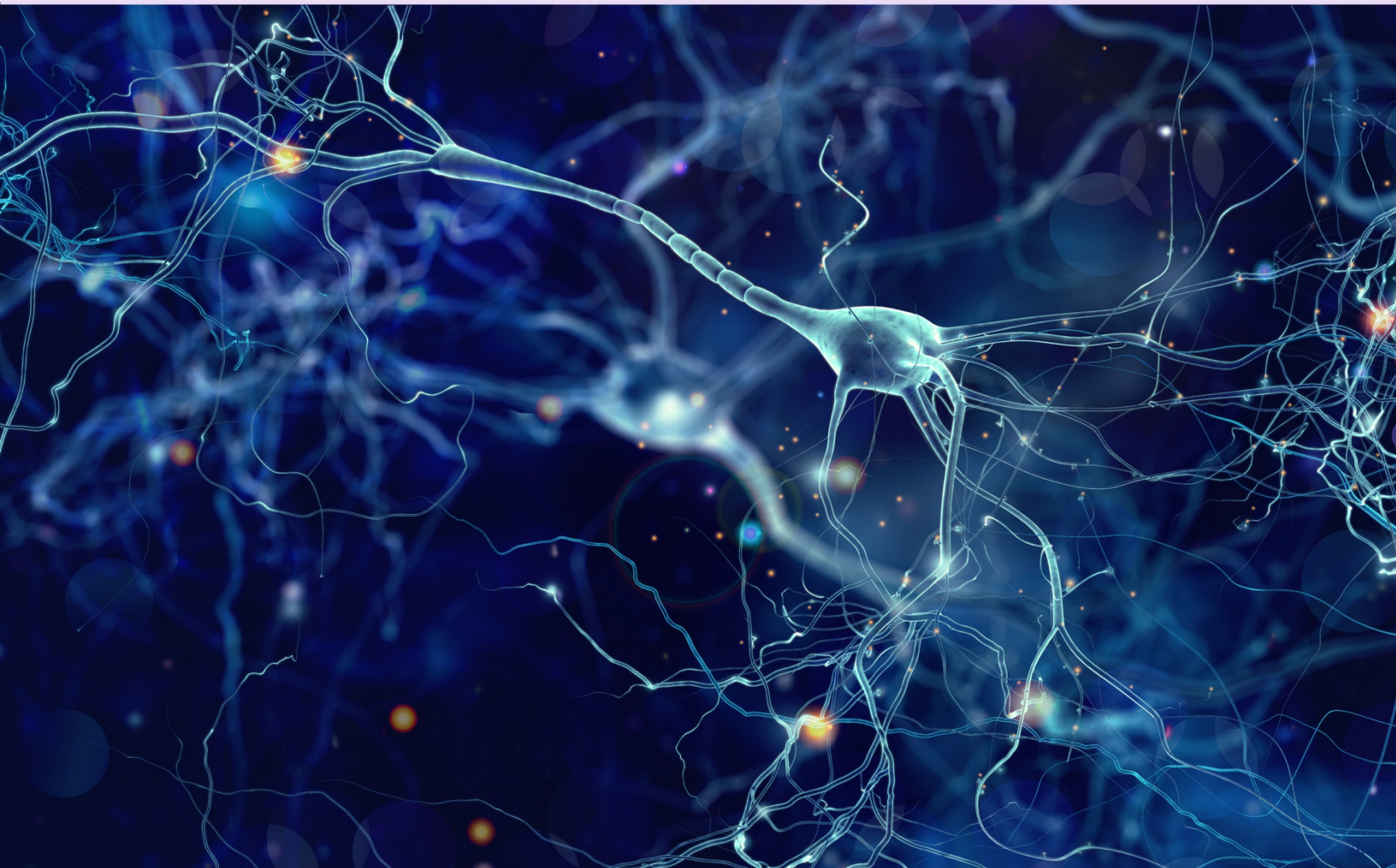


# INTRODUCTION TO MODERN CONVOLUTIONAL NETWORKS



Centrella, Christopher  
Franciscan University of Steubenville



## INTRODUCTION

In recent years, artificial intelligence has emerged as one of the fastest-growing technologies in the computer industry. Everything from cars to computer software, to retail store shopping sites, all use artificial intelligence or AI in their product. In our ever-dynamic computer-driven society, artificial intelligence is used everywhere, in every way, and is at the forefront of our most modern cutting-edge technologies. But what exactly is artificial intelligence, how does it differ from traditional computing, and what are some of the underlying artificial intelligence strategies used today? In this paper, I wish to investigate this topic in greater depth, first introducing artificial intelligence and then expounding upon this to discuss what is at the heart of modern computer image recognition, convolutional neural networks.

## WHAT IS ARTIFICIAL INTELLIGENCE?

First, what is AI? What does it mean to be “artificially intelligent”? Some people, with absolutely no knowledge of the subject, attribute artificial intelligence to playing the role of God in creating new living objects with arms and legs. Nothing could be further from the truth. Artificial intelligence does not involve any organic matter or living organism. Rather, artificial intelligence, at best is only a *simulation* of human intelligence.<sup>1</sup> Another way to describe it would be to say that artificial intelligence is the best attempt computers can make at imitating human behavior. This begs the question, in what ways do computers imitate human behavior? In *Artificial Intelligence for Dummies*, the authors describe four ways in which artificial intelligence attempts to mimic human behavior: Acting humanly, thinking humanly, thinking rationally, and acting rationally.<sup>2</sup>

Now there is an important distinction between acting humanly and acting rationally. John Paul Mueller and Luca Massaron, in the aforementioned book, explain

---

<sup>1</sup> (Mueller and Massaron 2022, Ch. 1)

<sup>2</sup> Ibid.

that acting rationally means that the system “always does the right thing based on the current information, given an ideal performance measure.” In contrast, acting humanly involves other variables that can influence a person’s behavior: instinct, intuition, and the like.<sup>3</sup> While any computer program acts rationally based on the person who programmed it, following a logic-based algorithm, or set of instructions,<sup>4</sup> using artificial intelligence allows the computer to act in a human way, which is much more than simply following a set of instructions. It includes simulating reason and intuition, even in very complex scenarios. Whereas with standard computer programming, every use case is handled explicitly by the programmer, with artificial intelligence the computer itself is making decisions.

There are several ways that artificial intelligence is able to mimic human behavior, but the primary way is through what is known as machine learning.<sup>5</sup> In machine learning, the computer abstracts functions from the data<sup>6</sup> to complete what is known as a model.<sup>7</sup> The skeleton of this model is created by the programmer, but the model is not finished until the individual *weights* are trained automatically, by the computer. The model represents a function, or a “deterministic mapping from a set of input values to one or more output values.”<sup>8</sup> Each unique situation is represented by an input value, and the result is described through an output value. After the model is complete, the computer *infers* from it to apply the model to the situation at hand.<sup>9</sup> This allows the computer to quickly apply the model to new situations and new circumstances without having to retrain, thus allowing for the speed and performance of the iPhone 13, or the Tesla Model 3.

---

<sup>3</sup> [Mueller and Massaron 2022, Ch. 1]

<sup>4</sup> [Kelleher 2019, 7]

<sup>5</sup> [Artascanchez and Joshi 2020, Ch. 1]

<sup>6</sup> [Kelleher 2019, 8]

<sup>7</sup> Ibid, 12-14.

<sup>8</sup> Ibid, 7.

<sup>9</sup> Ibid, 14.

## NEURAL NETWORKS

The scenario I described above represents a typical *neural network*, often referred to as a deep neural net. A neural network can be described as a mathematical model,<sup>10</sup> inspired by the operation of the human brain (see figure 1).<sup>11</sup> So the human brain

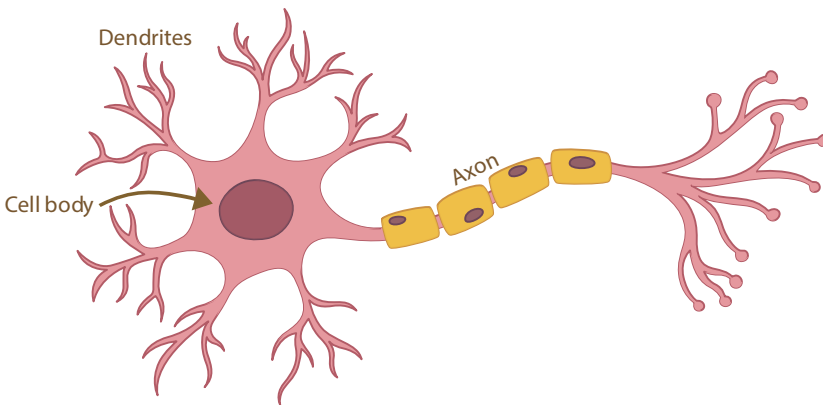


Figure 1. Brain Neuron  
Credit: Adobe Stock

consists of billions of tiny cells, called neurons.<sup>12</sup>

Each of these cells gets data from other cells via *dendrites* and then outputs data to additional cells via a long rod, known as the *axon*.<sup>13</sup> See figure 1. A neural network has all of these characteristics.

First, there are thousands of little computer “cells” which we refer to as *neurons*.<sup>14</sup> In neural networks, these neurons are grouped into layers,<sup>15</sup> where each output cell of a given layer, is connected to each input cell of the next layer. The layers consist of an input layer, “hidden” or intermediate layers, and one or more final output layers. In fact, the number of hidden or intermediate layers, and the output layer(s) is known as the *depth* of the network.<sup>16</sup> Now, each neuron has a unique *weight*, and the value of the neuron’s input is multiplied by the value of the weight, to calculate that neuron’s output.<sup>17</sup>

---

<sup>10</sup> (Kelleher 2019, 39)

<sup>11</sup> Ibid.

<sup>12</sup> Ibid, 65.

<sup>13</sup> Ibid, 66.

<sup>14</sup> Ibid, 67.

<sup>15</sup> Ibid, 67-8.

<sup>16</sup> Ibid, 68.

<sup>17</sup> Ibid, 70.

However, one last thing is needed, and that is the nonlinear activation function. So far, we have a bunch of layers, with a bunch of neurons; each layer takes in all the neurons from the previous layer, and then multiplies them by the appropriate given neuron values. The only issue is that this is a *linear* function.<sup>18</sup> A linear function means that if there is an increase or decrease in one thing, there is an increase or decrease in another thing at the same rate.<sup>19</sup> The problem is that in the real world, many things are not linear. If we limit an artificial intelligence model to a linear function, our model will automatically introduce errors, since it is constrained to only using linear functions, when our data uses a nonlinear function. To solve this problem, the concept of an *activation function* comes in.<sup>20</sup> An activation function applies a nonlinear mathematic function to the calculated value of the neuron, allowing the neuron's value to be nonlinear with the weight of the neuron, before it is input to the next layer.<sup>21</sup> Several activation functions exist, including the Sigmoid function, the SoftMax function, and the Rectified Linear Unit (ReLU) function.<sup>22</sup> The ReLU function is one of the best activation functions, because it reduces the problem of vanishing gradients, discussed later under the section "Backpropagation."<sup>23</sup>

## CALCULATING LOSS

After we choose our activation function and the number of neurons that we want our model to have, the model has to choose the right values for each of the neurons, known as *weights*<sup>24</sup>, or *parameters*.<sup>25</sup> These values are automatically generated by the

---

<sup>18</sup> [Kelleher 2019, 77-8]

<sup>19</sup> Ibid, 78.

<sup>20</sup> Ibid, 71, 76-9.

<sup>21</sup> Ibid.

<sup>22</sup> [Elgendy 2020, Ch. 2]

<sup>23</sup> Ibid.

<sup>24</sup> [Kelleher 2019, 70]

<sup>25</sup> [Elgendy 2020, Ch. 4]

neural network, and are the main way the system “learns.”<sup>26</sup> Now, how does the system know which are the correct parameters to set? This is through something called a loss function, or an error function.<sup>27</sup> The loss function calculates the loss or total error between the outputs generated using the model, and the outputs in the training data. To begin the process, we provide training data to the model, and our training data gives an expected output, which is compared to the output of the model.<sup>28 29</sup> After the first iteration of training is complete, the loss function calculates the total difference between the expected output and actual output, by calculating the sum of each individual difference.<sup>30 31</sup> Two of the most common loss functions are mean squared error and cross entropy.<sup>32</sup> The former does just what it sounds like; it squares the error for each output value, and then calculates the mean. Cross entropy uses a similar approach but with a logarithmic function.

## OPTIMIZATION

Now that we have a loss function, the next step is to choose an *optimization function*. The optimization function is used to determine the minimal loss; essentially, it *optimizes* for the minimum loss.<sup>33</sup> The optimization function will continually update the network so that the correct weights are chosen until the network’s loss is at a minimum.<sup>34</sup> The optimization function has to look at the current loss, compare this to the new loss, and then update the weights accordingly. However, this gets really difficult when you have thousands or millions of neurons, and hence the need for special

---

<sup>26</sup> [Elgendy 2020, Ch. 2]

<sup>27</sup> Ibid.

<sup>28</sup> [Kelleher 2019, 97]

<sup>29</sup> [Artasanchez and Joshi 2020, Ch. 3]

<sup>30</sup> [Kumar 2019, Ch. 4]

<sup>31</sup> [Elgendy 2020, Ch. 2]

<sup>32</sup> Ibid.

<sup>33</sup> Ibid.

<sup>34</sup> Ibid.

optimization algorithms.<sup>35</sup> One really popular optimization algorithm that can help overcome this problem is called batch gradient descent. Basically, you look at the slope or rate of change of the line tangent to a given point of the curve, also known as the *gradient*, and then keep adjusting the model until you can minimize the gradient.<sup>36</sup> Then we can specify the amount to change each time, known as the step size, which allows us to fine-tune the algorithm to work best for any given situation.<sup>37</sup> Today, better variations of this algorithm exist, such as stochastic gradient descent and mini-batch gradient descent. The former uses only one instance of the data that changes randomly, rather than all of the data, while the latter divides the training into sections containing multiple instances, but not all of the data at once.<sup>38</sup> Both of these generally perform better than the standard batch gradient descent.<sup>39</sup>

## BACKPROPAGATION

Once we have a loss function to determine the loss and we know the direction to adjust the weights given this loss, we still have to figure out how much to adjust each weight individually.<sup>40</sup> This is the function of the backpropagation algorithm. The backpropagation algorithm calculates the loss with respect to each specific weight, beginning at the last layer and working its way all the way up the network, until it reaches the initial input layer.<sup>41</sup> The algorithm accomplishes this functionality using partial derivatives. It distributes the loss among various weights by calculating the partial derivative relative to that weight, starting from the output layer and continuing until the start of the network.<sup>42</sup> However, one potential problem is that because the weights are calculated based on the partial derivative starting at the end of the network,

---

<sup>35</sup> [Elgendy 2020, Ch. 2]

<sup>36</sup> Ibid.

<sup>37</sup> Ibid.

<sup>38</sup> Ibid.

<sup>39</sup> Ibid.

<sup>40</sup> Ibid.

<sup>41</sup> Ibid.

<sup>42</sup> Ibid.

as the algorithm approaches the beginning, the value of the loss with respect to the weight diminishes until each weight is only given a very small effect. This is known as the problem of vanishing gradients.<sup>43</sup> Vanishing gradients makes choosing a good activation function much more complicated, because if the activation function grows to infinity as the value increases (or decreases), this will have a greater impact on the partial derivative and so will exacerbate the problem of vanishing gradients.<sup>44 45</sup>

## CONVOLUTIONAL NEURAL NETWORKS

Everything described thus far is at the heart of most modern artificial intelligence, specifically the field of deep learning, the most common artificial intelligence technology used today.<sup>46</sup> Basically, we have a sum of inputs that gather data, process that data through several hidden layers, applying both linear functions and nonlinear activation functions, and then finally output the data through a special output layer, which will be used to convert the result to the appropriate type.<sup>47</sup> Now in order for computers to sense different types of data, various add-ons and modifications to this model may be needed.<sup>48</sup> This paper describes the area of image recognition, also known as computer vision.<sup>49</sup> Computer vision, or the art of training a computer to observe patterns based on sight input,<sup>50</sup> makes use of a special type of neural network to enable image recognition and/or detection of visual objects, known as object detection.<sup>51</sup> This special type of neural network is called a convolutional neural network, or CNN.<sup>52</sup>

---

<sup>43</sup> [Elgendy 2020, Ch. 5]

<sup>44</sup> Ibid.

<sup>45</sup> [Elgendy 2020, Ch. 2]

<sup>46</sup> [Kelleher 2019, 2]

<sup>47</sup> Ibid, 76.

<sup>48</sup> Ibid, 159-60.

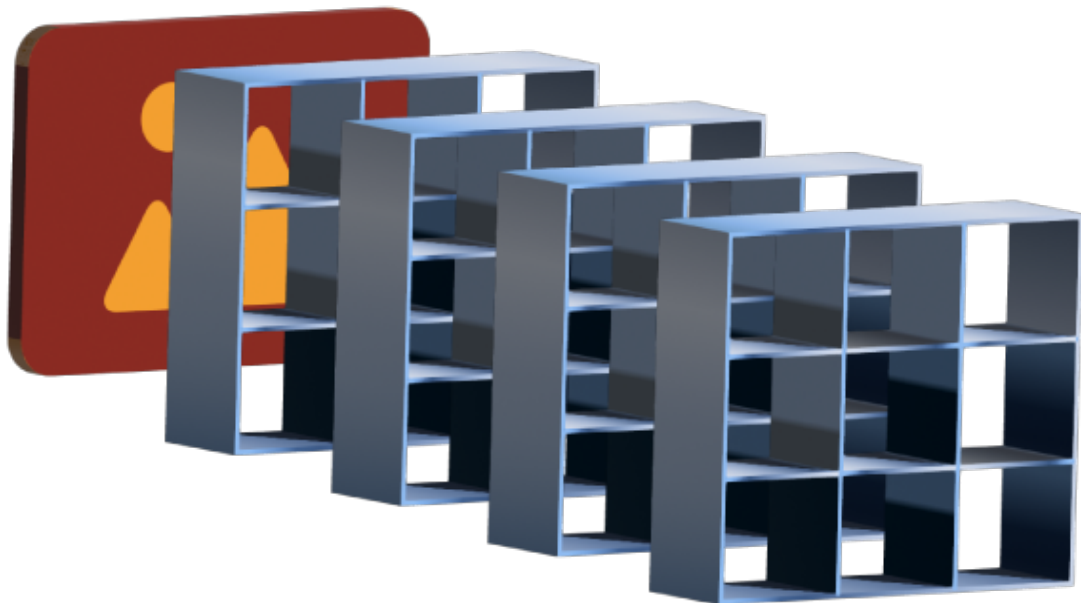
<sup>49</sup> [Elgendy 2020, Ch. 1]

<sup>50</sup> Ibid, Ch. 2.

<sup>51</sup> Ibid.

<sup>52</sup> Ibid, Ch. 3.

A convolutional neural network, or CNN, is a special type of neural net, which uses special sliding “layers,”<sup>53</sup> called convolutional layers, to slide over the object and thus be able to detect pixels at various locations of the image.<sup>54</sup> A convolutional layer consists of one or more *convolutional filters*<sup>55</sup>, where each filter contains  $M \times M$  neurons.<sup>56</sup> The filter, for example of size  $5 \times 5$  neurons, slides over the image and the result of the image pixel values, convolved with the filter, is passed into what is called a feature map, basically a map specifying the presence of various features or visual elements, at a given location [see figure 2].<sup>57</sup>



*Figure 2. Four convolutional layers, with feature maps of size 3 x 3  
Credit: Adobe Stock*

Each feature map searches for the presence of a very specific element, like a line or a curve<sup>58</sup> and as the model goes deeper into the network, the features become less

---

<sup>53</sup> [Elgendy 2020, Ch. 3]

<sup>54</sup> Ibid.

<sup>55</sup> Ibid.

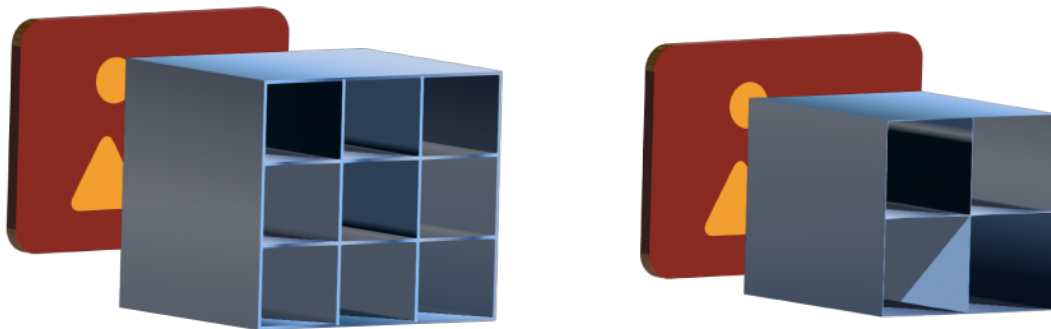
<sup>56</sup> Ibid.

<sup>57</sup> Ibid.

<sup>58</sup> Ibid.

general and more specific, until the system is able to detect the presence of specific objects.<sup>59</sup> Basically, layers are stacked against each other, until we have patterns within patterns within patterns, and the system is able to detect high-level, or human-friendly objects.<sup>60</sup> After all of these layers, there is one more layer. Known as a fully connected layer,<sup>61</sup> <sup>62</sup> this is just a normal layer and is used for the final step to classify the image; or for object detection, the last step before the network is brought into the next system.<sup>63</sup>

Now one problem with the process described above has to do with performance. Each convolutional filter stores pixel data for  $M \times M$  values and as we increase the filter count, the amount of data that has to be stored in memory (VRAM) during training grows to be outrageous. Thus comes the concept of *pooling layers* (see figure 3).<sup>64</sup> A *pooling layer* applies a statistical function to the output from the previous convolutional layer to shrink the size (width and height) of the previous layer.<sup>65</sup> So as we progress, even though the feature map becomes very deep, its width decreases, and so we are able to decrease the amount of memory required to train the model.<sup>66</sup> Another concept is that of *strides*



*Figure 3. Effect of Pooling Layers*  
Credit: Adobe Stock

---

<sup>59</sup> [Elgendy 2020, Ch. 3]

<sup>60</sup> Ibid.

<sup>61</sup> Ibid.

<sup>62</sup> See section on neural networks.

<sup>63</sup> [Elgendy 2020, Ch. 7]

<sup>64</sup> Ibid.

<sup>65</sup> Ibid.

<sup>66</sup> Ibid.

and *padding*. The stride specifies the amount to slide the convolutional filter each time, while padding specifies whether we should add additional empty data or 0s around the border of the image.<sup>67</sup> Nowadays, there is a push to use strides and padding to decrease the image size instead of pooling layers.<sup>68</sup>

## IMPROVING PERFORMANCE

The sections above describe the basic overview of a convolutional neural network. However, improving the performance of a convolutional network requires careful skill and observation, and is an ever-evolving field.<sup>69</sup> So what are some techniques for improving the performance of convolutional neural networks, in 2022? This section begins by discussing ways to measure performance, follows with a general introduction on designing accurate computer vision models and concludes by discussing some state-of-the-art techniques.

## MEASURING PERFORMANCE

In order to know how to improve our model, we first have to be able to measure performance.<sup>70</sup> While loss is used to measure model accuracy during training, how do we measure the performance of the trained model itself? There are several performance metrics we can use. Accuracy measures how many predictions are correct, compared with total predictions.<sup>71</sup> However, in some cases, this can be deceptive.<sup>72</sup> Elgendy gives an example of a rare disease that is only found 0.1% of the time. A test that always returned negative would be 99.9% accurate, but nevertheless not testing anything.<sup>73</sup> Recall, or sensitivity tells us how many the model marked as positive which were positive,

---

<sup>67</sup> [Elgendy 2020, Ch. 3]

<sup>68</sup> Ibid.

<sup>69</sup> Ibid, Ch. 4.

<sup>70</sup> Ibid.

<sup>71</sup> Ibid.

<sup>72</sup> Ibid.

<sup>73</sup> Ibid.

over the amount marked as positive and incorrectly marked as negative.<sup>74</sup> Precision or specificity, on the other hand, is used to measure how many the model marked as positive that were positive, over the total number marked as positive.<sup>75</sup> Finally, the F-score shows us both precision and recall and is calculated by the following equation<sup>76</sup>:  $\frac{2pr}{p+r}$ . Used together, we can find the best benchmarks for our model, depending on the situation.<sup>77</sup>

## GENERAL TIPS

All of the essential components of a deep learning model, like the loss, optimization, and backpropagation functions, as well as features specific to computer vision, such as convolutional layers, all have adjustable properties that we can tune to obtain optimal performance. These are known as hyperparameters,<sup>78</sup> and are set by the machine learning engineer, in contrast to parameters, which are set automatically by the system when the model is training.<sup>79</sup>

First and foremost, some of the hyperparameters that we can tune are the number of hidden layers and neurons in each layer.<sup>80</sup> While having more layers can produce better results, there are a couple of downsides. First, more layers or neurons can lead to longer training times, but also adding additional complexity will begin decreasing performance, a condition known as overfitting.<sup>81</sup> Generally, when our model starts underperforming when we add layers, we have two options. We can stop making the model more complex, or we can use additional techniques to prevent overfitting. One solution is to add dropout layers, which randomly turn off some of the neurons during

---

<sup>74</sup> (Elgendy 2020, Ch. 4)

<sup>75</sup> Ibid.

<sup>76</sup> Ibid.

<sup>77</sup> Ibid.

<sup>78</sup> (Kelleher 2020, 80)

<sup>79</sup> Ibid, 80, 82.

<sup>80</sup> (Elgendy 2020, Ch. 4)

<sup>81</sup> Ibid.

training.<sup>82</sup> Another tool we can use is called L2 Regularization. It adds a small penalty to the loss function that increases during each training iteration. This allows the model to update less frequently if there is a change, and so prevent overfitting.<sup>83</sup>

We can also choose which activation function to use. For example, the ReLU function will count negative values as zero, so a recent addition is to use a special type of ReLU function, called Leaky ReLU, so that negative numbers are accounted for.<sup>84</sup> We can then specify the rate at which the function “leaks.” We can also specify hyperparameters for the loss and optimization functions. For example, we can specify which loss function to use, or which mini-batch size and learning rates to use with mini-batch gradient descent.<sup>85</sup> Finally, we can choose to specify how many times to train the data. These are known as iterations or epochs. The more epochs, the more thorough our model, but also the longer and the greater the risk of overfitting our data.<sup>86</sup> One especially sneaky option is to set a special callback function, that either automatically decreases our learning rate or completely stops training once our loss hasn’t improved after the specified number of epochs.<sup>87</sup>

Finally, in addition to hyperparameters, there is one more thing, which has a huge impact on our data: preprocessing or preparing our training data. By conforming our data to some baseline, like the normal distribution, we can reduce the loss and help prevent overfitting.<sup>88</sup> One quick technique we should always use is to *normalize* our data.<sup>89</sup> Basically, we ensure that the data have small values all within the same range. This usually means that our data is converted to be in the range [0,1],<sup>90</sup> which is

---

<sup>82</sup> [Elgendy 2020, Ch. 3]

<sup>83</sup> Ibid.

<sup>84</sup> Ibid.

<sup>85</sup> Ibid.

<sup>86</sup> Ibid.

<sup>87</sup> Ibid.

<sup>88</sup> Ibid.

<sup>89</sup> Ibid.

<sup>90</sup> Ibid.

especially useful, since computers use the binary system. For convolutional neural networks, some preprocessing techniques we can use are augmenting different transformations to the images, converting images to grayscale, and resizing images to one size, prior to training.<sup>91</sup> These are also some examples of regularization techniques, which manipulate the data in some way to improve performance.<sup>92</sup>

## FEW-SHOT OBJECT DETECTION

Modern deep learning models perform well, even in the area of computer vision, when given very large amounts of data.<sup>93</sup> However, there is the catch: massive amounts of manually labeled, parsed data is required in order to train the model. While this can be feasible for some models through crowdsourcing sites like Amazon Turk<sup>94</sup>, it still remains a large hurdle before more widespread usage of CNN models is possible. One possible solution to this problem is Few-Shot Object Classification and Few-Shot Object Detection, collectively known as Few-Shot Learning.<sup>95</sup> Similar to how humans are able to quickly learn based on previous experience, Few-Shot Learning allows the system to learn how to identify objects after only a few examples, through utilizing previous knowledge.<sup>96</sup>

Few-Shot Object Detection operates on three essential components, the data, the model, and the algorithm. Prior knowledge is utilized for each of these categories to enable the system to learn based on only a few supervised examples.<sup>97</sup> <sup>98</sup> One technique is transfer learning, where a previously trained model is imported, but this can be

---

<sup>91</sup> [Elgendy 2020, Ch. 3]

<sup>92</sup> [Elgendy 2020, Ch. 4]

<sup>93</sup> [Antonelli et al. 2022]

<sup>94</sup> [Mueller and Massaron 2022, Ch. 11]

<sup>95</sup> [Antonelli et al. 2022]

<sup>96</sup> Ibid.

<sup>97</sup> Ibid.

<sup>98</sup> Whereas unsupervised learning is only based on finding patterns, supervised learning requires training data. See [Mueller and Massaron 2022].

difficult for object detection.<sup>99</sup> Another approach is distance metric learning, where samples are embedded to a low dimension space, and by comparing sample data and training data, the system is able to easily discover important features and distinguish between classes.<sup>100 101</sup> Finally, meta learning teaches the model how to learn, using what is known as meta knowledge. Basically, the system is able to extract helpful information and then use this information to improve the model.<sup>102</sup>

## ADDITIONAL REGULARIZATION TECHNIQUES

Regularization techniques were discussed earlier, in the section on general tips for improving performance. In particular, dropout layers and L2 regularization were mentioned, as well as data augmentation. However, in the past four years, numerous additional regularization technologies have been developed, a few of which are discussed here.<sup>103</sup>

## DATA AUGMENTATION

The first category of data regularization techniques are new data augmentation methods. One technique is Cutout, a form of data augmentation where a portion of the image is removed, before being fed into the network for training.<sup>104</sup> Because this is supervised data, the augmented data is labeled, and so it has the effect of adding additional blemished samples to the training dataset.<sup>105</sup> A similar technique is RandomErasing, where the portions that are removed are replaced with random noise, instead of whitespace.<sup>106</sup>

---

<sup>99</sup> (Antonelli et al. 2022)

<sup>100</sup> Ibid.

<sup>101</sup> (Dor Kedem 2021)

<sup>102</sup> Ibid.

<sup>103</sup> (Gonçalves Dos Santos et al. 2022)

<sup>104</sup> Ibid.

<sup>105</sup> Ibid.

<sup>106</sup> Ibid.

Furthermore, a new technology called AutoAugment, automatically determines the appropriate augmentations for a dataset, and obtained excellent results on the CIFAR-10, CIFAR-100 and ImageNet datasets, and could even be transferred to additional datasets.<sup>107</sup> Moreover, the biggest issue, the amount of time required to train, is largely fixed with a newer version called Fast AutoAugment.<sup>108</sup> Yet another technology, RandAugment looks at the 14 most common policies from AutoAugment and applies those, achieving remarkable results.<sup>109</sup> A few other augmentation regularization techniques are MixUp, CutMix, CutBlur, BatchAugment, and FixRes. MixUp and CutMix combine two images and help the model correctly decipher the probability of it being one of the two images.<sup>110</sup> CutBlur helps train super-resolution images as well as reconstruct data for blurred images.<sup>111</sup> BatchAugment splits the GPU memory when training, so half of the memory can be used for augmented images.<sup>112</sup> Finally, FixRes uses a larger resolution for testing than for training, which can improve performance by making the model more reliable and use less time.<sup>113</sup> Finally, performance gains can also be reached by combining augmentation methods.<sup>114</sup>

## INTERNAL STRUCTURE CHANGES

Another category of regularization techniques works by changing the internal structure of the neurons during training. For example, Dropout randomly shuts off some neurons when training each epoch.<sup>115</sup> A new version of Dropout, known as MaxDropout, shuts off neurons based on their activation above a given level. The more a neuron

---

<sup>107</sup> (Gonçalves Dos Santos et al. 2022)

<sup>108</sup> Ibid.

<sup>109</sup> Ibid.

<sup>110</sup> Ibid.

<sup>111</sup> Ibid.

<sup>112</sup> Ibid.

<sup>113</sup> Ibid.

<sup>114</sup> Ibid.

<sup>115</sup> Ibid.

exceeds a given threshold, the greater the probability of the neuron being turned off.<sup>116</sup> DropBlock removes correlated regions created by the CNN system, and TargetDrop can be combined with DropBlock to remove only appropriate regions.<sup>117</sup> Moreover, some additional strategies are AutoDrop, which automatically determines the best drop design, Shake-Shake, which assigns different weight values for each forward and backward (backpropagation) pass, and ShakeDrop, a smarter, more capable version of Shake-Shake.<sup>118</sup> Finally, a special technique, called Manifold MixUp, performs mix-up after training has begun, i.e., mixing up images after the input layer.<sup>119</sup>

## LABEL REGULARIZATION

Finally, some of the techniques mentioned above are also label regularizers. Namely, they control how the label data is presented, to help the model train more effectively.<sup>120</sup> For example, MixUp averages the values of the labels.<sup>121</sup> This is also known as label smoothing.<sup>122</sup> Other label smoothing techniques are also commonly used, like one-hot encoding, that reformats labels to help prevent overfitting.<sup>123</sup> <sup>124</sup> There are a few other techniques as well, which aim to improve label smoothing. Some of them are Two-Stage Label Smoothing, which only uses label smoothing for part of training, Structural Label smoothing, which uses Bayes' Theorem to calculate error, which is then used to determine smoothing values, and JoCor, which trains two networks at the same time and then compares the two.<sup>125</sup> Unfortunately, no regularizer is best for all circumstances. The

---

<sup>116</sup> (Gonçalves Dos Santos et al. 2022)

<sup>117</sup> Ibid.

<sup>118</sup> Ibid.

<sup>119</sup> Ibid.

<sup>120</sup> Ibid.

<sup>121</sup> Ibid.

<sup>122</sup> Ibid.

<sup>123</sup> Ibid.

<sup>124</sup> (Elgendy 2020, Ch. 3)

<sup>125</sup> (Gonçalves Dos Santos et al. 2022)

best approach is to attempt combining various regularizers, and adjusting them for the problem at hand, until the best results are obtained.<sup>126</sup>

## MODERN CNN BLOCKS

Lastly, in recent years there has been a huge development in CNN models, replacing standard convolution layers with special convolution blocks designed to improve performance.<sup>127</sup> One example is with the ResNet architecture, which uses special residual blocks to prevent the vanishing gradients problem that was discussed earlier.<sup>128</sup> Basically, there are “shortcut connections,” which enable “residual” connections, i.e. shortcuts to earlier layers in the network which are connected directly to the output.<sup>129</sup> More recently, EfficientNet uses multiple types of these Bottleneck Blocks or BBs to achieve both high accuracy and low computation cost.<sup>130 131</sup> This also enables the model to scale well even for very large networks.<sup>132</sup> One type of block used by EfficientNet is the point-wise convolution layer, which uses 1 x 1 filters, reducing the total number of operations.<sup>133</sup> Depth-wise convolution blocks are also employed, which divide the input (pixel) channel(s), into additional channels.<sup>134 135</sup> This helps improve performance since parallelization or parallel processing can be employed.

However, one setback we are facing is the issue of finding a CNN accelerator that improves performance for some of these recent, more efficient blocks.<sup>136</sup> For example,

---

<sup>126</sup> [Gonçalves Dos Santos et al. 2022]

<sup>127</sup> [Lee et al. 2022]

<sup>128</sup> [He et al. 2015]

<sup>129</sup> Ibid.

<sup>130</sup> [Lee et al. 2022]

<sup>131</sup> [Tan et al. 2019]

<sup>132</sup> Ibid.

<sup>133</sup> [Lee et al. 2022]

<sup>134</sup> Ibid.

<sup>135</sup> One input channel is used for grayscale images, while color images have three channels. See [Elgendy et al. 2020, Ch. 1].

<sup>136</sup> [Lee et al. 2022]

the depth-wise convolution blocks do not reuse the input feature maps between output channels.<sup>137</sup> However, this can be a problem, because it prevents the blocks from being sped up via conventional CNN accelerators.<sup>138</sup> One recent proposal is the MVP CNN Accelerator, which features matrix, vector, and processing units close to memory. Utilizing a concept similar to that used in Apple's M1 Chip<sup>139</sup>, this enables the accelerator to optimize the performance of both computationally intensive layers and memory-intensive layers, and thus be effective for all types of modern convolutional networks, at only a small performance penalty.<sup>140</sup>

## CONCLUSION

Modern artificial intelligence has gone a long way to become the backbone of computer-driven cutting-edge technologies. A simulation of human intelligence and based on the human brain, artificial intelligence seeks to imitate human behavior, and accomplishes this primarily through the use of what is known as a neural network or neural net. Containing millions of neurons, the system automatically finds the optimal values for these neurons through the use of loss, optimization and backpropagation functions. In the area of computer vision, artificial intelligence is applied primarily through the use of convolutional neural networks, which use a special convolution layer featuring a convolution filter that passes over the input and extracts data into feature maps. These feature maps start very generally and slowly grow more specific until the computer is able to recognize high-level objects.

Some ways of improving convolutional neural networks include the use of dropout layers, L2 Regularization, and normalization of data. Adjusting the hyperparameters allows us to fine-tune the model to perform best to our situation. Moreover, Few-Shot Learning allows the computer to easily learn new images and adapt to new scenarios

---

<sup>137</sup> (Lee et al. 2022)

<sup>138</sup> Ibid.

<sup>139</sup> (Apple 2020)

<sup>140</sup> Ibid.

through the use of previous knowledge, utilizing technologies like transfer learning and meta learning. Furthermore, numerous data regularization techniques exist, including several new data augmentation techniques, various methods that change the structure of the data, and lastly, various ways to obtain optimal performance through the use of label smoothing.

Finally, modern convolutional networks use special blocks instead of standard convolution layers, which help increase performance, including precision and training time. One state-of-the-art system utilizes a special accelerator that would enable these new blocks to be further optimized, just as conventional accelerators are able to optimize standard convolution layers. All of these techniques improve the performance of CNN networks, and are likely to be expanded upon in the near future, to pave way to new innovations and technologies. In the meantime, numerous artificial intelligence technologies and ever-growing state-of-the-art innovation make the current era one of the most exciting times to be a computer programmer, and one of the best for a machine learning engineer.

## REFERENCES

- Mueller, John Paul, and Luca Massaron. 2022. *Artificial Intelligence for Dummies. 2<sup>nd</sup> edition*. Hoboken, NJ: John Wiley & Sons, Inc.
- Artascanchez, Alberto, and Prateek Joshi. 2020. *Artificial Intelligence with Python. 2<sup>nd</sup> edition*. Birmingham, UK: Packt Publishing.
- John D. Kelleher. 2019. *Deep Learning*. Cambridge: The MIT Press.  
<https://doi.org/10.7551/mitpress/11171.001.0001>
- Mohamed Elgendy. 2020. *Deep Learning for Vision Systems*. Shelter Island, NY: Manning Publications.
- Rahul Kumar. 2019. *Machine Learning Quick Reference*. Birmingham, UK: Packt Publishing.
- Antonelli, Simone, Danilo Avola, Luigi Cinque, Donato Crisostomi, Gian Luca Foresti, Fabio Galasso, Marco Raoul Marini, Alessio Mecca, and Danielle Pannone. 2022. "Few-Shot Object Detection: A Survey." *ACM Computing Surveys*.  
<https://dl.acm.org/doi/pdf/10.1145/3519022>
- Dor Kedem. 2021. "PyData Global 2021 Introduction to Distance Metric Learning." Accessed April 25, 2022. [https://github.com/kedemdor/metric-learning-talk/blob/main/PyData\\_Global\\_2021\\_Introduction\\_to\\_Distance\\_Metric\\_Learning.ipynb](https://github.com/kedemdor/metric-learning-talk/blob/main/PyData_Global_2021_Introduction_to_Distance_Metric_Learning.ipynb)
- Gonçalves Dos Santos, Claudio Filipi, and João Paulo Papa. 2022. "Avoiding Overfitting: A Survey on Regularization Methods for Convolutional Neural Networks." *ACM Computing Surveys*. <https://dl.acm.org/doi/pdf/10.1145/3510413>
- Lee, Sunjung, Jaewan Choi, Wonkyung Jung, Byeongho Kim, Jaehyn Park, Hweesoo Kim, and Jung Ho Ahn. 2022. "MVP: An Efficient CNN Accelerator with Matrix, Vector, and Processing-Near-Memory Units." *ACM Computing Surveys*.  
<https://dl.acm.org/doi/pdf/10.1145/3497745>

He, Kaiming, Xiangyu Zhang, Shaoqing Ren, Jian Sun. 2015. "Deep Residual Learning for Image Recognition." <https://arxiv.org/pdf/1512.03385.pdf>

Tan, Mingxing, and Quoc V. Le. 2019. "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks." *International Conference on Machine Learning*. <https://arxiv.org/pdf/1905.11946.pdf>

Apple Press Release. 2020. "Apple Unleashes M1." Accessed April 25, 2022. <https://www.apple.com/newsroom/2020/11/apple-unleashes-m1/>