

Noise2Noise

Matteo Calafà, Paolo Motta, Thomas Rimbob

First Project for the Deep Learning (EE-559) course at EPFL Lausanne, Switzerland

Abstract—In this first part of the project, we analyse a possible denoising model aiming to restore images just by looking at corrupted examples, hence the name *Noise2Noise*. The based network architecture adopted is a U-Net, in which skip connections between the contracting and expansive paths are used.

I. INTRODUCTION

In many applications such as image denoising, image compression, and, in some cases, even image data generation, there is the necessity to explicitly model a high-dimensional signal. To this end, the use of an autoencoder allows a neural network to learn an efficient representation of unlabeled data by training the network to ignore signal “noise”. In the case of image denoising, the main idea is therefore to capture a small number of degrees of freedom that represent the physical context, and from these perform an efficient reconstruction. The auto-encoder therefore has two main building blocks. The first is an encoder, whose objective is to learn a lower-dimensional representation (encoding) for a higher-dimensional data. This is typically used for dimensionality reduction, and is achieved by training the network to capture the most important parts of the input image. The second is a decoder, whose objective is instead to reconstruct the image from its lower-dimensional representation.

The idea of the auto-encoder underlies the structure of the U-Net. This structure introduces *skip connections* between encoding and decoding layers, concatenating the states. This allows the U-Nets to use fine-grained details learned in the encoder part to reconstruct an image in the decoder part. Although the structure of the U-Net is mainly used for image segmentation tasks, in this paper we would like to show how the same structure can be adapted for an image denoising task, as already shown in [5]. As also stated in this paper, most image denoising models are trained using large numbers of pairs (\hat{x}_i, y_i) of noisy inputs \hat{x}_i and clean reference images y_i , but here the aim is instead to show how satisfactory results can be obtained even without clean samples, if the noise is additive and unbiased.

II. DATASET AND DATA AUGMENTATION

A. First Analysis of the Dataset

The dataset consists of 50'000 examples for the training set and 1000 examples for the test one. Each example is represented by a noisy pair of RGB images that are downsampled and pixelated. All images are 32×32 pixels in size. The networks proposed have the goal to reduce the effects of downsampling on unseen images.

In [Figure 1](#) we can observe two examples from the test set that we will reuse later in order to visually observe the results of our model.



Fig. 1: Noisy starting image (left) and target image (right)

B. Data Augmentation

A data augmentation technique is used to extend the dataset by adding transformed copies of already existing images. In doing so, we can increase the generalizability and robustness of our model, avoiding overfitting as well. In particular, we decided to add a horizontally and a vertically flipped version to our dataset. The motivation for this choice lies in the fact that our model must be able to reconstruct the main features of the image regardless of its possible orientation.

Other transformations, both geometric and colour-based, can be performed on the images, but in most cases they do not bring any real improvement other than a deterioration in performance when colours are changed.

III. MODELS AND METHODS

A. U-Net Architecture

In the original structure of the U-Net [6], the input and output images have different sizes. To be able to evaluate denoising quality with the same size, we reworked the model by changing the parameters of the convolution function. As mentioned earlier, the U-Net network can be divided into two parts. The first is the contracting path which uses a typical CNN architecture. Each block in the contracting path consists of two successive 3×3 convolutions with

padding followed by a *Leaky ReLU* activation unit and a *downsampling* layer. However, despite the fact that the downsampling layer is typically implemented through a max-pooling operation, we decided to replace this layer with a convolution with a larger stride. Indeed, max-pooling (or any kind of pooling) is a fixed operation and replacing it with a strided convolution can also be seen as learning the pooling operation, which increases the model’s expressiveness ability [8]. This structure is repeated several times. The main characteristic of U-Net comes in the second part, called the expansive path, in which each stage upsamples the feature map using 3×3 *transposed convolution* with a larger stride. Then, the feature map from the corresponding layer in the contracting path is *concatenated* onto the upsampled feature map. This is followed by two successive 3×3 convolutions with padding and a *Leaky ReLU* activation. At the final stage, an additional sequence of convolutions 3×3 is applied to reduce the feature map to the required number of channels and produce the denoised image. The overall result is a network with a U-shape and, more importantly, it propagates contextual information along the network, which allows it to properly reconstruct the context. Figure 4 illustrates the overall U-net architecture.

B. Residual U-Net Architecture

Along with the U-Net presented above, we also decided to implement a different structure of our network, combining the benefits of the U-Net structure with that of a Res-Net [2]. Instead of directly predicting the denoised image, the model predicts the residual noise of the corrupted image. Such a structure also makes it possible to use a network with many more layers without running into the vanishing gradient problem. The resulting final structure is called *Residual U-Net*, and although this architecture is also mainly used for image segmentation [9], we decided to adapt it for our denoising task.

The structure of the network is similar to that described for the U-Net, with the main differences being that each block is implemented as a residual block. Since within the residual block the input must be added to the output of the block, if the two quantities have discordant channel sizes, a 1×1 convolution is adopted to scale the number of input channels. Furthermore, it has been proven in the literature that *Batch Normalization* makes the training process smoother. However, it requires a sufficiently large batch size, and our choice of `batch_size = 30` may be too restrictive. For this reason, we also tried to make use of *Group Normalization* which, unlike batch normalization, does not require a very large number of batches, as it divides the channels into groups and normalizes the features within each group.

After experimenting without normalization, with batch normalization and with group normalization, we observed that the latter yielded much better results and therefore

decided to adopt it for our final model. Figure 2 illustrates the overall Residual U-Net architecture.

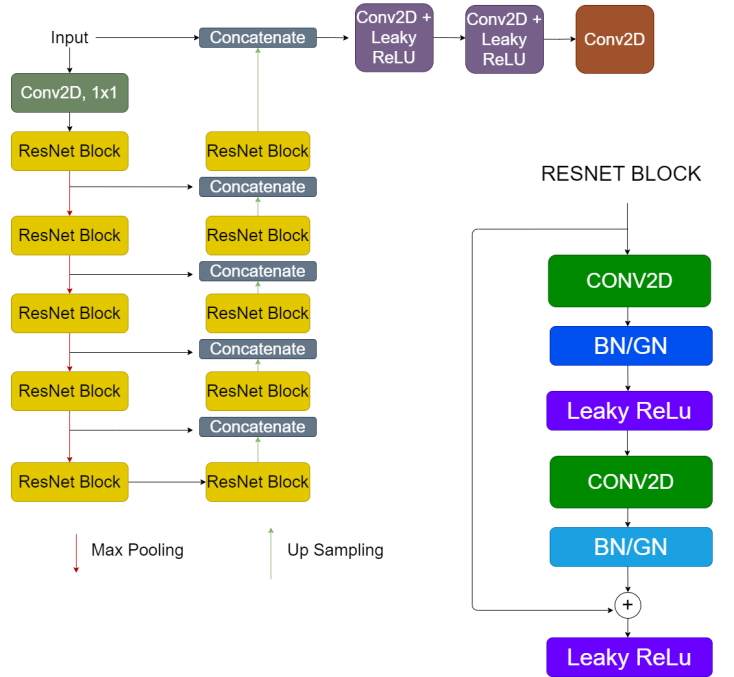


Fig. 2: Residual U-Net Architecture

C. Weight Initialization

In order to prevent layer activation outputs from exploding or vanishing gradients during training, an appropriate initialization of the weights is necessary and allows to achieve better performance. Among the most widely used are the Xavier normalization and He-et-al normalization when dealing with convolutional layers [3]. We experimented with both initializations for our networks and the former seems to lead to slightly better performances than the latter. We therefore stuck to this one in both models’ implementations.

IV. RESULTS

In this section we report the results obtained with our best model given by the U-Net architecture, together with the final performance of the Residual U-Net. The model’s performance is calculated through the PSNR between our prediction (denoised image I) and the target image T , given by:

$$\text{PSNR} := 20 * \log_{10} \left(\frac{\text{MAX}\{I\}}{\sqrt{\text{MSE}(I, T)}} \right)$$

where $\text{MAX}\{I\}$ is the maximum possible pixel value of the image and $\text{MSE}(I, T)$ is the mean square error between the images:

$$\text{MSE}(I, T) := \frac{1}{MN} \sum_i^{M-1} \sum_j^{N-1} \|I(i, j) - T(i, j)\|^2$$

Figure 3 illustrates the trend of the training loss and the PSNR on the test set during training, using 30 epochs.

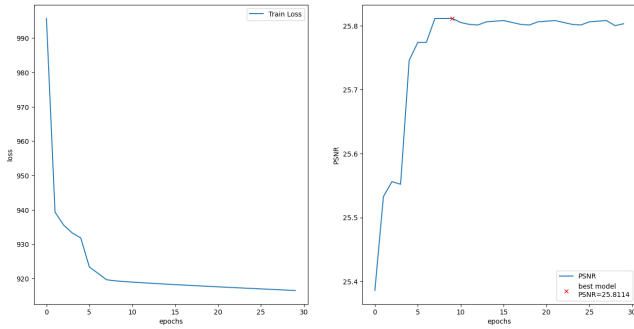


Fig. 3: Evolution of the loss (left) and PSNR (right)

Finally, we report in Figure 5 the results of our best model on the two pairs of test images previously presented.

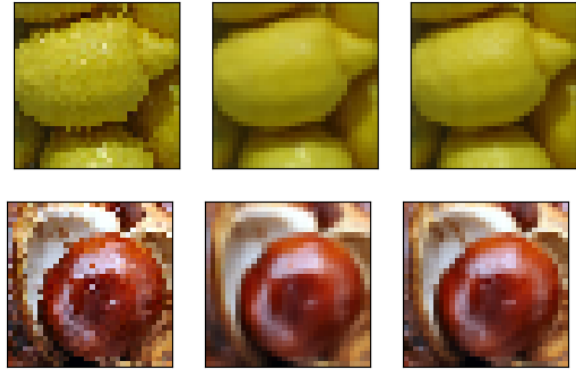


Fig. 5: Noisy input image (left), denoised output image (center) and target image (right)

The PSNR obtained on the test set with the best models are summarized in the following table.

Network Architecture	PSNR
U-Net	25.811
Residual U-Net	25.803

V. CONCLUSION

The obtained results confirm that the architecture of a U-Net can also be adapted for a denoising task. However, using the network on images with such a small resolution (32×32) does not allow to fully exploit its potential. The structure of the network also allows its use on larger images, on which its effectiveness could be tested. Finally, several variants have been built on the same basic architecture and can therefore be subsequently used and tested within this task [7].

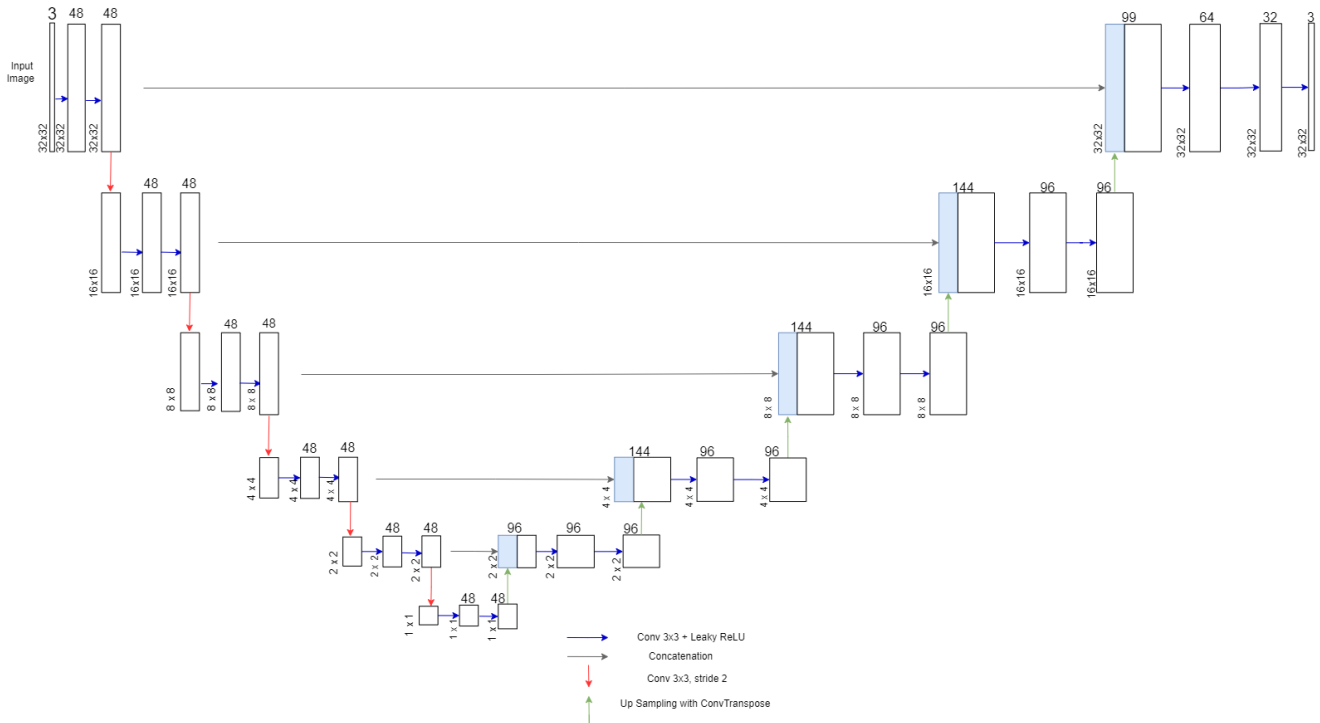


Fig. 4: U-Net Architecture

REFERENCES

- [1] Foivos I. Diakogiannis et al. “ResUNet-a: a deep learning framework for semantic segmentation of remotely sensed data”. In: *CoRR* abs/1904.00592 (2019). arXiv: [1904.00592](https://arxiv.org/abs/1904.00592). URL: <http://arxiv.org/abs/1904.00592>.
- [2] Kaiming He et al. “Deep Residual Learning for Image Recognition”. In: *CoRR* abs/1512.03385 (2015). arXiv: [1512.03385](https://arxiv.org/abs/1512.03385). URL: <http://arxiv.org/abs/1512.03385>.
- [3] Kaiming He et al. “Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification”. In: *CoRR* abs/1502.01852 (2015). arXiv: [1502.01852](https://arxiv.org/abs/1502.01852). URL: <http://arxiv.org/abs/1502.01852>.
- [4] Rina Komatsu and Tad Gonsalves. “Comparing U-Net Based Models for Denoising Color Images”. In: *AI* 1.4 (2020), pp. 465–486. ISSN: 2673-2688. DOI: [10.3390/ai1040029](https://doi.org/10.3390/ai1040029). URL: <https://www.mdpi.com/2673-2688/1/4/29>.
- [5] Jaakko Lehtinen et al. “Noise2Noise: Learning Image Restoration without Clean Data”. In: *CoRR* abs/1803.04189 (2018). arXiv: [1803.04189](https://arxiv.org/abs/1803.04189). URL: <http://arxiv.org/abs/1803.04189>.
- [6] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. “U-Net: Convolutional Networks for Biomedical Image Segmentation”. In: *CoRR* abs/1505.04597 (2015). arXiv: [1505.04597](https://arxiv.org/abs/1505.04597). URL: <http://arxiv.org/abs/1505.04597>.
- [7] Nahian Siddique et al. “U-Net and Its Variants for Medical Image Segmentation: A Review of Theory and Applications”. In: *IEEE Access* 9 (2021), pp. 82031–82057. DOI: [10.1109/ACCESS.2021.3086020](https://doi.org/10.1109/ACCESS.2021.3086020).
- [8] Jost Tobias Springenberg et al. “Striving for Simplicity: The All Convolutional Net”. In: *CoRR* abs/1412.6806 (2015).
- [9] Zhengxin Zhang, Qingjie Liu, and Yunhong Wang. “Road Extraction by Deep Residual U-Net”. In: *CoRR* abs/1711.10684 (2017). arXiv: [1711.10684](https://arxiv.org/abs/1711.10684). URL: <http://arxiv.org/abs/1711.10684>.

Manual implementation of a ConvNet

Matteo Calafà, Paolo Motta, Thomas Rimbob

Second project for the Deep Learning (EE-559) course at EPFL Lausanne, Switzerland

Abstract—This second part consists in the explicit implementation of the main blocks for a simple convolutional network. In other words, we aim at implementing blocks which could replace the PyTorch modules and allow to define a standard network as a sequence of those. While this does not present any practical utility because of the lower computational performance with respect to optimized PyTorch modules, this project’s goal is to clearly comprehend how convolutional networks work and some basic ideas about their implementation in the most common libraries.

I. INTRODUCTION

The goal of this project is to explicitly implement a few modules that are useful for the construction of a simple convolutional neural network. These modules are: Sigmoid, ReLU, 2D Convolutional Layer, Nearest Neighbor Upsampling. In addition to that, we aim to implement a standard optimizer (SGD) and a standard loss (MSE) in order to provide a fully functional Deep Learning pipeline. Finally, to define the network as a sequence of modules, we will define a class named `Sequential` presenting very similar characteristics with the one provided by PyTorch.

II. GENERAL ASPECTS ABOUT THE DEFINITION AND CONSTRUCTION OF THE NETWORK

A. Definition of the network

The network we aim to implement is shown in [Figure 1](#):

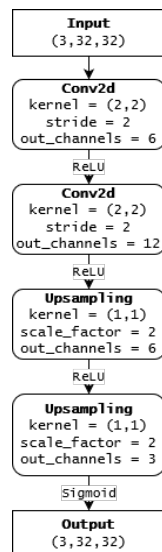


Fig. 1: Diagram of the network.

Upsampling can be defined as the composition of a Nearest Neighbor Upsampling and a 2D Convolutional Layer (respectively the equivalent of `UpsamplingNearest2d` and `Conv2d` in PyTorch).

B. The Module class

Every module has different inner tensors and will define three main methods:

- 1) `forward()`: to perform the forward pass.
- 2) `backward()`: to compute the gradient of the loss with respect to the input of the module, given the gradient with respect to the output.
- 3) `param()`: to return the module parameters and the gradient of the loss with respect to the same parameters (it is empty for parameterless modules, e.g ReLU).

In addition, we aim to implement the MSE loss as another class with the `forward()` and `backward()` methods (in this case, the latter will not accept any input) and the SGD optimizer with the `step()` method.

Finally we need a definition of the `Sequential` class to wrap the different blocks together and build the network. Its overridden `forward()` and `backward()` methods are explained in section [subsection II-C](#).

C. Basic functioning of the network

We stress again that every module/block contains various tensors which keep in memory the latest update of parameters such as weights, gradients, inputs and outputs. The training is based on the usual following steps:

- 1) Forward pass: from the current weights, each block computes the output from an input tensor. Therefore, the `Sequential` class can compute the output of the network by sequentially applying the `forward()` methods from all its building blocks. During this step, each block/module saves in memory its input and output.
- 2) Forward pass of the loss: from the output of the net, the MSE module can compute the error thanks to its `forward()` method.
- 3) Backward pass of the loss: the `backward()` method of the MSE module returns the gradient of the loss with respect to the output of the net.
- 4) Every block can compute the loss gradient with respect to its input from the gradient with respect to its output. Therefore, starting from the gradient returned by the loss, all the gradients can be recursively computed with the `backward()` methods from all the blocks, called

in reverse order. This is indeed the definition of the overridden `backward()` method in `Sequential`.

- 5) Optimization step: from the gradients with respects to the states, we can compute the gradients with respect to the weights/parameters through the `param()` methods of each module. Having the current values and gradients of all the weights and biases, the `SGD` class can perform the optimization step via its `step()` method.

III. OPERATING DEFINITIONS OF THE MODULES

In this section, we present how `Conv2d` and `Upsampling` explicitly perform their forward operation. The remaining blocks follow instead the standard definitions. Before proceeding, we explicit the following nomenclature:

- I and O are the input and output tensors of a generic layer.
- W and B are the weight and bias tensors of a convolutional layer.
- \mathcal{L} indicates the loss evaluated in a current state.
- B, C, N_X, N_Y are the batch size, the channels number and the x-y lengths of the image at a generic layer.

A. Nearest Neighbor Upsampling

The goal of this module is to increase the tensor's spatial size (based on a `scale_factor` parameter) by filling it with repeated values of the original tensor. In particular, let us assume that the original shape of the tensor is (B, C, N_x, N_y) . If `scale_factor` = k , the output tensor's shape should be (B, C, kN_x, kN_y) . The Nearest Neighbor Upsampling fills up the newly added $k \times k$ squares with the corresponding original values in the input tensor. We refer below a 2D example, inspired from the [PyTorch documentation page](#):

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \xrightarrow{k=2} \begin{bmatrix} 1 & 1 & 2 & 2 \\ 1 & 1 & 2 & 2 \\ 3 & 3 & 4 & 4 \\ 3 & 3 & 4 & 4 \end{bmatrix} \quad (1)$$

B. Conv2d as a linear layer

One can prove that convolution is in fact a linear operation. However, we a priori do not have access to suitable tensors for the linear products. Rather, they need a specific construction. We claim that there exist reshaped versions of the previously mentioned tensors (called $\tilde{W}, \tilde{B}, \tilde{O}$) such that:

$$\tilde{O} = \tilde{W} \otimes U + \tilde{B} \quad (2)$$

which, in components, it corresponds to:

$$\tilde{O}_{p,q,r} = \sum_m \tilde{W}_{q,m} U_{p,m,r} + \tilde{B}_{0,q,0} \quad (3)$$

where U is the unfolded version of I (see `PyTorch's` `unfold()` method). Therefore, the following forward and backward operations for the `Conv2d` block will be based on the equivalent linear operation in [Equation 2](#)

IV. COMPUTATION OF THE GRADIENTS

A. MSE loss

The definition of the MSE loss for a CNN is:

$$MSE(T, O) := \frac{1}{B \cdot C \cdot N_X \cdot N_Y} \cdot \sum_{B, C, N_X, N_Y} (T_{B, C, N_X, N_Y} - O_{B, C, N_X, N_Y})^2$$

where T and O are the target and output tensors. Therefore, the gradient can be easily computed as:

$$\frac{\partial \mathcal{L}}{\partial O} = \frac{2}{B \cdot C \cdot N_X \cdot N_Y} \cdot \sum_{B, C, N_X, N_Y} (T_{B, C, N_X, N_Y} - O_{B, C, N_X, N_Y})$$

B. Sigmoid

Since the forward operation is by definition $O = \sigma(I)$ (component-wise), where:

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

one can easily compute the gradient with respect to the input with:

$$\frac{\partial \mathcal{L}}{\partial I} = \frac{\partial \mathcal{L}}{\partial O} \odot \sigma'(I)$$

where \odot indicates the Hadamard product and the derivative of the sigmoid can be directly computed with:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

C. ReLU

By definition, the layer forward operation is $O = \text{ReLU}(I)$, where:

$$\text{ReLU}(x) := \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

Hence, as before, the gradient with respect to the input can be computed as:

$$\frac{\partial \mathcal{L}}{\partial I} = \frac{\partial \mathcal{L}}{\partial O} \odot H(I)$$

where, indeed, the Heavyside function H coincides with the ReLU subderivative.

D. Nearest neighbor upsampling

The general idea of the forward operation was already presented in [section III](#). To compute the backward pass, we instead need to sum the derivatives in each $k \times k$ square following intuitively the concept of *weight-sharing*. This intuition can also be mathematically demonstrated by the chain rule:

$$\frac{\partial \mathcal{L}}{\partial I_{i,j}} = \sum_{p,q} \frac{\partial \mathcal{L}}{\partial O_{p,q}} \frac{\partial O_{p,q}}{\partial I_{i,j}}$$

However, following [Equation 1](#), $\frac{\partial O_{p,q}}{\partial I_{i,j}}$ can only have value 1 when (p, q) are in the (i, j) k -square ($=: M_{i,j}^k$) and 0 otherwise. Therefore:

$$\frac{\partial \mathcal{L}}{\partial I_{i,j}} = \sum_{p,q \in M_{i,j}^k} \frac{\partial \mathcal{L}}{\partial O_{p,q}}$$

as expected.

E. Conv2d - gradient of states

Gradients with respects to the input I can be computed component-wise from Equation 3, which implies:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial U_{i,j,k}} &= \sum_{p,q,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,q,r}} \frac{\partial \tilde{O}_{p,q,r}}{\partial U_{i,j,k}} \quad (3) \\ \sum_{p,q,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,q,r}} \tilde{W}_{q,j} \delta_{p,i} \delta_{r,k} &= \sum_q \frac{\partial \mathcal{L}}{\partial \tilde{O}_{i,q,k}} \tilde{W}_{q,j} \end{aligned} \quad (4)$$

This component-wise operation can be computed with a direct tensor product, being careful to transpose some components. Moreover, notice that $\partial \mathcal{L} / \partial \tilde{O}$ can be easily obtained by reshaping $\partial \mathcal{L} / \partial O$. The last missing step is the computation of the gradient, not with respect to the unfolded input but to the very input instead. The `unfold()` operation creates a new tensor filled with the same original values but repeated multiple times. Therefore, to compute the last step, we once again need a *weigh-sharing* operation as in subsection IV-D. Fortunately, PyTorch provides a method called `fold()`, which consists in the inverse of the `unfold()` operation, with the additional multiplication for the repetition times. Therefore, we can directly go from $\partial \mathcal{L} / \partial U$ to $\partial \mathcal{L} / \partial I$ using `fold()`. Schematically:

$$W, \frac{\partial \mathcal{L}}{\partial O} \xrightarrow{\text{reshape}} \tilde{W}, \frac{\partial \mathcal{L}}{\partial \tilde{O}} \xrightarrow{(4)} \frac{\partial \mathcal{L}}{\partial U} \xrightarrow{\text{fold}} \frac{\partial \mathcal{L}}{\partial I}$$

F. Conv2d - gradient of weights and biases

With similar calculations, we get:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{W}_{i,j}} &= \sum_{p,q,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,q,r}} \frac{\partial \tilde{O}_{p,q,r}}{\partial \tilde{W}_{i,j}} \quad (3) \\ \sum_{p,q,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,q,r}} U_{p,j,r} \delta_{q,i} &= \sum_{p,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,i,r}} U_{p,j,r} \end{aligned} \quad (5)$$

which can be implemented with a standard tensor product if tensors are before properly translated and reshaped. Then,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial B_i} &= \frac{\partial \mathcal{L}}{\partial \tilde{B}_{0,i,0}} = \sum_{p,q,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,q,r}} \frac{\partial \tilde{O}_{p,q,r}}{\partial \tilde{B}_{0,i,0}} \quad (3) \\ \sum_{p,q,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,q,r}} \delta_{q,i} &= \sum_{p,r} \frac{\partial \mathcal{L}}{\partial \tilde{O}_{p,i,r}} \end{aligned} \quad (6)$$

which can be implemented with a sum over the first and third component. At this point, we show the sequential steps to get the final gradients:

$$\frac{\partial \mathcal{L}}{\partial O} \xrightarrow{\text{reshape}} \frac{\partial \mathcal{L}}{\partial \tilde{O}} \xrightarrow{(5)} \frac{\partial \mathcal{L}}{\partial \tilde{W}} \xrightarrow{\text{reshape}} \frac{\partial \mathcal{L}}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial O} \xrightarrow{\text{reshape}} \frac{\partial \mathcal{L}}{\partial \tilde{O}} \xrightarrow{(5)} \frac{\partial \mathcal{L}}{\partial B}$$

G. SGD

In order to actually update the network parameters from the gradients, we implement the SGD optimizer in the `SGD` class. In particular, this class provides a `step()` method, whose role is to loop through the model's modules and recursively update their parameters by calling their `update_params()` methods. The update rule is the SGD one, namely going in the opposite direction of the gradient with a tuneable step size:

$$p_{t+1} = p_t - \gamma \nabla_{p_t} \mathcal{L} \quad (7)$$

V. TRAINING, TESTING, SAVING AND LOADING

The `Model` class is the main interface for using our modules. It implements the sequential structure represented in Figure 1, stored in its `self.model` attribute, and defines the loss criterion (MSE) and optimizer (SGD). It also provides a `train()` method, responsible for running the full training loop. This method iterates on the input samples by batches, performs the prediction with `model.forward()`, computes the loss and the gradients with the `backward()` methods, and updates the parameters with the `optimizer.step()` method.

This class also provides support for saving and loading the model. Through the `load_pretrained_model()` and `save_pickle_state()` methods, the parameters and their gradients in each modules can be saved in and loaded from a pickle file `bestmodel.pth`.

VI. CONCLUSION

In order to assess the correctness of our implementation, these blocks have been compared to the PyTorch ones, both individually and sequentially. After careful examination, we concluded that the implementation was in fact correct, and attained an almost complete similarity with the library's performance. However, we point out that because of computational details, which are unrelated to the mathematics described in this report, it is very hard to actually get a fully precise implementation. Some examples are the floating point precision for digits, or different implementations of the same mathematical expressions, which have been proven to change results. In fact, this has been apparent when using an equivalent definition for the sigmoid function and a strict inequality instead of an inequality condition for the ReLU. Regarding the implementation of the CNN in Figure 1 for a Noise2Noise model, we obtained decent results and a PSNR of 19.5, by training the model for about two minutes on a CPU. While the score in the first part of the project is much higher, we remind that the proposed network as well as its optimizer are of very low complexity and could be more finely tuned to get better performances. Moreover, the higher computational times resulting from the manual implementation of the blocks instead of modules taken from optimized libraries prevented us from running large simulations. Despite this, we confirmed the correctness of the implementations and mathematical calculations.