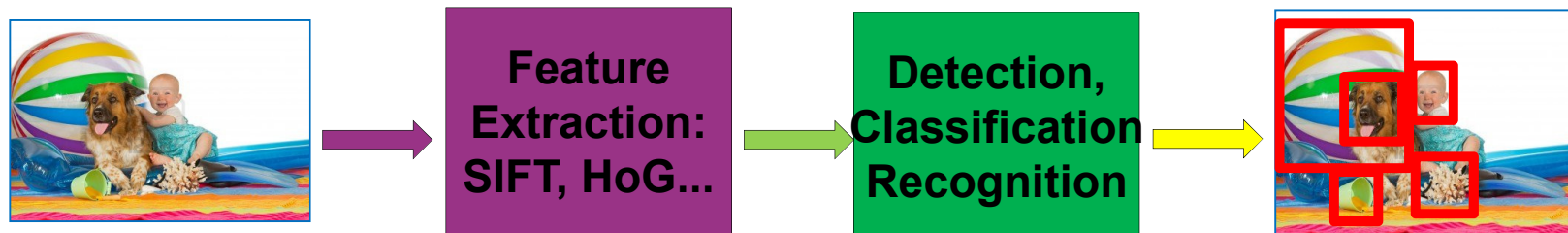Image Source: deeplearning.ai

# Computer Vision Problems

# Classical Computer Vision Pipeline

CV experts

1. Select / develop features: SURF, HoG, SIFT, RIFT, ...

2. Add on top of this Machine Learning for multi-class recognition and train classifier
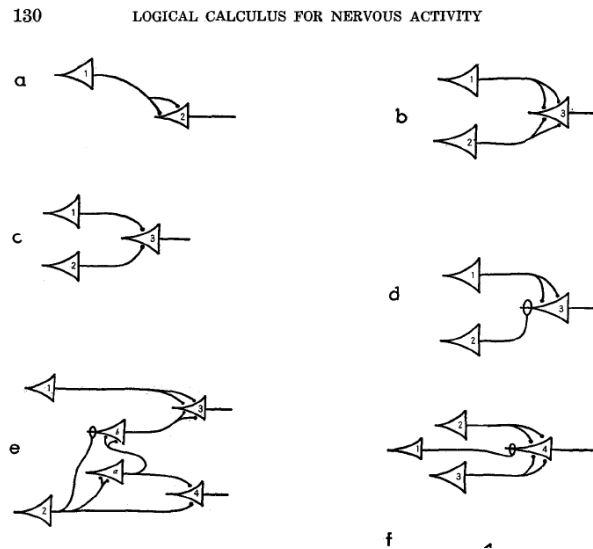


**Feature Extraction: SIFT, HoG...** → **Detection, Classification Recognition** →

Classical CV feature definition is domain-specific and time-consuming

# Neural Network



Warren McCulloch

Walter Pitts

A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY, Bulletin of Mathematical Biophysics, Vol. 5, pp. 115-133 (1943).

# Neural Network

Here x1 and x2 are normalized attribute value of data.

y is the output of the neuron , i.e the class label.

x1 and x2  values multiplied by weight values w1 and w2 are input to the neuron x.

Value of x1 is multiplied by a weight w1 and values of x2 is multiplied by a weight w2.

Given that



Fig1: an artificial neuron

- w1 = 0.5 and w2 = 0.5
- Say value of x1 is 0.3 and value of x2 is 0.8,

- So, weighted sum is :
- sum= w1 x x1 + w2 x x2 = 0.5 x 0.3 + 0.5 x 0.8 = 0.55

# Why We Need Multi Layer ?

Linear Separable:

$$x \vee y$$

$$x \wedge y$$

Linear inseparable:

$$x \veebar y$$

# Edge Detection



Vertical edges

Horizontal edges

How do we detect these edges

# Neural Network?

❑Suppose an image is of the size 68 X 68 X 3
  ◦ Input feature dimension then becomes 12,288

❑If Image size is of 720 X 720 X 3
  ◦ Input feature dimension becomes 1,555,200

❑Number of parameters will swell up to a HUGE number

❑Result in more computational and memory
requirements

# Another Application

Digit Recognition



$X_1,...,X_n \in \{0,1\}$ (Black vs. White pixels)

$Y \in \{5,6\}$ (predict whether a digit is a 5 or a 6)

# The Bayes Classifier

In class, we saw that a good strategy is to predict:

$$\arg\max_Y P(Y|X_1, \ldots, X_n)$$

- (for example: what is the probability that the image represents a 5 given its pixels?)

So ... how do we compute that?

# The Bayes Classifier

Use Bayes Rule!

Likelihood           Prior

$$P(Y|X_1, \ldots, X_n) = \frac{P(X_1, \ldots, X_n|Y)P(Y)}{P(X_1, \ldots, X_n)}$$

Normalization Constant

Why did this help?  Well, we think that we might be able to specify how features are "generated" by the class label

# The Bayes Classifier

Let's expand this for our digit recognition task:

$$P(Y = 5 | X_1, \ldots, X_n) = \frac{P(X_1, \ldots, X_n | Y = 5) P(Y = 5)}{P(X_1, \ldots, X_n | Y = 5) P(Y = 5) + P(X_1, \ldots, X_n | Y = 6) P(Y = 6)}$$

$$P(Y = 6 | X_1, \ldots, X_n) = \frac{P(X_1, \ldots, X_n | Y = 6) P(Y = 6)}{P(X_1, \ldots, X_n | Y = 5) P(Y = 5) + P(X_1, \ldots, X_n | Y = 6) P(Y = 6)}$$

To classify, we'll simply compute these two probabilities and predict based on which one is greater

# Model Parameters

For the Bayes classifier, we need to "learn" two functions, the likelihood and the prior

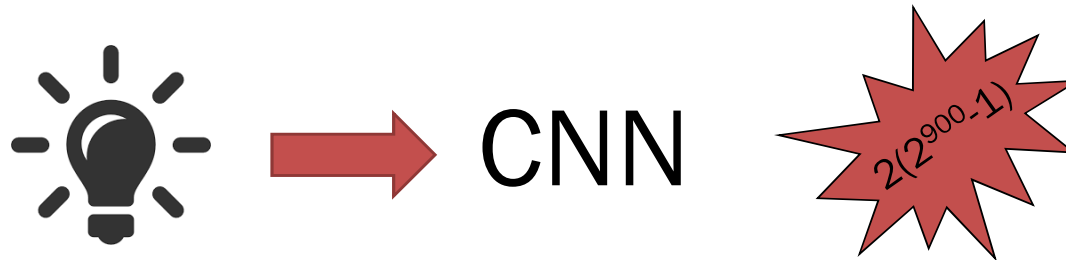How many parameters are required to specify the prior for our digit recognition example?

Just 1

# Model Parameters

How many parameters are required to specify the likelihood?

  ◦ (Supposing that each image is 30x30 pixels)



CNN

$2(2^{900}\text{-}1)$
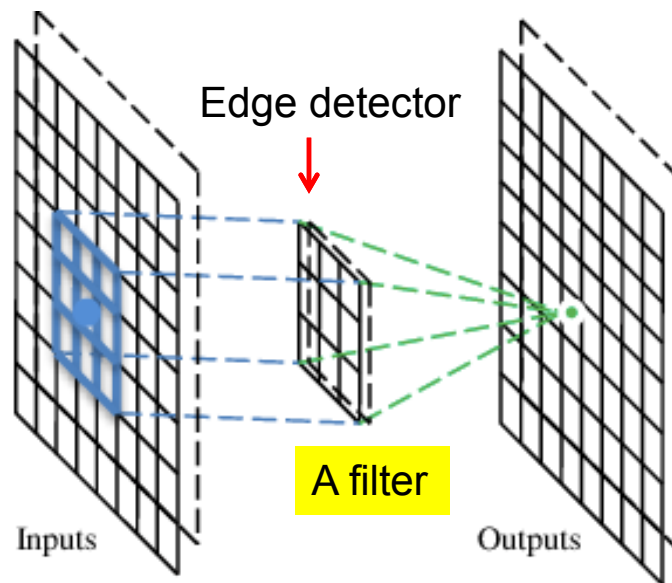
# Drive into CNN

In a convolutional network (ConvNet), there are basically three types of layers:
1. Convolution layer
2. Pooling layer
3. Fully connected layer

# A convolutional layer

A CNN is a neural network with some convolutional layers (and some other layers). A convolutional layer has a number of filters that does convolutional operation.

Edge detector

A filter

Inputs

Outputs

# Convolution

**These 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

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

Filter 2

⋮ ⋮

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

Image Source: internet

# 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 |

Dot product → 3  -1

6 x 6 image

# 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

# 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 |

# 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

Repeat this for each filter

| -1 | -1 | -1 | -1 |
|----|----|----|----|
| -1 | Feature Map | | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

Two 4 x 4 images
Forming 4 x 4 x 2
matrix

# Convolution over Volume

Color 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 |

Filter 1

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

Filter 2

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

# Convolution v.s. Fully Connected

Filter 1

6 x 6 image

fewer parameters!

1
2
3
4:
⋮

8
9
10:
⋮

1
3
14
15
16
⋮

3

Only connect to 9 inputs, not fully connected

Image Source: internet

Filter 1

6 x 6 image

Fewer parameters

Even fewer parameters

Shared weights

Image Source: internet

Suppose we have 10 filters applying on input (6 X 6 X 3), each of shape 3 X 3 X 3. What will be the number of parameters in that layer?

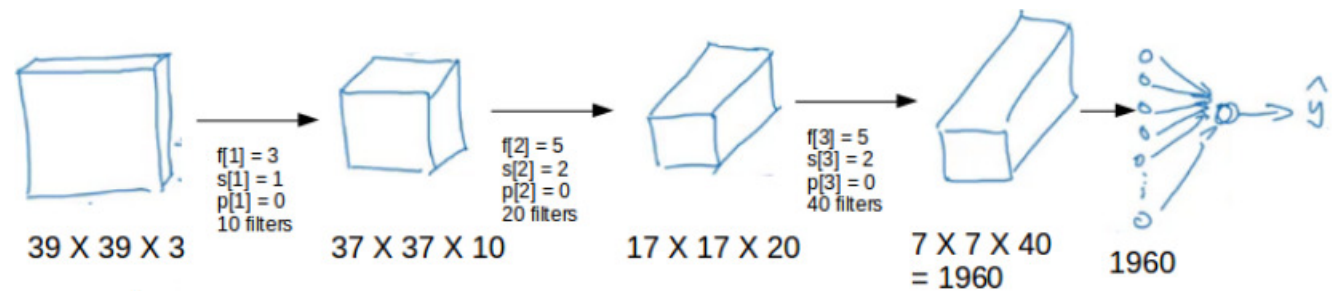- Number of parameters for each filter = 3*3*3 = 27

- There will be a bias term for each filter, so total parameters per filter = 28

- As there are 10 filters, the total parameters for that layer = 28*10 = 280

QUIZ

# Simple Convolutional Neural Network



39 X 39 X 3

f[1] = 3
s[1] = 1
p[1] = 0
10 filters

37 X 37 X 10

f[2] = 5
s[2] = 2
p[2] = 0
20 filters

17 X 17 X 20

f[3] = 5
s[3] = 2
p[3] = 0
40 filters

7 X 7 X 40
= 1960

1960

- **Size of feature vector : (n+2p-f)/s +1**
- **n : dimension of matrix**
- **p : size of padding**
- **f : size of filter**
- **s : size of stride**

Image Source: deeplearning.ai

The whole CNN

cat dog ……

Fully Connected Feedforward network

Flattened

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

Image Source: internet

# 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 |

# Why Pooling

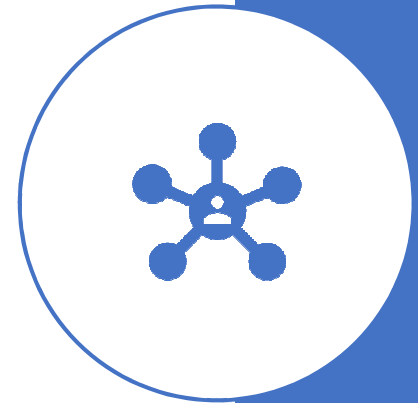- Subsampling pixels will not change the object
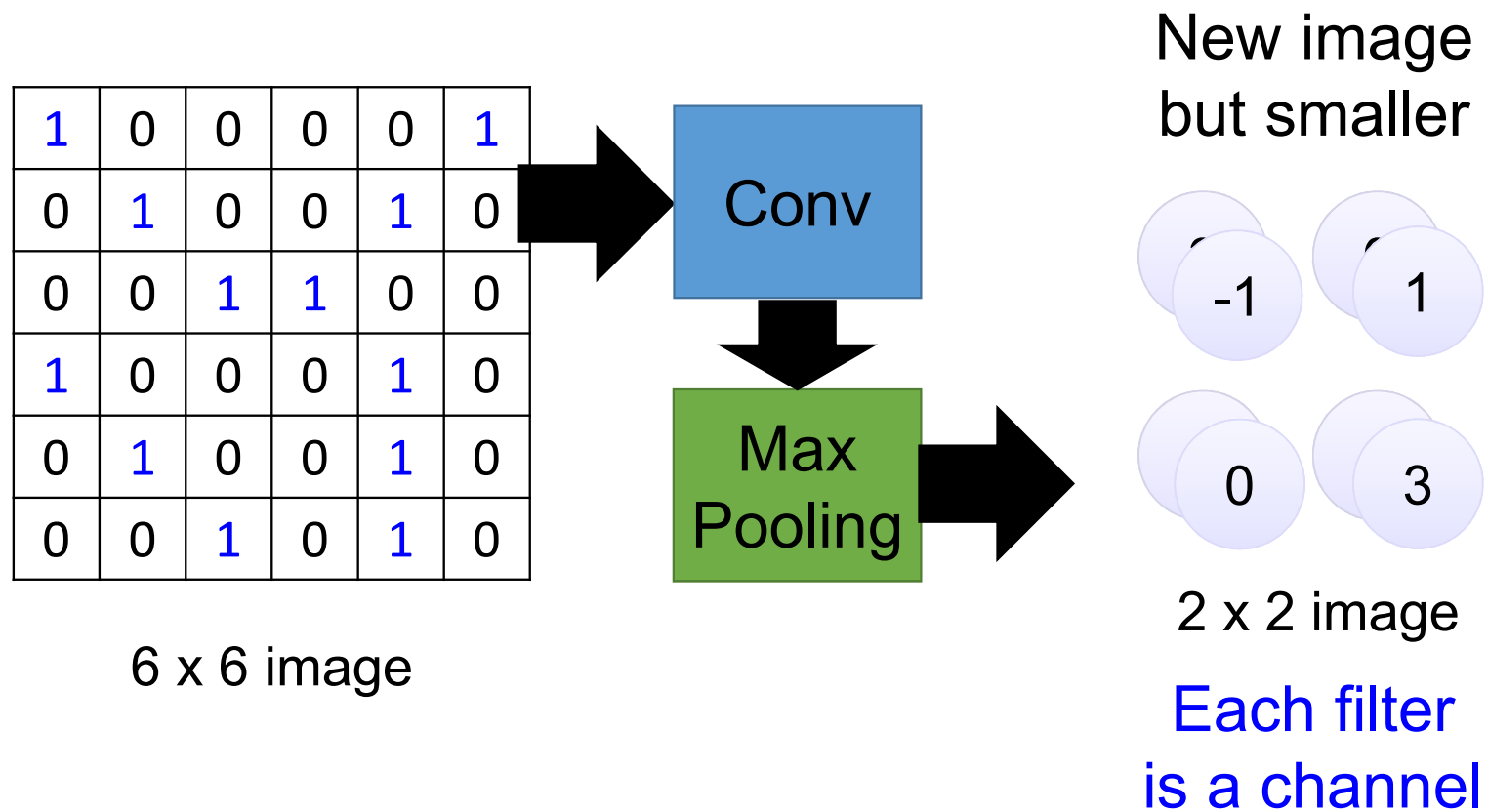
bird



bird

**Subsampling**

We can subsample the pixels to make image smaller → fewer parameters to characterize the image

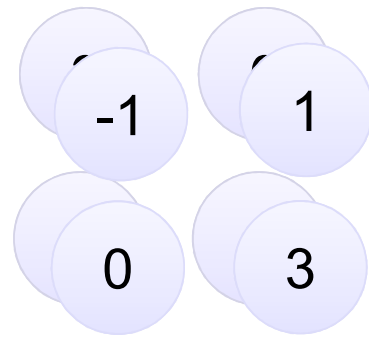# A CNN compresses a fully connected network in two ways:

- Reducing number of connections
- Shared weights on the edges
- Max pooling further reduces the complexity

# 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**

New image
but smaller

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

2 x 2 image

Each filter
is a channel
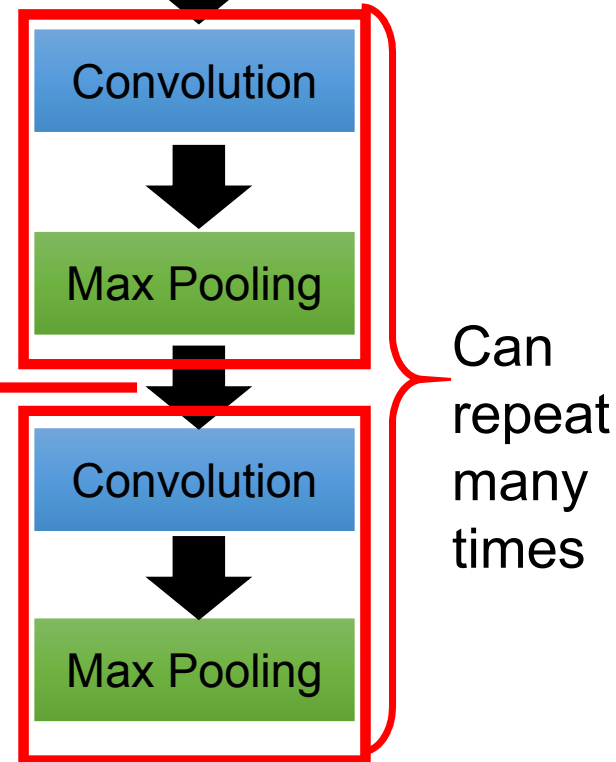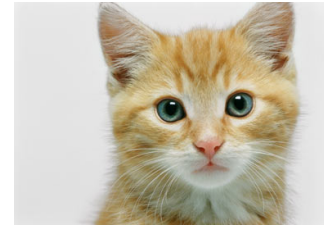
The whole CNN
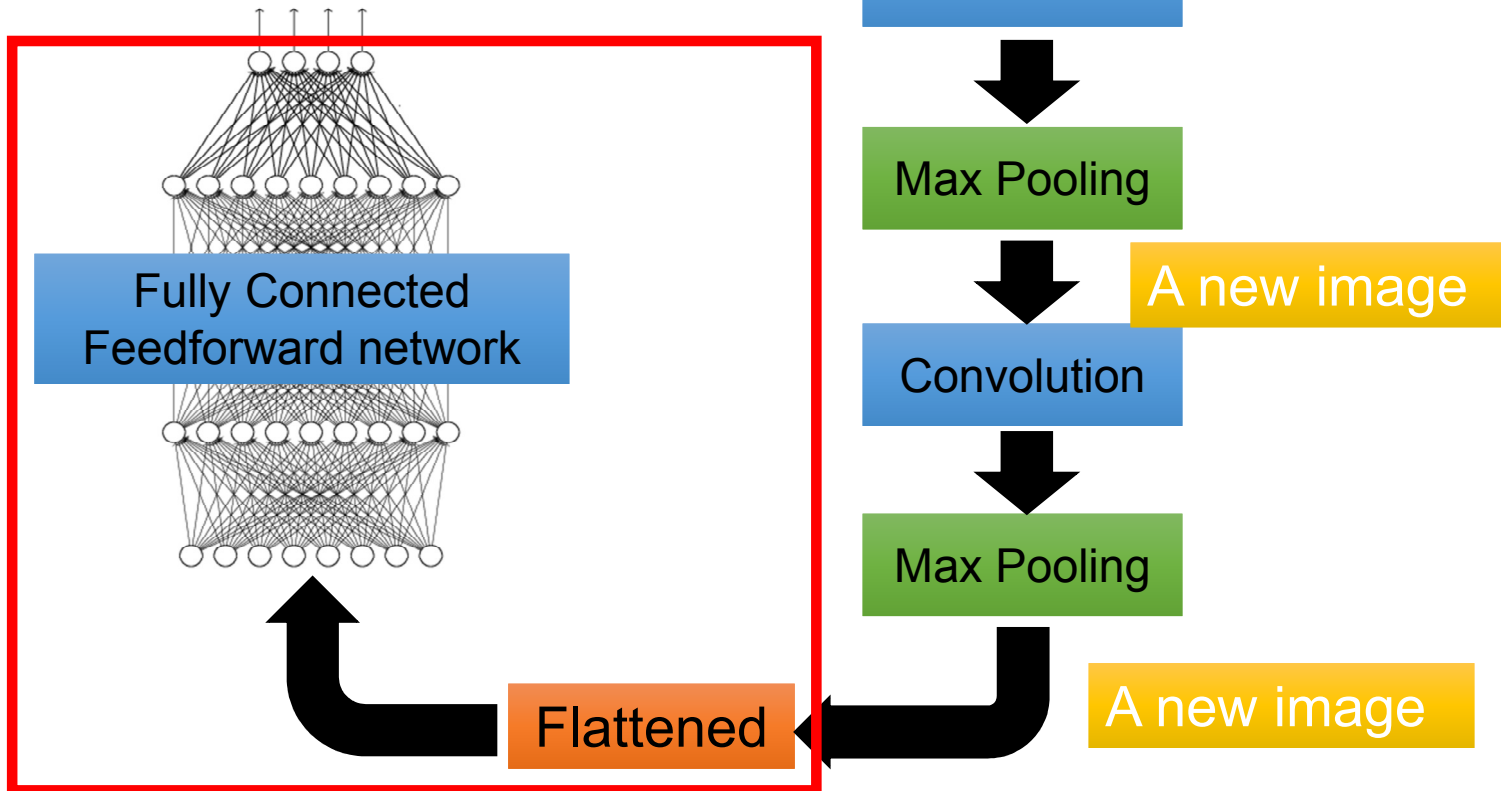
Smaller than the original image

The number of channels is the number of filters

The whole CNN

cat dog ……

Fully Connected Feedforward network

Flattened

Convolution

Max Pooling

A new image

Convolution

Max Pooling

A new image

Image Source: internet

# Flattening



Flattened

Fully Connected
Feedforward network

Image Source: internet

# Classic Networks

1. LeNet-5
2. AlexNet
3. VGG

# LeNet-5



- **Parameters:** 60k
- **Layers flow:** Conv → Pool → Conv → Pool → FC → FC → Output
- **Activation functions:** Sigmoid/tanh and ReLU

# AlexNet



- **Parameters:** 60 million
- **Activation functions:** ReLU

# VGG-16



224×224 ×3 → [CONV 64] ×2 → 224×224×64 → POOL → 112×112 ×64 → [CONV 128] ×2 → 112×112 ×128 → POOL → 56×56 ×128

→ [CONV 256] ×3 → 56×56 ×256 → POOL → 28×28 ×256 → [CONV 512] × 3 → 28×28 ×512 → POOL → 14×14×512

→ [CONV 512] ×3 → 14×14 ×512 → POOL → 7×7×512 → FC 4096 → FC 4096 → Softmax 1000

- **Parameters:** 138 million
- **Pool:** MAX with stride 2
- **CONV layer**: stride 1

# CNN in Keras

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

Input

1 x 28 x 28

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

Convolution

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

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

Max Pooling

25 x 13 x 13

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

Convolution

How many parameters for each filter?  225= 25x9  | 50 x 11 x 11

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

Max Pooling

50 x 5 x 5

Image Source: internet

# CNN in Keras

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

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

Flattened

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

1250

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'))
```