

Udacity: Follow Me Report

Shane Reynolds

May 31, 2018

Contents

1	Introduction	1
2	Network Architecture	2
2.1	Convolutional Neural Networks	4
2.2	Fully Convolutional Neural Networks	6
3	Network Training	6
3.1	Epochs	6
3.2	Learning Rate	6
3.3	Batch Size	6
4	Performance & Model Generalisation	6
5	Future Enhancements	6

1 Introduction

Computer vision is a subset of robotic perception - it has been defined as the development of autonomous systems which can perform tasks achieved by human visual systems (Huang, 1996). This means the acquisition of digital image data from an optical camera, and some type of interpretation of the acquired image. A simple example of a task that is routinely performed by a human visual system, which is sought for computer vision systems, is answering the question: *Is there a puppy in Figure 1?*, or *Where is the puppy in Figure 1?*.



Figure 1: Computer vision is interested in answering questions such as *Is there a puppy in the image?* or *Where is the puppy in the image?*

There are many sub-fields of computer vision such as scene reconstruction, event detection, video tracking, object recognition, 3D pose estimation, and motion estimation. This paper will focus on classification, using an approach called semantic segmentation. Shelhamer, Long and Darrell (2016) define semantic segmentation as a method of inference which is able to categorise fine image details.

This is achieved by classifying each pixel in the image, and labelling it with the class of its enclosing object or region. A Fully Convolutional Neural Network (FCN) is proposed as the architecture to implement semantic segmentation. An FCN model was trained and implemented on a robotic agent.



Figure 2: The UAV agent in the simulated environment.



Figure 3: The *hero* can be seen coloured in red, with the UAV agent following her.

The agent was tasked with the identification (and subsequent tracking) of an individual, known as the *hero*, in a 3D simulated environment built with Unity. The 3D simulated environment is a small city consisting of buildings, roads, elevated highways, and vegetation. The robotic agent is an unmanned aerial vehicle (UAV), as shown in Figure 2, which is fitted with a panning optical camera. The agent roams the simulation until it is able to locate the hero using the trained FCN, at which point the UAV will track the hero. The hero is a simulated person, as seen in Figure 3.

2 Network Architecture

FCNs are widely used for computer vision applications, and are a type of Artificial Neural Network (ANN). These ANNs are computational models which, once trained on a dataset, can be used to make classification predictions, or value estimations, based on a set of feature inputs that the model has been trained on. A typical fully connected feed-forward ANN consists of an input layer, one or more hidden layers, and an output layer, as shown in Figure 4. Hidden layers are made up of multiple nodes called neurons (perceptrons) which receive input from the previous layers in the model. Each of the neurons contain non-linear activation functions, such as a sigmoid or ReLU, which are loosely thought of as decision making elements.

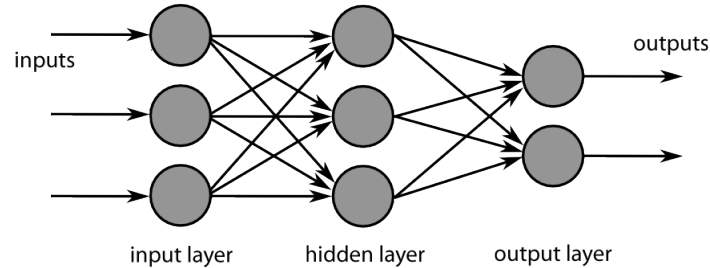


Figure 4: A feed-forward artificial neural network consists of an input layer, which receives feature inputs, some hidden layers, and an output layer for classification.

The model neurons are weighted, and changing a weight changes the neurons’s contribution to the model, which in turn affects the overall model output. Weight changes occur during model training, which uses large volumes of labelled data to adjust the weights. Hidden layers are important because they allow highly non-linear models to be constructed, providing an approach for estimating complex phenomena which may be difficult to model with classical approaches, or computationally intractable. Generally, the more hidden layers, the more non-linear the model. Network architectures with multiple hidden layers have become so wide spread that the term Deep Neural Network (DNN) was coined to describe feed-forward ANNs which use two or more hidden layers. It must be noted that whilst increased non-linearity may allow us to model more complex phenomenon, making the ANN deeper does not guarantee increased model performance. This is mainly due to the fact that deeper models may over-fit the data during training, resulting in a failure to generalise on test and validation data sets.

Table 1: Reproduced from Ciresan, Meier, Gambardella, and Schmidhuber (2010) - DNN architectures of varying size for classifying the MNIST data set, and the associated performance of each network.

Architecture (number of neurons in each layer)	Test Error Best Validation [%]	Best Test Error [%]	Simulation Time [min]	Weights [Millions]
1000, 500, 10	0.49	0.44	23.4	1.34
1500, 1000, 500, 10	0.46	0.40	44.2	3.26
2000, 1500, 1000, 500, 10	0.41	0.39	66.7	6.69
2500, 2000, 1500, 1000, 500, 10	0.35	0.32	114.5	12.11
$9 \times 1000, 10$	0.44	0.43	107.7	8.86

Fully connected feed-forward DNNs have proven effective in computer vision classification problems, such as optical character recognition. One of the most widely cited examples of this is a feed-forward DNN performing classification on the MNIST dataset. The MNIST dataset contains handwritten digits, from 0 to 9, and is considered a benchmark for measuring neural net classification performance for the optical character recognition problem. Table 1, taken from paper by Ciresan, Meier, Gambardella, and Schmidhuber (2010), shows a table of feed-forward DNNs with varying numbers of hidden layers, and hidden layer depth. Notably, every model listed presents an error rate of less than 1%. Generally, the deeper a network, the better the model’s predictive performance, although this is not always the case as previously outlined. Figure 5, taken from the same paper, shows a small set of the misclassified images. Despite the misclassification it can be seen that, in most cases, the handwritten digit bears a high resemblance to the predicted value, and that the second prediction is generally correct.

 1 ² 17	 1 ¹ 71	 9 ⁸ 98	 9 ⁹ 59	 9 ⁹ 79	 5 ⁵ 35	 8 ⁸ 23
 4 ⁹ 49	 5 ⁵ 35	 9 ⁴ 97	 9 ⁹ 49	 9 ⁴ 94	 0 ² 02	 3 ⁵ 35
 6 ⁶ 16	 9 ⁴ 94	 0 ⁰ 60	 6 ⁶ 06	 8 ⁶ 86	 1 ¹ 79	 1 ¹ 71
 9 ⁹ 49	 0 ⁰ 50	 5 ⁵ 35	 8 ⁸ 98	 9 ⁹ 79	 1 ⁷ 17	 1 ¹ 61
 2 ⁷ 27	 8 ⁸ 58	 2 ² 78	 6 ⁶ 16	 6 ⁵ 65	 9 ⁴ 94	 0 ⁰ 60

Figure 5: Misclassified hand written digits by the top performing DNN from Ciresan, Meier, Gambardella, and Schmidhuber (2010). The digit in the top right hand corner of the box is the observation label, and the two digits in the bottom right hand corner are the predictions from the DNN model.

Remarkably, DNNs are not considered state of the art for image classification problems - even with simple tasks like MNIST classification. This is due to model inefficiencies that arise from image variation in the spatial domain. The problem can be better understood by considering Figures 6

and 7. Figure 6 shows an image with a puppy on the left, and Figure shows an image with a puppy on the right. Suppose we create a simple model to classify whether an image has a puppy in it or not, and assume we train this model with lots of images like the one shown in Figure 6. If the trained model was then used to classify images like the one shown in Figure 7 it would perform poorly. This is because our model would have only learned to classify pixel features on the left side of the picture with puppies, which says nothing about identifying a puppy on the right hand side of an image. Put simply, there is no *translational invariance* in the model.

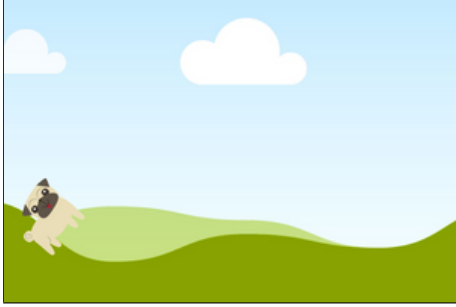


Figure 6: An image in which a puppy is located on the left hand side of the image.



Figure 7: An image in which a puppy is located on the right hand side of the image.

2.1 Convolutional Neural Networks

A Convolutional Neural Network (CNN) is a class of ANN, which has an underlying network structure which is better at learning shapes, edges, and colours meaning it is less reliant on the spatial location of a classification object in an image. Recall that vanilla feed-forward neural nets only have neuron connections from the previous layer, and there are no connections from neurons in the same layer - weights are not shared. In contrast, CNNs share neuron weights by using filters which are convolved over an input image. Consider a raw input image of say 32×32 pixels, with a depth of 3 colour channels, as shown in Figure 8. The convolving filter, which contains the model weights, is

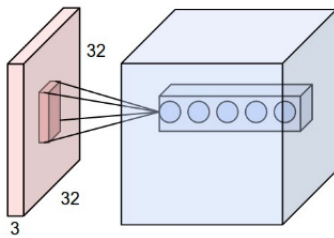


Figure 8: The raw image input is 32×32 , with a 3-channel depth (R,G,B). The filter is a 3×3 patch with the same depth as the input. The filter is convolved over the image using some stride - each convolution creates a single element output which forms part of the 2D activation map (i.e. the output). There are K filters convolved over the image, a parameter chosen as part of the architecture, and the output volume represents the stacked 2D activation maps.

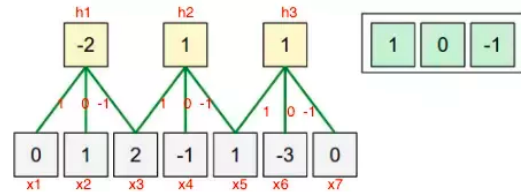


Figure 9: An simple example of weight sharing that takes place in a CNN. The white 1D array represents the input image, and the blue 1D array represents the filter. The filter (blue) is convolved across the input (white) using a stride of 2. The convolved output, which represents the activation map, can be seen in yellow. This architecture allows for the sharing of the weights in the model.

2.2 Fully Convolutional Neural Networks

FCN can be thought of in two distinct parts: encoders and decoders. This is shown in Figure XXXX. architectures avoid this problem by replacing the fully connected layers with 1×1 convolutional layers.

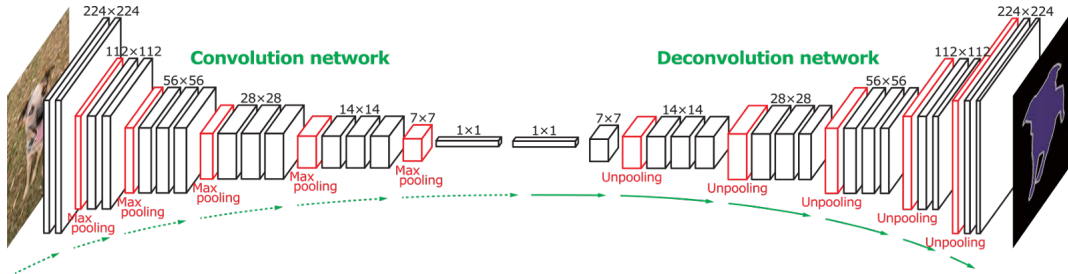


Figure 12: text

3 Network Training

3.1 Epochs

3.2 Learning Rate

3.3 Batch Size

4 Performance & Model Generalisation

5 Future Enhancements