# DEEP LEARNING SOLUTIONS FOR TANDEM-X-BASED FOREST CLASSIFICATION

*Antonio Mazza*[1], *Francescopaolo Sica*[2]

[1]DIETI, University Federico II, Via Claudio 21, 80125 Naples, Italy
[2]Microwaves and Radar Institute, DLR, Münchener Straße 20, 82234 Weßling, Germany

## ABSTRACT

In the last few years, deep learning (DL) has been successfully and massively employed in computer vision for discriminative tasks, such as image classification or object detection. This kind of problems are core to many remote sensing (RS) applications as well, though with domain-specific peculiarities. Therefore, there is a growing interest on the use of DL methods for RS tasks. Here, we consider the forest/non-forest classification problem with TanDEM-X data, and test two state-of-the-art DL models, suitably adapting them to the specific task. Our experiments confirm the great potential of DL methods for RS applications.

***Index Terms***— Deep Learning; Convolutional Neural Network (CNN); Vegetation Monitoring; Forest Classification; TanDEM-X.

## 1. INTRODUCTION

The monitoring of the state and health of forests is of primary importance for several reasons, such as the prevention of floods and landslides, the reduction of $CO_2$, or the preservation of biodiversity. Thanks to the wide availability of optical, multispectral and synthetic aperture radar (SAR) data from a variety of sensors, such phenomena can be observed on a global scale provided that adequate processing tools are available.

Optical images carry rich discriminative information about vegetation and are widely employed. The Normalized Difference Vegetation Index (NDVI) is a notable example of a standard and simple vegetation indicator that is extracted through a straightforward combination of spectral bands [1]. More specific indicators can also be derived from multispectral images, like the Enhanced Vegetation Index (EVI), more suited to discriminate canopy [2]. However, the use of optical data is severely undermined by their dependence on the weather conditions, which can be only partially mitigated through multitemporal processing and data fusion techniques [3, 4, 5]. On the contrary, SAR data are almost weather insensitive and carry precious information related to ground geometry and electromagnetic propagation [6]. In [7] SAR images obtained in different bands are combined for land cover classification. In [8, 9], the TanDEM-X forest/non-forest map is generated from TanDEM-X bistatic interferometric images, by linking the presence of vegetation to the retrieved InSAR volume decorrelation.

In this work we focus on this latter problem, and experiment with deep learning solutions based on some state-of-the-art models. Specifically, we define and train from scratch two DL architectures following the ResNet [10] and the DenseNet [11] models, respectively. These are adapted to the problem at hand and to the available dataset, also through the definition of suitable loss functions. Forest/non-forest classification maps obtained for a test area located in Pennsylvania confirm the great potential of DL approaches for RS classification tasks. In Section 2 we describe the approach and the details of the two proposed DL networks. Performance indicators, dataset, and experimental results are presented in Section 3. Finally, conclusions are drawn in Section 4.

## 2. PROPOSED DEEP LEARNING APPROACH

Deep learning models are characterized by an extremely large number of parameters to be trained, ranging from hundreds of thousands to billions, and organized in interconnected *layers* in order to generate a hierarchy of representations of the input. Convolutional Neural Networks (CNNs) are a popular family of DL models, particularly suited to solving image processing problems. In fact, under the assumptions of locality and shift-invariance, they adopt limited receptive fields and weight reuse, thereby ensuring a drastic reduction of the number of free parameters. In this work, we consider two state-of-the-art CNN models, ResNet [10] and DenseNet [11], which are particularly appealing as they can be reach a considerable depth avoiding vanishing gradient problems during training. Both solutions are modular, allowing to build a variety of different architectures, from a simple cascade structure of an arbitrary number of layers to multipath architectures differing in the layer definition. For both models, we describe here only the main functional aspects of interest for the present work, referring to the original papers [10, 11] for a thorough description of the network architecture.

In DenseNet, each layer is "densely" connected to all preceding ones. Therefore, the input of the $l$-th layer is obtained by concatenating the output features from all previous $l$-1 layers, not just the previous one. This approach, with direct

connections between each pair of layers, mitigates vanishing gradient and overfitting problems for large scale tasks. In ResNet, instead, the training phase is shortened by using stages (one or a few consecutive layers) whose output is the combination of the input (via skip connection) with the actual outcome of the trainable backbone. Although functionally unnecessary, skip connections have proven to speed-up the training [10, 12]

Here we propose for both models a cascade architecture with six convolutional layers with $3\times3$ kernels interleaved by ReLU (Rectified Linear Unit) activation functions [13]. Moreover, in order to output a classification probability map, an additional $1\times1$ convolutional output layer with a sigmoid activation function completes the network. The hyperparameters of the networks are summarized in Table 1. The 3-band input stack is formed by the SAR backscatter $\beta_0$, the interferometric coherence, and the local incidence angle, the latter obtained from the acquisition geometry and an external reference digital elevation model.

## 2.1. Training

In this work we exploit the same dataset used in [8], described in Section 3, including the ground-truth reference which is given in terms of density of forest in a squared area of $6\times6$ meters. In order to train the network we explored two different objective loss functions. The former combines two losses commonly used for classification and segmentation, based on cross-entropy ($L_{bce}$) and on the Jaccard distance ($L_J$). The latter includes also the $L_1$ norm, in order to reduce the absolute difference between the reference and the predicted density map. In formulas, the cross-entropy loss reads as

$$L_{bce} = -\frac{1}{N}\sum_n \left[y_n \log\left(\hat{y}_n\right) + (1-y_n)\log\left(1-\hat{y}_n\right)\right],$$
(1)

with $N$ being the number of pixels in a training batch, $y_n$ the class membership degree of pixel $n$ according to the ground-truth,[1] and $\hat{y}_n$ the membership estimated by the network. The Jaccard distance loss, which aims to maximize the overlap between the two maps [14], is defined as

$$L_J = 1 - \frac{\sum_n \left[y_n \cdot \hat{y}_n\right]}{\sum_n \left[y_n + \hat{y}_n - y_n \cdot \hat{y}_n\right]}.$$
(2)

Finally, the $L_1$ norm is

$$L_1 = \frac{1}{N}\sum_n |y_n - \hat{y}_n|.$$
(3)

The minimization of the loss function is performed using the ADAM algorithm [15], a gradient descend variant where the learning rate is updated at each iteration using estimates of low-order moments.

## 3. EXPERIMENTAL RESULTS

The region of interest of the available dataset is located in Pennsylvania (USA). We used 18707 tiles of $128\times128$ pixels for training (90%) and validation (10%). Tiles are grouped in mini-batches of 32 samples for the iterative optimization. For each configuration the initial learning rate was set to $10^{-4}$ and the training was carried out from scratch for 20 epochs.[2] Five large images not used for training, of about $1800\times1450$ pixels, were used to test the performance of the proposed methods. These latter were chosen to be representative of the different environmental contexts.

Two performance indicators related to classification accuracy and segmentation accuracy are considered. Following the methodology used in [8] we have chosen the accuray indicator ACC, a widespread quality index for binary classification problems, defined as

$$ACC = \frac{TP+TN}{TP+FP+TN+FN},$$
(4)

where TP, TN, FP, and FN count true positive, true negative, false positive, and false negative pixels, respectively. In addition, in order to measure performance from the perspective of segmentation, we also considered the Intersection-over-Union (IoU) indicator, which is the intersection between predicted and reference masks over their union. For binary masks it is given by

$$IoU = \frac{TP}{TP+FP+FN}.$$
(5)

Both indicators fall between 0 (worst case) and 1 (ideal prediction). Notice that the above definitions apply for binary images while the network is trained on probability values. Therefore, to make them suited to our problem we decided to properly threshold both reference and predicted maps to get the needed binary masks. To this end we followed the same criterion used in [8], maximizing the Pearson coefficient $\phi$ with respect to the threshold pair (one for prediction, one for reference). This coefficient is defined as

$$\phi = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{P \cdot RP \cdot RN \cdot N}},$$
(6)

where P[N] is the total number of positives[negatives] in the prediction map, RP is the number of positives in the reference (RP=TP+FN) and, conversely, RN is the number of negatives in the reference (RN=FP+TN). The Pearson coefficient is a sort of correlation coefficient which takes values between -1 and +1. The threshold pair that corresponds to the maximum value of $\phi$ is the optimal choice according to this criterion (see [8] for a deeper discussion). In our implementation we used part of the validation set to find the optimal thresholds.

By doing so we eventually collect the performance indicators gathered in Table 2 and Table 3, for IoU and ACC, respectively.

---

[1] We assume the ground-truth density map as membership degree.

[2] A pass on the whole training dataset.

| | ConvLayer 1 | ConvLayer 2 | ConvLayer 3 | ConvLayer 4 | ConvLayer 5 | ConvLayer 6 | ConvLayer 7 |
|---|---|---|---|---|---|---|---|
| Shape (ResNet) | $64\times3\times3\times3$ | $64\times64\times3\times3$ | $64\times64\times3\times3$ | $64\times64\times3\times3$ | $64\times64\times3\times3$ | $64\times64\times3\times3$ | $1\times64\times1\times1$ |
| Shape (DenseNet) | $64\times3\times3\times3$ | $64\times67\times3\times3$ | $64\times131\times3\times3$ | $64\times195\times3\times3$ | $64\times259\times3\times3$ | $64\times323\times3\times3$ | $1\times64\times1\times1$ |
| Activation | ReLU | ReLU | ReLU | ReLU | ReLU | ReLU | sigmoid |

**Table 1**. CNNs' hyper-parameters. Shape: # features $\times$ # channels $\times$ 2D support.

| | Region 1 | Region 2 | Region 3 | Region 4 | Region 5 |
|---|---|---|---|---|---|
| Baseline [8] | 0.4540 | 0.4592 | 0.4245 | 0.6897 | 0.4644 |
| ResNet | 0.5960 | 0.6242 | 0.6915 | 0.8128 | 0.5686 |
| ResNet ($+L_1$) | 0.6013 | **0.6362** | 0.7025 | **0.8354** | 0.5885 |
| DenseNet | **0.6062** | 0.6354 | **0.7087** | 0.8306 | **0.5946** |
| DenseNet ($+L_1$) | 0.5936 | 0.6105 | 0.6868 | 0.8175 | 0.5746 |

**Table 2**. Intersection-over-Union comparison. $+L_1$ marks models trained using the full loss $L = L_{bce} + L_J + L_1$, otherwise limited to the first two terms.

| | Region 1 | Region 2 | Region 3 | Region 4 | Region 5 |
|---|---|---|---|---|---|
| Baseline [8] | 0.6032 | 0.5964 | 0.6030 | 0.7636 | 0.8083 |
| ResNet | 0.8205 | 0.8290 | 0.8769 | 0.8766 | 0.9014 |
| ResNet ($+L_1$) | 0.8125 | 0.8259 | 0.8771 | **0.8901** | 0.9023 |
| DenseNet | 0.8231 | **0.8309** | **0.8841** | 0.8889 | 0.9097 |
| DenseNet ($+L_1$) | **0.8266** | 0.8241 | 0.8788 | 0.8829 | **0.9107** |

**Table 3**. Accuracy comparison.

The numerical results speak clearly in favor of the proposed DL solutions, with DenseNet slightly outperforming ResNet, on average. The use of the $L_1$ norm provides a negligible contribution, likely because the other two loss terms, based on cross-entropy and Jaccard distance, are directly related to ACC and IoU, respectively. It has to be remarked that for a fair comparison with the proposed methods, the baseline solution of [8] was used without masking any class, contrarily to what is done in the original formulation. Specifically, in [8] city and water classes are excluded by means of available masks, because forests, cities and water classes all exhibit a low volume correlation, the core feature proposed to classify forests.

For a further analysis of the performance of the proposed solutions we show some sample images to highlight merits and critical aspects of our proposal. In Fig.1, a case is shown where our proposals work fairly well. In this case, the baseline method also provides results that are coherent with the reference, but rather noisy. In Fig.2 the occurrence of a limited number of false positive (on bridges) can be easily observed for all DL methods. Moreover, some oversmoothing is also noticeable. On the other side, the baseline method falls in a typical failure case, where it is unable to discriminate among forests, water and man-made areas. Finally, Fig.3 shows a detail where all methods present many false negatives. However, the consistency between all predictions suggests that either a change in the scene occurred with respect to the refer-



**Fig. 1**. The mask produced by deep learning approaches are clean compared to the baseline.

ence, or radar data are unable to identify vegetated areas, in this case, due to a more complex backscattering mechanism.

## 4. CONCLUSIONS

In this work, we explored the use of deep learning methods for a forest/non-forest classification problem based on TanDEM-X data. Despite the limited amount of labeled data available for training, the proposed methods show very promising results, in terms of both objective numerical figures and subjective visual assessment. More accurate results, especially in fine details preservation, can be certainly obtained by using larger datasets for training, more sophisticated DL architectures, or additional hand-crafted features, which is the goal of our future work.

**Fig. 2**. False positives and oversmoothing for DL solutions and faliure of the baseline.



**Fig. 3**. False negatives for all.

## 5. REFERENCES

[1] A. Chakraborty, M.V.R. Seshasai, C. Sudhakar Reddy, and V.K. Dadhwal, "Persistent negative changes in seasonal greenness over different forest types of india using modis time series ndvi data (20012014)," *Ecological Indicators*, vol. 85, pp. 887 – 903, 2018.

[2] V. J. Pasquarella, C. E. Holden, and C. E. Woodcock, "Improved mapping of forest type using spectral-temporal landsat features," *Remote Sensing of Environment*, vol. 210, pp. 193 – 207, 2018.

[3] J. Inglada, et al., "Assessment of an operational system for crop type map production using high temporal and spatial resolution satellite optical imagery," *Remote Sensing*, vol. 7, no. 9, pp. 12356–12379, 2015.

[4] G. Scarpa, M. Gargiulo, A. Mazza, and R. Gaetano, "A CNN-Based Fusion Method for Feature Extraction from Sentinel Data," *Remote Sensing*, vol. 10, no. 2, 2018.

[5] A. Errico, C. V. Angelino, L. Cicala, D. P. Podobinski, G. Persechino, C. Ferrara, M. Lega, A. Vallario, C. Parente, G. Masi, R. Gaetano, G. Scarpa, D. Amitrano, G. Ruello, L. Verdoliva, and G. Poggi, "SAR/multispectral image fusion for the detection of environmental hazards with a gis," in *Proceedings of SPIE - The International Society for Optical Engineering*, 2014, vol. 9245.

[6] R. Gaetano, D. Amitrano, G. Masi, G. Poggi, G. Ruello, L. Verdoliva, and G. Scarpa, "Exploration of multitemporal COSMO-skymed data via interactive tree-structured MRF segmentation," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 7, no. 7, pp. 2763–2775, 2014.

[7] R. Hagensieker and B. Waske, "Evaluation of multi-frequency sar images for tropical land cover mapping," *Remote Sensing*, vol. 10, no. 2, 2018.

[8] M. Martone, P. Rizzoli, C. Wecklich, C. Gonzlez, J.-L. Bueso-Bello, P. Valdo, D. Schulze, M. Zink, G. Krieger, and A. Moreira, "The global forest/non-forest map from tandem-x interferometric SAR data," *Remote Sensing of Environment*, vol. 205, pp. 352 – 373, 2018.

[9] M. Martone, F. Sica, C. Gonzlez, J.-L. Bueso-Bello, P. Valdo, and P. Rizzoli, "High-resolution forest mapping from tandem-x interferometric data exploiting nonlocal filtering," *Remote Sensing*, vol. 10, pp. 1477, 2018.

[10] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, June 2016.

[11] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks.," in *CVPR*, 2017.

[12] G. Scarpa, S. Vitale, and D. Cozzolino, "Target-adaptive cnn-based pansharpening," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 9, pp. 5443–5457, Sep. 2018.

[13] A. Krizhevsky, I. Sutskever, and G. E Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[14] G. Csurka, D. Larlus, F. Perronnin, and F. Meylan, "What is a good evaluation measure for semantic segmentation?," in *BMVC*. Citeseer, 2013, vol. 27, p. 2013.

[15] Diederik P. Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014.