# CSCE 636 Neural Networks (Deep Learning)

Lecture 7: Deep Learning for Computer Vision

Anxiao (Andrew) Jiang

Based on the interesting lecture of Prof. Hung-yi Lee,
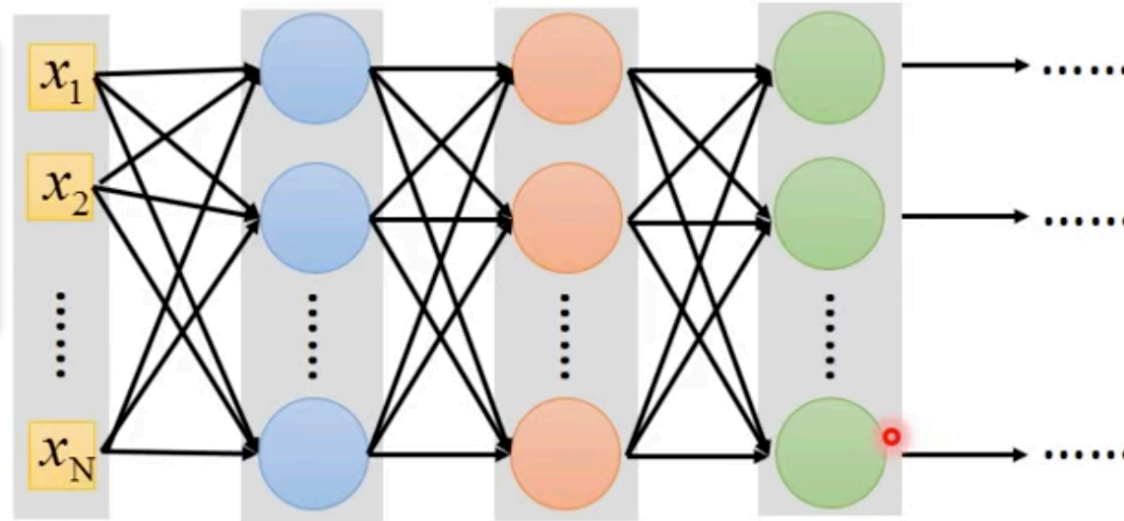https://www.youtube.com/watch?v=FrKWiRv254g&list=PLJV_el3uVTsPy9oCRY30oBPNLCo89yu49&index=19

# Convolutional Neural Network (CNN)
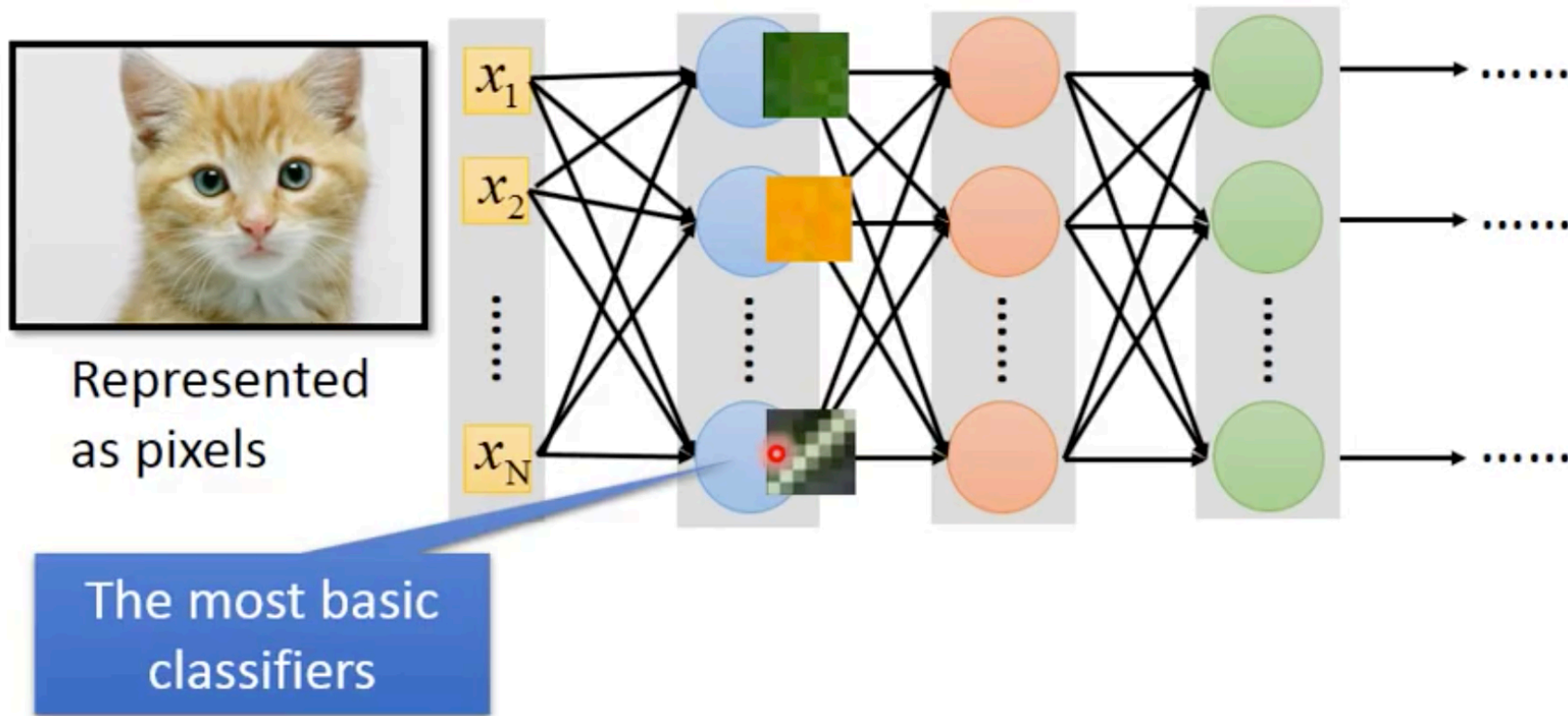
# Why CNN for Image? [Zeiler, M. D., *ECCV 2014*]



Represented as pixels

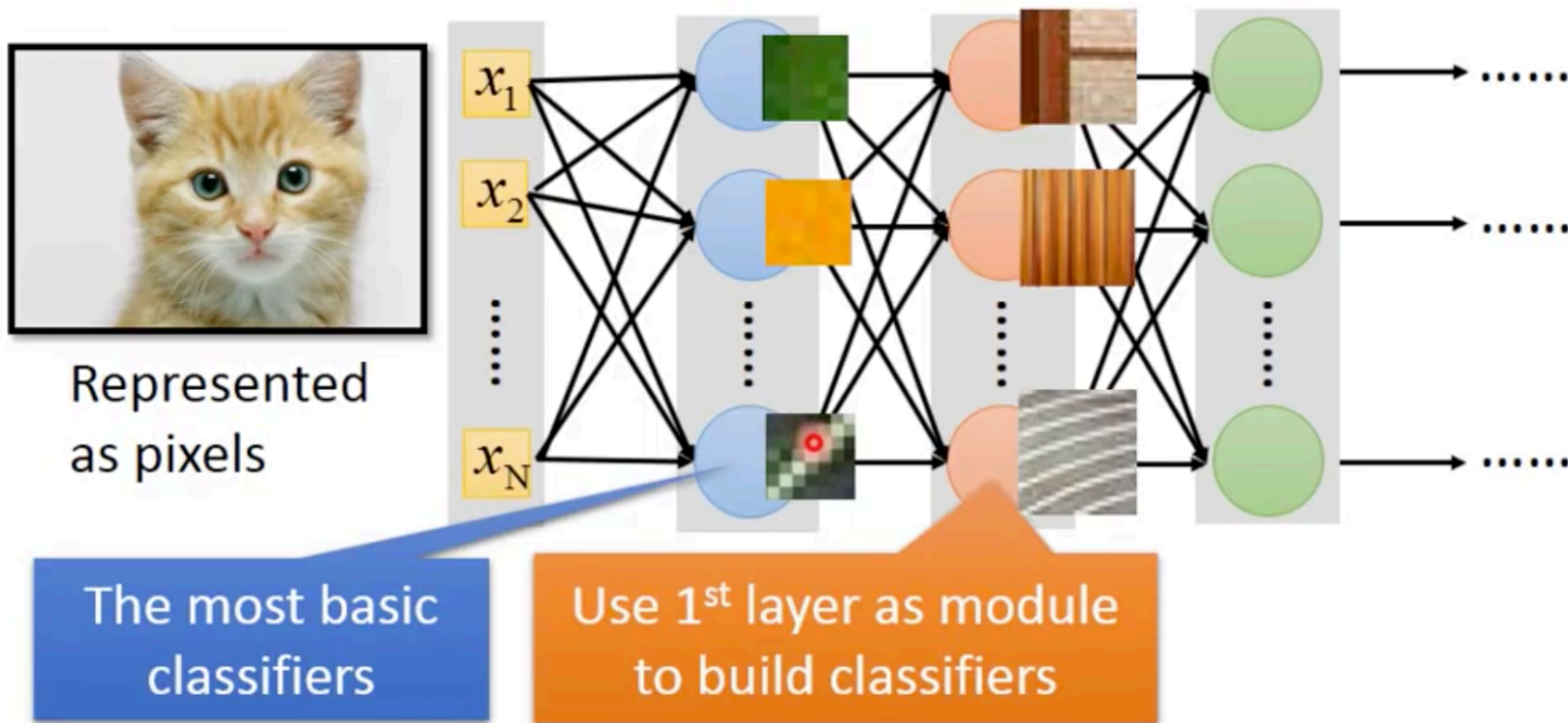# Why CNN for Image? [Zeiler, M. D., *ECCV 2014*]



Represented as pixels

The most basic classifiers

# Why CNN for Image?

[Zeiler, M. D., *ECCV 2014*]



Represented as pixels

The most basic classifiers

Use 1st layer as module to build classifiers

# Why CNN for Image?

[Zeiler, M. D., *ECCV 2014*]



Represented as pixels

$x_1$

$x_2$

$x_N$

The most basic classifiers

Use 1st layer as module to build classifiers

Use 2nd layer as module ......

Too many weights in a dense network!

# Why CNN for Image

- Some patterns are much smaller than the whole image.

A neuron does not have to see the whole image to discover the pattern.



有沒有某一個 pattern 出現

# Why CNN for Image

- Some patterns are much smaller than the whole image

A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters
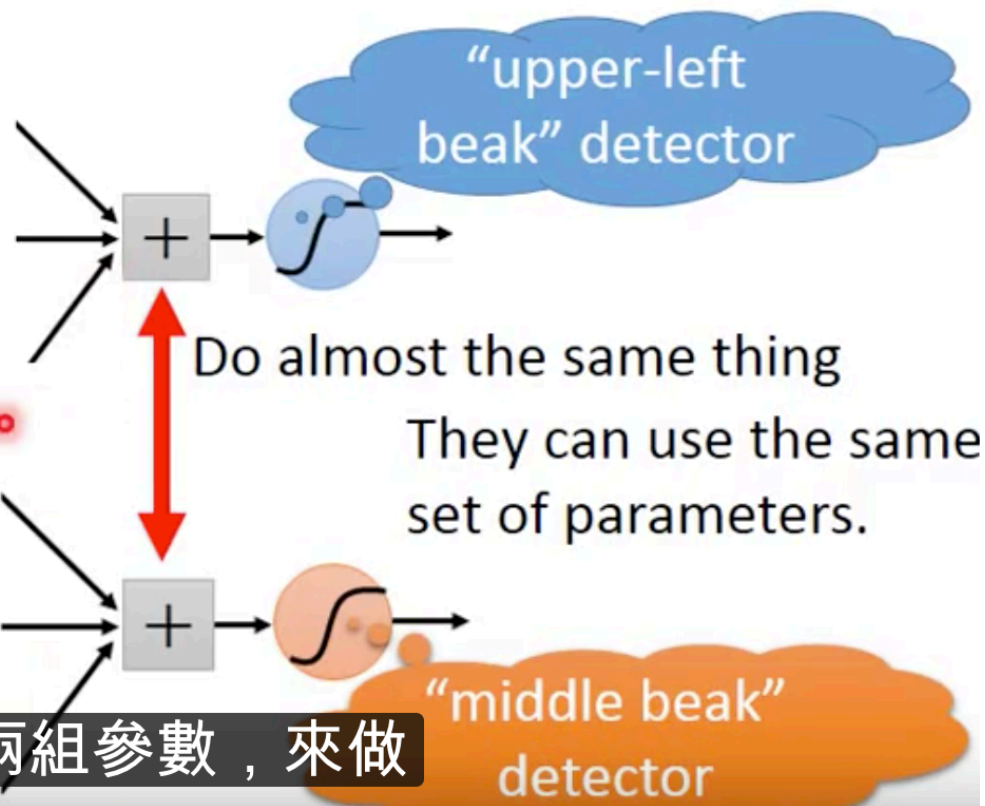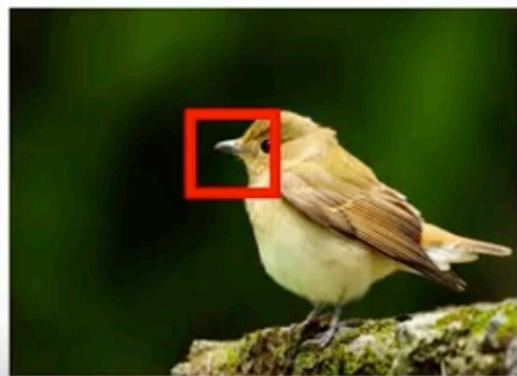


整張完整的圖

"beak" detector

# Why CNN for Image

- The same patterns appear in different regions.



同樣的這個 pattern

# Why CNN for Image

- The same patterns appear in different regions.

# Why CNN for Image
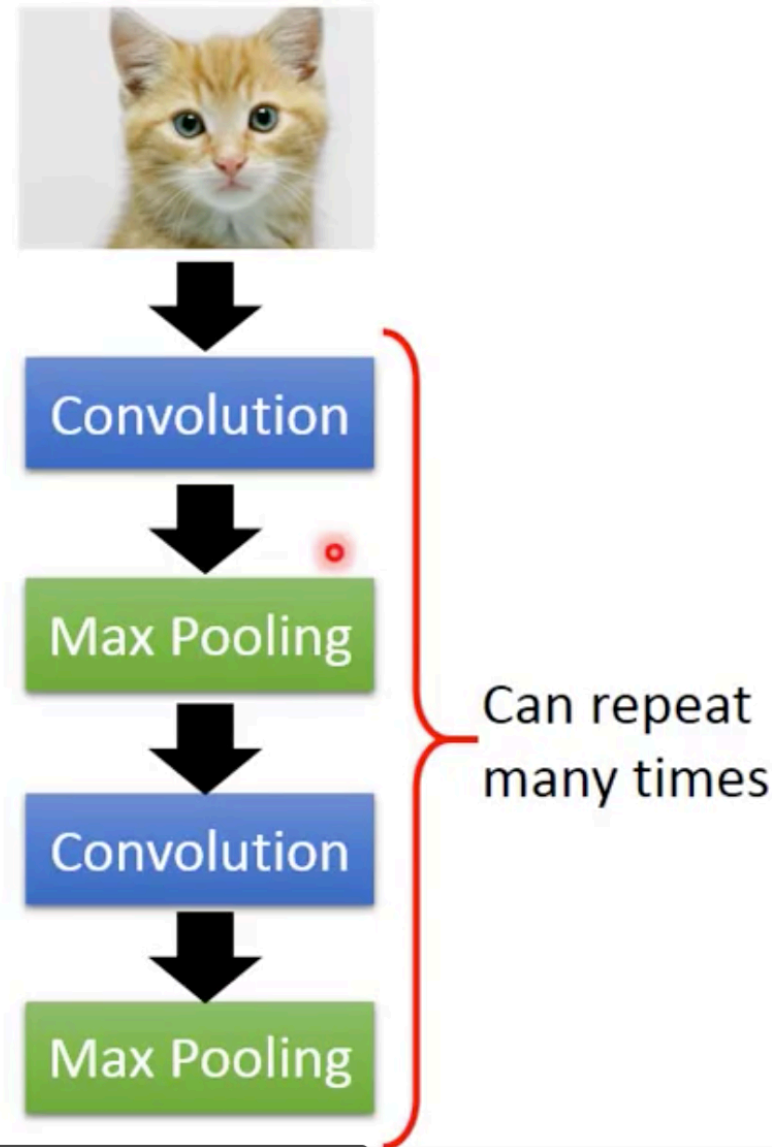
- Subsampling the pixels will not change the object

bird



subsampling

bird

We can subsample the pixels to make image smaller

# The whole CNN



Convolution

Max Pooling

Convolution

Max Pooling

Can repeat
many times

The whole CNN

cat dog ......

Fully Connected Feedforward network

Flatten

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

3 個對影像處理的觀察

# CNN – Convolution

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

# CNN – Convolution

| | | | | | |
|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

| | | |
|---|---|---|
| 1 | -1 | -1 |
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| | | |
|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

⋮

# CNN – Convolution

A filter corresponds to a set of neurons

**Those are the network parameters to be learned.**

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

| 1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1
Matrix

| -1 | 1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2
Matrix

# CNN – Convolution

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

How does a filter operate?

# CNN – Convolution

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Do inner product (dot product)

# CNN – Convolution

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

3

# CNN – Convolution

| 1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

3

# CNN – Convolution

Filter 1

|  |  |  |
|---|---|---|
| 1 | -1 | -1 |
| -1 | 1 | -1 |
| -1 | -1 | 1 |

stride=1

|  |  |  |  |  |  |
|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

( 3 )  ( -1 )

# CNN – Convolution

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

If stride=2

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

3    -3

We set stride = 1 in the following slides.

# CNN – Convolution

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

3    -1

6 x 6 image

# CNN – Convolution

| 1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

3    -1    -3

# CNN – Convolution

Filter 1

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

stride=1

6 x 6 image

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

3   -1   -3   -1

# CNN – Convolution

Filter 1

| 1  | -1 | -1 |
|----|----|----|
| -1 | 1  | -1 |
| -1 | -1 | 1  |

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

3    -1    -3    -1

-3

# CNN – Convolution

Filter 1

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

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

你就得到 -1

# CNN – Convolution



Filter 1

$$
\begin{array}{|c|c|c|}
\hline
1 & -1 & -1 \\
\hline
-1 & 1 & -1 \\
\hline
-1 & -1 & 1 \\
\hline
\end{array}
$$

stride=1

$$
\begin{array}{|c|c|c|c|c|c|}
\hline
1 & 0 & 0 & 0 & 0 & 1 \\
\hline
0 & 1 & 0 & 0 & 1 & 0 \\
\hline
0 & 0 & 1 & 1 & 0 & 0 \\
\hline
1 & 0 & 0 & 0 & 1 & 0 \\
\hline
0 & 1 & 0 & 0 & 1 & 0 \\
\hline
0 & 0 & 1 & 0 & 1 & 0 \\
\hline
\end{array}
$$

6 x 6 image

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

出現最大的值

# CNN – Convolution

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

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

一打 filter，它會有一大推的 filter

# CNN – Convolution

|   |   |   |
|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

**Do the same process for every filter**

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

比如說，這裡會有一個 filter 2

# CNN – Convolution

**Filter 2**

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Do the same process for every filter

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

再做 inner product 得到 -1

# CNN – Convolution

Filter 2

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Do the same process for every filter

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

再挪動一格，得到 1，就這樣，以此類推

# CNN – Convolution

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Do the same process for every filter

| -1 | -1 | -1 | -1 |
|----|----|----|----|
| -1 | -1 | -2 | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

# CNN – Convolution

Filter 2

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

stride=1

Do the same process for every filter

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Feature Map

| -1 | -1 | -1 | -1 |
|----|----|----|----|
| -1 |    |    | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

4 x 4 image

Created with Ever

# CNN – Colorful image

Colorful image

# Convolution v.s. Fully Connected



image

convolution

Fully-connected

$x_1$

$x_2$

$x_{36}$

有什麼關係？

Filter 1

6 x 6 image

Filter 1

6 x 6 image

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

6 x 6 image

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

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1

拉直變成這樣子 1
⋮

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| 1 | 0 | 0 | 0 | 0 | |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | |
| 0 | 0 | 1 | 1 | 0 | |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

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

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1

3

然後，你有一個 neuron 的 output 是 3

Filter 1

6 x 6 image

你看這個 filter

Created with E

Filter 1

6 x 6 image

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1

3

Fewer parameters!

Only connect to 9 input, not fully connected

你有 36 個 pixel 當作 input

Created with EverCam

Filter 1

6 x 6 image

Less parameters!

1: 1
2: 0
3: 0
4: 0
⋮
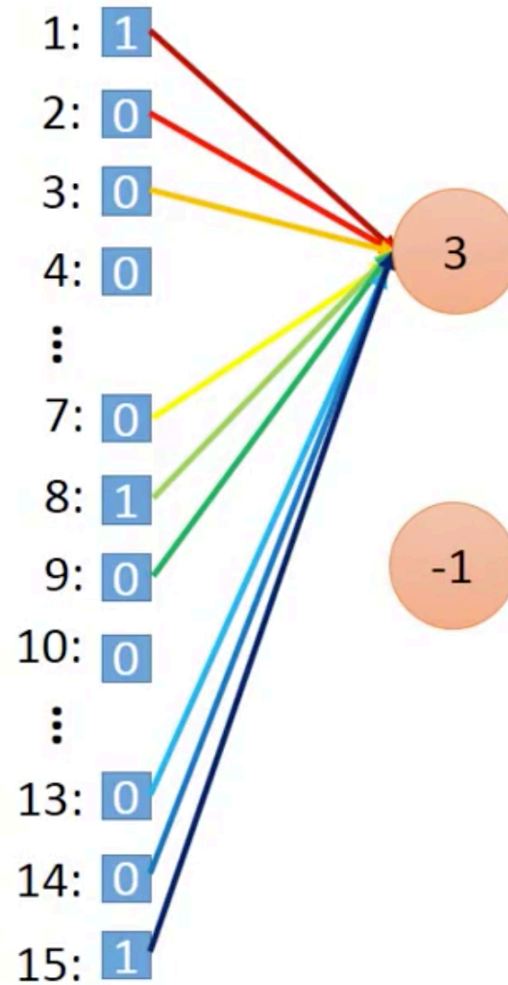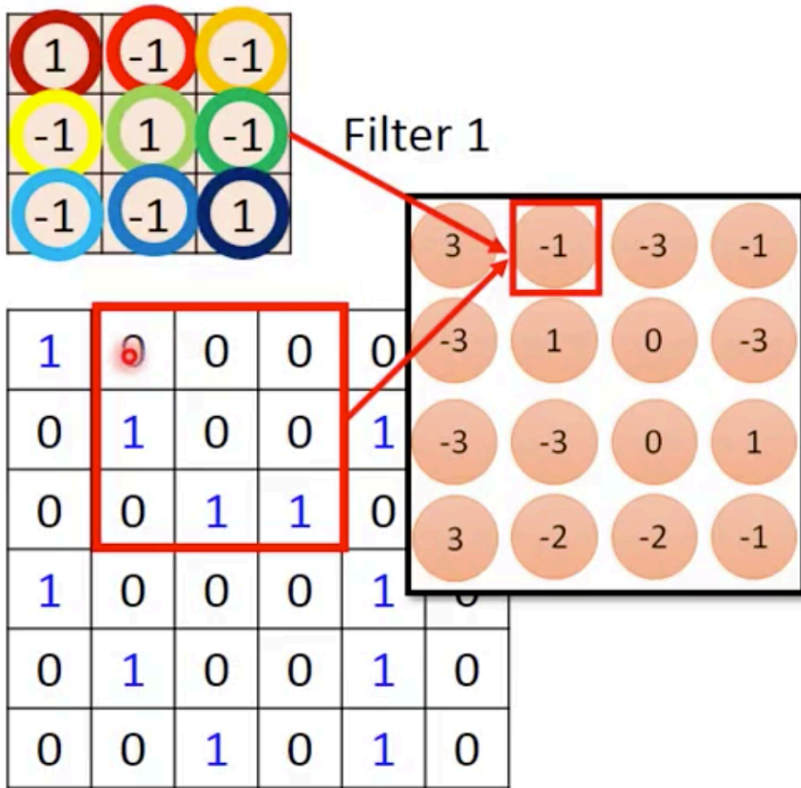7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
⋮

3

發生甚麼事呢？

Filter 1
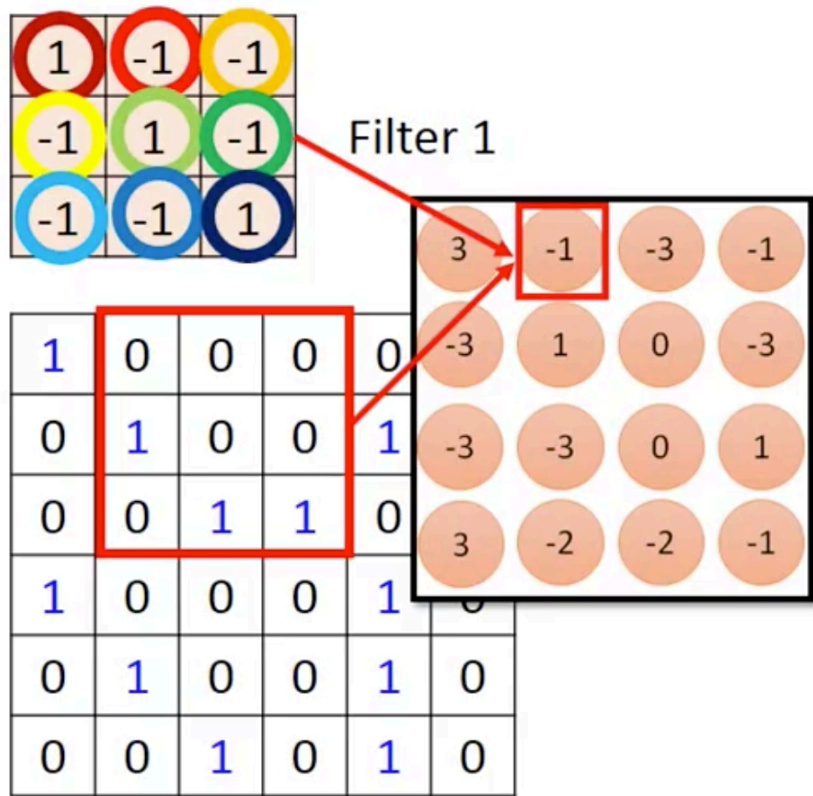
6 x 6 image

Less parameters!

那這個 neuron 連接到哪些 input 的 weight 呢

Filter 1

1  -1  -1
-1  1  -1
-1  -1  1

| 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

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

Less parameters!

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1

3

-1

這個框起來的地方，它正好就對應到 pixel 2, 3, 4

Filter 1

6 x 6 image

Less parameters!

1: 1
2: 0
3: 0
4: 0
...
7: 0
8: 1
9: 0
10: 0
...
13: 0
14: 0
15: 1

3

-1

How to train shared weight?
Average their gradients.

Shared weights

但是，當我們做這個

Created with EverCa

The whole CNN

cat dog ......

Fully Connected Feedforward network

Flatten

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

Max pooling 是甚麼呢？

# CNN – Max Pooling



Filter 1

Filter 2

所以，3, -1, 3, 1，我就選 3

# CNN – Max Pooling

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

| | |
|---|---|
| 3 | 0 |
| 3 | 1 |

| | |
|---|---|
| -1 | 1 |
| 0 | 3 |

它其實可以的

# CNN – Max Pooling



| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Conv

Max Pooling

| -1 | 1 |
|----|---|
| 0  | 3 |

2 x 2 image

Each filter is a channel

# The whole CNN



-1   1

0   3

**A new image**

Smaller than the original image

Convolution

↓

Max Pooling

↓

Convolution

↓

Max Pooling

Can repeat many times

# The whole CNN



-1    1

0    3

**A new image**

Smaller than the original image

The number of the channel is the number of filters

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

# The whole CNN

cat dog ......

Fully Connected
Feedforward network

Flatten

我們就繼續

Convolution

Max Pooling

A new image

Convolution

Max Pooling

A new image

# Flatten

## CNN in Keras

Only modified the **network structure** and **input format (vector -> 3-D tensor)**

input

↓

Convolution

↓

Max Pooling

↓

Convolution

↓

Max Pooling

## CNN in Keras

Only modified the **network structure** and **input format (vector -> 3-D tensor)**

```
model2.add( Convolution2D( 25,3,3,
            input_shape=(1,28,28) ) )
```

| 1 | -1 | |
|---|----|---|
| -1 | 1 | |
| -1 | -1 | |

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |
| -1 | 1 | -1 |

...... There are **25** **3x3** filters.

Input_shape = ( 1 , 28 , 28 )

1: black/weight, 3: RGB    28 x 28 pixels
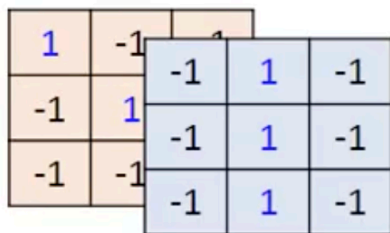
input

→ Convolution

→ Max Pooling

→ Convolution

→ Max Pooling

# CNN in Keras

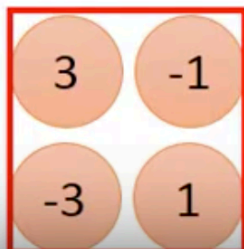Only modified the *network structure* and *input format (vector -> 3-D tensor)*

```
model2.add( Convolution2D( 25,3,3,
            input_shape=(1,28,28) ) )
```

There are **25** **3x3** filters.

Input_shape = ( 1 , 28 , 28 )

1: black/weight, 3: RGB     28 x 28 pixels

```
model2.add(MaxPooling2D((2,2)))
```

| 3 | -1 |
|---|----|
| -3 | 1 |

input

Convolution

Max Pooling

Convolution

Max Pooling

我們把 2*2 的這個

29:55 / 1:19:28

Created with EverCam

# CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D tensor)*

```
model2.add( Convolution2D( 25,3,3,
              input_shape=(1,28,28) ) )
```
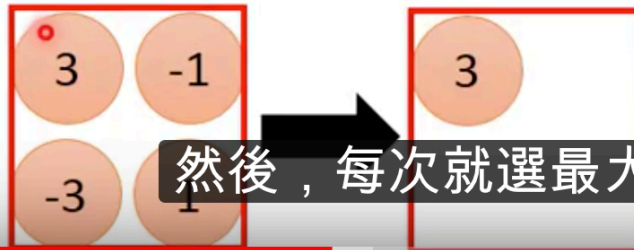
| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

...... There are 25 3x3 filters.

Input_shape = ( 1 , 28 , 28 )

1: black/weight, 3: RGB    28 x 28 pixels

```
model2.add(MaxPooling2D((2,2)))
```

| 3 | -1 |
|---|----|
| -3 | 1 |

→ | 3 |

**input**

↓

Convolution

↓

Max Pooling

↓

Convolution

↓

Max Pooling

然後，每次就選最大的那一個

## CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D tensor)*

input

1 x 28 x 28

```
model2.add( Convolution2D( 25,3,3,
         input_shape=(1,28,28) ) )
```

25 x 26 x 26

Convolution
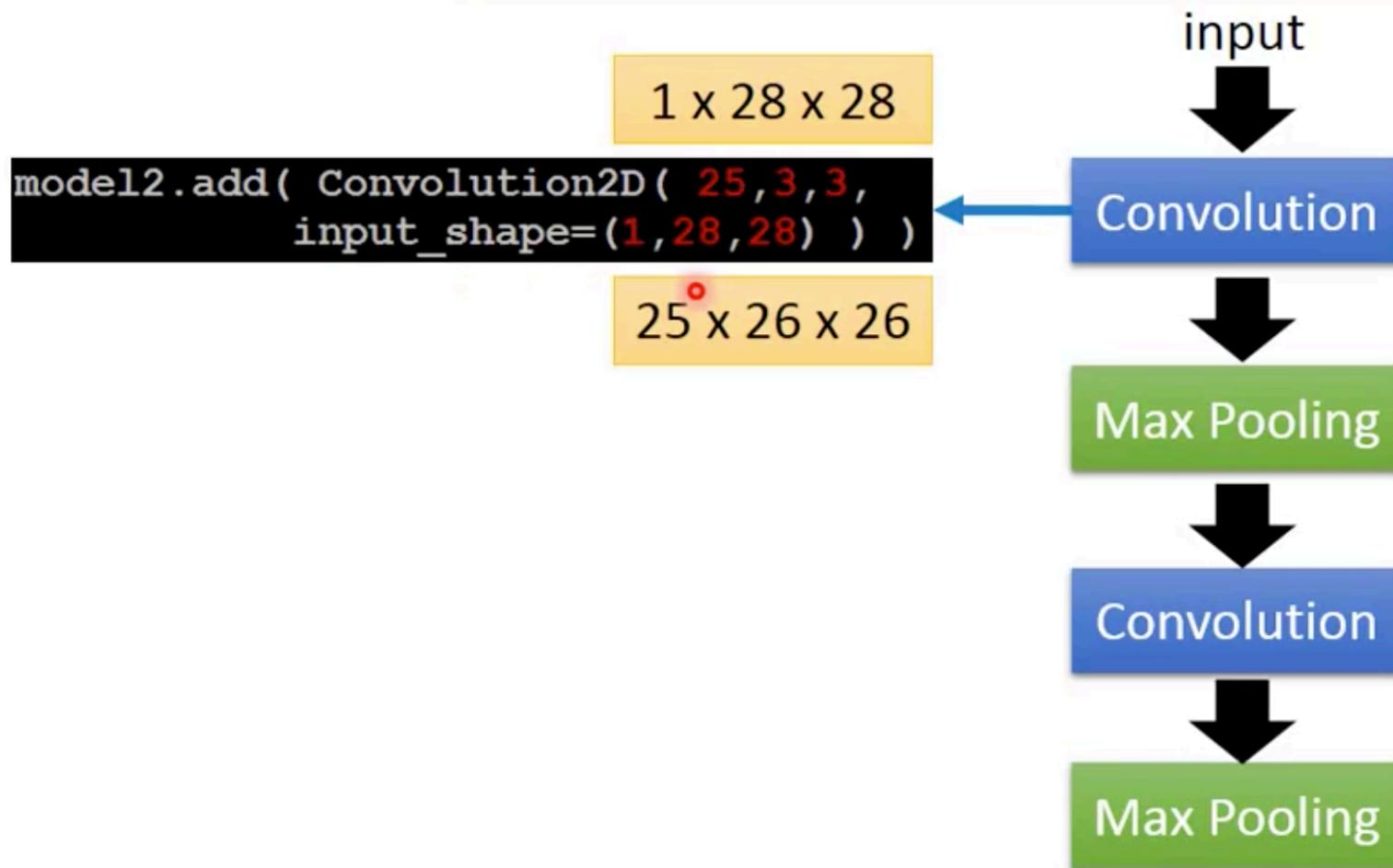
Max Pooling

Convolution

Max Pooling

# CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D tensor)*

input

↓

1 x 28 x 28

```
model2.add( Convolution2D( 25,3,3,
            input_shape=(1,28,28) ) )
```

Convolution

25 x 26 x 26
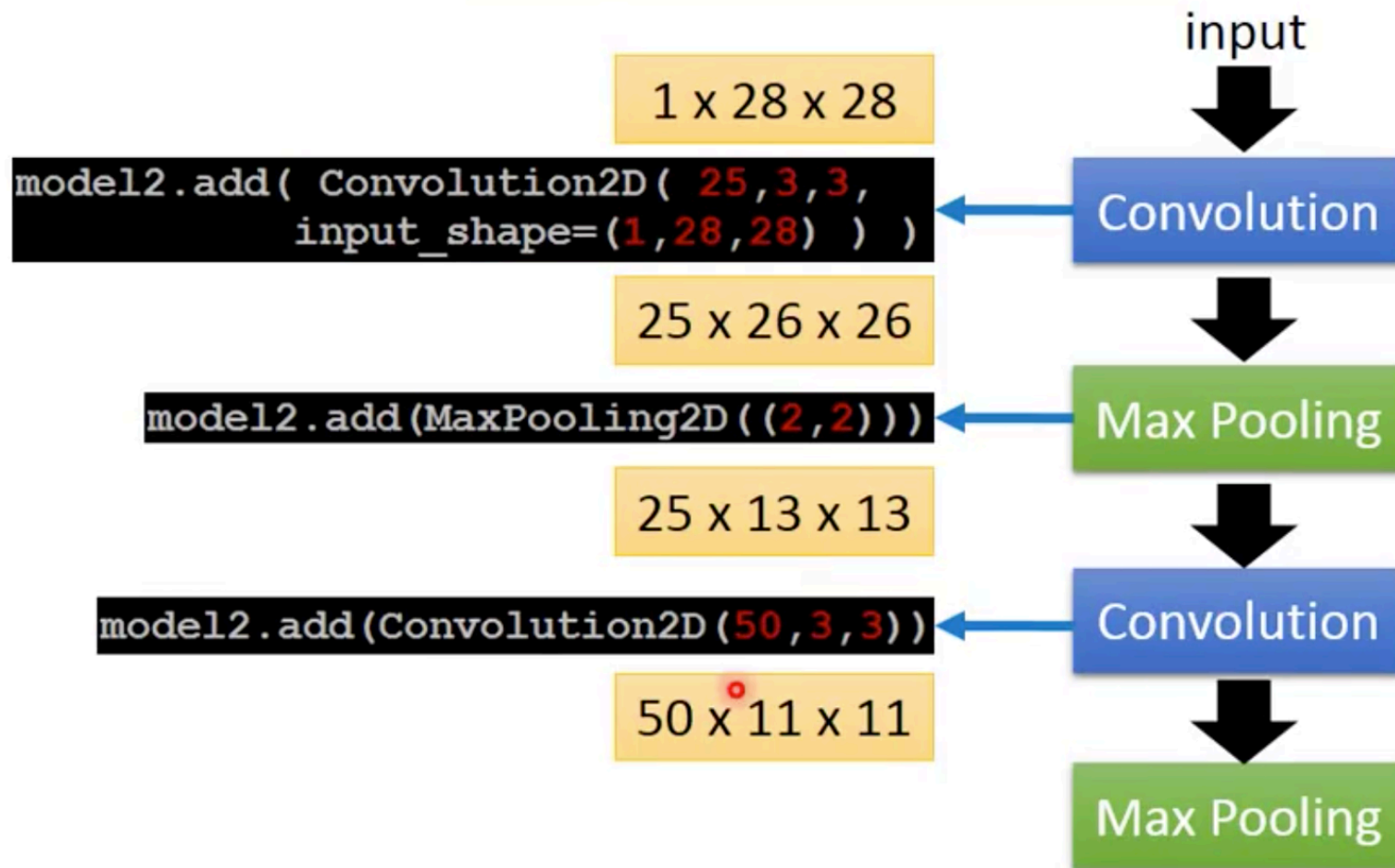
```
model2.add(MaxPooling2D((2,2)))
```

Max Pooling

25 x 13 x 13

Convolution

Max Pooling

# CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D tensor)*

input

↓

**Convolution**

```
1 x 28 x 28
```

```
model2.add( Convolution2D( 25,3,3,
            input_shape=(1,28,28) ) )
```

```
25 x 26 x 26
```

↓

**Max Pooling**

```
model2.add(MaxPooling2D((2,2)))
```

```
25 x 13 x 13
```

↓

**Convolution**

```
model2.add(Convolution2D(50,3,3))
```

```
50 x 11 x 11
```

↓

**Max Pooling**

# CNN in Keras

input

1 x 28 x 28

```
model2.add( Convolution2D( 25,3,3,
        input_shape=(1,28,28) ) )
```

Convolution

25 x 26 x 26

```
model2.add(MaxPooling2D((2,2)))
```

Max Pooling

25 x 13 x 13

```
model2.add(Convolution2D(50,3,3))
```

Convolution

50 x 11 x 11

```
model2.add(MaxPooling2D((2,2)))
```

Max Pooling

50 x 5 x 5

# CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D tensor)*

input

$\Downarrow$

1 x 28 x 28

```
model2.add( Convolution2D( 25,3,3,
        input_shape=(1,28,28) ) )
```
$\Leftarrow$ Convolution

How many parameters for each filter?  **9**   25 x 26 x 26

$\Downarrow$

```
model2.add(MaxPooling2D((2,2)))
```
$\Leftarrow$ Max Pooling

25 x 13 x 13

$\Downarrow$

```
model2.add(Convolution2D(50,3,3))
```
$\Leftarrow$ Convolution

How many parameters for each filter?  **225**   50 x 11 x 11

$\Downarrow$

```
model2.add(MaxPooling2D((2,2)))
```
$\Leftarrow$ Max Pooling

50 x 5 x 5

# CNN in Keras

**Only modified the *network structure* and *input format (vector -> 3-D tensor)***

input

1 x 28 x 28

Convolution

25 x 26 x 26

Max Pooling

25 x 13 x 13

Convolution

50 x 11 x 11

Max Pooling

50 x 5 x 5

Flatten 把它拉直

`model2.add(Flatten())`

# CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D tensor)*

output

**Fully Connected Feedforward network**

1250

input

| 1 x 28 x 28 |
| Convolution |
| 25 x 26 x 26 |
| Max Pooling |
| 25 x 13 x 13 |
| Convolution |
| 50 x 11 x 11 |
| Max Pooling |
| 50 x 5 x 5 |

`model2.add(Flatten())`

我們剛才 demo 的是一樣的

Created with EverC

# Deep Dream

• Given a photo, machine adds what it sees ......

然後，你把它的某一個 hidden layer

# Deep Dream

- Given a photo, machine adds what it sees ......



那 Deep Dream 還有一個進階的版本

# Deep Style
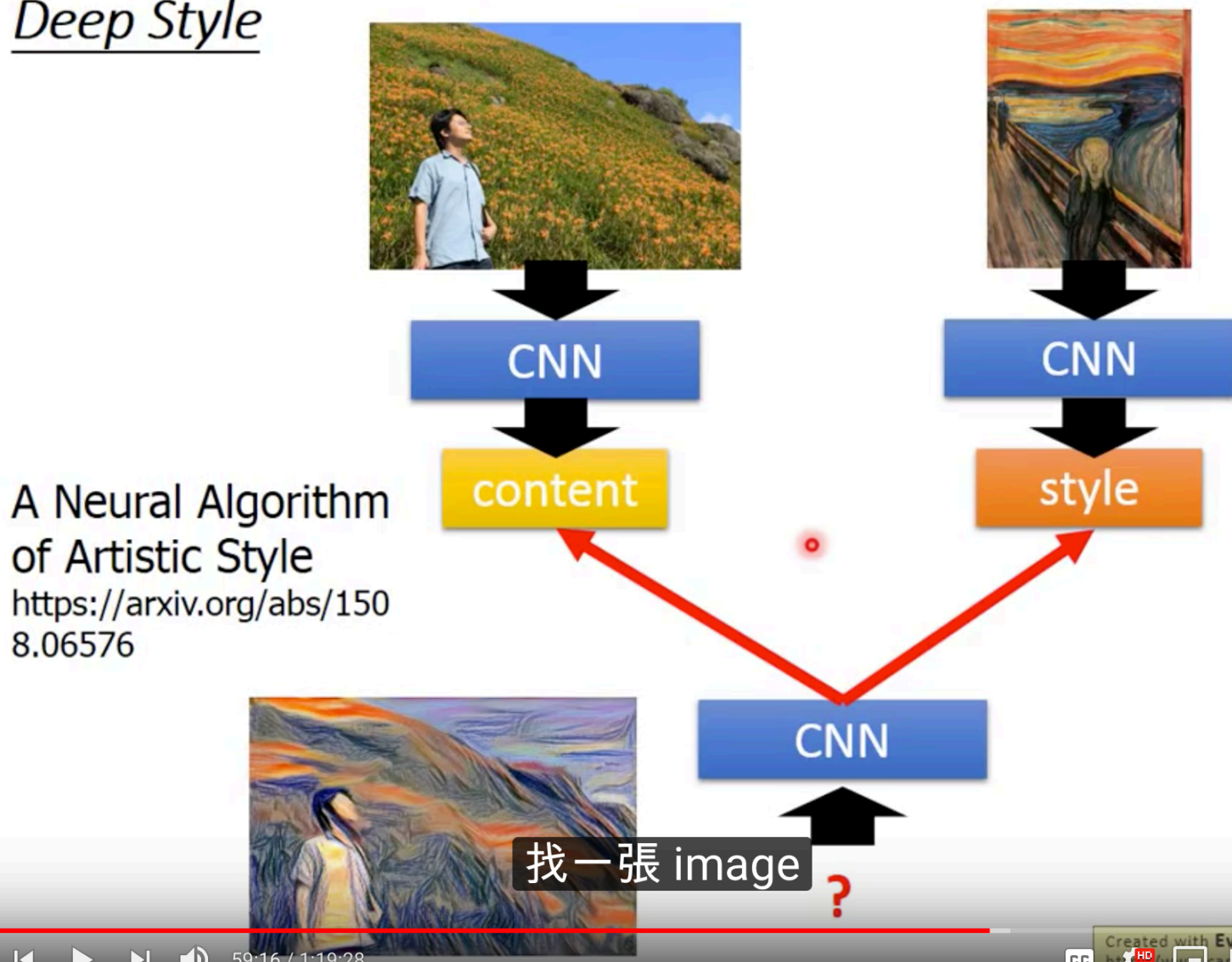
- Given a photo, make its style like famous paintings

# *Deep Style*



A Neural Algorithm
of Artistic Style
https://arxiv.org/abs/150
8.06576

# More Application: Playing Go



19 x 19 vector

Black: 1

white: -1

Network → Next move (19 x 19 positions)

19 x 19 vector

Fully-connected feedforward network can be used

But CNN performs much better.

# Why CNN for playing Go?

- Some patterns are much smaller than the whole image



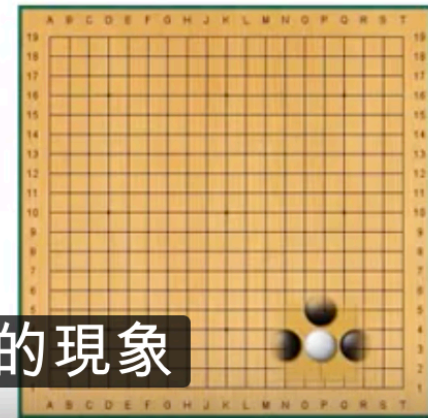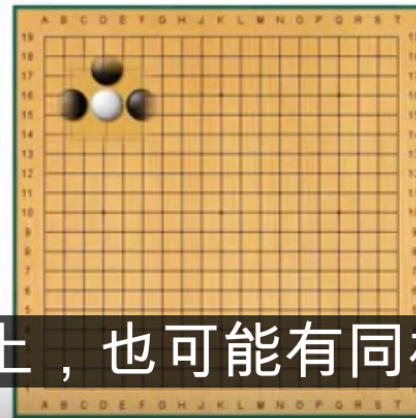- The same patterns appear in different regions.

# Why CNN for playing Go?

- Some patterns are much smaller than the whole image

  Alpha Go uses 5 x 5 for first layer



- The same patterns appear in different regions.



在圍棋上，也可能有同樣的現象

# Why CNN for playing Go?

- Subsampling the pixels will not change the object

**Max Pooling** How to explain this???

# Why CNN for playing Go?

- Subsampling the pixels will not change the object

➡️ **Max Pooling**  <span style="color:red">How to explain this???</span>

**Neural network architecture.** The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a $23 \times 23$ image, then convolves $k$ filters of kernel size $5 \times 5$ with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a $21 \times 21$ image, then convolves $k$ filters of kernel size $3 \times 3$ with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size $1 \times 1$ with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used $k = 192$ filters; Fig. 2b and Extended Data Table 3 additionally show the results of training with $k = 128, 256$ and 384 filters.

叫吃的狀態呢，等等

# Why CNN for playing Go?
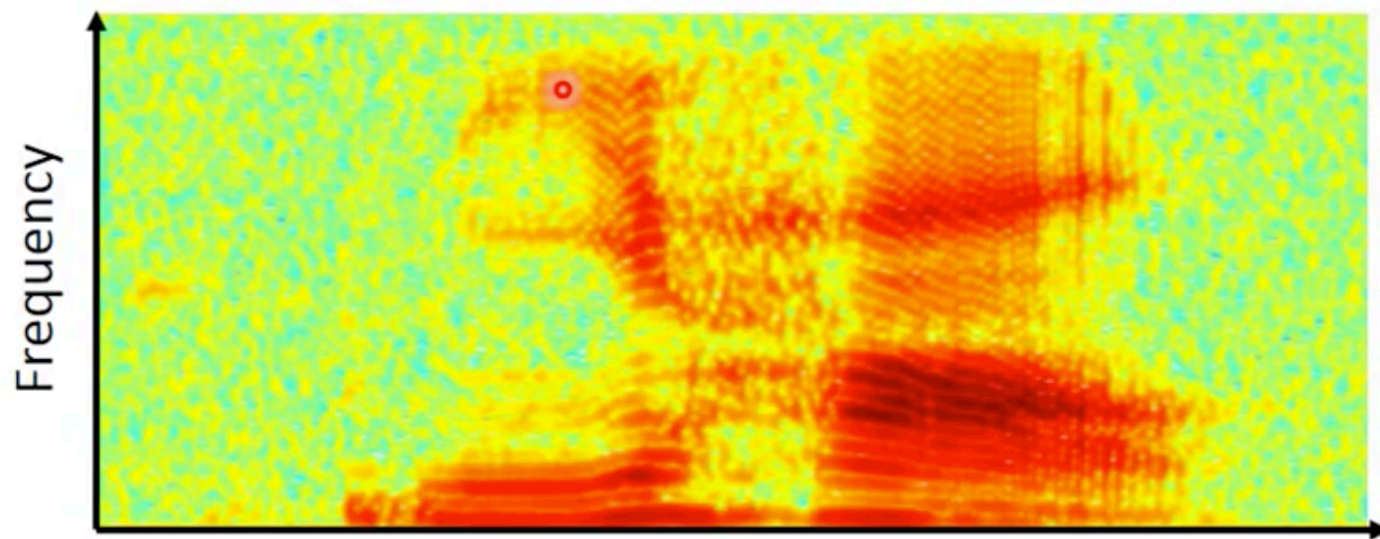
- Subsampling the pixels will not change the object

➡️ **Max Pooling** How to explain this???

**Neural network architecture.** The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a $23 \times 23$ image, then convolves $k$ filters of kernel size $5 \times 5$ with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a $21 \times 21$ image, then convolves $k$ filters of kernel size $3 \times 3$ with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size $1 \times 1$ with stride 1, with a different bias for each position, and applies a softmax function. The [...] Extended Data Table 3 additionally show the results of training with $k = 128$, 256 and 384 filters.

Alpha Go does not use Max Pooling ......
它是沒有用 Max Pooling 的

# More Application: Speech



**Spectrogram**

More Application: Speech

CNN

The filters move in the frequency direction.

Frequency

Spectrogram

它的寬，就跟我們 image 的寬是一樣的

# More Application: Text



sentence matrix $S \in \mathbb{R}^{d \times |s|}$ — convolutional feature map $C \in \mathbb{R}^{n \times |s| - m + 1}$ — pooled representation $c_{\text{pool}} \in \mathbb{R}^{1 \times n}$ — softmax

embedding dimension

$F \in \mathbb{R}^{d \times m}$

I love my new iphone :)

都用一個 vector 來表示

# To learn more ......

- The methods of visualization in these slides
  - https://blog.keras.io/how-convolutional-neural-networks-see-the-world.html
- More about visualization
  - http://cs231n.github.io/understanding-cnn/
- Very cool CNN visualization toolkit
  - http://yosinski.com/deepvis
  - http://scs.ryerson.ca/~aharley/vis/conv/