

1 Deep Learning Theory

The following section covers the most important background theory for the experiments in Section ?? . This includes an introduction to various types of neural networks, as well as an introduction to the optimization of such networks.

1.1 Feedforward Neural Networks

Feedforward neural networks are the most basic type of neural networks. The aim of a feed-forward neural network is to approximate some function f^* , by defining a mapping $\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta})$ and learning the parameters $\boldsymbol{\theta}$, that results in the best approximation of f^* . These models are called **feedforward** because there are no **feedback** connections in which the outputs of the model are fed back into itself. Instead, information flows through the function being evaluated from \mathbf{x} , through the intermediate computations used to define f , and finally to the output \mathbf{y} . Feedforward neural networks generally consists of multiple **layers**, arranged in a chain structure, with each layer being a function of the layer that preceded it [2].

1.1.1 Fully-connected Layers

The most simple type of layer found in a feedforward neural network is the **fully-connected layer**. The fully-connected layer usually consists of some learnable parameter matrix \mathbf{W} and learnable parameter vector \mathbf{b} , as well as a non-linear **activation function** g (which will be covered further in Section 1.4.1). In this case, the i 'th layer is defined as [2]

$$\mathbf{h}^{(i)} = \begin{cases} g^{(i)} (\mathbf{W}^{(i)\top} \mathbf{h}^{(i)} + \mathbf{b}^{(i)}) & \text{if } i > 1 \\ g^{(1)} (\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) & \text{if } i = 1 \end{cases} \quad (1)$$

1.1.2 Convolutional Layer

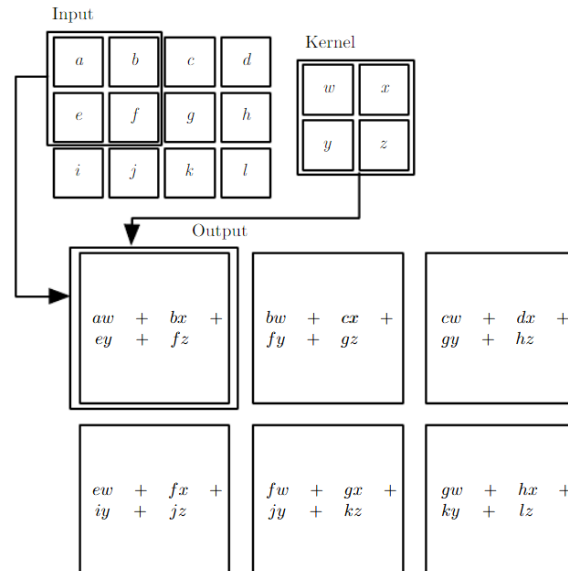


Figure 1: An example of applying a 2d kernel on an input [2].

A **convolutional layer** is a specialized kind of feedforward layer, usually used in analysis of time-series or image data [2]. If a network has at least one convolutional layer, it is called a **Convolutional neural network (CNN)**.

The convolutional layer consists of a set of **kernels**, each to be applied to the entire input vector, where each kernel is a learnable parameter matrix $k \times k$ [3]. Each kernel is applied on the input to produce a **feature map**. The kernels are applied to the input by "sliding" over the input (where the step size is called **stride**). Each $k \times k$ grid of the input is then used to compute the dot-product between the grid and each kernel, which is then placed in the corresponding feature map of each kernel, as visualized in Figure 1. [1]. To control the dimensions of the output, one might **pad** the sides with a constant value. Commonly, zero is used as the padding-value.

As seen in Figure 1, each kernel produces a linear combination of all pixel values in a neighbourhood defined by the size of the kernel. Thus, unlike a fully-connected layer, a convolutional layer captures the high correlation between a pixel and its neighbours. Further, by limiting the size of the kernel, the network will use much fewer parameters, than if a fully-connected layer would be used instead [2].

1.2 Recurrent Neural Networks

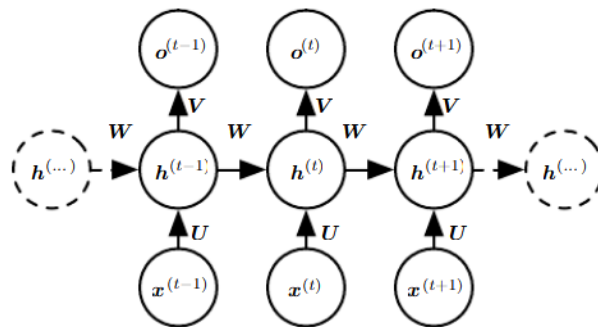


Figure 2: An illustration of an RNN [2].

Recurrent neural networks (RNNs) are a family of neural networks for processing sequential data. Figure 2 illustrates the general setup of such a network, which maps an input sequence of x values to a corresponding sequence of output o values. Generally, a RNN consists of three parts: (1) the input ($x^{(i)}$), (2) the hidden state ($h^{(i)}$), and (3) the output ($o^{(i)}$). During inference, the model maps each input value to an out value, in a sequential matter, where it first maps the first input value, then the second, then the third, and so forth. The network maps each input value to an output value by making use of the hidden state from the preceding step, where the first hidden state has to be initialized [2].

1.2.1 Convolutional Long Short-Term Memory

One common recurrent neural network unit is the **convolutional long short-term memory (ConvLSTM) cell**, which is an adaptation of the standard **long short-term memory (LSTM) cell** for sequences of images.

The idea of the ConvLSTM cell is to create paths through time that have derivatives that neither vanish nor explode. This is done by accumulating information over a long duration. Once that information has been used, the cell *forgets* the old state, which the ConvLSTM cell learns to decide when to clear. Each cell has the same inputs and outputs as an ordinary recurrent network, but also has more parameters and a system of gating units that controls the flow of information [2].

The ConvLSTM cell is defined as the following. Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_t$ be a sequence of input images, $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_t$ be a sequence of cell outputs, and $\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_t$ be a sequence of hidden states. Then, at each input step t , we compute [4]

$$\mathbf{i}_t = \sigma(\mathbf{W}_{xi} * \mathbf{X}_t + \mathbf{W}_{hi} * \mathbf{H}_{t-1} + \mathbf{W}_{ci} \circ \mathbf{C}_{t-1} + \mathbf{b}_i) \quad (2)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{xf} * \mathbf{X}_t + \mathbf{W}_{hf} * \mathbf{H}_{t-1} + \mathbf{W}_{cf} \circ \mathbf{C}_{t-1} + \mathbf{b}_f) \quad (3)$$

$$\mathbf{C}_t = \mathbf{f}_t \circ \mathbf{C}_{t-1} + \mathbf{i}_t \circ \tanh(\mathbf{W}_{xc} * \mathbf{X}_t + \mathbf{W}_{hc} * \mathbf{H}_{t-1} + \mathbf{b}_c) \quad (4)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{xo} * \mathbf{X}_t + \mathbf{W}_{ho} * \mathbf{H}_{t-1} + \mathbf{W}_{co} \circ \mathbf{C}_t + \mathbf{b}_o) \quad (5)$$

$$\mathbf{H}_t = \mathbf{o}_t \circ \tanh(\mathbf{C}_t) \quad (6)$$

1.2.2 Gated Recurrent Unit

1.3 Transformer

1.3.1 Encoder

1.3.2 Decoder

1.4 Training a Neural Network

1.4.1 Activation function