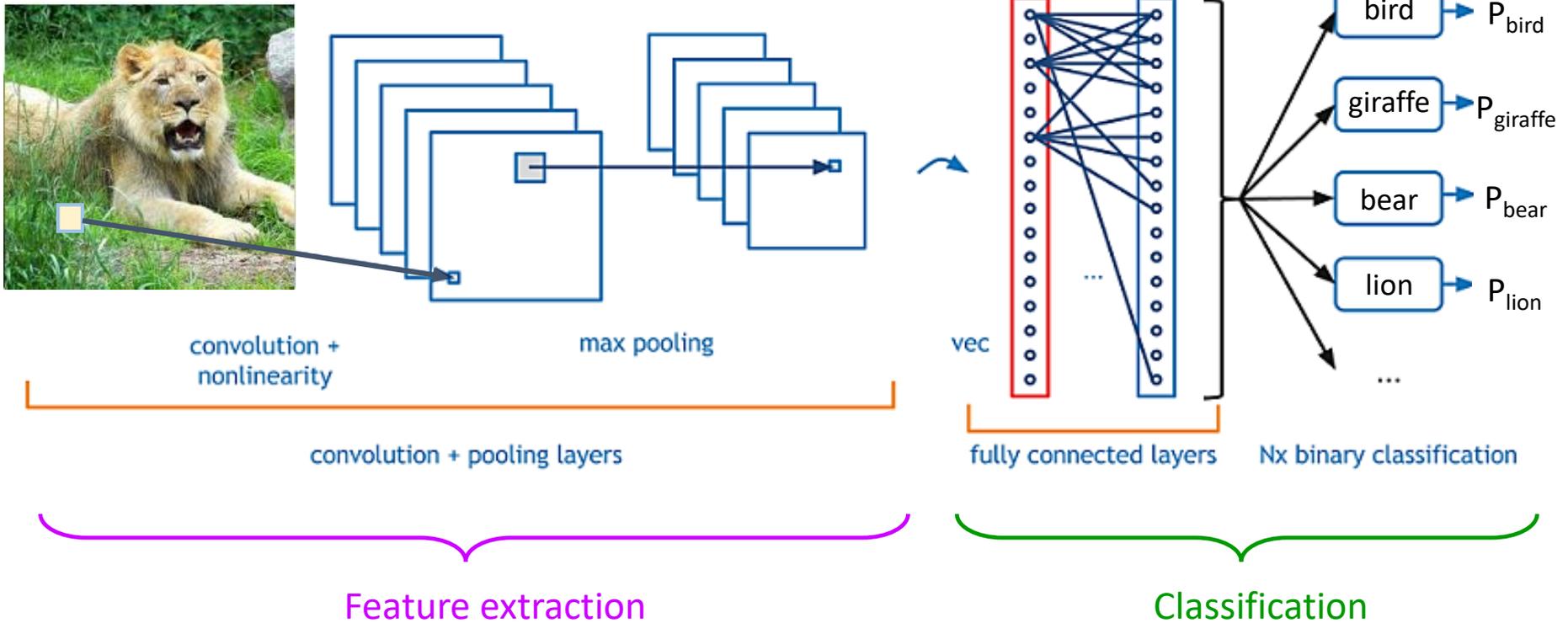


CNN

Convolutional Neural Network



CNN (ConvNet) for deep learning- intro

- CNN – Convolutional Neural Network

- ❑ A class of deep neural networks, most commonly applied to analyzing visual imagery.
- ❑ Applications in: image and video recognition, recommender systems, image classification, medical image analysis, and natural language processing.
 - ❖ Little pre-processing compared to other image classification algorithms.
 - ❖ The network **learns the filters** that in traditional algorithms were hand-engineered.
 - ✓ This independence from prior knowledge and human effort in feature design is a major advantage

CNN (ConvNet) for deep learning - intro

- ❑ The network employs a mathematical operation called **convolution**.
- ❑ CNN are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers.
- ❑ CNN have learnable parameter like neural network (weights, biases, etc).

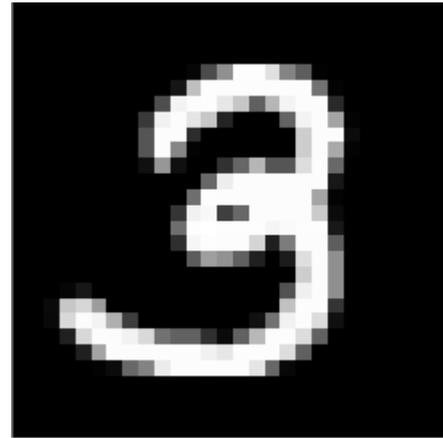
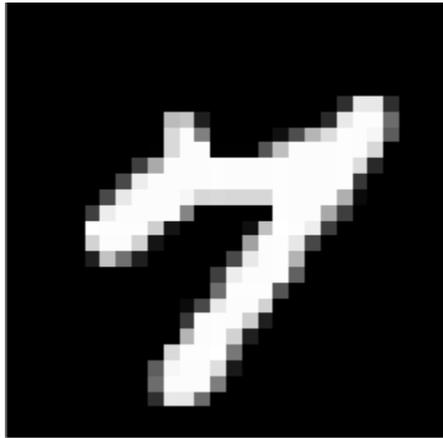
❖ Convolution

- a specialized kind of linear operation
- a mathematical operation on two functions (f and g) that produces a third function expressing how the shape of one (f) is modified by the other (g).
- defined as the integral of the product of the two functions after one (g) is reversed and shifted.

Why CNN?

If working with simple images, for example MNIST data set

- $28 \times 28 \times 1$ (b&w; 1 channel) = 784 features



The input layer in an ANN must have 784 neurons – can be manageable

If working with “real” images

- $224 \times 224 \times 3$ (3 channels) = 150,528 features



Needs 150,528 neurons!!!

Computationally inefficient!

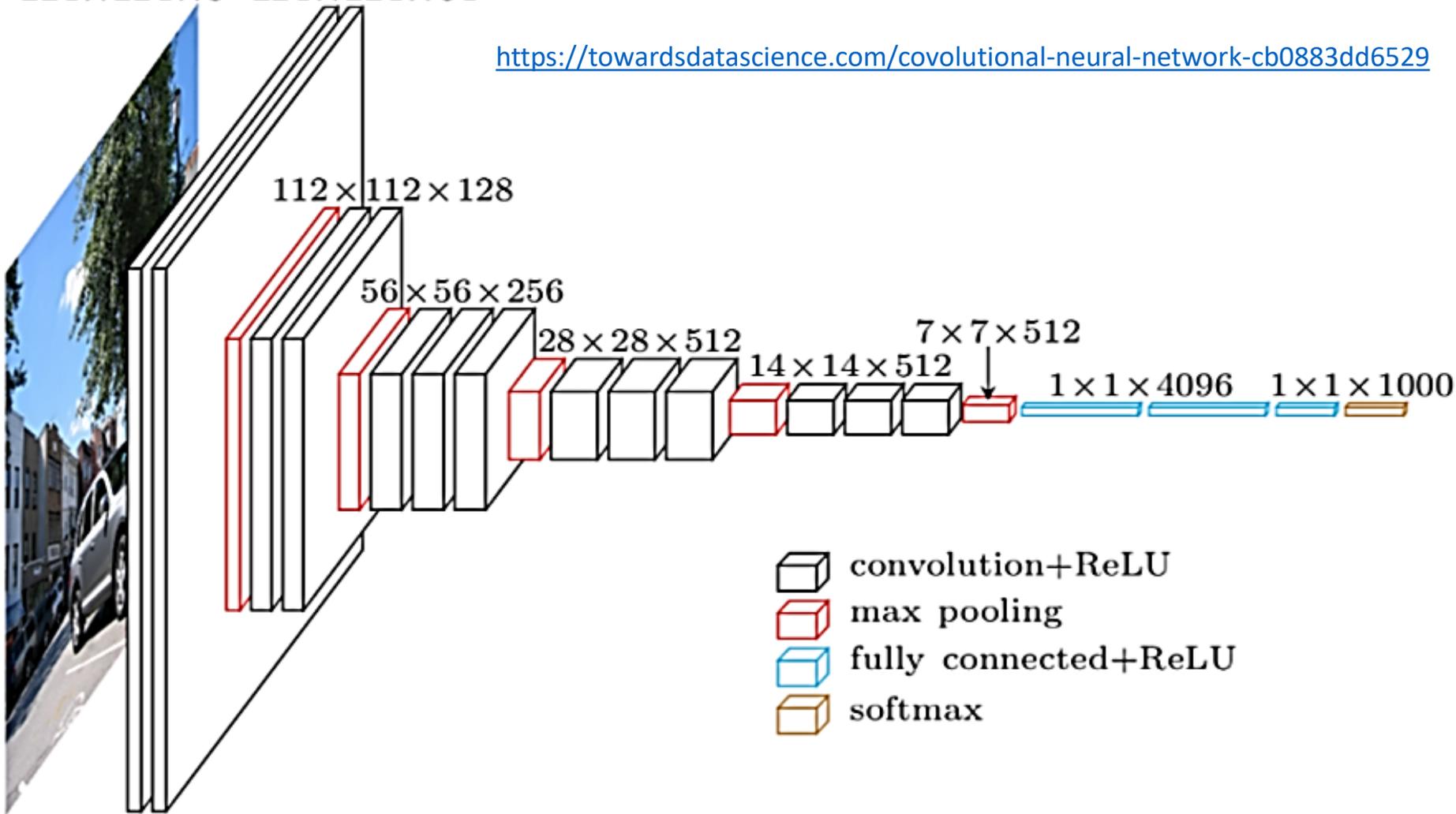
Applying convolution the input tensor dimension is finally reduced down to $1 \times 1 \times 1000$

Only 1000 neurons in the 1st layer of the feed forward neural network

Downsampling

$224 \times 224 \times 3$ $224 \times 224 \times 64$

<https://towardsdatascience.com/covolutional-neural-network-cb0883dd6529>

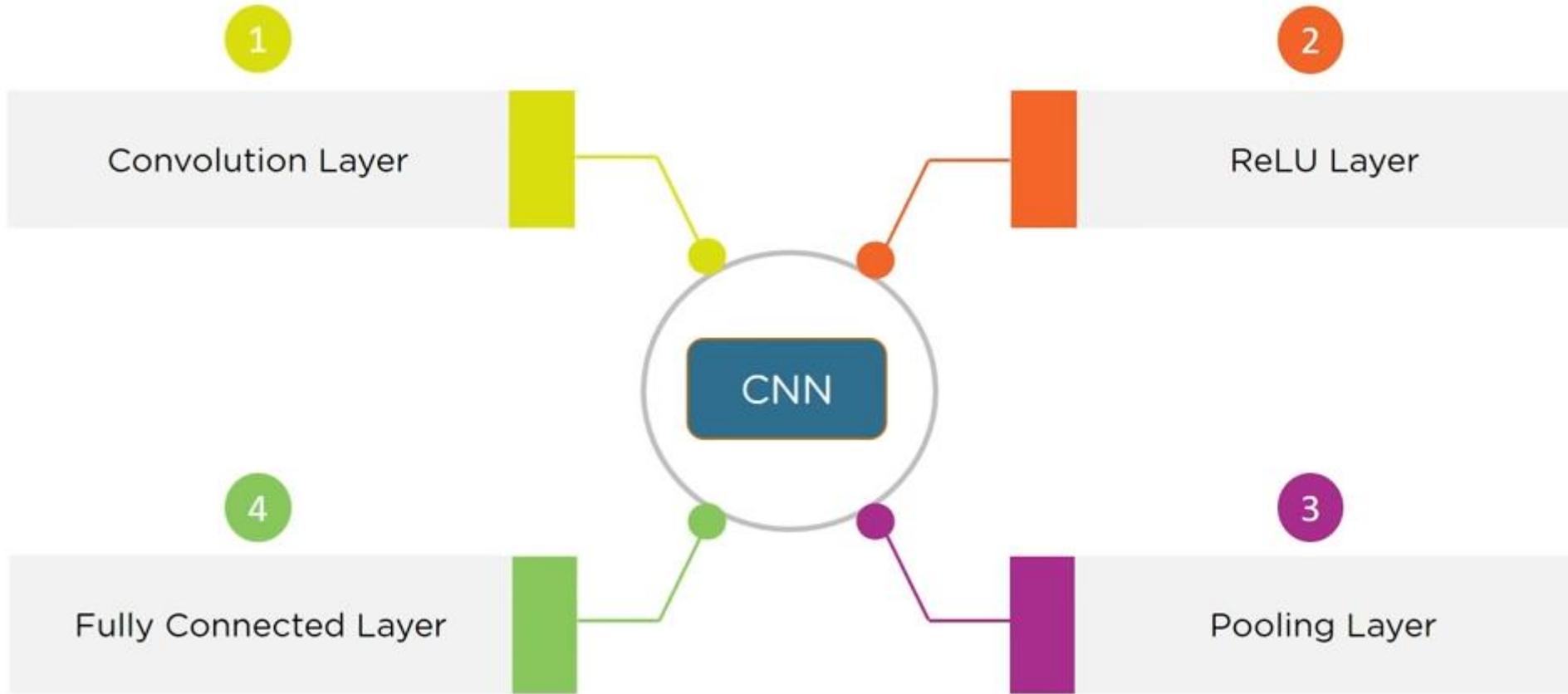


The hidden layers of a CNN typically consist of convolutional layers, pooling layers, fully connected layers, and normalization layers

□ The hidden layers of a CNN typically consist of

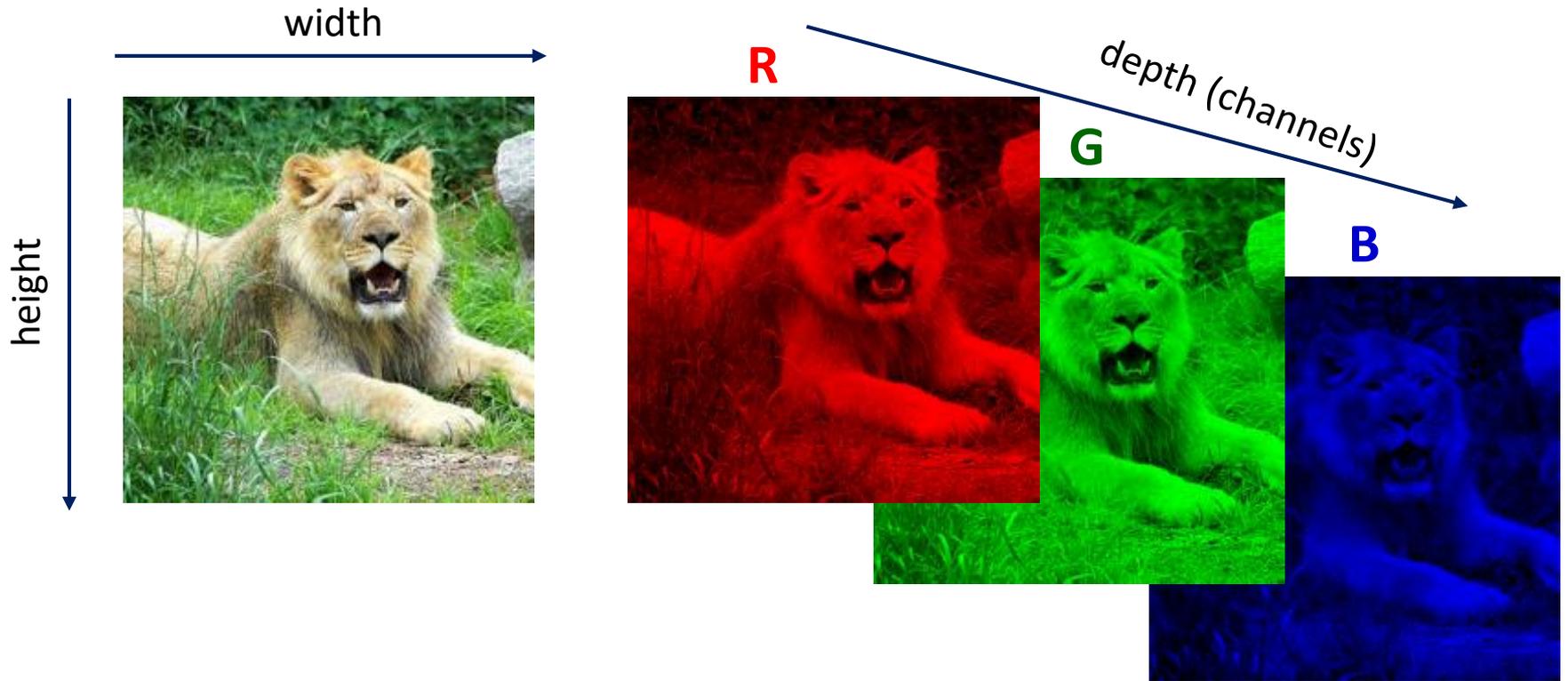
- convolutional layers + ReLU
- pooling layers
- fully connected layers
- normalization layers

Layers in CNN



https://www.youtube.com/watch?v=Jy9-aGMB_TE

Image representation



3D - tensor

Convolution

Convolution operates on two signals (in 1D) or two images (in 2D): you can think of one as the “input” signal (or image), and the other (called the kernel) as a “filter” on the input image, producing an output image.

Convolution takes two images as input and produces a third as output.

https://www.cs.cornell.edu/courses/cs1114/2013sp/sections/S06_convolution.pdf

Different kind of **information can be extracted**, according with the kernel (filter):

- Edge detection
- Sharpening
- Blurring
- Embossing

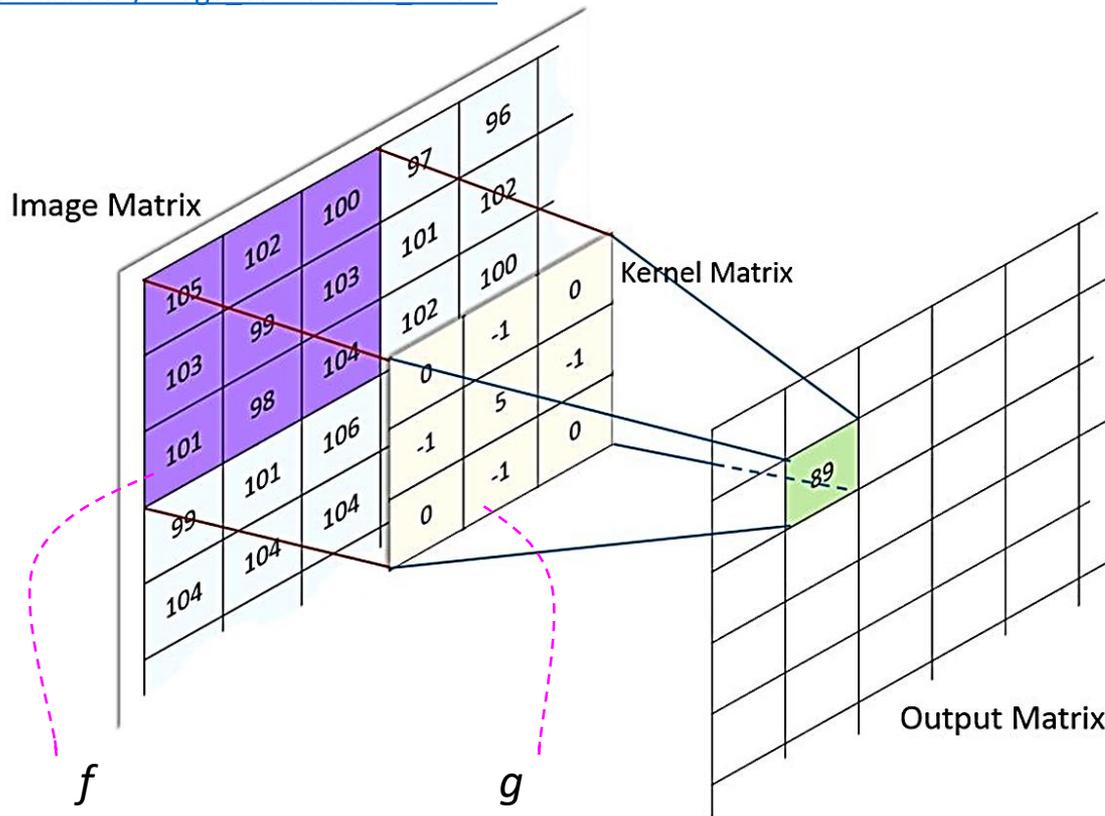
Convolution operation extract the **feature maps**

The values in all filters are training parameters of the CNN; they are learned during the training process

Convolution

Convolution is the process of adding each element of the image to its local neighbors, weighted by the kernel

http://machinelearninguru.com/computer_vision/basics/convolution/image_convolution_1.html



$$f * g$$

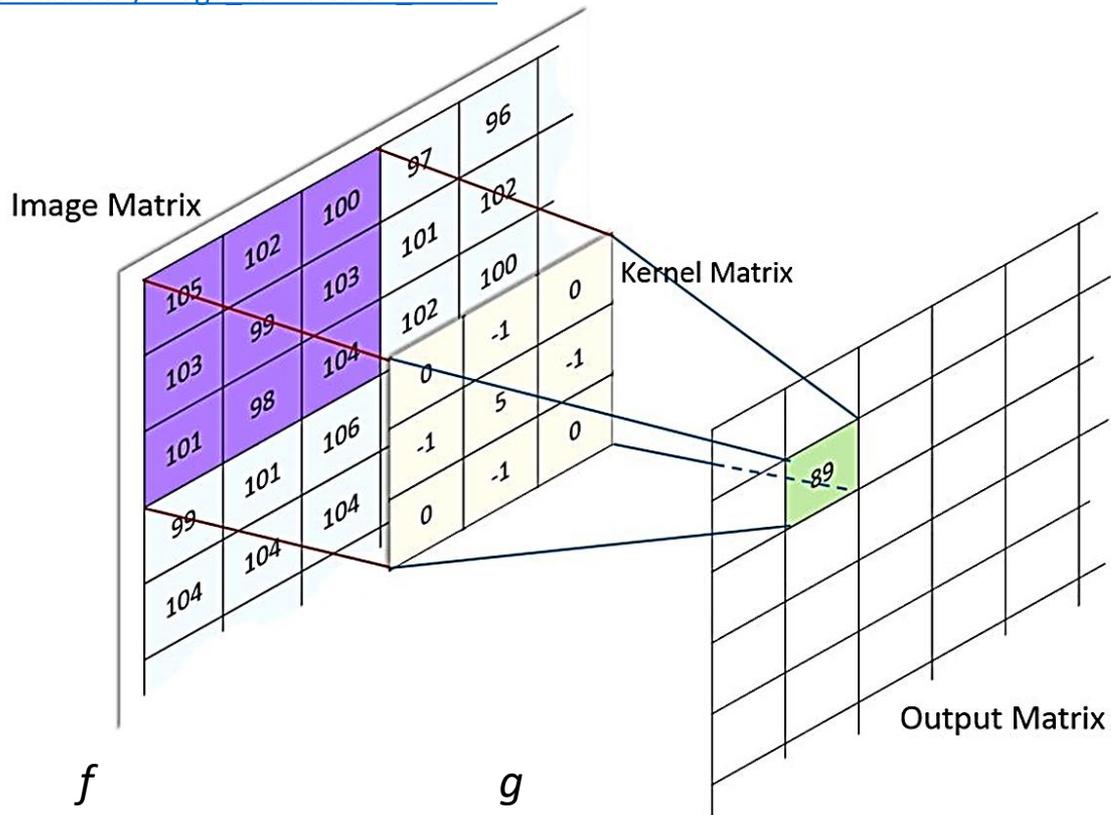
Dot product,
not matrix multiplication

Multiply the arrays
element wise and sum
the product

Convolution

Convolution is the process of adding each element of the image to its local neighbors, weighted by the kernel

http://machinelearninguru.com/computer_vision/basics/convolution/image_convolution_1.html



$f * g:$

$$\begin{aligned} & 105 * 0 + 102 * (-1) + 100 * 0 + \\ & 103 * (-1) + 99 * 5 + 103 * (-1) + \\ & 101 * 0 + 98 * (-1) + 104 * 0 = \quad \mathbf{89} \end{aligned}$$

Sharpening filter (kernel)

Sharpening an image increases the contrast between bright and dark regions to bring out features.

The sharpening process is basically the application of a **high pass filter** to an image.

Convolution

105	102	100	97	96
103	99	103	101	102
101	98	104	102	100
99	101	106	104	99
104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

	89			

Output Matrix

$$105 * 0 + 102 * -1 + 100 * 0 + 103 * -1 + 99 * 5 + 103 * -1 + 101 * 0 + 98 * -1 + 104 * 0 = 89$$

105	102	100	97	96
103	99	103	101	102
101	98	104	102	100
99	101	106	104	99
104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

	89	?		

Output Matrix

Convolution

105	102	100	97	96
103	99	103	101	102
101	98	104	102	100
99	101	106	104	99
104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

Kernel Matrix

	89			

Output Matrix

$$105 * 0 + 102 * -1 + 100 * 0 \\ + 103 * -1 + 99 * 5 + 103 * -1 \\ + 101 * 0 + 98 * -1 + 104 * 0 = 89$$

105	102	100	97	96
103	99	103	101	102
101	98	104	102	100
99	101	106	104	99
104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

Kernel Matrix

	89	111		

Output Matrix

$$102 * 0 + 100 * -1 + 97 * 0 \\ + 99 * -1 + 103 * 5 + 101 * -1 \\ + 98 * 0 + 104 * -1 + 102 * 0 = 111$$

Convolution

105	102	100	97	96
103	99	103	101	102
101	98	104	102	100
99	101	106	104	99
104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

Kernel Matrix

	89			

Output Matrix

$$105 * 0 + 102 * -1 + 100 * 0 \\ + 103 * -1 + 99 * 5 + 103 * -1 \\ + 101 * 0 + 98 * -1 + 104 * 0 = 89$$

Pixels on the border of image matrix?

105	102	100	97	96
103	99	103	101	102
101	98	104	102	100
99	101	106	104	99
104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

Kernel Matrix

?				
?	89	111		

Output Matrix

$$102 * 0 + 100 * -1 + 97 * 0 \\ + 99 * -1 + 103 * 5 + 101 * -1 \\ + 98 * 0 + 104 * -1 + 102 * 0 = 111$$

Convolution

Padding

The process of adding zeros to the input matrix symmetrically to maintain the dimension of output as in input.

(1 pixel padding here)

Padding depends on the dimension of the filter

0	0	0	0	0	0
0	105	102	100	97	96
0	103	99	103	101	102
0	101	98	104	102	100
0	99	101	106	104	99
0	104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

210	89	111		

Output Matrix

$$\begin{aligned}
 &0 * 0 + 105 * -1 + 102 * 0 \\
 &+ 0 * -1 + 103 * 5 + 99 * -1 \\
 &+ 0 * 0 + 101 * -1 + 98 * 0 = 210
 \end{aligned}$$

0	0	0	0	0	0
0	105	102	100	97	96
0	103	99	103	101	102
0	101	98	104	102	100
0	99	101	106	104	99
0	104	104	104	100	98

Image Matrix

0	-1	0
-1	5	-1
0	-1	0

320				
210	89	111		

Output Matrix

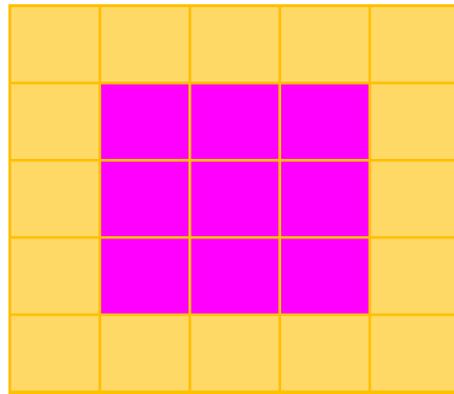
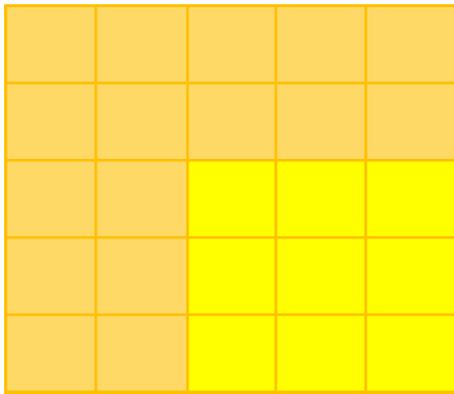
$$\begin{aligned}
 &0 * 0 + 0 * -1 + 0 * 0 \\
 &+ 0 * -1 + 105 * 5 + 102 * -1 \\
 &+ 0 * 0 + 103 * -1 + 99 * 0 = 320
 \end{aligned}$$

Convolution

Stride

Stride denotes how many steps we are moving in each steps in convolution.

Usually it is 1



Original gray image



Original gray image - histogram equalization



Sharpening filter

stride = 1

$$* \begin{array}{|c|c|c|} \hline 0 & -1 & 0 \\ \hline -1 & 5 & -1 \\ \hline 0 & -1 & 0 \\ \hline \end{array} =$$

Sharpen image



Sharpen image - histogram equalization



Original gray image

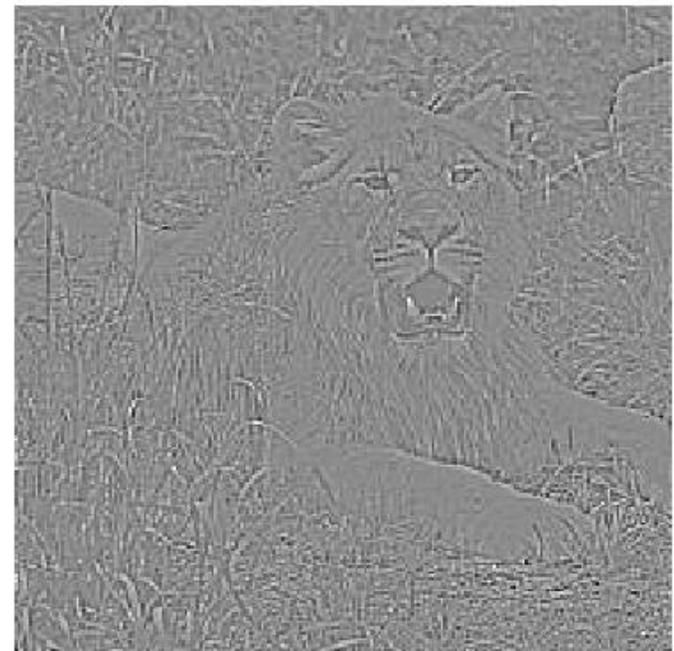


Edge detection filter

stride = 1

$$* \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} =$$

Edge detection image



Original gray image - histogram equalization



Edge detection image - histogram equalization



Original gray image



Original gray image - histogram equalization



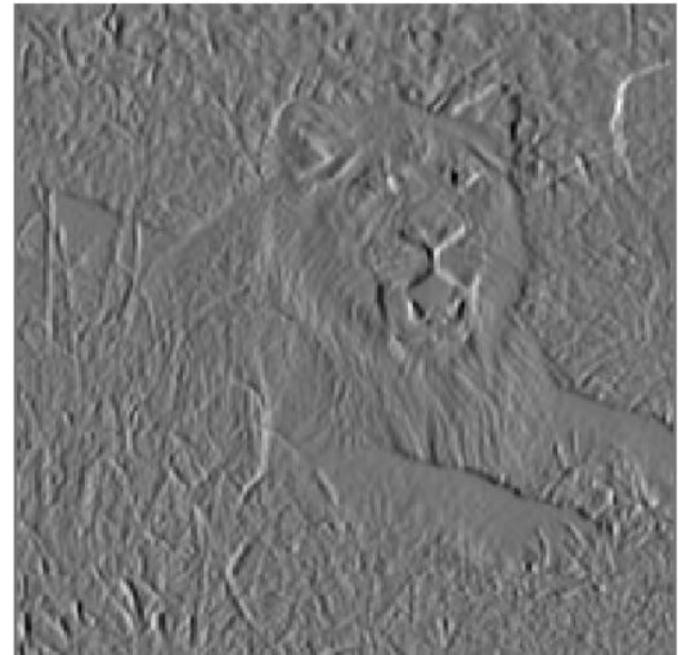
Edge detection filter

stride = 1

$$* \begin{matrix} \begin{matrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{matrix} = \end{matrix}$$

Sobel; horizontal changes (vertical edges)

Edge detection image

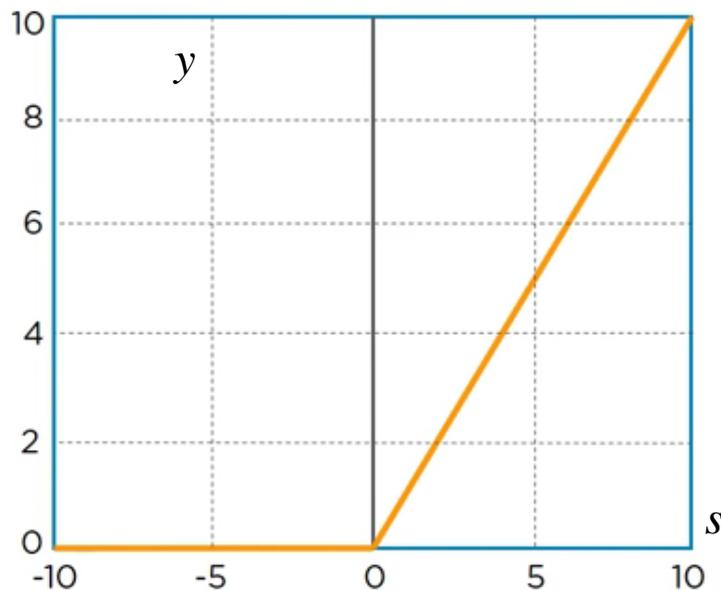


Edge detection image - histogram equalization

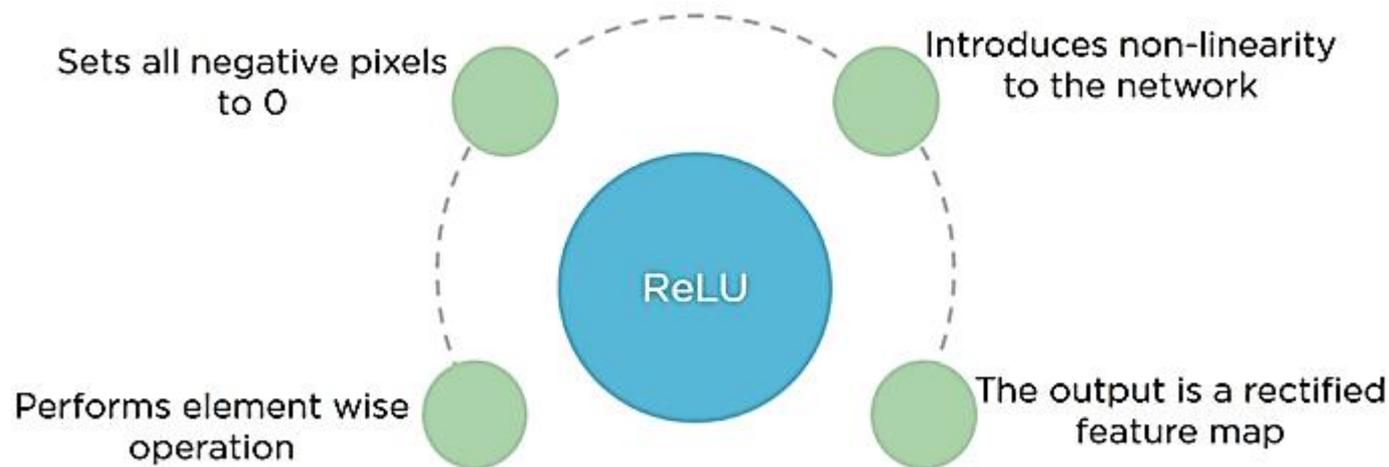


ReLU layer

Once the feature map is extracted by the convolution, the next operation is to apply the ReLU activation function



$$y = \max(0, s)$$

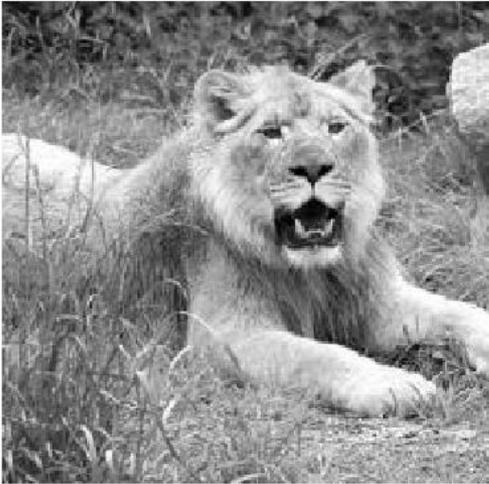


https://www.youtube.com/watch?v=Jy9-aGMB_TE

Convolution
Multiple filters

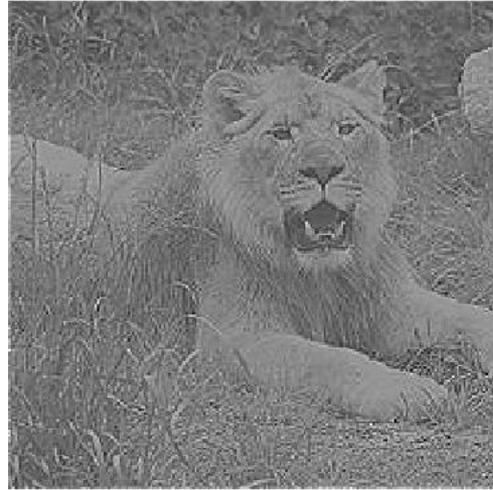


Original gray image

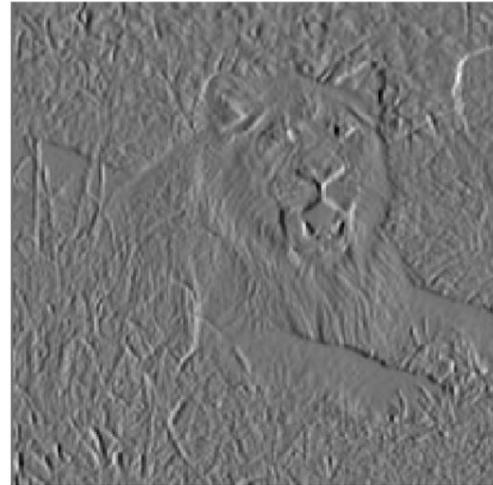


Input image

Sharpen image



Edge detection image



Input feature map

ReLU

Sharpen image + ReLU



Edge detection image + ReLU



Rectified feature map

Input image is scanned in **multiple** convolution and ReLU layers

Original gray image



First 5 columns and rows of the input image matrix:

```
[[33.6861 21.462 22.3216 32.8946 29.1734]
 [62.8323 38.6861 21.462 15.7388 12.1618]
 [53.7622 55.191 50.8323 39.2553 27.7504]
 [67.3392 79.3392 74.3372 57.8323 47.2553]
 [85.2028 79.8441 50.8421 29.3314 34.8323]]
```

Sharpen image



First 5 columns and rows of the image sharpen matrix:

```
[[ 84.1362 12.3258 36.2923 95.4401 78.5729]
 [188.1617 33.9352 -19.8438 -29.4384 -28.3799]
 [ 80.8206 54.1131 63.916 44.8381 15.6216]
 [121.9746 116.2651 133.6276 98.9823 78.1111]
 [196.6922 98.347 8.774 -31.7601 18.4111]]
```

Sharpen image + ReLU



First 5 columns and rows of the image sharpen+ReLU matrix:

```
[[ 84.1362 12.3258 36.2923 95.4401 78.5729]
 [188.1617 33.9352 0. 0. 0. ]
 [ 80.8206 54.1131 63.916 44.8381 15.6216]
 [121.9746 116.2651 133.6276 98.9823 78.1111]
 [196.6922 98.347 8.774 0. 18.4111]]
```

Pooling layer

Once the feature map is rectified by the ReLU activation function, the next operation is to **down-sampling** the images to **reduce the dimensionality** through a **pooling layer**

0	0	14	82
149	32	0	0
28	53	64	44
39	120	133	99

Rectified feature map

$$4 \times 4 = 16$$

max pooling



2 x 2 filter
stride 2

149	82
120	133

Pooled feature map

$$2 \times 2 = 4$$

Dimensionality reduction (given by the filter size and stride):

4 to 1; 4 times

Dimensionality reduction for 2x2 filter and stride 1?

Pooling layer

- Max pooling – how valuable is a feature in the area of the filter
- best (maximum feature value)

0	0	14	82
149	32	0	0
28	53	64	44
39	120	133	99

Rectified feature map

max pooling



2 x 2 filter
stride 2

149	82
120	133

Pooled feature map

$$2 \times 2 = 4$$

Pooling layers uses different filters to be able to identify different aspects in an image: edges, corners, body parts (eyes, ear, paw, fur, etc)

Flattening

Once the dimensionality is reduced to a manageable size, we have to connect further with the fully connected layer.

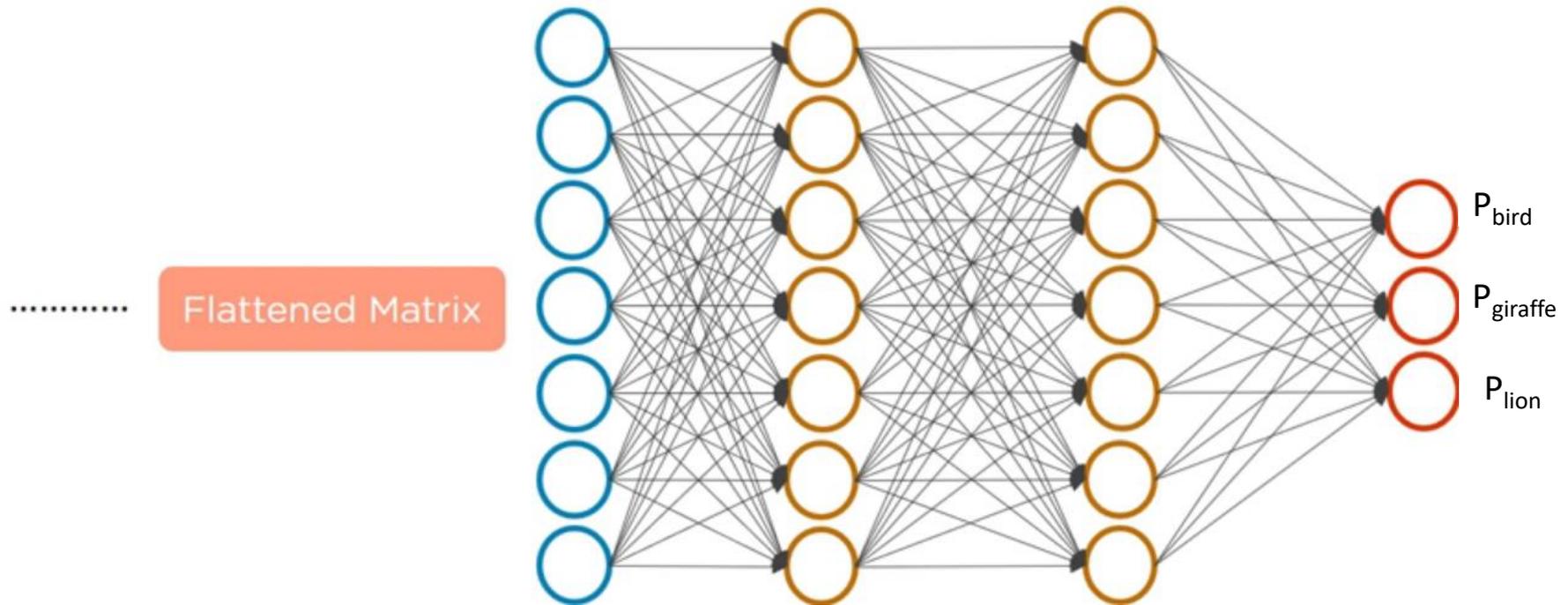
We need to convert the 2D dimensional array from pooled featured map to a single vector using the **flattening** operation



Pooled feature map

Fully connected layer

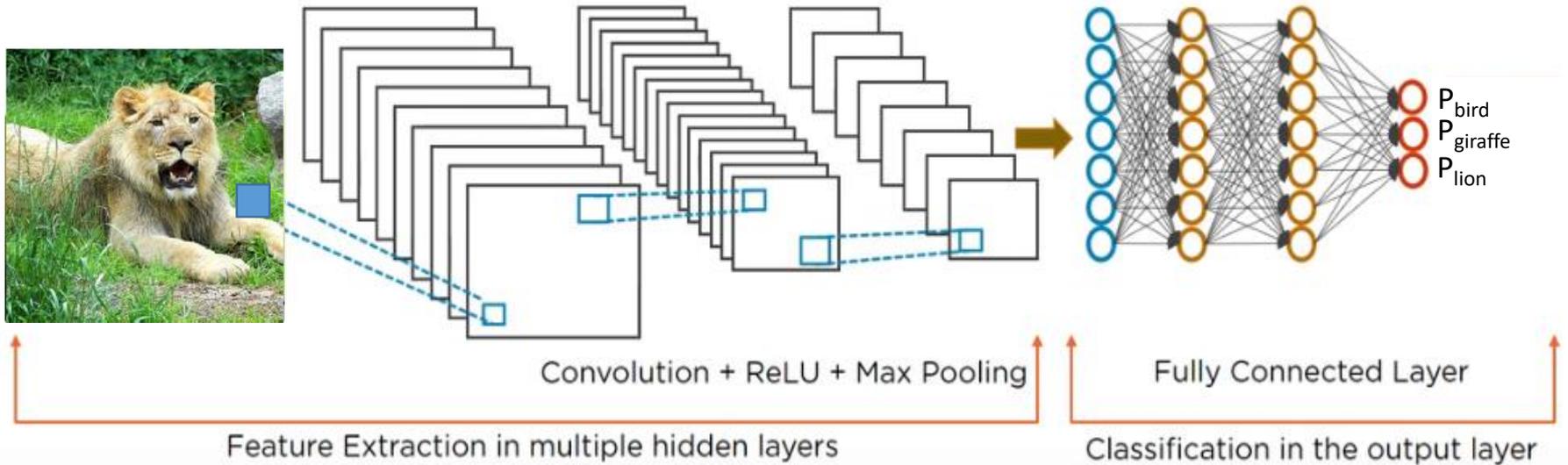
The vector (flattened 2D array) from the pooling layer is fed to the fully connected layer (conventional feed-forward ANN) to classify the image



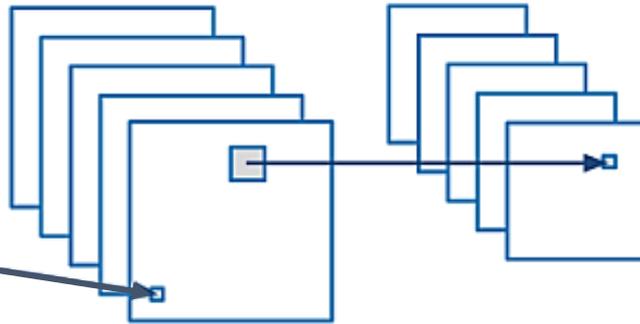
https://www.youtube.com/watch?v=Jy9-aGMB_TE

Fully connected layer

CNN – big picture



CNN – big picture

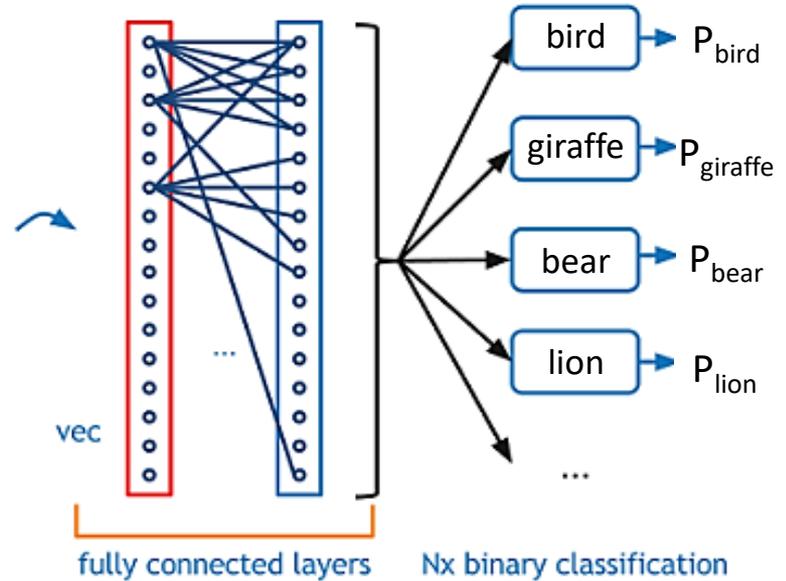


convolution +
nonlinearity

max pooling

convolution + pooling layers

Feature extraction



vec

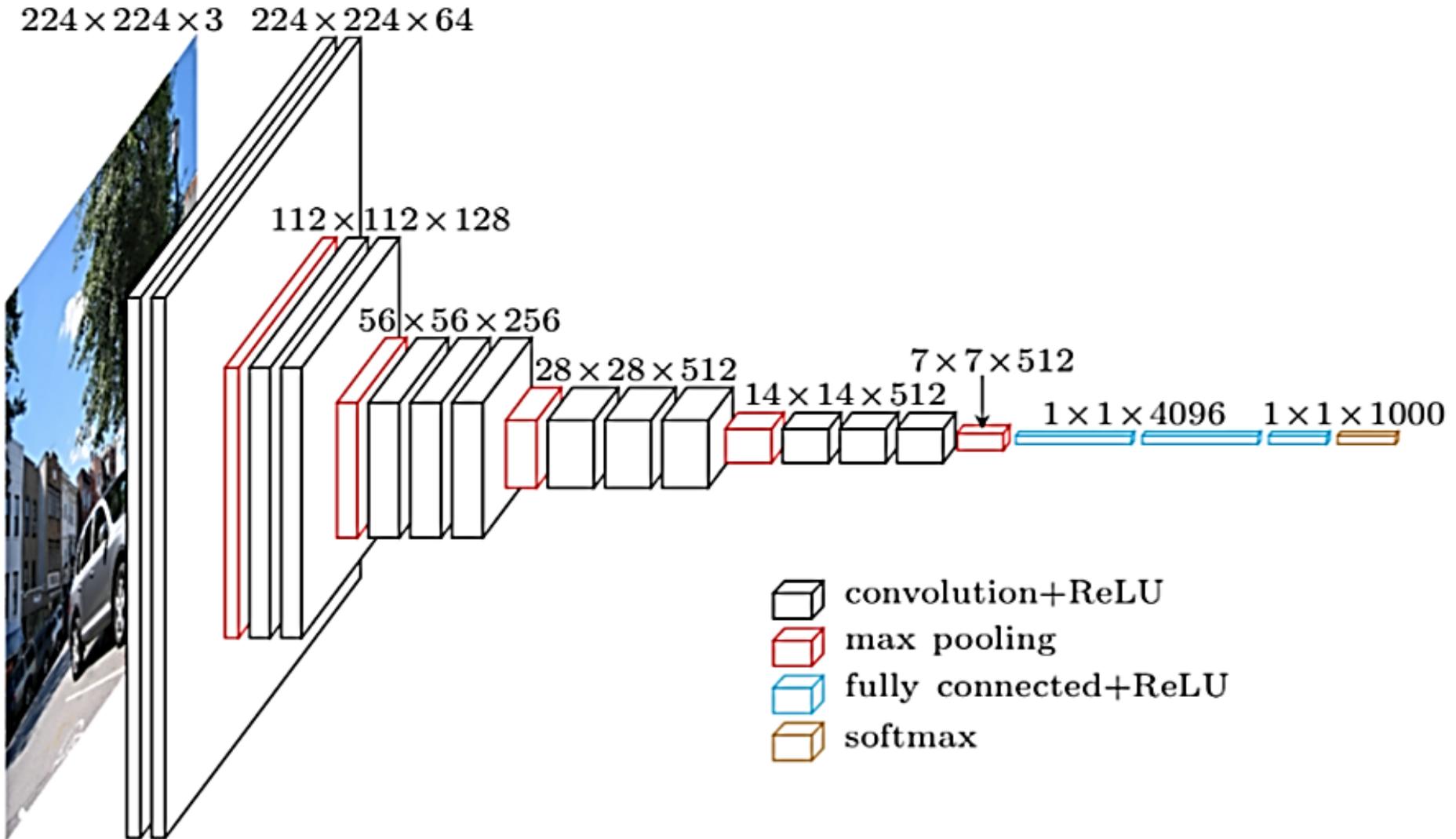
fully connected layers

$N \times$ binary classification

Classification

Adaptation after: <https://towardsdatascience.com/convolutional-neural-network-cb0883dd6529>

CNN – the big picture



Convolution Demo. Below is a running demo of a CONV layer. Since 3D volumes are hard to visualize, all the volumes (the input volume (in blue), the weight volumes (in red), the output volume (in green)) are visualized with each depth slice stacked in rows. The input volume is of size $W_1 = 5, H_1 = 5, D_1 = 3$, and the CONV layer parameters are $K = 2, F = 3, S = 2, P = 1$. That is, we have two filters of size 3×3 , and they are applied with a stride of 2. Therefore, the output volume size has spatial size $(5 - 3 + 2)/2 + 1 = 3$. Moreover, notice that a padding of $P = 1$ is applied to the input volume, making the outer border of the input volume zero. The visualization below iterates over the output activations (green), and shows that each element is computed by elementwise multiplying the highlighted input (blue) with the filter (red), summing it up, and then offsetting the result by the bias.

<http://cs231n.github.io/convolutional-networks/>

Input Volume (+pad 1) (7x7x3)

x[:, :, 0]						
0	0	0	0	0	0	0
0	2	1	1	1	0	0
0	2	1	1	2	2	0
0	2	1	1	0	0	0
0	1	2	0	2	2	0
0	0	2	0	1	0	0
0	0	0	0	0	0	0
x[:, :, 1]						
0	0	0	0	0	0	0
0	0	0	1	1	2	0
0	1	2	1	1	0	0
0	0	0	2	0	2	0
0	2	0	0	0	2	0
0	0	2	2	1	1	0
0	0	0	0	0	0	0
x[:, :, 2]						
0	0	0	0	0	0	0
0	0	2	0	2	1	0
0	0	0	2	1	2	0
0	0	0	1	0	2	0
0	1	2	2	0	2	0
0	2	0	1	1	1	0
0	0	0	0	0	0	0

Filter W0 (3x3x3)

w0[:, :, 0]		
1	-1	-1
-1	-1	-1
1	-1	1
w0[:, :, 1]		
-1	1	-1
0	1	1
1	-1	1
w0[:, :, 2]		
1	0	0
-1	1	0
-1	0	0
Bias b0 (1x1x1)		
b0[:, :, 0]		
1		

Filter W1 (3x3x3)

w1[:, :, 0]		
1	0	-1
1	1	0
-1	0	1
w1[:, :, 1]		
1	-1	1
-1	0	1
0	-1	0
w1[:, :, 2]		
0	0	0
-1	1	-1
-1	0	1
Bias b1 (1x1x1)		
b1[:, :, 0]		
0		

Output Volume (3x3x2)

o[:, :, 0]		
-2	2	1
-7	0	3
2	4	3
o[:, :, 1]		
0	0	-4
4	2	1
0	1	0

toggle movement

Input Volume (+pad 1) (7x7x3)

$x[:, :, 0]$

0	0	0	0	0	0	0
0	2	1	1	1	0	0
0	2	1	1	2	2	0
0	2	1	1	0	0	0
0	1	2	0	2	2	0
0	0	2	0	1	0	0
0	0	0	0	0	0	0

$x[:, :, 1]$

0	0	0	0	0	0	0
0	0	0	1	1	2	0
0	1	2	1	1	0	0
0	0	0	2	0	2	0
0	2	0	0	0	2	0
0	0	2	2	1	1	0
0	0	0	0	0	0	0

$x[:, :, 2]$

0	0	0	0	0	0	0
0	0	2	0	2	1	0
0	0	0	2	1	2	0
0	0	0	1	0	2	0
0	1	2	2	0	2	0
0	2	0	1	1	1	0
0	0	0	0	0	0	0

Filter W0 (3x3x3)

$w0[:, :, 0]$

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

$w0[:, :, 1]$

-1	1	-1
0	1	1
1	-1	1

$w0[:, :, 2]$

1	0	0
-1	1	0
-1	0	0

Bias b0 (1x1x1)

$b0[:, :, 0]$

1

Filter W1 (3x3x3)

$w1[:, :, 0]$

1	0	-1
1	1	0
-1	0	1

$w1[:, :, 1]$

1	-1	1
-1	0	1
0	-1	0

$w1[:, :, 2]$

0	0	0
-1	1	-1
-1	0	1

Bias b1 (1x1x1)

$b1[:, :, 0]$

0

Output Volume (3x3x2)

$o[:, :, 0]$

-2	2	1
-7	0	3
2	4	3

$o[:, :, 1]$

0	0	-4
4	2	1
0	1	0

toggle movement