

Computer Vision Problems

- Image Classification
- Object Detection
- Neural Style Transfer and more ...

Deep Learning on Large Images



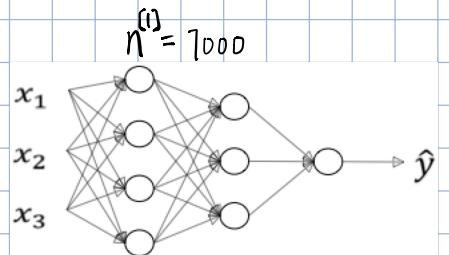
Low Resolution Image

$$64 \times 64 \times 3 = 12288, \quad W^{[1]} = (1000, 12288)$$



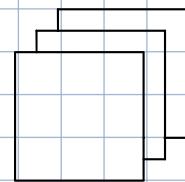
High Resolution Image

$$1000 \times 1000 \times 3 = 3 \times 10^6, \quad W^{[1]} = (1000, 3 \cdot 10^6)$$



- With 3 billions parameters, it's hard to get large amount of data to prevent overfitting.
- The memory consumption is so high that it's unrealistic to train it.

Edge Detection Example

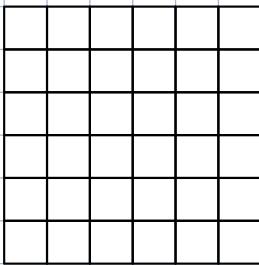


vertical edges



horizontal edges

Vertical Edge Detection

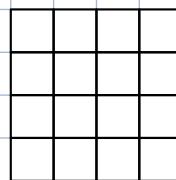


Convolution
*

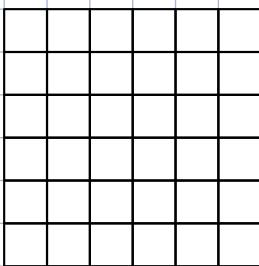
$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

Filter
Kernel

=



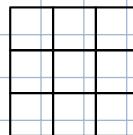
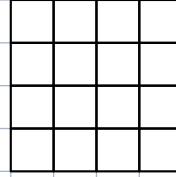
Horizontal Edge Detection



*

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

=



And other hand-crafted filters:

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

Sobel

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

Scharr

In convolutional neural network, it doesn't have hand-crafted filters.

Instead, it learns the nine parameters

$$\begin{array}{|c|c|c|} \hline w_1 & w_4 & w_7 \\ \hline w_2 & w_5 & w_8 \\ \hline w_3 & w_6 & w_9 \\ \hline \end{array}$$

. AMAZING!!!

Padding

Without Padding

$n \times n$

$$\begin{matrix} & \ast & \end{matrix} = \begin{matrix} (n-f+1) \times (n-f+1) \end{matrix}$$

First issue: the output size shrinks.

Second issue: the top-left corner pixel is used only once.

With Padding, $P = 1$

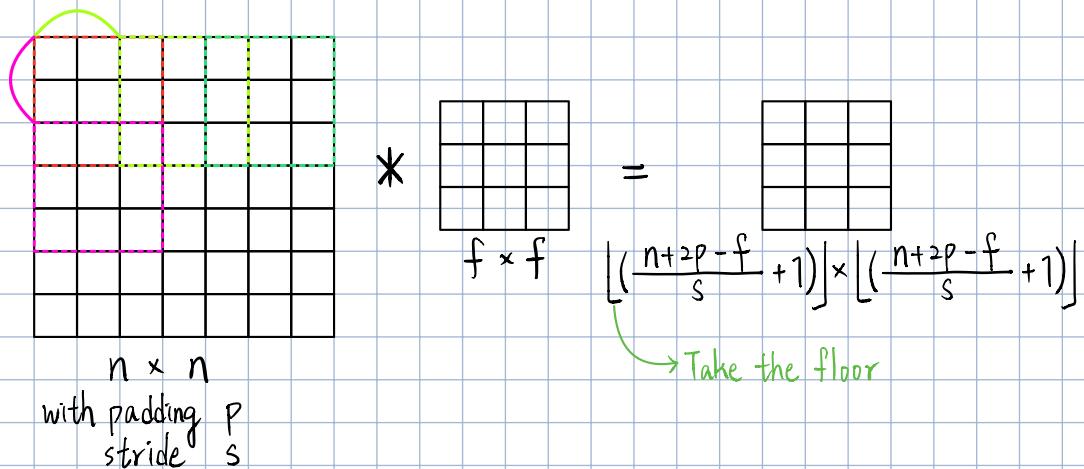
$(n+2p) \times (n+2p)$

$$\begin{matrix} & \ast & \end{matrix} = \begin{matrix} (n+2p-f+1) \times (n+2p-f+1) \end{matrix}$$

Convolution in Valid mode: No padding

Convolution in Same mode: Pad so that output size is the **Same** as the input size.

Strided Convolution



Cross-correlation v.s. Convolution

In theory, all convolution operation described above is actually cross-correlation.

The mathematical definition of convolution is the following:

Convolution \ast

3	4	5
1	0	2
-1	9	7

=

The kernel needs to be first flipped vertically and horizontally. And then does element-wise multiplication.

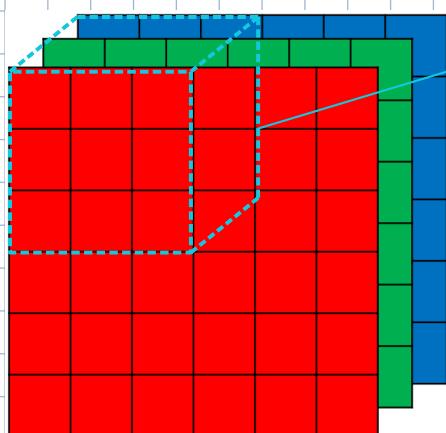
7	9	-1
2	0	1
5	4	3

However, many papers and Andrew don't bother call cross-correlation convolution.

Notice, in the signal processing world, convolution SHALL not be confused with cross-correlation.

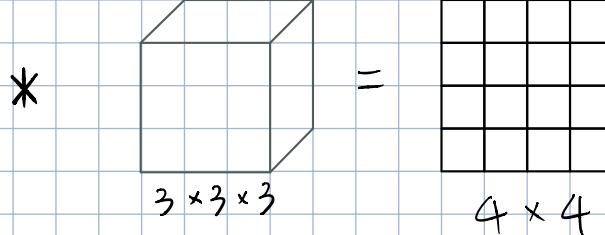
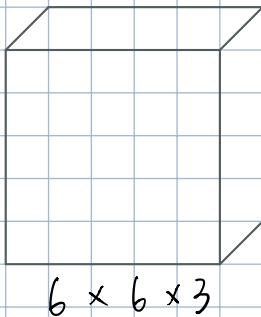
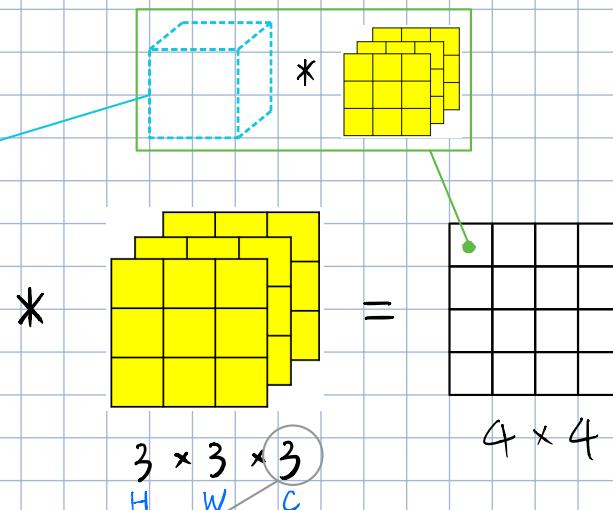
Convolutions Over Volume

Convolutions on RGB Image

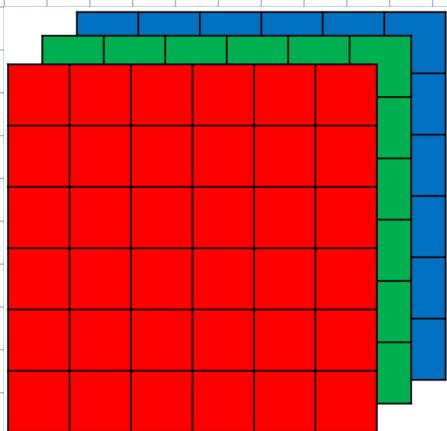


$6 \times 6 \times 3$
Height
Width
Channels

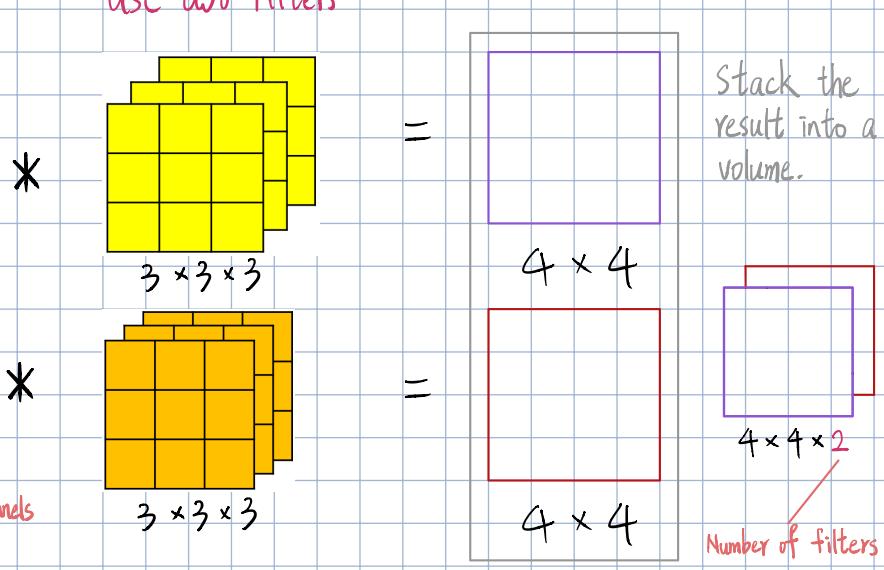
These two numbers have to be equal. (Channel-wise convolution)
The $3 \times 3 \times 3$ filter has 27 parameters to train.



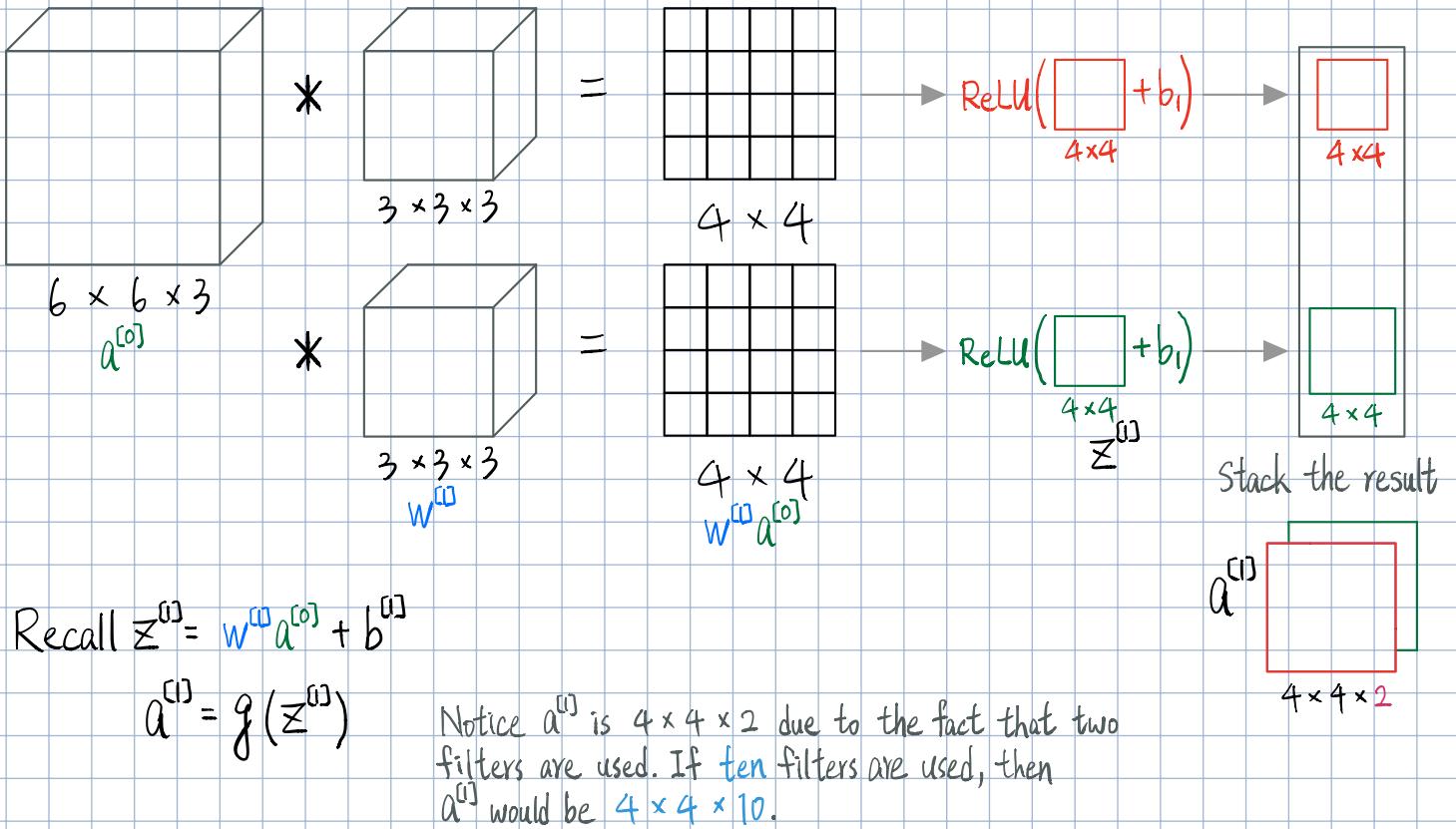
Using Multiple Filters



It's called number of channels or Depth.



One Layer of a Convolutional Network



Quick Quiz: If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?

Answer: $(3 \times 3 \times 3 + 1 \text{ bias}) \times 10 = 280 \text{ parameters}$

Summary of Notation: If layer ℓ is a convolution layer.

$f^{[\ell]}$ = filter size

$p^{[\ell]}$ = padding

$s^{[\ell]}$ = stride

$n_c^{[\ell]}$ = number of filters

Each filter is: $f^{[\ell]} \times f^{[\ell]} \times n_c^{[\ell-1]}$

Activations $a^{[\ell]}$: $n_H^{[\ell]} \times n_W^{[\ell]} \times n_C^{[\ell]}$

Weights: $f^{[\ell]} \times f^{[\ell]} \times n_c^{[\ell-1]} \times n_c^{[\ell]}$

bias: $n_c^{[\ell]}$, shape $(1, 1, 1, n_c^{[\ell]})$

Input: $n_H^{[\ell-1]} \times n_W^{[\ell-1]} \times n_C^{[\ell-1]}$ #channels in the previous $\ell-1$ layer

Output: $n_H^{[\ell]} \times n_W^{[\ell]} \times n_C^{[\ell]}$

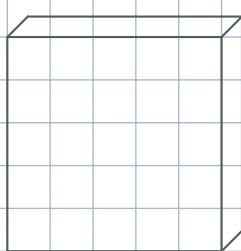
$$n_H^{[\ell]} = \left\lfloor \frac{n_H^{[\ell-1]} + 2p^{[\ell]} - f^{[\ell]}}{s^{[\ell]}} + 1 \right\rfloor \rightarrow \text{floor}$$

$$n_W^{[\ell]} = \left\lfloor \frac{n_W^{[\ell-1]} + 2p^{[\ell]} - f^{[\ell]}}{s^{[\ell]}} + 1 \right\rfloor$$

Batch $A^{[\ell]}$: $m \times n_H^{[\ell]} \times n_W^{[\ell]} \times n_C^{[\ell]}$
number of samples

Simple Convolutional Network Example

$39 \times 39 \times 3$ $\alpha^{[0]}$

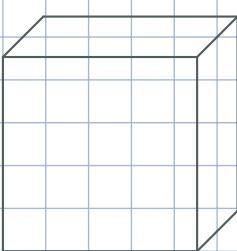


$l=0$

$$n_H^{[0]} = n_W^{[0]} = 39$$

$$n_C^{[0]} = 3$$

$37 \times 37 \times 10$ $\alpha^{[1]}$



$$\begin{aligned} f^{[1]} &= 3 \\ S^{[1]} &= 1 \\ P^{[1]} &= 0 \\ n_C^{[1]} &= 10 \end{aligned}$$

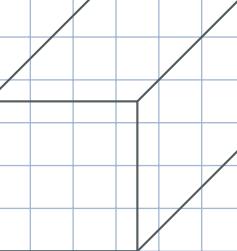
$l=1$

$$n_H^{[1]} = n_W^{[1]} = 37$$

$$n_C^{[1]} = 10$$

Example ConvNet

$17 \times 17 \times 20$ $\alpha^{[2]}$

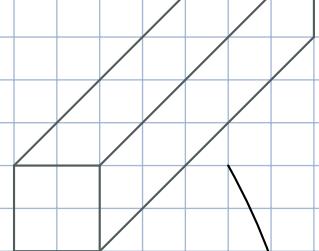


$l=2$

$$n_H^{[2]} = n_W^{[2]} = 17$$

$$n_C^{[2]} = 20$$

$7 \times 7 \times 40$ $\alpha^{[3]}$



$$\begin{aligned} f^{[3]} &= 5 \\ S^{[3]} &= 2 \\ P^{[3]} &= 0 \\ n_C^{[3]} &= 40 \end{aligned}$$

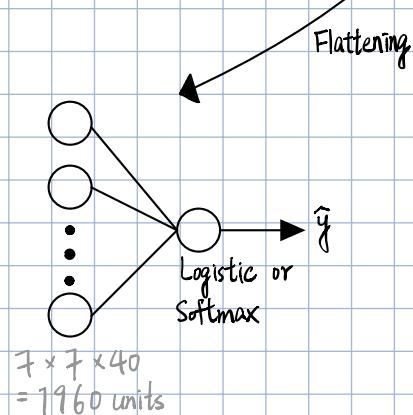
$l=3$

$$n_H^{[3]} = n_W^{[3]} = 7$$

$$n_C^{[3]} = 40$$

Types of layer in a convolutional network:

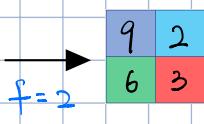
- Convolution (Conv)
- Pooling (Pool)
- Fully Connected (FC)



Pooling Layer : Max Pooling

stride

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2



filter size

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0

$$\begin{aligned} f &= 3 \\ s &= 1 \end{aligned}$$

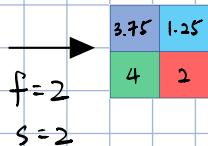
9	9	5
8	6	9

$$3 \times 3 \times n_C$$

$$5 \times 5 \times n_C$$

Pooling Layer : Average Pooling

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2



Doing pooling on 3D data is very similar to 2D.

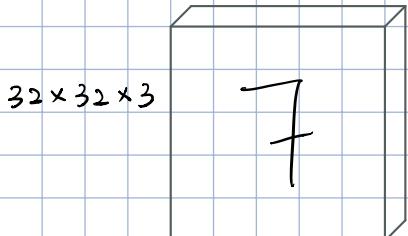
We just need to do pooling on each channel separately.

There are no parameters to learn in Pooling Layer.

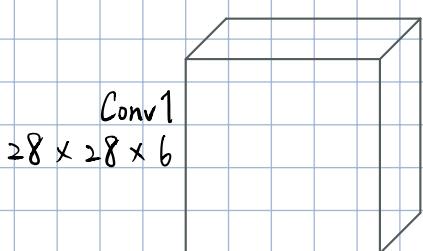
Convolutional Neural Network (CNN) Example

Inspired by LeNet-5

Goal: Build a neural network that can recognize 10 digits from 0 to 9.

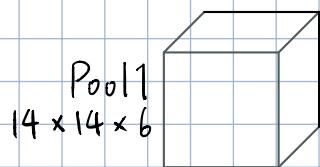


It's considered as the layer 1.



The pooling layer doesn't have weights to train. Thus, the convolutional layer and the pooling layer are considered as one layer.

MaxPooling
 $f=2$
 $s=2$

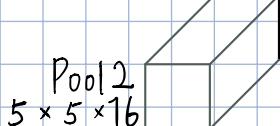


$f=5$
 $s=1$

Layer 2

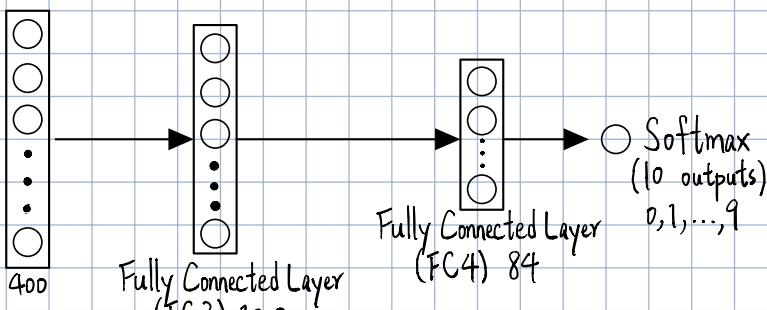
Conv2
 $10 \times 10 \times 16$

MaxPooling
 $f=2$
 $s=2$



Flattening
(No parameters here)

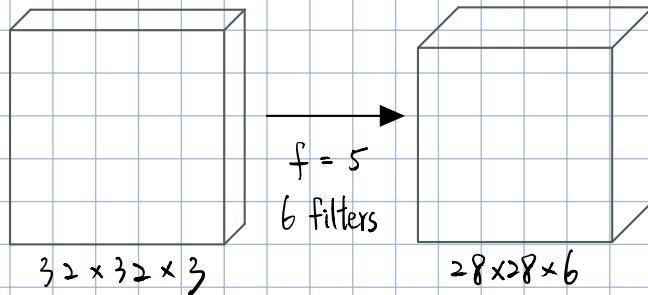
	Activation shape	Activation Size	# parameters
Input:	$(32, 32, 3)$	3,072	0
CONV1 ($f=5, s=1$)	$(28, 28, 6)$	<u>6,272</u> 2704	5x3x6+6 208 456
POOL1	$(14, 14, 6)$	<u>1,568</u> 1176	0
CONV2 ($f=5, s=1$)	$(10, 10, 16)$	<u>1,600</u>	416 416
POOL2	$(5, 5, 16)$	400	0
FC3	$(120, 1)$	120	120x480 48,000 48120
FC4	$(84, 1)$	84	84x120 10,080 10,084
Softmax	$(10, 1)$	10	841 850



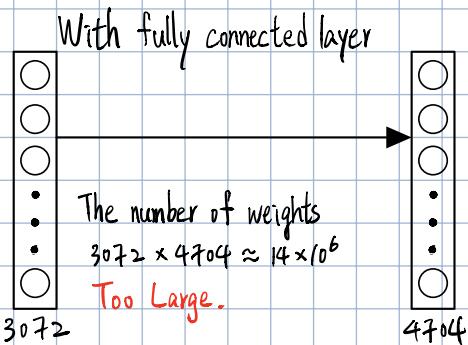
$$W^{[3]} : (120, 400)$$

$$b^{[3]} : (120,)$$

Why Convolutions?



With filters, there will only be $(5 \times 5 \times 3 + 1) \times 6$ weights to learn.



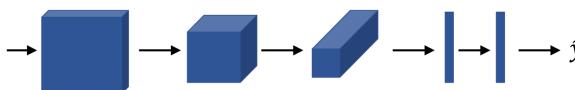
Quiz: The sparsity of connections and weight sharing are mechanisms that allow us to use fewer parameters in a CNN making it possible to train a network with smaller training set. Hmmm...

$\begin{matrix} & \ast & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} & = & \begin{matrix} \text{highlighted} & & \\ & & \end{matrix} \end{matrix}$

Parameter sharing: A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image. Only one filter/kernel is needed to run through the whole image to find features.

Sparsity of connections: In each layer, each output value depends only on a small number of inputs. The output pixel on the top-left corner is only depend on the 9 (3×3) out of 36 (6×6) pixels. The rest 27 ($36 - 9$) pixels do not have any effects on the output top-left corner pixel.

Putting it together



Training set: $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

Cost $J = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)})$. Use gradient descent to optimize parameters to reduce J .