



ARISTOTLE UNIVERSITY OF THESSALONIKI

Thesis Title

by

Iakovos Marios Tsouros

A thesis submitted in partial fulfillment for the
Graduate degree

in the
Faculty of Sciences
School of Physics

Supervising Professor: Dr. Panagiotis Argyrakis

Date

Declaration of Authorship

I, [author's name], declare that this thesis titled, [thesis title] and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

Inspiring quote goes here (optional)

Quote's attribution

ARISTOTLE UNIVERSITY OF THESSALONIKI

Abstract

Faculty of Sciences

School of Physics

Graduate Degree

by Iakovos Marios Tsouros

Abstract goes here.

Acknowledgements

Acknowledgements go here.

Contents

| | |
|---|-------------|
| Declaration of Authorship | i |
| Abstract | iii |
| Acknowledgements | iv |
| List of Figures | vii |
| List of Tables | viii |
| Abbreviations | ix |
| | |
| 1 Introduction | 1 |
| 1.1 Graphs | 1 |
| 1.1.1 Introduction | 1 |
| 1.1.2 Adjacency Matrix | 4 |
| 1.1.2.1 Adjacency List | 6 |
| 1.1.3 Graph Laplacian | 6 |
| 1.1.4 Edge weights | 7 |
| 1.1.5 Distance between nodes and shortest paths | 8 |
| 1.1.6 Node and edge properties | 9 |
| 1.2 Network Dynamics | 11 |
| | |
| 2 Neural Networks | 12 |
| 2.1 Historical Background | 12 |
| 2.1.1 Basics | 13 |
| 2.1.1.1 Summary | 13 |
| 2.1.1.2 Building Blocks | 14 |
| 2.2 Mathematics useful in ANNs | 15 |
| 2.2.1 Gradient Descent | 15 |
| 2.2.2 Hadamard Product | 16 |
| 2.3 Feedforward Networks | 16 |
| 2.3.1 Perceptrons, learning algorithms and simple NNs | 17 |
| 2.3.2 Backpropagation | 22 |
| 2.3.3 Common ANN Architectures | 26 |

| | | |
|-------------------------|---|---------------|
| 2.3.3.1 | Convolutional Neural Networks (CNNs) | 26 |
| 3 | Notation & Fundamentals | 28 |
| 4 | Basic Principles and Implementation Framework for an [Problem to be Solved | 29 |
| 5 | Implementation and Core Components of [Platform Title | 30 |
| 6 | Experimentation & Validation | 31 |
| 7 | Conclusions & Future Work | 32 |
| Bibliography | | 33 |

List of Figures

| | | |
|-----|--|----|
| 1.1 | An undirected pseudograph with labeled nodes and edges. | 2 |
| 1.2 | Two undirected multigraphs. | 4 |
| 1.3 | Simple example of an unordered graph with weighted edges | 7 |
| 1.4 | A graph with a maximum path of 3 (nodes 1 to 6). | 8 |
| 1.5 | Example graph of a small classroom with labeled edges and nodes | 10 |
| 2.1 | A simple neural network demonstrating the layered structure and and flow of data from input to output. | 13 |
| 2.2 | Gradient descent in three dimensional space | 15 |
| 2.3 | A simple perceptron | 17 |
| 2.4 | Plot of the sigmoid function, its first derivative and the Heaviside step function. | 19 |
| 2.5 | Input to a single neuron in a feedforward network. Here σ represents the activation function (a sigmoid function in this case) and the exponents represent layers. The activation of a layer can be conveniently represented in matrix form. Biases are added to the input of the node. | 20 |
| 2.6 | Noisy convergence of SGD compared to non-stochastic. Image source [1]. | 21 |
| 2.7 | Illustration of a convolutional layer with multiple feature maps. | 27 |
| 2.8 | Illustration of the pooling layer following the feature maps | 27 |

List of Tables

| | | |
|-----|---|----|
| 1.1 | Adjacency matrix for Figure 1.3a | 5 |
| 1.2 | Adjacency matrix for Figure 1.2b | 5 |
| 1.3 | Adjacency list for Figure 1.2b | 6 |
| 1.4 | A) Node properties B) Edge properties | 10 |
| 1.5 | Graph Properties | 10 |

Abbreviations

| Acronym | What (it) Stands For |
|---------|----------------------|
|---------|----------------------|

Dedication (optional)

Chapter 1

Introduction

1.1 Graphs

In this subsection, the main aspects of graph theory are briefly presented.

1.1.1 Introduction

In the real world, many problems can be described by a diagram connecting a set of points with lines, joining pairs of these points, or even creating loops on a single point. A simple example of that would be a set of points representing people with lines connecting acquaintances, or points representing atoms and lines representing chemical bonds, creating a representation of a molecule as a graph attribute. In the examples above, the only information contained is whether two points are associated, with the manner being disregarded. The concept of a graph consists of a mathematical abstraction of the above. [2]

Definition 1.1. Mathematically, in its simplest form, a **graph** is an ordered pair¹ $G = (V, E)$ of:

- V , a set of vertices (also known as nodes).
- $E \subseteq \{\{x, y\} | x, y \in V, x \neq y\}$, which is the set of **edges** which consists of unordered pairs of vertices that connect two nodes.

This type of object is called an **undirected simple graph** to avoid confusion with other types.

¹An ordered pair (a, b) is a pair of objects in which the order of appearance or insertion is significant; the ordered pair (a, b) is different than (b, a) unless $a = b$. An unordered pair is a set of the form a, b is a set having two elements with no relation between them and $a, b = b, a$.

Definition 1.2. A *graph* G is an ordered pair $(V(G), E(G))$ consisting of a set $V(G)$ of *vertices* (also called *nodes* or *points*) and a set $E(G)$, disjoint from $V(G)$ which consists of *edges* (also called *links* or *lines*) together with an incidence function ψ_G that associates with each edge of G an unordered pair of not necessarily distinct vertices of G . If e is an edge and u and v are vertices such that $\psi_G = u, v$ then e is said to *join* u and v , and the vertices u and v are called the *ends* of e . We denote the numbers of vertices and edges G by $u(G)$ and $e(G)$ which two parameters are called the *order* and *size* of G , respectively [2].

In short, we can define a **graph** as an ordered triple $G = (V, E, \phi_G)$ consisting of:

- V , a set of *vertices*
- E , a set of *edges*
- $\phi_G : E \rightarrow \{\{x, y\} | x, y \in V \text{ and } x \neq y\}$ an *incidence function* mapping every edge to an unordered pair of vertices - an edge associated with two distinct vertices. The incidence function is a function of the edges.

This type of object is called an *undirected multigraph*, to avoid confusion. Note, that the above definition of the *incidence function* does not allow for *loops* (mappings of an edge on the same vertex).

A *loop* is a an edge that allows a connection of a vertex to itself and a graph can be defined to either allow or disallow the presence of loops. Some authors allow for loops to exist on *multigraphs* [3], while other consider these kind of graphs to exist in a different category, called *pseudographs* [4]. Allowing loops requires modifying the incidence function so they can be supported. The new incidence function can be written as:

$$\phi_G : E \rightarrow \{\{x, y\} | x, y \in V\} \quad (1.1)$$

The example presented below should better illustrate clarify the definition (of a pseudograph).

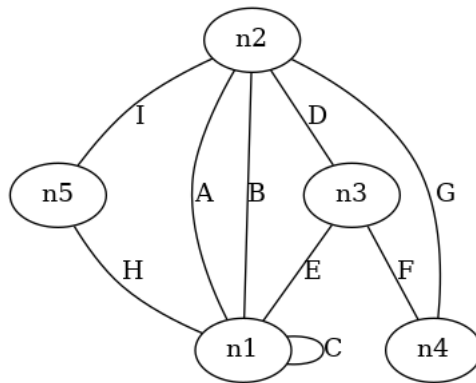


FIGURE 1.1: An undirected pseudograph with labeled nodes and edges.

Example 1.1.

For the graph presented in Figure 1.1 the following can be assumed:

$$G = (V(G), E(G))$$

and

$$V(G) = \{n_1, n_2, n_3, n_4, n_5\}$$

$$E(G) = \{A, B, C, D, E, F, G, H, I\}$$

and the incidence function is defined as:

$$\begin{aligned} \psi_G(A) &= n_1n_2, & \psi_G(B) &= n_1n_2, & \psi_G(C) &= n_1n_1, & \psi_G(D) &= n_2n_3, \\ \psi_G(E) &= n_1n_3, & \psi_G(F) &= n_3n_4, & \psi_G(G) &= n_2n_4, & \psi_G(H) &= n_1n_5, \\ \psi_G(I) &= n_2n_5 \end{aligned}$$

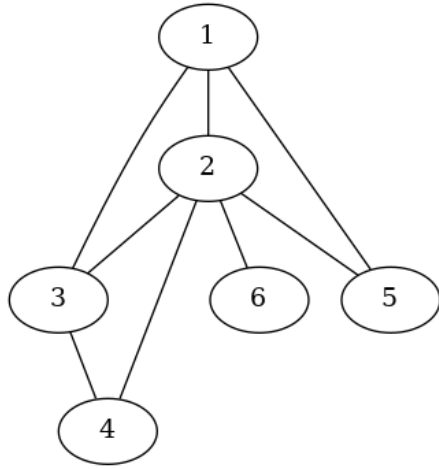
It should now be clear that with the newer definition of ϕ_G , self loops are now possible. Additionally, even though this was not prohibited by the previous definition, it is worth noting that a node can be connected to another with multiple edges (or multiedges), or that it can have zero connections to other nodes. Generally, V is assumed to be a non-empty set, but E can be empty.

It is now possible to define some characteristic attributes of graphs:

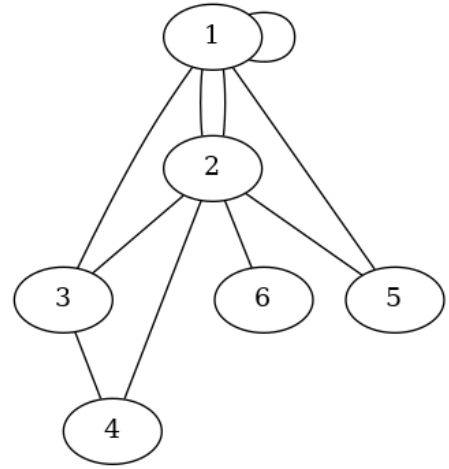
- $|V|$: the **order** of a graph is the number of its vertices.
- $|E|$: the **size** of a graph is the number of its edges.
- The **degree** (or **valency**) of a single node is the number of edges connected to it. The **degree** of a graph is the maximum number of edges connected to a single vertex that belongs to it.
- The edges of create a *homogenous relation*² \sim on the vertices of the graph that is called **adjacency relation**; for each edge (x, y) , its endpoints x, y are said to be **adjacent** to each other, denoted by $x \sim y$. This property will be particularly useful when the adjacency matrix is defined in the following section.

It can be inferred from the above definitions and attributes that for an undirected graph of order n , the maximum *degree* of a node is $n - 1$ and the maximum *size* of a graph is $n(n - 1)/2$.

²A **homogenous relation** (or **endorelation**) over a set X is a set of assignments (binary relations) over X and itself; i.e. it is a subset of the cartesian product $X \times X$



(A) Multigraph with no loops and multiple edges.



(B) Mutligraph with loops and multiple edges.

FIGURE 1.2: Two undirected multigraphs.

In this section only *undirected* graphs were considered, which are graphs with edges with no orientation. A whole other class of graph objects with edges which have orientation exists, called *directed graphs*. These kind of graphs objects are out of scope for this thesis and will not be presented.

1.1.2 Adjacency Matrix

Definition 1.3. The *adjacency matrix* is the fundamental mathematical representation of a graph. It is a square matrix, the elements of which represent which pair of nodes are *adjacent* or not. Thus, the adjacency matrix \mathbf{A} of a graph of order n is the $n \times n$ matrix with elements A_{ij} such that:

$$A_{ij} = \begin{cases} 1 & \text{if there exists at least one edge connecting } i \text{ and } j \\ 0 & \text{if there no edges connecting those edges directly.} \end{cases} \quad (1.2)$$

Considering the simple undirected graph of Figure 1.3a we can construct the following adjacency matrix:

$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

TABLE 1.1: Adjacency matrix for Figure 1.3a

For this simple network, which has no loops and only one edge connect two nodes, the diagonal matrix elements are always zero and the matrix is symmetric, as for each edge connecting i and j there is a representation for the j to i connection as well.

In a more complex case, such as the one presented in Figure 1.2b where loops and multiedges are present an adjacency matrix can still be constructed. In this case, a multiedge is represented by setting the value of the corresponding A_{ij} value equal to the multiplicity of the edge. In this case, $A_{12} = A_{21} = 2$.

For loops, the most common representation in the case of undirected graphs is to still set the value of the A_{ii} element equal to 2 (i.e. $A_{11} = 2$ in the example presented). Essentially, an edge of a loop has two ends that connect to the same node, thus the result [5, p. 68]. Additionally, defining the matrix in such manner, allows for better computations and is consistent with the definition of the representation of an edge connecting two nodes of an undirected graph [6, p. 108].

Thus, the adjacency matrix for the graph of Figure 1.2b is

$$A = \begin{pmatrix} 2 & 2 & 1 & 0 & 1 & 0 \\ 2 & 0 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

TABLE 1.2: Adjacency matrix for Figure 1.2b

Remark 1.4. Note that the degree of a node can be easily found by summing the values of the column or row of the adjacency matrix that correspond to said node.

1.1.2.1 Adjacency List

An alternative to the adjacency matrix is the *adjacency list*. An adjacency list is a collection of lists, one for each node i . Each list contains the labels of the nodes that i is connected to, and it is the most common method for storing networks on computers as it requires less space. It is also possible to represent edge attributes in an adjacency list by appending an extra column which holds these values. An example for the graph presented in Figure 1.2b with multiedges and loops:

| Node | Neighbors |
|------|-------------|
| 1 | 1,2,2,5,3 |
| 2 | 1,1,3,5,6,4 |
| 3 | 1,2,4 |
| 4 | 3,2 |
| 5 | 1,2 |
| 6 | 2 |

TABLE 1.3: Adjacency list for Figure 1.2b

Each edge of the network appears twice, thus for a network with m edges the size of the adjacency list would be $2m$, much smaller compared to the $n \times n$ matrix required to build an adjacency matrix. This is particularly useful in networks which are relatively *sparse*³, but have a high order.

1.1.3 Graph Laplacian

The graph laplacian is another representation matrix representation of a graph. In its simplest form, for a simple, undirected and unweighted network of order n , is a $n \times n$ matrix \mathbf{L} with elements:

$$L_{ij} = \begin{cases} k_i, & \text{if } i = j \\ -1, & \text{if } i \neq j \text{ and there exists an edge connecting } i \text{ and } j \\ 0, & \text{everywhere else} \end{cases}$$

where k_i is the degree of the node. Another way to write the graph Laplacian is as

$$\mathbf{L} = \mathbf{D} - \mathbf{A}$$

where \mathbf{D} is a diagonal matrix containing the degrees of the nodes.

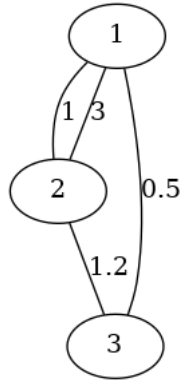
³*Sparse* networks are networks with a much lower number of edges than those possible.

In a similar manner the Laplacian matrix can be constructed for weighted networks, by replacing the the degree k_i of a node with the relevant matrix elements.

The Laplacian matrix has many applications in the study of dynamical systems, random walks and graph visualization. It has also found applications in graph neural networks, as its spectral decomposition allows the construction of low dimensional embeddings with applications in graph neural networks, such as ChebNet [7] which will be discussed in depth later.

1.1.4 Edge weights

So far, while presenting edges, we have considered graphs where connections between nodes represented binary relations between them; they either existed or they did not. In many situations when studying graphs, it is useful to represent edges as connections which carry some kind of attribute or value, commonly called *weight*. This weight could be any real number that fits a particular example, such as the distance between two airfields on an airline network, the kinship of connections on a social network (negative values can represent animosity and vice versa) or any type of relational attribute that can be quantified and characterizes the connection between nodes that belong in the same network[6, p. 109]. A simple example is presented in the figure below.



(A) Multigraph with no loops and multiple edges.

$$A = \begin{pmatrix} 0 & 4 & 0.5 \\ 4 & 0 & 1.2 \\ 0.5 & 1.2 & 0 \end{pmatrix}$$

(B) Corresponding adjacency matrix.

FIGURE 1.3: Simple example of an unordered graph with weighted edges

Generally, edges and nodes can hold any type of variable as values, such as vectors, the usefulness of which will become apparent when computer algorithms for graph representations and graph neural networks are discussed in later sections and chapters.

1.1.5 Distance between nodes and shortest paths

On a graph, any route that traverses nodes along the edges connecting them creates a sequence which is called a *walk*. Walks are not prohibited from traversing previously visited nodes and edges, but walks that do not intersect themselves are called *paths*.

Remark 1.5. Adjacency Matrix Powers Before continuing, it is worth mentioning that the powers of the adjacency matrix A^c directly provide the number of walks of length c among two nodes. If there is a connection between two nodes i and j , then $A_{ij} = 1$ else it is 0. Moving to the second power and taking for example an intermediate node k which might lie between i and j , the product $A_{ik}A_{kj}$ would be 1 if there is a node and 0 if there is not [6, p. 131]. Generalizing to walks of length c which traverse nodes i to j , their total number is:

$$N_{ij}^{(c)} = [A^c]_{ij}$$

The *shortest path* between two nodes is the shortest walk between those two nodes, the walk which traverses the least amount of nodes. In terms of edges, the least number that must be traversed is called *shortest distance* or often just “distance”. Mathematically, the shortest distance between two nodes i and j is the walk with the smallest value c where $[A^c]_{ij} > 0$

Example 1.2. For instance consider the following graph:

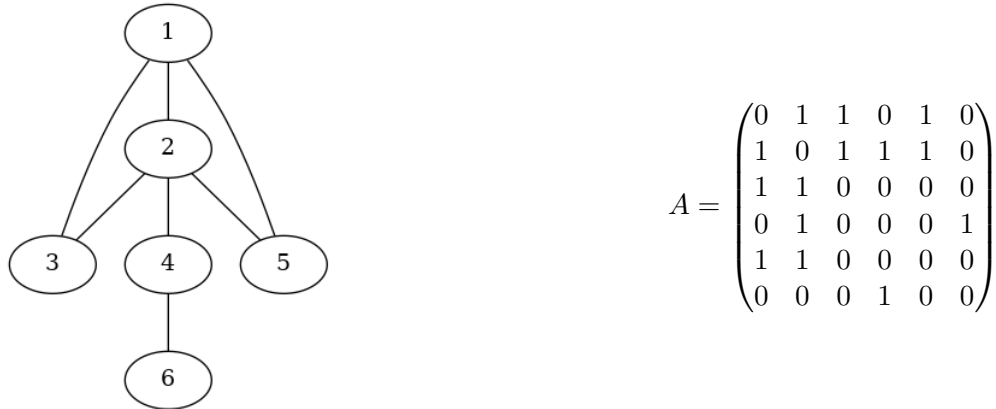


FIGURE 1.4: A graph with a maximum path of 3 (nodes 1 to 6).

In this example, it can be visually determined that the minimum walk between nodes 1 and 6 is of length 3. Indeed, raising the adjacency matrix to a power of three yields for these nodes:

...

$$N_{1,6}^{(2)} = \sum_{k=1}^n A_{1k}A_{k6} = [A^2]_{1,6} = 0$$

$$N_{1,6}^{(3)} = \sum_{k=1}^n A_{1k}A_{kl}A_{l6} = [A^3]_{1,6} = 1$$

1.1.6 Node and edge properties

As mentioned in Section 1.1.4, edges can hold more information than just the binary relations between nodes. In fact, this concept can be generalized to nodes and even whole graphs. When studying graphs and networks, especially when using computational methods for applications like complex dynamics, it is the most natural way to phrase data on a graph. The data can be any real number or even categorical data, such as colors.

The matrices which hold the information mentioned before are

- **V: vertex (or node) attributes** (i.e. a label or number of neighbors)
- **E: edge (or link) attributes** (i.e. a label or edge weight)
- **U: global (or master node) attributes** (i.e. number of nodes or number of paths of length 2)

Global attributes usually get their values as an aggregation of the attributes of the nodes and edges, and methods applied on them. For instance, a molecule might have chemical elements as node attributes, types of bonds as edge attributes and the toxicity of the molecule as graph attributes. An example of a classroom should better illustrate the concept.

Example 1.3. In this example, a classroom of 5 classmates is presented. Each student has a number of **node** attributes, their age, height and average grade (from F to A). **Edge** attributes between students hold information about their physical proximity when seated for class, and their kinship (as a real number between 0 – 1). Finally, **global** attributes consist of information about the classroom, such as total number of students and class's failure rate.

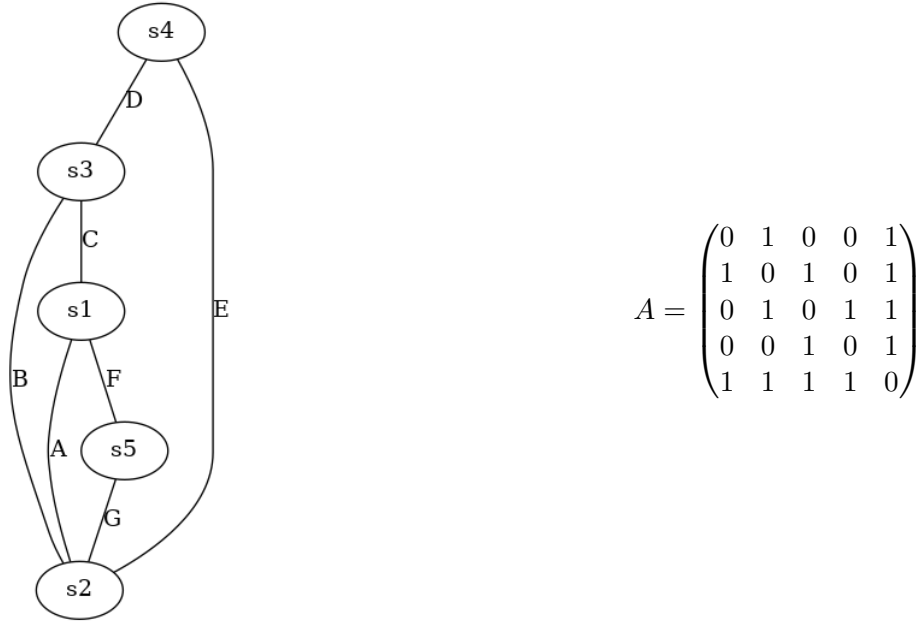


FIGURE 1.5: Example graph of a small classroom with labeled edges and nodes

| Node | Age | Height | Grade |
|------|------|--------|-------|
| s1 | 11.5 | 135 | C |
| s2 | 12 | 140 | B |
| s3 | 12 | 142 | A |
| s4 | 11.5 | 132 | A |
| s5 | 12 | 143 | B |

(A)

| Edge | Distance | Kinship |
|------|----------|---------|
| A | 1.5 | 0.5 |
| B | 1.2 | 0.5 |
| C | 0.5 | 0.8 |
| D | 2 | 0.35 |
| E | 2 | 0.3 |
| F | 0.5 | 0.9 |

(B)

TABLE 1.4: A) Node properties B) Edge properties

| | No. Nodes | Failure % |
|---|-----------|-----------|
| G | 5 | 0 |

TABLE 1.5: Graph Properties

1.2 Network Dynamics

Frequently, it is useful to consider cases where the status of the networks studied changes over time. This could mean that the topology of the network changes (i.e. nodes and edges are added or removed), the internal state of the network changes (the properties discussed in the previous section) or both.

The main concern of this thesis is with networks with a fixed topology, but with elements (nodes and edges) whose properties constitute dynamical quantities which can change over time. Following some dynamical rule, nodes can interact with their neighbors or change their own properties, with edges dictating which interactions are possible. In fact, the study of network dynamics combines graph theory with non-linear dynamics [8].

For many real world situations, a proper model of their processes consists of dynamical systems acting on networks. This can range from opinion spreading, epidemics, flow of electricity on grids, spread of packages on internet networks and many more systems whose dynamics are evolving on a network. In fact, many network processes can be understood

Chapter 2

Neural Networks

2.1 Historical Background

Artificial Neural Networks (ANN), or sometimes simply called **neural networks** is a class of computational models that mimic the way biological neural networks work, such as the human brain. Interest on the subject sparked after the seminal paper "*A Logical Calculus of the Ideas Immanent in Nervous Activity*" [9] by Warren McCulloch and Walter Pitts, where they proposed a computationally functional model of neural networks. Their suggestions showed that in principle, any function a digital computer can compute, a neural net should too. The models they described had weights and thresholds, but they lacked a training method.

The suggestions of McCulloch and Pitts lead Frank Rosenblatt to create the *Perceptron* in 1958 [10], a binary classifier algorithm based on supervised learning¹. Although initially promising, single layer perceptrons were not able to train on multiple classes of patterns and were eventually proven incapable of learning a XOR function² in the book *Perceptrons* [11], as the way they worked was by "separating" data linearly. This lead to a stagnation in machine learning research dubbed "AI winter", until the proposal of **backpropagation**³ by Paul John Werbos in 1975 [12].

A renewed interest in the field lead to the development of the Cresceptron [13] in 1992, a method of training large networks with pooling layers (**max-pooling**) and down-sampling. GPU⁴ usage made possible the training of larger networks, while new types of networks

¹Supervised learning is a machine learning training technique that optimizes a model based on examples input-output pairs.

²XOR (Exclusive or, $x \oplus y$) is a logical operation that is true only if its arguments differ.

³Backpropagation is a method of fine tuning a neural network based on the error rate obtained from previous runs of the program. It will be discussed in detail later in this thesis.

⁴GPU - Graphics Processing Units is a specialized electronic circuit, a central part of modern computers which excels in efficient computation of algorithms which process large blocks of data parallelly, thus exceling in machine learning applications.

emerged such as the **Recurrent Neural Networks (RNNs)**. **Convolutional Neural Networks** have recently proven to be far superior for image classification tasks.

In recent years, neural

2.1.1 Basics

2.1.1.1 Summary

One can think of ANNs as a directed graph, with a collection of nodes which are densely connected (called **artificial neurons**), transmitting signals to each other. These nodes are usually organized in sets of layers, with signals moving in one direction (i.e. *feed forward networks*) through weighted connections. Signals received on a single neuron are real numbers, and the output of a single neuron is the output of an aggregation of a non-linear function of the sum of its inputs. This function is called an *activation function* and its results are propagated to all of the nodes outgoing connections. Thus, each neuron can be thought as a simple processing unit. The weighted connections between nodes might have an excitatory or inhibitory effect, based on these weights which can be positive, negative or very close to zero, having no effect [14, Chap. 1]. While training, example data is passed through the input layer and gets radically transformed through the layers until it reaches the output layer. The weights and other trainable parameters are then adjusted until the training data consistently reaches satisfactory results.

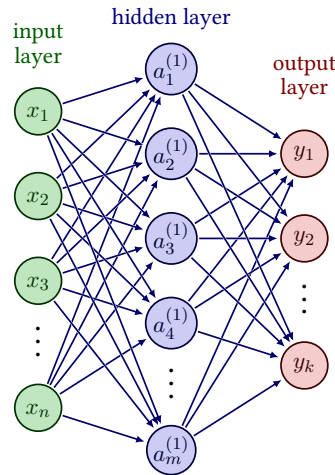


FIGURE 2.1: A simple neural network demonstrating the layered structure and flow of data from input to output.

2.1.1.2 Building Blocks

As mentioned before, the building blocks of ANNs are its artificial neurons organized into layers. In Figure 2.1 a basic ANN is presented, demonstrating the different layers that are typically present in a feedforward neural network ⁵ while the building blocks are widely considered [15] to be:

- **Input layer** This layer's purpose is to act as an entrypoint to the neural network and performs no computations. Data moves from this layer to the hidden layer succeeding it.
- **Hidden layer** Networks typically have one or more hidden layers, in which some computation takes place. Data moves from the input layer to a hidden layer where it is transformed and together with its weights moves to the next hidden layer or the output layer.
- **Output layer** Transformed data “exits” the network here, where it can pass through some function to reach the desired output format
- **Edges and Weights** Each node in a layer is connected to a set (usually all) of the nodes of the following layer with a weighted edge. Data from nodes $1 \dots n$ of layer k will be the input of node j of layer l , multiplied by a weight W_{ij} .
- **Activation Functions** An activation function takes as input some form of aggregate (usually the weighted sum) of the signals arriving at a node and produces an output. This function is typically nonlinear and differentiable for reasons which will be discussed later.
- **Learning** The learning process in a ANN involves modifying its weights and other learnable parameters to improve the accuracy of the result on the output layer. Learning usually involves a cost function which is evaluated on a predefined basis and adjustments are made accordingly. One of the most widespread learning techniques is *backpropagation*, where the error is propagated backwards through the network.

Along with the data from previous layers aggregated at a neuron, a bias is typically added which acts in the same way the intercept does in a linear equation. It adjusts the output of the activation function along with the weighted sum of the inputs to the neuron. It is also a trainable constant value provided to each node of a layer. Biases are node level parameters and do not depend on values provided by previous layers.

⁵Feedforward Neural Networks (FNNs) are the simplest type of neural networks, with the information moving only in one direction (“forward” through the layers), without any loops or cycles [15]. Different types of neural networks are discussed later.

2.2 Mathematics useful in ANNs

In this section some mathematical constructs which are widely used in machine learning and ANNs are briefly presented.

2.2.1 Gradient Descent

Gradient based optimizations methods are of great importance to NNs as they are most common method of training these models. They are tasked with minimizing or maximizing some function $f(x)$ which is oftenly called **objective function**. In NN training scenarios, where the goal is to minimize it the same function is also commonly called a **cost**, **loss** or **error** function[16].

Gradient descent is a method of minimizing a function $f(x)$ and finding a local minimum, given that its differentiable. The derivative $f'(x)$ is the slope of $f(x)$ at x , so it specifies how a change in x would provide a step towards the local minimum. For instance, for small values of ϵ , $f(x - \epsilon \text{sgn}(f'(x)))$ ⁶ is less than $f(x)$ and traversing the slope to ever decreasing values if possible through following the direction with the opposite sign of the derivate. This method is called **gradient descent** and it was first proposed by Cauchy [17] in 1847.

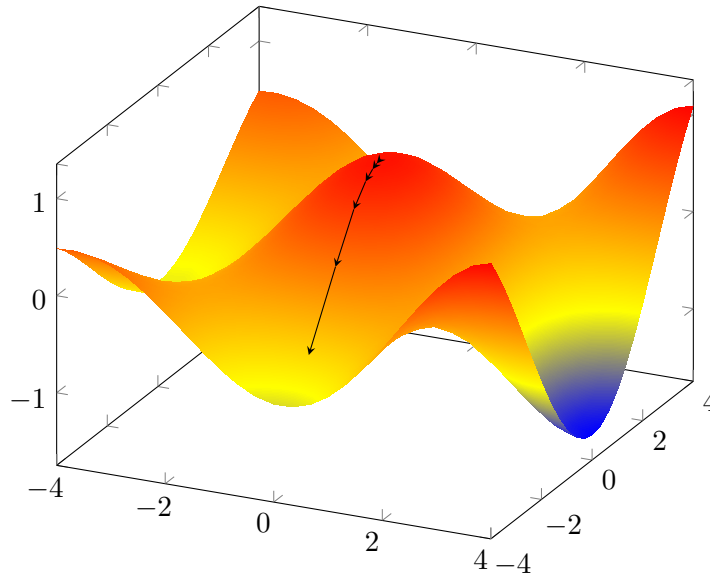


FIGURE 2.2: Gradient descent in three dimensional space

Definition 2.1. Consider a multi-variable function $F(x)$ which is **defined** and **differentiable** in a neighborhood of a point a . The value of $F(x)$ will decrease the fastest when moving from

⁶sgn is the *signum function*, a piecewise function which returns the sign of its input or 0 if input is 0. (i.e. $\text{sgn}(-1) = -1$ and $\text{sgn}(12) = 1$).

a towards the negative gradient direction of F , which is $-\nabla F(a)$. Thus when:

$$a_{n+1} = a_n - \gamma \nabla F(a_n) \quad (2.1)$$

then $F(a_n) \geq F(a_{n+1})$ and therefore the values of $F(a)$ move towards the local minimum. A sequence $x_0, x_1, x_2, \dots, x_m$ that follows the rule set by Equation 2.1 will lead to the monotonic sequence $F(x_0) \geq F(x_1) \geq F(x_2), \dots$ and the sequence x_n will converge to the local minimum. If some special choices are made for γ^7 the function is guaranteed to reach a local minimum. Additionally, if the function is *convex* local minima are global minima and gradient descent converges to a global solution.

For a function to be convex, a line segment connecting two of its points must lay on or above its curve. Mathematically for two points x_1 and x_2 , this can be expressed as

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2) \quad (2.2)$$

where λ is a location on a section line and $0 \leq \lambda \leq 1$.

2.2.2 Hadamard Product

The **Hadamard product**, also known as the **element-wise product** is a matrix operation in which two matrices of the same dimensions are multiplied to produce a matrix of equal dimensions, where its i, j elements are the product of elements i, j of the original two matrices. For two matrices A and B of dimensions $m \times n$ the Hamarand product can be written as:

$$(A \odot B)_{ij} = (A)_{ij}(B)_{ij} \quad (2.3)$$

2.3 Feedforward Networks

Feedforward neural networks (FNNs) are the exemplary of ANNs, and the first to be conceived and created by Rosenblatt in 1958 with the creation of the perceptron [10]. Their goal is to approximate some function f^* ; i.e. a classifier uses the function $y = f^*(x)$ to map the input x to some category y . The goal of a feedforward network would then be to define a mapping $y = f(x; \theta)$ and train in a way that the values of the parameter vector θ provide the best approximation of the function f^* .

⁷ γ in machine learning is also known as the “learning rate” and its one of the hypermaterees of NNs. Special choices made here include a selection of γ via line search which satisfies the Wolfe conditions or the Barzilai-Borwein method [18]

These networks are called feedforward as information flows from the input layer \mathbf{x} through intermediate computational layers used to define \mathbf{f} and finally to the output \mathbf{y} . There are no loops providing (called **feedback connections**) information from the output back into the input or other intermediate layers of the network.

Feedforward networks can be thought as a chain of functions, composing the final structure of the network. As an example, consider a network with three of these functions as $f(\mathbf{x}) = f^{(3)}(f^{(2)}(f^{(1)}(\mathbf{x})))$. In this case, the exponents denotes the layer, with $f^{(1)}$ being the first layer, $f^{(2)}$ the second etc. The total number of these functions is called the **depth** of the NN and the terminology “deep learning” arose from this layered structure.

During training the goal is to modify the parameters of the network in a way that $f(\mathbf{x})$ closely matches $f^*(\mathbf{x})$. The training set is composed of pairs of \mathbf{x} and labels $y \approx f^*(\mathbf{x})$ and the output of the network is evaluated at different training points. Training data defines the exact result expected from each input \mathbf{x} , a value as close as possible to y . The rest of the layers can have arbitrary behaviours as long as they transform the data in a way defined by the training goal. The learning algorithm can “use” them in any way that is useful to it, and the training data has no immediate effect on their behaviour. They are thus called “**hidden layers**” as they do not produce any meaningful result, related to the function [16, p. 160].

2.3.1 Perceptrons, learning algorithms and simple NNs

The perceptron serves as great precursory example to neural networks, as it introduces many concepts that are common in more complex NN paradigms. It is in fact a simple neuron, a computational unit, which takes a binary vector as input and produces a binary output.

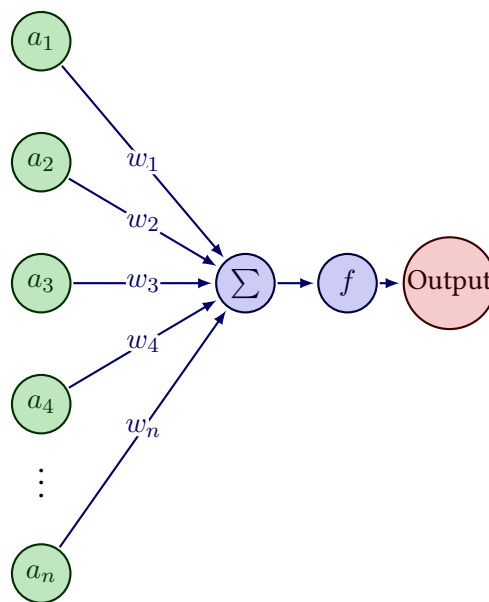


FIGURE 2.3: A simple perceptron

Inputs are multiplied with *weights*, and the output of the neuron is either 0 or 1, determined by whether the weighted sum of its inputs $\sum_i w_i a_i$ is greater than a threshold value *threshold*. All of these numbers are real numbers and a parameter of the neuron itself. The algebraic form of this can be written as:

$$\text{Output} = \begin{cases} 0 & \text{if } \sum_i w_i a_i \leq \text{threshold} \\ 1 & \text{if } \sum_i w_i a_i > \text{threshold} \end{cases}$$

By modifying the values of the threshold, the output can be changed. The above equation can be written in a more compact form by replacing the sums with the dot product of the weights and the input vector, $\sum_i w_i a_i = \mathbf{w} \cdot \mathbf{a}$. Moving the threshold to the left side of the equation and replacing it by what is referred to as the *bias* yields:

$$\text{Output} = \begin{cases} 0 & \text{if } \mathbf{w} \cdot \mathbf{a} + b \leq 0 \\ 1 & \text{if } \mathbf{w} \cdot \mathbf{a} + b > 0 \end{cases}$$

where $b = -\text{threshold}$. The bias can be seen as a measure of the tendency of the neuron to fire. Larger bias numbers makes the neuron easy to activate and output 1, while smaller ones require larger inputs and positive weights.

Remark 2.2. As mentioned before, perceptrons are able to solve problems which are linearly separable with the slope represented by the $\mathbf{w} \cdot \mathbf{a}$ term and the bias acting as the intercept. This excludes problems which are not, the most famous being the classic XOR gate, as there does not exist one single line capable of separating the predictions.

Multiple perceptrons can be combined in a network to compute any logical function, including the XOR gate. These kind of NNs are called *multilayer perceptrons*, and they are the basis modern artificial neural networks. As discussed before, these neurons (or nodes) are organized in layers and can have a depth based on the number of these layers. Typically, at least 3 layers are present.

A more practical activation function

So far, the way perceptrons transform data through a function has been described, but while they can produce sensible results as any other computing device, weights and biases had to be manually configured. It is possible to introduce a learning algorithm which does the tuning of these parameters automatically, in response to input-output (training data) pairs.

The perceptrons described used the Heaviside step function [19], often denoted with H . This function has a binary nature, with its output value being 0 or 1 based on a threshold, which presents a problem when fine tuning a perceptron or a NN based on them. Minor changes in a weight or bias value can completely alter the output and trigger a perceptron to flip, possibly changing the output of the whole network. A better choice for an activation function is the *sigmoid function* σ also called the *logistic function* which replaces the perceptron with the sigmoid neuron [20].

The sigmoid function is defined as:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{1}{1 + \exp(-\sum_j w_j x_j - b)} \quad (2.4)$$

It is easy to see the differences between these two functions when plotted:

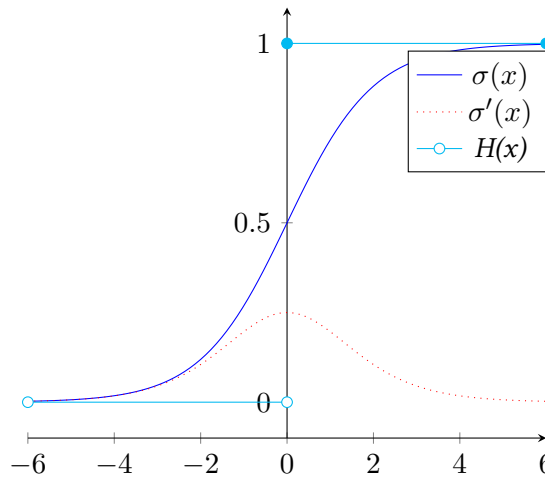


FIGURE 2.4: Plot of the sigmoid function, its first derivative and the Heaviside step function.

For very negative and very positive values of the inputs, the sigmoid approximates the behaviour of the step function, while for inputs around a neighborhood of 0 it differs by offering a smooth transition between 0 and 1 as an output. It also is a differentiable function. This properties allows an approximation of the output when the weights and biases change slightly as:

$$\Delta_{\text{output}} \approx \sum_j \frac{\partial \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \text{output}}{\partial b} \Delta b$$

This provides an estimation of how small changes in the biases and weights will affect the output of the sigmoid neuron, making choices for their values easy to calculate.

As stated before, all of the computations that take place on a neuron can be represented in a compact matrix form, as shown in Figure 2.5.

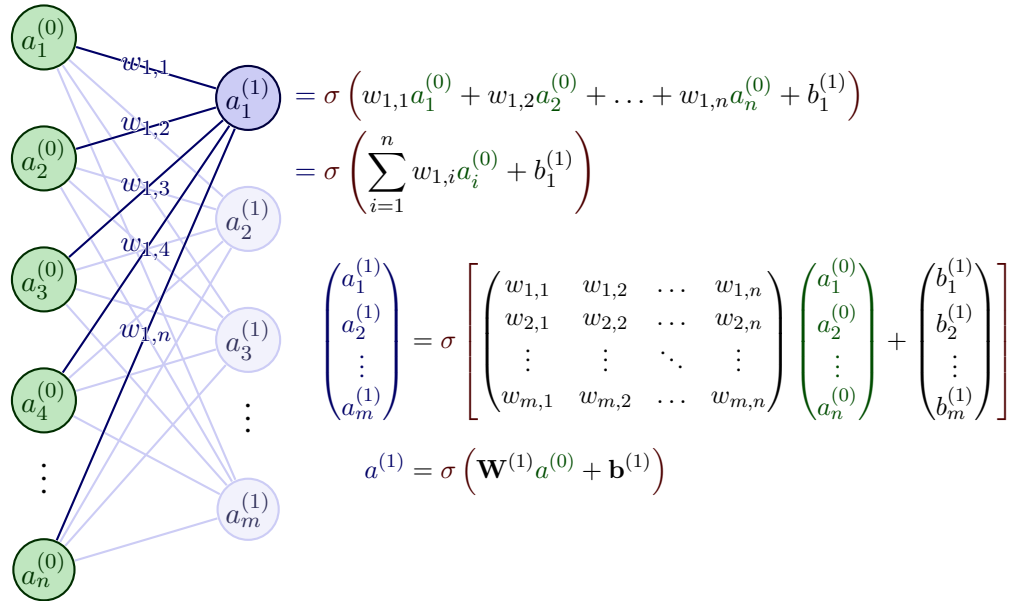


FIGURE 2.5: Input to a single neuron in a feedforward network. Here σ represents the activation function (a sigmoid function in this case) and the exponents represent layers. The activation of a layer can be conveniently represented in matrix form. Biases are added to the input of the node.

Learning

Before defining how training works for a neural network, it is imperative to describe an algorithm of quantifying far from the target output the network is with its current biases and weights. This algorithm is the *cost function* described in Section 2.2.1.

If for all inputs in a vector \mathbf{x} the desired output is known and is the output of some function $f^*(\mathbf{x})$ or $y(\mathbf{x})$ and the cost function can be defined as:

$$C(\mathbf{w}, \mathbf{b}) = \frac{1}{2n} \sum_x \|y(\mathbf{x}) - \mathbf{a}\|^2 \quad (2.5)$$

In this equation, \mathbf{w} and \mathbf{b} are the weights and biases of the network respectively and \mathbf{a} is a vector of the outputs when \mathbf{x} is the input.

This cost function is also known as *quadratic* or *mean squared error (MSE)* and it is one of the most commonly used loss functions. The goal of training is to minimize this function, and to achieve this gradient descent and variations of this method are used.

Applying gradient descent to Equation 2.5 yields:

$$\nabla C = \frac{1}{n} \sum_x \nabla C_x \quad (2.6)$$

where x are the training samples.

Stochastic Gradient Descent (SGD)

One of the most common variations is *Stochastic Gradient Descent (SGD)*. As training sets get larger, they become computationally expensive especially considering that in Equation 2.5, the total cost function $C = \frac{1}{n} \sum_x C_x$, averaging the costs of each training sample. The gradient of this function requires the computation of each ∇C_x for each sample and then averaging them as $\nabla C = \frac{1}{n} \sum_x \nabla C_x$ which can be prohibitive timewise.

SGD offers a solution to this problem, by drastically simplifying the above process; instead of finding the exact value of ∇C , it estimates its value by randomly picking one sample [21]. This stochastic approximation adds noise, but it has been proven to almost ensure convergence under mild conditions [22].

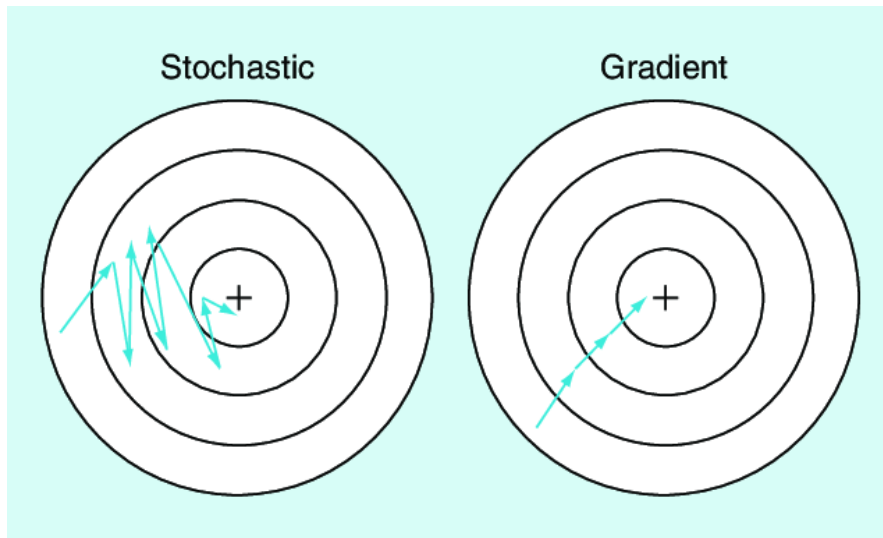


FIGURE 2.6: Noisy convergence of SGD compared to non-stochastic. Image source [1].

A smoother convergence by choosing a small number m of samples which is called a **mini-batch** and is the method most commonly used as it offers a compromise between speed and smoothness of convergence. In this case, Equation 2.6 becomes:

$$\nabla C \approx \frac{1}{m} \sum_{i=1}^m \nabla C_{X_j} \quad (2.7)$$

Training rules

We can then define a training rule for the network in terms of weights and biases and using the gradient of the cost function. Applying SGD on Equation 2.5 provides the values of weights w_k and biases b_l that minimizes it. An update rule for them can now be defined in terms of the

cost function:

$$\begin{aligned} w_k &\rightarrow w'_k = w_k - \frac{\eta}{m} \sum_j \frac{\partial C_{x_j}}{\partial w_k} \\ b_l &\rightarrow b'_l = b_l - \frac{\eta}{m} \sum_j \frac{\partial C_{x_j}}{\partial b_l} \end{aligned} \quad (2.8)$$

where η is a small positive value, the **learning rate**.

Having described the above, the process of training the NN can be summarized as:

1. The SGD is calculated for a random set number of training samples, the mini-batch.
2. Weights and biases are updated based on Equation ??.
3. The previous two steps are repeated until the training samples are exhausted.
4. When all samples are exhausted, the training **epoch** is complete, and the process starts again from the first step.
5. Training is usually said to be complete, when the value of the cost function is below a set threshold. The NN is now considered trained.

2.3.2 Backpropagation

Backpropagation is a method of computing the gradient of the loss function with respect to all the weights and biases of the network, based on single training examples, in contrast to more naive methods which compute the gradient on each weight or bias individually.

Even though these days the term *backpropagation* is sometimes used colloquially to refer to the whole process of the learning algorithm, strictly it is just an algorithm to compute the gradient, while SGD can be used for the training itself [16].

Historically, this method has been rediscovered multiple times, as early as 1960 [15] with its first implementation as software in 1970 [16, p. 229], albeit as a method for automatic differentiation with no mention to neural networks. In 1974, Werbos proposed its usage in training neural networks[23], but it was not until 1986 when Rumelhart, Hinton, and Williams [24] published their seminal paper “Learning representations by back-propagating errors” which experimentally showed that backpropagation can be used to train deep neural networks faster than any existing technique up to that point.

Ultimately the goal of backpropagation is to compute the partial derivatives $\frac{\partial C_X}{\partial w_{jk}^l}$ and $\frac{\partial C_X}{\partial b_j^l}$ — the partial derivatives of the cost function with respect to weights and biases in the network.

Matrix Notation

Before introducing the equations of backpropagation, it is important to present a matrix notation for quantities in a network which are used widely and unambiguously. Consider the equations presented in Figure 2.5. For a network with L layers, the connection of the k^{th} neuron in the $(l-1)^{th}$ layer with j^{th} node in l^{th} layer can be written as w_{jk}^l . In a similar fashion, the bias of the j^{th} node in the l^{th} layer can be represented as b_j^l and the *activation*⁸ of the same node is a_j^l . With this representations declared it is possible to define the matrix representation of a layer l as:

$$\mathbf{a}^l = \sigma(\mathbf{w}^l \mathbf{a}^{l-1} + \mathbf{b}^l) \quad (2.9)$$

where σ is the activation function. At this point the intermediate quantity $\mathbf{z}^l = \mathbf{w}^l \mathbf{a}^{l-1} + \mathbf{b}^l \rightarrow \mathbf{a}^l = \sigma(\mathbf{z}^l)$ commonly called the *weighted input* is worth declaring.

Fundamentals of Backpropagation

Backpropagation has two special requirements for the cost function. The first is that it must be able to be written as an average $C = \frac{1}{n} \sum_X C_X$ of the cost functions of individual training samples x . This ensures that after computing the partial derivatives $\frac{\partial C_X}{\partial w}$ and $\frac{\partial C_X}{\partial b}$ over the training set, it is possible to calculate $\frac{\partial C}{\partial w}$ and $\frac{\partial C}{\partial b}$ by averaging over the set.

The second requirement is that the cost function C can be written as a function of the outputs of the NN. For example, the MSE cost function described in Equation 2.5 for a single training example can be written as

$$C = \frac{1}{2} \|\mathbf{y} - \mathbf{a}^L\|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2 \quad (2.10)$$

which in fact can be regarded a function of the outputs \mathbf{a}^L alone, as the training examples $y(x)$ have fixed values. Note that L denotes the number of layers in the networks.

Equations of backpropagation

Before declaring the fundamental equations of backpropagation, an important intermediate quantity called the **error** δ_j^l must be declared. In a neuron j , during activation, a small quantity Δz_j^l is added as input and the neuron outputs $\sigma(z_j^l + \Delta z_j^l)$. This error is propagated through

⁸The activation of a neuron simply refers to the output of the activation function, taking as input the weighted sum of the inputs and the bias $a' = \sigma(wa + b)$

the network and finally the overall cost changes by $\frac{\partial C}{\partial z_j^l} \Delta z_j^l$. The error can be defined as:

$$\delta_j^l \equiv \frac{\partial C}{\partial z_j^l} \quad (2.11)$$

The fundamental equations which govern the backpropagation process are described below.

1. The first equation describes the error in the **output layer**:

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L) \quad (2.12)$$

or equivalently in matrix form:

$$\delta^L = \nabla_a C \odot \sigma'(z^L) \quad (2.13)$$

2. An equation for the error in terms of the error in the **next layer**:

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (2.14)$$

This equation is convenient as it propagates the error backwards in the network. Knowing the error in the $(l+1)^{th}$ layer and by transposing the weight matrix the error can then move backwards through the activation function in layer l . This process can be repeated all the way back.

3. An equation describing the rate of change of the cost with respect to **biases** in the network:

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \xrightarrow{2.14, 2.13} \frac{\partial C}{\partial b} = \delta \quad (2.15)$$

4. An equation describing the rate of change of the cost with respect to **weights** in the network:

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \quad (2.16)$$

Finally, Equation 2.15 and Equation 2.16 have been derived and can now be used in SGD or other learning algorithms.

The importance of backpropagation as a method arises from the fact that it is possible to compute all the partial derivatives of the cost function with just one forward and one backward pass through the network. Before its emergence, training a network required computing these derivatives for each weight and bias individually, which was prohibitively time consuming for larger networks.

Learning algorithm using backpropagation

The backpropagation algorithm can be represented in pseudocode as shown below.

Algorithm 1: Backpropagation Algorithm

Data: x the activation of the input layer a^1

Result: The gradient of C in terms of $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$ and $\frac{\partial C}{\partial b_j^l} = \delta_j^l$

begin

```

layers  $\leftarrow [2, 3, \dots, L]$ 
// Feedforward pass
for  $l$  in layers do
     $z^l \leftarrow w^l a^{l-1} + b^l$ 
     $a^l \leftarrow \sigma(z^l)$ 
end
// Output: Error vector
 $\delta^L \leftarrow \nabla_a C \odot \sigma'(z^L)$ 
// Backpropagate the error
layers_reverse  $\leftarrow [L-1, L-2, \dots, 2]$ 
for  $l$  in layers_reverse do
     $\delta^l \leftarrow ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$ 
end
// Output
 $\frac{\partial C}{\partial w_{jk}^l} \leftarrow a_k^{l-1} \delta_j^l$ 
 $\frac{\partial C}{\partial b_j^l} \leftarrow \delta_j^l$ 

```

end

Consequently, combining backpropagation and a learning algorithm such as a SGD, it is possible to train a network. For instance, in a SGD with m mini-batches case, the training rules, as described in Equation 2.8, can be represented with the following pseudo-code:

Algorithm 2: Training Rules using SGD with backpropagation

Data: Sets of training examples

Result: Updated weights and biases based on gradient descent of the cost function

begin

```

 $a^{x,l} \leftarrow$  set the activation for each sample
// Pass the activation to the backpropagation from Algorithm 1 and
get the error as output
 $\delta^{x,l} = ((w^{l+1})^T \delta^{x,l+1}) \odot \sigma'(z^{x,l})$ 
// Gradient Descent and updates for weights and biases
layers_reverse  $\leftarrow [L-1, L-2, \dots, 2]$ 
for  $l$  in layers_reverse do
     $w^{l'} \leftarrow w^l - \frac{\eta}{m} \sum_x \delta^{x,l} (a^{x,l-1})^T$ 
     $b^{l'} \leftarrow b^l - \frac{\eta}{m} \sum_x \delta^{x,l}$ 
end

```

end

2.3.3 Common ANN Architectures

So far, the models discussed were simple perceptrons and multi-layer perceptrons (MLPs), networks which are fully connected — each neuron in a layer is connected to all the neurons of the next layer (i.e. see Figure 2.1). This architecture does not take into account that data used as input might carry some spatial structure, as in the case of image recognition and computer vision. In this section, convolutional neural networks will be introduced briefly, which take into account the spatially local input patterns of some datasets, while at the same time reducing the number of free parameters allowing for deeper networks. While the origin of these types of networks can be traced to 1970, they were first established as a concept by LeCun in the paper “Gradient-based learning applied to document recognition” from Lecun et al. [25].

2.3.3.1 Convolutional Neural Networks (CNNs)

As mentioned before, CNNs take advantage of the spatial structure that some types of data have. For instance an image has a grid-topology, which is ignored in the case of simple or deep MLPs and as a consequence it inhibits learning by treating pixels which are spatially far, in the same manner. CNNs architecture solve this problem by taking into account the spatial characteristics of their input data. The basic concepts of CNNs are *local receptive fields*, *shared weights* and *pooling*. In this context, the input layer is considered to be 2-dimensional structure typically a square.

Local receptive fields

Each neuron of the hidden layer is connected to a small region of the input layer, typically a small square, which is called the *local receptive field* of the hidden neuron. Each of these connections has a trainable weight, and the neuron has a bias. The region is slid across the input layer, passing values to the hidden layer.

Shared weights and biases

All of the hidden neurons share the same weights and biases; with the activation of the i, k th neuron begin:

$$\sigma \left(b + \sum_{l=0}^{r_x} \sum_{m=0}^{r_y} w_{l,m} a_{j+l, k+m} \right) \quad (2.17)$$

where σ represents the neural activation function, $a_{x,y}$ are the activations of the input at those coordinates and r_x, r_y are the dimensions of the receptive field.

Sharing weights and biases allows for features to be detected from input data, in conjunction with their spatial coordinates. For this reason, the mapping of the input layer to the hidden

layer is called a *feature map*, and the shared weights and biases are often called a *kernel* or *filter*. Typically, CNNs have several of these feature maps, one for each feature they are configured (or learn) to detect.

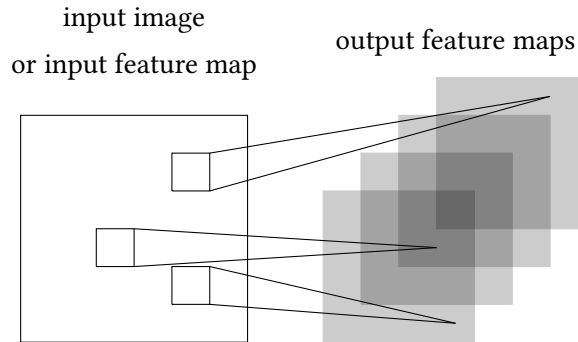


FIGURE 2.7: Illustration of a convolutional layer with multiple feature maps.

Pooling Layers Pooling layers are typically used immediately after the convolutional (hidden) layers. It reduces the dimensions of clusters of neurons into a single neuron in the next layer. The two most common pooling procedures are max-pooling which outputs the maximum activation from the pooling region and the average pooling which outputs the average value of the activations.

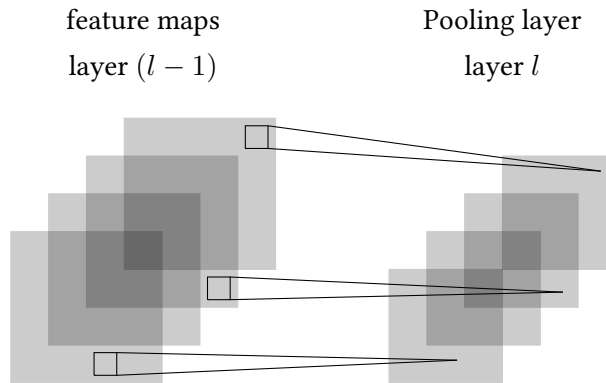


FIGURE 2.8: Illustration of the pooling layer following the feature maps

CNNs are important for understanding Graph neural networks as they introduce the concept of spatially aware data transformations, taking into account the topology of the input data. This, albeit with some changes, will be a critical concept when studying GNNs in the next chapter.

Chapter 3

Notation & Fundamentals

Chapter 4

]

Basic Principles and Implementation Framework for an [Problem to be Solved]

Chapter 5

]

Implementation and Core Components of [Platform Title]

Chapter 6

Experimentation & Validation

Chapter 7

Conclusions & Future Work

Bibliography

- [1] Kristy Carpenter, David Cohen, Juliet Jarrell, and Xudong Huang. Deep learning and virtual drug screening. *Future Medicinal Chemistry*, 10, 10 2018. doi: 10.4155/fmc-2018-0314.
- [2] U.S.R. Murty Adrian Bondy. *Graph theory*. Graduate texts in mathematics 244. Springer, 3rd corrected printing. edition, 2008. ISBN 1846289696; 9781846289699.
- [3] Bela Bollobas and Endre Szemerédi. Girth of sparse graphs. *Journal of Graph Theory*, 39: 194 – 200, 01 2002. doi: 10.1002/jgt.10023.
- [4] Ping Zhang Gary Chartrand. *A First Course in Graph Theory*. Dover Books on Mathematics. Dover Publications, 2012. ISBN 0486483681; 9780486483689.
- [5] A. Ashikhmin and A. Barg. *Algebraic Coding Theory and Information Theory: DIMACS Workshop, Algebraic Coding Theory and Information Theory, December 15-18, 2003, Rutgers University, Piscataway, New Jersey*. DIMACS series in discrete mathematics and theoretical computer science. American Mathematical Soc. ISBN 9780821871102. URL <https://books.google.gr/books?id=wp7XsCAm9EC>.
- [6] Mark Newman. *Networks*. Oxford University Press, 2 edition, 2018. ISBN 0198805098; 9780198805090.
- [7] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. 06 2016.
- [8] Jürgen Jost. *Dynamical Networks*, pages 35–62. Springer London, London, 2007. ISBN 978-1-84628-780-0. doi: 10.1007/978-1-84628-780-0_3. URL https://doi.org/10.1007/978-1-84628-780-0_3.
- [9] Warren S. McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The Bulletin of Mathematical Biophysics*, 5(4):115–133, 1943. doi: 10.1007/bf02478259.
- [10] Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65 6:386–408, 1958.

- [11] M. Minsky and S. Papert. *Perceptrons; an Introduction to Computational Geometry*. MIT Press, 1969. ISBN 9780262630221. URL <https://books.google.gr/books?id=Ow1OAQAAIAAJ>.
- [12] P.J. Werbos. *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*. Harvard University, 1975. URL <https://books.google.gr/books?id=z81XmgEACAAJ>.
- [13] J. Weng, N. Ahuja, and T.S. Huang. Cresceptron: a self-organizing neural network which grows adaptively. In *[Proceedings 1992] IJCNN International Joint Conference on Neural Networks*, volume 1, pages 576–581 vol.1, 1992. doi: 10.1109/IJCNN.1992.287150.
- [14] Kevin N. Gurney. *An introduction to neural networks*. 1997.
- [15] Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2014.09.003>. URL <https://www.sciencedirect.com/science/article/pii/S0893608014002135>.
- [16] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [17] A. Cauchy. Méthode générale pour la résolution des systèmes d'équations simultanées. *C. R. Acad. Sci*, 25:536–538, 1847.
- [18] Xiao Qi Yang Liqun Qi, Kok Lay Teo. *Optimization and control with applications*. Applied Optimization. Springer, 1 edition, 2005. ISBN 9780387242545; 0387242546.
- [19] Milton Abramowitz and Irene Stegun. *Abramowitz and Stegun*. Bracewell 200, 1972.
- [20] Michael A. Nielsen. *Neural Networks and Deep Learning*. Determination Press, 2015. URL <http://neuralnetworksanddeeplearning.com>.
- [21] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In Yves Lechevallier and Gilbert Saporta, editors, *Proceedings of COMPSTAT'2010*, pages 177–186, Heidelberg, 2010. Physica-Verlag HD. ISBN 978-3-7908-2604-3.
- [22] H. Robbins and D. Siegmund. A convergence theorem for non negative almost supermartingales and some applications**research supported by nih grant 5-r01-gm-16895-03 and onr grant n00014-67-a-0108-0018. In Jagdish S. Rustagi, editor, *Optimizing Methods in Statistics*, pages 233–257. Academic Press, 1971. ISBN 978-0-12-604550-5. doi: <https://doi.org/10.1016/B978-0-12-604550-5.50015-8>. URL <https://www.sciencedirect.com/science/article/pii/B9780126045505500158>.

- [23] Paul John Werbos. *The Roots of Backpropagation: From Ordered Derivatives to Neural Networks and Political Forecasting*. Wiley-Interscience, USA, 1994. ISBN 0471598976.
- [24] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323:533–536, 1986.
- [25] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.