

Arnav Goyal

Data Engineering & Machine Learning

WINTER 2024 / AISE 3010

Preface

The purpose of this document is to act as a comprehensive note for my understanding on the subject matter. I may also use references aside from the lecture material to further develop my understanding, and these references will be listed here.

This document should eventually serve as a standalone reference for learning or review of the subject matter. There is also a lot of organization within these documents, please refer to the table of contents within your PDF viewer for ease of navigation.

References

- Provided Course Material & Lecture Notes
- Neural Networks and Learning Machines - 3e - S. Haykin

The Perceptron

In this chapter we will be going over the architecture of the most basic **neural network (NN)**, called the **perceptron**¹.

The Neuron

The most basic unit of a perceptron is called the **neuron**². The structure of a neuron can be seen in Figure 1, each neuron possesses the following.

- Some **input features** $x_1 \dots x_n$
- A corresponding **weight** for each feature $w_1 \dots w_n$
- A **bias** b
- An **activation function** $\varphi(\cdot)$
- An **output** y

A neuron essentially performs the following operation, which can be seen from the signal flow graph.

$$y = \varphi \left(\sum_{i=1}^n w_i x_i + b \right)$$

¹ Also called Rosenblatt's perceptron

² More specifically, the McCulloch-Pitts model of a neuron

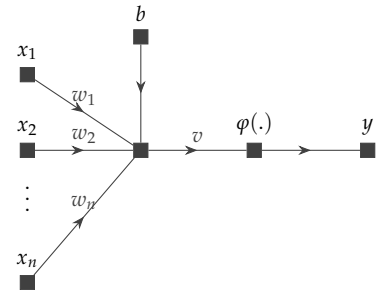


Figure 1: Signal Flow Graph of A Neuron

The Basic Perceptron

The perceptron was initially built with just a single neuron, and was limited to performing **binary classification** or **linear regression**.

- In binary classification, each training sample (X, y) contains some data $X = x_1, x_2, \dots$ and a label $y \in \{-1, +1\}$
- In linear regression the dependent variable is $y \in \mathbb{R}$

Overall, the perceptron was created to deal with binary classification problems.

- This means that both classes c_1 and c_2 need to be **linearly seperable**³
- The perceptron is essentially trying to find a **linear separator** between the two classes⁴, such that $\sum w_i x_i + b = 0$.

³ We need to be able to draw a hyperplane that clearly separates every instances from both classes, such that only instances from the same class are on a distinct side of the hyperplane

⁴ Given that the activation function, $\varphi(\cdot) \leftarrow \text{sgn}(\cdot)$

Multi-Layered Perceptrons

The most basic perceptron was basically a single-layered NN, but now we will discuss the **multi-layered perceptron (MLP)** which uses more layers, with possibly more neurons in each layer⁵.

Before I continue, lets clarify some terminology.

- The **input layer** is the first layer with neurons that each accept 1 feature. In total, an input layer with k neurons can accept k features.

⁵ We do this to overcome the perceptrons practical limitations

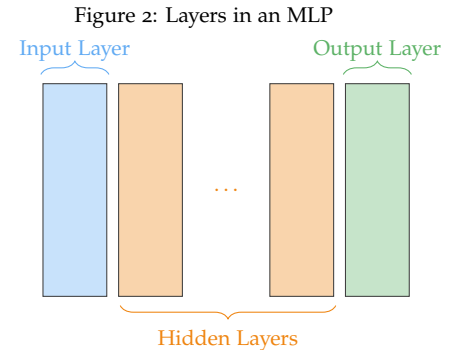
- The **output layer** is the last layer, containing neurons that spit out one output each. In total an output layer with h neurons produces an output with dimensionality of h .
- The **hidden layers** are the other layers sandwiched in between the input and output. When we say a "an n -layered network", we are referring to the amount of hidden layers in the network.

This is illustrated in Figure 2.

In order to construct a MLP we must ensure the following.

- The activation function $\varphi_i(.)$ of each neuron i is a **differentiable** and non-linear function.
- The network contains at least one layer that is **hidden** from the input and output nodes (hidden layers).
- The network exhibits a high degree of **connectivity** through synaptic weights of the network.

These conditions ensure that the network becomes complex, but also so complex that we struggle to understand the behavior of the network.



Batch & Iterative Learning

As we work our way up to understanding the training of the MLP we should account for the two methods of teaching the MLP, called **batch learning** and **iterative learning**⁶

- Batch learning consists of presenting a batch of training examples to the network at the same time
- Iterative learning consists of presenting one training example to the network at a time.

Batch learning has the advantages of providing an accurate-estimation of the gradient vector, and allowing a parallelization of the learning process. While iterative learning is simple to implement and can utilize redundancies in the dataset.

⁶ In the textbook this is called *on-line learning* but this was obviously confusing so I changed it for this note.

Backpropagation

Now we are going to dive into the complex algorithm called **backpropagation** which is the way modern NNs are trained. There are two cases for this algorithm

- When a neuron is an output neuron
- When a neuron isn't an output neuron

We will be starting with the former, but be warned! this is going to be a painful and math-intensive chapter.

The Error & Cost Function

Consider a neuron j that is located in the output layer of an MLP. It has internal structure as shown by Figure 3. We also define the **error** ε , that gives a measure of how far our current output \hat{y} , is from our desired output y . For the purposes of this chapter:

$$\varepsilon = y - \hat{y}$$

We should also define a **cost function** ℓ as the **total error energy** of the whole network.⁷

$$\ell = \frac{1}{2} \varepsilon^2$$

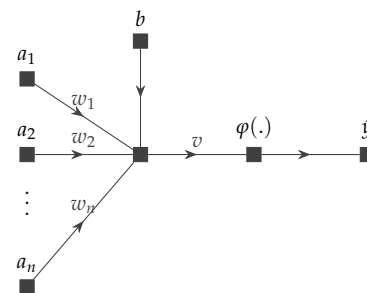


Figure 3: An Output Neuron

⁷ This is usually defined at some iteration n but this is left out for brevity.

The Output Neuron

We are concerned with finding out how the loss function ℓ changes with respect to a certain weight w_k . In other words we are tasked with finding the derivative.⁸ We can do this through the **chain rule**.

⁸ We are tasked with finding $\frac{\partial \ell}{\partial w_k}$

$$\frac{\partial \ell}{\partial w_k} = \frac{\partial \ell}{\partial \varepsilon} \cdot \frac{\partial \varepsilon}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_k}$$

We can break down this derivative even further using Figure 3, resulting in Equation 1. We can even extend this to see how we can change the biases, which is given in Equation 2.

$$\frac{\partial \ell}{\partial w_k} = \frac{\partial \ell}{\partial \varepsilon} \cdot \frac{\partial \varepsilon}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial v} \cdot \frac{\partial v}{\partial w_k} \quad (1)$$

$$\frac{\partial \ell}{\partial b_k} = \frac{\partial \ell}{\partial \varepsilon} \cdot \frac{\partial \varepsilon}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial v} \cdot \frac{\partial v}{\partial b_k} \quad (2)$$

Analytically solving Equation 1 gives the following result:

$$\frac{\partial \ell}{\partial w_k} = -\varepsilon a_k \varphi'(v)$$

Convolutional Neural Networks

The task of image classification is quite challenging given the structure of a MLP with **fully connected layers**. So far, these are the only layer types we have learned about. This chapter will start by introducing different layer types, their function, and how they are assembled to create a more efficient **image classifier**.

The CNN

A **Convolutional Neural Network (CNN)**, is a more efficient way to perform an image-classification task. Take a moment to consider the differences between an image as an input and a list of numbers. An image has a sense of *locality*, the pixels that make up a cat are usually together in a sense, and when we convert them into a list we lose that *local* information.

The CNN fixes this by training convolutional image filters (called kernels) to detect certain parts of the image, and eventually classify it.

The Conv Layer

The most basic building block of a **Convolutional Neural Network (CNN)** is the convolutional (conv) layer. A convolutional layer is defined by:

- Number of filters at this layer, f
- A **kernel** size⁹, $k \times k$
- A **stride**, s

⁹ We are considering a 2d kernel here

The kernel is the size of the convolutional filter, it's weights are what we train in this layer. The stride is how many steps to skip in each filter movement, in most conv layers it is set to 1. The number of filters determines the output dimensionality.

Consider a 32×32 RGB image, it would be represented with tensor¹⁰ of dimensionality: $32 \times 32 \times 3$. Suppose we have a conv layer defined generally above, We would get an output tensor of:

¹⁰ the last number is referred to as the number of **channels**

- $\frac{32}{s} \times \frac{32}{s} \times f$

Essentially we have converted a 3-channel tensor into a f channel tensor, this is how pytorch deals with this as well, by defining `in_channels=?` and `out_channels=?`

The Max-Pooling Layer