0	0	0	0
0	1	0	1	0	1	0
0	1	0	0	1	1	0
0	0	1	1	0	0	0
0	1	0	0	1	0	0
0	0	0	1	1	0	0
0	0	0	0	0	0	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

2

=

Output or feature map

# 如何避免image變太小

0	0	0	0	0	0	0
0	1	0	1	0	1	0
0	1	0	0	1	1	0
0	0	1	1	0	0	0
0	1	0	0	1	0	0
0	0	0	1	1	0	0
0	0	0	0	0	0	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

2	1
---	---

=

Output or feature map

# 如何避免image變太小

0	0	0	0	0	0	0
0	1	0	1	0	1	0
0	1	0	0	1	1	0
0	0	1	1	0	0	0
0	1	0	0	1	0	0
0	0	0	1	1	0	0
0	0	0	0	0	0	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

2	1	-1
---	---	----

=

Output or feature map

# 如何避免image變太小

0	0	0	0	0	0	0	0
0	1	0	1	0	1	0	
0	1	0	0	1	1	0	
0	0	1	1	0	0	0	
0	1	0	0	1	0	0	
0	0	0	1	1	0	0	
0	0	0	0	0	0	0	

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

$$\text{Feature Map Size} = n - f + 1 = m$$
$$\text{Feature Map Size} = 7 - 3 + 1 = 5$$

=

2	1	-1	2	2
-1	1	3	2	1
2	1	1	1	2
1	1	1	0	2
2	0	2	3	1

Output or feature map

# 為什麼需要padding?

- 當使用較深的模型時，避免feature map持續變小
- 增加角落的pixel的貢獻

只做一次convolution

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

好幾次convolution

# Stride

當 stride = 1

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2		

Output or feature map

# Stride

當 stride = 1

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	1	

Output or feature map

# Stride

當 stride = 1

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	1	-1
-1	1	3
2	1	1

Output or feature map

# Stride

當 stride = 2

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2
---

Output or feature map

# Stride

當 stride = 2

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	-1

Output or feature map

# Stride

當 stride = 2

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	-1
2	

Output or feature map

# Stride

當 stride = 2

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	-1
2	1

Output or feature map

## 使用不同stride

- 增加模型訓練效率
- 減少重複使用相同的pixel
- 控制feature map的大小

# 計算feature map大小

Stride = 2  
Padding = 0

$$\left(\frac{n+2p-f}{s} + 1\right) \times \left(\frac{n+2p-f}{s} + 1\right) = \left(\frac{5+(2 \times 0)-3}{2} + 1\right) \times \left(\frac{5+(2 \times 0)-3}{2} + 1\right) = 2 \times 2$$

1	0	1	0	1
1	0	0	1	1
0	1	1	0	0
1	0	0	1	0
0	0	1	1	0

Input image

\*

0	1	0
1	0	-1
0	1	0

Filter or Kernel

=

2	-1
2	1

Output or feature map

# 計算feature map大小

Stride = 1  
Padding = 1

0	0	0	0	0	0	0
0	1	0	1	0	1	0
0	1	0	0	1	1	0
0	0	1	1	0	0	0
0	1	0	0	1	0	0
0	0	0	1	1	0	0
0	0	0	0	0	0	0

$$(\frac{n+2p-f}{s} + 1) \times (\frac{n+2p-f}{s} + 1) = (\frac{5+(2 \times 1)-3}{1} + 1) \times (\frac{5+(2 \times 1)-3}{1} + 1) = 5 \times 5$$

\* =

2	1	-1	2	2
-1	1	3	2	1
2	1	1	1	2
1	1	1	0	2
2	0	2	3	1

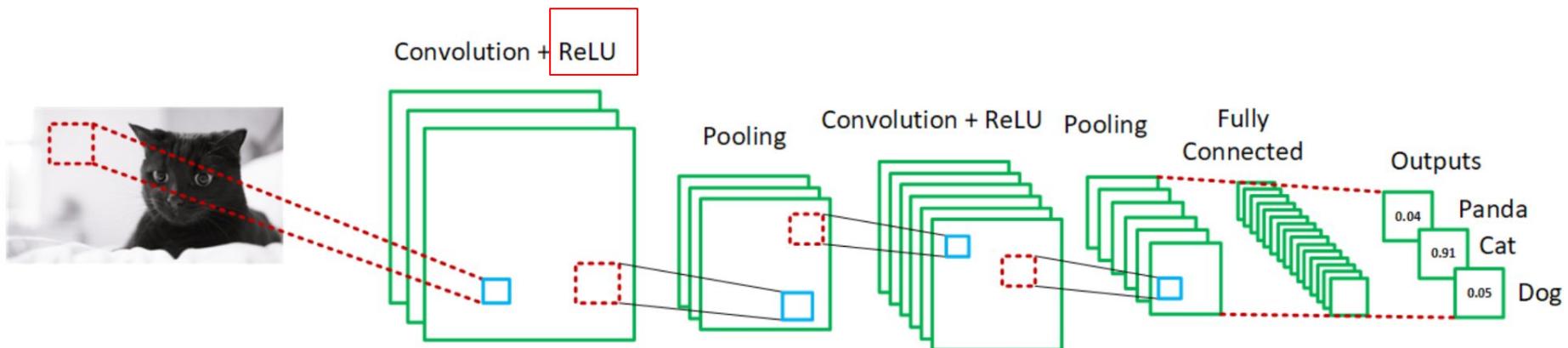
Input image

Filter or Kernel

Output or feature map

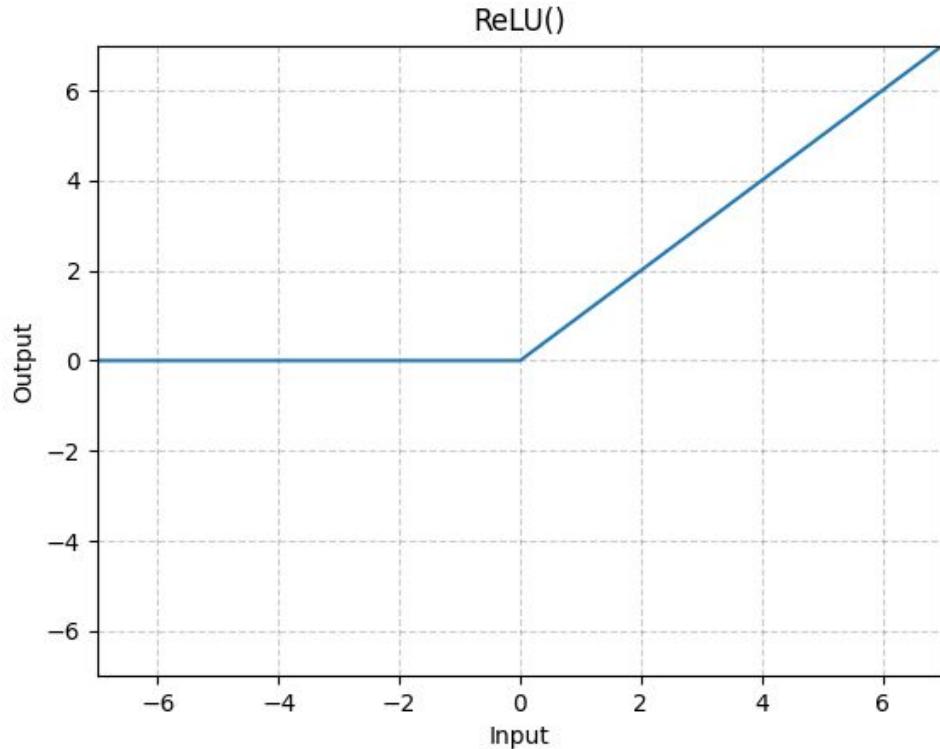
# Activation layer ReLU

# Activation layer ReLU



# ReLU

- 將負數轉為0
- 保留正數



# 對feature map使用ReLU

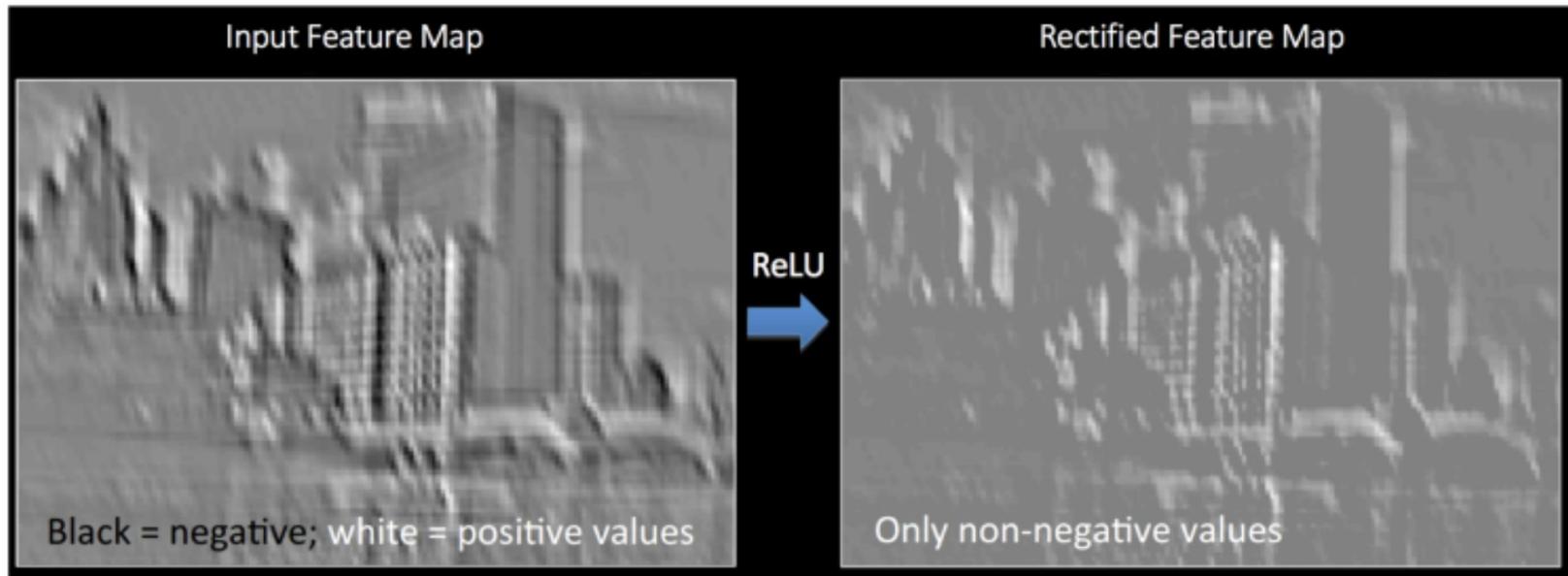
2	1	-1
-1	1	3
2	1	1



2	1	0
0	1	3
2	1	1

Output or feature map

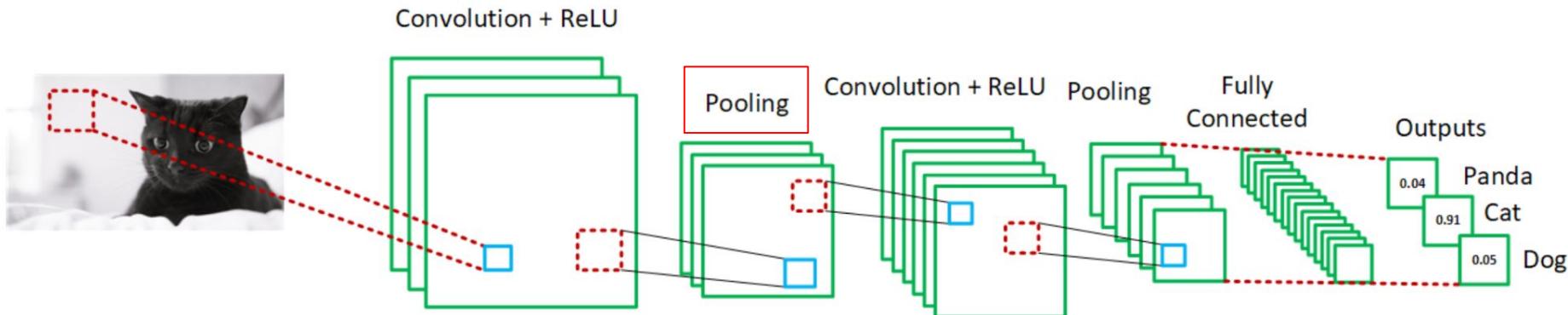
# 使用ReLU後的feature map



Source - [http://mlss.tuebingen.mpg.de/2015/slides/fergus/Fergus\\_1.pdf](http://mlss.tuebingen.mpg.de/2015/slides/fergus/Fergus_1.pdf)

# Pooling layer

# Pooling layer



# Pooling

- 減少feature map的大小
- 在不丟失重要資訊的情況下減少訓練參數
- 也叫Subsampling or Downsampling

# How max pooling works

4	123	1	34
56	99	222	253
45	122	165	12
21	187	133	124

**MaxPool Operation**



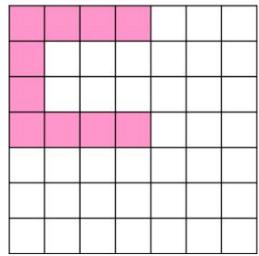
**Stride = 2**  
**Kernel = 2x2**

123	253
187	165

# About pooling

- 常使用的pooling layer: 2x2 kernels, strides = 2, no padding
- Feature map size =  $\frac{1}{2}$  input size
- 其他pooling方法: average pooling, sum pooling

# 為什麼pooling有用



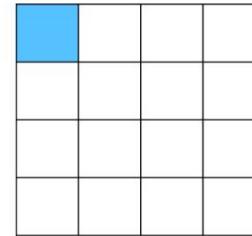
Convolution

\*

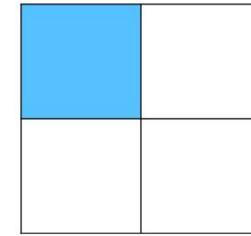
1	1	1	1
1	0	0	1
1	0	0	1
1	1	1	1



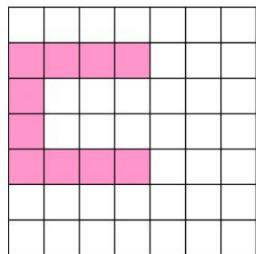
Feature Map



Max Pool



Shifting Our C down one pixel



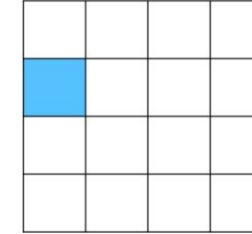
Convolution

\*

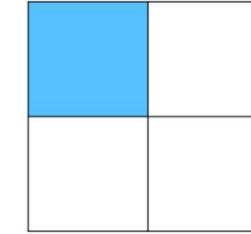
1	1	1	1
1	0	0	1
1	0	0	1
1	1	1	1



Feature Map



Max Pool

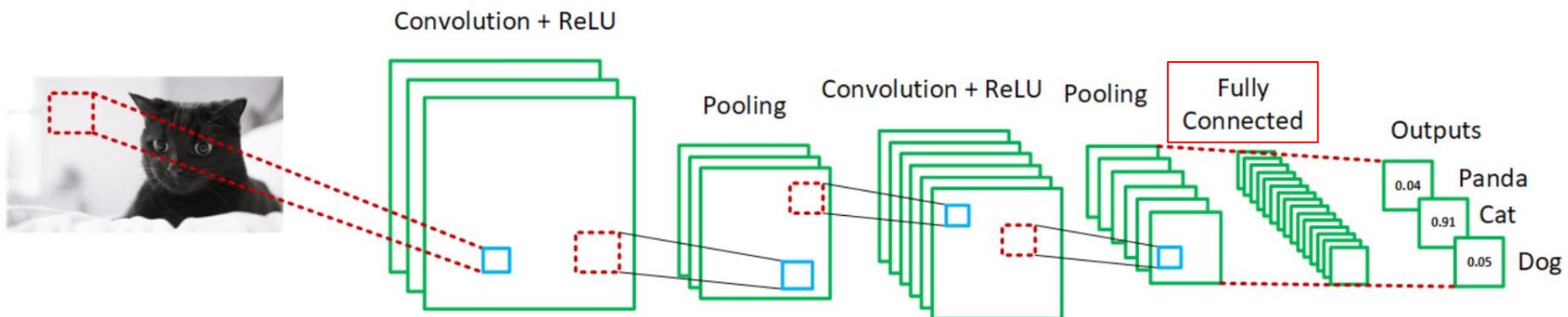


# 為什麼pooling有用

- 相鄰的pixel通常具有相似的特徵(尤其是初期的layer)
- 避免模型受到微小特徵差異的影響
- 在不丟失重要資訊的情況下減少訓練參數
- 太大的stride會丟失重要資訊

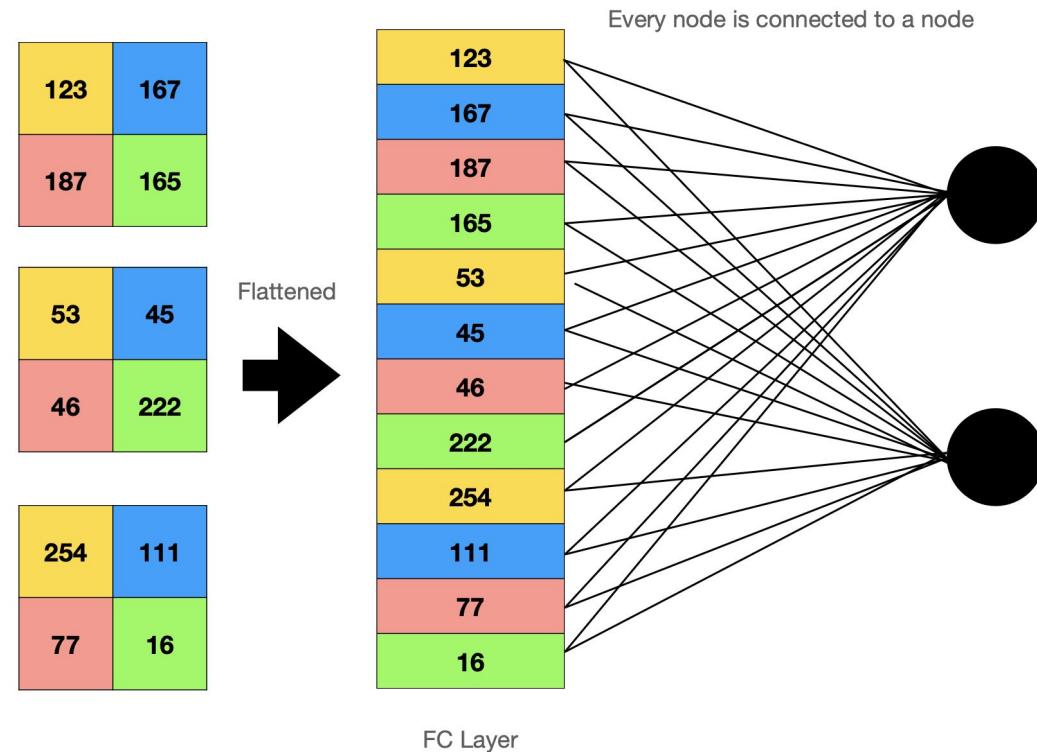
# Fully connected layers

# Fully connected layers



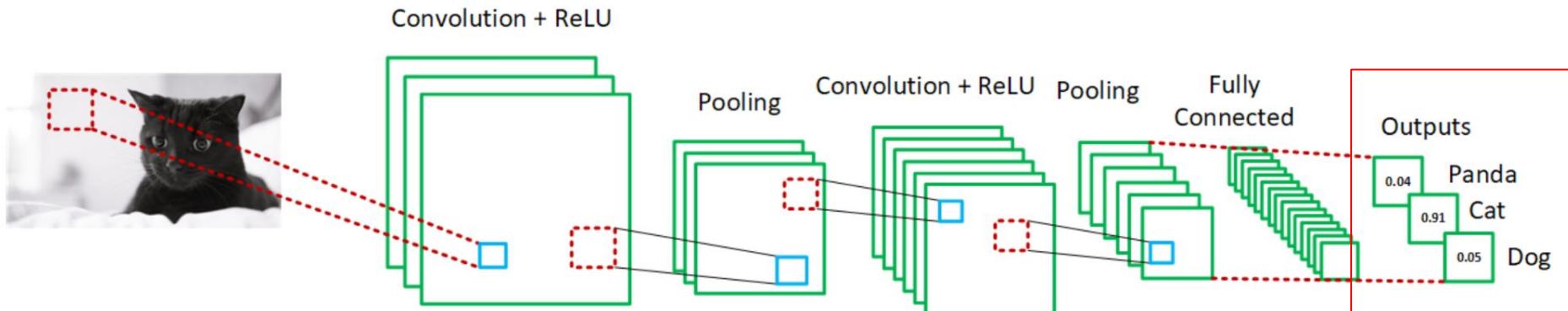
# Fully connected layers

- FC layer上的所有點都和下一層的每個點連接
- 也叫Dense layer

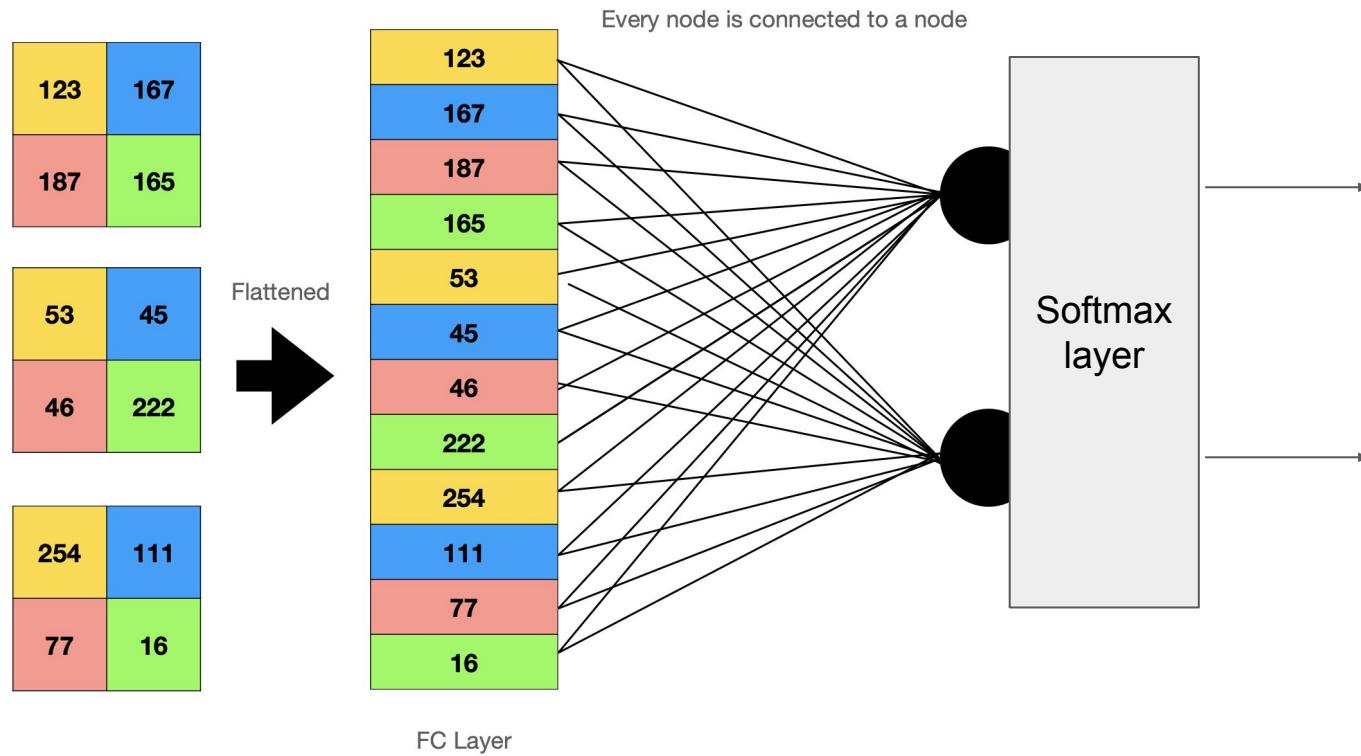


# Softmax layer

# Softmax layers



# Softmax layer



# Softmax layer

Logits Scores

2.0

1.0

0.1

Probabilities

0.7

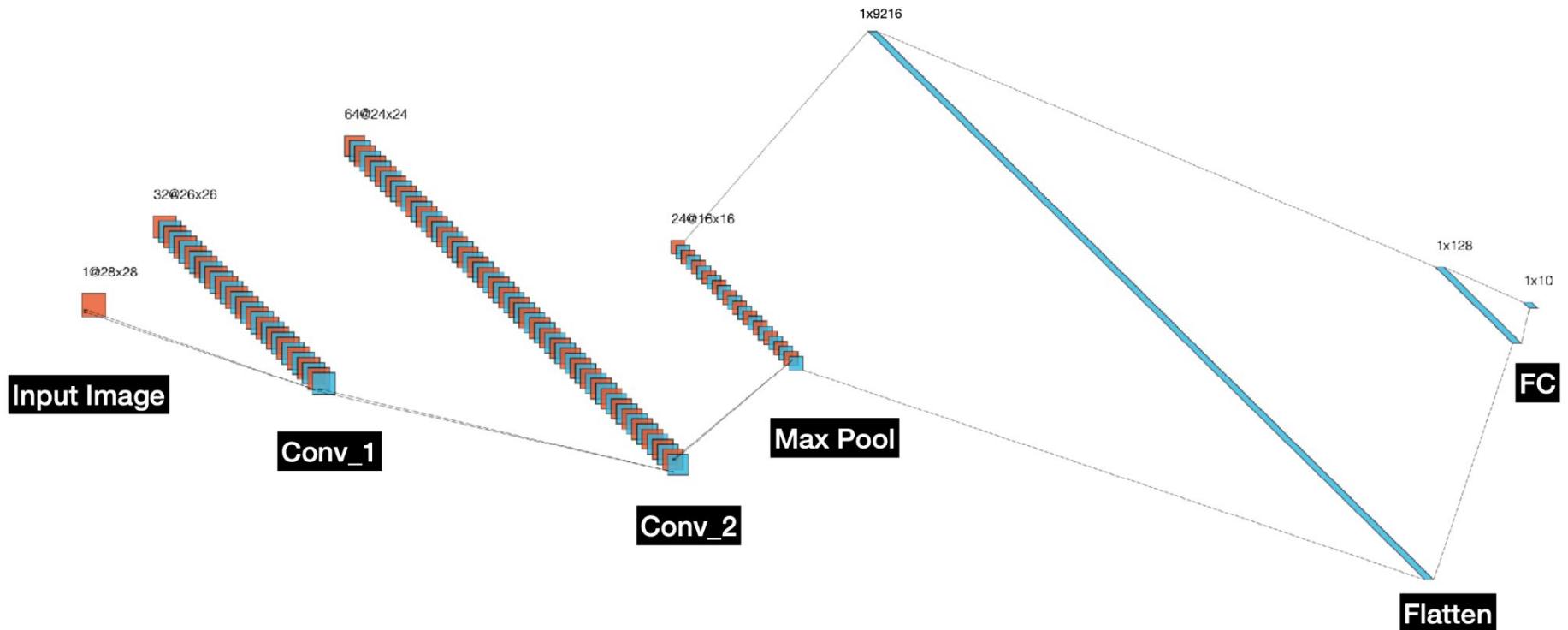
0.2

0.1

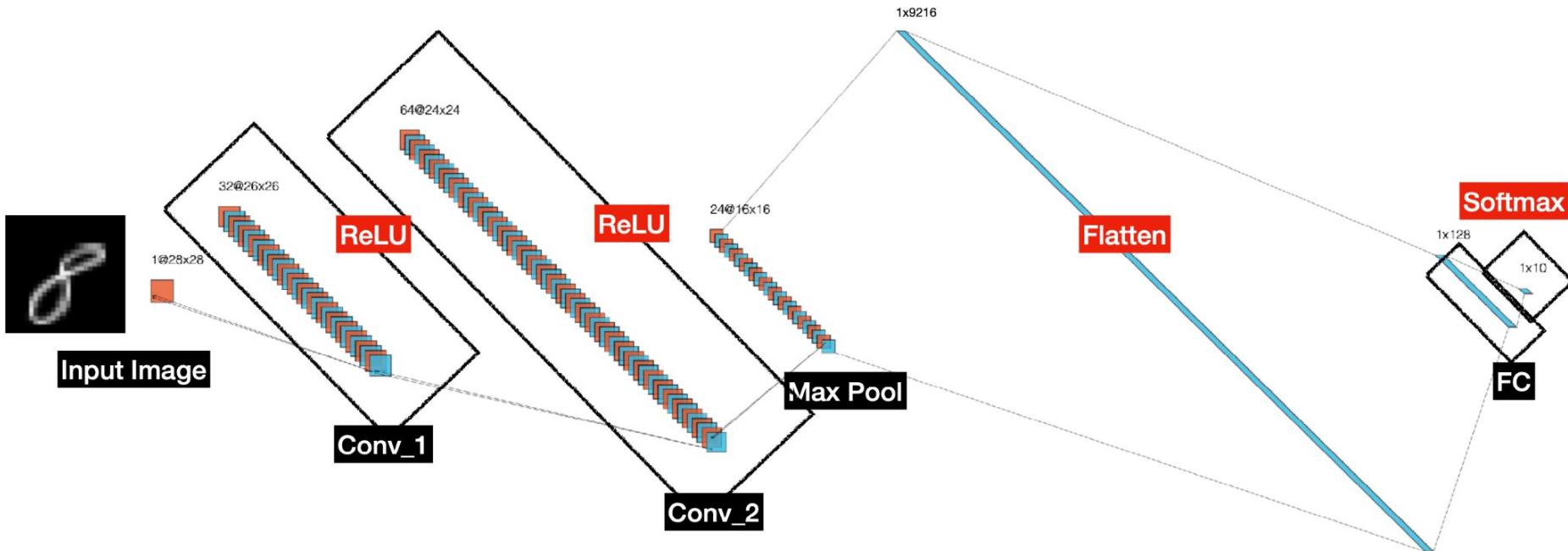
$$\text{softmax}(x)_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

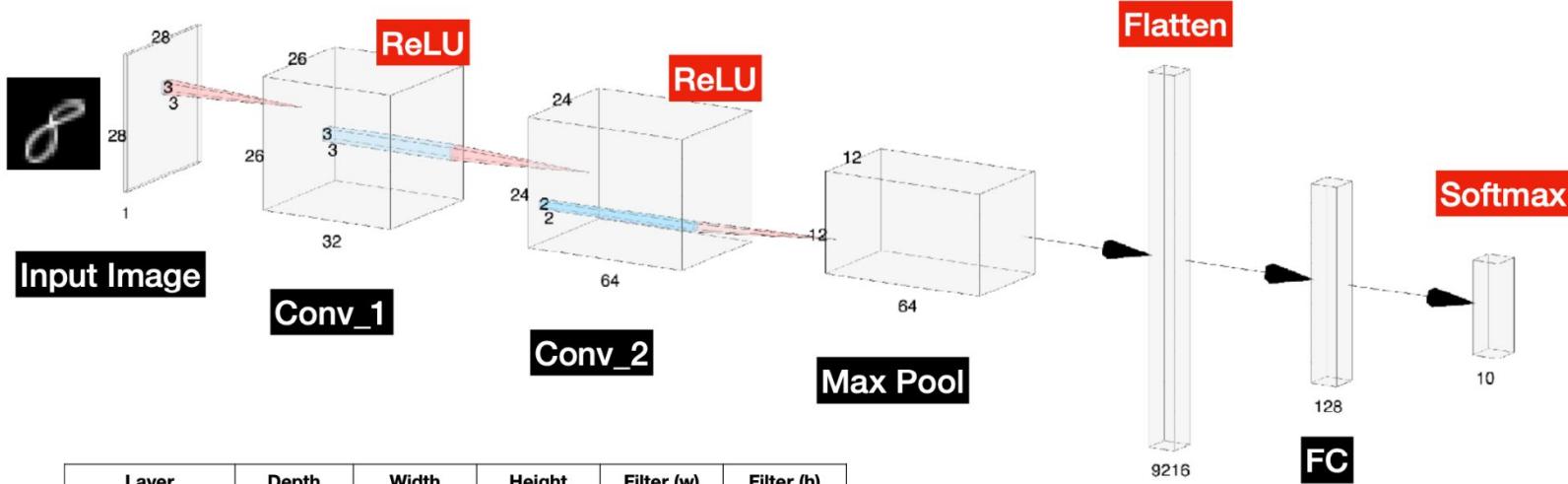
# Putting together

# 簡易CNN



# A 4 layer CNN for MNIST





### Notes:

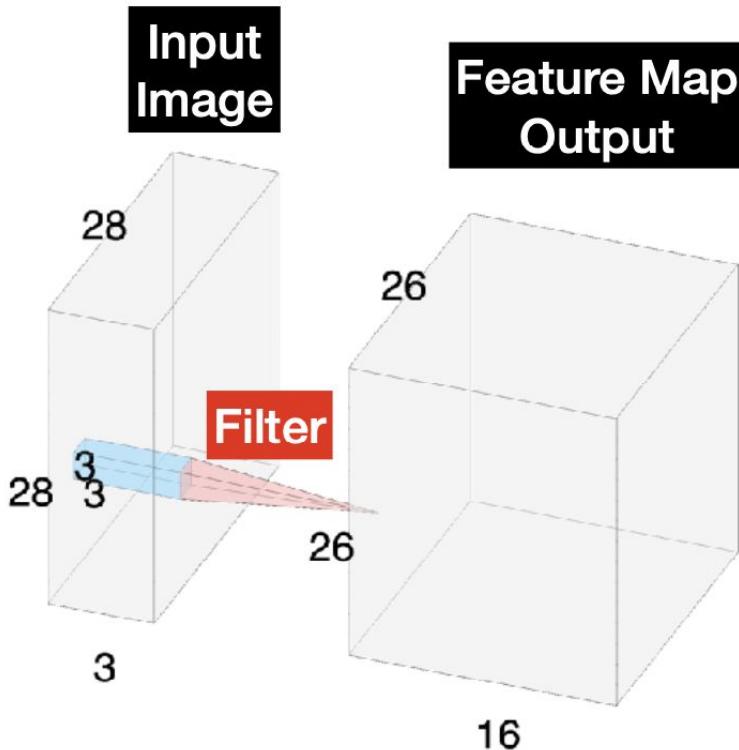
- We choose 32 & 64 Filters or Kernels for Conv\_1 & Conv\_2
- We choose to
- The Feature Maps are shown as Conv\_1 and Conv\_2
- Stride is 1
- Padding is 0 (not used)
- Max Pool Stride is 2

# 計算參數量

# 參數vs超參數

- 參數 (parameter)
  - 需要在訓練過程中讓模型去學習的變數
  - 又叫可學習的參數或權重(weights)
  - 隱藏層中的convolutional layer和fully connected layer都有權重
- 超參數 (hyperparameter)
  - 在訓練開始前由我們手動設定
  - 例如：節點數、activation function

# Conv layer上的參數量



Decided by kernel size, depth, and number of filters

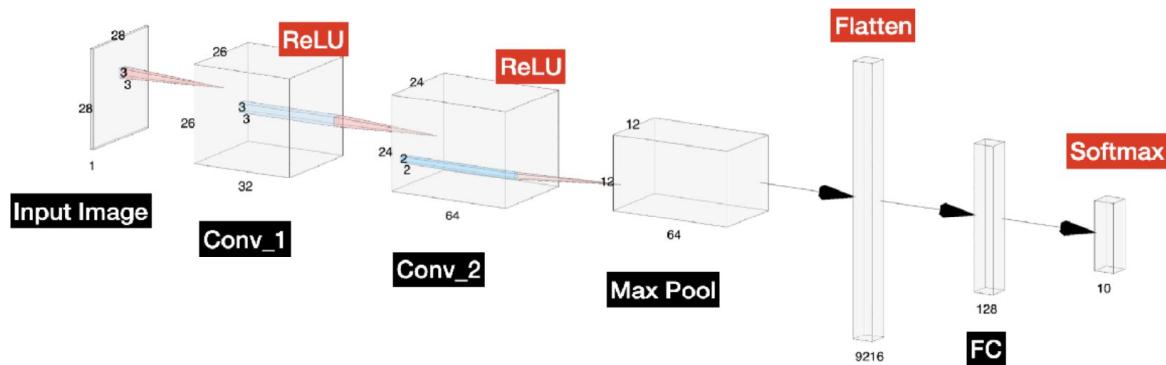
For one filter:

$$\begin{aligned} & (\text{weight} * \text{width} * \text{depth}) + \text{bias} \\ & = (3 * 3 * 3) + 1 = 10 \end{aligned}$$

For 16 filters:

$$\begin{aligned} & ((\text{weight} * \text{width} * \text{depth}) + \text{bias}) * N_f \\ & = ((3 * 3 * 3) + 1) * 16 = 160 \end{aligned}$$

# Simple CNN的參數量



Conv\_1:

$$((\text{Height} \times \text{Width} \times \text{Depth}) + \text{bias}) \times N_f$$
$$(3 \times 3 \times 1) + 1) \times 32 = 320$$

Conv\_2:

$$((\text{Height} \times \text{Width} \times \text{Depth}) + \text{bias}) \times N_f$$
$$(3 \times 3 \times 32) + 1) \times 64 = 18,496$$

沒有可訓練的參數的層：  
ReLU, Max Pool, Flatten

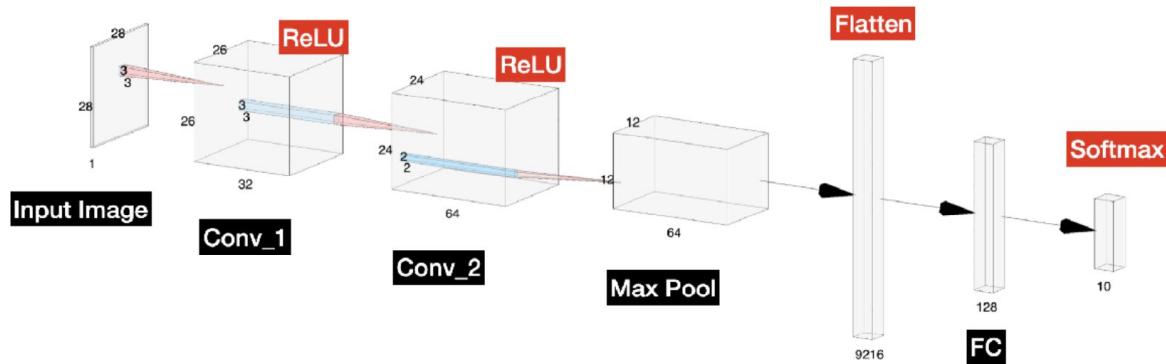
FC:

$$(\text{Length} + \text{bias}) \times N_{\text{nodes}}$$
$$(9216 + 1) \times 128 = 1,179,776$$

Final output (FC):

$$(\text{Length} + \text{bias}) \times N_{\text{nodes}}$$
$$(128 + 1) \times 10 = 1,290$$

# Simple CNN的參數量



Conv\_1:

$$((\text{Height} \times \text{Width} \times \text{Depth}) + \text{bias}) \times N_f$$
$$(3 \times 3 \times 1) + 1) \times 32 = 320$$

Conv\_2:

$$((\text{Height} \times \text{Width} \times \text{Depth}) + \text{bias}) \times N_f$$
$$(3 \times 3 \times 32) + 1) \times 64 = 18,494$$

沒有可訓練的參數的層：  
ReLU, Max Pool, Flatten

FC:

$$(\text{Length} + \text{bias}) \times N_{\text{nodes}}$$
$$(9216 + 1) \times 128 = 1,179,776$$

Final output (FC):

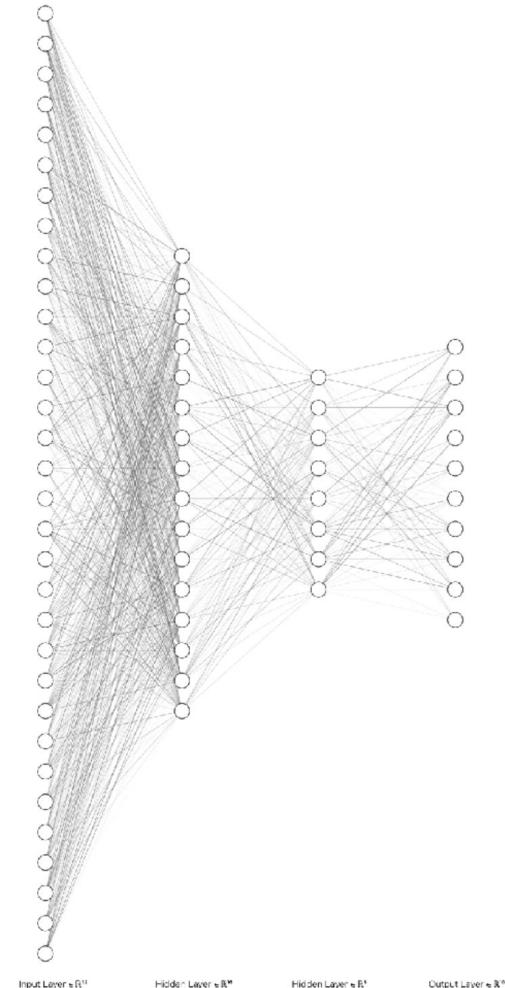
$$(\text{Length} + \text{bias}) \times N_{\text{nodes}}$$
$$(128 + 1) \times 10 = 1,290$$

$$\text{Total} = 320 + 18496 + 1179776 + 1290 = 1,199,882$$

為什麼CNN有用？

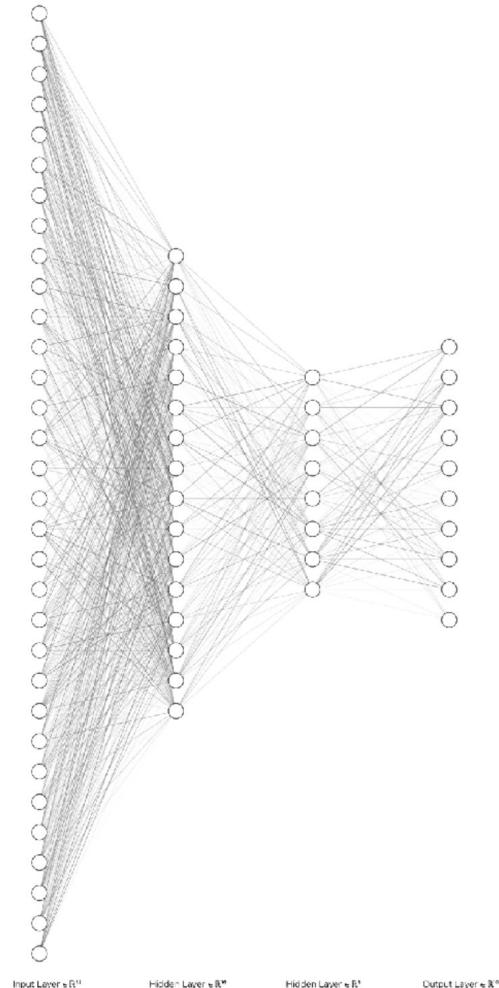
# 如果使用一般的神經網絡

- 一張 $28 \times 28$ 的灰階圖片有784個pixels
- 參照CNN裡的Conv\_1:  
32 filters for  $26 \times 26$  feature maps  
 $= 21,632$  nodes
- 參數量  $= 784 * 21632 = 16,959,488$  (一層)
- Conv\_1參數量:320



# 如果使用一般的神經網絡

- 無法擴展到資料量大的輸入  
Ex:  $320 \times 320$
- Overfitting

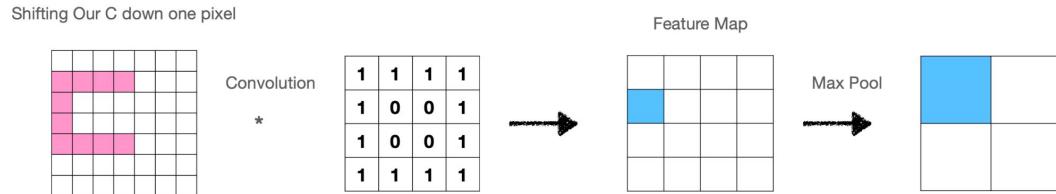
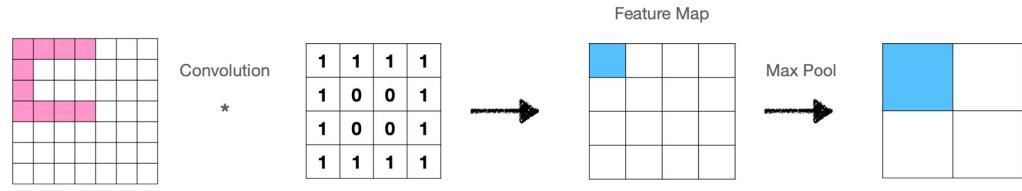


# CNN的優點

Parameter sharing:  
一個filter的參數可以用在  
image的不同地方

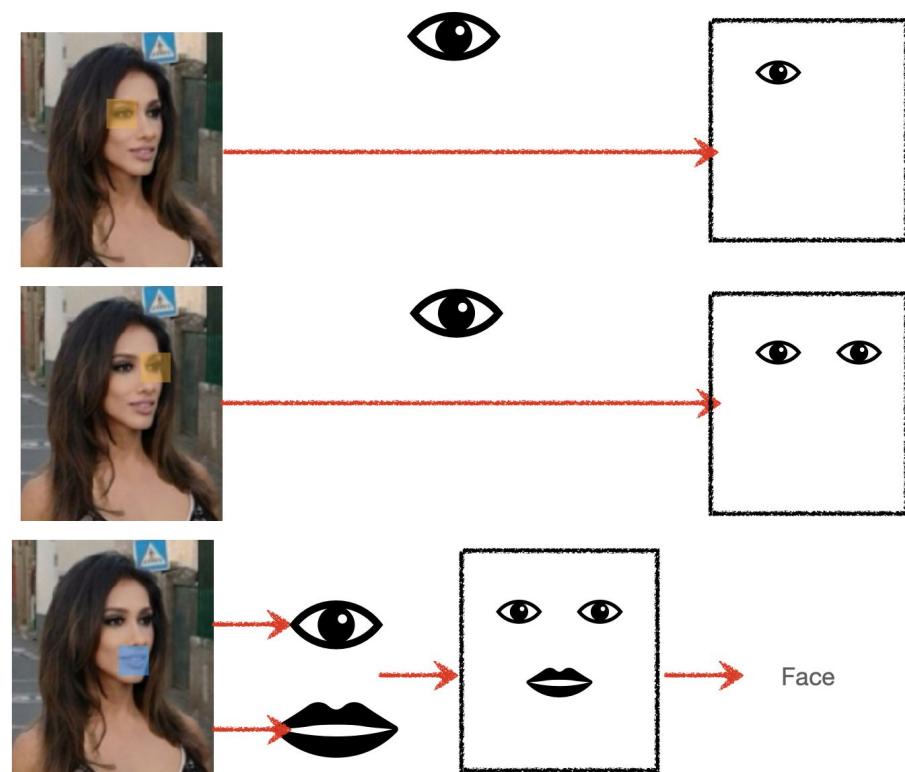
Sparsity of connection:  
在不丟失重要資訊的情況下捨  
棄一些pixels

Invariance:  
特徵可以出現在各個位置



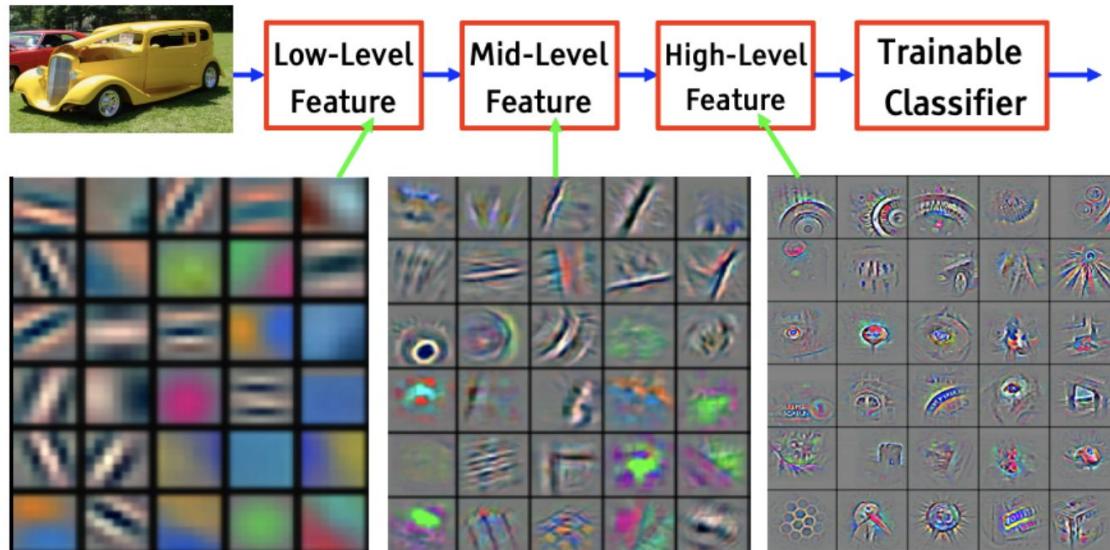
# CNN背後的假設

- 低階的特徵會靠比較近
- 特徵可以出現在不同位置
- 高階的特徵是由低階特徵組成



訓練CNN

# Conv filter學到了什麼



前面的層學低階的特徵  
Ex: 邊緣、線條

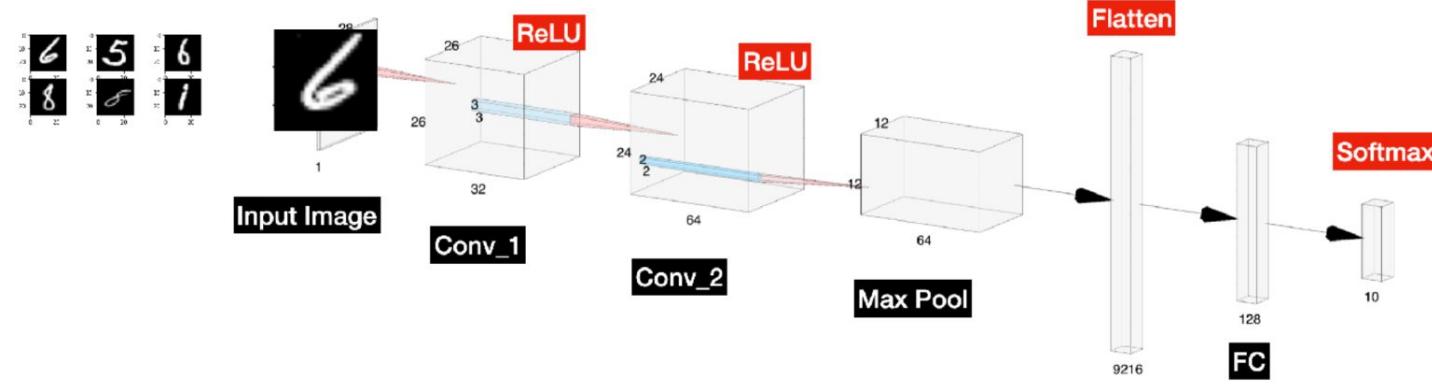
靠中間的層學中階的特徵  
Ex: 眼睛

後面的層學高階的特徵  
Ex: 鳥、花

# Filter是如何學習的？

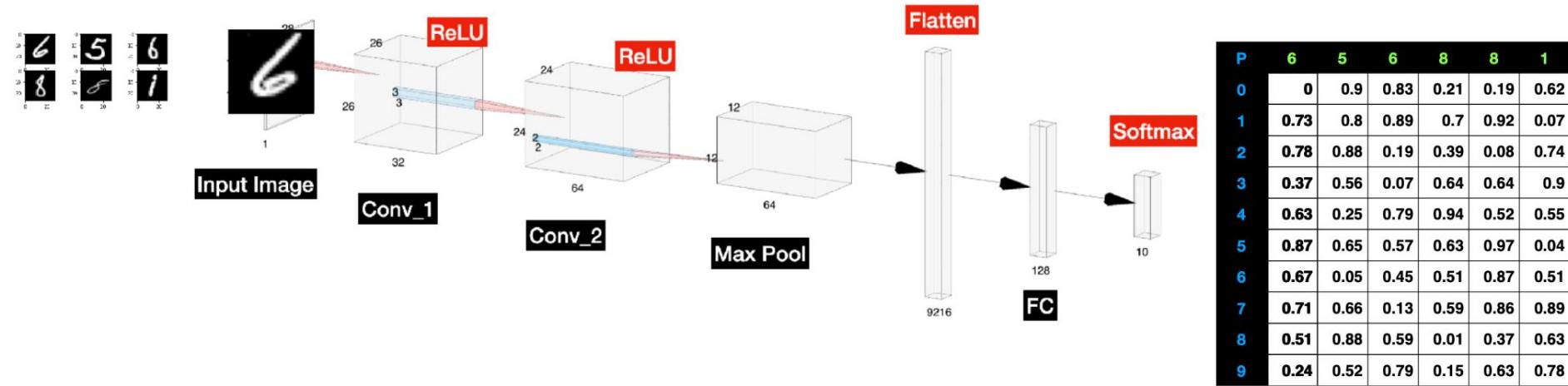
- 隨機初始權重 (weights)
- 將單張或多張圖片前向傳遞 (forward propagation)
- 計算誤差
- 使用反向傳遞 (backward propagation) 和梯度下降來更新權重
- 傳遞更多圖片 (one epoch)
- 反覆執行，直到訓練損失達到滿意程度

# 訓練過程



P	6	5	6	8	8	1
0	0	0.9	0.83	0.21	0.19	0.62
1	0.73	0.8	0.89	0.7	0.92	0.07
2	0.78	0.88	0.19	0.39	0.08	0.74
3	0.37	0.56	0.07	0.64	0.64	0.9
4	0.63	0.25	0.79	0.94	0.52	0.55
5	0.87	0.65	0.57	0.63	0.97	0.04
6	0.67	0.05	0.45	0.51	0.87	0.51
7	0.71	0.66	0.13	0.59	0.86	0.89
8	0.51	0.88	0.59	0.01	0.37	0.63
9	0.24	0.52	0.79	0.15	0.63	0.78

# 訓練過程



需要一個方法來告訴模型結果好壞

# Loss function

# Cross entropy

進一步細分

- Binary cross entropy -> binary class
- Categorical cross entropy -> multiclass

Class	Predicted Probabilities	Ground Truth
0	0.1	0
1	0.2	0
2	0.1	0
3	0.05	0
4	0.05	0
5	0.05	0
6	0.05	0
7	0.3	1
8	0.05	0
9	0.05	0

Cross entropy: 計算兩個機率分佈間的差異

$$L = -y \cdot \log(\hat{y})$$

正確值      預測值

Class	Predicted Probabilities	Ground Truth
0	0.3	0
1	0.6	1
2	0.1	0

- $L = -y \cdot \log(\hat{y})$
- $L = -(0 \times \log(0.3) + 1 \times \log(0.6) + 0 \times \log(0.1))$
- $L = -(0 + 1 \times -0.222 + 0) = 0.222$

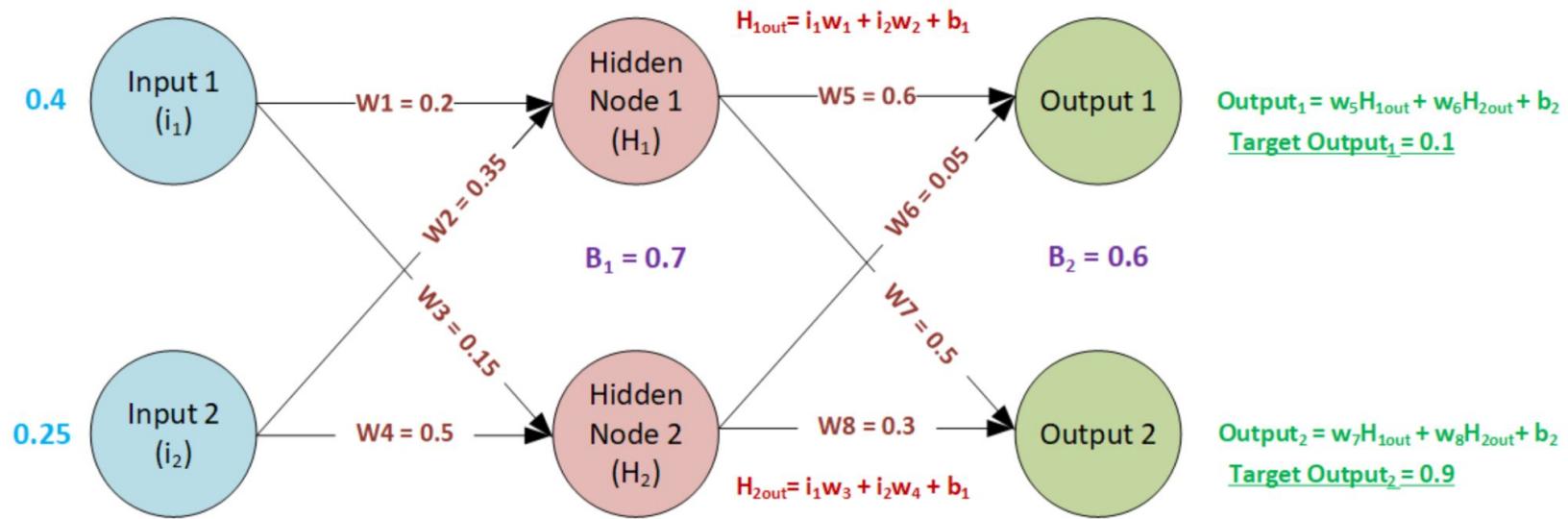
# 其他loss functions

- 迴歸任務 (regression): Mean square error (MSE)
  - Mean Square Error (MSE) = (目標 - 預測)<sup>2</sup>
- 更多：
  - L1, L2
  - Hinge Loss
  - Mean Absolute Error (MAE)

$$MSE = \frac{1}{n} \sum (Y_i - \hat{Y}_i)^2$$

# Backward propagation

# Backward propagation (反向傳遞)

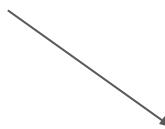


稍微增加  $w_5$  到 0.6001 或稍微減少到 0.5999, 哪個才能減少訓練損失

# Chain rule

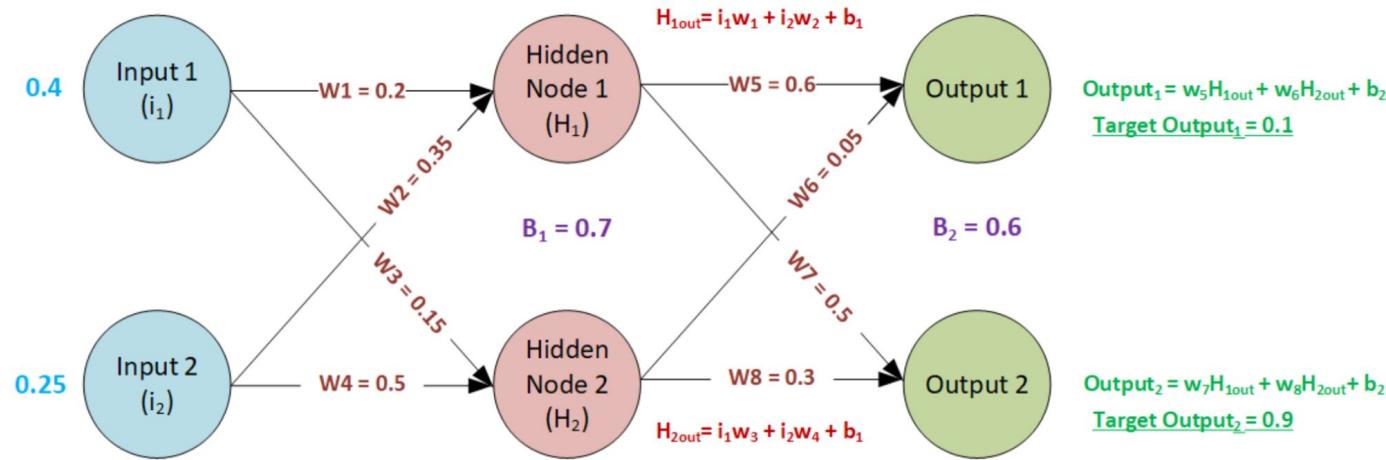
假設有兩個函式  $y = f(u)$ ,  $u = f(x)$ , 則  $y$  的導數為

$$\frac{dy}{dx} = \frac{dy}{du} \times \frac{du}{dx}$$



If  $y$  is the total loss and  $x$  is a weight, then  $dy/dx$  = the gradient of  $x$

# Example



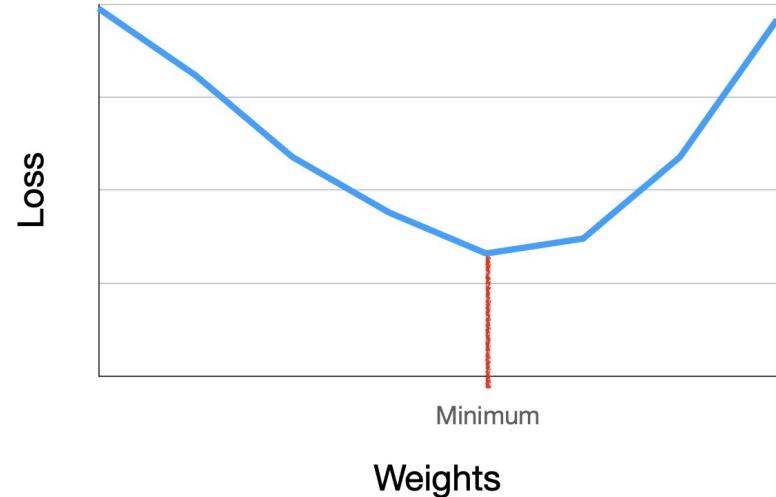
稍微改變  $W_5$  會對訓練損失影響多少？

$$\frac{dE_T}{dW_5} \quad E_T = \text{總訓練損失} \quad W_5 = W_{5old} - \lambda \times \frac{dE_T}{dW_5}$$

# Gradient descent

# Loss function

- 使用反向傳遞更新參數
  - $wx + b$
- 找到一組產生最低訓練損失的參數
- 更新的過程就是梯度下降 (gradient descent)

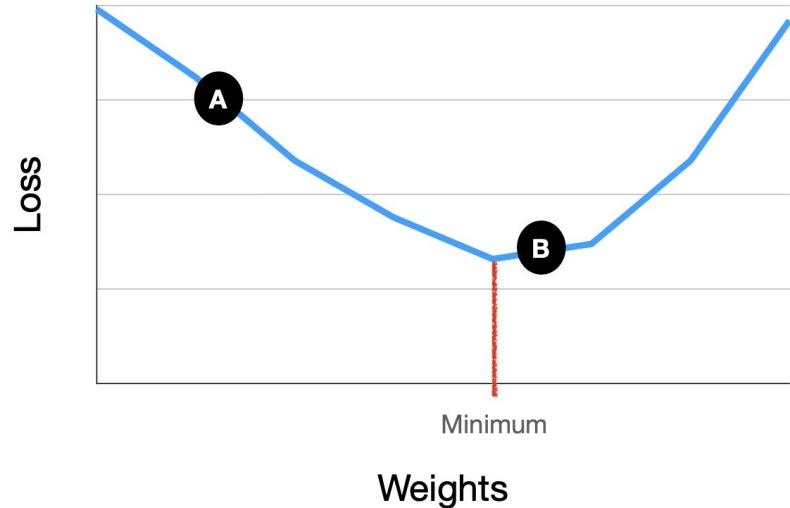


# Gradient descent

- 稍微改變W會對損失E影響多少？

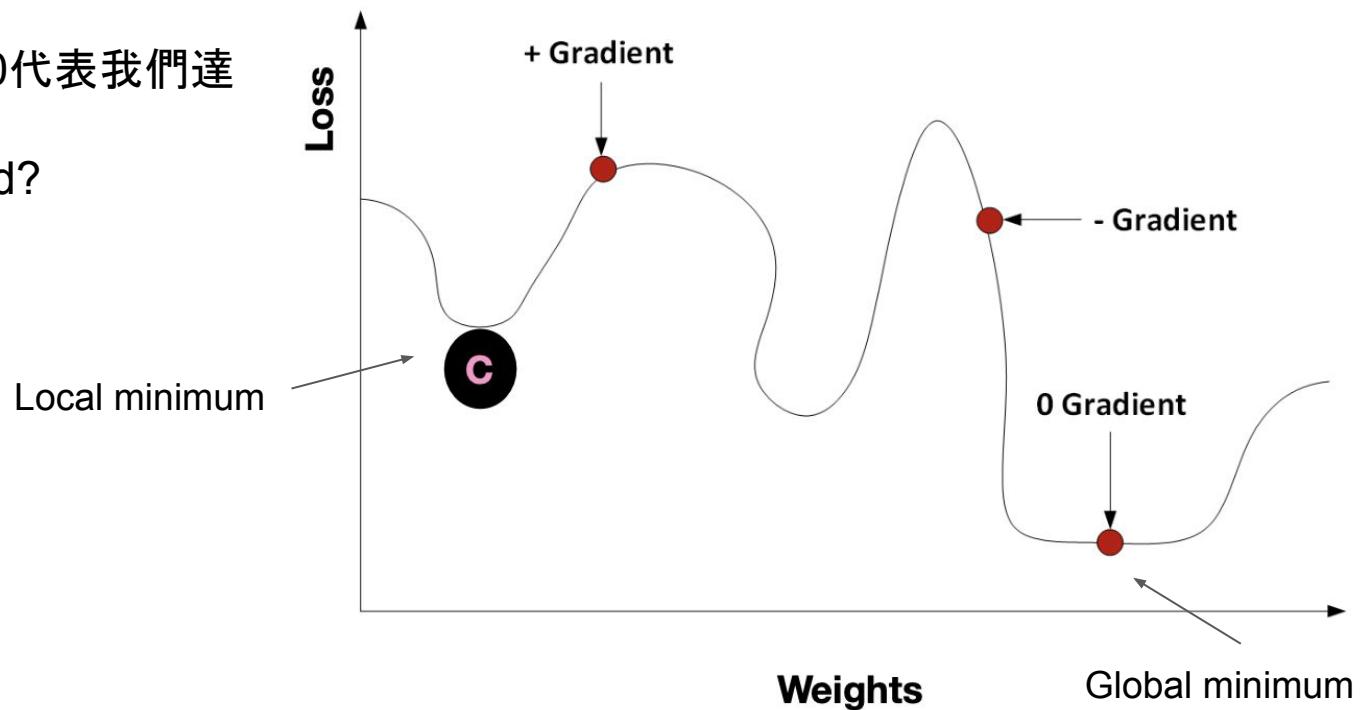
$$\text{Gradient} = \frac{dE}{dw}$$

- 如果gradient為正數，代表損失隨著權重增加而增加
- 如果gradient為負數，代表損失隨著權重增加而減少
- Gradient的值會告訴模型該往哪邊移動

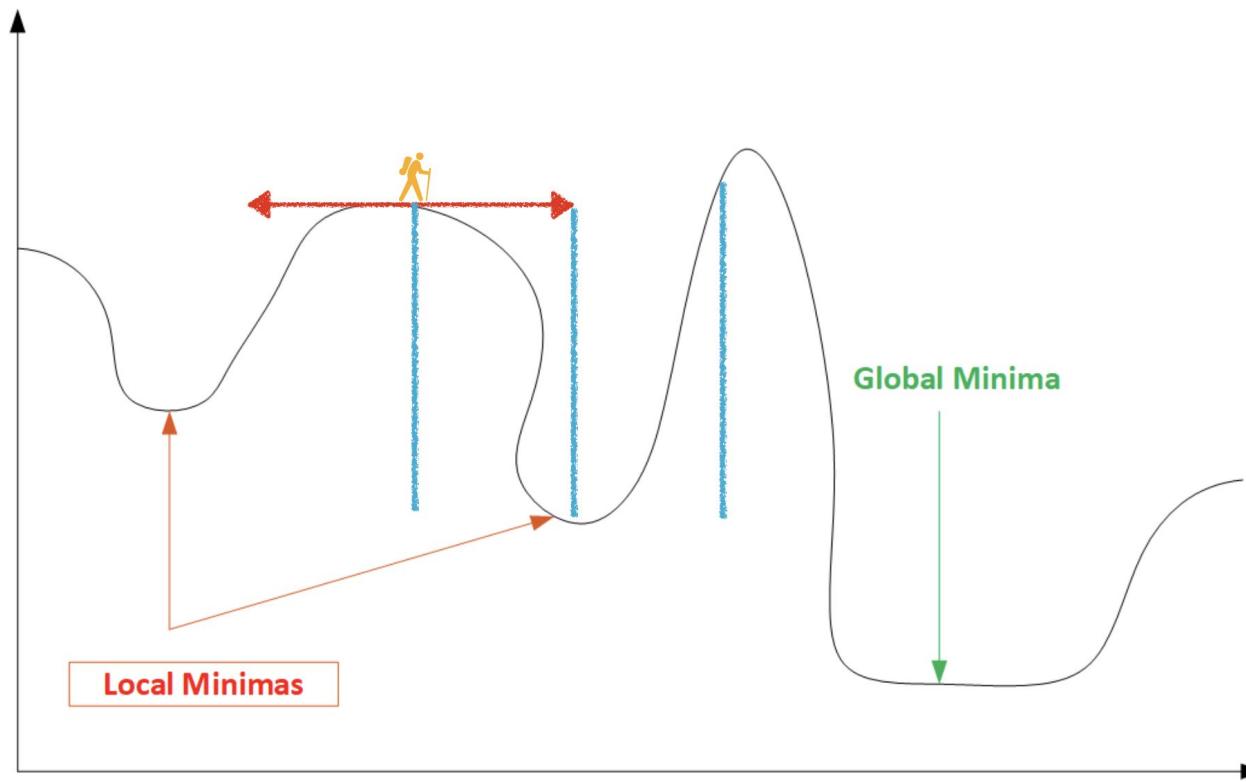


# If gradient becomes 0

- Gradient為0代表我們達到低點了
- Good or bad?



# Step size is important



# Learning rate

$$W_5 = W_{5old} - \lambda \times \frac{dE_T}{dW_5}$$

- $\lambda$  is learning rate
- Learning rate決定我們更新的步幅
- 找到最佳的learning rate可以避免模型卡在區域最佳解以及超過全域最佳解

# Gradient descent的方法

- Batch gradient descent - 每次看完全部資料才更新權重
  - 增加計算資源以及時間
- Stochastic gradient descent - 每次看完一筆資料就更新權重
  - 計算很快
  - 不穩定且有可能訓練很久
- Mini-batch gradient descent - 每次看完一批資料就更新權重
  - 增加訓練速度並達到global minimum
  - Batch size → 8~256

# Gradient descent的問題

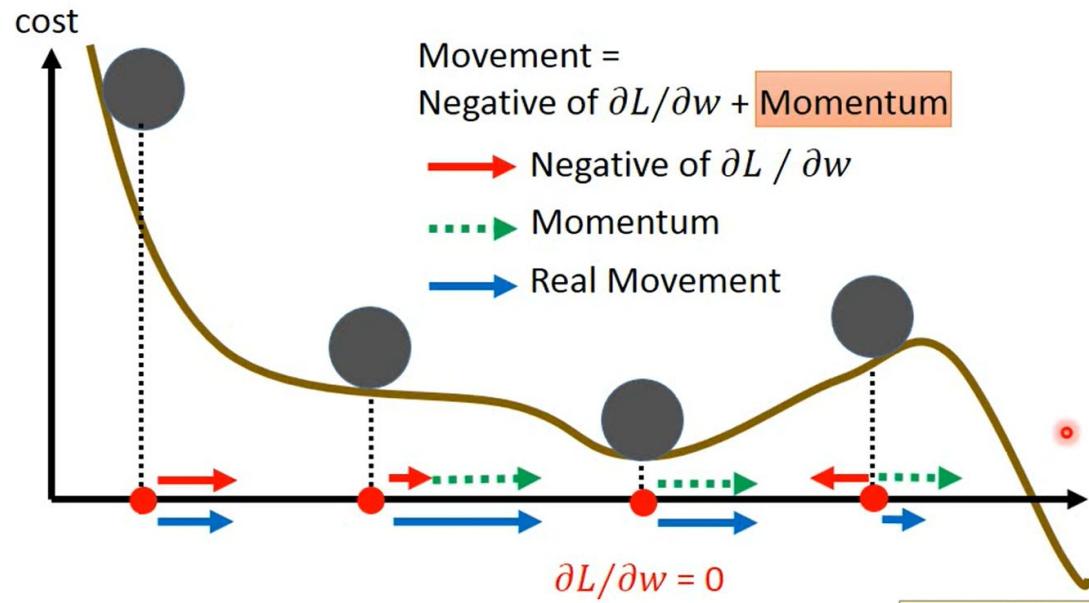
- 需要決定learning rate
- 使用固定的learning rate容易使模型卡在區域最佳解以及超過全域最佳解

# Solution

- Momentum
- Learning rate schedule

## Momentum

Still not guarantee reaching global minima, but give some hope .....



Created with EverCam.  
<http://www.camdem.com>

# Quiz next week

- Scope: CNN
- Paper based (bring a pen!)

# Implementation

Download notebook at:

<https://github.com/albert831229/nchu-computer-vision/tree/main/113/day>