

Convolutional Neural Network

Hung-yi Lee

Can the network be simplified by
considering the properties of images?

提出背景

图象分类可以用普通的 Fully-Connect 神经网络实现，即 DNN
但是一张图所有 pixel 都有输入方式，导致每个 pixel 都会有一个 w 被训练太耗效

解决了什么？

CNN 可以简化这件事情，类似于一种预处理，大大减少了 w 个数
所以 CNN 模型比 DNN 更简单

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

一些 pattern - 一般在图中只占很小一部分，所以没必要把整张图 pixel input 到一个神经元



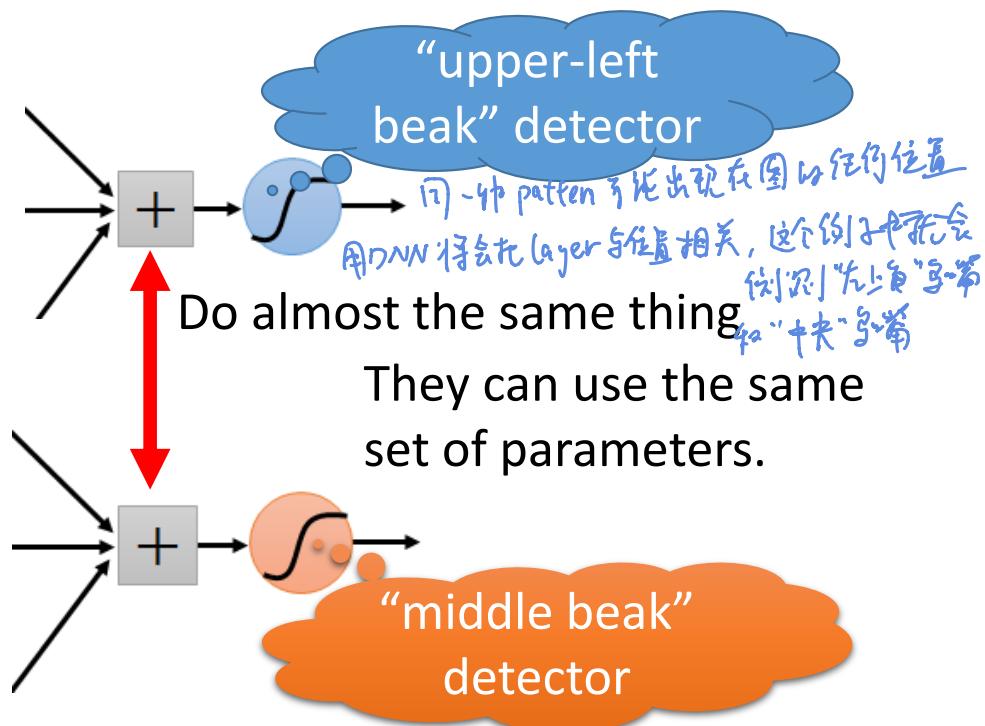
“beak” detector
鸟“嘴”

↓ 但若用 DNN, 全图所有 pixel 都需 input 到这个 layer

(一个 neuron 检测一个 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



bird



subsampling

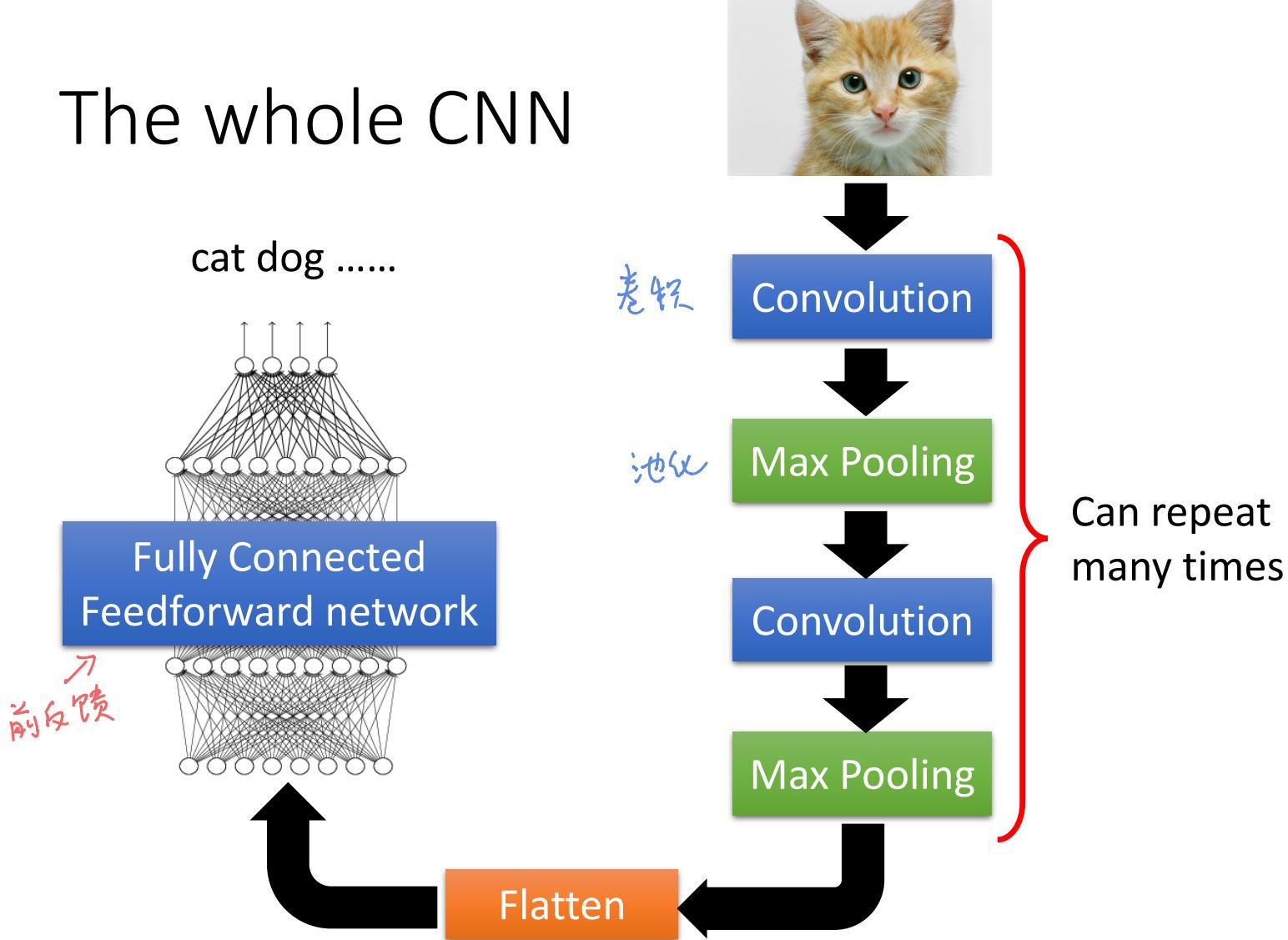
采样, resize圖, 減少 input

We can subsample the pixels to make image smaller



Less parameters for the network to process the image

The whole CNN



The whole CNN

Property 1

- Some patterns are much smaller than the whole image
取局部特征

Property 2

- The same patterns appear in different regions.
pattern 在不同区域

Property 3

- Subsampling the pixels will not change the object



Convolution



Max Pooling



Convolution



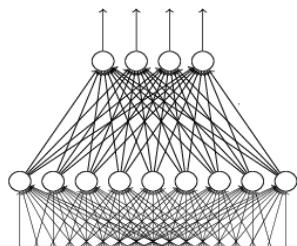
Max Pooling



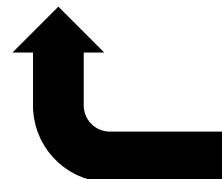
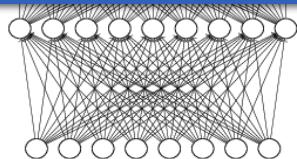
Can repeat many times

The whole CNN

cat dog



Fully Connected
Feedforward network



Flatten



Convolution

Max Pooling

Convolution

Max Pooling

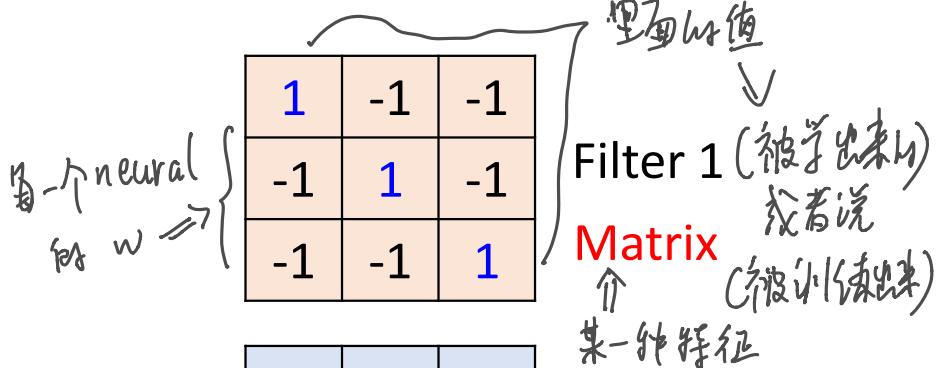
Can repeat
many times

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

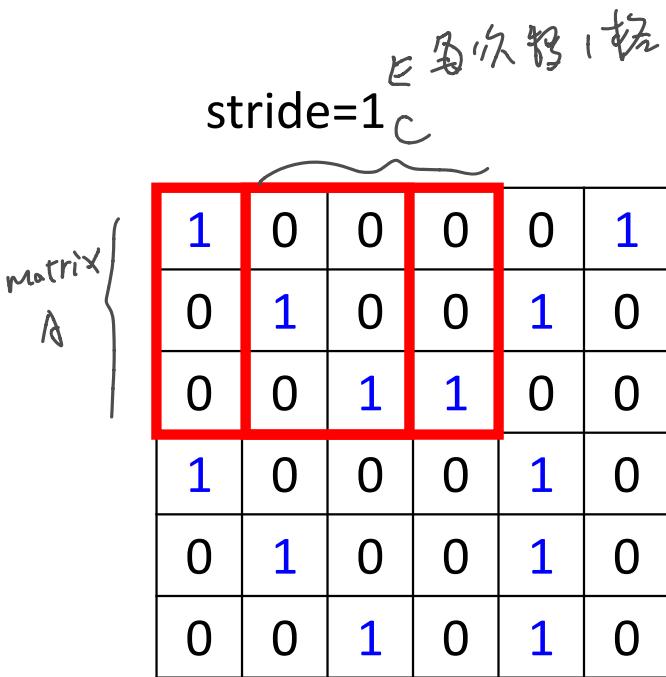
Those are the network parameters to be learned.



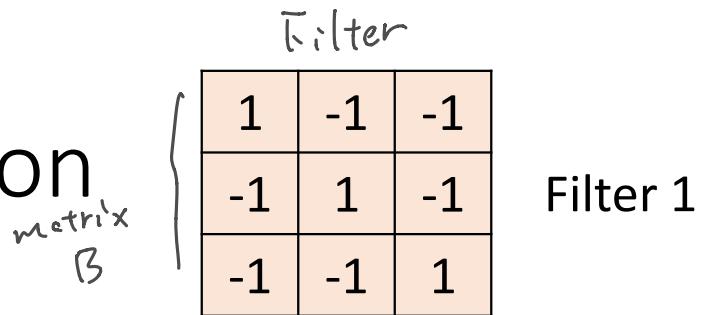
Property 1

Each filter detects a small pattern (3×3).

CNN – Convolution

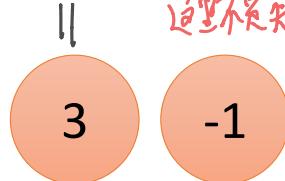


6 x 6 image



A · B 内积 (对应位置元素相乘, 加和 =)

逐项矩阵相乘



C · B 内积

CNN – Convolution

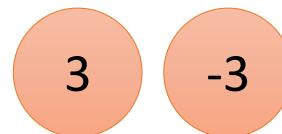
If stride=2 \leftarrow 每次移 2 格

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

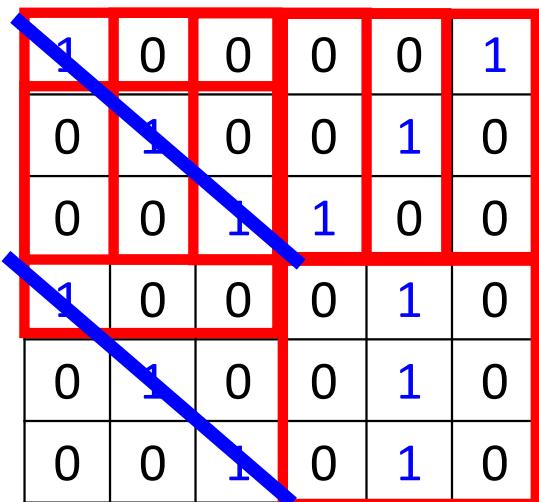
Filter 1



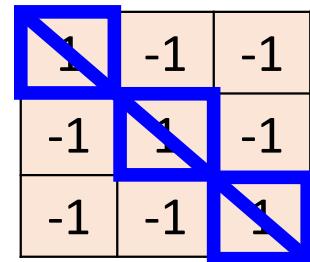
We set stride=1 below

CNN – Convolution

stride=1

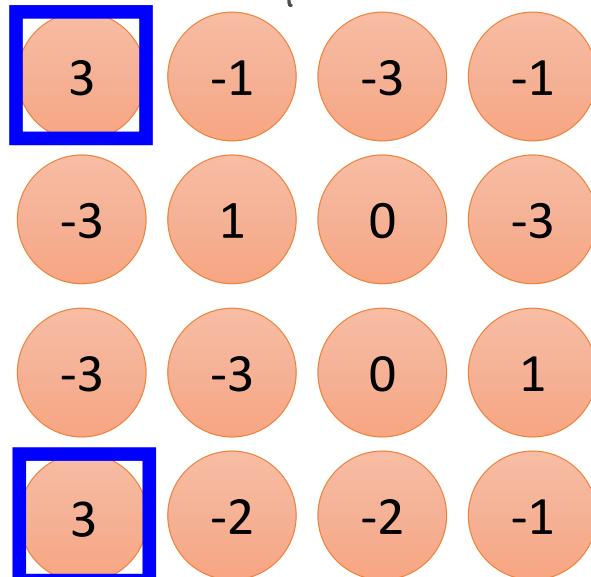


6 x 6 image



Filter 1

值越大，与 Filter 匹配度越高
(patch)



Property 2

CNN – Convolution

stride=1

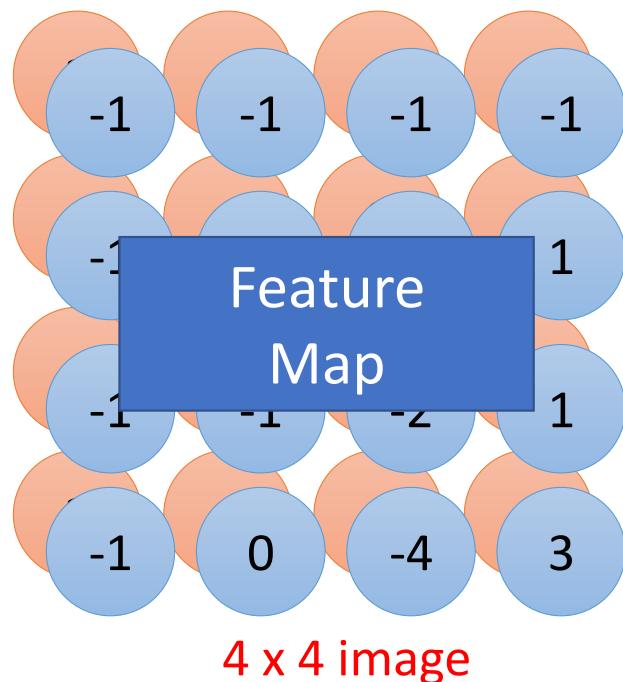
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

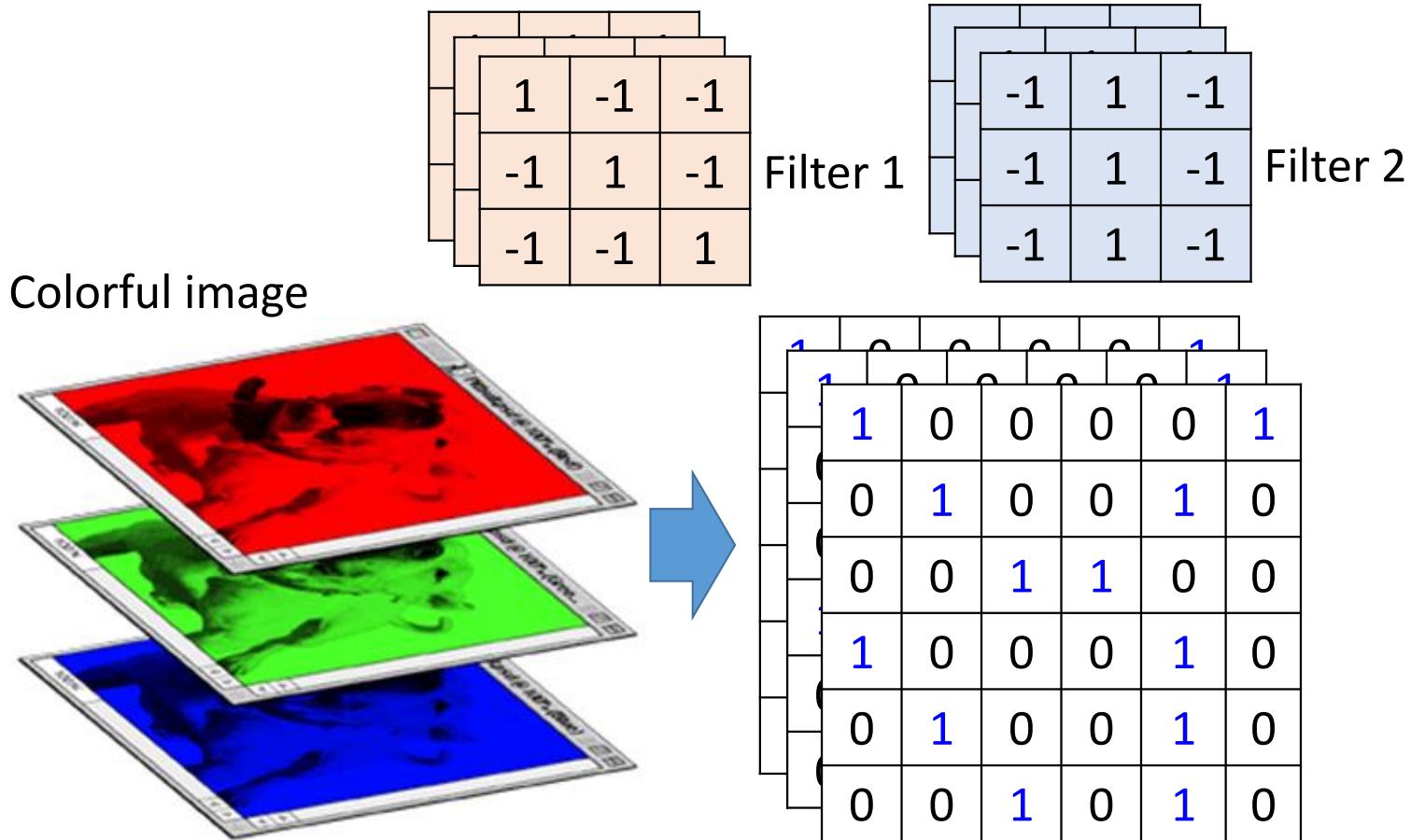
-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

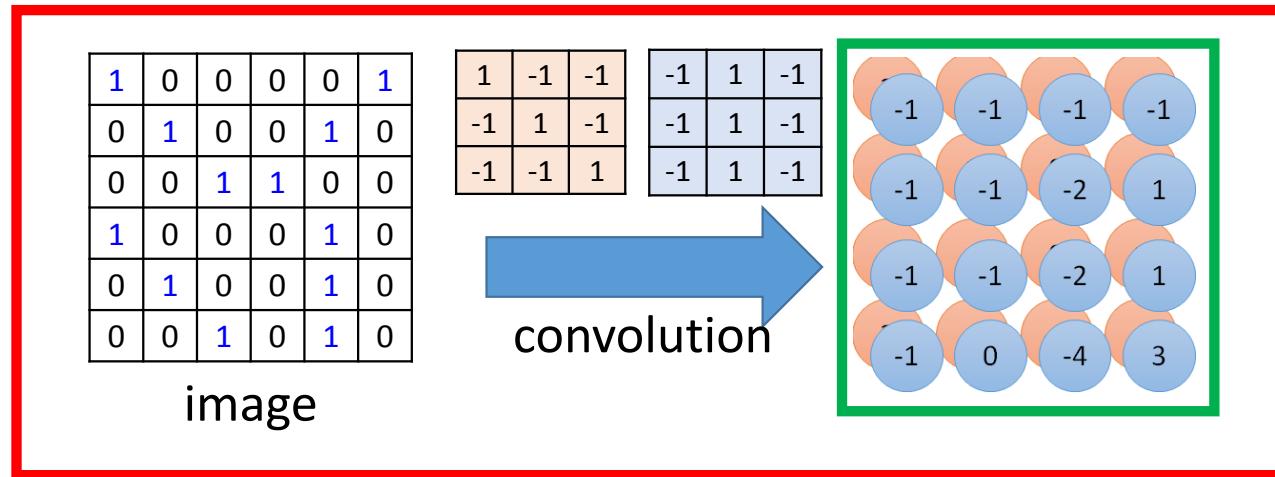
Do the same process for every filter



CNN – Colorful image

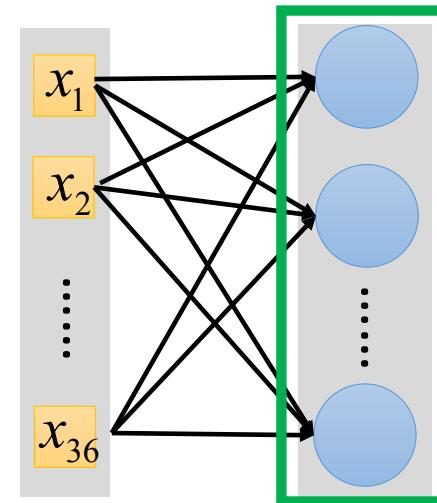


Convolution v.s. Fully Connected



Fully-
connected

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



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

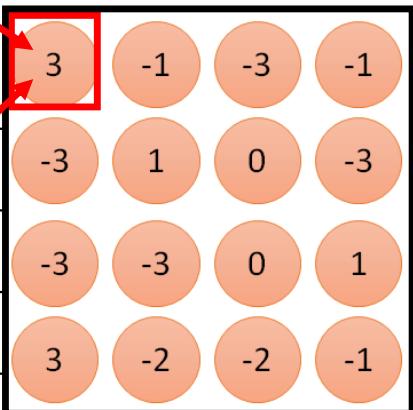
Filter 1

原数据太多，若直接对所有 pixel
运行神经网络 Fully Connected . 小能
太低 \Rightarrow 使用卷积 convolution 过滤。
不好的 pixel 被过滤了，

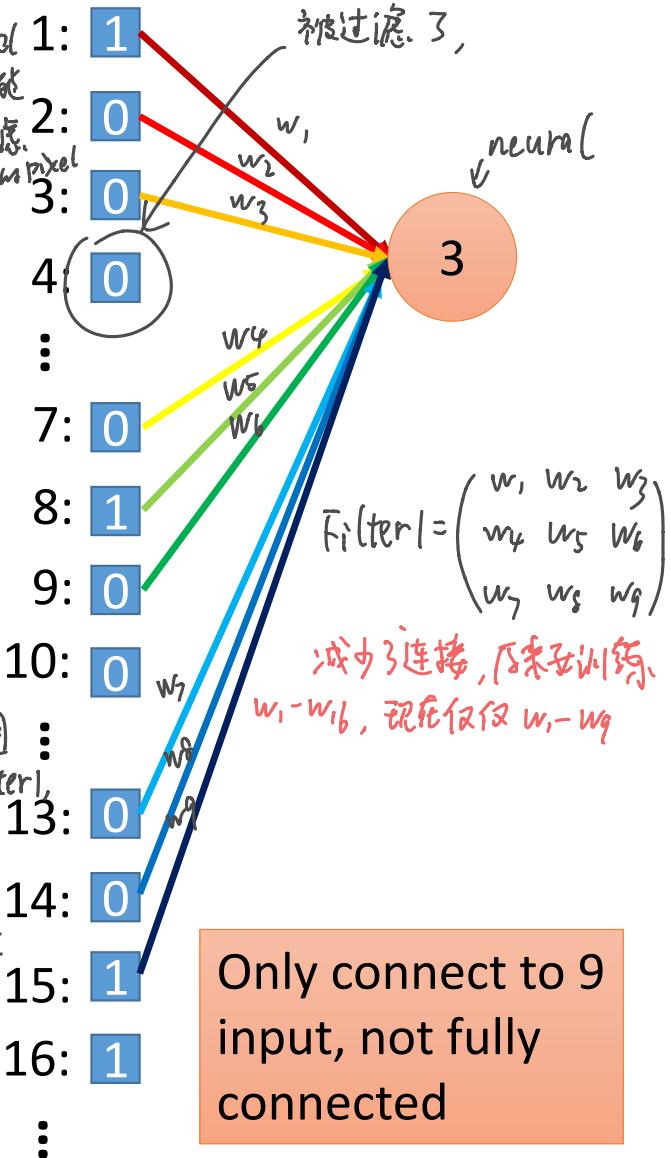
1	0	0	0	0	0
0	1	0	0	0	1
0	0	1	1	0	0
1	0	0	0	0	1
0	1	0	0	0	1
0	0	1	0	1	0

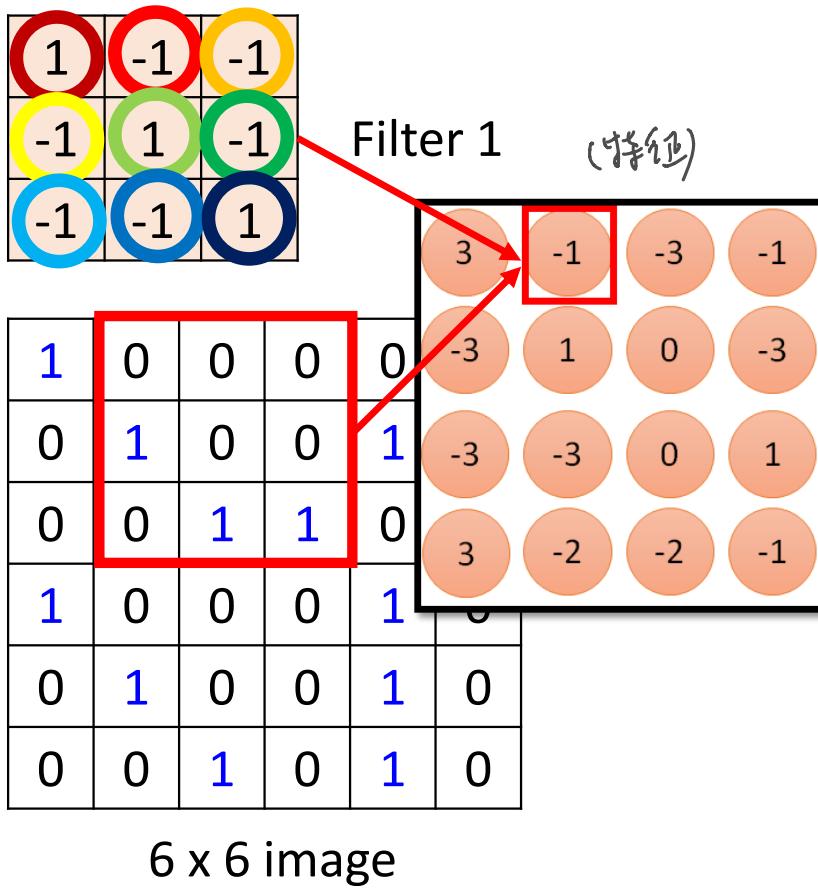
6 x 6 image

Less parameters!



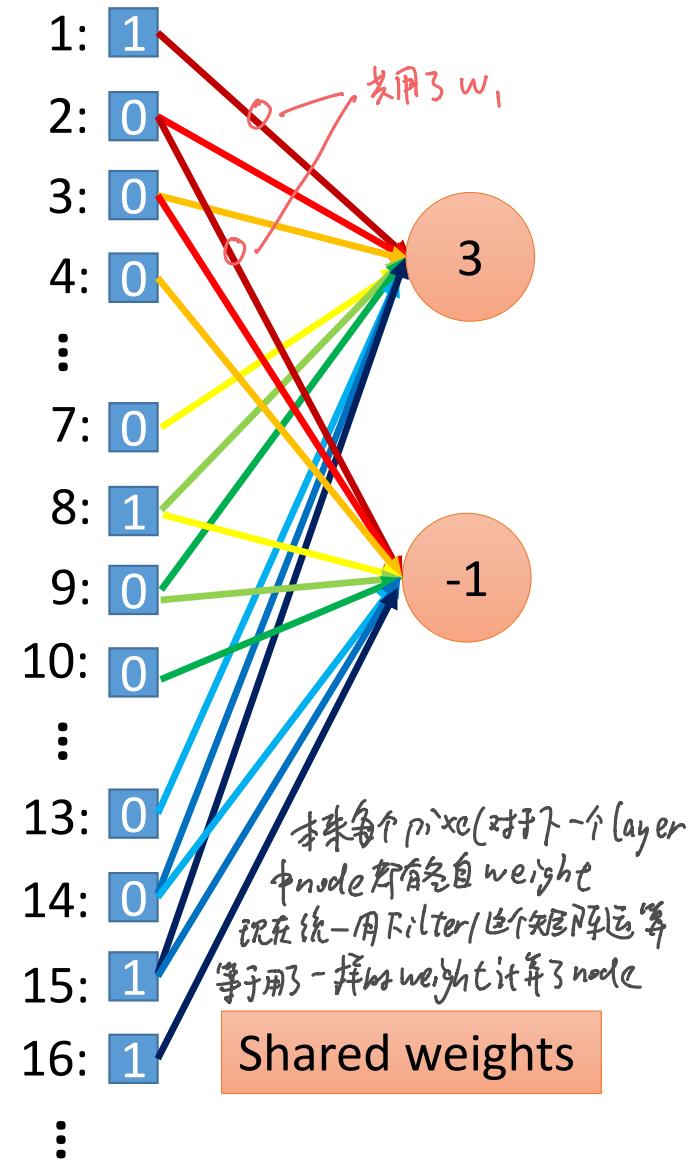
Filter 1 与左上角第一个 3x3 内核时，等价于下图：
中由 9 个像素 input 进 Filter 1，
变向实现了过滤多余元素，
并没有一次把整张图 input
进 Filter 1.





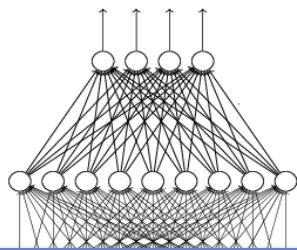
Less parameters!

Even less parameters!

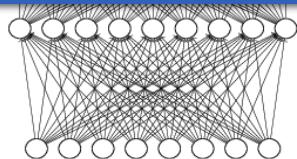


The whole CNN

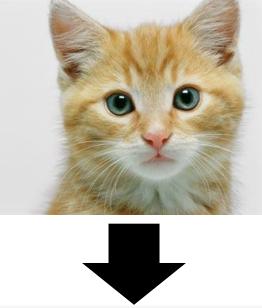
cat dog



Fully Connected
Feedforward network



Flatten



Max Pooling

Convolution

Max Pooling

Can repeat
many times

(进一步提取特征)

CNN – Max Pooling 池化

进一步缩小图片

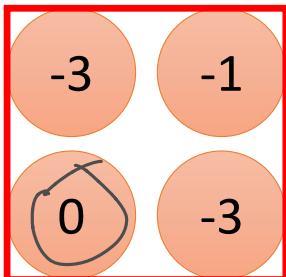
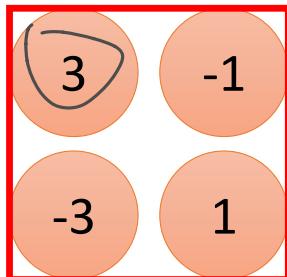
1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

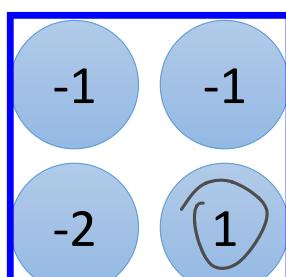
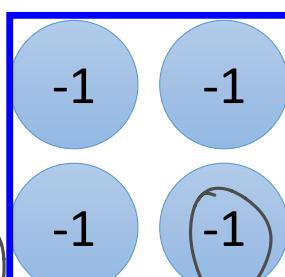
-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

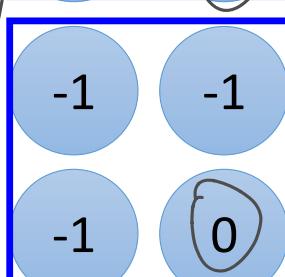
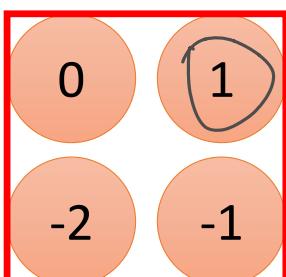
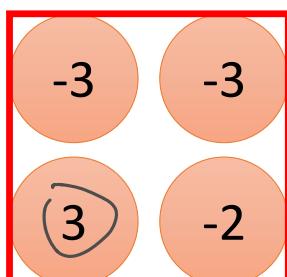
方式①，取组内最大



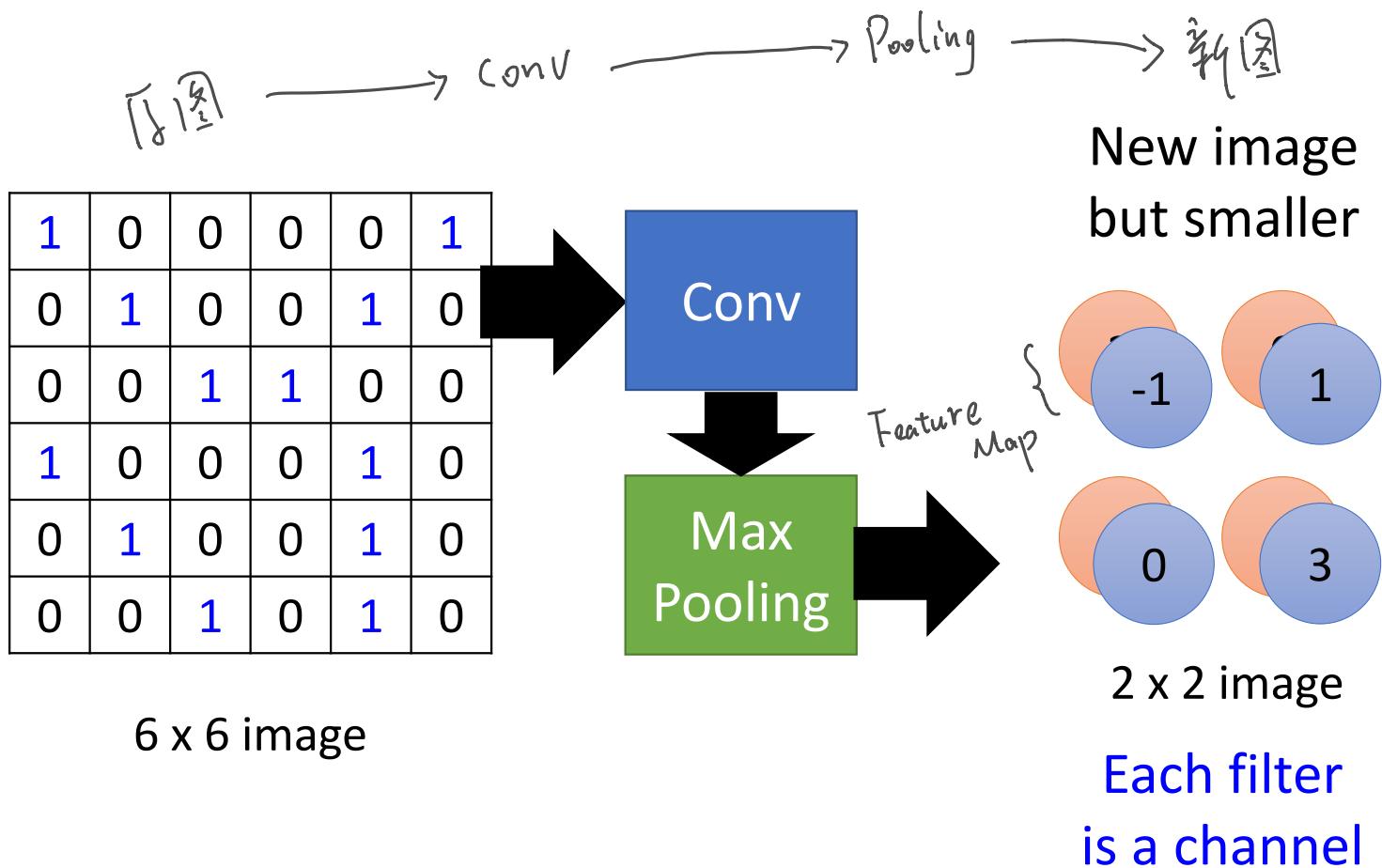
$$\Rightarrow \begin{pmatrix} 3 & 0 \\ 3 & 1 \end{pmatrix}$$



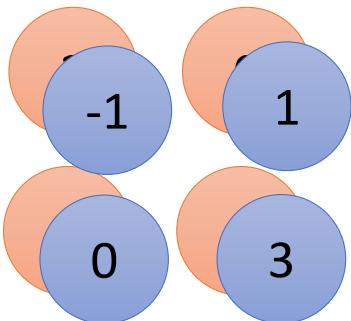
$$\Rightarrow \begin{pmatrix} -1 & 1 \\ 0 & 3 \end{pmatrix}$$



CNN – Max Pooling



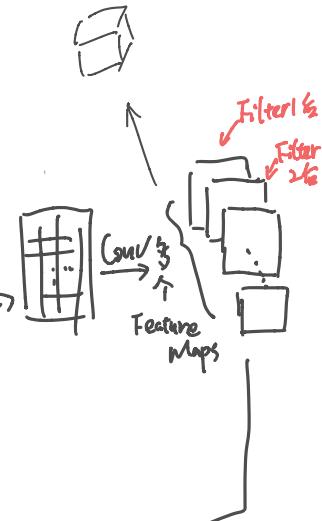
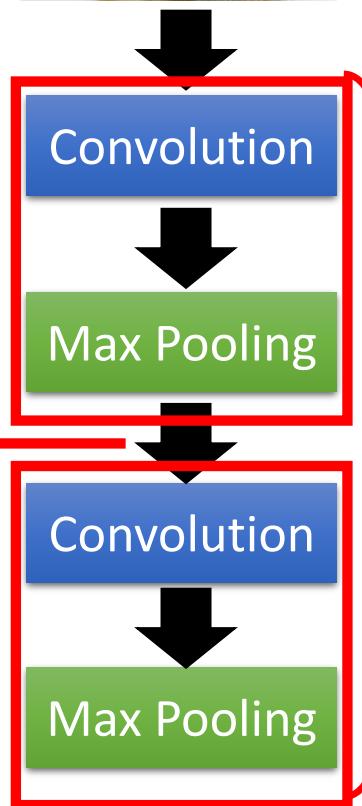
The whole CNN



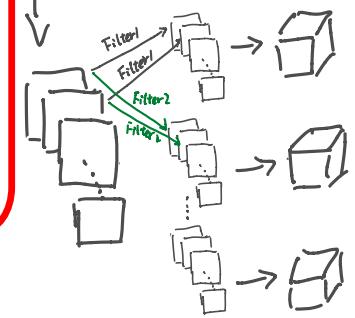
A new image

Smaller than the original image

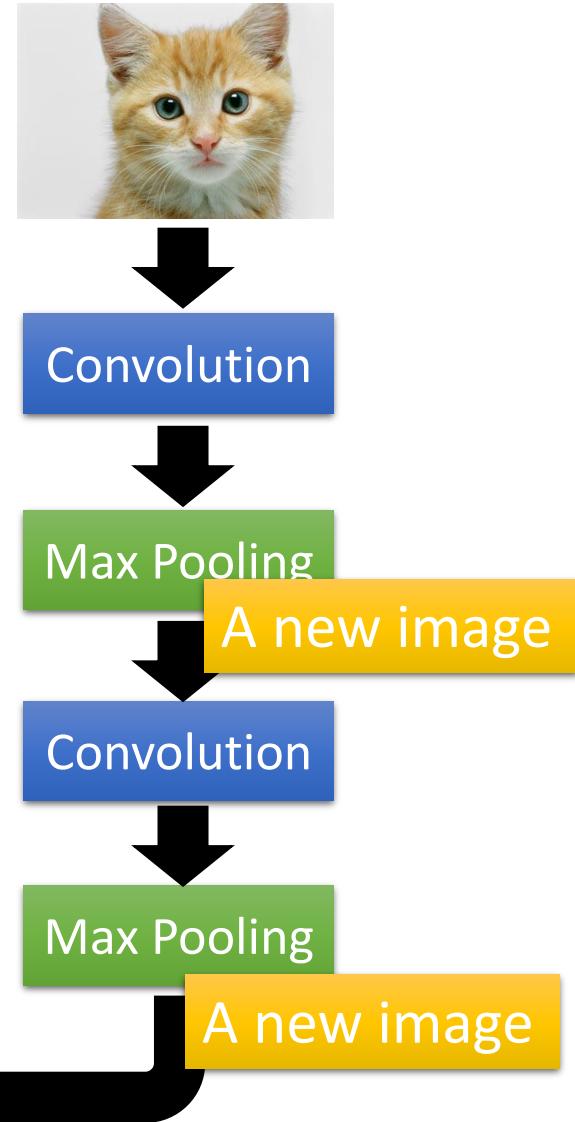
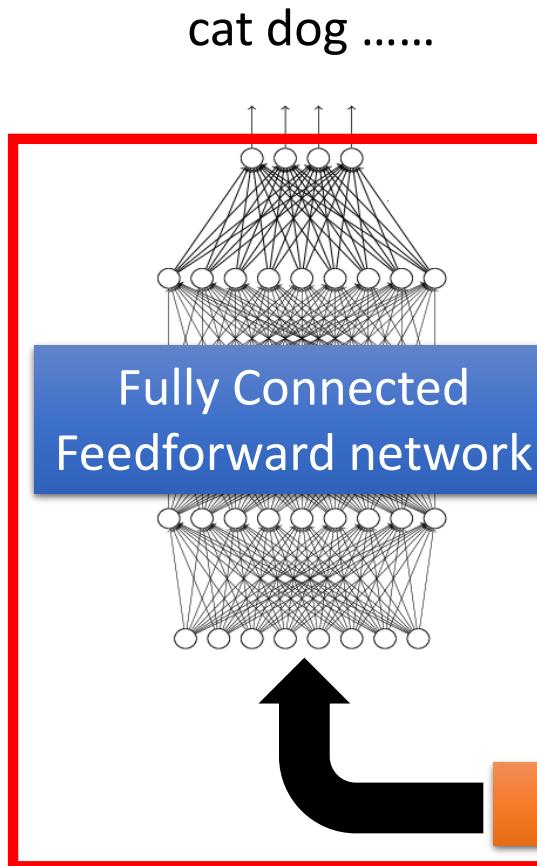
The number of the channel is the number of filters



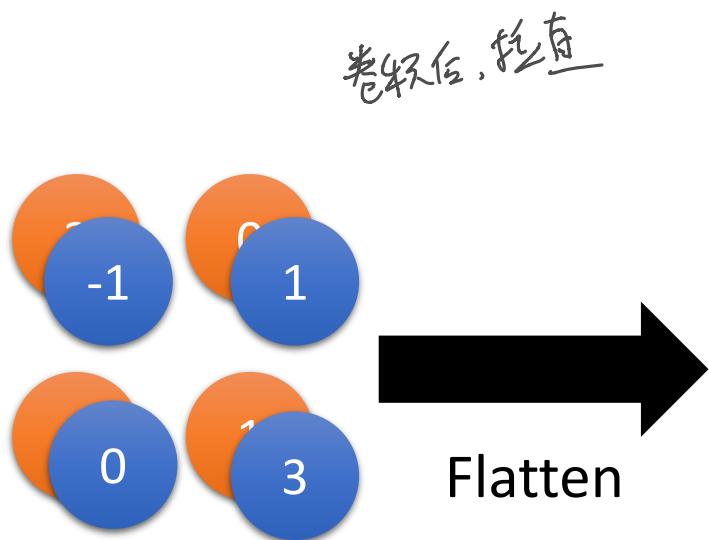
Can repeat many times
Conv/ $\xrightarrow{\text{Max Pool}}$ Conv



The whole CNN



Flatten

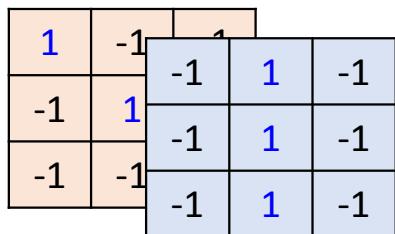


Fully Connected
Feedforward network

CNN in Keras

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

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

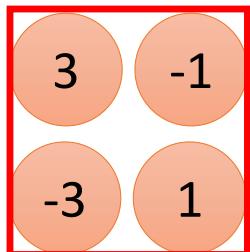


.....
There are 25
3x3 filters.

Input_shape = (28 , 28 , 1)

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

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



input
↓

Convolution



Max Pooling



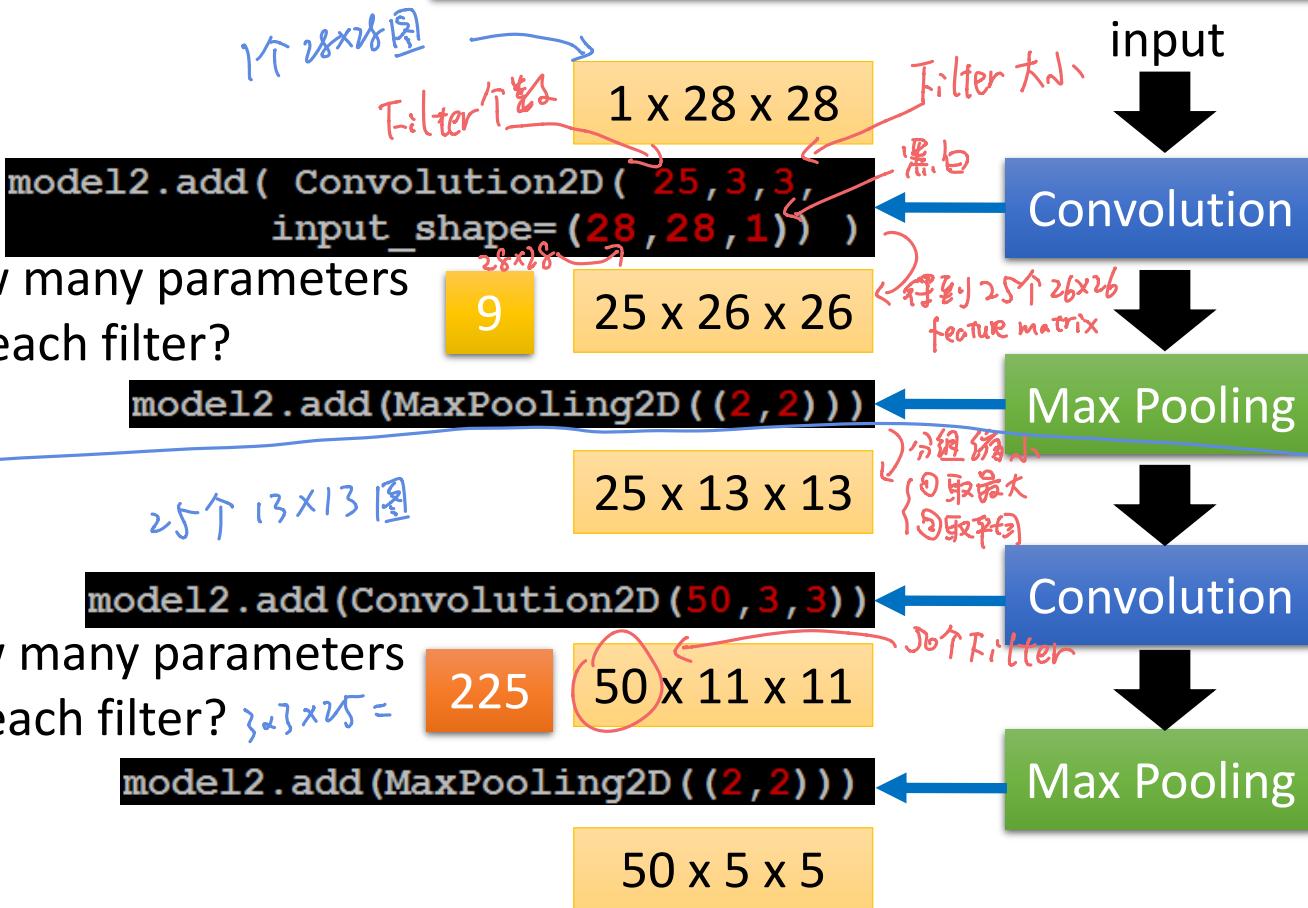
Convolution



Max Pooling

CNN in Keras

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



How many parameters for each filter?

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

25 个 13x13 圖

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

How many parameters for each filter? $3 \times 3 \times 25 =$

225

50 x 11 x 11

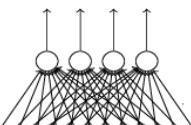
`model2.add(MaxPooling2D((2, 2)))`

50 x 5 x 5

CNN in Keras

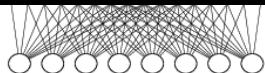
Only modified the *network structure* and *input format (vector -> 3-D tensor)*

output



Fully Connected
Feedforward network

```
model2.add(Dense(output_dim=100))  
model2.add(Activation('relu'))  
model2.add(Dense(output_dim=10))  
model2.add(Activation('softmax'))
```



1250

Flatten

```
model2.add(Flatten())
```

input

$1 \times 28 \times 28$

Convolution

$25 \times 26 \times 26$

Max Pooling

$25 \times 13 \times 13$

Convolution

$50 \times 11 \times 11$

Max Pooling

$50 \times 5 \times 5$

Live Demo

What does machine learn?



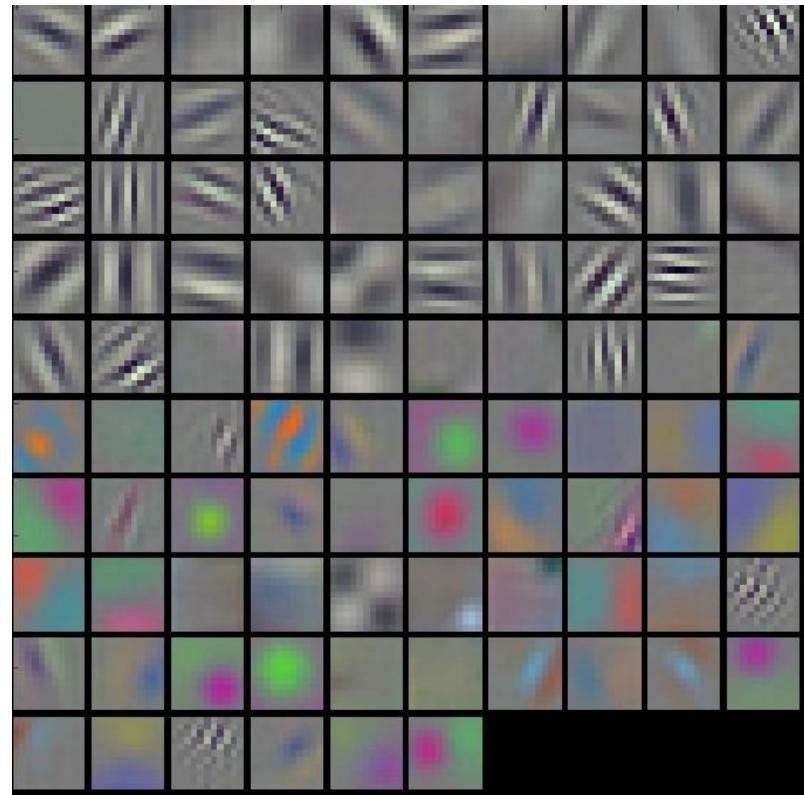
<http://newsneakernews.wpengine.netdna-cdn.com/wp-content/uploads/2016/11/rihanna-puma-creeper-velvet-release-date-02.jpg>

First Convolution Layer

- Typical-looking filters on the trained first layer

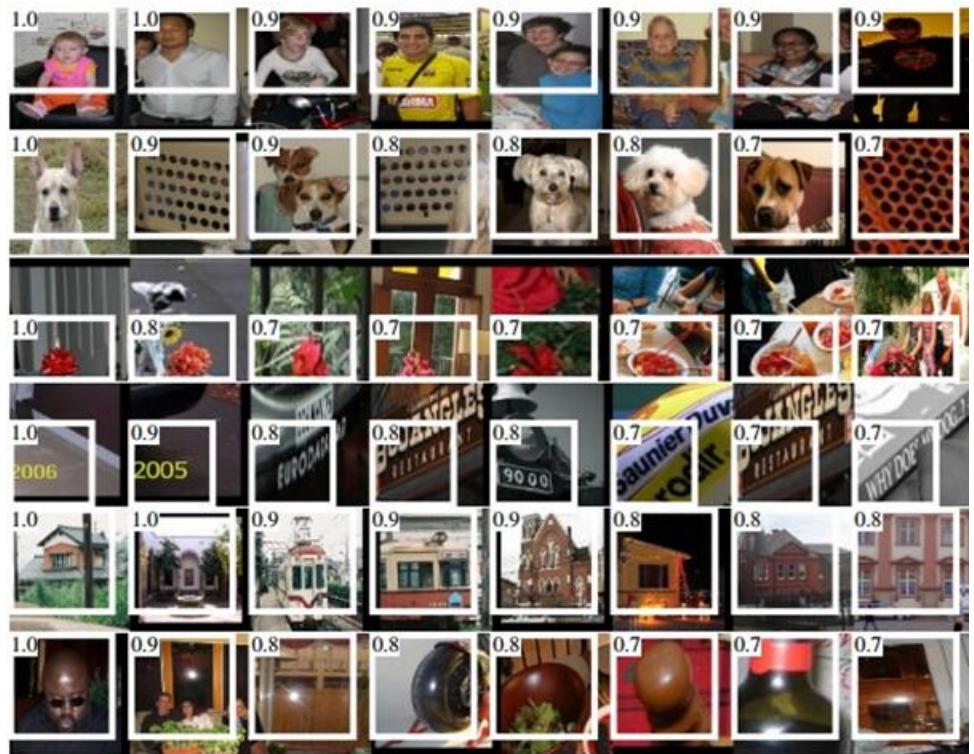
使用训练好的 Filter 反向
可视化，每个 Filter 最想看到什么

11 x 11
(AlexNet)



How about higher layers?

- Which images make a specific neuron activate



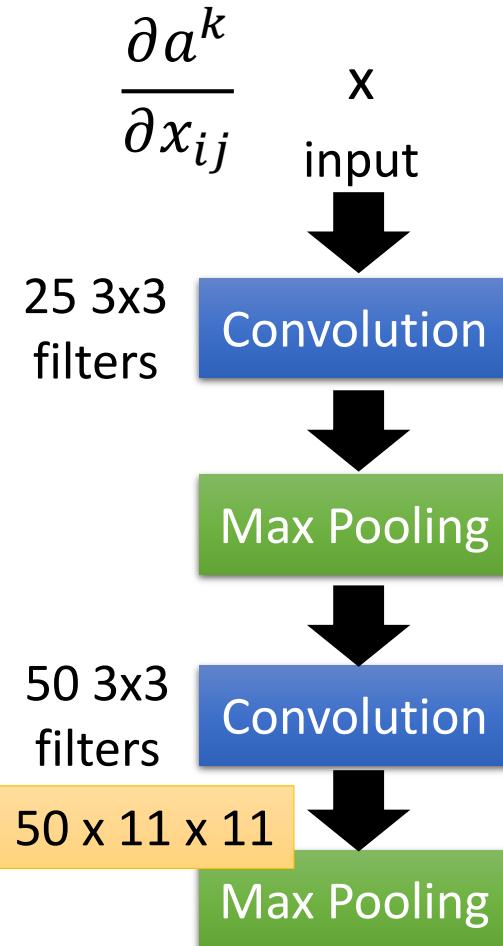
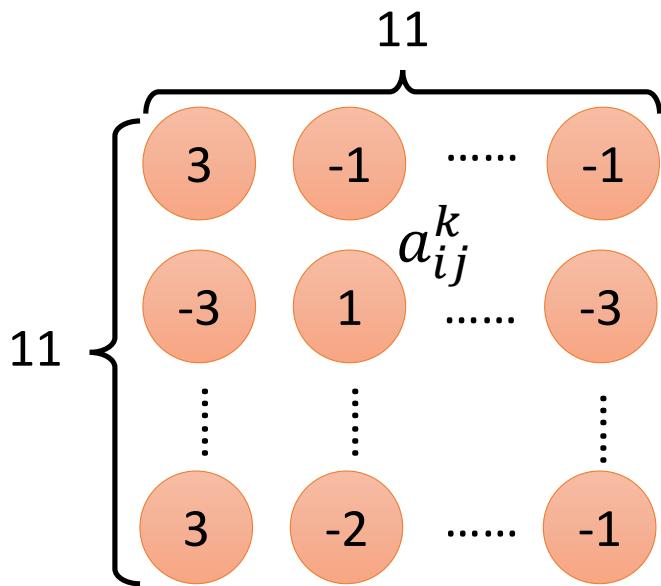
Ross Girshick, Jeff
Donahue, Trevor
Darrell, Jitendra Malik, "Rich
feature hierarchies for accurate
object detection and semantic
segmentation", CVPR, 2014

What does CNN learn?

The output of the k-th filter is a 11×11 matrix.

Degree of the activation of the k-th filter: $a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$

$$x^* = \arg \max_x a^k \text{ (gradient ascent)}$$

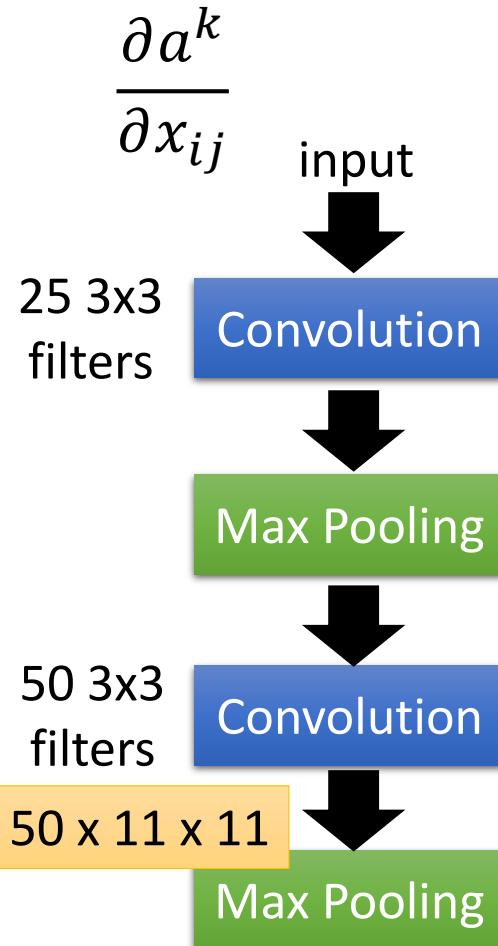
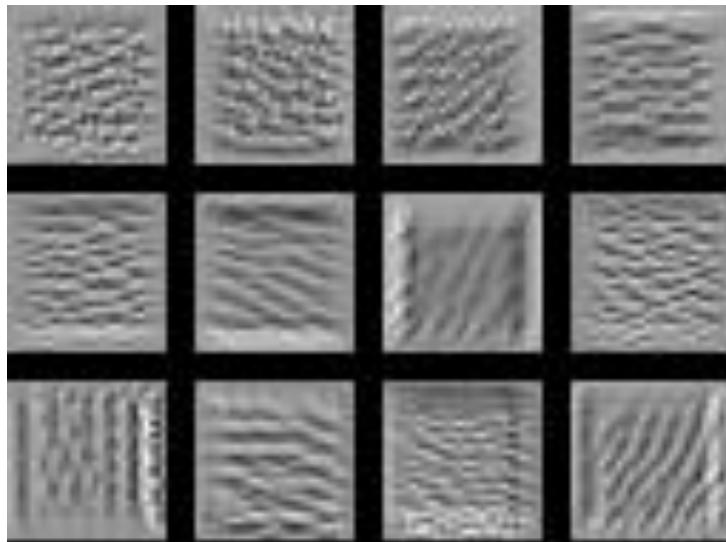


What does CNN learn?

The output of the k-th filter is a 11 x 11 matrix.

Degree of the activation of the k-th filter: $a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$

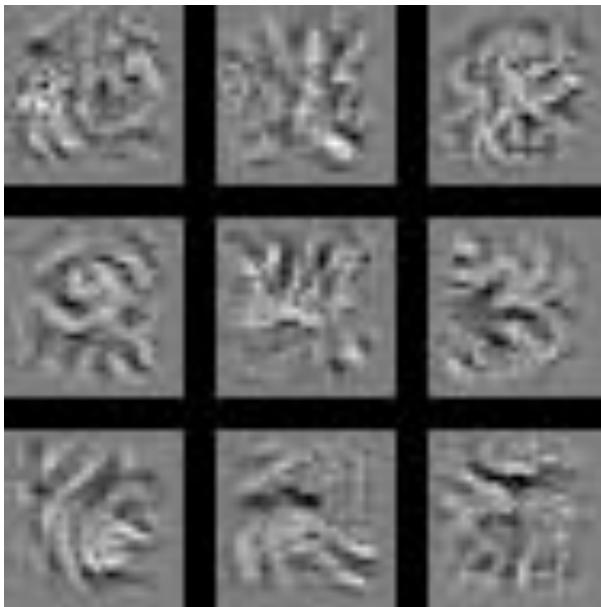
$$x^* = \arg \max_x a^k \text{ (gradient ascent)}$$



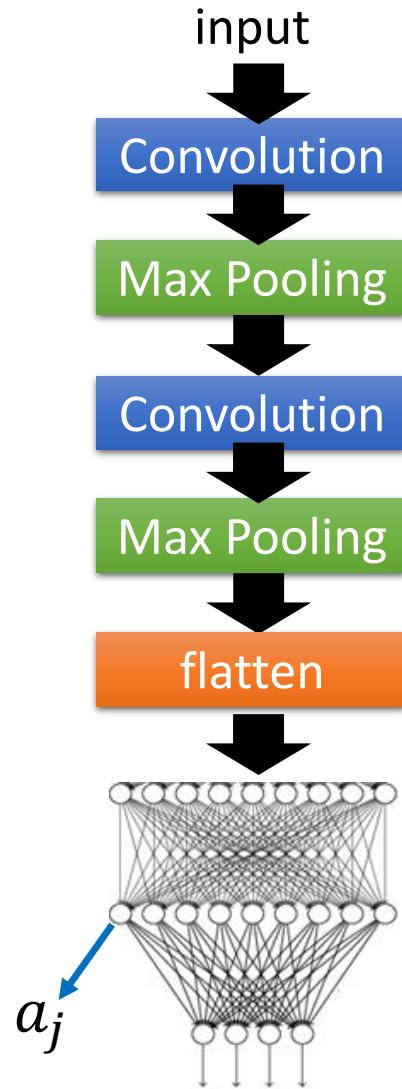
What does CNN learn?

Find an image maximizing the output of neuron:

$$x^* = \arg \max_x a_j$$



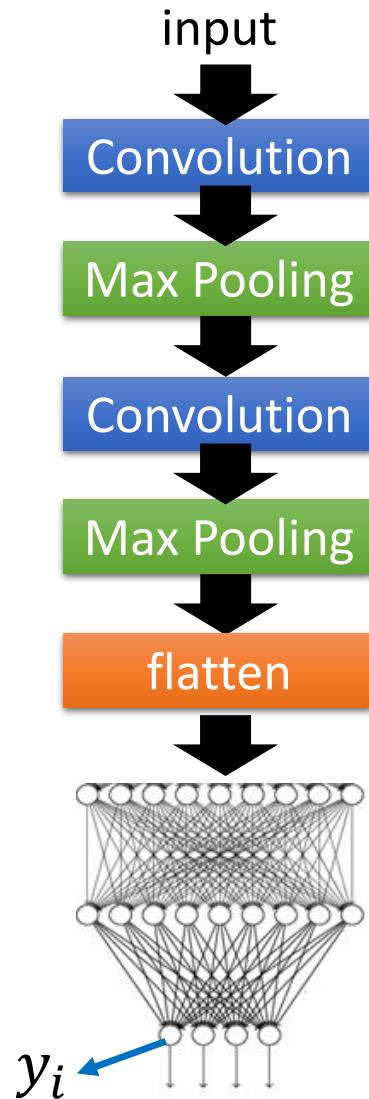
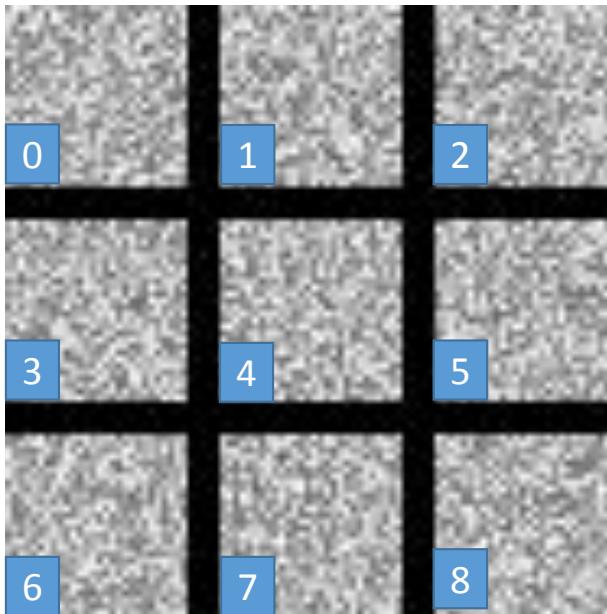
Each figure corresponds to a neuron



What does CNN learn?

$$x^* = \arg \max_x y^i$$

Can we see digits?

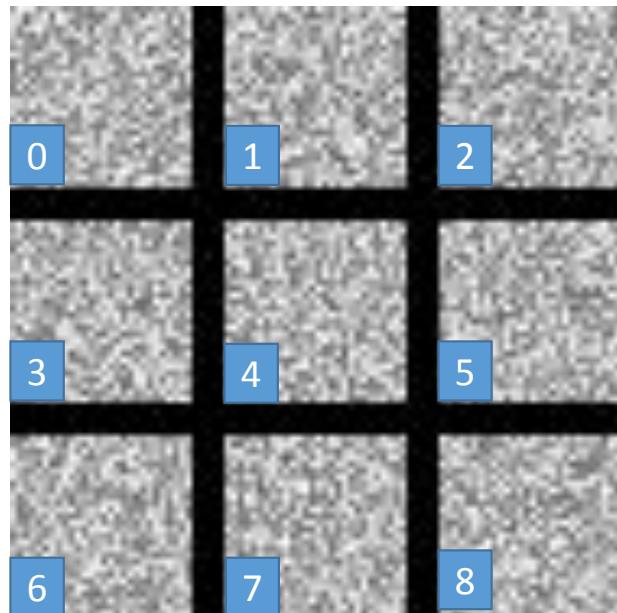


Deep Neural Networks are Easily Fooled

<https://www.youtube.com/watch?v=M2IebCN9Ht4>

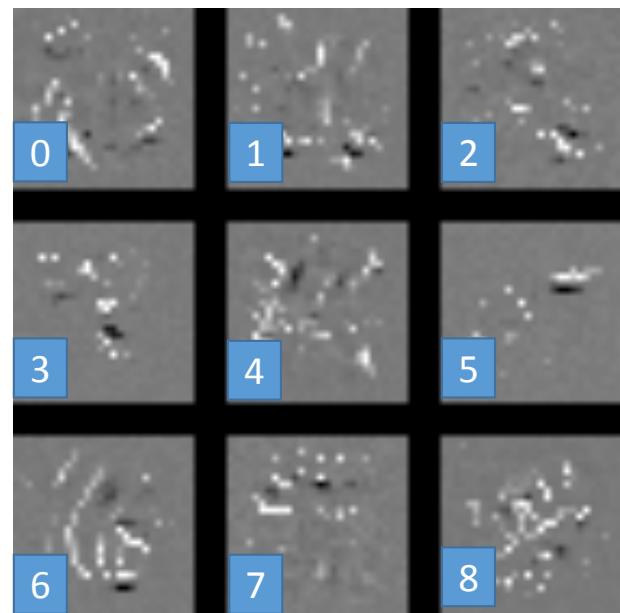
What does CNN learn?

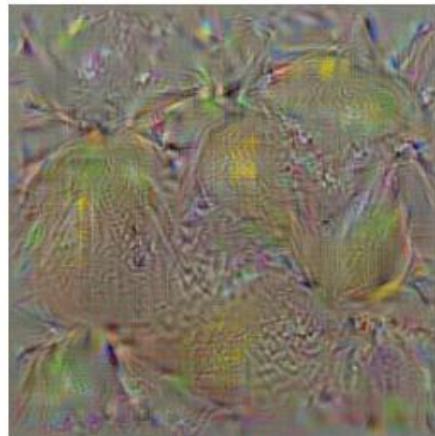
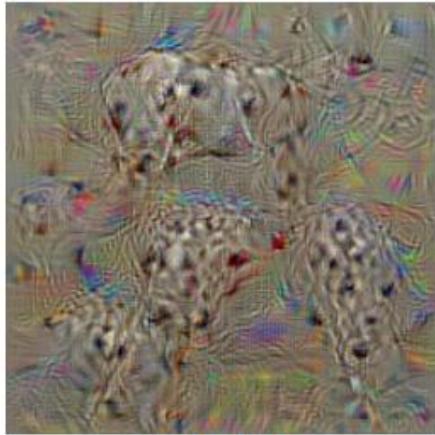
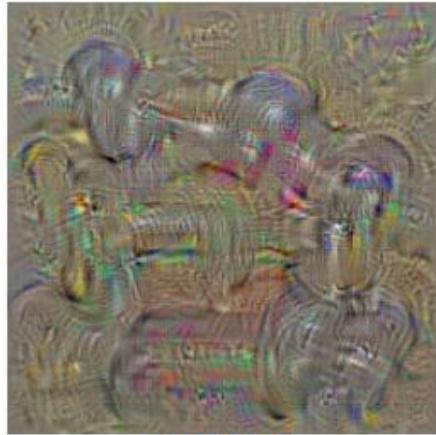
$$x^* = \arg \max_x y^i$$



Over all
pixel values

$$x^* = \arg \max_x \left(y^i - \sum_{i,j} |x_{ij}| \right)$$

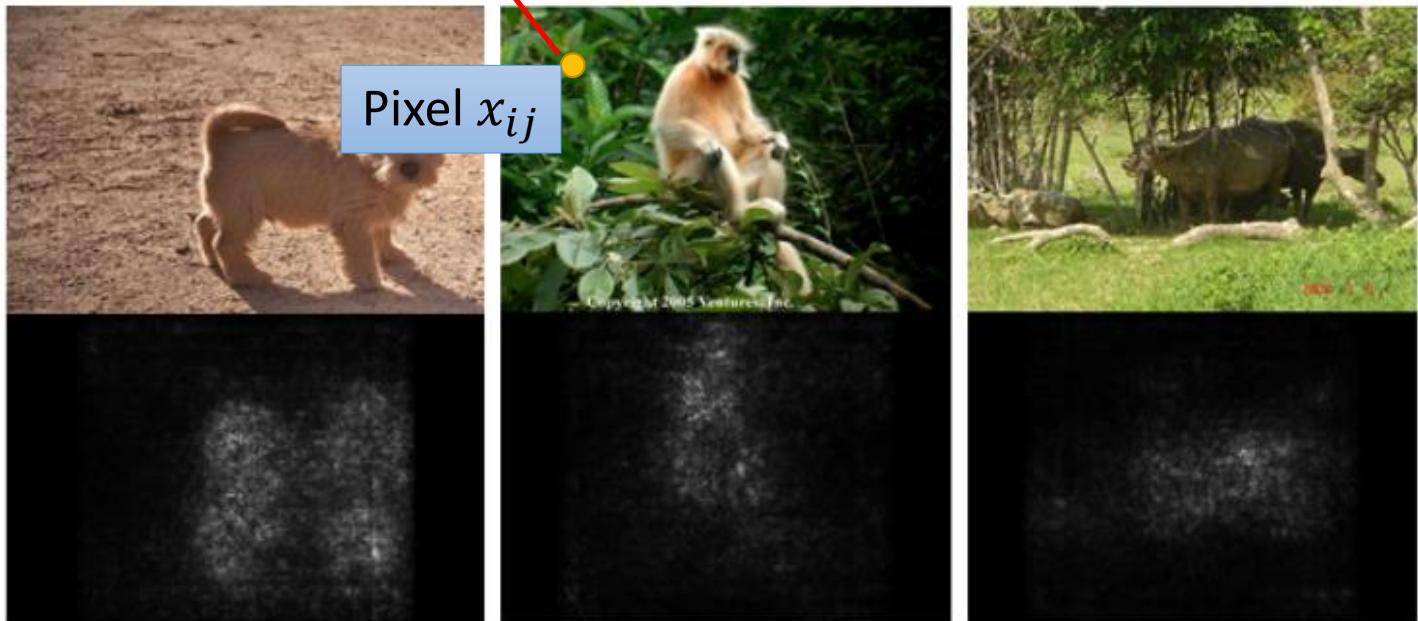




Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR, 2014

$$\left| \frac{\partial y_k}{\partial x_{ij}} \right|$$

y_k : the predicted class of the model



Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR, 2014



True Label: Pomeranian



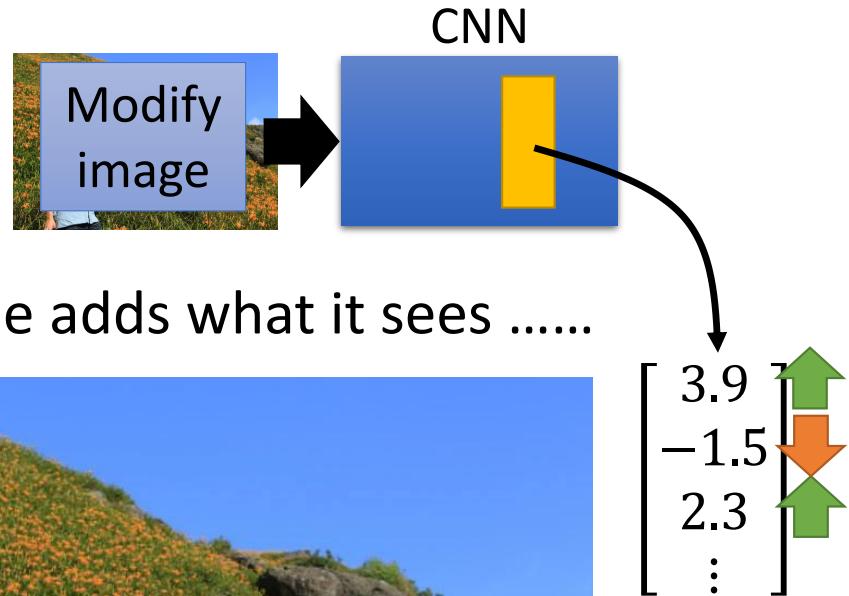
True Label: Car Wheel



True Label: Afghan Hound

Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014* (pp. 818-833)

Deep Dream



- Given a photo, machine adds what it sees



Deep Dream

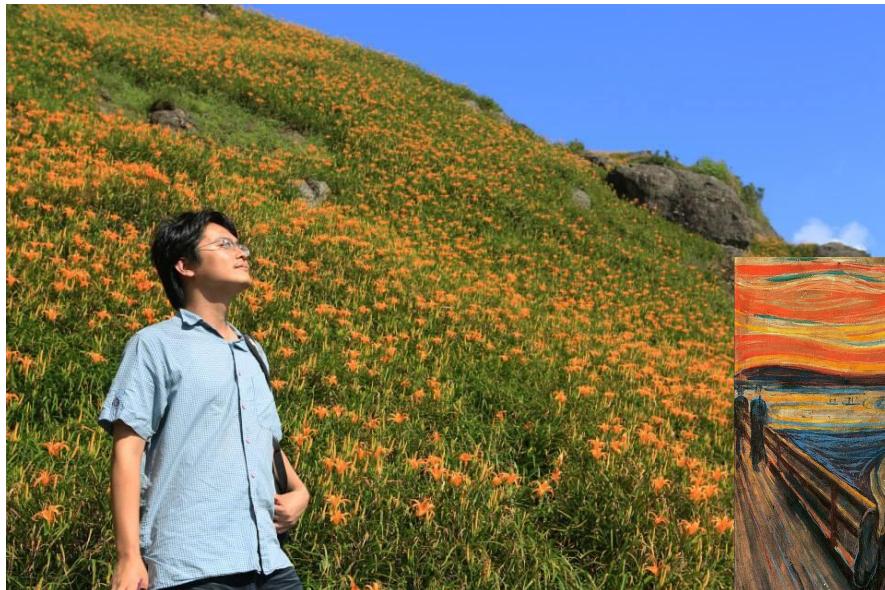
- Given a photo, machine adds what it sees



<http://deepdreamgenerator.com/>

Deep Style

- Given a photo, make its style like famous paintings



<https://dreamscopeapp.com/>

Deep Style

- Given a photo, make its style like famous paintings



<https://dreamscopeapp.com/>

Deep Style

A Neural
Algorithm of
Artistic Style

<https://arxiv.org/abs/1508.06576>



CNN

content



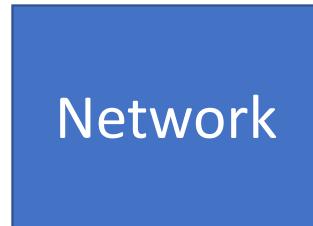
CNN

style

CNN

?

More Application: Playing Go



Next move
(19×19
positions)

19 x 19 matrix
(image)

Black: 1
white: -1
none: 0

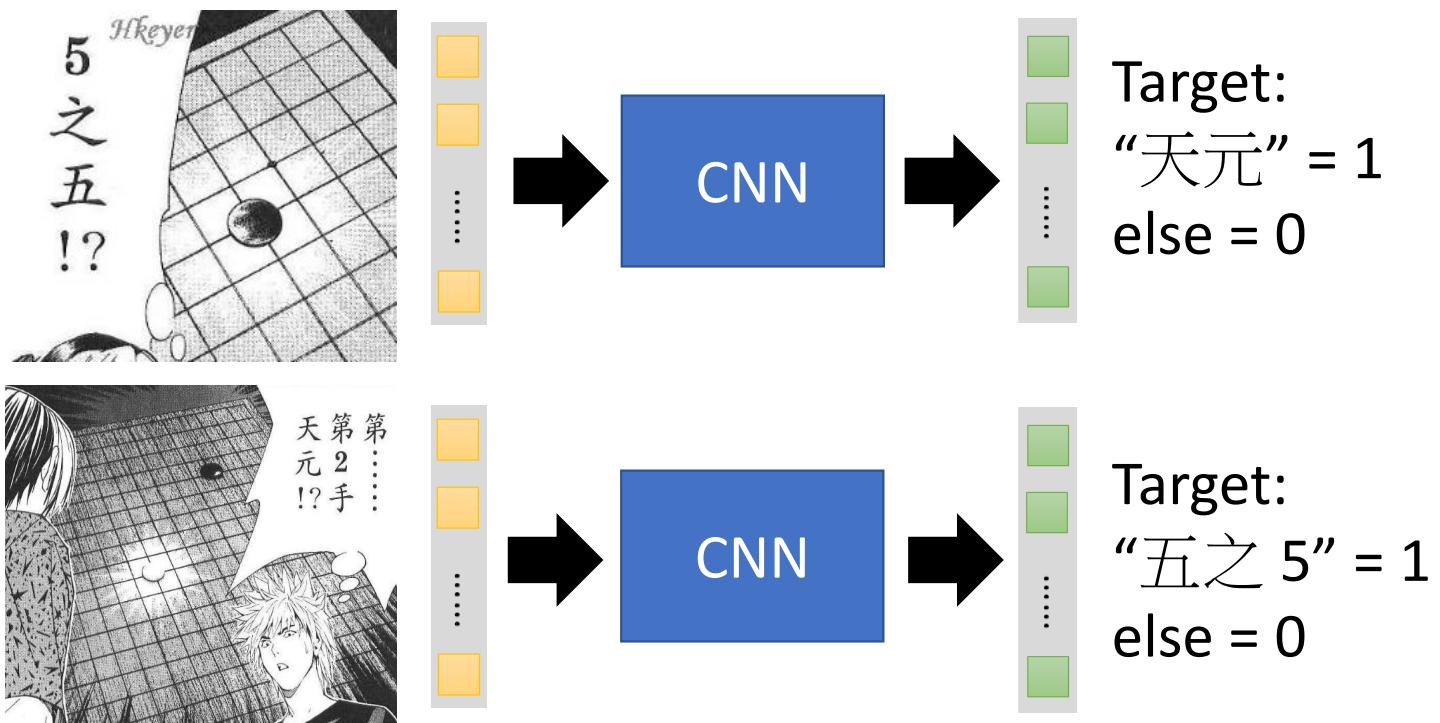
19 x 19 vector

Fully-connected feedforward
network can be used

But CNN performs much better.

More Application: Playing Go

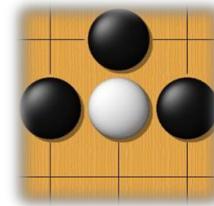
Training: record of previous plays 黑: 5之五 → 白: 天元 → 黑: 五之五 ...



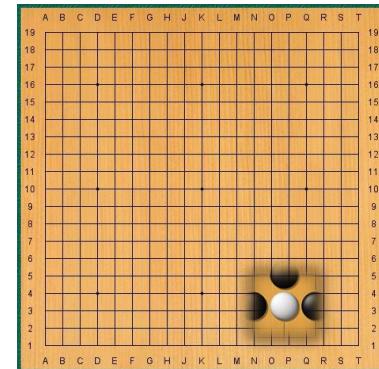
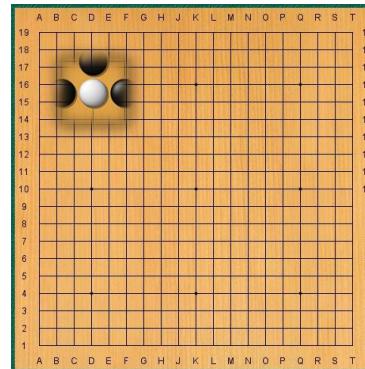
Why CNN for playing Go?

- 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.



Why CNN for playing Go?

- Subsampling the pixels will not change the object

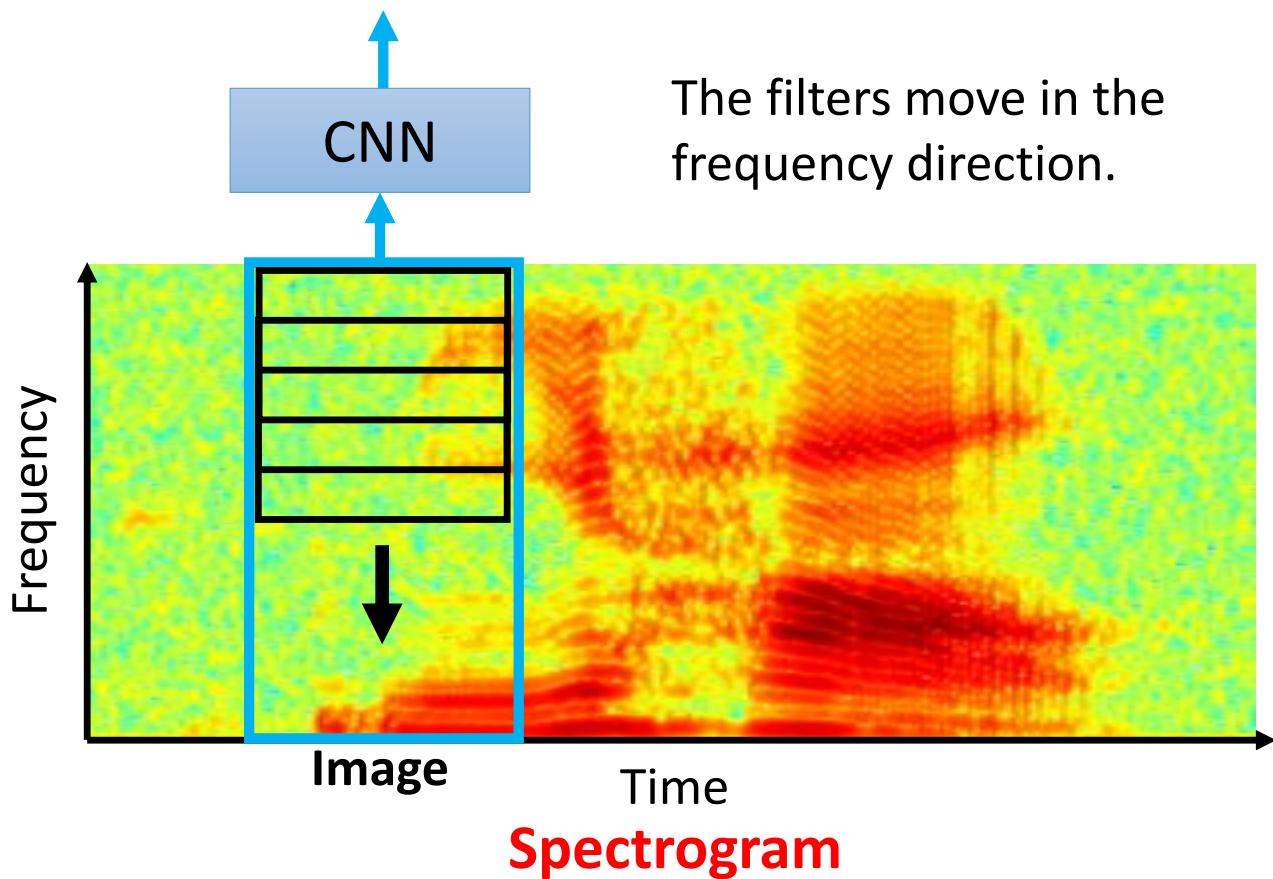


Max Pooling

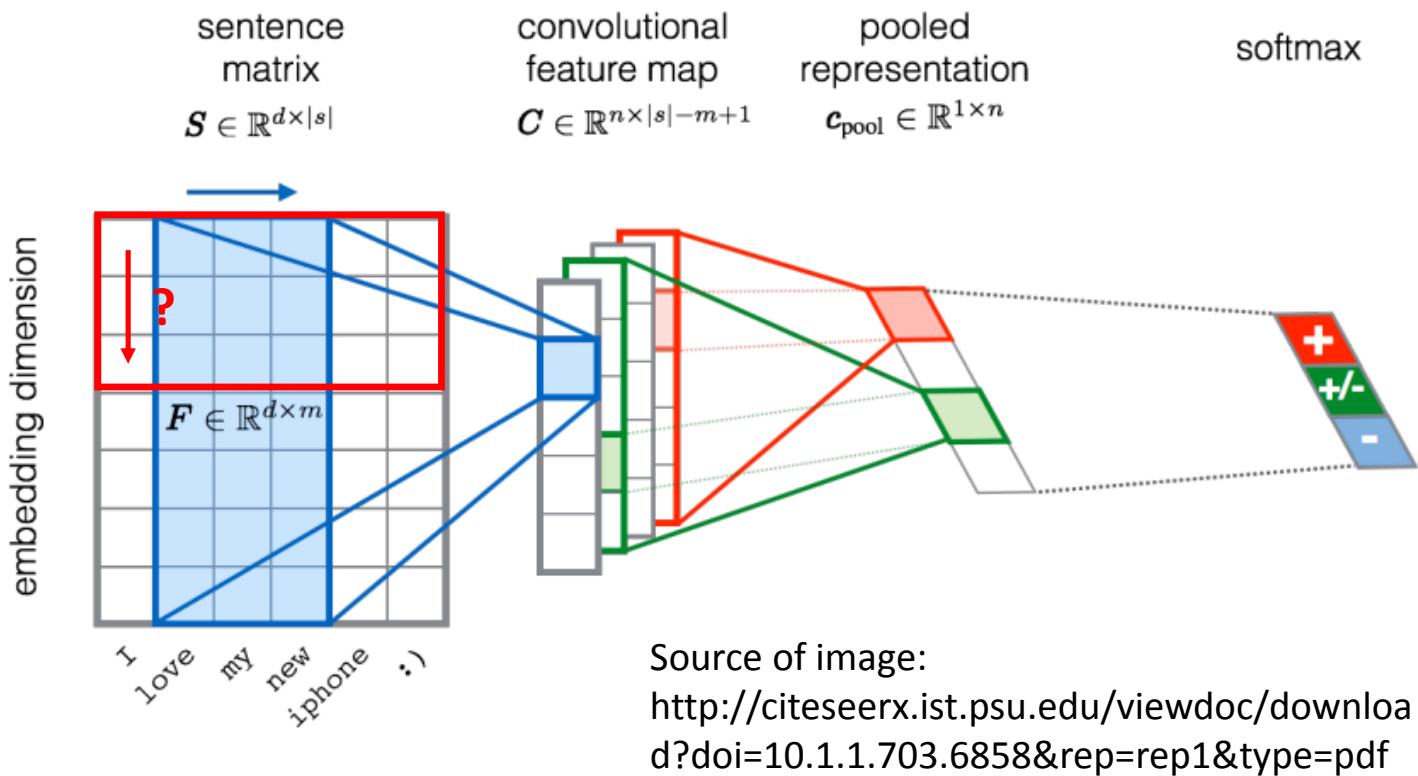
How to explain this???

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 23×23 image, then convolves k filters of kernel size 5×5 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 21×21 image, then convolves k filters of kernel size 3×3 with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size 1×1 with stride 1, with a different bias for each position, and applies a softmax function. The Alpha Go does not use Max Pooling Extended Data Table 3 additionally show the results of training with $k = 128, 256$ and 384 filters.

More Application: Speech



More Application: Text



Acknowledgment

- 感謝 Guobiao Mo 發現投影片上的打字錯誤