### DEEP LEARNING 2 FROM THEORY TO PRACTICE

#### Alexandre Vérine, Research Fellow, École Normale Supérieure Paris

Double Licence Intelligence Artificielle et Sciences des Organisations 3e année de Licence Université Paris-Dauphine, PSL

November 26, 2024

### SEMESTER SCHEDULE (TEMPORARY)

- ▶ 17/09: Fundamentals of Deep Learning
- ▶ 24/09: TP1 Classification Introduction to PyTorch
- ▶ 01/10: In a Deep Learning Model + Techniques to Improve Deep Learning Training
- ▶ 08/10: TP2 Autoencoders *Hyperparameter Tuning*
- ► 15/10: Advanced Deep Learning Techniques
- ▶ 22/10: TP3 Image Segmentation From CPU to GPU and Parallelization
- ▶ 29/10: No Class

### SEMESTER SCHEDULE (TEMPORARY)

#### ► 05/11: Graded Individual Practical Work

- ▶ 12/11: TP4 Deep Reinforcement Learning From Notebook to Script
- ▶ 19/11: TP5 Adversarial Attacks Importance of Git
- 26/11: Project Presentation Group Formation
- ▶ 03/12: Group Session Help with Projects
- ► 10/12: **Project Presentation**

## AI 101: FROM FUNDAMENTALS TO DEEP LEARNING

| 1 | Introc | luction to Artificial Intelligence                               |
|---|--------|--|
|   | 1.1    | Deep Learning in the AI family                                   |
|   | 1.2    | Representation Learning  |
| 2 | Neura  | al Networks Fundamentals   |
|   | 2.1    | Neurons  |
|   | 2.2    | Layers   |
|   | 2.3    | Activation Functions   |
| 3 | The M  | Aulti-layer Perceptron (MLP)                                     |
|   | 3.1    | The first Deep Learning Model                                    |
|   | 3.2    | Stochastic Gradient Descent                                      |
|   | 3.3    | Back-propagation   |
|   | 3.4    | Example : Image classification of handwritten digits from A to Z |

## IN A DEEP LEARNING MODEL : FROM NEURAL NETWORKS TO TRANSFORMER MODELS

| 1 | Conve  | olutional Neural Networks                 |
|---|--------|---|
|   | 1.1    | The Two dimensional Convolution           |
|   | 1.2    | CNN : Convolutional in a network Networks |
|   | 1.3    | CNN in practice: CIFAR 10                 |
| 2 | Recur  | rent Neural Networks                      |
|   | 2.1    | Recurrent Block                           |
|   | 2.2    | LSTM and GRU                              |
| 3 | Trans  | former and Attention Mechanism            |
|   | 3.1    | Self-Attention Mechanism                  |
|   | 3.2    | Transformers Model 123                    |
| 4 | TP2: 1 | Build and use an autoencoder              |
|   | 4.1    | Formal introduction of an autoencoder     |

## TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

| 1 | Techn | iques to Improve Deep Learning Training | 2 |
|---|-------|---|---|
|   | 1.1   | Data Augmentation                       | 3 |
|   | 1.2   | Learning Rate Scheduling                | 4 |
|   | 1.3   | Early Stopping                          | 5 |
|   | 1.4   | Gradient Clipping                       | 6 |
|   | 1.5   | Weight Initialization                   | 7 |
|   | 1.6   | Regularization                          | 9 |
|   | 1.7   | GPU Acceleration                        | 0 |

## DEEP LEARNING AND APPLICATIONS

| 1 | Learn      | ing to act with Deep Reinforcement Learning               |
|---|------------|---|
|   | 1.1<br>1.2 | Deep Q-Learning 142   The Cheese Game 144                 |
| 2 | Synth      | etic Data Generation with Generative Adversarial Networks |
|   | 2.1<br>2.2 | GANS Models    149      MNIST Generation    151           |
| 3 | Sentin     | nent Analysis with Transformers and GRU                   |
|   | 3.1<br>3.2 | Bert 154   Sentiment Analysis 157                         |
| 4 | Densi      | ty Estimation with Normalizing Flows                      |
|   | 4.1<br>4.2 | Estimating Density 162   Normalizing Flows 163            |
| 5 | Image      | Segmentation with U-Net                                   |
|   | 5.1<br>5.2 | Image Segmentation 175   U-Net Architecture 175           |

## MANAGING DEEP LEARNING PROJECTS

| 1 | Why Notebooks are not enough |
|---|------------------------------|
| 2 | Modularity                   |
| 3 | Local Environment Setup      |
| 4 | Running Experiments          |
| 5 | Monitoring Experiments       |
| 6 | Versioning your Code         |

## PROJECT PRESENTATION

| 1 | Mini l | Project 2: Adversarial Attacks   | 213 |
|---|--------|----------------------------------|-----|
|   | 1.1    | Principle of Adversarial Attacks | 213 |
|   | 1.2    | Attacks                          | 224 |

| 2 | <b>Final</b> | Project: Generative Adversarial Networks            | . 228 |
|---|--------------|---|-------|
|   | 2.1          | Introduction to Generative Models                   | . 228 |
|   | 2.2          | Precision and Recall in Generative Models           | . 233 |
|   | 2.3          | Generative Adversarial Networks                     | . 239 |
|   | 2.4          | Final Project: Tuning Quality and Diversity in GANS | . 243 |

## Part I

## AI 101: FROM FUNDAMENTALS TO DEEP LEARNING

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

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) = 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



Alexandre Vérine

Figure. 2D classification for different AI models.

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,*  $\mathcal{Y} \subset \mathbb{R}^m$ ,  $f : \mathcal{X} \to \mathcal{Y}$  *be a continuous function and*  $\sigma : \mathbb{R} \to \mathbb{R}$  *be a continuous real function. Then*  $\sigma$  *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

Multi-layer Perceptron







Alexandre Vérine

DEEP LEARNING 2

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.



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.





EXAMPLE OF REPRESENTATION LEARNING



Figure. MNIST : Layer 0

Figure. MNIST : Layer 1



EXAMPLE OF REPRESENTATION LEARNING



**Figure.** MNIST : Layer 0



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.



ALE 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  $f : \mathbb{R}^d \to \mathbb{R}^m$ :



ALEAN NEWFOR is a processing unit that receives input, performs a computation, and produces an output. Here, the inputs  $^{/247}$  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:



# 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:

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

The non-linear function  $\sigma_i$  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  $\sigma_i$  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



# NEURAL NETWORKS FUNDAMENTALS SOFTMAX

Softmax Function:

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

Final activation

 Use case : Multi-class Classification



### NEURAL NETWORKS FUNDAMENTALS Hyperbolic Tangent



# NEURAL NETWORKS FUNDAMENTALS RELU



### NEURAL NETWORKS FUNDAMENTALS Leaky Relu

- Leaky Rectified Linear Unit (Leaky ReLU):
  - $\sigma(x) = \max\{\alpha x, x\}$
- Intermediate activation



# NEURAL NETWORKS FUNDAMENTALS ELU



# NEURAL NETWORKS FUNDAMENTALS Swish



## 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.



DEEP LEARNING 2

## 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 a dataset  $\mathcal{D}$ . The model  $f_{\theta}$  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.

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

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  $\lambda$ 

- 1: Initialize parameters  $\theta$
- 2: while  $\theta$  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  $\theta$ :  $\nabla_{\theta} l(x_i, f_{\theta}(x_i))$
- 6: Update parameters  $\theta = \theta \lambda \nabla_{\theta} l(x_i, f(x_i))$
- 7: end for
- 8: end while
- 9: return  $\theta$

# 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  $\lambda$  and a batch size *b* 

- 1: Initialize parameters  $\theta$
- 2: Initialize the number of batches  $B = \lfloor \frac{N}{b} \rfloor$
- 3: while  $\theta$  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  $\theta$ :  $\frac{1}{B} \sum_{i=1}^{B} \nabla_{\theta} l(x_i, f_{\theta}(x_i))$
- 7: Update parameters  $\theta = \theta \lambda_{\overline{B}}^{1} \sum_{i=1}^{B} \nabla_{\theta} l(x_{i}, f(x_{i}))$
- 8: end for
- 9: end while
- 10: return  $\theta$

At every step *t* of the gradient descent, setting a learning rate  $\lambda$ , the parameter  $\theta$  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.



DEEP LEARNING 2

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:

 $l(f(x)) = l(\sigma_3 (A_3 x_2 + b_3))$ 

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

Let us look at the gradients with respect to  $A_3$ :

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

and therefore:

$$A_3 \leftarrow A_3 - \lambda l'(x_3)\sigma'_3 (A_3x_2 + b_3) x_2^T.$$

We need to keep in memory the latent values of *x*, i.e. *x*<sub>2</sub>.

Let us look at the gradients with respect to  $A_2$ :

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

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

We have to compute the gradient with respect to  $x_2$ :

$$\frac{\partial l}{\partial x_2} = \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial x_2} = l'(x_3) \frac{\partial \sigma_3 \left(A_3 x_2 + b_3\right)}{\partial x_2} = l'(x_3) \frac{\partial \left[A_3 x_2 + b_3\right]}{\partial x_2} \sigma'_3 \left(A_3 x_2 + b_3\right)$$
$$= l'(x_3) A_3^{\mathrm{T}} \sigma'_3 \left(A_3 x_2 + b_3\right)$$

Therefore:

$$A_{2} \leftarrow A_{2} - \lambda \left[ l'(x_{3}) A_{3}^{T} \sigma_{3}' \left( A_{3} x_{2} + b_{3} \right) \times \sigma_{2}' \left( A_{2} x_{1} + b_{2} \right) x_{1}^{T} \right]$$

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

#### **BACK-PROPAGATION**

We have to compute the gradient with respect to  $A_1$ :

$$\begin{split} \frac{\partial l}{\partial A_1} &= \frac{\partial l}{\partial x_1} \frac{\partial x_1}{\partial A_1} \\ &= \frac{\partial l}{\partial x_1} \frac{\partial \sigma_1 \left(A_1 x_0 + b_1\right)}{\partial A_1} \\ &= \frac{\partial l}{\partial x_1} \sigma_1' \left(A_1 x_0 + b_0\right) x_0^T, \end{split}$$

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

### **BACK-PROPAGATION**

Let us compute the gradient with respect to  $x_1$ :

$$\frac{\partial l}{\partial x_1} = \frac{\partial l}{\partial x_2} \frac{\partial x_2}{\partial x_1} = \frac{\partial l}{\partial x_2} \frac{\partial \sigma_2 \left(A_2 x_1 + b_2\right)}{\partial x_1} = \frac{\partial l}{\partial x_2} \frac{\partial \left[A_2 x_1 + b_2\right]}{\partial x_1} \sigma_2' \left(A_2 x_1 + b_2\right)$$
$$= \frac{\partial l}{\partial x_2} A_2^T \sigma_2' \left(A_2 x_1 + b_2\right)$$

Therefore:

$$A_{1} \leftarrow A_{1} - \lambda \left[ l'(x_{3}) A_{3}^{T} \sigma_{3}' \left( A_{3} x_{2} + b_{3} \right) A_{2}^{T} \sigma_{2}' \left( A_{2} x_{1} + b_{2} \right) \times \sigma_{1}' \left( A_{1} x_{0} + b_{1} \right) x_{0}^{T} \right]$$

### **BACK-PROPAGATION**

In other words, the update on the weights is:

$$\begin{aligned} A_{3} &\leftarrow A_{3} - \lambda l'(x_{3})\sigma_{3}' \left(A_{3}x_{2} + b_{3}\right)x_{2}^{T} \\ A_{2} &\leftarrow A_{2} - \lambda \left[l'(x_{3})A_{3}^{T}\sigma_{3}' \left(A_{3}x_{2} + b_{3}\right) \times \sigma_{2}' \left(A_{2}x_{1} + b_{2}\right)x_{1}^{T}\right] \\ A_{1} &\leftarrow A_{1} - \lambda \left[l'(x_{3})A_{3}^{T}\sigma_{3}' \left(A_{3}x_{2} + b_{3}\right)A_{2}^{T}\sigma_{2}' \left(A_{2}x_{1} + b_{2}\right) \times \sigma_{1}' \left(A_{1}x_{0} + b_{1}\right)x_{0}^{T}\right] \end{aligned}$$

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  $b_3$ :

$$\begin{aligned} \frac{\partial l}{\partial b_3} &= \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial b_3} \\ &= l'(x_3) \frac{\partial \sigma_3 \left(A_3 x_2 + b_3\right)}{\partial b_3} \\ &= l'(x_3) \sigma'_3 \left(A_3 x_2 + b_3\right) \frac{\partial \left[A_3 x_2 + b_3\right]}{\partial b_3} \\ &= \underbrace{l'(x_3)}_{\in \mathbb{R}} \underbrace{\sigma'_3 \left(A_3 x_2 + b_3\right)}_{\in \mathbb{R}^{k_i \times 1}} \end{aligned}$$

And thus :

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

Let's move on the second layer:

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

And thus :

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

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

For the first layer:

$$\frac{\partial l}{\partial b_1} = \frac{\partial l}{\partial x_1} \frac{\partial x_1}{\partial b_1}$$
$$= \frac{\partial l}{\partial x_1} \frac{\partial \sigma_1 \left(A_1 x_0 + b_1\right)}{\partial b_1}$$
$$= \frac{\partial l}{\partial x_1} \sigma_1' \left(A_1 x_0 + b_0\right)$$

And thus :

$$b_1 \leftarrow b_1 - \lambda \frac{\partial l}{\partial x_1} \sigma' \left( A_1 x_0 + b_1 \right)$$

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.

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.





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





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





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



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



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

# THE MULTI-LAYER PERCEPTRON (MLP) 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.





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



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





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





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





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

 $\mathsf{Example}:\mathsf{Image}\ \mathsf{classification}\ \mathsf{of}\ \mathsf{handwritten}\ \mathsf{digits}\ \mathsf{from}\ \mathsf{A}\ \mathsf{to}\ \mathsf{Z}$ 

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.

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 :

- ŷ ∈ ℝ<sup>N×K</sup> is the predicted probability distribution over K classes for N samples,
  y ∈ 0, 1<sup>N×K</sup> is the ground-truth one-hot encoded label matrix,

#### RECAP ON THE CROSS-ENTROPY LOSS



Model Predicted Distribution  $(\hat{y}_i)$ 



The cross-entropy loss for one sample is:

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

 $\mathsf{Example}:\mathsf{Image}\xspace$  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



 $\mathsf{Example}:\mathsf{Image}\ \mathsf{classification}\ \mathsf{of}\ \mathsf{handwritten}\ \mathsf{digits}\ \mathsf{from}\ \mathsf{A}\ \mathsf{to}\ \mathsf{Z}$ 



Alexandre Vérine

 $\mathsf{Example}:\mathsf{Image}\ \mathsf{classification}\ \mathsf{of}\ \mathsf{handwritten}\ \mathsf{digits}\ \mathsf{from}\ \mathsf{A}\ \mathsf{to}\ \mathsf{Z}$ 

With a interpretation tool such as SHAP:



## 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),Resdiual 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.





# 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.




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 :

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

KERNEL SIZE, PADDING AND STRIDE

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



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

KERNEL SIZE, PADDING AND STRIDE

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



KERNEL SIZE, PADDING AND STRIDE

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



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

KERNEL SIZE, PADDING AND STRIDE

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



KERNEL SIZE, PADDING AND STRIDE

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



**CNN** : CONVOLUTIONAL IN A NETWORK NETWORKS

#### We can represent a CNN as under this form:



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

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





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 take the maximum within a given sized sub-matrix. In practice, the matrix is size  $2 \times 2$  in order to reduce the dimension by 4 and doubling the scale.



### 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.



### 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.

The Net in 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  $64 \mapsto 10 + SoftMax$





88 / 247



DEEP LEARNING 2

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

- $\blacktriangleright \text{ Conv } 3 \times 32 \times 32 \mapsto 6 \times 28 \times 28 + \text{ReLU}$
- Max Pooling  $6 \times 28 \times 28 \mapsto 6 \times 14 \times 14$
- $\blacktriangleright Conv \ 6 \times 14 \times 14 \mapsto 16 \times 10 \times 10 + ReLU$
- Max Pooling  $16 \times 10 \times 10 \mapsto 16 \times 5 \times 5$
- ► Linear  $400 \mapsto 120 + \text{ReLU}$
- ► Linear  $120 \mapsto 84 + \text{ReLU}$
- ► Linear  $84 \mapsto 10 + SoftMax$





91 / 247



DEEP LEARNING 2

The Net is composed 6 convolutional layers and 3 linear layers:  $\blacktriangleright$  Conv 3 × 32 × 32  $\mapsto$  32 × 32 × 32 + BatchNorm2d + ReLU  $\blacktriangleright$  Conv 32  $\times$  32  $\times$  32  $\mapsto$  64  $\times$  32  $\times$  32 + ReLU • Max Pooling  $64 \times 32 \times 32 \mapsto 64 \times 16 \times 16$  $\blacktriangleright$  Conv 64 × 16 × 16  $\mapsto$  128 × 16 × 16 + BatchNorm2d + ReLU  $\blacktriangleright$  Conv 128 × 16 × 16  $\mapsto$  128 × 16 × 16 + ReLU Max Pooling  $128 \times 16 \times 16 \mapsto 128 \times 8 \times 8$  $\blacktriangleright$  Conv 128  $\times$  8  $\times$  8  $\mapsto$  256  $\times$  8  $\times$  8 + BatchNorm2d + ReLU  $\blacktriangleright$  Conv 256  $\times$  8  $\times$  8  $\mapsto$  256  $\times$  8  $\times$  8 + ReLU • Max Pooling  $256 \times 8 \times 8 \mapsto 256 \times 4 \times 4 + \text{DropOut } p = 0.05$  $\blacktriangleright$  Linear 4096  $\mapsto$  1024 + ReLU Linear  $1024 \mapsto 512 + \text{ReLU} + \text{DropOut } p = 0.05$  $\blacktriangleright$  Linear 512  $\mapsto$  10 + SoftMax We have added Batch Normalization to improve the training stability and Drop Out to reduce

DEEP LEARNING 2

overfitting.

# 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 BATCH NORMALIZATION

Batch normalization is a technique in deep learning that 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  $\mu_B$  and variance  $\sigma_B^2$  of the features in the mini-batch *B*.

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

3: Normalize each feature  $x_i$  in the mini-batch *B* using  $\mu_B$  and  $\sigma_B^2$ .

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

4: Scale and shift each normalized feature  $x_i$  using two learnable parameters  $\gamma$  and  $\beta$  respectively.

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

5: end for

#### Algorithm 1: Batch Normalization



97 / 247



DEEP LEARNING 2

### CONVOLUTIONAL NEURAL NETWORKS CNN in practice: CIFAR 10



DEEP LEARNING 2

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.



#### **INTUITION BEHIND CHANNELS**





Original Image



Best Match



Best Match



Best Match



Channel 26

Channel 31



Original Image



Best Match



Best Match

Best Match

Best Match

Best Match



Best Match

















Best Match



Best Match



Original Image

Original Image

Best Match

Best Match















Best Match





#### **INTUITION BEHIND CHANNELS**





Original Image



Best Match



Best Match



Best Match



Best Match



Original Image

Original Image

Original Image



Best Match



Best Match

Best Match

Best Match

DEEP LEARNING 2



Best Match





Best Match





#### Best Match





Alexandre Vérine

Channel 20

Channel 16

Best Match

Best Match











#### **INTUITION BEHIND CHANNELS**









Best Match



Best Match







Original Image

Original Image

Original Image

Channel 13

Channel 24



Best Match

Best Match

Best Match



Best Match

Best Match

Best Match



Best Match



Best Match











Best Match



Best Match



Best Match



Best Match DEEP LEARNING 2







#### **INTUITION BEHIND CHANNELS**





Original Image



Best Match



Best Match



Best Match



Original Image



Best Match



Best Match

Best Match



Best Match





Original Image



Original Image



Best Match



Best Match



Best Match







Best Match



Best Match



Best Match



Best Match







### **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.



## RECURRENT NEURAL NETWORKS

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.


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.



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.



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.



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.

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.



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.



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.

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

### TRANSFORMER AND ATTENTION MECHANISM

- Query (*Q*): Represents the query vector, which is used to calculate the attention scores. Intuitively, the query vector represents the token that we are interested in.
- Key (*K*): Represents the key vector, which is used to calculate the attention scores. The key vector helps to determine the importance of each token in the input sequence.
- Value (V): Represents the value vector, which is used to compute the weighted sum of the values. The value vector provides the information that is used to update the representation of the input sequence.

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 Qrepresenting 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.



DEEP LEARNING 2

#### TRANSFORMER AND ATTENTION MECHANISM TRANSFORMERS MODEL

More information about transformers and specific model architectures will be covered next semester in the course on Applied Deep Learning.

### TP2: BUILD AND USE AN AUTOENCODER

FORMAL INTRODUCTION OF AN AUTOENCODER

#### Definition

An autoencoder is a type of artificial neural network used to learn efficient codings of unlabeled data. It consists of two main components:

- An encoder function:  $encoder(x) : \mathbb{R}^d \to \mathbb{R}^m$ Maps an input *x* from the input space  $\mathbb{R}^d$  to a hidden representation space  $\mathbb{R}^m$ .
- A decoder function: *decoder*(*z*) : ℝ<sup>m</sup> → ℝ<sup>d</sup>
  Maps the hidden representation *z* back to the original input space ℝ<sup>d</sup>.

#### Goal

The primary goal of an autoencoder is to learn a representation (encoding) for a set of data, typically for the purpose of dimensionality reduction or feature learning. Through training, the autoencoder learns to compress the data from  $\mathbb{R}^d$  to  $\mathbb{R}^m$  (where m < d) and then reconstruct the data back to  $\mathbb{R}^d$  as accurately as possible. This process forces the autoencoder to capture the most important features of the data in the hidden representation *z*.

### TP2: BUILD AND USE AN AUTOENCODER

FORMAL INTRODUCTION OF AN AUTOENCODER



### TP2: BUILD AND USE AN AUTOENCODER

AUTOENCODERS AND UNSUPERVISED LEARNING

#### **Unsupervised Learning**

Autoencoders are a classic example of unsupervised learning. In unsupervised learning, the goal is to learn patterns from unlabelled data. Autoencoders learn to compress and decompress the input data without any explicit labels, aiming to capture the underlying structure of the data.

#### **Objective Function**

The learning process of an autoencoder is guided by the minimization of a loss function, typically involving a norm that measures the difference between the input and the reconstructed output. Formally, the objective is to minimize:

 $\min_{\theta} \mathbb{E}_{x \sim P} \left[ l(x, \operatorname{decoder}_{\theta}(\operatorname{encoder}_{\theta}(x))) \right]$ 

where *x* is the input data,  $\theta$  represents the parameters of the encoder and decoder, and *l* is a loss function.

### TP2: Build and use an autoencoder

AUTOENCODERS AND UNSUPERVISED LEARNING

#### **Unsupervised Learning**

Autoencoders are a classic example of unsupervised learning. In unsupervised learning, the goal is to learn patterns from unlabelled data. Autoencoders learn to compress and decompress the input data without any explicit labels, aiming to capture the underlying structure of the data.

#### **Objective Function**

The learning process of an autoencoder is guided by the minimization of a loss function, typically involving a norm that measures the difference between the input and the reconstructed output. Formally, the objective is to minimize:

$$\min_{\theta} \mathbb{E}_{x \sim P} \left[ \|x - \hat{x}\|_2^2 \right]$$

where *x* is the input data,  $\theta$  represents the parameters of the encoder and decoder and  $\hat{x} = \text{decoder}_{\theta}(\text{encoder}_{\theta}(x))$  in the reconstruction.

Reduce the size of the data to transfer: Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.

- Reduce the size of the data to transfer: Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.
- Denoise image: By learning to ignore the "noise" in the input data, autoencoders can reconstruct cleaner versions of noisy images, effectively removing the noise.



- Reduce the size of the data to transfer: Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.
- Denoise image: By learning to ignore the "noise" in the input data, autoencoders can reconstruct cleaner versions of noisy images, effectively removing the noise.
- Anomaly detection: Autoencoders can learn the normal patterns within data. Deviations from these patterns, when the reconstruction error is high, can indicate anomalies or outliers in the data.



- Reduce the size of the data to transfer: Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.
- Denoise image: By learning to ignore the "noise" in the input data, autoencoders can reconstruct cleaner versions of noisy images, effectively removing the noise.
- Anomaly detection: Autoencoders can learn the normal patterns within data. Deviations from these patterns, when the reconstruction error is high, can indicate anomalies or outliers in the data.



# TP2: BUILD AND USE AN AUTOENCODER YOUR TURN !

Get the TP2 on the course website and start working on it.

- https://www.alexverine.com
- ► Teaching
- ► Deep Learning II
- Lien Notebooks Python

### Part III

### 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

#### 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.



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

- 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  $2/(n_{in} + n_{out})$ , 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  $2/n_{in}$ , 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.

### Part IV

### DEEP LEARNING AND APPLICATIONS

# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING DEEP Q-LEARNING

Deep Q-Learning is a reinforcement learning algorithm that utilizes a neural network to approximate the optimal Q function, which is defined as the expected cumulative reward obtained by following a specific policy. The expected reward can be represented mathematically as follows:

$$R(\pi) = \sum_{t \leq T} E_{p^{\pi}}[\gamma^t r(s_t, a_t)],$$

where  $r(s_t, a_t)$  is the reward at time step  $t, \gamma$  is the discount factor, T is the final time step and

$$p^{\pi}(a_0, a_1, s_1, ..., a_T, s_T) = p(a_0) \prod_{t=1}^T \pi(a_t | s_t) p(s_{t+1} | s_t, a_t)$$
.

The Q function, Q(s, a), represents the expected cumulative reward obtained by taking action *a* in state *s*:

$$Q^{\pi}(s,a) = E_{p^{\pi}}[\sum_{t \le T} \gamma^{t} r(s_{t},a_{t}) | s_{0} = s, a_{0} = a].$$

The policy, represented by  $\pi(a|s)$ , is a probability distribution over actions given a state. The optimal Q function,  $Q^*(s, a)$ , can be found by solving the Bellman equation:

$$Q^*(s,a) = \mathbb{E}[R|s,a] = \mathbb{E}[r + \gamma \max_{a'} Q^{\pi}(s',a')|s,a].$$

DEEP LEARNING 2

### LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING DEEP Q-LEARNING

The loss in Deep Q-Learning method is the difference between the predicted Q-value and the target Q-value, which is the maximum expected reward obtained from the next state:

$$\mathcal{L}(\theta) = E_{s' \sim \pi^*(.|s,a)} \| r + \gamma \max_{a'} Q(s',a',\theta) - Q(s,a,\theta) \|^2.$$

This loss is used to update the parameters of the deep learning model in order to make the predictions more accurate. The target Q-value is typically computed as the reward obtained from taking an action in the current state, plus the maximum expected reward obtainable from the next state:

- 1. At the state  $s_t$ , select the action  $a_t$  with best reward using  $Q_t$  and store the results;
- 2. Obtain the new state  $s_{t+1}$  from the environment p;
- 3. Store  $(s_t, a_t, s_{t+1})$ ;
- 4. Obtain  $Q_{t+1}$  by minimizing  $\mathcal{L}$  from a batch recovered from previously stored results.

# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING THE MOUSE GAME

- A Mouse has to feed on food (red) and avoid poison (blue).
- It has a vision range of 2 squares. So it can see the 25 cells around.
- ► The reward for a cheese cell is 0.5, while the reward for eating poison is −1.

On this example, the mouse behaves randomly.

## LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING MODEL 1

#### Fully connected network:



## LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING MODEL 2

#### Convolutional network:



# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING $\epsilon$ -greedy Algorithm

The  $\epsilon$ -greedy algorithm is a common exploration strategy used in Deep Q learning. Balances exploration, where the agent tries out new actions and collects new data, and exploitation, where the agent uses the information it already has to select the action with the highest expected reward. The algorithm selects a random action with probability  $\epsilon$  and the action with the highest Q value with probability  $1 - \epsilon$ . The value of  $\epsilon$  decreases over time to gradually shift the focus from exploration.

### LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING

MODEL 2+ INCORPORATED  $\epsilon$ -EXPLORATION

### Convolutional network + $\epsilon$ -greedy Algorithm during training:

| Inr        | utl ovor | i      | input:   |                 | None, 5, 5, 2)]  |  |  |  |  |
|------------|----------|--------|--|-----------------|------------------|--|--|--|--|
| ш          | Juilayei | 0      | output:  |                 | None, 5, 5, 2)]  |  |  |  |  |
|            |          |        |  |                 |                  |  |  |  |  |
| _          |          |        | <u> </u>   |                 |                  |  |  |  |  |
| Conv2D     | in       | input: |  | (None, 5, 5, 2) |                  |  |  |  |  |
| Ľ          | 5011122  | ou     | output:  |                 | (None, 4, 4, 8)  |  |  |  |  |
|            |          |        | The second secon |                 |                  |  |  |  |  |
| Activation |          | i      | input:   |                 | (None, 4, 4, 8)  |  |  |  |  |
| 1          | cuvation | 0      | output:  |                 | None, 4, 4, 8)   |  |  |  |  |
|            |          |        |  |                 |                  |  |  |  |  |
|            |          |        | <b>*</b>   |                 |                  |  |  |  |  |
| Comu2D     | inj      | out:   | (None, 4, 4, 8)  |                 |                  |  |  |  |  |
|            | COIIV2D  |        | output:  |                 | (None, 3, 3, 16) |  |  |  |  |
|            |          |        |  |                 |                  |  |  |  |  |
| <b>_</b>   | latton   | inț    | input:   |                 | (None, 3, 3, 16) |  |  |  |  |
| ľ          | Flatten  |        | put:   | (None, 144)     |                  |  |  |  |  |
|            |          |        |  |                 |                  |  |  |  |  |
|            | Donco    | iı     | iput:  | (None, 144)     |                  |  |  |  |  |
|            | Dense    |        | output:  |                 | (None, 4)        |  |  |  |  |
|            |          |        | L.   |                 |                  |  |  |  |  |
| [          | Activati |        | inpu   | t:              | (None, 4)        |  |  |  |  |
|            | Activati |        | outpu  |                 | (None, 4)        |  |  |  |  |
|            |          |        |  |                 |                  |  |  |  |  |

# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS GANS MODELS

Generative Adversarial Networks (GANs) are a type of deep learning architecture composed of two neural networks, the generator and the discriminator, that are trained in a adversarial manner. The generator network is trained to generate fake data that appears similar to real data, while the discriminator network is trained to distinguish between real and fake data. The loss for GANs is defined as a min-max game, where the generator minimizes the loss function, and the discriminator maximizes it. *D* and *G* represent the discriminator and generator networks, respectively, and the goal is to find the optimal configuration for *D* and *G* such that the generated samples appear indistinguishable from real data.

# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS GANS MODELS

The loss of a (GAN) is defined as a minimax game between the generator and discriminator models. The generator aims to generate samples that are indistinguishable from real samples, while the discriminator aims to distinguish the generated samples from real samples. The loss function for the generator is defined as  $-\log(D(G(z)))$ , where *D* is the discriminator model and G(z) is the generator's output for a random noise vector *z*. The loss function for the discriminator is defined as  $\log(D(x)) + \log(1 - D(G(z)))$ , where *x* is a real sample.

$$\min_{D} \max_{G} \mathbb{E}_{x \sim p_{\text{real}}} \left[ \log(D(x)) \right] + \mathbb{E}_{z \sim q} \left[ \log(1 - D(G(z))) \right]$$

# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS MNIST GENERATION





Alexandre Vérine

DEEP LEARNING 2

# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS MNIST GENERATION



# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS MNIST GENERATION

# SENTIMENT ANALYSIS WITH TRANSFORMERS AND $\ensuremath{\mathsf{GRU}}$ bert

BERT (Bidirectional Encoder Representations from Transformers) is a pre-trained language model developed by Google in 2018. It is a transformer-based architecture that uses a masked language modeling task to generate a deep understanding of the contextual relationships between words in a sentence. BERT can be fine-tuned for various NLP tasks such as sentiment classification by adding a classification layer on top of its pre-trained representations. The model has achieved state-of-the-art performance in a wide range of NLP tasks, making it a popular choice for sentiment analysis.

A bidirectional transformer is a type of transformer architecture in natural language processing (NLP) where information from both past and future contexts is taken into consideration when making predictions. This is achieved by processing the input sequence in two directions, starting from the beginning and the end of the sequence, and concatenating the outputs to obtain the final representation. This allows the model to capture the context both in the forward and backward directions, providing a more comprehensive representation of the input sequence.

# SENTIMENT ANALYSIS WITH TRANSFORMERS AND $\ensuremath{\mathsf{GRU}}$ bert

BERT model consists of multiple transformer encoder blocks, with a self-attention mechanism, a feedforward neural network and layer normalization, stacked on top of each other. It also includes a positional encoding component to capture the relative position of tokens in a sequence, and a segment encoding component to differentiate between different sequences within the same input.



Bert

# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU training bert

Bert is trained on two tasks, the masked language model and the next sentence prediction. In the masked language model task, a portion of the input sequence is masked and the model must predict the original token based on its context. In the next sentence prediction task, the model receives a pair of sentences and must predict whether the second sentence follows the first one in the context of the input text. Both of these tasks are used to train Bert to understand the context of words in a sentence and how they relate to each other, allowing it to perform well on a wide range of natural language processing tasks.

| There  | is    | no    | curse  | in | Elvish, | Entish, | or | the | tongues | of | Men | for | this | treachery. |
|--------|-------|-------|--------|----|---------|---------|----|-----|---------|----|-----|-----|------|------------|
| Basic- | leve  | el m  | naskin | g  |         |         |    |     |         |    |     |     |      |            |
| There  | is    |       | curse  | in | Elvish, |         | or | the | tongues | of | Men | for | this |            |
| Entity | v-lev | vel 1 | maski  | ng |         |         |    |     |         |    |     |     |      |            |
| There  | is    |       |        | in | Elvish, | Entish, | or | the |         |    |     | for | this | treachery. |
| Phras  | e-le  | vel   | maski  | ng |         |         |    |     |         |    |     |     |      |            |
|        |       |       | curse  | in | Elvish, | Entish, | or |     |         |    |     | for | this | treachery. |

Sentence

# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU SENTIMENT ANALYSIS

Bert can be fine-tuned for sentiment analysis by adding a classifier layer on top of the pretrained Bert model. The layer is trained on a labeled sentiment analysis dataset to predict the sentiment of a given input sequence, which can be a sentence, paragraph, or document. Fine-tuning the model allows it to learn the specific nuances of sentiment in the target task and produce improved results compared to using the pre-trained model alone.





### Bert

(pre-trained)





| In | [117]: | predict | <pre>sentiment(model,</pre> | tokenizer, | "This | film | is | terrible") | ) |
|----|--------|---------|-----------------------------|------------|-------|------|----|------------|---|
|----|--------|---------|-----------------------------|------------|-------|------|----|------------|---|

Out[117]: 0.028073227033019066

- In [118]: predict\_sentiment(model, tokenizer, "This film is great")
- Out[118]: 0.9749133586883545
- In [119]: predict sentiment(model, tokenizer, "This lecture was quite boring")
- Out[119]: 0.0537508986890316
- In [120]: predict\_sentiment(model, tokenizer, "This lecture was challenging")
- Out[120]: 0.6184227466583252
- In [121]: predict sentiment(model, tokenizer, "This lecture was dense but interesting")

Out[121]: 0.7105910181999207

### DENSITY ESTIMATION WITH NORMALIZING FLOWS Estimating Density

In Machine Learning, the task of *density estimation* consists in estimating the probability density function of a random variable from a set of observations. This is a fundamental problem in statistics and machine learning, with applications in a wide range of fields, including anomaly detection, clustering, and generative modeling. There are several methods for density estimation, including parametric models, non-parametric models, and deep learning models. In this section, we will focus on a specific type of deep learning model called Normalizing Flows, which is used for density estimation and generative modeling.



Alexandre Vérine

### DENSITY ESTIMATION WITH NORMALIZING FLOWS Overview

A Normalizing Flow is usually seen as:

- ▶ a generative model,
- ▶ a bijective mapping,
- ▶ an *invertible neural network*,
- ► a density estimator.

#### DENSITY ESTIMATION WITH NORMALIZING FLOWS OVERVIEW



Point to point

#### DENSITY ESTIMATION WITH NORMALIZING FLOWS OVERVIEW



### DENSITY ESTIMATION WITH NORMALIZING FLOWS

MATHEMATICAL FRAMEWORK

#### **Normalizing Flow**

A Normalizing Flow is a bijective function between a data space  $\mathcal{X}$  and a latent space  $\mathcal{Z}$ , both subset of  $\mathbb{R}^d$ .

 $\begin{array}{cccc} F: & \mathcal{X} & \longmapsto & \mathcal{Z} \\ & x & \longmapsto & z = F(x) \end{array}$ 

#### **Data and Latent Distributions**

In theory, a NF maps a target distribution *P*, ie the data distribution to a simple latent distribution *Q*. Usually, *Q* is set to be a Normal Gaussian multivariate distribution  $\mathcal{N}(0_d, I_d)$ . *p* and *q* are respectively the probability densities of *P* and *Q*.

# DENSITY ESTIMATION WITH NORMALIZING FLOWS HOW DOES IT WORK ?

In practice, the mapping is *not perfect*.  $P^*$  induces a distribution Q and similarly, the latent distribution Q induces  $\hat{P}$ , which is the learned distribution. The forward pass F is called the *Normalizing* direction while the inverse pass  $F^{-1}$  is called the *Generative* direction.



Figure. 1D Normalizing Flow process.

### DENSITY ESTIMATION WITH NORMALIZING FLOWS INDUCED PROBABILITIES ?

#### **Change of Variable Formula**

For a bijective and continuous fonction *F* and a latent distribution *Q*, the distribution induced by *Q* and *F* is defined through the *change of variable formula*:

$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \operatorname{Jac}_F(x)| \, q(F(x)). \tag{1}$$

## DENSITY ESTIMATION WITH NORMALIZING FLOWS INDUCED PROBABILITIES ?

 $\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \operatorname{Jac}_F(x)| q(F(x)).$ PX FQ  $\mathcal{Z}$ 

### DENSITY ESTIMATION WITH NORMALIZING FLOWS

INDUCED PROBABILITIES ?

 $\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \operatorname{Jac}_F(x)| q(F(x)).$ 



#### DENSITY ESTIMATION WITH NORMALIZING FLOWS INDUCED PROBABILITIES ?

 $\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \operatorname{Jac}_F(x)| q(F(x)).$ 



### DENSITY ESTIMATION WITH NORMALIZING FLOWS

INDUCED PROBABILITIES ?

 $\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \operatorname{Jac}_F(x)| q(F(x)).$ 



Alexandre Vérine

## DENSITY ESTIMATION WITH NORMALIZING FLOWS INDUCED PROBABILITIES ?

-----

 $\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \operatorname{Jac}_F(x)| q(F(x)).$ 



# DENSITY ESTIMATION WITH NORMALIZING FLOWS DENSITY ESTIMATION

To perfom density estimation:

- 1. Draw  $x \sim P^*$ ,
- 2. Compute F(x) and  $|\det \operatorname{Jac}_F(x)|$ ,
- 3. Compute  $\widehat{p}(x) = q(F(x)) |\det \operatorname{Jac}_F(x)|$ .



Figure. 1D Normalizing Process of Density Estimation.

#### DENSITY ESTIMATION WITH NORMALIZING FLOWS Data Generation

To perform data generation:

- 1. Draw  $z \sim Q$ ,
- 2. Compute  $x = F^{-1}(x)$ .



Figure. 1D Normalizing Flow process of Generation.
# DENSITY ESTIMATION WITH NORMALIZING FLOWS LOG-LIKELIHOOD

### Loss

The objective is to approximate  $P^*$  with  $\hat{P}$ . We can minimize the Kullback-Leiber Divergence :

$$\theta = \operatorname*{arg\,min}_{\theta} \mathcal{D}_{\mathrm{KL}}(P^* \| \widehat{P}).$$

This is equivalent to maximizing the log likelihood :

$$\theta = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{X}} \left[ \log \widehat{p}(x) \right].$$

# DENSITY ESTIMATION WITH NORMALIZING FLOWS LOG-LIKELIHOOD

$$\mathcal{D}_{\mathrm{KL}}(P^* \| \widehat{P}) = \int_{\mathcal{X}} p^*(x) \log\left(\frac{p^*(x)}{\widehat{p}(x)}\right) dx$$

$$\mathbf{nll} = -\mathbb{E}_{x \sim \mathcal{X}} \left[ \log \widehat{p}(x) \right]$$

# DENSITY ESTIMATION WITH NORMALIZING FLOWS LEARNING STEPS



**Figure.** Learning Process for a 1D Normalizing Flow.

# IMAGE SEGMENTATION WITH U-NET

Image segmentation is the process of partitioning an image into multiple segments or regions based on the characteristics of the pixels. It is a fundamental task in computer vision and has applications in various fields, including medical imaging, autonomous driving, and satellite image analysis. There are several methods for image segmentation, including thresholding, clustering, and deep learning-based approaches. In this section, we will focus on a deep learning model called U-Net, which is commonly used for image segmentation tasks.



DEEP LEARNING 2

## Image Segmentation with U-Net

APPLICATIONS OF IMAGE SEGMENTATION

Image segmentation has a wide range of applications in computer vision and image processing, including:

- Medical Imaging: Segmentation of organs, tumors, and other structures in medical images.
- Autonomous Driving: Segmentation of objects such as cars, pedestrians, and road signs in images captured by autonomous vehicles.
- Satellite Image Analysis: Segmentation of land cover types, buildings, and other features in satellite images.
- Object Detection: Segmentation of objects in images to localize and classify them.
- Image Editing: Segmentation of objects for image editing tasks such as background removal and image compositing.

# IMAGE SEGMENTATION WITH U-NET

Image segmentation can be broadly classified into two types: semantic segmentation and instance segmentation.

- Semantic Segmentation: Semantic segmentation assigns a class label to each pixel in an image, such as road, car, person, etc. The goal is to partition the image into semantically meaningful regions.
- Instance Segmentation: Instance segmentation goes a step further than semantic segmentation by distinguishing between different instances of the same class. It assigns a unique label to each object instance in the image.



**Original Image** 



Semantic Segmentation



Instance Segmentation

## IMAGE SEGMENTATION WITH U-NET U-NET ARCHITECTURE

U-Net is a convolutional neural network architecture designed for image segmentation tasks. It consists of an encoder-decoder structure with skip connections that allow the model to capture both local and global features in the input image.



3 main loss functions are used for image segmentation tasks:

- Categorical Cross-Entropy Loss: Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.
- Dice Loss: A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.
- Shaped aware Loss: A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

3 main loss functions are used for image segmentation tasks:

Categorical Cross-Entropy Loss: Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} y_{i,c} \log(\hat{y}_{i,c})$$

- Dice Loss: A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.
- Shaped aware Loss: A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

3 main loss functions are used for image segmentation tasks:

- Categorical Cross-Entropy Loss: Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.
- Dice Loss: A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.

$$\mathcal{L} = 1 - \frac{2\sum_{i=1}^{N} y_i \hat{y}_i}{\sum_{i=1}^{N} y_i + \sum_{i=1}^{N} \hat{y}_i}$$

Shaped aware Loss: A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

3 main loss functions are used for image segmentation tasks:

- Categorical Cross-Entropy Loss: Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.
- Dice Loss: A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.
- Shaped aware Loss: A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

$$\mathcal{L} = -w(x) \left[ \frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} y_{i,c} \log(\hat{y}_{i,c}) \right]$$

with w(x) a weight function that assigns higher weights to pixels near the object boundaries.

# Part V

# MANAGING A DEEP LEARNING PROJECT

Jupyter or Colab Notebooks are great for prototyping, showcasing, and data exploration. However, they are not enough for managing a deep learning project. In particular, they lack the following features:

- Reproducibility: Notebooks are not reproducible. They are not designed to be run multiple times with the same results.
- Versioning: Notebooks are not versioned. It is hard to track changes and revert to versions.
- **Modularity:** Notebooks are not modular. It is hard to reuse code across different notebooks.
- **Scalability:** Notebooks are not scalable. They are not designed to run on multiple machines or in the cloud.
- Monitoring: Notebooks are not monitored. It is hard to track the progress of a long-running experiment.
- Collaboration: Notebooks are not collaborative. It is hard to work with multiple people on the same notebook.
- **Deployment:** Notebooks are not deployable. They are not designed to be run in production.
- Performance: Notebooks are not performant, depending on the amount of log, they can be up to 50% slower than a script.

## WHY NOTEBOOKS ARE NOT ENOUGH Why Scripts are Better

Directly running Python script is a better alternative to Jupyter Notebooks for a deep learning project. They allow:

- **Support for Scheduling and Automation:** Scripts can be scheduled to run at specific times or intervals.
- **Support for Versioning:** Scripts can be versioned using Git or other version control systems.
- **Support for Scalability:** Scripts can be run on multiple machines or in the cloud.
- **Support for Command Line Arguments:** Scripts can take command line arguments for customization.
- Support for Ressource Management: Scripts can be run in the background without blocking the terminal or integrate cluster management tools (e.g. SLURM).
- **Support for Logging:** Scripts can log information to files for monitoring and debugging.
- Support for Modularity: Scripts can be modularized into functions and classes.

# WHY NOTEBOOKS ARE NOT ENOUGH PROGRAM

In this lecture, we will learn how to manage a deep learning project using Python scripts. We will cover the following topics:

- Project Structure: We will learn how to organize a deep learning project into directories and files for better modularity and reusability.
- **Dependency Management:** We will learn how to manage dependencies using a virtual environment and a requirements file.
- **Running code as a script:** We will learn how to run Python code as a script from the command line.
- **Logging:** We will learn how to log information to files for monitoring and debugging.

## MODULARITY Why Code Should Be Modular

Writing modular code is essential for managing complex deep learning projects. It offers several advantages:

- Reusability: Modular code allows functions and classes to be reused across different parts of the project, reducing redundancy.
- **Maintainability:** With a clear structure, modular code is easier to maintain and update, especially in large projects.
- **Testing and Debugging:** Smaller, independent modules make it easier to test and debug specific parts of the codebase.
- Collaboration: A modular design enables multiple team members to work on different parts of the codebase simultaneously without conflicts.
- **Scalability:** Modular code can scale more effectively, as each component can be optimized or extended individually.

## MODULARITY Organizing Your Project



#### **Explanation:**

- main.py: Entry point for running experiments
- models.py, agents.py, etc.: Modular files for specific tasks
- utils.py: Utility functions and classes
- requirements.txt: List of dependencies
- experiments/: Separate folder for experiment outputs, logs, and saved models

## MODULARITY Key Principles of Modularity

To build a maintainable and scalable project, follow these modularity principles:

- Separation of Concerns: Each module should have a single responsibility (e.g., separate data processing, model architecture, and training code).
- **Encapsulation:** Keep functionality within defined boundaries to avoid cross-module dependencies.
- Configurable Parameters: Use configuration files for parameters, ensuring that code can be adapted to new experiments without modification.
- Interface Design: Define clear inputs and outputs for each module, making it easier to swap or update components.

### MODULARITY WRITING MODULAR CODE IN PYTHON

Best practices for maintaining modularity in Python:

- **Use Functions and Classes:** Encapsulate related functionality in functions and classes for better organization.
- **Separate Configurations:** Store experiment parameters and settings in YAML or JSON files instead of hardcoding.
- Follow a Naming Convention: Consistent, descriptive names for modules, functions, and variables improve readability.
- **Documentation:** Document modules with docstrings and comments to clarify each part's role and dependencies.

### MODULARITY Defining the NeuralNet Class

### Code example:

```
import torch.nn as nn

class NeuralNet(nn.Module):
    def __init__(self):
        super(NeuralNet, self).__init__()
        self.block = None

    def forward(self, x):
        x = self.block(x)
        return x
```

## MODULARITY Extending NeuralNet for Different Architectures

Modularity allows for easy extension of base classes. Here, we extend NeuralNet to create specialized architectures:

```
class DenseNet(NeuralNet):
    def init (self, input size, hidden size, output size):
        super(DenseNet, self). init ()
       self.blocks = nn.ModuleList([nn.Linear(input size, hidden size),
                                        nn.ReLU(),
                                        nn.Linear(hidden size, output size)])
class ConvNet (NeuralNet):
    def init (self, input channels, num classes):
        super(ConvNet, self). init ()
        self.conv1 = nn.Conv2d(input channels=16,
                                kernel size=3.
                                stride=1.
                                padding=1)
       self.fc = nn.Linear(16 * 28 * 28, num classes)
        self.blocks = nn.ModuleList([self.conv1, nn.ReLU(), self.fc])
```

### LOCAL ENVIRONMENT SETUP Why Set UP a Local Environment?

Setting up a local environment is crucial for deep learning projects for several reasons:

- Dependency Isolation: Keeps project-specific dependencies separate from system-wide packages, preventing version conflicts.
- **Reproducibility:** Ensures consistent dependencies, allowing for reproducible experiments and easier collaboration.
- Project Portability: Makes it simpler to share the project with others, as they can replicate the exact environment. Tools for Creating Local Environments:
  - virtualenv: Lightweight environment creation tool
  - conda: Environment management with additional package management capabilities

## LOCAL ENVIRONMENT SETUP Why Use Virtualenv?

virtualenv offers several advantages over conda for local environments:

- **Lightweight:** virtualenv is faster to create and has a smaller footprint compared to conda.
- **Python Version Agnostic:** Works with different versions of Python, regardless of the system's Python version.
- Flexible Integration: Compatible with pip and various Python packaging tools, making it easy to install dependencies from requirements.txt.

### **Creating a Virtual Environment:**

# Install virtualenv if not already installed
pip install virtualenv

# Create a new virtual environment
virtualenv my\_env

## LOCAL ENVIRONMENT SETUP

ESSENTIAL COMMANDS FOR VIRTUALENV

Once the virtual environment is created, here are the key commands to manage it:

#### Activate the Environment:

# On Windows
my\_env\Scripts\activate

# On macOS/Linux
source my\_env/bin/activate

### Deactivate the Environment:

deactivate

### Install Dependencies:

# Install a package
pip install <package\_name>

# Install from a requirements file
pip install -r requirements.txt

### Tip: Always activate the environment before installing packages.

## LOCAL ENVIRONMENT SETUP Why Use Conda?

conda provides several advantages over virtualenv for local environments:

- **Built-In Package Management:** conda can install packages directly.
- Cross-Language Support: Supports packages beyond Python, such as R or C libraries, making it versatile for data science projects.

**Easier Dependency Resolution:** conda handles complex dependencies automatically, reducing conflicts.

### **Creating a Conda Environment:**

```
# Create a new environment
conda create --name my_env python=3.8
```

## LOCAL ENVIRONMENT SETUP

ESSENTIAL COMMANDS FOR CONDA

Once the conda environment is created, here are the key commands to manage it:

#### Activate the Environment:

conda activate my\_env

#### **Deactivate the Environment:**

conda deactivate

#### Install Dependencies:

# Install a specific package
conda install <package\_name>

# Install from an environment file
conda env update --file environment.yml

Tip: Use environment.yml files for reproducible setups when sharing your project.

INTRODUCTION TO ARGPARSE

The argparse library in Python allows for flexible command-line argument parsing, which is essential for configuring scripts without modifying code.

- **Customizable Scripts:** Enable users to specify parameters at runtime.
- **Improves Code Reusability:** Parameters can be adjusted dynamically, making the script adaptable across tasks.

• User-Friendly Interface: Provides help messages and default values, making it easier to use scripts.

Example:

```
import argparse
parser = argparse.ArgumentParser(description="My script description")
```

ADDING ARGUMENTS

Using argparse, we can add different types of arguments with various options:

#### **Explanation:**

- ► -epochs: Integer argument with a default value of 10.
- ▶ -lr: Float argument with a default value of 0.001.
- batch-size: Integer argument to specify the batch size.

EXAMPLE USAGE OF ARGPARSE

Once arguments are defined, we can access them via args:

# Accessing arguments
print(f"Train for {args.epochs} epochs with batch size {args.batch\_size}")

Example of running the script with arguments:

```
python my_script.py
```

python my\_script.py --epochs 20 --lr 0.01 --batch-size 64

```
python my_script.py --epochs 200 --lr 1e-5
```

This allows the script to be configured dynamically without changing the code.

RUNNING PYTHON COMMANDS IN THE TERMINAL (MAC/LINUX)

Running Python scripts from the terminal on Mac and Linux provides flexibility and resource control:

Basic Command:

python my\_script.py --epochs 10 --lr 0.001

- Background Execution: Running a script in the background frees up the terminal for other tasks. python my\_script.py --epochs 10 --lr 0.001 &
- **Monitoring Background Jobs:** Use jobs to list background jobs and fg to bring them to the foreground.
- Importance of Background Execution: Ideal for long-running tasks, especially on servers or remote machines.

RUNNING PYTHON COMMANDS IN THE TERMINAL (WINDOWS)

Running Python scripts from the Command Prompt or PowerShell on Windows is straightforward:

**Basic Command:** 

python my\_script.py --epochs 10 --lr 0.001

#### **•** Running in Background:

- Use pythonw to launch a new command window: pythonw my\_script.py --epochs 10 --lr 0.001
- Alternatively, use start /B for background execution without opening a new window: start /B python my\_script.py --epochs 10 --lr 0.001

**Importance of Background Execution:** Allows other terminal tasks and is useful for long-running processes.

SCHEDULING JOBS WITH CRON AND ALTERNATIVES

Cron is a task scheduler for Unix-based systems that runs scripts or commands at scheduled intervals.

- Cron Syntax: The format is min hour day month day\_of\_week command.
- ► Example:

# Run a script every day at midnight
0 0 \* \* \* /usr/bin/python /path/to/my\_script.py --epochs 10

#### **Alternative Tools:**

- ► MacOS: Use launchd, a built-in tool for scheduling tasks.
- Windows: Use the Task Scheduler, which provides a GUI for managing scheduled tasks.

Why Scheduling Matters: Automating tasks saves time and allows for consistent, timely execution of regular jobs (e.g., retraining a model).

## MONITORING EXPERIMENTS Why Monitoring Matters

Monitoring experiments is essential in deep learning projects for several reasons:

- **Track Progress:** Provides insights into training dynamics, allowing for better optimization and tuning.
- Diagnose Issues Early: Identify potential issues (e.g., vanishing gradients, overfitting) quickly, saving time and resources.
- Reproducibility and Accountability: Logs and metrics enable reproducibility and allow other collaborators to understand the experiment setup and results.
- Model Comparisons: Facilitates comparison of different runs and hyperparameters to identify the best-performing model.

# MONITORING EXPERIMENTS

TOOLS FOR MONITORING EXPERIMENTS

There are various tools for effective monitoring and tracking:

- **Basic Logging:** Use the logging module to track metrics, errors, and messages within code.
- **JSONL Files:** Store logs as structured data, allowing for detailed tracking and easy parsing.
- **TensorBoard:** Provides a visual interface for monitoring metrics, loss curves, and other parameters in real-time.
- Weights & Biases: A powerful tool for experiment tracking and visualization, allowing for collaborative analysis and sharing.

These tools support different needs, from lightweight solutions (logging) to more robust experiment tracking (TensorBoard, Weights & Biases).

# MONITORING EXPERIMENTS

The logging module in Python is a flexible way to capture and track key information during training:

- Purpose: Logs provide structured, timestamped messages about code execution, helping with debugging and monitoring.
- **Logging Levels:** logging supports different log levels:
  - **INFO:** General messages about code execution and key steps.
  - WARNING: Messages for potential issues that are not critical.
  - **ERROR:** Logs errors that may stop part of the execution.
- **Formats:** Logs can include timestamps, filenames, and message content for better traceability.

## MONITORING EXPERIMENTS Setting UP a Logger

Here's how to set up a basic logger in Python:

```
import logging
# Configure the logger
logging.basicConfig(
    filename='experiment.log',
    level=logging.INFO,
    format='%(asctime)s - %(levelname)s - %(message)s'
)
# Log messages
logging.info("Training started")
logging.warning("Batch size might be too large")
logging.error("Out of memory error")
```

**Explanation:** 

- ▶ filename: Saves logs to a file for permanent record.
- level: Sets the minimum log level (INFO, WARNING, ERROR) for capture.
- ► format: Configures the log format with timestamp, level, and message.
### MONITORING EXPERIMENTS Saving Metrics and Plots

For effective monitoring, it's important to save, rather than print, key metrics and visualizations:

- Metrics Tracking: Save loss, accuracy, and other key metrics to log files or structured formats (e.g., JSON, JSONL) for later analysis.
- Plotting Loss Curves: Generate and save plots of loss, accuracy, and other metrics during training, rather than printing them in the terminal.
- **Tools for Visualization:** Use tools like matplotlib to save plots, or real-time dashboards such as:
  - **TensorBoard:** Visualize training metrics live for a smoother debugging experience.
  - Weights & Biases: Track metrics, visualize experiments, and compare runs.

Tip: Automate saving plots and metrics to ensure nothing is missed during training.

### VERSIONING YOUR CODE Why Git Matters

Git is essential for managing code in deep learning projects for several reasons:

- Version Control: Git tracks changes over time, enabling easy rollback to previous versions and enhancing reproducibility.
- Collaboration: Git allows multiple contributors to work on the same project, with tools to merge code and resolve conflicts.
- Backup and Security: Code stored on remote servers (e.g., GitHub, GitLab) provides a secure, backed-up version of the project.
- Experiment Tracking: Branching allows you to test experimental features or new ideas without affecting the main code.

### VERSIONING YOUR CODE How Git Works

Git is a distributed version control system where every version of the repository is saved locally and remotely (e.g., on GitHub):

- Local Repository: Changes are tracked locally on your machine, allowing for offline work and complete version history.
- Remote Repository: A server (e.g., GitHub, GitLab) stores your code and history, making it accessible for collaboration.

**Commit History:** Every change is saved as a "commit," building a history of modifications for easy review. Each commit is like a snapshot of the project at a specific point in time, allowing for seamless version tracking.

### VERSIONING YOUR CODE Essential Git Commands

Here are some fundamental Git commands to manage your repository:

git pull: Fetch updates from the remote repository and integrate them into your local repository. git pull origin main

**git add:** Stage changes for the next commit.

**git rm:** Remove a file from the repository.

git rm <file\_name>

#### VERSIONING YOUR CODE COMMITTING AND PUSHING CHANGES

Committing and pushing changes are key steps in Git:

- git commit: Save a snapshot of the staged changes to the local repository with a message describing the change. git commit -m "Describe the change"
- git push: Upload your commits from the local repository to the remote repository. git push origin main

Regular commits with clear messages help track progress and make it easier to understand code changes.

#### VERSIONING YOUR CODE BRANCHING AND PULL REQUESTS

Git branches enable experimentation and collaboration:

**Branches:** Create separate lines of development to test features without affecting the main codebase.

| git | branch new-feature   | # | Create | а  | brand | ch     |
|-----|----------------------|---|--------|----|-------|--------|
| git | checkout new-feature | # | Switch | to | the   | branch |

- Merging: Combine changes from a branch back into the main branch once the feature is ready.
  - git checkout main git merge new-feature
- Pull Requests (PRs): PRs are used on platforms like GitHub to review and approve code before merging it into the main branch.

#### VERSIONING YOUR CODE GIT BRANCHING WORKFLOW

The following diagram represents a typical Git workflow with branching, committing, and pushing changes.



### Part VI

## **PROJECT PRESENTATION**

# MINI PROJECT 2: Adversarial Attacks A dataset



### MINI PROJECT 2: ADVERSARIAL ATTACKS

A DECISION BOUNDARY



## MINI PROJECT 2: ADVERSARIAL ATTACKS A CLASSIFIER



### MINI PROJECT 2: ADVERSARIAL ATTACKS

CHOOSING A DATA POINT



### MINI PROJECT 2: ADVERSARIAL ATTACKS

PERTURBING THE DATA POINT



#### MINI PROJECT 2: ADVERSARIAL ATTACKS Adversarial Attacks

What if  $\delta$  is imperceptible ?

### MINI PROJECT 2: ADVERSARIAL ATTACKS

ADVERSARIAL ATTACKS IN IMAGE RECOGNITION



Source : Explaining and Harnessing Adversarial Examples, Goodfellow et al, ICLR 2015.

### MINI PROJECT 2: ADVERSARIAL ATTACKS

ADVERSARIAL ATTACKS IN IMAGE RECOGNITION



Figure. Adversarial traffic signs (Sitawarin, Bhagoji et al., 2018)

# MINI PROJECT 2: ADVERSARIAL ATTACKS DEFINITIONS

#### To be imperceptible, the norm of the perturbation is bounded

We define an  $\epsilon \in \mathbb{R}$  such that  $\|\delta\|_p \leq \epsilon$ . In practice, we use  $\ell_2$  and  $\ell_{\infty}$  norm to bound the perturbation.

#### Generating a adversarial example

Let  $f : \mathbb{R}^d \to \mathcal{Y}$  be a classifier. Given an example  $x \in \mathcal{X} \subset \mathbb{R}^d$  and its true label  $y \in \mathcal{Y}$ , the goal is to find  $\delta \in \mathbb{R}^d$  such that :

**Untargeted attacks**  $\|\delta\|_p \le \epsilon \text{ and } f(x+\delta) \ne y$ 

## **Targeted attacks** $\|\delta\|_{v} \le \epsilon$ and $f(x + \delta) = t$ with $t \ne y$

### MINI PROJECT 2: ADVERSARIAL ATTACKS

Generating an adversarial example with  $\ell_2\text{-}\text{NORM}$ 



### MINI PROJECT 2: ADVERSARIAL ATTACKS

Generating an adversarial example with  $\ell_\infty\text{-}\mathsf{NORM}$ 



### MINI PROJECT 2: ADVERSARIAL ATTACKS FGSM Attack

#### FGSM

The Fast Gradient Sign Method (FGSM) is an attack scheme that uses the gradients of the neural network to create adversarial examples, it is defined as:

 $x_{adv} = x + \epsilon \cdot \operatorname{sign}(\nabla_x L(\theta, x, y))$ 

Source: Explaining and Harnessing Adversarial Examples, Goodfellow et. al, ICLR 2015.

# MINI PROJECT 2: Adversarial Attacks $\ell_2$ -PGD attack

#### $\ell_2\text{-PGD}$

 $\ell_2$ -PGD is an iterative method similar to  $\ell_\infty$ -PGD, but it constrains the perturbation to an  $\ell_2$ -norm ball. The iteration is defined as follows:

- 1.  $x_0 \leftarrow x$
- 2. repeat *n* times :

 $x_{t+1} = \prod_{B_2(x,\epsilon)} \left( x_t + \eta \nabla_x L_\theta(x_t, y) \right)$ 

Source:

Towards Deep Learning Models Resistant to Adversarial Attacks, Madry et. al, ICLR 2018.

### MINI PROJECT 2: ADVERSARIAL ATTACKS $\ell_2$ -PGD Attack



Alexandre Vérine

# MINI PROJECT 2: Adversarial Attacks $\ell_{\infty}$ -PGD attack

#### $\ell_\infty\text{-PGD}$

 $\ell_\infty\text{-}PGD$  is an iterative method that constructs the perturbed data as follows :

1.  $x_0 \leftarrow x$ 

2. repeat *n* times :

 $x_{t+1} = \prod_{B_{\infty}(x,\epsilon)} (x_t + \eta sign(\nabla_x L_{\theta}(x_t, y)))$ 

Source:

Towards Deep Learning Models Resistant to Adversarial Attacks, Madry et. al, ICLR 2018.











• Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .



• Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}$  that approximate *P*:



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}_G$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .
  - 2. Take *a generator model G* represented by a neural network. Take  $\hat{P}_G = G \# Q$ .

FRAMEWORK

- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\widehat{P}_G$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .
  - 2. Take *a generator model G* represented by a neural network. Take  $\hat{P}_G = G \# Q$ .
  - 3. Compute  $G^{\text{opt}}$  that minimize *a dissimilarity measure D* between *P* and  $\hat{P}_G$ :

$$G^{\text{opt}} = \operatorname*{arg\,min}_{G} D(P, \widehat{P}_{G})$$



- Assumption: There is an unknown *target distribution* P in  $\mathcal{X} \subset \mathbb{R}^d$ .
- Goal: Learn a *parameterized distribution*  $\hat{P}$  that approximate *P*:
  - 1. Consider a distribution Q in *a latent space*  $\mathcal{X} \subset \mathbb{R}^m$ , usually  $\mathcal{N}(0, I_m)$ .
  - 2. Take *a generator model G* represented by a neural network. Take  $\hat{P} = G \# Q$ .
  - 3. Compute  $G^{\text{opt}}$  that minimize *a dissimilarity measure D* between *P* and  $\widehat{P}$ :

$$G^{\text{opt}} = \operatorname*{arg\,min}_{G} D(P, \widehat{P})$$

#### GENERATIVE MODELS IN PRACTICE



Midjourney v5 (2023)



Prompt: *A dog playing with a child*.

2

#### GENERATIVE MODELS IN PRACTICE



Midjourney v5 (2023)



Prompt: *A dog playing with a child*.

 $\neq$
## GENERATIVE MODELS



 $\widehat{P} \neq P$ 

## GENERATIVE MODELS



### Low Diversity

### GENERATIVE MODELS



### $\label{eq:precision} Precision \ \text{and} \ Recall \ \text{for Generative Models}$

METRICS TO EVALUATE QUALITY AND DIVERSITY

To assess models, we use the notion of Precision and Recall, inspired from Information Retrieval:

## Quality



### PRECISION AND RECALL FOR GENERATIVE MODELS

METRICS TO EVALUATE QUALITY AND DIVERSITY

To assess models, we use the notion of Precision and Recall, inspired from Information Retrieval:



What proportion of generated samples are realistic?

### PRECISION AND RECALL FOR GENERATIVE MODELS

METRICS TO EVALUATE QUALITY AND DIVERSITY

To assess models, we use the notion of Precision and Recall, inspired from Information Retrieval:



What proportion of generated samples are realistic? What proportion of real samples can be generated?

## PRECISION AND RECALL FOR GENERATIVE MODELS FOR FINITE SUPPORT



# PRECISION AND RECALL FOR GENERATIVE MODELS FOR FINITE SUPPORT

#### Definition 2.1 (Support-Based Precision and Recall - [6].)

For any distributions  $P \in \mathcal{P}(\mathcal{X})$  and  $\widehat{P} \in \mathcal{P}(\mathcal{X})$ , we say that the distribution P has precision  $\overline{\alpha}$  at recall  $\overline{\beta}$  with respect to  $\widehat{P}$  if

$$\bar{\alpha} = \widehat{P}(\operatorname{Supp}(P)) \quad et \quad \bar{\beta} = P(\operatorname{Supp}(\widehat{P})).$$
 (2)

Precision for finite support is the proportion of generated data that lies on the support of the real data:

$$\bar{\alpha} = \widehat{P}(\operatorname{Supp}(P)).$$



Precision for finite support is the proportion of generated data that lies on the support of the real data:

$$\bar{\alpha} = \widehat{P}(\operatorname{Supp}(P)).$$



Recall for finite support is the proportion of the support of the real data that is covered by the generated data:

$$\bar{\beta} = P(\operatorname{Supp}(\widehat{P})).$$



Recall for finite support is the proportion of the support of the real data that is covered by the generated data:

$$\bar{\beta} = P(\operatorname{Supp}(\widehat{P})).$$



# PRECISION AND RECALL FOR GENERATIVE MODELS IN PRACTICE

MNIST Dataset [7]



# PRECISION AND RECALL FOR GENERATIVE MODELS IN PRACTICE



Precision: 0.54 Recall: 0.91

Precision: 0.80 Recall: 0.70

# PRECISION AND RECALL FOR GENERATIVE MODELS FOR LLMS

On open-ended generation, the quality and diversity of LLMs can also be evaluated using Precision and Recall: [2]



#### GENERATIVE ADVERSARIAL NETWORKS ORIGINAL FRAMEWORK

- Let  $G : \mathcal{Z} \to \mathcal{X}$  be a generator model parameterized by a neural network.
- Let  $D : \mathcal{X} \to [0, 1]$  be a discriminator model parameterized by a neural network.

The original GAN framework [3] is defined by the following optimization problem:

$$\min_{G} \max_{D} \mathbb{E}_{x \sim P} \left[ \log D(x) \right] + \mathbb{E}_{x \sim \widehat{P}_{G}} \left[ \log(1 - D(x)) \right].$$
(3)

#### GENERATIVE ADVERSARIAL NETWORKS ORIGINAL FRAMEWORK

- Let  $G : \mathcal{Z} \to \mathcal{X}$  be a generator model parameterized by a neural network.
- Let  $D : \mathcal{X} \to [0, 1]$  be a discriminator model parameterized by a neural network.

The original GAN framework [3] is defined by the following optimization problem:

$$\min_{G} \max_{D} \mathbb{E}_{x \sim P} \left[ \log D(x) \right] + \mathbb{E}_{x \sim \widehat{P}_{G}} \left[ \log(1 - D(x)) \right].$$
(3)

# TUNING PRECISION AND RECALL IN GENERATIVE MODELS TRUNCATION

Hard Trunctation Karras et al. [4]

Soft Trunctation Kingma and Dhariwal [5]

# FINAL PROJECT: GENERATIVE ADVERSARIAL NETWORKS HARD TRUNCATION



**Figure.** From left to right:  $\psi = 0.0$ ,  $\psi = 0.3 \ \psi = 0.7 \ \psi = 1.0$ .



# FINAL PROJECT: GENERATIVE ADVERSARIAL NETWORKS SOFT TRUNCATION



(a)  $\psi = 0.04$ 

(b)  $\psi = 0.5$ 

(c)  $\psi = 1.0$ 

(d)  $\psi = 2.0$ 

Figure. Soft-Truncation on BigGAN. Source:[1].

#### FINAL PROJECT TUNING QUALITY AND DIVERSITY IN GANS

- Goal: Implement a GAN architecture to generate samples from a given dataset.
- Dataset: MNIST
- Code: Use the repository given in the Project
- Enrolling: GitHub Classroom https://classroom.github.com/a/lj2\_ipKW

### Possible improvements

- ► f-GANs
  - 1. f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization
  - 2. Precision-Recall Divergence Optimization for Generative Modeling with GANs and Normalizing Flows
- ► WGAN
  - 1. Wasserstein GAN
- Rejection Sampling
  - 1. Discriminator Rejection Sampling
  - 2. Metropolis-Hastings Generative Adversarial Networks
  - 3. Optimal Budgeted Rejection Sampling for Generative Models
- Latent Rejection sampling
  - 1. Latent reweighting, an almost free improvement for GANs
- Gradient ascent
  - 1. Discriminator optimal transport
  - 2. Refining Deep Generative Models via Discriminator Gradient Flow
  - 3. Your GAN is Secretly an Energy-based Model and You Should use Discriminator Driven Latent Sampling
- Classifier guidance generation
  - 1. MMGAN: Generative Adversarial Networks for Multi-Modal Distributions

#### FINAL PROJECT: GENERATIVE ADVERSARIAL NETWORKS FINAL PROJECT: TUNING OUALITY AND DIVERSITY IN GANS

Specifics:

- ▶ For fair comparison, use the same *generator* architecture as the one in the project.
- Your code will be pulled and test every 24 hours (at 02:00 AM).
- ▶ You have to choose one of the methods (or any other method that makes sense to you) and implement it.

During the presentation, you will have to:

- Present the paper you chose.
- Present how you adapted the method to the project.
- Explain your results.
- Explain the limitations of your method.

You will be graded:

- ▶ on the quality of the presentation.
- on how well you understood the method.
- on how well you adapted the method to the project.

You will not be graded on:

- Your rank on the platform.
- ► The quality of your code.

### FINAL PROJECT: GENERATIVE ADVERSARIAL NETWORKS

The work is:

- ► Individual.
- Due on the 10th of December at 8:15 AM.

The report and the presentation:

- The report should be a PDF file to push on the Github repository.
- ► The report should be at most 5 pages long. (excluding references and title page)
- ► The report should be written in English.
- ▶ The presentation should be a PDF file to push on the Github repository.
- ► The presentation must be 15 minutes long (hard limit) + 5 minutes for questions.
- The presentation should be in English or French.

### **REFERENCES I**

- [1] Brock, A., Donahue, J., and Simonyan, K. (2019). Large Scale GAN Training for High Fidelity Natural Image Synthesis. arXiv:1809.11096 [cs, stat].
- [2] Bronnec, F. L., Verine, A., Negrevergne, B., Chevaleyre, Y., and Allauzen, A. (2024). Exploring Precision and Recall to assess the quality and diversity of LLMs. *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*. arXiv:2402.10693 [cs].
- [3] Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative Adversarial Networks. In 27th Conference on Neural Information Processing Systems (NeurIPS 2014). arXiv: 1406.2661.
- [4] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., and Aila, T. (2020). Analyzing and Improving the Image Quality of StyleGAN. arXiv:1912.04958 [cs, eess, stat].
- [5] Kingma, D. P. and Dhariwal, P. (2018). Glow: Generative Flow with Invertible 1x1 Convolutions. In 32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Montréal, Canada., volume 31.
- [6] Kynkäänniemi, T., Karras, T., Laine, S., Lehtinen, J., and Aila, T. (2019). Improved Precision and Recall Metric for Assessing Generative Models. In 33rd Conference on Neural Information Processing Systems (NeurIPS 2019), Vancouver, Canada. arXiv: 1904.06991.
- [7] Yann LeCun, Corinna Cortes, and Burges, C. (2010). MNIST handwritten digit database. ATT Labs, 2.