

CSE 4553
Machine Learning
Lecture 8: Convolutional Neural Network (CNN)

Winter 2024

Prof. Dr. Hasan Mahmud | hasan@iut-dhaka.edu

Deep Learning

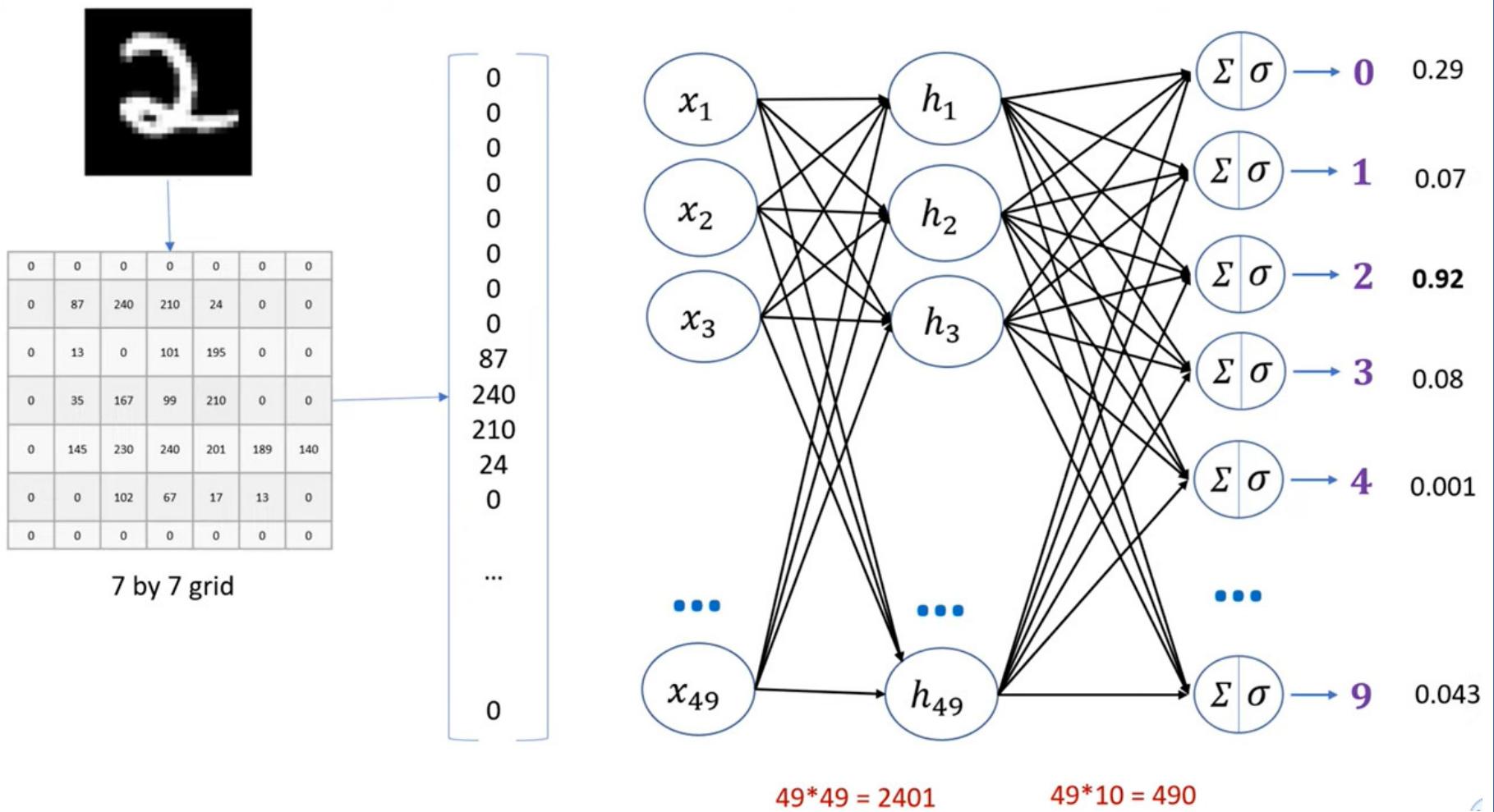
- Deep learning is supervised learning technique that contains a stack of hidden layers in addition to input and output layers of ANN.
- Deep Neural Network are capable of extracting deep features out of the data; hence the name deep learning.
- Some deep learning algorithms are:
 - Convolution Neural Network (CNN)
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)

ANN in handwritten digit recognition

- ANN can handle these variations of input of the hand written digit '9'



ANN in handwritten digit recognition



Problems with ANN



Image size = 1920 x 1080 X 3

First layer neurons = $1920 \times 1080 \times 3 \sim 6 \text{ million}$

↳

Hidden layer neurons = Let's say you keep it $\sim 4 \text{ million}$

Weights between input and hidden layer = $6 \text{ mil} * 4 \text{ mil}$
 $= 24 \text{ million}$

- Too much computation
- Treats local pixels same as pixels far apart
- Sensitive to location of an object in an image

Introduction to CNN

- In 2012, neural nets grew to prominence as Alex Krizhevsky used them to win that year's ImageNet competition (basically, the annual Olympics of computer vision), dropping the classification error record from 26% to 15%, an astounding improvement at the time.
- Deep learning-based applications:
 - Facebook uses neural nets for their automatic tagging algorithms
 - Google for their photo search
 - Amazon for their product recommendations
 - Pinterest for their home feed personalization
 - Instagram for their search infrastructure

Introduction to CNN...

- When we see an image or just when we look at the world around us, most of the time we are able to immediately characterize the scene and give each object a label, all without even consciously noticing.
- These skills of being able to quickly recognize patterns, generalize from prior knowledge, and adapt to different image environments are ones that we do not share with our fellow machines.



Introduction to CNN...

- Image classification is the task of taking an input image and outputting a class (a cat, dog, etc) or a probability of classes that best describes the image.
- The way that we see images around us is different than the way machine sees it.



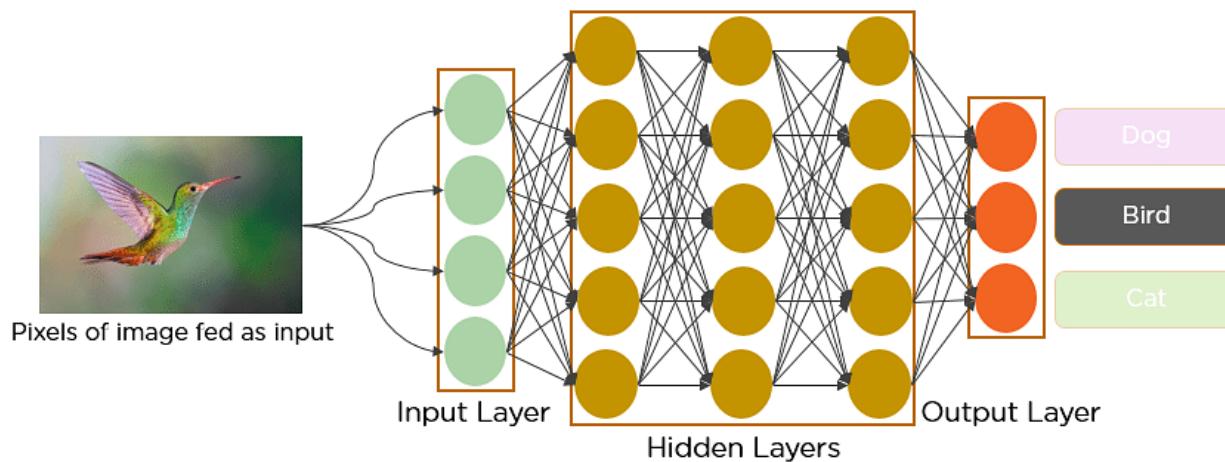
What we see

0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

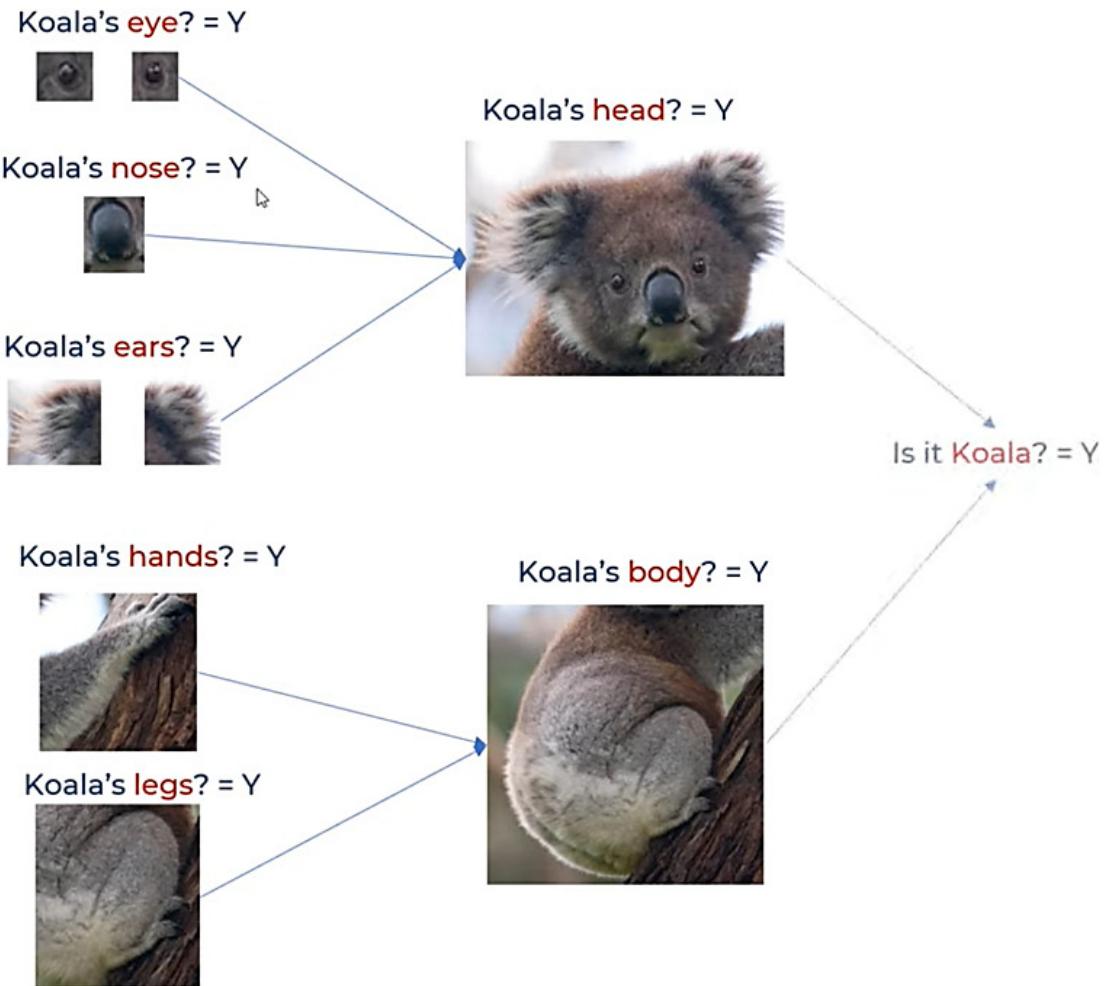
What a computer sees

Convolution Neural Network (CNN)

- It is a multilayered feed-forward neural network with a special architecture to detect complex features in data.
- CNNs are able to classify images by detecting features, similar to how the human brain detects features to identify objects.
- Images are high-dimensional vectors. It would take a huge amount of parameters to characterize the network. To address this problem, bionic convolutional neural networks are proposed to reduced the number of parameters and adapt the network architecture specifically to vision tasks.

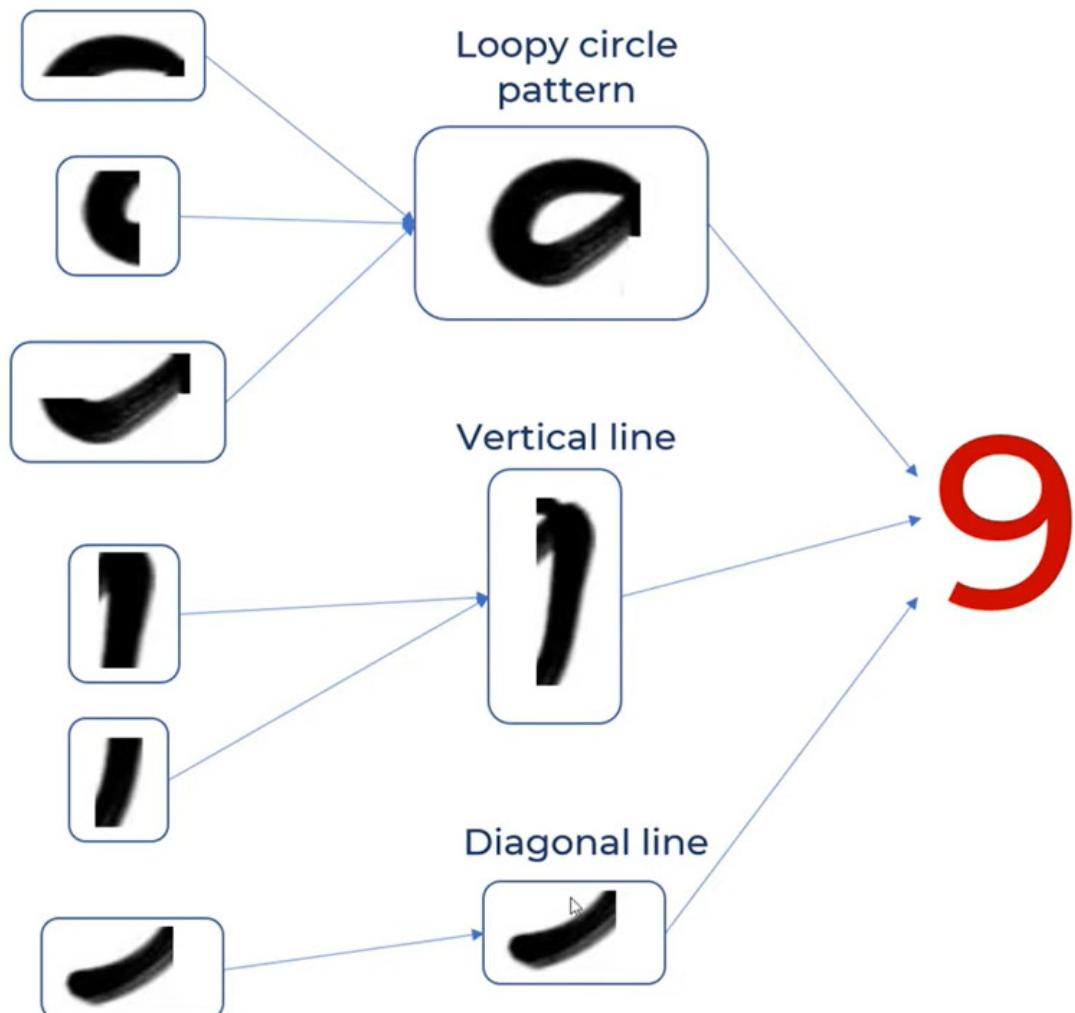


Conceptual understanding of CNN

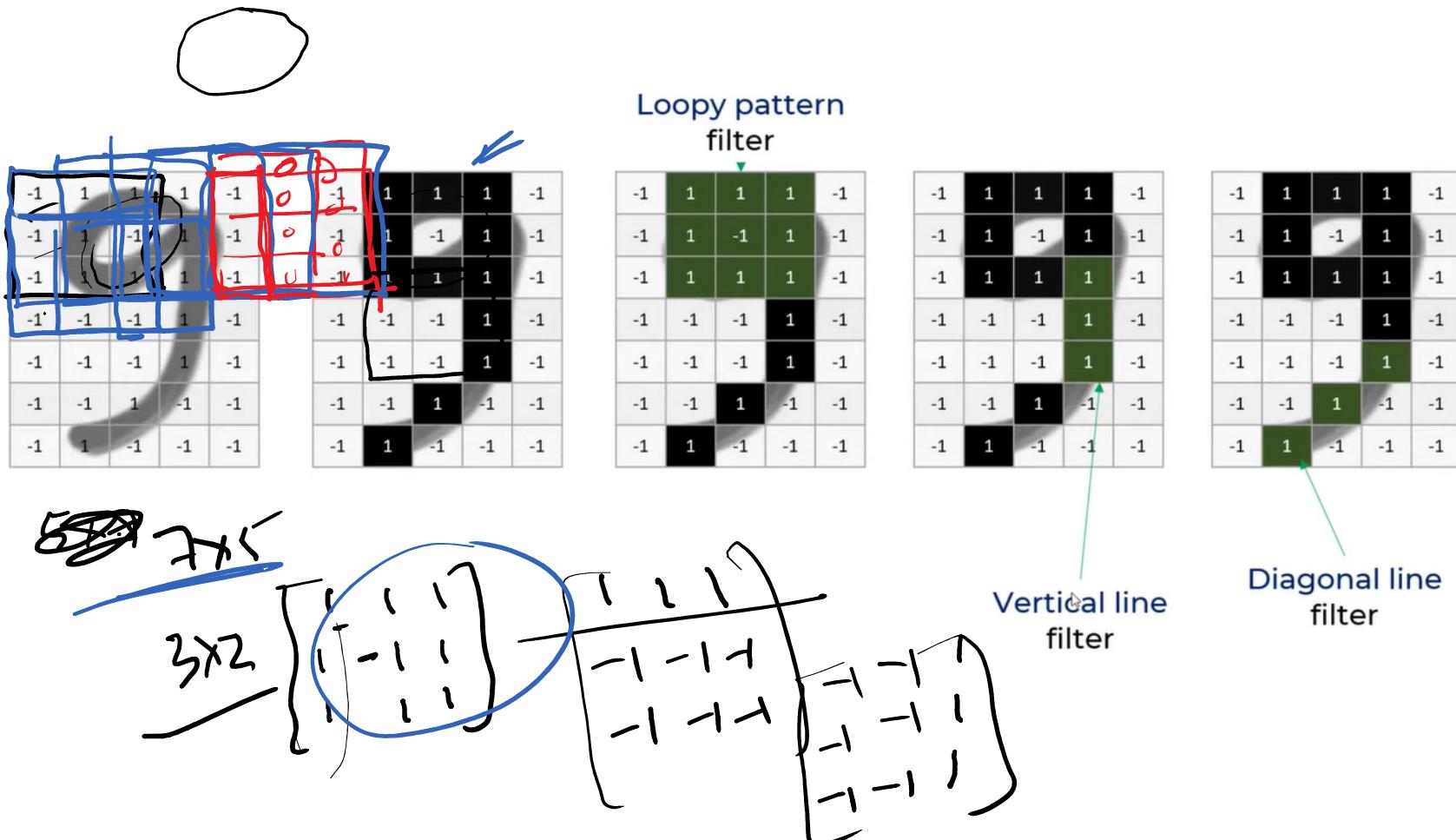


Conceptual understanding of CNN

g

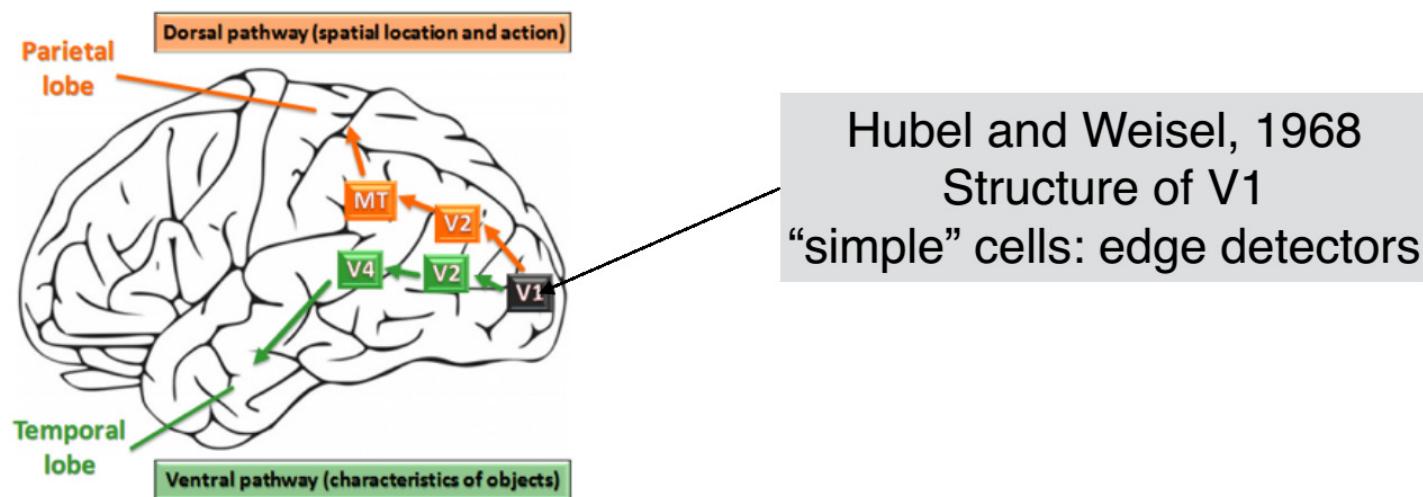


Conceptual understanding of CNN



CNN

- ◆ A CNN unit contains the following layers:
 1. Convolutional layer containing a set of filters
 2. Pooling layer
 3. Non-linearity
- ◆ Deep CNN: a stack of multiple CNN units
 - Inspired by the human visual system (V1, V2, V3)

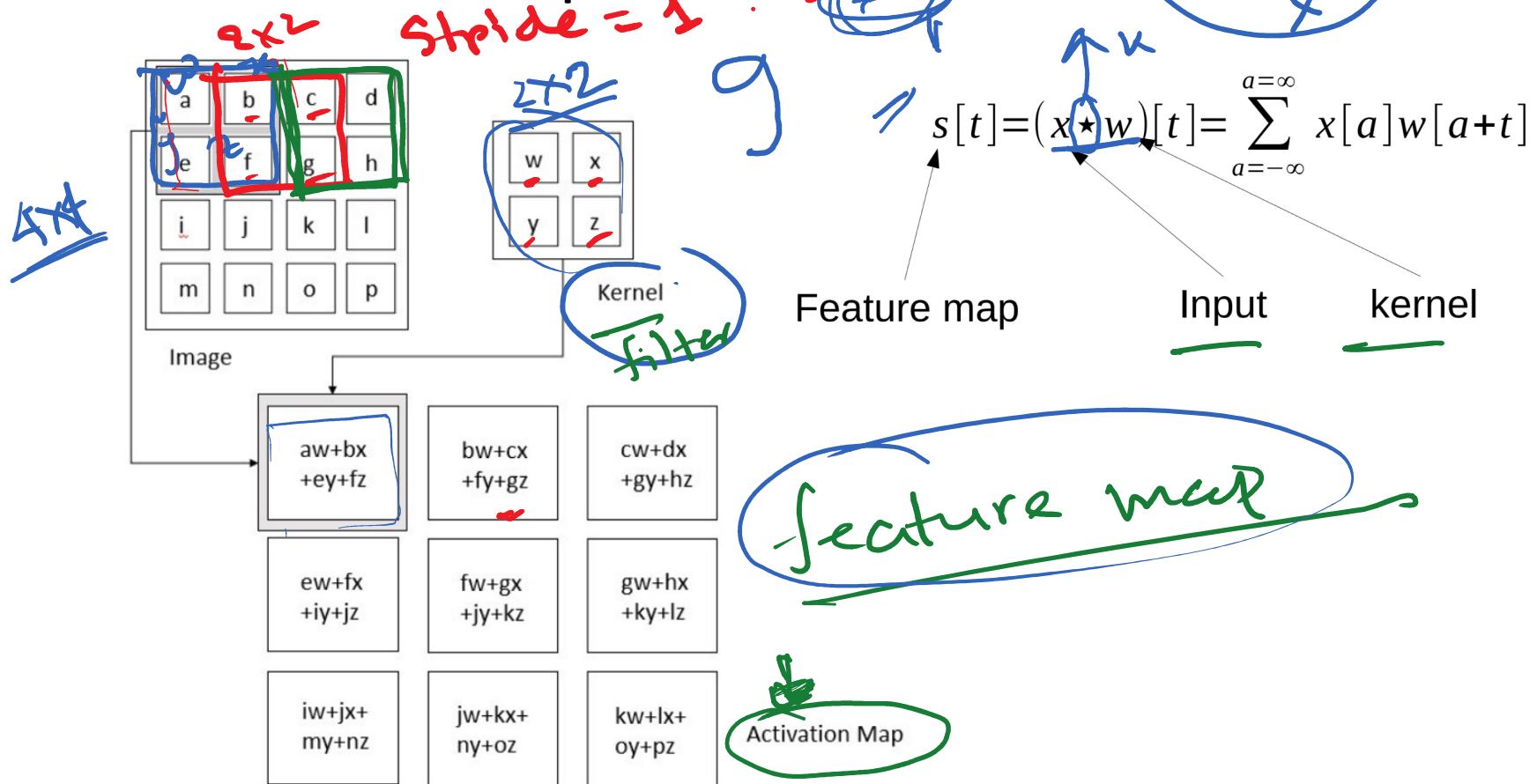


Typical steps of CNN

- **Step 1: Convolution**
- **Step 1b: ReLU Layer**
- **Step 2: Pooling**
- **Step 3: Flattening**
- **Step 4: Full Connection**

Convolution

- Convolution operation



Convolution...

$$-1+1+1-1-1-1+1+1 = -1 \rightarrow -1/9 = -0.11$$



7x5

*

1	1	1
1	-1	1
1	1	1

3×3

-0.11		

Feature / Activation

Convolution...

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

*

1	1	1
1	-1	1
1	1	1

-0.11	1	



Convolution...

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

*

1	1	1
1	-1	1
1	1	1

-0.11	1	-0.11



Convolution...

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

*

1	1	1
1	-1	1
1	1	1

-0.11	1	-0.11
-0.55		



Convolution...

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

*

1	1	1
1	-1	1
1	1	1

-0.11	1	-0.11
-0.55	0.11	



Convolution...

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

*

1	1	1
1	-1	1
1	1	1

-0.11	1	-0.11
-0.55	0.11	-0.33



Convolution...

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

*

1	1	1
1	-1	1
1	1	1

$$f = \max(0, z)$$

-0.11	1	-0.11
-0.55	0.11	-0.33
-0.33	0.33	-0.33
-0.22	-0.11	-0.22
-0.33	-0.33	-0.33

Feature Map
OR Activation Map

Convolution...

9 * Loopy pattern detector =

1	1	1
1	-1	1
1	1	1

6 * Loopy pattern detector =

1	1	1
1	-1	1
1	1	1

8 * Loopy pattern detector =

1	1	1
1	-1	1
1	1	1

96 * Loopy pattern detector =

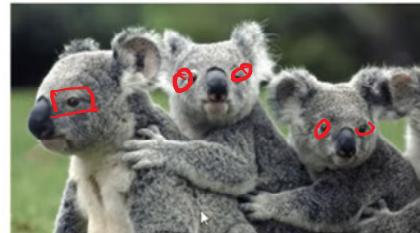
1	1	1
1	-1	1
1	1	1

Shared weight and sliding window

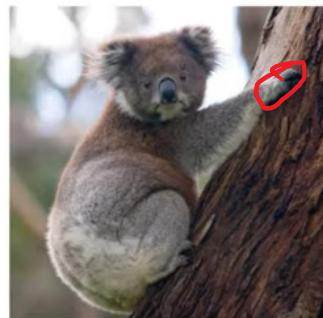
Convolution...



$$* \begin{array}{c} \text{eye} \\ \text{detector} \\ \hline \text{---} \\ \text{---} \end{array} = \boxed{\begin{array}{cc} 1 & 1 \end{array}}$$



$$* \begin{array}{c} \text{eye} \\ \text{detector} \\ \hline \text{---} \\ \text{---} \end{array} = \boxed{\begin{array}{cccc} \checkmark & & & \\ & 1 & 1 & \\ & & \dots & \end{array}}$$



$$* \begin{array}{c} \text{hands} \\ \text{detector} \\ \hline \text{---} \\ \text{---} \end{array} = \boxed{\begin{array}{c} \checkmark \end{array}}$$

Location invariant: It can detect eyes in any location of the image

spatial



$$* \begin{array}{c} \text{eye} \\ \text{detector} \\ \hline \text{---} \\ \text{---} \end{array} = \boxed{\begin{array}{cc} 1 & 1 \end{array}}$$

Convolution...

Loopy pattern detector

$$\begin{matrix} \text{9} & * & \begin{matrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \end{matrix} & = & \begin{matrix} & & 1 & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix} \end{matrix}$$

Vertical line detector

$$\begin{matrix} \text{9} & * & \begin{matrix} -1 & 1 & -1 \\ -1 & 1 & -1 \\ -1 & 1 & -1 \end{matrix} & = & \begin{matrix} & & & & \\ & & & & \\ & & 1 & & \\ & & & & \\ & & & & \end{matrix} \end{matrix}$$

Diagonal line detector

$$\begin{matrix} \text{9} & * & \begin{matrix} -1 & -1 & 1 \\ -1 & 1 & -1 \\ 1 & -1 & -1 \end{matrix} & = & \begin{matrix} & & & & \\ & & & & \\ & & & 1 & & \\ & & & & & \\ & & & & & \end{matrix} \end{matrix}$$

Feature Maps

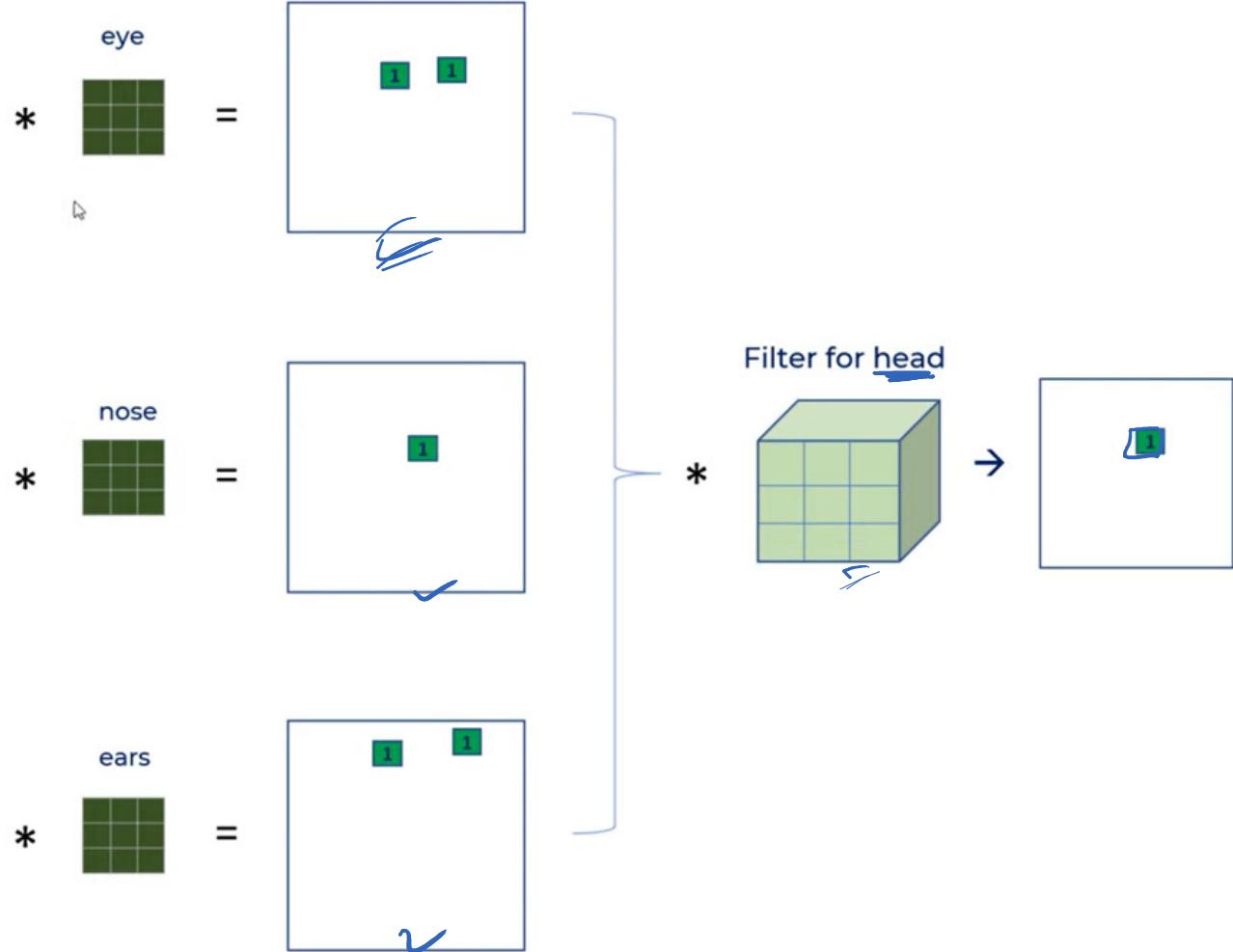
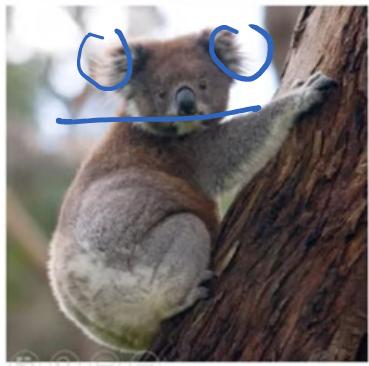
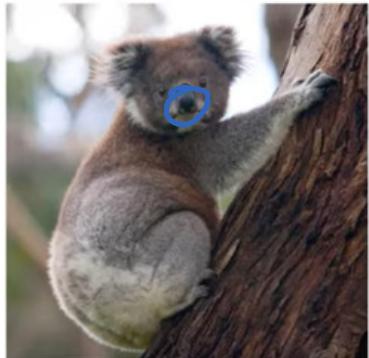
Filters

$$\begin{matrix} \text{9} & * & \begin{matrix} 1 & 1 & 1 & 1 & 1 \\ -1 & -1 & 1 & 1 & 1 \\ 1 & 1 & -1 & 1 & 1 \\ -1 & 1 & -1 & 1 & 1 \\ 1 & 1 & -1 & 1 & 1 \end{matrix} & = & \begin{matrix} & & 1 & & \\ & & & & \\ & & 1 & & \\ & & & & \\ & & & & \end{matrix} \end{matrix}$$

3D-filter

3D-feature map.

Convolution...



Convolution...



$$* \begin{matrix} \text{eye} \\ \text{grid} \end{matrix} = \begin{matrix} \text{grid} \\ \text{grid} \end{matrix}$$



$$* \begin{matrix} \text{nose} \\ \text{grid} \end{matrix} = \begin{matrix} \text{grid} \\ \text{grid} \end{matrix}$$



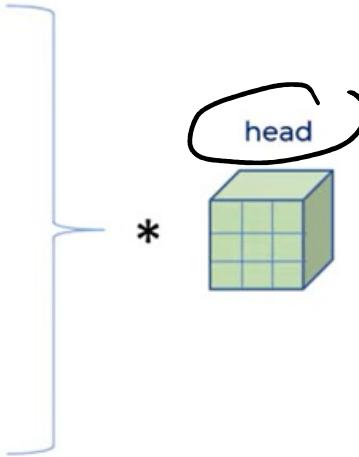
$$* \begin{matrix} \text{ears} \\ \text{grid} \end{matrix} = \begin{matrix} \text{grid} \\ \text{grid} \end{matrix}$$



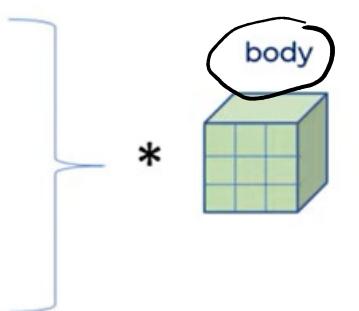
$$* \begin{matrix} \text{hands} \\ \text{grid} \end{matrix} = \begin{matrix} \text{grid} \\ \text{grid} \end{matrix}$$



$$* \begin{matrix} \text{legs} \\ \text{grid} \end{matrix} = \begin{matrix} \text{grid} \\ \text{grid} \end{matrix}$$

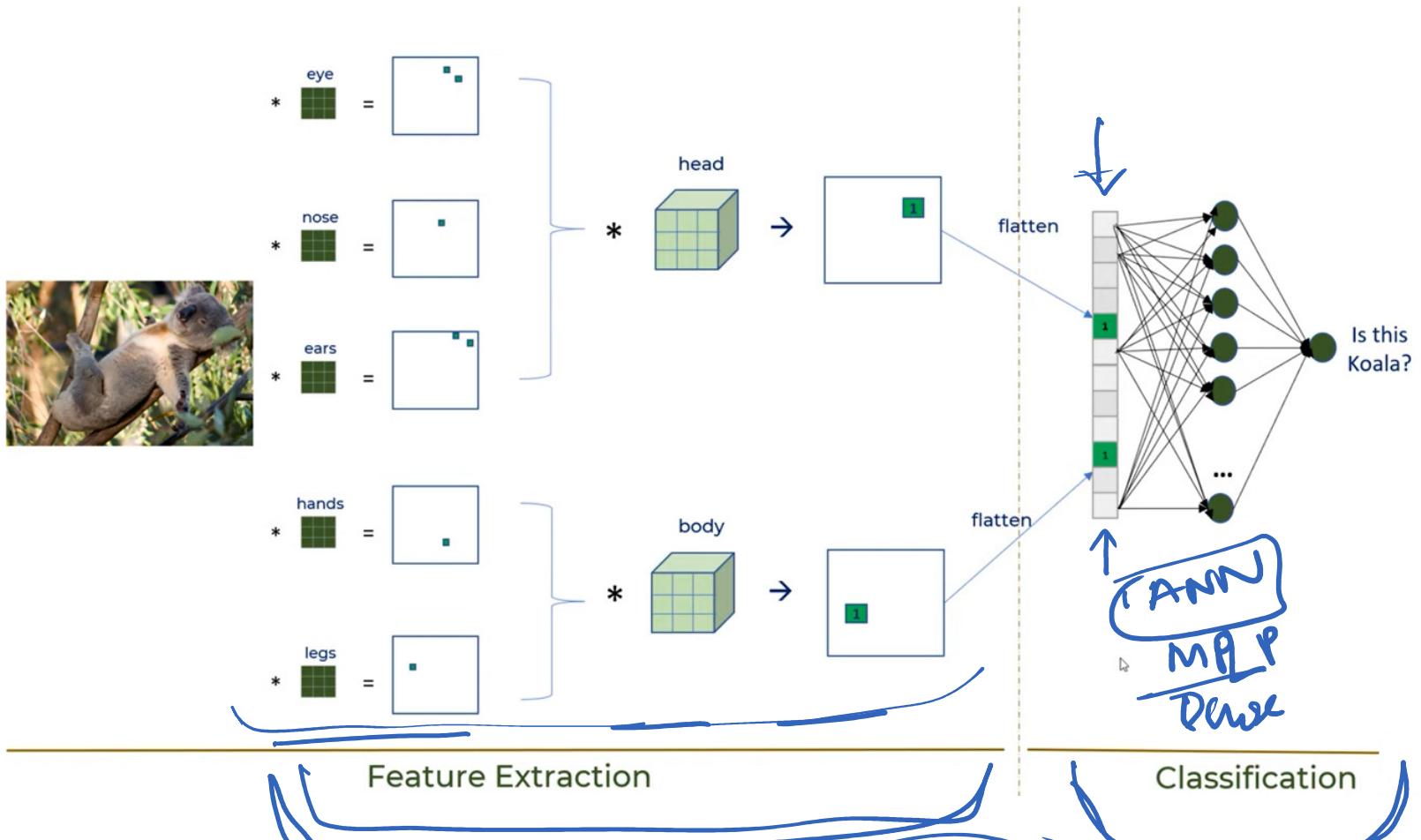


$$\rightarrow \begin{matrix} \text{grid} \\ \text{grid} \end{matrix}$$



$$\rightarrow \begin{matrix} \text{grid} \\ \text{grid} \end{matrix}$$

Convolution...



Convolution...

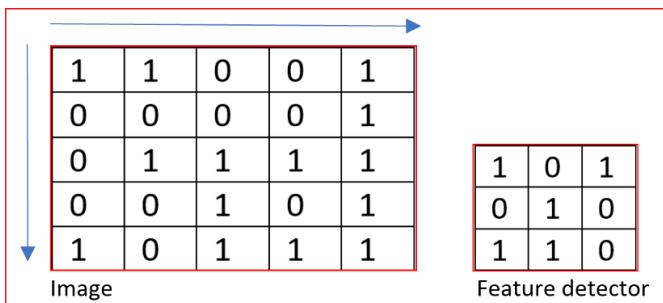
- Kernel: A kernel is a small matrix whose contents are based upon the operations to be performed. The kernel is the filter which is convolved with the original image to detect important patterns.

<i>Original</i>	<i>Gaussian Blur</i>	<i>Sharpen</i>	<i>Edge Detection</i>
$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$
			

- Stride: Stride defines by what step does the kernel move, for example stride of 1 makes kernel slide by one row/column at a time and stride of 2 moves kernel by 2 rows/columns.

Convolution...

- let's take an image of 5 by 5 matrix with 3 channels(RGB) , a feature detector of 3 by 3 with 3 channels (RGB) and scan the feature detector over the image by 1 stride.
- Dimension of the feature map as a function of the input image size(W), feature detector size(F), Stride(S) and Zero Padding on image(P) is: $(W - F + 2P)/S + 1$**



- Input image size W in our case is 5.
- Feature detector or receptive field size is F , which in our case is 3
- Stride (S) is 1, and the amount of zero padding used (P) on the image is 0.
- so, our feature map dimension will $(5-3+0)/1 + 1 = 3$.
- so feature map will a 3*3 matrix with three channels(RGB).**

<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>1</td><td>1</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> </table>	1	1	0	0	1	0	0	0	0	1	0	1	1	1	1	0	0	1	0	1	1	0	1	1	1	$(1*1 + 0*1 + 1*0 + 0*0 + 1*0 + 1*0 + 1*0 + 1*1 + 0*1) = 2$
1	1	0	0	1																						
0	0	0	0	1																						
0	1	1	1	1																						
0	0	1	0	1																						
1	0	1	1	1																						
<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>1</td><td>1</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> </table>	1	1	0	0	1	0	0	0	0	1	0	1	1	1	1	0	0	1	0	1	1	0	1	1	1	$2 \quad (1*1 + 0*0 + 1*0 + 1*0 + 1*0 + 1*0 + 1*1 + 1*1 + 0*1) = 3$
1	1	0	0	1																						
0	0	0	0	1																						
0	1	1	1	1																						
0	0	1	0	1																						
1	0	1	1	1																						
<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>1</td><td>1</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> </table>	1	1	0	0	1	0	0	0	0	1	0	1	1	1	1	0	0	1	0	1	1	0	1	1	1	$2 \quad 3 \quad (1*0 + 0*0 + 1*1 + 0*0 + 1*0 + 0*1 + 1*1 + 1*1 + 0*1) = 3$
1	1	0	0	1																						
0	0	0	0	1																						
0	1	1	1	1																						
0	0	1	0	1																						
1	0	1	1	1																						
<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>1</td><td>1</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> </table>	1	1	0	0	1	0	0	0	0	1	0	1	1	1	1	0	0	1	0	1	1	0	1	1	1	$2 \quad 3 \quad 3 \quad (1*0 + 0*0 + 1*0 + 0*0 + 1*1 + 0*1 + 1*0 + 1*0 + 1*1) = 2$
1	1	0	0	1																						
0	0	0	0	1																						
0	1	1	1	1																						
0	0	1	0	1																						
1	0	1	1	1																						
<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>1</td><td>1</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> </table>	1	1	0	0	1	0	0	0	0	1	0	1	1	1	1	0	0	1	0	1	1	0	1	1	1	$2 \quad 3 \quad 3 \quad 3$
1	1	0	0	1																						
0	0	0	0	1																						
0	1	1	1	1																						
0	0	1	0	1																						
1	0	1	1	1																						

Continuing to scan through input matrix and the final feature map will be

2	3	3
2	2	3
2	4	4

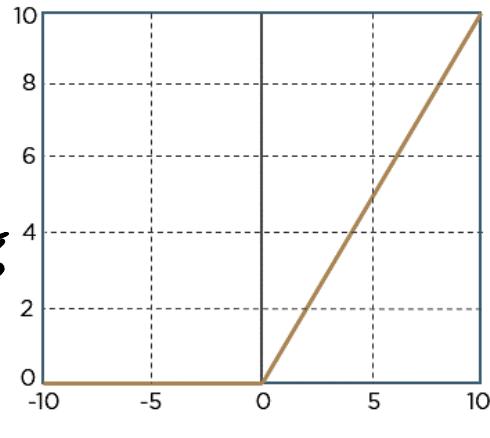
Padding

0	0	0	0	0	0	0
0	1	1	0	0	1	0
0	0	0	0	0	1	0
0	0	1	1	1	1	0
0	0	0	1	0	1	0
0	1	0	1	1	1	0
0	0	0	0	0	0	0

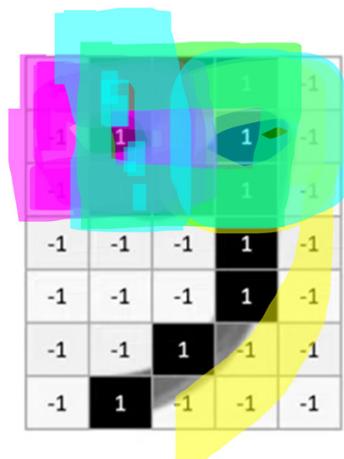
- Losing the edge pixels
- Padding with zero value pixels
- Reflection padding

Rectified Linear Unit (ReLU)

- ① shared weights
- ② sliding window



$$R(z) = \max(0, z)$$



*

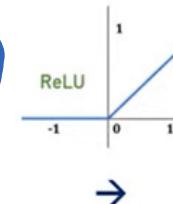
Loopy pattern filter

1	1	1
1	-1	1
1	1	1

→

-0.11	1	-0.11
-0.55	0.11	-0.33
-0.33	0.33	-0.33
-0.22	-0.11	-0.22
-0.33	-0.33	-0.33

ReLU



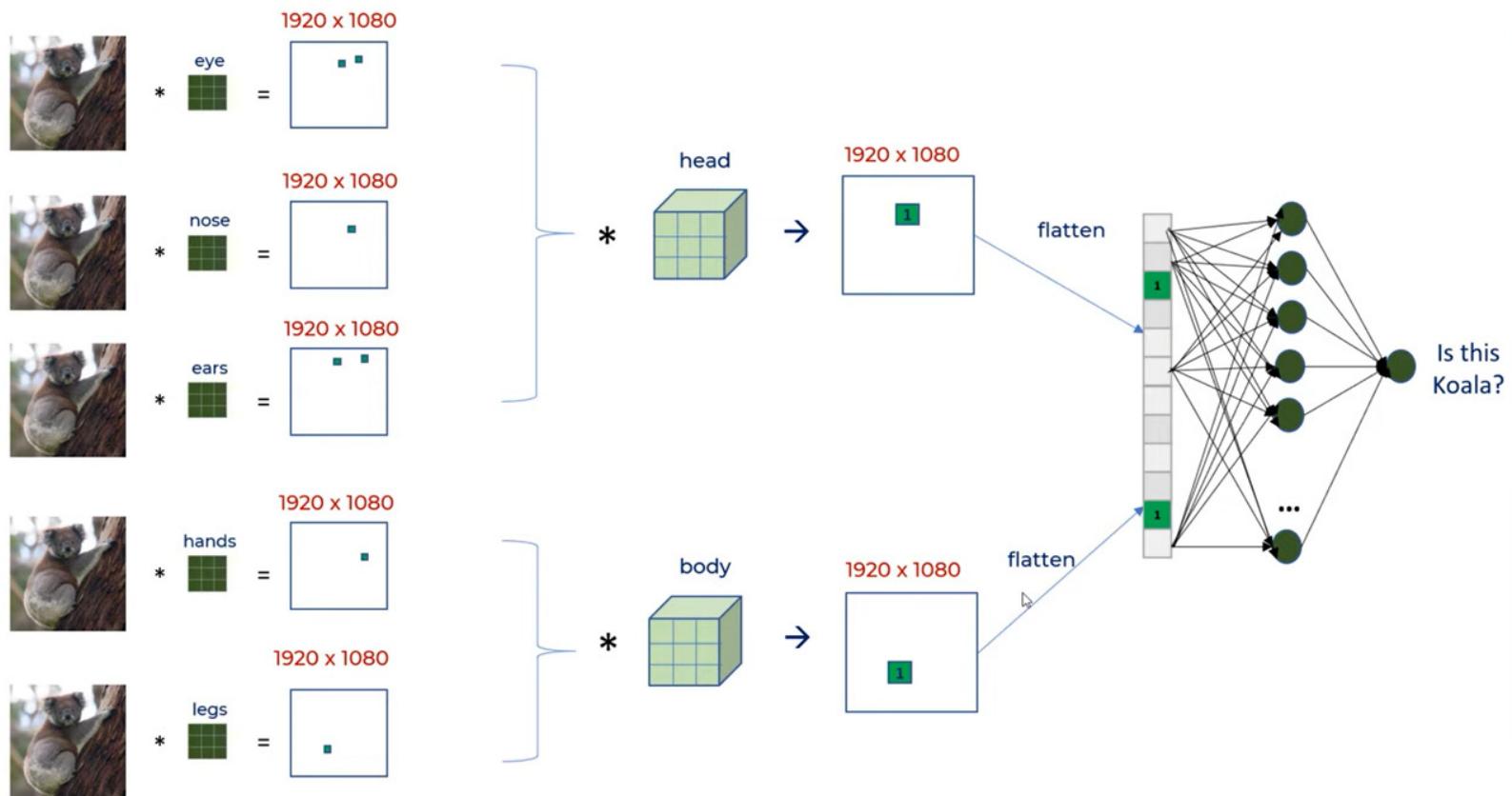
→

0	1	0
0	0.11	0
0	0.33	0
0	0	0
0	0	0

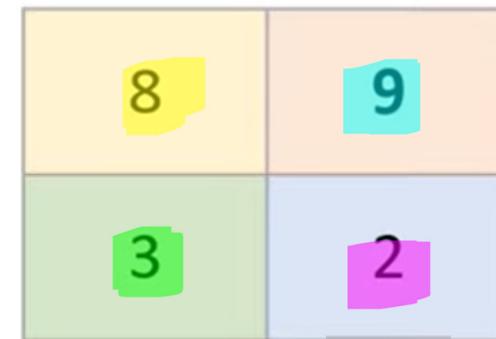
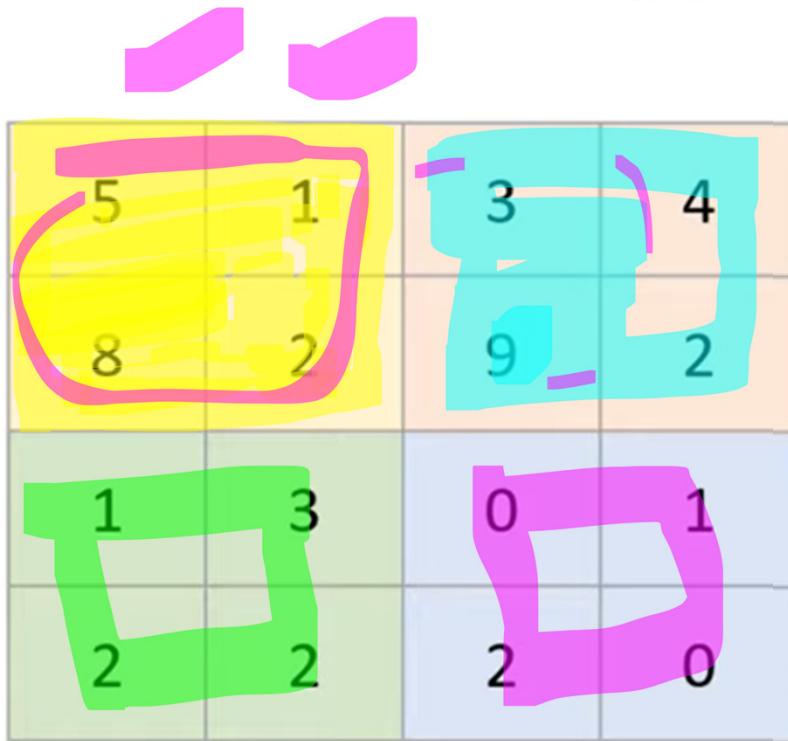
Pooling

- Reduces dimensions and computations
- Reduce overfitting as there are less parameters
- Model is tolerant towards spatial/translations, distortion
- Types of pooling
 - Mean pooling
 - Max pooling
 - Average Pooling
 - Sum pooling

Pooling

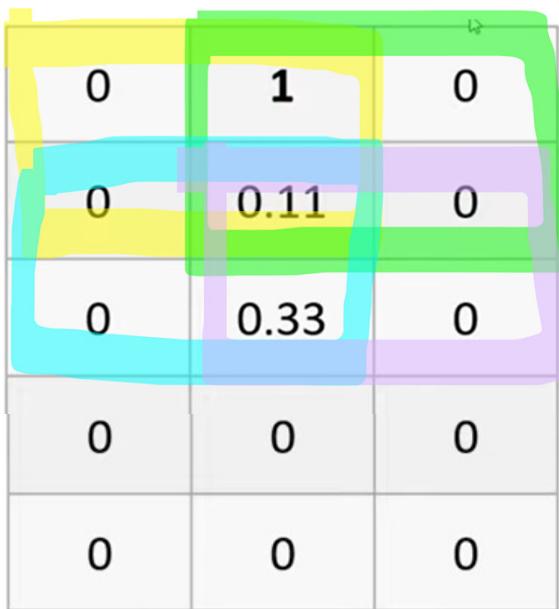


Pooling...



2 by 2 filter with stride = 2

Pooling...



1	1
0.33	0.33
0.33	0.33
0	0

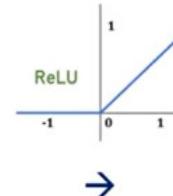
2 by 2 filter with stride = 1

Position invariant property using pooling

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

Loopy pattern filter
 $\begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$

-0.11	1	-0.11
-0.55	0.11	-0.33
-0.33	0.33	-0.33
-0.22	-0.11	-0.22
-0.33	-0.33	-0.33



0	1	0
0	0.11	0
0	0.33	0
0	0	0
0	0	0

Max pooling
 \rightarrow

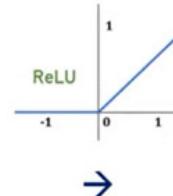
1	1
0.33	0.33
0.33	0.33
0	0

Shifted 9 at different position

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

Loopy pattern filter
 $\begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$

1	-0.11	-0.11
0.11	-0.33	0.33
0.33	-0.33	-0.33
-0.11	-0.55	-0.33
-0.55	-0.33	-0.55



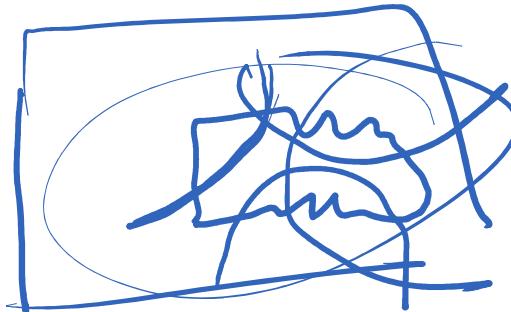
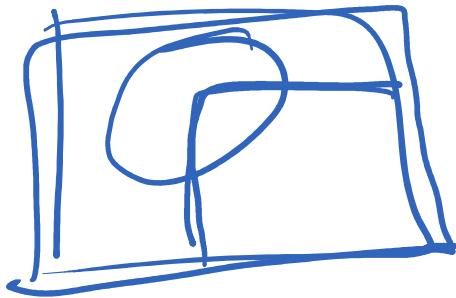
1	0	0
0.11	0	0.33
0.33	0	0
0	0	0
0	0	0

Max pooling
 \rightarrow

1	0.33
0.33	0.33
0.33	0
0	0

2. Pooling

- enables CNN to detect features like lighting and angle orientation at the images.
- spatial invariant concept
- applies non-linear downsampling on the activation map.



Max pooling

29	15	28	184
0	10	70	38
12	12	7	2
12	12	45	6

2×2

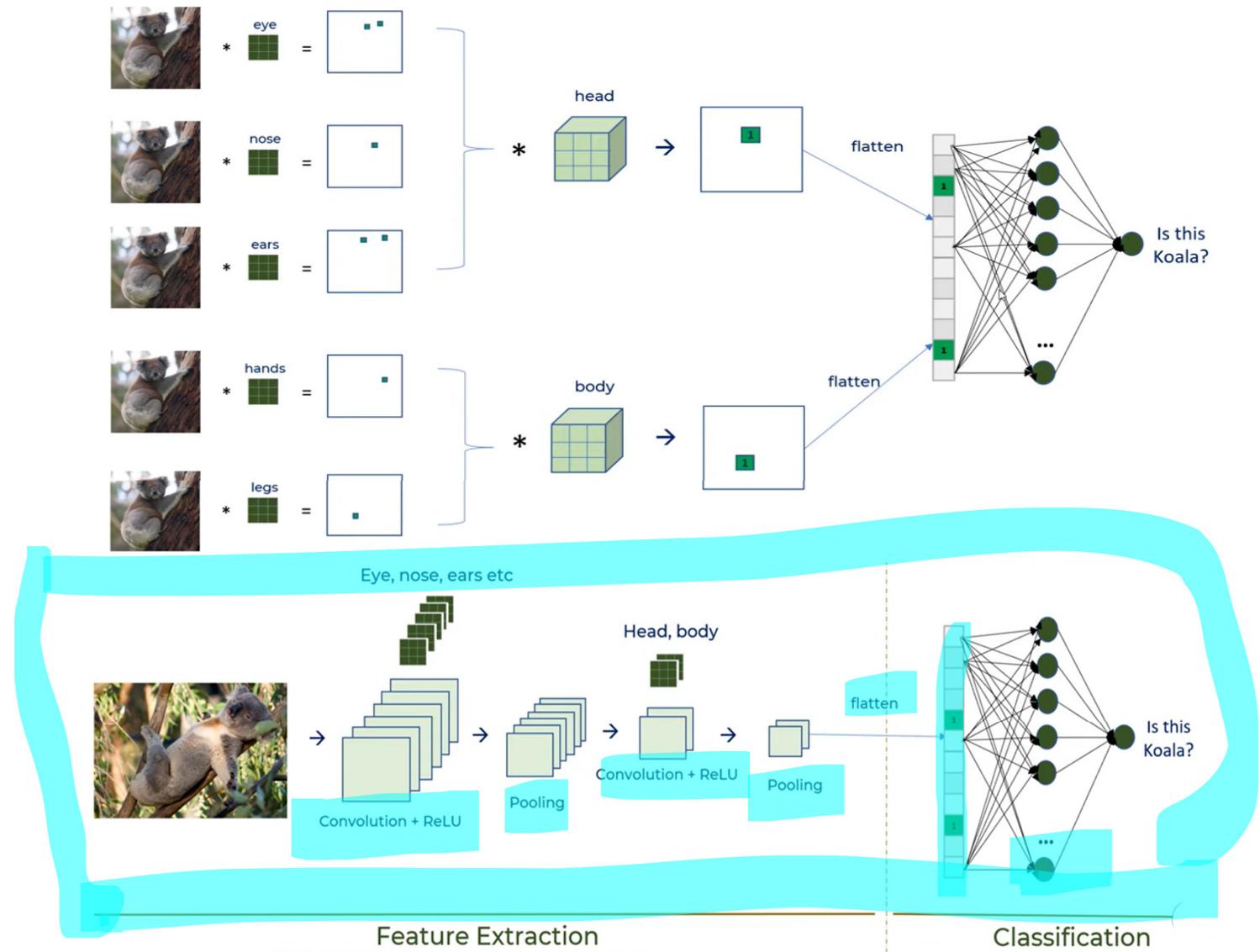
100	184
12	45

Avg - Pooling

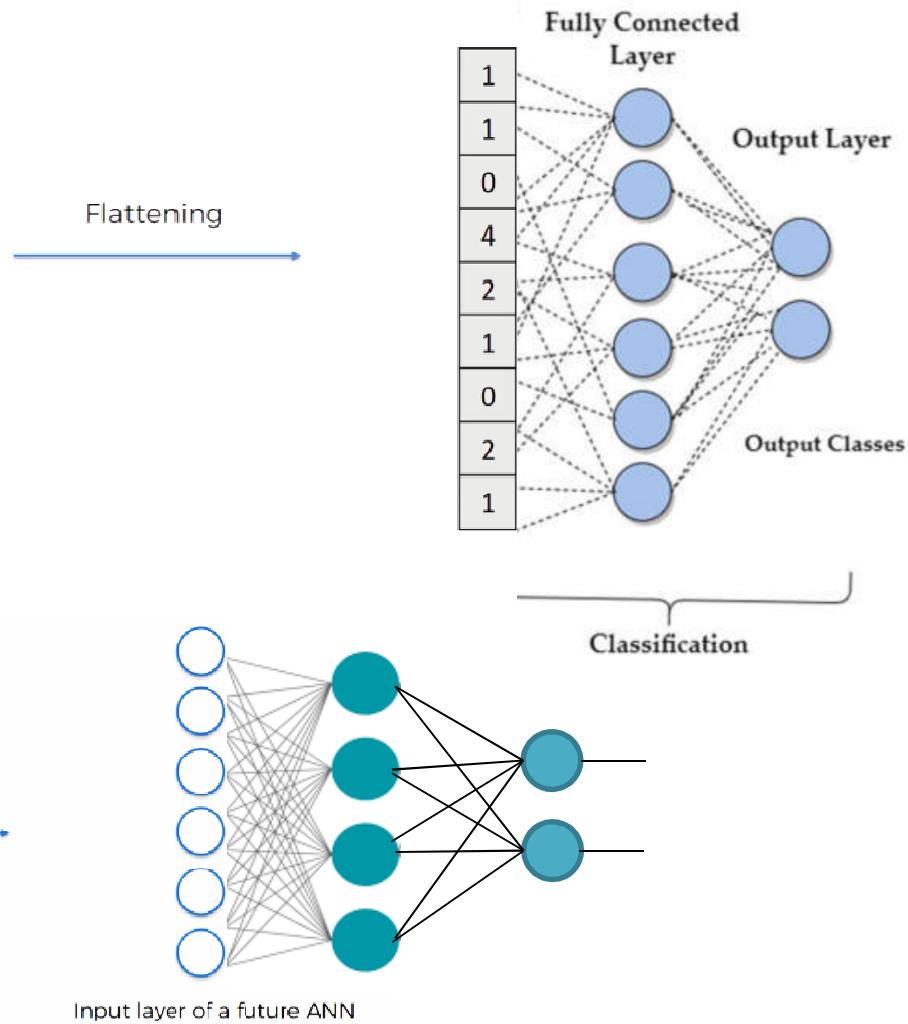
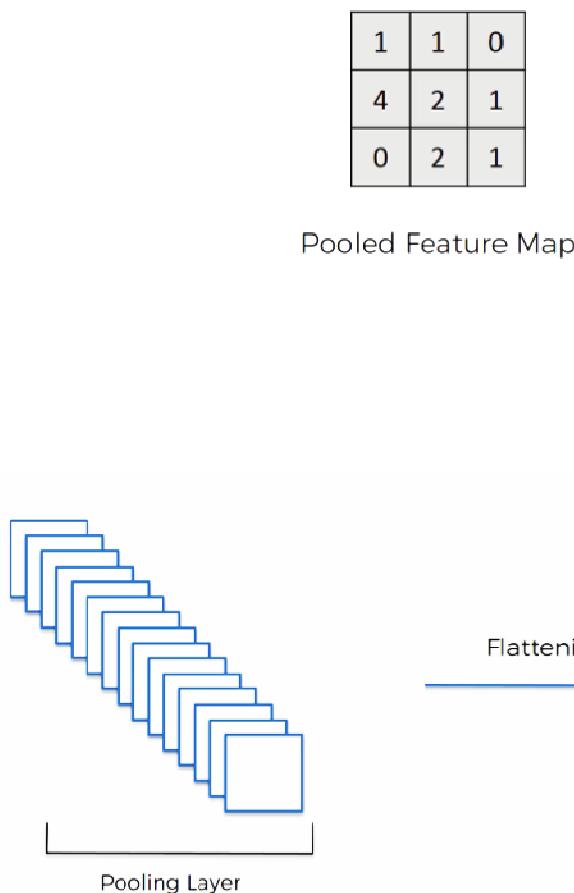
2×2

36	80
12	15

Spatial invariant



Flattening and fully connected layer



③ Flattening

1.	1	0.
4	2	1.
0	2	1

flattening

gradient
descent

→ omitted
parameters

final connection

e_{zi}

Cx

$$Cx = \sum e_{zi}$$

ARIMA

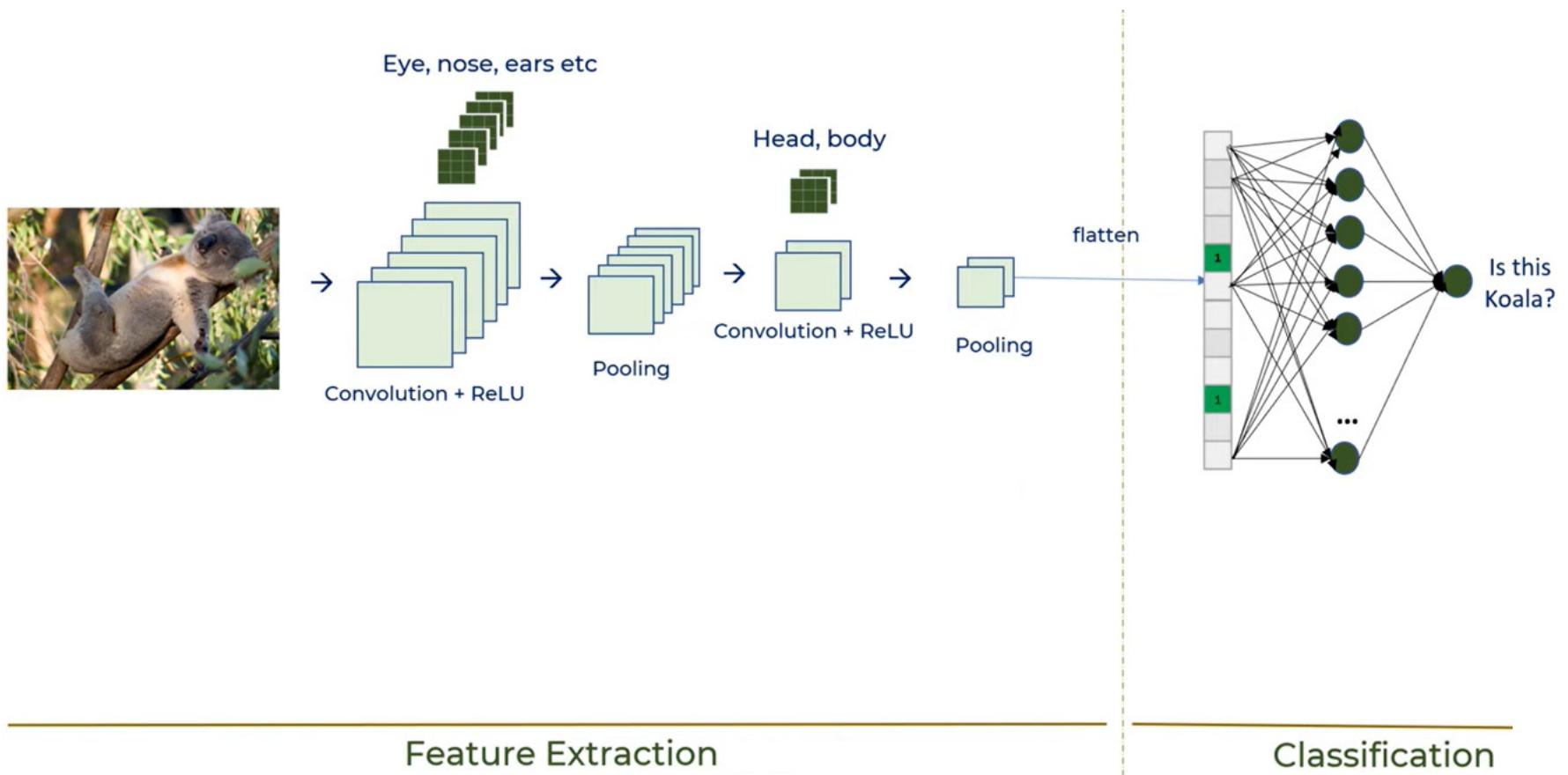
0.1 Soft. Mat
function

CSE 4553 | Machine Learning | Winter 2024

Deep Fake error / loss
2000



Complete structure of CNN



Calculations of the CNN parameters

CONV 1
Input Size ($W_1 \times H_1 \times D_1$) = 28 x 28 x 1
<ul style="list-style-type: none">• Requires four hyperparameter:<ul style="list-style-type: none">○ Number of kernels, $k = 16$○ Spatial extend of each one, $F = 5$○ Stride Size, $S = 1$○ Amount of zero padding, $P = 2$
<ul style="list-style-type: none">• Outputting volume of $W_2 \times H_2 \times D_2$<ul style="list-style-type: none">○ $W_2 = (28 - 5 + 2(2)) / 1 + 1 = 28$○ $H_2 = (28 - 5 + 2(2)) / 1 + 1 = 28$○ $D_2 = k$
Output of Conv 1 ($W_2 \times H_2 \times D_2$) = 28 x 28 x 16

POOL 1
Input Size ($W_2 \times H_2 \times D_2$) = 28 x 28 x 16
<ul style="list-style-type: none">• Requires two hyperparameter:<ul style="list-style-type: none">○ Spatial extend of each one, $F = 2$○ Stride Size, $S = 2$
<ul style="list-style-type: none">• Outputting volume of $W_3 \times H_3 \times D_2$<ul style="list-style-type: none">○ $W_3 = (28 - 2) / 2 + 1 = 14$○ $H_3 = (28 - 2) / 2 + 1 = 14$
Output of Pool 1 ($W_3 \times H_3 \times D_2$) = 14 x 14 x 16

CONV 2
Input Size ($W_3 \times H_3 \times D_2$) = 14 x 14 x 16
<ul style="list-style-type: none">• Requires four hyperparameter:<ul style="list-style-type: none">○ Number of kernels, $k = 32$○ Spatial extend of each one, $F = 5$○ Stride Size, $S = 1$○ Amount of zero padding, $P = 2$
<ul style="list-style-type: none">• Outputting volume of $W_4 \times H_4 \times D_3$<ul style="list-style-type: none">○ $W_4 = (14 - 5 + 2(2)) / 1 + 1 = 14$○ $H_4 = (14 - 5 + 2(2)) / 1 + 1 = 14$○ $D_3 = k$
Output of Conv 2 ($W_4 \times H_4 \times D_3$) = 14 x 14 x 32

POOL 2
Input Size ($W_4 \times H_4 \times D_3$) = 14 x 14 x 32
<ul style="list-style-type: none">• Requires two hyperparameter:<ul style="list-style-type: none">○ Spatial extend of each one, $F = 2$○ Stride Size, $S = 2$
<ul style="list-style-type: none">• Outputting volume of $W_5 \times H_5 \times D_3$<ul style="list-style-type: none">○ $W_5 = (14 - 2) / 2 + 1 = 7$○ $H_5 = (14 - 2) / 2 + 1 = 7$
Output of Pool 2 ($W_5 \times H_5 \times D_3$) = 7 x 7 x 32

FC Layer
Input Size ($W_5 \times H_5 \times D_3$) = 7 x 7 x 32
Output Size (Number of Classes) = 10

Convolution

- Connections sparsity reduces overfitting
- Conv + Pooling gives location invariant feature detection
- Parameter sharing

ReLU

- Introduces nonlinearity
- Speeds up training, faster to compute

Pooling

- Reduces dimensions and computation
- Reduces overfitting
- Makes the model tolerant towards small distortion and variations