

## PoliMi Dropouts

### AN2DL - Second Homework Report

Mattia Piccinato, Matteo Salari, Davide Salonic, Federica Topazio

mattia.piccinato, matteosalari, dadopres, topaz

249496, 252201, 249519, 244935

## 1 Introduction

The goal of this challenge is to address a semantic segmentation task on Mars terrain imagery, classifying each pixel in 64x128 grayscale images into five terrain classes: Background, Soil, Bedrock, Sand, and Big Rocks. The official leaderboard score uses the Mean Intersection over Union (IOU) metric for classes 1 to 4.

$$\frac{1}{|C|} \sum_{c \in C} \frac{\mathbf{1}(y = c) \wedge \mathbf{1}(\hat{y} = c)}{\mathbf{1}(y = c) \vee \mathbf{1}(\hat{y} = c)}, \quad (1)$$

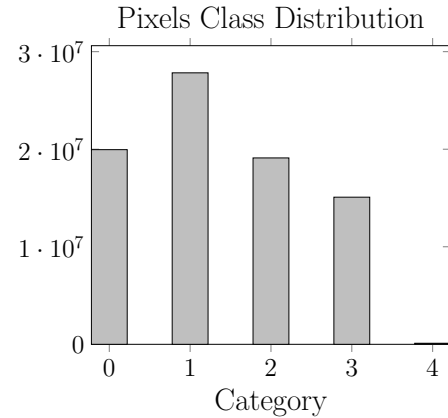
This report outlines the development of our deep learning model, built entirely from scratch without pretrained components, and provides a detailed analysis of the data that guided our decisions.

## 2 Dataset

### 2.1 Analysis and Pre-Processing

The dataset comprised 2615 Martian surface images, including 110 poorly informative images that we promptly removed as it is possible to see in the Python notebook `DataInspection&Cleaning`.

Analyzing the class distribution, we explored strategies such as selectively augmenting the least frequent class2.1. Interestingly, models trained with the original distribution consistently performed better, improving *MeanIOU* by 1%-3%, likely due to align with the test set distribution.



### 2.2 Data Augmentation

We implemented data augmentation to enhance dataset robustness. The final augmentation pipeline, involving quadrupling the dataset with three additional flips per image, yielded the best performance, improving *MeanIOU* by 1%-6% over other methods. Flips preserved dataset characteristics without degrading performance, as opposed to transformations like rotations, zoom, or noise, which consistently underperformed by 1%-3% down in the adopted metric.

### 2.3 Other Methods

Extensive trials tested various transformations individually, in pairs, and with probabilistic combinations. Techniques to increase "Big Rocks" representation included aggressive augmentation, but none

outperformed the flip-based pipeline. Ultimately, this simple approach proved to be optimal without extensive tuning.

## 3 Models

### 3.1 Architecture Exploration

Our model uses a UNet architecture implemented in Keras [RFB15] in our notebook `MultipathUNet`, trained from scratch as per assignment requirements. It follows a four-pathway encoder-decoder structure to capture diverse Mars terrain features.

The network incorporates several key components:

- Multiple parallel encoder paths, which are concatenated in the bottleneck layer:
  - A path incorporating conventional convolutions and max pooling;
  - A residual path featuring skip connections to enhance gradient flow;
  - A global context path designed to capture broad contextual information;
  - A multiscale path with dilated convolutions for multi-resolution feature process.
- Feature fusion with Squeeze-and-Excitation blocks in the bottleneck layer, to effectively merge the extracted features;
- In the decoder, we use concatenation from the corresponding blocks in the encoder to reconstruct the predicted mask;
- Batch normalization layers integrated throughout the network;

For the loss function we used Focal Loss ( $\gamma = 0.5$ ) with class weights [0.0 0.0359 0.0523 0.0663 9.2720], found as the inverse frequency of the pixels per each class (weight for Background class set to 0 manually). This combination helps balance between pixel-wise accuracy and segmentation quality.

$$\mathcal{L}_{\text{focal}} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C w_c \cdot y_{i,c} \cdot (1 - \hat{y}_{i,c})^\gamma \cdot \log(\hat{y}_{i,c}), \quad (2)$$

The final layer employs a Softmax activation to output probability distributions across the five terrain classes for each pixel in the 64x128 input images.

### 3.2 Hyperparameter Tuning

Given the high dimensionality of the hyperparameter space, we adopted a two-phase approach to hyperparameter optimization. Initially, we conducted a manual search focusing on a subset of critical parameters that would most significantly impact model performance. Following this preliminary optimization, we implemented an automated hyperparameter exploration using the `Optuna` library [ASY<sup>+</sup>19], enabling a "distributed" version of the Optuna's Bayesian optimization algorithm throughout the team members.

However, due to the extensive search space, we believe that more optimal hyperparameter configurations could be discovered with additional computational resources and time. We provide a detailed example of our hyperparameter tuning process in our notebook `optuna.ipynb`.

### 3.3 Other methods

In parallel with our main approach, we experimented PSPnet, Manet, DeepLabv3+ and SegFormer, eventually implemented using the PyTorch library [PGM<sup>+</sup>19], which all yielded between 48% and 53% in MeanIOU. Despite proving effectiveness in various semantic segmentation tasks, it yielded inferior results compared to our Keras UNet implementation. Moreover, we implemented a binary classifier as an additional UNet specifically trained to recognize class 4, which proved particularly challenging to predict.

Find an example of implementation using PyTorch library in our notebook `PSPNet`.

We also experimented different loss functions, especially hybrid one combining Sparse Categorical Crossentropy, Dice, Focal or Tversky losses but we never managed to out perform a pure Focal loss. Also different forms of regularization didn't improve results.

## 4 Training

### 4.1 Methods

We decided to use a training and validation split of 0.8 and 0.2, respectively. After experimenting with various splits, this configuration consistently

yielded the best results by ensuring sufficient training data while retaining a meaningful validation set for hyperparameter tuning and a test set for final evaluation.

To further optimize training, we have incorporated callbacks such as **EarlyStopping** and **ReduceLROnPlateau**, which dynamically adjusted hyperparameters like learning rates. More into detail, **EarlyStopping** monitored validation accuracy metric to halt training when improvements plateaued, preventing overfitting; we tried with several other metrics, such as Recall and Loss, but results did not vary significantly.

The batch size was set to the biggest value which could be handled by your development environment, and we found that it could be either set to 128 or 256.

## 4.2 Optimizer

To optimize the model, we employed the AdamW optimizer with a weight decay of 1e-3, coupled with a ReduceLROnPlateau scheduler that reduces the learning rate by a factor of 0.5 after 25 epochs without improvement.

Across different model architectures, our experiments consistently demonstrated that AdamW optimizer, when training with optimized initial learning rates and weight decay parameters, consistently outperformed other widely used optimizers such as standard Adam, SGD, RMSprop and NAdam.

## 4.3 Challenge-Driven submissions

As mentioned in the introduction, the leaderboard metric ignores class 0 (Background). To improve our score, we recalibrated the model to prioritize distinguishing the four target classes (1 to 4), as shown in loss 2, since these classes directly impacted the scoring metric. This adjustment led to a significant score increase, from 0.53 to 0.70, greatly im-

proving our leaderboard position. In a calculated effort to enhance our results and climb the leaderboard, we deliberately deviated from best practices in deep neural network training. Specifically, we performed the data split after applying augmentation, knowingly introducing bias into the validation process. This unorthodox approach allowed us to feed more information into the training phase. While acknowledging the methodological compromise, this strategy yielded a notable improvement of nearly 2% in our performance metrics.

## 5 Results & Conclusions

In Table 1, we listed some of our most relevant submissions (not in chronological order). It shows in each row the model used, some of its main characteristics and the result obtained on the public test set. Accordingly with what stated before, the model featuring multiple encoders with a pure focal loss resulted to be the best for this task. *Multipath UNet* proved to be our best model, as it combines multiple features from the latest state-of-the-art models, offering a powerful combination of advanced techniques.

## 6 Contributions

In our collaborative project, each team member played a crucial role. Federica Topazio worked mainly on literary review and PSPNet, Davide Saltonico focused on hyperparameter tuning, Matteo Salari worked on different models and Mattia Piccinato played a pivotal role in the Data Augmentation domain. We care to clarify that even though we mainly focused on diverse topics, each member was involved in each single choice.

Model	Notes	MeanIOU
PSPNet	Focal&Dice with se-resnext	0.53
Multipath UNet	Focal&Dice Squeeze&Exitation bottleneck	0.70
Multipath UNet	Previous but with CrossEntropyLoss	0.71
<b>Multipath UNet</b>	<b>Previous but with pure Focal Loss</b>	<b>0.74</b>

Table 1: Best submissions

## References

- [ASY<sup>+</sup>19] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A Next-generation Hyperparameter Optimization Framework. *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2623–2631, 2019.
- [PGM<sup>+</sup>19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems 32*, pages 8024–8035, 2019.
- [RFB15] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.