# Introduction to Computer Vision CNNs

A special
mathematical operation

# Convolutional
# Neural Network

It's another class of neural networks used to
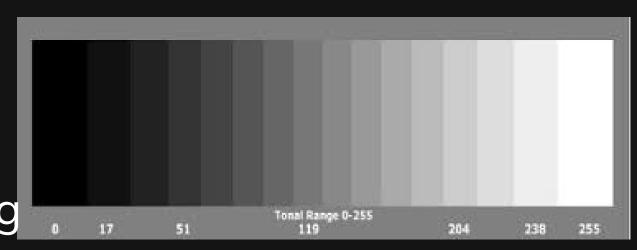get important information in images

# CONTENT

- **What are Grayscale/RGB images?**
- **The Convolution Operation**
- **Layers of a Convolutional Neural Network**
  - **Convolution**
  - **Activation Functions**
  - **Pooling**
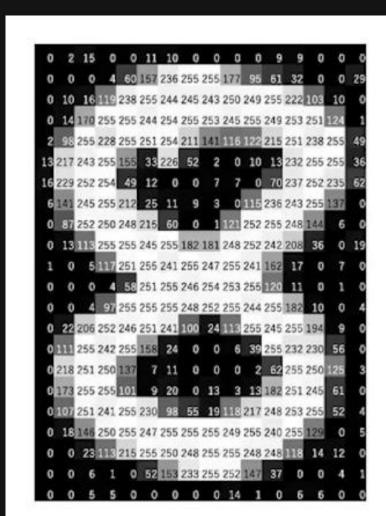  - **Fully Connected Network**
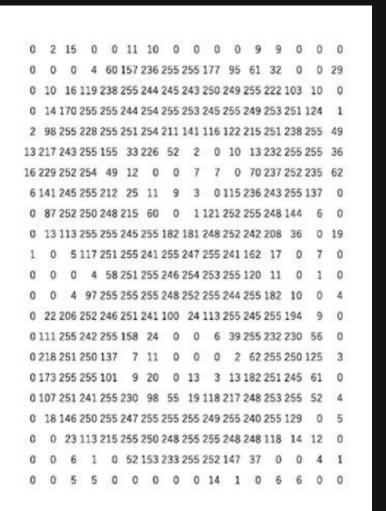
# HOW DO COMPUTERS PROCESS IMAGES?

## We convert it into a grayscale/RGB image

Each pixel of the image has a different intensity
of the colour white. Namely, this intensity ranges from
0 to 255. Thus we convert each pixel into its corresponding
numerical value.
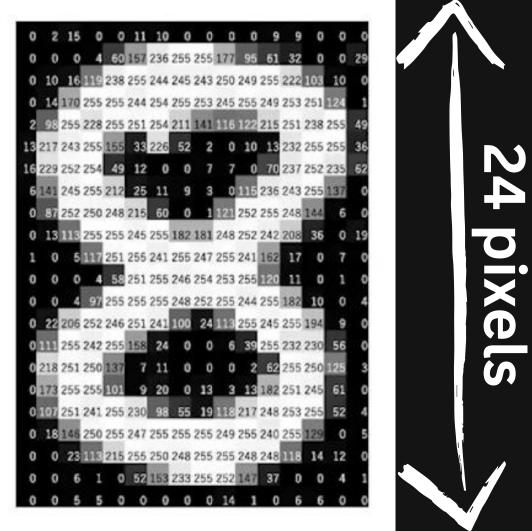
# RED GREEN BLUE IMAGES(RGB)

In most cases we deal with colour images which can be represented simply with their red, blue and green parts. Just like in Grayscale images, each of the pixels in the Red, Green and Blue images are converted into their corresponding numerical values depending on the intensity of the colour in the image



Colour Image = Red + Green + Blue

# WHY?

Neural networks will take the numerical value from each pixel as an input. That's 24x16=384 nodes in the input layer of our neural network!
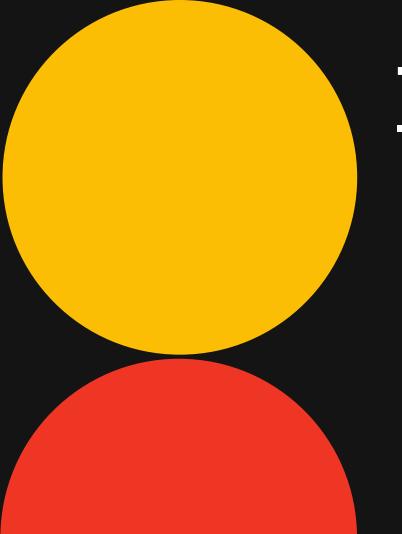


24 pixels

16 pixels

Each of the input node will be connected to each node of the hidden layer. If we assume a hidden layer of 36 nodes, we need 384*36 = 13,824 weights and 13,824 biases!. This takes a huge amount of time to compute. Thus we need to decrease the amount of inputs without losing out the important features of the image.

Now for a different image of 8, the previous network will not be as accurate because instead of capturing the features of the image, its been trained to assign weights and biases for each pixel.

Hence traditional neural networks are ineffective for image detection.

The solution? ⟶ CNNs

# THE CONVOLUTION OPERATION

## Kernel

A kernel is a small matrix which extracts the required features from the given image. The kernel is much smaller than the input image and have different kernels for different tasks like blurring, sharpening or edge detection.

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

- Typically kernels of size 3x3
- The values of the kernels will be updated by back propagation

a small part from the image matrix

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

**kernel**

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

\*

=

$3(0)+3(1)+2(2)+$
$0(2)+0(2)+1(0)$   **= 12**
$3(0)+1(1)+2(2)$

a small part from the image matrix

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

**kernel**

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

*

=

$3(0)+2(1)+1(2)+$

$0(2)+1(2)+3(0)$   = **12**

$1(0)+2(1)+2(2)$

| 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 |

# THE CONVOLUTION OPERATOR

$$\begin{bmatrix} 4 & 9 & 0 \\ 8 & 6 & 1 \\ 7 & 2 & 34 \end{bmatrix} * \begin{bmatrix} 1 & 0 & 1 \\ 0 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} = \begin{array}{l} 4(1) + 9(0) + 0(1) \\ + 8(0) + 6(2) + 1(0) \\ + 7(1) + 2(0) + 34(1) \end{array}$$
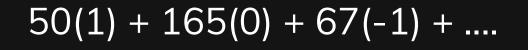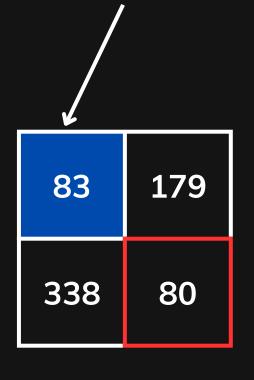
$$= \quad 57$$

# CONVOLUTION OF AN IMAGE

$50(1) + 165(0) + 67(-1) + ....$

| | | | |
|---|---|---|---|
| 50 | 165 | 67 | 0 |
| 94 | 23 | 88 | 12 |
| 178 | 56 | 90 | 64 |
| 234 | 204 | 78 | 123 |

*

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

=

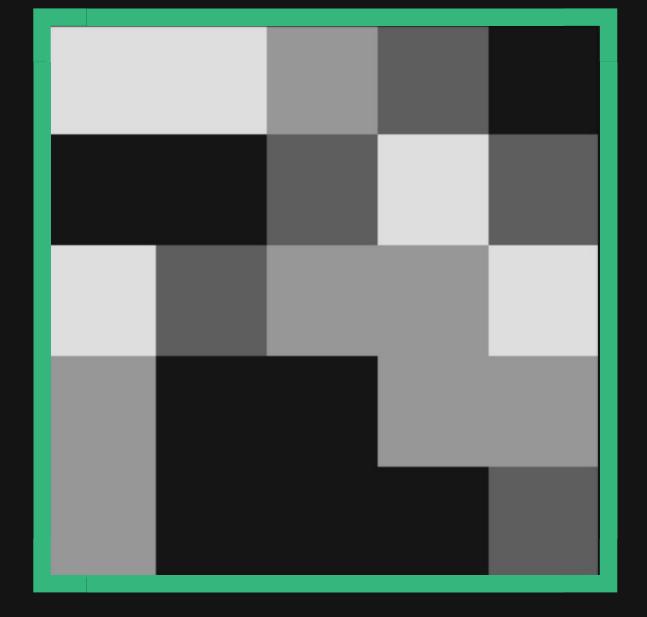| | |
|---|---|
| 83 | 179 |
| 338 | 80 |

Image

Kernel

The new image

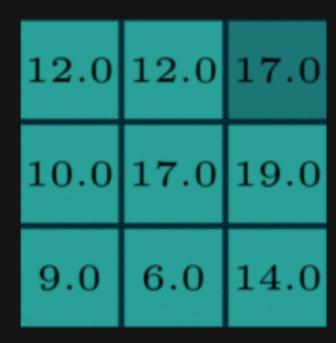# WHY IS CONVOLUTION NEEDED FOR IMAGE DETECTION?

By Convolving the Image matrix with the given kernel we extract the information from each pixel along with the influence its neighboring pixels have. Doing this allows us to extract the features of the image.



by applying the convolution we are focusing on the central pixel and extract its info along with it neighbours

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

# FILTERS VS KERNELS

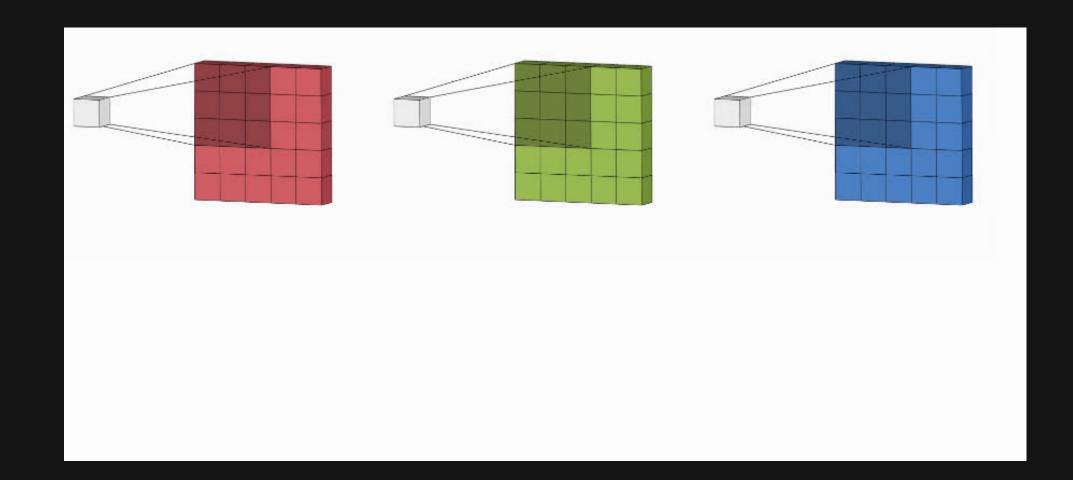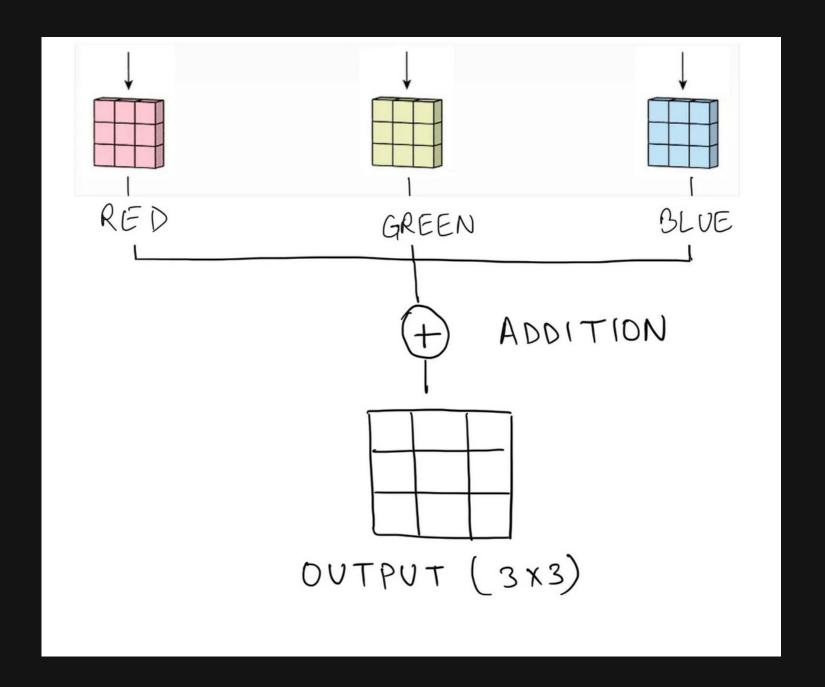Often in computer visiion , the terms kernel and filter are used interchangeably but there is a slight difference between them. Filters are groups of kernels and a reused in convoling RGB images, which have three channels (red, green, blue) instead of just one. IN each filter the kernels might be the same or sifferennt according to the features we need to extract.

The resultant three matrices are then combined using matrix addition to get a single output.

Remember that this is only for RGB images. in a gray scale a image filer is equal to the kernel because we only have one channel of input

# Few more terminologies in CNNs

Stride is the no. of steps we move the kernel each time we convolve.

Normally we use a stride of 1, so every possible section of the image is taken.



Stride = 1



Stride = 3

# STRIDE

a different stride value might be
taken according to the need of the CNN model

# PADDING

During the convolution operation we extract the information of each pixel by applying a kernel over it such the required pixel is at the middle of the kernel.
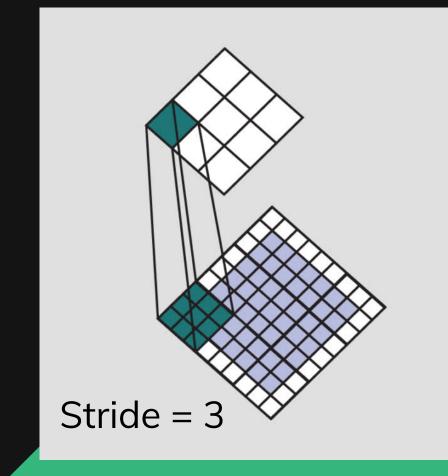
However, we end up losing the data on the edges of the image because we are unable to apply the kernel over those pixels.

| $3_0$ | $3_1$ | $2_2$ | 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 |

we retain the information of these pixels but lose out of the edge pixels

## THE SOLUTION?

### Padding the image

Padding is adding a layer of zeros around the image. It will help in convolving the entire image without distorting any information as the zeros will have no effect on the convolution

**kernel**

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

*

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

=

| 6 | 14 | 17 | 11 | 3 |
|---|----|----|----|---|
| 14 | 12 | 12 | 17 | 11 |
| 8 | 10 | 17 | 19 | 13 |
| 11 | 9 | 6 | 14 | 12 |
| 6 | 4 | 4 | 6 | 4 |

**thus by padding the image we have retained its size and the information of the edges**

# PADDING

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | 50 | 165 | 67 | 0 | 0 |
| 0 | 94 | 23 | 88 | 12 | 0 |
| 0 | 178 | 56 | 90 | 64 | 0 |
| 0 | 234 | 204 | 78 | 123 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |

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

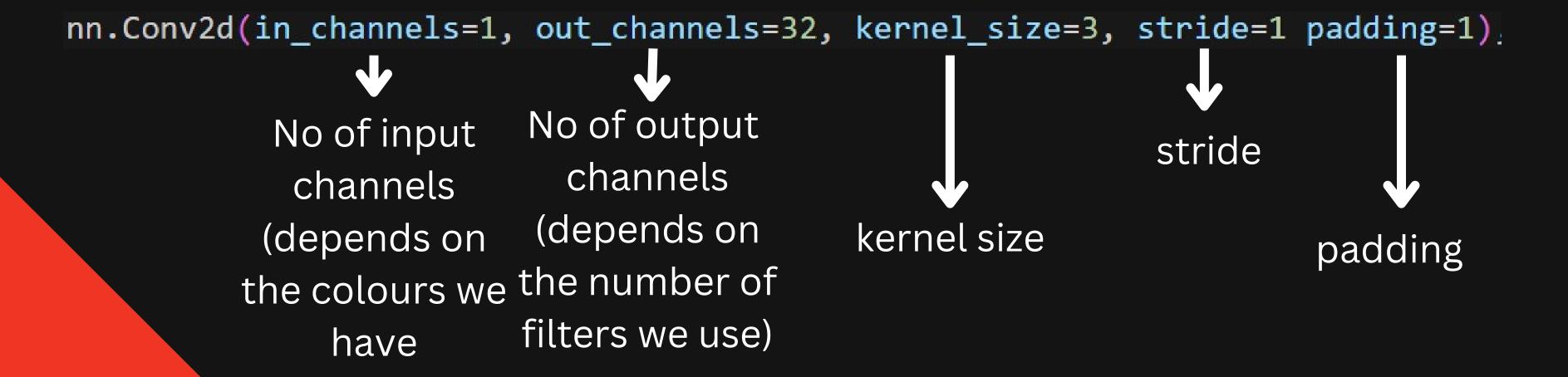| -353 | 28 | 341 | 222 |
|---|---|---|---|
| -267 | 83 | 179 | 333 |
| -343 | 338 | 80 | 346 |
| -472 | 400 | 162 | 246 |

We add zeros around the image

The kernel is the same

Now the size is the same and entire information of the image is retained

# Congratulations! You just learnt the basics of a Convolutional Layer!

**Before we move forward, let's look at syntax of convolutions in PyTorch:**

```python
nn.Conv2d(in_channels=1, out_channels=32, kernel_size=3, stride=1 padding=1)
```

No of input channels (depends on the colours we have)

No of output channels (depends on the number of filters we use)

kernel size

stride

padding

# Also here's a simple formula

It's just for reference

dont worry about memorizing this

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

where
o is the output size
i is the input size
k is kernel size
p is padding
s is stride

# Some Examples of Kernels

# GAUSSIAN BLUR

$$\begin{bmatrix} .075 & .124 & .075 \\ .124 & .204 & .124 \\ .075 & .124 & .075 \end{bmatrix}$$

Gaussian blurring is a technique that blurs a pixel giving more weight to middle pixel. Its similar to the gaussian curve. Blurring the image reduces computation time by filering out noise(unwanted data) in the image while reataing the context

Blurring the image retained important information
but also got rid of unnecessary noise
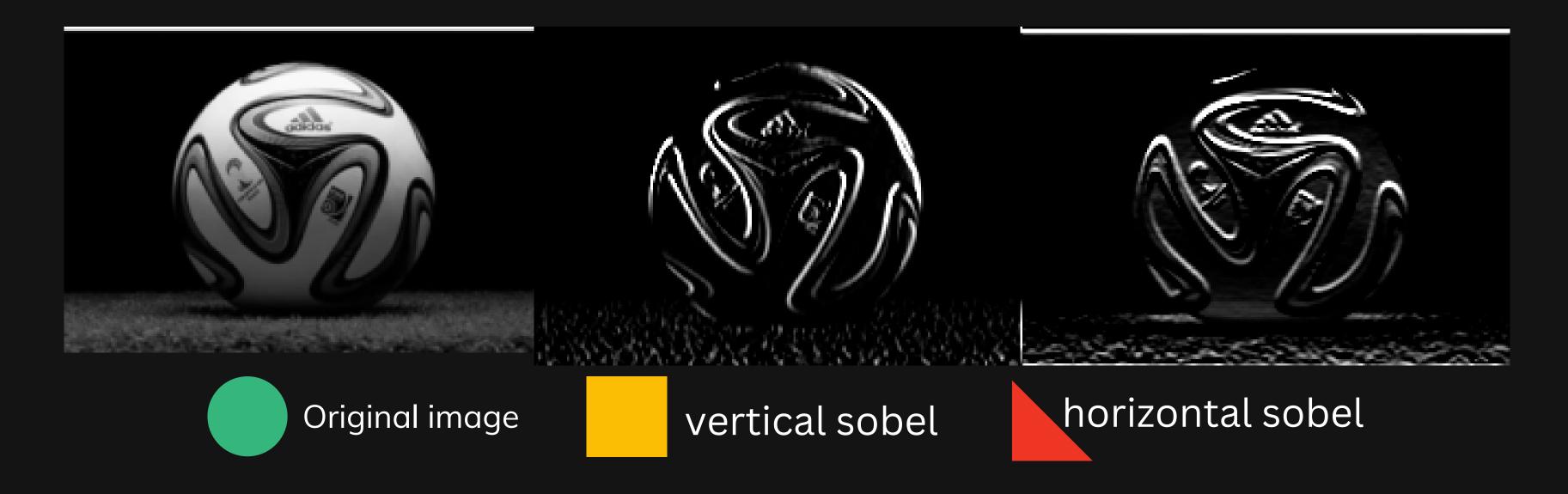
# THE SOBEL OPERATOR

### Vertical Edges

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

### Horizontal Edges

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

The kernel shown before is actually the Sobel operator - it's used to find the edges in an image.
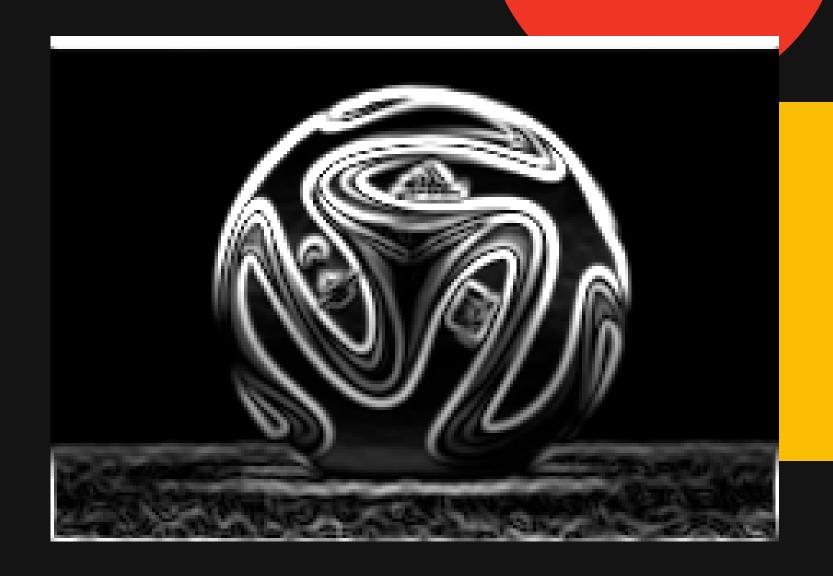
# SOBEL OPERATOR



● Original image   ■ vertical sobel   ◢ horizontal sobel

We then combine both outputs to get our final result

Gaussian blur also reduces noise, as you can see when you compare the edges detected in both images



Sobel
(Horizontal and Vertical Combined)

Gaussian Blur + Sobel

# Other Layers of the CNN

Images by default have intensities ranging from 0 to 255, but neural networks usually work best if inputs are normalized, i.e, between 0 and 1. This helps in reducing the learning time and helps in preventing overfitting of the model

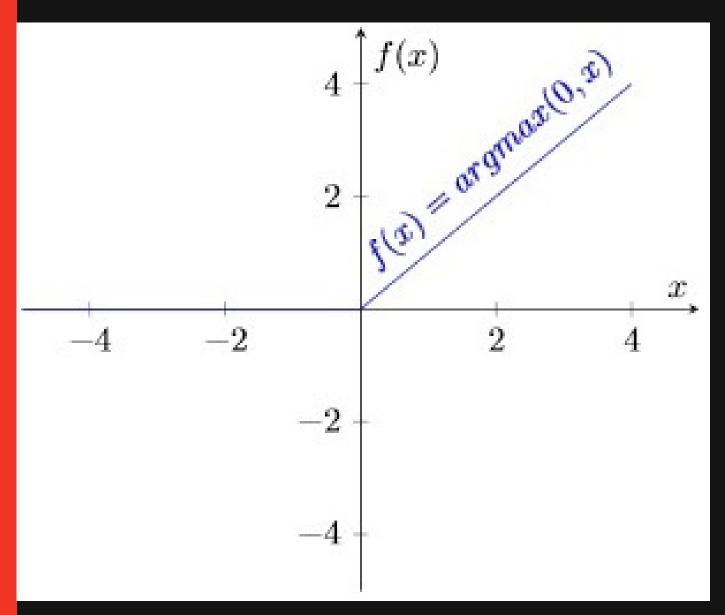So, we divide our numpy array before applying convolutions.

# NORMALIZATION

## PyTorch code

```
transforms.Normalize((0,), (1,))
```

# ACTIVATION FUNCTIONS

Just like in artificial neural networks, we can make use of activation functions like ReLU here.

After we apply convolution and batch normalization to an image, we can apply an activation function to each pixel.

# Relu



**f(x) = argmax(0,x) where**
**x is the value of each pixel in the image array**

nn.ReLU().

# POOLING

We reduce the dimensionality of the convolved image by pooling it. Normally, a stride of the size of the block is used (so there's no overlap)
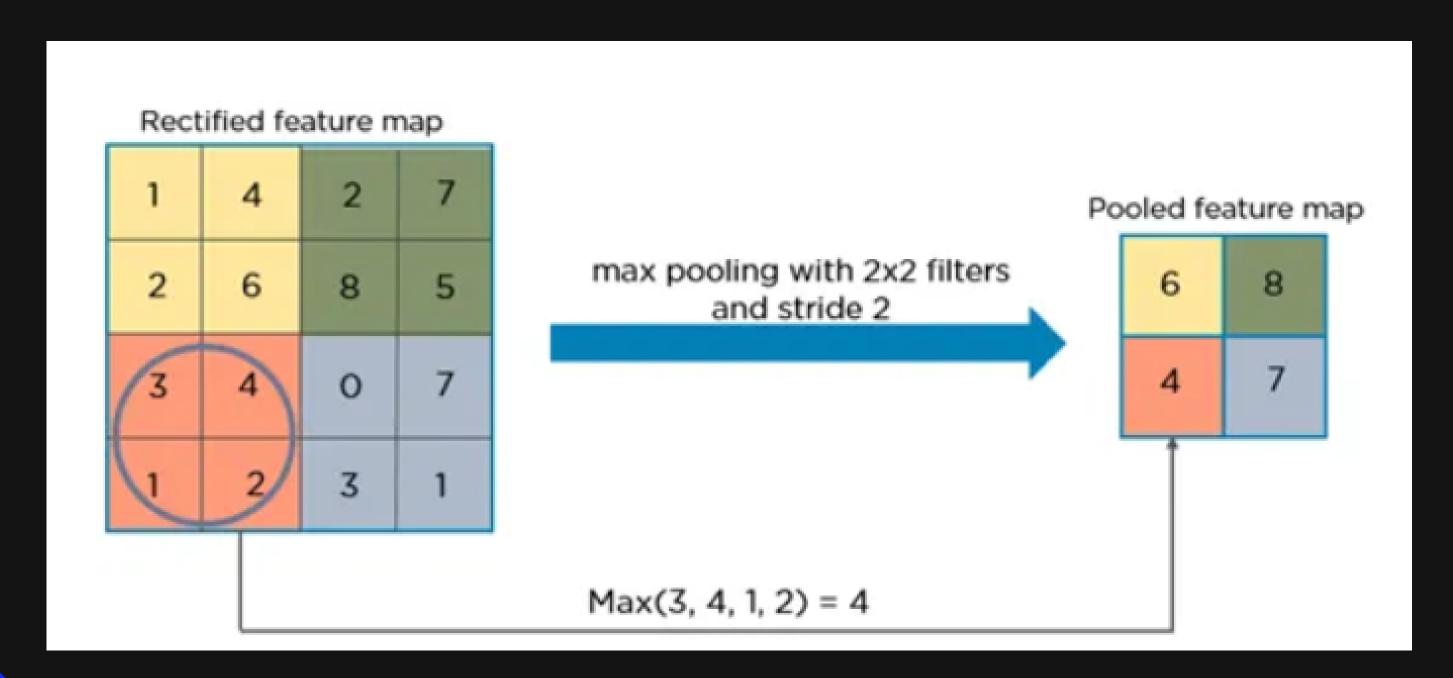
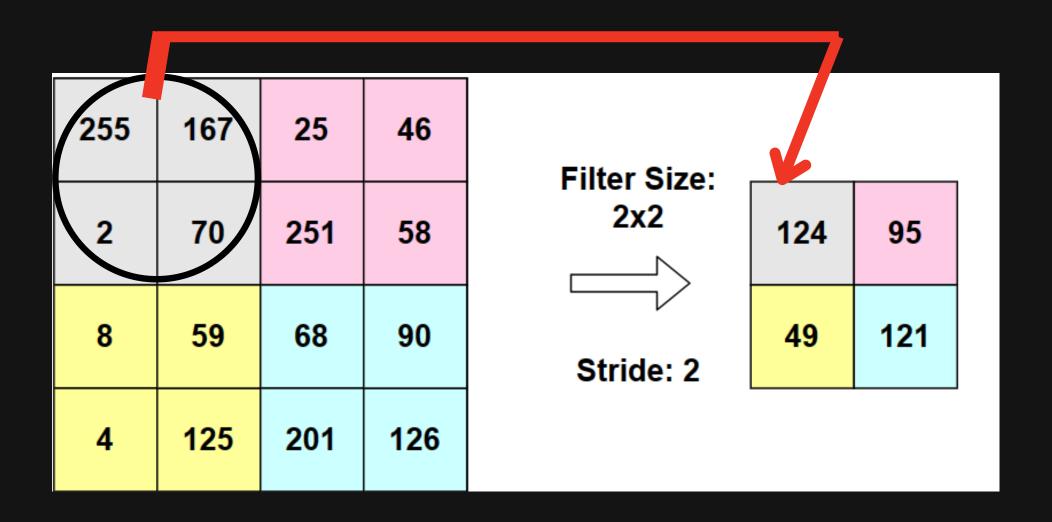This greatly reduces the size of the image, and also removes unnecessary detail (preventing overfitting)

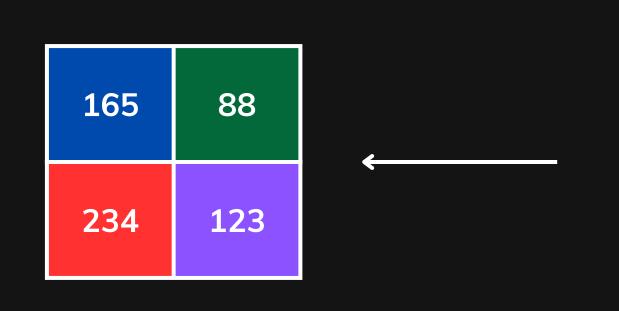## Max Pooling

## Average Pooling

# MAX POOLING



Max pooling is a pooling operation that selects the maximum element from the region of the feature map covered by the filter. Thus, the output after max-pooling layer would be a feature map containing the most prominent features of the previous feature map
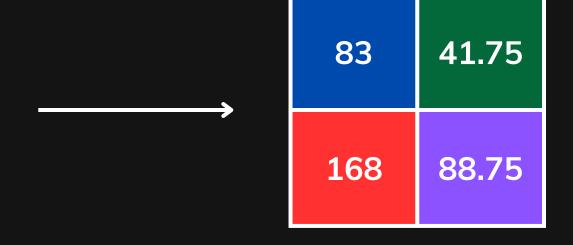
# AVERAGE POOLING

Average pooling gives us the average of the features in a filter size

# TYPES OF POOLING

| | |
|---|---|
| 50 | 165 | 67 | 0 |
| 94 | 23 | 88 | 12 |
| 178 | 56 | 90 | 64 |
| 234 | 204 | 78 | 123 |

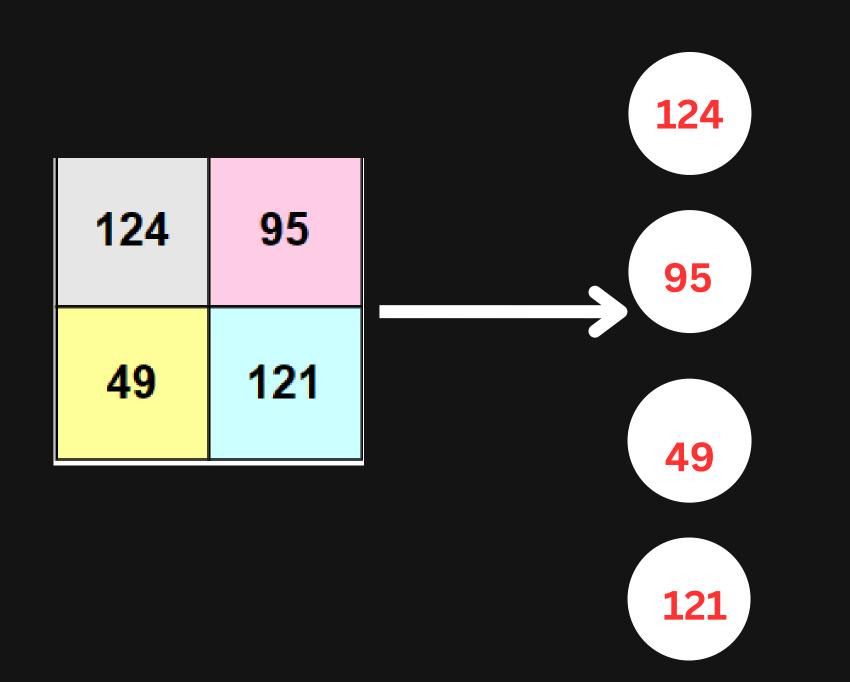| | |
|---|---|
| 165 | 88 |
| 234 | 123 |

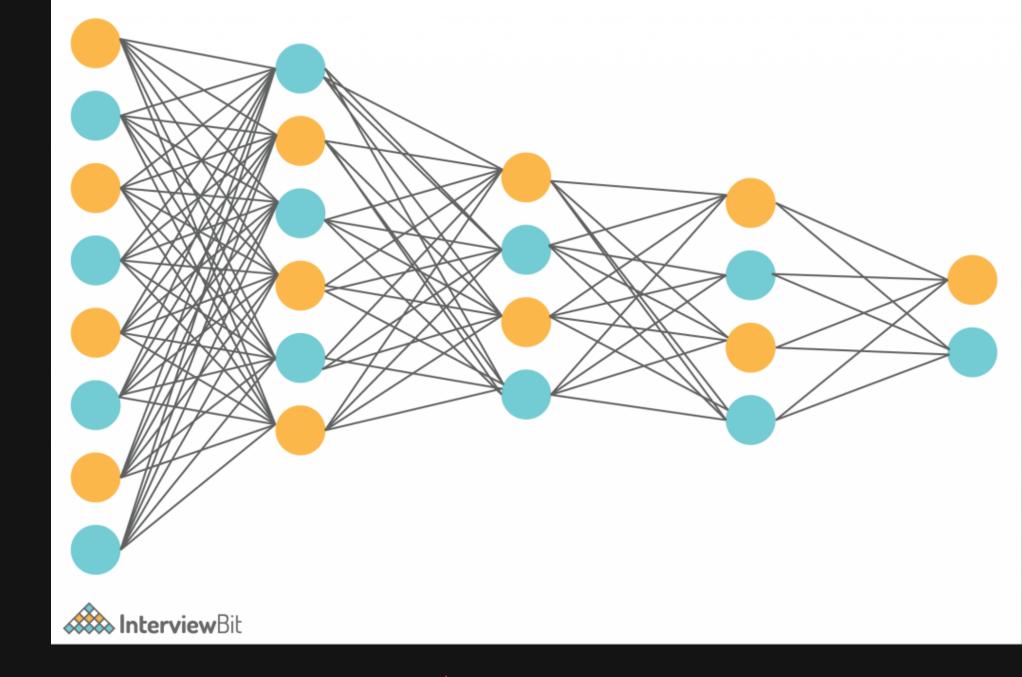| | |
|---|---|
| 83 | 41.75 |
| 168 | 88.75 |

Max Pooling

Average Pooling

# FLATTENING

The output of the final pooling layer is fed into a fully connected neural network. To do this we "flatten" the outputs

# FULLY CONNECTED LAYER

**Finally we input our results into a neural network to run the classification**



Convolutions just let the network learn features of the image. We then use these in a normal network to get our results
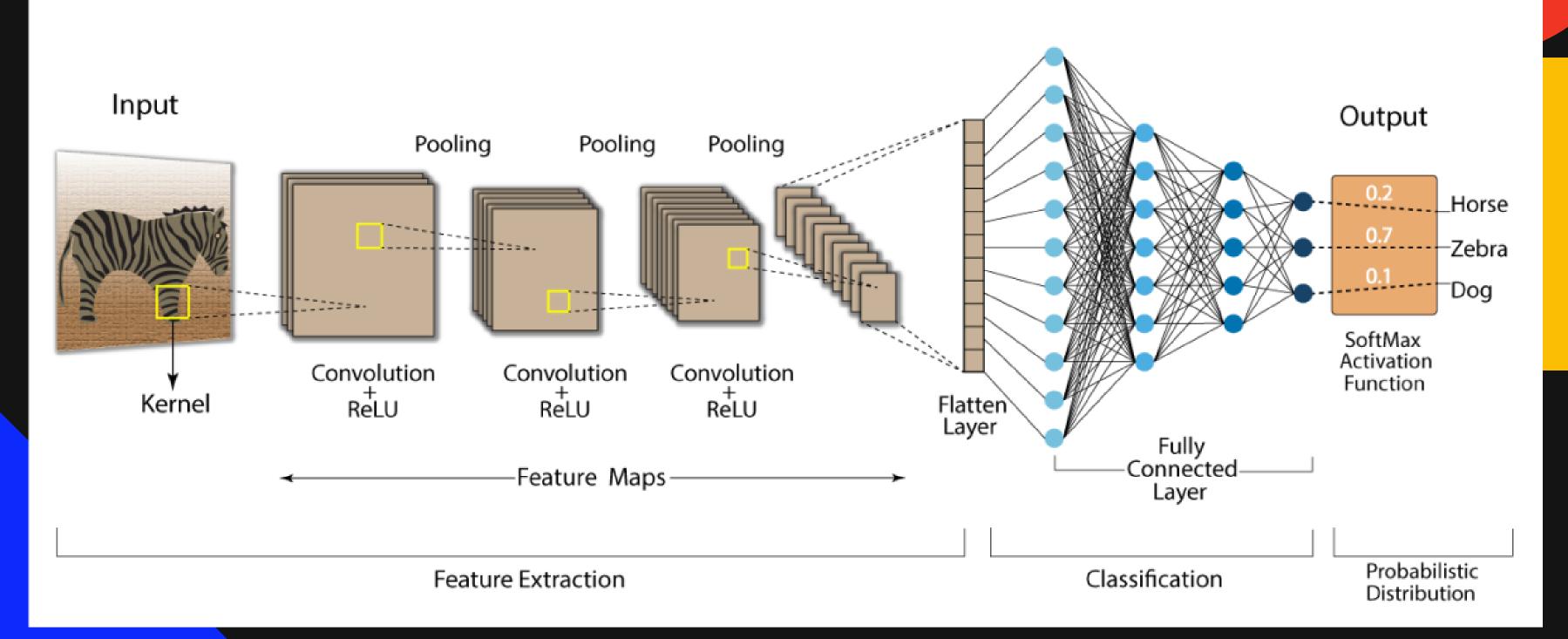
First, we take our image, which is a matrix, and flatten it. This acts as the input layer for our artificial neural network

Next, we add hidden layers and output layers as in a normal neural network.

# THE ENTIRE NETWORK



Convolution Neural Network (CNN)

Input

Pooling   Pooling   Pooling

Kernel

Convolution + ReLU

Convolution + ReLU

Convolution + ReLU

Flatten Layer

Fully Connected Layer

Output

0.2 — Horse
0.7 — Zebra
0.1 — Dog

SoftMax Activation Function

Feature Maps

Feature Extraction

Classification
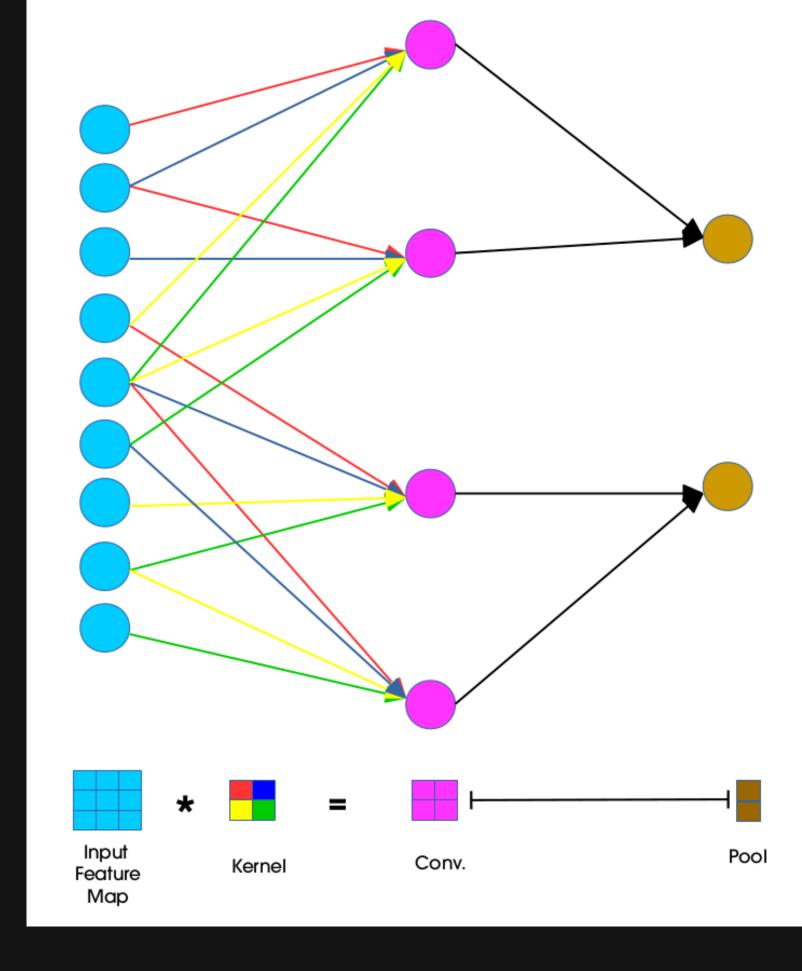
Probabilistic Distribution

# BACKPROPAGATION



Backpropagation is applied on the kernels in CNNs.

CNNs can be considered a partially connected neural network

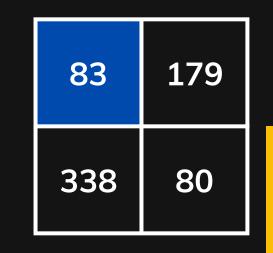~~We~~ Pytorch can do the standard backpropagation by only considering connected neurons

$W_{11}$

| 50 | 165 | 67 | 0 |
|----|-----|----|----|
| 94 | 23 | 88 | 12 |
| 178 | 56 | 90 | 64 |
| 234 | 204 | 78 | 123 |

*

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

=

| 83 | 179 |
|----|-----|
| 338 | 80 |

**Our target is to update W11. to do that lets see which pixels it affects.**

$$\frac{\partial L}{\partial W_{11}} = (y_{11} - t_{11})x_{11} + (y_{12} - t_{12})x_{12} + (y_{21} - t_{21})x_{21} + (y_{22} - t_{22})x_{22}$$

where x11, x12 so on are from the image like 50, 165
t11,t12 are the targets for each output pixel
y11 ,y12 etc are the present values, that is 83, 179