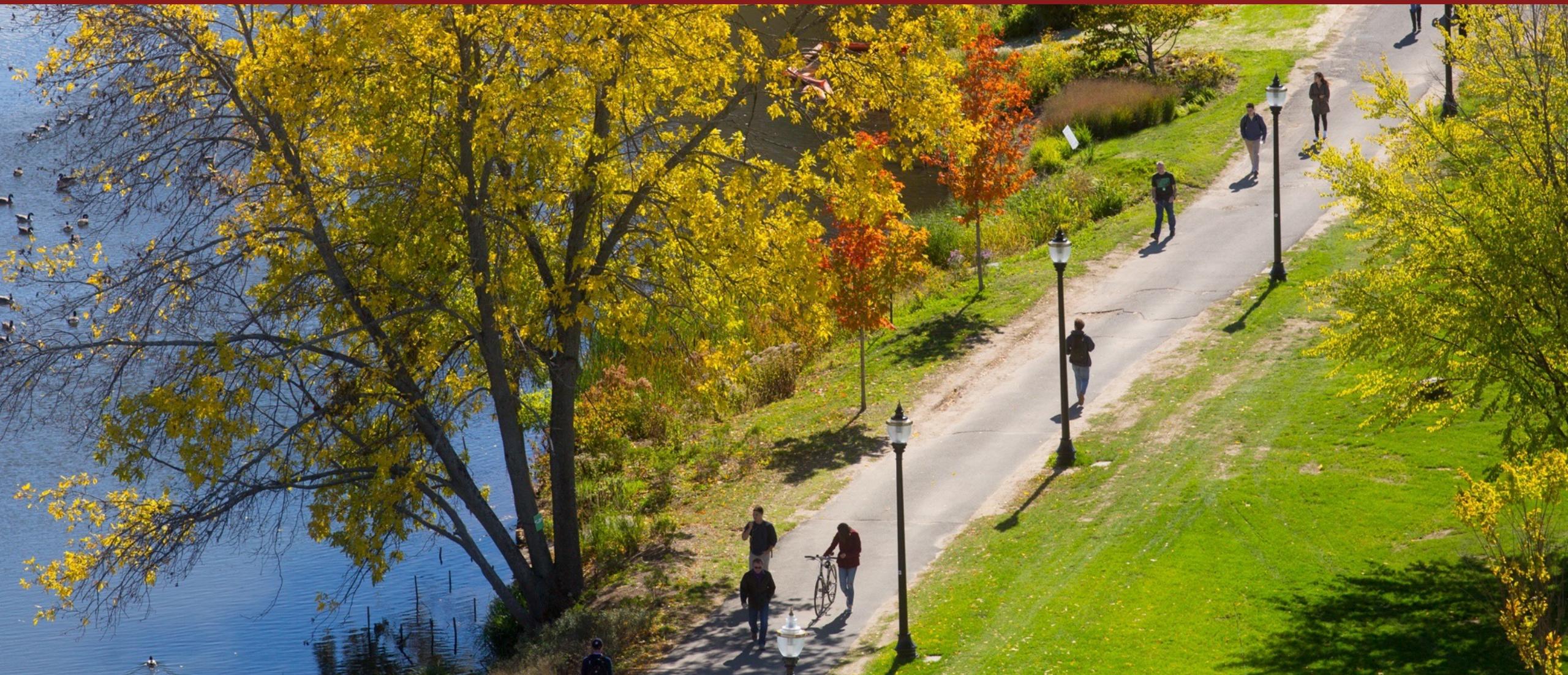


Digital Image Processing ECE 566

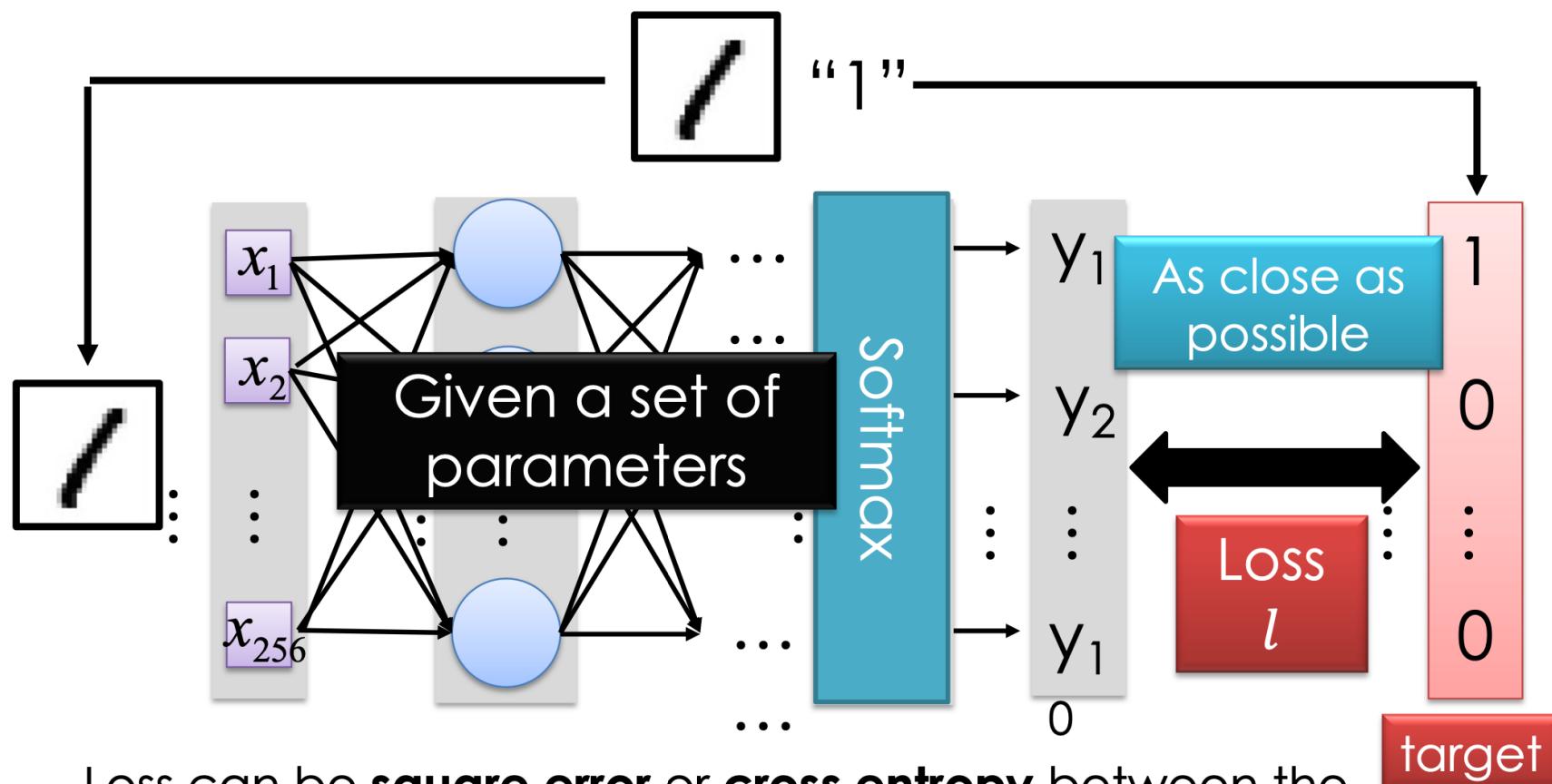
Ahmad Ghasemi

Department of Electrical and Computer Engineering



Loss

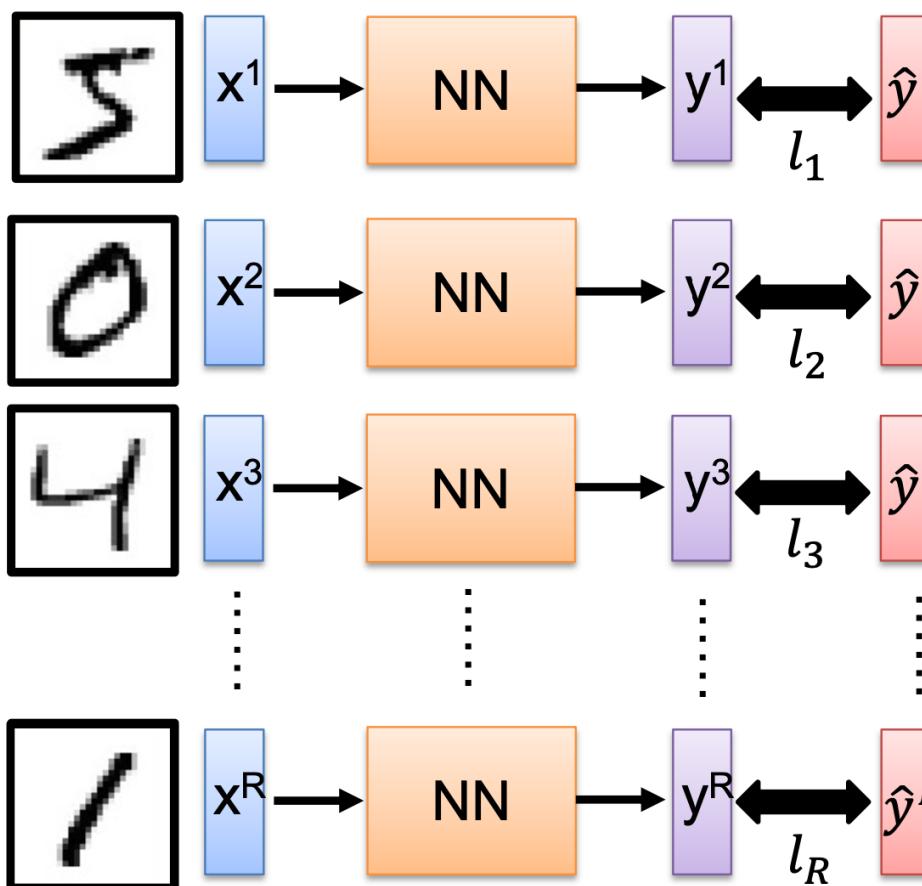
A good function should make the loss of all examples as small as possible.



Loss can be **square error** or **cross entropy** between the network output and target

Total Loss

For all training data ...



Total Loss:

$$L = \sum_{r=1}^R l_r$$

As small as possible

Find a function in function set that minimizes total loss L

Find the network parameters θ^* that minimize total loss L

Three steps for Deep Learning

Step 1: define a set of function



Step 2: goodness of function



Step 3: pick the best function

Testing or Inference

Convolutional Neural Network (CNN)

Widely used in image
processing and computer vision

Why CNN for Images?

- Some patterns are much smaller than the whole image

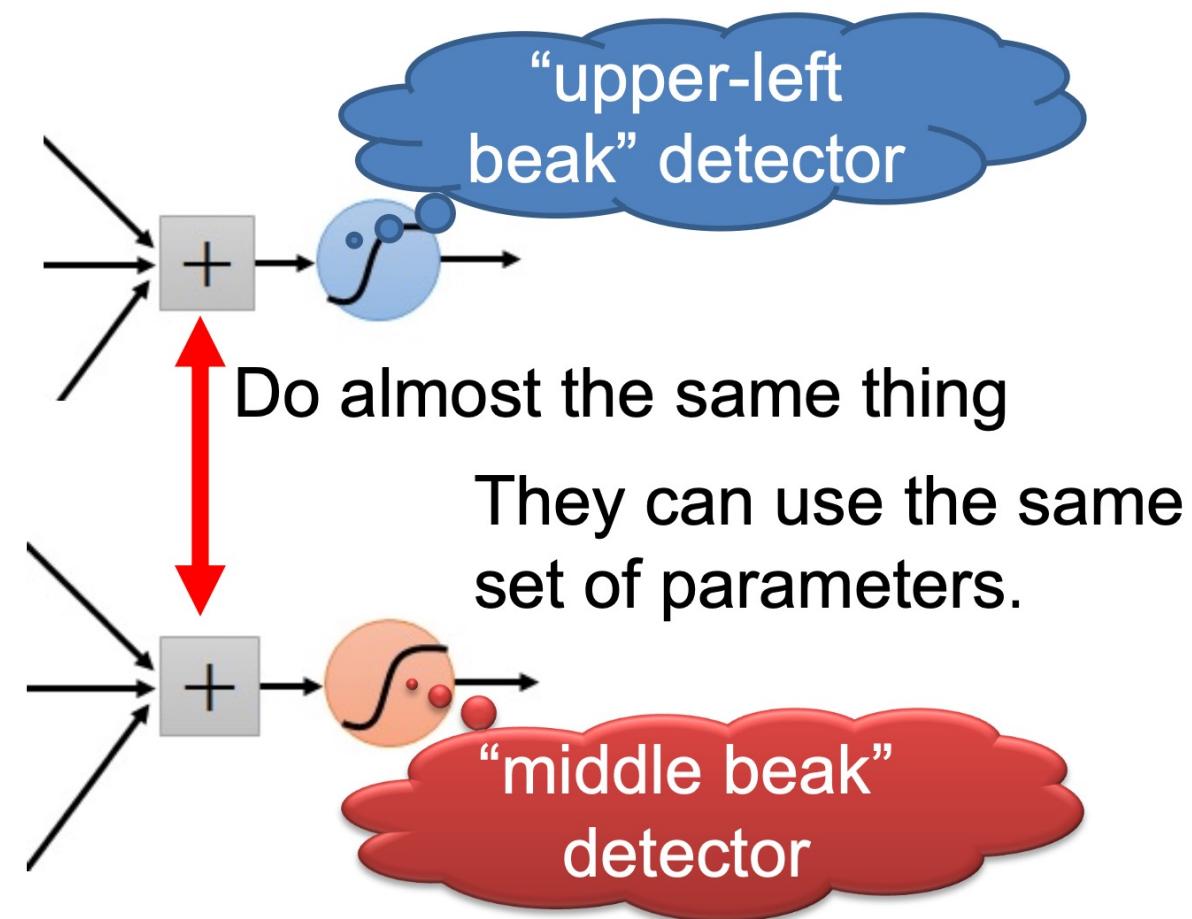
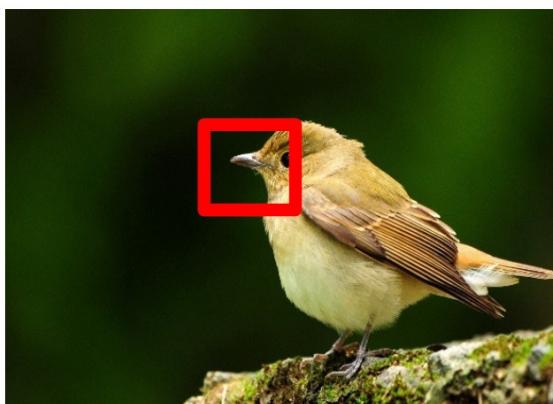
A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters



Why CNN for Images?

- The same patterns appear in different regions.



Why CNN for Images?

- Subsampling the pixels will not change the object

bird



subsampling

bird



We can subsample the pixels to make image smaller



Less parameters for the network to process the image

Why CNN for Images?

Property 1

- Some patterns are much smaller than the whole

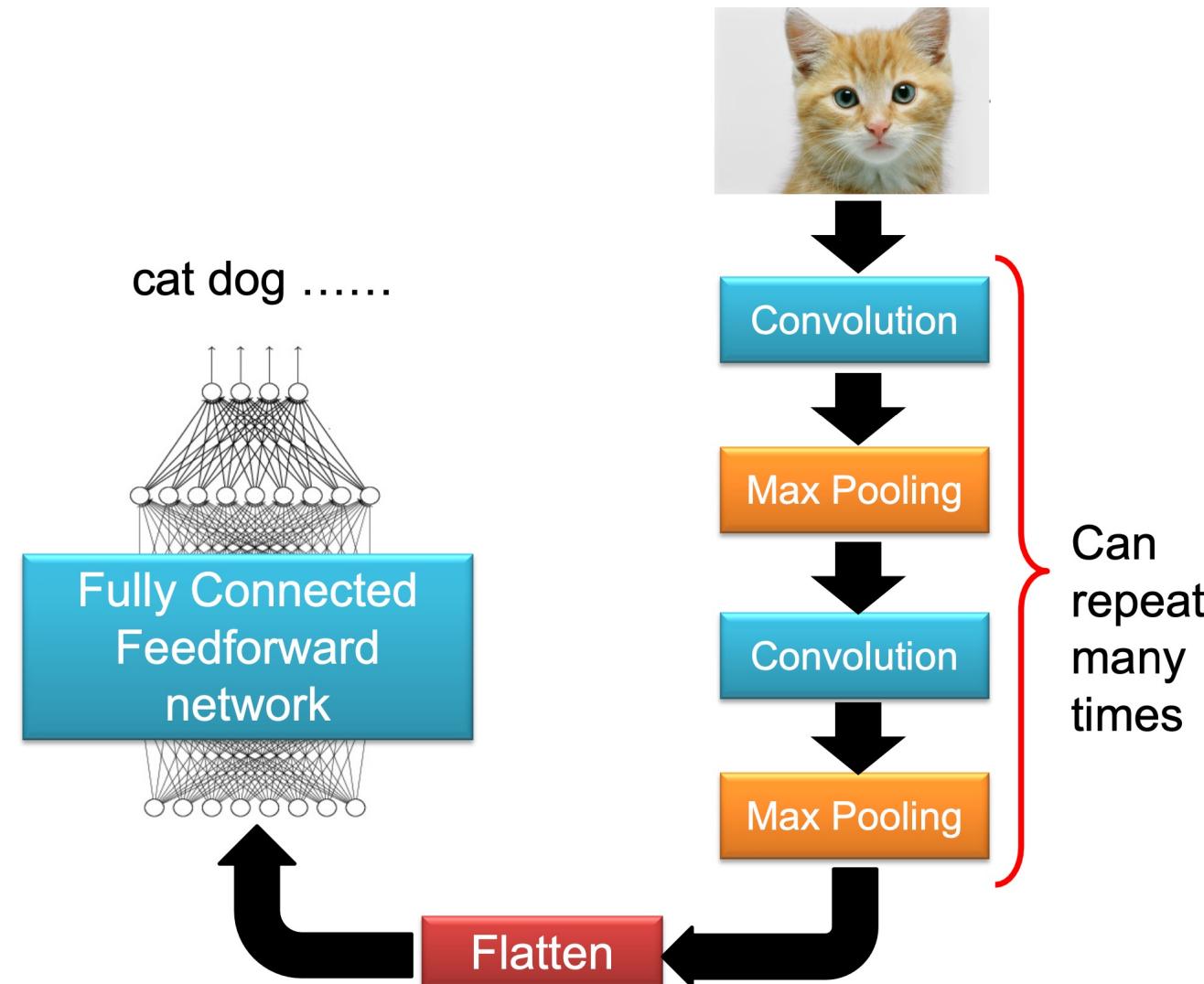
Property 2

- The same patterns appear in different regions.

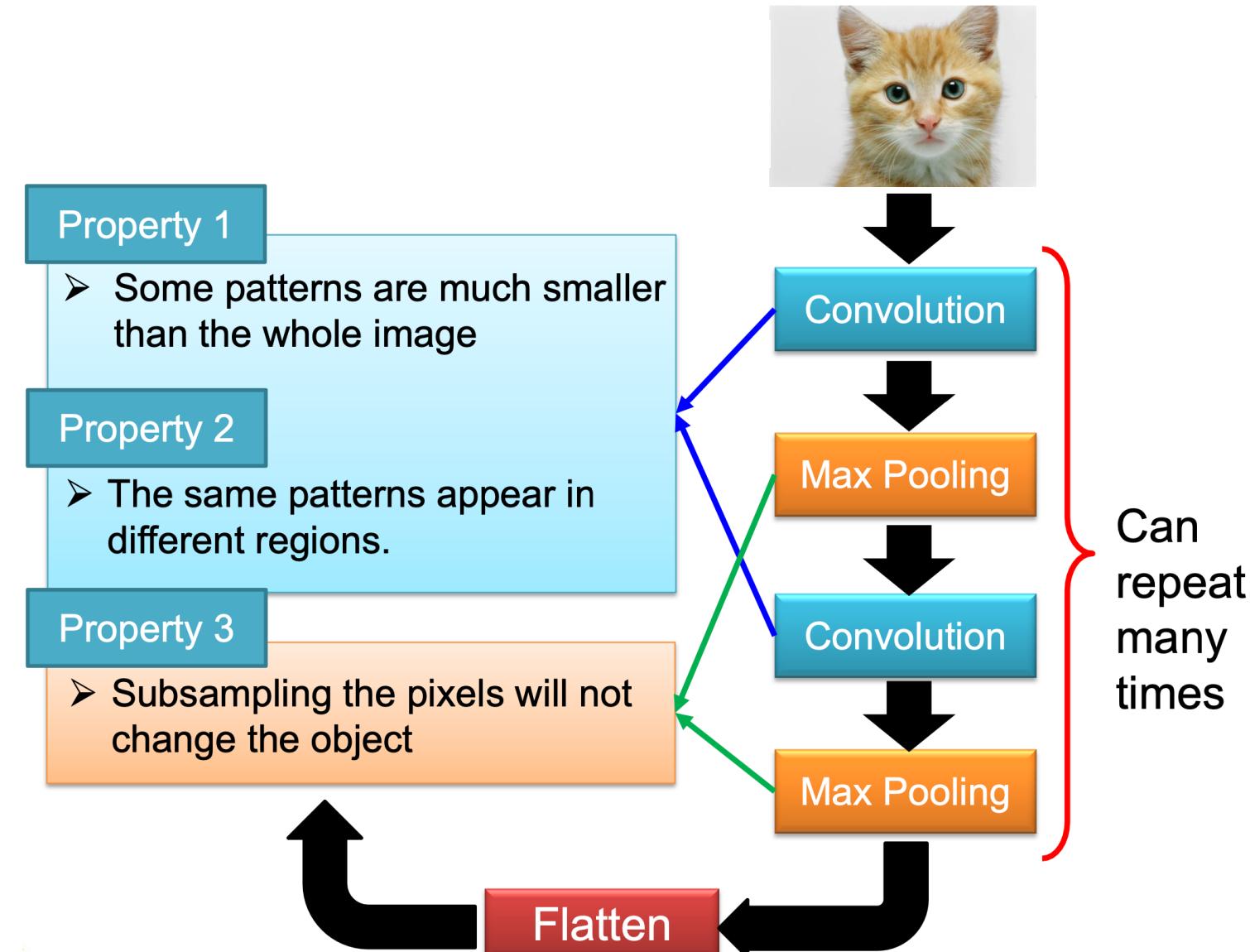
Property 3

- Subsampling the pixels will not change the object

The Whole CNN



The Whole CNN



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/kernel 1

Matrix

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

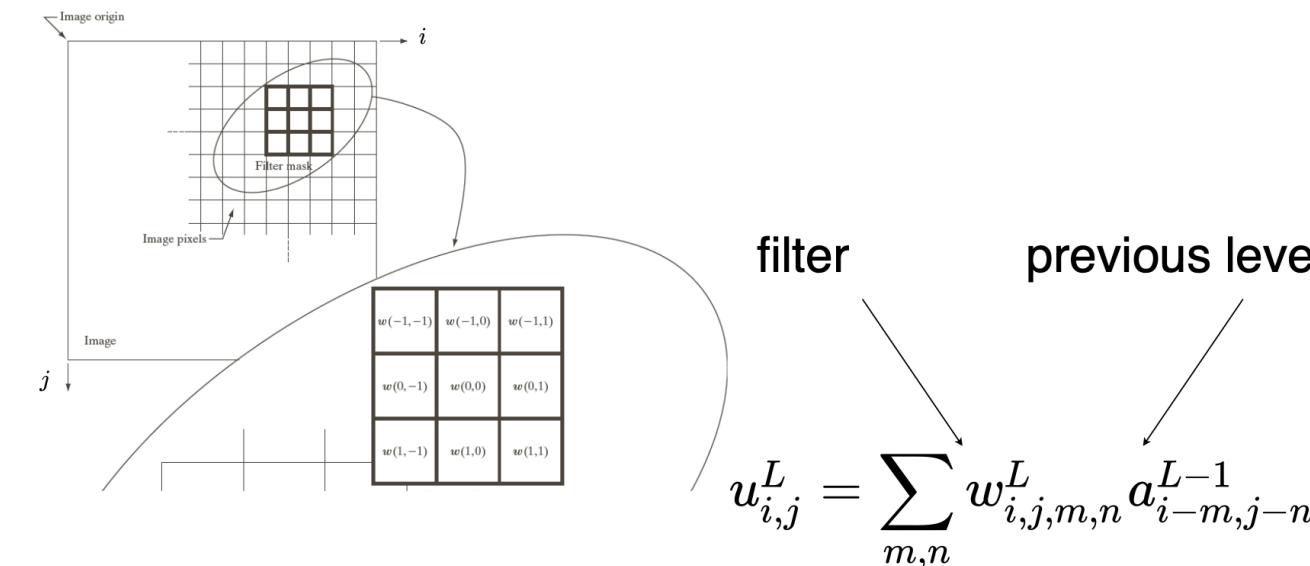
Filter 2

Matrix

⋮ ⋮

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

CNN: Convolution



Filter/Kernel 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

3 -1

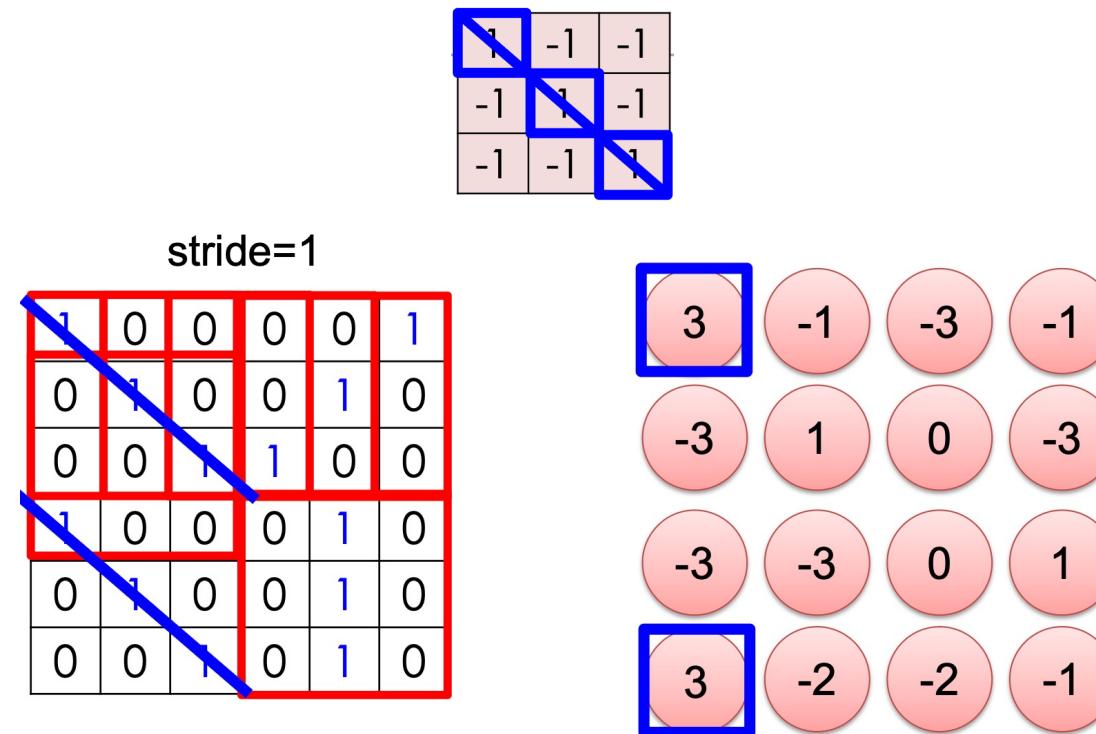
If stride=2

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

3 -3

We set stride=1 below

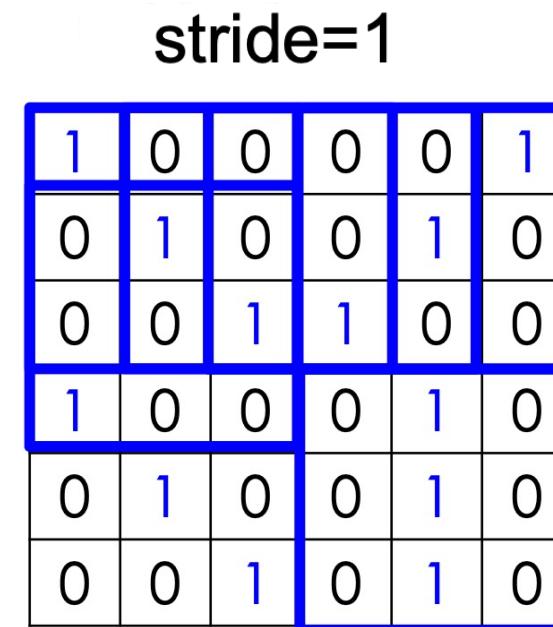
CNN: Convolution



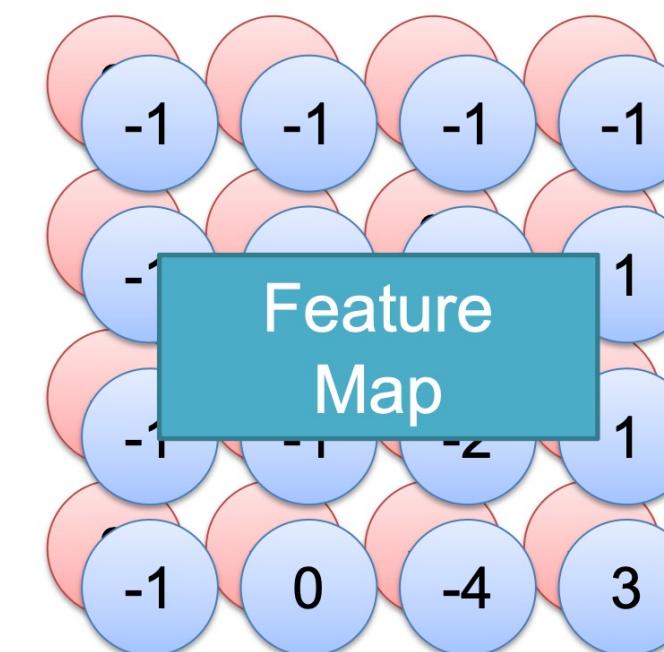
CNN: Convolution

Filter 2

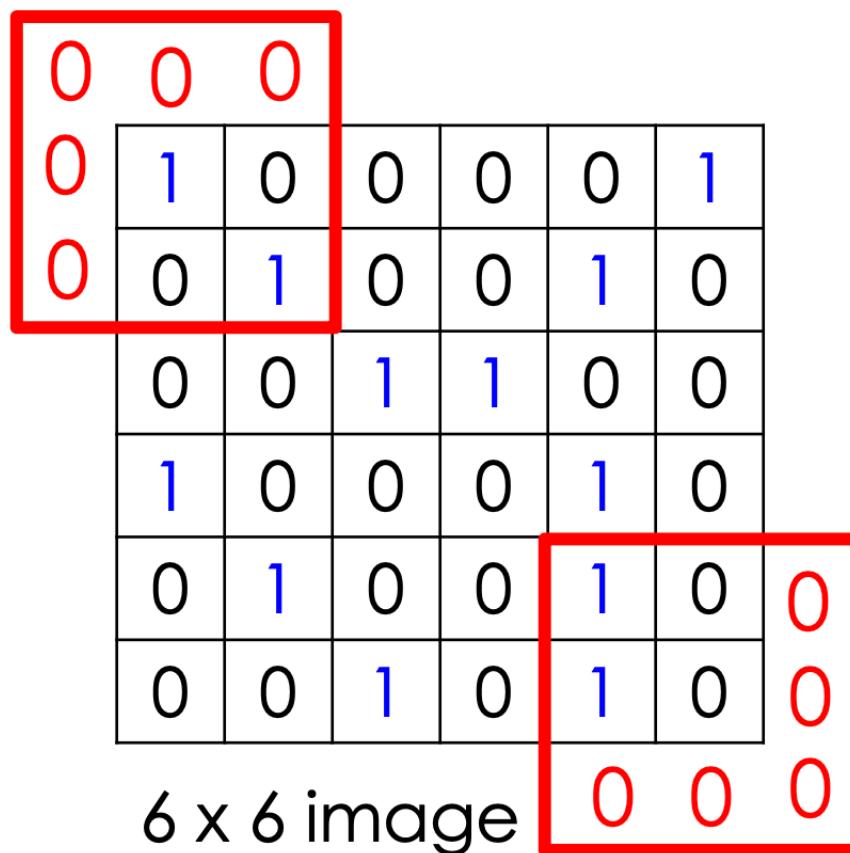
-1	1	-1
-1	1	-1
-1	1	-1



Do the same process
for every filter



CNN: Zero Padding

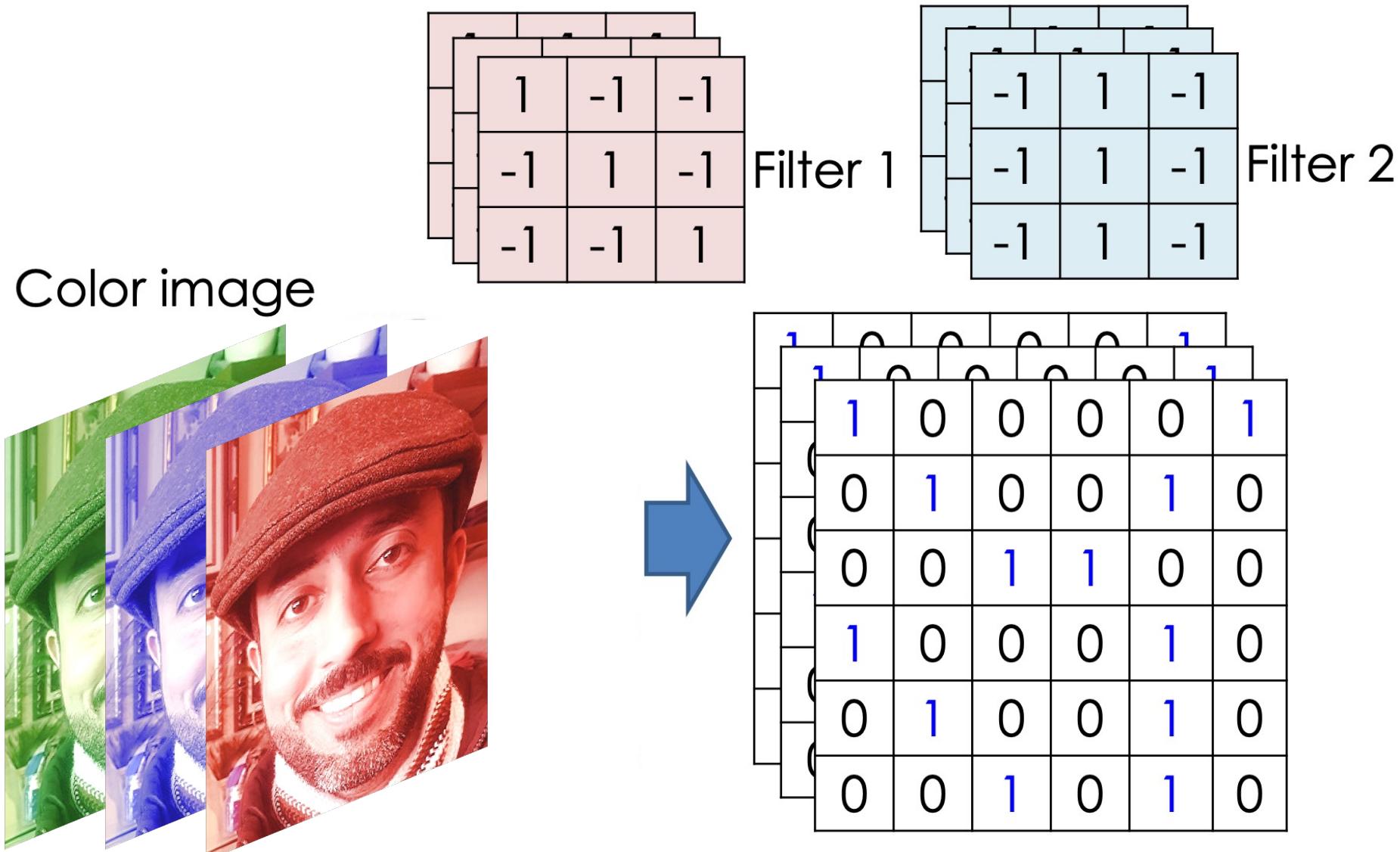


You will get another 6×6 image in this way

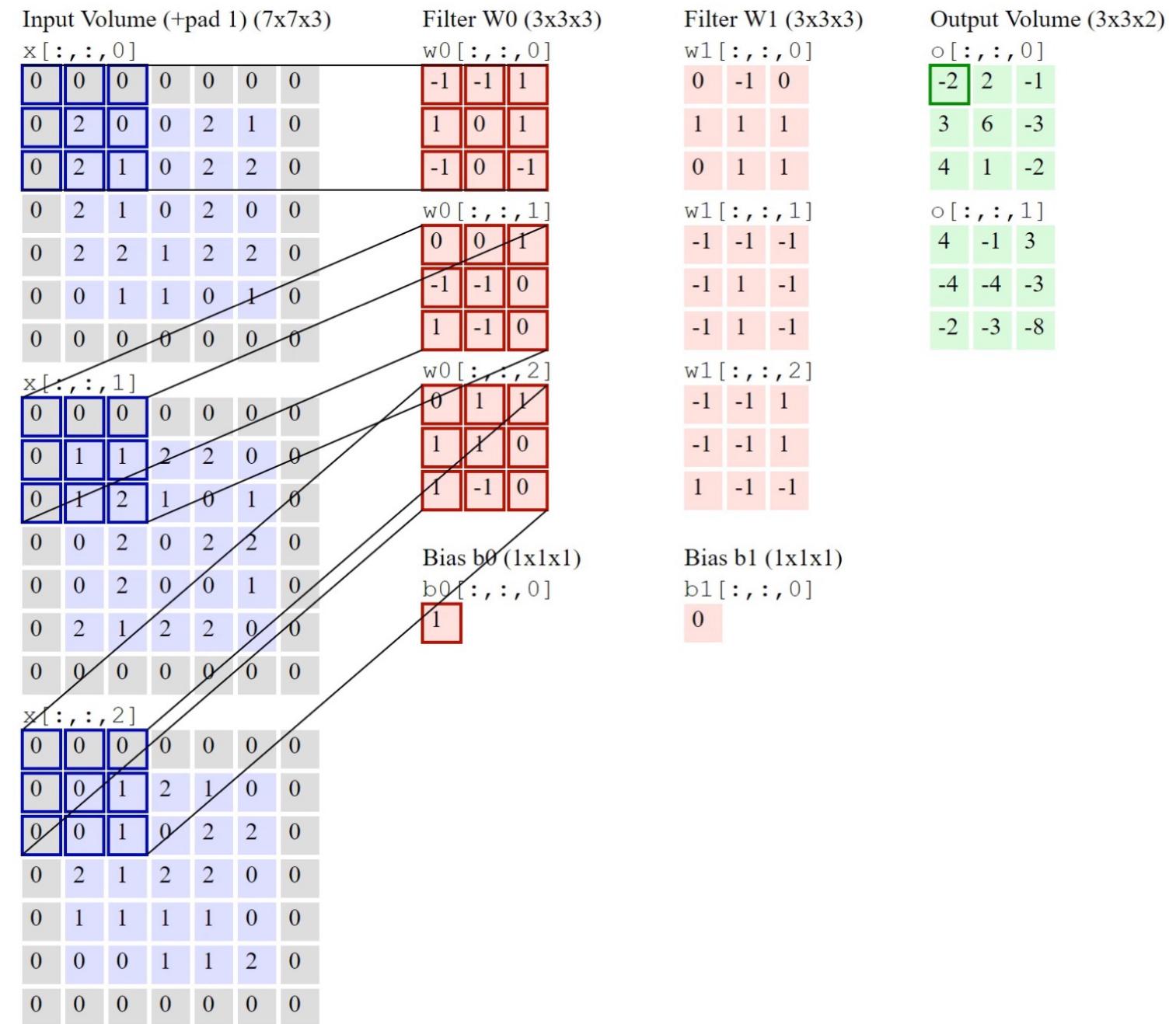


Zero padding

CNN: Color Image



CNN: Color Image



Convolutional Network: Example

- Convolution network is a sequence of these layers

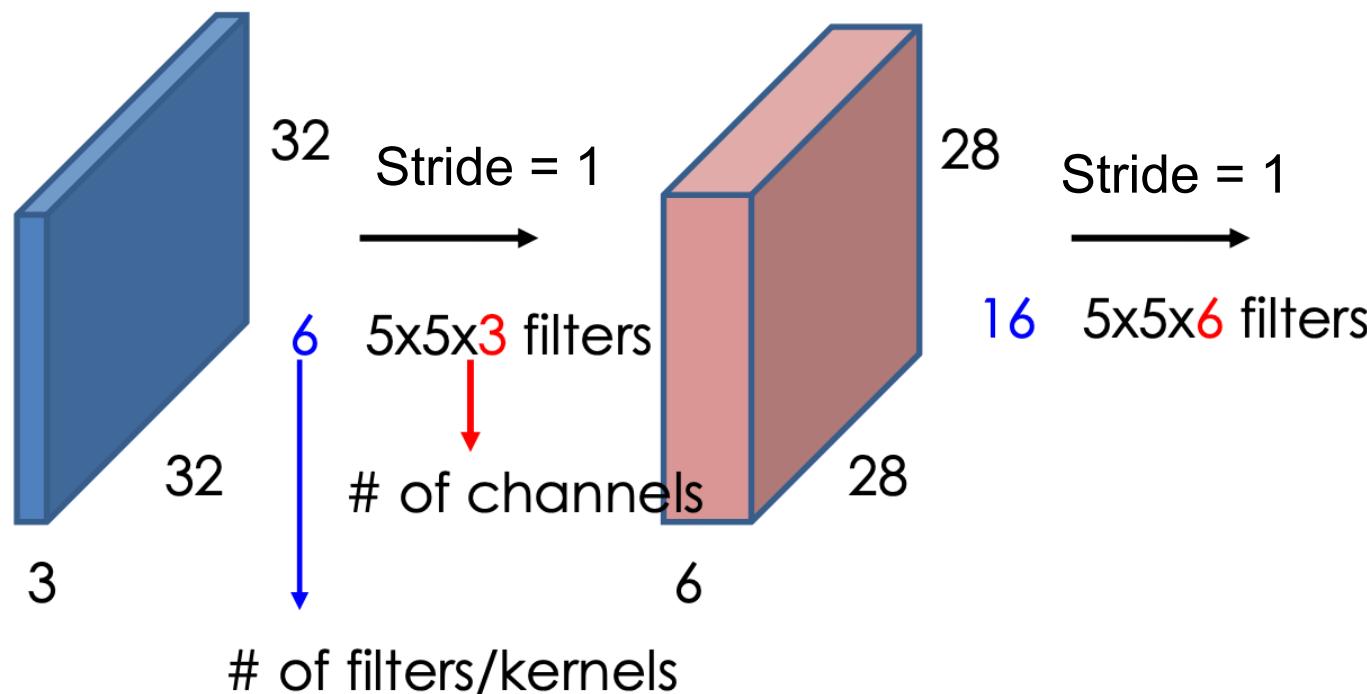
$$\text{Output size} = \frac{\text{Input size} - \text{Filter size}}{\text{Stride}} + 1$$



Convolutional Network: Example

- Convolution network is a sequence of these layers

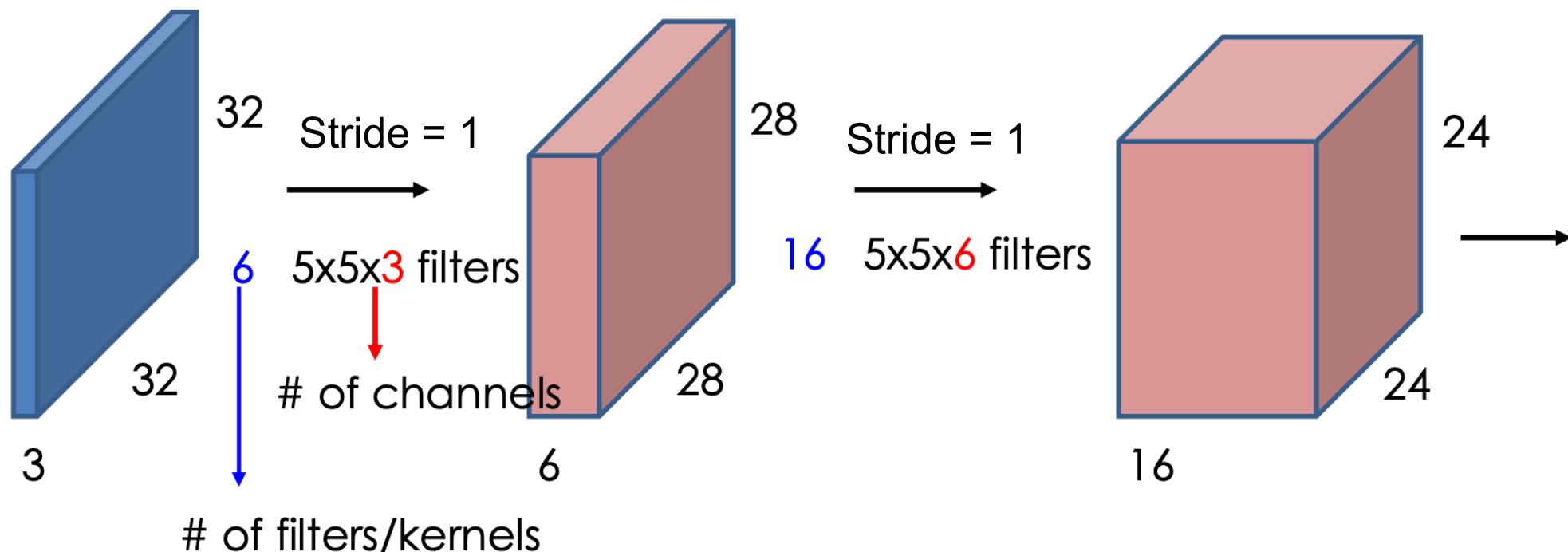
$$\text{Output size} = \frac{\text{Input size} - \text{Filter size}}{\text{Stride}} + 1$$



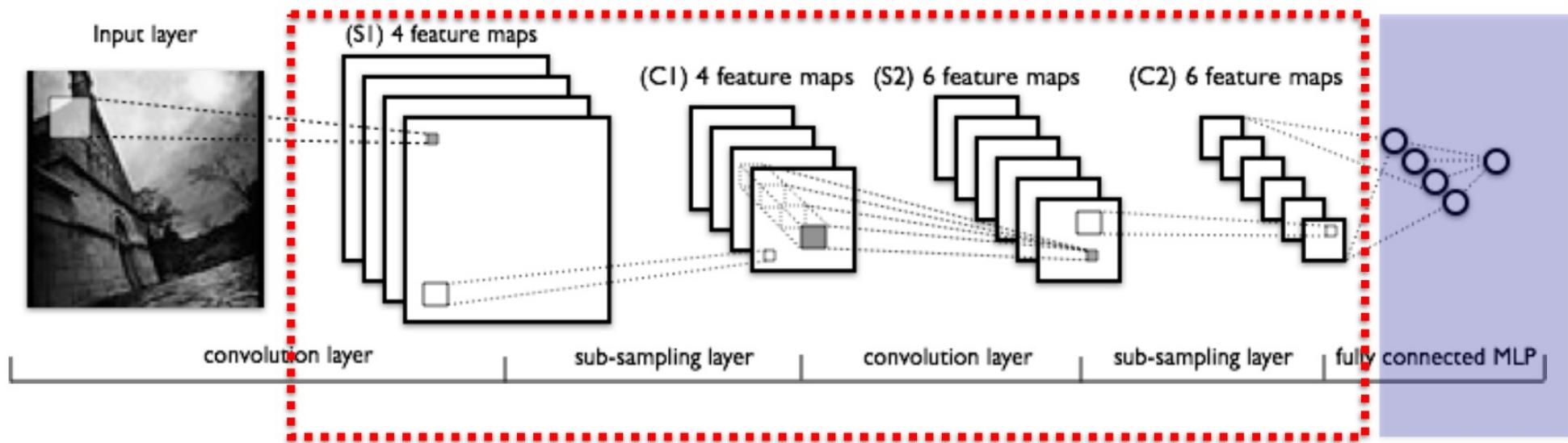
Convolutional Network: Example

- Convolution network is a sequence of these layers

$$\text{Output size} = \frac{\text{Input size} - \text{Filter size}}{\text{Stride}} + 1$$

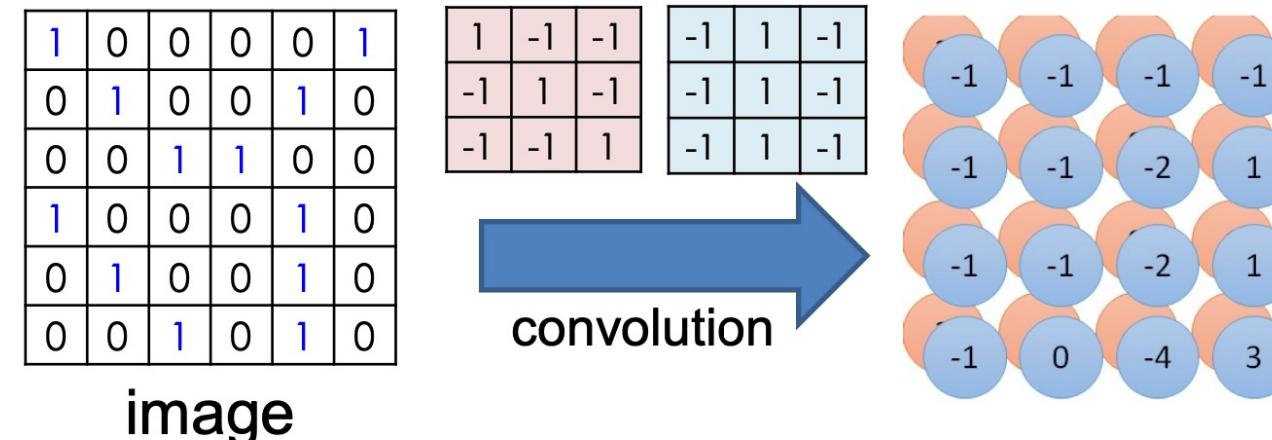


Convolutional Network

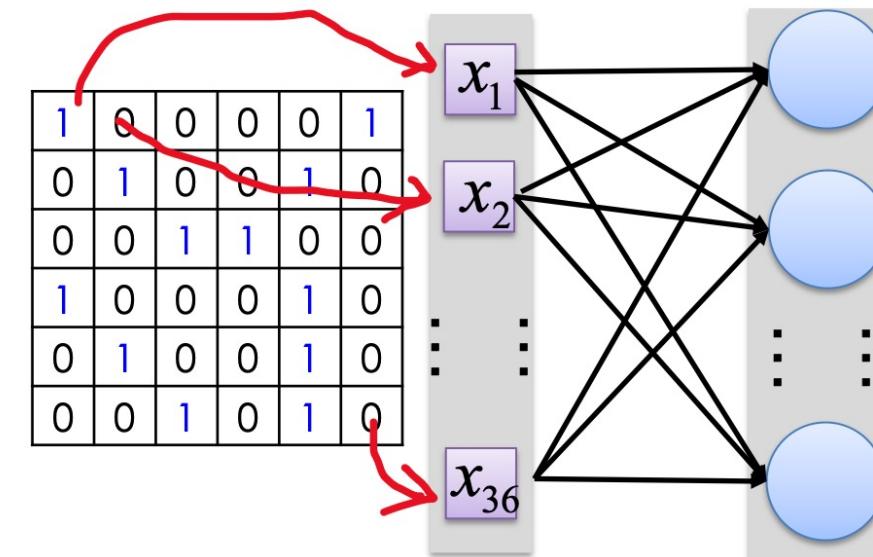


**Feature extraction using
convolutional layers**

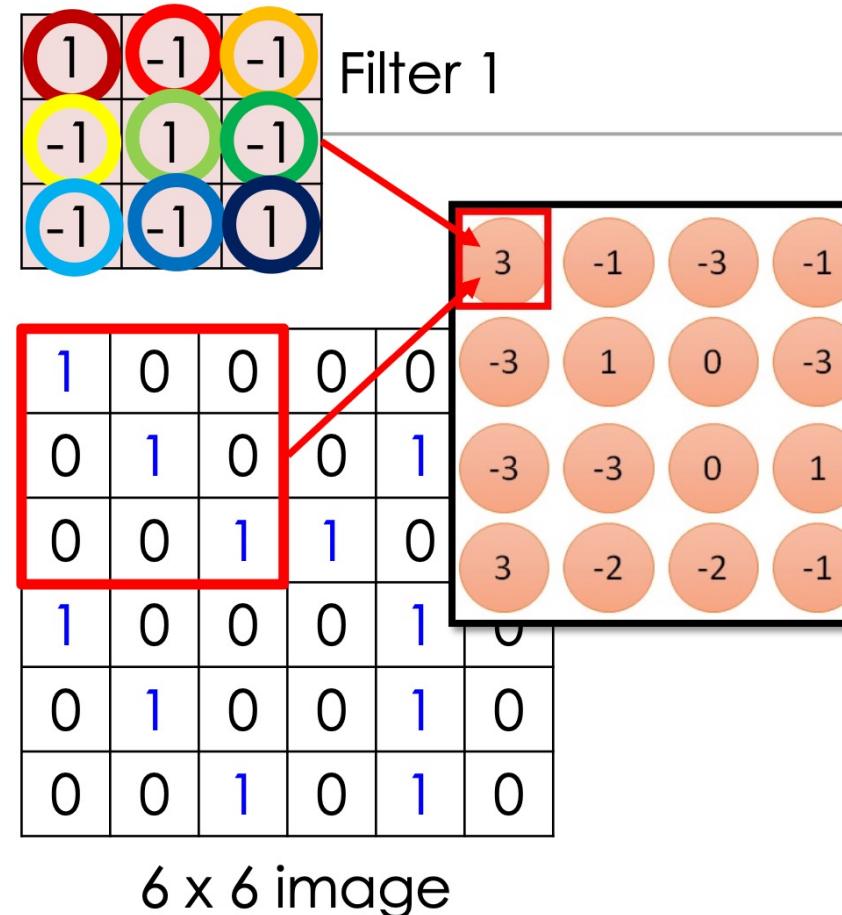
Convolution vs. Fully Connected



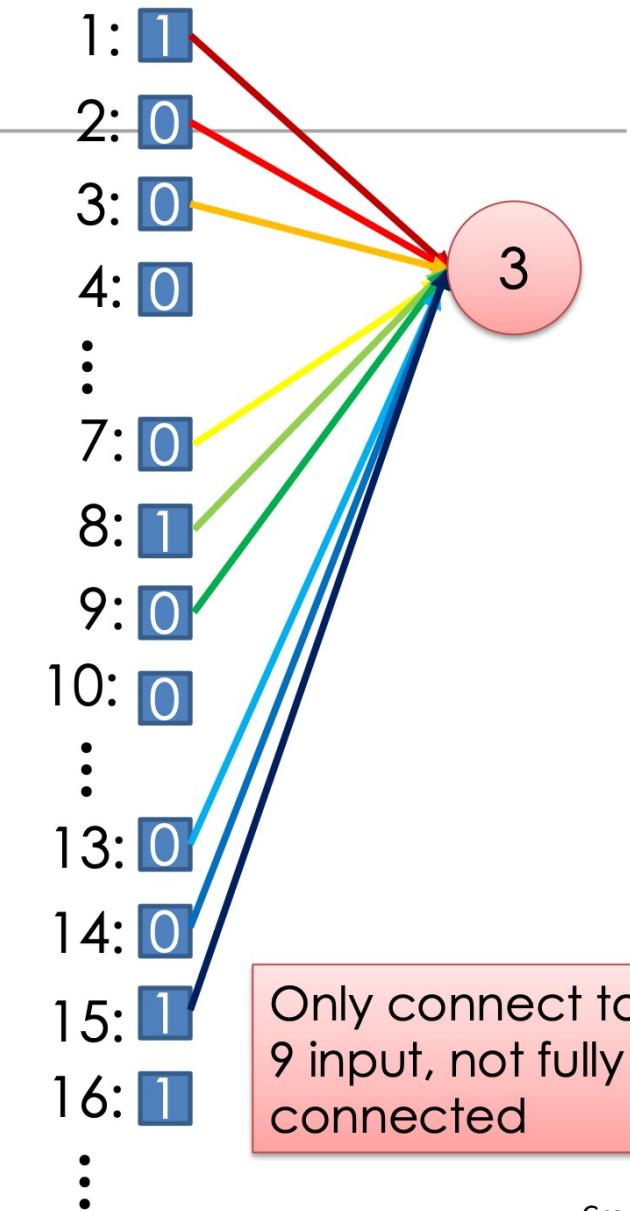
Fully-connected



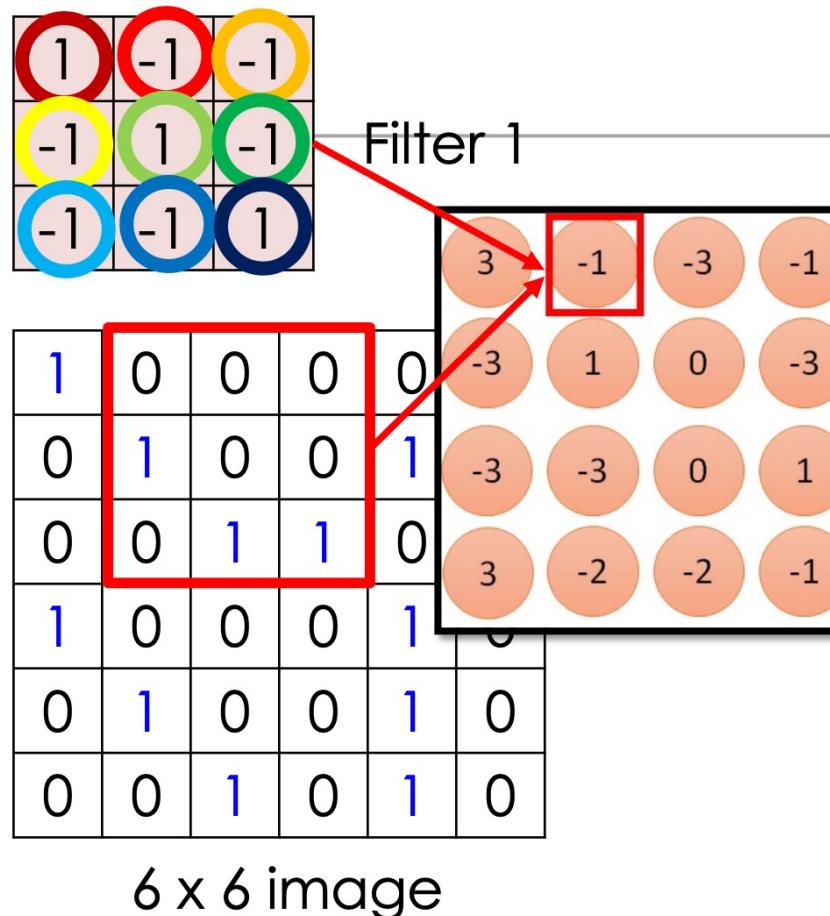
Convolution vs. Fully Connected



Less parameters!

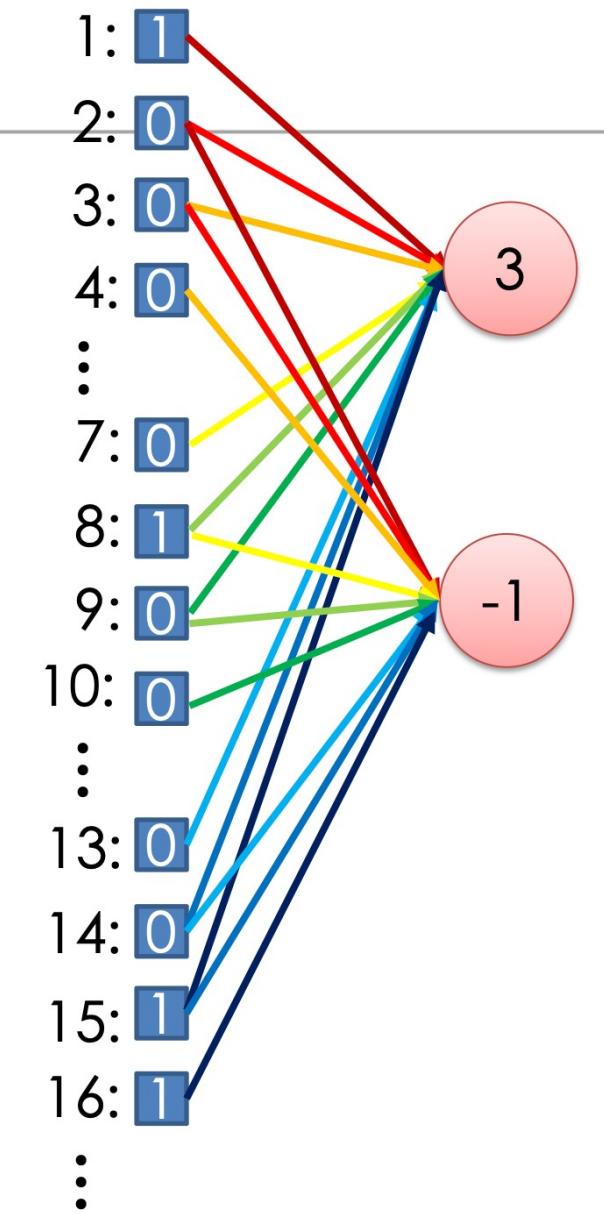


Convolution vs. Fully Connected



Less parameters!

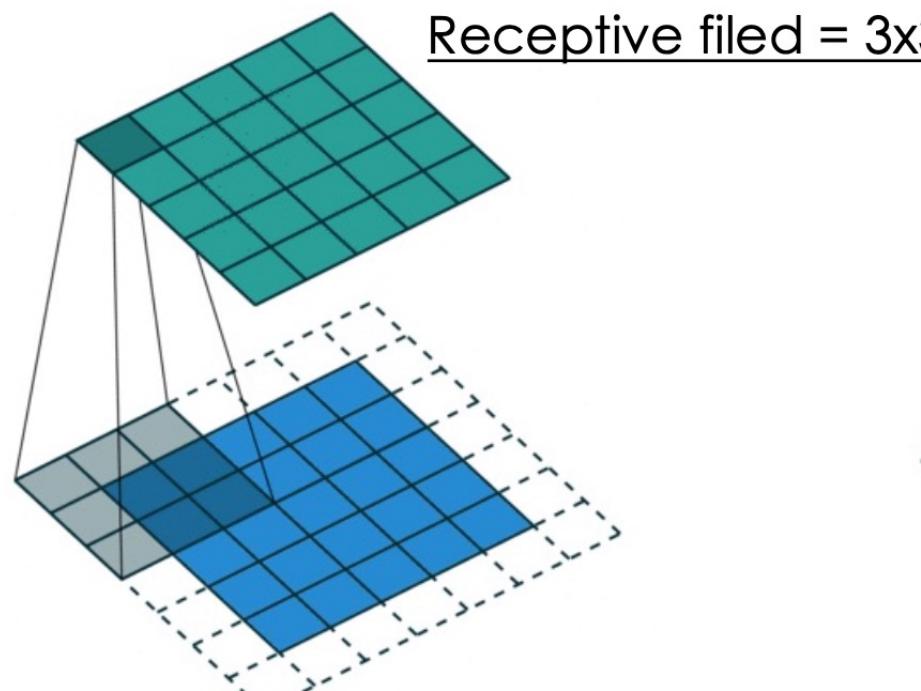
Even less parameters!



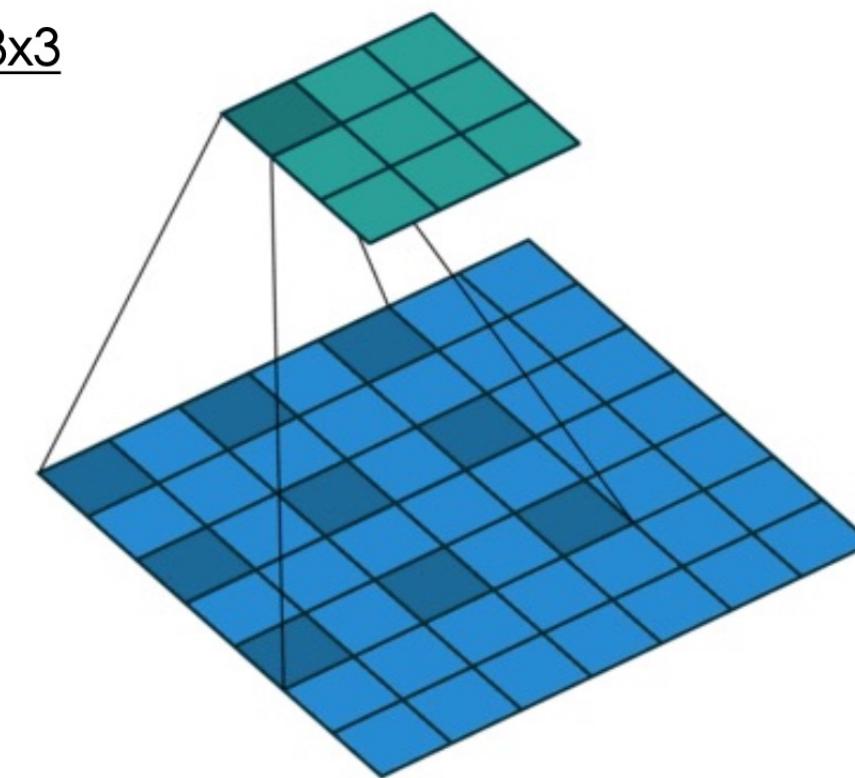
Variants of Convolution Operation

- Dilated/Atrous Convolution

Dilation rate = 2
Receptive field = 5x5
With only 9 parameters

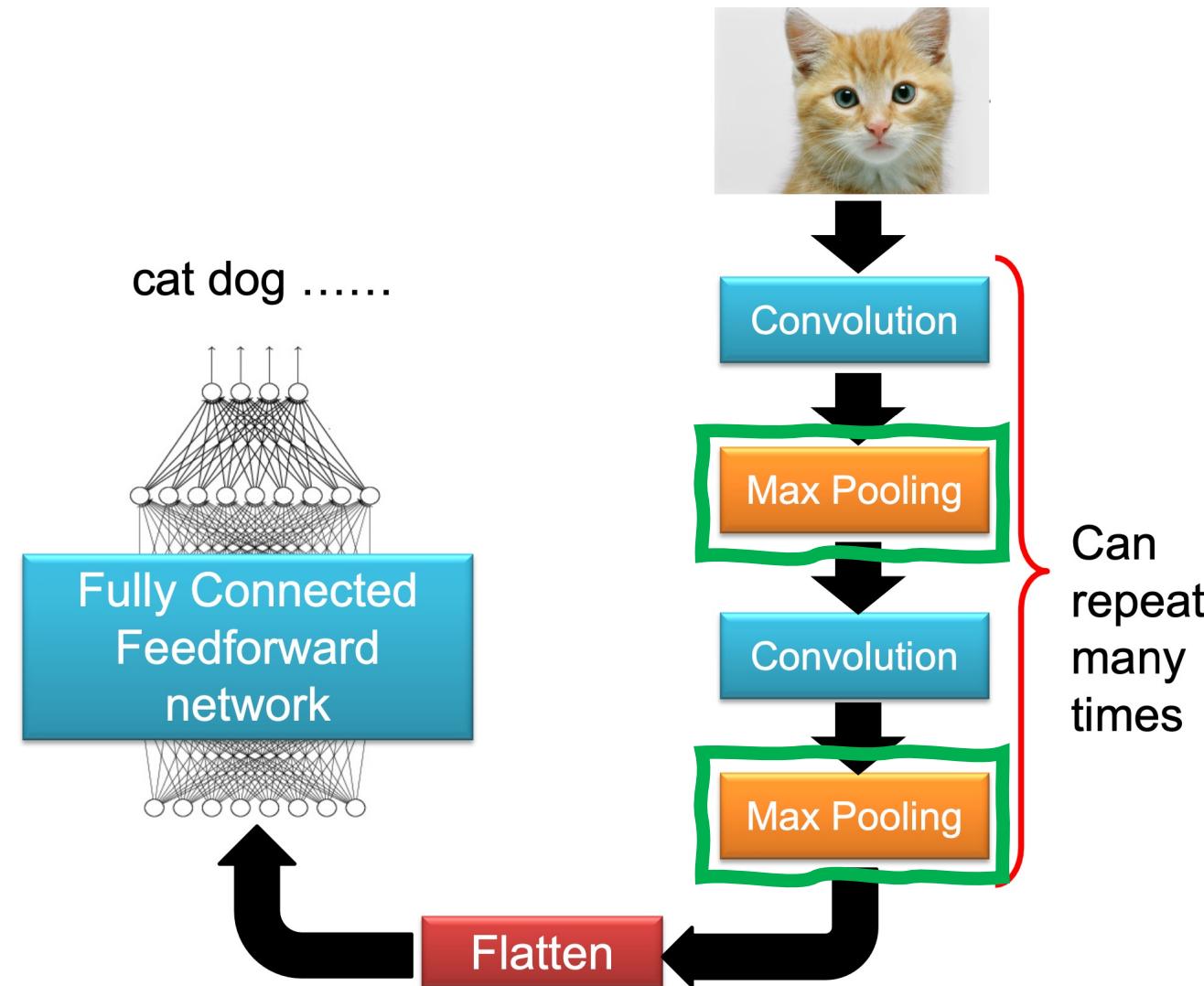


Normal convolution



Effective in segmentation task

The Whole CNN



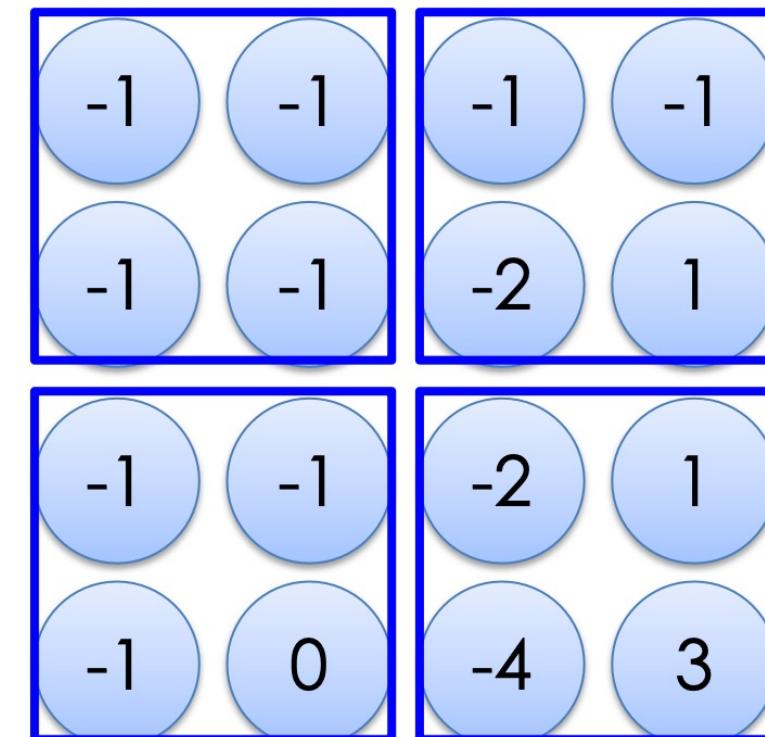
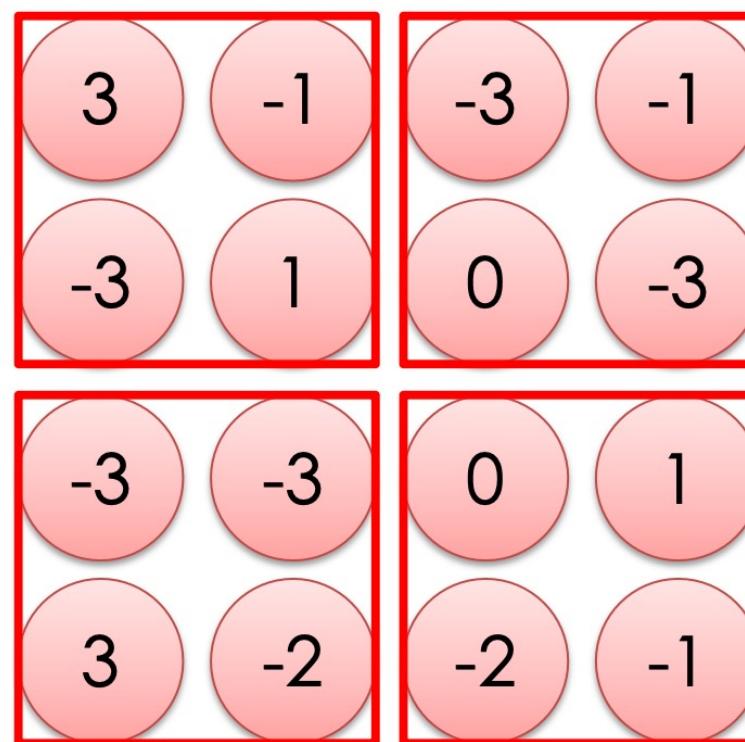
CNN: Max Pooling

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

Filter 1

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

Filter 2



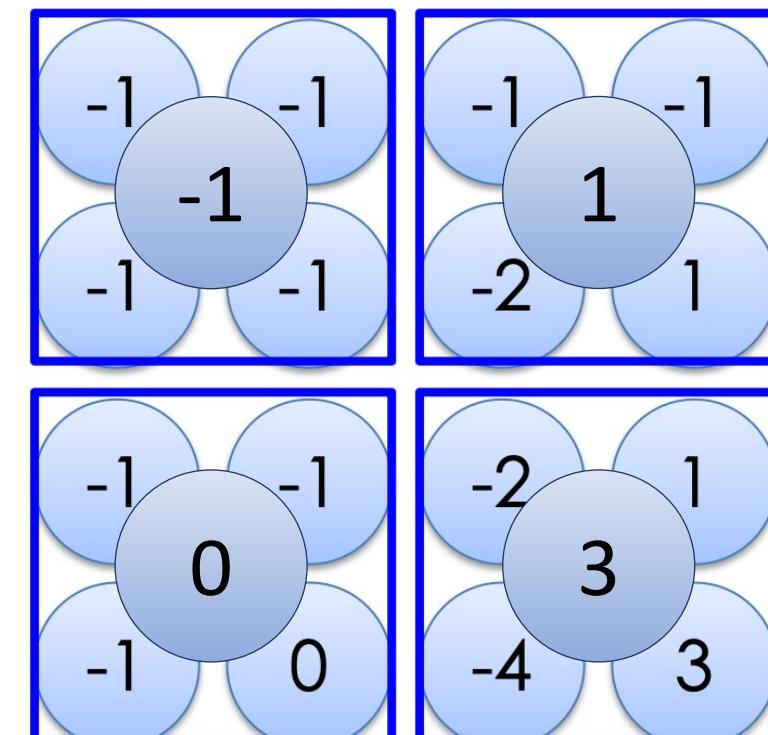
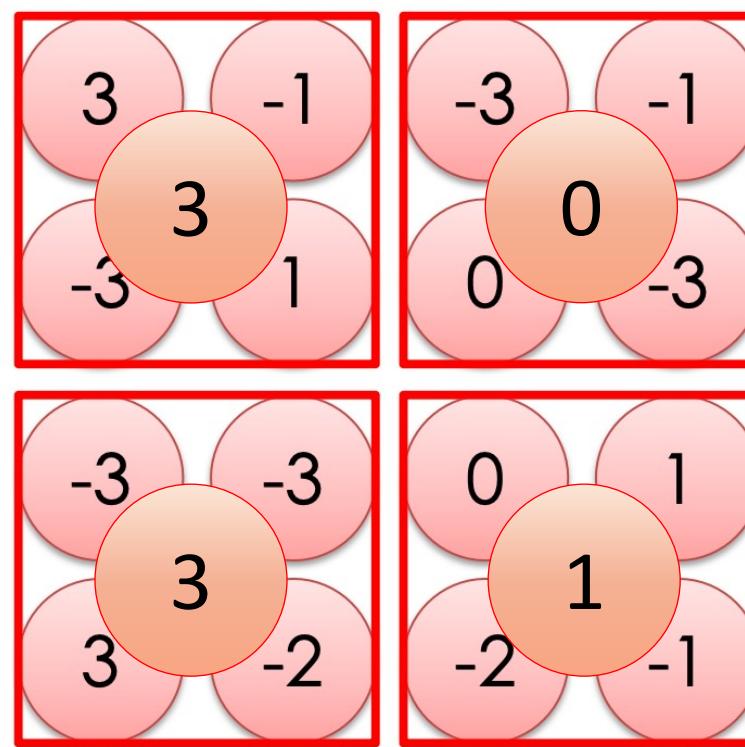
CNN: Max Pooling

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

Filter 1

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

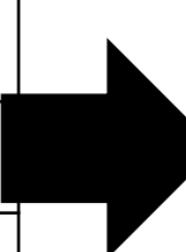
Filter 2



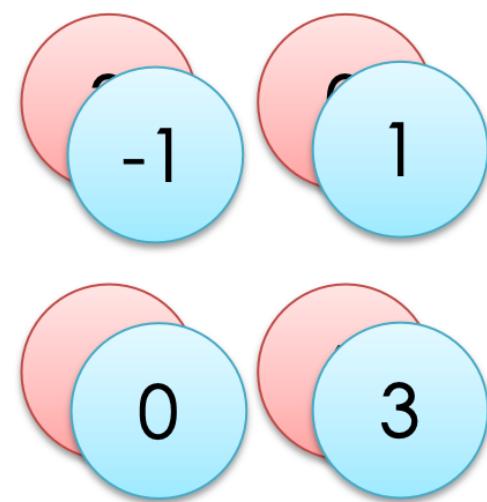
CNN: Max Pooling

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

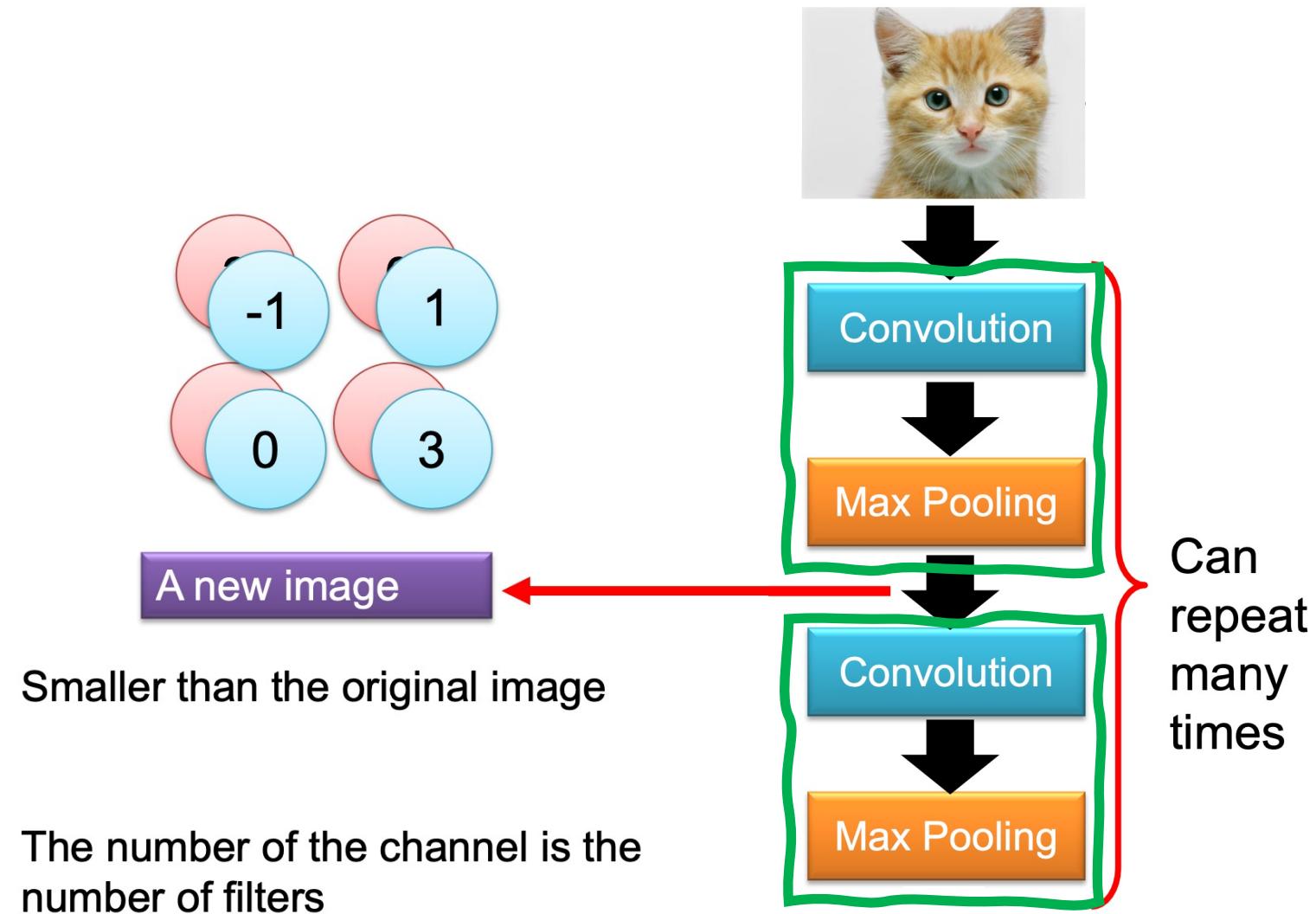


New image
but smaller



2 x 2 image

The Whole CNN



Why Pooling?

- It progressively reduce the spatial size of the representation
 - ➔ reduce the number of parameters and computation in the network.

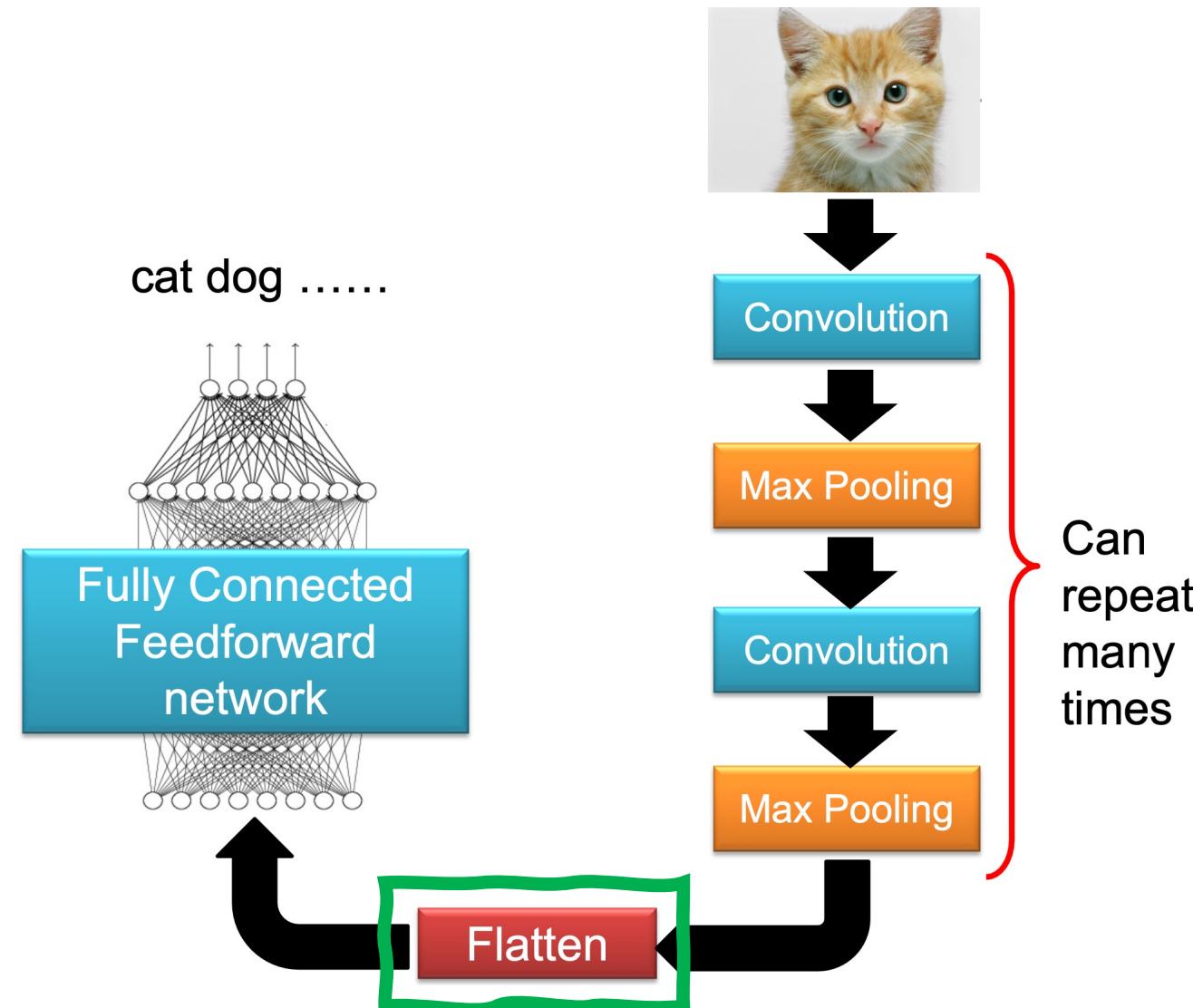
Why Pooling?

- It progressively reduce the spatial size of the representation
 - ➔ reduce the number of parameters and computation in the network.
- Another way to reduce feature map size?
 - ➔ convolution with stride, e.g., 2.

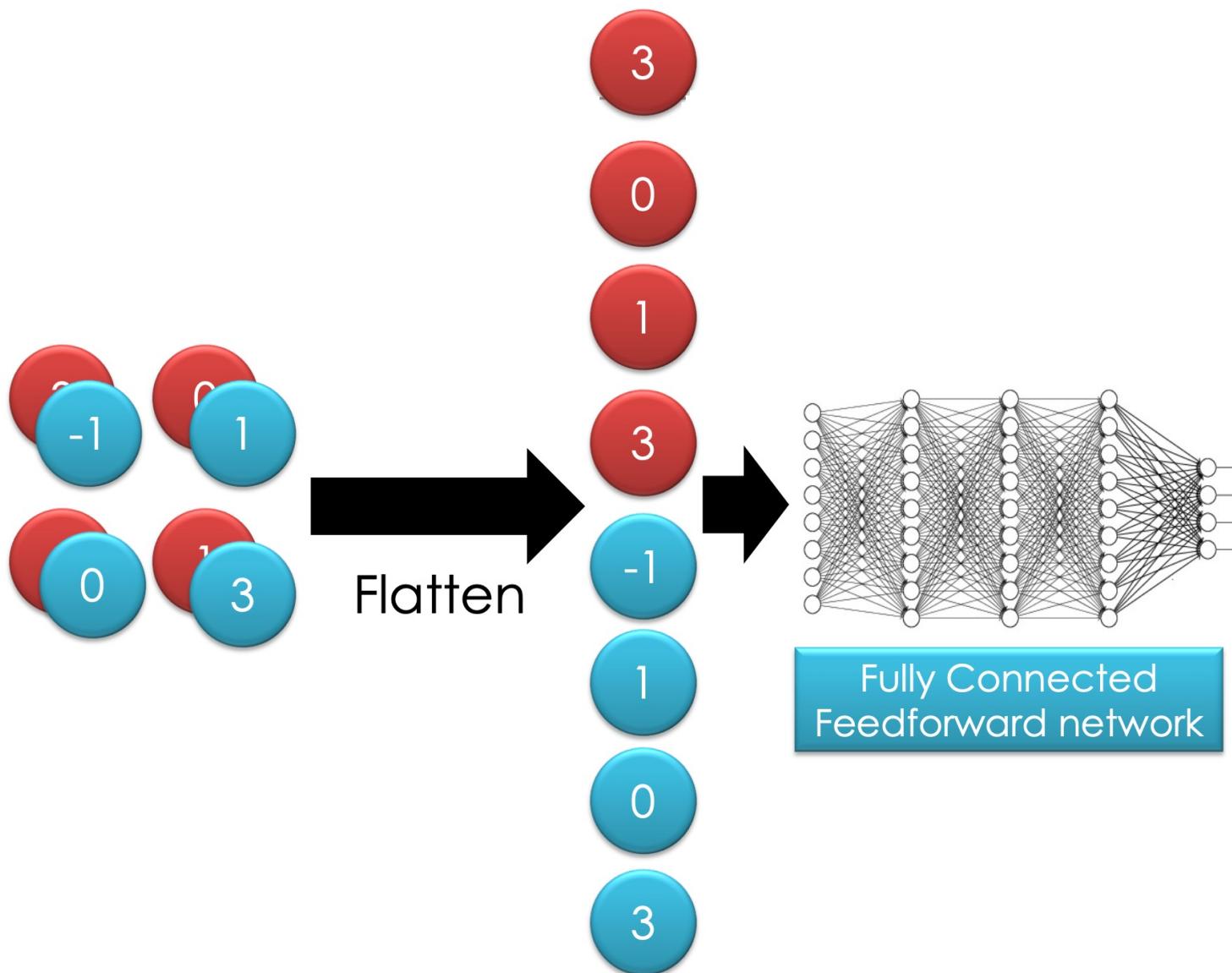
Why Pooling?

- It progressively reduce the spatial size of the representation
 - ➔ reduce the number of parameters and computation in the network.
- Another way to reduce feature map size?
 - ➔ convolution with stride, e.g., 2.
- If instead of taking the maximum, using the **average** will be the **average pooling**.

The Whole CNN

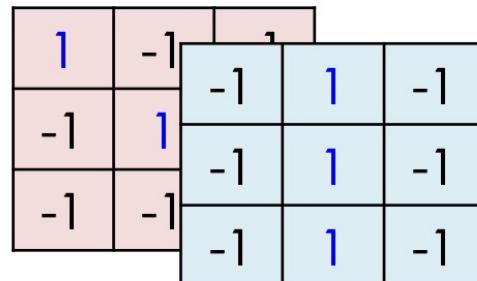


Flatten



CNN in Keras

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

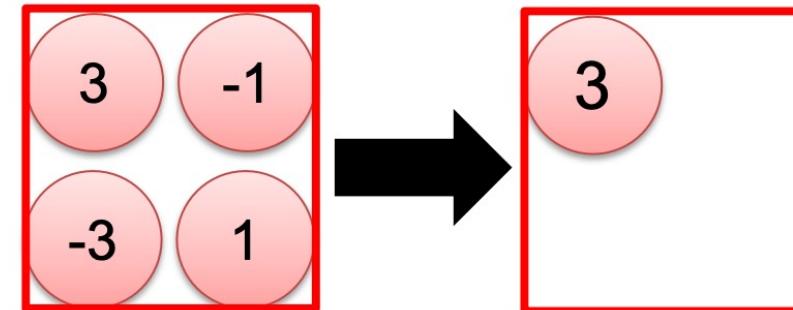


There are
25 3x3
filters.

Input_shape = (1 , 28 , 28)

1: grayscale, 3: RGB 28 x 28 pixels

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



input

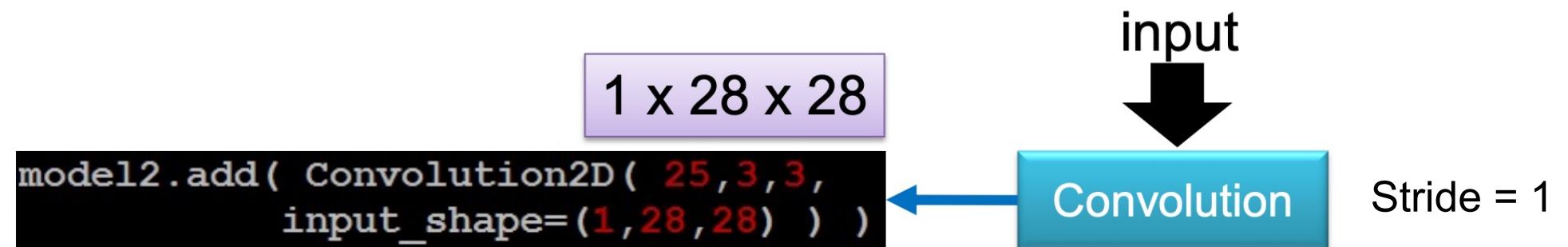
Convolution

Max Pooling

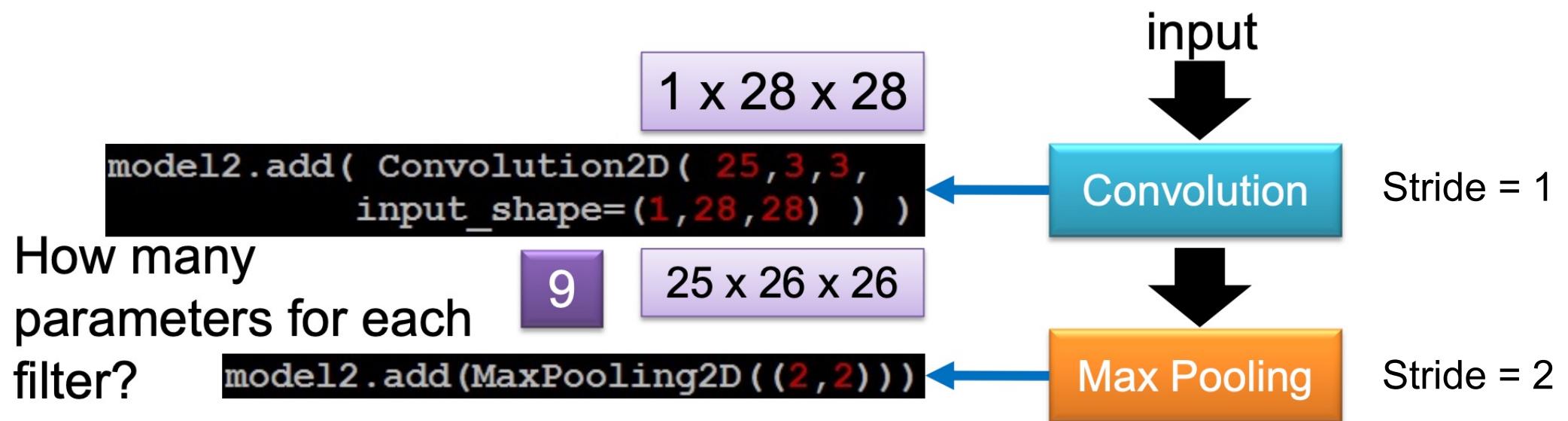
Convolution

Max Pooling

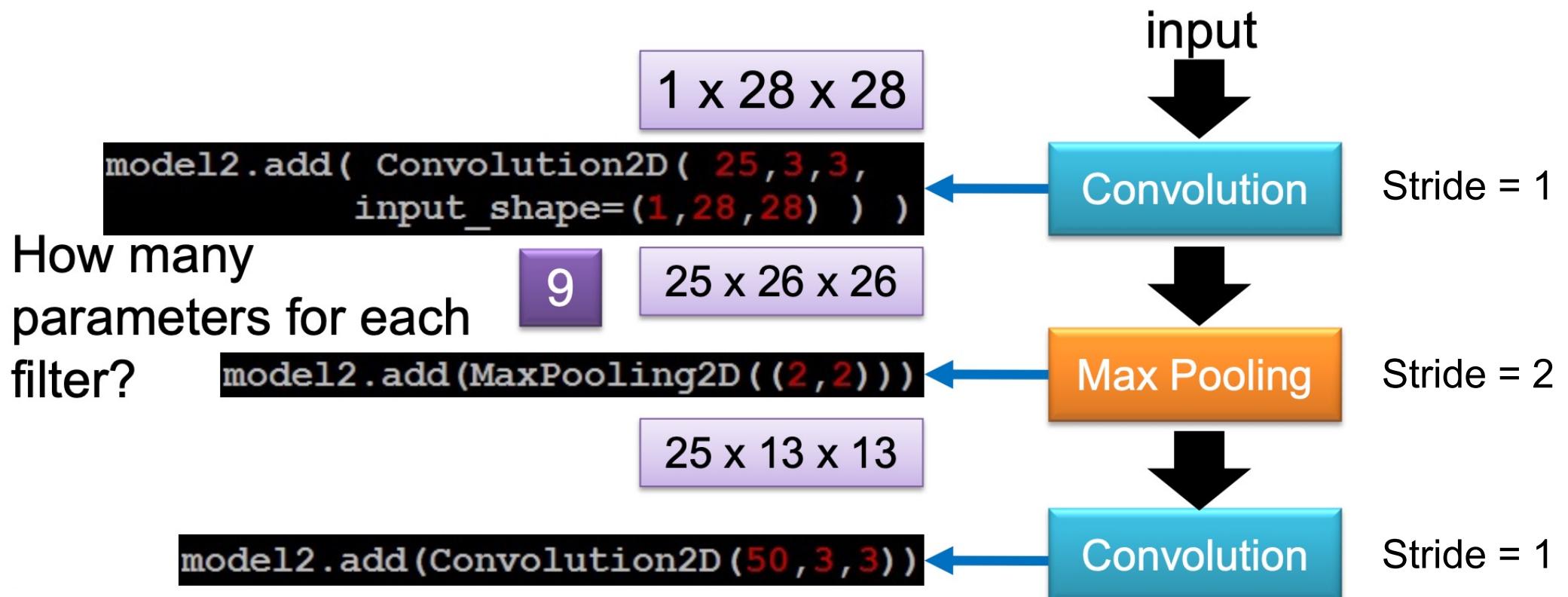
CNN in Keras



CNN in Keras



CNN in Keras



CNN in Keras

How many parameters for each filter?

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

1 x 28 x 28

9

25 x 26 x 26

How many parameters for each filter?

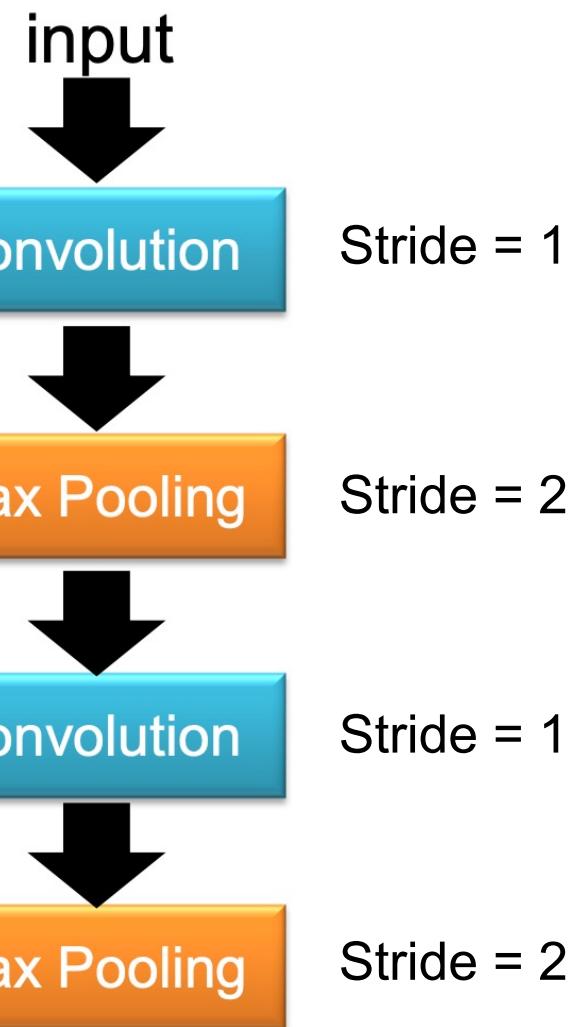
```
model2.add(Convolution2D(50, 3, 3))
```

225

50 x 11 x 11

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

50 x 5 x 5



CNN in Keras

output

Fully Connected
Feedforward

```
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

$50 \times 5 \times 5$

Loss Function

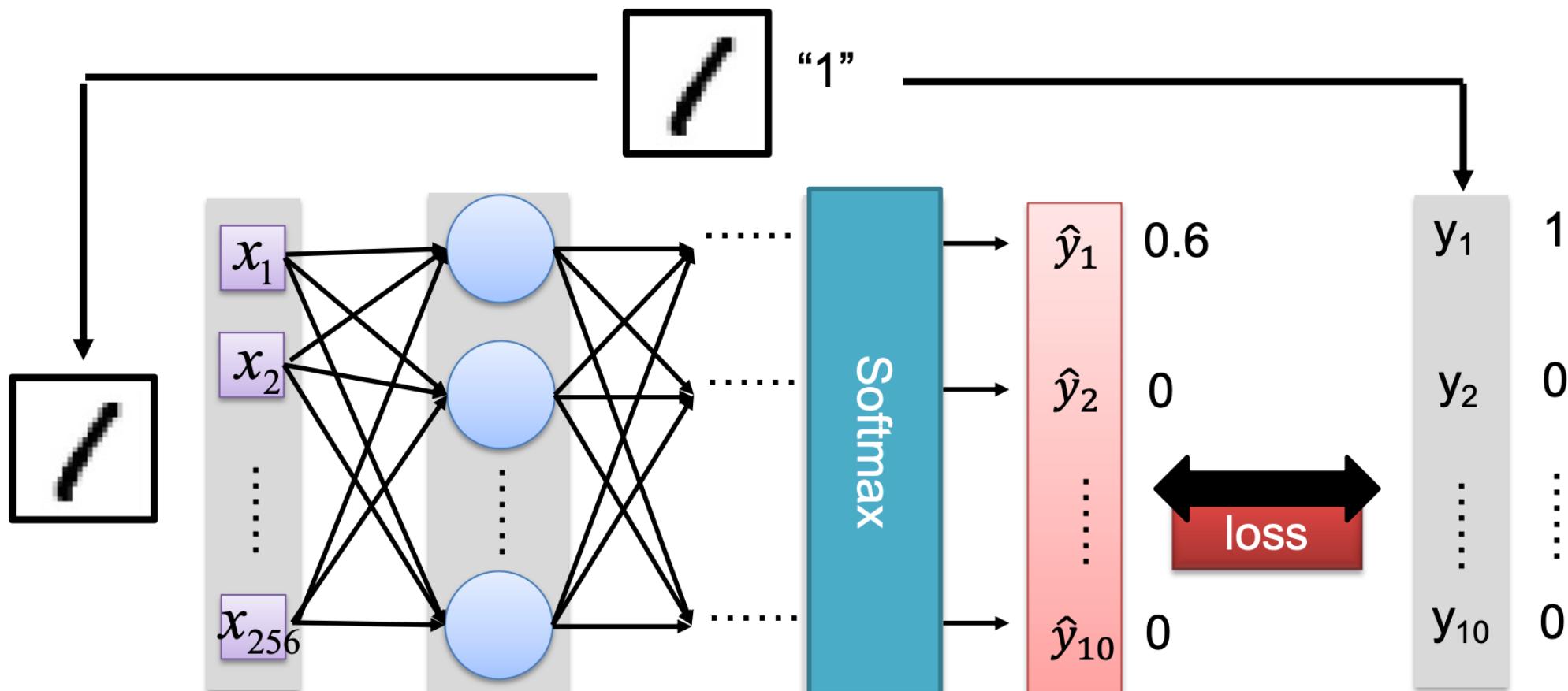
- Way to define how good the network is performing
 - In terms of prediction
- Network training (Optimization)
 - Find the best network parameters to minimize the loss

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(\mathbf{x}_i, W), y_i)$$

Total loss for a training set
N samples

input Ground truth
Network parameters
Loss function
network

Proper Loss Selection



Square
Error

$$\sum_{i=1}^{10} (y_i - \hat{y}_i)^2$$

Target
(one hot encoding)

Cross-entropy Loss

- Binary case (we know this from logistic regression)

$$-(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$$

- General form

K classes

$$-\sum_{i=1}^K y_i \log(\hat{y}_i)$$

Label \mathbf{y} is a one-hot vector

$$y_i \in [0,1]$$

Cross-entropy Loss

- Multi-class classification

Cross-Entropy loss for one sample

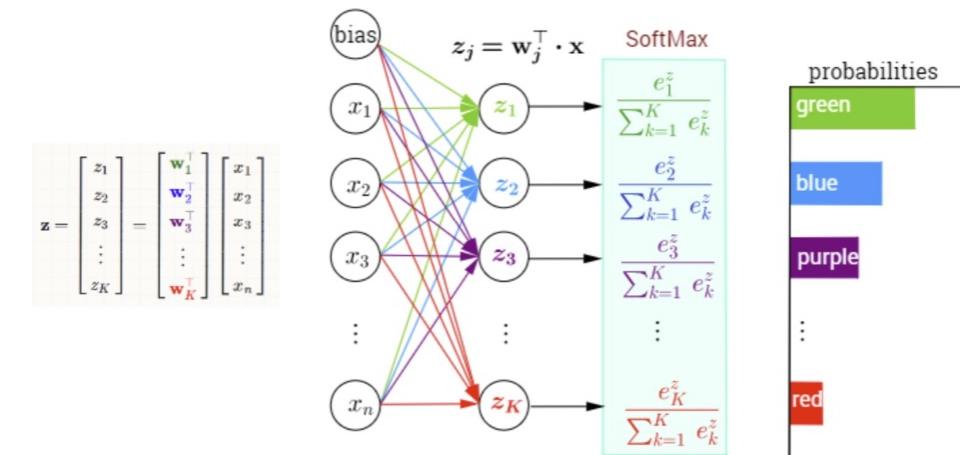
$$-\sum_{i=1}^K y_i \log(\hat{y}_i)$$

Label \mathbf{y} is a one-hot vector

$$(y_i) \in [0,1]$$

computed	targets	correct?

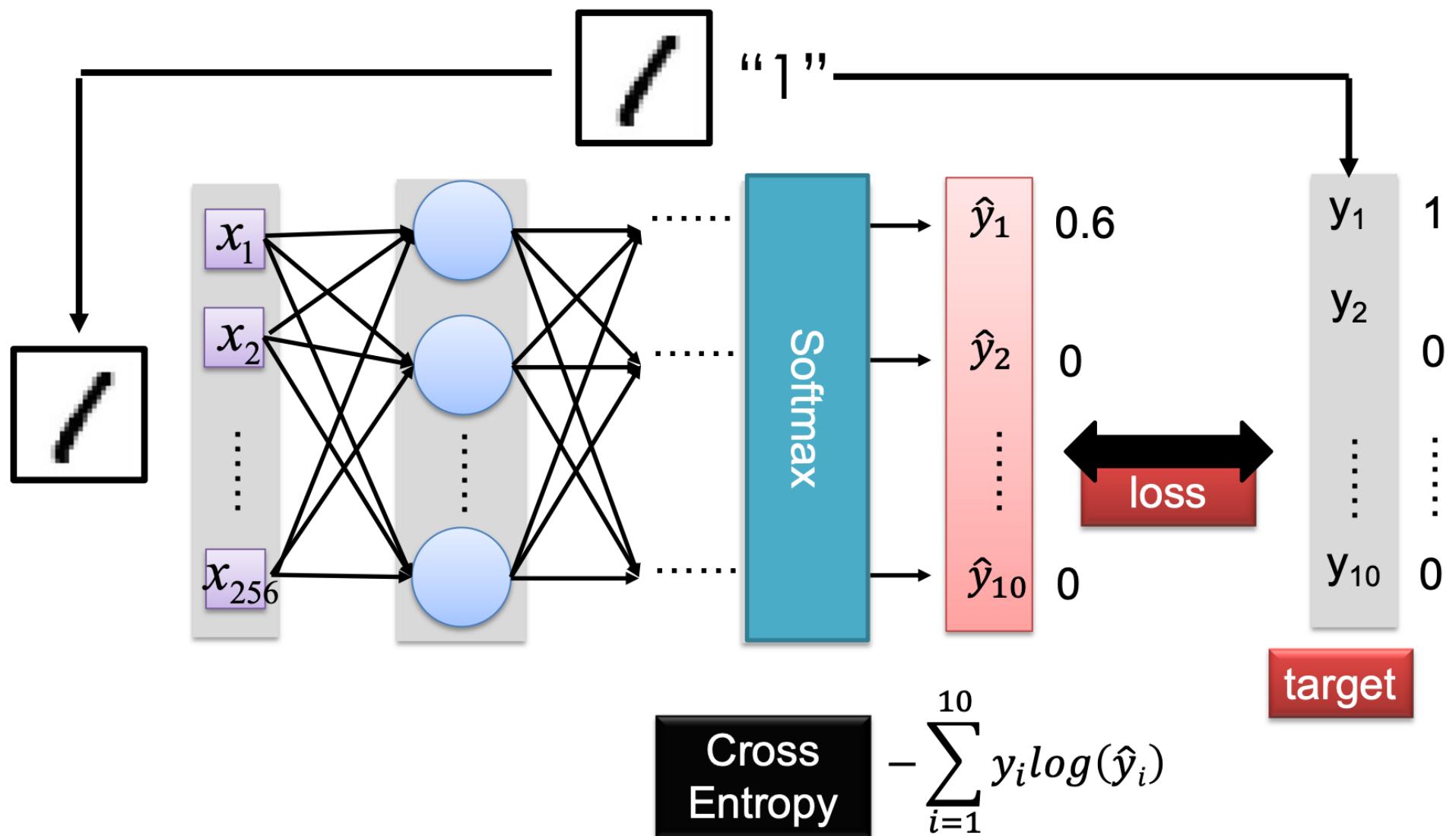
0.3 0.3 0.4	0 0 1 (democrat)	yes
0.3 0.4 0.3	0 1 0 (republican)	yes
0.1 0.2 0.7	1 0 0 (other)	no



$$\sigma(j) = \frac{\exp(\mathbf{w}_j^\top \mathbf{x})}{\sum_{k=1}^K \exp(\mathbf{w}_k^\top \mathbf{x})} = \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)}$$

$$- ((\log(0.3) \times 0) + (\log(0.3) \times 0) + (\log(0.4) \times 1)) = -\log(0.4)$$

Cross-entropy Loss



Proper Loss Selection

Square Error

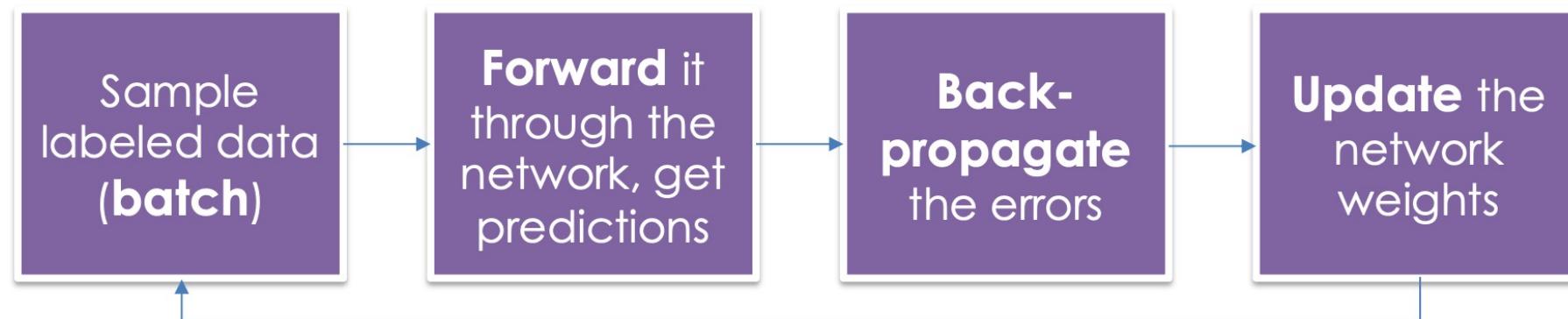
```
model.compile(loss='mse',  
              optimizer=SGD(lr=0.1),  
              metrics=['accuracy'])
```

Cross Entropy

```
model.compile(loss='categorical_crossentropy',  
              optimizer=SGD(lr=0.1),  
              metrics=['accuracy'])
```

Several alternatives: <https://keras.io/objectives/>

Training



Optimize (min. or max.) **objective/cost function $J(\theta)$**

Generate **error signal** that measures difference between predictions and target values

Use error signal to change the **weights** and get more accurate predictions

Subtracting a fraction of the **gradient** moves you towards the **(local) minimum of the cost function**

Gradient Descent

Optimize (min. or max.) **objective/cost function $J(\theta)$**

Use error signal to change the **weights** and get more accurate predictions

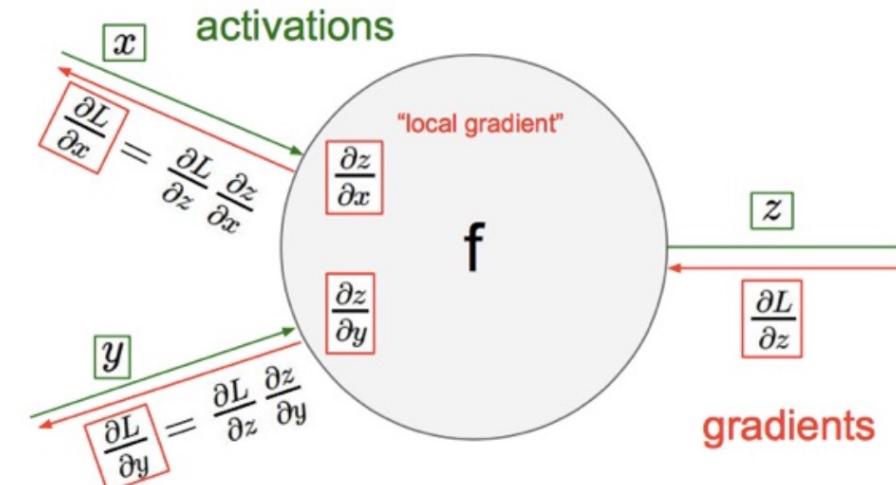
Subtracting a fraction of the **gradient** moves you towards the **(local) minimum of the cost function**

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} - \alpha \frac{d}{d\theta_j^{\text{old}}} J(\theta) \quad \text{Update each element of } \theta$$

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \nabla_{\theta} J(\theta)$$

learning rate

Matrix notation for all parameters

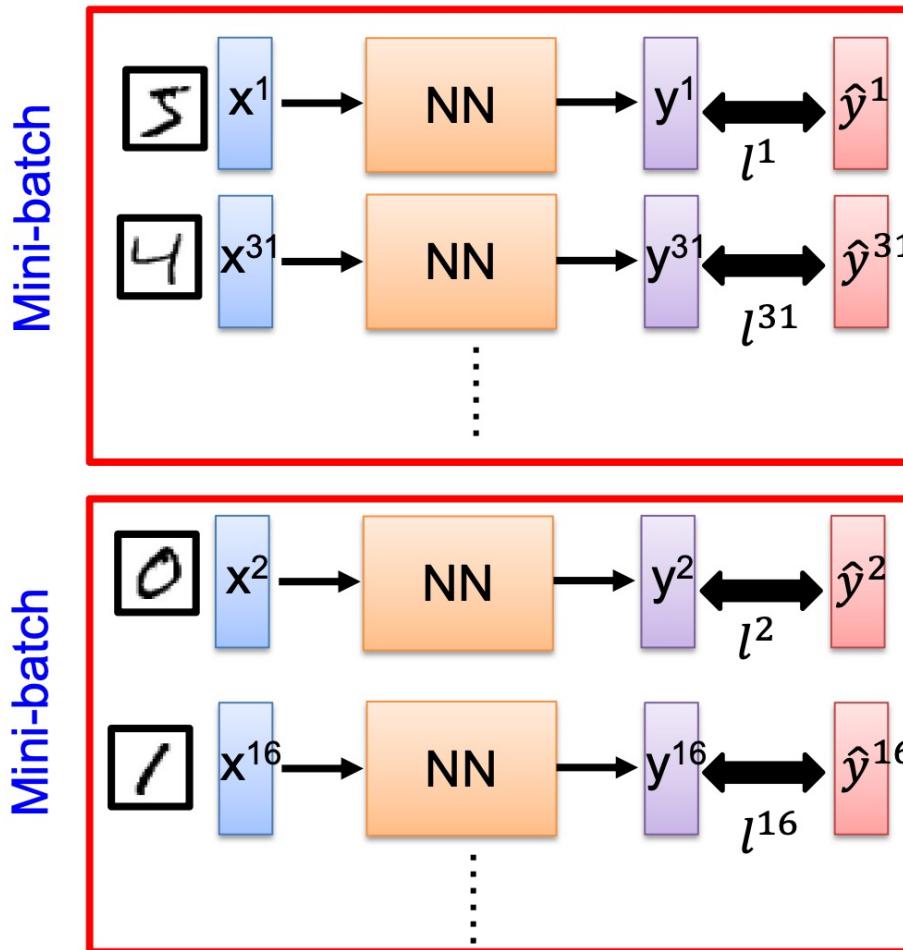


Recursively apply **chain rule** though each node

Stochastic Gradient Descent

- Gradient over entire dataset is impractical
- Better to take quick, noisy steps
- Estimate gradient over a **mini-batch** of examples

Mini Batch



➤ Randomly initialize network parameters

➤ Pick the 1st batch

$$L' = l^1 + l^{31} + \dots$$

Update parameters once

➤ Pick the 2nd batch

$$L'' = l^2 + l^{16} + \dots$$

Update parameters once

➤ Until all mini-batches have been picked

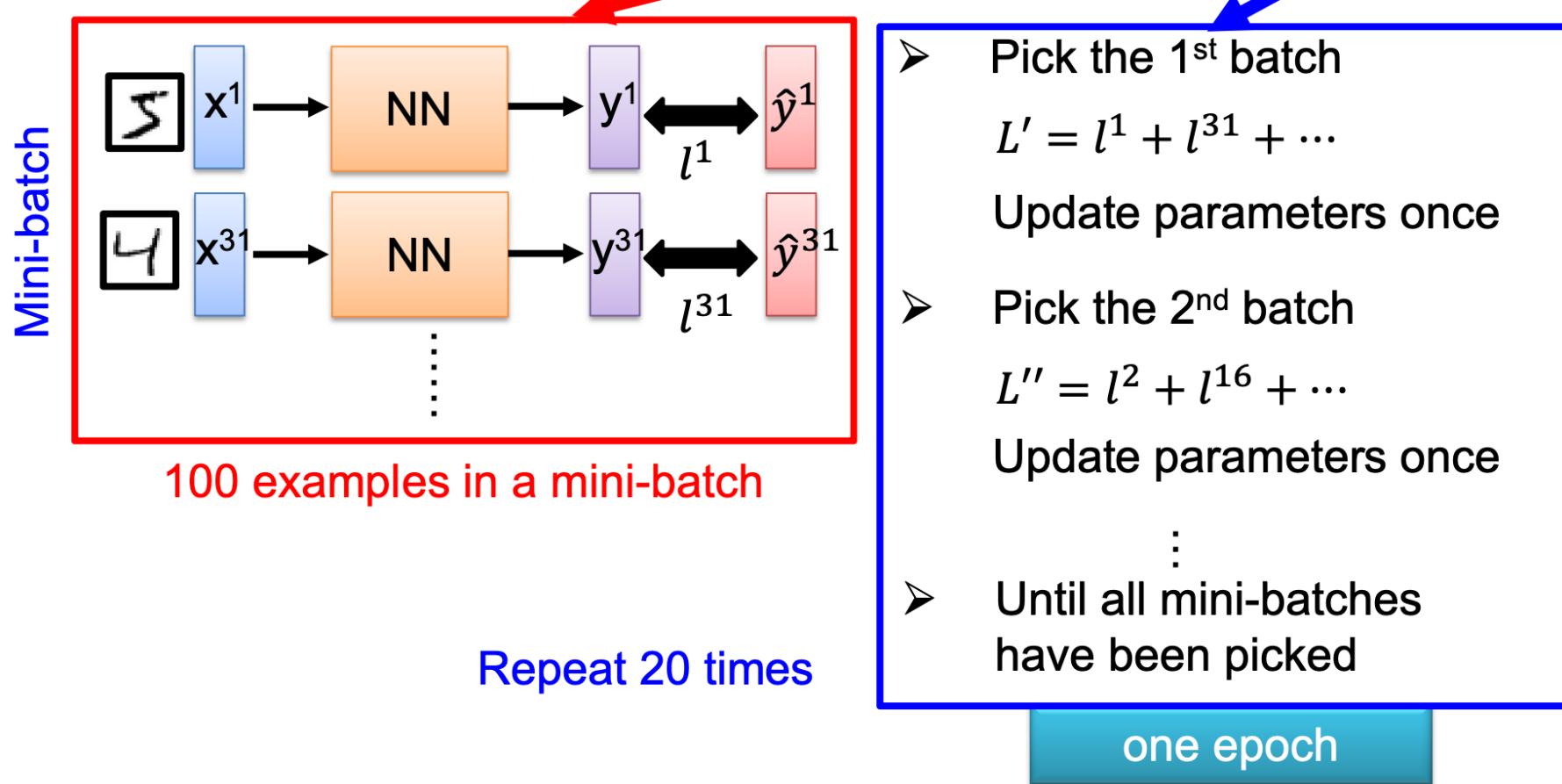
one epoch

Repeat the above process

We do not really minimize total loss!

Mini Batch

```
model.fit(x_train, y_train, batch size=100, nb epoch=20)
```



Other Optimizers

Adagrad [John Duchi, JMLR'11]

RMSprop <https://www.youtube.com/watch?v=O3sxAc4hxZU>

Adadelta [Matthew D. Zeiler, arXiv'12]

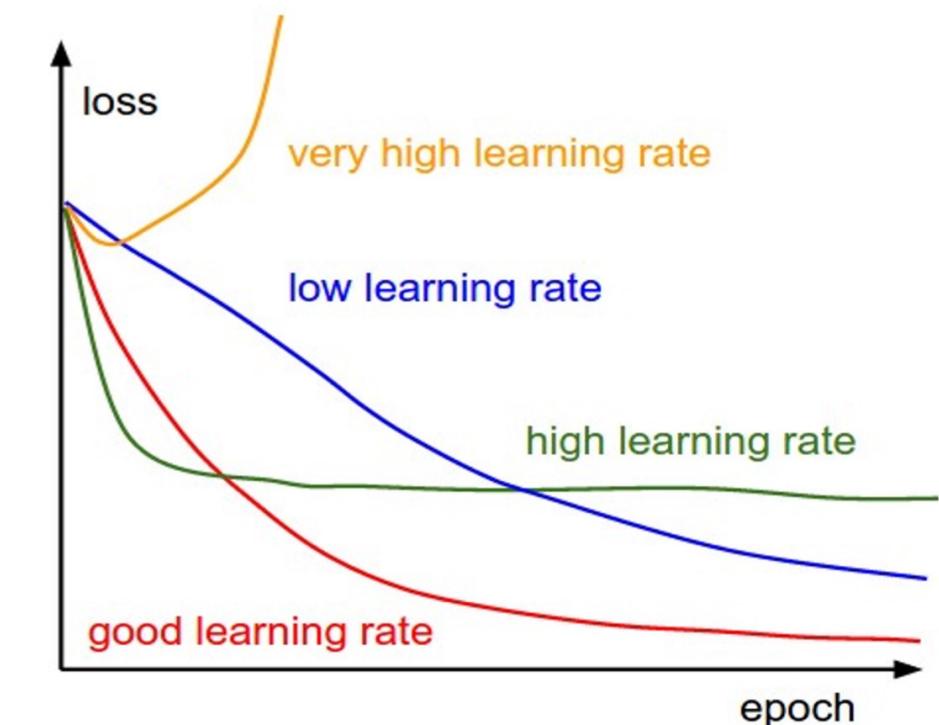
“No more pesky learning rates” [Tom Schaul, arXiv'12]

AdaSecant [Caglar Gulcehre, arXiv'14]

Nadam http://cs229.stanford.edu/proj2015/054_report.pdf

Best Learning Rate?

- Too big = diverge, too small = slow convergence
- No “one learning rate to rule them all”
- Start from a high value and keep cutting by half if model diverges
- Learning rate schedule: decay learning rate over time



Batch Normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Batch Normalization

```
1 # example of batch normalization for an cnn
2 from keras.layers import Dense
3 from keras.layers import Conv2D
4 from keras.layers import MaxPooling2D
5 from keras.layers import BatchNormalization
6 ...
7 model.add(Conv2D(32, (3,3), activation='relu'))
8 model.add(Conv2D(32, (3,3), activation='relu'))
9 model.add(BatchNormalization())
10 model.add(MaxPooling2D())
11 model.add(Dense(1))
12 ...
```