

LUS Images classification with uncertainty detection and image similarity

Davide Modolo (#229297)

Abstract

The aim of this project was to find an alternative way for lung ultrasound images classification. The developed model is built by three main components and outputs the predicted score of a frame. The entire project can be found on GitHub [1].

Introduction

We were tasked to try an alternative method to create a classifier model that predicts the score of a LUS image.

As explained in the article by S. Roy et al. [2], LUS images are scored as:

- 0: no artifact in the picture;
- 1: at least one vertical artifact (B-line);
- 2: small consolidation below the pleural surface;
- 3: wider hyperechogenic area ($< 50\%$) below the pleural surface.

Frames are taken from videos taken using ultrasound probes and are taken in a maximum of 14 different spots (N on front and M on the back of the patient), as explained in the article by Ali et al. [2].

The first idea that came to my mind was to build something that could be used by doctors. The existing methodology consists in scoring all 14 different spots and summing their values. If the result is $\leq 24/42$, the patient can be left going home because it indicates low probability of worsening.

My idea was to retrieve similar images when the scoring of a specific frame was not sure in order to “help” with the decision.

So, my model can be summarized as follows: multi-class frame classifier, an uncertainty detection model that tries to understand if the first model is confident enough, a similarity module recalled when the first model is not confident that retrieves similar images and analyzes them.

Unfortunately, we will see this idea is probably not effective and results are not encouraging.

1. Data

We have been given a partial dataset from the San Matteo hospital, consisting of 11 patients for a total of $\sim 47k$ frames.

The dataset score distribution is shown in Figure 1; at a first glance it could seem to be almost balanced (with only the score 1 that has less frames), but in reality many patients are inherently unbalanced (the score distribution for each patient is shown in Figure 2).

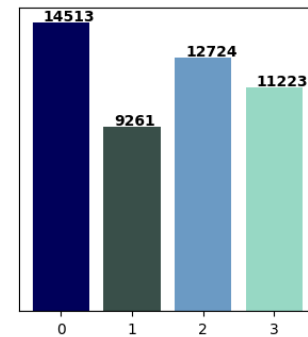


Figure 1. Number of frames for each score in the entire dataset.

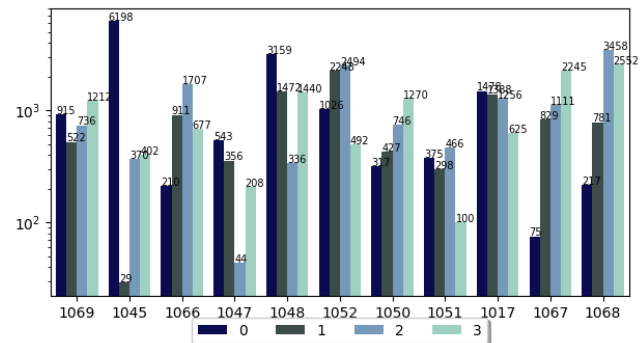


Figure 2. Number of frames for each score for each patient (log scale for better visualization).

1.1 Augmentation

Using the raw dataset got me overfitting even after the first epoch. To address this issue I implemented some transformations taken from the article [2].

In specific, each transformation is activated with a probability of 50%. The set of my augmentation function is:

- affine transformations (translation = $\pm 15\%$, rotation = $\pm 15^\circ$, scaling $\pm 45\%$, and shearing = $\pm 4.5^\circ$)
- multiplication with a constant ($\pm 45\%$)
- Gaussian blurring ($\sigma = 3/4$)
- horizontal flipping ($p = 0.5$)

1.2 Data splitting

Having 11 patients available, my idea was to use 8 of them to train the model and the remaining 3 for testing. This was due to the fact that using a portion of frame for a patient in test and another in train easily leads to overfitting. Even dividing by exams would not been effective since different exams for the same patients still have big correlation.

The first attempt was to test with a k -fold approach and then choose the best configuration, but having 165 combinations with $\sim 4h$ per combination it was unfeasible.

So, to balance the dataset I computed the standard deviation within scores for each 8-patients combination and selected the one with lowest std (Figure 3), resulting in the following division: [0 : 7888, 1 : 7540, 2 : 7189, 3 : 7592]; the problem now was with the test set, that resulted to be very unbalanced. After different attempts to balance both sets, I decided to just select an equal number of images for each score from the training patients set of frames to use in the `test_model` method (still, confusion matrices on this report are built using the entire available number of frames).

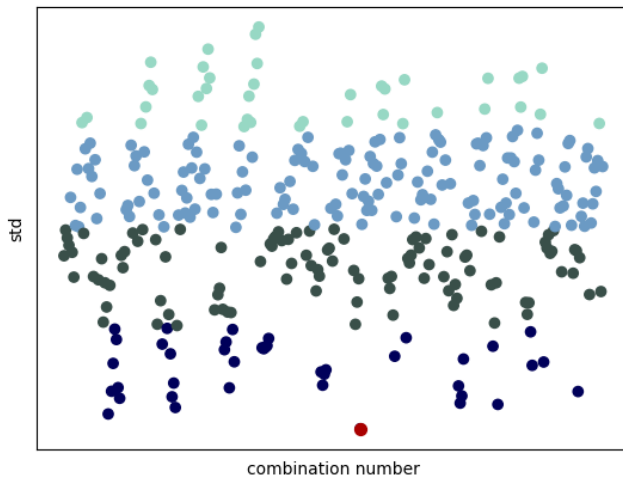


Figure 3. Standard deviation within number of frames per scores of every combination of 8 patients, the red one is the minimum (and so, it is the selected combination).

2. Multi-class classifiers

The first module of my project consists in a deep learning classifier that predicts the score from a frame.

Different pre-trained models have been tested with several different values for my hyperparameters. The training part has been made several times in order to find a model that didn't overfit in the first epoch or didn't stuck in a local minima that always gave one single score.

Frames are very similar, using models too big could get overfitting and using models too small could get no good generalization capability.

2.1 ResNet18

ResNet (Residual Network) is a network introduced by K. He et al. [3] trained on the ImageNet dataset [4].

There are different version of this model based on the number of layers. Looking for a “small” model, ResNet18 was the smallest one and so it has been selected for test.

After many runs, I was able to achieve an accuracy of $\sim 60\%$ in my test set before overfitting. Confusion matrix on the test set can be seen in Figure 4, resulting in an accuracy class-wise that can be seen on Table 1.

0	1	2	3
55.52%	66.88%	40.43%	77.03%

Table 1. Accuracy class-wise of the fine-tuned ResNet18.

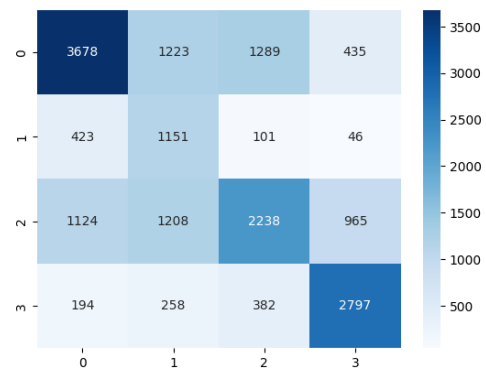


Figure 4. Confusion matrix of the fine-tuned ResNet18.

2.2 VGG16

VGG (Visual Geometry Group) is a convolutional Neural Network built by K. Simonyan, A. Zisserman [5]. It has been trained on a subset of the ImageNet dataset.

Similarly to ResNet, VGG is available with 16 and 18 layers. For the same reasons as above, VGG16 has been selected and tested.

Independently from my fine tuning tries, VGG16 started memorizing the training data even in the first epoch (even having less parameters than ResNet18).

2.3 SqueezeNet

SqueezeNet is a model developed by F. N. Iandola et al. [6] in 2016.

Following the idea to find a compact model, I found out this variation of AlexNet that is still capable of very good performance while requiring less parameters. It has been trained on ImageNet.

SqueezeNet gave me the best results in the early stage of the project, but after refining the fine-tuning of the ResNet, I was able to improve performance by not using SqueezeNet.

2.4 Built-from-scratch model

I even tried building from scratch a Convolutional Neural Network. I tried different combinations of Convolutional

layers but result were very poor, resulting in a path I didn't follow deeper.

3. Binary classifiers

The second goal of my project was find a way to understand if the first module is confident.

The first idea was to use some threshold, but during class presentation we decided it was more interesting to find something that learns the behaviour of the net in both correct and wrong results. Also, max softmax value are very similar between correct and wrong predictions, meaning that a threshold approach would not work.

This created the idea to use the softmax value of the net and classify it in a binary way.

3.1 SVC

bla bla bla

3.2 Deep model

bla bla bla

4. Image similarity

The third module of my project consists in an image similarity model.

The idea was to help the doctor by retrieving ...

4.1 Near Duplicate Image Search

bla bla bla

4.2 t-SNE

bla bla bla

4.2.1 Embedding

bla bla bla

4.2.2 Raw Images

bla bla bla

4.2.3 Behavior

bla bla bla Even if this representation is purely based on the first classifier softmax values (and not on visually similar images) it got the best results.

5. Performance analysis

Very very slow due to image similarity.

Bad performance maybe to: binary gets wrong when is correct. Or binary gets 0 and the t-SNE is wrong in exactly those.

6. Conclusions

bla bla bla slow tsne low acc

6.1 Future works

bla bla bla

References

- [1] GitHub repository with the project. [Online]. Available: <https://github.com/davidemodolo/Lung-Ultrasound-Image-Classifer>
- [2] "Deep Learning for Classification and Localization of COVID-19 Markers in Point-of-Care Lung Ultrasound." [Online]. Available: <https://ieeexplore.ieee.org/document/9093068>
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015.
- [4] ImageNet site. [Online]. Available: <https://www.image-net.org/>
- [5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2015.
- [6] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 10.5mb model size," 2016.

7. Data