

DIVE INTO DEEP LEARNING: A COMPREHENSIVE INTRODUCTION

FROM AI FUNDAMENTALS TO CUTTING-EDGE DEEP LEARNING TECHNIQUES

**Alexandre Vérine,
Research Fellow, ENS-PSL
Université PSL**

Executive Master IASD
Université Paris-Dauphine, PSL

January 27, 2025

AI 101: FROM FUNDAMENTALS TO DEEP LEARNING

1	Introduction to Artificial Intelligence	5
1.1	Deep Learning in the AI family	5
1.2	Representation Learning	10
2	Neural Networks Fundamentals	15
2.1	Neurons	16
2.2	Layers	18
2.3	Activation Functions	20
3	The Multi-layer Perceptron (MLP)	29
3.1	The first Deep Learning Model	30
3.2	Stochastic Gradient Descent	31
3.3	Back-propagation	34
3.4	Example : Image classification of handwritten digits from A to Z	56

DEEP LEARNING IN ACTION: FROM NEURAL NETWORKS TO TRANSFORMER MODELS

1 Convolutional Neural Networks	65
1.1 The Two dimensional Convolution	66
1.2 CNN : Convolutional in a network Networks	74
1.3 CNN in practice: CIFAR 10	81
2 Recurrent Neural Networks	106
2.1 Recurrent Block	107
2.2 LSTM and GRU	109
3 Transformer and Attention Mechanism	119
3.1 Self-Attention Mechanism	120
3.2 Transformers Model	124

TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

1 Techniques to Improve Deep Learning Training	126
1.1 Data Augmentation	127
1.2 Learning Rate Scheduling	128
1.3 Early Stopping	129
1.4 Gradient Clipping	130
1.5 Weight Initialization	131
1.6 Regularization	133
1.7 GPU Acceleration	134

Part I

AI 101: FROM FUNDAMENTALS TO DEEP LEARNING

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

DEEP LEARNING IN THE AI FAMILY

In general, among all the class of AI algorithms, we make the difference between 3 sub-categories :

- ▶ **Artificial Intelligence** : human designed program and...
- ▶ **Machine Learning** : human designed features with learned mapping such as Support Vector Machine, Kernels methods, Logistic Regression and ...
- ▶ **Deep Learning**: Learned features with learned mapping such as Multilayer Perceptron, Convolutional Networks, ...

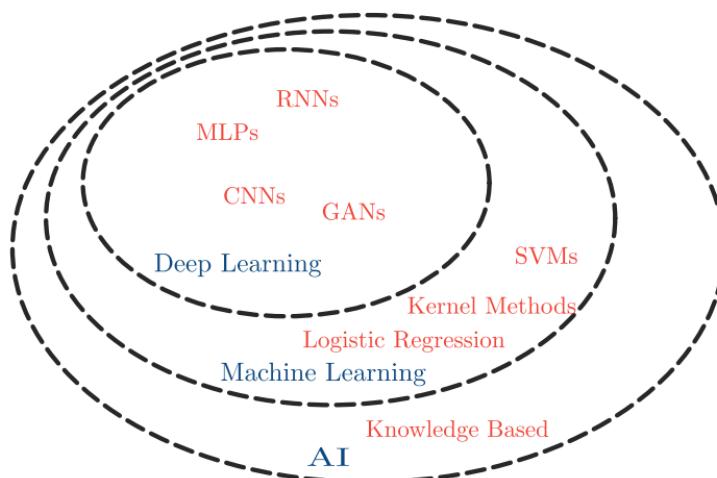


Figure. Subsets of Artificial Intelligence

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

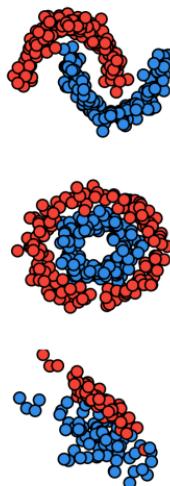
DEEP LEARNING IN THE AI FAMILY

In the field of Artificial Intelligence, the fundamental objective is to find a function f that can perform a desired task. This function can either be set by a human or can be learned through training.

For example, in the context of a binary classification task, the goal is to determine $f(x)$ such that $f(x) = 0$ when the label of x is 0 and $f(x) = 1$ when its label is 1. The choice of AI model impacts the expressivity of the function f .

For example, a logistic regression model uses a linear function to make decisions, where $f(x) = \text{sgn}(Ax + b)$. The expressivity of the model can be increased by using more complex functions, such as polynomials or radial basis functions.

Input data



INTRODUCTION TO ARTIFICIAL INTELLIGENCE

CLASSIFICATION TASK

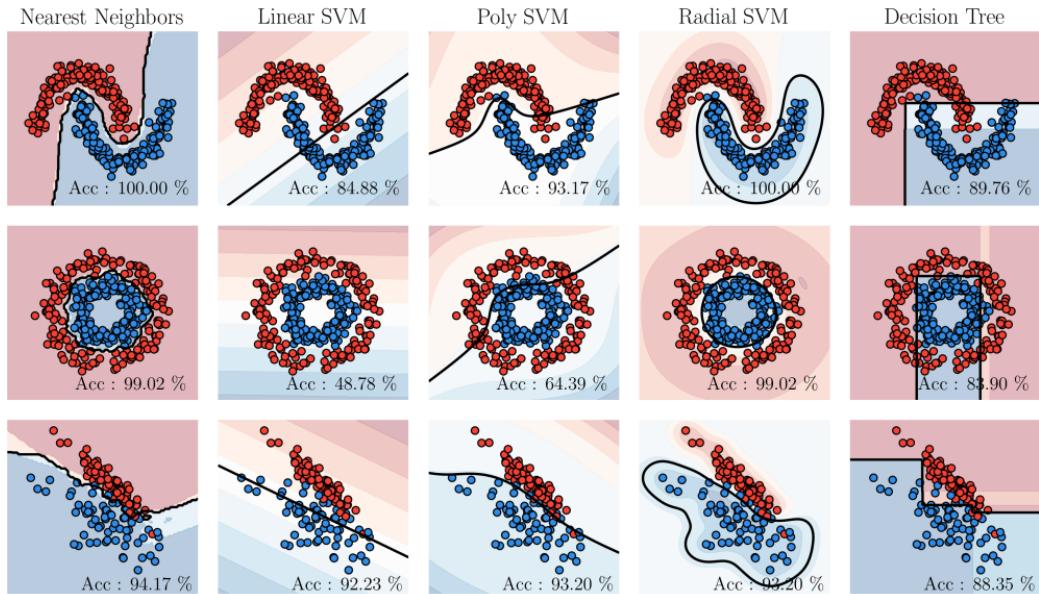


Figure. 2D classification for different AI models.

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

THE UNIVERSAL APPROXIMATION THEOREM

The Universal Approximation Theorem is a fundamental result in the field of artificial neural networks. It states that a deep learning model can approximate any function.

Theorem 1 (Universal Approximation Theorem)

Let $\mathcal{X} \subset \mathbb{R}^d$ be compact, $\mathcal{Y} \subset \mathbb{R}^m$, $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a continuous function and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous real function.

Then σ is not polynomial if and only if for every $\epsilon > 0$, there exist $k \in \mathbb{N}$, $A \in \mathbb{R}^{k \times d}$, $b \in \mathbb{R}^k$ and $C \in \mathbb{R}^{m \times k}$ such that

$$\sup_{x \in \mathcal{X}} \|f(x) - g(x)\| \leq \epsilon$$

where $g(x) = C \times \sigma(Ax + b)$.

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

CLASSIFICATION TASK

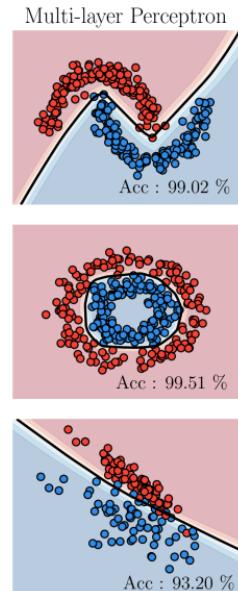


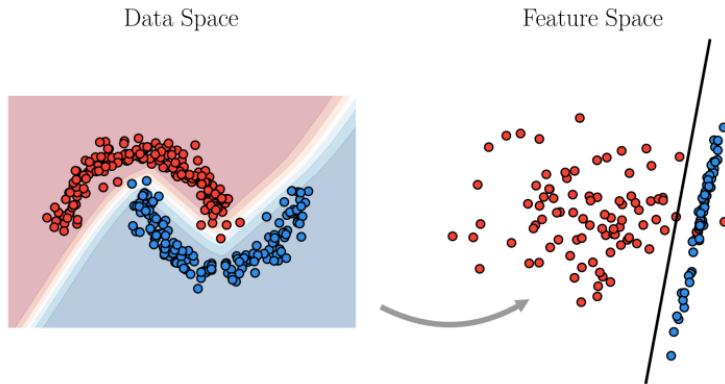
Figure. 2D classification for small Neural Network.

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

REPRESENTATION LEARNING

How does deep learning work in practice ?

Deep learning is a subset of representation learning that uses deep neural networks to learn meaningful representations of data. In deep learning, representations are learned through a hierarchy of nonlinear transformations, where each layer of the network builds upon the previous one to extract increasingly abstract and higher-level features from the input data.



INTRODUCTION TO ARTIFICIAL INTELLIGENCE

EXAMPLE OF REPRESENTATION LEARNING

Consider the task of recognizing objects in images. A traditional approach would be to hand-engineer features such as edge detectors and color histograms that can be fed into a classifier.

However, with deep learning representation learning, the model learns to automatically discover these features from the data. The network might start by learning simple features such as edges and color blobs in the first layer, then build upon these to learn more complex features such as parts of objects in subsequent layers, until finally, the final layer outputs a probability distribution over classes of objects.

In this way, deep learning of representation enables the model to automatically learn a rich and meaningful representation of the data, without the need for manual feature engineering.

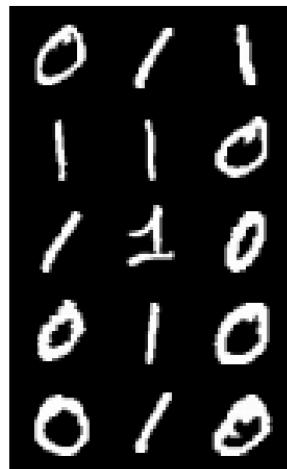


Figure. MNIST

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

EXAMPLE OF REPRESENTATION LEARNING

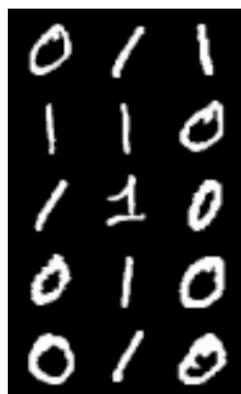


Figure. MNIST : Layer 0

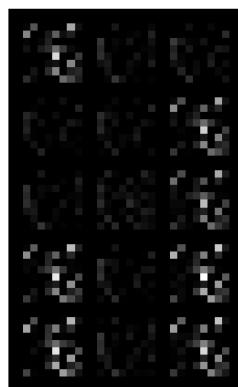


Figure. MNIST : Layer 1

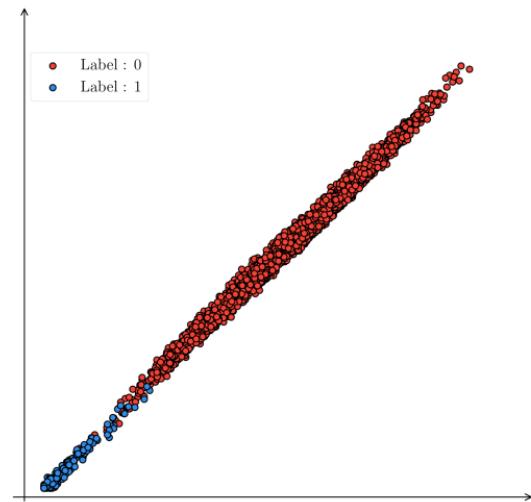


Figure. MNIST : Layer 2

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

EXAMPLE OF REPRESENTATION LEARNING



Figure. MNIST : Layer 0

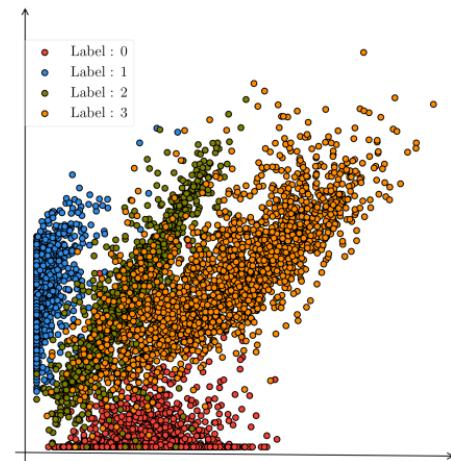


Figure. MNIST : Layer 2

INTRODUCTION TO ARTIFICIAL INTELLIGENCE

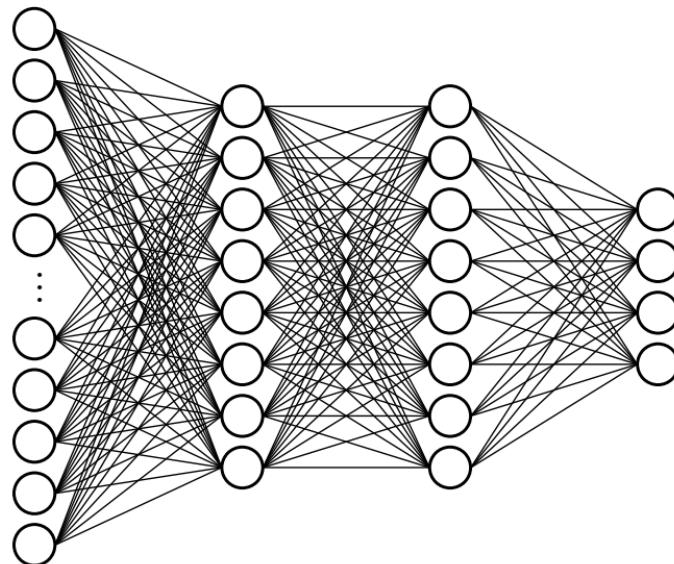
DEEP LEARNING AND NEURAL NETWORKS

Ok, Deep Learning is a model that learns a good representation of the feature. But how?

- ▶ How does it work ?
- ▶ How can we build a model ?
- ▶ How does it learn ?

NEURAL NETWORKS FUNDAMENTALS

Typically, a neural network is defined as a computational model composed of interconnected nodes, organised into layers, that perform transformations on input data.

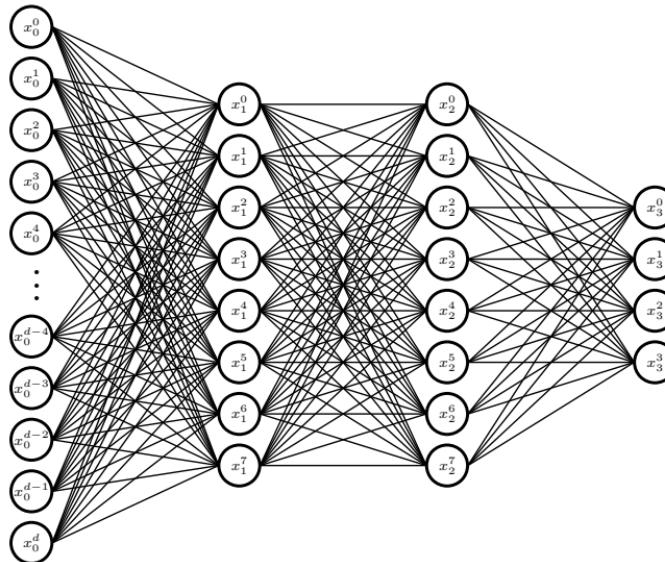


Let's see what the interconnected nodes, the layers and the transformations are.

NEURAL NETWORKS FUNDAMENTALS

NEURONS

If we consider that the Neural Network is a function $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$:



A **Neuron** is a processing unit that receives input, performs a computation, and produces an output. Here, the inputs are x_{i-1} and the output is x_i^k .

NEURAL NETWORKS FUNDAMENTALS

NEURONS

For example, with an image dataset, the image can be flattened:

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.99	0.91	0.02	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.99	0.45	0.18	0.66	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.99	0.07	0.00	0.00	0.99	0.00	0.00	0.00	0.00
0.00	0.00	0.30	0.44	0.00	0.00	0.00	0.99	0.00	0.00	0.00	0.00
0.00	0.00	0.33	0.00	0.00	0.00	0.99	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.33	0.99	0.99	0.77	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

$$\in [0, 1]^{d/2 \times d/2}$$

$$x_0 = [0.00, 0.00, \dots, 0.00, 0.99, 0.07 \dots, 0.00, 0.00] \in [0, 1]^d$$

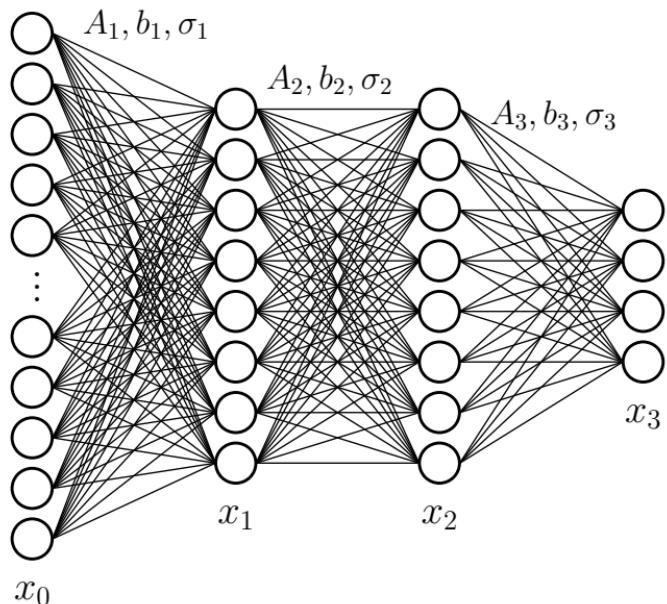
NEURAL NETWORKS FUNDAMENTALS

LAYERS

A layer i is defined by a matrix $A_i \in \mathbb{R}^{k_{i-1} \times k_i}$, a vector $b_i \in \mathbb{R}^{k_i}$ and a nonlinear function $\sigma_i : \mathbb{R} \mapsto \mathbb{R}$. The transformation made by a layer is:

$$x_i = \sigma_i(A_i x_{i-1} + b_i).$$

The non-linear function σ_i the **activation function**.



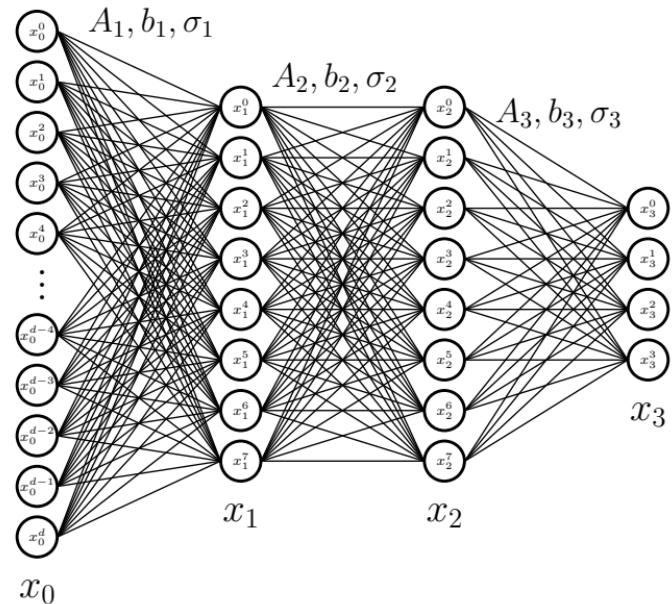
NEURAL NETWORKS FUNDAMENTALS

LAYERS

A layer i is defined as a matrix $A_i \in \mathbb{R}^{k_{i-1} \times k_i}$, a vector $b_i \in \mathbb{R}^{k_i}$ and a nonlinear function $\sigma_i : \mathbb{R} \mapsto \mathbb{R}$. The transformation made by a layer is:

$$x_i^k = \sigma_i \left(\sum_{l=1}^{k_i} [A_i]_{l,k} x_{i-1} + [b_i]_k \right).$$

The non-linear function σ_i the activation function.



NEURAL NETWORKS FUNDAMENTALS

ACTIVATION FUNCTIONS

The activation functions play a crucial role in the implementation of deep neural networks, as they allow them to approximate any continuous function, as stated by the Universal Approximation Theorem. We can list some activation function that are commonly used :

- ▶ Linear
- ▶ Sigmoid
- ▶ Hyperbolic Tangent
- ▶ Rectified Linear Unit (ReLU)
- ▶ Leaky Rectified Linear Unit (Leaky ReLU)
- ▶ Exponential Linear Unit (ELU)
- ▶ Sigmoid-Weighted Linear Unit (Swish)
- ▶ Softmax

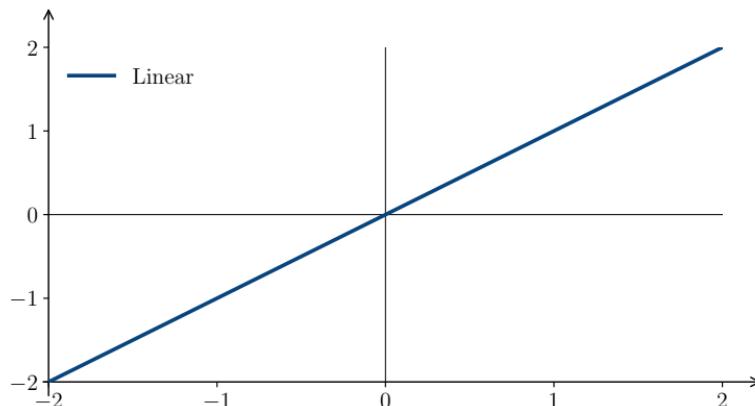
NEURAL NETWORKS FUNDAMENTALS

LINEAR

- ▶ Linear activation Function:

$$\sigma(x) = x$$

- ▶ Final activation
- ▶ Use case : Regression



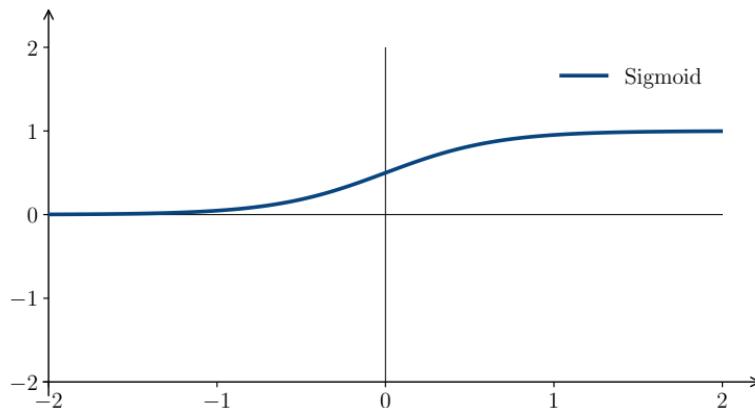
NEURAL NETWORKS FUNDAMENTALS

SIGMOID

- ▶ Sigmoid Function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- ▶ Final activation
- ▶ Use case : Classification



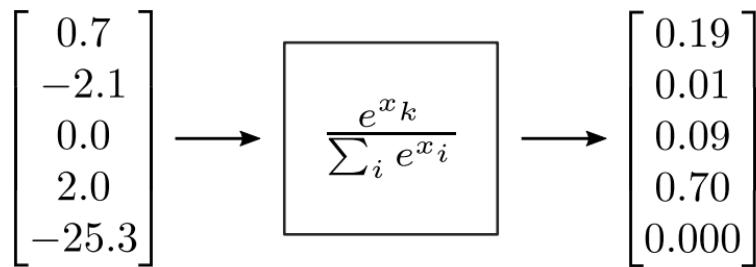
NEURAL NETWORKS FUNDAMENTALS

SOFTMAX

- ▶ Softmax Function:

$$\sigma(x_k) = \frac{e^{x_k}}{\sum_{i=1}^k e^{x_i}}$$

- ▶ Final activation
- ▶ Use case : Multi-class Classification



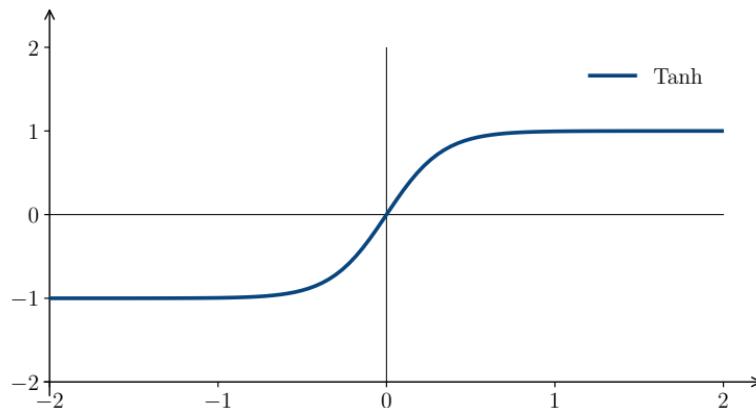
NEURAL NETWORKS FUNDAMENTALS

HYPERBOLIC TANGENT

- ▶ Hyperbolic Tangent

$$\sigma(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- ▶ Final activation
- ▶ Use case : Generative task



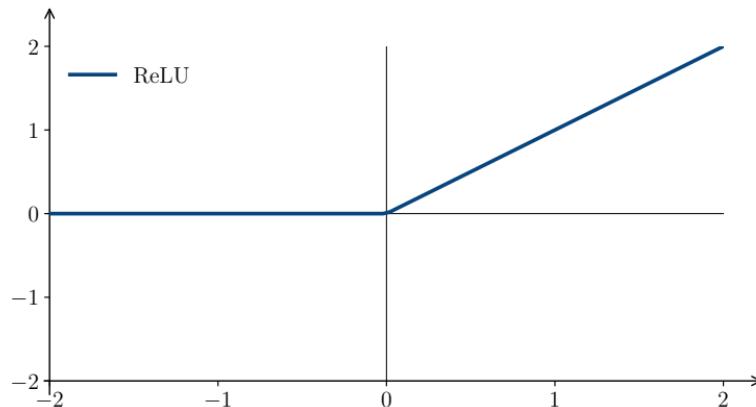
NEURAL NETWORKS FUNDAMENTALS

ReLU

- ▶ Rectified Linear Unit (ReLU):

$$\sigma(x) = \max\{0, x\}$$

- ▶ Intermediate activation



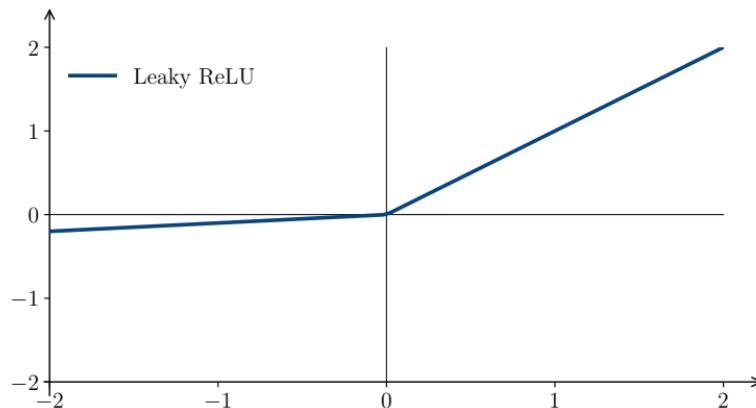
NEURAL NETWORKS FUNDAMENTALS

LEAKY RELU

- ▶ Leaky Rectified Linear Unit (Leaky ReLU):

$$\sigma(x) = \max\{\alpha x, x\}$$

- ▶ Intermediate activation



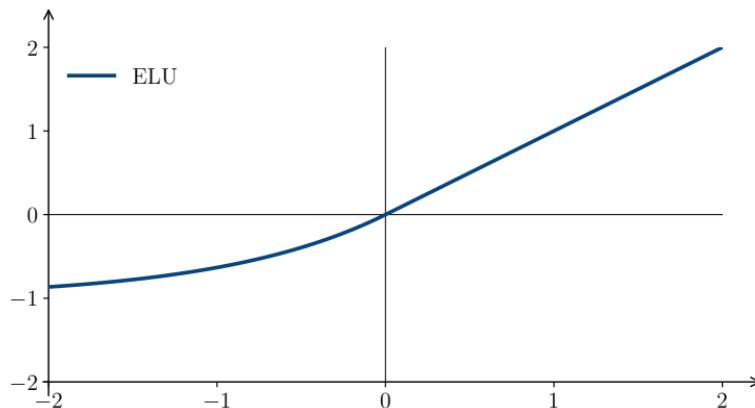
NEURAL NETWORKS FUNDAMENTALS

ELU

- ▶ Exponential Linear Unit (ELU):

$$\sigma(x) = \begin{cases} \alpha(e^x - 1) & \text{if } x < 0, \\ x & \text{if } x \geq 0. \end{cases}$$

- ▶ Intermediate activation



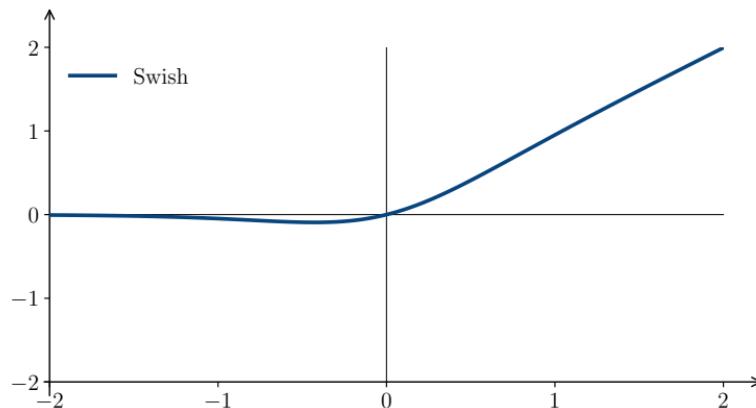
NEURAL NETWORKS FUNDAMENTALS

SWISH

- ▶ Sigmoid-Weighted Linear Unit (Swish):

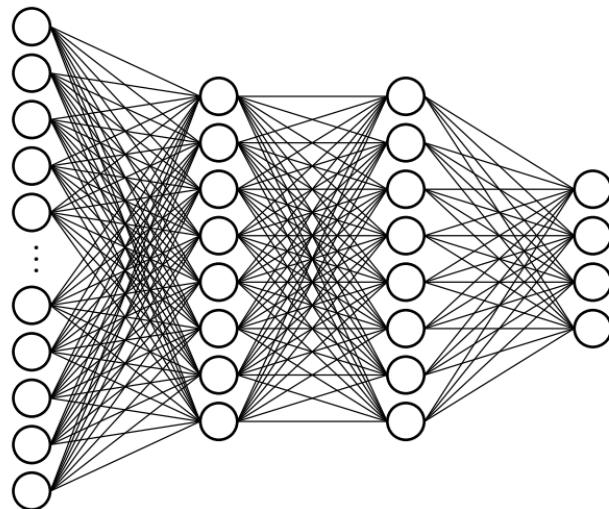
$$\sigma(x) = \frac{x}{1 + e^{-x}}$$

- ▶ Intermediate activation



THE MULTI-LAYER PERCEPTRON (MLP)

Having discussed the structure of a neural network, we will proceed to examine the process of training a model for a specific task. As an illustration, we will consider the example of a Multilayer Perceptron. The two intermediate activation functions are ReLUs and the final activation is a softmax to perform multi-class classification on MNIST. We will consider only 4 classes.

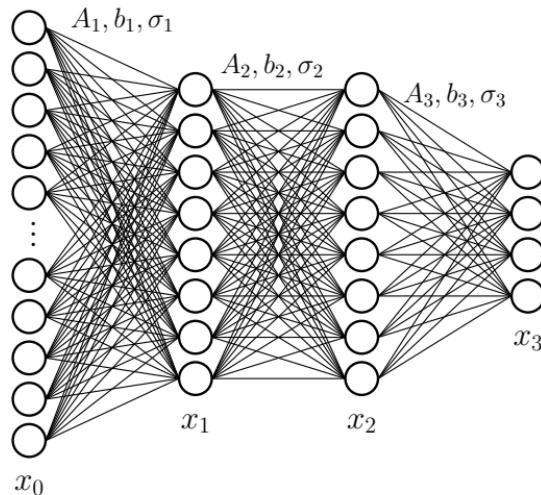


THE MULTI-LAYER PERCEPTRON (MLP)

THE FIRST DEEP LEARNING MODEL

To introduce the training process, we will consider a 3 layers MLP trained to minimise a loss \mathcal{L} over a given dataset \mathcal{D} . The model f_θ is parameterised by a vector $\theta = \{A_1, A_2, A_3, b_1, b_2, b_3\}$:

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta, \mathcal{D})$$



THE MULTI-LAYER PERCEPTRON (MLP)

STOCHASTIC GRADIENT DESCENT

Stochastic gradient descent (SGD) is widely used in deep learning instead of traditional gradient descent due to its efficiency and faster convergence rate. SGD updates the model parameters after computing the gradient of the loss function with respect to each parameter using only a single randomly selected sample. This leads to a faster convergence rate and improved optimization compared to traditional gradient descent, which uses the entire training dataset to compute the gradient at each iteration.

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta, \mathcal{D}) = \arg \min_{\theta} \mathbb{E}_{x \sim \mathcal{D}} [l(x, f_{\theta}(x))]$$

THE MULTI-LAYER PERCEPTRON (MLP)

STOCHASTIC GRADIENT DESCENT

Theoretically the algorithm is the following:

Require: Given a loss function l , a dataset $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ and a learning rate λ

- 1: Initialize parameters θ
- 2: **while** θ has not converged **do**
- 3: **for** $i = 1$ to N **do**
- 4: Randomly select x_i from the dataset
- 5: Compute gradient of the loss with respect to θ : $\nabla_{\theta} l(x_i, f_{\theta}(x_i))$
- 6: Update parameters $\theta = \theta - \lambda \nabla_{\theta} l(x_i, f_{\theta}(x_i))$
- 7: **end for**
- 8: **end while**
- 9: **return** θ

THE MULTI-LAYER PERCEPTRON (MLP)

SGD IN MINI-BATCH

In practice the algorithm is modified to use mini-batches of data instead of single samples. This is done to improve the stability of the optimization process and reduce the variance of the gradient estimates. The algorithm is as follows:

Require: Given a loss function l , a dataset $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$, a learning rate λ and a batch size b

- 1: Initialize parameters θ
- 2: Initialize the number of batches $B = \lfloor \frac{N}{b} \rfloor$
- 3: **while** θ has not converged **do**
- 4: **for** $i = 1$ to B **do**
- 5: Randomly select a mini-batch of b samples from the dataset
- 6: Compute gradient of the loss with respect to θ : $\frac{1}{B} \sum_{i=1}^B \nabla_\theta l(x_i, f_\theta(x_i))$
- 7: Update parameters $\theta = \theta - \lambda \frac{1}{B} \sum_{i=1}^B \nabla_\theta l(x_i, f_\theta(x_i))$
- 8: **end for**
- 9: **end while**
- 10: **return** θ

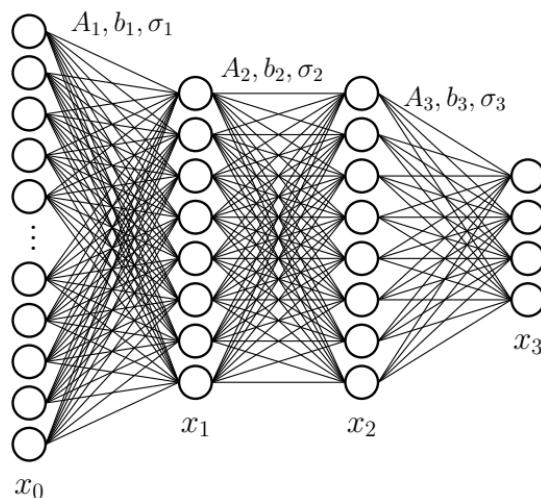
THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

At every step t of the gradient descent, setting a learning rate λ , the parameter θ is updated as:

$$\theta_{t+1} = \theta_t - \lambda \nabla_\theta l(f(x_i), y_i)$$

But $\theta = \{A_1, A_2, A_3, b_1, b_2, b_3\}$ and the gradient is computed with respect to each parameter.



THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

First we will consider a single data point x , the loss will depend on the output only: $l(f(x))$.

f is a layered composed function. Let us focus on the last layer:

$$f(x) = x_3 = \sigma_3(A_3x_2 + b_3)$$

Therefore:

$$l(f(x)) = l(\sigma_3(A_3x_2 + b_3))$$

To minimise the loss, we have to act on A_3 , b_3 and x_2 .

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Let us look at the gradients with respect to A_3 :

$$\begin{aligned}\frac{\partial l}{\partial A_3} &= \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial A_3} = l'(x_3) \frac{\partial \sigma_3(A_3 x_2 + b_3)}{\partial A_3} = l'(x_3) \sigma'_3(A_3 x_2 + b_3) \frac{\partial [A_3 x_2 + b_3]}{\partial A_3} \\ &= \underbrace{l'(x_3)}_{\in \mathbb{R}} \underbrace{\sigma'_3(A_3 x_2 + b_3)}_{\in \mathbb{R}^{k_i \times 1}} \underbrace{x_2^T}_{\in \mathbb{R}^{1 \times k_{i-1}}}\end{aligned}$$

and therefore:

$$A_3 \leftarrow A_3 - \lambda l'(x_3) \sigma'_3(A_3 x_2 + b_3) x_2^T.$$

We need to keep in memory the latent values of x , i.e. x_2 .

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Let us look at the gradients with respect to A_2 :

$$\begin{aligned}\frac{\partial l}{\partial A_2} &= \frac{\partial l}{\partial x_2} \frac{\partial x_2}{\partial A_2} \\ &= \frac{\partial l}{\partial x_2} \frac{\partial \sigma_2(A_2x_1 + b_2)}{\partial A_2} \\ &= \frac{\partial l}{\partial x_2} \sigma'_2(A_2x_1 + b_2) \frac{\partial [A_2x_1 + b_2]}{\partial A_2} \\ &= \frac{\partial l}{\partial x_2} \sigma'_2(A_2x_1 + b_2) x_1^T\end{aligned}$$

which depends on $\frac{\partial l}{\partial x_2}$, we need to compute it.

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

We have to compute the gradient with respect to x_2 :

$$\begin{aligned}\frac{\partial l}{\partial x_2} &= \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial x_2} = l'(x_3) \frac{\partial \sigma_3(A_3x_2 + b_3)}{\partial x_2} = l'(x_3) \frac{\partial [A_3x_2 + b_3]}{\partial x_2} \sigma'_3(A_3x_2 + b_3) \\ &= l'(x_3) A_3^T \sigma'_3(A_3x_2 + b_3)\end{aligned}$$

Therefore:

$$A_2 \leftarrow A_2 - \lambda \left[l'(x_3) A_3^T \sigma'_3(A_3x_2 + b_3) \times \sigma'_2(A_2x_1 + b_2) x_1^T \right]$$

The update of A_2 depends on $l'(x_3)$,

BACK-PROPAGATION

We have to compute the gradient with respect to A_1 :

$$\begin{aligned}\frac{\partial l}{\partial A_1} &= \frac{\partial l}{\partial x_1} \frac{\partial x_1}{\partial A_1} \\ &= \frac{\partial l}{\partial x_1} \frac{\partial \sigma_1(A_1 x_0 + b_1)}{\partial A_1} \\ &= \frac{\partial l}{\partial x_1} \sigma'_1(A_1 x_0 + b_1) x_0^T,\end{aligned}$$

which depends on $\frac{\partial l}{\partial x_1}$, we need to compute it.

BACK-PROPAGATION

Let us compute the gradient with respect to x_1 :

$$\begin{aligned}\frac{\partial l}{\partial x_1} &= \frac{\partial l}{\partial x_2} \frac{\partial x_2}{\partial x_1} = \frac{\partial l}{\partial x_2} \frac{\partial \sigma_2 (A_2 x_1 + b_2)}{\partial x_1} = \frac{\partial l}{\partial x_2} \frac{\partial [A_2 x_1 + b_2]}{\partial x_1} \sigma'_2 (A_2 x_1 + b_2) \\ &= \frac{\partial l}{\partial x_2} A_2^T \sigma'_2 (A_2 x_1 + b_2)\end{aligned}$$

Therefore:

$$A_1 \leftarrow A_1 - \lambda \left[l'(x_3) A_3^T \sigma'_3 (A_3 x_2 + b_3) A_2^T \sigma'_2 (A_2 x_1 + b_2) \times \sigma'_1 (A_1 x_0 + b_1) x_0^T \right]$$

BACK-PROPAGATION

In other words, the update on the weights is:

$$A_3 \leftarrow A_3 - \lambda l'(x_3) \sigma'_3 (A_3 x_2 + b_3) x_2^T$$

$$A_2 \leftarrow A_2 - \lambda \left[l'(x_3) A_3^T \sigma'_3 (A_3 x_2 + b_3) \times \sigma'_2 (A_2 x_1 + b_2) x_1^T \right]$$

$$A_1 \leftarrow A_1 - \lambda \left[l'(x_3) A_3^T \sigma'_3 (A_3 x_2 + b_3) A_2^T \sigma'_2 (A_2 x_1 + b_2) \times \sigma'_1 (A_1 x_0 + b_1) x_0^T \right]$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

If we look at the update of the different biases, we can easily compute the different gradient and see the updates. First, let us compute the gradient with respect to b_3 :

$$\begin{aligned}\frac{\partial l}{\partial b_3} &= \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial b_3} \\ &= l'(x_3) \frac{\partial \sigma_3(A_3x_2 + b_3)}{\partial b_3} \\ &= l'(x_3) \sigma'_3(A_3x_2 + b_3) \frac{\partial [A_3x_2 + b_3]}{\partial b_3} \\ &= \underbrace{l'(x_3)}_{\in \mathbb{R}} \underbrace{\sigma'_3(A_3x_2 + b_3)}_{\in \mathbb{R}_i^{k \times 1}}\end{aligned}$$

And thus :

$$b_3 \leftarrow b_3 - \lambda l'(x_3) \sigma'(A_3x_2 + b_3)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Let's move on the second layer:

$$\begin{aligned}\frac{\partial l}{\partial b_2} &= \frac{\partial l}{\partial x_2} \frac{\partial x_2}{\partial b_2} \\ &= \frac{\partial l}{\partial x_2} \frac{\partial \sigma_2 (A_2 x_1 + b_2)}{\partial b_2} \\ &= \frac{\partial l}{\partial x_2} \sigma'_2 (A_2 x_1 + b_2)\end{aligned}$$

And thus :

$$b_2 \leftarrow b_2 - \lambda \frac{\partial l}{\partial x_2} \sigma' (A_2 x_1 + b_2)$$

We need to back-propagate the term $\frac{\partial l}{\partial x_2}$ computed for the first layer.

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

For the first layer:

$$\begin{aligned}\frac{\partial l}{\partial b_1} &= \frac{\partial l}{\partial x_1} \frac{\partial x_1}{\partial b_1} \\ &= \frac{\partial l}{\partial x_1} \frac{\partial \sigma_1(A_1x_0 + b_1)}{\partial b_1} \\ &= \frac{\partial l}{\partial x_1} \sigma'_1(A_1x_0 + b_1)\end{aligned}$$

And thus :

$$b_1 \leftarrow b_1 - \lambda \frac{\partial l}{\partial x_1} \sigma'(A_1x_0 + b_1)$$

We need to back-propagate the term $\frac{\partial l}{\partial x_1}$ computed for the second layer which has been computed with $\frac{\partial l}{\partial x_2}$ back-propagated from the first layer.

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

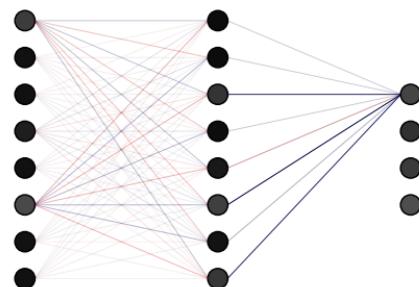
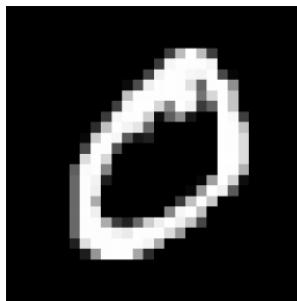
To update the weights, we need to compute the gradient of the loss with respect to the output of the network, and then **back-propagate** the gradient of the loss with respect to each activation, the $\frac{\partial l}{\partial x_i}$, through the network to compute the gradients with respect to the weights and biases of each layer.

THE MULTI-LAYER PERCEPTRON (MLP)

LAST LAYER

We can plot the current state of the network for a given input.

The red lines show positive values for A_i , the blue lines represent negative values for A_i . The level of transparency is proportional to the previous neurons.

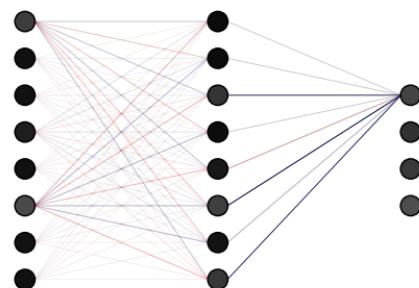
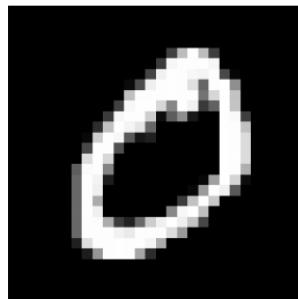


$$x_3^1 = \sigma_3 \left(A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \dots + A_3^{1,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

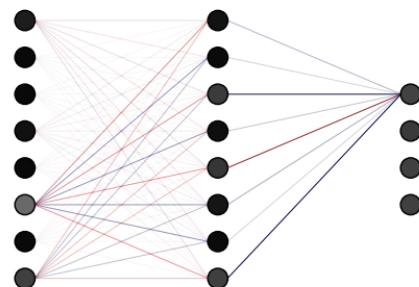
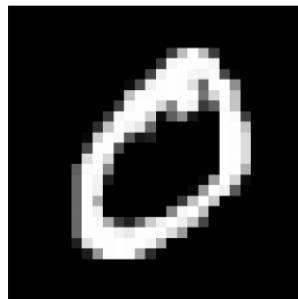


$$x_3^1 = \sigma_3 \left(A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

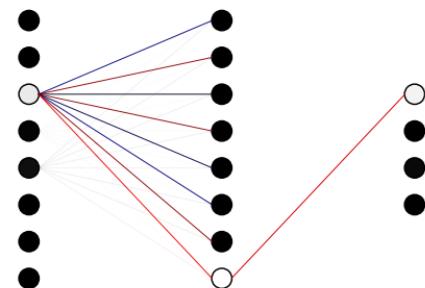
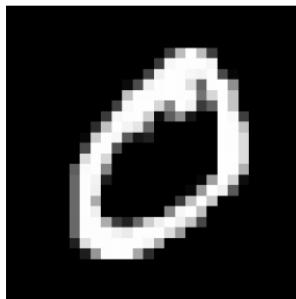


$$x_3^1 = \sigma_3 \left(A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

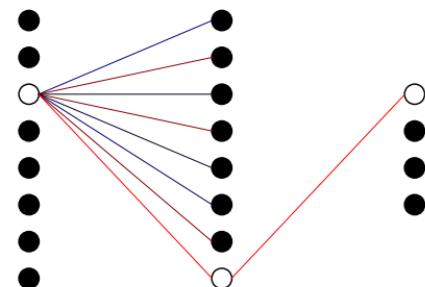
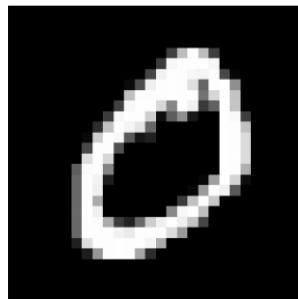


$$x_3^1 = \sigma_3 \left(A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.



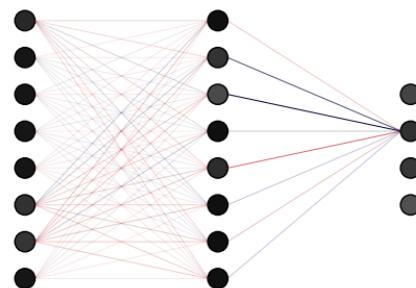
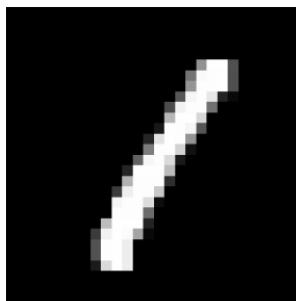
$$x_3^1 = \sigma_3 \left(A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

LAST LAYER

We can plot the current state of the network for a given input.

Red lines show positive values of A_i , Blue lines represent negative values of A_i . The level of transparency is proportional to the previous neurons.

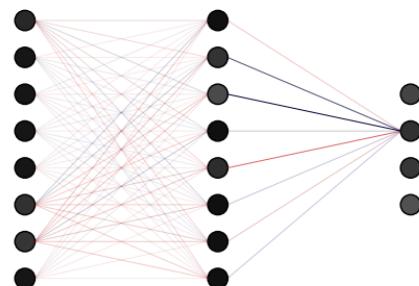
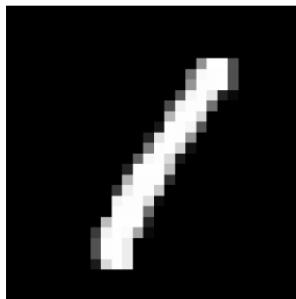


$$x_3^2 = \sigma_3 \left(A_3^{2,1}x_2^1 + A_3^{2,2}x_2^2 + \dots + A_3^{2,8}x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

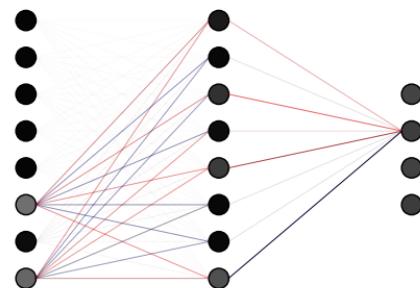
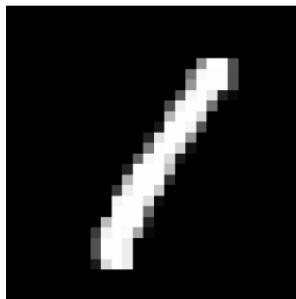


$$x_3^2 = \sigma_3 \left(A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

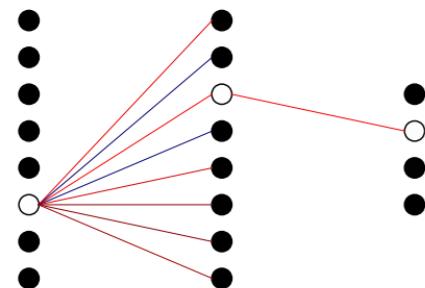
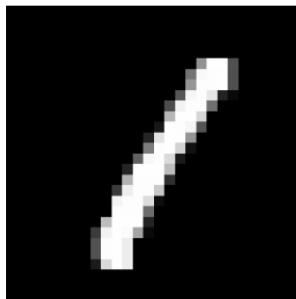


$$x_3^2 = \sigma_3 \left(A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

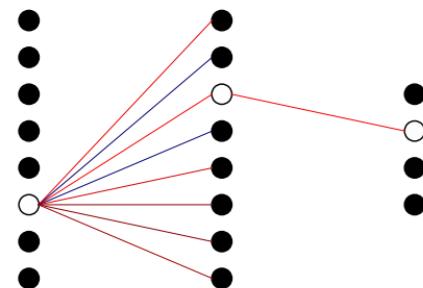
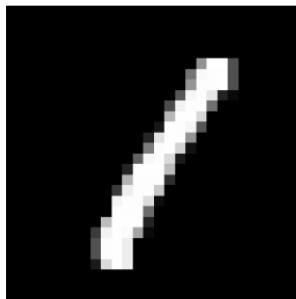


$$x_3^2 = \sigma_3 \left(A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

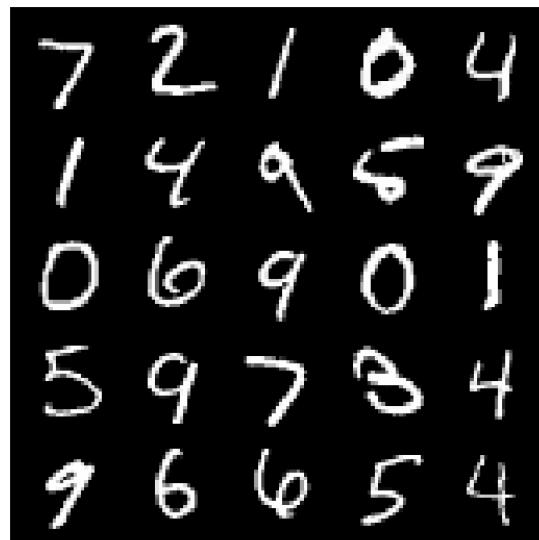


$$x_3^2 = \sigma_3 \left(A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

Having discussed the theory behind Artificial Neural Networks and the training process, we will now proceed to demonstrate a comprehensive end-to-end example of image classification on MNIST.



THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

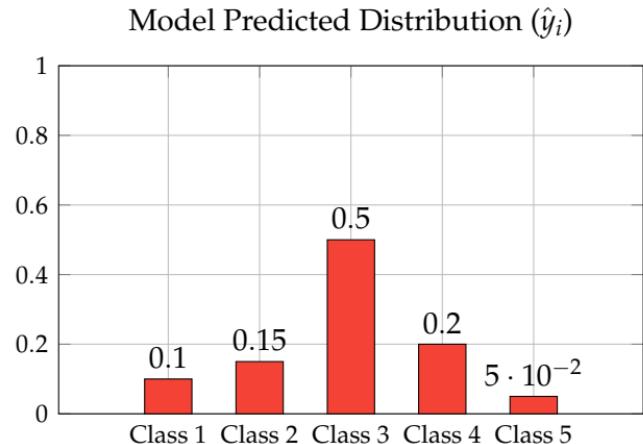
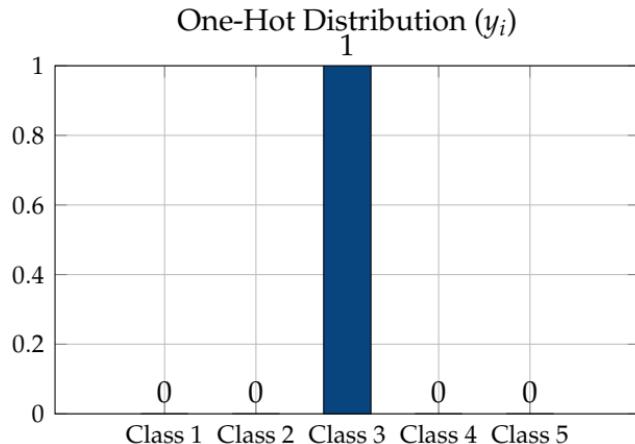
- ▶ Input shape : $1 \times 28 \times 28$.
- ▶ Number of Classes : 10.
- ▶ Number of training samples (x, y) : 60000.
- ▶ Number of evaluating samples: 10000.
- ▶ Loss : cross-entropy

$$L(\hat{y}, y) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij})$$

where :

- $\hat{y} \in \mathbb{R}^{N \times K}$ is the predicted probability distribution over K classes for N samples,
- $y \in \{0, 1\}^{N \times K}$ is the ground-truth one-hot encoded label matrix,

RECAP ON THE CROSS-ENTROPY LOSS



The cross-entropy loss for one sample is:

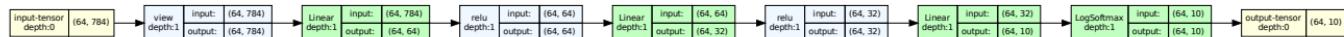
$$l(\hat{y}_i, y_i) = - \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij}).$$

THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

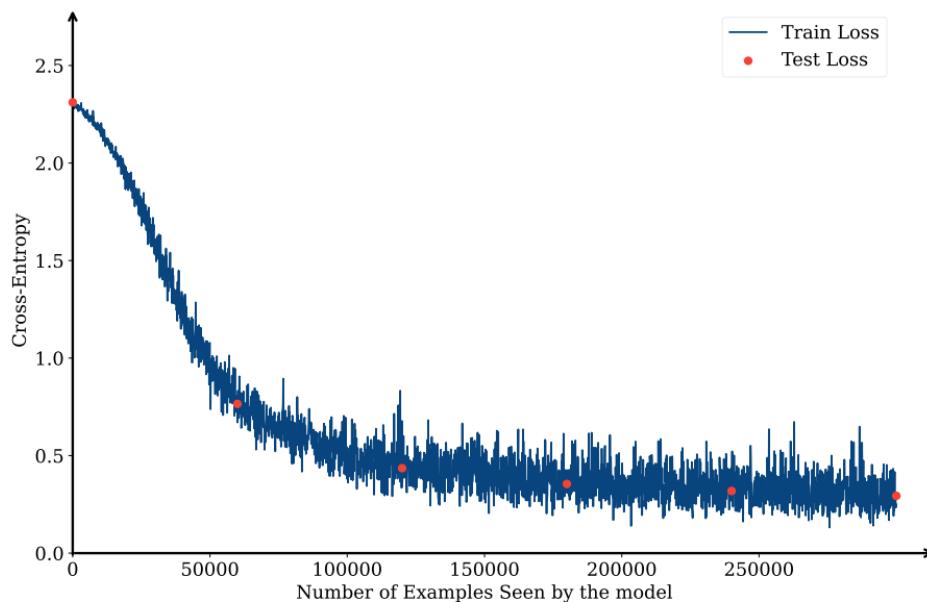
We build a 3 layers network.

- ▶ Batch size : 64
- ▶ Learning rate : 0.01
- ▶ Intermediate activation : ReLU
- ▶ Final activation : Softmax
- ▶ Number of epochs : 12
- ▶ Number of trained parameters: 52.6k



THE MULTI-LAYER PERCEPTRON (MLP)

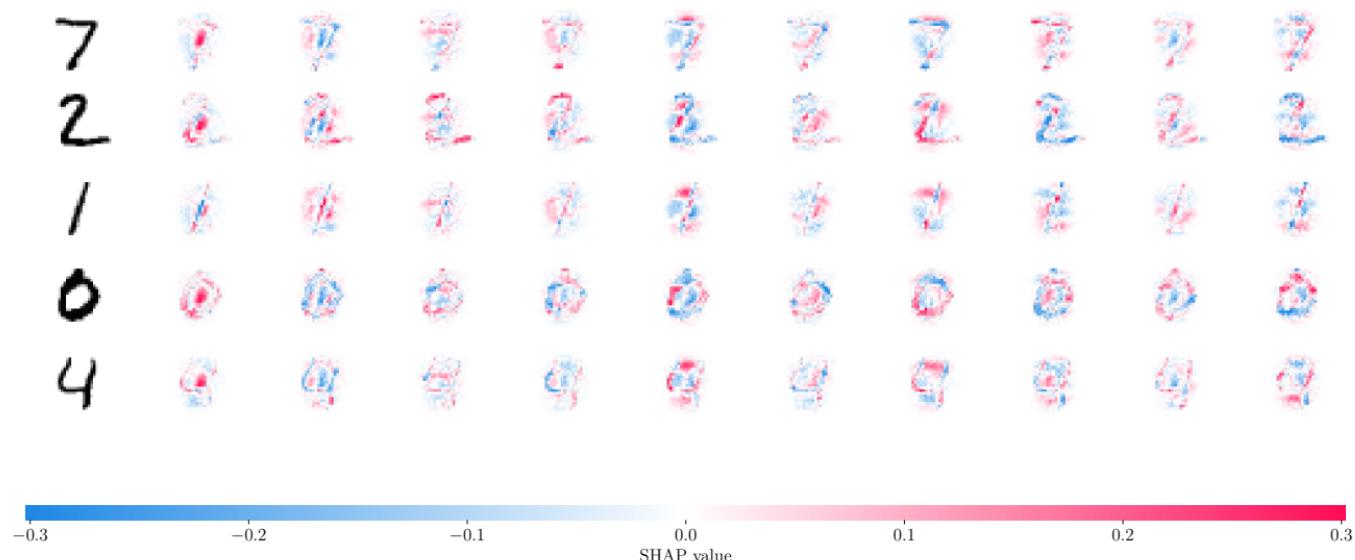
EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z



THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

With an interpretation tool such as SHAP:



TP1: THE MULTI-LAYER PERCEPTRON (MLP)

THE FIRST DEEP LEARNING MODEL

Link to the notebook (ipynb): [TP1.ipynb](#)

Link to the notebook (html): [TP1.html](#)

Part II

DEEP LEARNING IN ACTION: FROM NEURAL NETWORKS TO TRANSFORMER MODELS

Now that we have an understanding of the training procedure for Artificial Neural Networks, we shall examine several widely-utilized structures within the literature of Neural Networks, including Convolutional Neural Networks (CNN), Residual Networks (ResNet), Recurrent Neural Networks (RNN), and Transformers.

CONVOLUTIONAL NEURAL NETWORKS

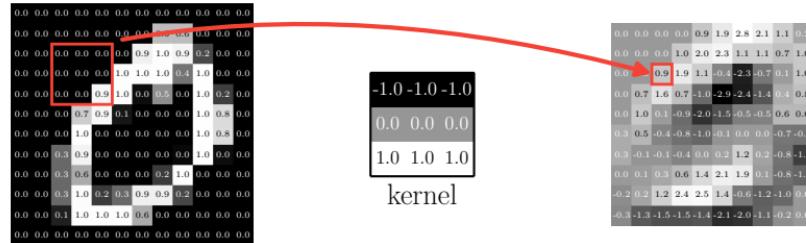
In the field of image processing, the Convolution Operators are widely considered as the most favoured approach. While it has been demonstrated that Dense blocks, or Linear blocks, are capable of accurately classifying images in the case of the MNIST dataset, the need for convolutional transformations arises when addressing wider and more intricate datasets.

CONVOLUTIONAL NEURAL NETWORKS

THE TWO DIMENSIONAL CONVOLUTION

A 2D convolution in a neural network context can be mathematically represented as a sliding window operation where a filter (also called kernel) w of size $k \times k$ is applied to each $k \times k$ sub-matrix of the input matrix x . The operation can be defined as the element-wise multiplication of the filter w and the sub-matrix followed by summing the results, i.e.

$$y_{i,j} = \sum_{m=1}^k \sum_{n=1}^k w_{m,n} \cdot x_{i+m,j+n}$$

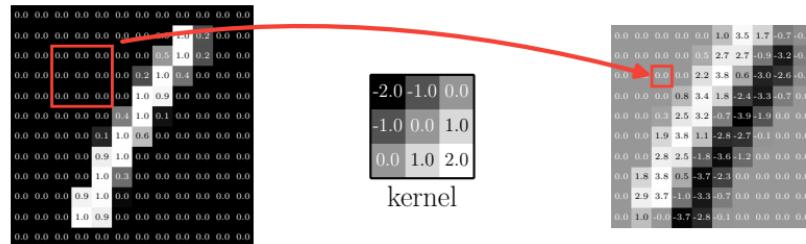


CONVOLUTIONAL NEURAL NETWORKS

THE TWO DIMENSIONAL CONVOLUTION

A 2D convolution in a neural network context can be mathematically represented as a sliding window operation where a filter (also called kernel) w of size $k \times k$ is applied to each $k \times k$ submatrix of the input matrix x . The operation can be defined as the element-wise multiplication of the filter w and the submatrix followed by summing the results, i.e.

$$y_{i,j} = \sum_{m=1}^k \sum_{n=1}^k w_{m,n} \cdot x_{i+m,j+n}$$



CONVOLUTIONAL NEURAL NETWORKS

KERNEL SIZE, PADDING AND STRIDE

In every Deep Learning library, the Conv2D block takes three parameters in argument:

- ▶ the Kernel's size,
- ▶ the Stride,
- ▶ the Padding.

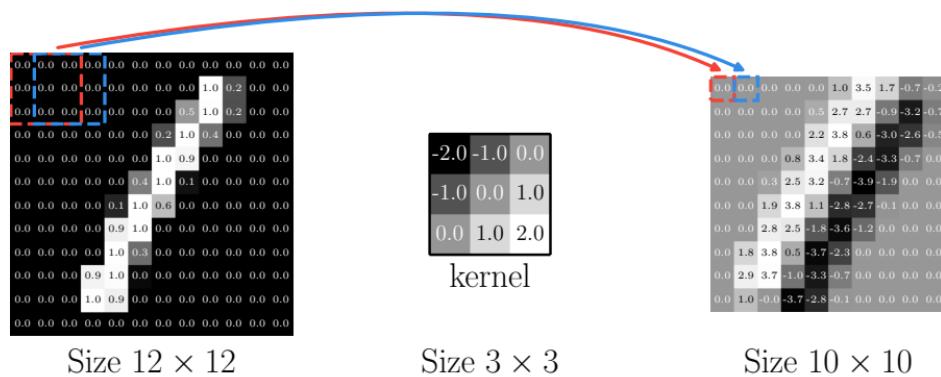
The size out the output is :

$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

CONVOLUTIONAL NEURAL NETWORKS

KERNEL SIZE, PADDING AND STRIDE

Kernel size: 3, Padding: 0, Stride: 1

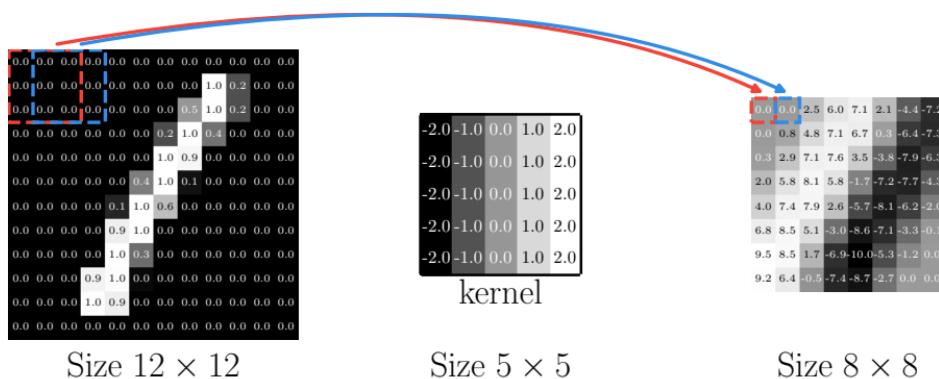


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

CONVOLUTIONAL NEURAL NETWORKS

KERNEL SIZE, PADDING AND STRIDE

Kernel size: 5, Padding: 0, Stride: 1

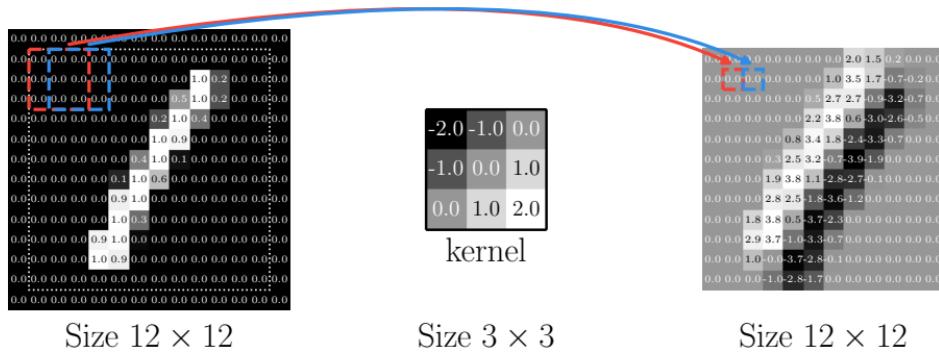


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

CONVOLUTIONAL NEURAL NETWORKS

KERNEL SIZE, PADDING AND STRIDE

Kernel size: 3, Padding: 1, Stride: 1. Padding mode can be 'zeros', 'reflect', 'replicate' or 'circular'.

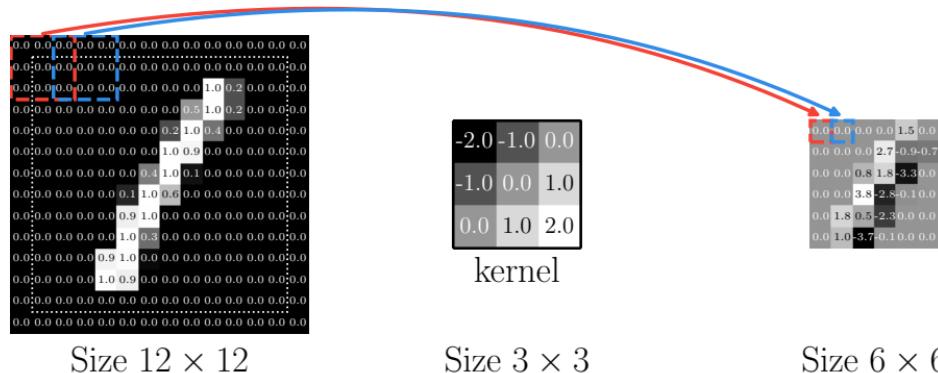


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

CONVOLUTIONAL NEURAL NETWORKS

KERNEL SIZE, PADDING AND STRIDE

Kernel size: 3, Padding: 1, Stride: 2

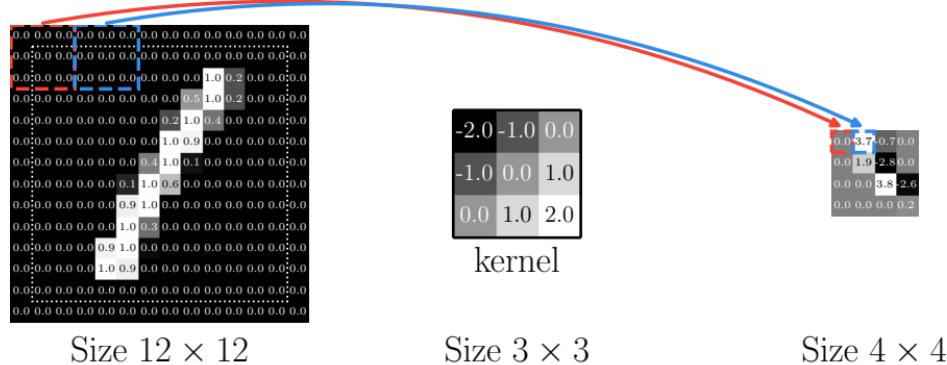


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

CONVOLUTIONAL NEURAL NETWORKS

KERNEL SIZE, PADDING AND STRIDE

Kernel size: 3, Padding: 1, Stride: 3

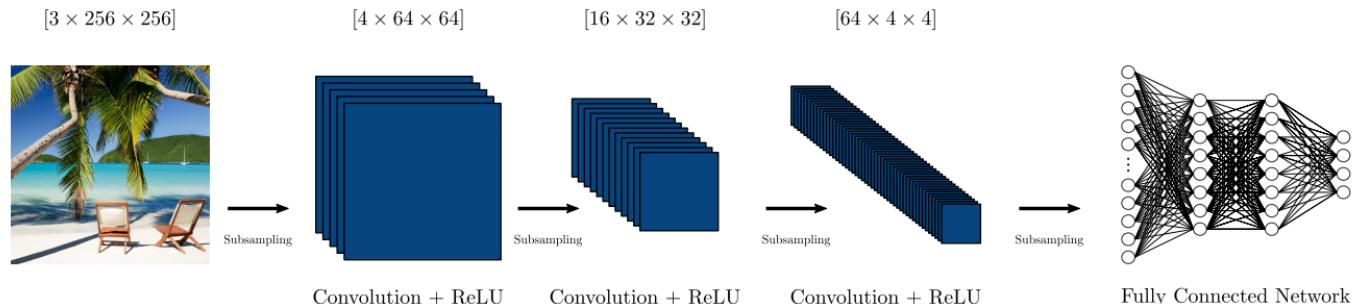


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

CONVOLUTIONAL NEURAL NETWORKS

CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

We can represent a CNN as under this form:

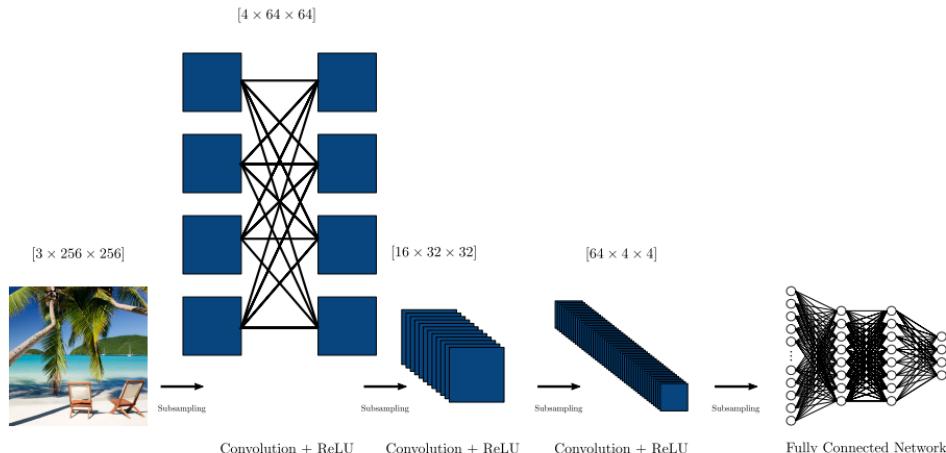


CONVOLUTIONAL NEURAL NETWORKS

CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

Usually, the output of a convolutional block is linear combination of the Convolutional output of every previous channels and a bias:

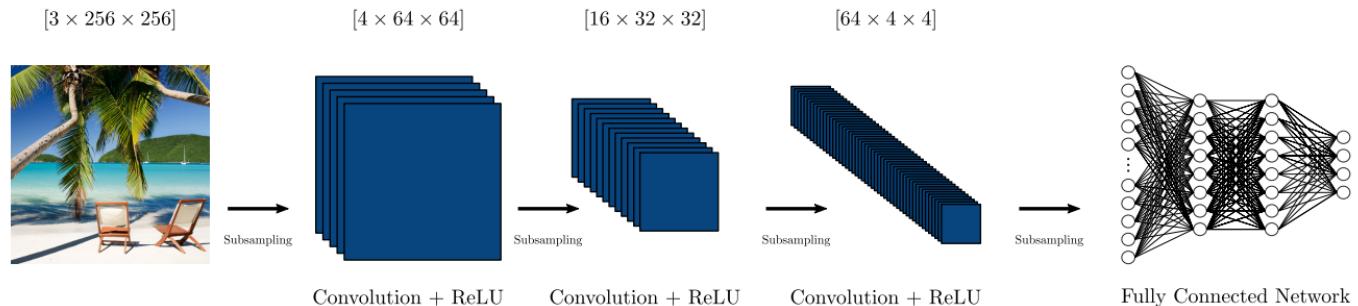
$$\text{out}_{i,j}(c_{\text{out}}) = \text{bias}(c_{\text{out}}) + \sum_{k=0}^{|c_{\text{in}}|-1} \text{Conv}(\text{input}(k), \text{kernel}_k)_{i,j}$$



CONVOLUTIONAL NEURAL NETWORKS

CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

In practice, we split the image into multiple channels : the three channels RGB to begin with. Then we apply convolutional operation on different scales and then we use a fully connected tail. To change the scale we can use different sub-sampling : Max pooling, Average pooling or Invertible pooling.



CONVOLUTIONAL NEURAL NETWORKS

MAX POOLING

Max Pooling takes the maximum within a given sized sub-matrix. In practice, the matrix is size 2×2 in order to reduce the dimension by 4 and doubling the scale.

0.2	1.0	0.3	0.8	0.1	0.8	0.6	0.9
0.7	0.3	0.3	0.5	0.4	0.3	0.3	0.4
0.2	0.6	0.7	0.9	0.9	0.1	0.3	0.5
0.8	0.7	0.1	0.3	0.3	0.6	0.9	0.5
0.7	0.1	0.1	0.2	0.8	0.4	0.9	0.7
0.9	0.1	0.8	0.9	0.6	0.3	0.1	0.9
0.8	0.7	0.2	0.7	0.0	0.6	0.9	0.5
0.9	0.5	0.1	0.2	0.0	0.1	0.1	0.4

Max Pooling
Subsampling

1.0	0.8	0.8	0.9
0.8	0.9	0.9	0.9
0.9	0.9	0.8	0.9
0.9	0.7	0.6	0.9

CONVOLUTIONAL NEURAL NETWORKS

AVERAGE POOLING

The Average pooling takes the average value within the sub-matrix.

0.2	1.0	0.3	0.8	0.1	0.8	0.6	0.9
0.7	0.3	0.3	0.5	0.4	0.3	0.3	0.4
0.2	0.6	0.7	0.9	0.9	0.1	0.3	0.5
0.8	0.7	0.1	0.3	0.3	0.6	0.9	0.5
0.7	0.1	0.1	0.2	0.8	0.4	0.9	0.7
0.9	0.1	0.8	0.9	0.6	0.3	0.1	0.9
0.8	0.7	0.2	0.7	0.0	0.6	0.9	0.5
0.9	0.5	0.1	0.2	0.0	0.1	0.1	0.4

Average Pooling
Subsampling

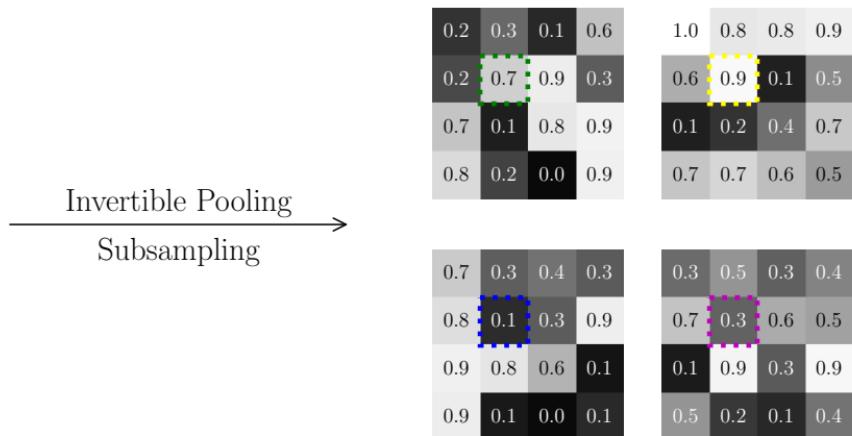
0.5	0.5	0.4	0.6
0.6	0.5	0.5	0.5
0.4	0.5	0.5	0.6
0.7	0.3	0.2	0.5

CONVOLUTIONAL NEURAL NETWORKS

INVERTIBLE POOLING

For Invertible Networks, we can use Invertible Pooling, aka Squeeze. It preserves the information contained in the channels and keeps the dimension constant.

0.2	1.0	0.3	0.8	0.1	0.8	0.6	0.9
0.7	0.3	0.3	0.5	0.4	0.3	0.3	0.4
0.2	0.6	0.7	0.9	0.9	0.1	0.3	0.5
0.8	0.7	0.1	0.3	0.3	0.6	0.9	0.5
0.7	0.1	0.1	0.2	0.8	0.4	0.9	0.7
0.9	0.1	0.8	0.9	0.6	0.3	0.1	0.9
0.8	0.7	0.2	0.7	0.0	0.6	0.9	0.5
0.9	0.5	0.1	0.2	0.0	0.1	0.1	0.4



CONVOLUTIONAL NEURAL NETWORKS

CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

Convolutional Neural Networks are more suitable for image processing compared to fully connected networks due to their ability to efficiently handle the spatial relationships between pixels in an image. This is achieved through the use of convolutional layers that apply filters to small portions of an image, rather than fully connected layers that process the entire image as a single vector. Additionally, the shared weights in convolutional layers allow for learning of hierarchical features, reducing the number of parameters in the network and increasing its ability to generalize to new images.

CONVOLUTIONAL NEURAL NETWORKS

CNN IN PRACTICE: CIFAR 10

- ▶ Input shape : $3 \times 32 \times 32$.
- ▶ Number of Classes : 10.
- ▶ Number of training samples (x, y): 50000.
- ▶ Number of evaluating samples: 10000.



CONVOLUTIONAL NEURAL NETWORKS

CNN IN PRACTICE: CIFAR 10

We will compare three different models:

- ▶ Model 1 : Fully Connected Neural Network with 3.4 million parameters.
- ▶ Model 2 : CNN with 62 thousand parameters.
- ▶ Model 3 : Wider and longer CNN with 5.8 million parameters.

CONVOLUTIONAL NEURAL NETWORKS

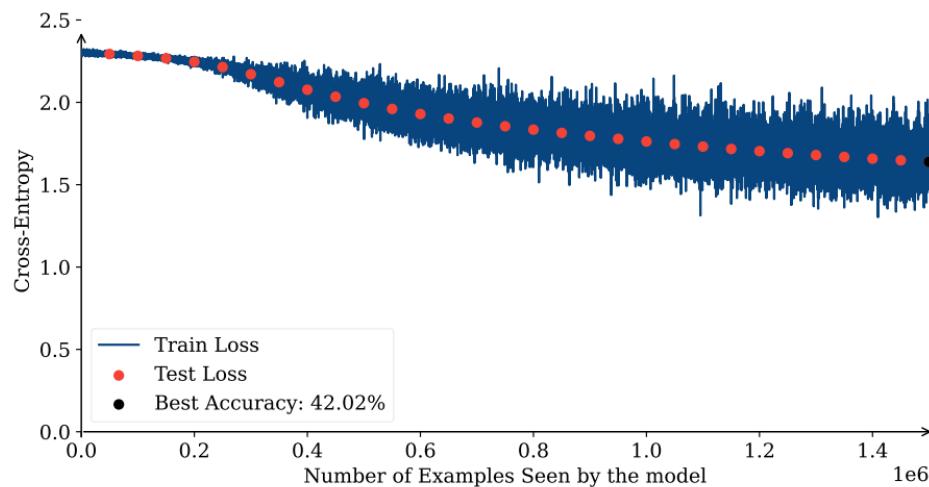
MODEL 1 : FULLY CONNECTED NEURAL NETWORK

The Net is composed of 4 linear layers with ReLU activations:

- ▶ Linear $3072 \mapsto 1024 + \text{ReLU}$
- ▶ Linear $1024 \mapsto 256 + \text{ReLU}$
- ▶ Linear $256 \mapsto 64 + \text{ReLU}$
- ▶ Linear $64 \mapsto 10 + \text{SoftMax}$

CONVOLUTIONAL NEURAL NETWORKS

MODEL 1



CONVOLUTIONAL NEURAL NETWORKS

USING A BETTER OPTIMIZATION ALGORITHM ?

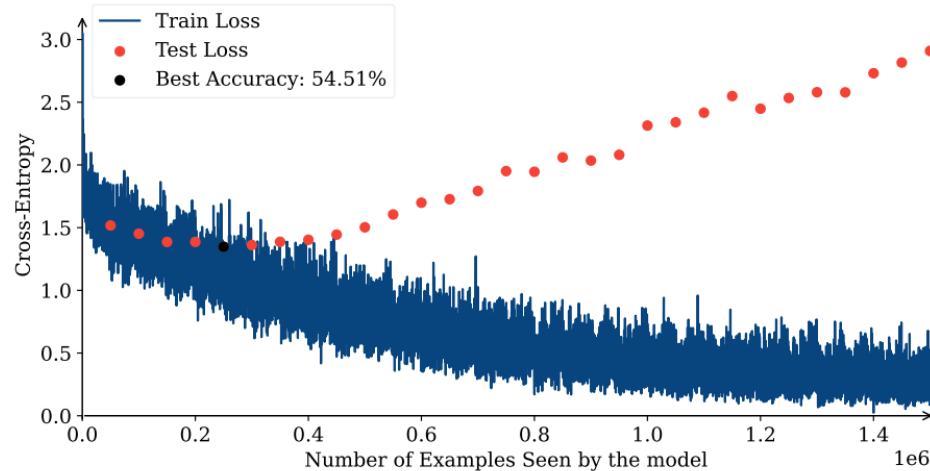
Adam: Adam is an optimization algorithm that can be used instead of the classical stochastic gradient descent procedure to update network weights iteratively based on training data. It combines the advantages of two other extensions of stochastic gradient descent, AdaGrad and RMSProp, and is designed to work well for a wide range of applications:

$$\begin{aligned}\theta_{t+1} &= \theta_t - \frac{\lambda}{\sqrt{\hat{v}_t + \varepsilon}} \hat{m}_t \\ \hat{m}_t &= \beta_1 m_{t-1} + (1 - \beta_1) \nabla_\theta l(\theta_t) \\ \hat{v}_t &= \beta_2 v_{t-1} + (1 - \beta_2) \nabla_\theta l(\theta_t)^2\end{aligned}$$

- ▶ λ is the learning rate.
- ▶ \hat{m}_t is the moving average of the gradient. It allows to smooth the gradient step.
- ▶ By dividing the gradient by the square root of the moving average of the squared gradient, we can adapt the learning rate for each parameter to improve the convergence speed.

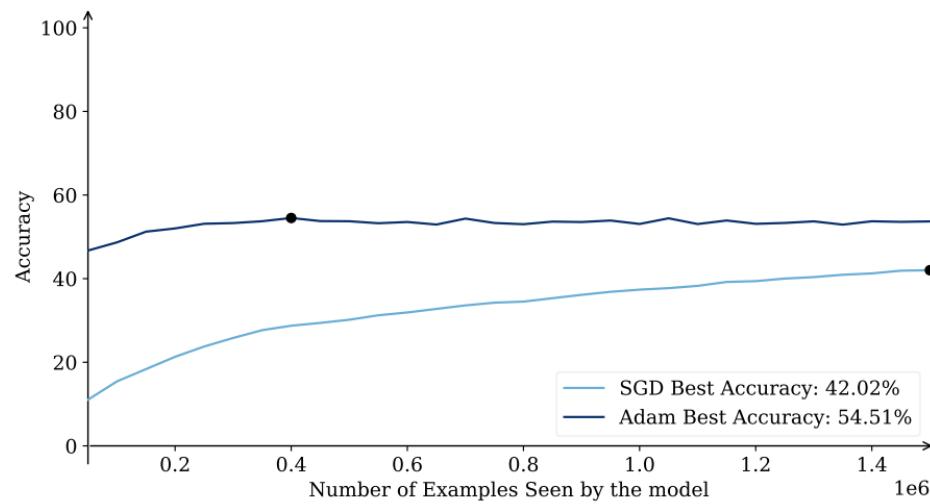
CONVOLUTIONAL NEURAL NETWORKS

MODEL 1



CONVOLUTIONAL NEURAL NETWORKS

MODEL 1



CONVOLUTIONAL NEURAL NETWORKS

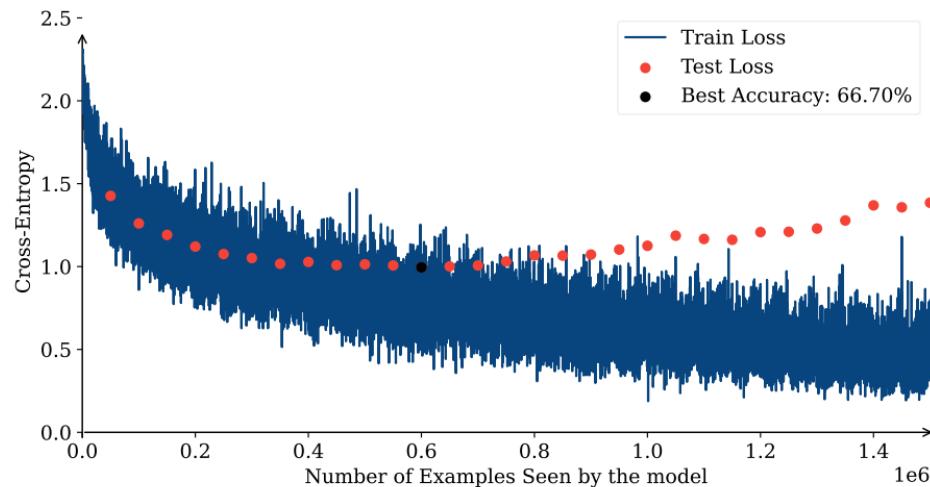
MODEL 2

The Net is composed 2 convolutional layers and 2 linear layers:

- ▶ Conv $3 \times 32 \times 32 \mapsto 6 \times 28 \times 28$ + ReLU
- ▶ Max Pooling $6 \times 28 \times 28 \mapsto 6 \times 14 \times 14$
- ▶ Conv $6 \times 14 \times 14 \mapsto 16 \times 10 \times 10$ + ReLU
- ▶ Max Pooling $16 \times 10 \times 10 \mapsto 16 \times 5 \times 5$
- ▶ Linear $400 \mapsto 120$ + ReLU
- ▶ Linear $120 \mapsto 84$ + ReLU
- ▶ Linear $84 \mapsto 10$ + SoftMax

CONVOLUTIONAL NEURAL NETWORKS

MODEL 2



CONVOLUTIONAL NEURAL NETWORKS

MODEL 3

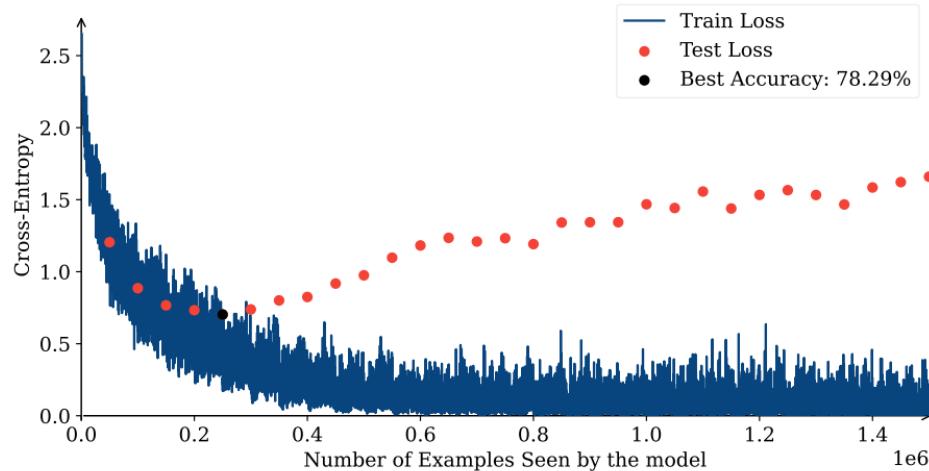
The Net is composed 6 convolutional layers and 3 linear layers:

- ▶ Conv $3 \times 32 \times 32 \mapsto 32 \times 32 \times 32$ + BatchNorm2d + ReLU
- ▶ Conv $32 \times 32 \times 32 \mapsto 64 \times 32 \times 32$ + ReLU
- ▶ Max Pooling $64 \times 32 \times 32 \mapsto 64 \times 16 \times 16$
- ▶ Conv $64 \times 16 \times 16 \mapsto 128 \times 16 \times 16$ + BatchNorm2d + ReLU
- ▶ Conv $128 \times 16 \times 16 \mapsto 128 \times 16 \times 16$ + ReLU
- ▶ Max Pooling $128 \times 16 \times 16 \mapsto 128 \times 8 \times 8$
- ▶ Conv $128 \times 8 \times 8 \mapsto 256 \times 8 \times 8$ + BatchNorm2d + ReLU
- ▶ Conv $256 \times 8 \times 8 \mapsto 256 \times 8 \times 8$ + ReLU
- ▶ Max Pooling $256 \times 8 \times 8 \mapsto 256 \times 4 \times 4$ + DropOut $p = 0.05$
- ▶ Linear $4096 \mapsto 1024$ + ReLU
- ▶ Linear $1024 \mapsto 512$ + ReLU + DropOut $p = 0.05$
- ▶ Linear $512 \mapsto 10$ + SoftMax

We have added Batch Normalization to improve the training stability and Drop Out to reduce overfitting.

CONVOLUTIONAL NEURAL NETWORKS

MODEL 3 WITHOUT BATCH NORMALIZATION AND DROP OUT



CONVOLUTIONAL NEURAL NETWORKS

BATCH NORMALIZATION

Batch normalization is used to normalize the activations of a layer within a batch of data.

- ▶ This helps to prevent the problem of vanishing or exploding gradients and also speeds up the training process.
- ▶ By normalizing the activations, batch normalization helps to stabilize the distribution of the inputs to each layer, reducing the covariate shift and allowing the network to learn more effectively.

CONVOLUTIONAL NEURAL NETWORKS

BATCH NORMALIZATION

- 1: **for** each x_i in a mini-batch B of size b **do**
- 2: Compute the mean μ_B and variance σ_B^2 of the features in the mini-batch B .

$$\mu_B = \frac{1}{b} \sum_i x_i \quad \text{and} \quad \sigma_B^2 = \frac{1}{m} \sum_i (x_i - \mu_B)^2$$

- 3: Normalize each feature x_i in the mini-batch B using μ_B and σ_B^2 .

$$\bar{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}$$

- 4: Scale and shift each normalized feature x_i using two learnable parameters γ and β respectively.

$$y_i = \gamma \bar{x}_i + \beta$$

- 5: Update the moving average of the mean and variance for inference:

$$\mu = (1 - \gamma)\mu + \gamma\mu_B \quad \text{and} \quad \sigma^2 = (1 - \gamma)\sigma^2 + \gamma\sigma_B^2$$

- 6: **end for**

Algorithm 1: Batch Normalization during training

CONVOLUTIONAL NEURAL NETWORKS

BATCH NORMALIZATION

At inference:

1: **for** each x_i **do**

2: Normalize each feature x_i using the moving average of the mean and variance:

$$\bar{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \varepsilon}}$$

3: Scale and shift each normalized feature x_i using two learnable parameters γ and β respectively.

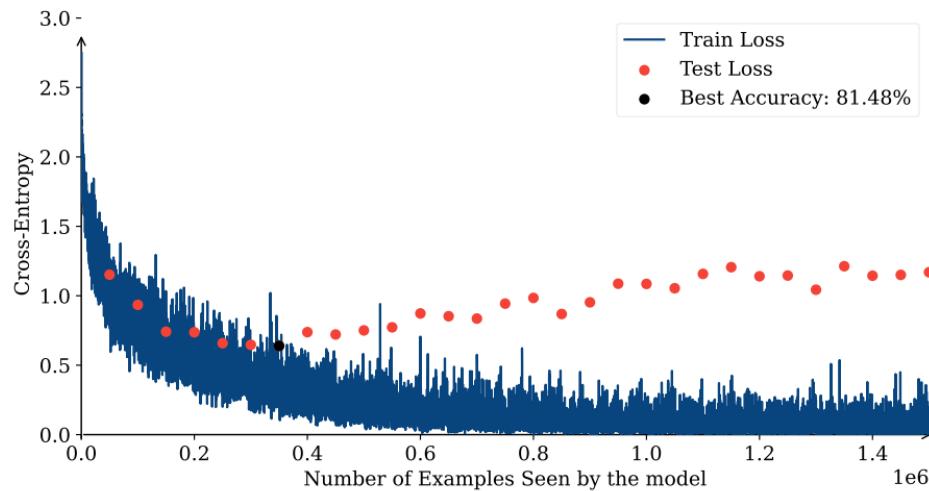
$$y_i = \gamma \bar{x}_i + \beta$$

4: **end for**

Algorithm 2: Batch Normalization at inference

CONVOLUTIONAL NEURAL NETWORKS

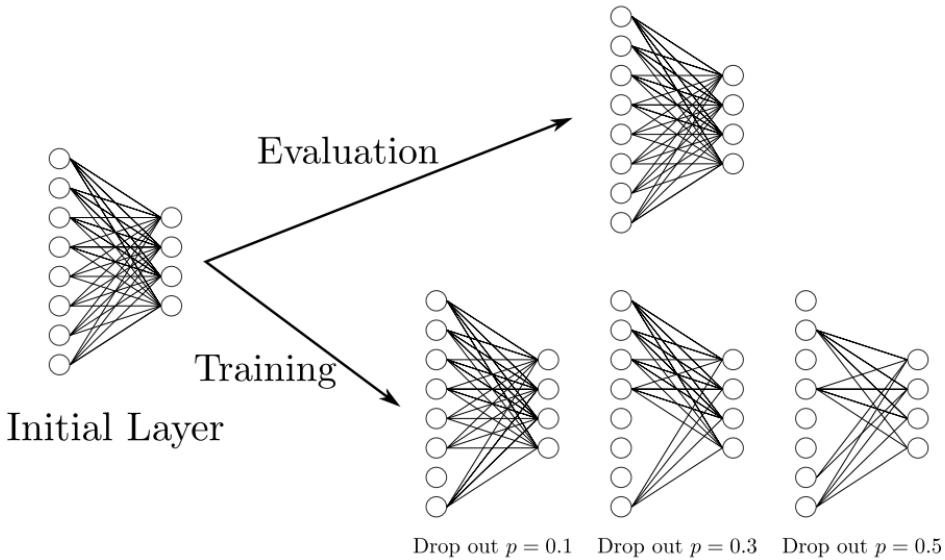
MODEL 3 WITHOUT DROP OUT



CONVOLUTIONAL NEURAL NETWORKS

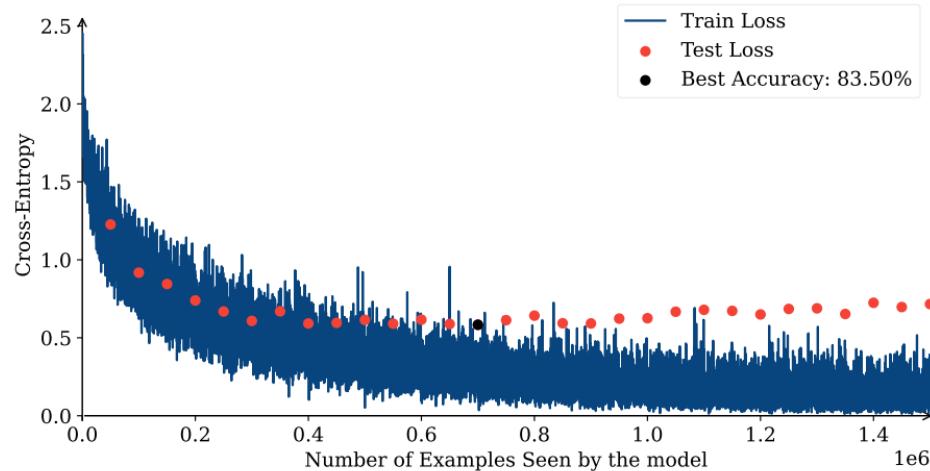
DROP OUT

Dropout is a regularization technique in neural networks where *during training*, a portion of the nodes are randomly "dropped out" or ignored during each iteration. This helps prevent over-fitting by preventing the model from relying too heavily on any one node. The result is a more robust and generalizable model that can better handle unseen data.



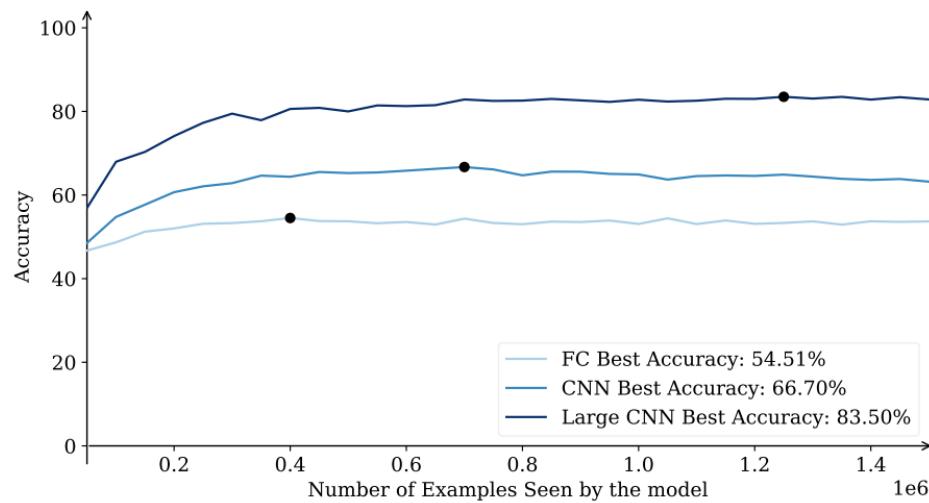
CONVOLUTIONAL NEURAL NETWORKS

MODEL 3



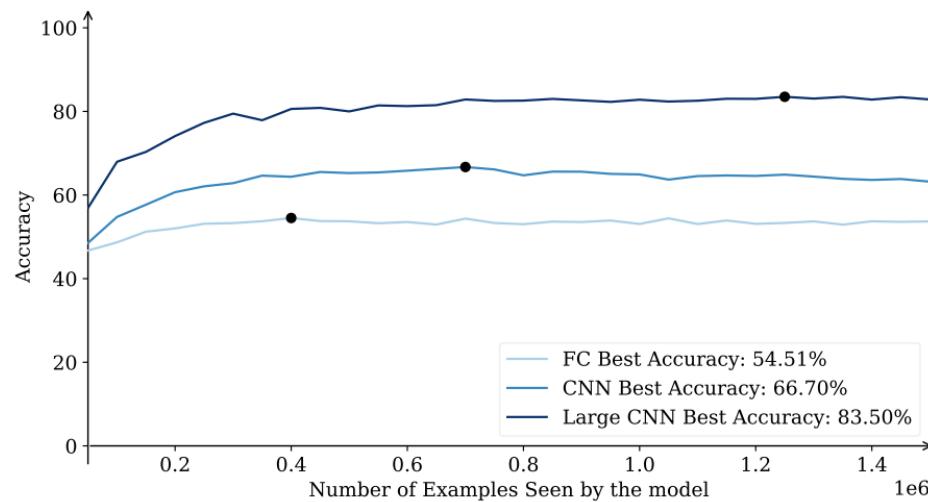
CONVOLUTIONAL NEURAL NETWORKS

MODEL 3



CONVOLUTIONAL NEURAL NETWORKS

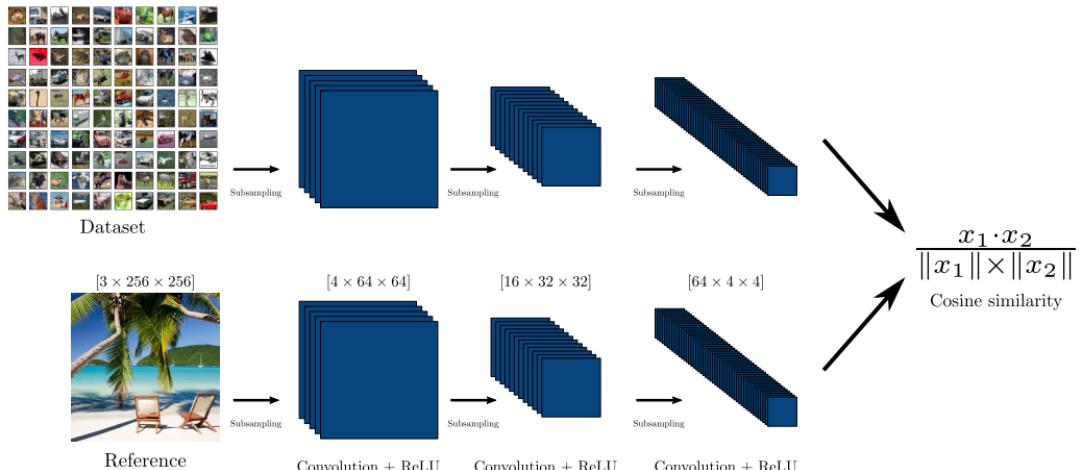
CNN IN PRACTICE: CIFAR 10



CONVOLUTIONAL NEURAL NETWORKS

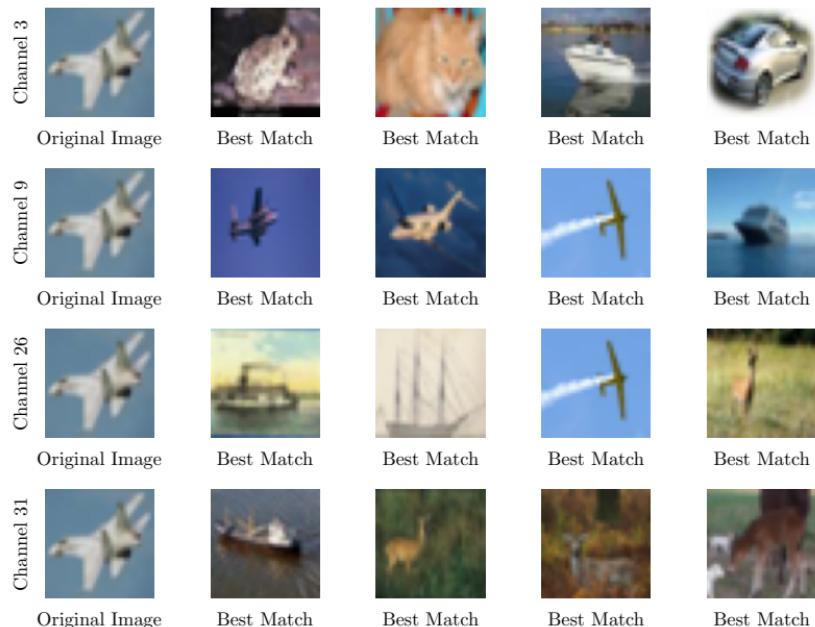
INTUITION BEHIND CHANNELS

To examine the information captured by different channels in a Neural Network, we can compare their output on a dataset. For a given input x , we can compute the similarity between the output of a specific channel and the same channel for other images in the dataset.



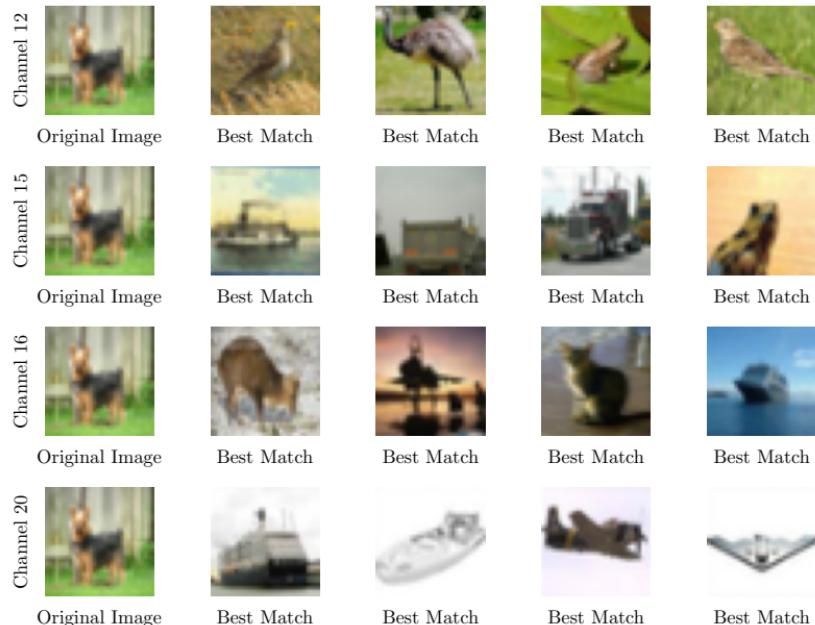
CONVOLUTIONAL NEURAL NETWORKS

INTUITION BEHIND CHANNELS



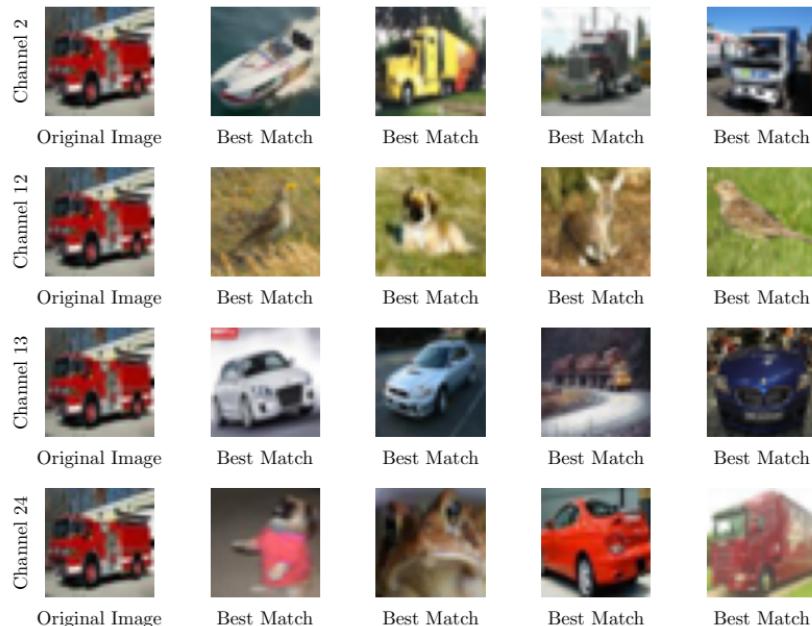
CONVOLUTIONAL NEURAL NETWORKS

INTUITION BEHIND CHANNELS



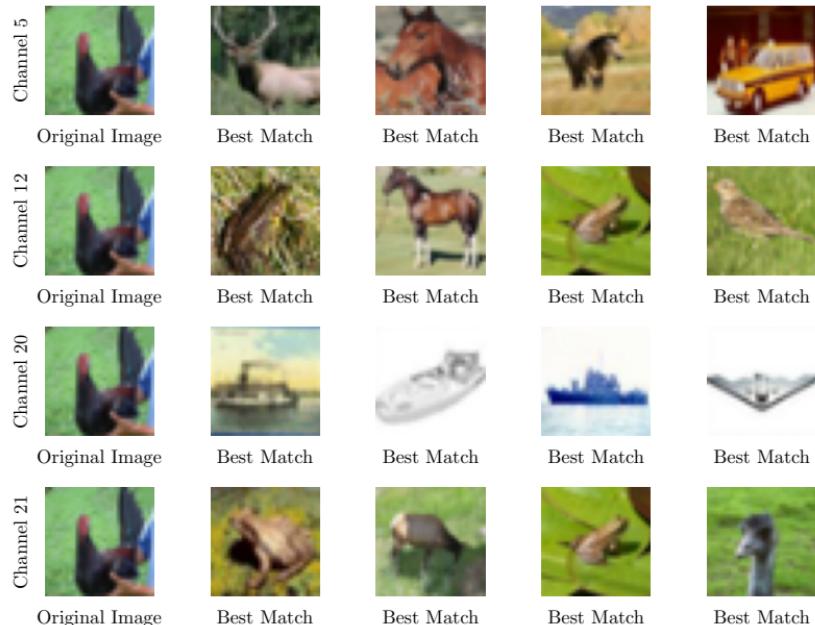
CONVOLUTIONAL NEURAL NETWORKS

INTUITION BEHIND CHANNELS



CONVOLUTIONAL NEURAL NETWORKS

INTUITION BEHIND CHANNELS



TP2: THE CONVOLUTIONAL NEURAL NETWORK

INTRODUCTION TO CNN AND CIFAR-10 DATASET

Link to the notebook (ipynb): TP2.ipynb

Link to the notebook (html): TP2.html

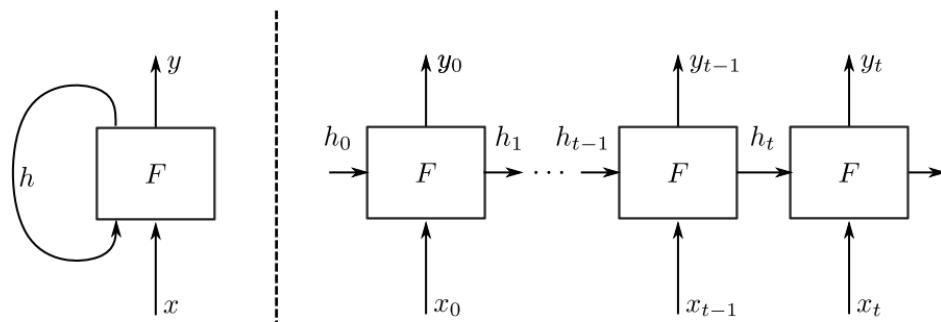
RECURRENT NEURAL NETWORKS

Recurrent Networks (RNNs) are a type of neural network that are specifically designed to handle sequential data, whereas CNNs are more suited for image and grid-like data. The main difference between RNNs and CNNs lies in the way they process data, with RNNs considering the sequence of elements and their interdependencies, while CNNs focus on capturing local patterns within the input.

RECURRENT NEURAL NETWORKS

RECURRENT BLOCK

A Recurrent Network is a type of neural network that contains a loop mechanism, allowing previous outputs to be used as inputs for future computations. This creates a form of memory that allows the network to process sequential data with variable-length sequences.



Rolled

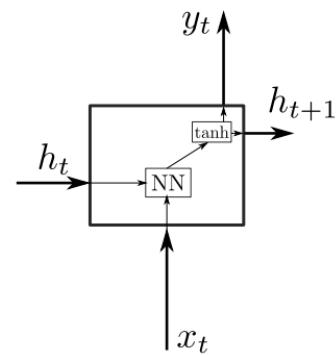
Unrolled

RECURRENT NEURAL NETWORKS

RECURRENT BLOCK

Some of the limitations of Vanilla RNNs:

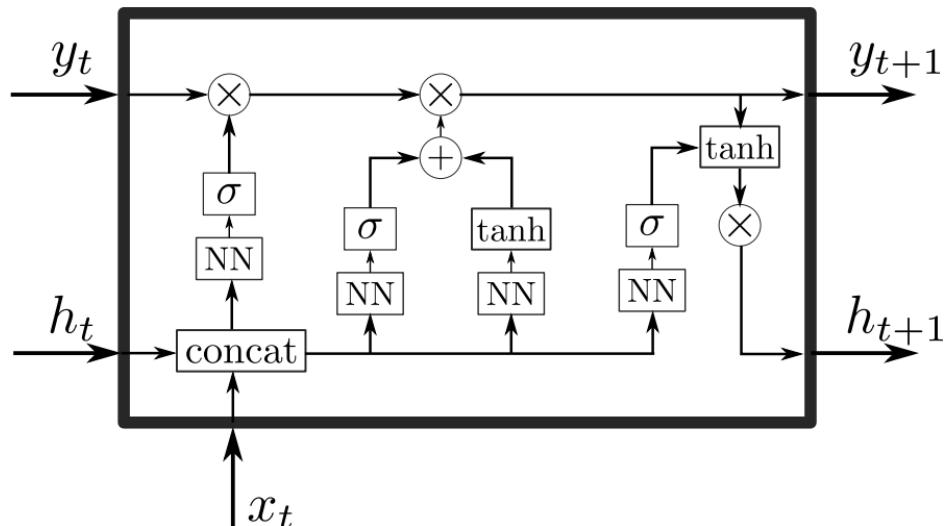
- ▶ Vanishing gradient problem: The gradient signals used to update the weights during training can become very small, making it difficult to train RNNs effectively.
- ▶ Exploding gradient problem: On the other hand, gradients can become too large and cause numeric instability, making it difficult to train RNNs effectively.
- ▶ Short-term memory: Vanilla RNNs have difficulty retaining information over long periods of time, making them unsuitable for tasks that require remembering information from previous time steps.
- ▶ Computational limitations: RNNs can be computationally intensive, making it difficult to apply them to large sequences of data.
- ▶ Difficulty with parallelization: The sequential nature of RNNs can make it difficult to take advantage of parallel processing to speed up training and inference.



RECURRENT NEURAL NETWORKS

LSTM

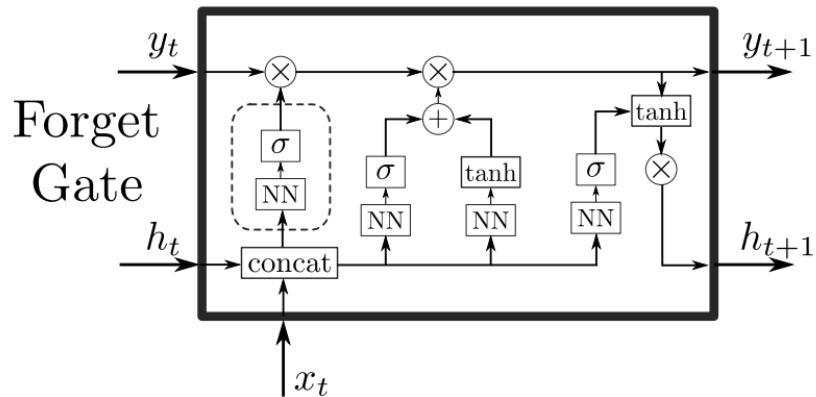
Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the vanishing gradient problem and difficulty in learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating cells, which allows them to selectively preserve information from previous time steps.



RECURRENT NEURAL NETWORKS

LSTM

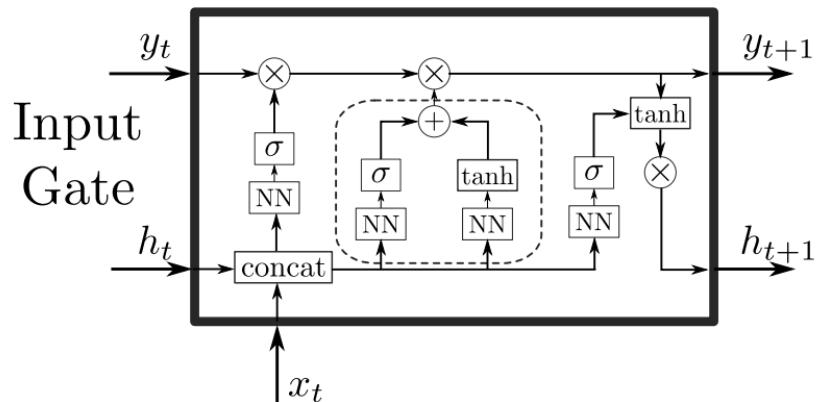
Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the problem of vanishing gradients and the difficulty of learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating cells, which allows them to selectively preserve information from previous time steps.



RECURRENT NEURAL NETWORKS

LSTM

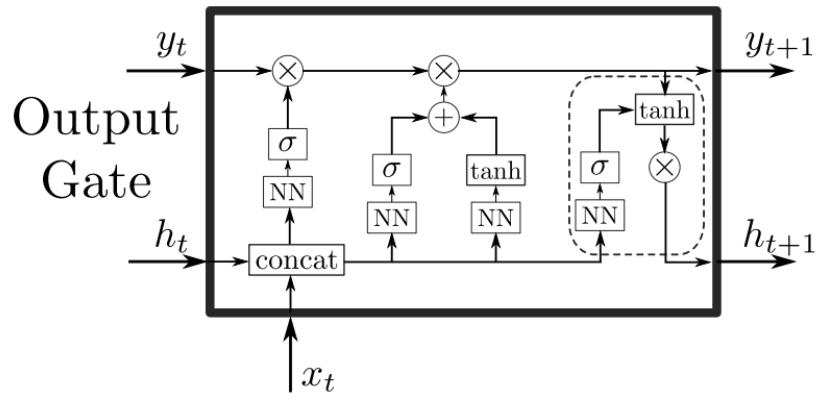
Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the vanishing gradient problem and difficulty in learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating the cells, which allows them to selectively preserve information from previous time steps.



RECURRENT NEURAL NETWORKS

LSTM

Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the vanishing gradient problem and difficulty in learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating the cells, which allows them to selectively preserve information from previous time steps.



RECURRENT NEURAL NETWORKS

LIMITS OF LSTM

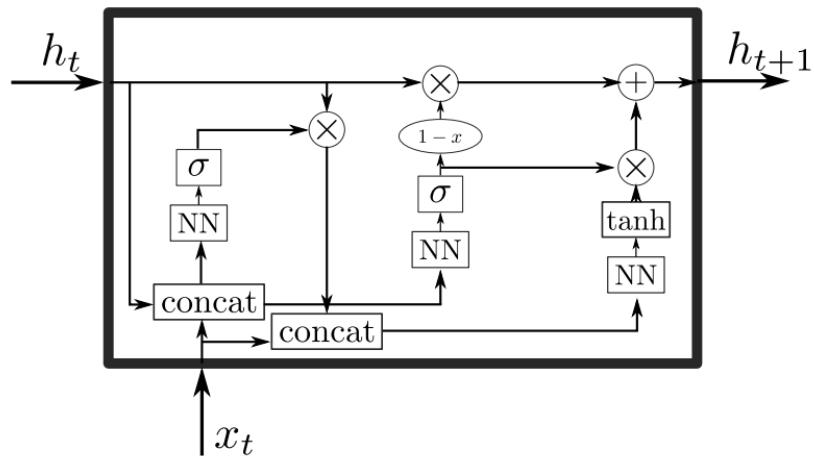
Limitations of LSTM RNNs:

- ▶ High computational cost: LSTMs are computationally more expensive compared to other traditional neural network models due to the presence of multiple gates and their sequential processing nature.
- ▶ Vanishing Gradient Problem: LSTMs, like any other RNNs, are prone to the vanishing gradient problem when the sequences are too long, making it difficult for the model to learn long-term dependencies.
- ▶ Overfitting: LSTMs are complex models and are more susceptible to overfitting compared to simple feedforward networks.
- ▶ Difficult to parallelize: Due to the sequential nature of LSTMs, they are difficult to parallelize and can take longer to train.
- ▶ Gradient Explosion: LSTMs can also suffer from the gradient explosion problem, where the gradients can become too large and cause numerical instability during training.

RECURRENT NEURAL NETWORKS

GRU

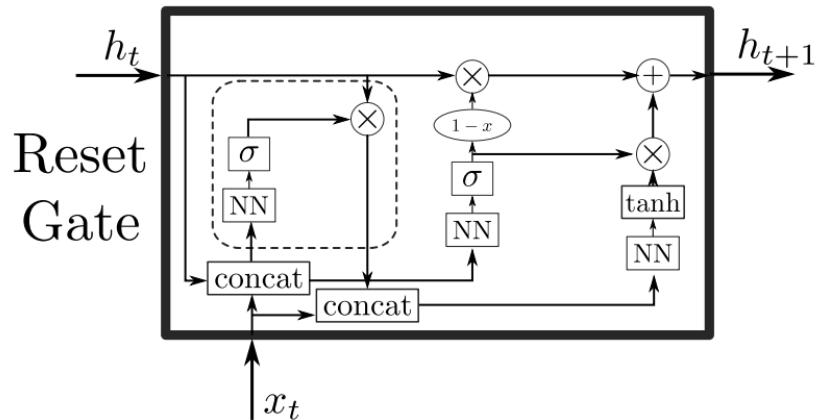
GRU blocks, or Gated Recurrent Units, are a type of recurrent neural network architecture that are similar to LSTMs in their function and ability to process sequential data. GRUs were introduced as a simplification of LSTMs, with the aim of reducing the number of parameters in the network and improving computational efficiency. GRUs achieve this by merging the forget and input gates in LSTMs into a single update gate, effectively combining the two operations in a single step.



RECURRENT NEURAL NETWORKS

GRU

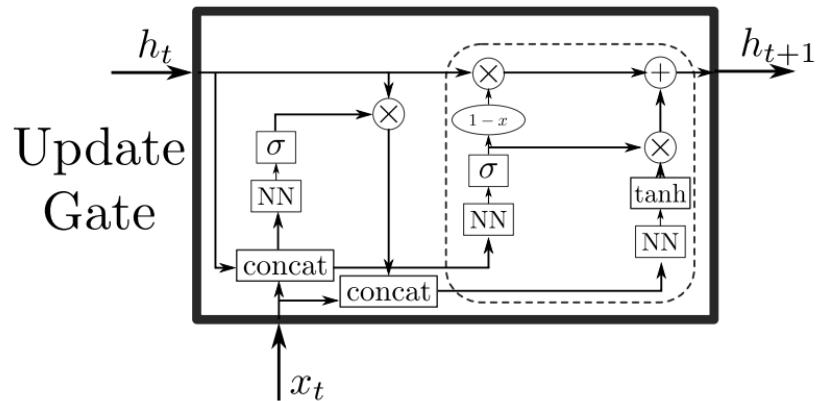
GRU blocks, or Gated Recurrent Units, are a type of recurrent neural network architecture that are similar to LSTMs in their function and ability to process sequential data. GRUs were introduced as a simplification of LSTMs, with the aim of reducing the number of parameters in the network and improving computational efficiency. GRUs achieve this by merging the forget and input gates in LSTMs into a single update gate, effectively combining the two operations in a single step.



RECURRENT NEURAL NETWORKS

GRU

GRU blocks, or Gated Recurrent Units, are a type of recurrent neural network architecture that are similar to LSTMs in their function and ability to process sequential data. GRUs were introduced as a simplification of LSTMs, with the aim of reducing the number of parameters in the network and improving computational efficiency. GRUs achieve this by merging the forget and input gates in LSTMs into a single update gate, effectively combining the two operations in a single step.



RECURRENT NEURAL NETWORKS

LSTM AND GRU

Limitations of GRU RNNs:

- ▶ Computational complexity: GRUs are more computationally efficient than LSTMs but still more complex than feedforward neural networks.
- ▶ Long-term dependencies: GRUs may struggle with capturing long-term dependencies in sequences, although they perform better in this regard than vanilla RNNs.
- ▶ Vanishing gradient problem: GRUs can still be affected by the vanishing gradient problem that plagues all RNN models. This problem makes it difficult for the model to learn from long sequences.
- ▶ Non-stationary data: GRUs may struggle with nonstationary data, where the statistical properties of the data change over time.

RECURRENT NEURAL NETWORKS

APPLICATION OF RNNs

Applications of RNNs:

- ▶ Natural language processing (NLP): Using RNNs for text classification, language translation, and text generation.
- ▶ Time-series prediction: Using RNNs to make predictions based on sequential data, such as stock prices and weather patterns.
- ▶ Speech recognition: Using RNNs for speech-to-text conversion.

TRANSFORMER AND ATTENTION MECHANISM

Transformers and Attention Mechanisms are relatively recent developments in the field of deep learning, which have become popular for processing sequential data, such as natural language processing (NLP) tasks. Unlike Recurrent Neural Networks (RNNs) which process sequential data by repeatedly applying the same set of weights to the inputs over time, Transformers and Attention Mechanisms use self-attention mechanisms to dynamically weight the importance of different elements in the sequence. This enables Transformers to better capture the long-range dependencies between elements in the sequence, leading to improved performance on NLP tasks.

TRANSFORMER AND ATTENTION MECHANISM

SELF-ATTENTION MECHANISM

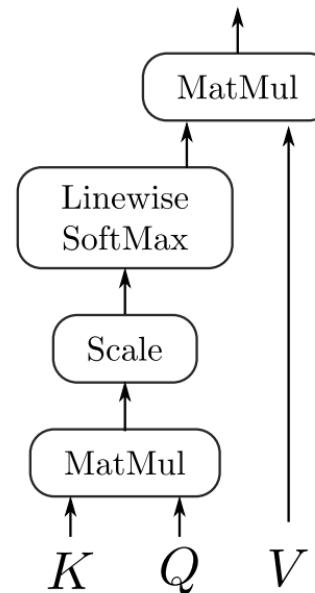
Self-attention mechanism in transformers is a method of calculating the weight of each input token in a sequence with respect to every other token in the same sequence, resulting in a representation of the input sequence in which the most relevant tokens have the highest weight. Mathematically, the self-attention mechanism can be represented as a dot product between the query (Q), key (K) and value (V) matrices, obtained from the input sequence, followed by a softmax activation to obtain the attention scores. These scores are then used to compute a weighted sum of the value matrix to produce the final representation.

$$\text{Attention}(Q, K, V) = \text{Softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad \text{where } Q \in \mathbb{R}^{m \times d_k}, K \in \mathbb{R}^{n \times d_k}, V \in \mathbb{R}^{n \times d_v}$$

TRANSFORMER AND ATTENTION MECHANISM

SELF-ATTENTION MECHANISM

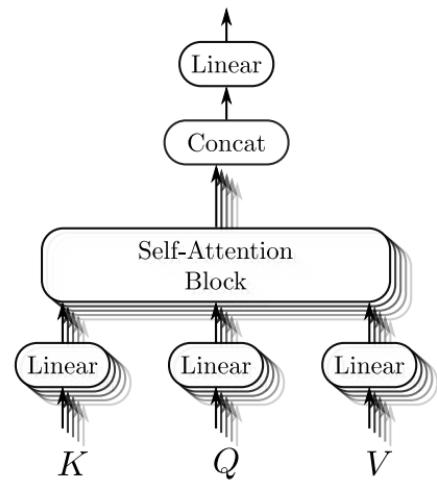
The intuition behind these matrices is to provide a way to measure the similarity between each query and key, which is then used to determine the contribution of the values to the final result. The resulting weighted sum of the values represents the output of the self-attention mechanism, capturing the relationships between different parts of the input sequence.



TRANSFORMER AND ATTENTION MECHANISM

MULTI-HEAD ATTENTION

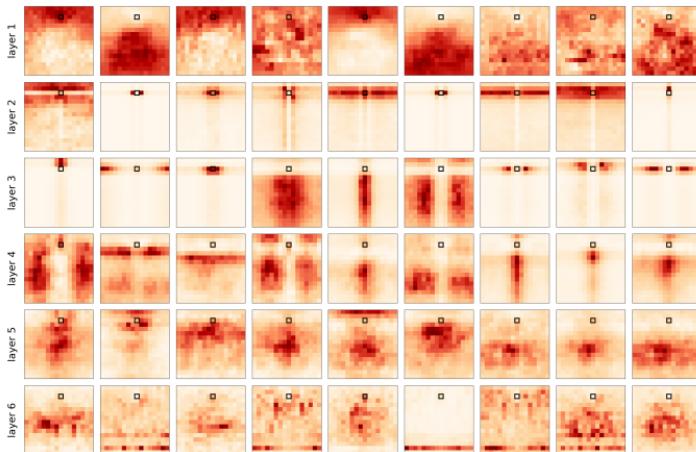
In Multi-head Attention, the self-attention mechanism is performed multiple times in parallel with different weight matrices, before being concatenated and once again projected, leading to a more robust representation of the input sequence. The intuition behind the three matrices (Q , K , V) remains the same as in self-attention, with Q representing the query, K the key and V the value. Each head performs an attention mechanism on the input sequence, capturing different aspects and dependencies of the data, before being combined to form a more comprehensive representation of the input.



TRANSFORMER AND ATTENTION MECHANISM

VISUALIZING MULTI-HEAD ATTENTION

Visualizing Self-Attention for
Image:
Link



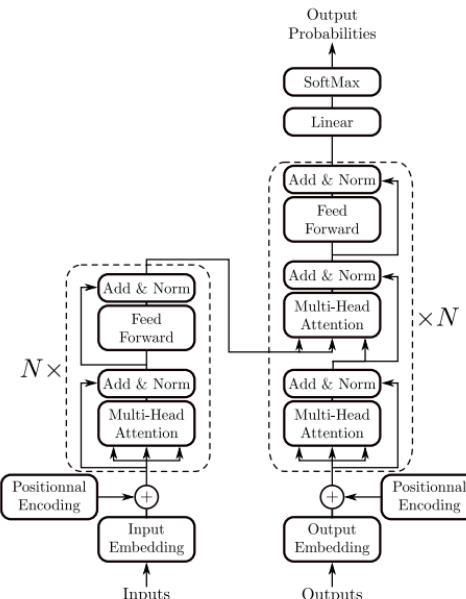
TRANSFORMER AND ATTENTION MECHANISM

TRANSFORMERS MODEL

Transformers are neural network models that use an encoder-decoder architecture. The encoder takes the input sequence and converts it into a continuous hidden representation, which is then passed to the decoder to generate the output sequence. The architecture of the transformer model is designed to allow the model to process the entire sequence in parallel, rather than processing one element at a time like in traditional RNNs.

Training of transformers involves optimizing a loss function that measures the difference between the model predictions and the true outputs. This loss function is usually based on the cross entropy between the predicted and true sequences.

The encoder-decoder mechanism is commonly referred to as the seq2seq mechanism.



Part III

TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

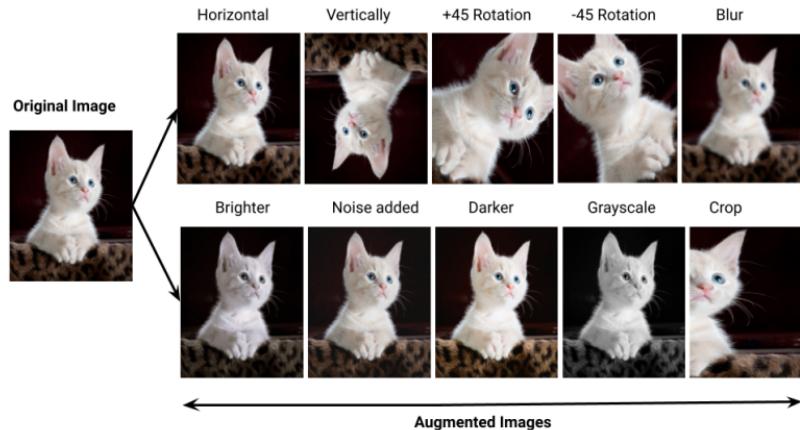
TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

- ▶ **Batch Normalization** - Seen
- ▶ **Dropout** - Seen
- ▶ **Data Augmentation**
- ▶ **Learning Rate Scheduling**
- ▶ **Early Stopping**
- ▶ **Gradient Clipping**
- ▶ **Weight Initialization**
- ▶ **Regularization**
- ▶ **GPU Acceleration**

TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

DATA AUGMENTATION

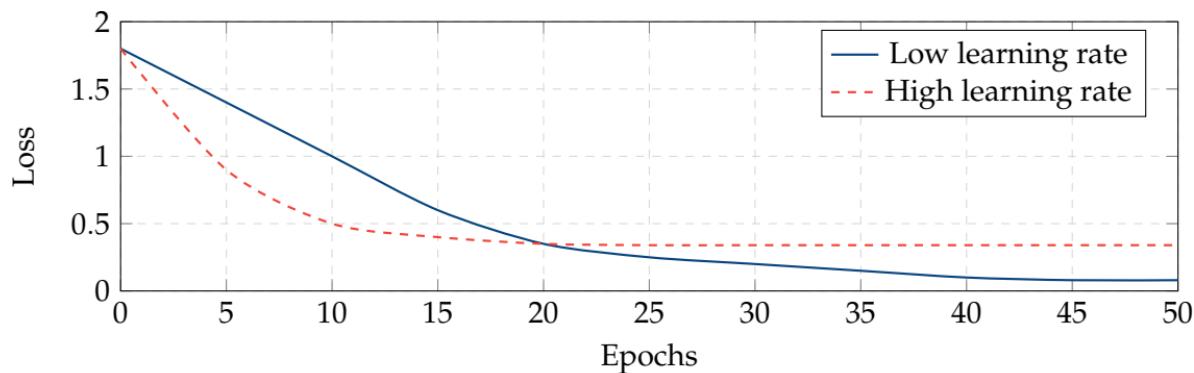
- ▶ Data Augmentation is a technique to increase the diversity of your training set by applying random (but realistic) transformations to the training images.
- ▶ The goal is to train a model that is robust to these transformations.
- ▶ For example, you can randomly rotate, scale, and flip the images in your training set.
- ▶ This helps expose the model to different aspects of the data and reduce overfitting.



TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

LEARNING RATE SCHEDULING

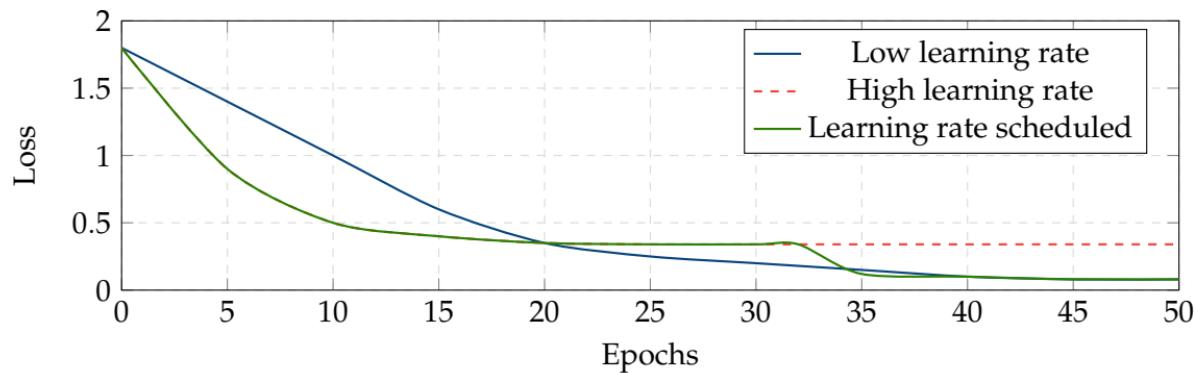
- ▶ The learning rate is one of the most important hyperparameters to tune for your deep learning model. The learning rate determines how quickly the model learns the optimal weights and how refined the gradient descent process is.
- ▶ If the learning rate is too high, the model may not converge or converge to a higher loss. If it is too low, the model may take too long to train.



TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

LEARNING RATE SCHEDULING

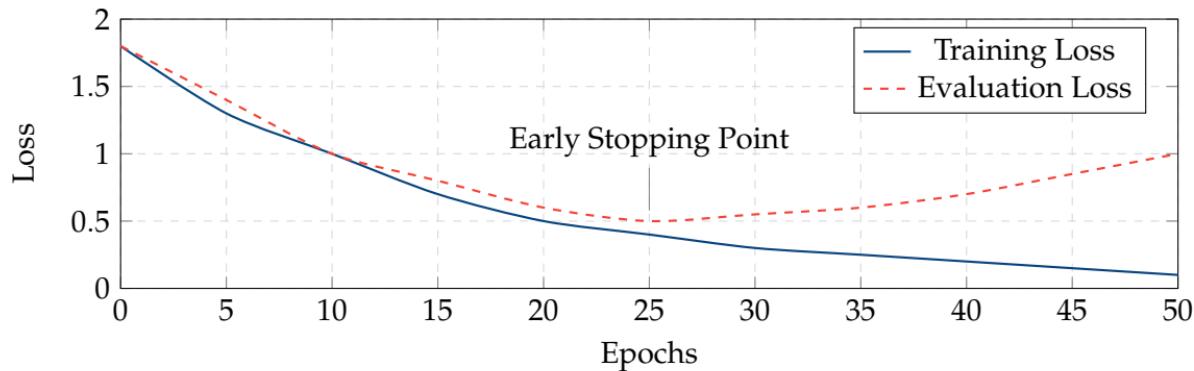
- ▶ The learning rate is one of the most important hyperparameters to tune for your deep learning model. The learning rate determines how quickly the model learns the optimal weights and how refined the gradient descent process is.
- ▶ If the learning rate is too high, the model may not converge or converge to a higher loss. If it is too low, the model may take too long to train.
- ▶ Learning rate scheduling is a technique to adjust the learning rate during training. For example, you can start with a high learning rate and then decrease it over time.



TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

EARLY STOPPING

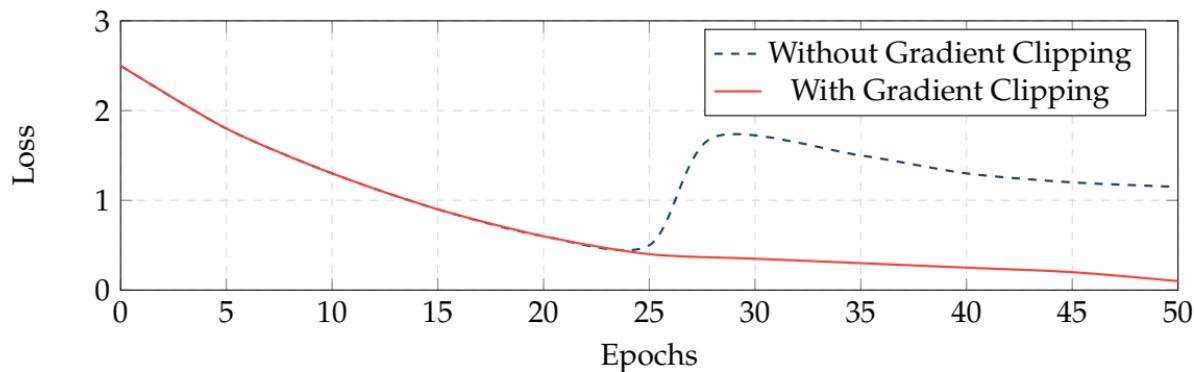
- ▶ Early stopping is a technique to prevent overfitting by stopping the training process when the model's performance on the validation set starts to degrade.
- ▶ The idea is to monitor the validation loss during training and stop training when the validation loss stops decreasing.



TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

GRADIENT CLIPPING

- ▶ Gradient clipping is a technique to prevent exploding gradients during training.
- ▶ Exploding gradients occur when the gradients of the loss function with respect to the model's parameters are too large.
- ▶ This can cause the model to diverge and fail to learn.
- ▶ Gradient clipping involves scaling the gradients if their norm exceeds a certain threshold.



TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

WEIGHT INITIALIZATION

- ▶ Weight initialization is a technique to set the initial values of the weights in the model.
- ▶ The initial values of the weights can have a significant impact on the training process and the final performance of the model.
- ▶ If the weights are initialized too small, the model may not learn effectively. If they are initialized too large, the model may not converge.
- ▶ Common weight initialization techniques include Xavier/Glorot initialization and He initialization.

TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

WEIGHT INITIALIZATION

- ▶ Xavier/Glorot initialization: The weights are initialized from a normal distribution with mean 0 and variance $2/(n_{\text{in}} + n_{\text{out}})$, where n_{in} and n_{out} are the number of input and output units, respectively. It helps prevent the gradients from vanishing or exploding during training by ensuring that the gradients have a similar scale.
- ▶ He initialization: The weights are initialized from a normal distribution with mean 0 and variance $2/n_{\text{in}}$, where n_{in} is the number of input units. It is commonly used for ReLU activation functions.

TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

REGULARIZATION

- ▶ Regularization is a technique to prevent overfitting by adding a penalty term to the loss function that discourages the model from learning complex patterns that may not generalize well.
- ▶ L1 regularization adds a penalty term to the loss function that is proportional to the absolute value of the weights. It encourages sparsity in the weights.
- ▶ L2 regularization adds a penalty term to the loss function that is proportional to the square of the weights. It encourages the weights to be small.

TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

GPU ACCELERATION

CPUs and GPUs are very different in terms of architecture and performance. CPUs are more suited for general-purpose computing tasks, while GPUs are optimized for parallel processing of simple operations, making them ideal for deep learning tasks.

- ▶ GPUs are much faster than CPUs for deep learning tasks because they have many more cores and can perform many more operations in parallel.
- ▶ Deep learning frameworks like PyTorch and TensorFlow are designed to take advantage of GPUs to accelerate the training process.
- ▶ Only the forward and backward passes of the model are executed on the GPU. The data loading and preprocessing are still done on the CPU.