

RESEARCH

Introduction to Focus Areas - Project 2

Maximilian Otto^{*}, Dominik Bannwitz, Mariana Steffens and Florian Herzler

^{*}Correspondence:

maximilian.otto@fu-berlin.de

Department of Mathematics and
Computer Science, Free University
of Berlin, Takustraße 9, 14195
Berlin, Germany

Full list of author information is
available at the end of the article

[†]Equal contributor

Abstract

Goal of the project: Classification of images into benign or malignant.

Methods used in the project: 3 different residual neural networks in PyTorch.

Main results of the project: ResNet18 performed better than the more complicated models.

Personal key learnings: Implementation of residual neural networks in PyTorch.

Estimated working hours: 54

Possible improvements: Fine tuning of hyperparameters, downsampling to balance classes, further pre-processing of images

Keywords: Classifiers; Deep Learning; Breast Cancer; ResNet; ResNeXt

1 Scientific Background

Classification in data science is used to classify data into classes, in our case classifying images of tissue samples into benign or malignant tumor classes. A trained classification model can be used, to predict the classes of new, similar data. The used data is unstructured, as it is not stored in tabular but in image form. The information contained in the images can not be accessed directly, but has to be extracted by more complex methods. Hence, the data needs to be split in a train and test set to train a prediction model and evaluate their overall performance. Also, the data can be transformed beforehand to artificially obtain more training data.

2 Goal

The goal of the authors was to provide a publicly available data set of breast cancer histopathology images and perform automated classification into benign or malignant images as a supporting tool for pathologists. Similar to the paper [1] by *Spanhol et al. 2015*, we were tasked with developing and evaluating three different Deep Learning (DL) topologies that work directly on the images of the provided data set. Our goal was to correctly classify as many images of the labeled *test* data as possible as one of the binary outcomes: benign (B) or malignant (M).

3 Data and Preprocessing

Data

The provided publicly available data set^[1] consists of 7909 images of breast tissue of 82 patients, classified into different classes of benign (*adenosis*, *fibroadenoma*, *phylloides tumor*, *tubular adenoma*) and malignant (*ductal carcinoma*, *lobular carcinoma*, *mucinous carcinoma*, *papillary carcinoma*) tumors. The images were obtained in surgical open biopsy (SOB) and classified and labeled as one of the above tumor types by a pathologist. For each type, *Spanhol et al.* provide different microscopic magnifying factors of the tissues (Table 1).

Magnification	Benign	Malignant	Total
40X	652	1,370	1,995
100X	644	1,437	2,081
200X	623	1,390	2,013
400X	588	1,232	1,820
Total of Images	2,480	5,429	7,909

Table 1 Distribution and amounts of breast cancer images with different magnification factors.

We focus our work on just one subset of the data (*fold2*)^[2].

We chose magnification factor 200X and thus ended up with all of the images of factor 200X re-ordered in the corresponding benign or malignant sub folder in a 70 to 30 percent split between training and testing data, similar to *Spanhol et al.*

Preprocessing

Spanhol et al. also provide a Python script^[3] to distribute the images to training and testing folders. Because the authors ran five trials, we modified this script to only re-order the images once, creating a **train** and a **test** directory with sub directories **B** and **M** (benign and malignant) for one magnifying factor. No changes were made to the original data, as there were no missing values in the image data set and we did not balance the sample size of the classes.

Several data augmentations were performed, which got executed during the runtime, but could also be associated as a part of the preprocessing. After downsizing the images to 224x224 pixels, we used randomized horizontal and vertical flips, a random rotation of each image of $(-20, 20)$ degrees, and normalized the images by the mean and standard deviation of all pixels within the train set. The data augmentations are used to simulate more data than provided to tackle an immediate overtraining and give the model a chance to learn more generalized features.

4 Methods

The Python library of choice was PyTorch^[2]. From that library three different networks were used: ResNet-18^[3], wide-ResNet-50^[4] and ResNeXt-50_32xd4^[5]. All three are residual neural networks, which provide an architecture similar to classical Convolutional Neural Networks (CNN). ResNet-networks usually have between 18

^[1]http://www.inf.ufpr.br/vri/databases/BreaKHis_v1.tar.gz

^[2]www.inf.ufpr.br/lesoliveira/download/TBME-00608-2015-R2-preprint.pdf

^[3]<http://www.inf.ufpr.br/lesoliveira/download/mkfold.tar.gz>

and 152 convolutional layers, but can support up to thousands of layers. Shallow CNNs with that many layers would run into the "vanishing gradient problem", which occurs when the algorithm tries to find small weights and by that the loss function is set to the minimal value. With too many layers this leads to a very small gradient until it disappears and no further activation of the following neurons can be registered. A residual network solves this problem by adding *skip-connections*, which add a residue of the previously executed block to the output of a current residual block. Each block consists of multiple convolutional layers, followed by activation functions and batch-normalizations. By using this method, they tackle the "vanishing gradient problem" and can outperform traditional CNNs. That is why we have chosen residual networks as our preferred, underlying architecture. Wide-ResNet und ResNeXt are variations of residual networks, which come with their own advantages and disadvantages. ResNeXt repeats a build-block of ResNet. This enables the ResNeXt to split the input of each block into different channels and analyze spatial relationships, due to the individual pathways within each block. Wide-ResNets simply add more channels to the ResNet-architecture and transform the data into even higher dimensions. Thereby, our Wide-ResNet-50 had about 60 million trainable parameters compared to 25 million of the used ResNeXt-50. The main hyperparameters for this networks remained the same for all models: number of classes, learning rate, the batch size, drop out parameter, and optimizer (Table 2).

<i>Parameter</i>	<i>fValue</i>	<i>Short description</i>
Number of epoch	5	How many epochs the network runs through
Number of classes	2	How many output classes to classify
Learning rate	0.0005	Rate with which the weights are modulated
Batch size	32	Image per iteration
Drop out rate	0.4	Probability of eliminating a neuron
Optimizer	ADAMAX	Gradient descent method

Table 2 Main hyperparameters for training our neural networks

We used the same testing to training split ratio as *Spanhol et al.* with 70% training data to 30% testing data. We tested different combinations of hyperparameters and modified the learning rate with an exponential learning rate scheduler throughout training. The ResNet-18 got trained for 15 epochs. To shorten the training time, we used pre-trained networks. They were previously trained on *ImageNet* [6] and we re-trained all layers. The fully connected layer for classification got replaced with three fully connected layers (# outputs: 256, 100, 2) and *ReLU* as an activation function and a drop-out layer in between.

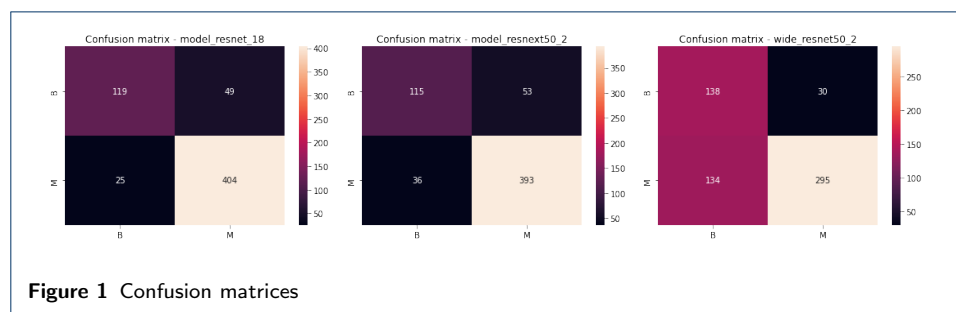
The evaluation consisted of displaying the precision and recall for both classes B and M in a confusion matrix. Additionally, we created heatmaps of the confusion matrices. We also kept track of accuracy and loss (*CrossEntropyLoss*) for both training and test sets and validated the results with the weighted *F1-Score*.

5 Results

To evaluate the three residual neural networks precision and recall as well as the *F1-Score* were determined. Figure 1 shows the confusion matrices for the predictions of the three residual networks on the test data. From these it can be recognised, that the ResNet18 correctly predicted the most malignant tumors on the images with an amount of 404 out of 429 in the test data.

ResNext50 correctly classified only a bit less then ResNet18 with 393, but Wide-ResNet50 was only able to predict 295 of the malignant tumors out of the 429 correctly.

For the prediction of the benign data, it can be recognised in Figure 1 that this was less successful than the classification of malignant tissue data. The ResNet18 was only able to predict 119 out of 168 of the benign tumors on the pictures correctly. The ResNext50 network was again a little less successful than the ResNet18 by predicting 115 of the benign tumors, but the Wide-ResNet50 outperformed the two other networks by correctly predicting 138 of the 168 benign tumors.



ResNet-18

The ResNet18 network reached the highest *F1-Score* out of the three trained models with 87% accuracy. The precision on the benign data lays at 83% and the recall at 71%. The malignant data could be predicted with a precision of 89% and a recall of 94% (See table 3).

	precision	recall	f1-score	support
B	0.83	0.71	0.76	168
M	0.89	0.94	0.92	429
weighted avg	0.87	0.88	0.87	597

Table 3 Precision, recall and *F1-Score* of the ResNet18 network

ResNeXt-50

With a f1-score of 85% the ResNeXt-50 network had the second highest score. When predicting the benign data, it only reached a precision of 76% and a recall of 68%. The prediction for the malignant data works a lot better with a precision of 88% and recall of 92% (see Table 4 on the facing page).

	precision	recall	f1-score	support
B	0.76	0.68	0.72	168
M	0.88	0.92	0.90	429
weighted avg	0.85	0.85	0.85	597

Table 4 Precision, recall and *F1-Score* of the ResNeXt50 network

Wide ResNet-50

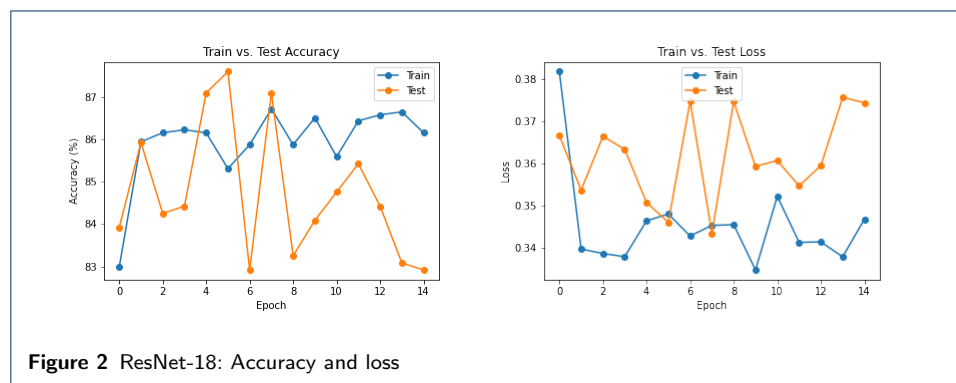
The wide ResNet-50 reached only a *F1-Score* of 74% and performed the worst out of our three residual networks. The prediction for the benign tissue images reached only a precision of 51% and a recall of 82%. For the malignant it looks mirrored. The prediction for malignant tissue samples yields a precision of 91%, which is the highest one of all tree networks, but only a recall of 69% (see Table 5).

	precision	recall	f1-score	support
B	0.51	0.82	0.63	168
M	0.91	0.69	0.78	429
weighted avg	0.80	0.73	0.74	597

Table 5 Precision, recall and *F1-Score* of the Wide ResNet50 network

6 Discussion

Analyzing the performance of the three residual models the first thing what can be compared is the accuracy and loss function on the training vs. the test data over the different epochs. The ResNeXt-50 and wide ResNet-18 are both getting a higher accuracy and a smaller loss on the training data then on the test data (see appendix figure 3 on the following page & figure 4 on the next page). This is a sign of overtraining, which fits in well with the fact that the two networks are significantly more complicated than ResNet-18. The latter also seems to overtrain in some epochs, but on the other hand the performance on the test data was even better then on the training data (see Figure 2). It seems that simpler model architectures suffice.



Improvements

The original data set contains pictures with four different magnifying factors. Since the 200X images performed best in the original paper we also used them to train

our networks. One way to overcome the over training problem could be to use instead all the pictures. Another ratio of training and test data set size (80 / 20 split) could also lead to more stable training results.

Fixing the class imbalance problem would decrease the amount of overtraining for only one class, because in the training phase the model would see the same number of images of benign and malignant tissue.

Fine tuning the hyperparameters can always help to improve the performance of a model. Our learning rate is quite small, but tests with 0.1 instead of 0.005 produced a worse performance. We would suggest to further tune the parameters. Batch size and number of epochs can also be modified, but these might not lead to an improved performance.

7 Appendix

Workload distribution:

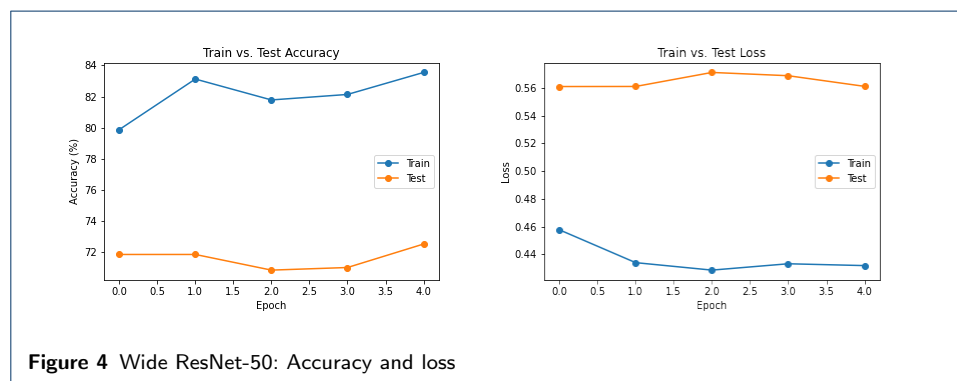
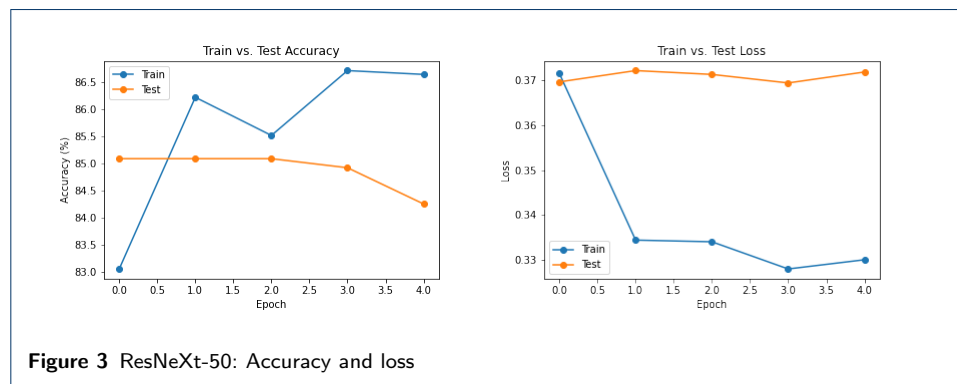
Dominik Bannwitz (Bioinformatics) - contributed 16 hours

Florian Herzler (Bioinformatics) - contributed 16 hours

Maximilian Otto (Bioinformatics) - contributed 16 hours

Mariana Steffens (Data Science) - contributed 6 hours

Figures



References

1. Spanhol, F.A., Oliveira, L.S., Petitjean, C., Heutte, L.: A Dataset for Breast Cancer Histopathological Image Classification. *IEEE Transactions on Biomedical Engineering* **63**(7), 1455–1462 (2016). doi:[10.1109/TBME.2015.2496264](https://doi.org/10.1109/TBME.2015.2496264). Accessed 2022-11-18
2. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E.Z., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. *CoRR abs/1912.01703* (2019). [1912.01703](https://arxiv.org/abs/1912.01703)
3. He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. arXiv. Number: arXiv:1512.03385 arXiv:1512.03385 [cs] version: 1 (2015). doi:[10.48550/arXiv.1512.03385](https://doi.org/10.48550/arXiv.1512.03385). <http://arxiv.org/abs/1512.03385> Accessed 2022-11-18
4. Zagoruyko, S., Komodakis, N.: Wide Residual Networks. arXiv. Number: arXiv:1605.07146 arXiv:1605.07146 [cs] version: 4 (2017). doi:[10.48550/arXiv.1605.07146](https://doi.org/10.48550/arXiv.1605.07146). <http://arxiv.org/abs/1605.07146> Accessed 2022-11-18
5. Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated Residual Transformations for Deep Neural Networks. arXiv. Number: arXiv:1611.05431 arXiv:1611.05431 [cs] version: 2 (2017). doi:[10.48550/arXiv.1611.05431](https://doi.org/10.48550/arXiv.1611.05431). <http://arxiv.org/abs/1611.05431> Accessed 2022-11-18
6. Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition, pp. 248–255 (2009). doi:[10.1109/CVPR.2009.5206848](https://doi.org/10.1109/CVPR.2009.5206848)