# 人工智能技术及应用
Artificial Intelligence and Application

# Classification

# Classification

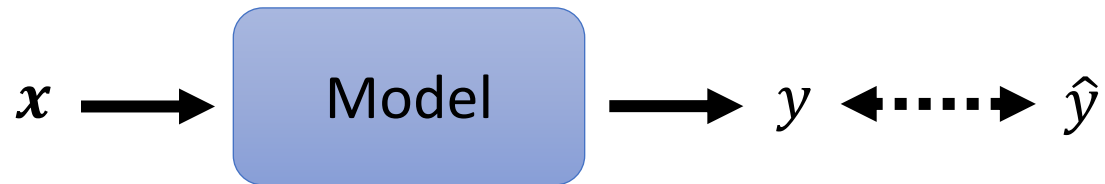$x$ → **Function** → Class n

- Credit Scoring
  - Input: income, savings, profession, age, past financial history ……
  - Output: accept or refuse
- Medical Diagnosis
  - Input: current symptoms, age, gender, past medical history ……
  - Output: which kind of diseases
- Handwritten character recognition

  Input:   output: 金

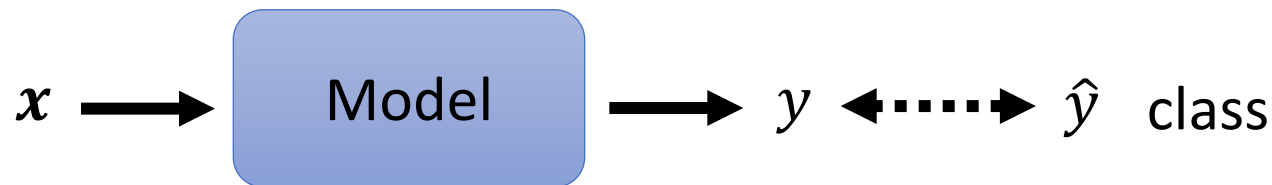- Face recognition
  - Input: image of a face, output: person
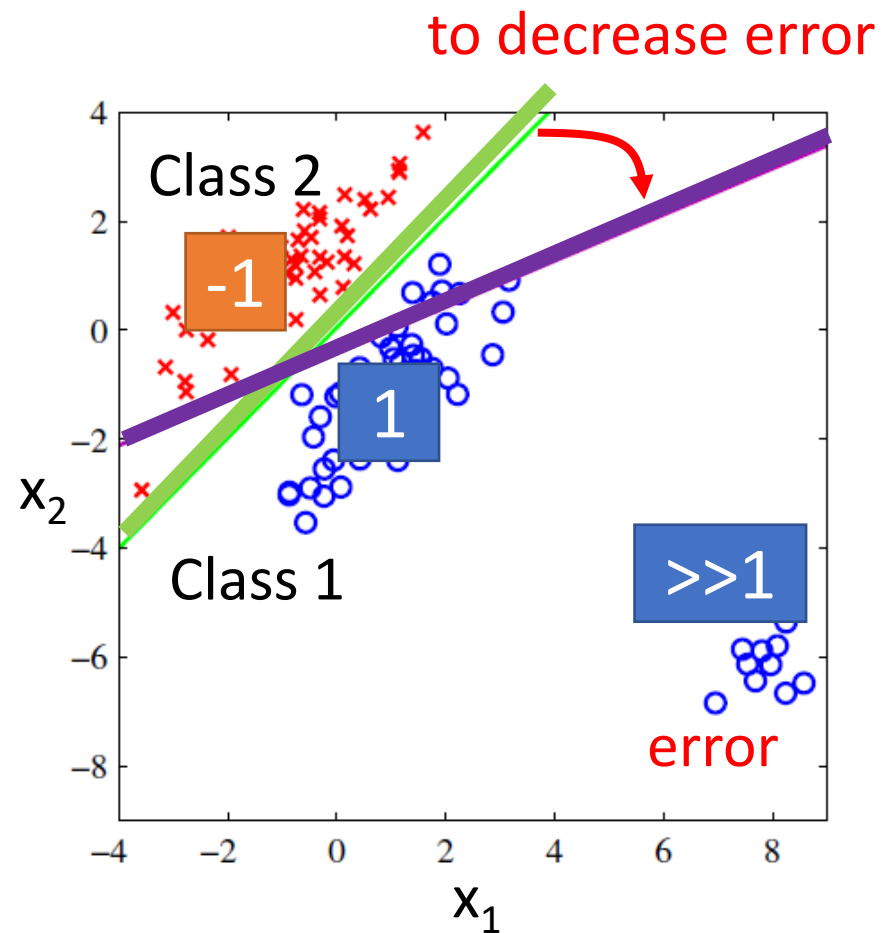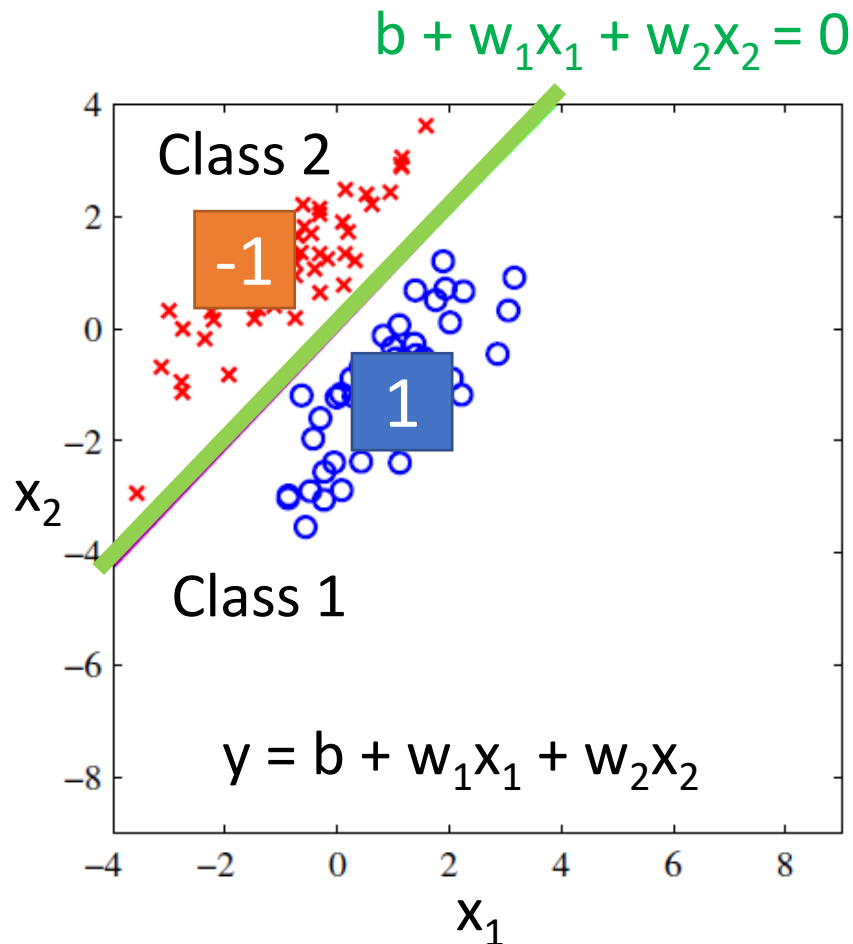
# Classification as Regression?

- Regression

$$x \longrightarrow \boxed{\text{Model}} \longrightarrow y \longleftarrow\cdots\longrightarrow \hat{y}$$

- Classification as regression?

$$x \longrightarrow \boxed{\text{Model}} \longrightarrow y \longleftarrow\cdots\longrightarrow \hat{y} \quad \text{class}$$

1 = class 1
2 = class 2
3 = class 3

different?

similar?

$b + w_1x_1 + w_2x_2 = 0$

to decrease error

Class 2

-1

1

Class 1

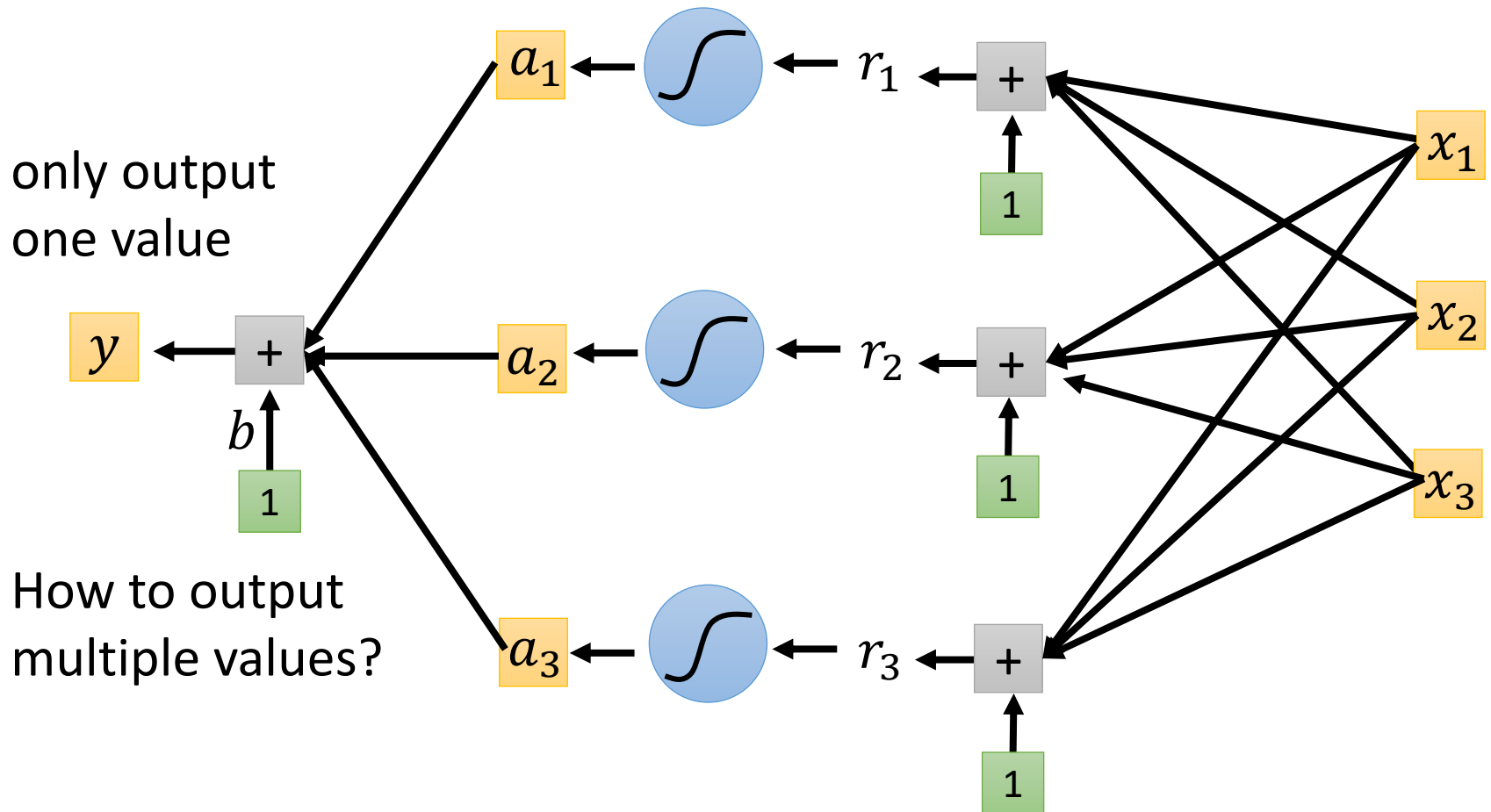$y = b + w_1x_1 + w_2x_2$

Class 2

-1

1

Class 1

>>1

error

Penalize to the examples that are "too correct" …   (Bishop, P186)

- Multiple class: Class 1 means the target is 1; Class 2 means the target is 2; Class 3 means the target is 3 …… problematic

# Class as one-hot vector

$$\hat{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Class 1    Class 2    Class 3
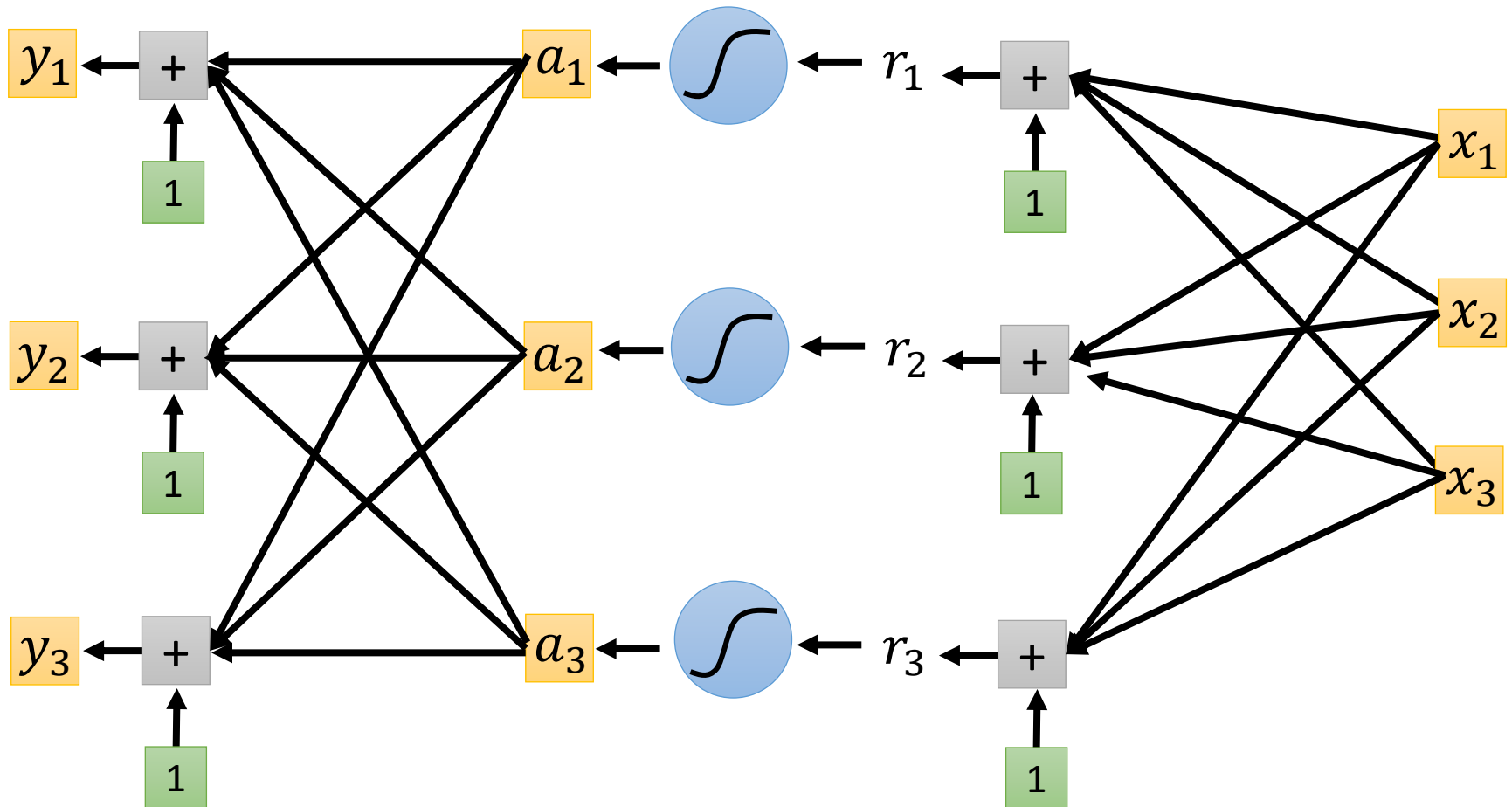
only output
one value

How to output
multiple values?

# Class as one-hot vector

$$\widehat{\boldsymbol{y}} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Class 1      Class 2      Class 3

## *Regression*

label

$$\hat{y} \longleftrightarrow y = b + \boldsymbol{c}^T \, \sigma( \boldsymbol{b} + W \boldsymbol{x} )$$

feature

## *Classification*

feature

$$\boldsymbol{y} = \boldsymbol{b}' + W' \, \sigma( \boldsymbol{b} + W \boldsymbol{x} )$$

label $\widehat{\boldsymbol{y}} \longleftrightarrow \boldsymbol{y}' = softmax( \boldsymbol{y} )$

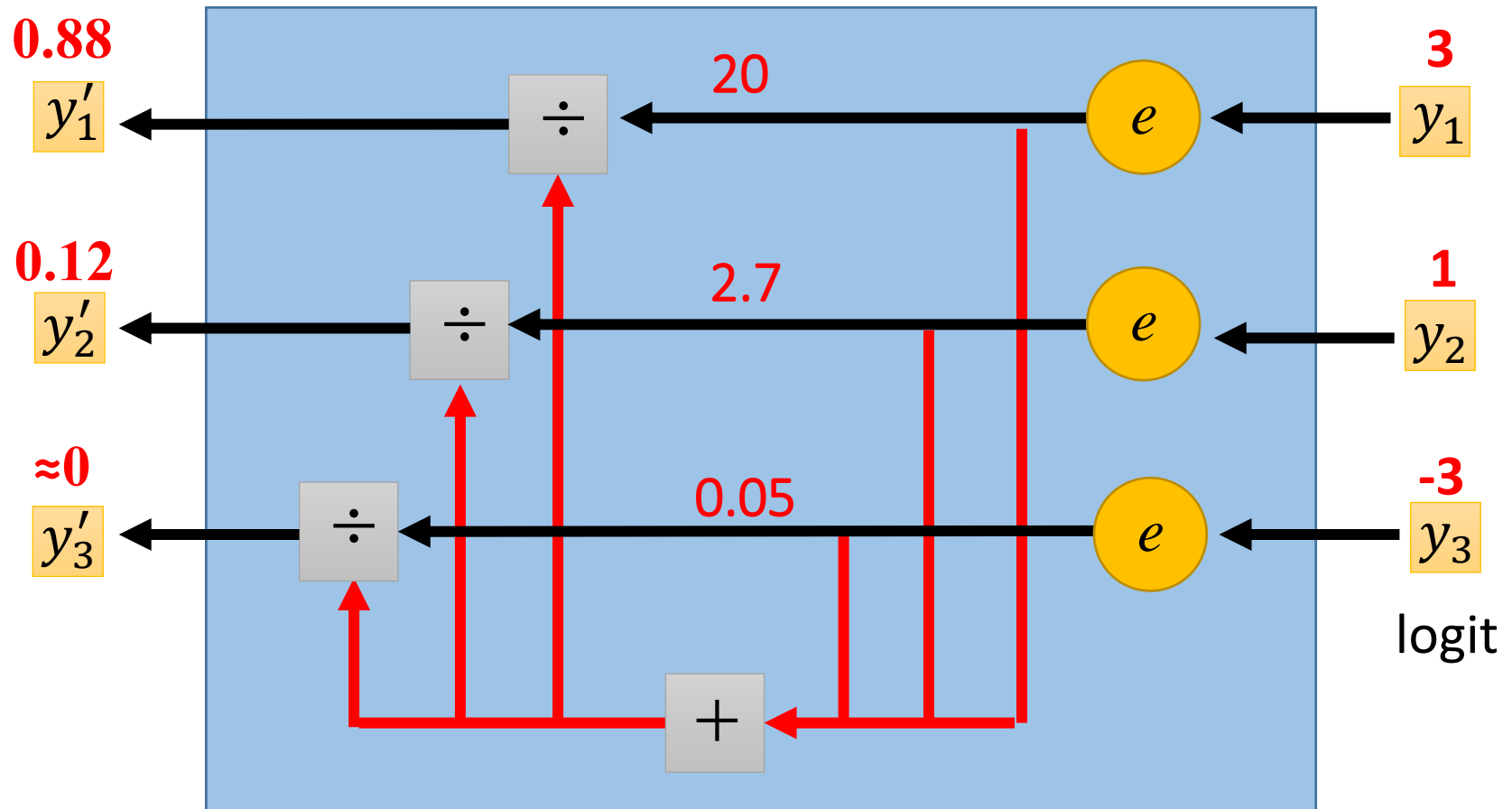0 or 1     Make all values between 0 and 1     Can have any value

# Soft-max

$$y_i' = \frac{exp(y_i)}{\sum_j exp(y_i)}$$

- $1 > y_i' > 0$
- $\sum_i y_i' = 1$

*__Softmax__*

How about **binary classification**? ☺



0.88 $y_1'$ ← ÷ ← 20 ← $e$ ← 3 $y_1$

0.12 $y_2'$ ← ÷ ← 2.7 ← $e$ ← 1 $y_2$

≈0 $y_3'$ ← ÷ ← 0.05 ← $e$ ← -3 $y_3$

logit

+

# *Multi-class Classification* (3 classes as example)

$z_1 = w^1 \cdot x + b_1$

$x$

$z_2 = w^2 \cdot x + b_2$

$z_3 = w^3 \cdot x + b_3$

Softmax

$y$

$y_1$

$y_2$

$y_3$

Cross Entropy

$$-\sum_{i=1}^{3} \hat{y}_i \ln y_i$$

$\hat{y}$

$\hat{y}_1$

$\hat{y}_2$

$\hat{y}_3$

target

If x ∈ class 1

$$\hat{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$-\ln y_1$

If x ∈ class 2

$$\hat{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$-\ln y_2$

If x ∈ class 3
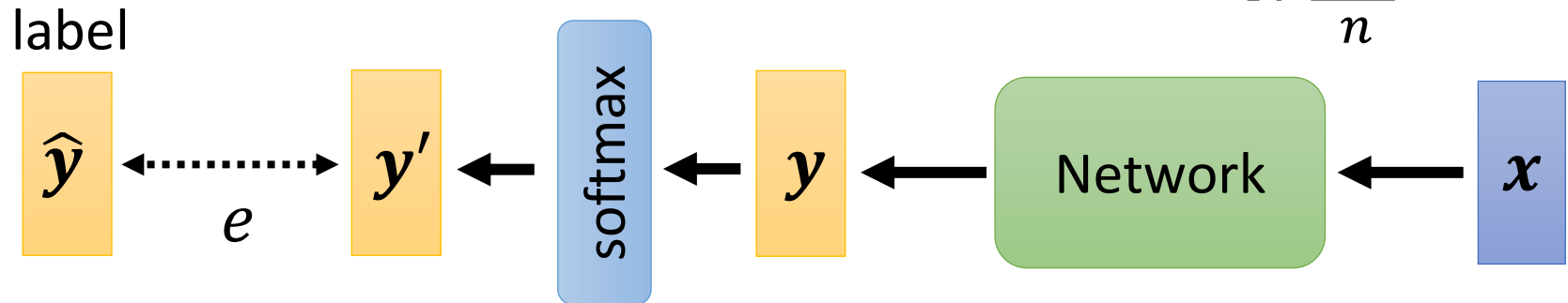
$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$-\ln y_3$

# Loss of Classification

$$L = \frac{1}{N} \sum_n e_n$$

label

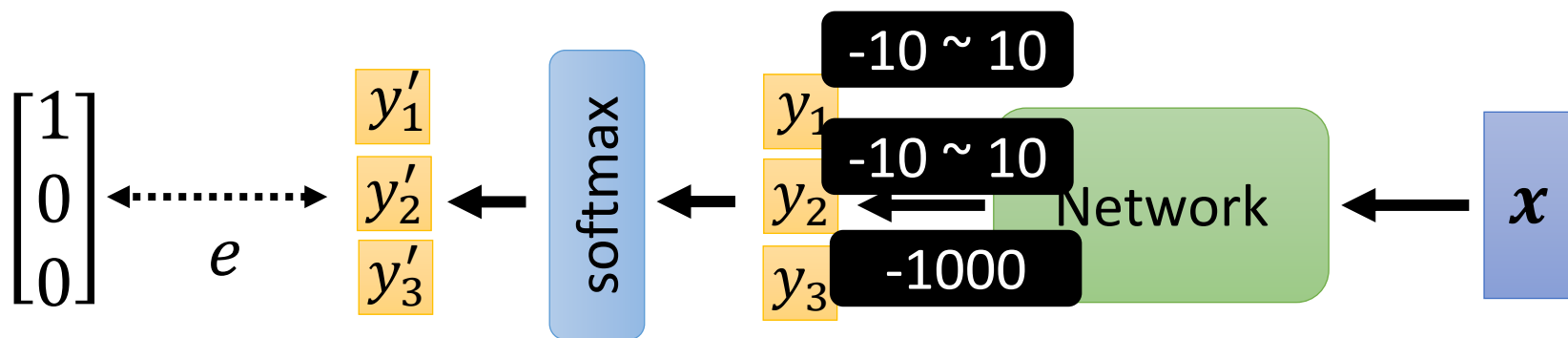$$\hat{\boldsymbol{y}} \quad \overset{e}{\longleftrightarrow} \quad \boldsymbol{y'} \longleftarrow \text{softmax} \longleftarrow \boldsymbol{y} \longleftarrow \text{Network} \longleftarrow \boldsymbol{x}$$

*Mean Square Error (MSE)* $\quad e = \sum_i (\hat{\boldsymbol{y}}_i - \boldsymbol{y'}_i)^2$

*Cross-entropy* **WIN** $\quad e = -\sum_i \hat{\boldsymbol{y}}_i \ln \boldsymbol{y'}_i$

**Minimizing cross-entropy** is equivalent to **maximizing likelihood**.

**Changing the loss function can change the difficulty of optimization.**

Convolutional Neural Network (CNN)

# Image Classification



$$\begin{bmatrix} \vdots \\ 0.2 \\ 0.7 \\ 0.1 \\ \vdots \end{bmatrix} \quad \begin{matrix} \text{dog} \\ \text{cat} \\ \text{tree} \end{matrix} \quad \begin{bmatrix} \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix}$$
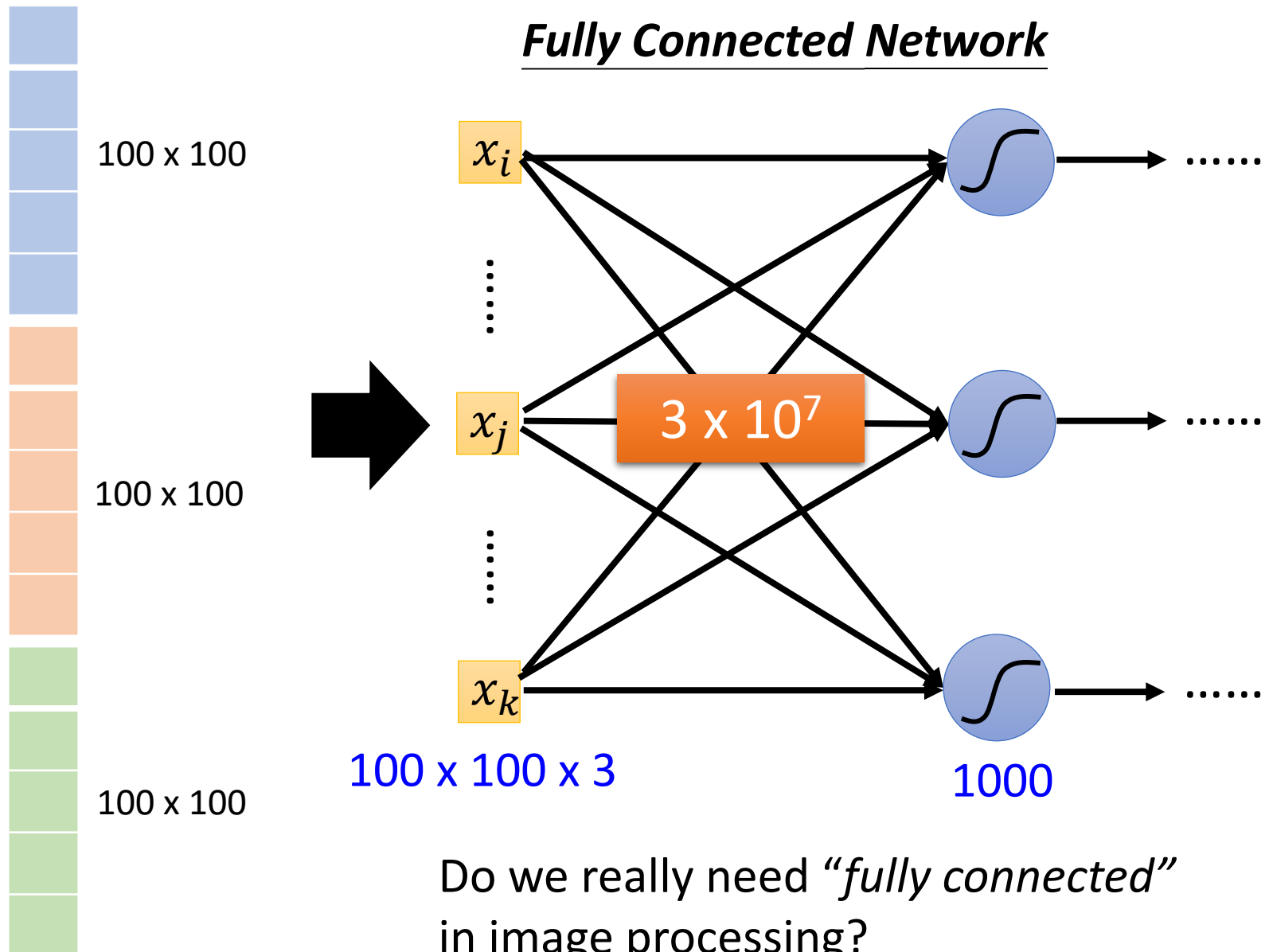
Model

100 x 100

$y'$ ◀┅┅┅▶ $\hat{y}$

Cross entropy

(All the images to be classified have the same size.)

# Image Classification



3 channels

3-D tensor

100

100

100 x 100

100 x 100

100 x 100

100 x 100

value represents intensity

**_Fully Connected Network_**

$x_i$

$x_j$

$3 \times 10^7$

$x_k$

100 x 100

100 x 100

100 x 100

100 x 100 x 3
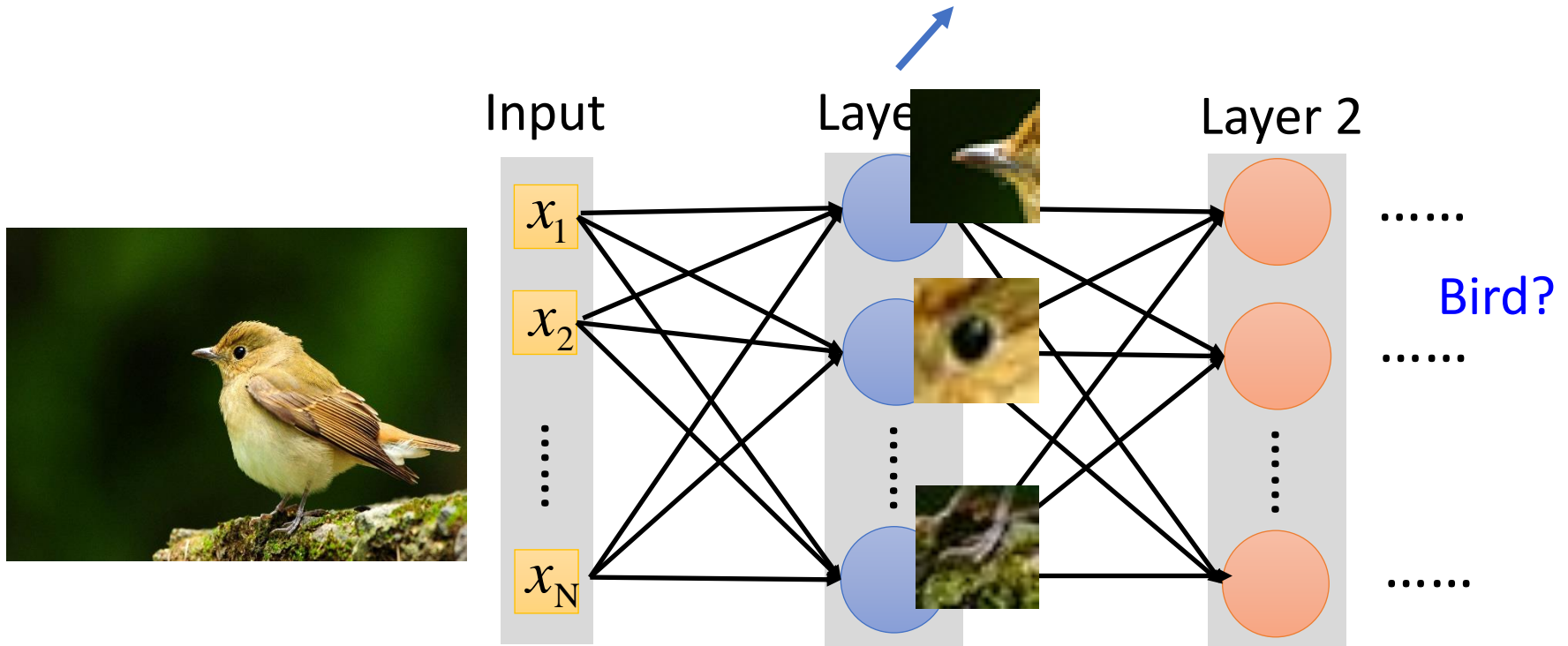
1000

Do we really need "*fully connected*" in image processing?

16

# Observation 1
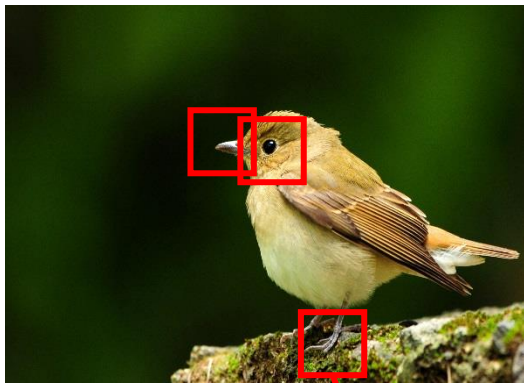
Identifying some critical patterns



Perhaps human also identify birds in a similar way ... ☺

# Observation 1

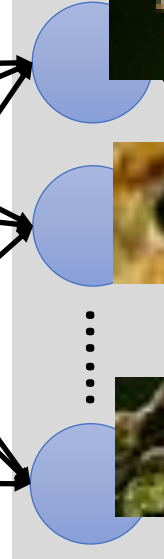A neuron does not have to see the whole image.
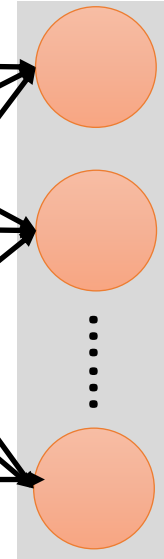
Need to see the whole image?



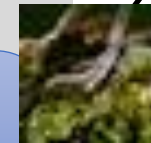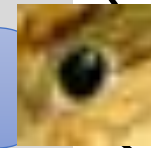Input     Layer     Layer 2

$x_1$

$x_2$

$x_N$

basic detector

advanced detector

...... 

bird

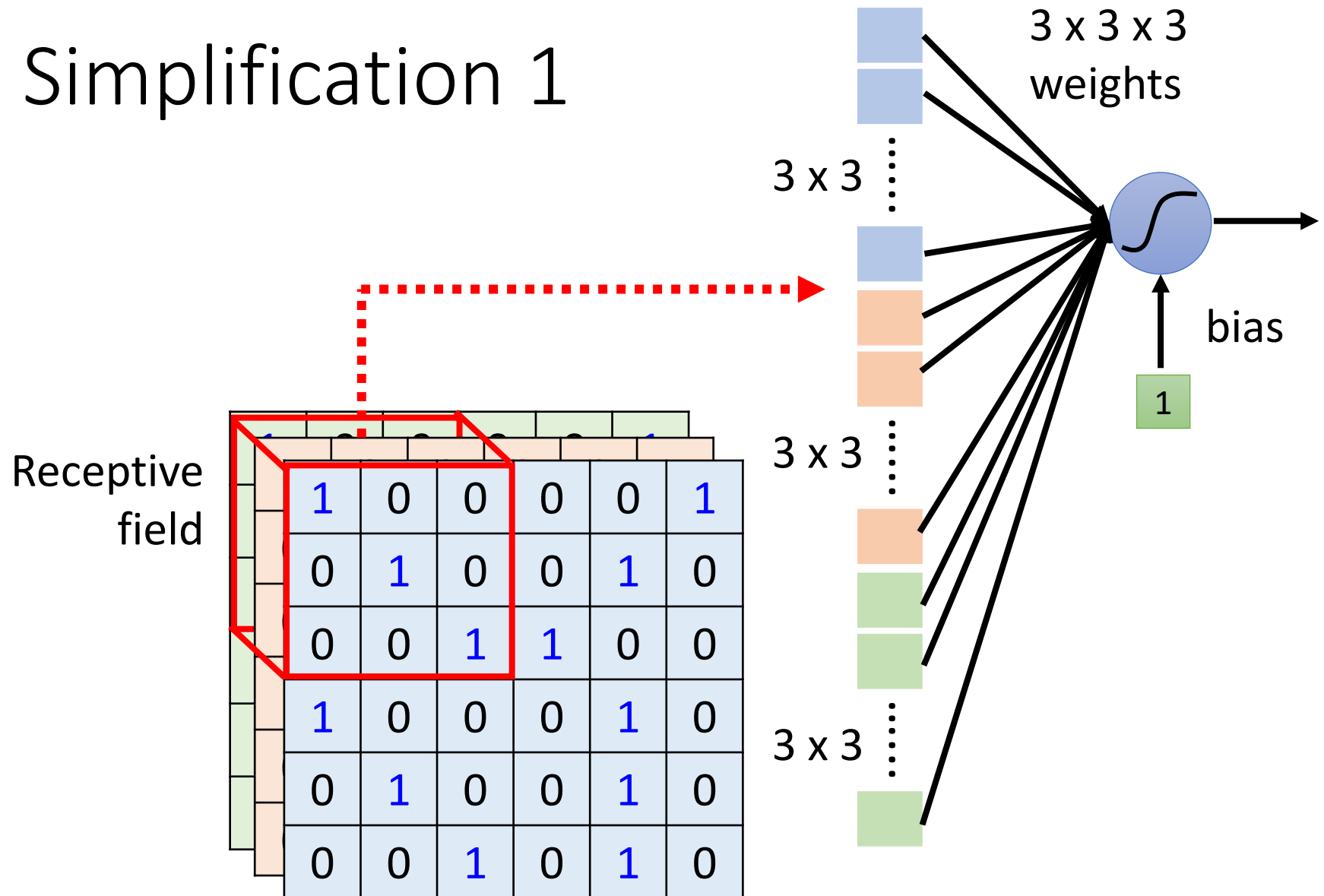...... 

Some patterns are much smaller than the whole image.

Connecting to small region with less parameters

# Simplification 1
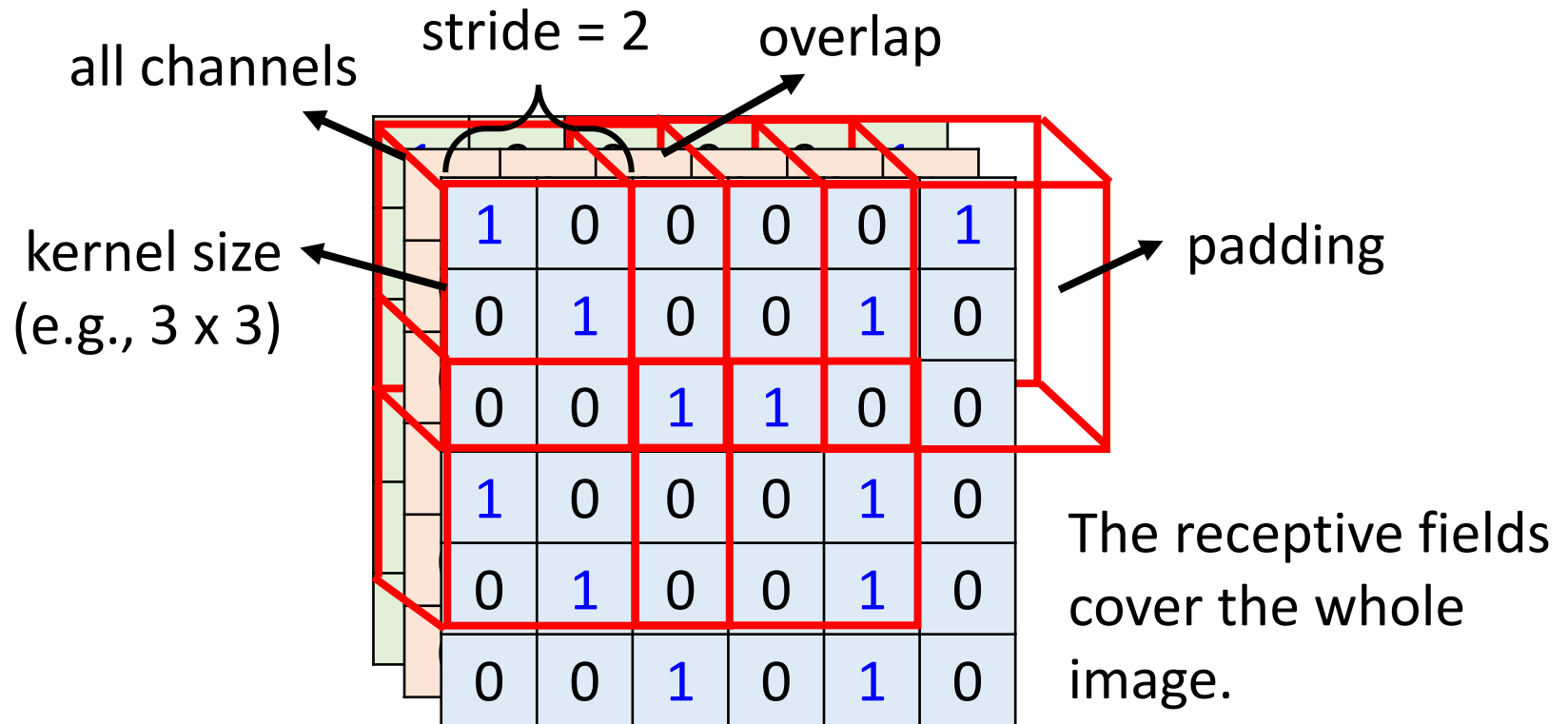


Receptive field

3 x 3 x 3 weights

3 x 3

3 x 3

3 x 3

bias

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 |

19

# Simplification 1

- Can different neurons have different sizes of receptive field?
- Cover only some channels?
- Not square receptive field?

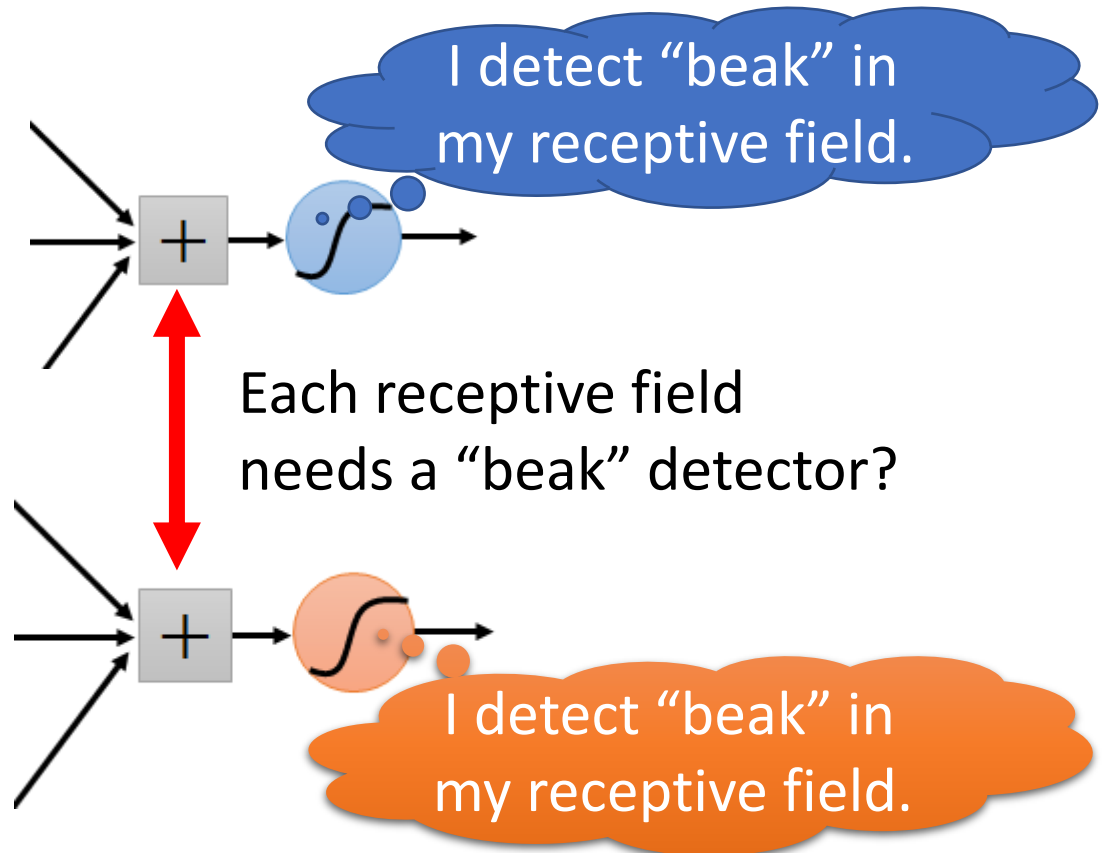3 x 3 x 3 weights

Receptive field

the same receptive field

Can be overlapped

| 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 |

# Simplification 1 – Typical Setting

Each receptive field has a set of neurons (e.g., 64 neurons).



stride = 2

overlap

all channels

kernel size (e.g., 3 x 3)

padding

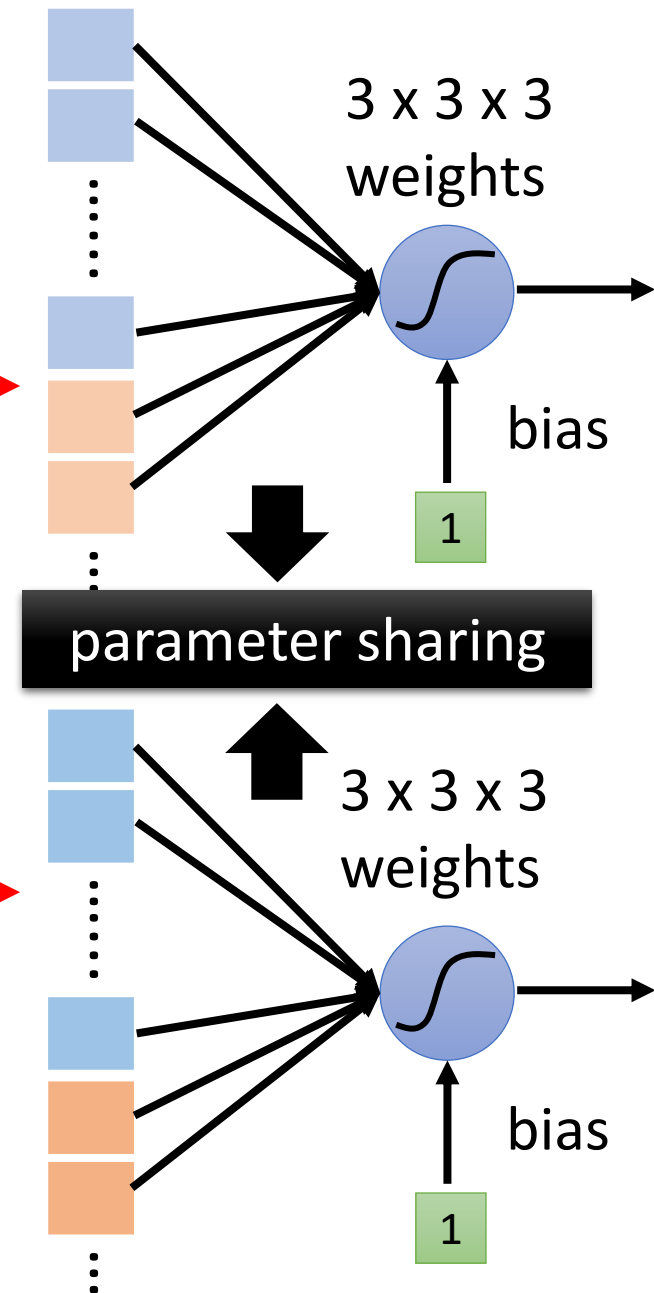| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

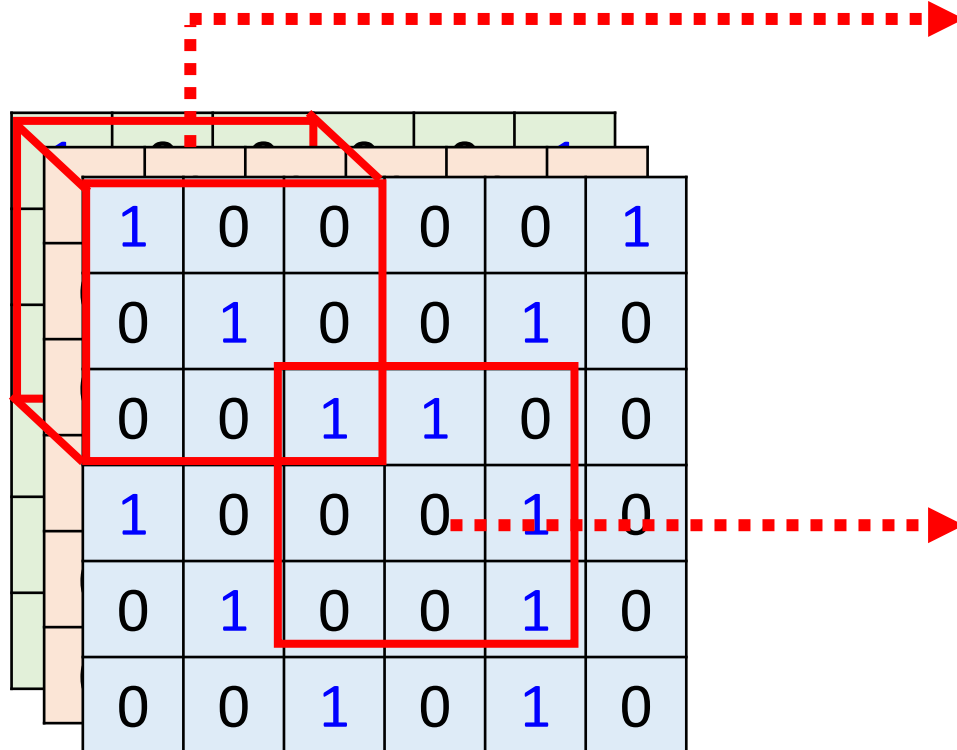The receptive fields cover the whole image.

# Observation 2

- The same patterns appear in different regions.

# Simplification 2



3 x 3 x 3 weights

bias

1

**parameter sharing**

3 x 3 x 3 weights

bias

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 |

# Simplification 2

$\sigma(\textcolor{red}{w_1}x_1 + \textcolor{orange}{w_2}x_2 + \cdots)$

$x_1$

$x_2$

$\textcolor{red}{w_1}$

$\textcolor{orange}{w_2}$

bias

1

$\sigma(\textcolor{red}{w_1}x_1' + \textcolor{orange}{w_2}x_2' + \cdots)$

$x_1'$

$x_2'$

$\textcolor{red}{w_1}$

$\textcolor{orange}{w_2}$

bias

1



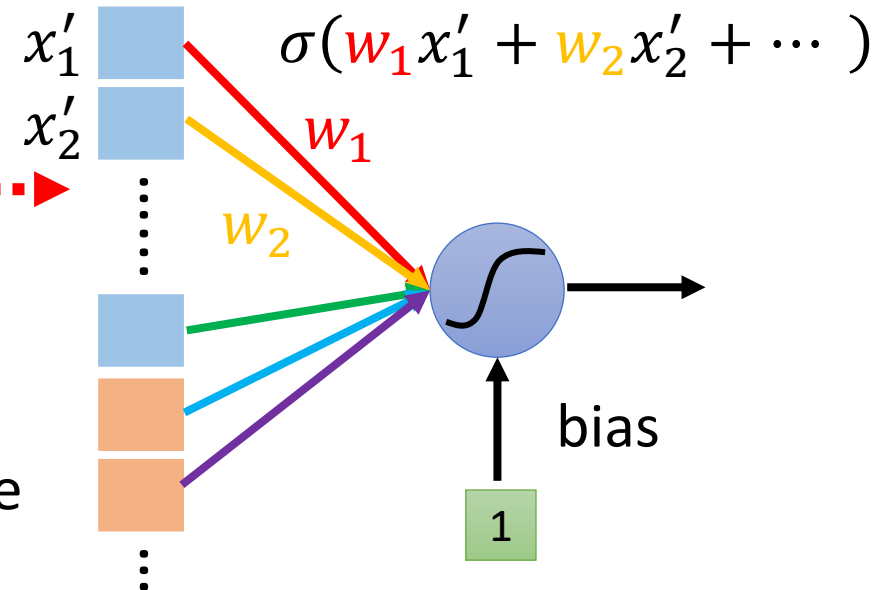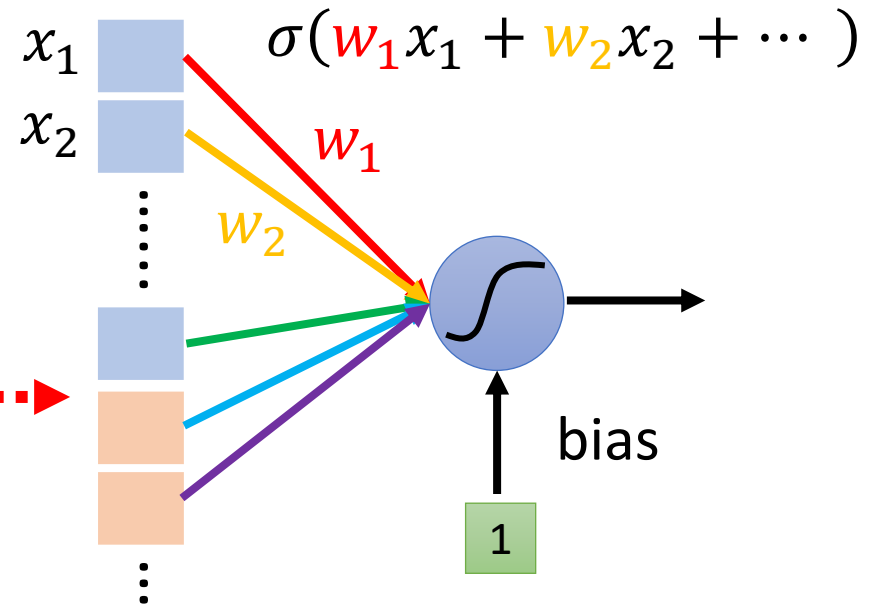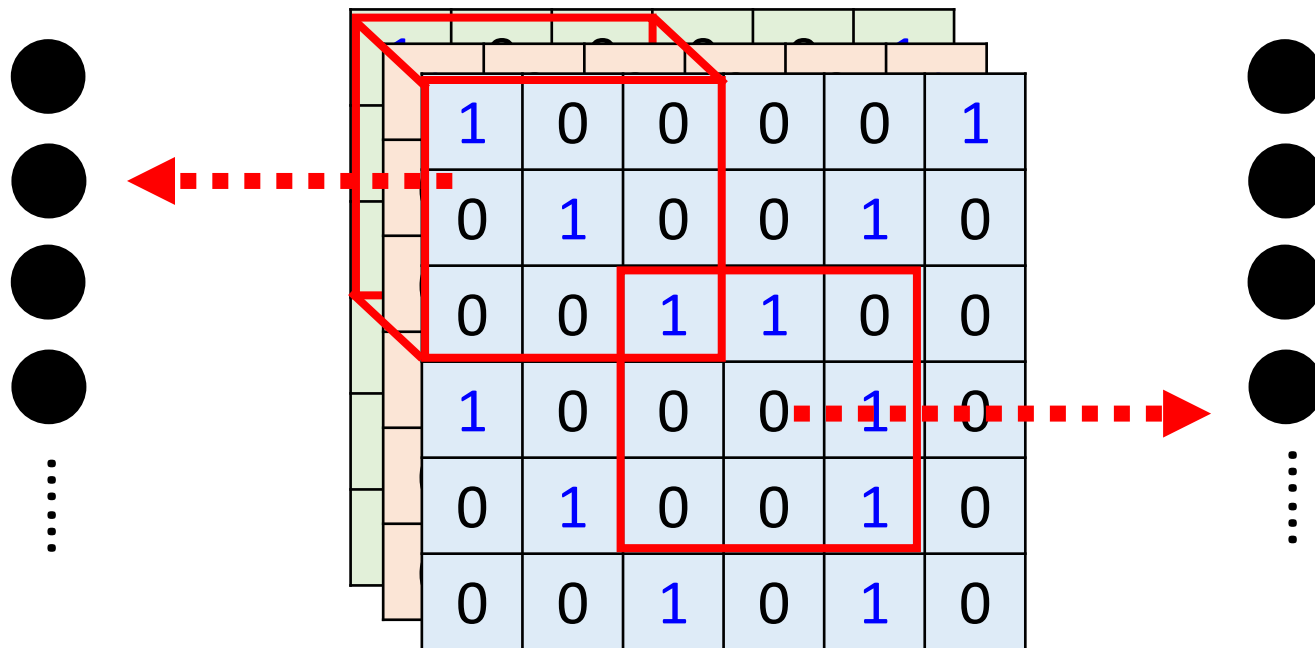| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

Two neurons with the same receptive field would not share parameters.
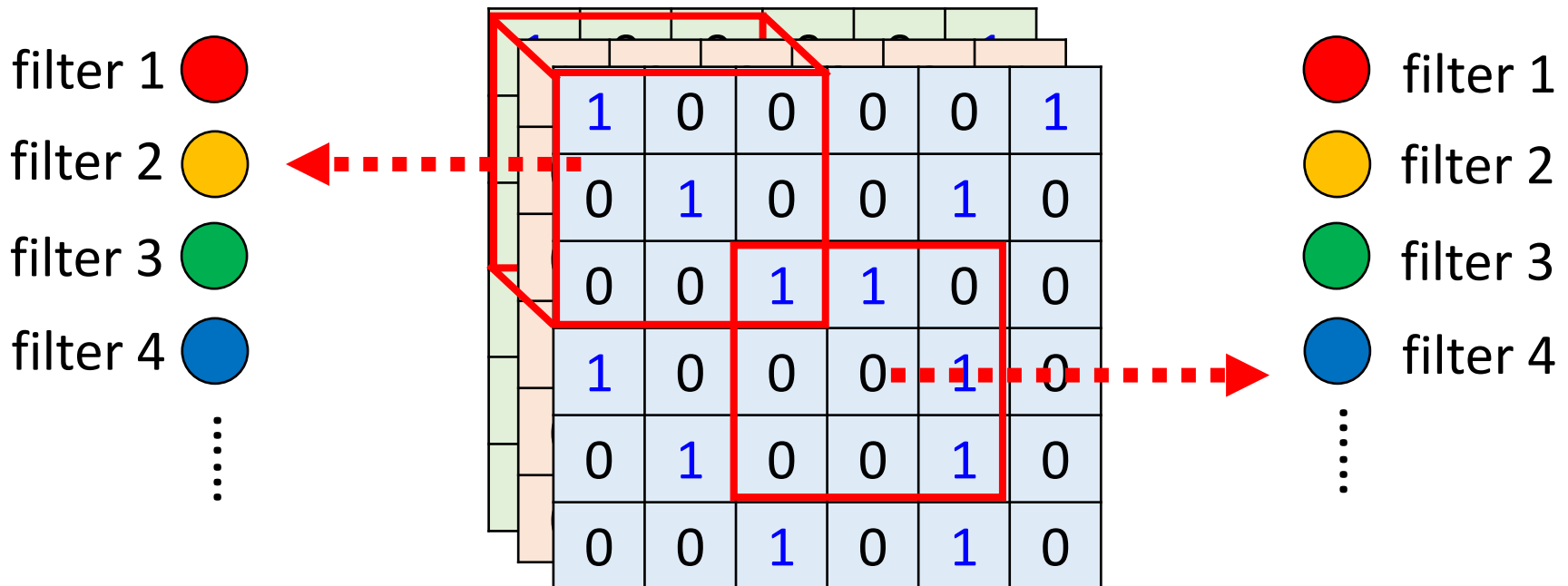
# Simplification 2 – Typical Setting

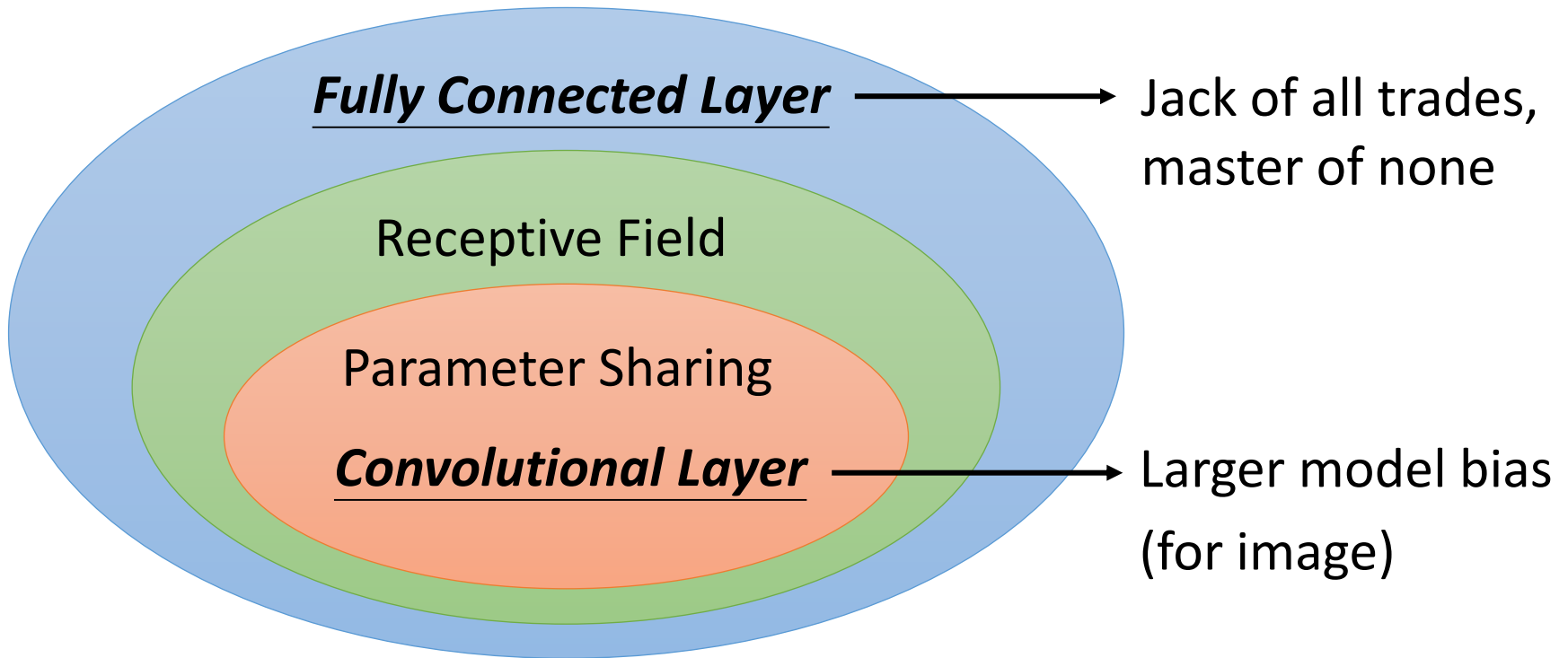Each receptive field has a set of neurons (e.g., 64 neurons).

# Simplification 2 – Typical Setting

Each receptive field has a set of neurons (e.g., 64 neurons).

Each receptive field has the neurons with the same set of parameters.
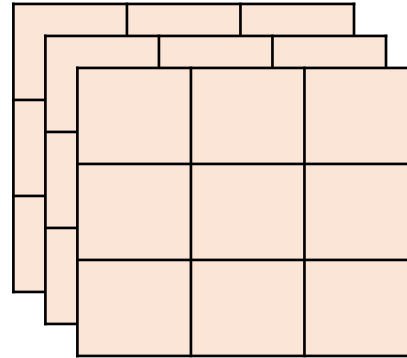
# Benefit of Convolutional Layer

Fully Connected Layer → Jack of all trades, master of none

Receptive Field

Parameter Sharing

Convolutional Layer → Larger model bias (for image)

- Some patterns are much smaller than the whole image.
- The same patterns appear in different regions.

# Convolutional Layer



channel = 3  (colorful)

channel = 1  (black and white)

Filter 1
3 x 3 x channel
tensor

Filter 2
3 x 3 x channel
tensor

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

28

# CNN – 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

⋮

**Property 1**  Each filter detects a small pattern (3 x 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

Filter 1

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 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 below

# CNN – Convolution

Filter 1

stride=1

6 x 6 image

Property 2

# Convolutional Layer

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 |

# CNN – Colorful image



Filter 1

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

Filter 2

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

Colorful 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 |

# Convolutional Layer



64 filters

Convolution

Convolution

"Image" with 64 channels

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

# *Multiple Convolutional Layers*

64 filters

**Convolution**

**Convolution**



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

"Image" with 64 channels

Filter:
3 x 3 x 64

64

# Multiple Convolutional Layers



64 filters

Convolution

Convolution

⋮

# *Convolution v.s. 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 |

image

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

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

convolution

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

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 |

$x_1$
$x_2$
$\vdots$
$x_{36}$

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
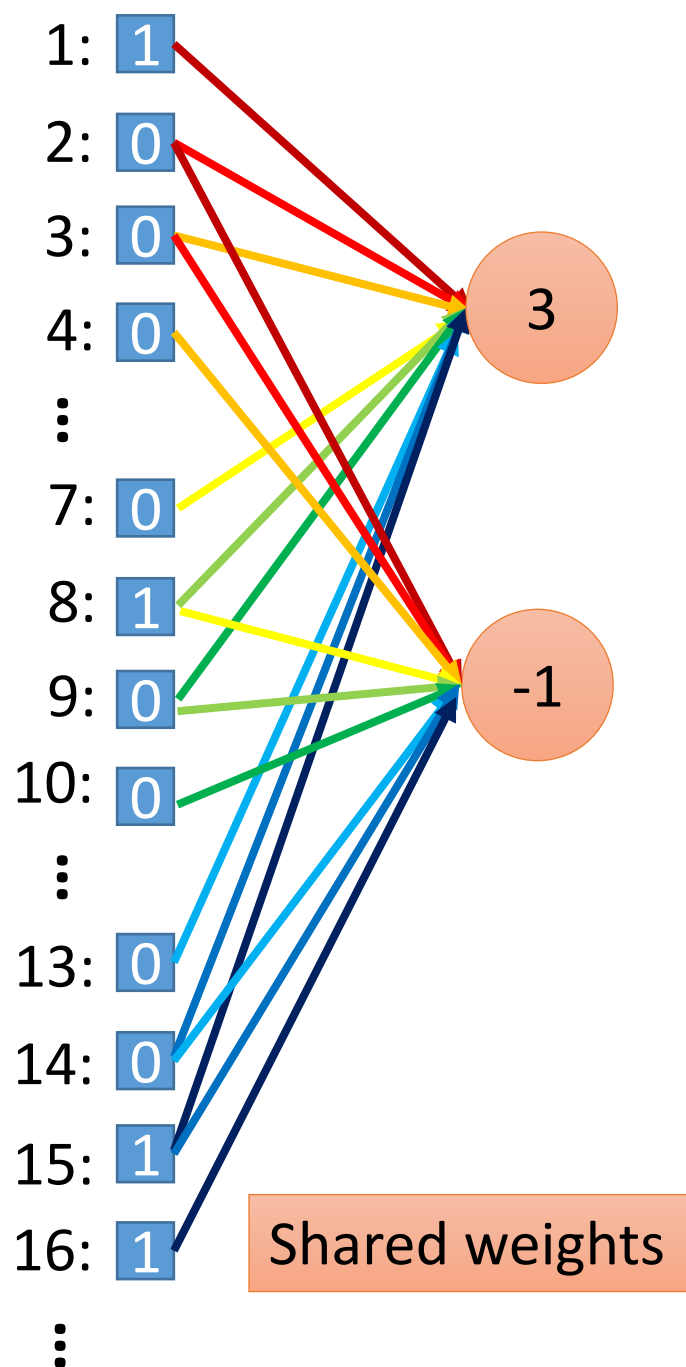15: 1
16: 1
⋮

3

Only connect to 9 input, not fully connected

Filter 1

6 x 6 image

Less parameters!

Even 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

-1

Shared weights

# Comparison of Two Stories
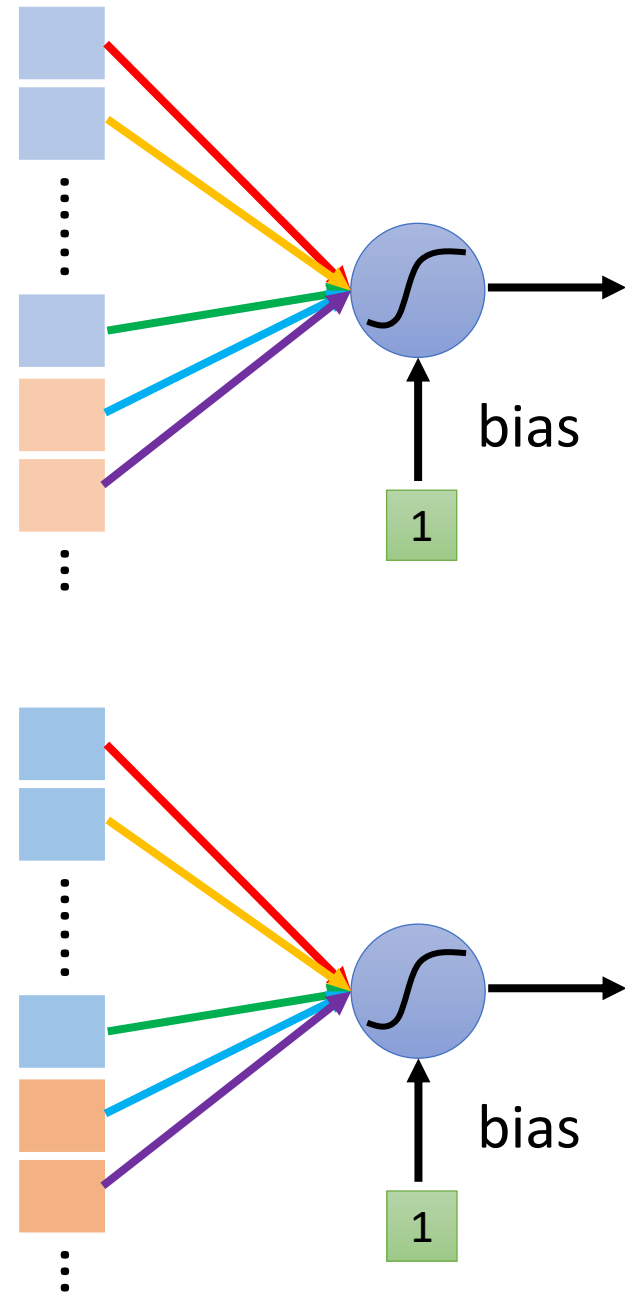


Receptive field

Filter
3 x 3 x channel tensor

(ignore bias in this slide)

41

The neurons with different receptive fields **share the parameters**.

| 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 |

bias

1

bias

1

**Each filter convolves over the input image.**

42

# Convolutional Layer

| ***Neuron Version Story*** | ***Filter Version Story*** |
|---|---|
| Each neuron only considers a receptive field. | There are a set of filters detecting small patterns. |
| The neurons with different receptive fields share the parameters. | Each filter convolves over the input image. |

They are the same story.

# Observation 3

- Subsampling the pixels will not change the object

bird



bird

subsampling



We can subsample the pixels to make image smaller

Less parameters for the network to process the image
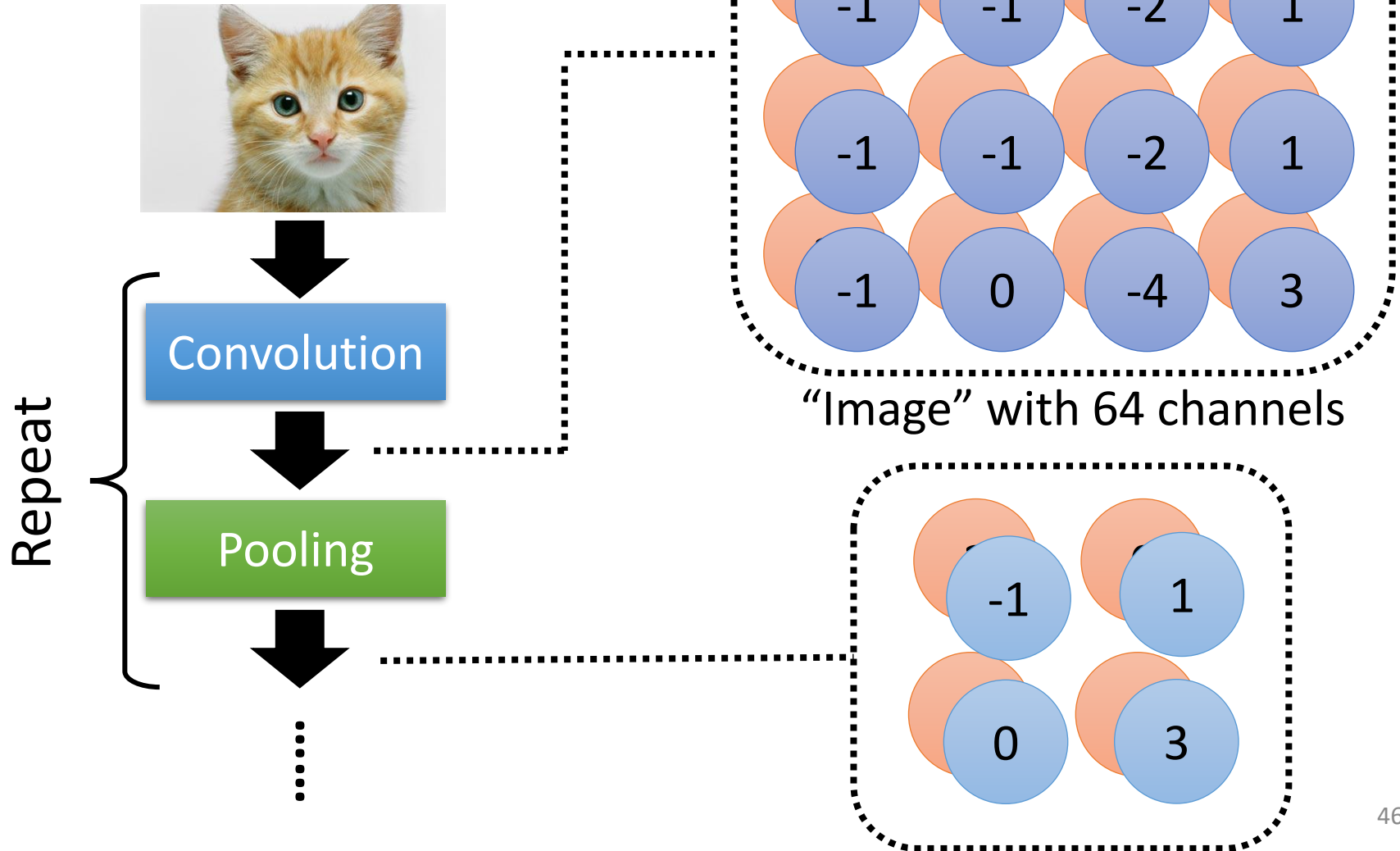
# Pooling – Max Pooling

| | | |
|---|---|---|
| 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 |
| 3 | -2 |

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

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

| | |
|---|---|
| -1 | -1 |
| -2 | 1 |

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

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

# Convolutional Layers + Pooling



"Image" with 64 channels

# The whole CNN



cat dog ......

softmax

Fully Connected
Feedforward network

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat
many times

Flatten

# The whole CNN



**Property 1**
➢ Some patterns are much smaller than the whole image

**Property 2**
➢ The same patterns appear in different regions.

**Property 3**
➢ Subsampling the pixels will not change the object
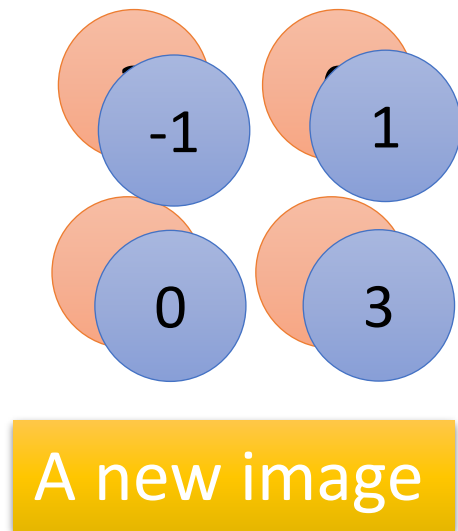
Convolution

Max Pooling

Convolution

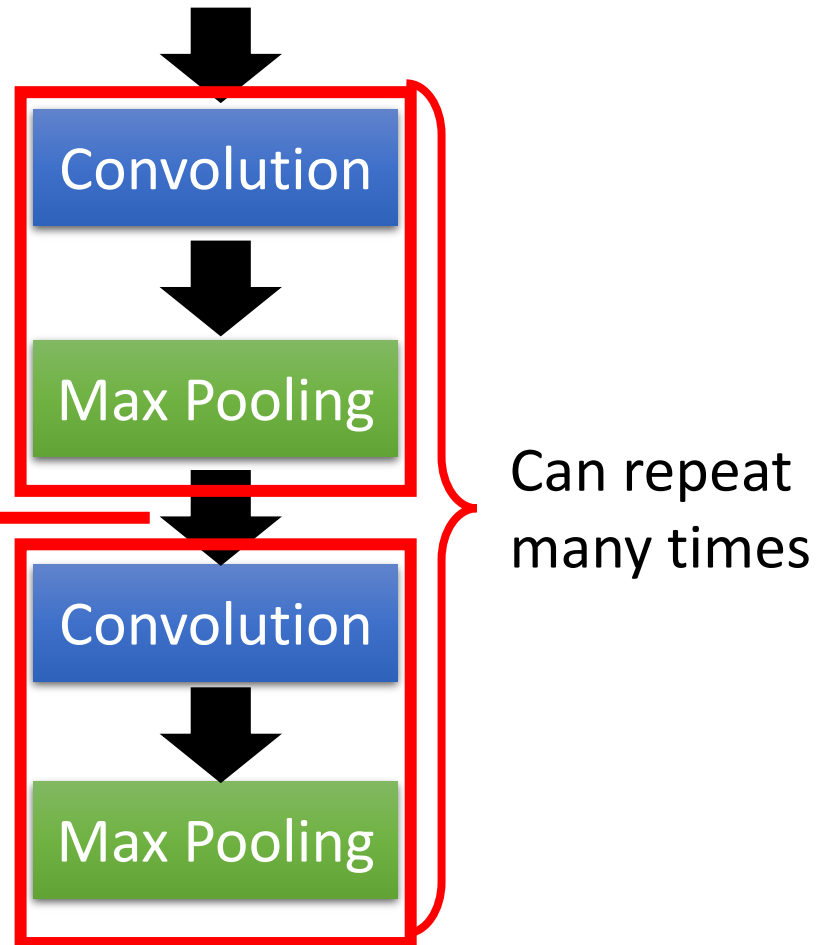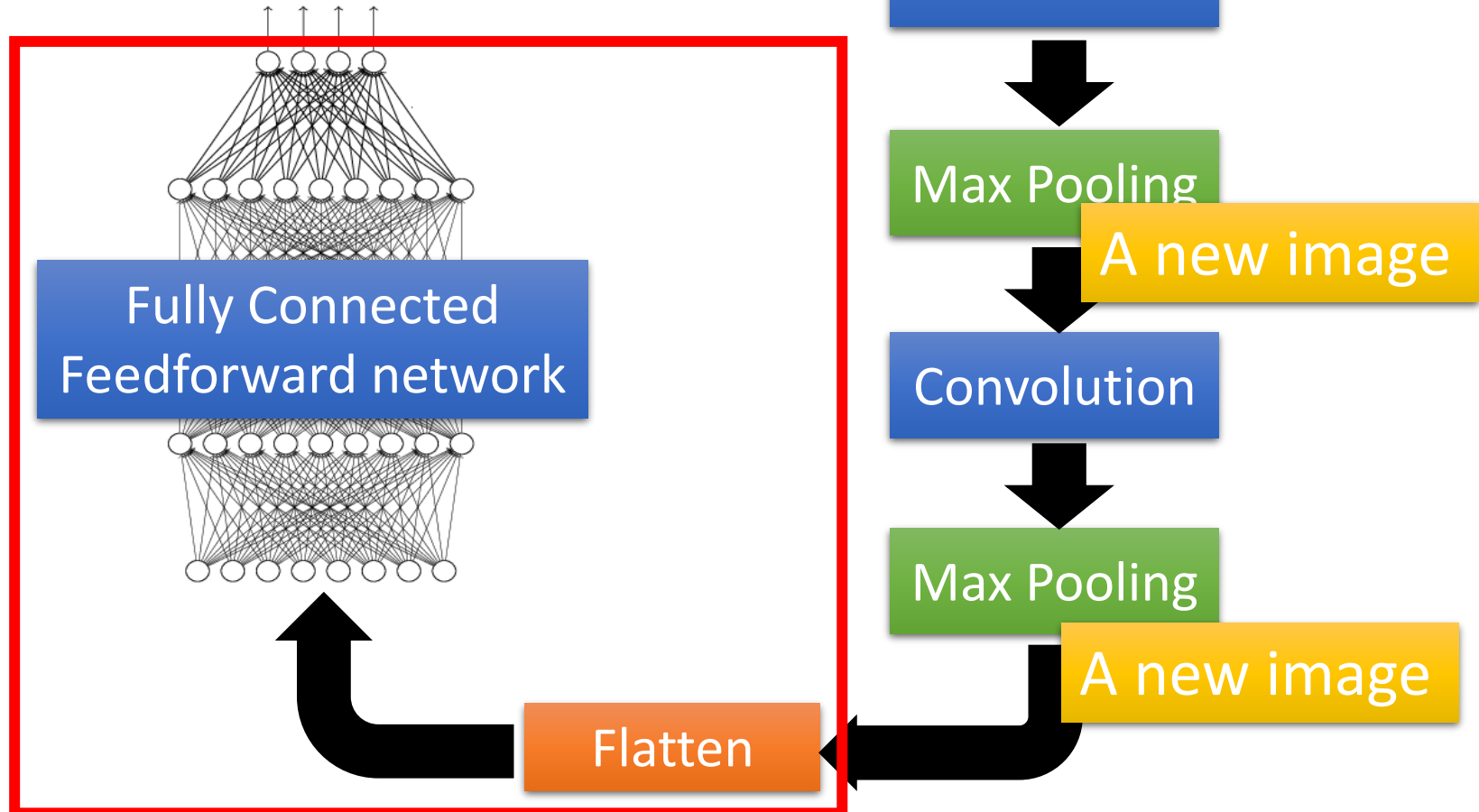Max Pooling

Flatten

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

# Flatten



Flatten

Fully Connected
Feedforward network

# Application: Playing Go



19 x 19 matrix (image)

Network → Next move (19 x 19 positions)

19 x 19 classes

48 channels in Alpha Go

Black: 1
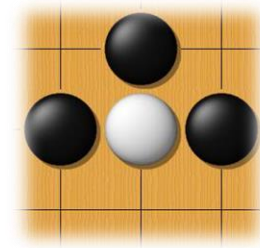
white: -1

none: 0

Fully-connected network can be used
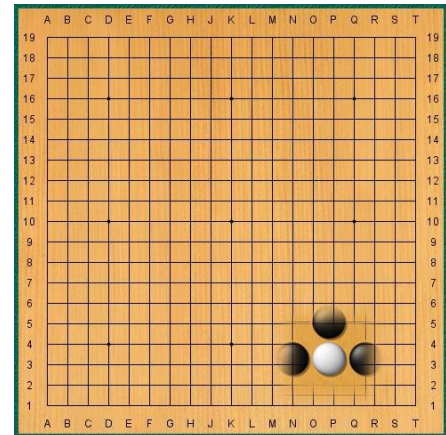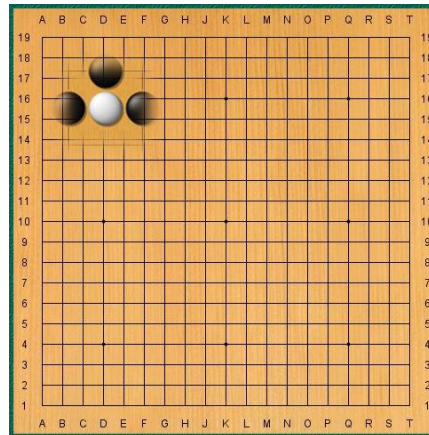
But CNN performs much better.

# Why CNN for Go playing?

- 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 Go playing?

- Subsampling the pixels will not change the object

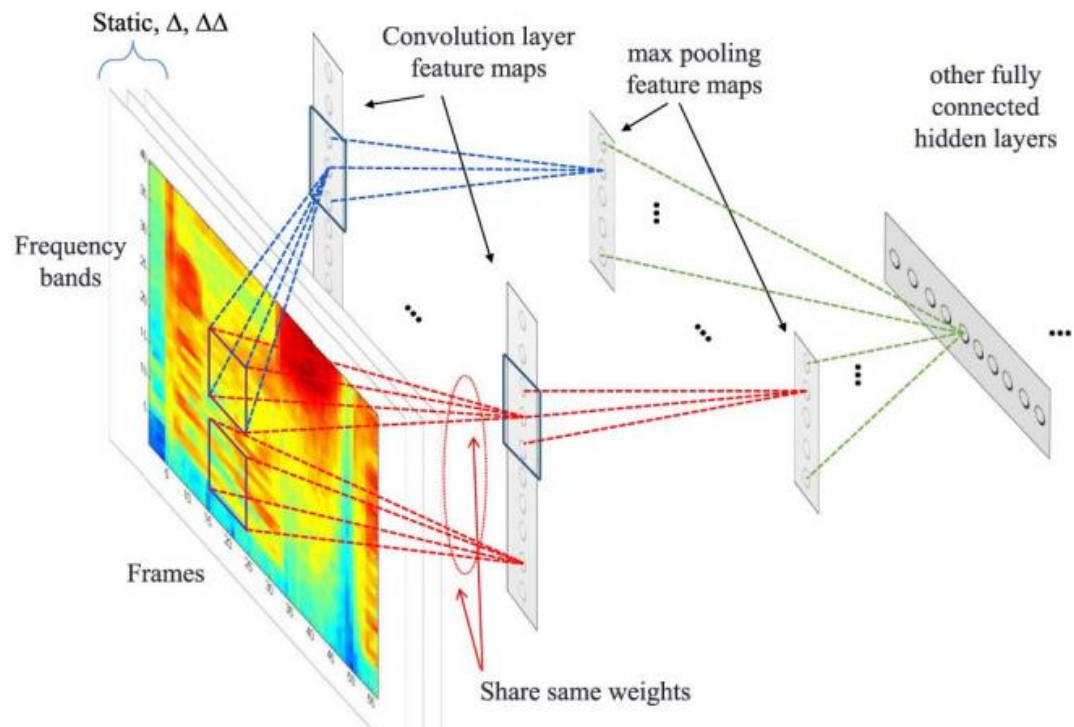Pooling    How to explain this???

**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 match version of AlphaGo used $k = 192$ filters; Fig. 2b and Extended Data Tabl 256 and 384 filters

Alpha Go does not use Pooling ......

# *More Applications*



## Speech

https://dl.acm.org/doi/10.1109/TASLP.2014.2339736

## Natural Language Processing

https://www.aclweb.org/anthology/S15-2079/