

Table of content

Reference to the code..... 1

1. Contained files and folders..... 1

2. Prerequisites for using the given code..... 2

 2.1 Needed packages 2

 2.2 Prerequisites to images 3

3. Analyze foci images with a pre-trained model 3

 3.1 Setting parameters..... 3

 3.2 Parameters used for the provided pre-trained image 4

 3.3 Analyzing images and output files 5

4. Train model on labeled images 5

 4.1 Setting parameters..... 5

 4.2 Setting the model 7

 4.3 Train model and output files..... 8

5. License 8

Reference to the code

If the code provided here or modifications of it are used for scientific publications please refer to it as published in: T. Hohmann et al, Evaluation of machine learning models for automatic detection of DNA double strand breaks after irradiation using a γ H2AX foci assay, PLOS One, 2020

1. Contained files and folders

The following folders are provided:

- 1) Application
- 2) Training
- 3) Sample Images

The folder “Application” contains 3 “.py” files called AnalyseImage.py, FociDetctImgAnalysis.py and GetFociData.py. FociDetctImgAnalysis.py and GetFociData.py contain helper functions for the main file AnalyseImage.py, which is used for analyzing images with a pre-trained model. Furthermore, this folder contains a sub-folder called “Trained Model”, containing all used scalings and the pre-trained model used in T. Hohmann et al, Evaluation of machine learning models for automatic detection of DNA double strand breaks after irradiation using a γ H2AX foci assay, PLOS One, 2020. The model file is called “MLP_Reduced_Trained”.

The folder “Training” contains 4 “.py” files called TrainModels.py, FociDetctImgAnalysis.py, GetBestParams.py and GetFociData.py. FociDetctImgAnalysis.py, GetBestParams.py and GetFociData.py contain helper functions for the main file TrainModels.py, which is used for training models on previously labeled foci images.

The folder “Sample Images” contains 2 identical images of nuclei, foci and manually marked foci, as well as all sub-folders and files created during model-training and application on these two files (see next chapters).

2. Prerequisites for using the given code

2.1 Needed packages

The code provided was written in Python 3 and relies on the following Python packages:

- NumPy
- Pandas
- SciPy
- os
- sys
- pickle
- csv
- scikit-learn package v0.20.3
- sci-kit image package v0.15

2.2 Prerequisites to images

To use the given code several pre-requisites for the images have to be fulfilled. The code in its current form is accepting “.tif”, “.png”, “.jpg” and “.bmp” files that can either be rgb or grayscale images. As the code reads any image of the given format in the given folder it should not contain additional unwanted images. Additional sub-folders do not pose obstacles. Image stacks are not supported.

For correct working of the code, images of foci and the nucleus have to be provided.

3. Analyze foci images with a pre-trained model

The provided code allows users to apply a pre-trained model to their foci images. As the given approach reduces analysis to a region of interest defined via the nuclear region, images showing the nucleus need to be provided. For starting the procedure open the AnalyseImages.py file in the “Application” folder with appropriate software (e.g. Spyder).

Together with the source code two identical example images, their manual labeling, as well as all results created by training a model and applying it to those images are included.

3.1 Setting parameters

In lines 16-21 of the source code you can specify the folders containing the images of the foci (line 17), the nuclei (line 19) and the folder specifying the path to the model data (line 21). In line 23 the exact name of the model file needs to be given. If the provided pre-trained model is used, input the name “MLP_Reduced_Trained”. This is also illustrated in Figure 1.

In lines 25-35 specific image analysis parameters are set. In line 27 a minimal size in pixel for nuclei is set. Objects in nuclei images smaller than this threshold get removed. In line 29 the same property is set for the foci images. In lines 32 and 35 the color channel for the nuclear or foci images are set, assuming those images are rgb. The value 0 corresponds to the red channel, 1 to the green channel and 2 to the blue channel. For gray scale images these properties are not used and can be chosen arbitrarily. This is also illustrated in Figure 1.

In lines 37-44 properties and sizes of filters are set. Please refer to “T. Hohmann et al, Evaluation of machine learning models for automatic detection of DNA double strand breaks after irradiation using a γ H2AX foci assay, PLOS One, 2020” and the scikit-learn documentation for more information. This is also illustrated in Figure 1.

Line 48 introduces a rescale factor, increasing or reducing image size in each dimension by the given factor. This might be useful to adopt a model to images of different resolution to reach a roughly constant foci size between test and training images. Notably, the higher the rescale factor the higher the computation time.

```
1 #-*- coding: utf-8 -*-
2 """
3 Created on Fri Mar 29 10:39:10 2019
4
5 @author: Tim Hohmann et al. - "Evaluation of machine learning models for
6 automatic detection of DNA double strand breaks after irradiation using a gH2AX
7 foci assay", PLOS One, 2020
8 """
9
10
11 # main file to analyse new images using a pre trained model
12
13 #####
14 # Parameters and file path that have to be set manually:
15
16 # directory containing the foci images:
17 im_path_foci = "D:\\Sample Images\\foci"
18 # directory containing the dapi images:
19 im_path_dapi = "D:\\Sample Images\\dapi"
20 # path with the information for the trained model:
21 model_path = "D:\\Sample Images\\trained Models"
22 # file name of the model
23 model_name = "MLP_Trained"
24
25 # Parameters:
26 # min area of nucleus
27 min_area = 4000
28 # min foci area:
29 min_area_foci = 16
30 # color channel of nucleus. 0 = red, 1 = green, 2 = blue. for grayscale images
31 # this value is ignored.
32 nuc_chan = 2
33 # color channel of foci. 0 = red, 1 = green, 2 = blue. for grayscale images
34 # this value is ignored.
35 foci_chan = 1
36
37 # specify filter sizes. for best results these should be identical to the ones
38 # used for model training.
39 # used filter sizes
40 filt_range = [2,10,30]
41 # scaling range for frangi filter
42 sc_range = list(range(2,11))
43 #frequency range for gabor filter
44 freq = [0.08,0.16,0.2]
45
46 # adjust image size (should be adjusted to match resolution of training images)
47 #image rescale factor:
48 rescale_factor = 1.8
49
```

Figure 1: Sample parameter setting for model application.

3.2 Parameters used for the provided pre-trained image

For the pre-trained model provided together with the source code the following settings were used:

`filt_range = [2,10,30]`

`sc_range = list(range(2,5))`

freq = [0.08,0.16,0.2]

Pixelsize of the images in μm : 1px \approx 0.12 μm .

3.3 Analyzing images and output files

After setting all these parameters the program can be run and will analyze all images. Thereby, it will create a sub-folder called “Single Nuclei” within the folder containing the nucleus images. This folder contains binary images showing the nucleus segmentation.

Furthermore, a sub-folder named “Model Name”_Trained_Results will be created containing the foci images – with foci now displayed in green – with the respective classification of each pixel, shown in white. Additionally, a .csv file named “Model Name”_Trained_Results.csv is created which contains the number of pixels covered by foci relative to the nucleus area for each nucleus and image.

4. Train model on labeled images

The provided code allows users to train a model using manually labeled foci images. As the given approach reduces analysis to a region of interest defined via the nuclear region, images showing the nucleus need to be provided. For starting the procedure open the AnalyseImages.py file in the “Application” folder with appropriate software (e.g. Spyder).

Together with the source code two identical example images, their manual labeling, as well as all results created by training a model and applying it to those images are included.

4.1 Setting parameters

In lines 45-50 of the source code you can specify the folders containing the images of the foci (line 46), images of the manually labeled foci (line 48) and the nuclei (line 50). In line 49 the name of the model file is specified, which is later used for savings. This is also illustrated in Figure 2.

In lines 12-22 specific image analysis parameters are set. In line 13 a minimal size in pixel for nuclei is set. Objects in nuclei images smaller than this threshold get removed. In lines 16 and 19 the color channel for the nuclear or foci images are set, assuming those images are rgb. The value 0 corresponds to the red channel, 1 to the green channel and 2 to the blue channel. For gray scale images these properties are not used and can be chosen arbitrarily. In line 22 the color used for the manual markings is defined. As this is fixed to be red, green or blue it is highly recommended to use one of these colors for the markings. This is also illustrated in Figure 2.

Line 27 introduces a rescale factor, increasing or reducing image size in each dimension by the given factor. This property is useful for downscaling images of very high resolution to improve calculation time. Line 34 can also be used to speed up the training process, as the sampling property sets what proportion of the data of each image is used for training. This is especially useful when using many training images to speed up the training process. Notably, both rescaling and sampling generally affect the accuracy of model training negatively and thus these parameters have to be chosen with care.

Another strategy employed here for improving training time is that of dimension reduction. Therefore, linear principal component (PCA) analysis is used. After PCA the obtained principal components are linearly de-correlated and higher order components thus do not add significant “new” information (in the linear PCA sense) and thus can likely be omitted. The parameter in line 31 sets the number of kept principal components by defining the cumulative variance (“information”) they contain. Higher numbers keep more components, lower less. Consequently, reducing this value speeds up the training process, but may lead to a loss of information.

In lines 36-44 properties and sizes of filters are set. Please refer to “T. Hohmann et al, Evaluation of machine learning models for automatic detection of DNA double strand breaks after irradiation using a γ H2AX foci assay, PLOS One, 2020” and the scikit-learn documentation for more information. This is also illustrated in Figure 2.

```

1 #- coding: utf-8 -*-
2 """
3 Created on Sun Feb 24 19:08:39 2019
4 @author: Tim Hohmann et al. - "Evaluation of machine learning models for
5 automatic detection of DNA double strand breaks after irradiation using a gH2AX
6 foci assay", PLOS One, 2020
7 """
8 # main file for training machine learning models using previously labeled data
9 #####
10 # Parameters and file path that have to be set manually:
11 # Parameters:
12 # min area of nucleus
13 min_area = 4000
14 # color channel of nucleus. 0 = red, 1 = green, 2 = blue. for grayscale images
15 # this value is ignored.
16 nuc_chan = 2
17 # color channel of foci. 0 = red, 1 = green, 2 = blue. for grayscale images
18 # this value is ignored.
19 foci_chan = 1
20 # color channel of marked image. 0 = red, 1 = green, 2 = blue. corresponds to
21 # the color of the markings in the manually labeled foci images
22 mark_chan = 0
23
24 # adjust image size - might be useful to save calculation time. needs to be
25 # identical for foci and nucleus images
26 # image rescale factor:
27 rescale_factor = 1.0
28
29 # take only those PCA components cumulatively explaining var_max of the variance
30 # 0 < var_max <= 1.
31 var_max = 0.95
32 # randomly sample a proportion of the training data from each image (0 < sampling <= 1).
33 # speeds up training process if smaller than 1
34 sampling = 1
35
36 # used filter sizes
37 filt_range = [2,3,4,5,8,10,15,20,25,30,35]
38 # scaling range for frangi filter
39 sc_range = list(range(2,11))
40 # frequency range for gabor filter
41 freq = [0.08,0.10,0.13,0.16,0.2]
42
43 # Name used for saving the trained model and related images:
44 model_name = "MLP"
45 # directory containing the foci images:
46 im_path_foci = "D:\\Sample Images\\foci"
47 # directory containing the labeled images:
48 im_path_foci_marked = "D:\\Sample Images\\foci marked"
49 # directory containing the nucleus images:
50 im_path_dapi = "D:\\Sample Images\\dapi"
51

```

Figure 2: Sample parameter setting for model training.

4.2 Setting the model

In the current version of the code a multi-layer perceptron model is used for training. To change this go to lines 137 and the following (Figure 3) and change it accordingly. The variable “parameters” describes multiple properties of the model that can be set and changed. In the current version it defines multiple “alpha-values” and “batch sizes”. The values set there are tested using a grid search (line 162/163) to find best suited settings.

Additional code for sample models – all part of the scikit-learn package – are also provided and just need to be uncommented (e.g. random forest in lines 141-146, Figure 3). For more information on available models, their parameters etc. the reader is referred to the scikit-learn documentation.

```

129 #####
130 # This part is for model training assuming "get_labeled_data" was ran successfully
131
132 # get trainingsdata:
133 x_train_transformed, y_train, idx, s1, s2, p, mnm = image_pipeline(x_data, y_data, removed_im = [], var_max = var_max, sampling = sampling)
134
135 # Chose the model to train:
136 # neural network:
137 model = MLPClassifier(alpha=0.1, batch_size = 2000, learning_rate = "adaptive", \
138                       learning_rate_init = 0.1, max_iter = 300, tol = 10**-4, early_stopping = True)
139 parameters = {'batch_size': [100, 500, 1000, 2000, 3000, 4000, 5000], 'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 1]}
140
141 # random forest:
142 clf = RandomForestClassifier(criterion = "entropy", min_weight_fraction_leaf = 0.005, \
143                             n_estimators = 15, max_depth = 50, min_samples_leaf = 10, min_samples_split = 100, n_jobs = -1)
144 # model = AdaBoostClassifier(base_estimator = clf, n_estimators=5)
145 # parameters = {'base_estimator__min_weight_fraction_leaf': [0.0001, 0.001, 0.005], \
146 #               'base_estimator__n_estimators': [5, 10, 15, 20], 'base_estimator__min_samples_leaf': [10, 20, 100]}
147
148 # complement naive bayes:
149 clf = ComplementNB(alpha = 0.0, norm = True)
150 # model = AdaBoostClassifier(base_estimator = clf, n_estimators=15)
151 # parameters = {'base_estimator__alpha': [0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06], 'base_estimator__norm': [True, False]}
152
153 # support vector machine:
154 # linear svm
155 clf = LinearSVC(penalty = "l2", Loss = "hinge", C = 2, class_weight = "balanced", max_iter = 5000)
156 # model = AdaBoostClassifier(base_estimator = clf, n_estimators=5, algorithm="SAMME")
157 # parameters = {"base_estimator__C": [0.1, 0.3, 0.6, 1, 2, 3]}
158
159
160
161 # get best model parameters:
162 clf = GridSearchCV(model, parameters, cv = 3)
163 clf.fit(x_train_transformed, y_train)
164

```

Figure 3: Sample model for model training.

4.3 Train model and output files

After setting all these parameters the program can be run and will analyze all images. First the model will be trained on all but one images and tested on the remaining one. Afterwards, it will create a sub-folder called “Results Model Validation” within the folder containing the foci images, containing classified test images, with identified foci marked in white. Furthermore, an image_statistics.txt file is created, containing the image number, the confusion matrix, mean accuracy of classification, precision and recall for each image tested with the leave-one-out strategy.

As a last step the model is trained on the whole dataset and a folder called “Trained Models” containing a folder “Model Name” is created. Inside the folder all necessary model and scaling data is stored. Notably, the saved model is called “Model_name”_Trained.

5. License

Copyright (c) 2019, Tim Hohmann

Kommentar [OTM1]: @Faramarz: Wen muss man da eigentlich angeben?

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.