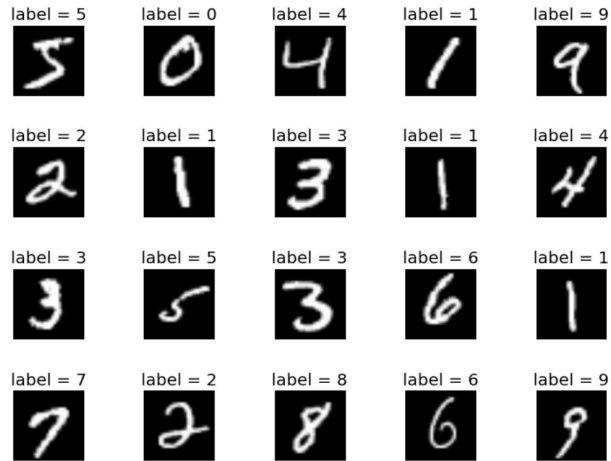


APS360: Applied Fundamentals of Deep Learning

Week 4: Convolutional Neural Network - Part I

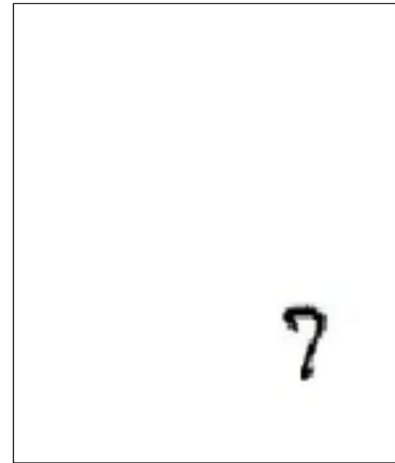
Last Time: MNIST with ANN

MNIST dataset is nicely preprocessed, centred, etc.



MNIST Data (Cleaned)

What happens if the new sample is not preprocessed?



New Data Sample (Raw)

What if we make our deep NN bigger?

Input: 200×200 pixels \rightarrow 1st hidden layer: 500 \rightarrow 2nd hidden layer: 200

Q: How many weights are there?

What if we make our deep NN bigger?

Input: 200×200 pixels \rightarrow 1st hidden layer: 500 \rightarrow 2nd hidden layer: 200

Q: How many weights are there?

Using a large fully connected network has some down sides

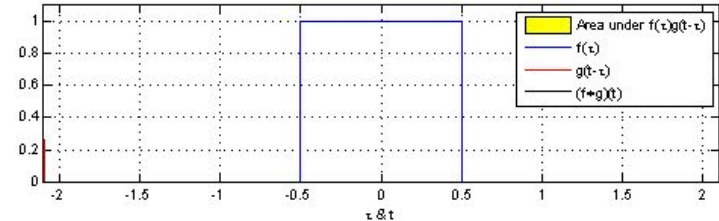
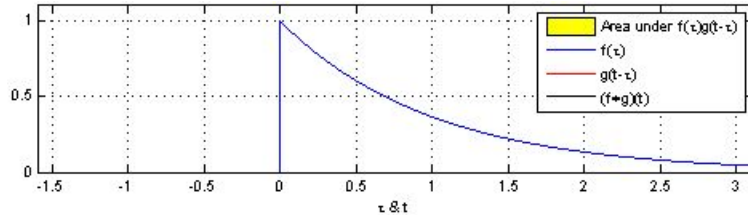
- **Computational complexity grows** \rightarrow harder to train, **more data** to **generalize**
- **Bad inductive bias** \rightarrow Ignores geometry of image data
- **Not flexible** \rightarrow Different image sizes require different models

Convolution Operator

Convolution operator in One-Dimension

Convolution is a mathematical operation on two functions f and g that expresses how the shape of one is modified by the other.

$$(f * g)[n] = \sum_{k=-\infty}^{\infty} f[k] g[n - k]$$



Convolution in 2D for images

Convolution of Image **I** with filter kernel **K**

1. Multiply each pixel in **I** in range of kernel by the corresponding element of kernel **K**
2. Sum all these products and write to a new 2d array
3. Slide kernel across all areas of the image until you reach the ends.

$$y[m, n] = I[m, n] * K[m, n] = \sum_{j=-\infty}^{\infty} \sum_{i=-\infty}^{\infty} I[i, j].K[m - i, n - j]$$

Convolution in 2D for images

Convolution of Image I with filter kernel K

1. Multiply each pixel in I in range of kernel by the corresponding element of kernel K
2. Sum all these products and write to a new 2d array
3. Slide kernel across all areas of the image until you reach the ends.

$$y[m, n] = I[m, n] * K[m, n] = \sum_{j=-\infty}^{\infty} \sum_{i=-\infty}^{\infty} I[i, j].K[m - i, n - j]$$

I

| | | | | |
|---|---|---|---|---|
| 7 | 2 | 3 | 3 | 8 |
| 4 | 5 | 3 | 8 | 4 |
| 3 | 3 | 2 | 8 | 4 |
| 2 | 8 | 7 | 2 | 7 |
| 5 | 4 | 4 | 5 | 4 |

$*$

K

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

$=$

y

| | | |
|---|--|--|
| 6 | | |
| | | |
| | | |

$7 \times 1 + 4 \times 1 + 3 \times 1 +$
 $2 \times 0 + 5 \times 0 + 3 \times 0 +$
 $3 \times -1 + 3 \times -1 + 2 \times -1$
 $= 6$

Convolution in 2D for images

Convolution of Image **I** with filter kernel **K**

1. Multiply each pixel in **I** in range of kernel by the corresponding element of kernel **K**
2. Sum all these products and write to a new 2d array
3. Slide kernel across all areas of the image until you reach the ends.

$$y[m, n] = I[m, n] * K[m, n] = \sum_{j=-\infty}^{\infty} \sum_{i=-\infty}^{\infty} I[i, j].K[m - i, n - j]$$

I

| | | | | |
|---|---|---|---|---|
| 7 | 2 | 3 | 3 | 8 |
| 4 | 5 | 3 | 8 | 4 |
| 3 | 3 | 2 | 8 | 4 |
| 2 | 8 | 7 | 2 | 7 |
| 5 | 4 | 4 | 5 | 4 |

K

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

*

y

| | | |
|---|--|--|
| 6 | | |
| | | |
| | | |

=

$7 \times 1 + 4 \times 1 + 3 \times 1 + 2 \times 0 + 5 \times 0 + 3 \times 0 + 3 \times -1 + 3 \times -1 + 2 \times -1 = 6$

I

| | | | | |
|---|---|---|---|---|
| 7 | 2 | 3 | 3 | 8 |
| 4 | 5 | 3 | 8 | 4 |
| 3 | 3 | 2 | 8 | 4 |
| 2 | 8 | 7 | 2 | 7 |
| 5 | 4 | 4 | 5 | 4 |

*

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

*

| | | |
|---|----|--|
| 6 | -9 | |
| | | |
| | | |

=

$2 \times 1 + 5 \times 1 + 3 \times 1 + 3 \times 0 + 3 \times 0 + 2 \times 0 + 3 \times -1 + 8 \times -1 + 8 \times -1 = -9$

Convolution in 2D for images

| | | | | |
|---|---|---|---|---|
| 7 | 2 | 3 | 3 | 8 |
| 4 | 5 | 3 | 8 | 4 |
| 3 | 3 | 2 | 8 | 4 |
| 2 | 8 | 7 | 2 | 7 |
| 5 | 4 | 4 | 5 | 4 |

*

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

=

| | | |
|---|--|--|
| 6 | | |
| | | |
| | | |

$$\begin{aligned} &7 \times 1 + 4 \times 1 + 3 \times 1 + \\ &2 \times 0 + 5 \times 0 + 3 \times 0 + \\ &3 \times -1 + 3 \times -1 + 2 \times -1 \\ &= 6 \end{aligned}$$

Convolution in 2D for images

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

K

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

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

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

Convolution in 2D for images

What effect would the following filter have on an image?



* 1/9

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

Convolution in 2D for images

What effect would the following filter have on an image?

Blurring averages out pixel intensities in an image.

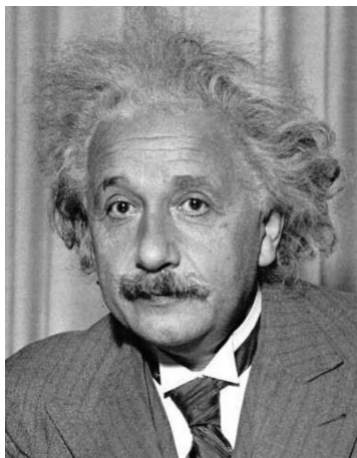


$$* \frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} \rightarrow$$



Convolution in 2D for images

What effect would the following filter have on an image?



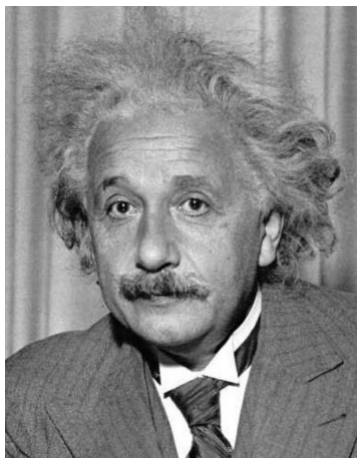
*

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

Convolution in 2D for images

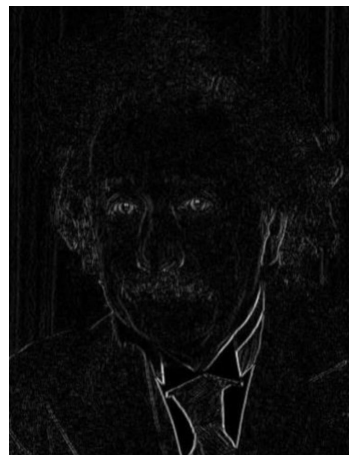
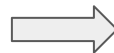
What effect would the following filter have on an image?

Vertical edge detector



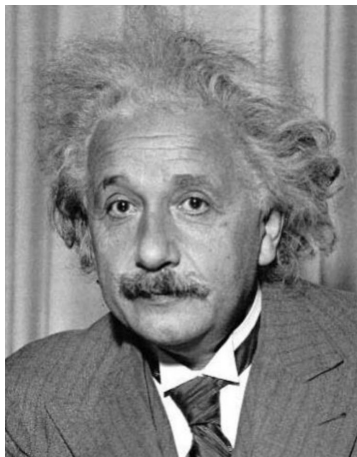
*

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



Convolution in 2D for images

What effect would the following filter have on an image?



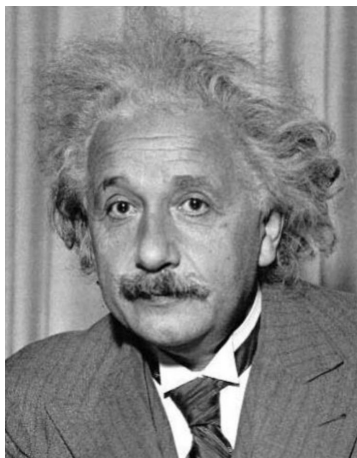
*

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

Convolution in 2D for images

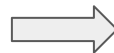
What effect would the following filter have on an image?

Horizontal edge detector



*

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



Convolution in 2D for images

What effect would the following filter have on an image?



*

$$\begin{pmatrix} 0 & 0 & 3 & 2 & 2 & 2 & 3 & 0 & 0 \\ 0 & 2 & 3 & 5 & 5 & 5 & 3 & 2 & 0 \\ 3 & 3 & 5 & 3 & 0 & 3 & 5 & 3 & 3 \\ 2 & 5 & 3 & -12 & -23 & -12 & 3 & 5 & 2 \\ 2 & 5 & 0 & -23 & -40 & -23 & 0 & 5 & 2 \\ 2 & 5 & 3 & -12 & -23 & -12 & 3 & 5 & 2 \\ 3 & 3 & 5 & 3 & 0 & 3 & 5 & 3 & 3 \\ 0 & 2 & 3 & 5 & 5 & 5 & 3 & 2 & 0 \\ 0 & 0 & 3 & 2 & 2 & 2 & 3 & 0 & 0 \end{pmatrix}$$

Convolution in 2D for images

What effect would the following filter have on an image?

Blob detector → Regions that differ in properties, such as brightness or color, compared to surrounding regions



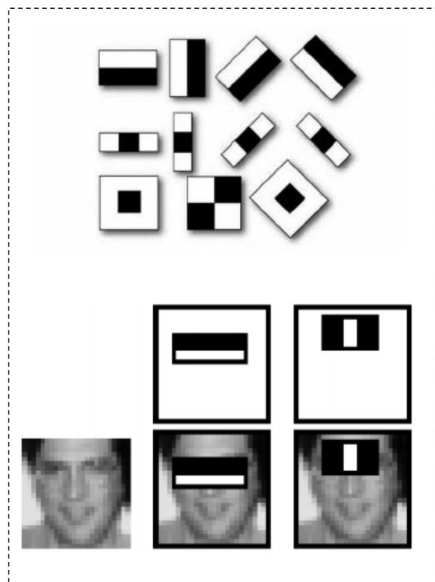
*

$$\begin{pmatrix} 0 & 0 & 3 & 2 & 2 & 2 & 3 & 0 & 0 \\ 0 & 2 & 3 & 5 & 5 & 5 & 3 & 2 & 0 \\ 3 & 3 & 5 & 3 & 0 & 3 & 5 & 3 & 3 \\ 2 & 5 & 3 & -12 & -23 & -12 & 3 & 5 & 2 \\ 2 & 5 & 0 & -23 & -40 & -23 & 0 & 5 & 2 \\ 2 & 5 & 3 & -12 & -23 & -12 & 3 & 5 & 2 \\ 3 & 3 & 5 & 3 & 0 & 3 & 5 & 3 & 3 \\ 0 & 2 & 3 & 5 & 5 & 5 & 3 & 2 & 0 \\ 0 & 0 & 3 & 2 & 2 & 2 & 3 & 0 & 0 \end{pmatrix}$$

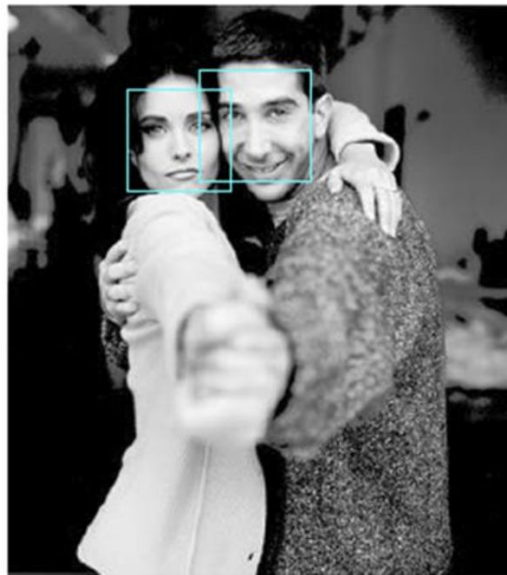


Convolution in 2D for images

Even more complicated kernels for detecting faces



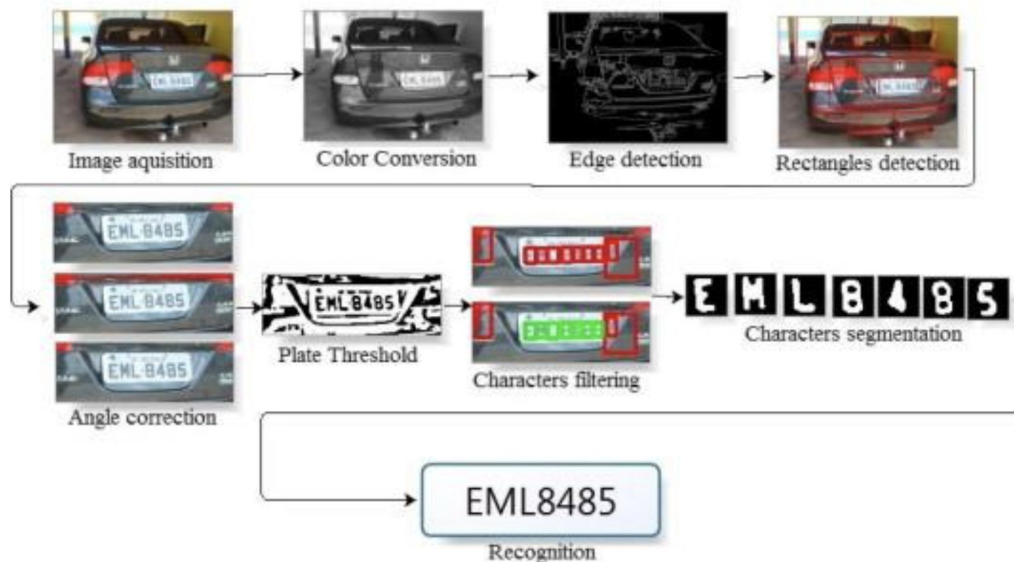
Convolutional kernels



Convolution in 2D for images

Where the kernels come from? They used to be **hand-crafted**

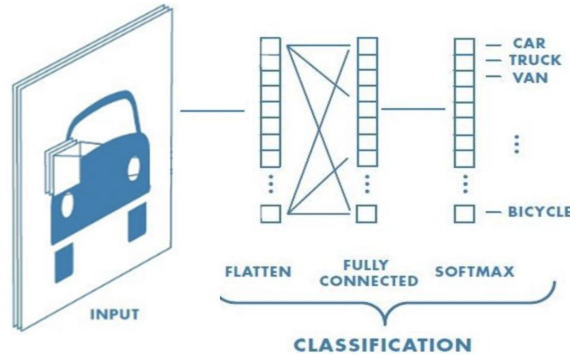
Classic computer vision → **multi-stage feature (kernel) engineering**



Convolutional Neural Networks

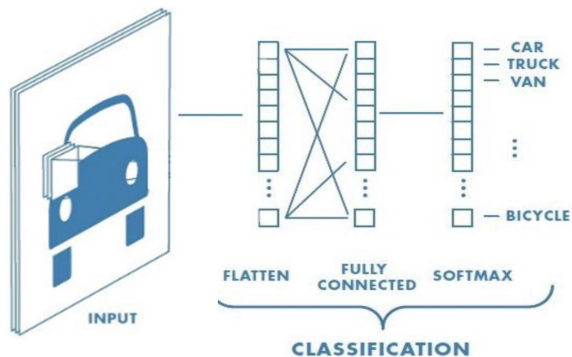
Convolutional Neural Networks

MLP → Requires that we preprocess the input!

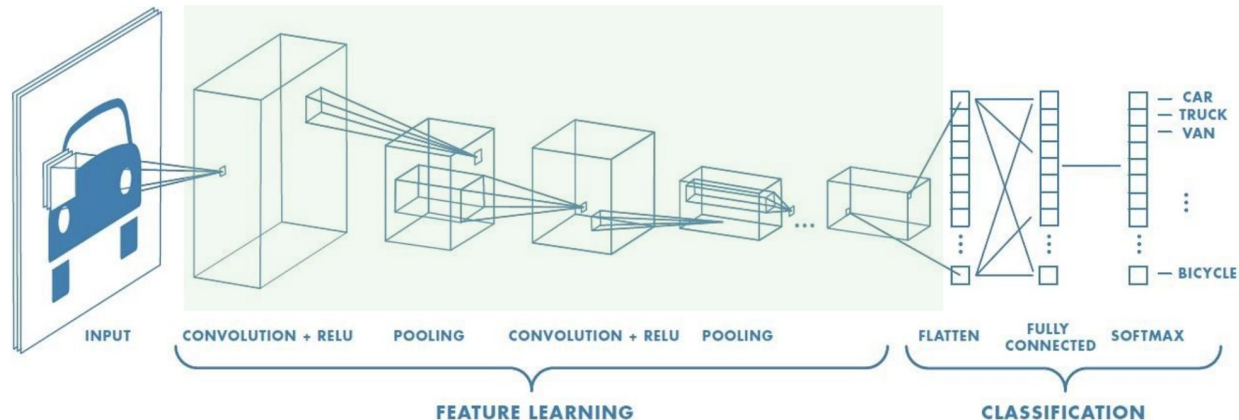


Convolutional Neural Networks

MLP → Requires that we preprocess the input!



CNN → Apply convolution to image tensors.



Forward & Backward pass

- Initialize the kernels randomly
- In forward pass, convolve the image with the kernel
- In backward pass, update the kernel using gradients

input

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

kernel (weights)

| | | |
|-----------------|-----------------|-----------------|
| W ₀₀ | W ₀₁ | W ₀₂ |
| W ₁₀ | W ₁₁ | W ₁₂ |
| W ₂₀ | W ₂₁ | W ₂₂ |

➤ The kernel or filter contains the **trainable weights**. In this picture, the kernel size is 3 × 3

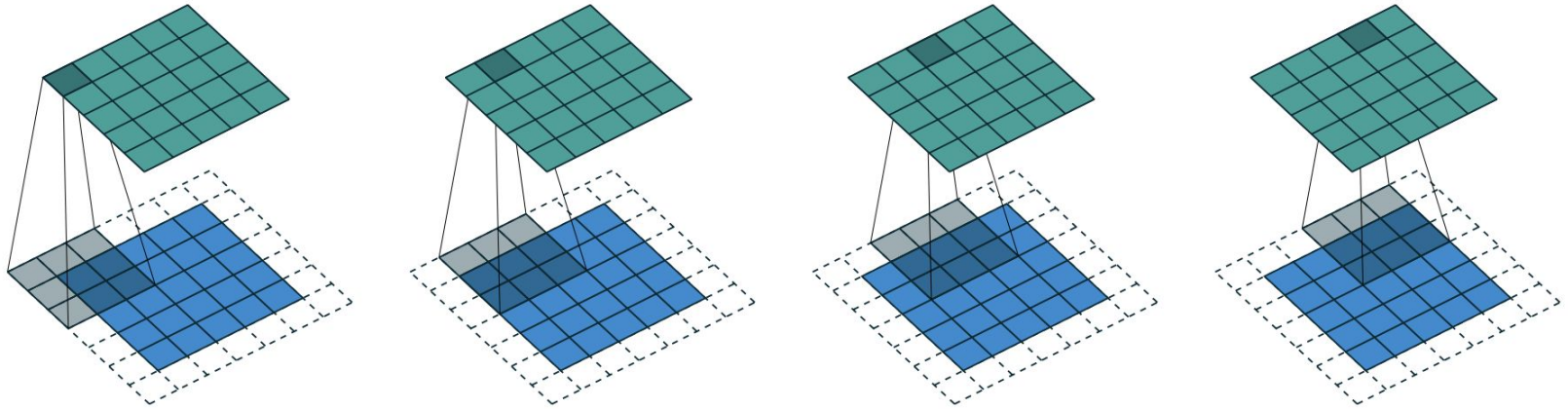
output

| | | |
|------|------|------|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

Zero Padding

Adding zeros around the border of the image before convolution. Why?

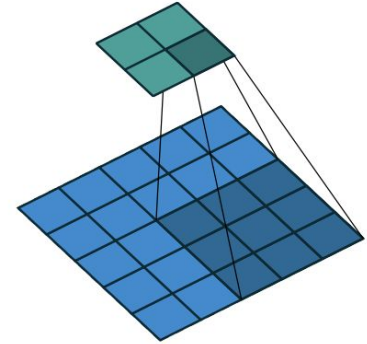
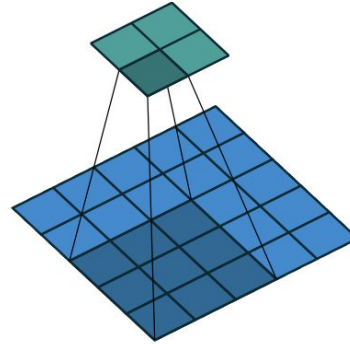
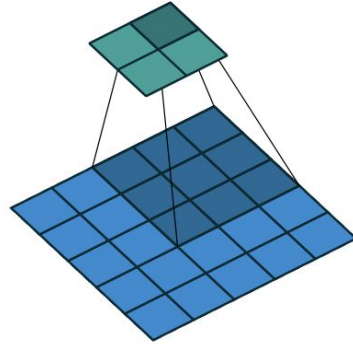
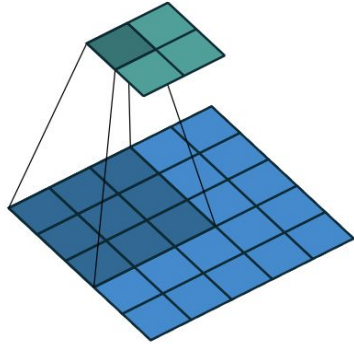
- Keep width and height consistent with the previous layer
- Keep the information around the border of the image



Stride

Distance between two consecutive positions of the kernel

- Allows us to control the output resolution



Computing the output size

For each dimension of an input image with

- Image dimension of size i
- Kernel of size k
- Padding of size p
- Stride of size s

The size of output dimension is computed by:

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$$

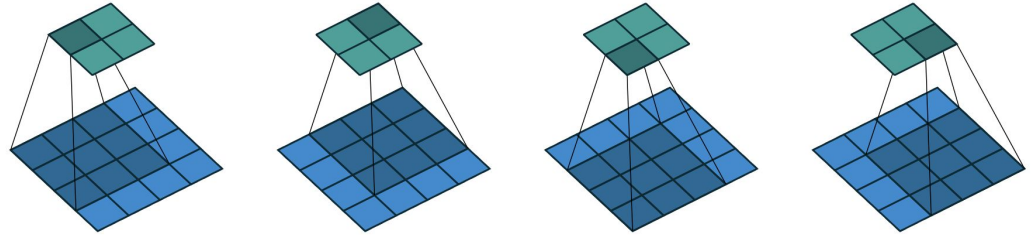
Example

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$$

- No padding, unit stride \rightarrow (the output size $<$ the input size)

$$i = 4, k = 3, s = 1 \text{ and } p = 0$$

$$o = \left\lfloor \frac{4 + 2 \times 0 - 3}{1} \right\rfloor + 1 = 2$$



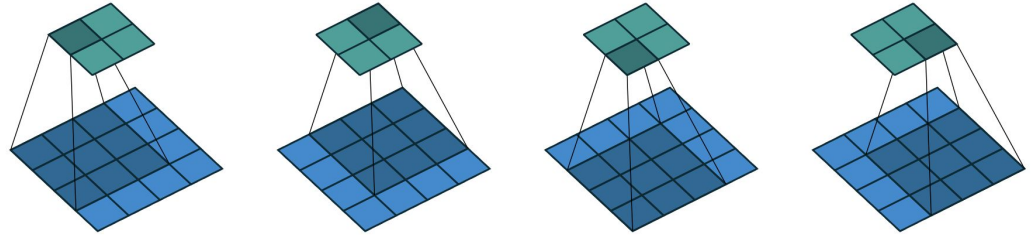
Example

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$$

- No padding, unit stride \rightarrow (the output size $<$ the input size)

$$i = 4, k = 3, s = 1 \text{ and } p = 0$$

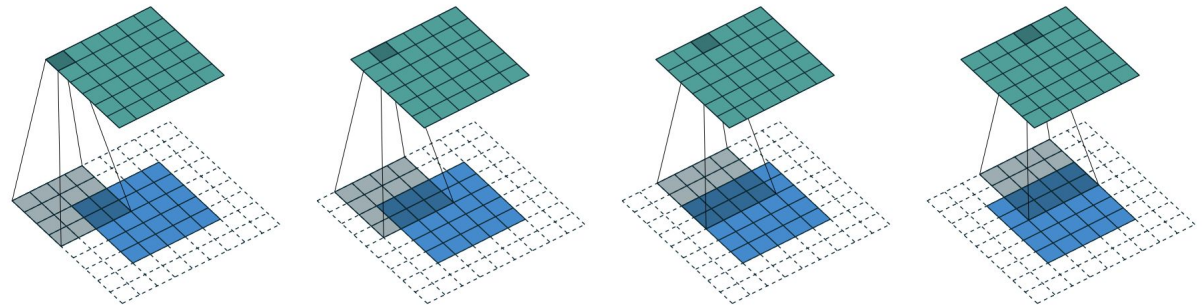
$$o = \left\lfloor \frac{4 + 2 \times 0 - 3}{1} \right\rfloor + 1 = 2$$



- Arbitrary padding, unit strides

$$i = 5, k = 4, s = 1 \text{ and } p = 2$$

$$o = \left\lfloor \frac{5 + 2 \times 2 - 4}{1} \right\rfloor + 1 = 6$$



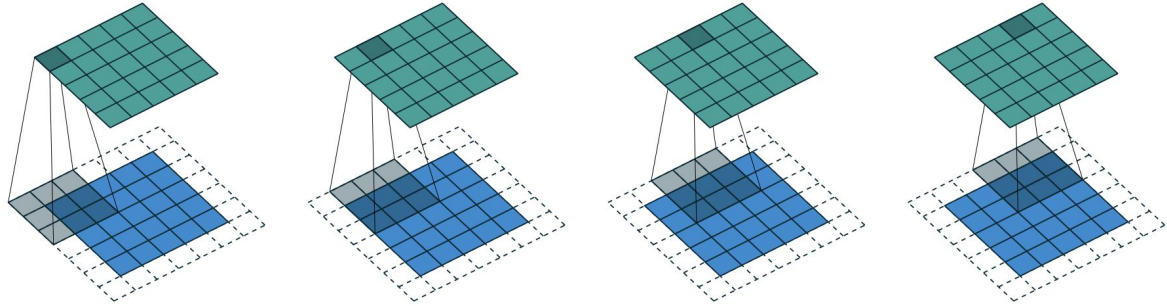
Example

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$$

- Half padding, unit strides \rightarrow (the output size be the same as the input size)

$i = 5, k = 3, s = 1$ and $p = 1$

$$o = \left\lfloor \frac{5 + 2 \times 1 - 3}{1} \right\rfloor + 1 = 5$$



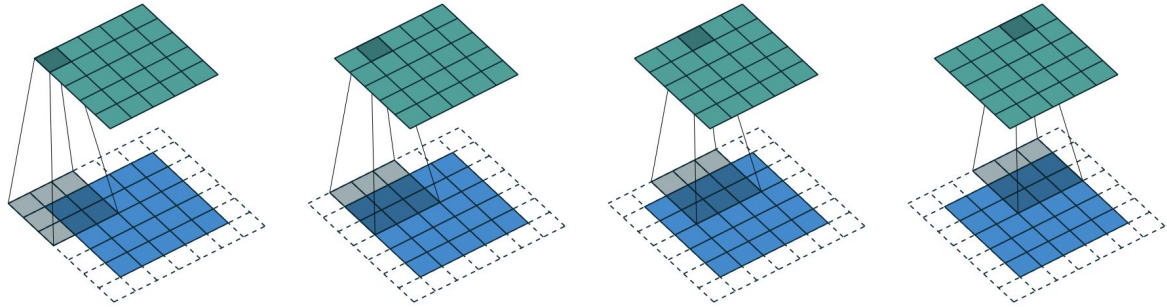
Example

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$$

- Half padding, unit strides \rightarrow (the output size be the same as the input size)

$i = 5, k = 3, s = 1$ and $p = 1$

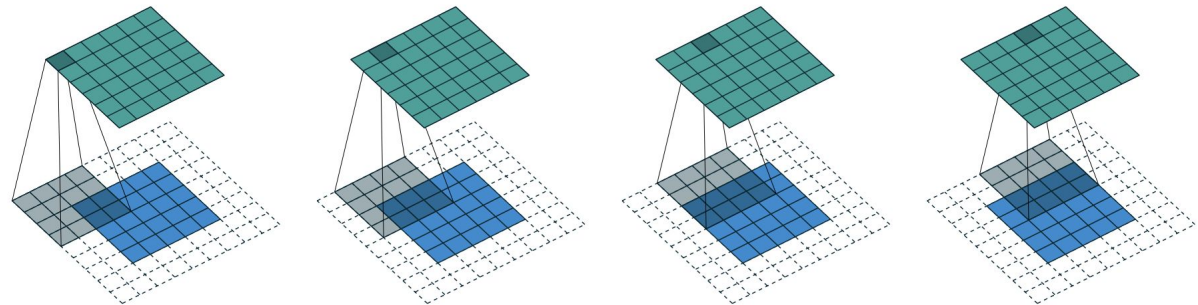
$$o = \left\lfloor \frac{5 + 2 \times 1 - 3}{1} \right\rfloor + 1 = 5$$



- Full padding, unit strides \rightarrow (the output size \geq the input size)

$i = 5, k = 3, s = 1$ and $p = 2$

$$o = \left\lfloor \frac{5 + 2 \times 2 - 3}{1} \right\rfloor + 1 = 6$$



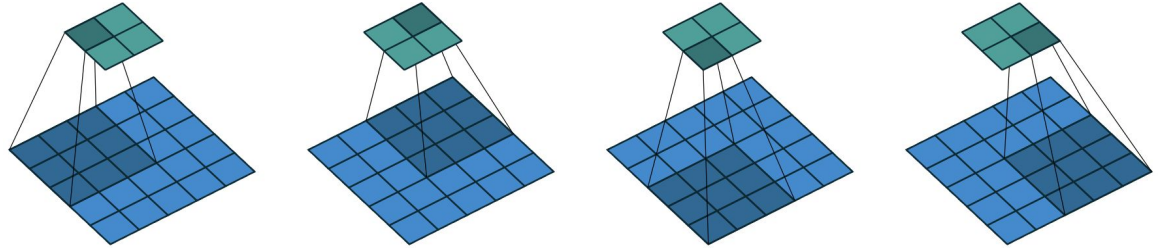
Example

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$$

- No padding, arbitrary strides

$i = 5$, $k = 3$, $s = 2$ and $p = 0$

$$o = \left\lfloor \frac{5 + 2 \times 0 - 3}{2} \right\rfloor + 1 = 2$$



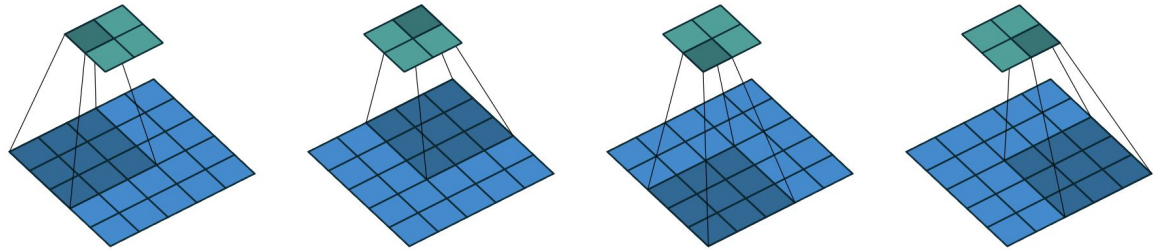
Example

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$$

- No padding, arbitrary strides

$i = 5, k = 3, s = 2$ and $p = 0$

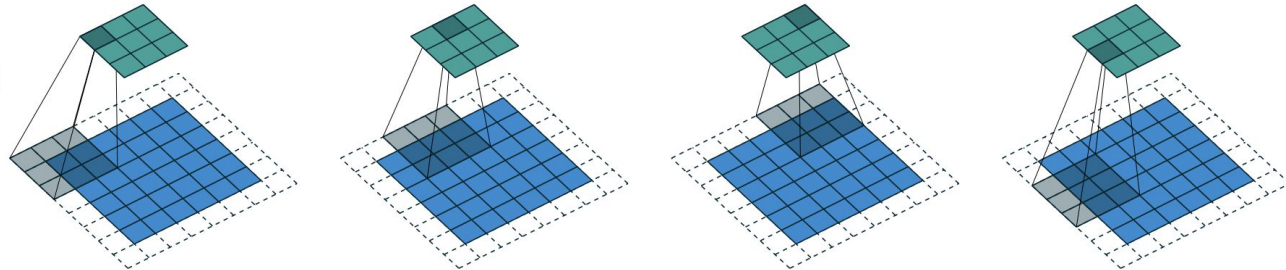
$$o = \left\lfloor \frac{5 + 2 \times 0 - 3}{2} \right\rfloor + 1 = 2$$



- Arbitrary padding and strides

$i = 6, k = 3, s = 2$ and $p = 1$

$$o = \left\lfloor \frac{6 + 2 \times 1 - 3}{2} \right\rfloor + 1 = 3$$



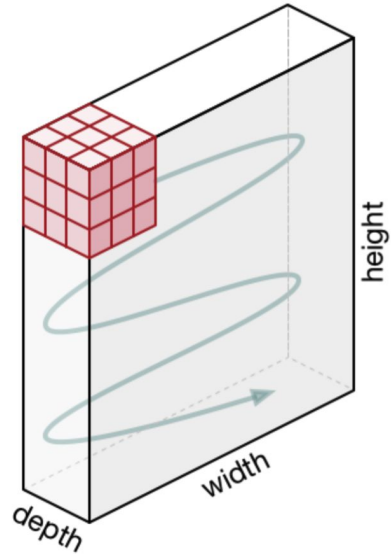
CNN on RGB



What if we have 3
colour channels?

RGB

Convolution on RGB input



Convolution in RGB

kernel dim = $3 \times k \times k$
image dim = $3 \times i \times i$

The kernel becomes a 3-dimensional tensor!

Convolution on RGB input

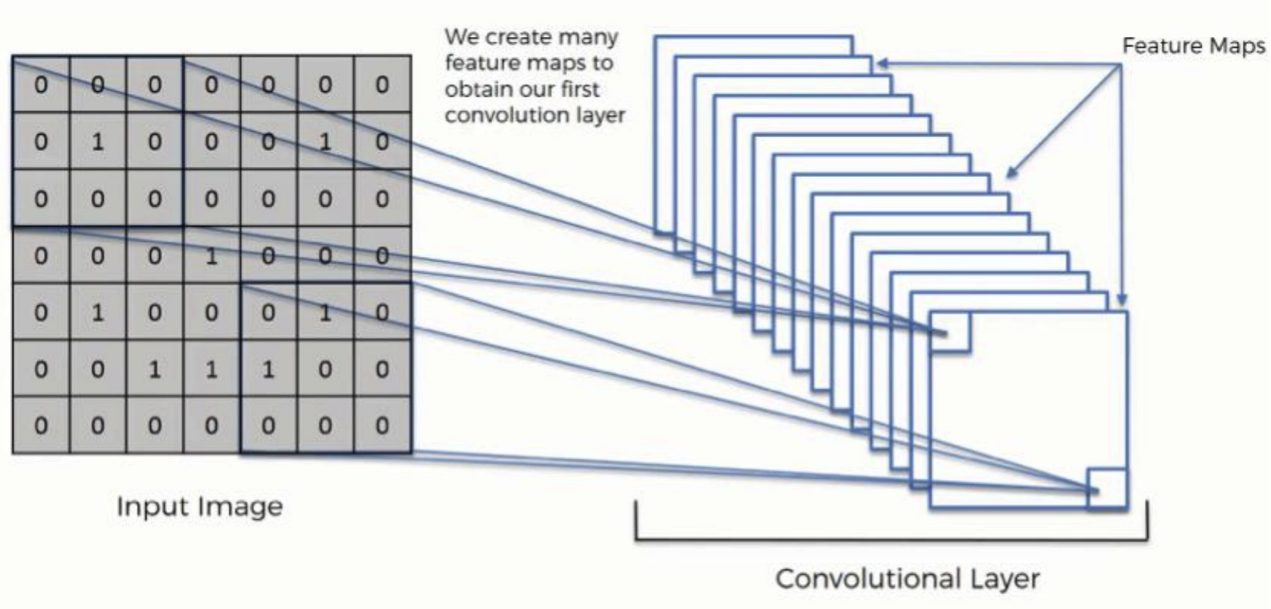
Colour input image: $3 \times 28 \times 28$

Convolution kernel: $3 \times 3 \times 3$

Questions

- How many trainable weights are there?
- What if we want to detect many **features** ?

Expanding Feature Maps



stride = 2
 Padding = 1

R

G

B

Input Volume (+pad 1) (7x7x3)

x[:, :, 0]

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

x[:, :, 1]

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

x[:, :, 2]

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

Filter W0 (3x3x3)

w0[:, :, 0]

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

w0[:, :, 1]

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

w0[:, :, 2]

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

Bias b0 (1x1x1)

b0[:, :, 0]

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

Filter W1 (3x3x3)

w1[:, :, 0]

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

w1[:, :, 1]

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

w1[:, :, 2]

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

Bias b1 (1x1x1)

b1[:, :, 0]

| |
|---|
| 0 |
|---|

Output Volume (3x3x2)

o[:, :, 0]

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

o[:, :, 1]

| | | |
|----|----|----|
| -2 | -3 | -5 |
| -2 | -1 | -2 |
| -3 | -4 | 2 |

Result of adding 27 numbers + bias

Bias term is often ignored in the example convolutions, but we can assume it's always there and adds a constant to each result

Convolution of RGB (3-Channel) image with Filter "W0"

R

Input Volume (+pad 1) (7x7x3)

$x[:, :, 0]$

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

G

$x[:, :, 1]$

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

B

$x[:, :, 2]$

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

Filter W0 (3x3x3)

$w0[:, :, 0]$

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

$w0[:, :, 1]$

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

$w0[:, :, 2]$

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

Bias b0 (1x1x1)

$b0[:, :, 0]$

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

Filter W1 (3x3x3)

$w1[:, :, 0]$

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

$w1[:, :, 1]$

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

$w1[:, :, 2]$

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

Bias b1 (1x1x1)

$b1[:, :, 0]$

| |
|---|
| 0 |
|---|

Output Volume (3x3x2)

$o[:, :, 0]$

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

$o[:, :, 1]$

| | | |
|----|----|----|
| -2 | -3 | -5 |
| -2 | -1 | -2 |
| -3 | -4 | 2 |

Convolution Interactive Demo:
<https://deeplizard.com/resource/pavq7noze2>

Pink values are weights learned by convolutional layer, everything else is input/output

Convolution of RGB (3-Channel) image with Filter "W1"

R

Input Volume (+pad 1) (7x7x3)

$x[:, :, 0]$

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

Filter W0 (3x3x3)

$w0[:, :, 0]$

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

Filter W1 (3x3x3)

$w1[:, :, 0]$

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

Output Volume (3x3x2)

$o[:, :, 0]$

| | | |
|---|----|----|
| 2 | 3 | 3 |
| 3 | 7 | 3 |
| 8 | 10 | -3 |

G

$x[:, :, 1]$

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

$w0[:, :, 1]$

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

$w1[:, :, 1]$

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

$o[:, :, 1]$

| | | |
|----|----|----|
| -8 | -8 | -3 |
| -3 | 1 | 0 |
| -3 | -8 | -5 |

B

$x[:, :, 2]$

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

$w0[:, :, 2]$

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

$w1[:, :, 2]$

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

Bias $b0$ (1x1x1)

$b0[:, :, 0]$

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

Bias $b1$ (1x1x1)

$b1[:, :, 0]$

| |
|---|
| 0 |
|---|

Convolution on RGB input

Colour input image: $3 \times 28 \times 28$

Convolution kernels: $5 \times 3 \times 8 \times 8$

Questions

- How many input channels are there?
- How many output channels are there?
- How many trainable weights are there?

Pooling Operator

Consolidating Information

In a neural network with fully-connected layers, we reduced the number of units before the final output layer.

Q: Why?

Consolidating Information

In a neural network with fully-connected layers, we reduced the number of units before the final output layer.

Q: Why?

To be able to **consolidate information** , and remove information not useful for the current task

Consolidating Information

In a neural network with fully-connected layers, we reduced the number of units before the final output layer.

Q: Why?

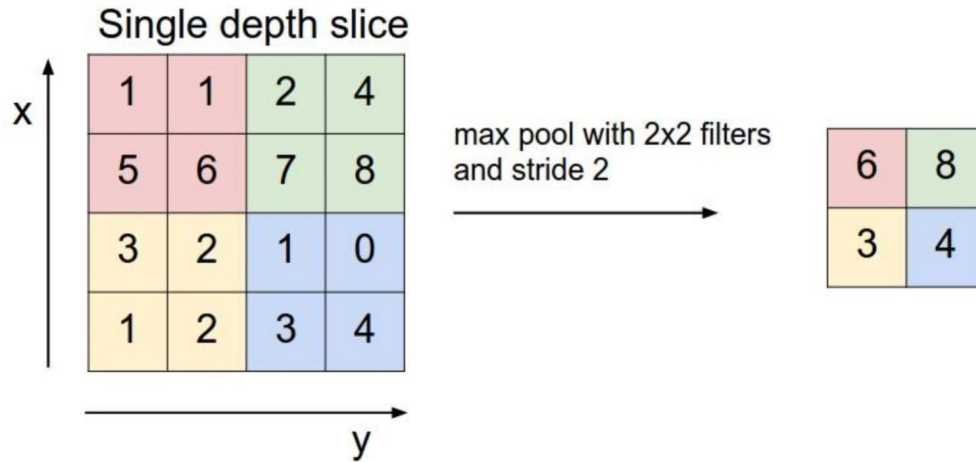
To be able to **consolidate information** , and remove information not useful for the current task

Q: How can we consolidate information in a neural network with convolutional layers?

Strided convolutions , **max pooling** , **average pooling**

Max pooling

pooling layers provide invariance to small translations of the input.

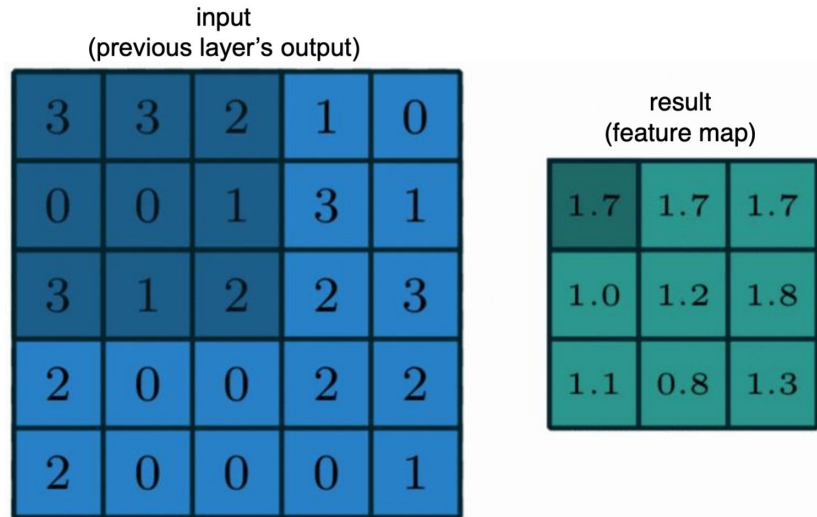


The size of output dimension is computed by:

$$o = \left\lfloor \frac{i - k}{s} \right\rfloor + 1$$

- Image dimension of size i
- Kernel of size k
- Stride of size s

Average pooling

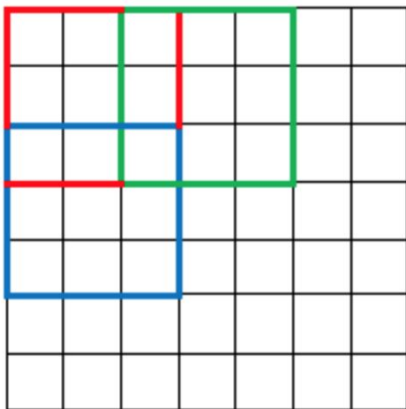


Max pooling generally works better

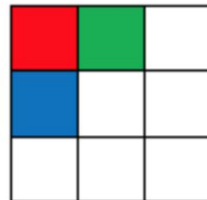
Strided convolution

- More recently people are doing away with pooling operations, **using strided convolutions instead** .
- Shift the kernel by s (e.g. $s=2$) when computing convolution.

7 x 7 Input Volume

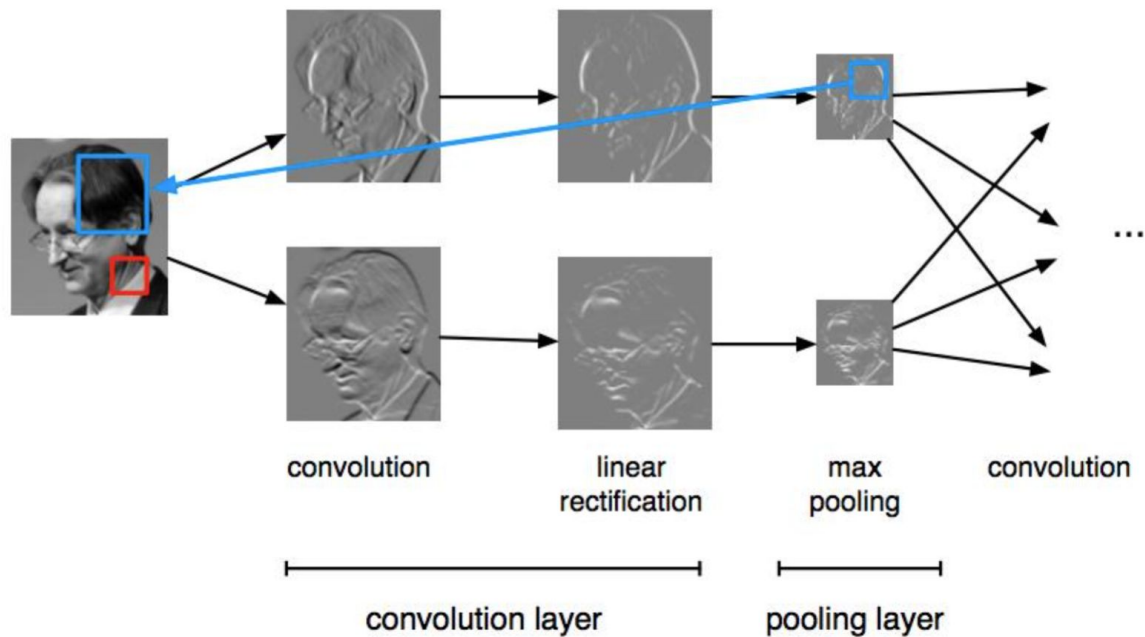


3 x 3 Output Volume



PyTorch Implementation

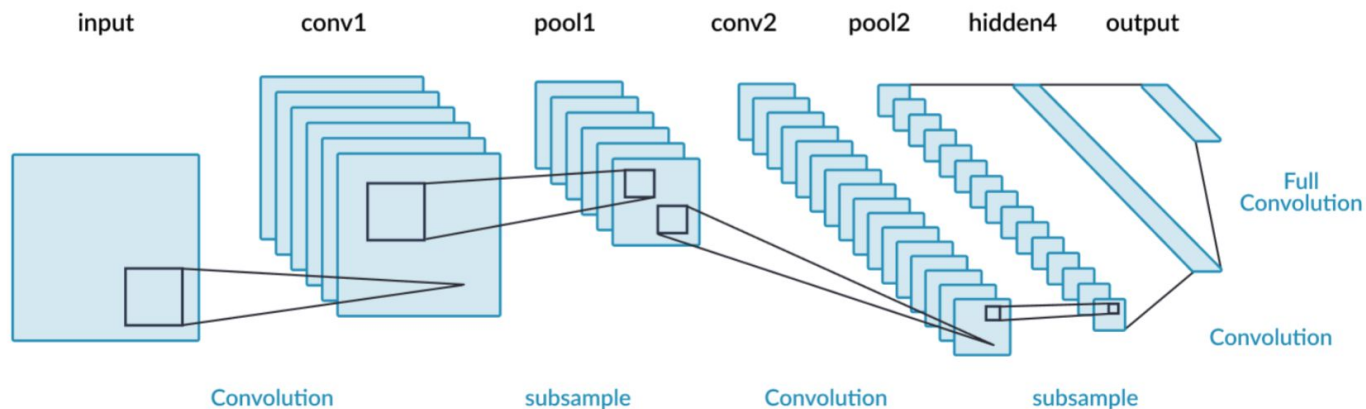
Putting it all together



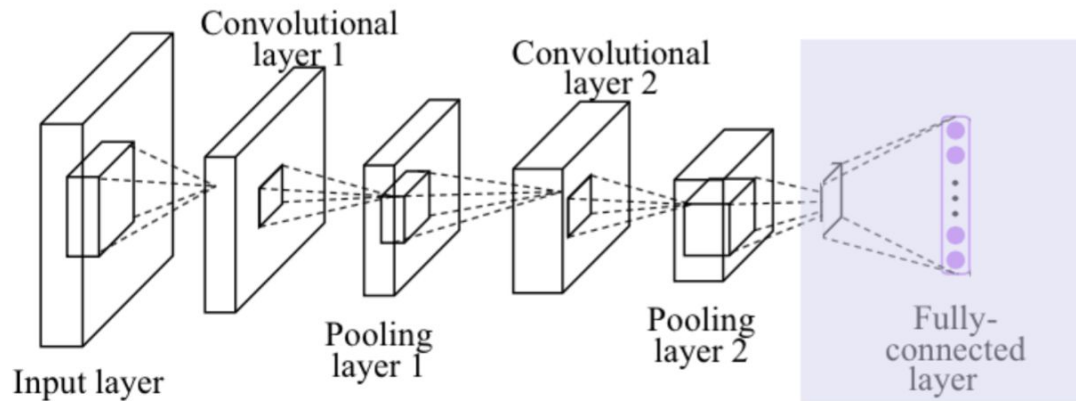
CNN Architecture Blueprint

As we go through the CNN network layer by layer:

- The **filter depth increases**
- The feature map **height and width decreases**

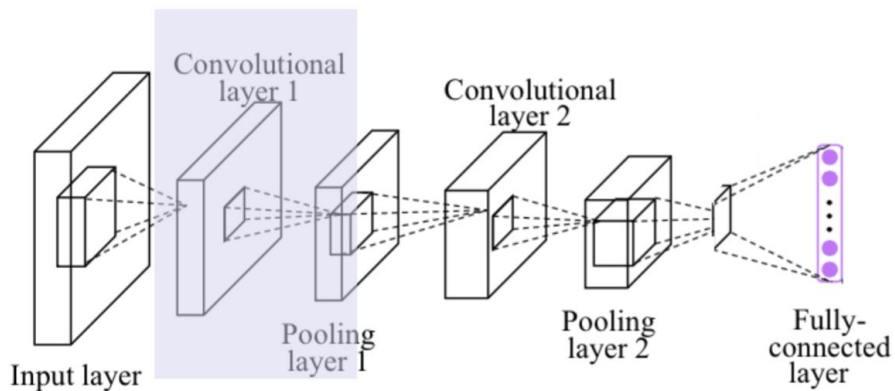


CNN in PyTorch: Recall Linear



```
self.fc = nn.Linear(n, 10)
```

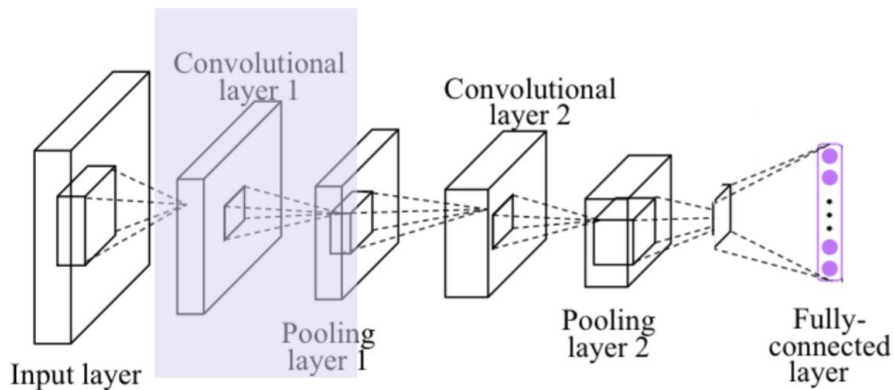
CNN in PyTorch: Conv2D



```
self.conv1 = nn.Conv2d(in_channels = 3,  
                        out_channels = 7,  
                        kernel_size= 5,  
                        stride = 1,  
                        padding = 1)
```

CNN in PyTorch: Conv2D

```
# non-square kernels and unequal stride and with padding  
m = nn.Conv2d(16, 33, (3, 5), stride=(2, 1), padding=(4, 2))
```



```
self.conv1 = nn.Conv2d(  
    in_channels = 3, }  
    out_channels = 7, }  
    kernel_size = 5, }  
    stride = 1, }  
    padding = 1) }  
}
```

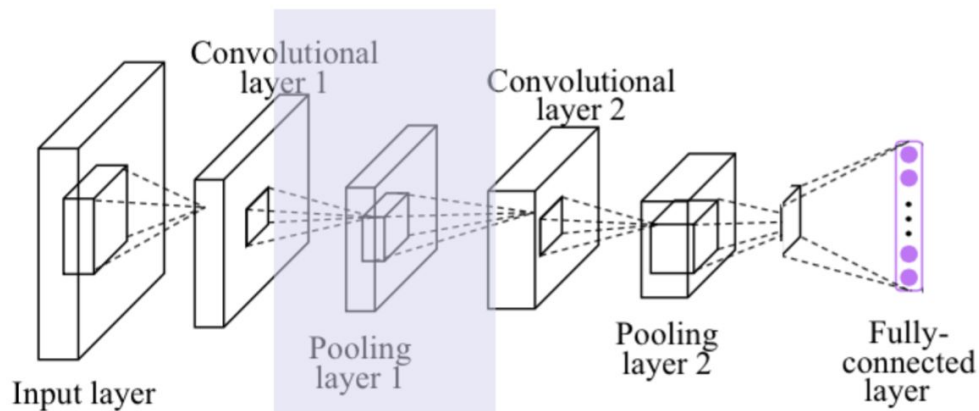
int or tuple ←

int

Default → 1

Default → 0

CNN in PyTorch: MaxPool2d



```
self.pool = nn.MaxPool2d(kernel_size = 2,  
                           stride = 2)
```

CNN in PyTorch

```
class LargeNet(nn.Module):
    def __init__(self):
        super(LargeNet, self).__init__()
        self.name = "large"
        self.conv1 = nn.Conv2d(3, 5, 5)
        self.pool = nn.MaxPool2d(2, 2)
        self.conv2 = nn.Conv2d(5, 10, 5)
        self.fc1 = nn.Linear(10 * 5 * 5, 32)
        self.fc2 = nn.Linear(32, 10)

    def forward(self, x):
        x = self.pool(F.relu(self.conv1(x)))
        x = self.pool(F.relu(self.conv2(x)))
        x = x.view(-1, 10 * 5 * 5)
        x = F.relu(self.fc1(x))
        x = self.fc2(x)
        return x
```

Questions?