# Convolutional Operation



## Vertical Edge Detection



$$\begin{array}{|c|c|c|c|c|c|}
\hline
3 & 0 & 1 & 2 & 7 & 4 \\
\hline
1 & 5 & 8 & 9 & 3 & 1 \\
\hline
2 & 7 & 2 & 5 & 1 & 3 \\
\hline
0 & 1 & 3 & 1 & 7 & 8 \\
\hline
4 & 2 & 1 & 6 & 2 & 8 \\
\hline
2 & 4 & 5 & 2 & 3 & 9 \\
\hline
\end{array}$$

6×6

"convolution"

$$* \quad \begin{array}{|c|c|c|}
\hline
1 & 0 & -1 \\
\hline
1 & 0 & -1 \\
\hline
1 & 0 & -1 \\
\hline
\end{array}$$

3×3

Filter / kernel

$$= \quad \begin{array}{|c|c|c|c|}
\hline
-5 & -4 & 0 & 8 \\
\hline
-10 & \cdot & & \\
\hline
 & & & \\
\hline
 & & & \\
\hline
\end{array}$$

4×4

Python :   conv-forward
tensorflow : tf. nn. conv2d
Keras :   Conv 2D

## How it works:

vertical edge detected

doesn't matter what is in the middle

bright pixels on left     dark pixels on right

negative if values flipped
* can take abs value

$$\begin{array}{|c|c|c|c|c|c|}
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
\end{array}$$

6×6

$$* \quad \begin{array}{|c|c|c|}
\hline
1 & 0 & -1 \\
\hline
1 & 0 & -1 \\
\hline
1 & 0 & -1 \\
\hline
\end{array}$$

3×3

$$= \quad \begin{array}{|c|c|c|c|}
\hline
0 & 30 & 30 & 0 \\
\hline
0 & 30 & 30 & 0 \\
\hline
0 & 30 & 30 & 0 \\
\hline
0 & 30 & 30 & 0 \\
\hline
\end{array}$$



Andrew Ng

edge detected

# Horizontal Edge Detection

$$
\begin{array}{|c|c|c|c|c|c|}
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
10 & 10 & 10 & 0 & 0 & 0 \\
\hline
0 & 0 & 0 & 10 & 10 & 10 \\
\hline
0 & 0 & 0 & 10 & 10 & 10 \\
\hline
0 & 0 & 0 & 10 & 10 & 10 \\
\hline
\end{array}
\quad * \quad
\begin{array}{|c|c|c|}
\hline
1 & 1 & 1 \\
\hline
0 & 0 & 0 \\
\hline
-1 & -1 & -1 \\
\hline
\end{array}
\quad = \quad
\begin{array}{|c|c|c|c|}
\hline
0 & 0 & 0 & 0 \\
\hline
30 & 10 & -10 & -30 \\
\hline
30 & 10 & -10 & -30 \\
\hline
0 & 0 & 0 & 0 \\
\hline
\end{array}
$$

horizontal filter/kernel

$6 \times 6$
$n \times n$

$3 \times 3$
$f \times f$

$4 \times 4$
$(n-f+1) \times (n-f+1)$

Other filter options:

→
$$
\begin{array}{|c|c|c|}
\hline
1 & 0 & -1 \\
\hline
2 & 0 & -2 \\
\hline
1 & 0 & 1 \\
\hline
\end{array}
$$
Sobel filter: puts emphasis on the central row/pixel

→
$$
\begin{array}{|c|c|c|}
\hline
3 & 0 & -3 \\
\hline
10 & 0 & -10 \\
\hline
3 & 0 & -3 \\
\hline
\end{array}
$$
Scharr filter: puts emphasis on the central row/pixel

\# flip to get horizontal edge detection

→
$$
\begin{array}{|c|c|c|}
\hline
W_1 & W_2 & W_3 \\
\hline
W_4 & W_5 & W_6 \\
\hline
W_7 & W_8 & W_9 \\
\hline
\end{array}
$$
Learn these as parameters such that this gives us a good edge detector

- Learning is even more robust
- helps learn these features its trying to detect
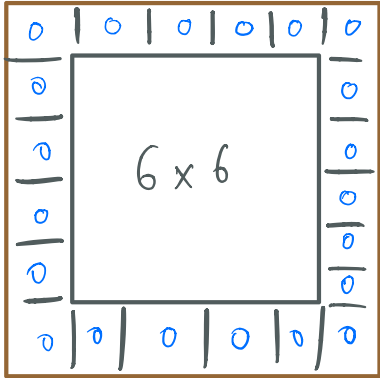- can learn edges at angles too

# Padding

$6 \times 6$ mat $\quad * \quad$ <sup>convolutional</sup> $\quad 3 \times 3$ mat $\quad = \quad 4 \times 4$ mat

problems
→ shrinking output
→ throwing away info from edges

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | | | | | 0 |
| 0 | | $6 \times 6$ | | | 0 |
| 0 | | | | | 0 |
| 0 | | | | | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |

$* \quad 3 \times 3$ mat $\quad = \quad 6 \times 6$ mat

$(n + 2p - f + 1) \times (n + 2p - f + 1)$

$8 \times 8$ mat <u>padded</u>

- pad with zeros
- $p$ = padding amount = 1 (in the example)

## Valid & Same convolutions

Valid → no padding applied

$$n \times n \quad * \quad f \times f = (n - f + 1) \times (n - f + 1)$$

Same → pad so the output is same as the input size

$$(n + 2p - f + 1) \times (n + 2p - f + 1) \quad * \quad f \times f = (n + 2p - f + 1) \times (n + 2p - f + 1)$$

★ $f$ is usually odd

# Strided Convolutions



stride = 2

6 × 6

dropped because of floor

dimensions –

$$n \times n \quad * \quad f \times f \quad \longrightarrow \quad \left\lfloor \frac{n + 2p - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n + 2p - f}{s} + 1 \right\rfloor$$
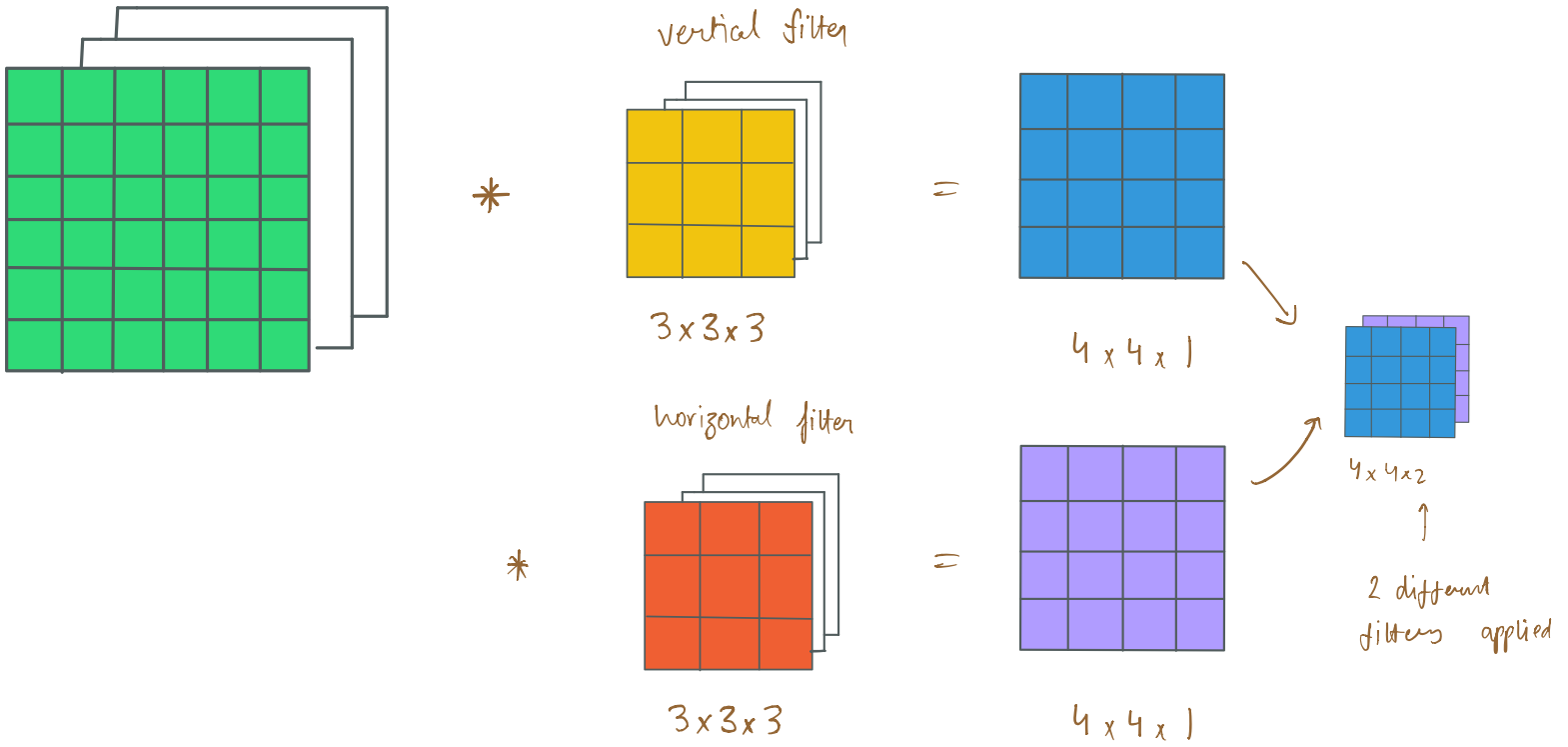
Padding P          stride S

# Convolutions Over Volume



convolve

*

=

6 × 6 × 3

height    width    # channels

3 × 3 × 3

4 × 4 × 1

edges in red channel

| R | | | G | | | B | | |
|---|---|---|---|---|---|---|---|---|
| 1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | -1 | 0 | 0 | 0 | 0 | 0 | 0 |

# Multiple Filters



vertical filter

$3 \times 3 \times 3$

$4 \times 4 \times 1$

$4 \times 4 \times 2$

2 different filters applied

horizontal filter

$3 \times 3 \times 3$

$4 \times 4 \times 1$
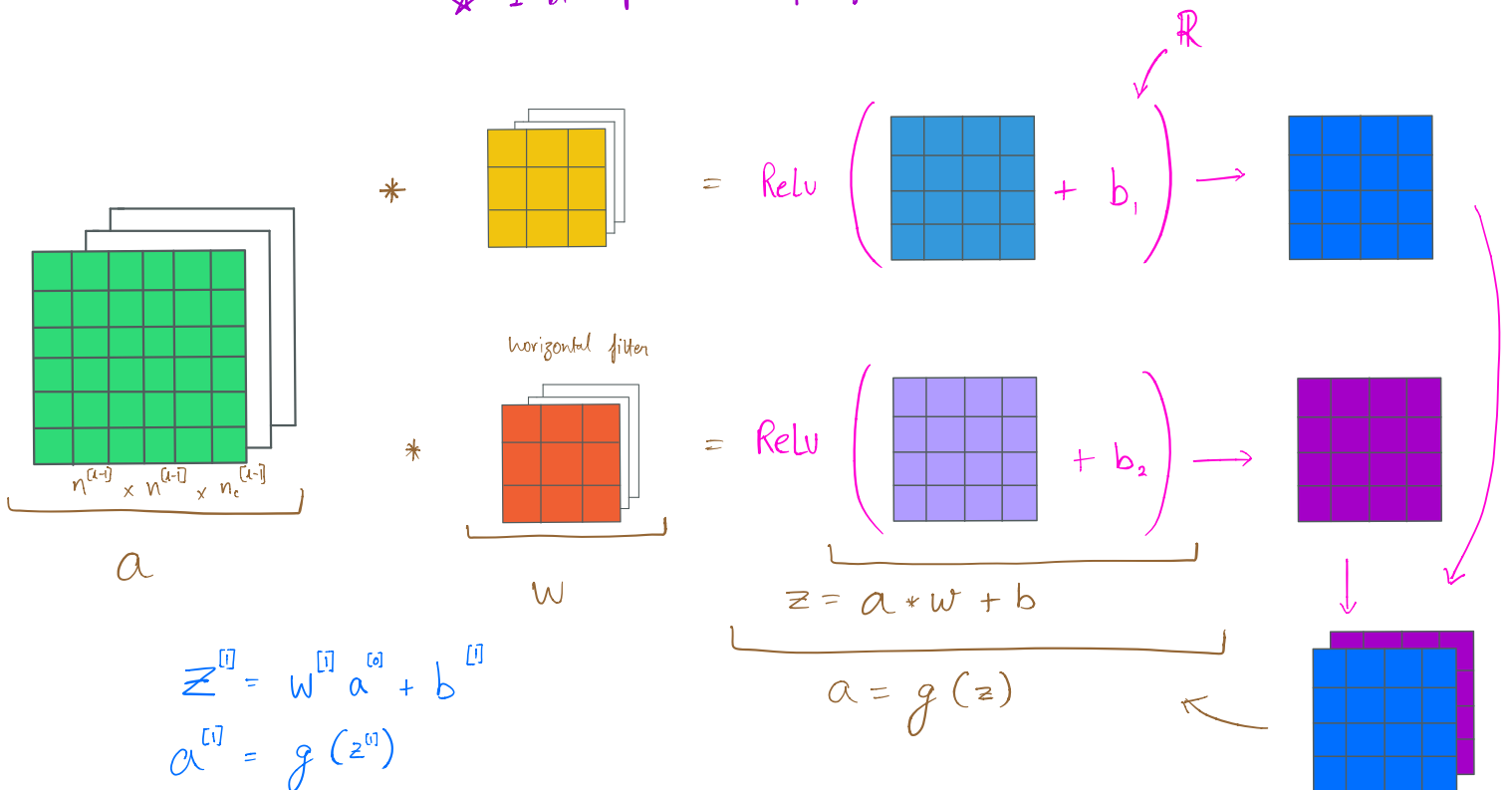
dimension:

$$n \times n \times n_c \quad * \quad f \times f \times n_c \rightarrow \frac{n-f+1}{4} \times \frac{n-f+1}{4} \times n_c$$

# filters
depth

★ 1 bias parameter per filter

$\mathbb{R}$

$= \text{Relu} \left( \phantom{xxx} + b_1 \right) \rightarrow$

horizontal filter

$= \text{Relu} \left( \phantom{xxx} + b_2 \right) \rightarrow$

$\underbrace{n^{[l-1]} \times n^{[l-1]} \times n_c^{[l-1]}}$

$a$

$W$

$\underbrace{z = a * w + b}$

$a = g(z)$

$$z^{[l]} = W^{[l]} a^{[0]} + b^{[l]}$$
$$a^{[l]} = g(z^{[l]})$$

derivatives —
$$dA \mathrel{+}= \sum_{h=0}^{n_H} \sum_{w=0}^{n_w} W_c \times dZ_{hw}$$

$$dW_c \mathrel{+}= \sum_{h=0}^{n_H} \sum_{h=0}^{n_w} a_{slice} \times dZ_{hw}$$

$$db = \sum \sum dZ_{hw}$$



$6 \times 6 \times 3$          $4 \times 4 \times 2$
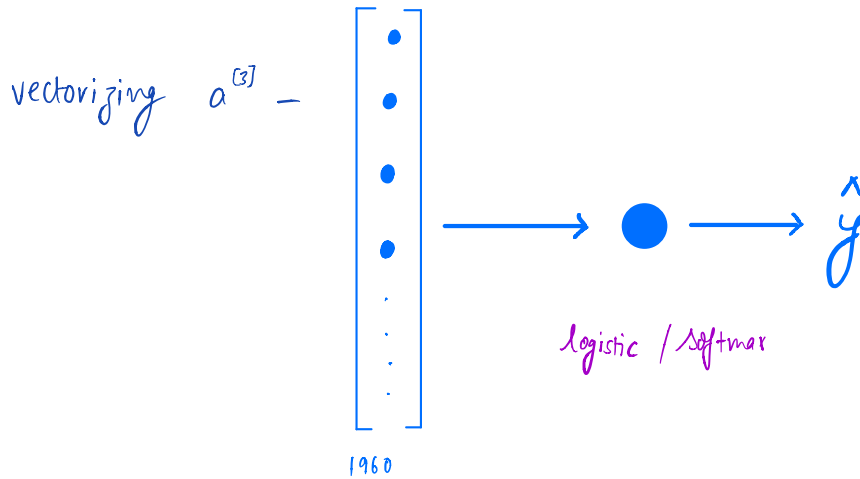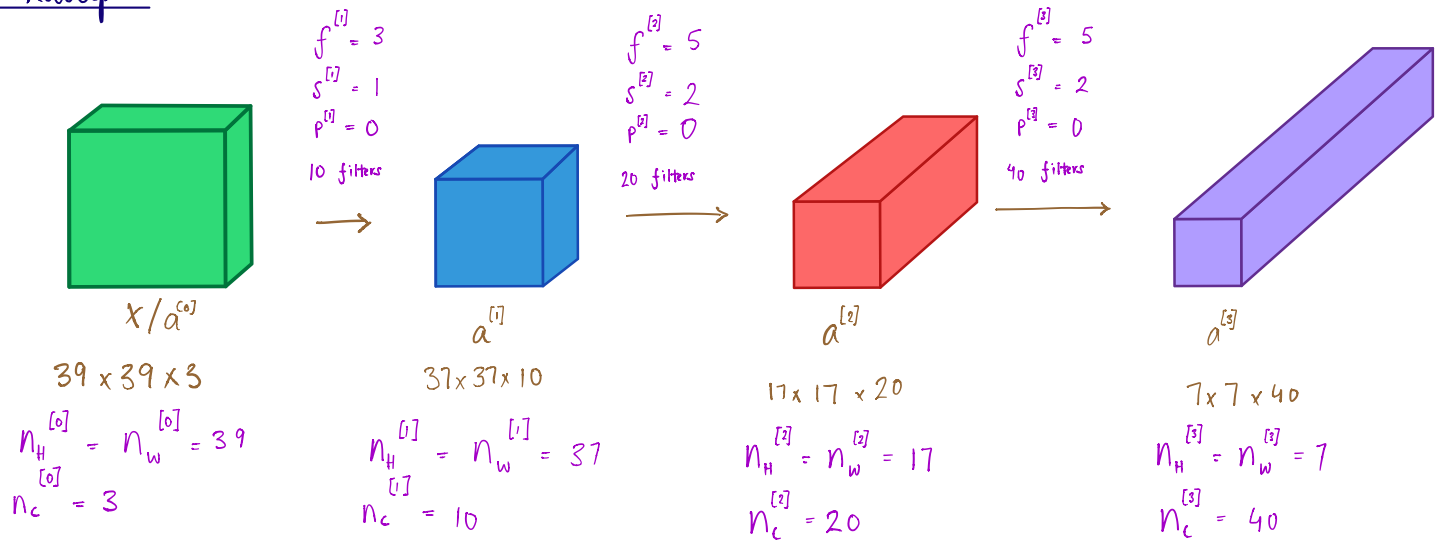
$a^{[0]}$          $a^{[i]}$

## Notation

If layer 1 is a convolution layer —

- $f^{[l]}$ = filter size

- $p^{[l]}$ = padding

- $s^{[l]}$ = stride

- input — $n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$

- output — $n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

- $n^{[l]} = \left\lfloor \dfrac{n^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$

- filter — $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

- activations - $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

- weights - $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

  depth from last layer $\rfloor$      $\lfloor$ # of filters in layer $l$

- bias - $n_c^{[l]} \rightarrow 1, 1, 1, n_c^{[l]}$

## Example



$x/a^{[0]}$     $a^{[1]}$     $a^{[2]}$     $a^{[3]}$

$39 \times 39 \times 3$     $37 \times 37 \times 10$     $17 \times 17 \times 20$     $7 \times 7 \times 40$

$n_H^{[0]} = n_W^{[0]} = 39$    $n_H^{[1]} = n_W^{[1]} = 37$    $n_H^{[2]} = n_W^{[2]} = 17$    $n_H^{[3]} = n_W^{[3]} = 7$

$n_c^{[0]} = 3$     $n_c^{[1]} = 10$     $n_c^{[2]} = 20$     $n_c^{[3]} = 40$

$f^{[1]} = 3$, $s^{[1]} = 1$, $p^{[1]} = 0$, 10 filters

$f^{[2]} = 5$, $s^{[2]} = 2$, $p^{[2]} = 0$, 20 filters

$f^{[3]} = 5$, $s^{[3]} = 2$, $p^{[3]} = 0$, 40 filters

vectorizing $a^{[3]}$ —

1960

logistic / softmax

$\hat{y}$

Types of layers in a convolutional network –

$\rightarrow$ convolution (CONV)

$\rightarrow$ pooling (POOL)

$\rightarrow$ fully connected (FC)

# Pooling Layers

## Max Pooling



hyperparameters

$f = 2$
$s = 2$

getting the max
from each color

- Some set of features or activations from previous layer, then max pooling detects a particular feature and preserves it due to the high number

→ no parameters, only hyperparameters ⇒ fixed computation

example:



5x5

$f = 3$
$s = 1$

3x3

$$\frac{n + 2p - f}{s} + 1 = 3$$

# Average Pooling



$f = 2$
$s = 2$

$7 \times 7 \times 1000 \longrightarrow 1 \times 1 \times 1000$

→ helps make the size of the matrix smaller by taking the averages

# Example of a full Convolutional Network



| | Activation shape | Activation Size | # parameters |
|---|---|---|---|
| Input: | (32,32,3) | 3,072 | 0 |
| CONV1 (f=5, s=1) | (28,28,8) | 6,272 | 608 |
| POOL1 | (14,14,8) | 1,568 | 0 |
| CONV2 (f=5, s=1) | (10,10,16) | 1,600 | 3216 |
| POOL2 | (5,5,16) | 400 | 0 |
| FC3 | (120,1) | 120 | 48120 |
| FC4 | (84,1) | 84 | 10164 |
| Softmax | (10,1) | 10 | 850 |

# Why Convolutions

CONV nets have comparatively way less parameters

↳ Parameter Sharing - A feature detector like vertical/horizontal edge detectors are useful in all parts of the image

↳ Sparsity of connections - in each layer, each output value depends only on a small number of inputs

As we take a picture suppose RGB, and we apply a filter/kernel to it, we are able to make that picture go through a filter like edge detections and leave out other features of the image. The image we get is particular to that filter. As we apply more filters/kernels to the image, we can stack these output images and we get images that have certain filters applied to them and these images tend to be smaller in width and height as compared to the actual images.
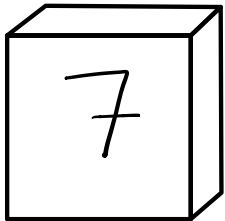
We can view this as a neural network if we try to stack all the pixels of the original picture and connect it to the weights which are the values of the pixels in the filter. The concept of convolutional neural networks is that we connect a lot of these pixels values / input units to the weights / hidden units but not all of them to each one of the hidden units. The missing connections or the less connections help us to end up with an easier network with less parameters than a traditional neural network where all input units are connected to every one of the hidden units.

As use this filter to get another image that we feed again to another convolutional network, this helps us separate or individually feed these features to further layers which eventually help us identify what the image is. We train the parameters, i.e the pixel values of the filter, to make new filters that help us detect features important in identifying the image. As we apply the filter we get a part or a feature of the image important for our classification.
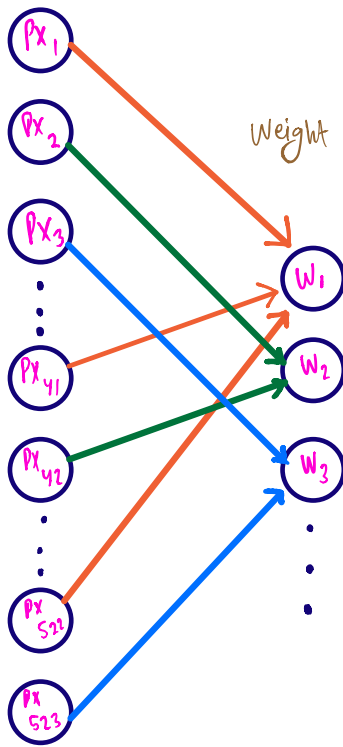


What pooling does is it helps us to kind of average out and make the size of the height and width smaller to prevent having a large number of parameters. It is like a control check on the network and helps us keep the number of parameters in control.

Pixel matrix

Orignal image

Weight values

Trained filter

$Px_1$
$Px_2$
$Px_3$
$Px_{41}$
$Px_{42}$
$Px_{521}$
$Px_{523}$

$W_1$
$W_2$
$W_3$

| $W_1$ | $W_2$ | $W_3$ |
|-------|-------|-------|
| $W_4$ | $W_5$ | $W_6$ |
| $W_7$ | $W_8$ | $W_9$ |