

**SPATIAL ENHANCEMENT OF REMOTELY
SENSED IMAGES USING CONVOLUTIONAL
NEURAL NETWORKS**

Ugo Palatucci

July 2020



Supervisor:

Prof. Restaino Rocco

Contents

Introduction	3
1 Pansharpening state of the art	5
1.1 CS	7
1.2 MRA	12
1.3 Quality Assessment	14
2 Pansharpening applications of Deep Learning	20
2.1 Neural Networks	22
2.1.1 Perceptron	22
2.1.2 Activation Function	23
2.1.3 Layers	24
2.2 The Deep Learning Paradigm	26
2.3 Pansharpening Applications	27
3 Proposed Solution	31
3.1 Loss Function Issue	31

3.2	Automatic Differentiation	33
3.3	Tensorflow and custom loss function implementation	36
3.3.1	Loss Function	38
4	Experimental results	40
4.1	Description	40
4.2	Experimental Settings	43
4.3	Results	46
	Conclusions	56
A	Python implementation of no-reference quality indexes	65
A.1	QNR	65
A.2	HQNR	67

Introduction

Pansharpening refers to a particular data fusion issue where two images, one panchromatic (PAN) and one multispectral (MS), representing the same area can be combined to enhance the peculiarity of both. The panchromatic image is acquired with a wide spectrum sensor that can have a higher spatial resolution compared to a multispectral one. However, the sensor cannot acquire different bands. A multispectral sensor, instead, collects several channels at a lower spatial resolution. As physical constraints occur, the creation of a sensor specialized in both types of resolution would not be possible. For this reason, the fusion of both images can be the only possibility to have a new image with higher spatial and spectral resolution. Pansharpening is an urgent topic for remote sensing; indeed, the result can be used upstream of another process such as change detection [1], object recognition [2], visual image analysis and scene interpretation [3].

Based on a convolutional neural network, a new pansharpening method has been proposed recently [4]. Using degraded versions of the PAN and MS images, the network weights were trained to fuse the images, optimizing a

reference index with the gradient descendent algorithm: the mean squared error. After the training, the same weights were applied to fuse the original PAN and MS images.

The purpose of this work is to improve the current methodology using the original PAN and MS images for a training with a *no-reference* index like Quality with no reference (QNR) or hybrid QNR (HQNR) indexes. Furthermore, the intention is to remove the error added using the degraded images that do not reflect the models were the pansharpening algorithm should be utilised. In [4] the code was written in python 2.7 using the Theano library. The Theano project has been not updated since 2018 [5]. For this reason, the project is incompatible with new Cuda libraries and NVIDIA drivers and many issues to run the training process on the GPU have been founded. To solve that, a new software has been created using TensorFlow. Tensorflow is a more modern and popular Deep Learning library maintained by Google that allows more compatibility with the newest libraries and a larger community, helping the development in case of uncommon errors. Moreover, Tensorflow implements the *automatic differentiation* [6], a critical feature for this research. *Automatic differentiation* allows writing differentiable functions and subsequently using them for the core algorithm of the neural network's backpropagation training: the gradient descendent algorithm.

Chapter 1

Pansharpening state of the art

According to the existing literature, the pansharpening techniques are divided into two main areas: *component substitution* (CS) and *multiresolution analysis* (MRA). The techniques belonging to the first class consist in representing the MS and PAN in a different domain that can entirely split the spatial information from the spectral information. In this domain, the spatial information part of the MS image can be replaced with the PAN image. After this substitution, the MS image can be back-transformed in the original domain. Clearly, the less the PAN is correlated with the replaced component, the more distortion is introduced. The most famous techniques of this class are intensity-hue-saturation (IHS) [7, 8], in which the images are represented in the IHS domain, principal component analysis (PCA) [9, 10] and Gram-Schmidt (GS) spectral sharpening [11]. On the one side, those techniques preserve the PAN spatial information. On the other side they can produce a high spectral

distortion. This is because PAN and MS are obtained in spectral ranges that only partially overlap.

The second class of techniques, MRA, are based on the introduction of spatial details, extracted from the PAN image through a multigiscale decomposition, into the up-sampled version of the MS. This approach promises a better spectral fidelity but often present spatial distortions.

The lack of the reference image is the principal issue in the evaluation of the pansharpening methods. When a couple of images are fused, the result cannot be compared with anything else. The sensors used for the acquisition cannot reach alone both spatial and spectral resolution of the result. As the two models are different, the result cannot be compared with another image acquired with a different sensor. For this reason, there is no universal measure of quality for the pansharpening. The scientific community common practice is to use the verification criteria that were proposed in the most credited work [12]. This study defines two properties to use for the evaluation of the fused product: consistency and synthesis. The first means that the original MS image should be obtained with a degradation of the fused result. The second property describe that the fused image should preserve both the features of each band and the mutual relations among them. The definitions of an algorithm that accomplishes these properties and of an index that can guarantee the correct evaluation are an open problem. But, no matter what index is decided to use, the unavailability of a reference image is a huge problem and a visual inspection is always mandatory.

Accordingly, there are two techniques that can be used for the quality assessment. The first is to degrade both the images given in input to the pansharpening algorithm and use the original MS image as a reference for the result evaluation. The downside of this method is the assumption of invariance between scales, which justifies that the same algorithm operates similarly at reduced scale. The cited hypothesis is not always verified as documented in [9, 12]. A second technique is the use of an index that does not require a reference image.

1.1 CS

The CS family is based on converting the MS image into a domain in which the spatial and spectral pieces of information can be better separated. In this domain, the component containing the spatial information can be replaced by the PAN image. The greater the correlation between the PAN image and the replaced component, the lower the distortion introduced by the fusion. For this reason, the histogram matching of the PAN with the component that contains the spatial part of the MS information is preliminarily performed. After the substitution, the data can be represented in the original space with an inverse transformation. This approach is applied to the whole image in the same way. Techniques of this category have high fidelity regarding the fusion of spatial details and are fast and easy to implement. But as the acquisition spectrum of the sensors used to produce the PAN and MS image differ each other, the

process may produce significant spectral distortions [13, 14]. In the studies [15, 16, 17, 18, 19], it was shown that, when a linear transformation is used, the substitution and fusion can be obtained without the explicit forward and backward transformation of the images but with a precise injection scheme. This scheme can be formalized according to the following equation:

$$\widehat{MS}_k = \widetilde{MS}_k + g_k(P - I_L), \quad k = 1, \dots, N \quad (1.1)$$

in which k indexes the spectral bands, g_k are the injection gains, \widehat{MS}_k is the k -th band of the pansharpened image, \widetilde{MS}_k is the k -th band of the MS image interpolated to the PAN scale and I_L is the intensity component derived from the MS image according to the relation:

$$I_L = \sum_{i=1}^N w_i \widetilde{MS}_i \quad (1.2)$$

The weight vector $w = [w_1, w_2, \dots, w_k]$ is the first row of the forward transformation matrix and depends on the spectral overlap among MS channels and PAN.

The CS approach procedure is illustrated in Fig. 1.1. Four important steps can be noticed: 1) interpolation of MS image for matching the PAN scale; 2) calculation of I_L using Eq. (1.2); 3) histogram matching between PAN and intensity component; 4) details injection according to Eq. (1.1).

The various CS techniques such as IHS [7, 8], PCA [10, 1] and GS [11, 15] define different $w_{k,i}$ and g_k .

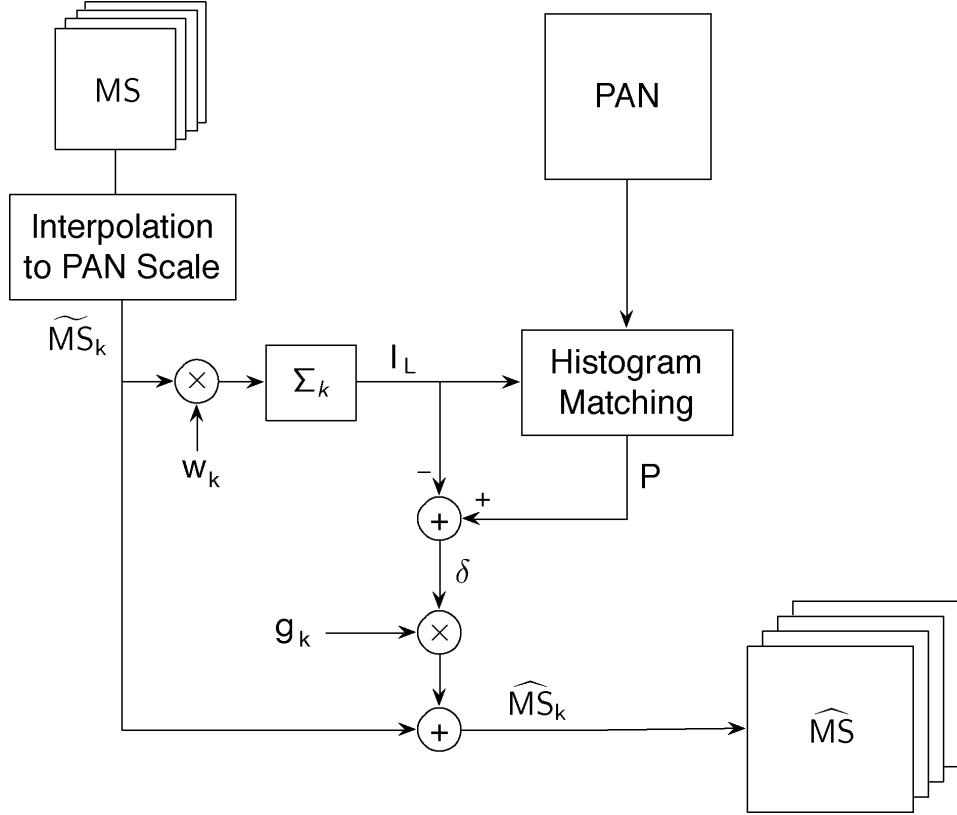


Figure 1.1: Flowchart of CS approach [20].

In the IHS pansharpening method the IHS transformation is used. This is the major limitation of this technique because it can transform only RGB images and often the MS image has 4 or also 8 and more bands. As a workaround, the authors of paper [18] has proved that GIHS, a generalization of the IHS transformation for more bands, can be formulated for any arbitrary set of nonnegative spectral weights as described in the following equation:

$$\widehat{MS}_k = \widetilde{MS}_k + \left(\sum_{i=1}^N w_i \right)^{-1} (P - I_L), \quad k = 1, \dots, N \quad (1.3)$$

in which w_i are all equal to $1/N$ [7].

With the injection gains defined such that:

$$g_k = \frac{\widetilde{MS}_k}{I_L}, \quad k = 1, \dots, N \quad (1.4)$$

\widehat{MS}_k can be calculated as

$$\widehat{MS}_k = \widetilde{MS}_k \cdot \frac{P}{I_L} \quad (1.5)$$

which is the well-known Brovey Transform.

In the PCA pansharpening method, it is used the PCA transformation, also called Karhunen-Loeve transform. It is a linear transformation that can be implemented for a multidimensional image, so it is not limited as the IHS method, and consists into the projection of all the components along the eigenvectors of the covariance matrix. This means that each component is orthogonal and statistically uncorrelated from the others. The hypothesis introduced in this step is that the spatial information is concentrated in the first component, the component with the larger eigenvalue. The PCA can be implemented by using Eq. (1.1), in which w is the first row of the forward transformation matrix; g is the first column of the backward transformation matrix.

The GS transformation is a common technique used to orthogonalize a set of vectors in linear algebra. First of all, the \widetilde{MS} bands are organized in vectors to obtain a two dimensional matrix in which the columns are constituted by the bands organized as vectors. The mean of each band is subtracted from all the

columns. The orthogonalization procedure is used to create a low-resolution version of the PAN image, i.e., I_L . The last step is the replacement of I_L with the histogram matched PAN before the inverse transformation. GS is a generalization of PCA in which PC1 may be any component and the remaining ones are calculated to be orthogonal with PC1. Also the GS procedure can be described by Eq. (1.1) if g_k is defined as:

$$g_k = \frac{\text{cov}(\widetilde{MS}_k, I_L)}{\text{var}(I_L)}, \quad k = 1, \dots, N \quad (1.6)$$

in which $\text{cov}(\cdot, \cdot)$ is the covariance between two images and $\text{var}(\cdot)$ is the variance. There are several versions of this technique that differ on how the I_L is created. The simplest way is to set $w_i = 1/N$. This version is called GS mode 1 [11]. In [15], also an *adaptive* version of this mode was proposed. It is called GSA and the I_L is generated by a weighted average of the MS bands. Another technique defined in [11] and called GS mode 2 suggests to generate the I_L by applying a low pass filter to the PAN image. This last step leads the GS mode 2 that belongs to the MRA class of techniques.

Another noteworthy technique is described in [19] and introduces the concept of *partial replacement* of the intensity component. An intensity image is created for every band of the MS from the PAN image; it is calculated with the following equation:

$$P^{(k)} = CC(I_L, \widetilde{MS}_k) \cdot P + (1 - CC(I_L, \widetilde{MS}_k)) \cdot \widetilde{MS}'_k \quad (1.7)$$

in which \widetilde{MS}'_k is the k -th MS band histogram-matched to PAN and CC is the correlation coefficient. I_L is defined using in Eq. (1.2) a vector w obtained with a linear regression of \widetilde{MS}'_k on P_L , the degraded version of the PAN. The injection gains are the result of:

$$g_k = \beta \cdot CC(P_L^{(k)}, \widetilde{MS}_k) \cdot \frac{std(\widetilde{MS}_k)}{\frac{1}{N} \sum_{i=1}^N std(\widetilde{MS}_i)} L_k \quad (1.8)$$

β is empirically tuned and is a factor that normalizes the high frequencies. $P_L^{(k)}$ is a low-pass-filtered version of $P^{(k)}$, and L_k is an adaptive factor that removes the local spectral instability error between the synthetic component image and the MS band defined as:

$$L_k = 1 - \left| 1 - CC(I_L, \widetilde{MS}_k) \frac{\widetilde{MS}_k}{P_L^{(k)}} \right|. \quad (1.9)$$

1.2 MRA

In the MRA class of techniques, the pansharpened image is defined as:

$$\widehat{MS}_k = \widetilde{MS}_k + g_k(P - P_L), \quad k = 1, \dots, N. \quad (1.10)$$

$P - P_L$ is the operation performed to obtain the high-frequency details of the PAN image. The algorithm to create the P_L and the chosen g_k weights differentiate the MRA pansharpening techniques of this class.

However, in general, all the techniques follow the algorithm described in

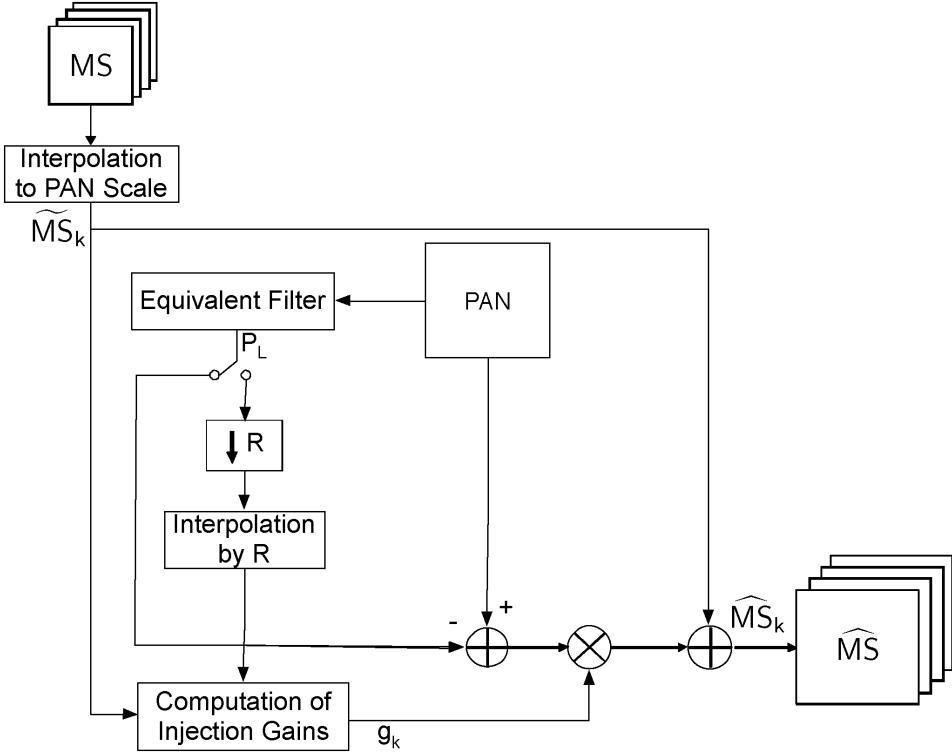


Figure 1.2: Flowchart of MRA approach [20]

Fig. 1.2. First of all, the MS image is interpolated to the PAN scale. The second step is to calculate P_L , the low pass version of PAN obtained by means of an equivalent filter. The vector of injection weights g_k can be computed using the \widetilde{MS}_k in combination with P_L . Interpolation is less crucial in MRA respect to CS methods. A method to produce the P_L image consists in applying a low pass filter h_{LP} to the PAN image P . So Eq. (1.2) can be rewritten as:

$$\widehat{MS}_k = \widetilde{MS}_k + g_k(P - P * h_{LP}), \quad k = 1, \dots, N \quad (1.11)$$

where $*$ is the convolution operation. A more general method to obtain the

P_L is called *Pyramidal Decompositions* and the number of filterings can be one or more. A filter type that proves to be a good choice is a Gaussian filter that closely matches the sensor MTF. A noteworthy option is the MTF-GLP with a context-based decision (MTF-GLP-CBD) [21] where the injection gains are defined as follows:

$$g_k = \frac{\text{cov}(\widehat{MS}_k, P_L^{(k)})}{\text{var}(P_L^{(k)})} \quad (1.12)$$

It is context-based because it can be applied on nonoverlapping image patches to improve the quality of the final product.

1.3 Quality Assessment

As explained above, the lack of a reference image is the main limitation. The community has proposed two assessment procedures as a workaround. The first procedure consists in using the images at a lower spatial resolution and keeping the original MS image as a reference. However, the output of an algorithm can have different performance at different scales, as it is showed in [22]. This happens because the performance assessment depends intrinsically on the image scale, mostly in case of pansharpening methods that apply spatial filters. The second procedure consists in using non-reference quality indexes. Both types of procedures require also a visual inspection for spectral distortions and spatial details.

The Wald's protocol is composed by three requirements:

1. \widehat{MS}_k degraded to the original MS scale should be as identical as possible

to the MS_k .

2. The fused image \widehat{MS}_k should be as identical as possible to the MS_k that the sensor would acquire at the highest resolution
3. The MS set of synthetic images $\widehat{MS} = \{\widehat{MS}\}_{k=1,\dots,N}$ should be as identical as possible to the MS set of images $HRMS = \{HRMS\}_{k=1,\dots,N}$ that the corresponding sensor would observe at the highest resolution.

In the previous definitions, HRMS is the reference image. For a reduced-resolution assessment, the filter choice for the downsampling is crucial in the validation. A bad filter choice results in an image degradation that does not reflect the sensor characteristics at a lower scale. So the algorithm, after the degradation, is applied to images that reflect the wrong sensor model. This means that the same algorithm can have a more different result at the original and lower scales. On the contrary, with a good filter that preserves the sensor characteristics at the lower resolution, the algorithm has a higher possibility to reflect the quality of the original resolution. Indeed, the filter used for the MS degradation should simulate the transfer function of the remote sensor and so, it should match the sensor's MTF [23]. Similarly, the PAN image has to be degraded in order to contain the details that would have been seen if the image were acquired at the reduced resolution. The fused image obtained from the degraded PAN and MS, can be evaluated by means of different indexes using the MS as a reference image.

The *Spectral Angle Mapper* (SAM) is a vector measure that is useful to eval-

uate the spectral distortion. In simple terms, denoting by $I_{(n)} = [I_{1,n}, \dots, I_{N,n}]$ a pixel vector of the MS image, with N bands, the SAM between the corresponding pixel vectors of two images is defined as:

$$SAM(I_i, J_i) = \arccos \left(\frac{\langle I_i, J_i \rangle}{\|I_i\| \|J_i\|} \right) \quad (1.13)$$

$\langle I_i, J_i \rangle$ is the scalar product and $\|I_i\|$ is the vector l_2 -norm. Applying this equation to every pixel results in a so-called SAM map. Averaging all the pixel of the SAM map returns the SAM index for the whole image. The optimal value of the SAM index is 0.

The *Root Mean Square Error* (RMSE) is used to calculate the spatial/radiometric distortions. It is defined as:

$$RMSE(I, J) = \sqrt{E[(I - J)^2]} \quad (1.14)$$

where $E[\cdot]$ is the arithmetic mean. The ideal value of RMSE is zero and is achieved if and only if $I = J$. But it is not an efficient index because it is not considered the error for each band, but is global. So, to better measure the error for each band, the *Erreur Relative Globale Adimensionnelle de Synthèse* (ERGAS) is used. The ERGAS index evaluates the RMSE error with a different weight for each band.

$$ERGAS = \frac{100}{R} \sqrt{\frac{1}{N} \sum_{k=1}^N \left(\frac{RMSE(I_k, J_k)}{\mu(I_k)} \right)^2} \quad (1.15)$$

Obviously, the ERGAS is composed of a sum of RMSEs, so the optimal value is also 0.

Another important index is the *Universal Image Quality Index* (UIQI) or also called Q-index, proposed in [24]. Its expression is:

$$Q(I, J) = \frac{\sigma_I J}{\sigma_I \sigma_J} \frac{2\bar{I}\bar{J}}{\bar{I}^2 + \bar{J}^2} \frac{2\sigma_I \sigma_J}{(\sigma_I^2 + \sigma_J^2)} \quad (1.16)$$

where σ_{IJ} is the covariance of I and J , and \bar{I} is the mean of I . The first fraction represents an estimation of the covariance, the second is a difference in the mean luminance and the third is the difference in the mean contrast. The Q-index varies in the range $[-1, 1]$ with 1 as the optimal value.

Q4 is an extension of the UIQI for images with 4 bands [25]. Let a, b, c and d denote the radiance values of the given image pixel in four bands, and construct the quaternions:

$$z_A = a_A + ib_A + jc_A + kd_A \quad (1.17)$$

$$z_B = a_B + ib_B + jc_B + kd_B \quad (1.18)$$

The Q4 is defined as :

$$Q4 = \frac{4|\sigma_{z_A z_B}| \cdot |\bar{z}_A| \cdot |\bar{z}_B|}{(\sigma_{z_A}^2 + \sigma_{z_B}^2 (|z_A|^2 + |z_B|^2))} \quad (1.19)$$

where \bar{z} is the complex conjugate. If eventually, the bands are more than 4,

the Q4 can be replaced with Q average.

An index used for the validation at full-resolution is the *Quality with no reference* (QNR) index [26]. It is defined by the following equation:

$$QNR = (1 - D_\lambda)^\alpha (1 - D_S)^\beta \quad (1.20)$$

wherein α and β are two coefficients which can be tuned to weight more the spectral or the spatial distortion, respectively.

The maximum theoretical value of the index is 1 and is reached when D_λ and D_S are 0. The spectral distortion is calculated with D_λ using this equation:

$$D_\lambda = \sqrt[p]{\frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j=1, j \neq i}^N |d_{i,j}(MS, \widehat{MS}|^p)} \quad (1.21)$$

where $d_{i,j}(MS, \widehat{MS}) = Q(MS_i, MS_j) - Q(\widehat{MS}_i, \widehat{MS}_j)$, \widehat{MS} is the fused image and p is a parameter typically set to one [26]. The objective is to create an image with the same spectral features of the original MS image.

The spatial distortion is calculated by:

$$D_S = \sqrt[q]{\frac{1}{N} \sum_{i=1}^N |Q(\widehat{MS}_i, P) - Q(MS_i, P_{LP})|^q} \quad (1.22)$$

where P_{LP} is the low-resolution PAN image at the same scale of the MS image and q is usually set to one [26].

Khan protocol [27] extends the consistency property of Wald's protocol.

The pansharpened image is considered as a sum of a lowpass term plus a high pass term. The lowpass term is the original interpolated low resolution MS image and the highpass term corresponds to the spatial details extracted from the PAN and injected into the MS image. A Gaussian model of the sensor's MTF is used to build the filters. The similarity between the lowpass component and the original MS image can be calculated using the Q4 index or any other similarity measure for images with more bands. The similarity between PAN and the spatial component is measured as the average of UIQIs calculated using the PAN and each band of MS. The same similarity is calculated also between the original MS and the degraded version of the PAN.

The QNR and the spectral distortion of Khan's protocol can be combined to yield another quality index, the *Hybrid QNR* (HQNR) [28], defined as

$$HQNR = (1 - D_{\lambda}^{(K)})^{\alpha} (1 - D_s)^{\beta} \quad (1.23)$$

in which α and β are coefficients tuned to weight more the spectral or the spatial distortion, as in Eq. (1.20) and $D_{\lambda}^{(K)}$ is given by

$$D_{\lambda}^{(K)} = 1 - Q4(\widehat{M}_L, M) \quad (1.24)$$

\widehat{M}_L indicates the fused image degraded to the resolution of the original MS image.

Chapter 2

Pansharpening applications of Deep Learning

Early works in the field of *Deep Learning* (DL) have been made in the 1940s starting from the concept of *Perceptron* [29] and afterwards in the 60s with the invention of backpropagation, the most commonly used algorithm in the present day to train a *Neural Network*. The Neural Network is a model constituted by several neurons, that are eventually organized into different layers and have the ability to learn, extract and distinguish different features from the data given into input. In order to obtain these capabilities, the network should be subjected to a training phase that requires a huge amount of computational power. As in the early 60s the computers were no powerful enough and there was not a great data availability, DL had no reason to be implemented as today.

This approach is characterized by different phases. One of them, the training, is a critical phase in which the model learns from new data and can differ for the type of issues that the model should solve. In most of the cases, the model gives a prediction and the result is compared with a target. Initially, the error is calculated between the output of the prediction and the target. After that, when the error is propagated in all the neurons of the model, the weights of all the neurons are modified so that the next prediction for the same input will be much similar to the target output. How this error is calculated and what means "similar" is established by the loss function that calculates the error. The propagation uses a method called *Gradient Descent* algorithm, that allows the network to understand how to change the weights and minimize the loss. In order to use the gradient descent, the loss function must be differentiable so that the gradient operator can be applied more times. Many layers the net have, many times the algorithm apply the gradient to the error.

LeCun (1989) for the first time used a backpropagation algorithm to train a neural network to classify handwritten digits.

Nowadays, the data collected through internet, the incredible amount of computational power exhibited by data centers and GPUs, the performance of this type of algorithm and a large number of applicative fields, have encouraged the growth of Machine Learning (ML) and in particular, DL.

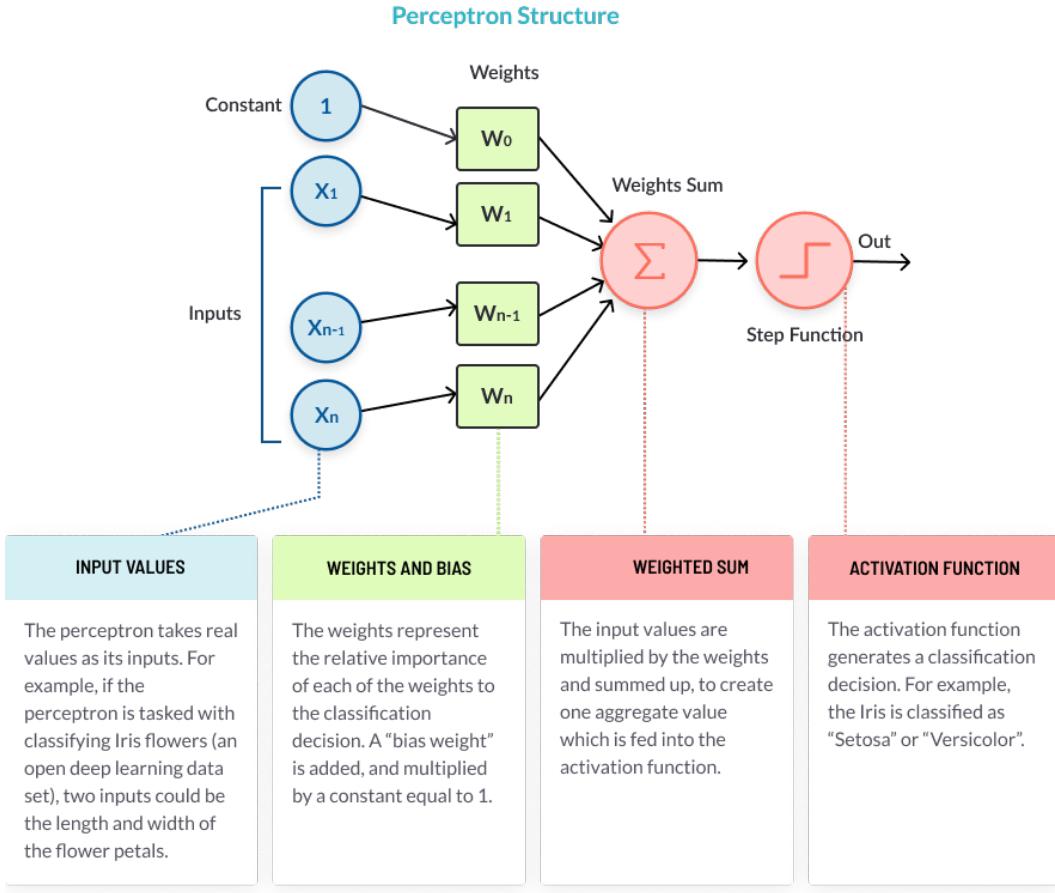


Figure 2.1: Perceptron in detail [30]

2.1 Neural Networks

2.1.1 Perceptron

Essentially a perceptron is an element in which the output is a result of a weighted sum of the inputs. There is always only one output but it can be more inputs as depicted in Fig. 2.1. A perceptron uses only one linear or non-linear activation function. The following paragraph will discuss what is an activation function and how it works. With a non-linear transfer function, perceptron

can build a nonlinear plane separating data points of different classes. In 1969, it was proved that the perceptron itself may fail in certain simple tasks, as, for instance, the separation of a plane described by the XOR function [31]. However, 3 perceptron organized in 2 layers can separate an XOR plane. This opened up the development of multi-layer models and subsequently, for training optimization for specific applications, the creation of different types of layers.

2.1.2 Activation Function

The activation function is used to transform the weighted data (input multiplied by weights) in outputs. By the use of a non-linear transformation, we create new relation between the points. This consents to the ML model to create increasingly complex features with every layer. Features of many layers that use pure linear transformations can be reproduced by a single layer that uses a non-linear function. The most common activation function is the so called ReLu that implements the relation $y = \max(0, x)$, where x is the input and y the output. The gradient is 1 when the value of x is positive, and 0 when negative. This means that during the training, negative gradients will not update the weights. A gradient equal to 1 means that the training will be much faster compared with other activation functions like logistic sigmoid.

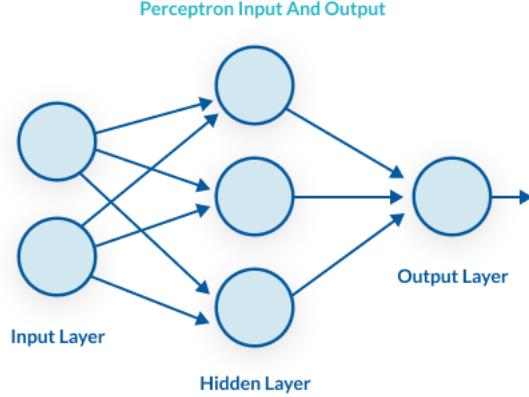


Figure 2.2: Multilayer architecture example [30]

2.1.3 Layers

The DL network can have different type of layers that differ in how the perceptrons are organized and how they modify their weights during the training. Fig. 2.2 shows an example of a multilayer model. The first experimented layer was the dense layer where all neurons are connected with all neurons of the next layer. Every neuron has his own weights and considering all the neurons connections, this type of layer generates promptly a large number of weights. This requires a large dataset and a long time for the training.

The convolutional layer gave an important improvement in computer vision applications. This layer can apply different filters to the data at the same time.

Commonly, several convolutional layers at the start of a neural network are used to extract features from an image. In a first phase, the DL model is able to extract features in the data with complexity that grows with the model dimension, as the Fig. 2.3 shows. In a second phase, the features can be

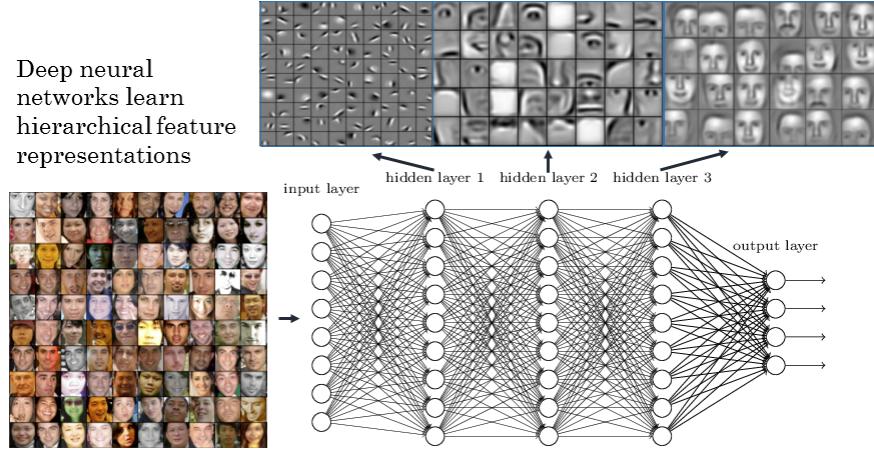


Figure 2.3: Feature extraction of a deep model [34]

analyzed by different layers, for instance dense layers can elaborate features for classification or for other tasks. The filter weights of the convolutional layer are learned in the training and applied to the whole input where the neurons share the same weights. For this reason, respect to a dense layer, the convolutional layer has a lower weights count, is much faster to train and allocates a minor amount of memory.

To recognize patterns among a sequence of data such as video or audio, *Long short-term memory* (LSTM) [32] and *Gated Recurrent Unit* (GRU) [33] are used and, unlike standard layers, have feedback connections that allow the model to memorize data for subsequent inputs. Other layers like the pooling layer, max-pooling layer and dropout were created to optimize the training.

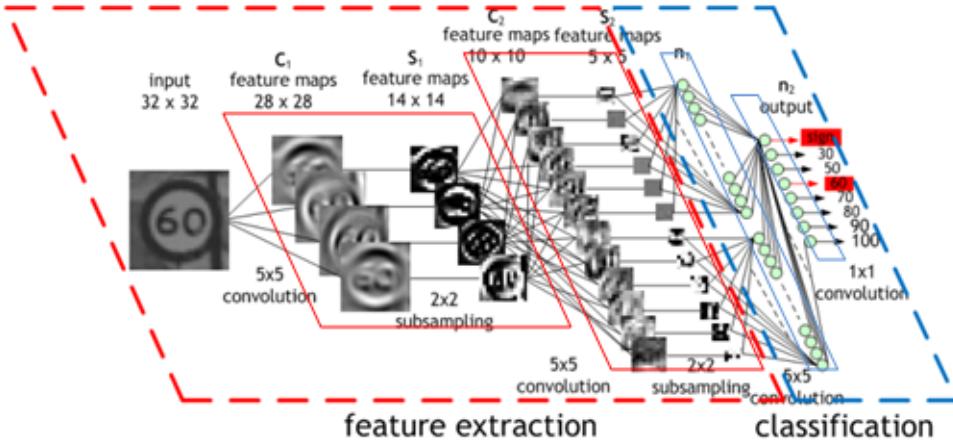


Figure 2.4: Image by Maurice Peemen

2.2 The Deep Learning Paradigm

Deep Learning [35] is a subfield of ML focalized in the study of deep neural networks. The model architecture is made in such a way to extract features from the data and recognize patterns from them. This is the main difference between the DL and other ML techniques. Other ML techniques are focused on learning from handcrafter features. The extraction process of this features should be tuned for the specific data structure and requires a high knowledge of the particular context.

In DL, the net is trained also for feature extraction, but it requires a more complex architecture. Fig. 2.4 illustrates a structure of a DL model.

Moreover, it gave a strong impulse to the development of very efficient algorithms in the computer vision field. Indeed, nowadays there are a lot of tools that simplify the use of complex models trained to recognize hundreds

of objects in an image.

The so called fine-tuning technique allows net to recognize a different set of objects. The purpose of fine-tuning is to reset the weights of final layers in a way that the model learned on some given features, can accomplish another task. With this technique it will be just necessary to train final layers instead of the entire model for days and using large datasets.

2.3 Pansharpening Applications

Recently, a new pansharpening method based on a convolutional neural network has been proposed [4]. To accomplish that, it has been specialized a network built for super-resolution [36]. From that model, other three different models have been tuned for GeoEye1, IKONOS and WorldView-2 sensors with the aim of predicting the sensor characteristics and giving better performance.

An important issue in this field is that it is difficult to found good images because of the high cost for high-resolution images. For this reason, it cannot be possible to train deep network. The choice of the authors was to use three convolutional layers, illustrated in Fig. 2.5. The advantage of using only convolutional layers is that the input can have different sizes.

The training phase showed in Fig. 2.6 uses a reference approach in which the images are downsampled and fused. After the downgrande, the P_L image is attached to the downgraded MS in order to provide input to the model. The MSE error between the original MS and the network output are calculated and

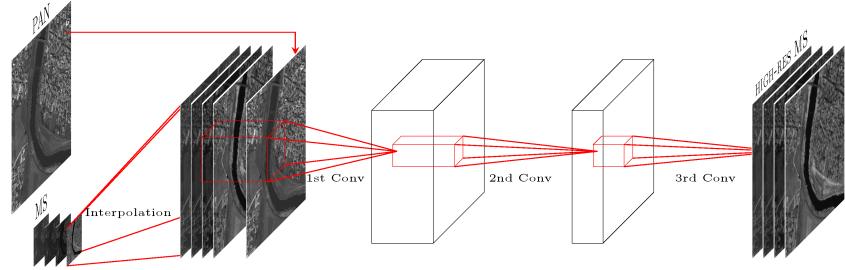


Figure 2.5: CNN architecture for pansharpening[4]

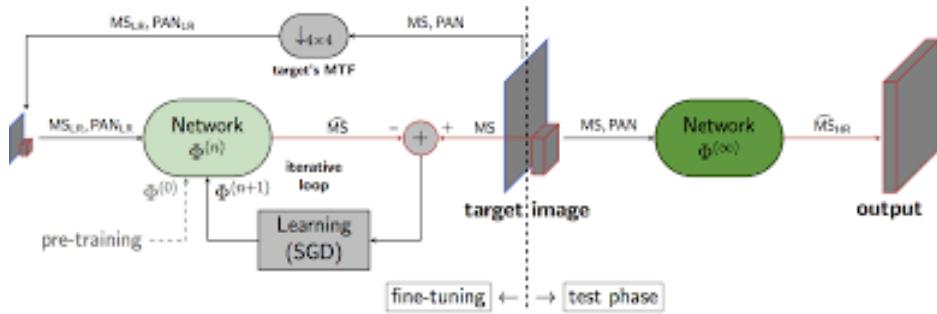


Figure 2.6: Training and Test[37]

this error is used for the training of the weights.

Radiometric relevant indexes has been used to improve the results [4]. It was showed that the network avoids learning the provided indexes. Another important step introduced in [37] was to run a fast session of fine-tuning before the application of the model. With the fine-tuning, the net can specialize the weights for the specific environment described by the downgraded image and, subsequently, complete the fusion on the original images, with better results. However, as described in the conclusions of [37], performance with this approach is not as much relevant as *no-reference* approach that have a major impact to performance. Another important disadvantage of the method is that when the model works with misaligned images, it is trained with images

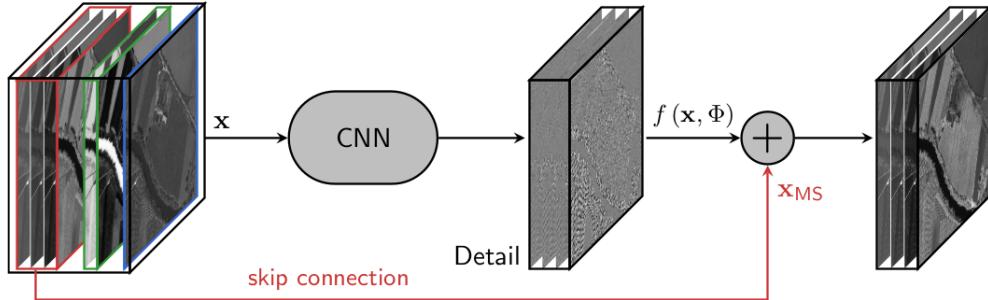


Figure 2.7: Residual-based version [38]

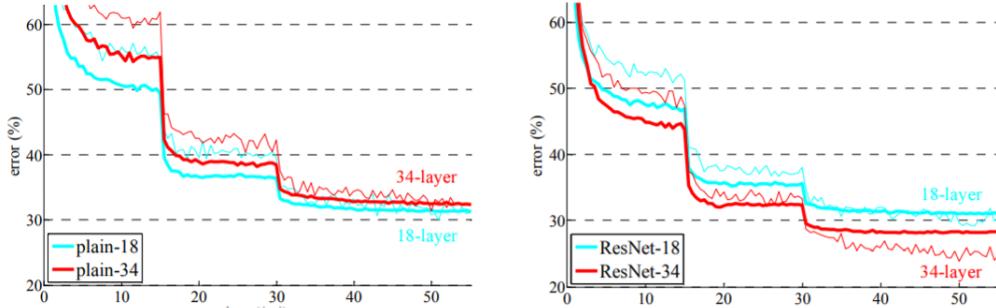


Figure 2.8: Performance comparison between residual (right) and non residual (left) networks [39]

that have a fraction of the misalignment, depending on the scale ratio between PAN and MS.

In another paper [38], the authors has tried different structures and strategies to improve the network released in [4]. An important result was achieved with a residual version of the original network. The structure is showed in Fig. 2.7.

The residual version was not trained to reproduce the whole image, but just the high-pass component. Indeed, the low-pass component is represented

by the multispectral image provided by input. This means that the network reconstructs just the missing parts.

Residual-based structures were used in other papers [40, 41], in particular in the DL in [42, 39].

In conclusion, it was observed that it is easier to train a neural network for differences instead of reproducing an ouput really similar to the input.

The main results of [39] are shown in Fig. 2.8.

Chapter 3

Proposed Solution

The aim of this chapter is to discuss and evaluate the difference between the training of the reduced resolution approach proposed in [4] and the training of the *no-reference* approach developed in the research. The analysis will cover the *automatic differentiation* and the implementation in Tensorflow.

3.1 Loss Function Issue

To assess the training set for the *reduced resolution* (RR) method described in [4], the algorithm synthesized in the Fig. 3.1 has been implemented.

Both PAN and MS are downsampled into two images called MS_L and PAN_L . The MS_L is upsampled to match the PAN_L size. The images represent the input of the model that compute the error by the use of MS. To update the model weights, a learning rate scheduler is used: the *Stochastic Gradient Descent* (SGD) [43]. A learning rate scheduler consent to modify the learning

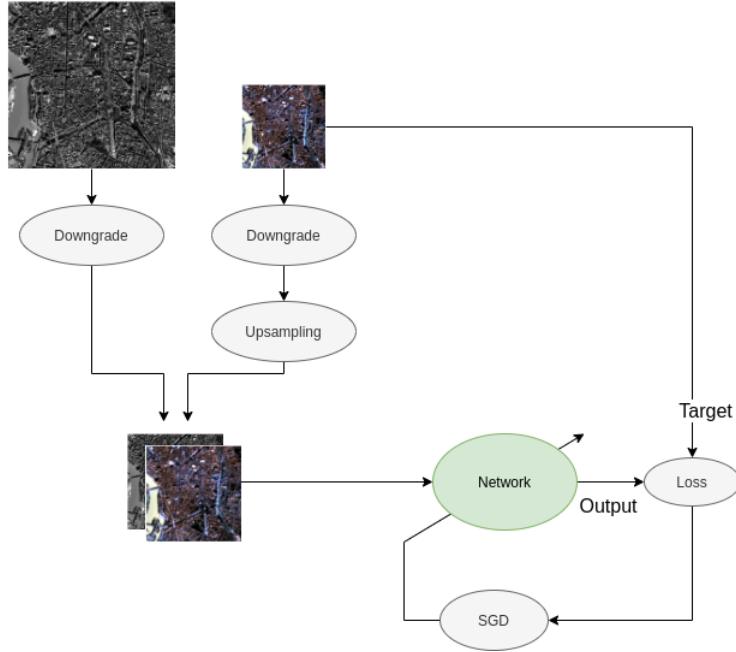


Figure 3.1: Graph of the training algorithm released by [4].

rate during the training in order to optimize performance and reach easier the minimum error. Other details about learning rate and learning rate scheduler can be find in the paragraph 4.1

Instead of this procedure, we decided to use the *no-reference* approach (NOREF) illustrated in Fig. 3.2. As the figure shows, the images are not downsampled. The great advantage is that the model is trained with original images and not with the downgraded version. Also in this case, the MS is upsampled to match the PAN size. Without the upsampling, it will not be possible to stack the images and apply the 3D filters of the convolutional layers.

In this case, the loss function does not use a target. This is because the model is trained with the images at the highest possible resolution. Indeed,

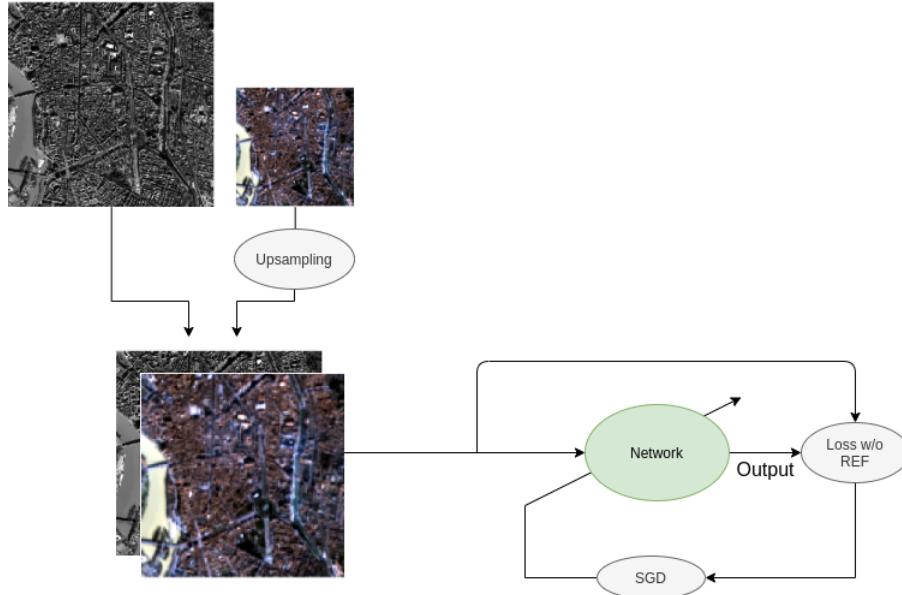


Figure 3.2: Graph of the training algorithm developed.

the loss computation is done using a *no-reference* index. For this reason, this approach belongs into the *unsupervised learning* field. Oppositely to *supervised learning*, unsupervised learning is a subfield of ML where the training process is not leaded by a target.

3.2 Automatic Differentiation

Automatic differentiation is a set of techniques to efficiently calculate the gradient of a function by a computer. Every computer program executes a sequence of elementary operations. There are different kinds of techniques to differentiate a function: *automatic differentiation*, *numerical differentiation* and *symbolic differentiation*. The Fig. 3.3 shows an example of all techniques

that evaluate the derivative of the same function. The various approach are described in the following.

Numerical differentiation is the finite approximation of the derivatives.

$$\frac{df(x)}{dx} \approx \frac{f(x+h) - f(x)}{h} \quad (3.1)$$

In this technique, the computer evaluates the function described in Eq. (3.1) at some sample points. As described in [44], this type of differentiation is easy to code. However, the $O(n)$ cost of the evaluation in n dimensions is so high that ML, where n can be millions or also billions, would not be accessible.

Symbolic differentiation is a technique that uses the chain rule, the product rule and other rules to split the expression in known primitive derivative to obtain the result. Eq. (3.2) shows how product rule and sum rule are applied:

$$\begin{aligned} \frac{d}{dx}(f(x) + g(x)) &= \frac{d}{dx}f(x) + \frac{d}{dx}g(x) \\ \frac{d}{dx}(f(x)g(x)) &= \left(\frac{d}{dx}f(x)\right)g(x) + f(x)\left(\frac{d}{dx}g(x)\right) \end{aligned} \quad (3.2)$$

Techniques that belong to this class are really difficult to convert into a computer programs.

Automatic differentiation systems explicitly split the operations in simpler ones and build a computation graph. Each node of the graph have some attributes like values, primitive operations and parents.

Jacobian matrixes such as

$$J_f = \begin{bmatrix} \frac{dy_1}{dx_1} & \cdots & \frac{dy_1}{dx_n} \\ \vdots & \ddots & \\ \frac{dy_m}{dx_1} & & \frac{dy_m}{dx_n} \end{bmatrix} \quad (3.3)$$

are used to represent the derivative of each output y respect to each input x of a graph node.

For each primitive operation, it must be defined a *Vector-Jacobian Product* (VJP). Combining all the VJP nodes, the result would be the value of the gradient.

Automatic differentiation provide different modes like *forward mode* and *reverse mode*. In forward mode, *automatic differentiation* and *symbolic differentiation* are equivalent, as described in the paper [45]. They both apply the chain rule and other differentiation rules and actually create expression graphs. However, the first one was created specifically for the computer manipulation and includes numerical values. Indeed, the *automatic differentiation* can handle also control flow statements like "if", "while" and "loops". The second one, *symbolic differentiation*, operates on mathematical expressions and symbols.

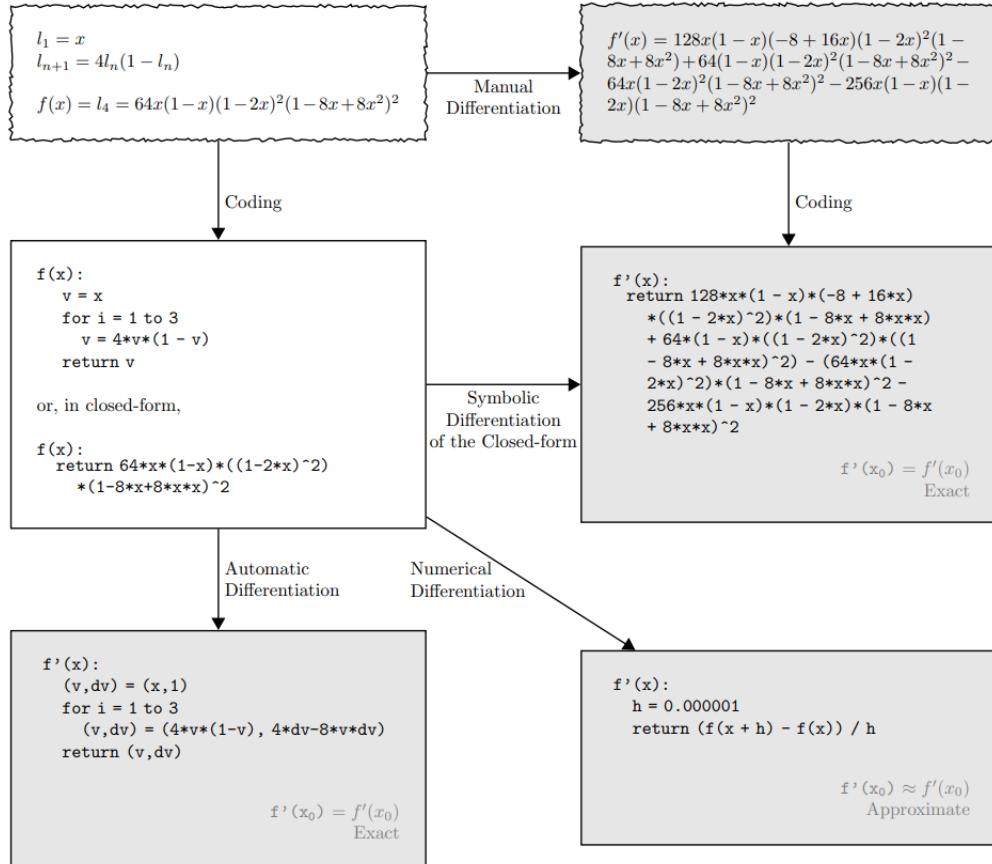


Figure 3.3: Example with all type of differentiation [44]

3.3 Tensorflow and custom loss function implementation

Tensorflow is the most famous ML and DL tool that implements the *automatic differentiation*. To build a custom loss function, Tensorflow provide APIs. These APIs define all the operations between *Tensors* that are the basic datatype in Tensorflow. The Tensorflow APIs used in the research code are showed in the Table 3.1. The external package *tensorflow probability* APIs

used are showed in Table 3.2.

After the definition of the custom loss function, at runtime Tensorflow translate the Python code in C/C++ for efficiency purpose, compile it, and builds the computation graphs. Also, Tensorflow for a couple of years, integrate on it Keras, another framework that simplifie the creation of the model with all common layers well-defined into **Classes**. To create a model it is necessary to create a **Model** class. It is possible to add a layer using `model.add(Layer(args))`. **Layer** is the corresponding class of the layer and args are the arguments like the input shape and the number of filters the layer should have. Firstly, using such a widespread framework has the advantage of a large community that supports challenging situations; secondly, it will be a larger compatibility on Operating Systems, GPUs and hardware in general.

Table 3.1: Table of tensorflow API used in the thesis.

<code>abs(...);</code>	Computes the absolute value of a tensor.
<code>add(...);</code>	Returns $x + y$ element-wise.
<code>cast(...)</code> :	Casts a tensor to a new type.
<code>constant(...)</code> :	Creates a constant tensor from a tensor-like object.
<code>divide(...);</code>	Computes Python style division of x by y .
<code>expand_dims(...)</code> :	Returns a tensor with an additional dimension inserted at index axis.
<code>multiply(...);</code>	Returns an element-wise $x * y$.
<code>ones(...)</code> :	Creates a tensor with all elements set to one.

<i>pad(...)</i> :	Pads a tensor.
<i>reduce_mean(...)</i> ;	Computes the mean of elements across dimensions of a tensor.
<i>reduce_std(...)</i> ;	Computes the standard deviation of elements across dimensions of a tensor.
<i>reduce_sum(...)</i> ;	Computes the sum of elements across dimensions of a tensor.
<i>sqrt(...)</i> ;	Computes element-wise square root of the input tensor.
<i>square(...)</i> ;	Computes square of x element-wise.
<i>squeeze(...)</i> :	Removes dimensions of size 1 from the shape of a tensor.

Table 3.2: Table of tensorflow probability API used in the thesis

<i>stats.covariance(...)</i> ;	Sample covariance between observations indexed by event_axis.

3.3.1 Loss Function

At the beginning, to build the loss function, it was used a QNR approximation

$$f(x) = 1 - \widetilde{QNR} \quad (3.4)$$

where \widetilde{QNR} is the approximated QNR.

The D_s and D_λ were calculated using the Qavg, that is the average of Q indexes computed for each band. The Q index was calculated considering the whole image and not dividing it into small patches. The result shows a series of negative values instead of values between 0 and 1 as the original images format. This is why the performance was unsatisfactory (results are illustrated in the Chapter 4). Indeed, the model minimizes the loss functions (namely, it maximizes the QNR, considering 3.4) no matter the sense of output values. After some tests, it was considered to solve this problem using the D_{sreg} described in [46] instead of D_s and migrate to an approximation of HQNR instead of the QNR. Implementing the exact HQNR was unfeasible due to the difficulty on manipulating quaternions and hypercomplex numbers in TensorFlow in such a way to keep the function differentiable. In order to gain better performance, different versions of D_λ indexes have been implemented, according to the diverse approaches proposed by the scientific community.

Chapter 4

Experimental results

4.1 Description

The implemented code has been written in Python 3.7 with the use of a file `main.py` as a entrypoint. It is compatible with Linux, Windows and MacOS and it can run on dedicated GPU as well as on CPU.

By the use of the standard `argparse` library, all the possible parameters can be codified on terminal. This allows to run, with a bash file, all the experiments in series without changing the code.

At the beginning, the QNR function has been developed with approximations. In the Appendix A.1 the most important functions used for the QNR evaluation are presented. As showed in the Appendix A.1, the Q index has been obtained as an average of the Q calculated band-by-band instead of being calculated in patches as the original paper of the QNR described. In order to

obtain the spectral consistency measure D_λ , for each band of MS and fused image, it has been necessary the calculation of Q between bands, as formalized in Eq. (1.22). According to Eq. (1.21), D_s has been calculated as the Q between MS bands and PAN_L and also between the fused image and the PAN. This gives a quality measurements of the details inserted into the image.

Due to the poor performance obtained via the QNR, it has been decided to use an approximation of HQNR. Approximations are due to the lack of Tensorflow hypercomplex API. Instead of using Q4, for calculating D_λ , the Qavg used for the QNR has been adopted. The D_λ is thus calculated as $1 - Q_{avg}$, where Q_{avg} is computed between the fused image filtered through a Gaussian filter and the MS.

Moreover, we employed as spatial distortion index the D_{sreg} [46], an implementation of the D_s that takes advantages from the *coefficient of determination* calculation. The coefficient of determination, also called *R-squared* (R^2), is a statistical measure that represents the proportion of the variance for a dependent variable that is predictable from the independent variables.

The D_{sreg} is obtained as $1 - R^2$, in which the coefficient of determination is calculated between the fused image and the PAN. The implementation is illustrated in the Appendix A.2. Depending on the sensor used for the image acquisition, the filter would be different to match the sensor MTF. The filters were created with the Matlab Toolbox, exported in the .mat format and imported in the project at runtime. It was decided to use no padding after the filtering process and cut the pixels affected by border effects from the MS

during the Q calculation.

An important choice regarded the learning rate. The learning rate is multiplied with the gradient of the error for each layer and it will determines how quickly the error decrease.

If the learning rate is too high, the error will drop down quickly, but difficulties will be encountered in finding minimum. With a too low learning rate, it can take long time to reach the convergence or the algorithm gets stuck in a local minimum. According to the community, the best approach is to use a learning rate scheduler like Adam [47] or SGD [43]. At the start of the training, the model can have a high learning rate so that it can explore all the loss function. After that, the learning rate will be gradually decreased to avoid oscillations. The learning rate is one of the parameters that requires some experiments to be tuned. The best case is to reach the maximum possible quality with a constant increase in performance during the training. The main index of performance was Q2N. At the moment, Q2N is the best index used with reference image. It is important in the training to avoid a bell shape in the index graph. Bell shape indicates that the Q2N increases really fast at the beginning. After reaching the pick value of optimum, it decreases really fast. This characteristic would made the training difficult in real applications. The reason why is that the calculation would be impossible due to the unavailability of the reference image.

Without the bell shape during the training, the model can only have an increase in performance.

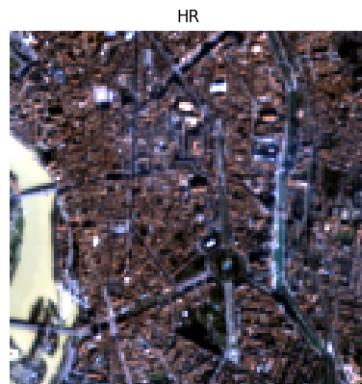


Figure 4.1: Multispectral and Panchromatic Toulouse images used for all the tests

4.2 Experimental Settings

All types of training have been executed in a test image illustrated with a true color representation in Fig. 4.1 and identified as Toulouse dataset. The Toulouse dataset is a couple of IKONOS images acquired over the city of Toulouse, France, in 2000. The MS sensor acquires an image of 512×512 pixels in the visible and near-infrared spectrum range in 4 different bands and with

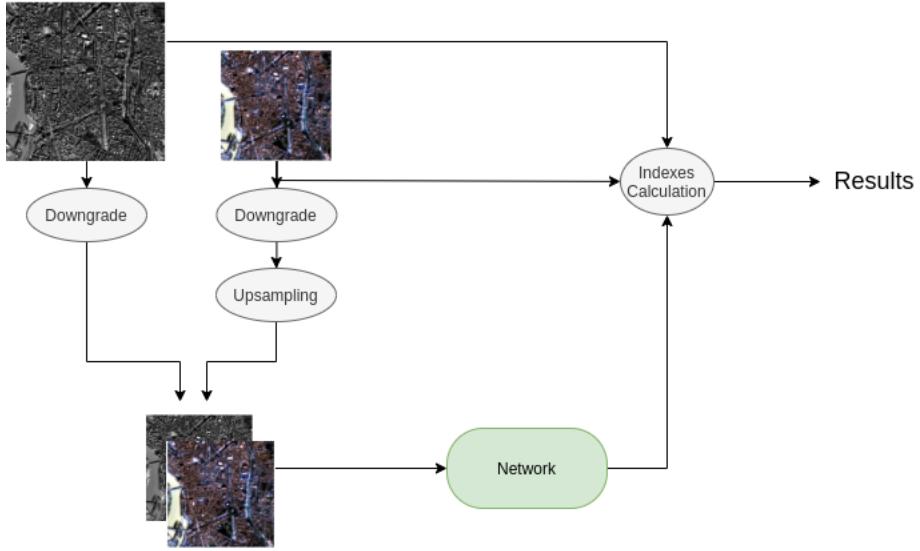


Figure 4.2: Validation process

a 4m spatial sampling interval (SSI). The panchromatic image is constituted by an image of 2048×2048 pixels with a 1m SSI.

DL models should not be specialized too much for the training set otherwise, the model could have low performance for new data not included in the training set. To avoid this situation, it is necessary to create a validation set in which there are different data. Performance on the validation set are more similar to real world performance. The validation process is described in the Fig. 4.2.

To generate a reference image for the quality assessment, MS and PAN has been downsampled. The original MS has been used as Ground Truth (GT).

At every epoch, the downsampled images PAN_L and MS_L were used to produce a fused image with the updated weights and compared with the GT. This means that the images will constitute the validation set. Reference in-

dexes like SAM, ERGAS and Q together with *no-reference* indexes QNR and HQNR were calculated between the output of the model and the GT. Those indexes has been obtained with the MatLab toolbox functions provided by [20] and MatLab engine for Python [48]. The procedure can give a real performance assessment of the whole model. The training set have always better performance as the model optimizes the error of these data. It needs to be highlighted that only the performance of the validation set are relevant.

For the training in RR approach, as described in the Paragraph 3.1, the PAN_L and MS_L are downsampled again as the model would have a different target from GT. Training the model with the GT as target have no valid applications. In real world, the model would not have a GT. Additionally, the validation set would have no relevant performance as the model is trained with the same data.

It has been noticed that, for some images, the model that uses the fine-tuning with the reference approach has worse performance than the original one. The reason is that the model is fine-tuned with the downgraded images.

The characteristics of the sensors, in general, do not satisfy the scale invariance hypothesis and downsampling the input does not guarantee similar performance like using the original images as described also in the previous chapters.

In the NOREF approach, the images are not downsampled again as the model use a *no-reference* loss function.

To compare the different methods, the model has been also trained with

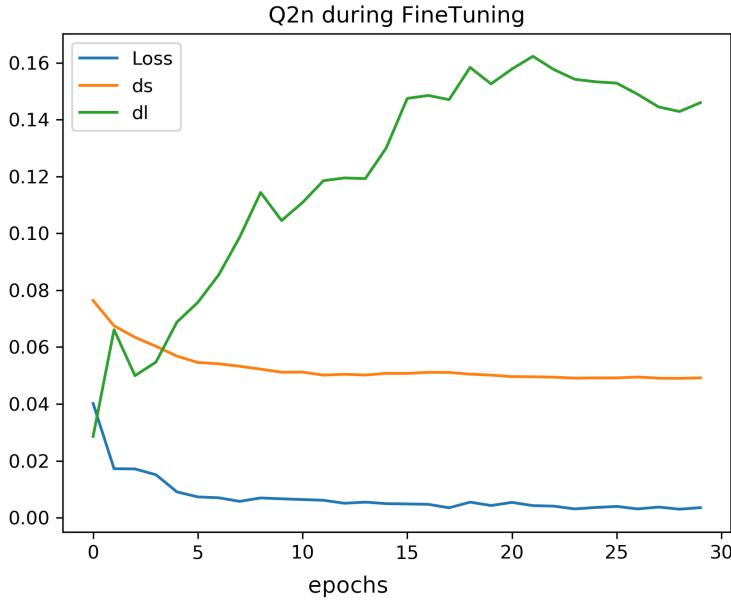


Figure 4.3: Values of Loss, DS and DL of QNR during the training.

the GT image itself using the MSE loss function. The purpose is to compare the previous described methods with the best possible solution; however it is not applicable in a real approach.

4.3 Results

The use of the classical QNR function, whose formula is given by Eq. (1.20)), has produced no appreciable results. During the training, while the error continues to decrease, the image quality indexes computed through the Matlab Toolbox get worse, as showed in Fig. 4.3, 4.4 and 4.5, even the QNR itself. The reason is related to the fact that, in order to minimize the error, it has been generated an image with negative values. An input with negative values is not

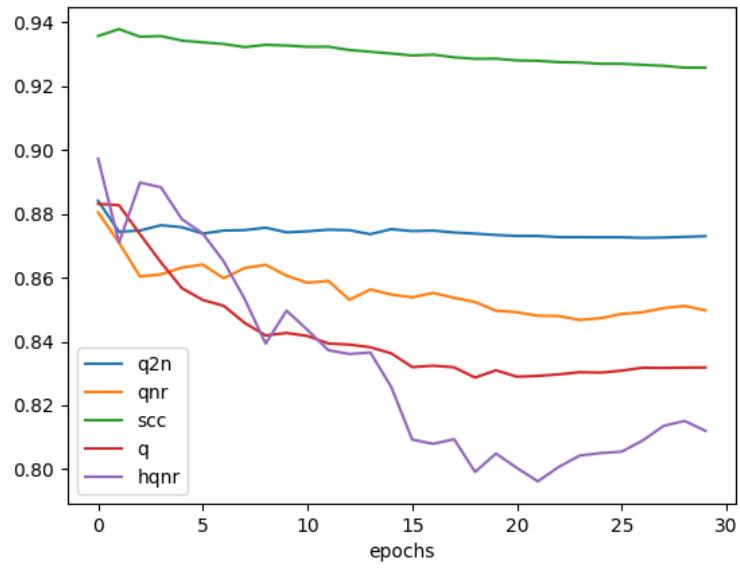


Figure 4.4: Values of Q2N, QNR, SCC, Q and HQNR during the training.

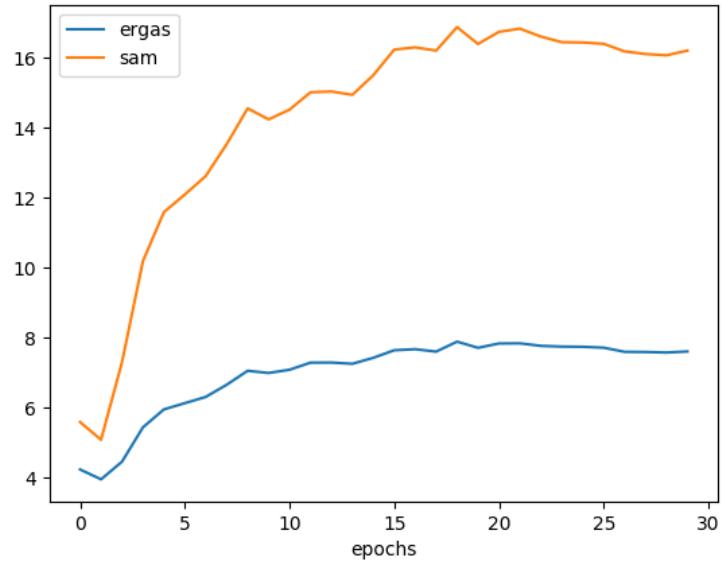


Figure 4.5: Values of ERGAS and SAM during the training.

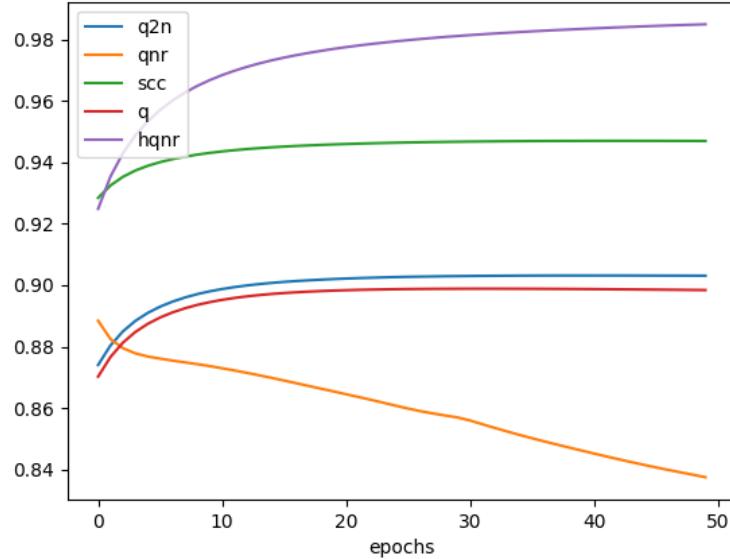


Figure 4.6: Values of Q2N, QNR, SCC, Q and HQNR during the training using the loss function with HQNR.

expected by any of the indexes used. This creates a discrepancy between the QNR used in the loss function and the QNR of the Matlab Toolbox. Indeed, with inputs having values between 0 and 1, the QNR developed in Tensorflow and the Matlab one returns the same results. In other words, the loss tends to reach his minimum value even though this is not reflected in an increase of the quality indexes.

After those tests, it has been decided to change the loss function using another *no-reference* index, the HQNR, whose expression is reported in Eq. (1.23). However, this index is more complex than the QNR, as it involves quaternions and operations with hypercomplex numbers.

For this reason, it was applied an approximated version of HQNR that is

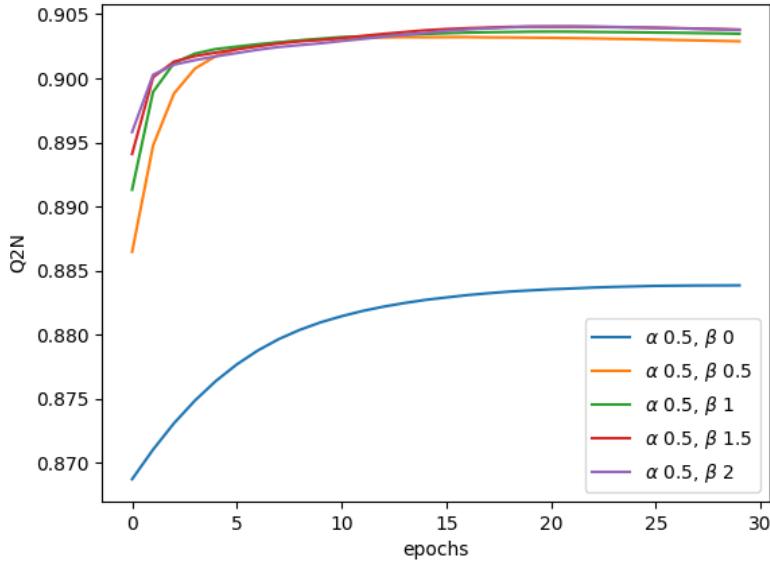


Figure 4.7: Learning curves obtained by varying β with a 0.5 step in the range $[0, 2]$ and α set to 0.5.

fully described in the previous chapter. The results obtained for the Ikonos Toulouse dataset are shown in Fig. 4.6.

With this type of loss function, it was recorded a performance improvement in terms of all the indexes, except the QNR. The latter, as evidenced also by the previous results, does not seem to be a good index for the quality assessment.

To explore the behavior of the training algorithm for different values of the coefficients α and β that are present in Eq. (1.23), several tests have been carried out by firstly modifying the parameters with steps of 0.5. An interesting trend has been noticed with α equal to 0.5 and results of those tests are available in Fig. 4.7. Furthermore, tests with α set to 0.5 and β ranging from 0.5 to 2 with 0.25 steps have been carried out. The results are

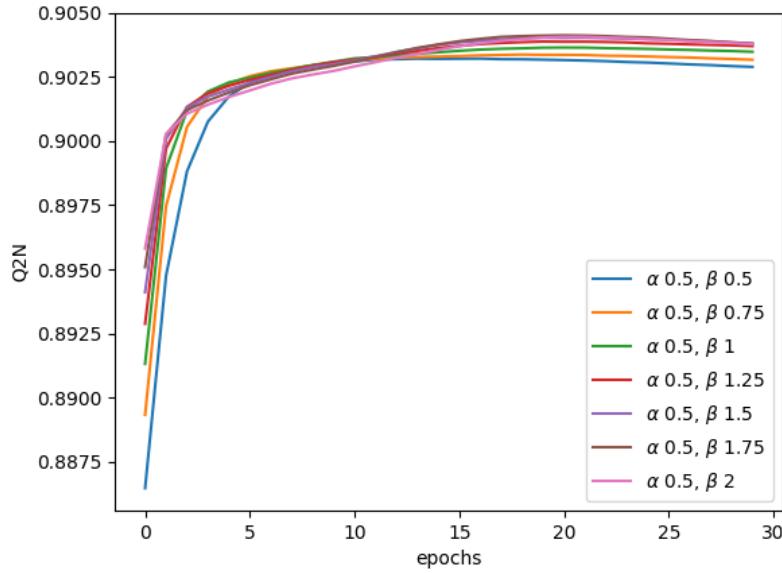


Figure 4.8: Learning curves obtained by varying β with a 0.25 step in the range [0.5, 2] and α set to 0.5.

showed in Figs. 4.8 and 4.9. The close-up provided from the latter figure allows to conclude that the best results are obtained with $\alpha = 0.5$ and $\beta = 1.75$.

After these experiments aimed at selecting the optimal values for the proposed approach, a comparison with existing techniques has been carried out. In particular, the GT method uses the *ground truth* during the training and is thus introduced as a reference since it cannot be employed in the practice; instead, RR denotes the method used in [4] with reduced resolution training. The method proposed in this work that uses full-resolution images and is described in the previous paragraphs is here denoted as NOREF.

Fig. 4.10 shows the comparison between the three different methodologies. The graph indicates that the Q2N trend (in green) of the GT approach has

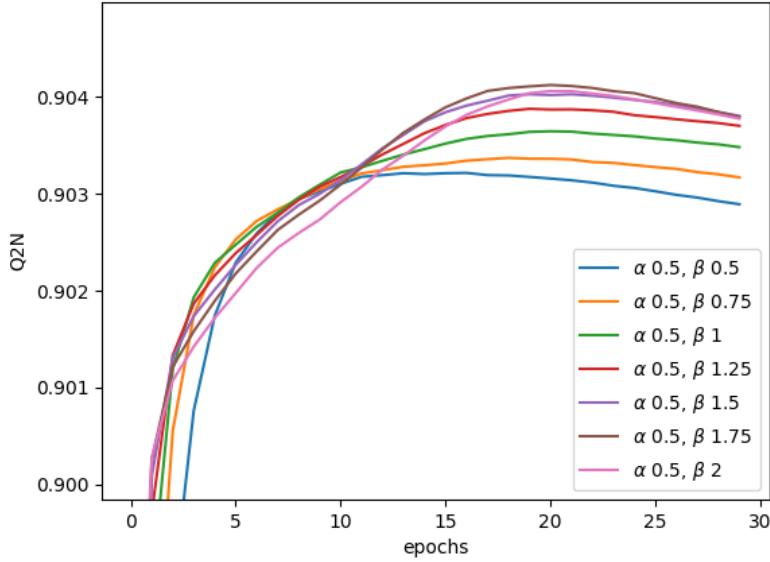


Figure 4.9: Zoom of Fig. 4.8 with Q2N in the range [0.9, 0.905].

a high margin of performance over the others. The expectation was that the model would reach a level close to 1. However, with the actual architecture, the model cannot achieve higher performance. Worst performance has been obtained by RR (in blue) where the trend increases in the first two epochs and then remains stuck below 0.89. The NOREF method (in yellow) yields a continuous growth and remains constant from the 30th epoch onwards. The model has better performance than RR despite his recent application.

In Figs. 4.12, 4.13, 4.14 the images obtained with GT, NOREF and RR methods are illustrated respectively. The images show some border effects noticeable primarily on the left border. Those effects are introduced since the prediction of the model (high pass compo-

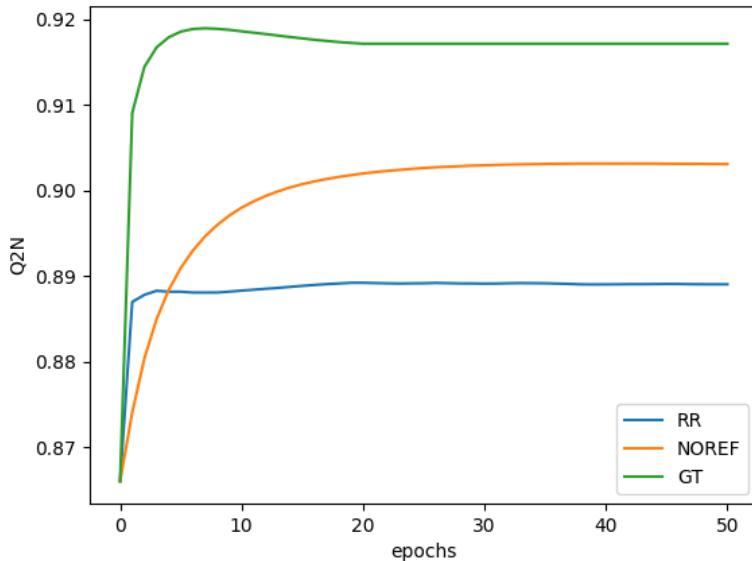


Figure 4.10: Q2N experiment results with GT, RR and NOREF.

ment) that have a reduced size is added to the MS image used for the image (low pass component). The convolutional layers of the network reduce the size of the output layer-by-layer depending on the used kernel. In particular, the model used for the experiments decreases the image by 6 pixels per side.

For a better visual inspection, the differences of NOREF and RR with the ground truth are showed in the Fig. 4.15. The figure emphasise with red and yellow the significant errors and with blue the areas the negligible ones. The RR difference image shows some red spots on the bottom and on the center but in the NOREF these areas are processed in a better way.



Figure 4.11: Ground Truth used for validation



Figure 4.12: Result of GT method



Figure 4.13: Result of NOREF method



Figure 4.14: Result of RR method

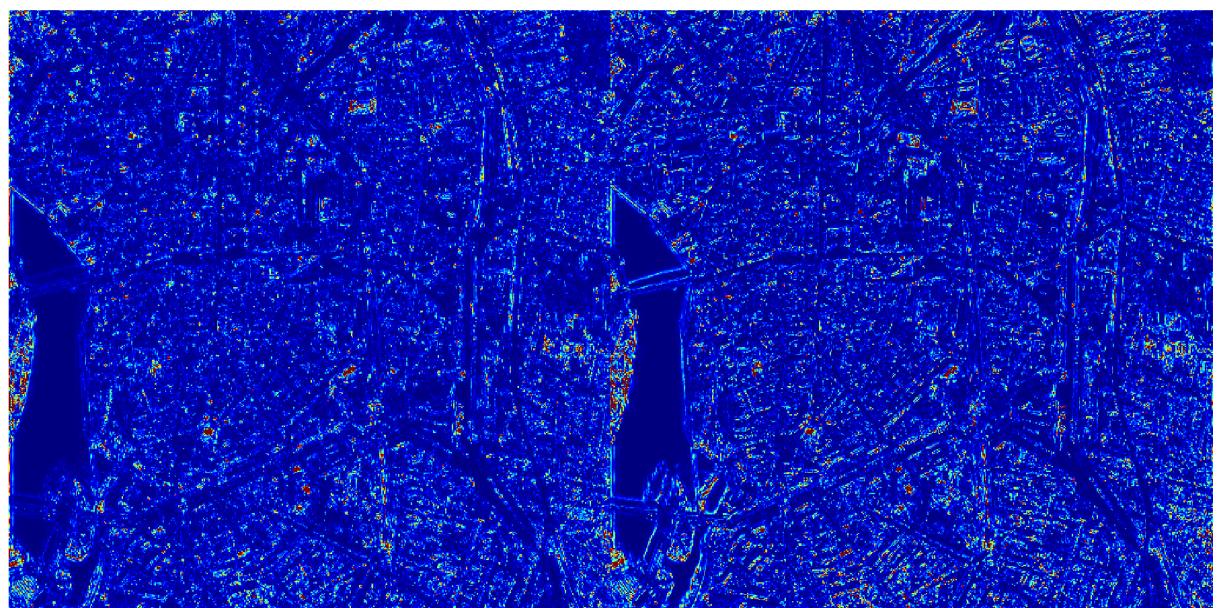


Figure 4.15: Images representing the channels sum of the difference between NOREF and RR with ground truth

Conclusions

In this work, we started from the recently proposed [37] CNN-based pan-sharpening method and, using modern and more compatible tools, explored a *no-reference* approach. This opened numerous issue regarding the loss function to choose. Without a target for the model, the loss function became the critical aspect more than before.

In the current study, we implemented new training and validation processes that do not make use of a reference image (NOREF) and compared them with existing and reference approaches. It was demonstrated that the use of the QNR index as full resolution loss is not appropriate, since it does not provide valuable results in this settings. On the contrary, we developed an approximate implementation of the HQNR index that obtains better performance than the widespread training operating at reduced resolution (RR). It was also shown, by resorting to an ideal training scheme based on the ground truth (GT), that there is room for further improvements.

It is arguable that the use of a hybrid RR and NOREF approach can have even better performance as it inherits the best aspects of both methods.

Indeed, a hybrid model would use reduced resolution images with a robust reference index and, at the same time, images at the original resolution would compensate any problems due to the different behavior of the algorithms across scales.

Furthermore, different architectures can be considered. In fact, an architecture with only convolutional layers includes a limited number of filters and every filter is applied to the whole image. This implicates that a single filter is trained at each level by optimize the performance on the whole image instead on a particular area. Images can have different environments (like rural or urban) and the model could benefit by neurons divided to specialize different patches.

This results can be achieved with different techniques. At first glance, it would be possible to generate multiple small models, which are fine-tuned for patches. Furthermore, changing the model layers would deliver images processed in a specialized way. To this extent, in a dense layer, all neurons have their own specific weights that would be trained for the exact pixels in input. Indeed, the actual limit in *pansharpening neural networks* (PNN) is treating images without considering the different morphology of them.

For these reasons, the current research lay the foundation for future development in the field of pansharpening. A *no-reference* model with a more implemented and specialized image processing can overcome the limits of the model in use.

Bibliography

- [1] C. Souza Jr. L. Firestone L. M. Silva and D. Roberts. *Mapping forest degradation in the Eastern Amazon from SPOT 4 through spectral mixture models.* Remote Sens. Environ., 2003.
- [2] A. Mohammadzadeh A. Tavakoli and M. J. Valadan Zoej. *Road extraction based on fuzzy logic and mathematical morphology from pansharpened IKONOS images.* Photogramm. Rec., 2006.
- [3] F. Laporterie-Déjean H. de Boissezon G. Flouzat and M.-J. Lefèvre-Fonollosa. *Thematic and statistical evaluations of five panchromatic/-multispectral fusion methods on simulated PLEIADES-HR images.* Inf. Fusion, 2005.
- [4] Giuseppe Masi Davide Cozzolino Luisa Verdoliva and Giuseppe Scarpa. *Pansharpening by Convolutional Neural Networks.* Remote Sensing, 2016.
- [5] Pascal Lamblin. *MILA and the future of Theano.* <https://groups.google.com/forum/#topic/theano-users/7Poq8BZutbY>, 2018.

- [6] tensorflow.org. *Introduction to Gradients and Automatic Differentiation*.
<https://www.tensorflow.org/guide/autodiff?hl=da>, 2020.
- [7] W. Carper T. Lillesand and R. Kiefer. *The use of intensity-hue-saturation transformations for merging SPOT panchromatic and multi-spectral image data*,. Photogramm. Eng. Remote Sens., 1990.
- [8] P. S. Chavez Jr. S. C. Sides and J. A. Anderson. *Comparison of three different methods to merge multiresolution and multispectral data: Landsat TM and SPOT panchromatic*,. Photogramm. Eng. Remote Sens., 1991.
- [9] C. Thomas and L. Wald. *Analysis of changes in quality assessment with scale*. Proc. 9th Int. Conf. Inf. Fusion, 2006.
- [10] V. P. Shah N. H. Younan and R. L. King. *An efficient pan-sharpening method via a combined adaptive-PCA approach and contourlets*. IEEE Trans. Geosci. Remote Sens., 2008.
- [11] C. A. Laben and B. V. Brower. *Process for enhancing the spatial resolution of multispectral imagery using pan-sharpening*,. U.S. Patent 6 011 875, 2000.
- [12] L. Wald T. Ranchin and M. Mangolini. *Fusion of satellite images of different spatial resolutions: Assessing the quality of resulting images*. Photogramm. Eng. Remote Sens, 1997.
- [13] L. Wald C. Thomas, T. Ranchin and J. Chanussot. *Synthesis of multispectral images to high spatial resolution: A critical review of fusion*

methods based on remote sensing physics. IEEE Trans. Geosci. Remote Sens., 2008.

- [14] M. Vega R. Molina I. Amro, J. Mateos and A. K. Katsaggelos. *A survey of classical methods and new trends in pansharpening of multispectral images.* EURASIP J. Adv. Signal Process., 2011.
- [15] B. Aiazzi S. Baronti and M. Selva. *Improving component substitution pansharpening through multivariate regression of MS+Pan data.* IEEE Trans. Geosci. Remote Sens., 2007.
- [16] T.-M. Tu S.-C. Su H.-C. Shyu and P. S. Huang. *A new look at IHS-like image fusion methods.* Inf. Fusion, 2001.
- [17] T.-M. Tu P. S. Huang C.-L. Hung and C.-P. Chang. *A fast intensity-hue-saturation fusion technique with spectral adjustment for IKONOS imagery.* IEEE Trans. Geosci. Remote Sens., 2004.
- [18] W. Dou Y. Chen X. Li and D. Sui. *A general framework for component substitution image fusion: An implementation using fast image fusion method.* Comput. Geosci., 2007.
- [19] J. Choi K. Yu and Y. Kim. *A new adaptive component-substitution based satellite image fusion by using partial replacement.* IEEE Trans. Geosci. Remote Sens., 2011.
- [20] Gemine Vivone Luciano Alparone Jocelyn Chanussot Mauro Dalla Mura andrea Garzelli Giorgio A. Licciardi Rocco Restaino and Lucien Wald. *A*

Critical Comparison Among Pansharpening Algorithms. IEEE Transactions on Geoscience and Remote Sensing, 2015.

- [21] L. Alparone et al. Comparison of three different methods to merge multiresolution and multispectral data: Landsat tm and spot panchromatic. 2007.
- [22] F. Laporterie-Déjean H. de Boissezon G. Flouzat and M.-J. Lefèvre-Fonollosa. *Thematic and statistical evaluations of five panchromatic/multispectral fusion methods on simulated PLEIADES-HR images.* Inf. Fusion, 2005.
- [23] B. Aiazzi L. Alparone S. Baronti A. Garzelli and M. Selva. *MTF-tailored multiscale fusion of high-resolution MS and Pan imagery.* Photogramm. Eng. Remote Sens., 2006.
- [24] Z. Wang and A. C. Bovik. *A universal image quality index.* IEEE Signal Process. Lett., 2002.
- [25] S. Garzelli A. Alparone, L. Baronti and Nencini F. *A global quality measurement of pan-sharpened multispectral imagery.* IEEE Geosci. Remote Sens., 2004.
- [26] L Alparone et al. *Multispectral and panchromatic data fusion assessment without reference.* Photogramm. Eng. Remote Sens., 2008.

- [27] Khan M. M. Alparone L. and Chanussot J. *Pansharpening quality assessment using the modulation transfer functions of instruments*. IEEE Trans. Geosci. Remote Sens., 2009.
- [28] B. Aiazza L. Alparonea S. Barontia R. Carl A. Garzelliia L. Santurri. *Full scale assessment of pansharpening methods and data products*. Photogramm. Eng. Remote Sens., 2008.
- [29] Warren S. McCulloch and Walter Pitts. *A Logical Calculus of Ideas Immanent in Nervous Activity*. University of Chicago, 1943.
- [30] MissingLink. *Perceptrons and Multi-Layer Perceptrons: The Artificial Neuron at the Core of Deep Learning*.
<https://missinglink.ai/guides/neural-network-concepts/perceptrons-and-multi-layer-perceptrons-the-artificial-neuron-at-the-core-of-deep-learning>.
- [31] Marvin Minsky and Seymour Papert. *Perceptrons: an introduction to computational geometry*. 1969.
- [32] Sepp Hochreiter and Jurgen Schmidhuber. Long short-term memor. 1997.
- [33] Cho Kyunghyun van Merriënboer Bart Gulcehre Caglar Bahdanau Dzmitry Bougares Fethi Schwenk Holger Bengio Yoshua. Learning phrase representations using rnn encoder-decoder for statistical machine translation. 2014.

- [34] RSIPvision. <https://www.rsipvision.com/exploring-deep-learning>.
- [35] Bengio Yoshua LeCun Yann Hinton and Geoffrey. Deep learning. *Nature* 521, 2015.
- [36] K. Tang X Dong C. Loy C. He. *Image super-resolution using deep convolutional networks*. IEEE Trans. Pattern Anal. Mach. Intell., 2006.
- [37] Giuseppe Scarpa Sergio Vitale and Davide Cozzolino. *CNN-based pansharpening of multi-resolution remote-sensing images*. Conference: 2017 Joint Urban Remote Sensing Event (JURSE).
- [38] Giuseppe Scarpa Sergio Vitale and Davide Cozzolino. *Target-adaptive CNN-based pansharpening*. 2017.
- [39] K He X Zhang S Ren and J Sun. *Deep residual learning for image recognition*. IEEE Conf on Computer Vision and Pattern Recognition (CVPR), 2016.
- [40] R Zeyde M Elad and M Protter. *On single image scale-up using sparse-representations*. Curves and Surfaces Springer, 2012.
- [41] R Timofte V De and L V Gool. *Anchored neighborhood regression for fast example-based super-resolution*. IEEE Int Conf on Computer Vision (ICCV), 2013.

- [42] J K L J Kim and K M Lee. *Accurate image super-resolution using very deep convolutional networks*. IEEE Conf on Compute Vision and Pattern Recognition (CVPR), 2016.
- [43] George Dahl Ilya Sutskever1, James Martens. On the importance of initialization and momentum in deep learning. 2013.
- [44] Atilim Gunes Baydin Barak A Pearlmutter Alexey Andreyevich Radul and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. 2018.
- [45] Sören Laue. *On the Equivalence of Forward Mode Automatic Differentiation and Symbolic Differentiation*. 2019.
- [46] Luciano Alparone Andrea Garzelli Gemine Vivone. *Spatial Consistency for Full-Scale Assessment of Pansharpening*. IEEE International Geoscience and Remote Sensing Symposium, 2018.
- [47] Jimmy Ba Diederik P. Kingma. Adam: A method for stochastic optimization. 2014.
- [48] Matlab. <https://it.mathworks.com/help/matlab/matlab-engine-for-python.html>.

Appendix A

Python implementation of no-reference quality indexes

A.1 QNR

```
import tensorflow as tf
import tensorflow_probability as tfp

def q_index(y_true, y_pred):
    two = tf.constant(2.0, tf.float32)

    cov_b = tfp.stats.covariance(y_true, y_pred, [0, 1], None)
    true_b_std = tf.math.reduce_std(y_true, [0, 1])
    pred_b_std = tf.math.reduce_std(y_pred, [0, 1])

    true_b_mean = tf.cast(\n        tf.reduce_mean(y_true, [0, 1]), \n        tf.float32)
    pred_b_mean = tf.cast(\n        tf.reduce_mean(y_pred, [0, 1]), \n        tf.float32)

    q1_b = cov_b / (true_b_std * pred_b_std)
    q2_b = (two * true_b_mean * pred_b_mean) / \
            (tf.square(true_b_mean) + tf.square(pred_b_mean))
    q3_b = (two * true_b_std * pred_b_std) / \
            (tf.square(true_b_std) + tf.square(pred_b_std))

    q_b = q1_b * q2_b * q3_b
    return tf.reduce_mean(q_b)
```

```

def d_lambda(ms, f, p, b):
    result = tf.constant(0.0, tf.float32)
    for l in range(b-1):
        for r in range(l+1, b):
            result += tf.abs(\n
                tf.cast(\n
                    q_index(f[:, :, l:l+1], f[:, :, r:r+1]),\n
                    tf.float32) - \
                tf.cast(\n
                    q_index(ms[:, :, l:l+1], ms[:, :, r:r+1]),\n
                    tf.float32)\n
            )**p

    b = tf.constant(b, tf.float32)

    s = ( b * ( b - tf.constant(1.0, tf.float32) ) ) \
        / tf.constant(2.0, tf.float32)

    result = result / s
    result = result ** (1.0/p)
    return result

def d_s(ms, fused, pan, pan_degraded, q, b):
    result = tf.constant(0.0, tf.float32)

    for l in range(b):
        result += tf.abs(\n
            tf.cast(\n
                q_index(fused[:, :, l:l+1], pan),\n
                tf.float32) - \
            tf.cast(\n
                q_index(ms[:, :, l:l+1], pan_degraded),\n
                tf.float32)\n
        ))**q

    b = tf.constant(b, tf.float32)

    result = result / b

    r = result **(1./q)
    return r

def qnr(fused, ms, pan, pan_degraded, alpha, beta, p, q, bands):

```

```

a = (1-d_lambda(ms, fused, p=p, b=bands))**alpha
b = (1-d_s(ms, fused, pan, pan_degraded, q, b=bands))**beta
return a*b

```

A.2 HQNR

```

import tensorflow as tf

def filter_image(origin_image, kernel):
    image = tf.expand_dims(origin_image, 0)
    kernel = tf.expand_dims(kernel, -1)
    output = tf.nn.depthwise_conv2d(\ 
        image, kernel, strides=(1, 1, 1, 1), padding="VALID")
    output = tf.squeeze(output)
    return output

def gaussian_filtered_image(image, sensor):
    if sensor == "WV2":
        gauss_kernel = hWV2
    elif sensor == "WV3":
        gauss_kernel = hWV3
    elif sensor == "GeoEye1":
        gauss_kernel = hGE1
    else:
        gauss_kernel = hIK

    return filter_image(image, gauss_kernel)

def d_s_reg(ms, pan):
    ms = tf.reshape(ms, (ms.shape[0] * ms.shape[1], ms.shape[2]))
    pan = tf.reshape(pan, (pan.shape[0] * pan.shape[1], 1))

    alpha = tf.matmul(tf.linalg.pinv(ms), pan)
    fi = tf.matmul(ms, alpha)
    return 1 - r_squared(pan, fi)

def d_lambda(ms, fused, p, b, sensor):
    fused_filtered = gaussian_filtered_image(fused, sensor)
    return 1 - q_index(fused_filtered, ms[R:-R, R:-R, :])

def hqnr(f, ms, pan, alpha, beta, bands, sensor):
    a = (1 - d_lambda(\ 
        ms, f, b=bands, sensor=sensor))**alpha
    b = (1 - d_s_reg(f, pan)) ** beta

```

```
return a*b
```