
Skin Lesion Segmentation and Classification using U-Net

Frederico Guth and Teofilo E. deCampos*

Departamento de Ciência da Computação,
Universidade de Brasília (UnB), Brasília-DF, Brazil, CEP 70910-900
fredguth@fredguth.com, t.decampos@oxfordalumni.org

Abstract

This paper describes our approach for the ISIC Challenge 2018 - Skin Lesion Analysis Towards Melanoma Detection. For Part 1 - Lesion Segmentation, we developed a U-Net based convolutional neural network pretrained with the ImageNet dataset and applied several data augmentation and hyperparameters optimization strategies, obtaining threshold Jaccard index of 77.5%. For Part 3 - Lesion Classification, we developed an ensemble strategy that achieved an 84.6% accuracy, which is in the upper bound of human specialists performance range.

1 Introduction

According to the World Health Organization, between 2 and 3 million non-melanoma skin cancers and 132,000 melanoma skin cancers occur globally each year [12]. Despite representing less than 6.5% of all skin cancers, melanomas are the most dangerous type, accounting for approximately 75% of all skin cancer related deaths [12, 3].

Early detection is critical to increase survival expectancy and visual inspection still is the most common diagnostic technique.

Deep convolutional neural networks (CNNs) already exceed human performance in visual classification[4]. In some areas of oncology, such as histological image analysis, CNNs have also proven to match the performance of experts, e.g. [14]. In an attempt to improve the scalability of diagnostic expertise, CNNs have been developed to locate and classify skin cancers in images with dermatologist-level accuracy [3].

Dermoscopy is a technique for examination of skin lesions that, with proper training, increase diagnostic accuracy from 60% (unaided expert visual inspection) to 75%-84%[1]. The International Skin Imaging Collaboration (ISIC) has a large-scale publicly accessible dataset of more than 20,000 dermoscopy images and host an annual benchmark challenge on dermoscopic image analysis since 2016. The challenge comprise 3 tasks of lesion analysis: Segmentation, Dermoscopic feature extraction and Classification. We worked on Segmentation, identifying the lesion region in a dermoscopic image, and Classification, identifying from the image which kind of lesion it is (Nevus, Melanoma, Pigmented Benign Keratosis, Basal Cell Carcinoma, Actinic Keratosis, Vascular Lesion or Dermatofibroma).

In this paper, we describe our approach for the ISIC Challenge 2018:

- **Section 2** describes our methodology for **lesion segmentation**;
- **Section 3** describes our methodology for **lesion classification**;

*<http://cic.unb.br/~teodecampos>

2 Lesion Segmentation

2.1 Segmentation Model: Unet34

Introduced in 2015, U-Net is an encoder-decoder architecture designed for biomedical image segmentation[10], with great results in image segmentation problems [9]. In an U-Net, the output is an image with the same dimension of the input, but with one channel. The encoder path is a typical CNN, where each downsampling step doubles the number of feature channels. But what makes this architecture unique is the decoder path, where each upsampling step input is a concatenation of the output of the previous step with the output of the corresponding (same height) downsampling step. This strategy enables precise localization with a very simple network.

ResNet is a very successful architecture in several computer vision tasks [15]. It mitigates the degradation problem that happens when very deep networks starts converging. Instead of learning a direct mapping $H(x) = y$, it learns the residual function $F(x) = H(x) - x$, which can be reframed into $H(x) = F(x) + x = y$, where $F(x)$ is a stack of non-linear layers and x is the identity function(input=output). The formulation of $F(x) + x$ can be implemented by feedforward neural networks with “shortcut connections”. ResNet34, specifically, is composed of an initial convolutional layer, 16 blocks of 2 layers and a final fully connected layer.

Unet34 is the idea of using a pretrained ResNet34 model as an U-Net encoder path [7]. First, every step from the adaptive pooling onwards is removed, keeping only ResNet backbone. Then we save the output of results of the initial layer, 3rd , 8th, and 14th blocks (of 16 in total). During the upsampling we concatenate the output of those with the outputs of upsampling steps. We used Adam optimizer and Binary Cross Entropy with Logits as the loss function.

2.2 Data and Augmentation

We used “ISIC 2018: Skin Lesion Analysis Towards Melanoma Detection” grand challenge datasets [2, 13] and no additional external data. The Segmentation dataset comprises 2597 training images and 101 validation images acquired with a variety of dermatoscope types, from different anatomic sites, from sample of patients of different institutions. There are more benign lesions than malignant, but an over-representation of malignancies. Mask images are encoded as grayscale 8-bit PNGs, where each pixel is either 0, background, or 255, lesion.

All images were first resized to 128x128 pixels, 256x256 and 512x512; and preprocessed to adjust color balance. Random transformations on input images to augment the dataset were made: dihedral transformation, rotation (up to 44 degrees), zooming (up to 1.05), flipping and random lightening changes. The official training dataset was then split in 3-folds of training and validation datasets.

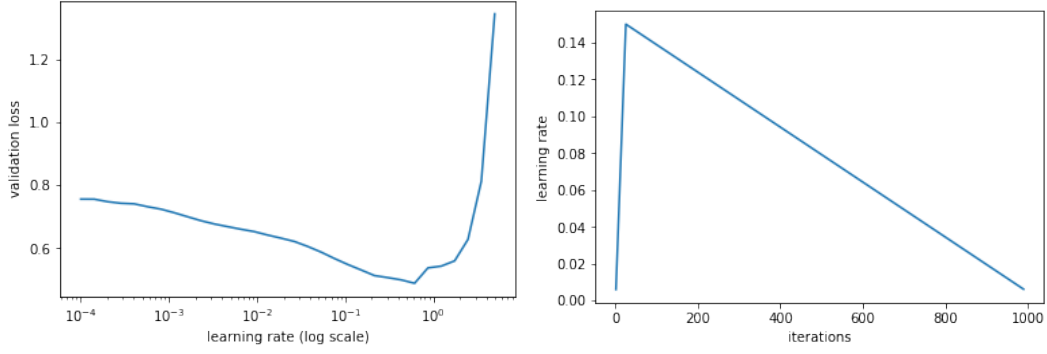
2.3 Segmentation Experiments and Training

Our training strategy was to first train the model with 128x128 images and transfer this learning to train the same model with images with 256 x 256 images. We would suggest using the same strategy to go from 256x256 to 512 x 512 images, but due to GPU memory constraints, we did not manage to accomplish this last step.

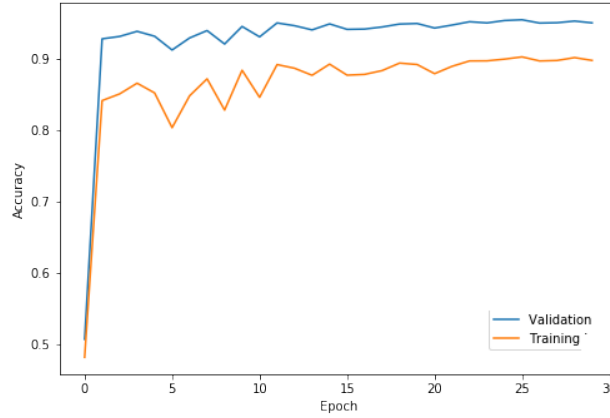
The training procedure was the same on 128x128 and 256x256:

1. Freeze the first layer group.
2. Define the optimal learning rate with the method proposed by [11] and implemented by [7], where one batch is trained with different learning rates, starting at very low and linearly increasing it at every iteration and plotting a chart of the learning rate versus loss (see figure 1a).
3. We use the 1 cyclical learning rate policy (figure 1b), also proposed by [11], to obtain training convergence in only 30 epochs (superconvergence).
4. Unfreeze the model, keeping only the batch normalization layers frozen, and repeat steps 2 and 3.

It would be advisable to use a loss function more similar to the evaluation criteria. As jaccard is not differentiable, one could use a soft jaccard variation[9]. However, in our preliminary trials the



(a) Learning Rate vs Validation Loss (Learning Rate Optimization) (b) Iterations vs Learning Rate (cyclical learning rate policy)



(c) Epoch vs Accuracy (superconvergence)

Figure 1: Learning rate optimization strategy

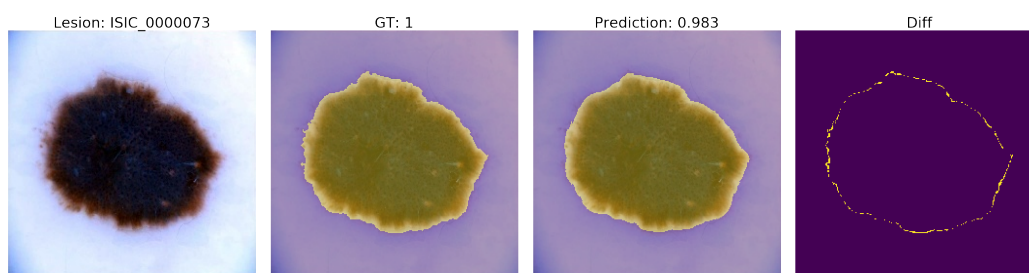
implemented soft jaccard did not improve over the Binary Cross Entropy with Logits loss function and we decided to keep the later.

We developed two strategies for Task 1: *BestDice*: This strategy just predicts the input with the model that presented the best dice index on the validation of our splitted training dataset. *Ensemble*: We used the 3-folds of our training dataset and trained with BestDice model and ensemble than to give a prediction.

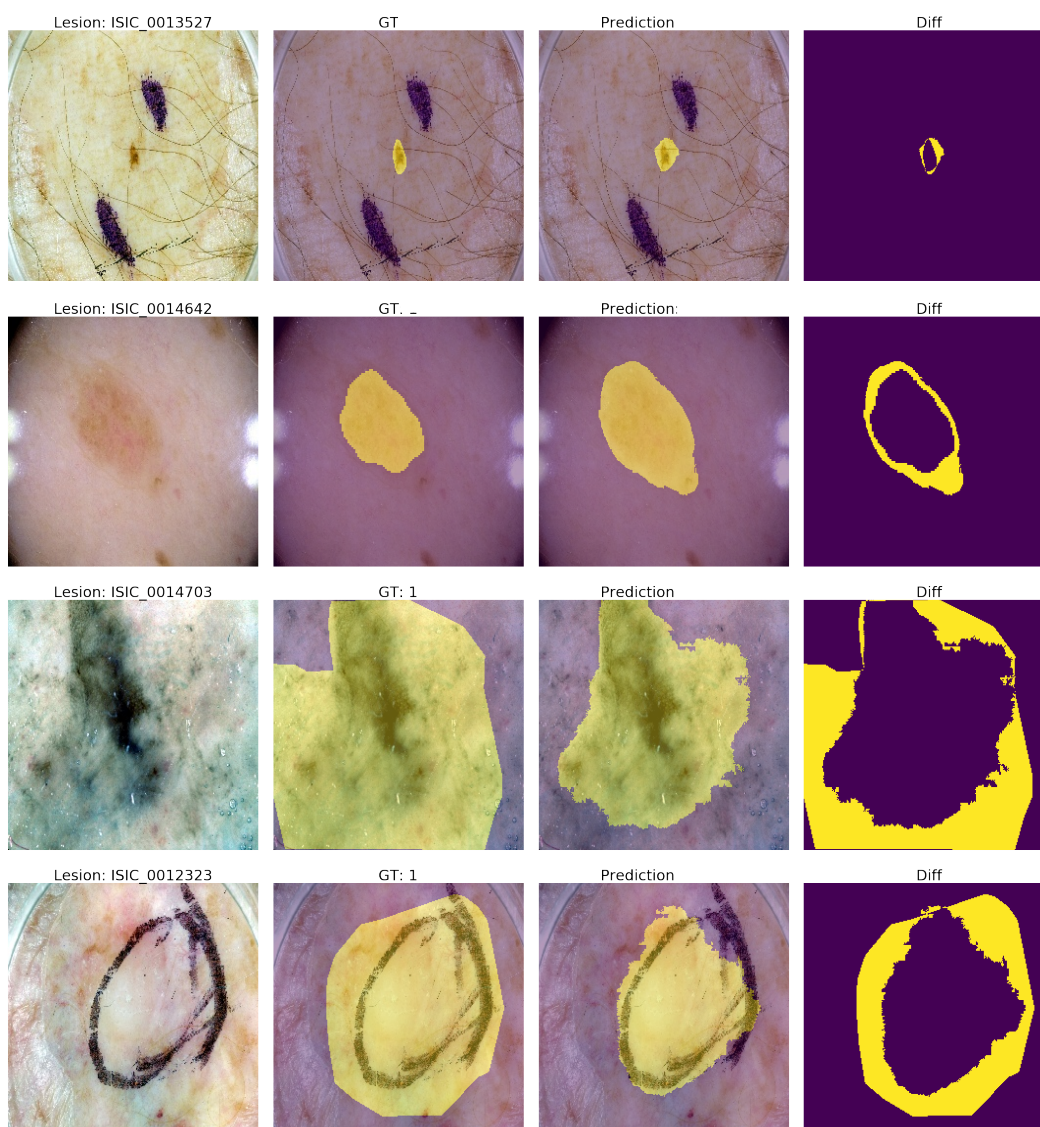
2.4 Segmentation Results

The best result on our validation set was obtained using the Ensemble strategy. It achieves a 85.39% Jaccard index and 78.43% Threshold Jaccard index (with cut at 65%). Surprisingly, it did not score so well with the online score and given official validation set, scoring 71.4%. The small size of the official validation set and the threshold at 65% may be the reason for that (one change in prediction from 64% or 66% Jaccard makes a huge difference in the average threshold jaccard index). The BestDice strategy achieved an online score of 75.5% with the official validation set. Anyway, the top ranked participant in 2017 achieved an average Jaccard Index of 76.5%, which should be compared with our 85.39% score.

Visually the best segmentations are almost identical to the ground truth (see Figure 2a). But we can learn even more from our mistakes. Analysing the worse segmentations there are cases where as non specialists is hard to say if the algorithm was wrong or the ground truth was; there are cases where our algorithm got confused by the pen marker or the glass used by the doctor; and it is clear that in general it doesn't do a good job when the lesion is small relative to the overall image.



(a) Good Result Sample



(b) Bad Result Samples

Figure 2: Qualitative assessment of segmentation result

3 Lesion Classification

3.1 Classification Model

We used ResNet50, an architecture which has been very succesful in image classification problems. We tried using ResNet101 and Resnext101, but due to the lack of computer resources they didn't suit the fast iteration process our research needed. We used the fastai implementation of ResNet50, which has some improvements over the original model proposed by [6]. We customized the ResNet50 model pretrained with the Imagenet dataset [5], by first, removing all the steps from the adaptive pooling onwards keeping only ResNet backbone as a feature extractor. Then we plugged a very simple network consisting of:

- BatchNorm1d (1024, momentum=0.1)
- Dropout (p=0.25)
- Linear (input=1024, output=512)
- ReLU
- BatchNorm1d (1024, eps=1e-5, momentum=0.1)
- Dropout (p=0.5)
- Linear (input=512, output=7)
- LogSoftmax

We used Adam optimizer and the negative log likelihood loss function.

3.2 Data and Augmentation

We used "ISIC 2018: Skin Lesion Analysis Towards Melanoma Detection" grand challenge datasets [2, 13]. The dataset comprises 10015 images and one CSV file containing ground truth response for each image with one of the 7 lesion types: Nevus, Melanoma, Pigmented Benign Keratosis, Basal Cell Carcinoma, Actinic Keratosis, Vascular Lesion or Dermatofibroma.

When we used a pretrained Task1 model, we considered that external data was used. All images were first resized to 224x224 and a color balance copy of the dataset was made. Random transformations on input images to agument the dataset were made: dihedral transformation, rotation (up to 35 degrees), zooming (up to 1.05), flipping and random lightning changes. The official training dataset was then splited in 3-folds of training and validation datasets.

3.3 Classification: Experiments and Training

Priliminar results showed that the color balanced dataset was performing worse than the original and we kept the later. The training procedure was similar to Task 1 (2.3), but we used a cosine annealing learning rate strategy (see figure 3) instead of the one cycle policy (both proposed by [11]). We then tried 3 strategies:

SnapEsem: We saved snapshots of our model during training and then made an ensemble of those snapshots, as proposed by [8].

CropBest: We applied our segmentation algorithm from Task1 to the dataset of Task3, cropping the images by the lesion mininum bounding rectangle. We used a pretrained model from SnapEsem and retrained it using the cropped images.

KFoldCrop: We repeated CropBest strategy to 3 different folds and ensemble them.

On all strategies we applied Test Time Augmentation (TTA) with 5 augmentations per input to improve results.

3.4 Classification Results

Although the best result within our validation set was obtained using the CropBest strategy, we believe that the KfoldCrop can be more general. Interesting to point that all strategies had better accuracy than dermatologists which have an accuracy of 75% to 84%[1].

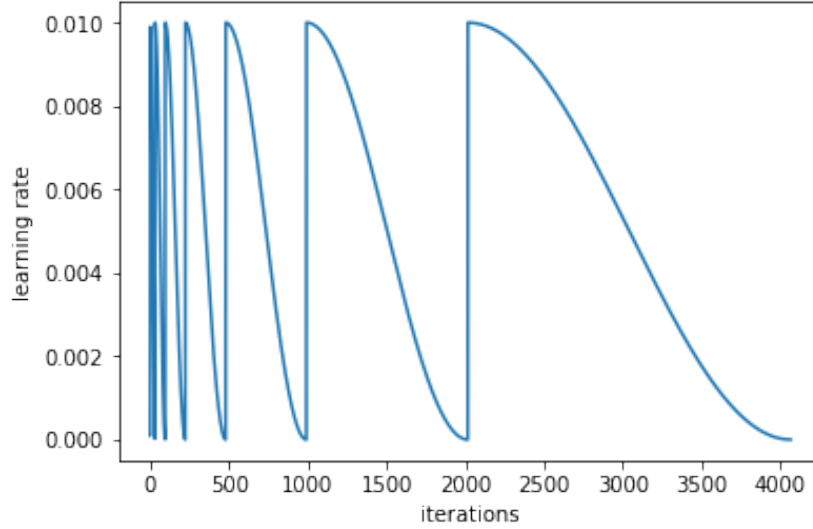


Figure 3: Cosine annealing learning rate strategy.

Table 1: Classification Results.

	Accuracy (%)	Online Score
SnapEsem	89.71	84.6
CropBest	93.26	76.4
KfoldCrop	92.31	75.4

Better than dermatologist specialists (diagnostic accuracy of 75%-84%).

4 Computational Resources

We used Paperspace GPU Cloud service for running all our experiments. The instance server had: 8 cores CPU with 30GB of RAM, Quadro P4000 8 core GPU with 8GB of RAM. The code was developed in PyTorch and Fast.ai[7] on Jupyter Notebooks.

5 Conclusion and Future Work

This paper describes our approach for the ISIC Challenge 2018. For the Lesion Segmentation Task, our model Unet34 achieves 77.5% in the competition threshold jaccard index, which is better than previous results in the competition. For the Lesion Classification, we developed an ensemble strategy that achieves 84.6% accuracy, which is on the upper bound of dermatologists accuracy range.

Future work could improve in many ways. For lesion segmentation, we could repeat the experiment using other architectures as feature extractors. Test time augmentation could be applied, it made difference in our classification algorithm and was not implemented due to the difficulties of inverting the augmentations to later average. We suggest to implement a preprocessing pipeline for removing the pen marks in images and maybe a 2 step process can be developed to roughly finding where the lesion is and cropping a part of the input image to prevent the bad results due to lesion size. For lesion classification, one could different ResNet variations as backbone (ResNet18, ResNet101, ResNext101, etc) and/or using inception. We didn't have time to use the snapshots of the cropped images and just used the final model in the kfold, and this may show better results.

References

- [1] Noel C. F. Codella, David Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin K. Mishra, Harald Kittler, and Allan Halpern. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (ISIC). *CoRR*, abs/1710.05006, 2017.
- [2] Noel C. F. Codella, David Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin K. Mishra, Harald Kittler, and Allan Halpern. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (ISBI), hosted by the international skin imaging collaboration (ISIC). *CoRR*, abs/1710.05006, 2017.
- [3] Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639):115–118, jan 2017.
- [4] Fei-Fei Li and Jia Deng. Where have we been? where are we going? Slides of Talk presented at the ILSVRC – Beyond ImageNet Large Scale Visual Recognition Challenge – workshop in conjunction with CVPR, July 26 2017.
- [5] M. Guillaumin and V. Ferrari. Large-scale knowledge transfer for object localization in ImageNet. In *Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2012.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [7] Jeremy Howard et al. The fast.ai deep learning library, lessons, and tutorials. <https://github.com/fastai/fastai>, 2018.
- [8] Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E. Hopcroft, and Kilian Q. Weinberger. Snapshot ensembles: Train 1, get M for free. *CoRR*, abs/1704.00109, 2017.
- [9] Vladimir Iglovikov, Sergey Mushinskiy, and Vladimir Osin. Satellite imagery feature detection using deep convolutional neural network: A kaggle competition. *CoRR*, abs/1706.06169, 2017.
- [10] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015.
- [11] Leslie N. Smith. A disciplined approach to neural network hyper-parameters: Part 1 - learning rate, batch size, momentum, and weight decay. *CoRR*, abs/1803.09820, 2018.
- [12] Bernard Stewart. *World cancer report 2014*. International Agency for Research on Cancer, WHO Press, World Health Organization, Lyon, France and Geneva, Switzerland, 2014.
- [13] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Sci. Data*, 5:180161, 2018.
- [14] Mitko Veta, Paul J. van Diest, Stefan M. Willems, Haibo Wang, Anant Madabhushi, Angel Cruz-Roa, Fabio Gonzalez, Anders B.L. Larsen, Jacob S. Vestergaard, Anders B. Dahl, Dan C. Ciresan, Jurgen Schmidhuber, Alessandro Giusti, Luca M. Gambardella, F. Boray Tek, Thomas Walter, Ching-Wei Wang, Satoshi Kondo, Bogdan J. Matuszewski, Frederic Precioso, Violet Snell, Josef Kittler, Teofilo E. de Campos, Adnan M. Khan, Nasir M. Rajpoot, Evdokia Arkoumani, Mianela M. Lacle, Max A. Viergever, and Josien P.W. Pluim. Assessment of algorithms for mitosis detection in breast cancer histopathology images. *Medical Image Analysis*, 20(1):237 – 248, 2015.
- [15] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Computer Vision and Pattern Recognition CVPR*, pages 5987–5995, Honolulu, HI, USA, July 21-26 2017.