

# An Overview on Deep Neural Networks

By Hamed Araab

Supervisor: Dr. Marzieh Zarinbal

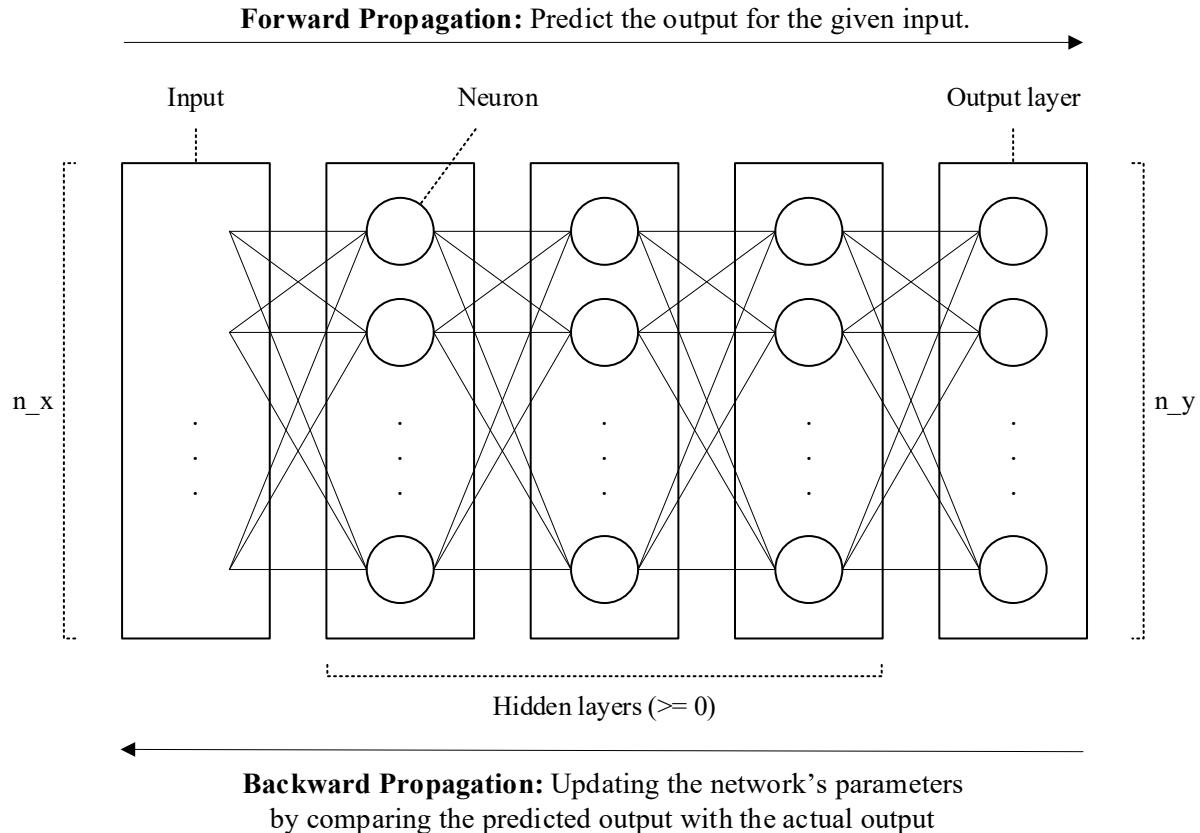
Winter 2024

## Table of Contents

1	Introduction.....	3
2	Common Tasks in Deep Learning.....	3
2.1	Classification.....	3
2.1.1	Binary.....	3
2.1.2	Multi-Class.....	4
2.1.3	Multi-Label .....	4
2.2	Regression.....	4
3	Comments on Some Notations.....	4
4	Components of a Simple DNN .....	4
4.1	Data .....	4
4.1.1	Sizes .....	4
4.1.2	Objects .....	5
4.2	Learnable Parameters.....	5
4.3	Hyperparameters .....	5
4.3.1	Examples.....	5
5	Forward Propagation.....	6
5.1	Mathematical Representation.....	6
5.1.1	Objects .....	6
5.1.2	Equations.....	7
6	Backward Propagation (Batch Gradient Descent) .....	7
6.1	Mathematical Representation.....	7
6.1.1	Objects .....	7
6.1.2	Equations.....	8
6.2	Activation Functions' Derivatives .....	8
6.3	Loss Functions' Derivatives.....	9
7	Different Types of Gradient Descent .....	9
7.1	Batch .....	9
7.2	Stochastic .....	9
7.3	Mini-Batch .....	9
8	Appendix.....	10
8.1	Some Mathematical Notations and Their NumPy Equivalents.....	10

# 1 Introduction

Deep Neural Networks (DNNs) are a type of artificial neural networks that stand out for having multiple layers. This multi-layer structure allows them to model complex relationships and create hierarchical representations in datasets. In each layer of a DNN, there are nodes or neurons connected to each other. These connections have weights, and there's also a bias term. The bias term makes the model more flexible, helping it capture subtle patterns in the data.



# 2 Common Tasks in Deep Learning

## 2.1 Classification

### 2.1.1 Binary

Binary classification involves the process of categorizing items into one of two distinct classes or categories. The primary objective is to assign each item to either Class 1 or Class 0. To illustrate, consider the task of classifying items as either a "cat" (Class 1) or "not cat" (Class 0):

→ [1]

→ [0]

→ [0]

### 2.1.2 Multi-Class

Multi-class classification involves the assignment of items into more than two classes or categories. The primary objective is to categorize items across multiple classes rather than just two. To illustrate, consider a scenario where we aim to classify items into distinct categories such as "cat," "dog," and "neither":

$$\text{cat} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{dog} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\text{neither} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\text{elephant} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

### 2.1.3 Multi-Label

Multi-label classification is a classification task where each input is capable of belonging to multiple classes simultaneously. In this scenario, the goal is to identify and assign multiple labels to each input. For instance, we may seek to find out whether each item contains a "cat," "dog," or "elephant":

$$\text{cat} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{elephant, cat} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\text{cat, dog} \rightarrow \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

$$\text{bird} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{horse} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

## 2.2 Regression

Regression is a type of machine learning task where the goal is to predict a continuous numerical value rather than a class label. It involves learning a mapping function from input features to a continuous output variable. For example, we can predict house prices based on features such as square footage, number of bedrooms, and location.

## 3 Comments on Some Notations

You can see the explanations for some notations that you are going to encounter later on:

- $\mathbb{Z}_{>0}$  denotes positive integers.
- $\mathbb{R}_{\geq 0}$  denotes non-negative real numbers.
- Superscript  $[l]$  denotes the  $l^{th}$  layer.
- $\odot$  denotes element-wise product of two matrices.
- $\oslash$  denotes element-wise division of two matrices.
- $:$  denotes tensor double contraction.

## 4 Components of a Simple DNN

### 4.1 Data

#### 4.1.1 Sizes

- $m \in \mathbb{Z}_{>0}$  The **number** of the items in the dataset

- $n_x \in \mathbb{Z}_{>0}$  The **size** of the **input** for each item
- $n_y \in \mathbb{Z}_{>0}$  The **size** of the **output** for each item

#### 4.1.2 Objects

- $X \in \mathbb{R}^{n_x \times m}$  The input **matrix**
- $Y \in \mathbb{R}^{n_y \times m}$  The actual output **matrix**

### 4.2 Learnable Parameters

- $W^{[l]} \in \mathbb{R}^{n_h^{[l]} \times n_h^{[l-1]}}$ : The weight **matrix** of the  $l^{th}$  layer
- $b^{[l]} \in \mathbb{R}^{n_h^{[l]} \times 1}$ : The bias **vector** of the  $l^{th}$  layer

### 4.3 Hyperparameters

- $L \in \mathbb{Z}_{>0}$  The **number** of the layers in the network (hidden + output)

- $n_h^{[l]} \in \mathbb{Z}_{>0}$  The **number** of the units (neurons) in the  $l^{th}$  layer  
It is also possible to denote  $n_x = n_h^{[0]}$  and  $n_y = n_h^{[L]}$ .

- $\sigma^{[l]}(Z^{[l]} \in \mathbb{R}^{n_h^{[l]} \times m}) \rightarrow \mathbb{R}^{n_h^{[l]} \times m}$  The activation **function** of the  $l^{th}$  layer
- $J(Y \in \mathbb{R}^{n_y \times m}, \hat{Y} \in \mathbb{R}^{n_y \times m}) \rightarrow \mathbb{R}_{\geq 0}$  The **loss (error) function** of the model for all training examples
- $\alpha \in (0,1]$  The **learning rate** of the model

#### 4.3.1 Examples

##### 4.3.1.1 Activation Functions

- Rectified Linear Unit (ReLU):  

$$\sigma^{[l]}(Z^{[l]})_{i,j} = \max\{0, Z_{i,j}^{[l]}\}$$
  
 Usually applied to **the output layer in regression**  
 One of the common functions in **the hidden layers of modern models**
- Sigmoid:  

$$\sigma^{[l]}(Z^{[l]})_{i,j} = \frac{1}{1+\exp(-Z_{i,j}^{[l]})}$$
  
 Usually applied to **the output layer in binary or multi-label classification**
- Softmax:  

$$\sigma^{[l]}(Z^{[l]})_{i,j} = \frac{\exp(Z_{i,j}^{[l]})}{\sum_{k=1}^{n_h^{[l]}} \exp(Z_{k,j}^{[l]})}$$

Usually applied to **the output layer** in **multi-class classification**

#### 4.3.1.2 Loss Functions

- Squared of Squared Errors (SSE):  $J(Y, \hat{Y}) = \frac{1}{2} \sum_{i=1}^{n_y} \sum_{j=1}^m (\hat{Y}_{i,j} - Y_{i,j})^2$   
Commonly used in **regression**  
The sum is divided by two for easier derivative calculation.
- Binary Cross Entropy (BCE):  $J(Y, \hat{Y}) = - \sum_{i=1}^{n_y} \sum_{j=1}^m (Y_{i,j} \log(\hat{Y}_{i,j}) + (1 - Y_{i,j}) \log(1 - \hat{Y}_{i,j}))$   
Commonly used in **binary or multi-label classification**
- Cross Entropy (CE):  $J(Y, \hat{Y}) = - \sum_{i=1}^{n_y} \sum_{j=1}^m (Y_{i,j} \log(\hat{Y}_{i,j}))$   
Commonly used in **multi-class classification**

#### 4.3.1.3 Summary

Task	Output Layer Activation Function	Loss Function
Regression	ReLU or Linear	SSE
Binary of Multi-Label Classification	Sigmoid	BCE
Multi-Class Classification	Softmax	CE

## 5 Forward Propagation

Forward propagation is a fundamental step in how deep neural networks work. It includes sending input data through the network's layers to produce an output. In each layer, the input is multiplied by weights, and the bias term is added. Activation functions bring in non-linearities to capture complex patterns effectively. The output from one layer becomes the input for the next layer, continuing through the network until the final layer makes the model's prediction.

### 5.1 Mathematical Representation

#### 5.1.1 Objects

- $Z^{[l]} \in \mathbb{R}^{n_h^{[l]} \times m}$   
The pre-activation **matrix** of the  $l^{th}$  layer  
Pre-activation values are also referred to as **logits**.
- $A^{[l]} \in \mathbb{R}^{n_h^{[l]} \times m}$   
The activation **matrix** of the  $l^{th}$  layer
- $\hat{Y} \in \mathbb{R}^{n_y \times m}$   
The predicted output **matrix**

### 5.1.2 Equations

- $A^{[0]} = X$
- $Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]}$   $\forall l = 1, 2, \dots, L$   
 $b$  is added to each column of  $W^{[l]}A^{[l-1]}$ . This is called broadcasting.
- $A^{[l]} = \sigma^{[l]}(Z^{[l]})$   $\forall l = 1, 2, \dots, L$
- $\hat{Y} = A^{[L]}$

## 6 Backward Propagation (Batch Gradient Descent)

Backward propagation, or backpropagation, is a crucial step in training deep neural networks. It helps the network learn from mistakes. In backpropagation, we calculate the error in the model's output by comparing what it predicted with the actual values we wanted. We then move backward through the network, one layer at a time. For each layer, we adjust the weights and biases in the opposite direction of the error gradient. This is done through an iterative process, guided by optimization algorithms like gradient descent. It helps the neural network improve its parameters, reducing prediction errors. Backward propagation is essential for the network to get better at making accurate predictions on new, unseen data. This makes deep learning models successful in various applications.

## 6.1 Mathematical Representation

### 6.1.1 Objects

- $C \in \mathbb{R}_{\geq 0}$  The **cost** of the predictions
- $\frac{\partial C}{\partial \hat{Y}} \in \mathbb{R}^{n_y \times m}$
- $\frac{\partial A^{[l]}}{\partial Z^{[l]}} \in \mathbb{R}^{n_h^{[l]} \times m \times n_h^{[l]} \times m}$
- $\frac{\partial Z^{[l]}}{\partial W^{[l]}} \in \mathbb{R}^{m \times n_h^{[l-1]}}$
- $\frac{\partial Z^{[l]}}{\partial b^{[l]}} \in \mathbb{R}^{m \times 1}$
- $\frac{\partial Z^{[l]}}{\partial A^{[l-1]}} \in \mathbb{R}^{n_h^{[l-1]} \times n_h^{[l]}}$
- $\frac{\partial C}{\partial A^{[l-1]}} \in \mathbb{R}^{n_h^{[l-1]} \times m}$
- $\frac{\partial C}{\partial W^{[l]}} \in \mathbb{R}^{n_h^{[l]} \times n_h^{[l-1]}}$
- $\frac{\partial C}{\partial b^{[l]}} \in \mathbb{R}^{n_h^{[l]} \times 1}$

### 6.1.2 Equations

#### 6.1.2.1 Calculating the Cost of the Predictions

- $\mathcal{C} = \frac{1}{m} J(Y, \hat{Y})$

#### 6.1.2.2 Calculating the Partial Derivatives

- $\frac{\partial \mathcal{C}}{\partial \hat{Y}} = \frac{1}{m} \frac{\partial J}{\partial \hat{Y}}$
- $\frac{\partial \mathcal{C}}{\partial A^{[L]}} = \frac{\partial \mathcal{C}}{\partial \hat{Y}}$
- $\frac{\partial A^{[l]}}{\partial Z^{[l]}} = \sigma'^{[l]}(Z^{[l]}) \quad \forall l = L, L-1, \dots, 1$
- $\frac{\partial Z^{[l]}}{\partial W^{[l]}} = A^{[l-1], T} \quad \forall l = L, L-1, \dots, 1$
- $\frac{\partial Z^{[l]}}{\partial b^{[l]}} = 1_{m \times 1} \quad \forall l = L, L-1, \dots, 1$
- $\frac{\partial Z^{[l]}}{\partial A^{[l-1]}} = W^{[l], T} \quad \forall l = L, L-1, \dots, 1$

#### 6.1.2.3 Applying the Chain Rule

- $\frac{\partial \mathcal{C}}{\partial W^{[l]}} = \left( \frac{\partial \mathcal{C}}{\partial A^{[l]}} : \frac{\partial A^{[l]}}{\partial Z^{[l]}} \right) \frac{\partial Z^{[l]}}{\partial W^{[l]}} = \left( \frac{\partial \mathcal{C}}{\partial A^{[l]}} : \sigma'^{[l]}(Z^{[l]}) \right) A^{[l-1], T} \quad \forall l = L, L-1, \dots, 1$
- $\frac{\partial \mathcal{C}}{\partial b^{[l]}} = \left( \frac{\partial \mathcal{C}}{\partial A^{[l]}} : \frac{\partial A^{[l]}}{\partial Z^{[l]}} \right) \frac{\partial Z^{[l]}}{\partial b^{[l]}} = \left( \frac{\partial \mathcal{C}}{\partial A^{[l]}} : \sigma'^{[l]}(Z^{[l]}) \right) 1_{m \times 1} \quad \forall l = L, L-1, \dots, 1$
- $\frac{\partial \mathcal{C}}{\partial A^{[l-1]}} = \frac{\partial Z^{[l]}}{\partial A^{[l-1]}} \left( \frac{\partial \mathcal{C}}{\partial A^{[l]}} : \frac{\partial A^{[l]}}{\partial Z^{[l]}} \right) = W^{[l], T} \left( \frac{\partial \mathcal{C}}{\partial A^{[l]}} : \sigma'^{[l]}(Z^{[l]}) \right) \quad \forall l = L, L-1, \dots, 1$

#### 6.1.2.4 Updating the Learnable Parameters

- $\Delta W^{[l]} = -\alpha \frac{\partial \mathcal{C}}{\partial W^{[l]}} \quad \forall l = L, L-1, \dots, 1$
- $\Delta b^{[l]} = -\alpha \frac{\partial \mathcal{C}}{\partial b^{[l]}} \quad \forall l = L, L-1, \dots, 1$

## 6.2 Activation Functions' Derivatives

$$\sigma'^{[l]}(Z^{[l]})_{i,j,u,v} = \frac{\partial}{\partial Z^{[l]}_{u,v}} \sigma^{[l]}(Z^{[l]})_{i,j}$$

- ReLU:  $\sigma'^{[l]}(Z^{[l]})_{i,j,u,v} = \begin{cases} 1 & i = u \text{ and } j = v \text{ and } Z^{[l]}_{u,v} > 0 \\ 0 & \text{otherwise} \end{cases}$

- Sigmoid:

$$\sigma'^{[l]}(Z^{[l]})_{i,j,u,v} = \begin{cases} A_{i,j}^{[l]}(1 - A_{i,j}^{[l]}) & i = u \text{ and } j = v \\ 0 & \text{otherwise} \end{cases}$$

- Softmax:

$$\sigma'^{[l]}(Z^{[l]})_{i,j,u,v} = \begin{cases} A_{i,j}^{[l]}(1 - A_{i,j}^{[l]}) & i = u \text{ and } j = v \\ -A_{i,j}^{[l]}A_{u,j}^{[l]} & i \neq u \text{ and } j = v \\ 0 & j \neq v \end{cases}$$

## 6.3 Loss Functions' Derivatives

- SSE:

$$\frac{\partial \mathcal{L}}{\partial \hat{Y}_{i,j}} = \hat{Y}_{i,j} - Y_{i,j} \Rightarrow \frac{\partial \mathcal{L}}{\partial \hat{Y}} = \hat{Y} - Y$$

- BCE:

$$\frac{\partial \mathcal{L}}{\partial \hat{Y}_{i,j}} = \frac{\hat{Y}_{i,j} - Y_{i,j}}{\hat{Y}_{i,j}(1 - \hat{Y}_{i,j})} \Rightarrow \frac{\partial \mathcal{L}}{\partial \hat{Y}} = (\hat{Y} - Y) \oslash (\hat{Y} \odot (1 - \hat{Y}))$$

- CE:

$$\frac{\partial \mathcal{L}}{\partial \hat{Y}_{i,j}} = -\frac{Y_{i,j}}{\hat{Y}_{i,j}} \Rightarrow \frac{\partial \mathcal{L}}{\partial \hat{Y}} = -Y \oslash \hat{Y}$$

## 7 Different Types of Gradient Descent

### 7.1 Batch

Batch gradient descent involves computing the gradient of the entire dataset's cost function with respect to the model parameters in each iteration. This method provides a stable and accurate estimate of the gradient but can be computationally expensive, especially for large datasets.

### 7.2 Stochastic

Stochastic gradient descent takes a different approach by updating the model parameters after evaluating the gradient for each individual data point in the training set. While this can be computationally more efficient, it introduces higher variability in the updates and may result in noisy convergence. Despite the increased variance, stochastic gradient descent is particularly useful for large datasets and online learning scenarios.

### 7.3 Mini-Batch

Mini-batch gradient descent strikes a balance between the previous two approaches. It involves dividing the dataset into smaller batches and computing the gradient based on each mini-batch. This method combines the computational efficiency of stochastic gradient descent with the stability of batch gradient descent. The mini-batch size is a hyperparameter that can be tuned to find the right balance between efficiency and accuracy, making mini-batch gradient descent a popular choice in many machine learning applications.

## 8 Appendix

### 8.1 Some Mathematical Notations and Their NumPy Equivalents

- $AB$   $\equiv$   $A @ B$
- $A:B$   $\equiv$   $\text{np.tensordot}(A, B)$
- $A \odot B$   $\equiv$   $A * B$
- $A \oslash B$   $\equiv$   $A / B$