《数据挖掘技术》

# ⭐ **CHX5 Neural Network**

➡️ Create by *Wang JingHui*

➡️ Version: *4.0*

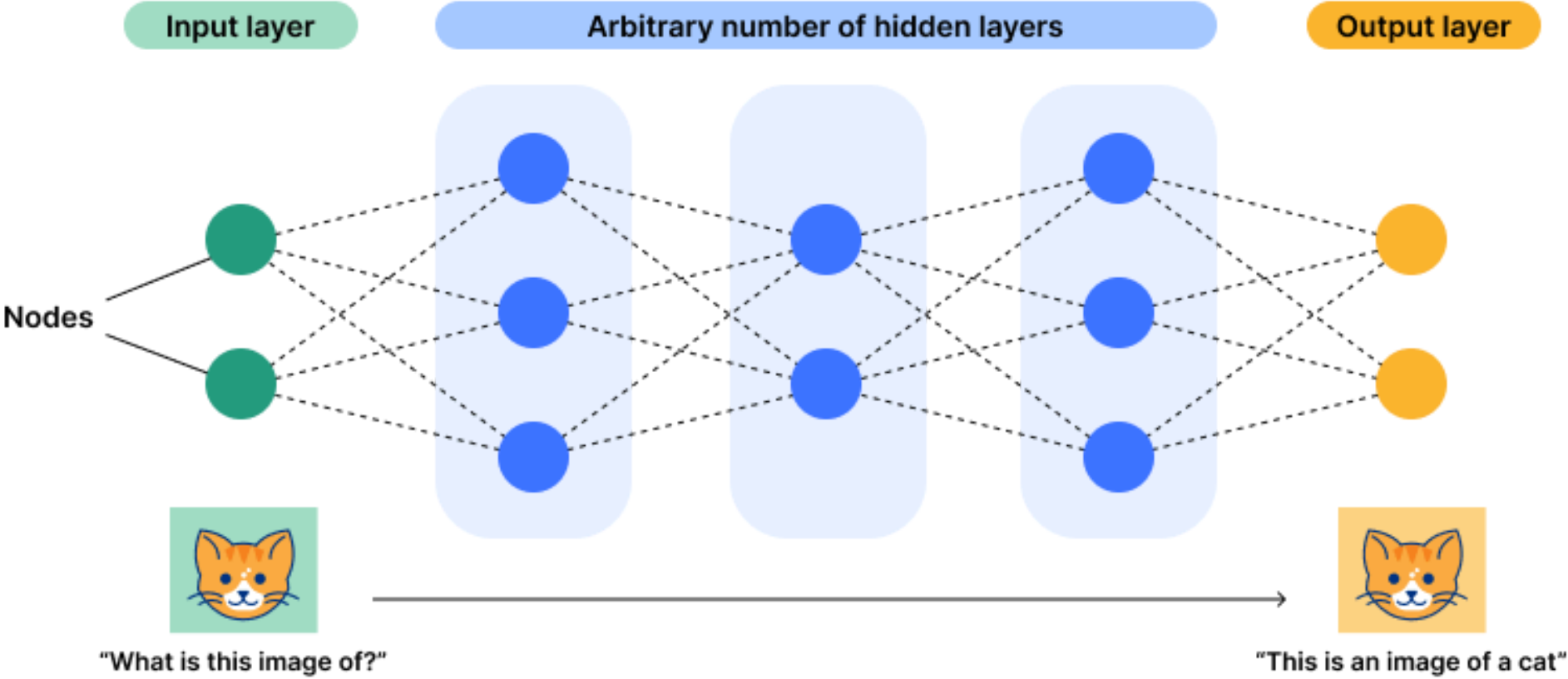# Contents

# ⚛ What is a neural network?

- A neural network is a machine learning program, or model, that makes decisions in a manner similar to the human brain, by using processes that mimic the way biological neurons work together to identify phenomena, weigh options and arrive at conclusions.

- Every neural network consists of layers of nodes, or artificial neurons—an input layer, one or more hidden layers, and an output layer. Each node connects to others, and has its own associated weight and threshold. If the output of any individual node is above the specified threshold value, that node is activated, sending data to the next layer of the network. Otherwise, no data is passed along to the next layer of the network.

- https://www.geeksforgeeks.org/backpropagation-in-machine-learning/

# Neural network

Input layer

Arbitrary number of hidden layers

Output layer

Nodes

"What is this image of?"

"This is an image of a cat"
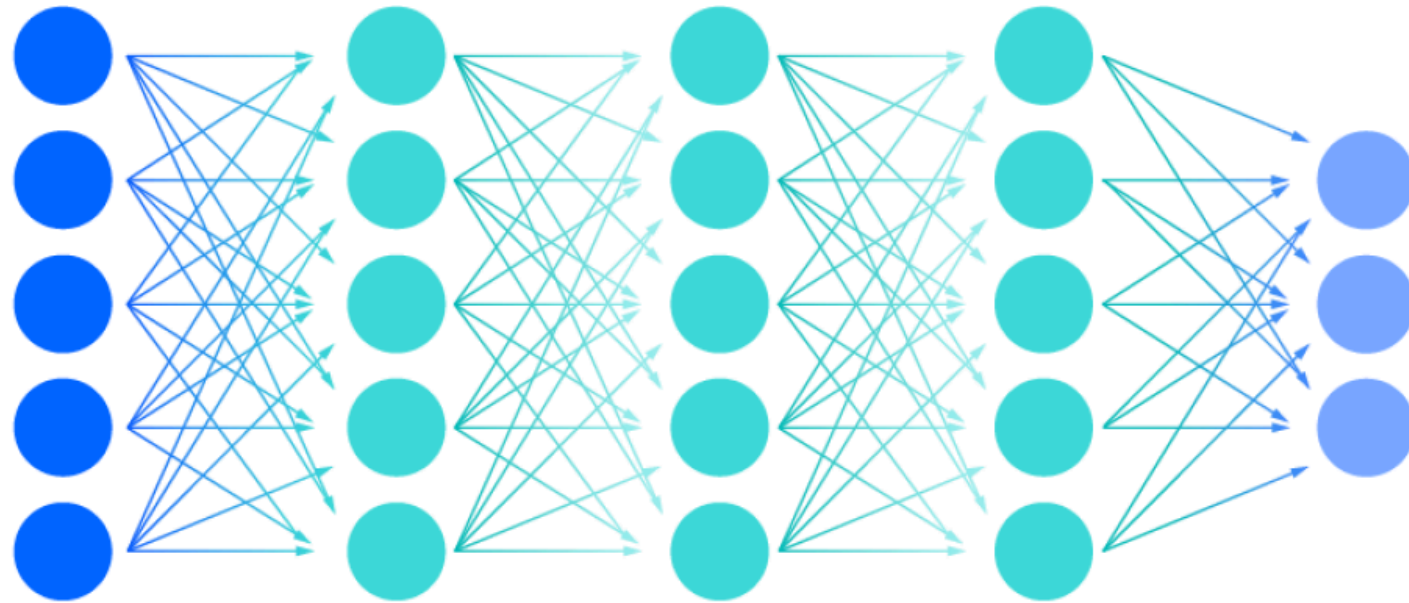
# Deep neural network

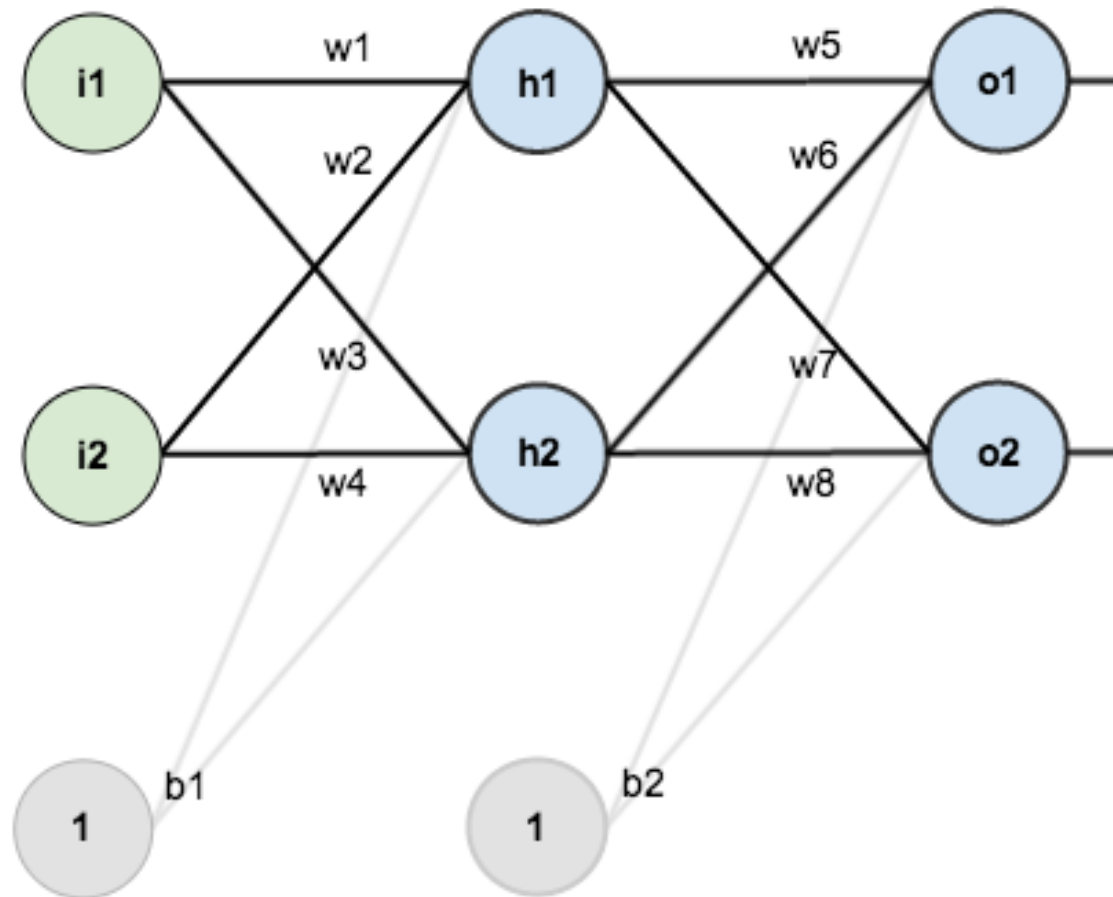Input layer          Multiple hidden layer          Output layer

# ⚛️ What is backpropagation?

- In machine learning, backpropagation is an effective algorithm used to train artificial neural networks, especially in feed-forward neural networks.

- Backpropagation is an iterative algorithm, that helps to minimize the cost function by determining which weights and biases should be adjusted. During every epoch, the model learns by adapting the weights and biases to minimize the loss by moving down toward the gradient of the error. Thus, it involves the two most popular optimization algorithms, such as gradient descent or stochastic gradient descent.

- Computing the gradient in the backpropagation algorithm helps to minimize the cost function and it can be implemented by using the mathematical rule called chain rule from calculus to navigate through complex layers of the neural network.
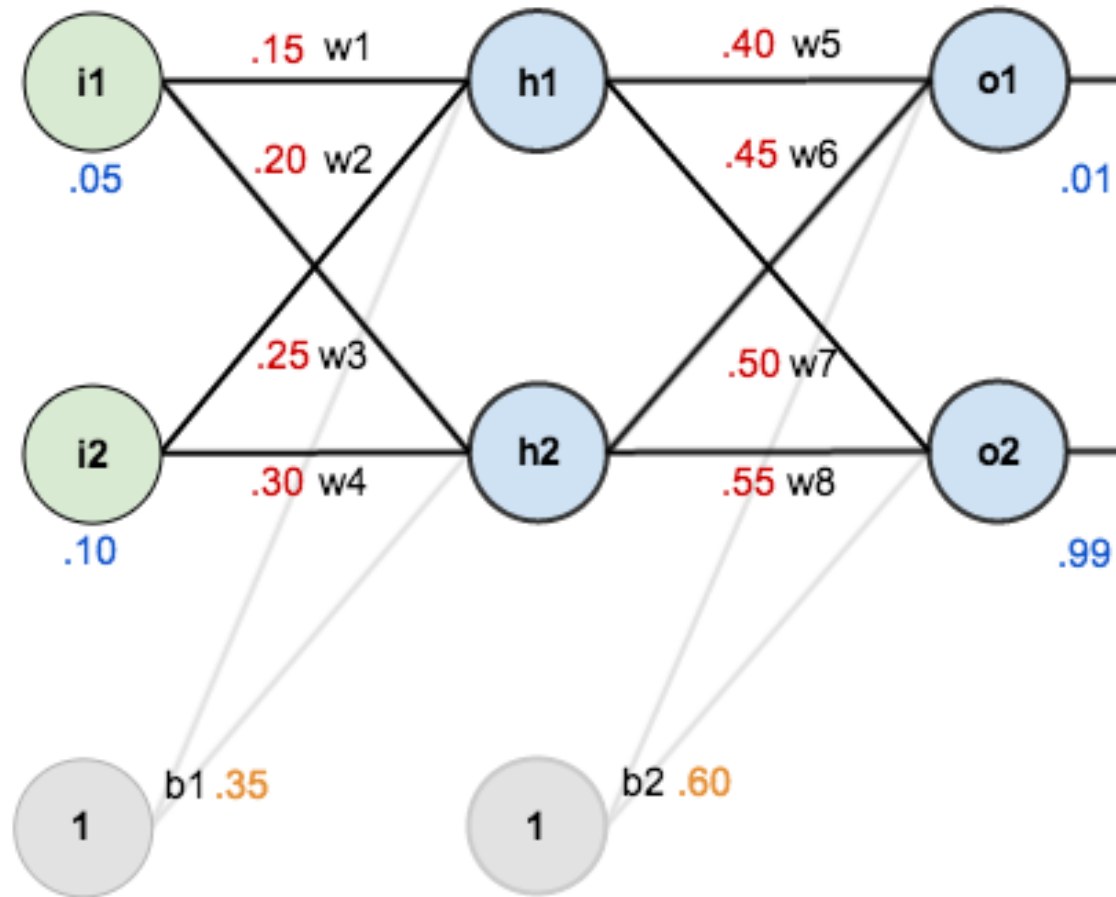
- Backpropagation is a common method for training a neural network. There is no shortage of papers online that attempt to explain how backpropagation works, but few that include an example with actual numbers. This post is my attempt to explain how it works with a concrete example that folks can compare their own calculations to in order to ensure they understand backpropagation correctly.

- https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/

For this tutorial, we're going to use a neural network with two inputs, two hidden neurons, two output neurons. Additionally, the hidden and output neurons will include a bias.

Here's the basic structure:

In order to have some numbers to work with, here are the initial weights, the biases, and training inputs/outputs:

The goal of backpropagation is to optimize the weights so that the neural network can learn how to correctly map arbitrary inputs to outputs.

For the rest of this tutorial we're going to work with a single training set: given inputs 0.05 and 0.10, we want the neural network to output 0.01 and 0.99.

# The Forward Pass

To begin, lets see what the neural network currently predicts given the weights and biases above and inputs of 0.05 and 0.10. To do this we'll feed those inputs forward though the network.

We figure out the total net input to each hidden layer neuron, squash the total net input using an activation function (here we use the logistic function), then repeat the process with the output layer neurons.

Total net input is also referred to as just net input by some sources.
Here's how we calculate the total net input for $h_1$:

$$net_{h1} = w_1 * i_1 + w_2 * i_2 + b_1 * 1$$
$$net_{h1} = 0.15 * 0.05 + 0.2 * 0.1 + 0.35 * 1 = 0.3775$$

We then squash it using the logistic function to get the output of $h_1$:

$$out_{h1} = \frac{1}{1 + e^{-net_{h1}}} = \frac{1}{1 + e^{-0.3775}} = 0.593269992$$

Carrying out the same process for $h_2$ we get:

$$out_{h2} = 0.596884378$$

We repeat this process for the output layer neurons, using the output from the hidden layer neurons as inputs.

Here's the output for $o_1$:

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$
$$net_{o1} = 0.4 * 0.593269992 + 0.45 * 0.596884378 + 0.6 * 1 = 1.105905967$$
$$out_{o1} = \frac{1}{1 + e^{-net_{o1}}} = \frac{1}{1 + e^{-1.105905967}} = 0.75136507$$

And carrying out the same process for $o_2$ we get:

$$out_{o2} = 0.772928465$$

## Calculating the Total Error

We can now calculate the error for each output neuron using the squared error function and sum them to get the total error:

$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

The $\frac{1}{2}$ is included so that exponent is cancelled when we differentiate later on. The result is eventually multiplied by a learning rate anyway so it doesn't matter that we introduce a constant here [1].

For example, the target output for $o_1$ is 0.01 but the neural network output 0.75136507, therefore its error is:

$$E_{o1} = \frac{1}{2}\left(target_{o1} - out_{o1}\right)^2 = \frac{1}{2}(0.01 - 0.75136507)^2 = 0.274811083$$

Repeating this process for $o_2$ (remembering that the target is 0.99) we get:

$$E_{o2} = 0.023560026$$

The total error for the neural network is the sum of these errors:

$$E_{total} = E_{o1} + E_{o2} = 0.274811083 + 0.023560026 = 0.298371109$$

# The Backwards Pass

Our goal with backpropagation is to update each of the weights in the network so that they cause the actual output to be closer the target output, thereby minimizing the error for each output neuron and the network as a whole.
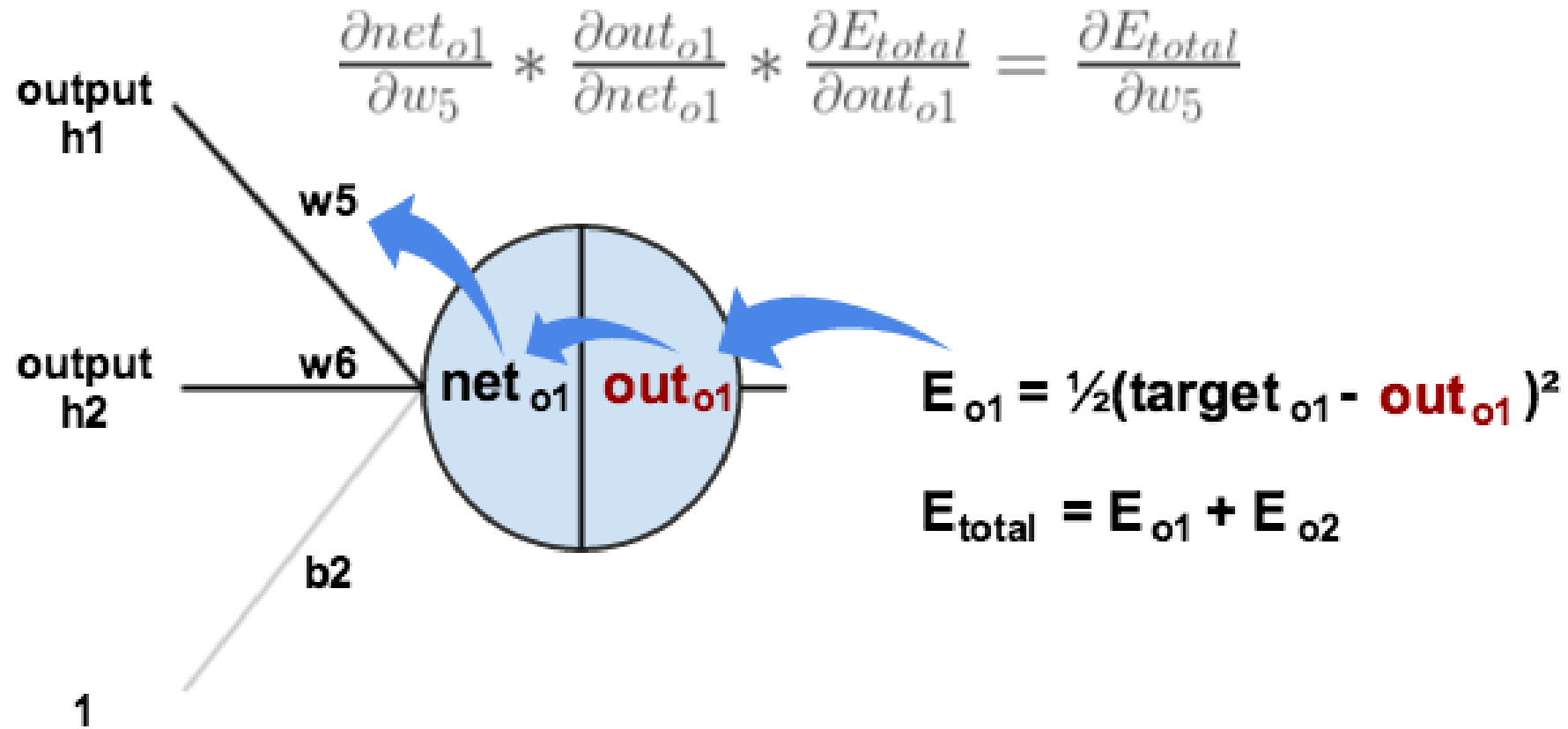
## Output Layer

Consider $w_5$. We want to know how much a change in $w_5$ affects the total error, aka $\frac{\partial E_{total}}{\partial w_5}$.

- $\frac{\partial E_{total}}{\partial w_5}$ is read as "the partial derivative of $E_{total}$ with respect to $w_5$". You can also say "the gradient with respect to $w_5$".

- By applying the chain rule we know that:

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

Visually, here's what we're doing:

$$\frac{\partial net_{o1}}{\partial w_5} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial E_{total}}{\partial out_{o1}} = \frac{\partial E_{total}}{\partial w_5}$$



$$E_{o1} = \frac{1}{2}(target_{o1} - out_{o1})^2$$

$$E_{total} = E_{o1} + E_{o2}$$

We need to figure out each piece in this equation.

First, how much does the total error change with respect to the output?

$$E_{total} = \frac{1}{2}\left(target_{o1} - out_{o1}\right)^2 + \frac{1}{2}\left(target_{o2} - out_{o2}\right)^2$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = 2 * \frac{1}{2}\left(target_{o1} - out_{o1}\right)^{2-1} * -1 + 0$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = -\left(target_{o1} - out_{o1}\right) = -\left(0.01 - 0.75136507\right) = 0.74136507$$

$-(target - out)$ is sometimes expressed as $out - target$

When we take the partial derivative of the total error with respect to $out_{o1}$, the quantity $\frac{1}{2}\left(target_{o2} - out_{o2}\right)^2$ becomes zero because $out_{o1}$ does not affect it which means we're taking the derivative of a constant which is zero.

Next, how much does the output of $o_1$ change with respect to its total net input?

The partial derivative of the logistic function is the output multiplied by 1 minus the output:

$$out_{o1} = \frac{1}{1 + e^{-net_{o1}}}$$

$$\frac{\partial out_{o1}}{\partial net_{o1}} = out_{o1}(1 - out_{o1}) = 0.75136507(1 - 0.75136507) = 0.186815602$$

Finally, how much does the total net input of o1 change with respect to w_5?

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\frac{\partial net_{o1}}{\partial w_5} = 1 * out_{h1} * w_5^{(1-1)} + 0 + 0 = out_{h1} = 0.593269992$$

Putting it all together:

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

$$\frac{\partial E_{total}}{\partial w_5} = 0.74136507 * 0.186815602 * 0.593269992 = 0.082167041$$

You'll often see this calculation combined in the form of the delta rule:

$$\frac{\partial E_{total}}{\partial w_5} = -(target_{o1} - out_{o1}) * out_{o1}(1 - out_{o1}) * out_{h1}$$

Alternatively, we have $\frac{\partial E_{total}}{\partial out_{o1}}$ and $\frac{\partial out_{o1}}{\partial net_{o1}}$ which can be written as $\frac{\partial E_{total}}{\partial net_{o1}}$, aka $\delta_{o1}$ (the Greek letter delta) aka the node delta.

We can use this to rewrite the calculation above:

$$\delta_{o1} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} = \frac{\partial E_{total}}{\partial net_{o1}}$$

$$\delta_{o1} = -(target_{o1} - out_{o1}) * out_{o1}(1 - out_{o1})$$

Therefore:

$$\frac{\partial E_{total}}{\partial w_5} = \delta_{o1} out_{h1}$$

Some sources extract the negative sign from $\delta$ so it would be written as:

$$\frac{\partial E_{total}}{\partial w_5} = -\delta_{o1} out_{h1}$$

To decrease the error, we then subtract this value from the current weight (optionally multiplied by some learning rate, eta, which we'll set to 0.5):

$$w_5^+ = w_5 - \eta * \frac{\partial E_{total}}{\partial w_5} = 0.4 - 0.5 * 0.082167041 = 0.35891648$$

Some sources use $\alpha$ (alpha) to represent the learning rate, others use $\eta$ (eta), and others even use $\epsilon$ (epsilon).

We can repeat this process to get the new weights $w_6$, $w_7$, and $w_8$:

$$w_6^+ = 0.408666186$$
$$w_7^+ = 0.511301270$$
$$w_8^+ = 0.561370121$$

We perform the actual updates in the neural network after we have the new weights leading into the hidden layer neurons (ie, we use the original weights, not the updated weights, when we continue the backpropagation algorithm below).

# Hidden Layer

Next, we'll continue the backwards pass by calculating new values for $w_1$, $w_2$, $w_3$, and $w_4$.
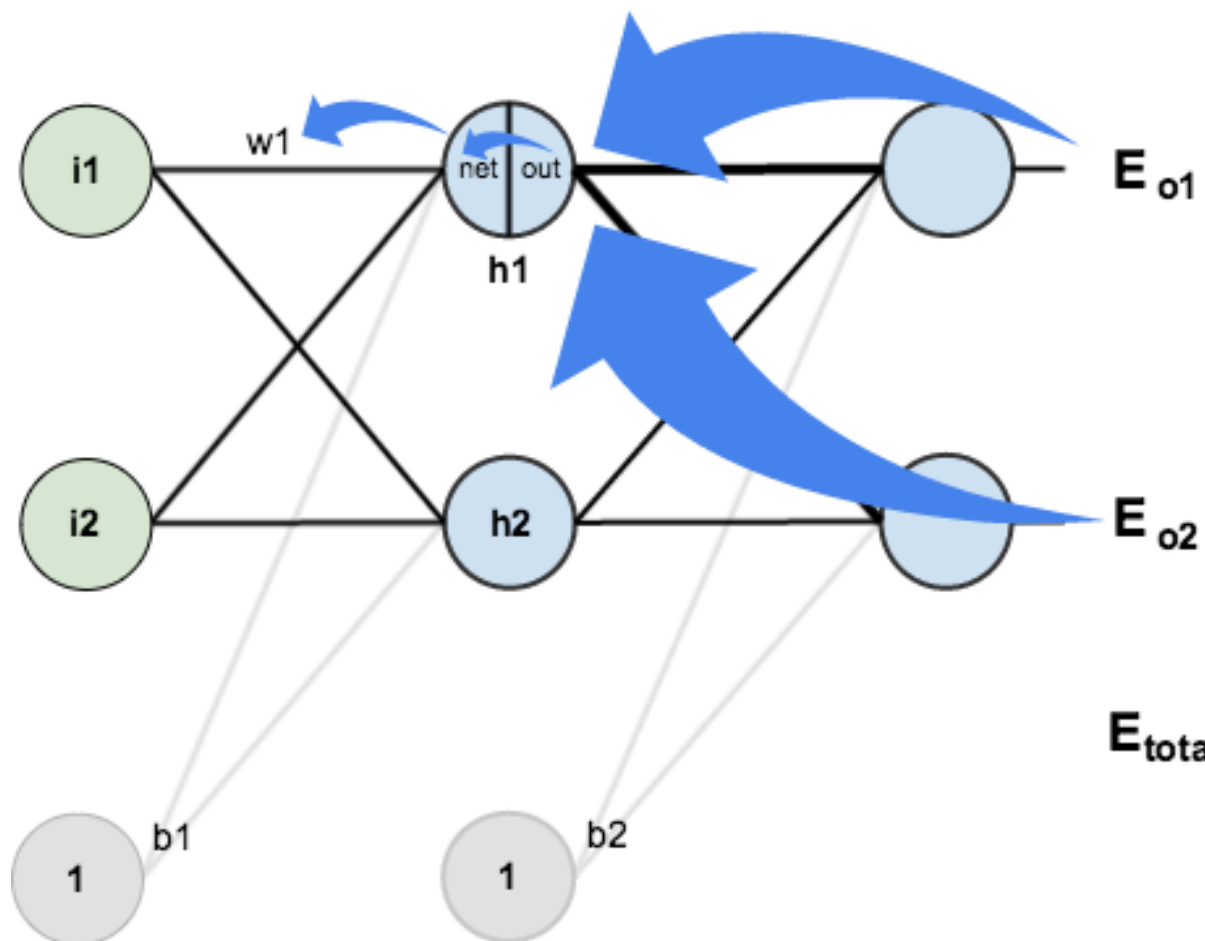
Big picture, here's what we need to figure out:

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

$E_{total} = E_{o1} + E_{o2}$

We're going to use a similar process as we did for the output layer, but slightly different to account for the fact that the output of each hidden layer neuron contributes to the output (and therefore error) of multiple output neurons. We know that $out_{h1}$ affects both $out_{o1}$ and $out_{o2}$ therefore the $\frac{\partial E_{total}}{\partial out_{h1}}$ needs to take into consideration its effect on the both output neurons:

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

Starting with $\frac{\partial E_{o1}}{\partial out_{h1}}$:

$$\frac{\partial E_{o1}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial out_{h1}}$$

We can calculate $\frac{\partial E_{o1}}{\partial net_{o1}}$ using values we calculated earlier:

$$\frac{\partial E_{o1}}{\partial net_{o1}} = \frac{\partial E_{o1}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} = 0.74136507 * 0.186815602 = 0.138498562$$

And $\frac{\partial net_{o1}}{\partial out_{h1}}$ is equal to $w_5$:

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\frac{\partial net_{o1}}{\partial out_{h1}} = w_5 = 0.40$$

Plugging them in:

$$\frac{\partial E_{o1}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial out_{h1}} = 0.138498562 * 0.40 = 0.055399425$$

Following the same process for $\frac{\partial E_{o2}}{\partial out_{h1}}$, we get:

$$\frac{\partial E_{o2}}{\partial out_{h1}} = -0.019049119$$

Therefore:

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}} = 0.055399425 + -0.019049119 = 0.036350306$$

Now that we have $\frac{\partial E_{total}}{\partial out_{h1}}$, we need to figure out $\frac{\partial out_{h1}}{\partial net_{h1}}$ and then $\frac{\partial net_{h1}}{\partial w}$ for each weight:

$$out_{h1} = \frac{1}{1 + e^{-net_{h1}}}$$

$$\frac{\partial out_{h1}}{\partial net_{h1}} = out_{h1}(1 - out_{h1}) = 0.59326999(1 - 0.59326999) = 0.241300709$$

We calculate the partial derivative of the total net input to h_1 with respect to w_1 the same as we did for the output neuron:

$$net_{h1} = w_1 * i_1 + w_3 * i_2 + b_1 * 1$$

$$\frac{\partial net_{h1}}{\partial w_1} = i_1 = 0.05$$

Putting it all together:

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial w_1} = 0.036350306 * 0.241300709 * 0.05 = 0.000438568$$

You might also see this written as:

$$\frac{\partial E_{total}}{\partial w_1} = (\sum_o \frac{\partial E_{total}}{\partial out_o} * \frac{\partial out_o}{\partial net_o} * \frac{\partial net_o}{\partial out_{h1}}) * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial w_1} = (\sum_o \delta_o * w_{ho}) * out_{h1}(1 - out_{h1}) * i_1$$

$$\frac{\partial E_{total}}{\partial w_1} = \delta_{h1} i_1$$

We can now update $w_1$:

$$w_1^+ = w_1 - \eta * \frac{\partial E_{total}}{\partial w_1} = 0.15 - 0.5 * 0.000438568 = 0.149780716$$

Repeating this for $w_2$, $w_3$, and $w_4$

$$w_2^+ = 0.19956143$$
$$w_3^+ = 0.24975114$$
$$w_4^+ = 0.29950229$$

Finally, we've updated all of our weights! When we fed forward the 0.05 and 0.1 inputs originally, the error on the network was 0.298371109. After this first round of backpropagation, the total error is now down to 0.291027924. It might not seem like much, but after repeating this process 10,000 times, for example, the error plummets to 0.0000351085. At this point, when we feed forward 0.05 and 0.1, the two outputs neurons generate 0.015912196 (vs 0.01 target) and 0.984065734 (vs 0.99 target).

# ⚛️ What is a Convolutional Neural Network (CNN)?

A Convolutional Neural Network (CNN), also known as ConvNet, is a specialized type of deep learning algorithm mainly designed for tasks that necessitate object recognition, including image classification, detection, and segmentation. CNNs are employed in a variety of practical scenarios, such as autonomous vehicles, security camera systems, and others.

- https://www.datacamp.com/tutorial/introduction-to-convolutional-neural-networks-cnns

Convolutional neural networks were inspired by the layered architecture of the human visual cortex, and below are some key similarities and differences:

- Hierarchical architecture: Both CNNs and the visual cortex have a hierarchical structure, with simple features extracted in early layers and more complex features built up in deeper layers. This allows increasingly sophisticated representations of visual inputs.

- Local connectivity: Neurons in the visual cortex only connect to a local region of the input, not the entire visual field. Similarly, the neurons in a CNN layer are only connected to a local region of the input volume through the convolution operation. This local connectivity enables efficiency.
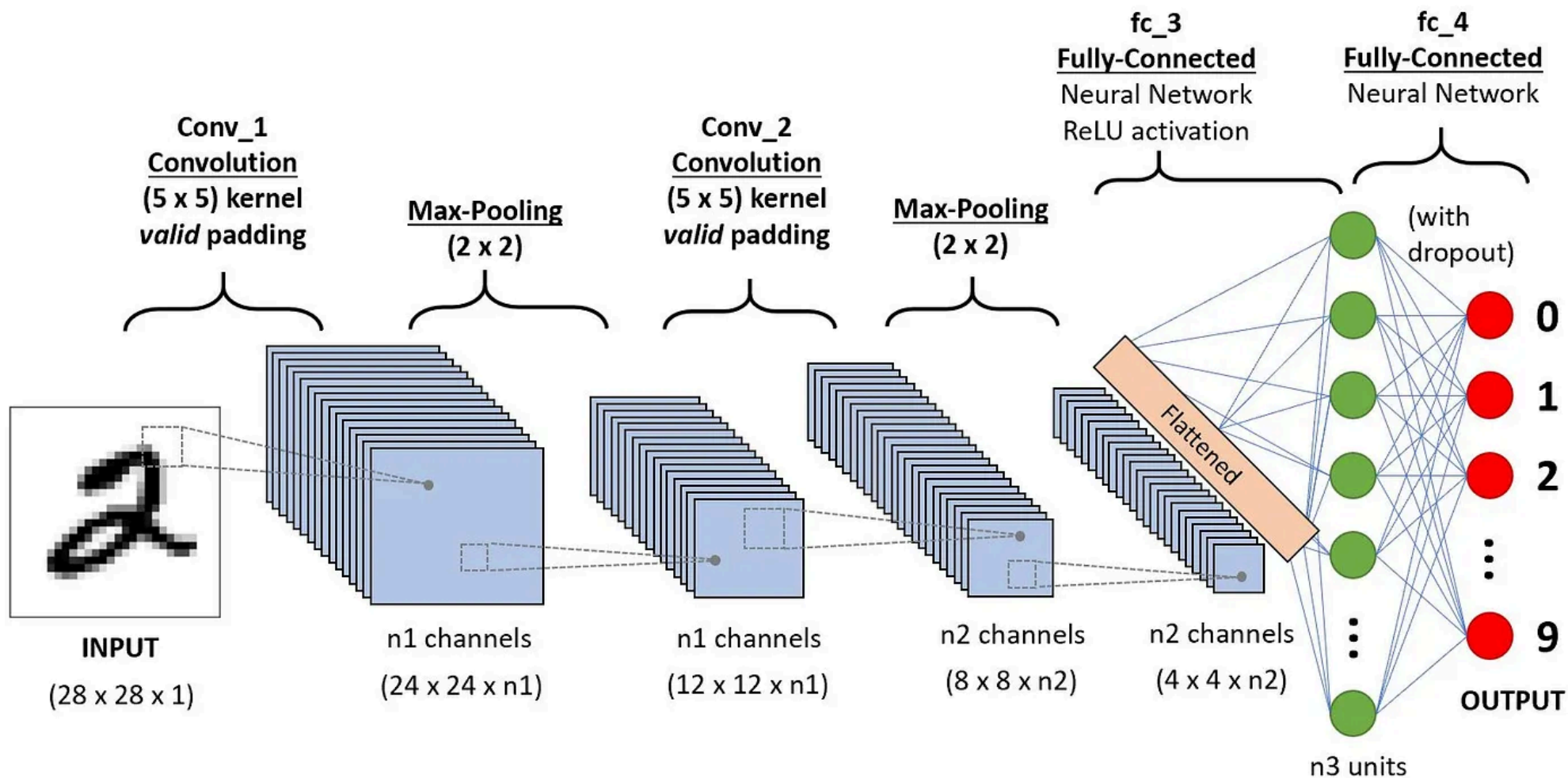
- Translation invariance: Visual cortex neurons can detect features regardless of their location in the visual field. Pooling layers in a CNN provide a degree of translation invariance by summarizing local features.

- Multiple feature maps: At each stage of visual processing, there are many different feature maps extracted. CNNs mimic this through multiple filter maps in each convolution layer.

- Non-linearity: Neurons in the visual cortex exhibit non-linear response properties. CNNs achieve non-linearity through activation functions like ReLU applied after each convolution.
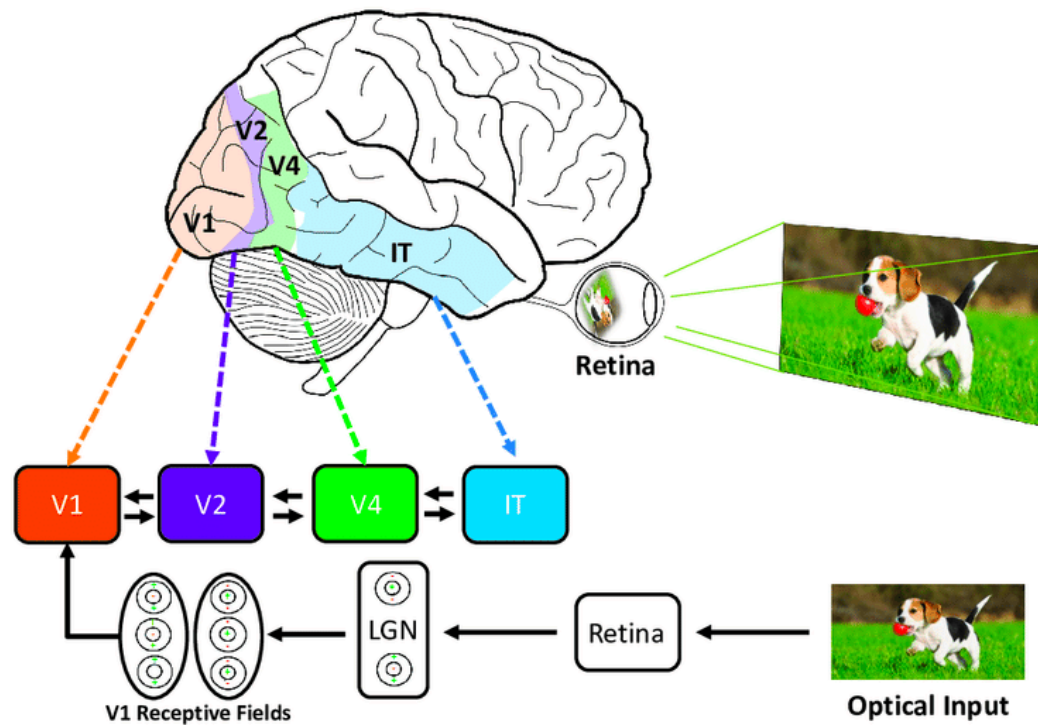
**Key Components of a CNN**

The convolutional neural network is made of four main parts.

They help the CNNs mimic how the human brain operates to recognize patterns and features in images:
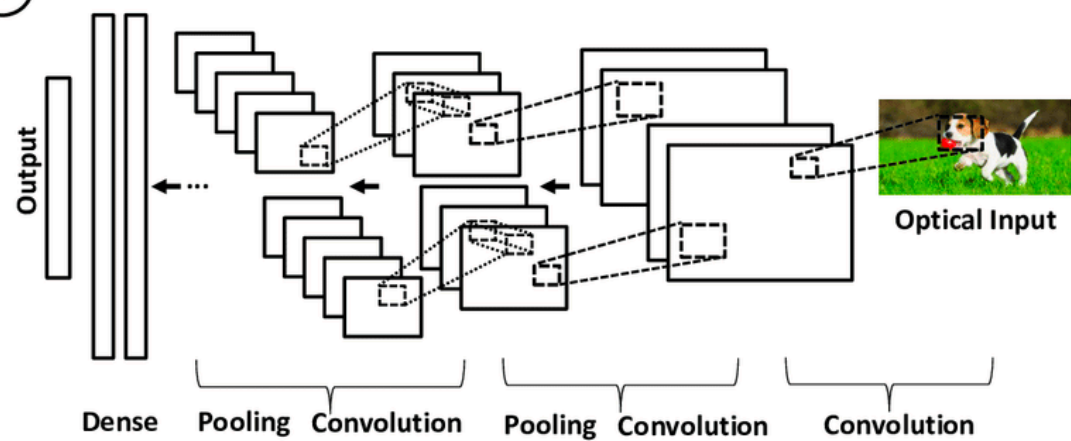
- Convolutional layers

- Rectified Linear Unit (ReLU for short)

- Pooling layers

- Fully connected layers

**Conv_1**
Convolution
(5 x 5) kernel
*valid* padding

**Max-Pooling**
(2 x 2)

**Conv_2**
Convolution
(5 x 5) kernel
*valid* padding

**Max-Pooling**
(2 x 2)

**fc_3**
**Fully-Connected**
Neural Network
ReLU activation

**fc_4**
**Fully-Connected**
Neural Network

(with dropout)

Flattened

INPUT

(28 x 28 x 1)

n1 channels

(24 x 24 x n1)

n1 channels

(12 x 12 x n1)

n2 channels

(8 x 8 x n2)

n2 channels

(4 x 4 x n2)

n3 units

0
1
2

9

OUTPUT

**a**

V2
V4
V1
IT

Retina

V1 | V2 | V4 | IT

V1 Receptive Fields

LGN

Retina

Optical Input

**b**

Output

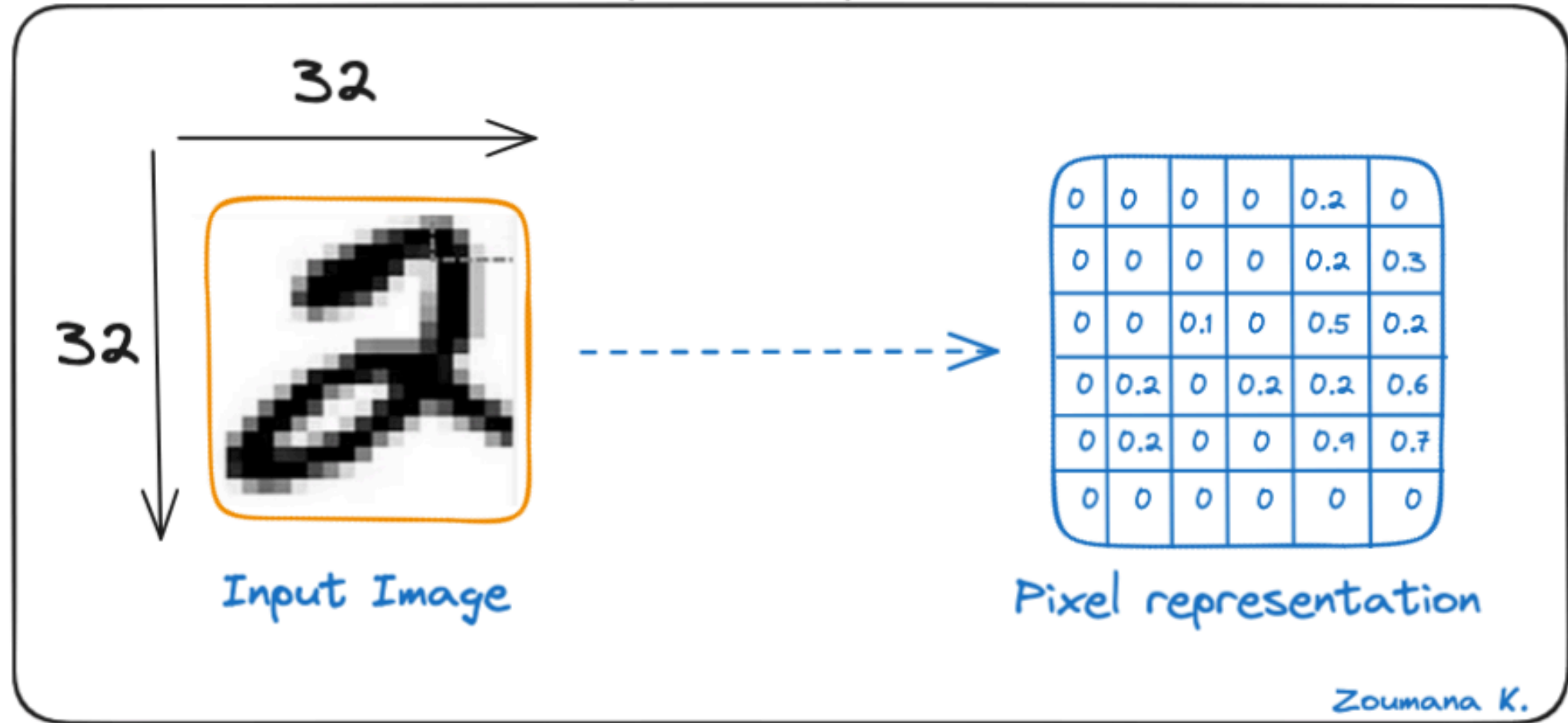Dense    Pooling    Convolution    Pooling    Convolution    Convolution
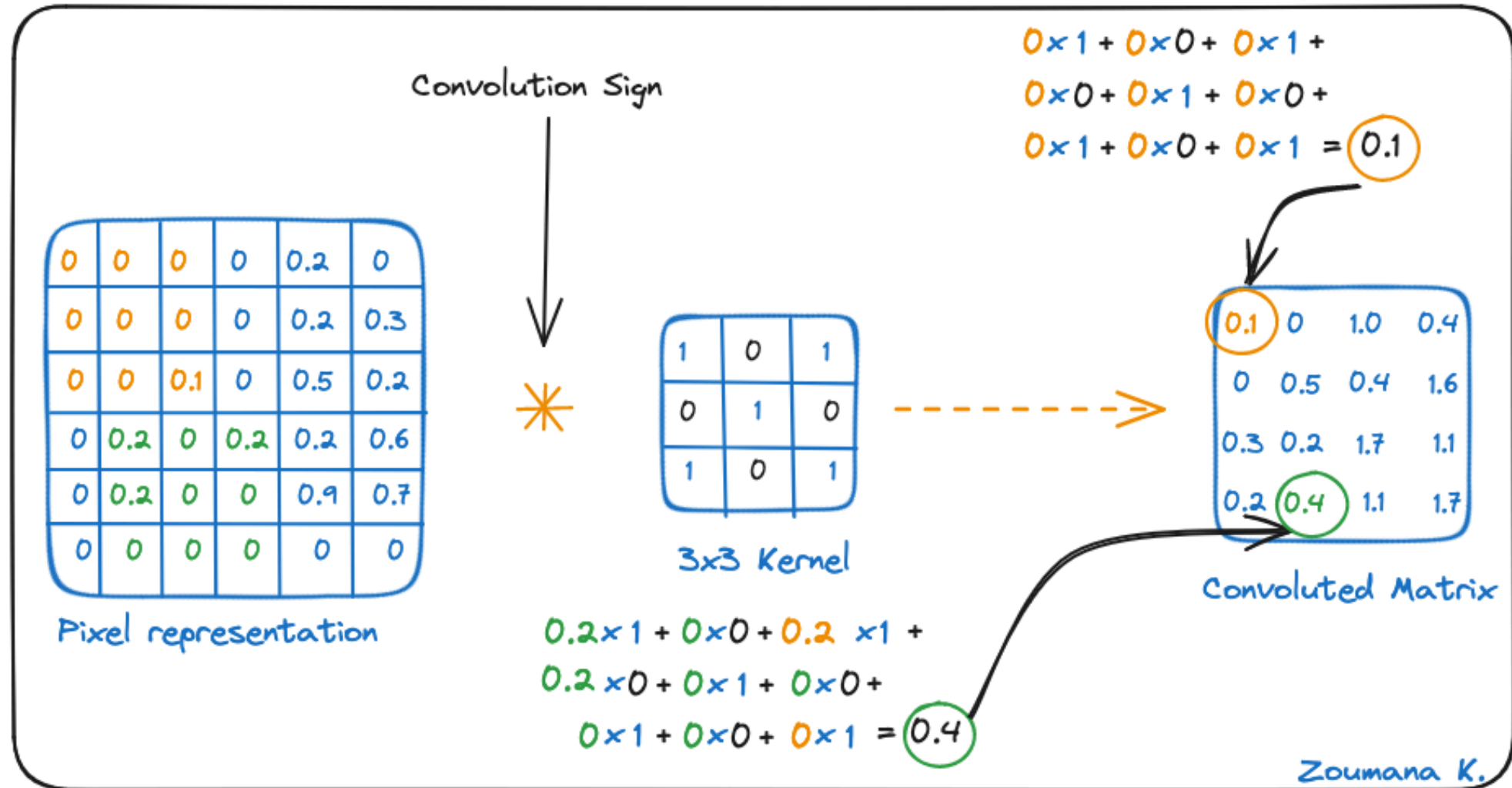
Optical Input

## Convolution layers

This is the first building block of a CNN. As the name suggests, the main mathematical task performed is called convolution, which is the application of a sliding window function to a matrix of pixels representing an image. The sliding function applied to the matrix is called kernel or filter, and both can be used interchangeably.

Illustration of the input image and its pixel representation

32

32

Input Image

Pixel representation

| 0 | 0 | 0 | 0 | 0.2 | 0 |
|---|---|---|---|-----|---|
| 0 | 0 | 0 | 0 | 0.2 | 0.3 |
| 0 | 0 | 0.1 | 0 | 0.5 | 0.2 |
| 0 | 0.2 | 0 | 0.2 | 0.2 | 0.6 |
| 0 | 0.2 | 0 | 0 | 0.9 | 0.7 |
| 0 | 0 | 0 | 0 | 0 | 0 |

Zoumana K.

Application of the convolution task using a stride of 1 with 3x3 kernel

**Activation function**

A ReLU activation function is applied after each convolution operation. This function helps the network learn non-linear relationships between the features in the image, hence making the network more robust for identifying different patterns. It also helps to mitigate the vanishing gradient problems.

**Pooling layer**

The goal of the pooling layer is to pull the most significant features from the convoluted matrix. This is done by applying some aggregation operations, which reduce the dimension of the feature map (convoluted matrix), hence reducing the memory used while training the network. Pooling is also relevant for mitigating overfitting.

The most common aggregation functions that can be applied are:

Max pooling, which is the maximum value of the feature map
Sum pooling corresponds to the sum of all the values of the feature map
Average pooling is the average of all the values.

Application of max pooling with a stride of 2 using 2x2 filter

Convoluted Matrix

| 0.1 | 0 | 1.0 | 0.4 |
|-----|-----|-----|-----|
| 0 | 0.5 | 0.4 | 1.6 |
| 0.3 | 0.2 | 1.7 | 1.1 |
| 0.2 | 0.4 | 1.1 | 1.7 |

Max Pooling
2x2 filter
Stride=2

Pooling Result
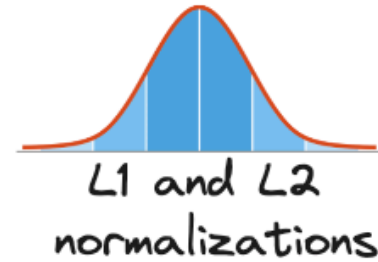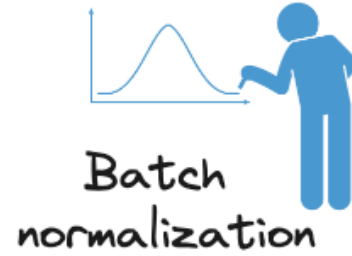
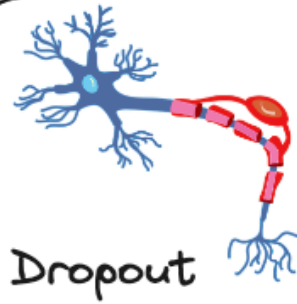| 0.5 | 1.6 |
|-----|-----|
| 0 | 0.5 |

Zoumana K.

**Fully connected layers**

These layers are in the last layer of the convolutional neural network, and their inputs correspond to the flattened one-dimensional matrix generated by the last pooling layer. ReLU activations functions are applied to them for non-linearity.

Finally, a softmax prediction layer is used to generate probability values for each of the possible output labels, and the final label predicted is the one with the highest probability score.

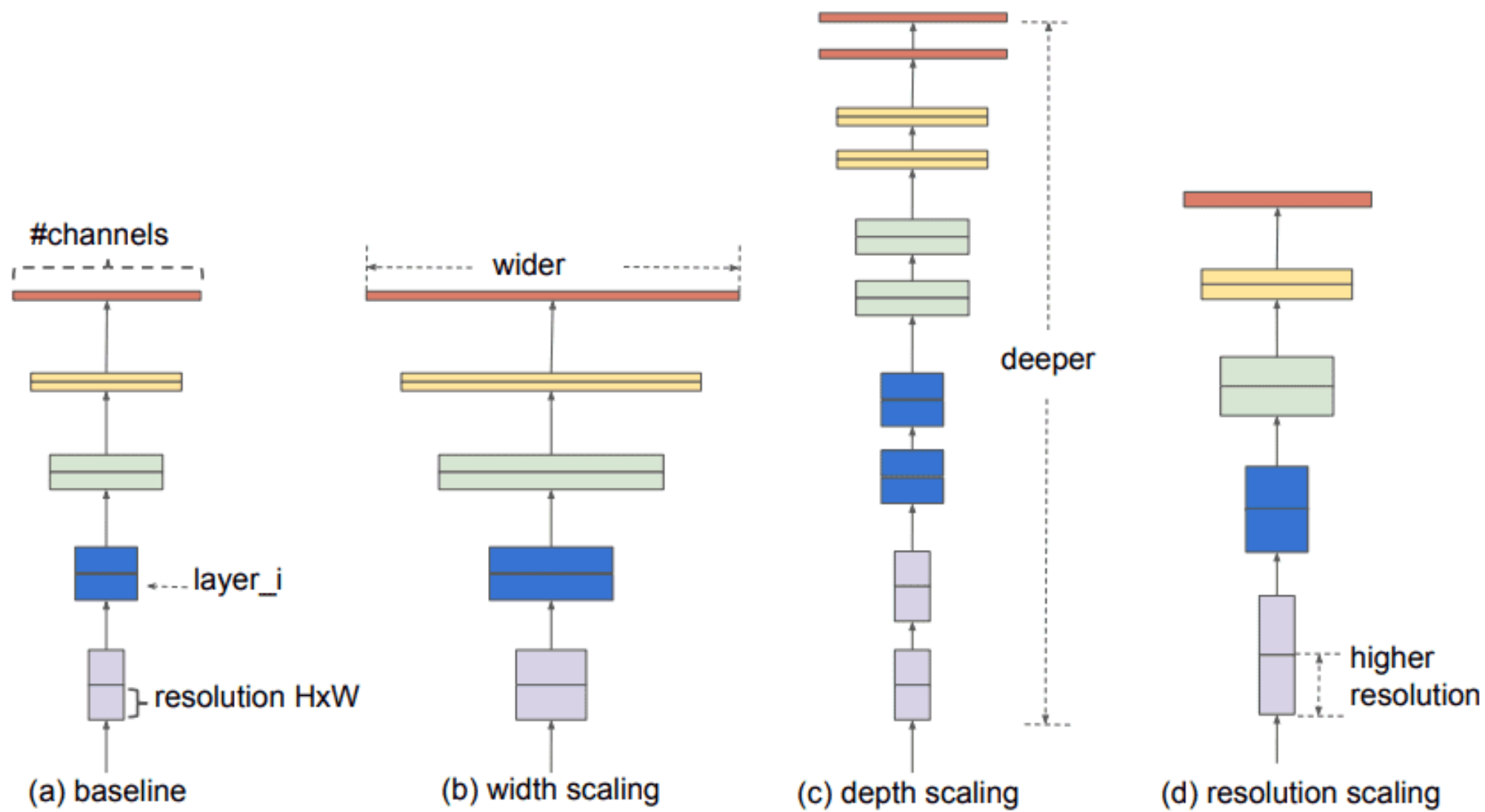# Overfitting and Regularization in CNNs

# 7 Strategies to Mitigate Overfitting in CNNs

Dropout

Batch normalization

Pooling Layers

Early stopping

Noise injection

L1 and L2 normalizations

Data augmentation

Zoumana K.

# ⚛️ Neurnal Network Summary

## Terminology

- A wider network means more feature maps (filters) in the convolutional layers

- A deeper network means more convolutional layers

- A network with higher resolution means that it processes input images with larger width and depth (spatial resolutions). That way the produced feature maps will have higher spatial dimensions.

- https://theaisummer.com/cnn-architectures/

(a) baseline     (b) width scaling     (c) depth scaling     (d) resolution scaling

## AlexNet: ImageNet Classification with Deep Convolutional Neural Networks (2012)

Alexnet [1] is made up of 5 conv layers starting from an 11x11 kernel. It was the first architecture that employed max-pooling layers, ReLu activation functions, and dropout for the 3 enormous linear layers. The network was used for image classification with 1000 possible classes, which for that time was madness.

Now, you can implement it in 35 lines of PyTorch code:

```python
class AlexNet(nn.Module):
    def __init__(self, num_classes: int = 1000) -> None:
        super(AlexNet, self).__init__()
        self.features = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(64, 192, kernel_size=5, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(192, 384, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(384, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(256, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
        )

        self.avgpool = nn.AdaptiveAvgPool2d((6, 6))
        self.classifier = nn.Sequential(
            nn.Dropout(),
            nn.Linear(256 * 6 * 6, 4096),
            nn.ReLU(inplace=True),
            nn.Dropout(),
            nn.Linear(4096, 4096),
            nn.ReLU(inplace=True),
            nn.Linear(4096, num_classes),
        )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        x = self.features(x)
        x = self.avgpool(x)
        x = torch.flatten(x, 1)
        x = self.classifier(x)
        return x
```
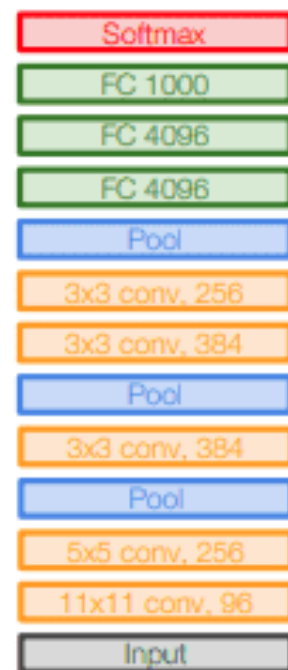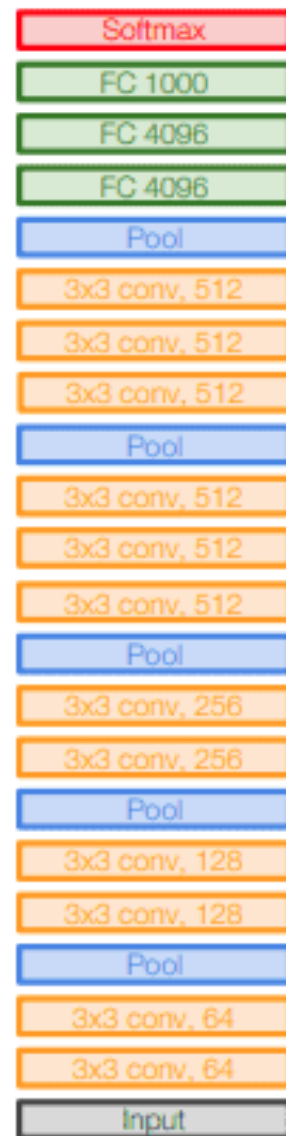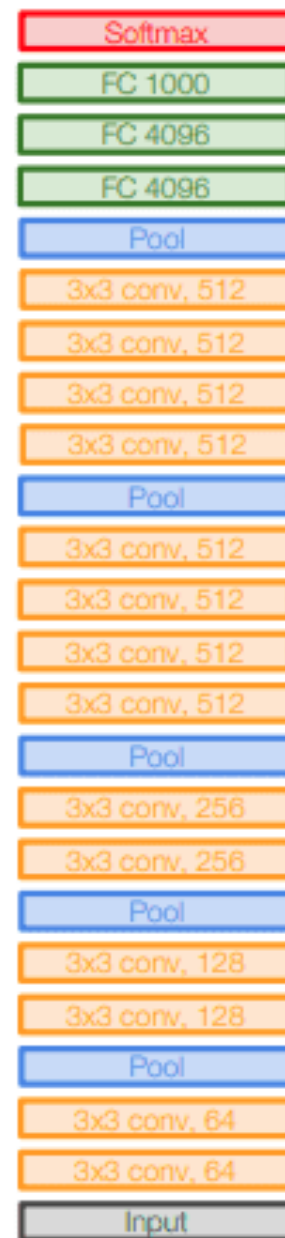
## VGG (2014)

The famous paper "Very Deep Convolutional Networks for Large-Scale Image Recognition" [2] made the term deep viral. It was the first study that provided undeniable evidence that simply adding more layers increases the performance. Nonetheless, this assumption holds true up to a certain point. To do so, they use only 3x3 kernels, as opposed to AlexNet. The architecture was trained using 224 × 224 RGB images.
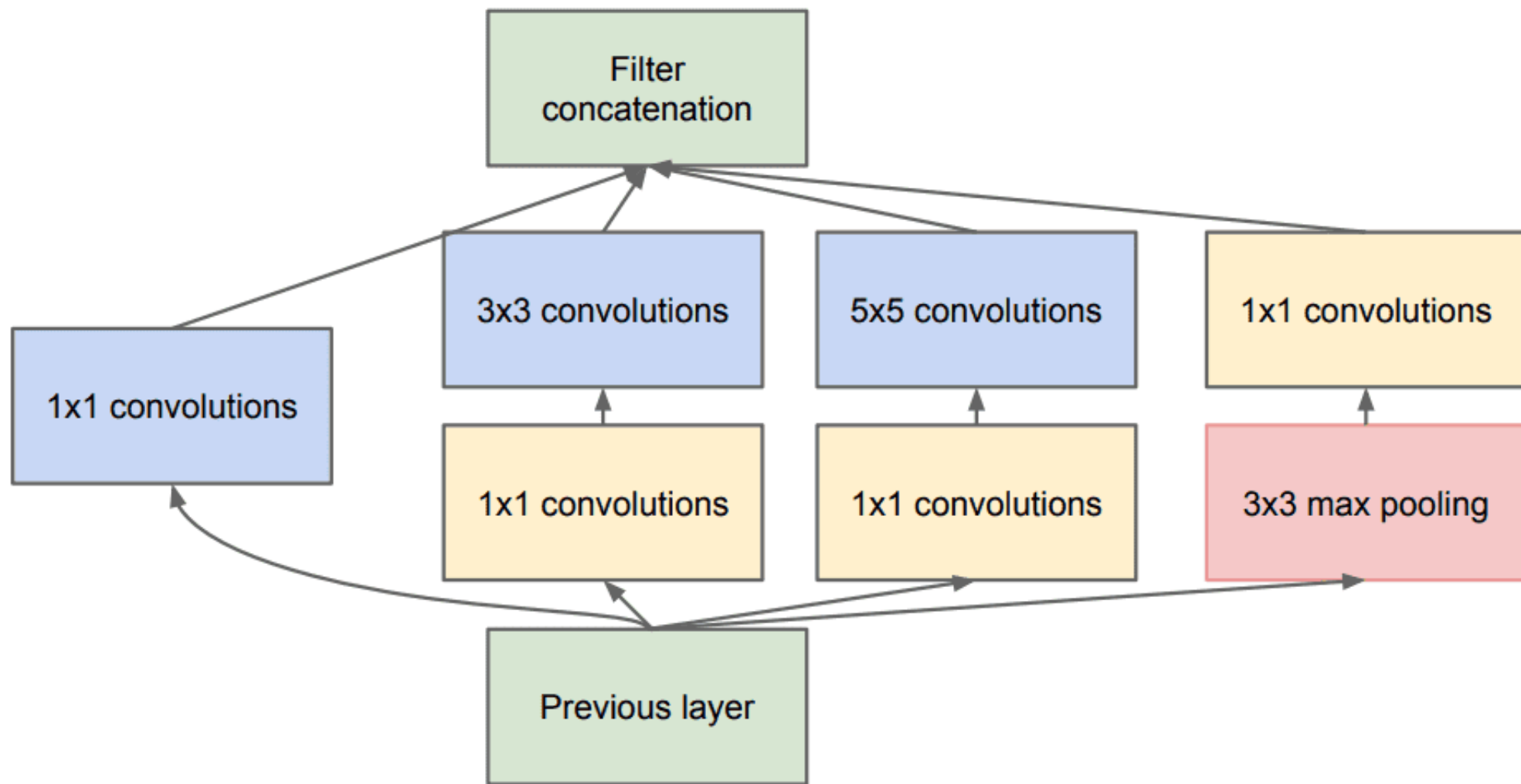
| AlexNet | VGG16 | VGG19 |

## InceptionNet/GoogleNet (2014)

After VGG, the paper "Going Deeper with Convolutions" by Christian Szegedy et al. was a huge breakthrough.

Motivation: Increasing the depth (number of layers) is not the only way to make a model bigger. What about increasing both the depth and width of the network while keeping computations to a constant level?

Filter concatenation

3x3 convolutions    5x5 convolutions    1x1 convolutions

1x1 convolutions

1x1 convolutions    1x1 convolutions    3x3 max pooling

Previous layer

```python
import torch
import torch.nn as nn

class InceptionModule(nn.Module):
    def __init__(self, in_channels, out_channels):
        super(InceptionModule, self).__init__()
        relu = nn.ReLU()
        self.branch1 = nn.Sequential(
                nn.Conv2d(in_channels, out_channels=out_channels, kernel_size=1, stride=1, padding=0),
                relu)

        conv3_1 = nn.Conv2d(in_channels, out_channels=out_channels, kernel_size=1, stride=1, padding=0)
        conv3_3 = nn.Conv2d(out_channels, out_channels, kernel_size=3, stride=1, padding=1)
        self.branch2 = nn.Sequential(conv3_1, conv3_3,relu)

        conv5_1 = nn.Conv2d(in_channels, out_channels=out_channels, kernel_size=1, stride=1, padding=0)
        conv5_5 = nn.Conv2d(out_channels, out_channels, kernel_size=5, stride=1, padding=2)
        self.branch3 = nn.Sequential(conv5_1,conv5_5,relu)

        max_pool_1 = nn.MaxPool2d(kernel_size=3, stride=1, padding=1)
        conv_max_1 = nn.Conv2d(in_channels, out_channels=out_channels, kernel_size=1, stride=1, padding=0)
        self.branch4 = nn.Sequential(max_pool_1, conv_max_1,relu)

    def forward(self, input):
        output1 = self.branch1(input)
        output2 = self.branch2(input)
        output3 = self.branch3(input)
        output4 = self.branch4(input)
        return torch.cat([output1, output2, output3, output4], dim=1)

model = InceptionModule(in_channels=3,out_channels=32)
inp = torch.rand(1,3,128,128)
print(model(inp).shape)
```
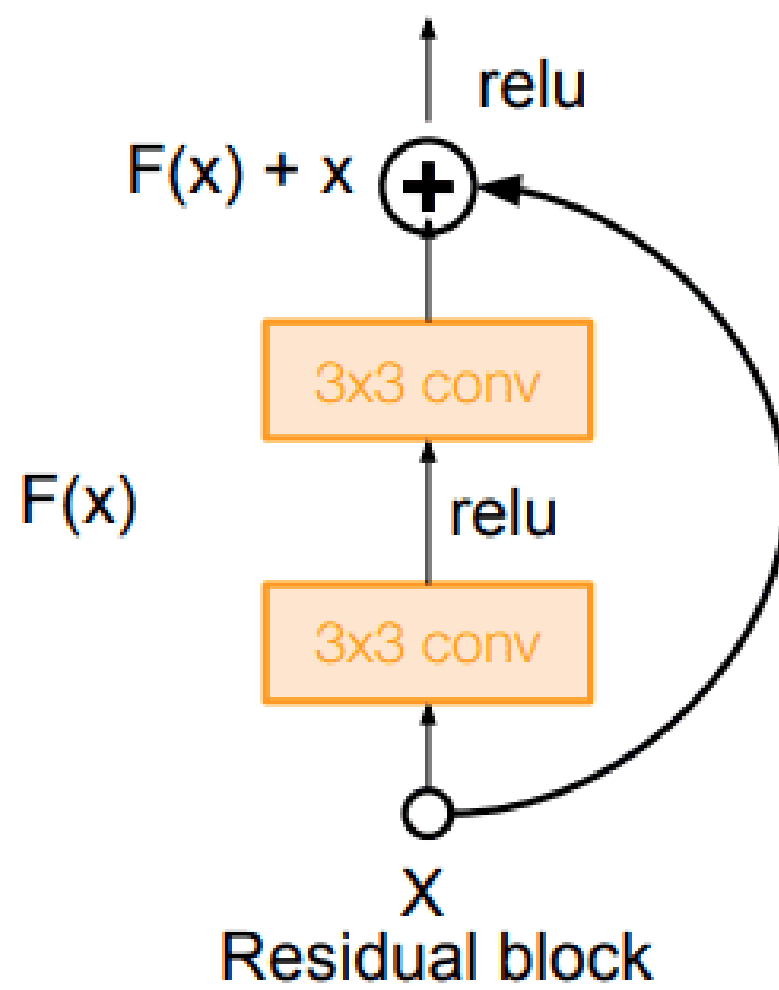
## Inception V2, V3 (2015)

Later on, in the paper "Rethinking the Inception Architecture for Computer Vision" the authors improved the Inception model based on the following principles:

- Factorize 5x5 and 7x7 (in InceptionV3) convolutions to two and three 3x3 sequential convolutions respectively. This improves computational speed. This is the same principle as VGG.

- They used spatially separable convolutions. Simply, a 3x3 kernel is decomposed into two smaller ones: a 1x3 and a 3x1 kernel, which are applied sequentially.

- The inception modules became wider (more feature maps).

- They tried to distribute the computational budget in a balanced way between the depth and width of the network.

- They added batch normalization.

- Later versions of the inception model are InceptionV4 and Inception-Resnet.

**ResNet: Deep Residual Learning for Image Recognition (2015)**

All the predescribed issues such as vanishing gradients were addressed with two tricks:

- batch normalization and

- short skip connections

relu

F(x) + x

3x3 conv

F(x)

relu

3x3 conv

X
Residual block

DenseNet: Densely Connected Convolutional Networks (2017)

Big Transfer (BiT): General Visual Representation Learning (2020)

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (2019)
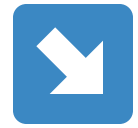
# The ImageNet dataset

The ImageNet dataset contains 14,197,122 annotated images according to the WordNet hierarchy. Since 2010 the dataset is used in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), a benchmark in image classification and object detection. The publicly released dataset contains a set of manually annotated training images. A set of test images is also released, with the manual annotations withheld. ILSVRC annotations fall into one of two categories: (1) image-level annotation of a binary label for the presence or absence of an object class in the image, e.g., "there are cars in this image" but "there are no tigers," and (2) object-level annotation of a tight bounding box and class label around an object instance in the image, e.g., "there is a screwdriver centered at position (20,25) with width of 50 pixels and height of 30 pixels". The ImageNet project does not own the copyright of the images, therefore only thumbnails and URLs of images are provided.

| Model name | Number of parameters [Millions] ImageNet Top 1 | Accuracy | Year |
|---|---|---|---|
| AlexNet | 60 M | 63.3 % | 2012 |
| Inception V1 | 5 M | 69.8 % | 2014 |
| VGG 16 | 138 M | 74.4 % | 2014 |
| VGG 19 | 144 M | 74.5 % | 2014 |
| Inception V2 | 11.2 M | 74.8 % | 2015 |
| ResNet-50 | 26 M | 77.15 % | 2015 |
| ResNet-152 | 60 M | 78.57 % | 2015 |
| Inception V3 | 27 M | 78.8 % | 2015 |

| Model name | Number of parameters [Millions] ImageNet Top 1 | Accuracy | Year |
|---|---|---|---|
| DenseNet-121 | 8 M | 74.98 % | 2016 |
| DenseNet-264 | 22M | 77.85 % | 2016 |
| BiT-L (ResNet) | 928 M | 87.54 % | 2019 |
| NoisyStudent EfficientNet-L2 | 480 M | 88.4 % | 2020 |
| Meta Pseudo Labels | 480 M | 90.2 % | 2021 |

# ⚛️ 参考

[1] https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/

[2] https://www.datacamp.com/tutorial/introduction-to-convolutional-neural-networks-cnns

[3] https://theaisummer.com/cnn-architectures/

# ⬊ **Enjoy your machine learning!**

**https://github.com/wjssx/Statistical-Learning-Slides-Code**

E-mail: csr_dsp@sina.com