



大数据导论

Introduction to Big Data



大数据深度学习基础

叶允明

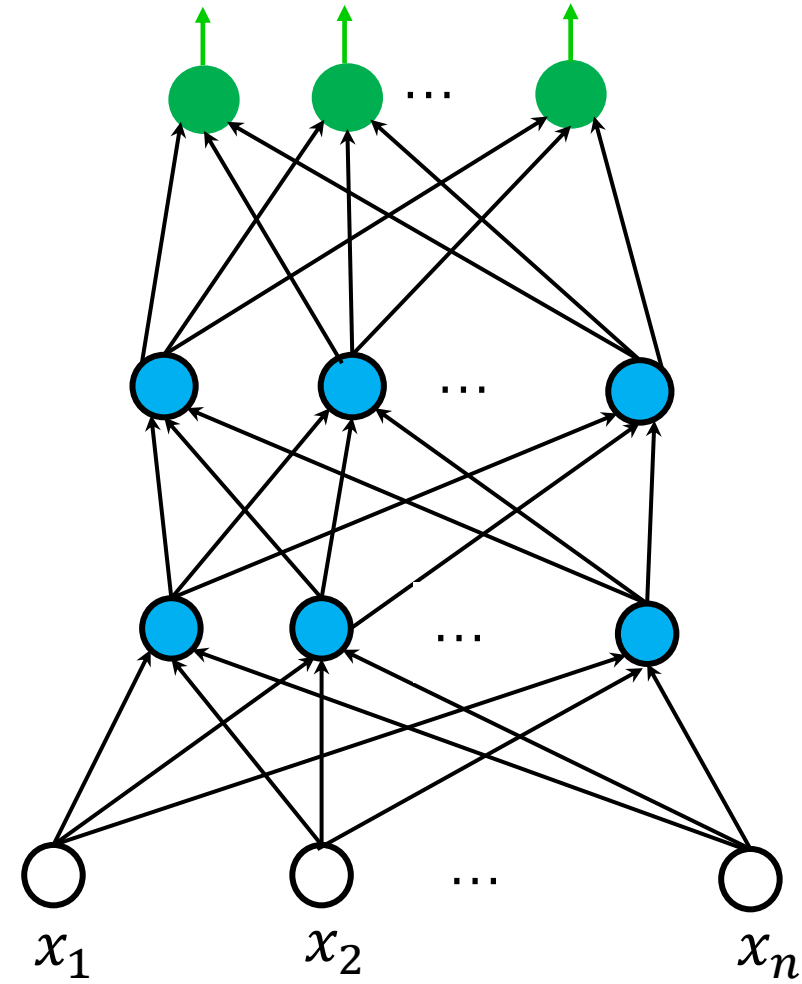
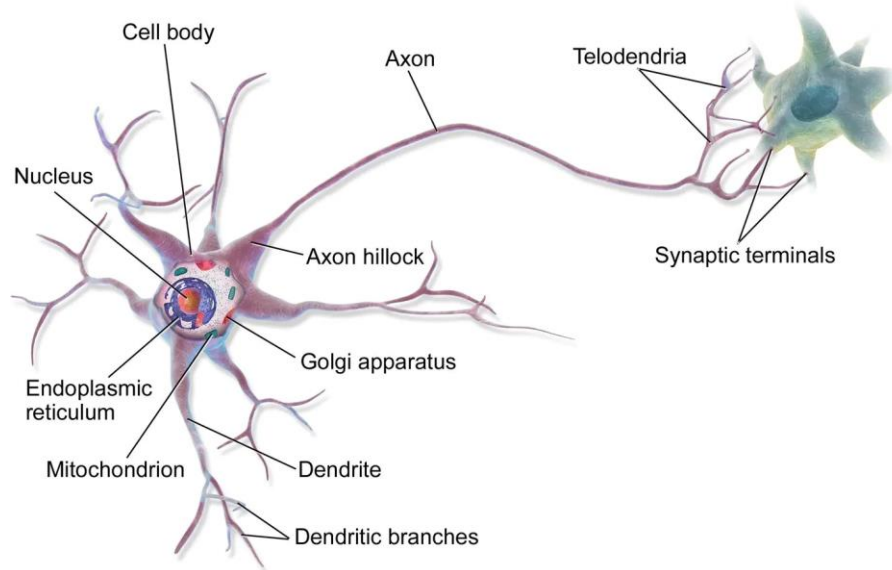
计算机科学与技术学院
哈尔滨工业大学（深圳）

目录

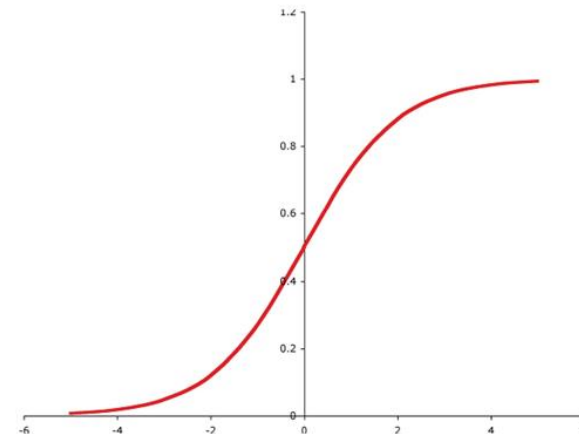
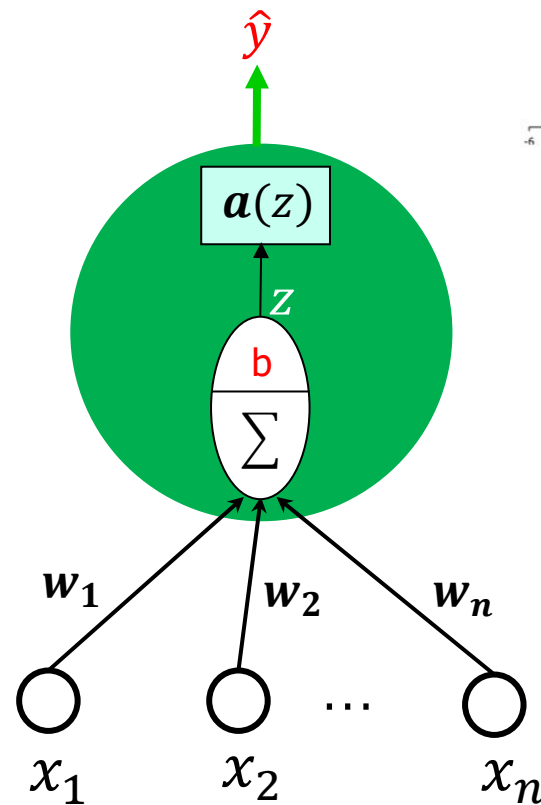
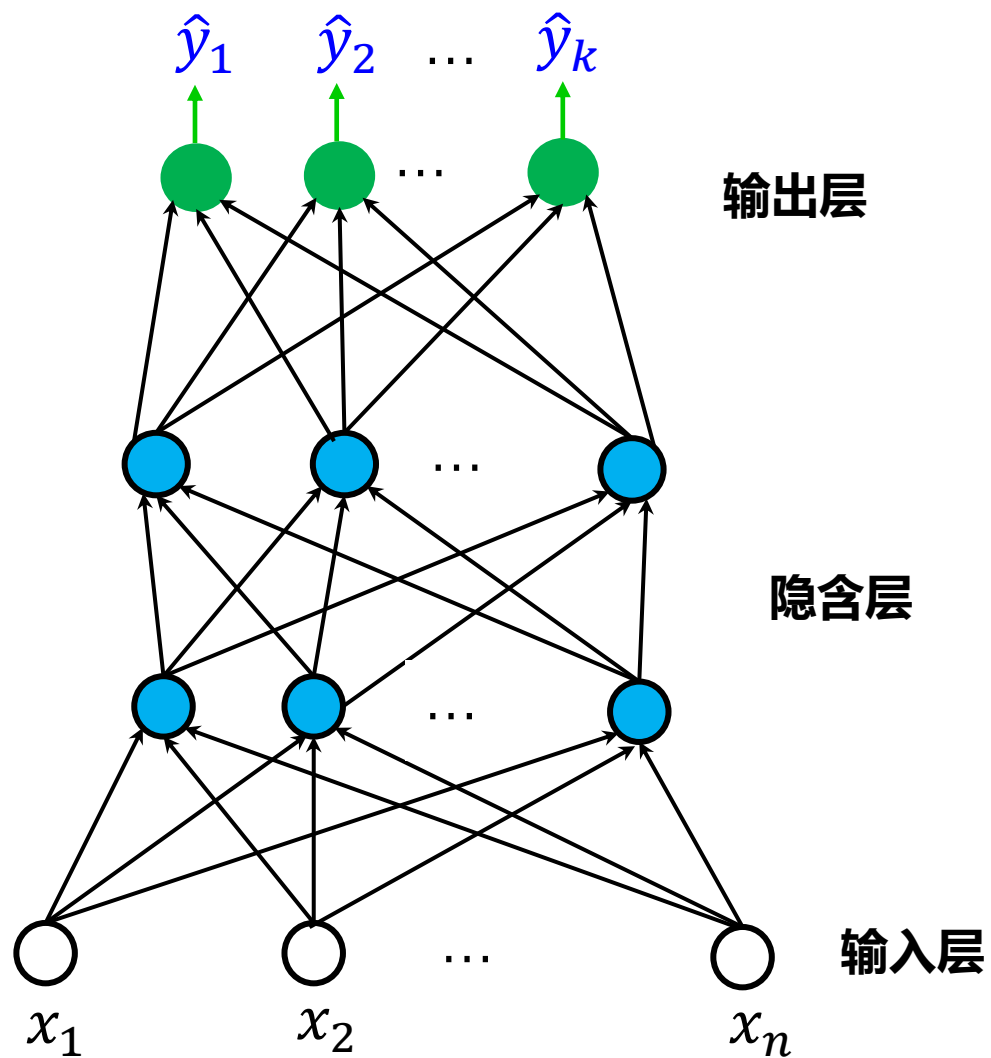
- 神经网络与深度学习基础
- 如何学习出更好的神经网络

神经网络与深度学习基础

人工神经网络(Artificial Neural Network)

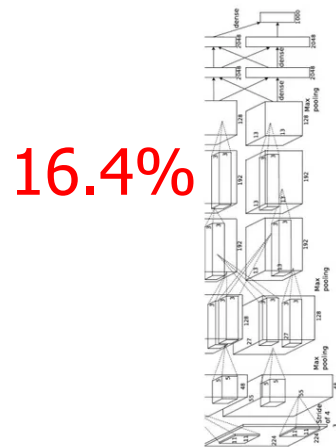


神经网络的常用结构



$$a(z) = \frac{1}{1 + e^{-z}}$$

深度学习：深度神经网络



16.4%

AlexNet (2012)

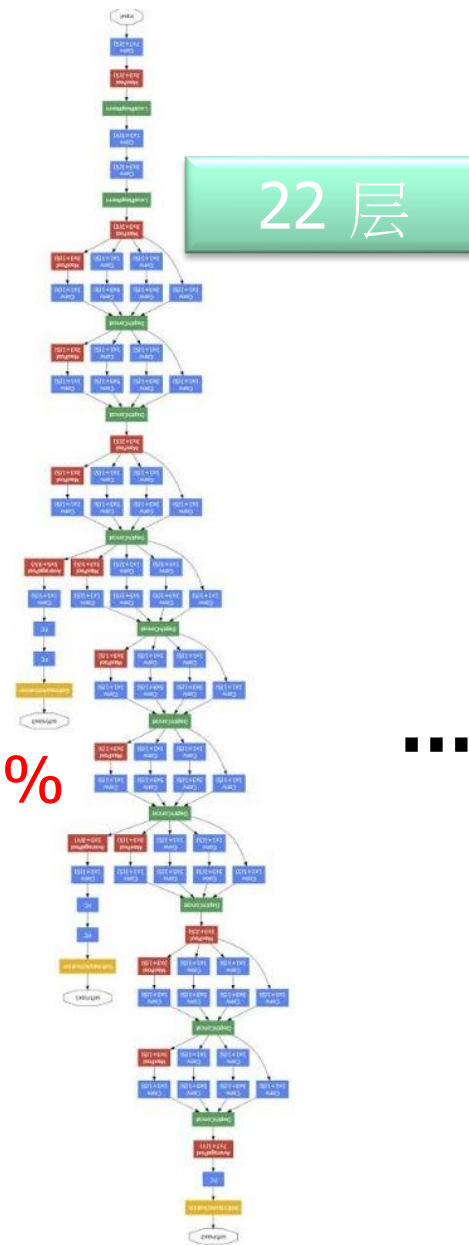
7.3%

VGG (2014)



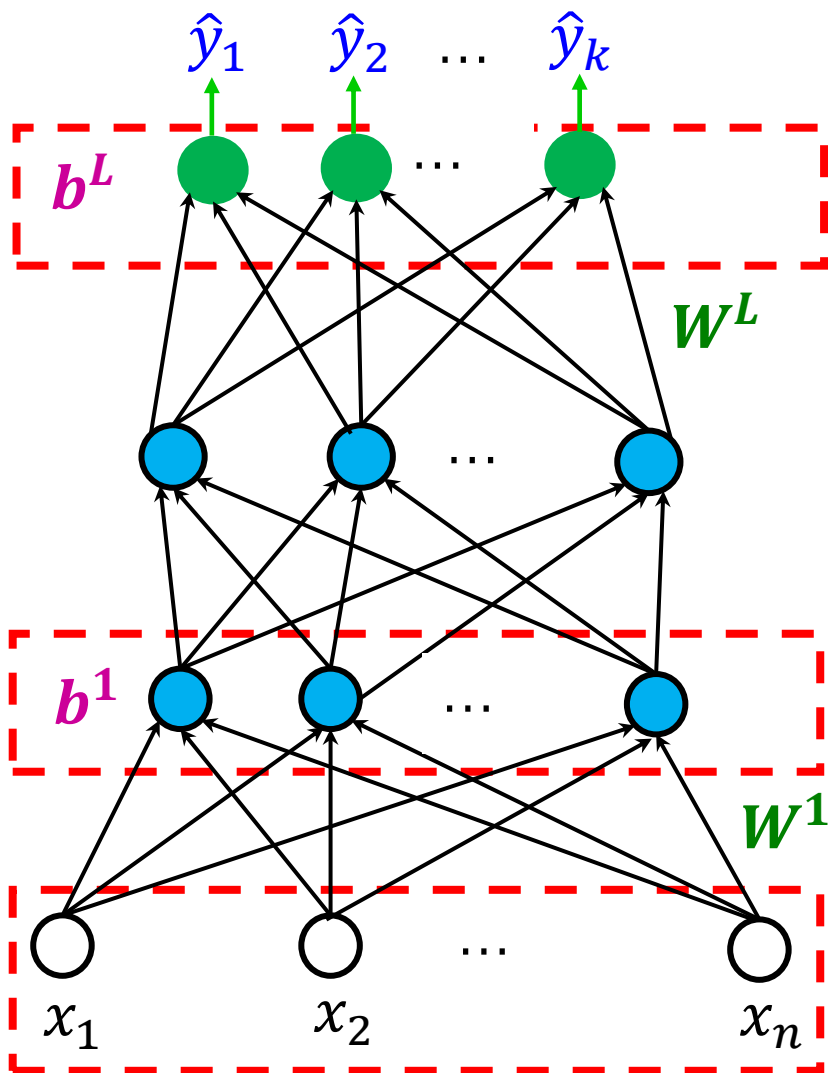
6.7%

GoogleNet (2014)



据推测：**GPT-4**
在**120层**中总共
包含了**1.8万亿**
参数！

全连接前馈神经网络的计算过程



$$a^L = \sigma(W^L a^{L-1} + b^L)$$

d^L is the total number of neurons of L layer

$$a^1 = \sigma(W^1 a^0 + b^1)$$

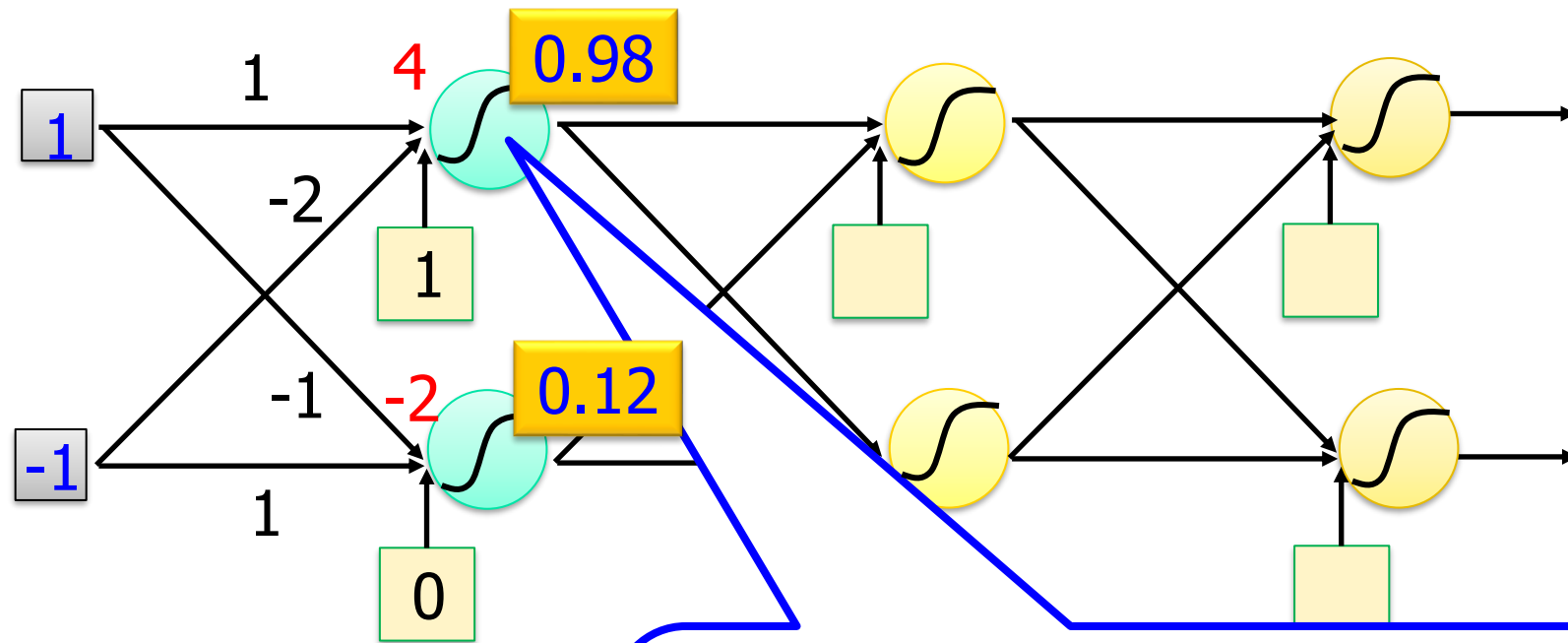
a^0

$$W^1 = \begin{bmatrix} (w_1^1)^T \\ (w_2^1)^T \\ \vdots \\ (w_i^1)^T \\ \vdots \\ (w_{d^1}^1)^T \end{bmatrix} = \begin{bmatrix} w_{11}^1 & \dots & w_{1d^0}^1 \\ \dots & w_{ij}^1 & \dots \\ w_{d^1d^0}^1 & \dots & w_{d^1d^0}^1 \end{bmatrix}$$

$$a(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

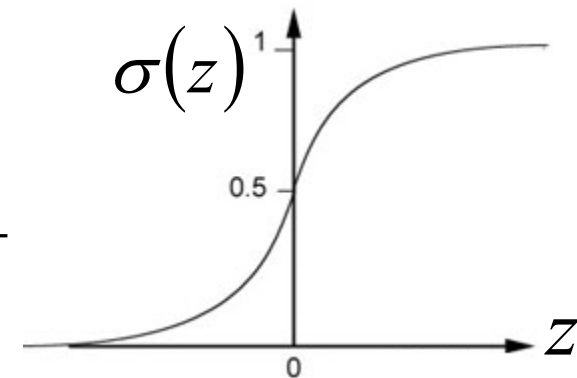
$$b^1 = \begin{bmatrix} b_1^1 \\ b_2^1 \\ \vdots \\ b_i^1 \\ \vdots \\ b_{d^1}^1 \end{bmatrix}$$

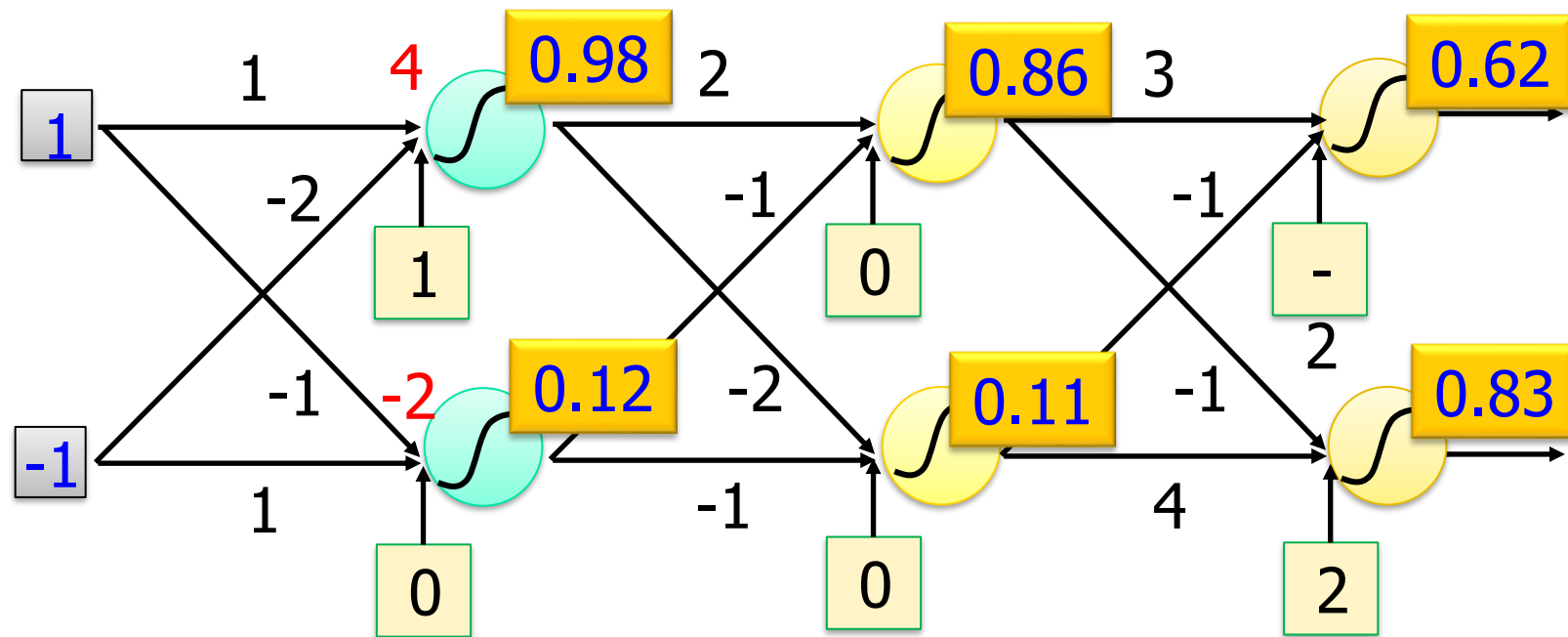
$$w_i^1 = \begin{bmatrix} w_{i1}^1 \\ w_{i2}^1 \\ \vdots \\ w_{ij}^1 \\ \vdots \\ w_{id^0}^1 \end{bmatrix}$$



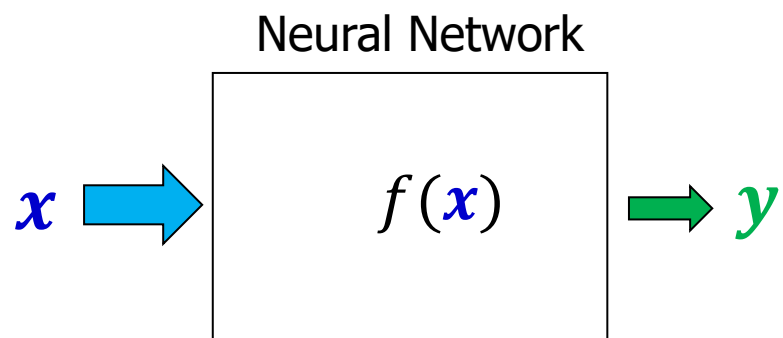
Sigmoid
Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

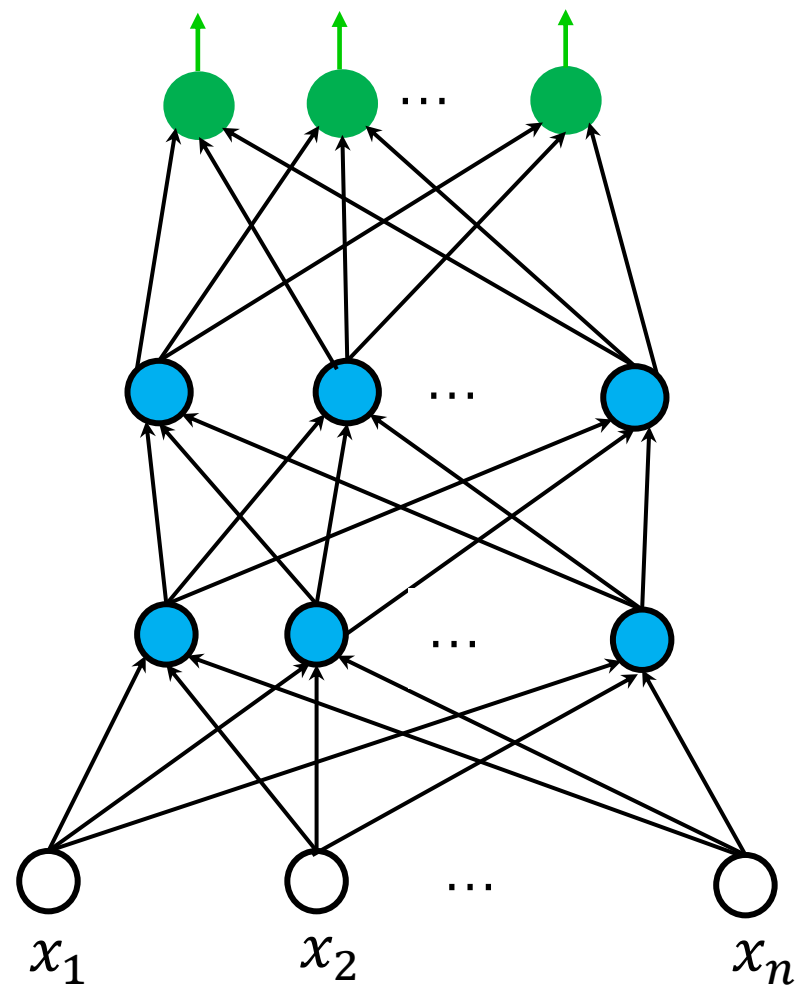




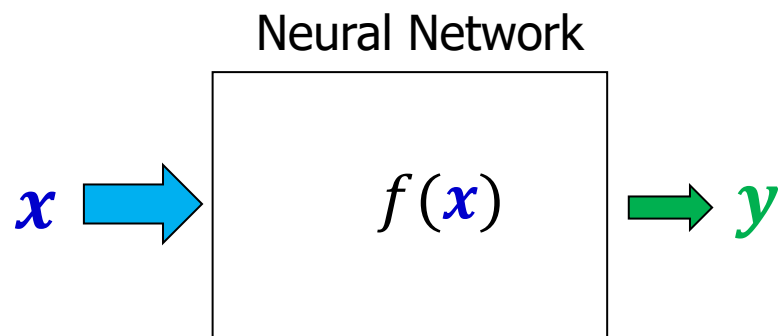
神经网络的本质



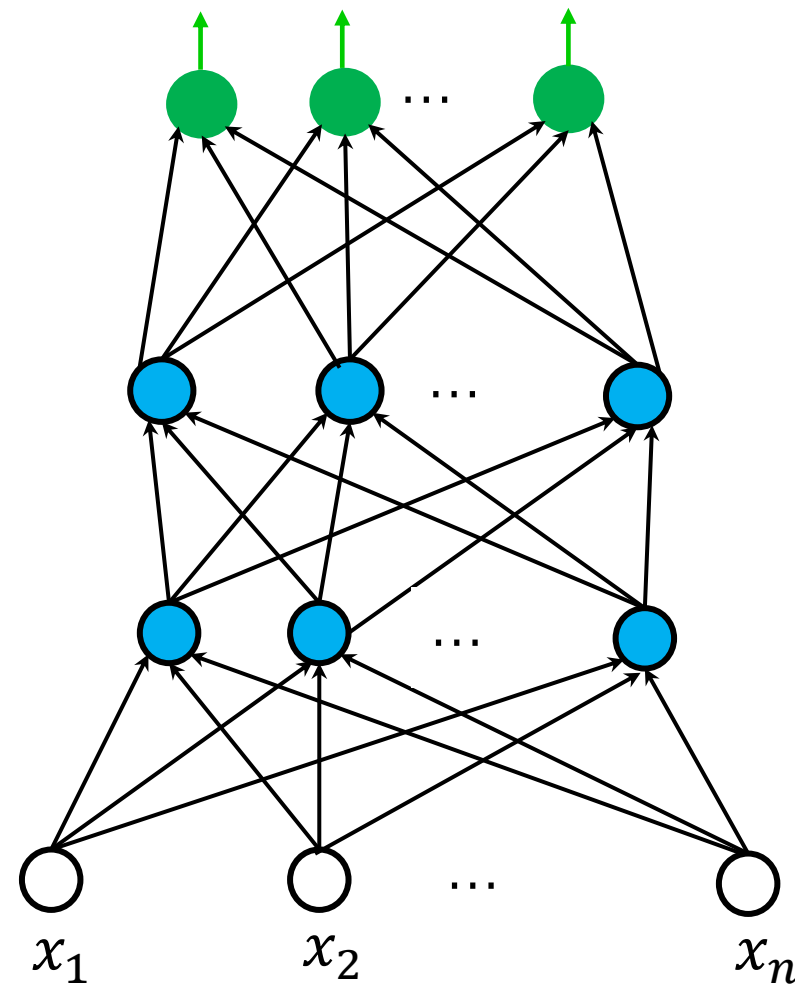
$$f(x) = f_n(f_{n-1}(\dots f_2(f_1(x)) \dots))$$



神经网络能做什么？



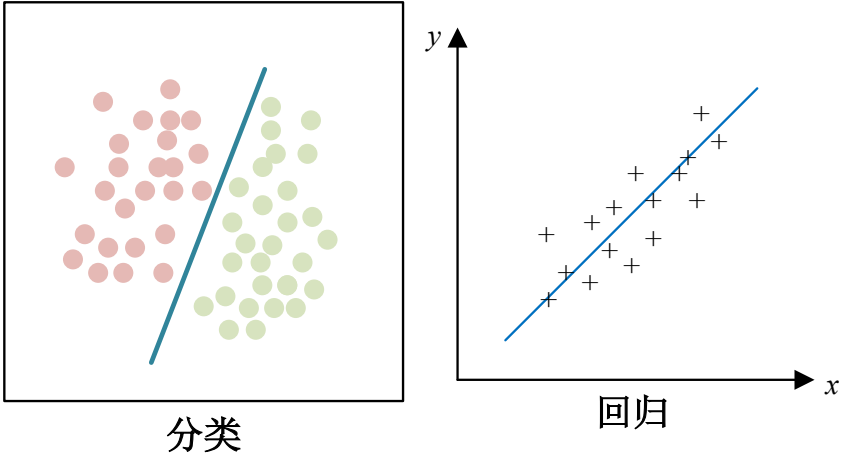
- 表征学习 (Representation learning)
- 执行决策 (Decision making)



为什么神经网络具有普适的应用价值

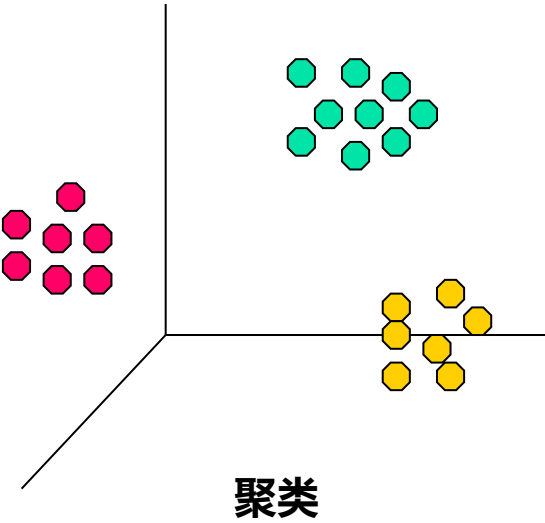
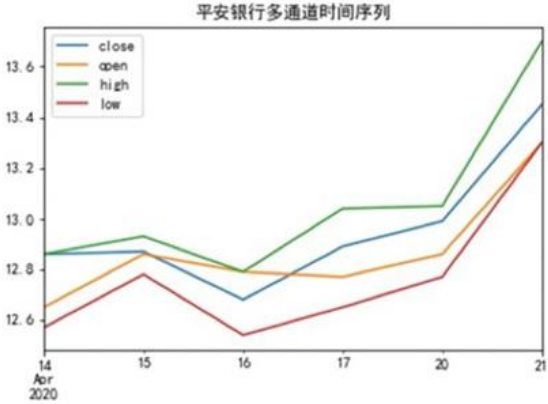
- 事物的向量化表示具有通用性
- 各种应用任务都可转化为函数 $f(x)$ 拟合问题

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



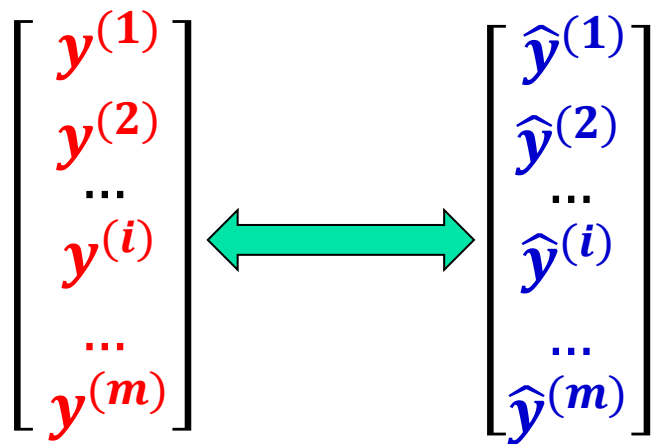
据中央气象台消息，今年第11号台风“轩岚诺”（超强台风级）的中心今天（8月31日）早晨5点钟位于日本冲绳县那霸市偏东方向约340公里的西北太平洋洋面上，就是北纬26.1度、东经131.1度，中心附近最大风力17级以上（62米/秒），中心最低气压为915百帕，七级风圈半径220~230公里，十级风圈半径70公里，十二级风圈半径40公里。

预计，“轩岚诺”将以每小时30公里左右的速度向西偏南转西南方向移动，9月1~2日将在琉球群岛以东洋面停滞或回旋，而后转向北偏西方向移动，3日夜间移入东海东南部海面。未来4~5天“轩岚诺”将维持超强台风级的强度，最大强度可达17级以上（62~70米/秒）。



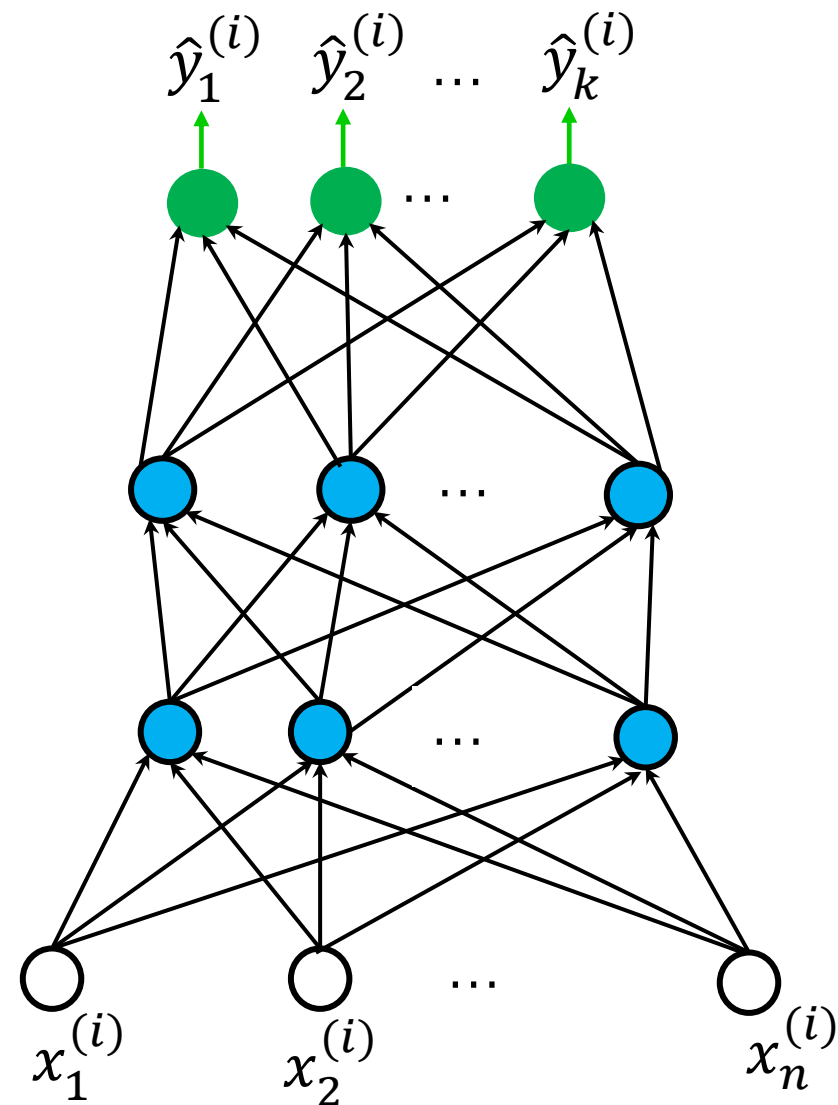
如何学习出神经网络模型: loss function

- Given $\mathbb{X} = \{\langle \mathbf{x}^{(1)}, y^{(1)} \rangle, \langle \mathbf{x}^{(2)}, y^{(2)} \rangle \dots, \langle \mathbf{x}^{(m)}, y^{(m)} \rangle\}$
- Translate $y^{(i)}$ to be a k-dim one-hot vector $\mathbf{y}^{(i)}$



Minimize loss:

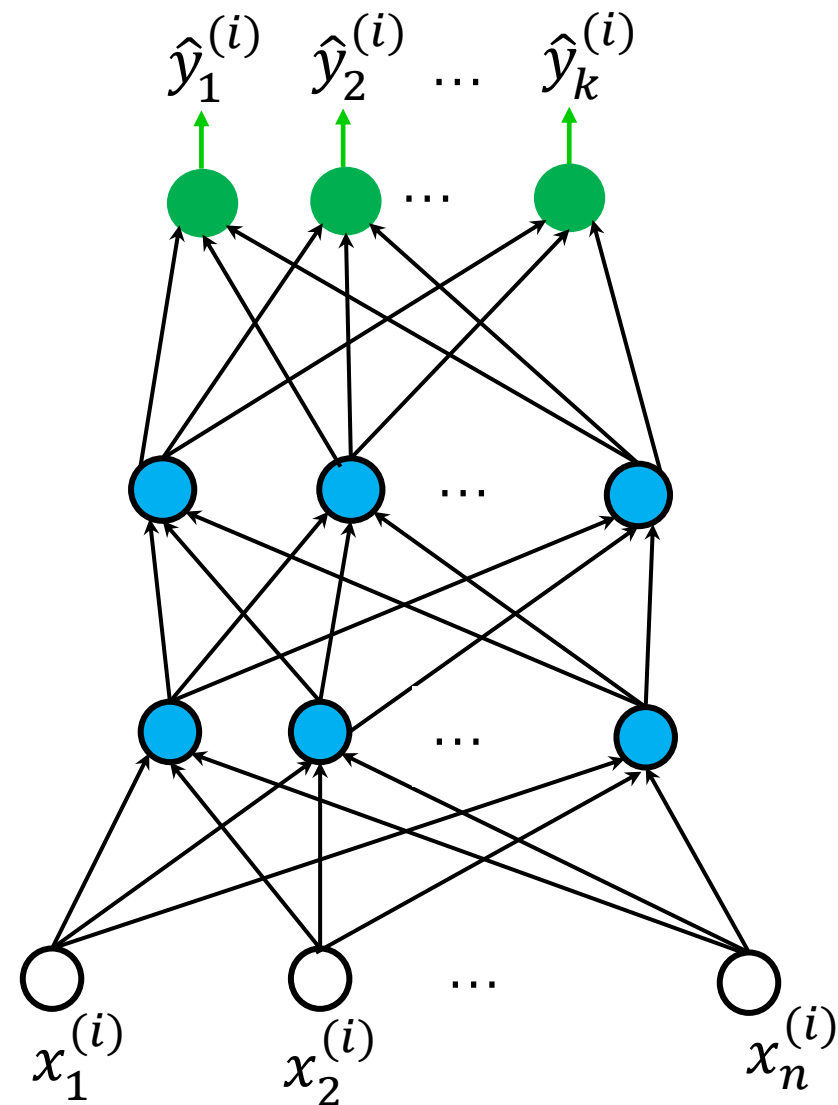
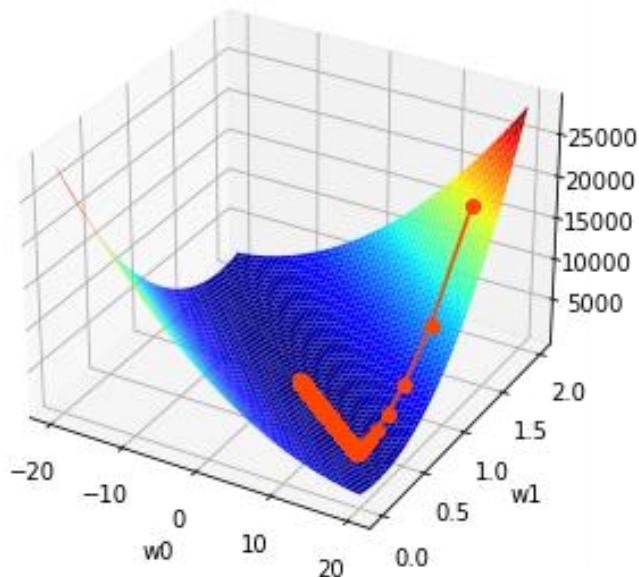
$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \|\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)}\|^2$$



如何学习出神经网络模型： 优化问题

$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \|\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)}\|^2$$

$$\begin{aligned} \hat{\mathbf{y}}^{(i)} &= f(\mathbf{x}^{(i)}) \\ &= \sigma(W^L \sigma(W^{L-1} \dots \sigma(W^2 \sigma(W^1 \mathbf{a}^0 + \mathbf{b}^1) + \mathbf{b}^2) \dots + \mathbf{b}^L) \end{aligned}$$



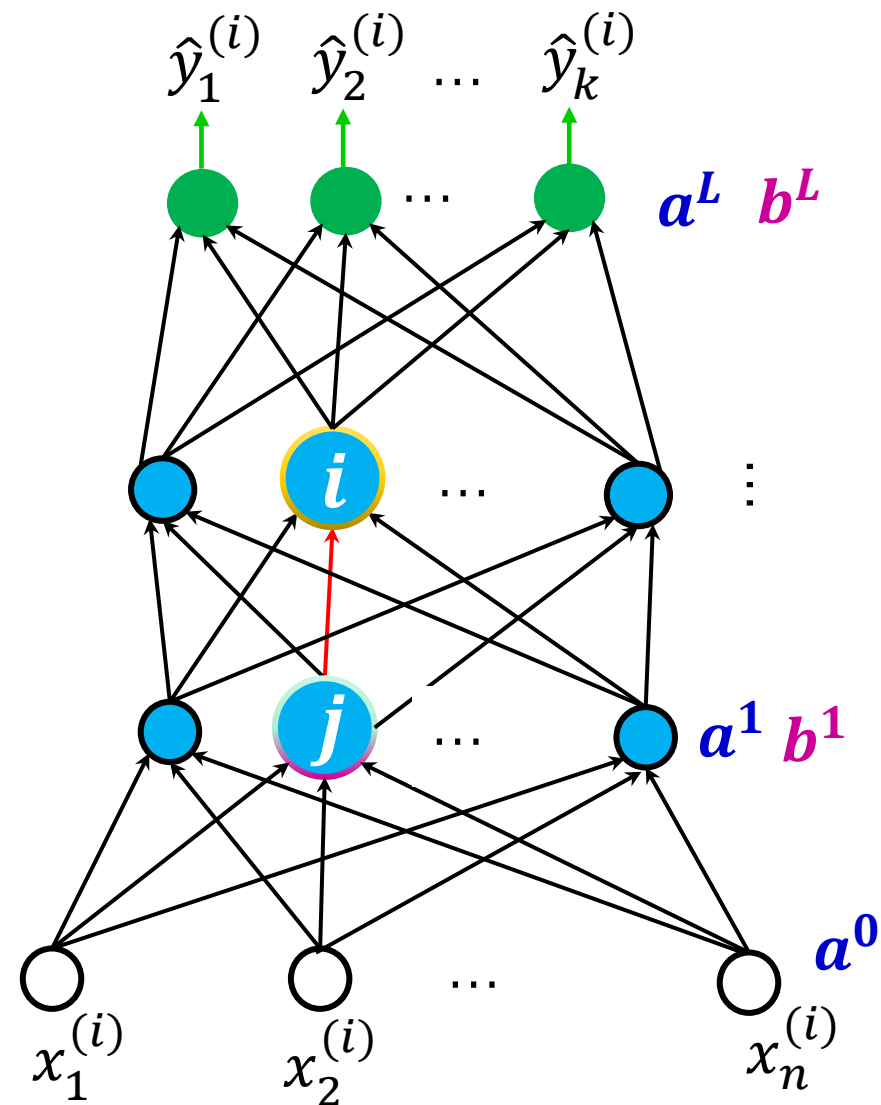
如何学习出神经网络模型：梯度下降法

$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \|\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)}\|^2$$

$$\left[\begin{array}{ccc} w_{11}^1 & \dots & w_{1d^0}^1 \\ \dots & w_{ij}^1 & \dots \\ w_{d^1d^0}^1 & \dots & w_{d^1d^0}^1 \end{array} \right] \dots \quad \nabla L = \left[\begin{array}{c} \frac{\partial L}{\partial w_{11}^1} \\ \frac{\partial L}{\partial w_{12}^1} \\ \vdots \\ \frac{\partial L}{\partial w_{ij}^r} \\ \vdots \end{array} \right]$$

$$\left[\begin{array}{c} b_1^1 \\ b_2^1 \\ \dots \\ b_i^1 \\ \dots \\ b_{d^1}^1 \end{array} \right] \dots$$

$$w_{ij}^{r(new)} = w_{ij}^{r(old)} - \eta \frac{\partial L}{\partial w_{ij}^r}$$



$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2$$

$$a_i^r = a(z_i^r):$$

$$z_i^r = \sum_{j=1}^{d^{r-1}} w_{ij}^r \cdot a_j^{r-1} + b_i^r$$

$$\frac{\partial L}{\partial w_{ij}^r} = ?$$

$$\hat{\mathbf{y}} = f(\mathbf{x})$$

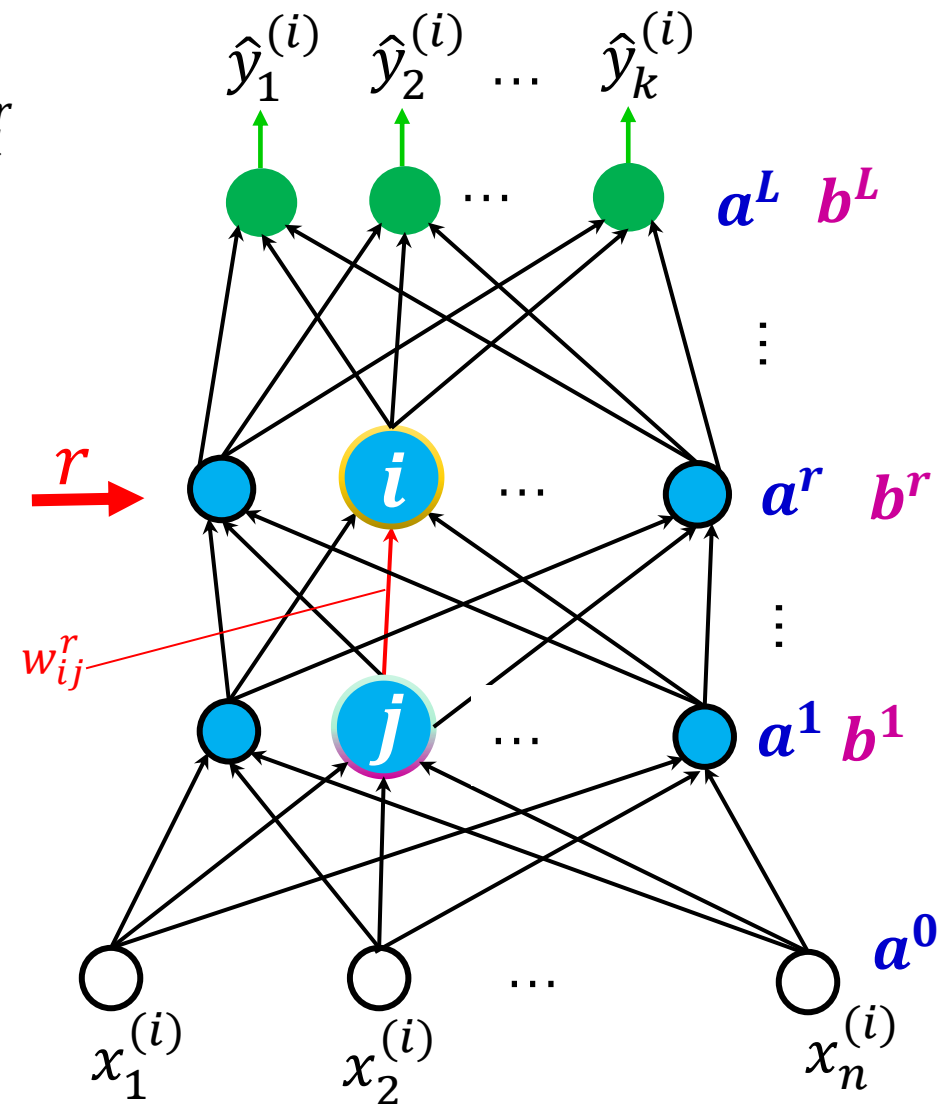
$$a(W^L a(W^{L-1} \dots a(W^r \dots a(W^1 \mathbf{x} + \mathbf{b}^1) \dots + \mathbf{b}^r) \dots + \mathbf{b}^L)$$

$$w_{ij}^r \cdot a_j^{r-1} \longrightarrow z_i^r \longrightarrow a_i^r \dots a^L = \hat{\mathbf{y}} \longrightarrow L(\cdot)$$

Chain rule for derivation:

$$y = f(u), \quad u = g(x)$$

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



BP (Back Propagation) 算法: basic idea

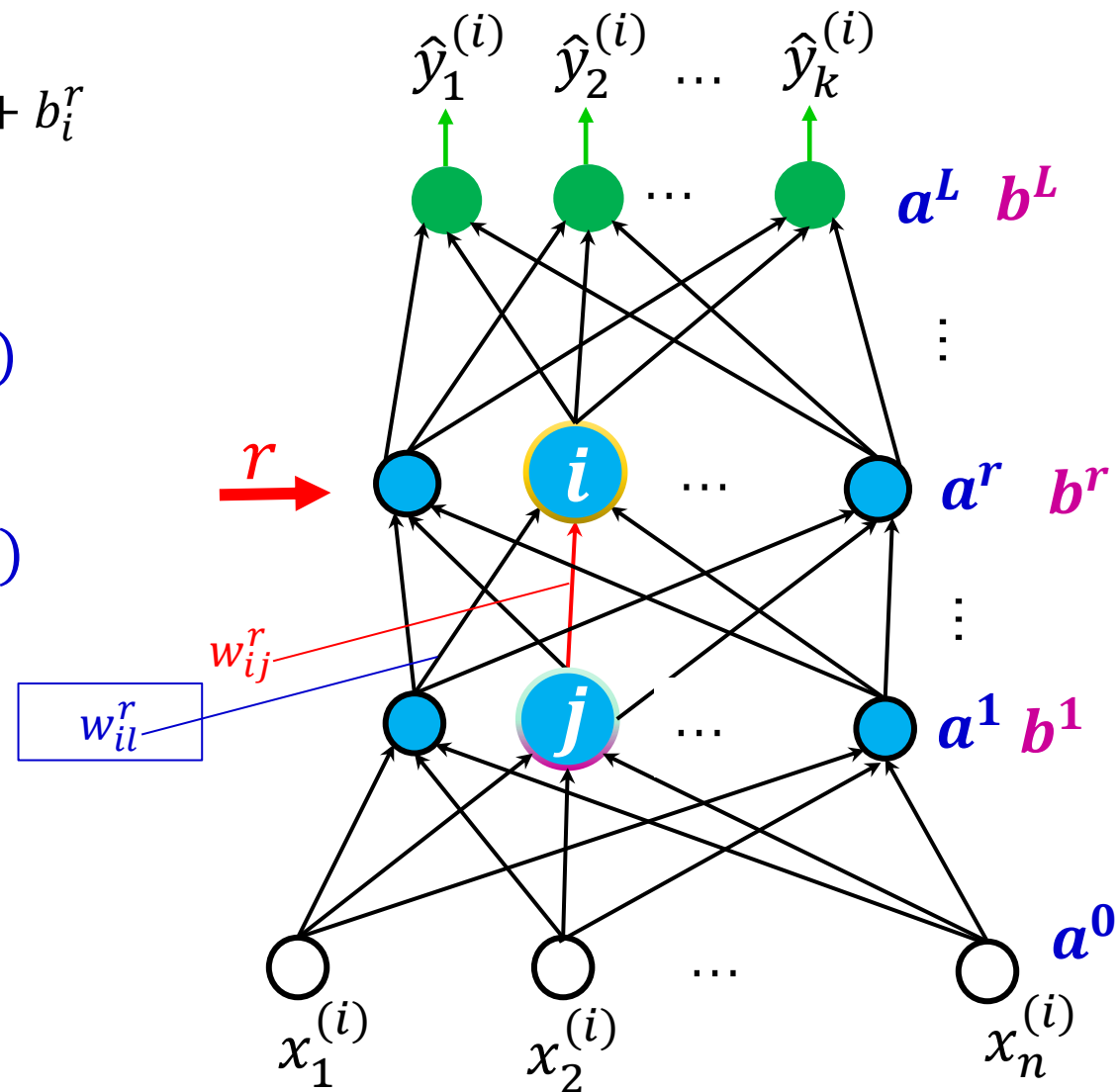
$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2 \quad z_i^r = \sum_{j=1}^{d^{r-1}} w_{ij}^r \cdot a_j^{r-1} + b_i^r$$

$$w_{ij}^r \cdot a_j^{r-1} \longrightarrow z_i^r \longrightarrow a_i^r \quad \dots \quad \mathbf{a}^L = \hat{\mathbf{y}} \longrightarrow L(\cdot)$$

$$w_{il}^r \cdot a_l^{r-1} \longrightarrow z_i^r \longrightarrow a_i^r \quad \dots \quad \mathbf{a}^L = \hat{\mathbf{y}} \longrightarrow L(\cdot)$$

$$\frac{\partial L}{\partial w_{ij}^r} = \frac{\partial L}{\partial \mathbf{a}^L} \dots \dots \frac{\partial \mathbf{a}^{r+1}}{\partial a_i^r} \frac{\partial a_i^r}{\partial z_i^r} \frac{\partial z_i^r}{\partial w_{ij}^r}$$

$$\frac{\partial L}{\partial b_i^r} = \frac{\partial L}{\partial \mathbf{a}^L} \dots \dots \frac{\partial \mathbf{a}^{r+1}}{\partial a_i^r} \frac{\partial a_i^r}{\partial z_i^r} \frac{\partial z_i^r}{\partial b_i^r}$$



$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2$$

$$z_i^r = \sum_{j=1}^{d^{r-1}} w_{ij}^r \cdot a_j^{r-1} + b_i^r$$

$$\frac{\partial L}{\partial w_{ij}^r} = \underbrace{\frac{\partial L}{\partial \mathbf{a}^L} \cdots \cdots \frac{\partial \mathbf{a}^{r+1}}{\partial a_i^r} \frac{\partial a_i^r}{\partial z_i^r} \frac{\partial z_i^r}{\partial w_{ij}^r}}_{\delta_i^r}$$

$$\delta_i^r = \frac{\partial L}{\partial z_i^r}$$

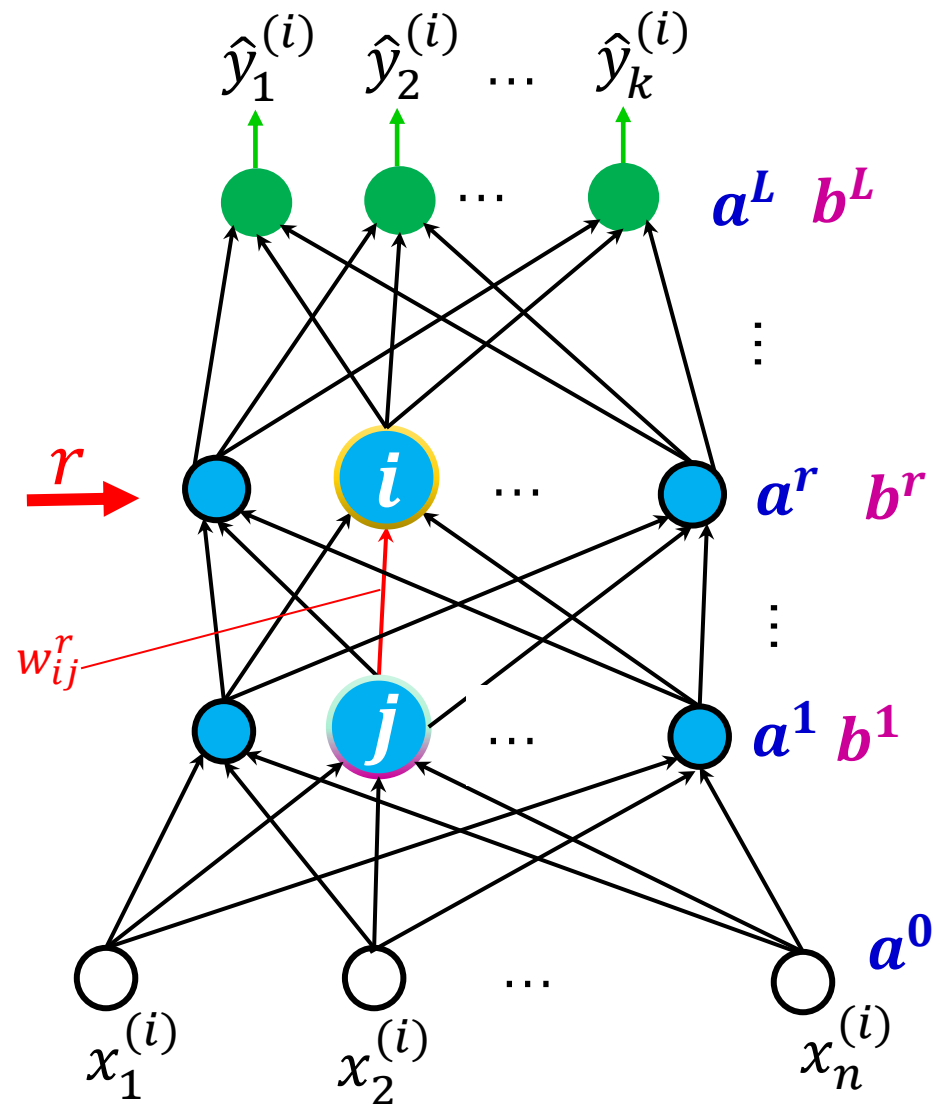
$$\frac{\partial L}{\partial w_{ij}^r} = \delta_i^r \times a_j^{r-1}$$

$$\frac{\partial z_i^r}{\partial w_{ij}^r} = a_j^{r-1}$$

$$\frac{\partial a_i^r}{\partial z_i^r} :$$

$$\downarrow \text{a: } \sigma = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial a_i^r}{\partial z_i^r} = a_i^r (1 - a_i^r)$$



$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2$$

$$z_i^r = \sum_{j=1}^{d^{r-1}} w_{ij}^r \cdot a_j^{r-1} + b_i^r$$

$$\frac{\partial L}{\partial w_{ij}^r} = \frac{\partial L}{\partial \mathbf{a}^L} \cdots \cdots \underbrace{\frac{\partial \mathbf{a}^{r+1}}{\partial a_i^r} \frac{\partial a_i^r}{\partial z_i^r} \frac{\partial z_i^r}{\partial w_{ij}^r}}_{\delta_i^r}$$

$$\delta_i^r = \frac{\partial L}{\partial z_i^r}$$

$$\frac{\partial L}{\partial w_{ij}^r} = \delta_i^r \times a_j^{r-1}$$

Output Layer($r = L$):

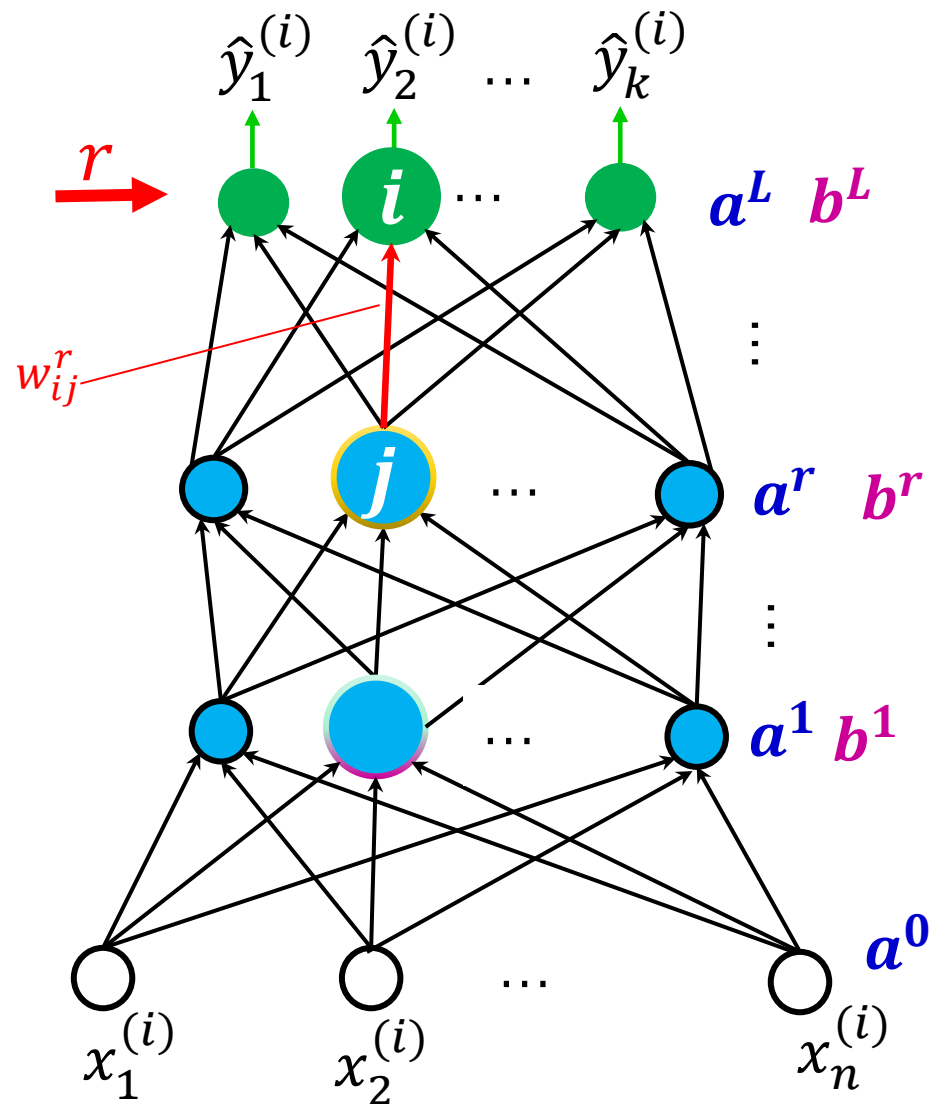
$$\delta_i^r = -(\mathbf{y}_i - \mathbf{a}_i^L) \mathbf{a}_i^L (1 - \mathbf{a}_i^L)$$

$$\frac{\partial z_i^r}{\partial w_{ij}^r} = a_j^{r-1}$$

$$\frac{\partial a_i^r}{\partial z_i^r} :$$

$$\downarrow \sigma = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial a_i^r}{\partial z_i^r} = a_i^r (1 - a_i^r)$$



$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2$$

$$z_i^r = \sum_{j=1}^{d^{r-1}} w_{ij}^r \cdot a_j^{r-1} + b_i^r$$

$$\frac{\partial L}{\partial w_{ij}^r} = \frac{\partial L}{\partial \mathbf{a}^L} \cdots \cdots \underbrace{\frac{\partial \mathbf{a}^{r+1}}{\partial a_i^r} \frac{\partial a_i^r}{\partial z_i^r} \frac{\partial z_i^r}{\partial w_{ij}^r}}_{\delta_i^r}$$

$$\delta_i^r = \frac{\partial L}{\partial z_i^r}$$

$$\frac{\partial L}{\partial w_{ij}^r} = \delta_i^r \times a_j^{r-1}$$

Output Layer ($r = L$):

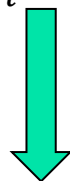
$$\delta_i^r = -(\mathbf{y}_i - \mathbf{a}_i^L) \mathbf{a}_i^L (1 - \mathbf{a}_i^L)$$

Hidden Layer ($r < L$):

$$\delta_i^r = \mathbf{a}_i^r (1 - \mathbf{a}_i^r) \sum_{s \in \text{Next}(i)} \delta_s^{r+1} w_{si}^{r+1}$$

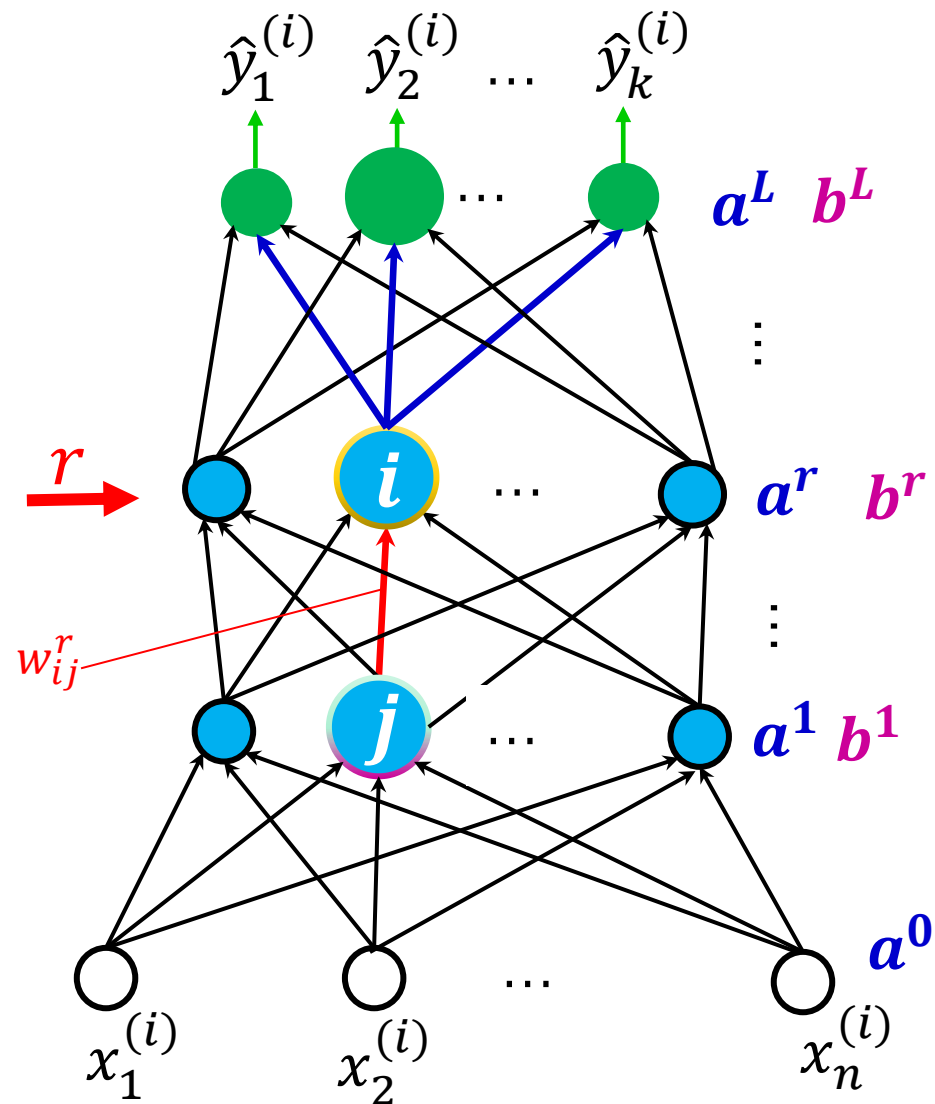
$$\frac{\partial z_i^r}{\partial w_{ij}^r} = a_j^{r-1}$$

$$\frac{\partial a_i^r}{\partial z_i^r} :$$



$$a: \sigma = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial a_i^r}{\partial z_i^r} = a_i^r (1 - a_i^r)$$



批处理反向传播 (BP with Batch Training)

$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2$$

- Given $\mathbb{X} = \{ \langle \mathbf{x}^{(1)}, y^{(1)} \rangle, \langle \mathbf{x}^{(2)}, y^{(2)} \rangle \dots, \langle \mathbf{x}^{(m)}, y^{(m)} \rangle \}$, and hyper parameters

➤ Translate $y^{(i)}$ to be a k-dim one-hot vector $\mathbf{y}^{(i)}$ for classification task

1. Initialize network weights with small random values

$$\frac{\partial L}{\partial w_{ij}^r} = \delta_i^r \times a_j^{r-1}$$

2. **Do**

① Initialize Loss: $L(W, b) = 0$

$$\delta_i^r = -(\mathbf{y}_i - a_i^L) a_i^L (1 - a_i^L)$$

② **For each sample** $\langle \mathbf{x}^{(N)}, y^{(N)} \rangle$

a. **Compute the output of** $\hat{\mathbf{y}}^{(N)} = f(\mathbf{x}^{(N)})$

$$\delta_i^r = a_i^r (1 - a_i^r) \sum_{s \in \text{Next}(i)} \delta_s^{r+1} w_{si}^{r+1}$$

b. **Update the loss:** $L(W, b) = L(W, b) + \frac{1}{2} (y^{(N)} - \hat{\mathbf{y}}^{(N)})^2$

③ **Calculate the gradient with** $L(W, b)$ **and current weights** $w_{ij}^{r(\text{old})}$

④ **Update the weights:** $w_{ij}^{r(\text{new})} = w_{ij}^{r(\text{old})} - \eta \frac{\partial L}{\partial w_{ij}^r}$

3. **Until stopping criteria satisfied**

随机梯度下降算法 (BP with SGD)

$$L(W, b) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2$$

- Given $\mathbb{X} = \{ \langle \mathbf{x}^{(1)}, y^{(1)} \rangle, \langle \mathbf{x}^{(2)}, y^{(2)} \rangle \dots, \langle \mathbf{x}^{(m)}, y^{(m)} \rangle \}$, and hyper parameters

➤ Translate $y^{(i)}$ to be a k-dim one-hot vector $\mathbf{y}^{(i)}$ for classification task

1. Initialize network weights with small random values

$$\frac{\partial L}{\partial w_{ij}^r} = \delta_i^r \times a_j^{r-1}$$

2. Do

① For each sample $\langle \mathbf{x}^{(N)}, y^{(N)} \rangle$

$$\delta_i^r = -(\mathbf{y}_i - a_i^L) a_i^L (1 - a_i^L)$$

a. Compute the output of $\hat{\mathbf{y}}^{(N)} = f(\mathbf{x}^{(N)})$

$$\delta_i^r = a_i^r (1 - a_i^r) \sum_{s \in \text{Next}(i)} \delta_s^{r+1} w_{si}^{r+1}$$

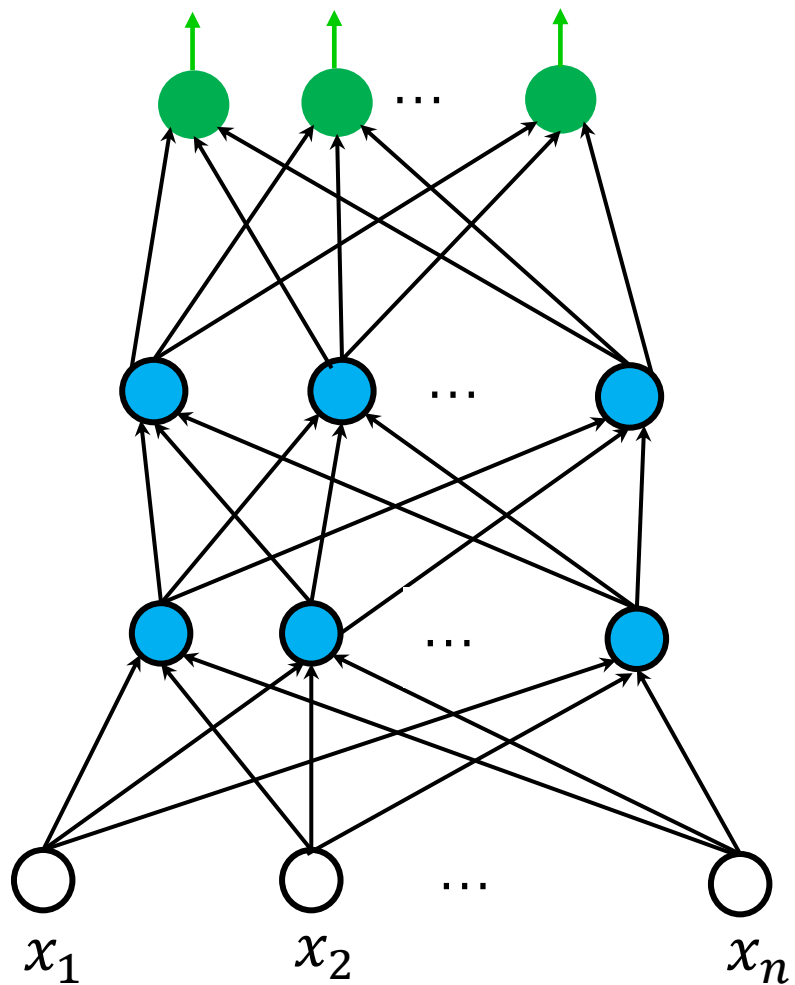
b. Compute the loss of $\mathbf{x}^{(N)}$: $l(W, b) = \frac{1}{2} (y^{(N)} - \hat{\mathbf{y}}^{(N)})^2$

c. Calculate the gradient with $l(W, b)$ and current weights $w_{ij}^{r(\text{old})}$

d. Update the weights: $w_{ij}^{r(\text{new})} = w_{ij}^{r(\text{old})} - \eta \frac{\partial L}{\partial w_{ij}^r}$

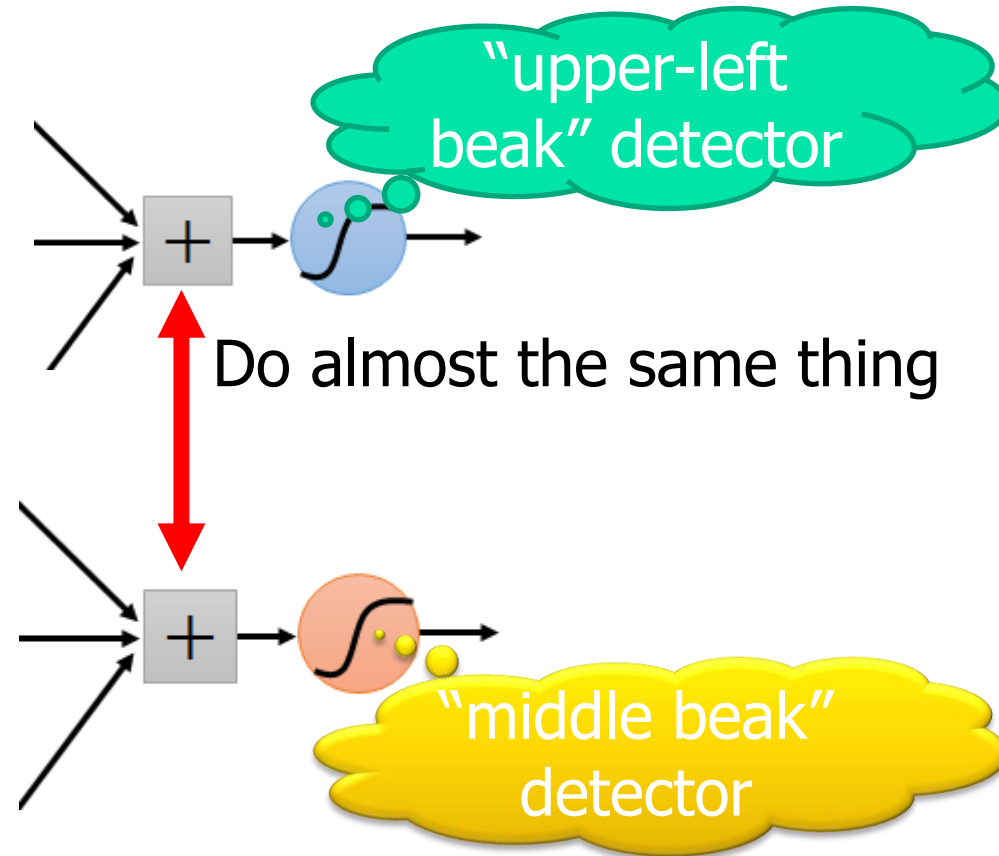
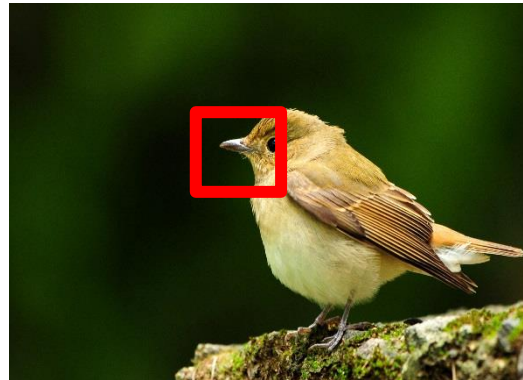
3. Until stopping criteria satisfied

全连接网络的问题



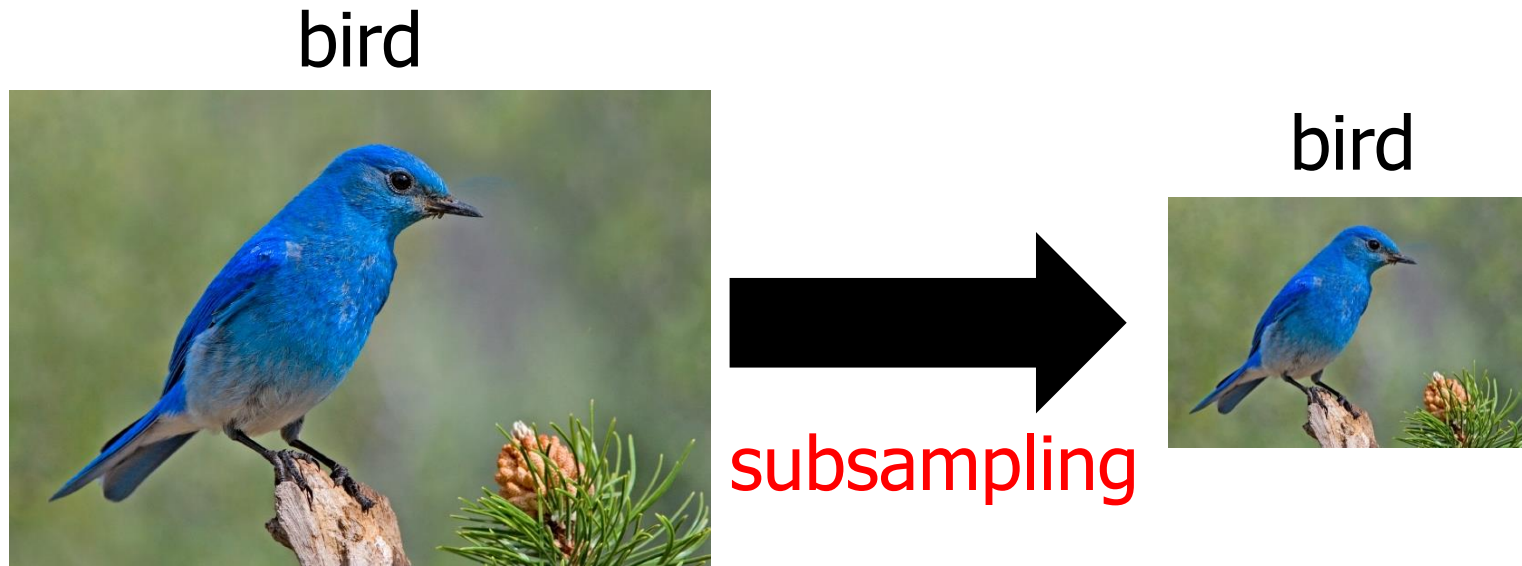
卷积神经网络 (CNN) : 可挖掘局部模式

- The same patterns appear in different regions.



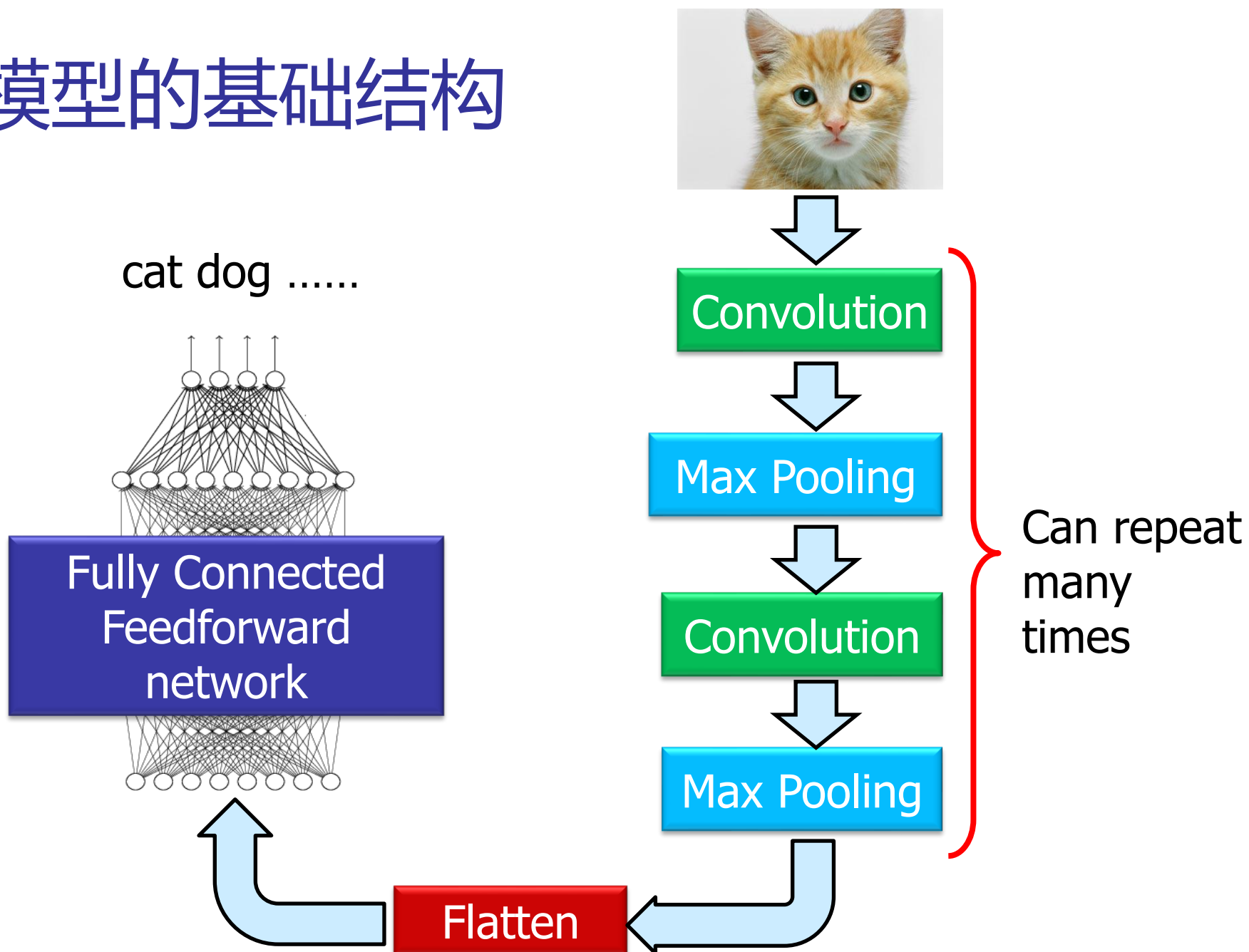
卷积神经网络 (CNN) : 可挖掘局部模式

- Subsampling the pixels will not change the object



We can subsample the pixels to make image smaller
Less parameters for the network to process the image

CNN模型的基础结构



CNN – 卷积 (Convolution)

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

⋮ ⋮

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

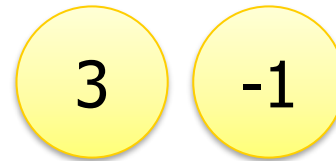
stride=1

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 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

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

Filter 1

3 -3

We set stride=1 below

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 1

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

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

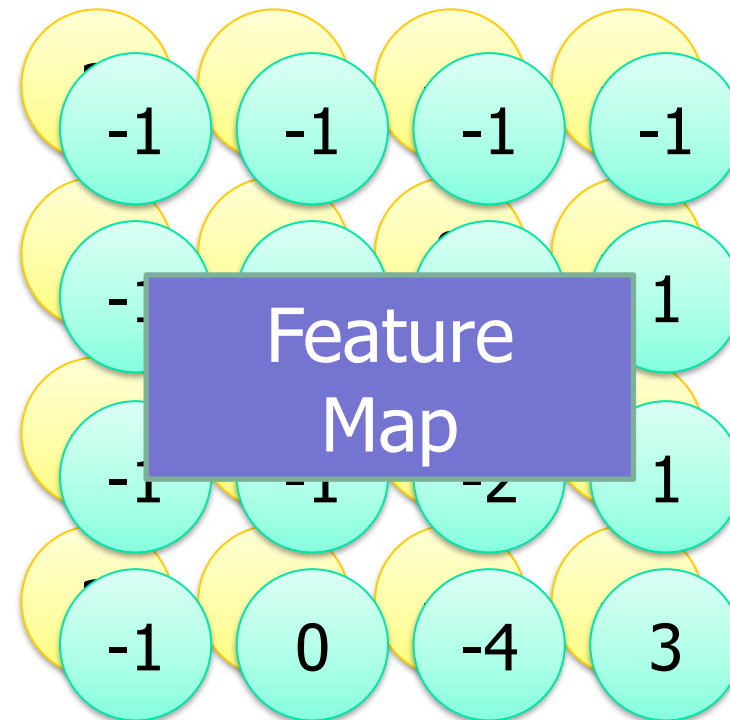
Filter 2

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

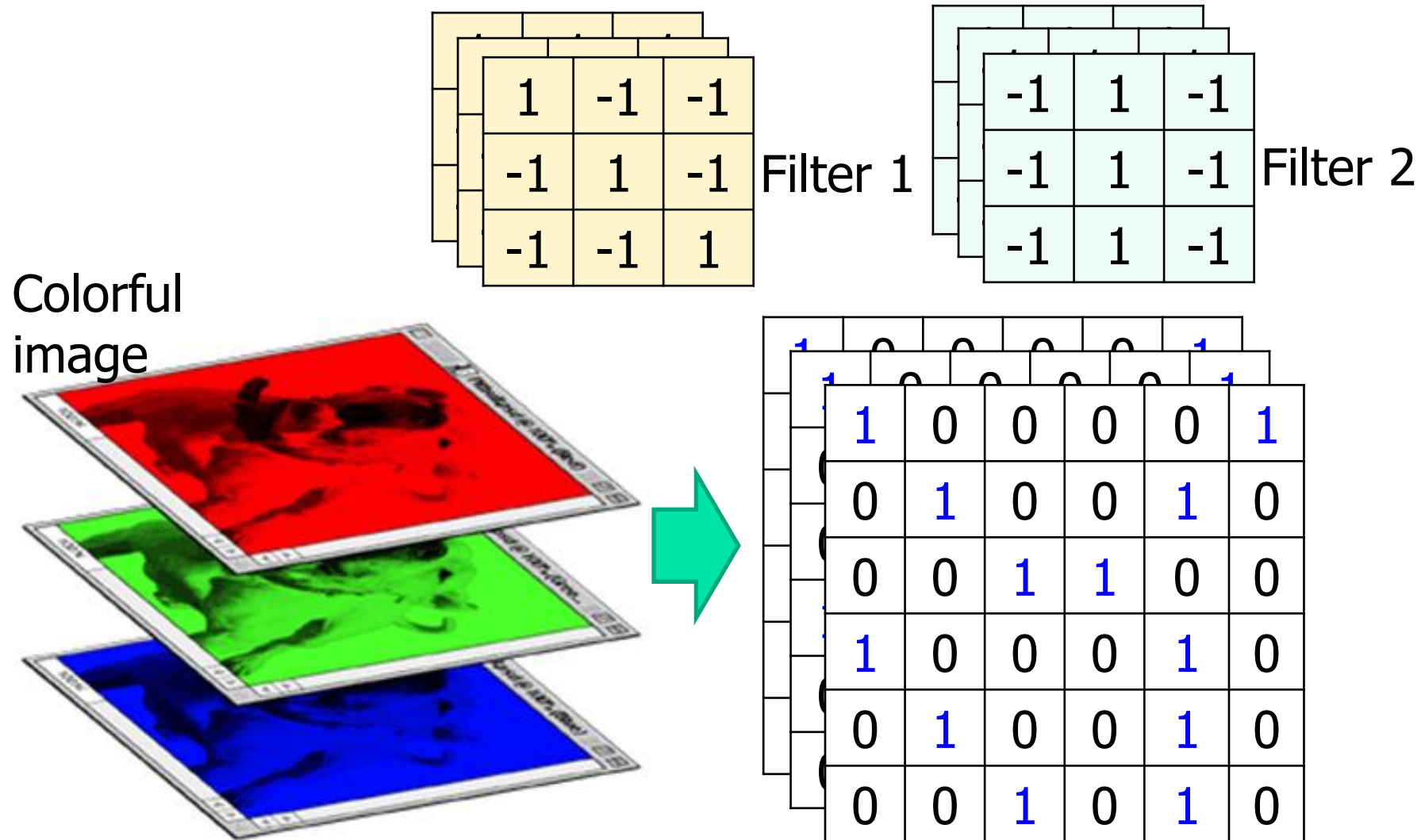
6 x 6 image

Do the same process
for every filter

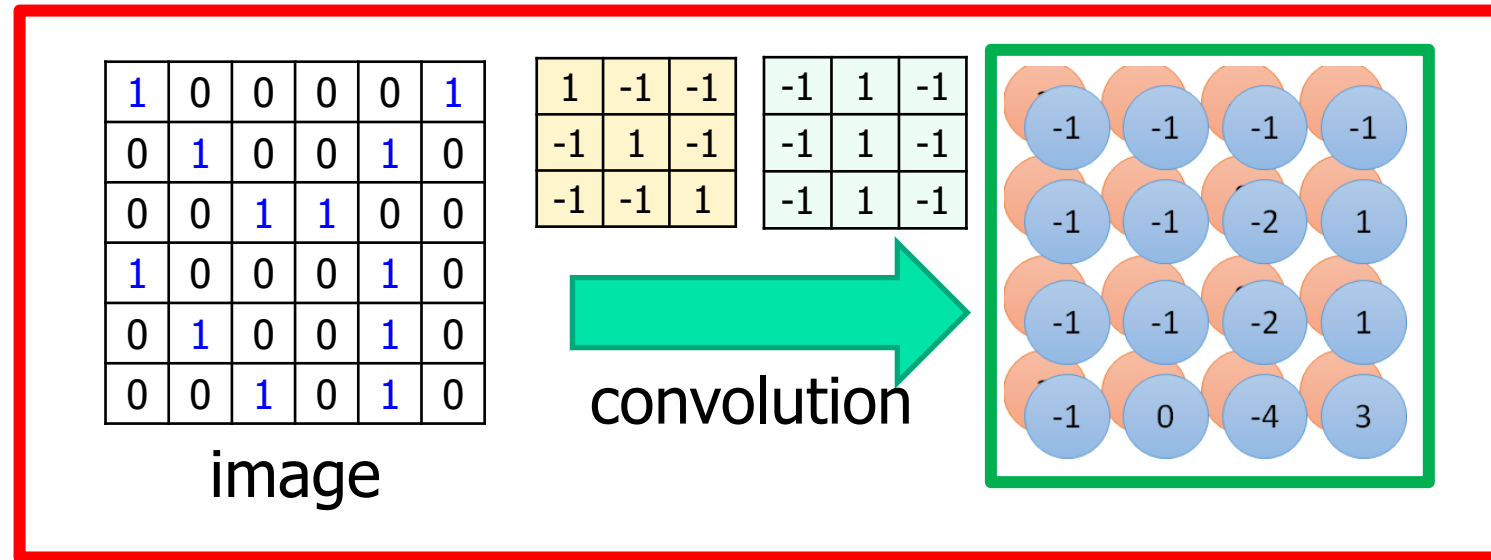


4 x 4 image

CNN – Colorful image

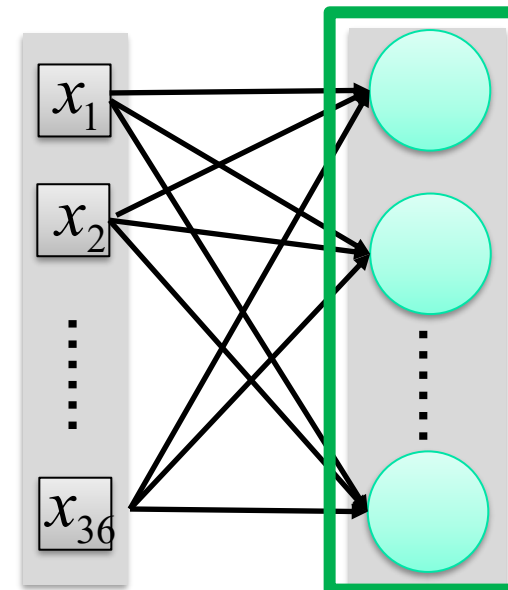


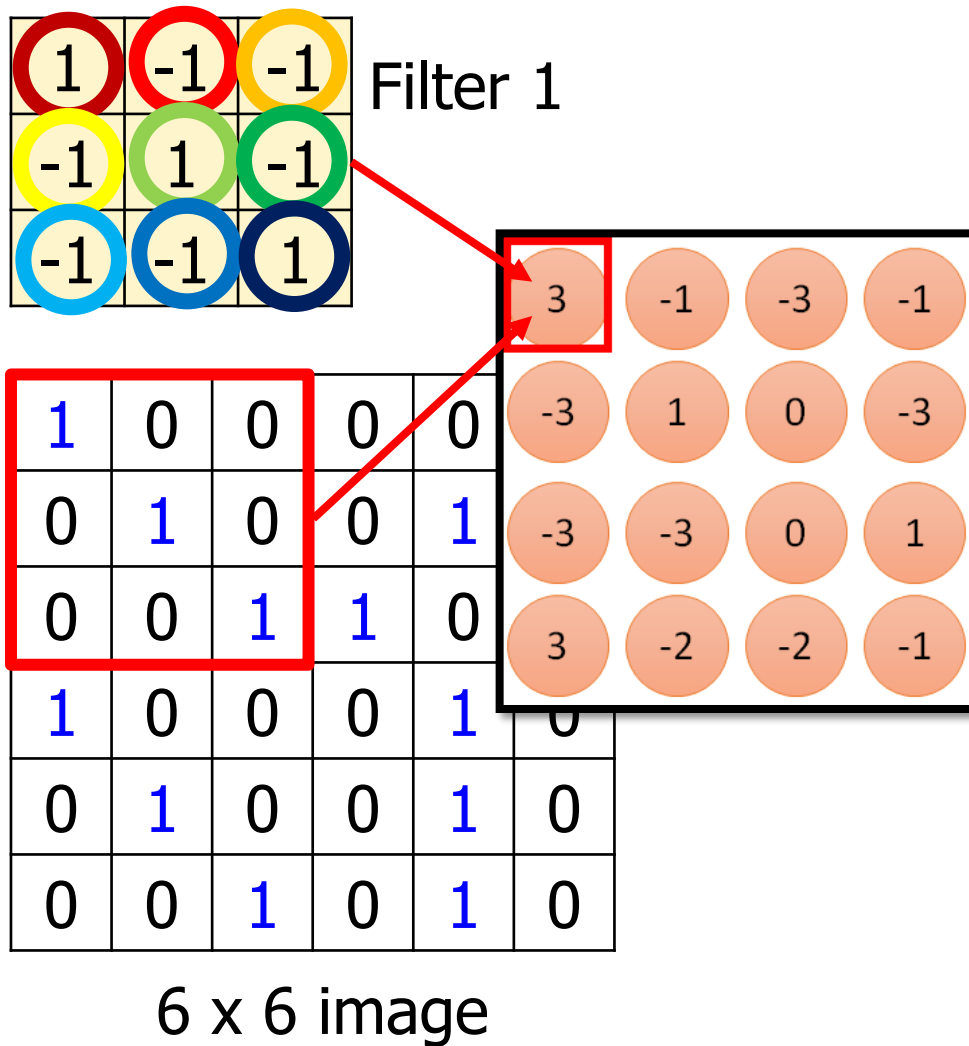
Convolution v.s. Fully Connected



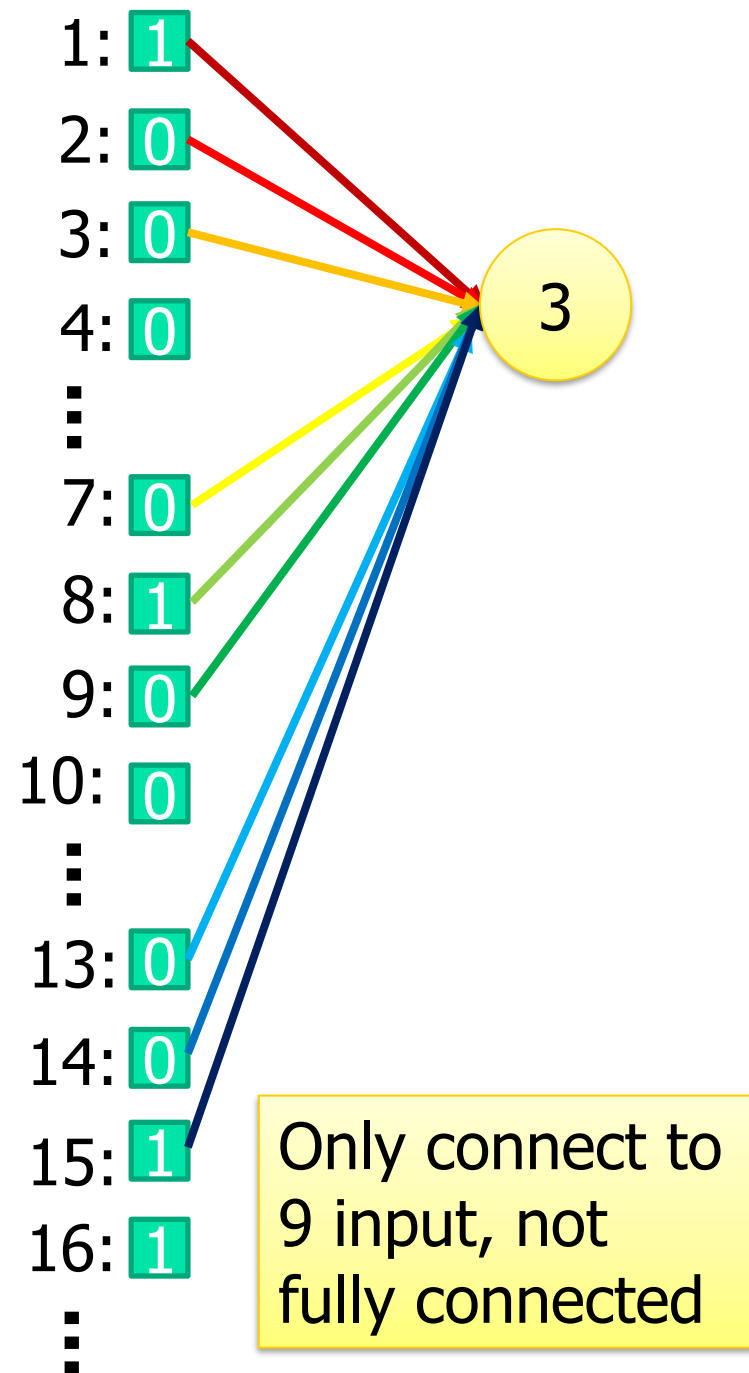
Fully-
connected

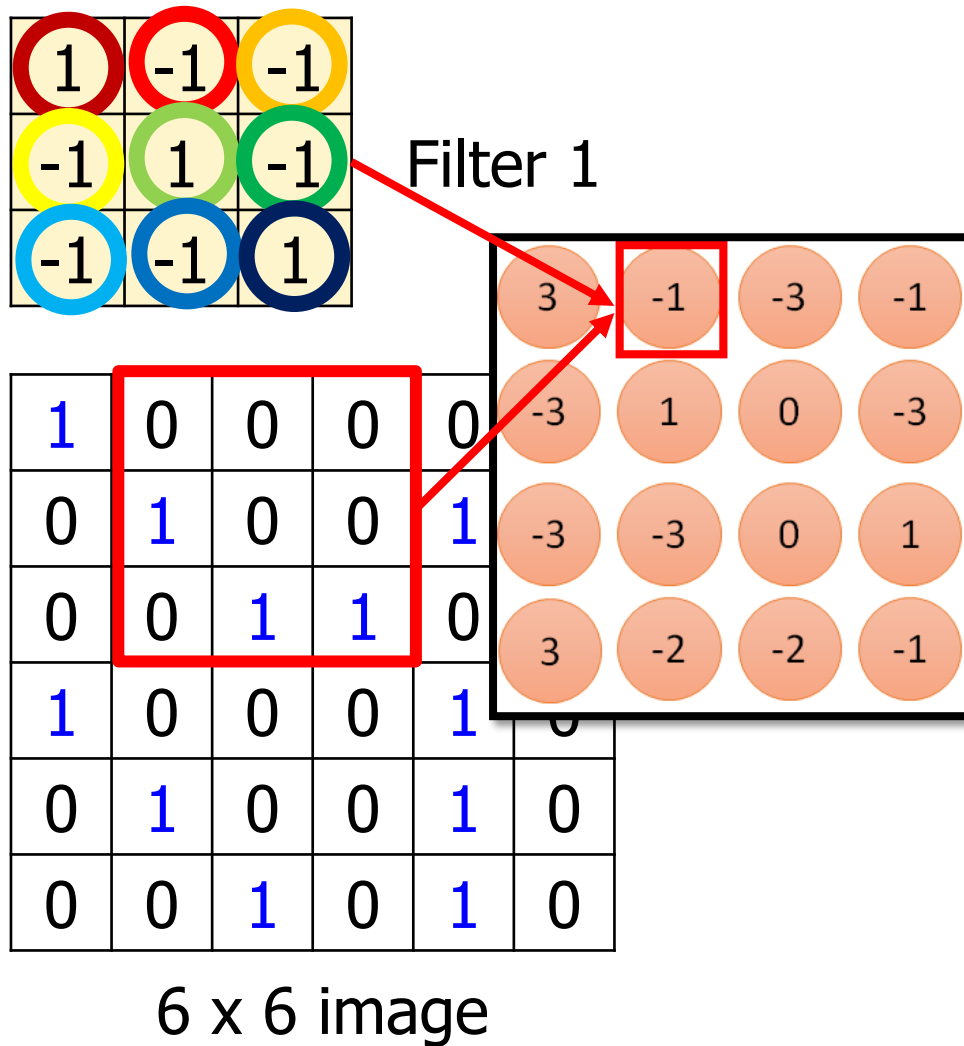
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0





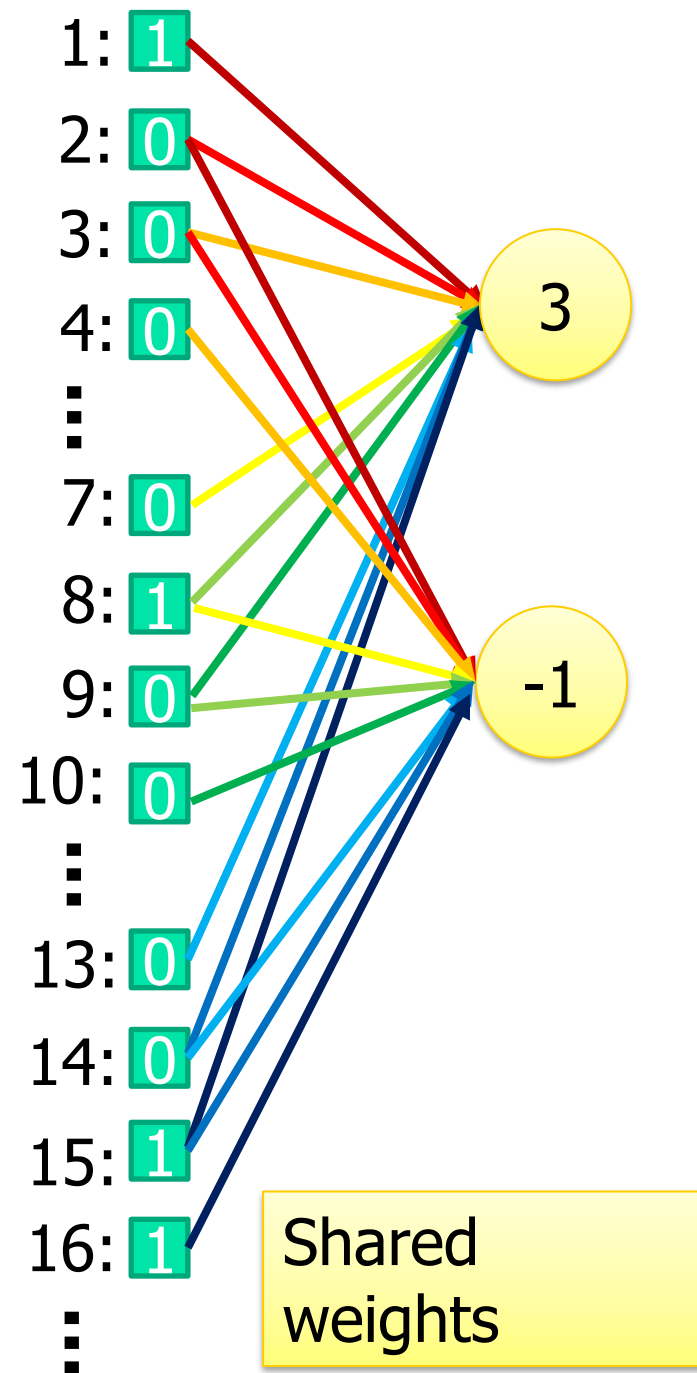
Less parameters!





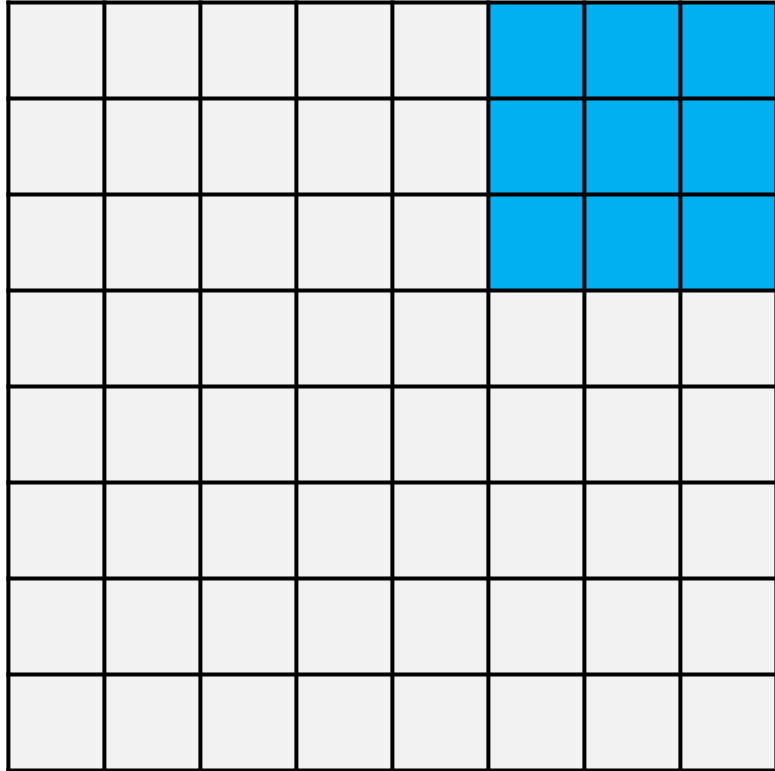
Less parameters!

Even less parameters!

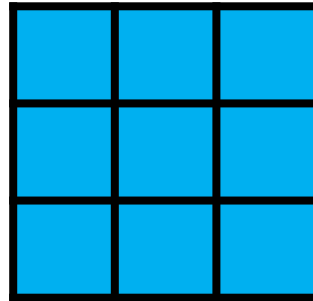


Convolution: 填充 (padding)

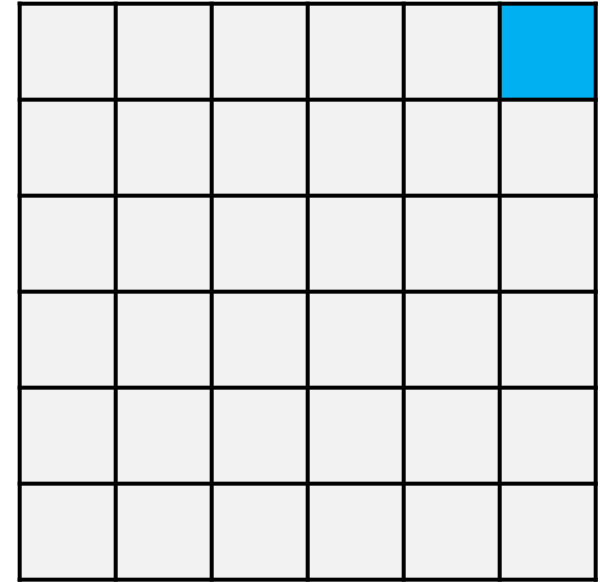
input: $H \times W = 8 \times 8$



filter



output: $H \times W = 6 \times 6$



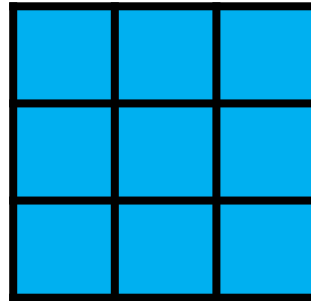
$$H_{\text{out}} = H_{\text{in}} - K_h + 1$$

Convolution: 填充 (padding)

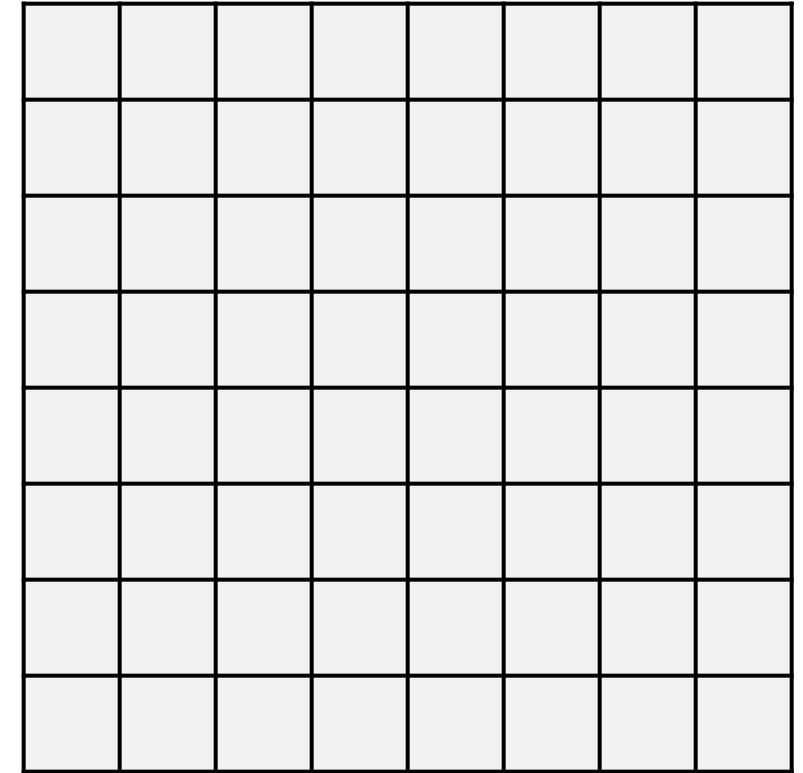
input: 8×8 , + pad

0	0	0	0	0	0	0	0	0	0
0									0
0									0
0									0
0									0
0									0
0									0
0									0
0									0
0									0
0	0	0	0	0	0	0	0	0	0

filter



output: $H \times W = 8 \times 8$

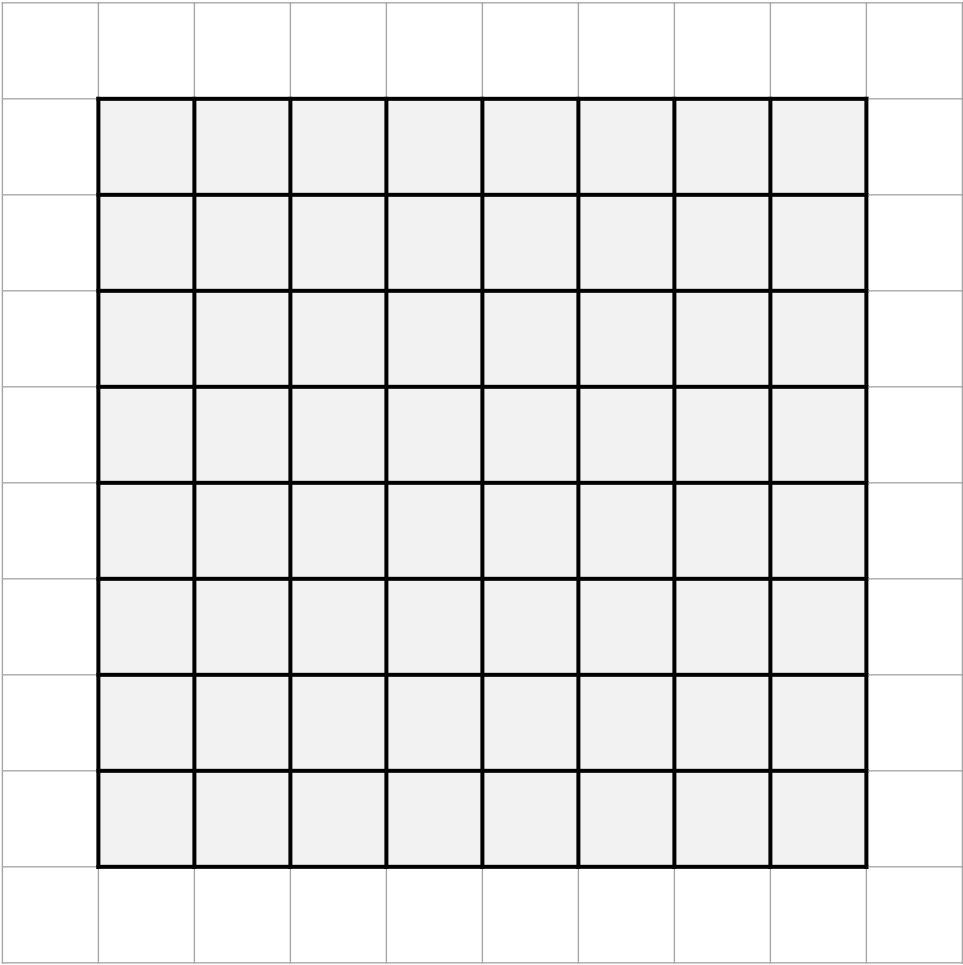


$$H_{\text{out}} = H_{\text{in}} + 2\text{pad}_h - K_h + 1$$

- $\text{pad} = \lfloor \text{kernel_size} / 2 \rfloor$
- maintains feature map size

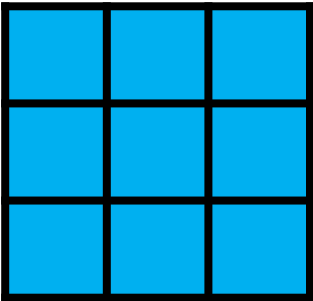
Convolution: 步长 (stride)

input

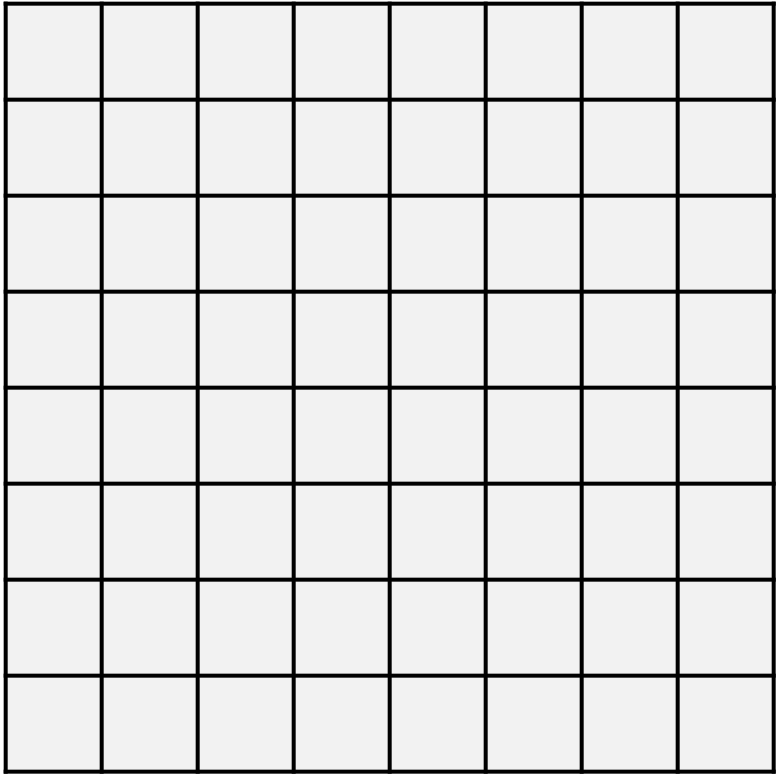


stride = 2

filter

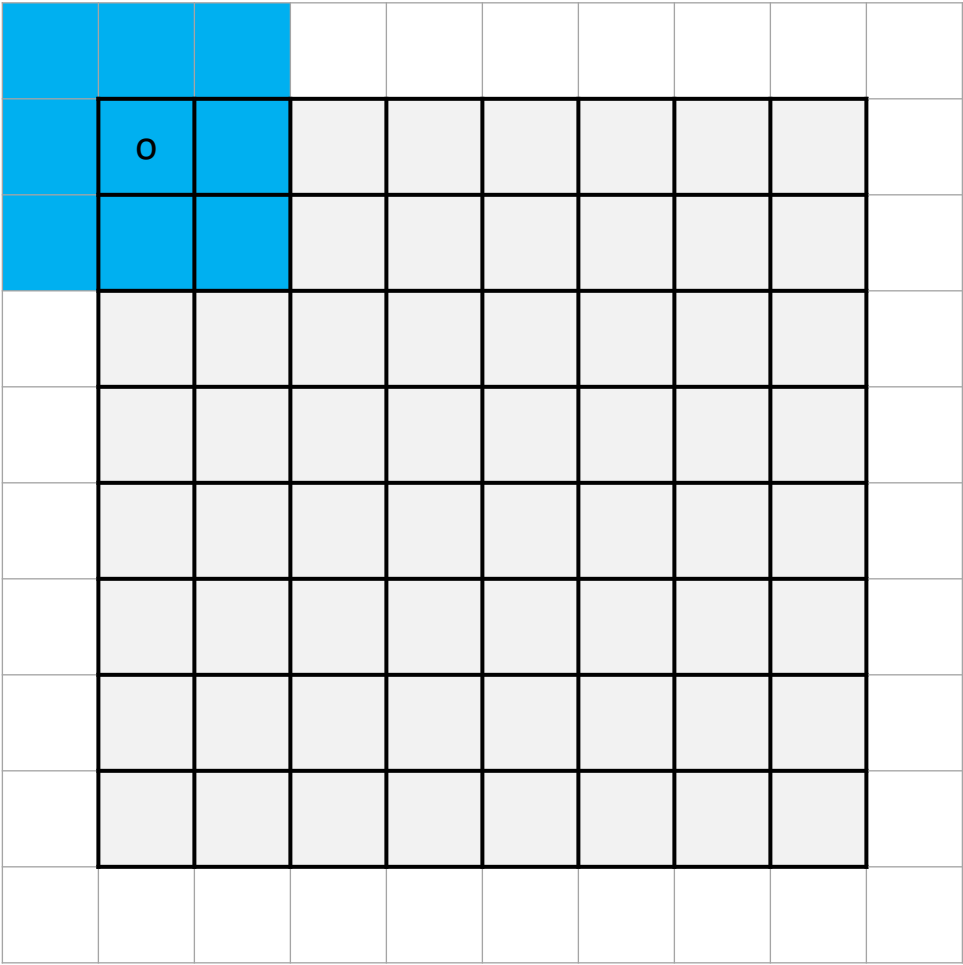


output



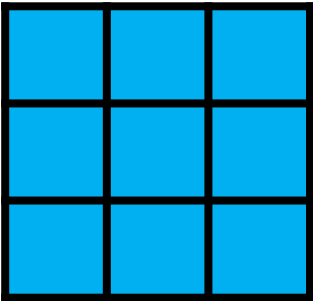
Convolution: 步长 (stride)

input

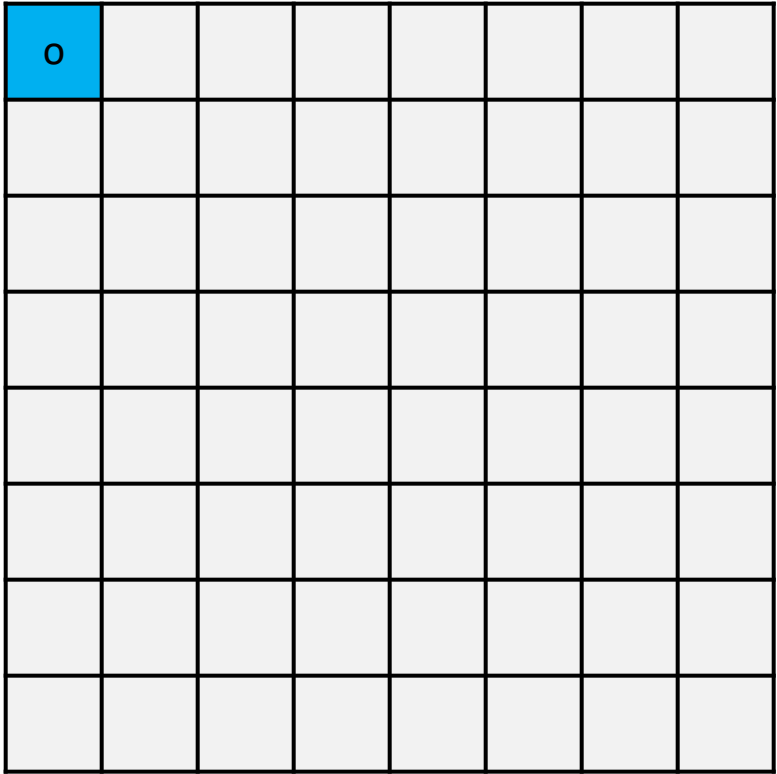


stride = 2

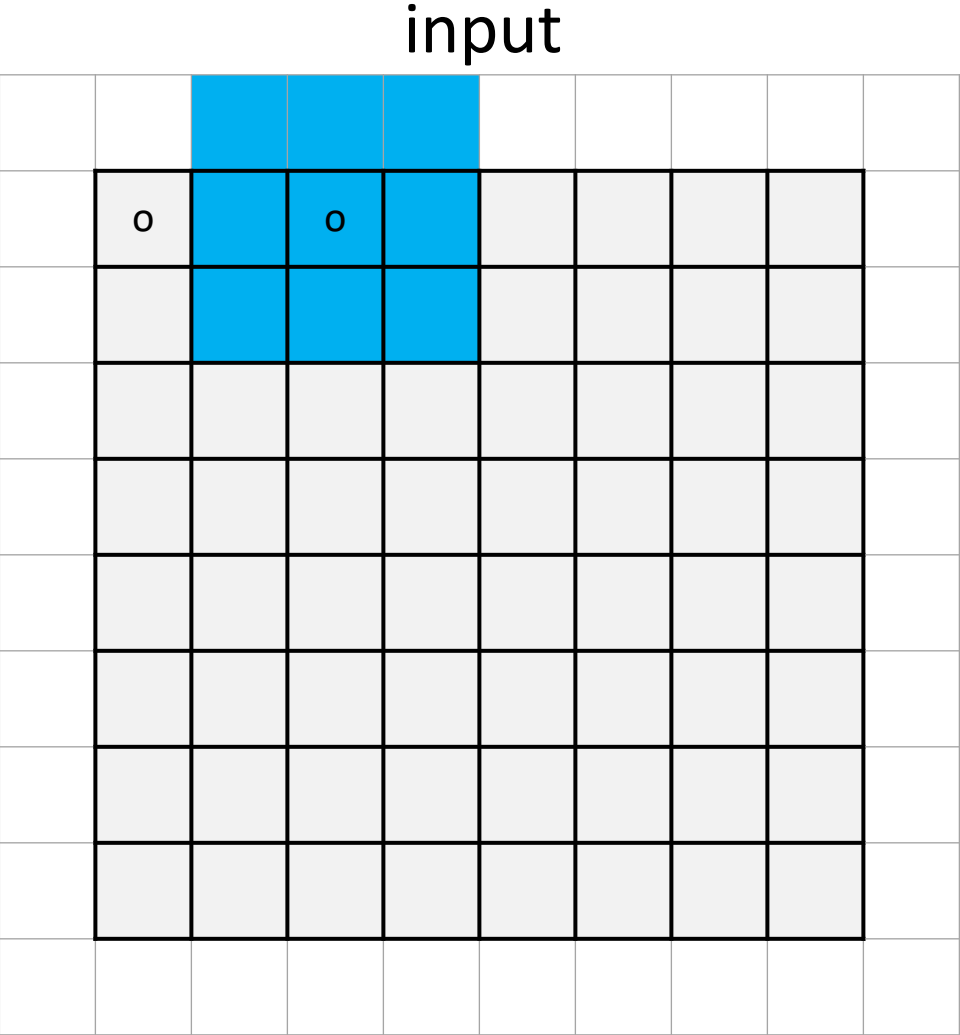
filter



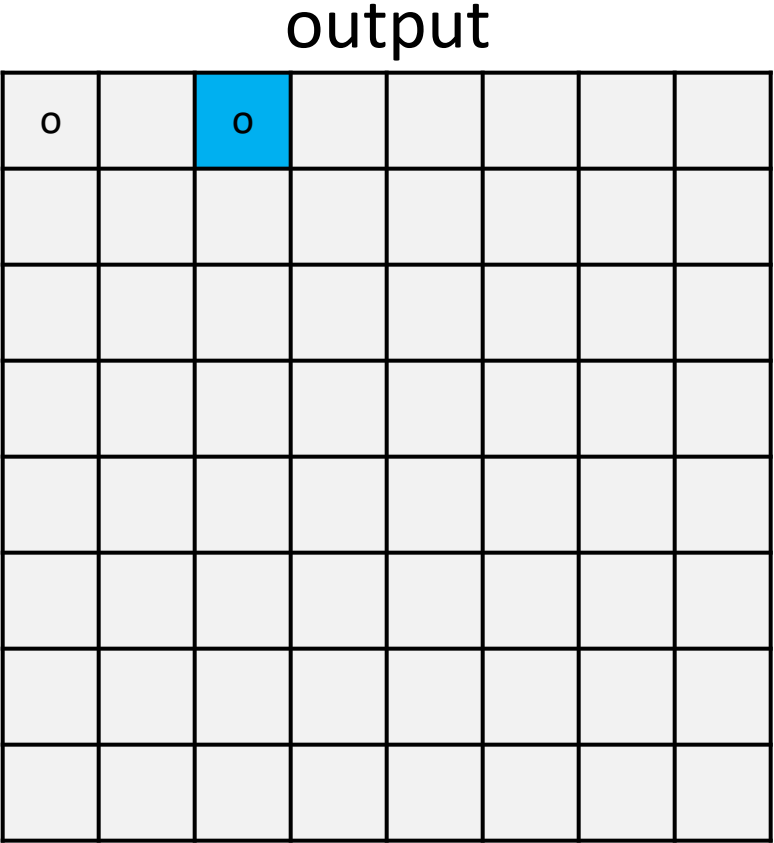
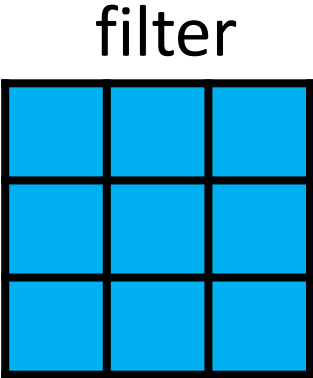
output



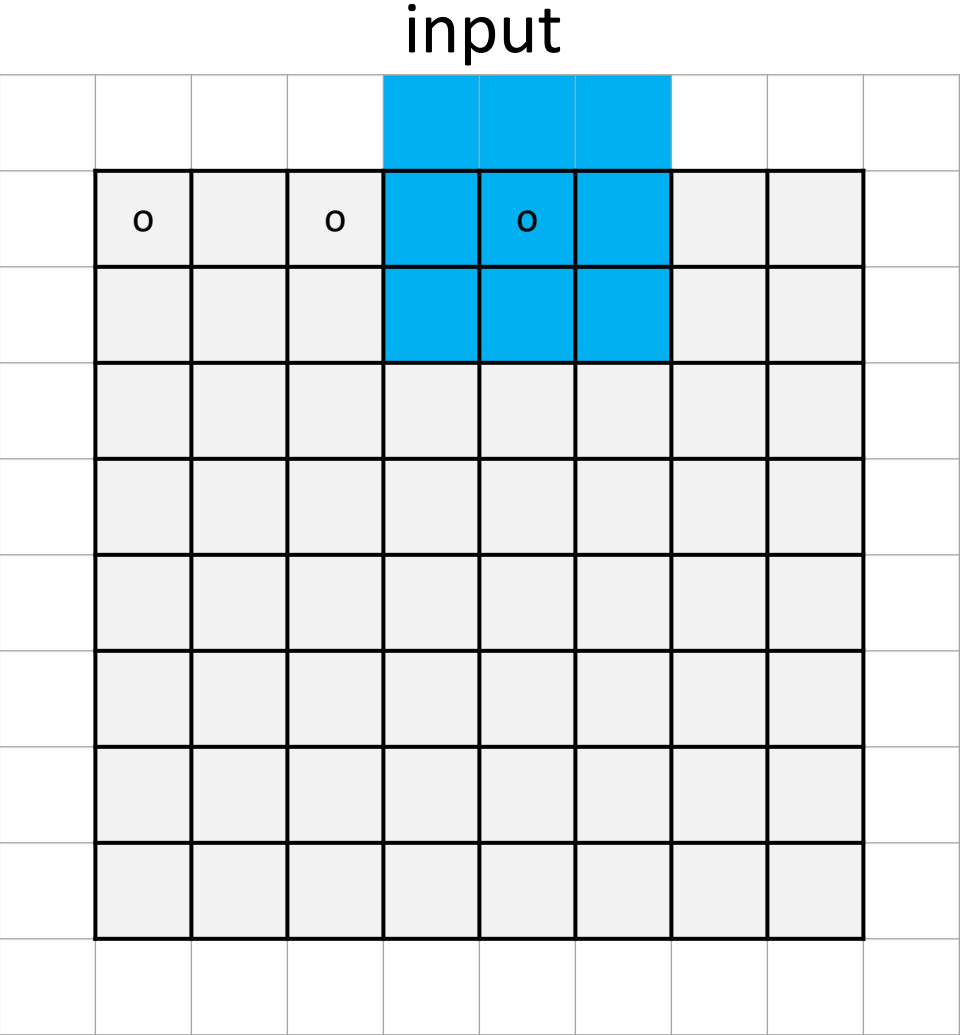
Convolution: 步长 (stride)



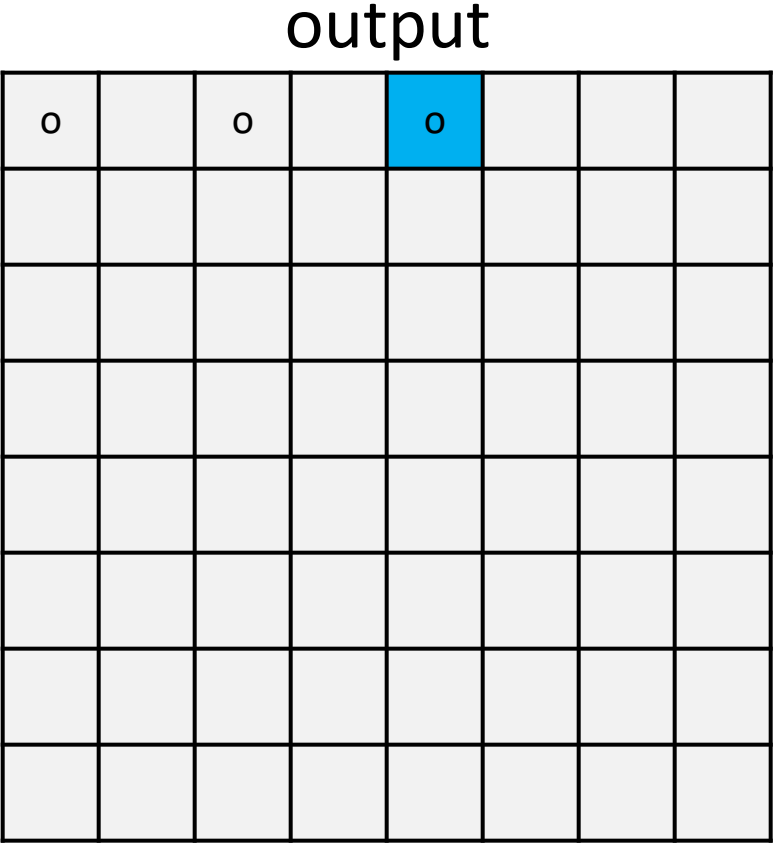
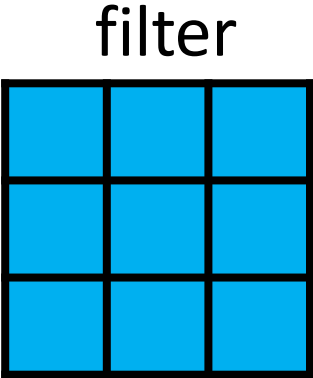
stride = 2



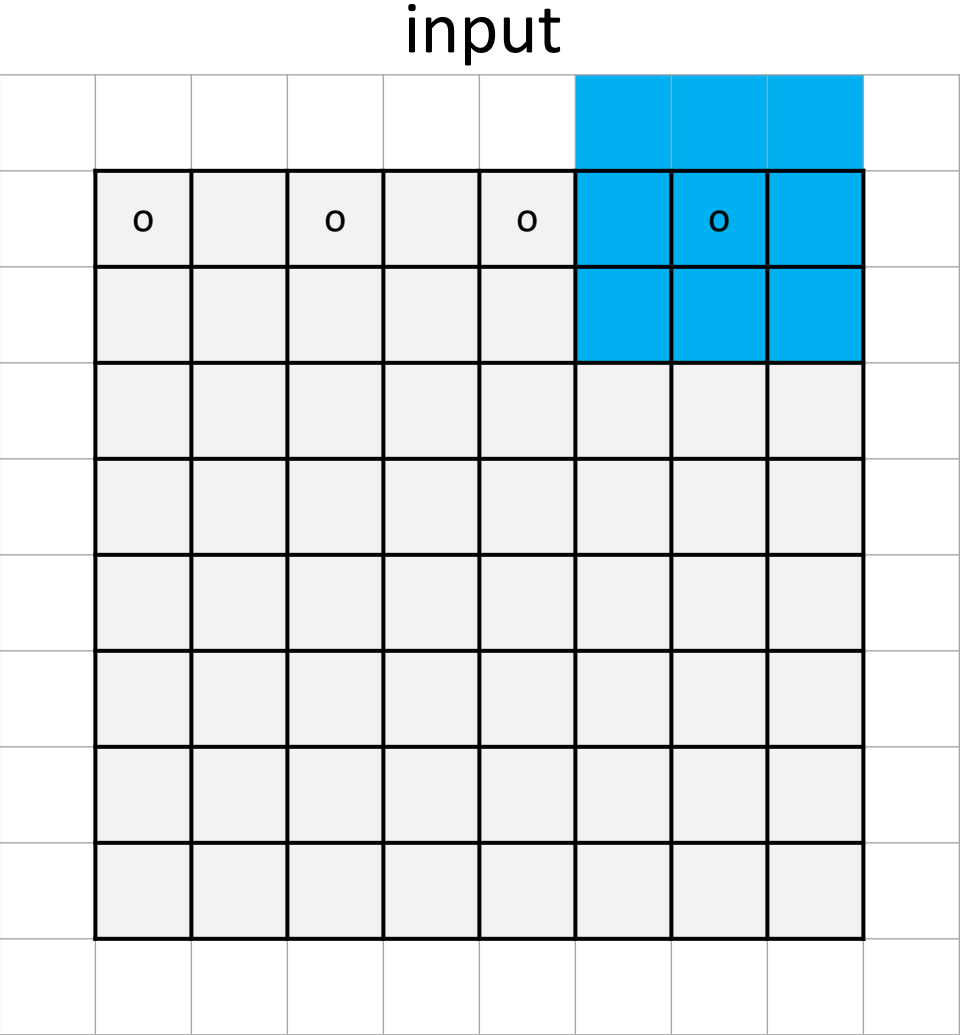
Convolution: 步长 (stride)



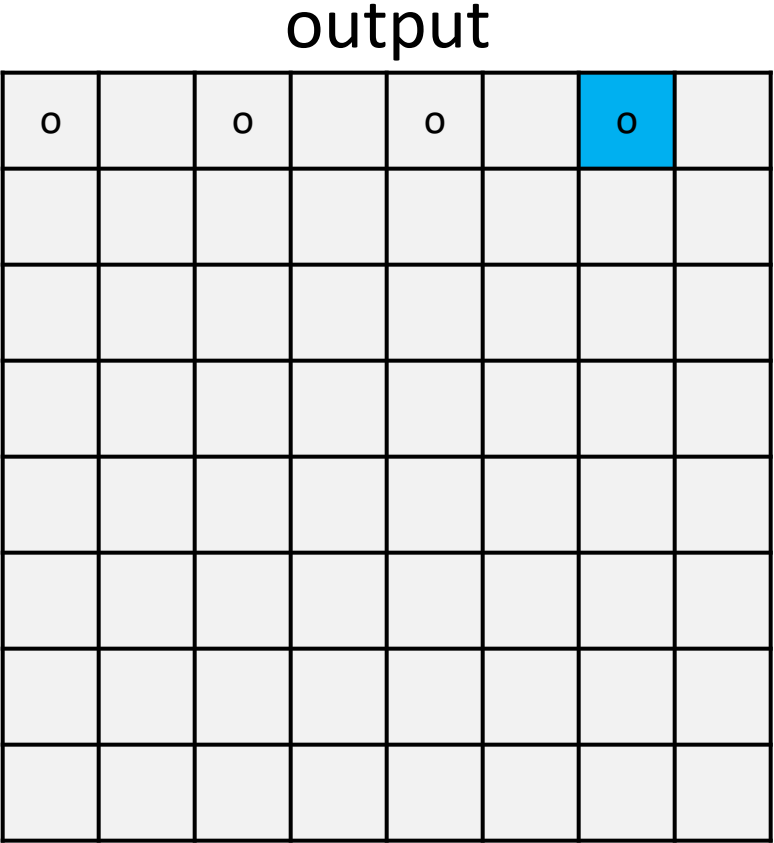
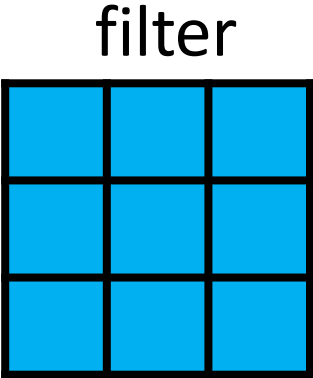
stride = 2



Convolution: 步长 (stride)

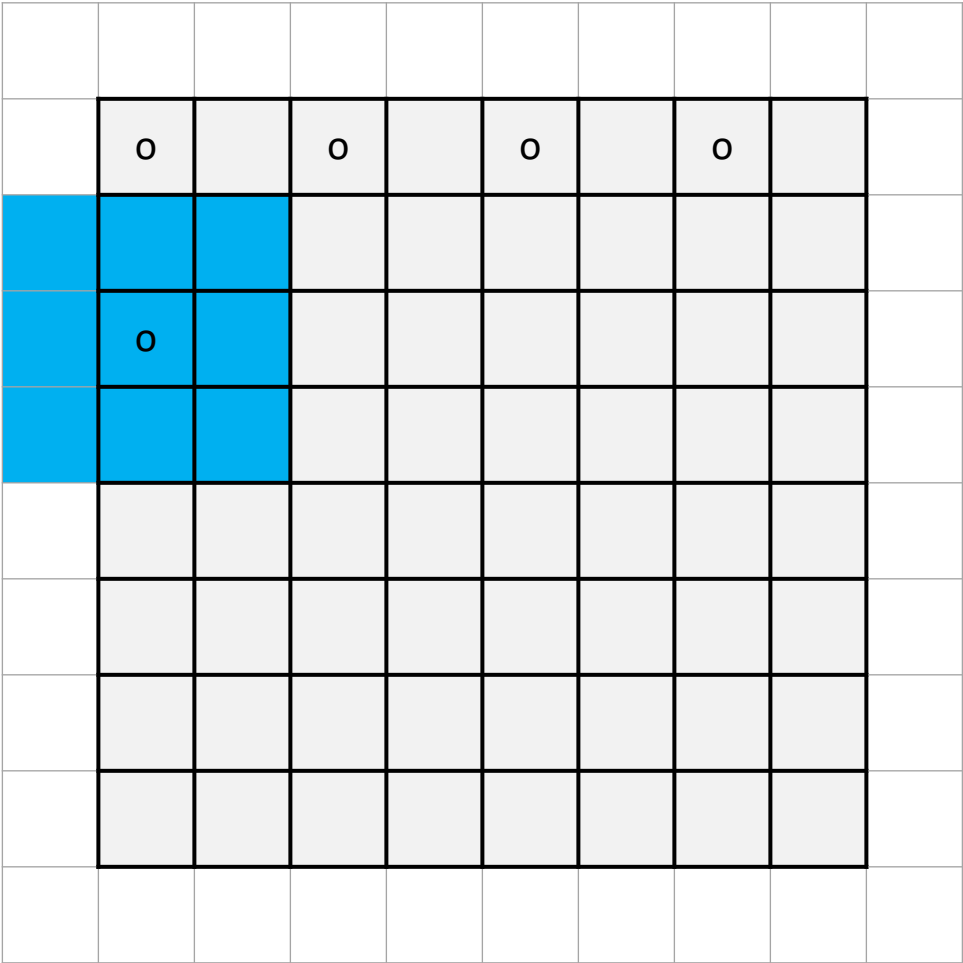


stride = 2



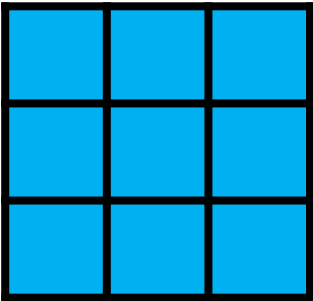
Convolution: 步长 (stride)

input

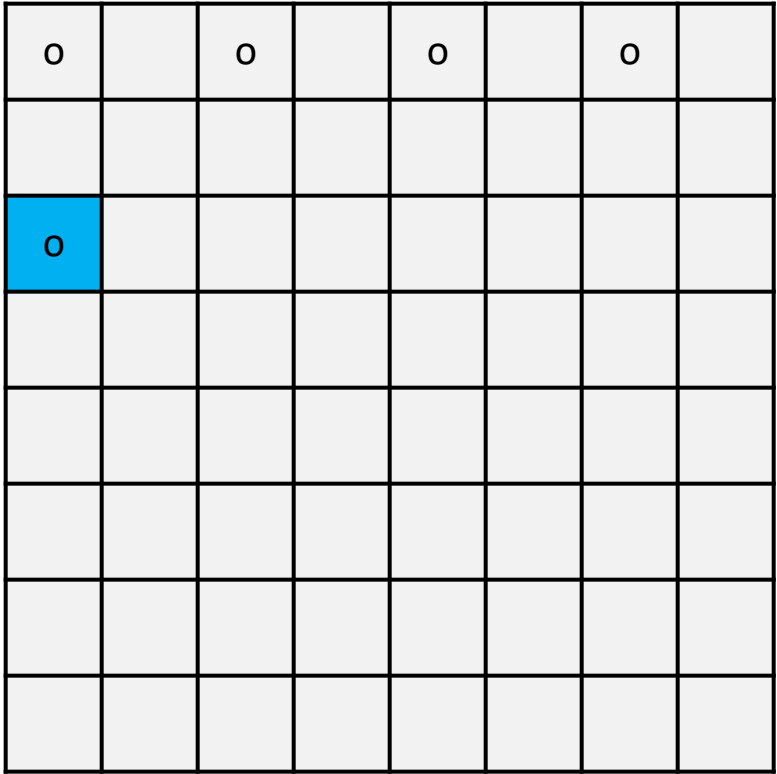


stride = 2

filter



output



Convolution: 步长 (stride)

input

	o		o		o		o		
	o		o		o		o		
	o		o		o		o		
	o		o		o		o		
	o		o		o		o		

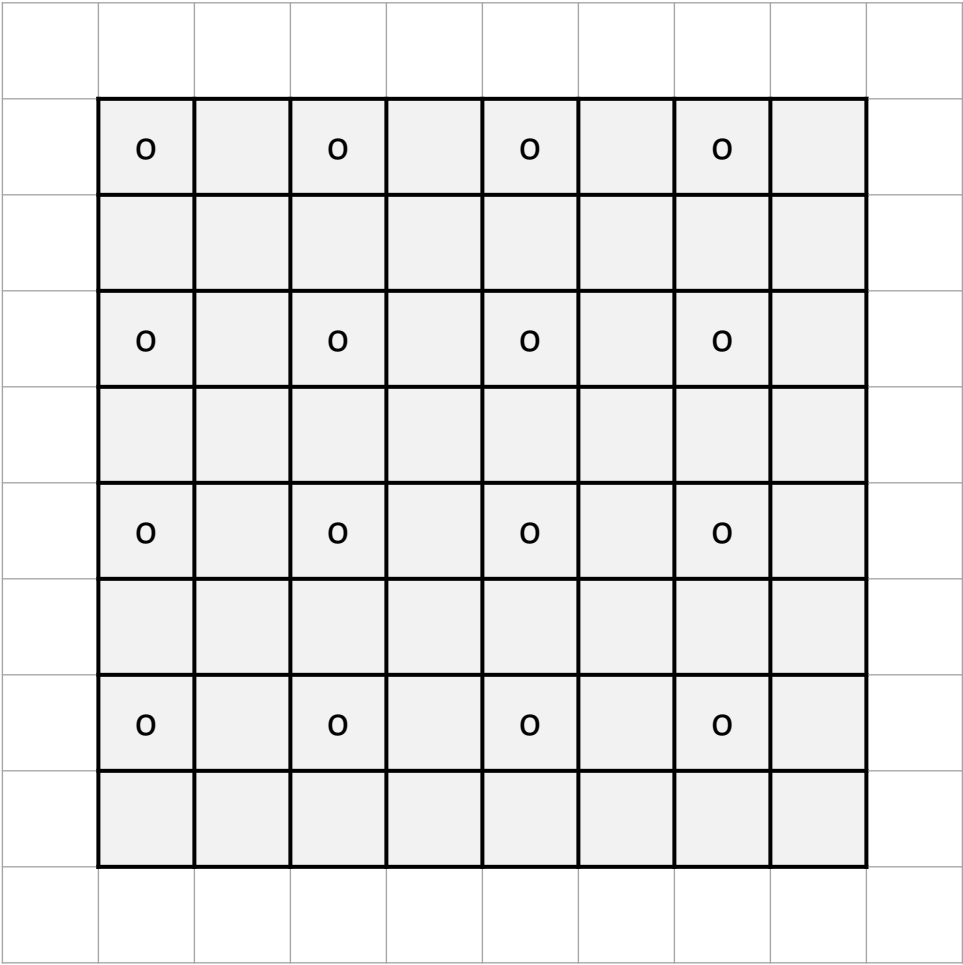
filter

output

o		o		o		o	
o		o		o		o	
o		o		o		o	
o		o		o		o	
o		o		o		o	

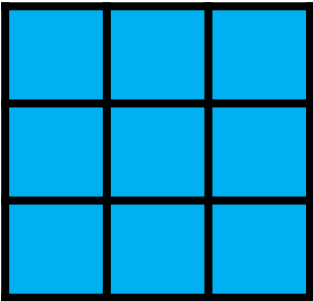
Convolution: 步长 (stride)

input

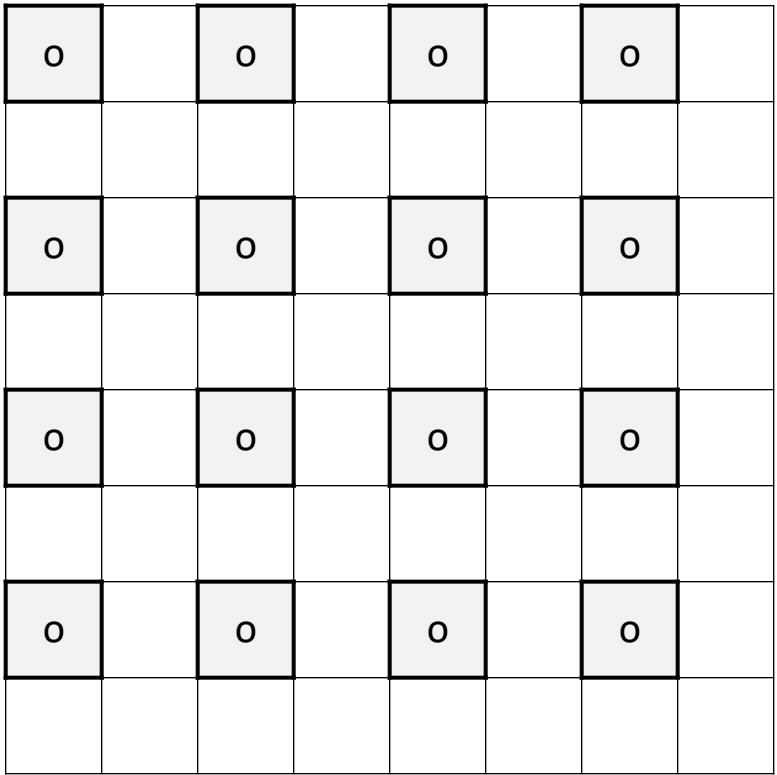


stride = 2

filter



output: $H \times W = 4 \times 4$



Convolution: 步长 (stride)

input

	o		o		o		o		
	o		o		o		o		
	o		o		o		o		
	o		o		o		o		

stride = 2

filter

- reduces feature map size
- compress and abstract

output: $H \times W = 4 \times 4$

o	o	o	o
o	o	o	o
o	o	o	o
o	o	o	o

$$H_{\text{out}} = \lfloor (H_{\text{in}} + 2\text{pad}_h - K_h) / \text{str} \rfloor + 1$$

CNN – Max Pooling

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

Filter 1

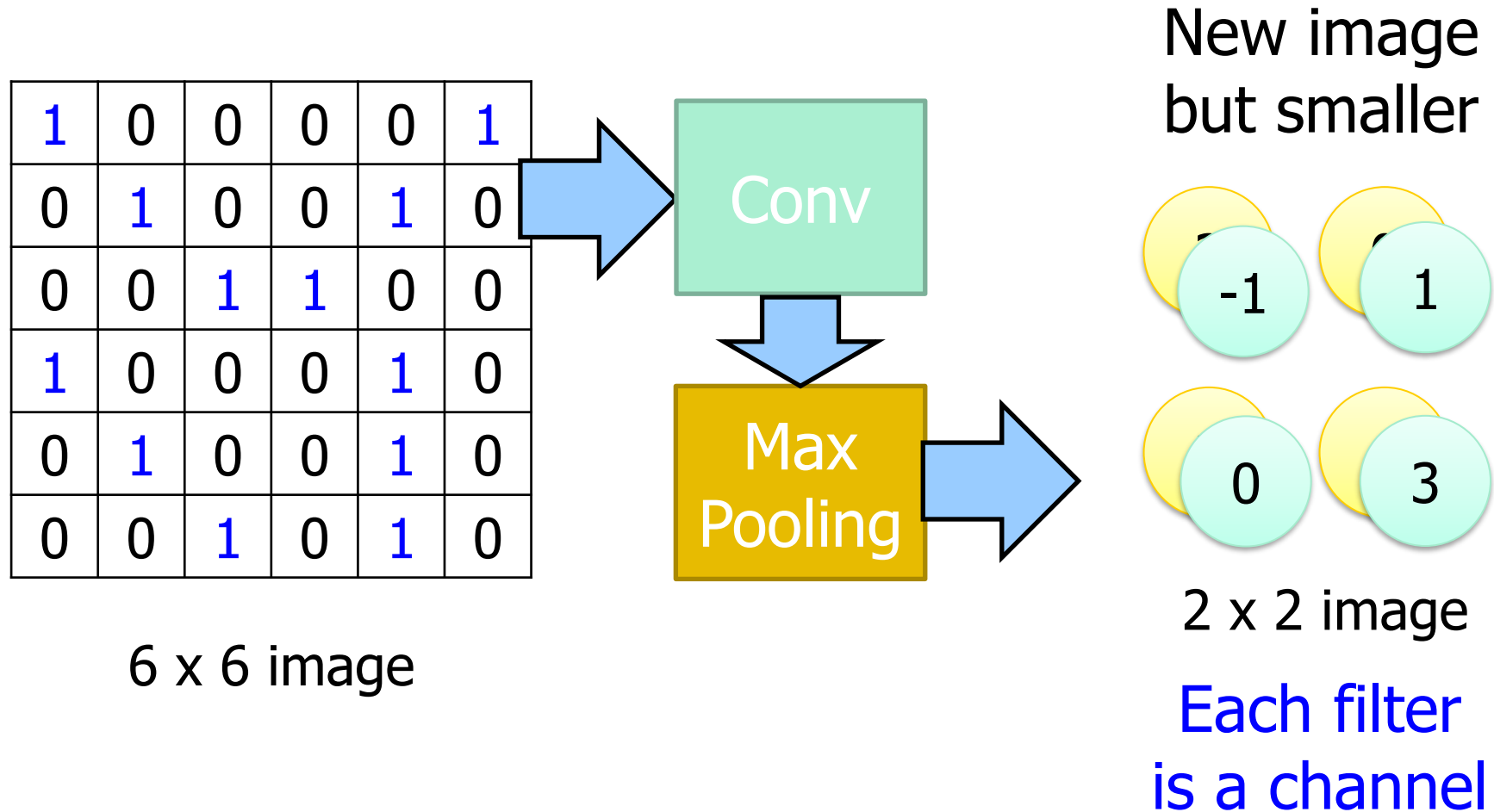
-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

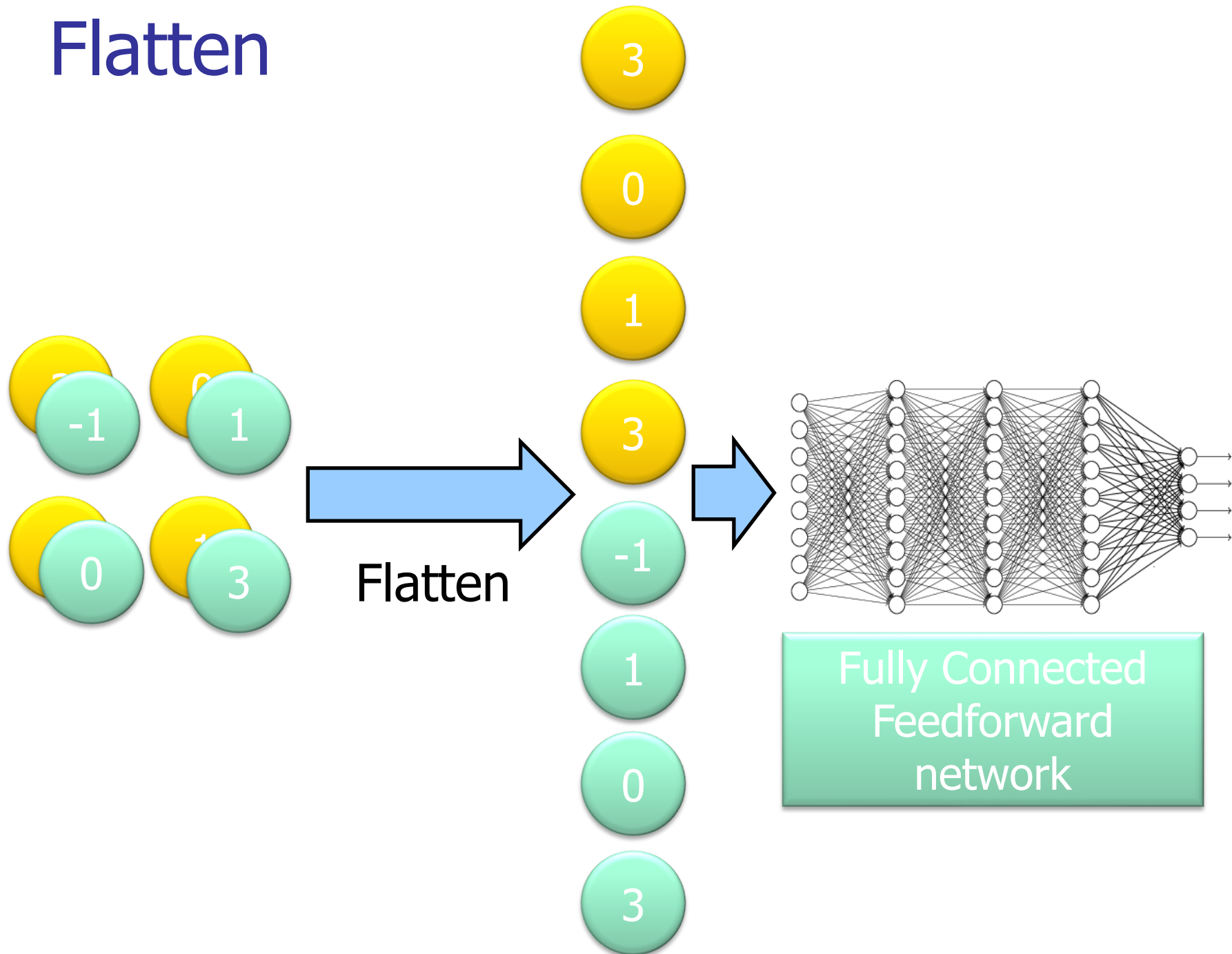
3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

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

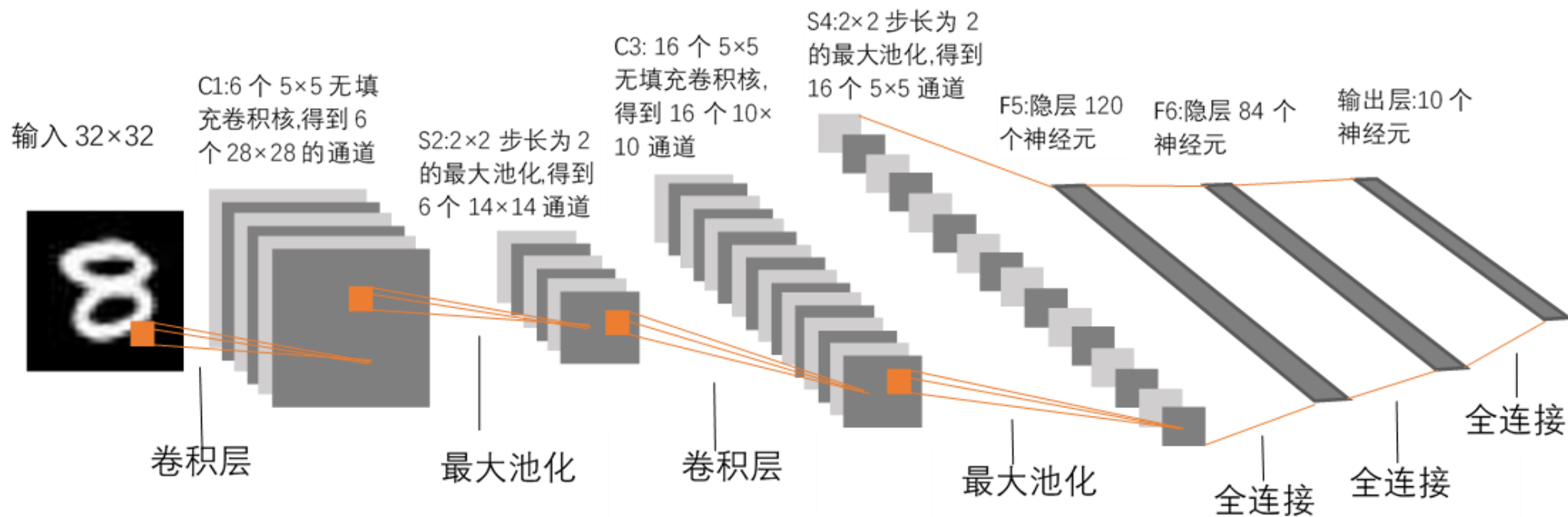
CNN – Max Pooling



Flatten

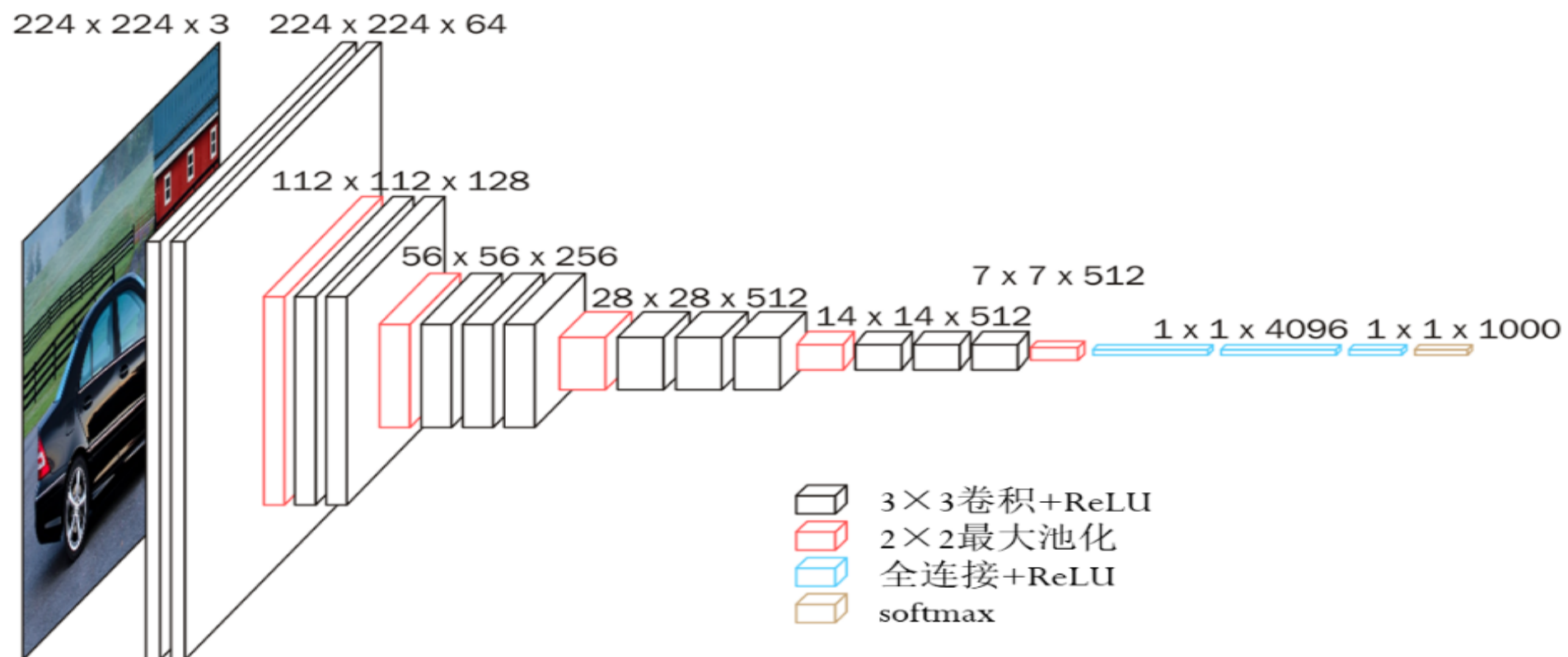


LeNet

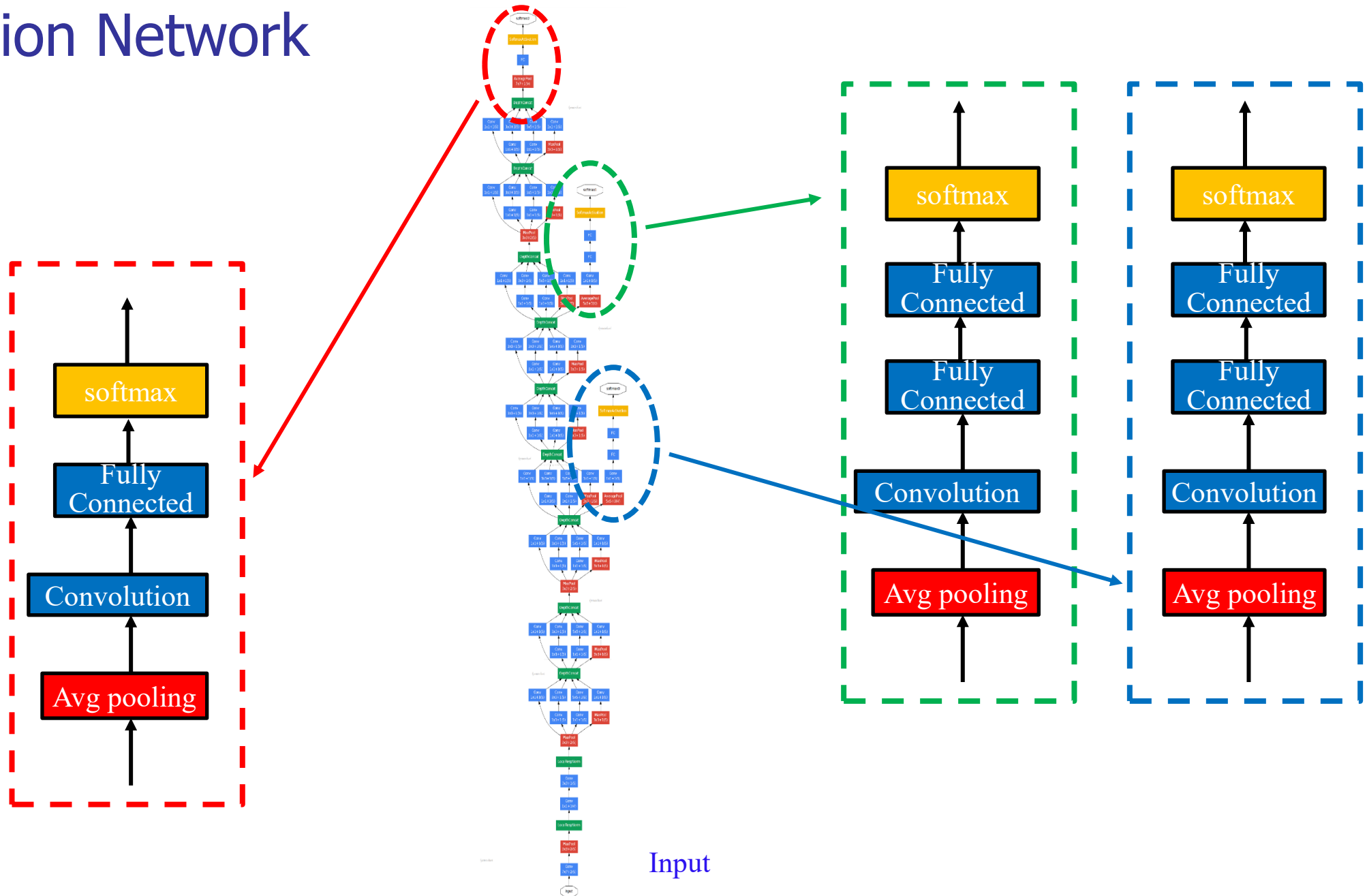


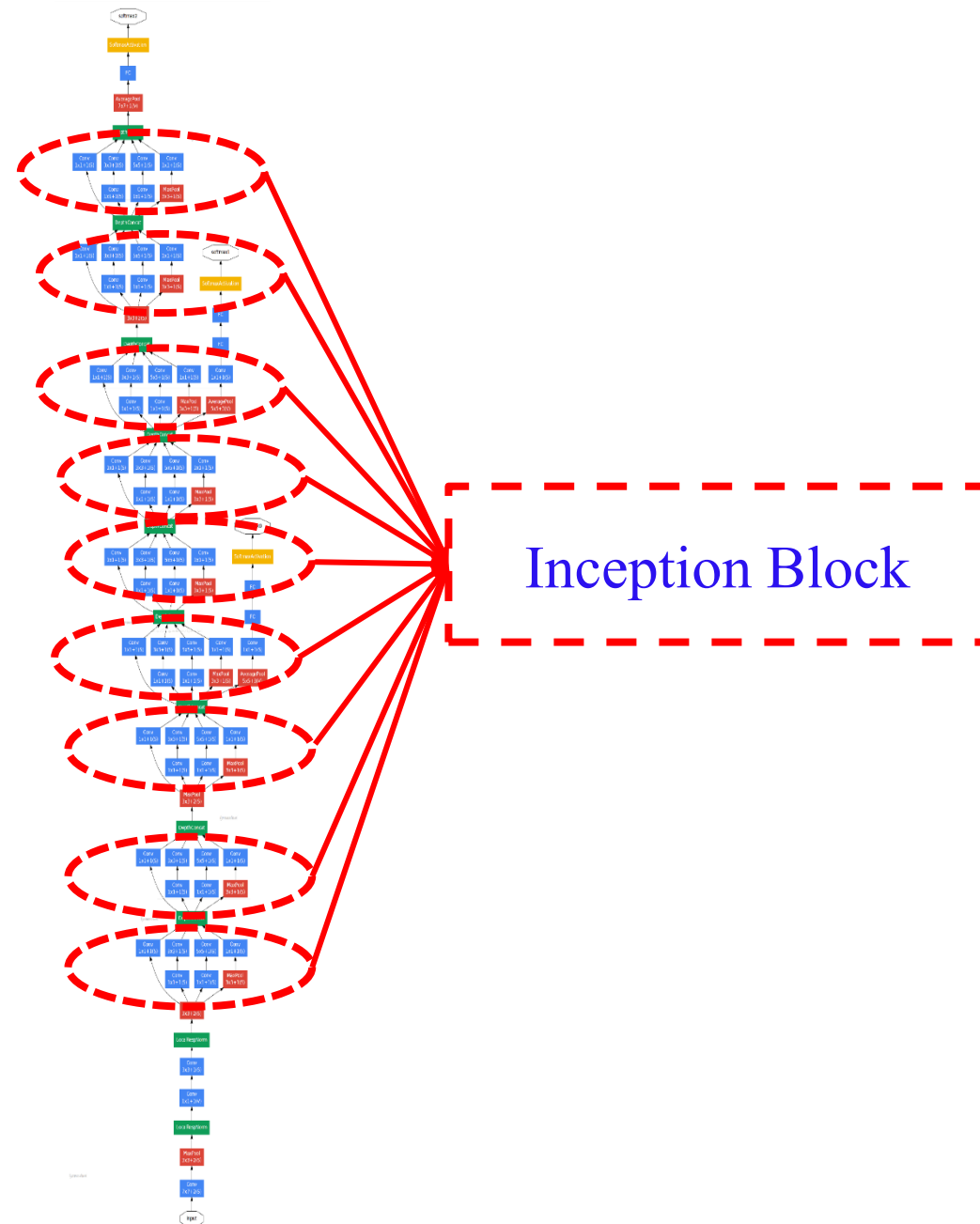
This slide and the following seven slides are from: <https://mp.weixin.qq.com/s/sdrIgGOo6DLTGZ3uth6I0>

VGG-16

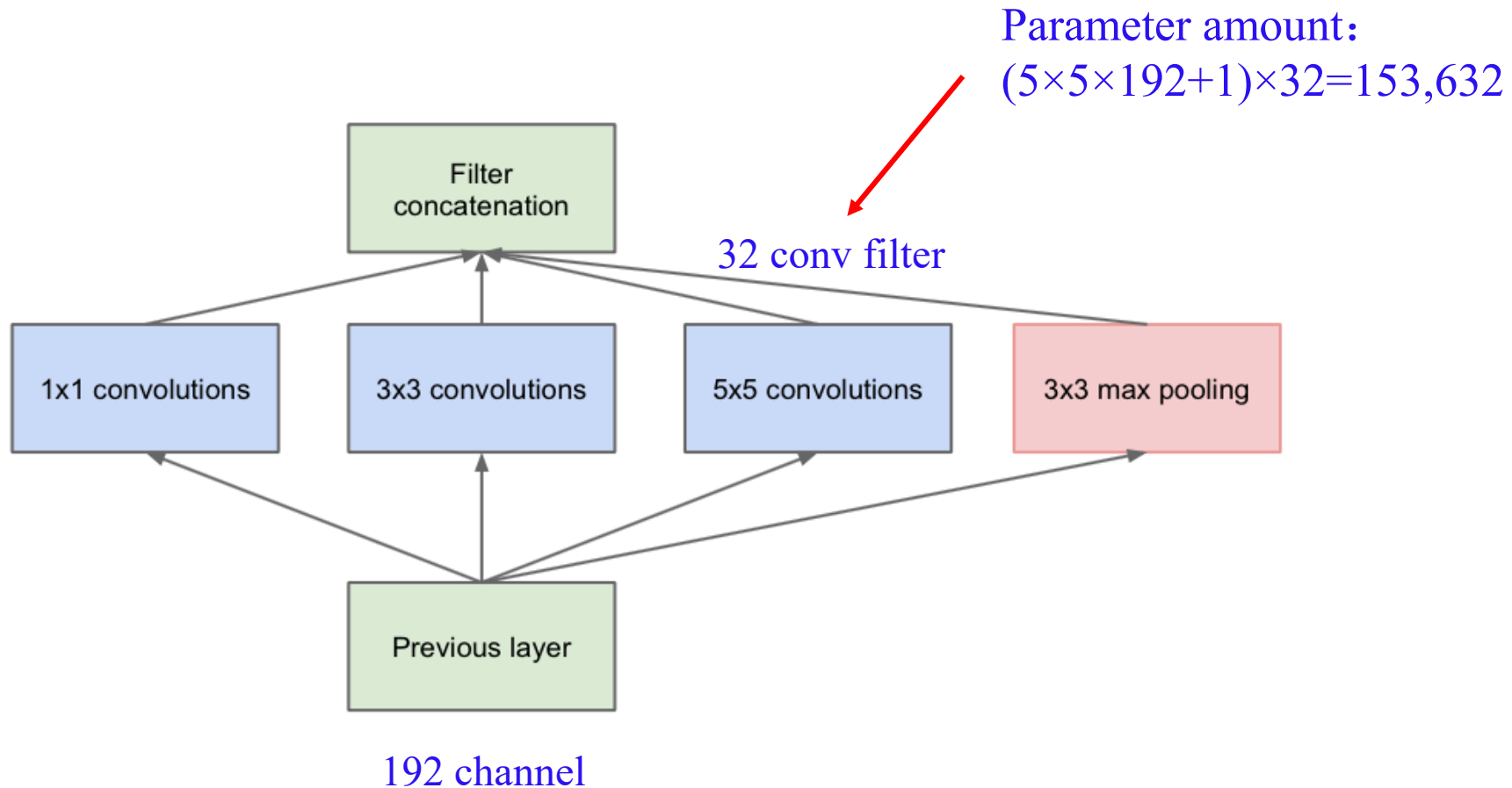


Inception Network

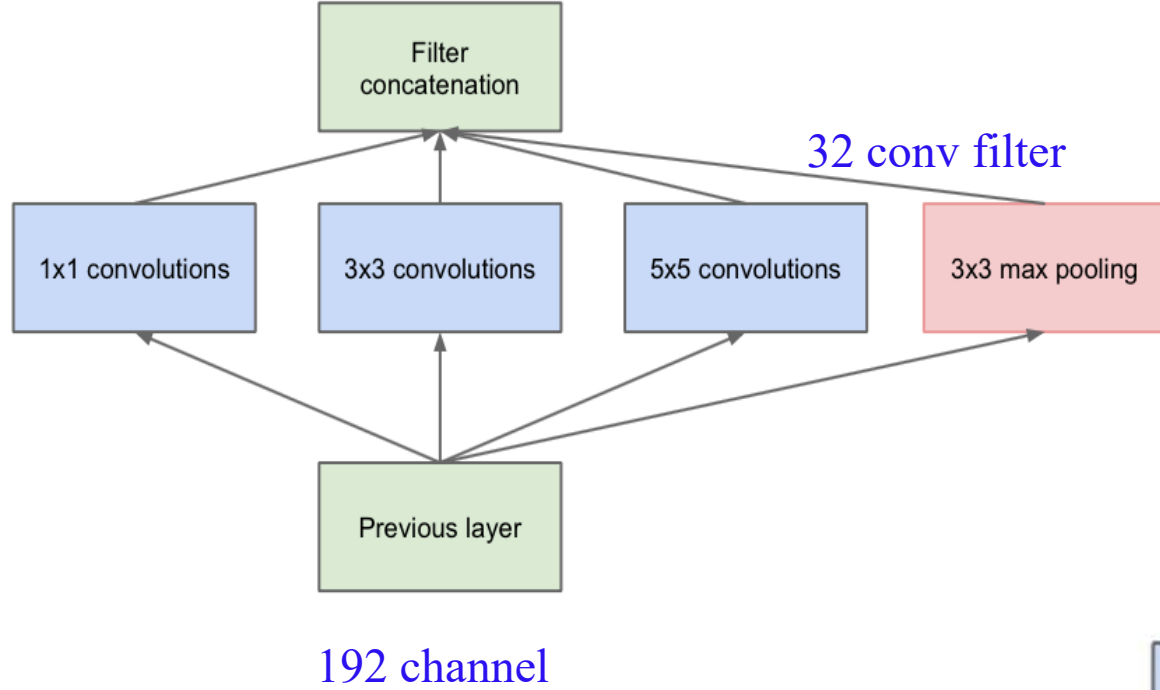




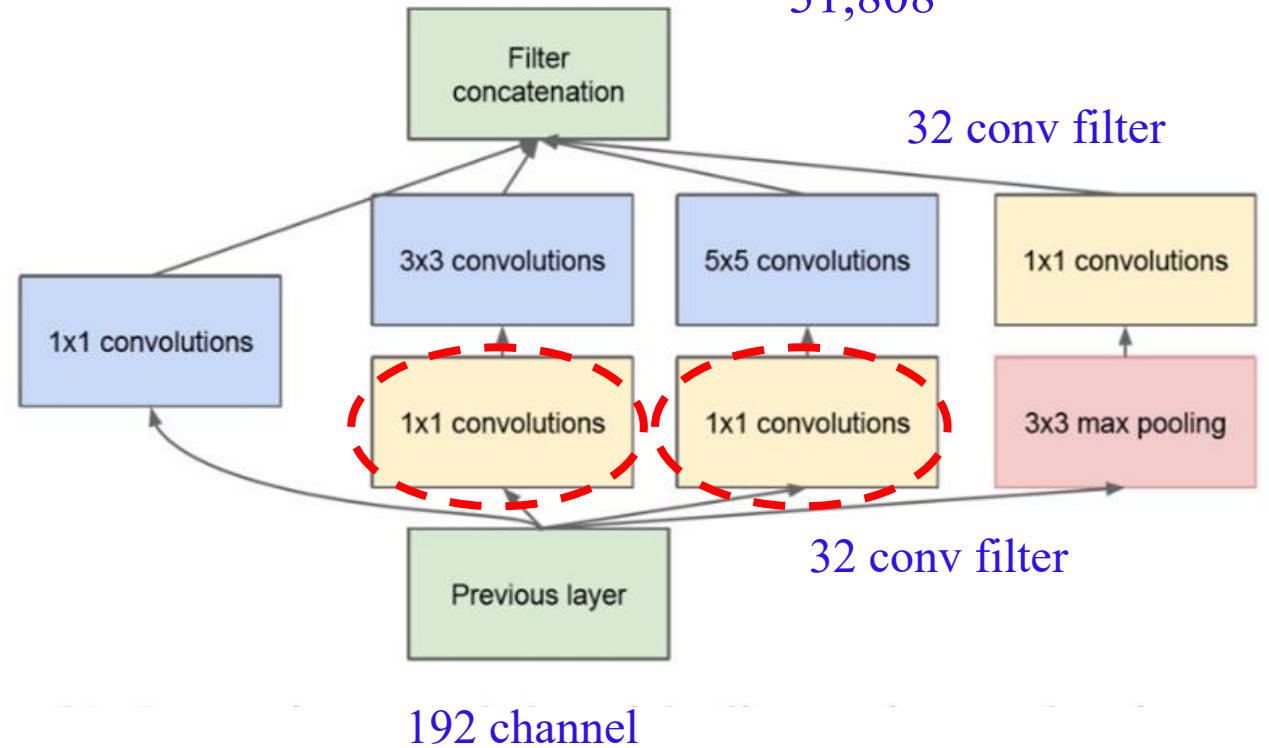
Inception Block

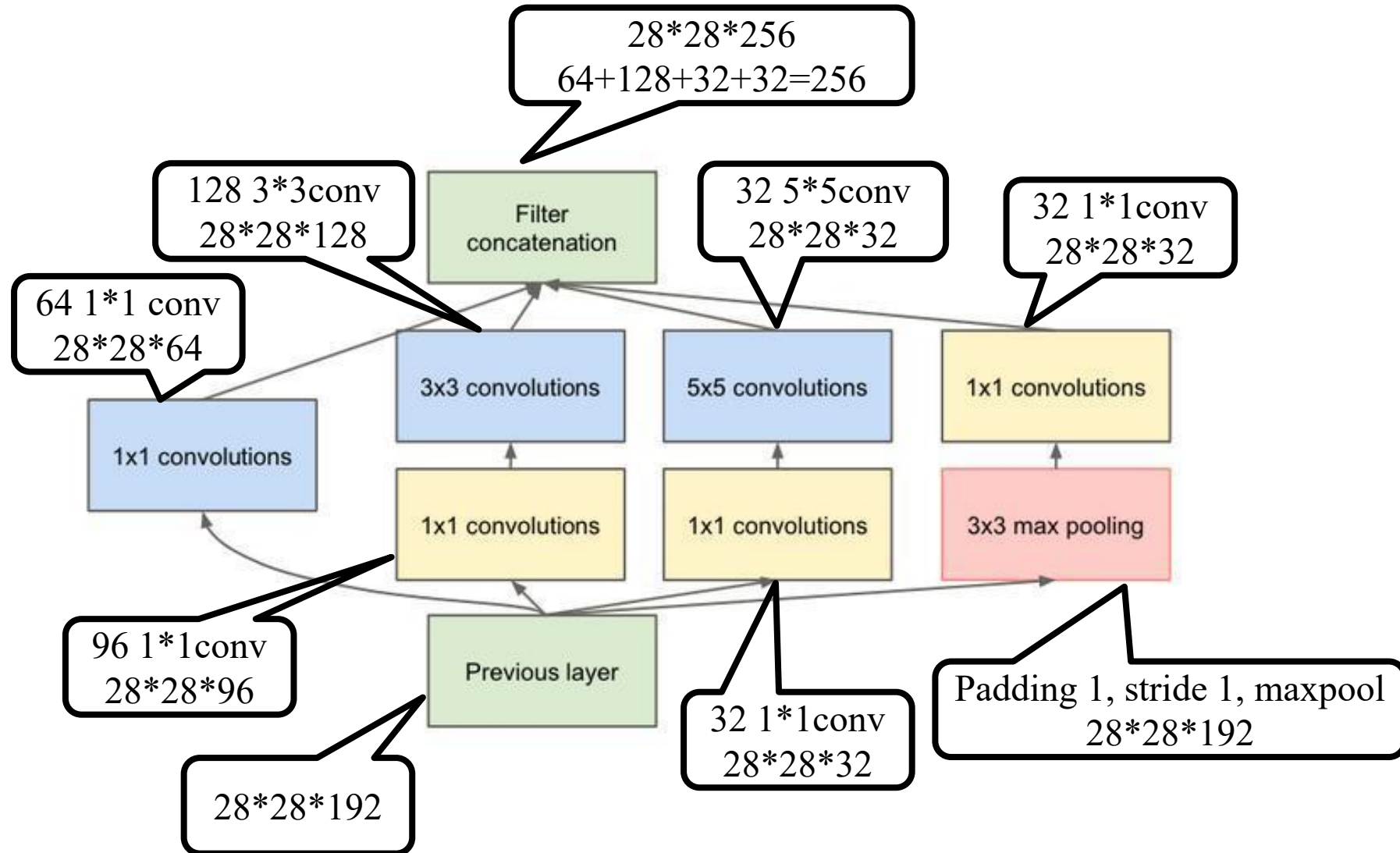


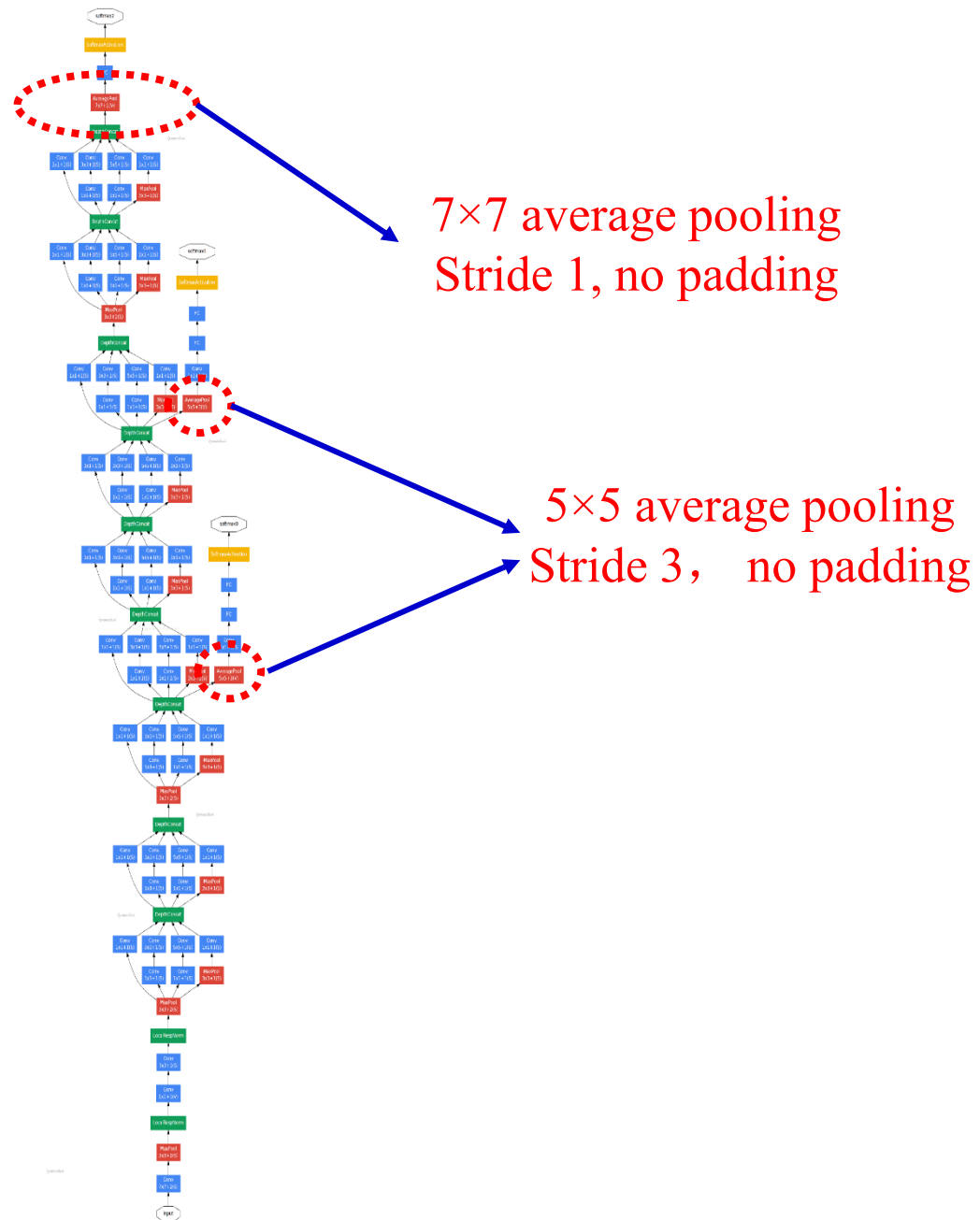
Parameters:
 $(5 \times 5 \times 192 + 1) \times 32 = 153,632$



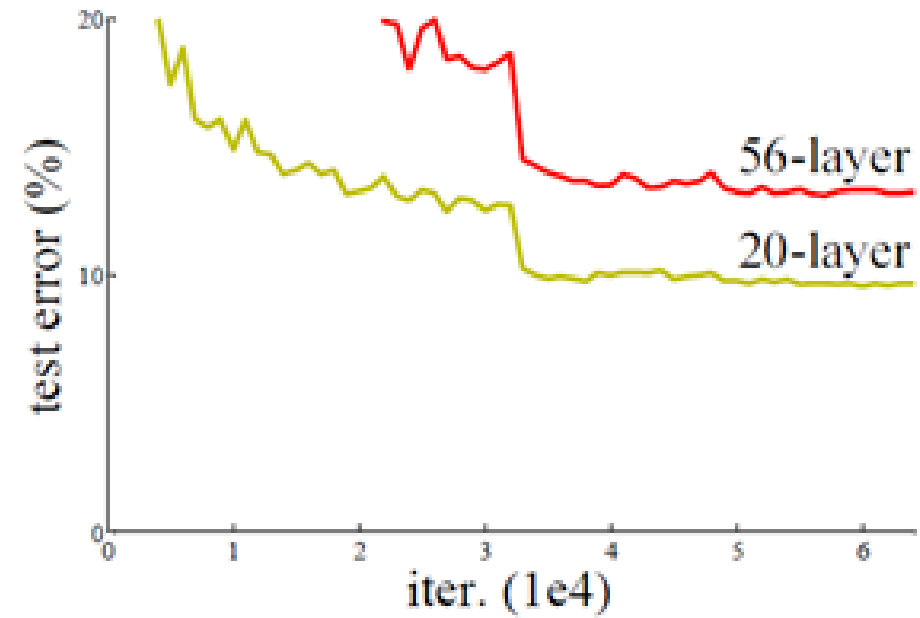
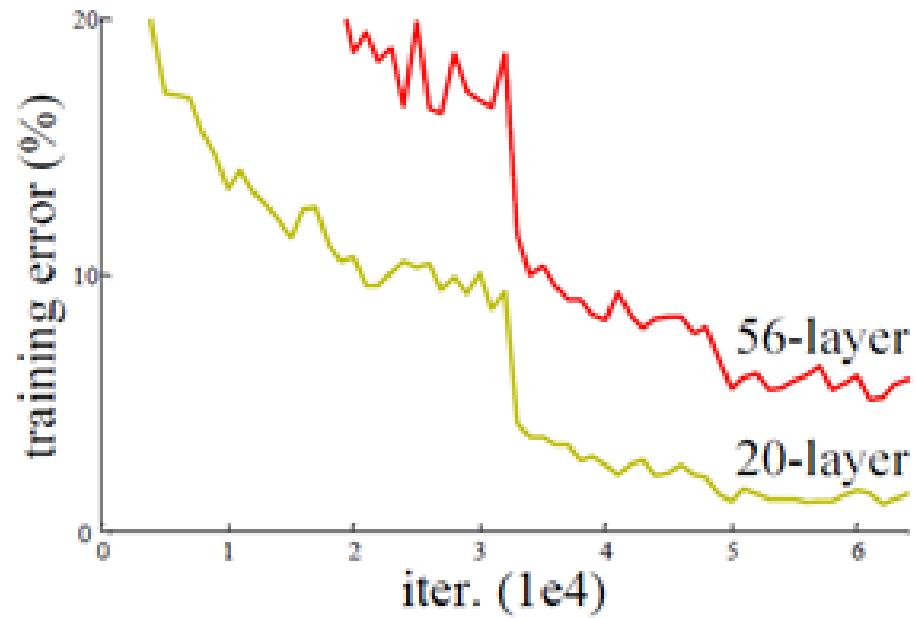
Parameters:
 $(1 \times 1 \times 192 + 1) \times 32 +$
 $(5 \times 5 \times 32 + 1) \times 32$
 $= 31,808$





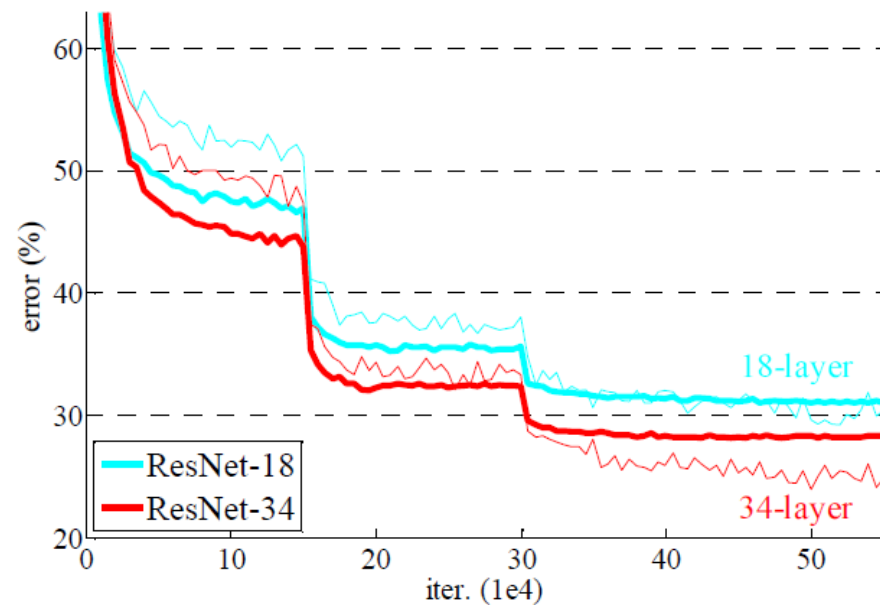
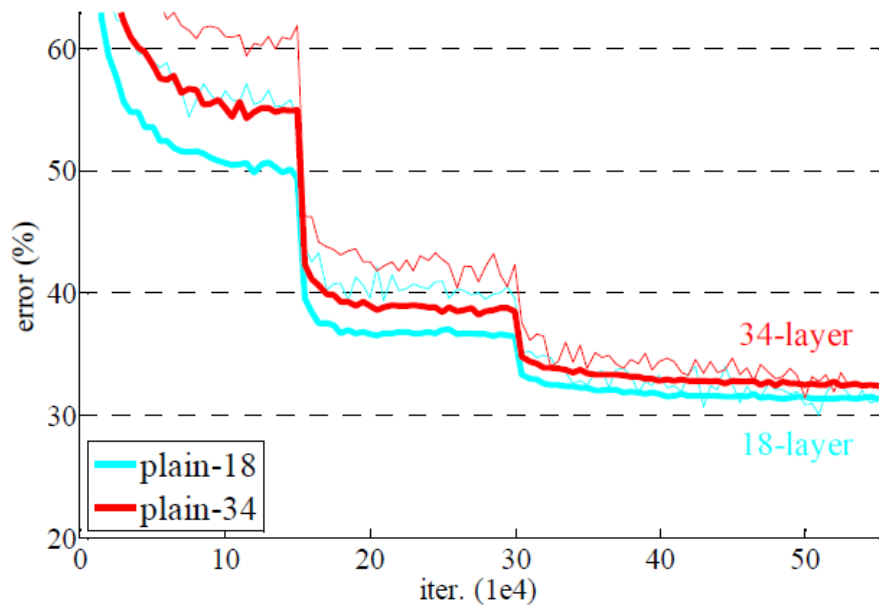


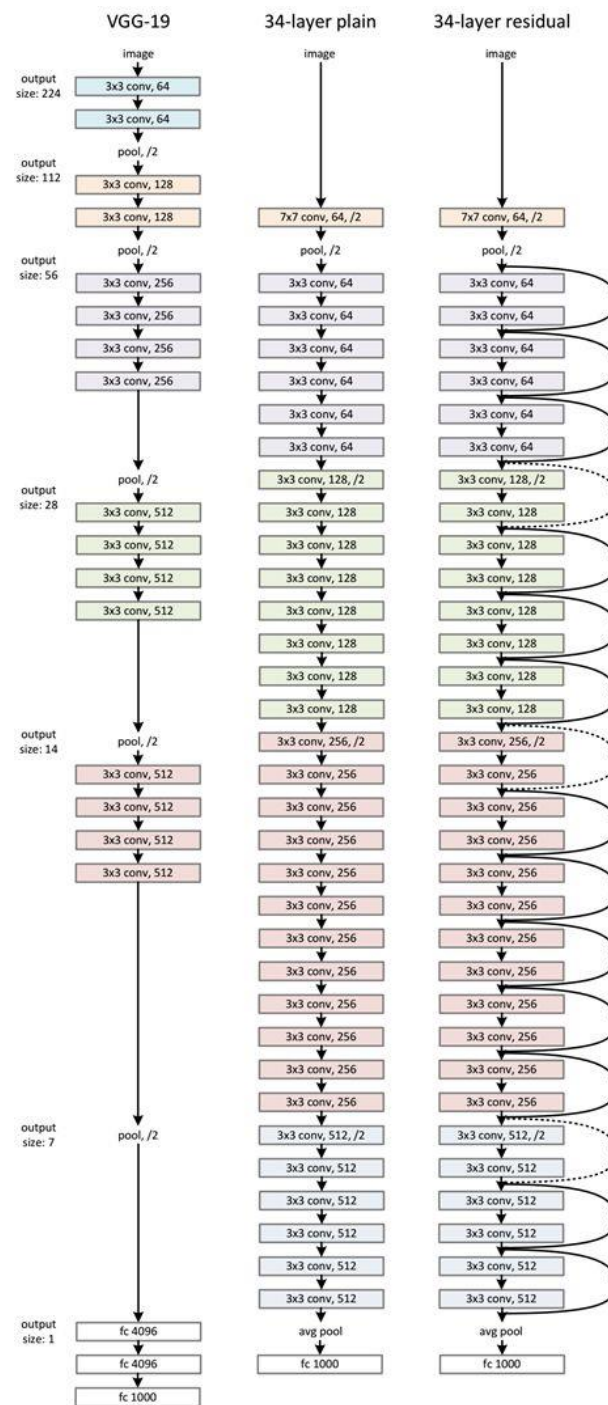
Degradation Problem in Deep Networks



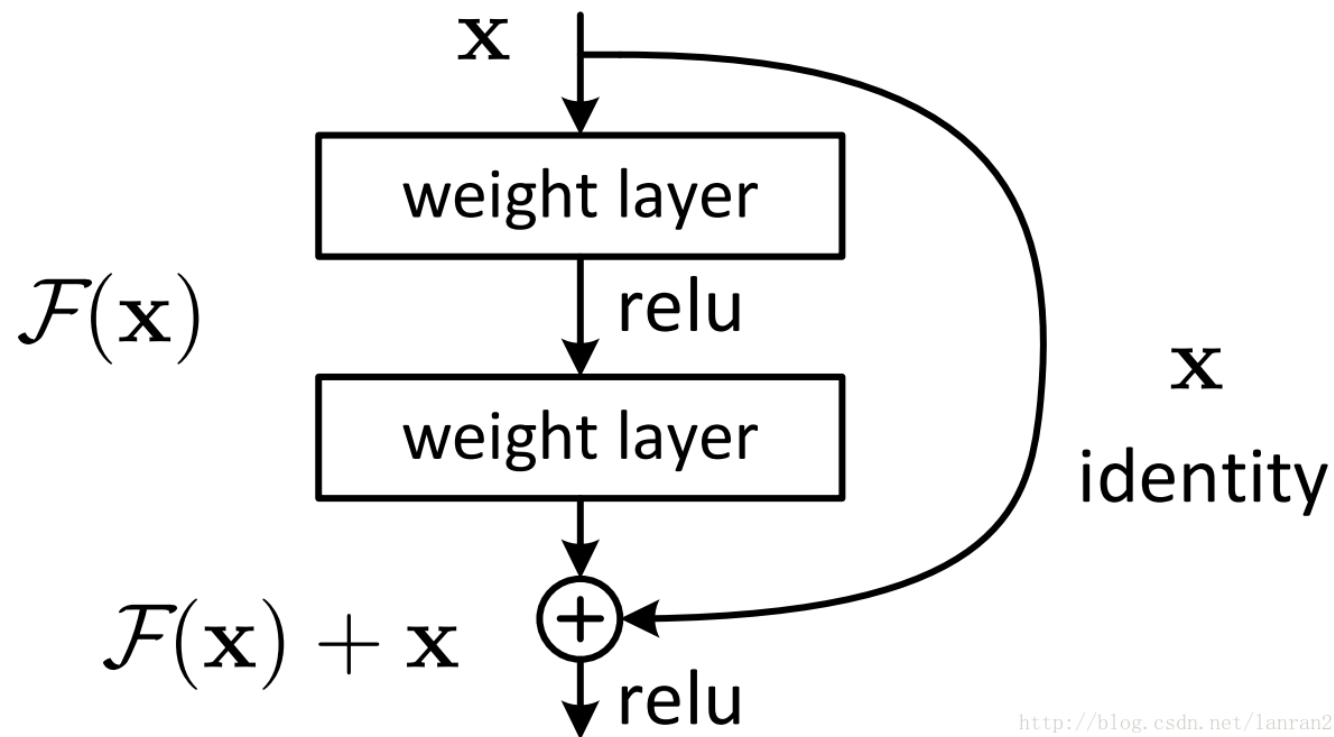
Deep Residual Learning

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. Deep Residual Learning for Image Recognition. CVPR 2016.

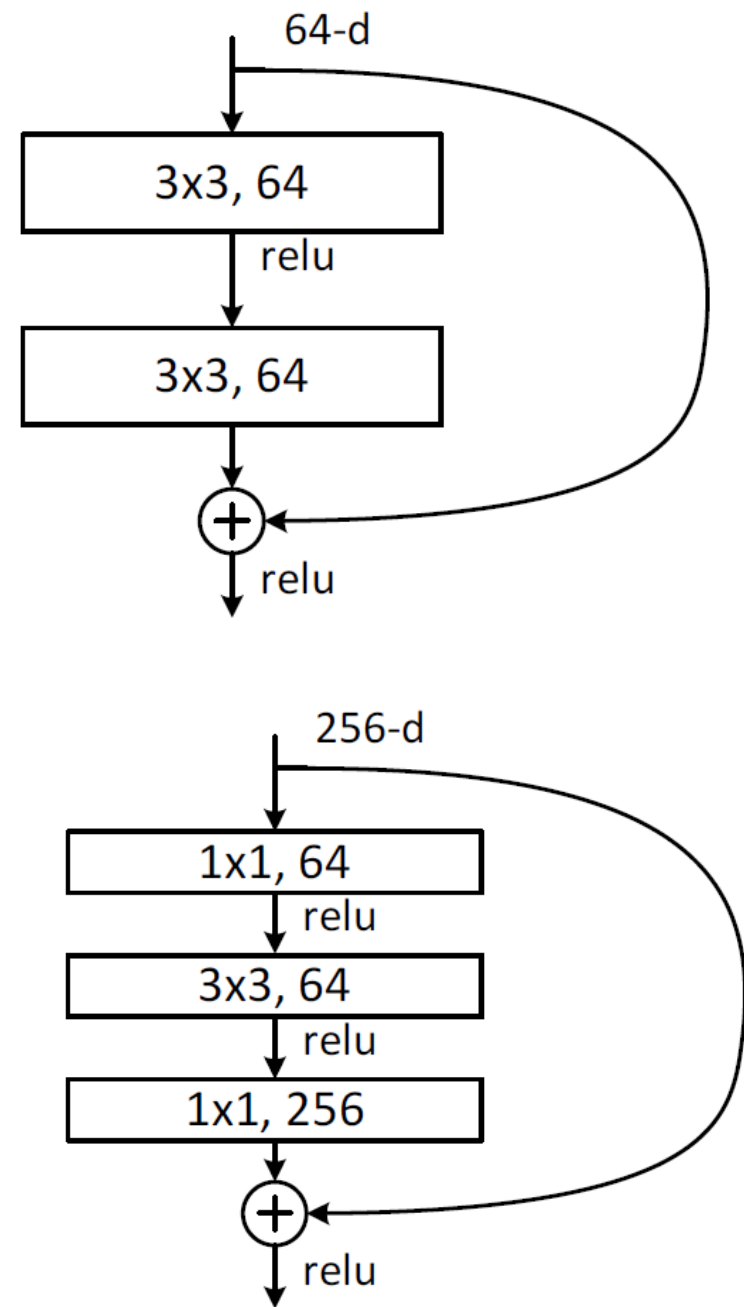




Identity Mapping by Shortcuts



<http://blog.csdn.net/lanran2>



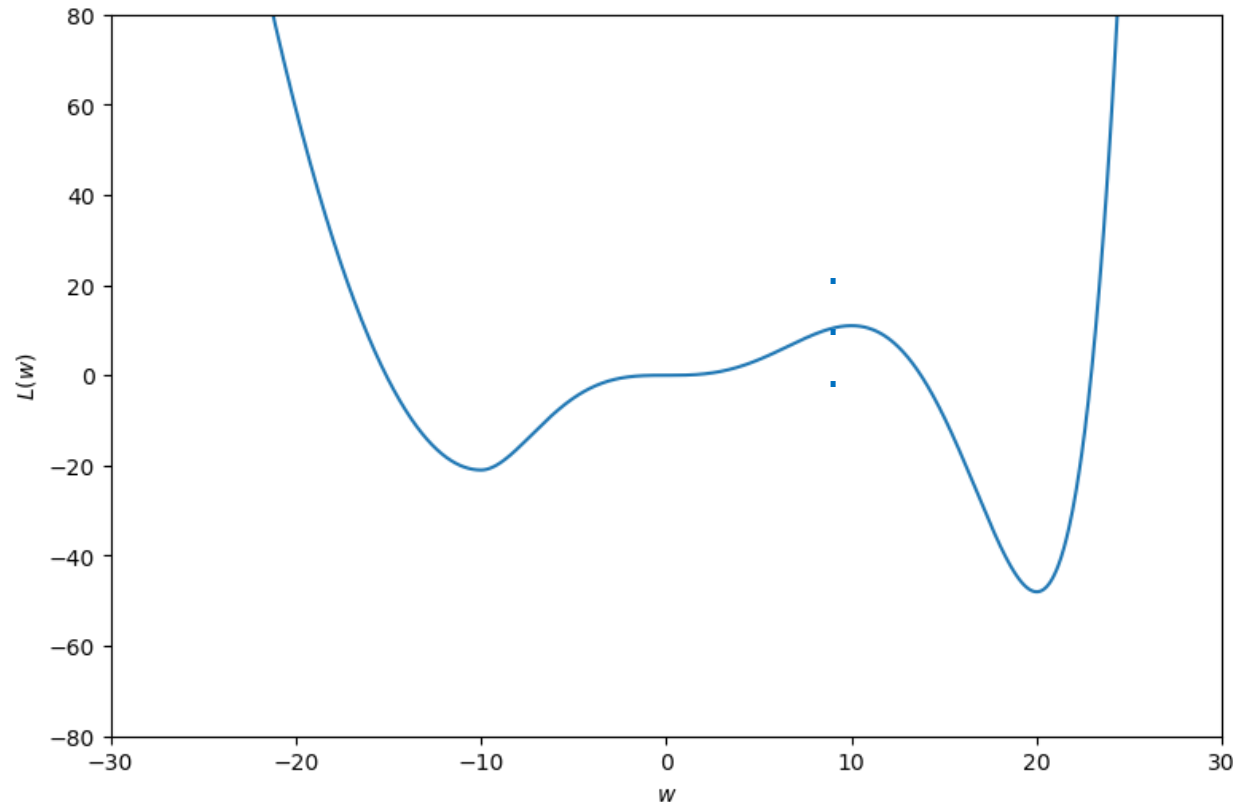
Some Techniques for Learning Better Deep Networks

Techniques for Learning Better NNs

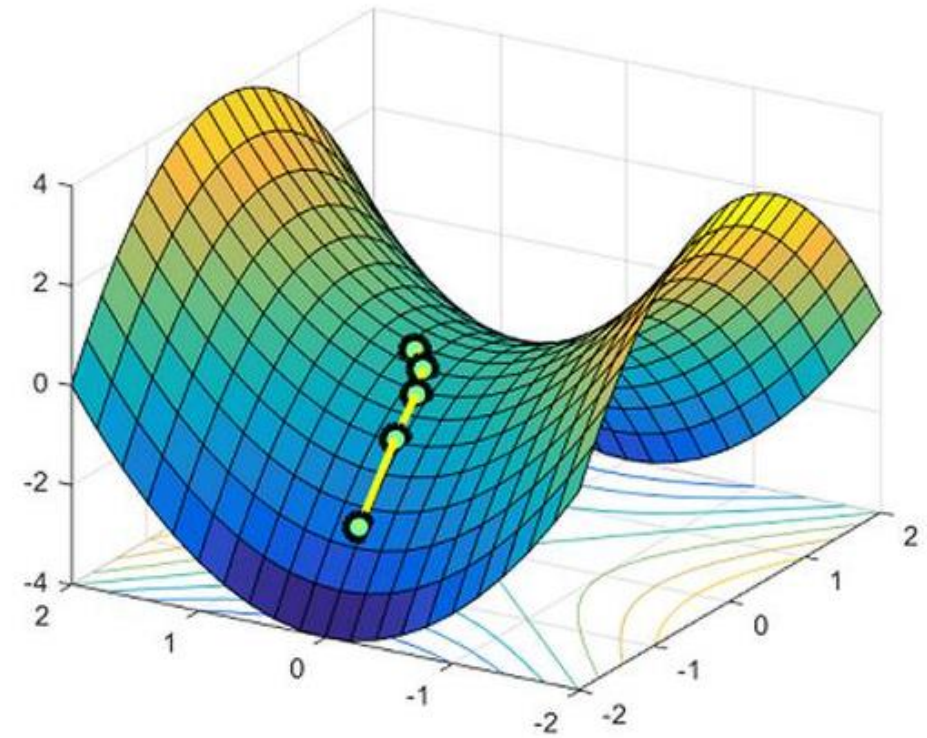
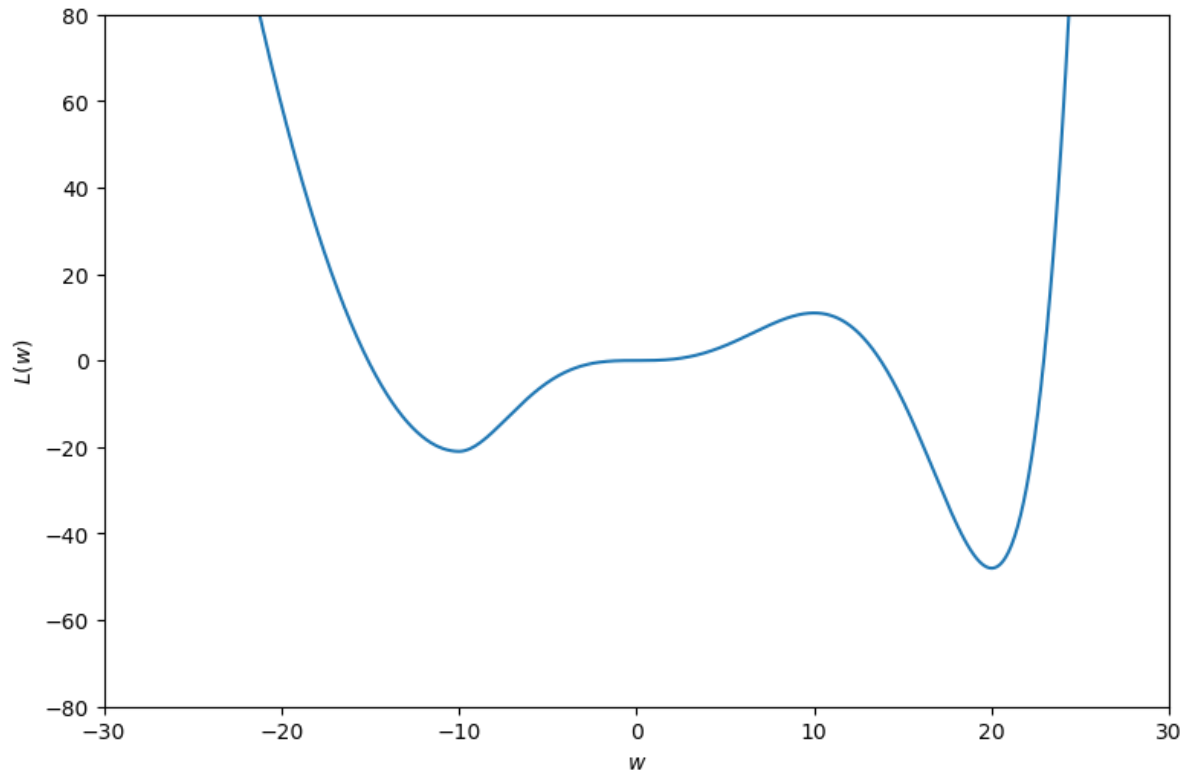
- Design good network structures
- Design good loss?
- Optimization techniques
- Others
 - Data preprocessing
 - Data augmentation

Remember Two Facts

- A NN is a complicated function
 - We can study it from function properties
- Optimizing a NN is usually a non-convex problem



Stationary Point and Saddle Points in High-dimensional Space



Flat Minima in NN

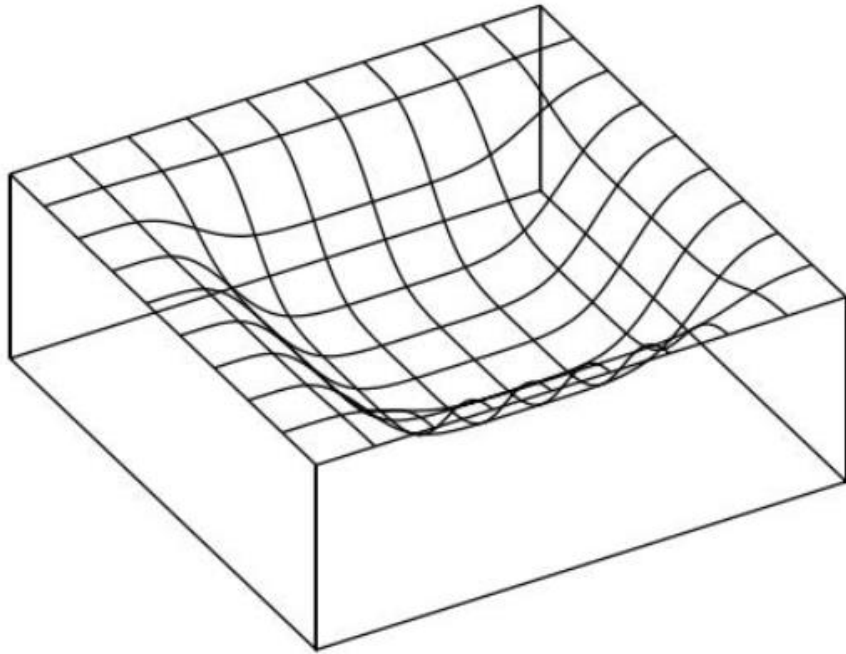


Figure 1: *Example of a “flat” minimum.*

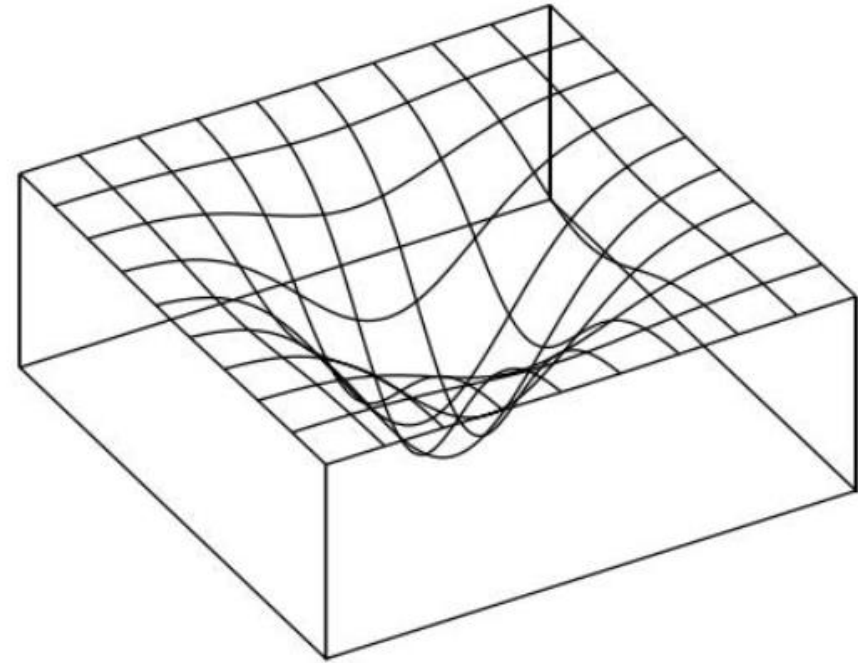
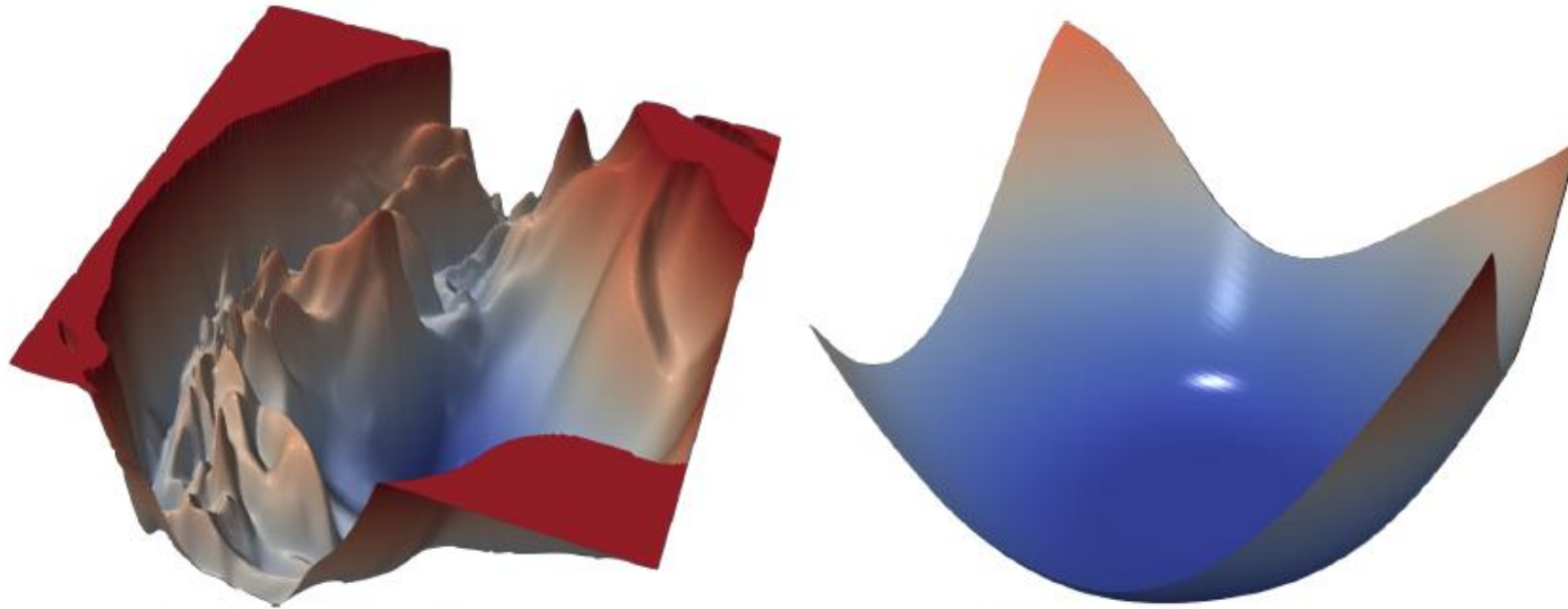


Figure 2: *Example of a “sharp” minimum.*

Loss Landscape of NN



(a) ResNet-110, no skip connections

(b) DenseNet, 121 layers

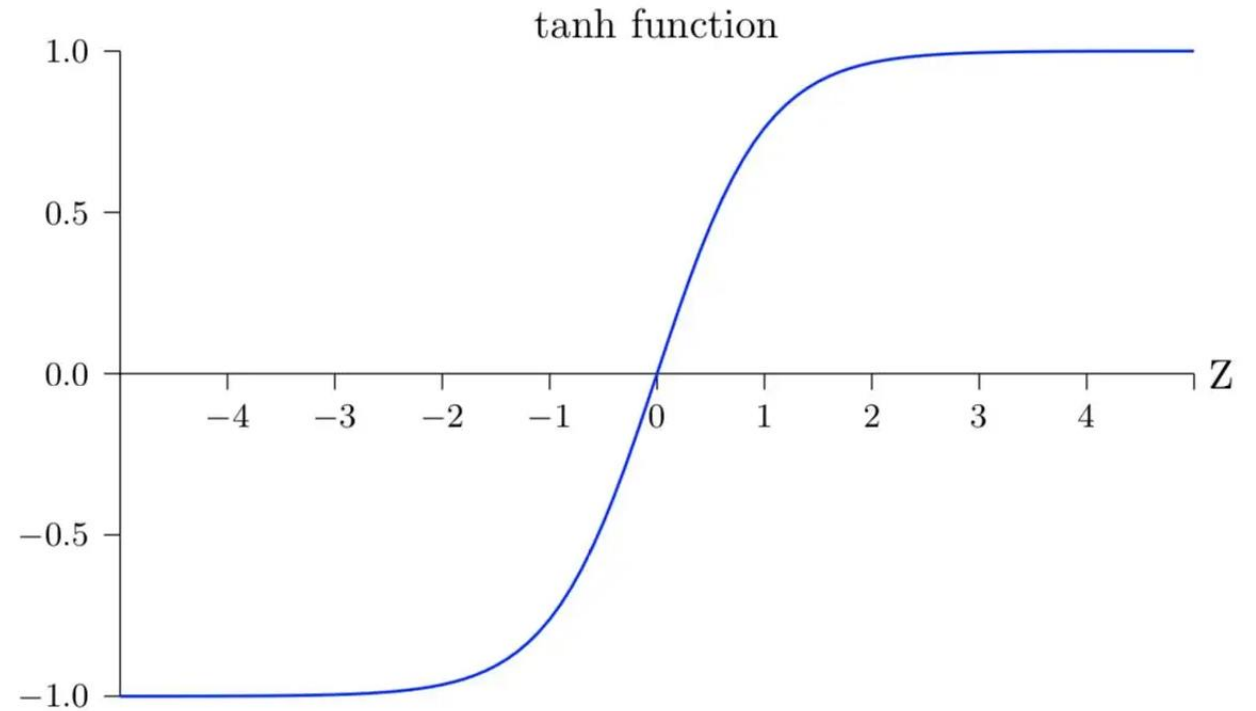
Figure 4: The loss surfaces of ResNet-110-noshort and DenseNet for CIFAR-10.

Network Design Issues

- How many layers?
- How many neurons in a layer?
- Connection structure:
 - other than fully-connected
- Activation Functions

Activation Function: Hyperbolic Tangent

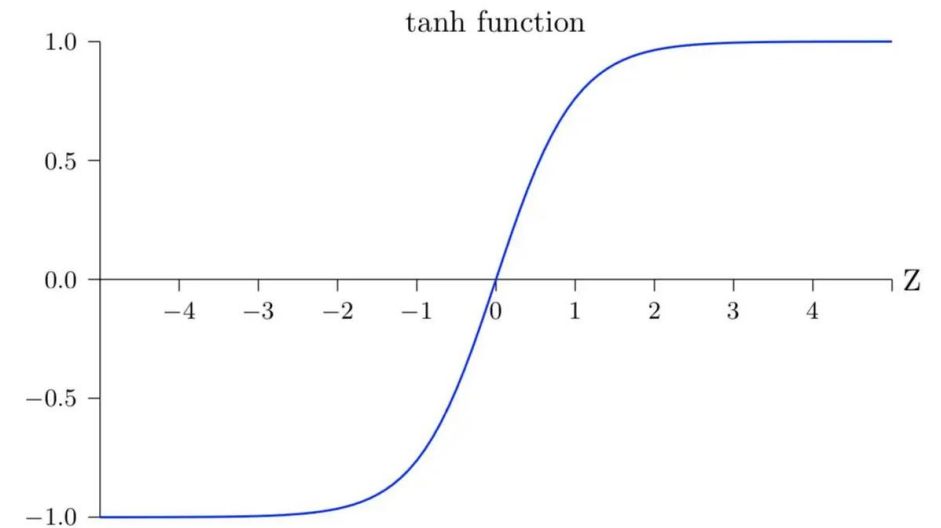
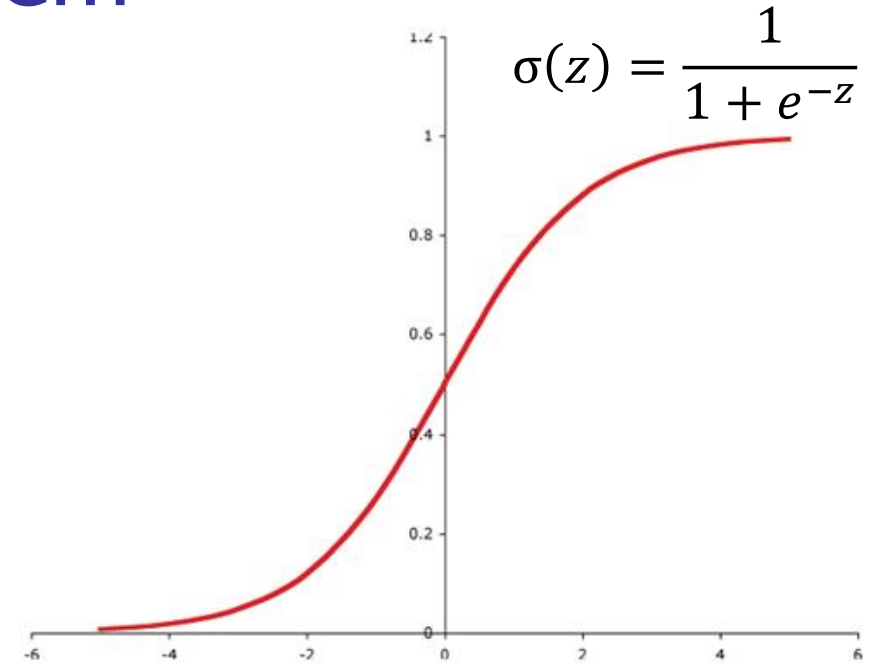
$$a(z) = \tanh(z)$$
$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Gradient Vanishing Problem

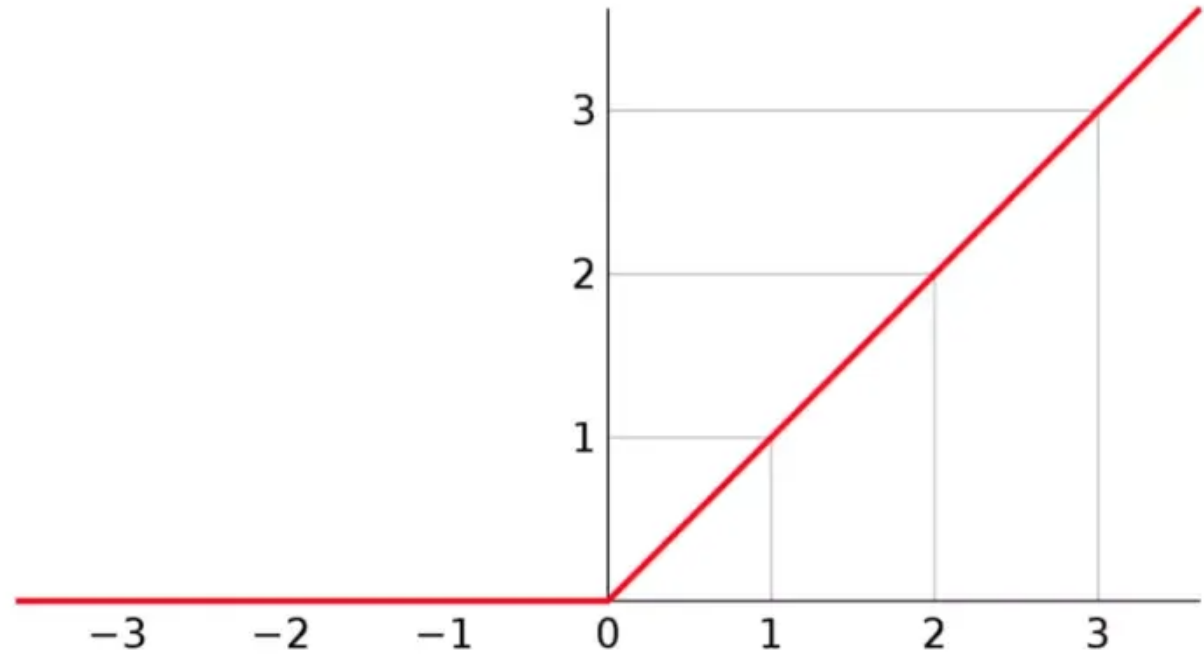
$$\frac{\partial L}{\partial w_{ij}^r} = \frac{\partial L}{\partial \mathbf{a}^L} \cdots \cdots \frac{\partial \mathbf{a}^{r+1}}{\partial a_i^r} \frac{\partial a_i^r}{\partial z_i^r} \frac{\partial z_i^r}{\partial w_{ij}^r}$$

Sigmoid function are now often used as gate unit



Activation Function: ReLU

$$a(z) = \text{ReLU}(z)$$
$$= \max(0, z)$$



Design Good Loss Functions

- Different tasks should have specific loss functions

- MSE

- Cross entropy

- Focal loss

-

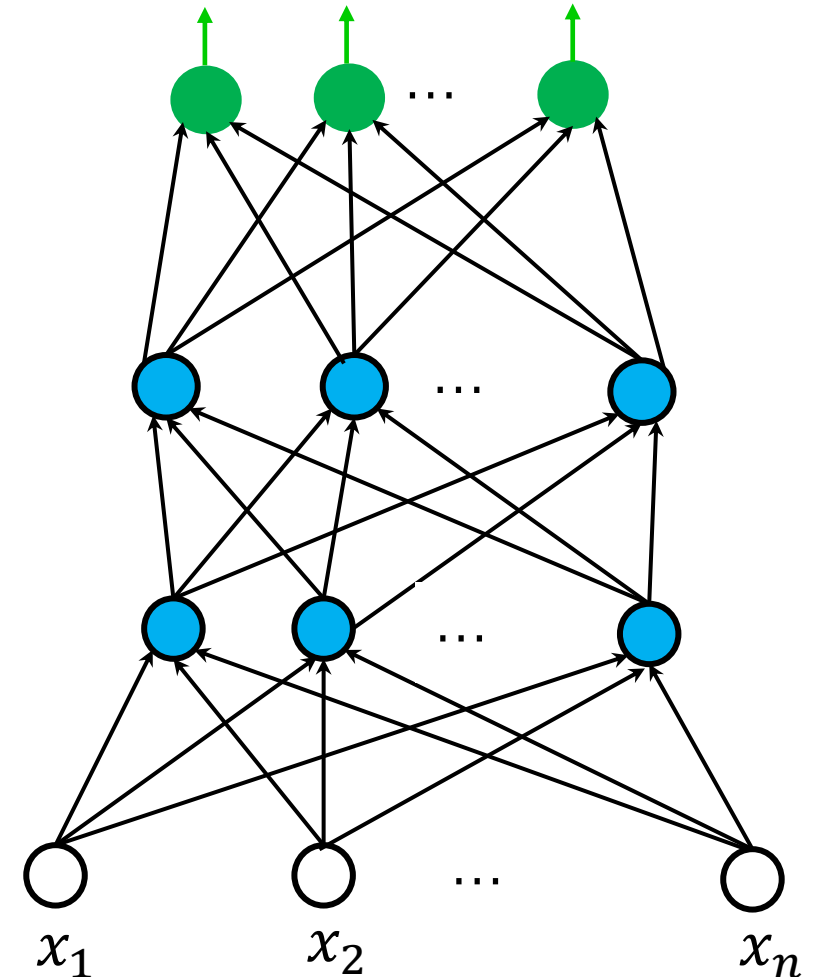
$$L(W) = \frac{1}{2} \sum_{i=1}^m \| \mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \|^2$$

$$L(W) = -\frac{1}{m} \sum_{i=1}^m \mathbf{y}^{(i)} \ln(\hat{\mathbf{y}}^{(i)})$$

- Adding regularization terms: L1, L2,

Parameter Initialization Techniques

- Initialize weights with Small random numbers
 - Set bias to zero initially
- Layer-wise pretraining
- Pretrained model+finetune



Employ Mini-batch instead of SGD

- Bias and Variance
- Small mini-batch usually don't change the bias expectation of SGD
- Mini-batch can fully exploit the GPU power
- Big mini-batch can have lower variance without big loss of bias expectation
 - Can have bigger learning rate
 - Correspondingly, with small mini-batch, learning rate should also be small
- Can set the learning rate linearly with mini-batch size

Gradient Updating Optimization

- The problem of gradient estimation error

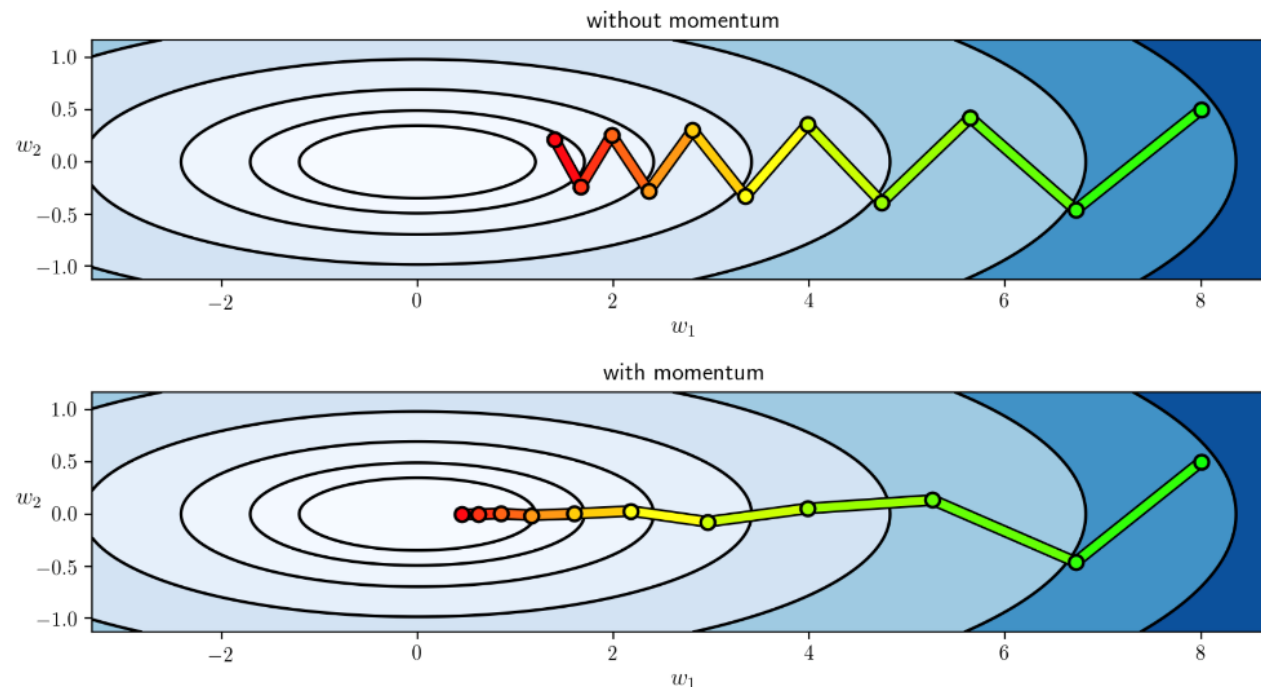
- Momentum+

Adaptive Learning Rate

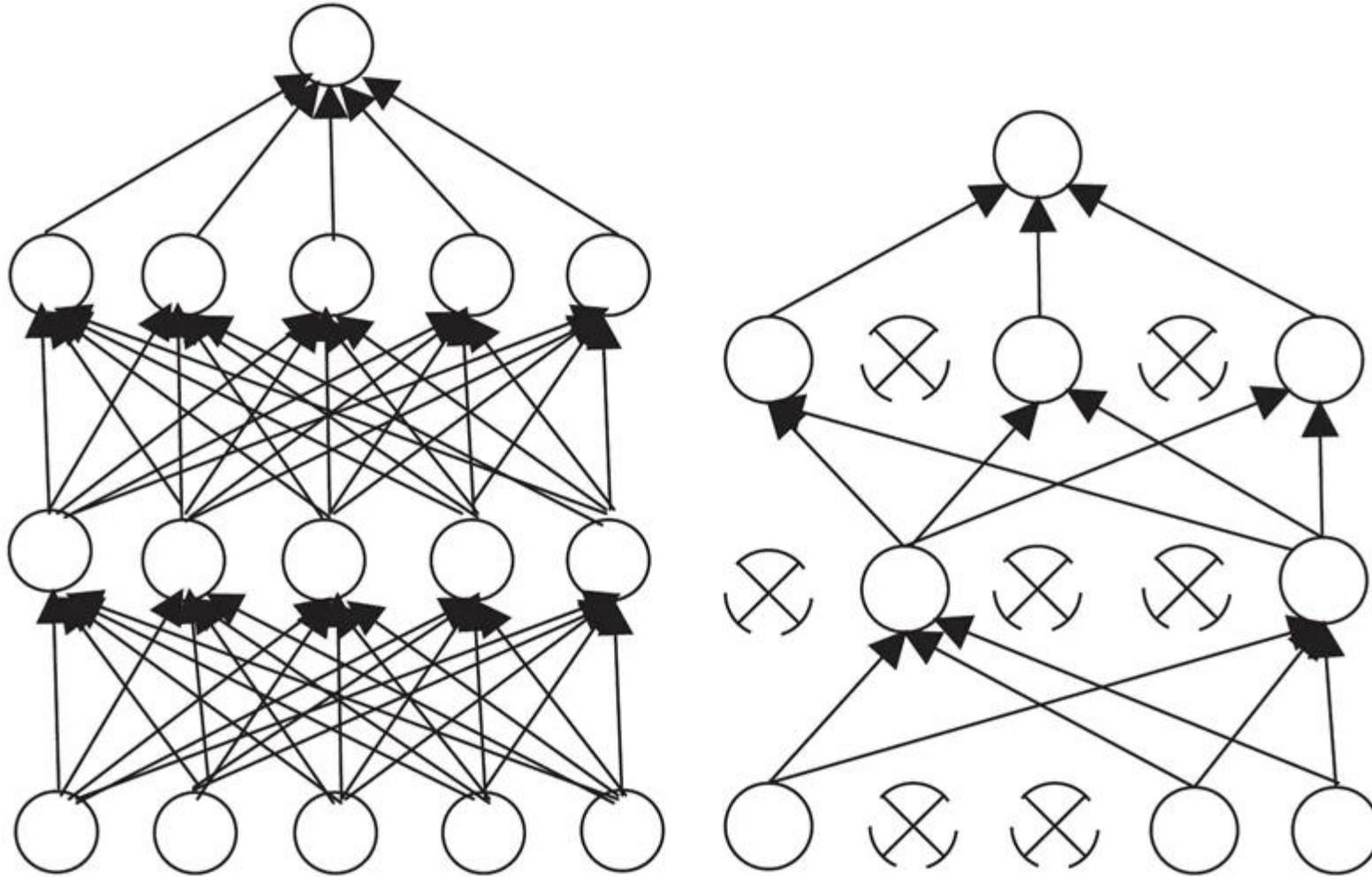
- Adam

$$\Delta\theta_t = \rho\Delta\theta_{t-1} - \alpha\mathbf{g}_t$$

- Gradient clipping



Dropout



Other Classical Techniques

- Batch normalization
- Layer normalization
- Weight decay
- Early stopping
- Hyperparameter optimization
 - Grid search
 - Bayesian optimization

Acknowledgements

- Some text, figures and formulations are from WWW.
- This lecture is distributed for nonprofit purpose.

Thank You for Your Attention

Contact me at: yym@hit.edu.cn

Tel: 26033008, 13760196623

Address: Rm.1402, H# Building

