



Diabetic Retinopathy Detection

Detect diabetic retinopathy to
stop blindness before it's too
late





TABLE OF CONTENTS

01.

Definition

- 1.1 Project Overview
- 1.2 Dataset and Inputs
- 1.3 Problem Statement

02.

Analysis

- 2.1 Data Exploration
- 2.2 Exploratory Visualization

03.

Methodology

- 3.1 Data Preprocessing
- 3.2 Algorithms and Techniques
- 3.3 Evaluation Metrics

04.

Implementation

- 4.1 Pre-training
- 4.2 Fine-tuning

05.

Results

- Summary
- Explain Results

06.

References

- Medical Reference
- Technical Reference





01.

Definition

- 1.1 Project Overview
- 1.2 Dataset and Inputs
- 1.3 Problem Statement





Project Overview

Diabetic retinopathy (DR) is one of the leading causes of vision loss. According to a recent study from International Diabetes Federation, the global prevalence of DR among the individuals with diabetes for the period from 2015 to 2019 was at more than 25%. According to the World Health Organization, more than 300 million people worldwide have diabetes, and the disease prevalence has been rising rapidly in developing countries.



Project Overview

Early detection and treatment are crucial steps towards preventing DR. Currently, detecting DR is a time-consuming process. The screening procedure requires a trained clinical expert to examine the fundus photographs of the patient's retina. This creates delays in diagnosis and treatment of the disease. Automated evaluation of retina photographs can speed up the efficiency and coverage of the DR screening programs. This is especially relevant for developing countries, which often lack qualified medical stuff to perform the diagnosis.

The project aims at developing a deep learning model for predicting the severity of DR disease based on the patient's retina photograph.

Previous research has explored the usage of deep learning for detecting DR and concluded that convolutional neural networks (CNNs) have high potential in this task . The Asia Pacific Tele

Ophthalmology Society (APTOS) has launched two Kaggle competitions with a goal of promoting the use of deep learning for DR detection and boosting the development of automated detection systems.



Data sets and Inputs

The project leverages two data

Main Data

Main data set used for modeling and evaluation our model.

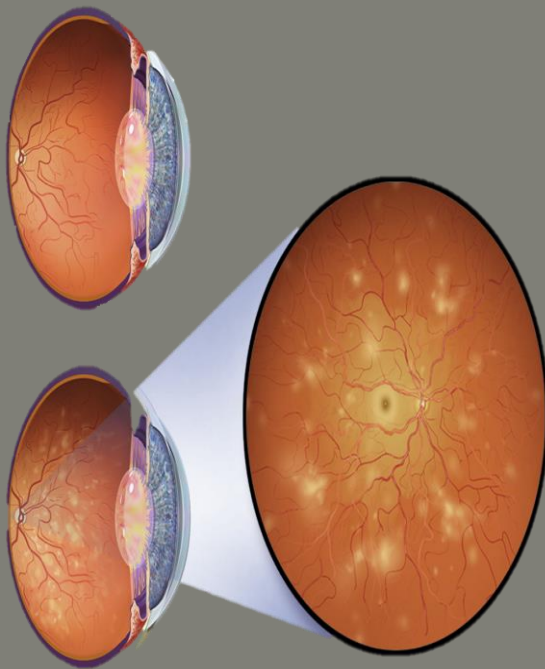
Supplementary data

Supplementary data set for pre-training model before use it on the main data.



Main Data set

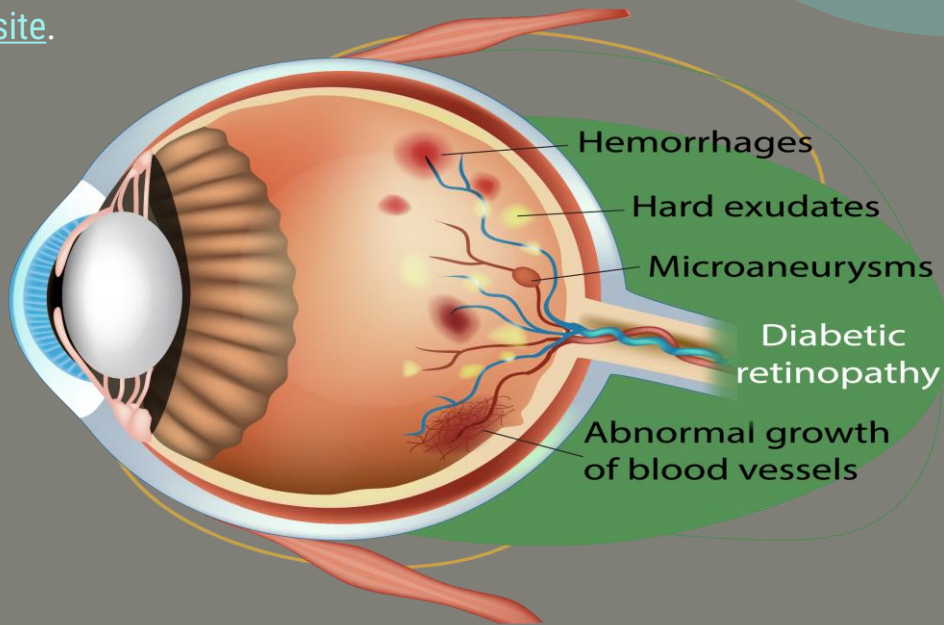
The main data set is provided by APTOS. It has been employed in the APTOS 2019 Blindness Detection competition on Kaggle and is available for the download at the [competition website](#).



The data set includes 3,662 labeled retina images of clinical patients. The images are taken using a fundus photography technique. The images are labeled by a clinical expert. The labels indicate the severity of DR from 0 to 4 in accordance with the scale provided in Section 1.2. The data also include a test set with 1,928 images whose labels are not explicitly disclosed.

Supplementary data

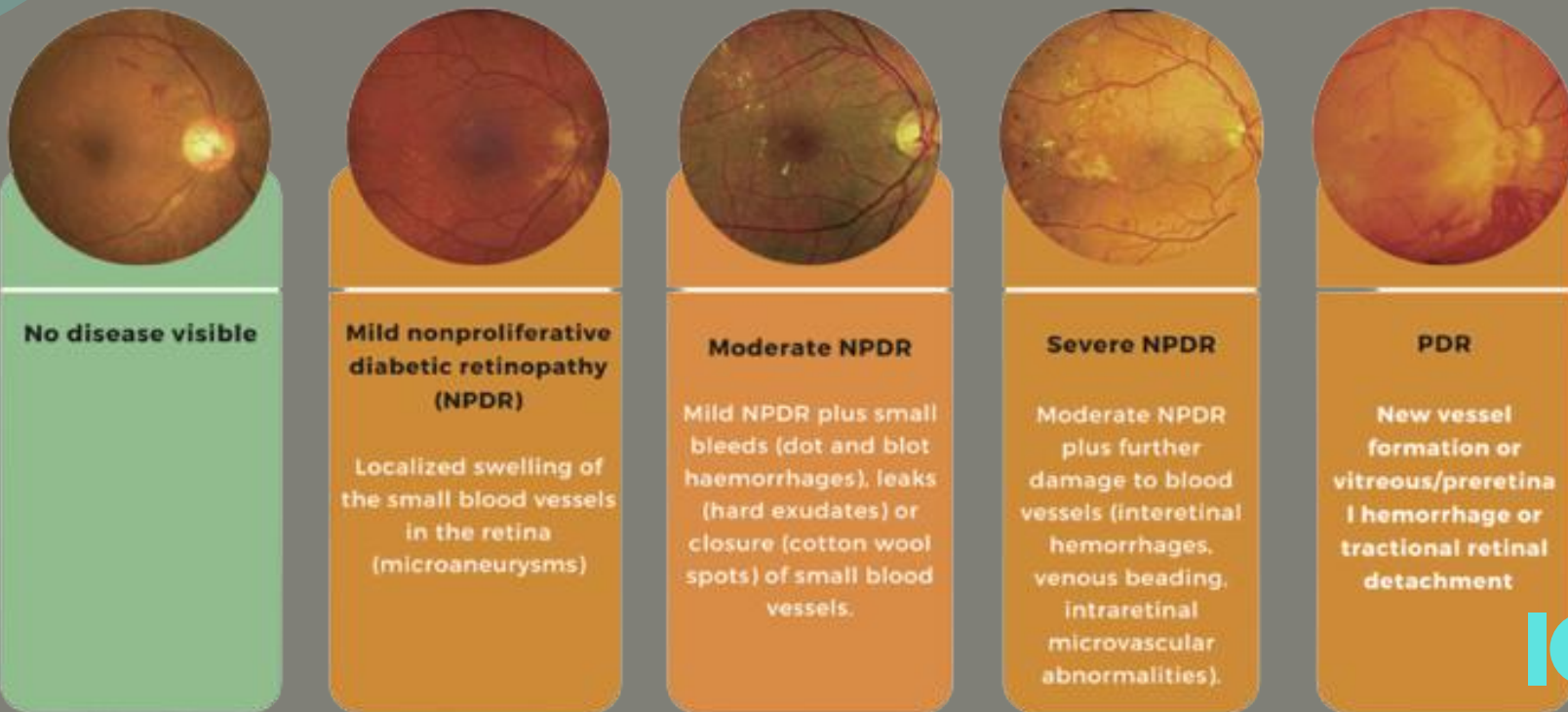
The supplementary data set features 35,126 retina images labeled by a clinician using the same scale as the main data set. The data set has been used in the 2015 Diabetic Retinopathy Detection competition and is available for the download at the [competition website](#).



Problem Statement



The project aims at developing a deep learning model for predicting the severity of the DR disease based on the patient's retina photograph. Severity of DR determines the necessary actions for prevention and/or required treatment. Therefore, we consider the multiple classification task, in which five disease stages are distinguished based on its severity:





02. *Analysis*

- 2.1 Data Exploration
- 2.2 Exploratory Visualization



Data Exploration

The data set is collected from multiple clinics in India using a variety of different camera models, which creates discrepancies in the image resolution, aspect ratio and other parameters. This is demonstrated in Figure 1 below, where we plot the histograms of image width, height and aspect ratio. A high variety in the image parameters requires to be accounted for during the data preprocessing. All retina images are colored images with three RGB channels

Analysis

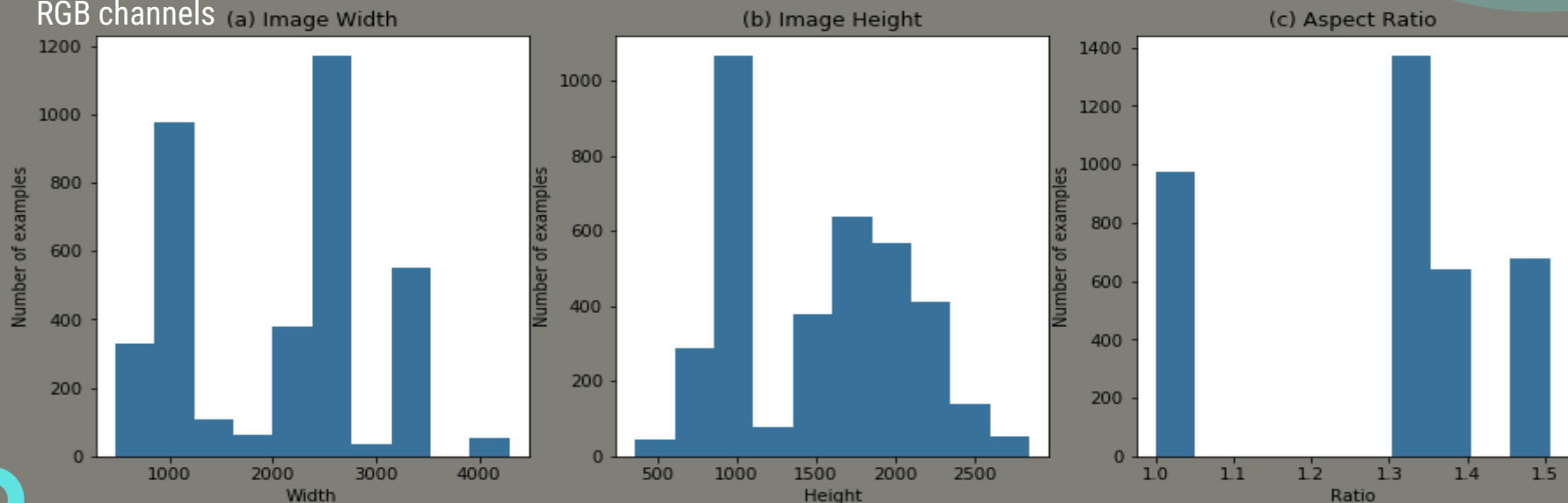


Figure 1. Image Size Distribution

Data Exploration

The data set is imbalanced. Figure 2 illustrates the class proportions in the labeled data. 49% of images are healthy patients. The remaining 51% are pictures with different stages of DR. The least common class is 3 (severe stage) with only 5% of the total examples. The overall sample size of the labeled data set is 3,662 images.

Class Distribution

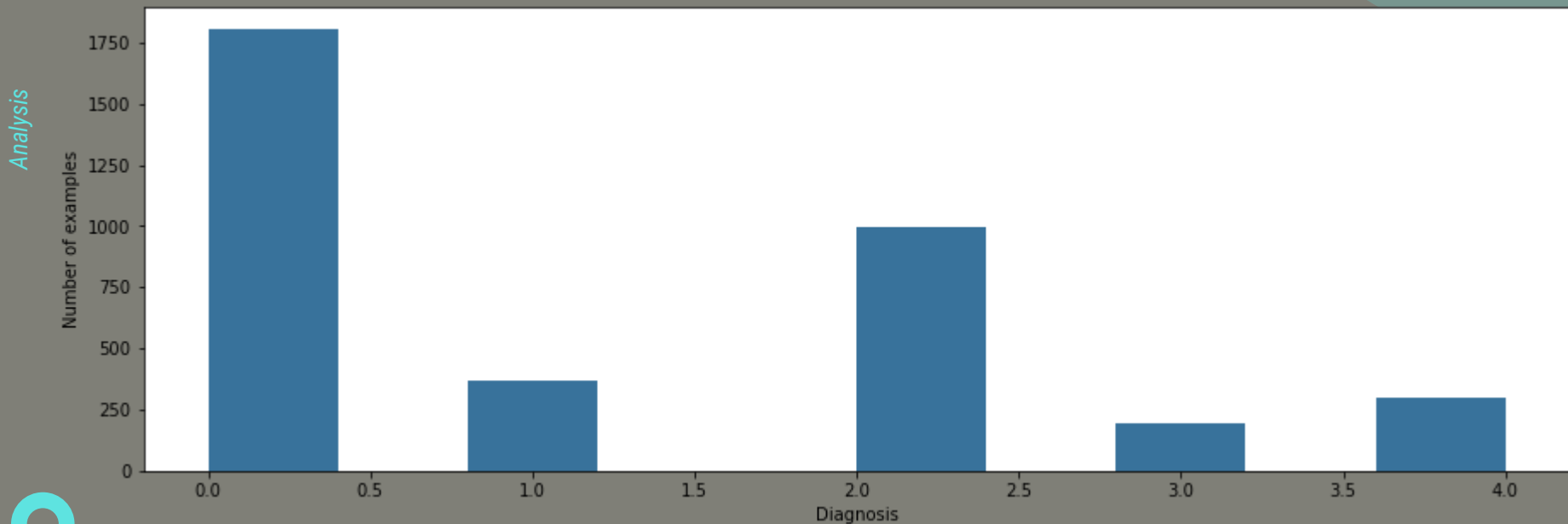


Figure 2. Class Distribution

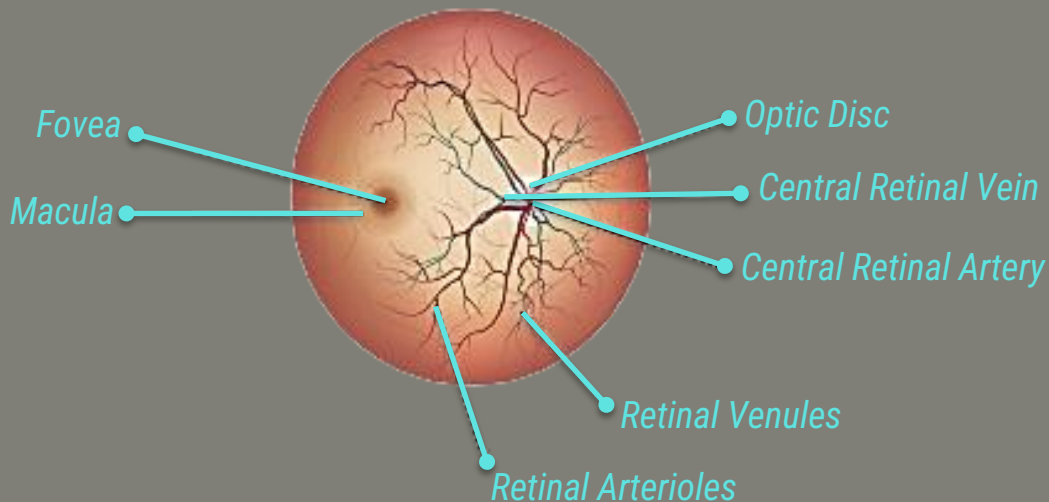


Exploratory Visualization

In practice, the severity of DR is diagnosed by the presence of different visual cues on the retina photographs. This includes signs like abnormal blood vessels, hard exudates and so-called «cotton wool» spots. Figure 3 illustrates the exemplary retina photograph with present DR signs used to guide doctors.

Analysis

Normal Retina



Diabetic Retinopathy

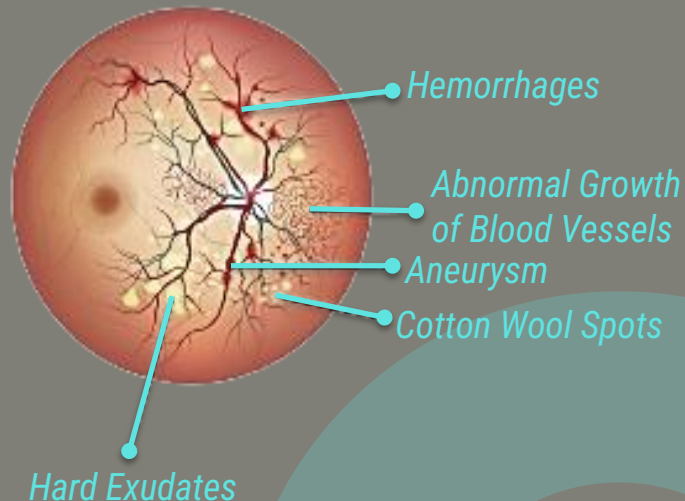


Figure 3. Illustration of the Visual Signs of Normal, DR

Exploratory Visualization

Figure 4 visualizes a batch of sample images from the main data set. The labels are provided in the image titles. Comparing images of different classes, we can see the presence of exudates and «cotton wool» spots on some of the retina photographs of sick patients. The illustration further emphasizes the difference between the retina photographs in terms of aspect ratio, lighting conditions and camera quality. We will discuss ways to tackle these discrepancies in Section 3.1.

Analysis

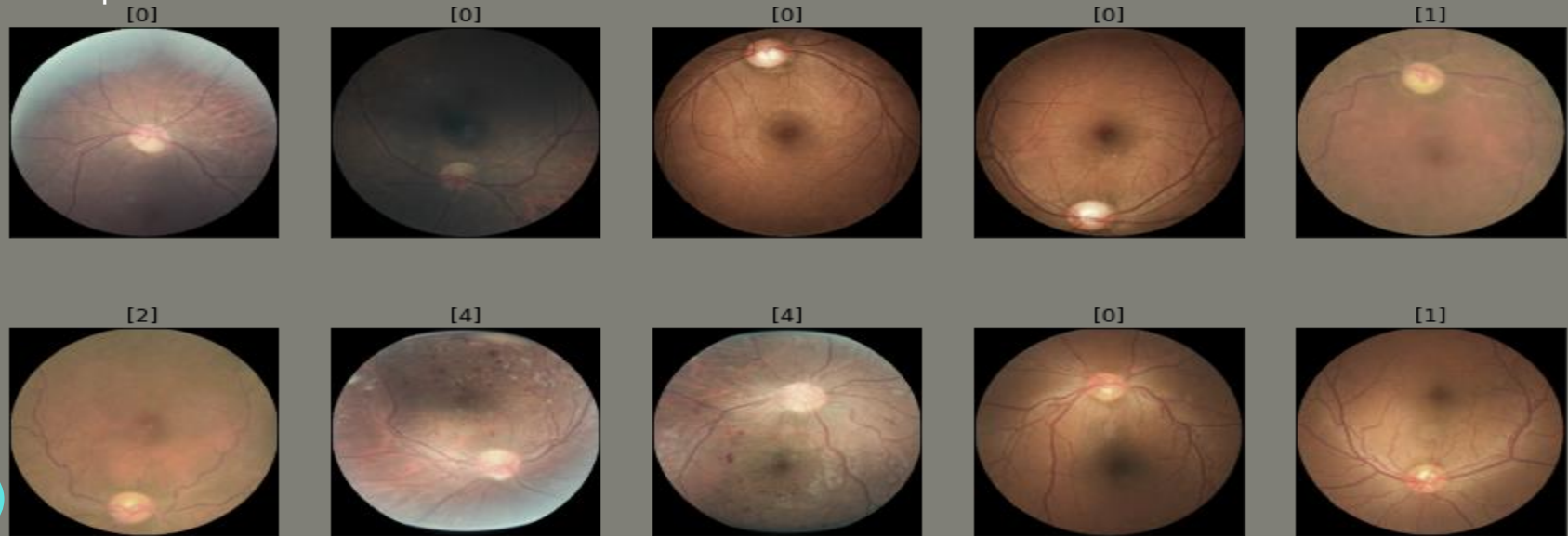


Figure 4. Data Illustration



03.

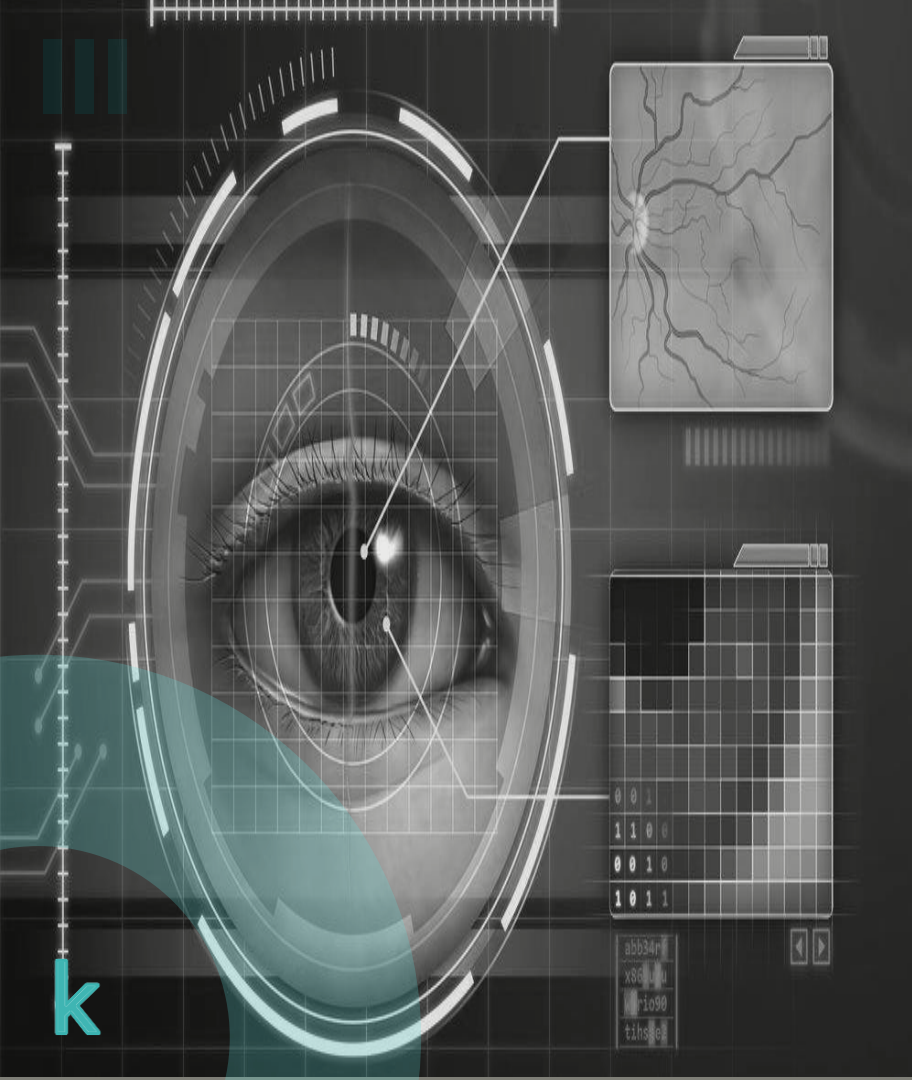
Methodology

3.1 Data Preprocessing

3.2 Algorithms and Techniques

3.3 Evaluation Metrics

Methodology



k

Data Preprocessing

As illustrated in Section 2, the retina photographs are characterized by a high variety of image parameters. This can have a strong impact on the performance of the classification model. We make multiple steps towards standardizing the images.

Visual inspection of retina images suggests that images taken from cameras with different aspect ratios result in some images having large black areas around the eye. The black areas do not contain information relevant for prediction and can be cropped. However, we see that the size of black areas varies from one image to another. To address this, we develop a cropping function that converts the image to grayscale and marks black areas based on the pixel intensity.

Next, we find a mask of the image by selecting rows and columns in which all pixels exceed the intensity threshold. This helps to remove vertical or horizontal rectangles filled with black like the ones observed in the upper-right image in Figure 4. After removing the black stripes, we resize the images to the same height and width.

Another issue is the eye shape. Depending on the image parameters, some eyes appear to have a circular form, whereas others look like ovals. As the size and shape of items located in the retina determine the disease severity, it is crucial to standardize the eye shape as well.

To do so, We develop another cropping function that makes a circular crop of a particular radius around the center of the image. Finally, we correct for the lightning and brightness discrepancies by smoothing the images using a Gaussian filter from the cv2 image preprocessing library. Figure 5 depicts the example batch of eye images after preprocessing. Comparing the images to the ones presented in Figure 4 we can see that the apparent discrepancies between the photographs are now fixed.

Data Preprocessing

To reduce overfitting and improve the generalization ability of the model, we apply data augmentations during model training. We use the following augmentations:

- Random rotation of the image by up to 360 degrees.
- Random horizontal flip of the image with a probability of 0.5.
- Random vertical flip of the image with a probability of 0.5.

Methodology

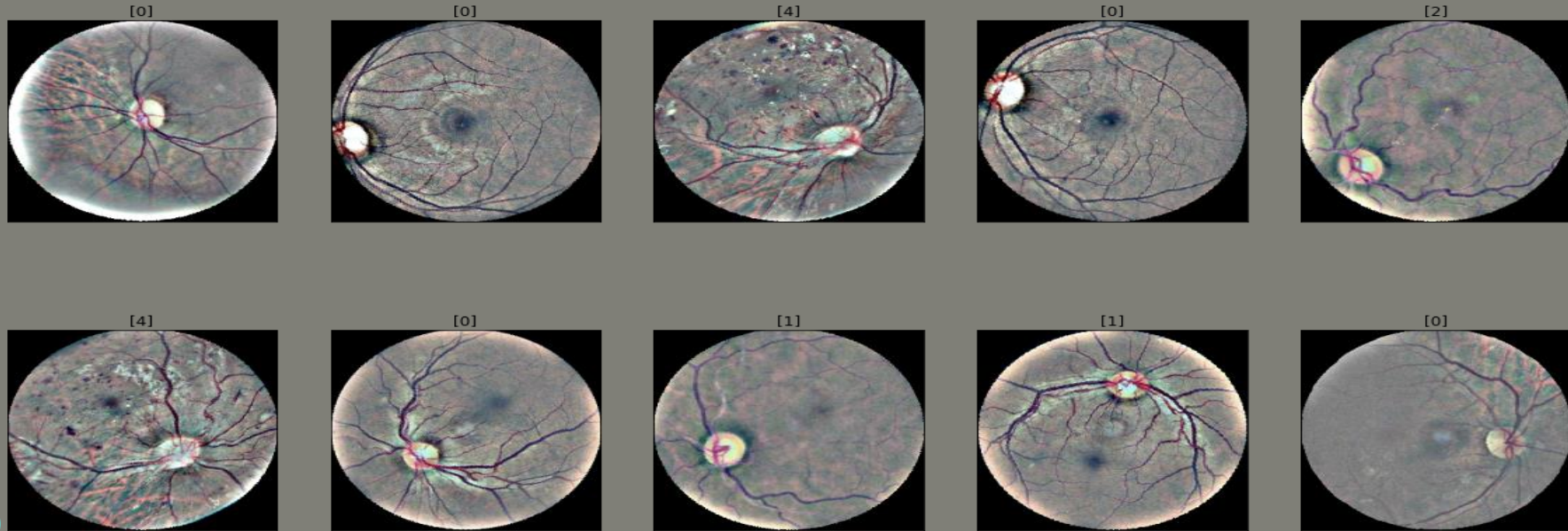


Figure 5. Data Illustration After Preprocessing



Algorithms and Techniques

This project develops a CNN-based deep learning model for the DR severity classification task. CNNs achieve state-of-the-art performance in computer vision tasks such as image classification. Recent medical research also shows a high potential of applying CNNs in the DR classification task. Furthermore, using CNNs allows to take advantage of transfer learning techniques, which is relevant given a small sample size of the considered data set. The project employs the EfficientNet architecture for the CNN classifier. EfficientNet is one of the recent state-of-the-art CNN models in image classification. EfficientNet encompasses 8 architecture variants (B0 to B7) which differ in the model complexity and default image size. The architecture of EfficientNet B7 is presented in Figure 6. We test multiple EfficientNet neural network architectures and use the one that demonstrates the best performance.

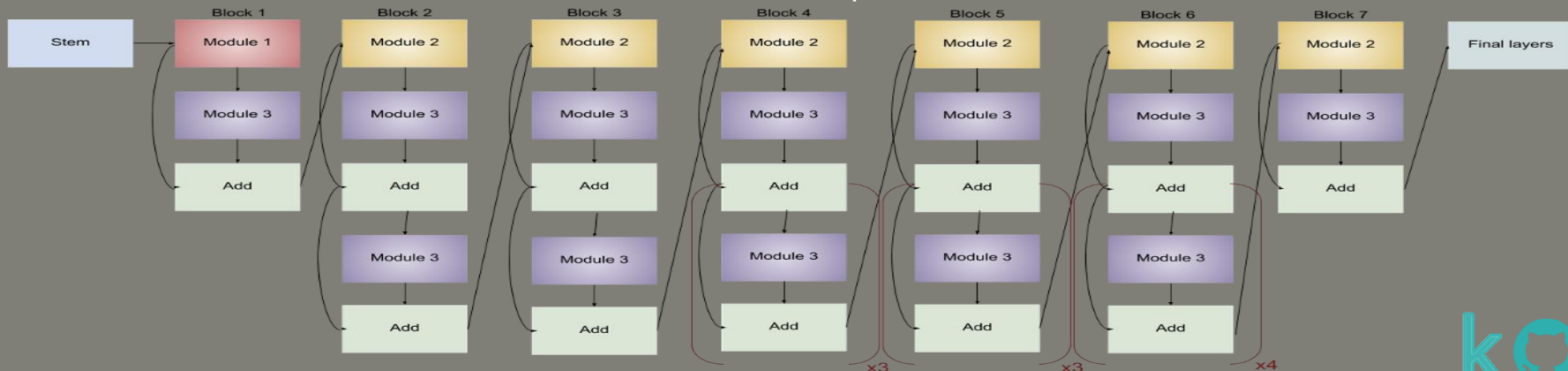


Figure 6. EfficientNet B7 Architecture.

Evaluation Metrics

The project considers DR detection as an ordinal classification task. In such a setting, a suitable evaluation metric is the Cohen's Kappa, which measures the agreement between the actual and predicted labels. The metric varies from 0 (random agreement) to 1 (perfect agreement). The Kappa is computed as follows.

First, one constructs three matrices:

- I. The matrix of observed scores X .
- II. The matrix of expected scores based on chance agreement M .
- III. The weight matrix W .

Next, the Kappa is computed as:

$$\kappa_w = \frac{\sum_{i=1}^k \sum_{j=1}^k w_{ij} x_{ij}}{\sum_{i=1}^k \sum_{j=1}^k w_{ij} m_{ij}}$$

where x_{ij} and m_{ij} are elements in the observed and predicted matrices, respectively, whereas the w_{ij} are weights. We will use Kappa with quadratic weights to stronger penalize larger errors.

04.

Implementation

The models are implemented in Python 3.7 using PyTorch library. All implementations are provided in Jupyter notebooks in the «codes» folder. This section discusses the implementation process of pre-training, fine-tuning and inference stages of the modeling pipeline.

[GitHub Repository](#) [Kaggle Notebook](#)



- 4.1 Pre-training
- 4.2 Fine-tuning
- 4.3 Inference

Pre-training

The first stage is pre-training a CNN model on the supplementary 2015 data set. First, we instantiate the EfficientNet B7 model. We initialize the model weights using the parameters pre-trained on ImageNet by downloading the model architecture and weights in the PyTorch format. The convolutional part of the network responsible for feature extraction outputs a tensor with 1792 features. To adapt the trained EfficientNet to our task, we replace the last fully-connected classification layer with a (1792, 5) fully-connected layer.

Since we are dealing with a multiple classification problem, we use cross-entropy as a loss function.

We implement `nn.CrossEntropyLoss()` function which combines logsoftmax and negative log-likelihood loss and applies them to the output of the last network layer. Next, we define optimizer for the network. We use adaptive learning rate optimization algorithm (Adam optimizer) with a starting learning rate of 0.001. During training, we use a learning rate scheduler, which multiplies the learning rate by 0.5 after every 5 epochs. This helps to make smaller changes to the network weights when we are getting closer to the optimum. Weights of all network layers are optimized during the pre-training. We create a dataloader object using a batch size of 20 and the image size of 256. The choice of these parameters is a trade-off between a predictive performance and resource capacity. During training, we iteratively load a batch of training examples from the supplementary data set and conduct forward and backward passes with the network. All data from the supplementary data set is used as the training sample.

After each training epoch, we validate the model on the validation examples. We extract class scores from the last fully-connected layer and predict the image class corresponding to the highest score. We use images from the main 2019 data set as a validation sample. This allows us to validate the network on the images from the main data set. We train the network for 15 epochs, tracking the validation loss and Cohen's kappa. If the kappa does not increase for 5 consecutive epochs, we terminate the training process and save model weights for the epoch which is associated with the highest validation kappa.

Pre-training

Figure 7 visualizes the training and validation loss over the training epochs as well as the Cohen's kappa on the validation sample. The results suggest that the cross-entropy loss on validation sample reaches a minimum already after 3 epochs and does not improve with further training. At the same time, kappa continues to increase up to the 15th epoch. Since we use kappa to evaluate the quality of our solution, we save model weights after 15 epochs.

Implementation

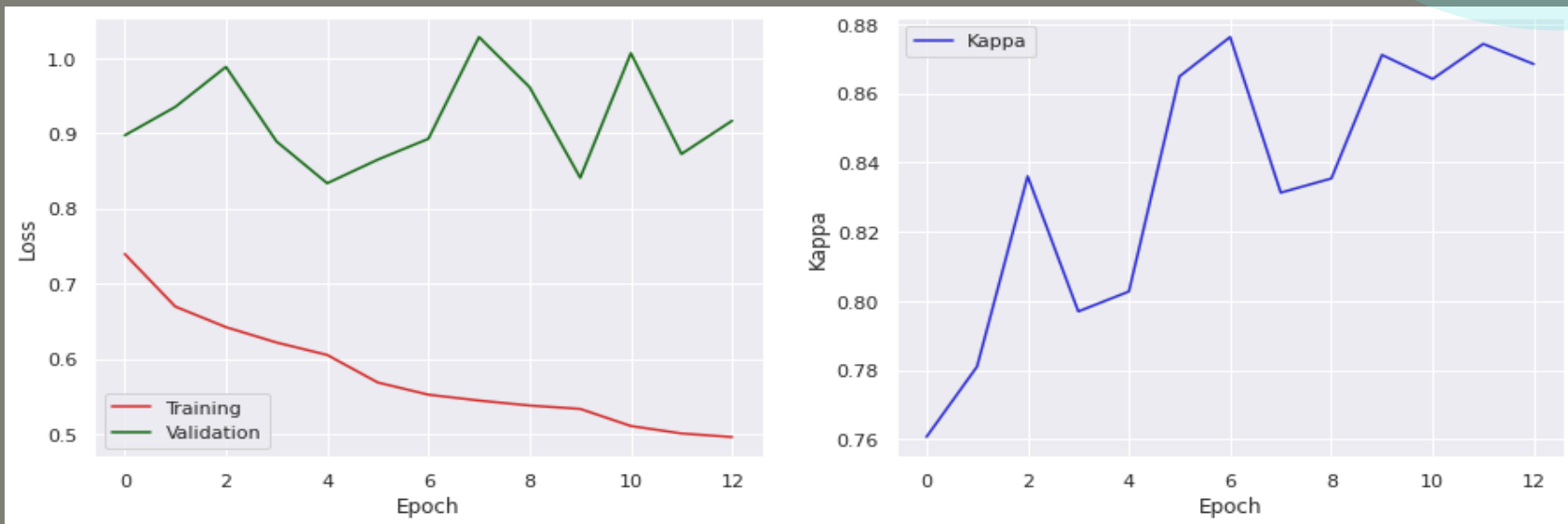


Figure 7. Pre-Training: Validation Performance Dynamics



Pre-training

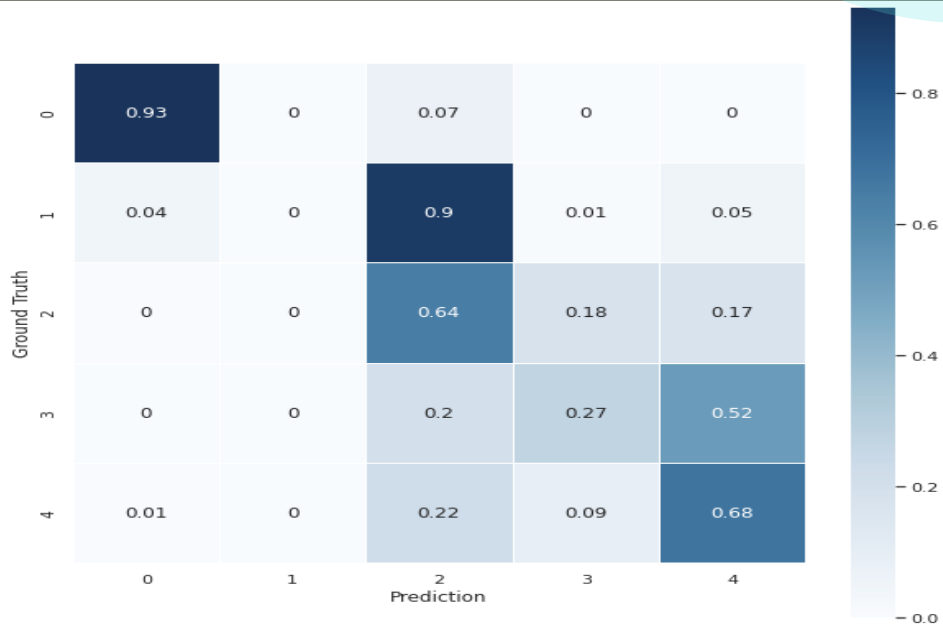


Figure 8 visualizes the confusion matrix of the trained model on validation set. The numbers in the table cells are percentages. According to the results, the model does a poor job in distinguishing the mild and moderate stages of DR: 90% of images with mild DR are classified as moderate DR by the model. The best performance 93% is observed for classifying the retina photographs of healthy people. Overall, we see that the model tends to confuse nearby severity stages but rarely misclassifies the proliferate and mild stages.

OOF loss = 1.0280
 Classification Report in Test: OOF kappa = 0.8313

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.99	0.93	0.96	1805
1	0.00	0.00	0.00	370
2	0.53	0.64	0.58	999
3	0.20	0.27	0.23	193
4	0.40	0.68	0.51	295
accuracy			0.70	3662
macro avg	0.43	0.51	0.46	3662
weighted avg	0.68	0.70	0.68	3662



Implementation
 Figure 8. Classification Report, Confusion Matrix

Fine-tuning

The fine-tuning on the main data set is performed using cross-validation with 4 folds. To ensure that we have enough examples of each class in the training sample, we perform cross-validation with stratification. On each iteration, we instantiate an EfficientNet B7 model with the same architecture as described in the previous section. Next, we load the saved weights from the model pre-trained on the supplementary data set. We freeze all network layers except for the last fully-connected classification layer. The weights in this layer are fine-tuned during training.

As on the pre-training stage, we use Adam optimizer with a learning rate of 0.001 and implement a learning rate scheduler which multiplies the learning rate by 0.5 after every 5 epochs. We also track the model performance on the validation fold and stop the training if kappa does not increase for 5 consecutive epochs. This process is repeated for each of the 4 folds, and the best model weights are saved for each combination of the training folds.

After running the cross-validation loop, we aggregate the model performance including validation loss and kappa over all out-of-fold examples from validation folds. This provides the estimate of the model performance. Based on the obtained kappa values, we test multiple variants of the model and iteratively refine the proposed solution. The details on the refinement process are provided in Section 5.

The results indicate that after pre-training, the model converges quickly during fine-tuning. The best validation performance is obtained after 3 to 7 training epochs depending on a fold.

Fine-tuning

This provides the estimate of the model performance. Based on the obtained kappa values, we test multiple variants of the model and iteratively refine the proposed solution. The details on the refinement process are provided in Section 5. The results indicate that after pre-training, the model converges quickly during fine-tuning. The best validation performance is obtained after 3 to 7 training epochs depending on a fold.

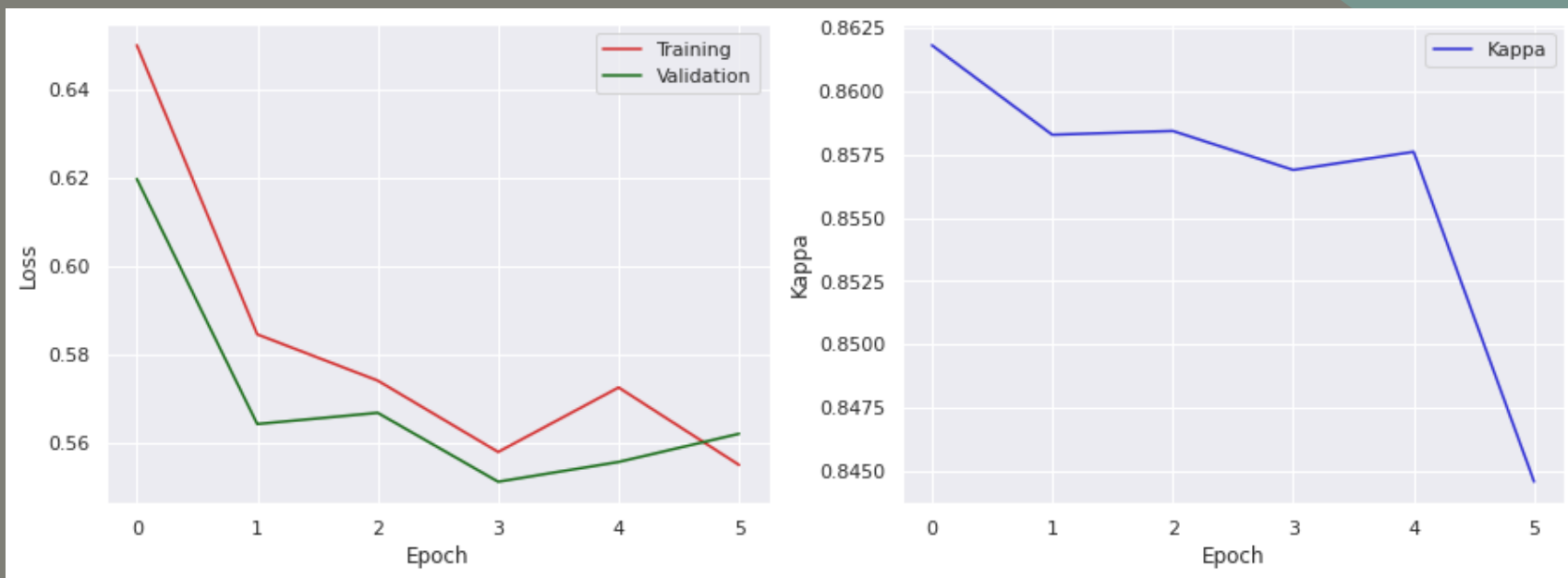


Figure 9. Fine-Tuning: Validation Performance Dynamics



Fine-tuning



Figure 10 visualizes the confusion matrix of the best model after fine tuning. The confusion matrix illustrates the advantages of the fine-tuned model over the pre-trained CNN and indicates a better performance in classifying moderate stages of the DR. However, we also observe that the model classifies too many examples as a mild stage of DR (class = 2). We try to address this issue on the inference stage by manipulating the thresholds.

Classification Report in Test: OOF kappa = 0.8742

	precision	recall	f1-score	support
0	0.97	0.99	0.98	1805
1	0.59	0.23	0.33	370
2	0.63	0.91	0.74	999
3	0.26	0.03	0.05	193
4	0.62	0.46	0.53	295
accuracy			0.80	3662
macro avg	0.61	0.52	0.53	3662
weighted avg	0.77	0.80	0.76	3662

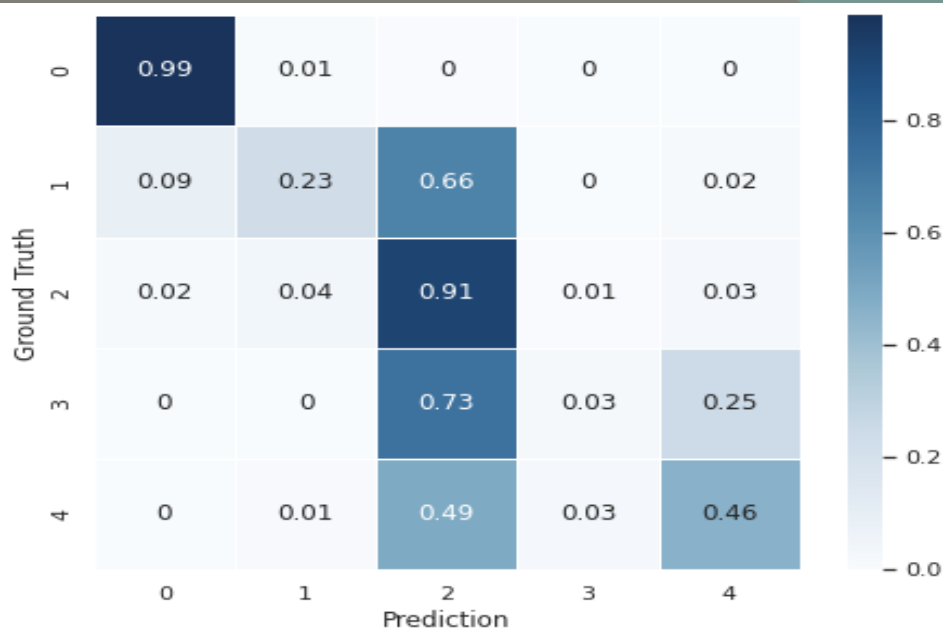


Figure 10. Classification Report, Confusion Matrix

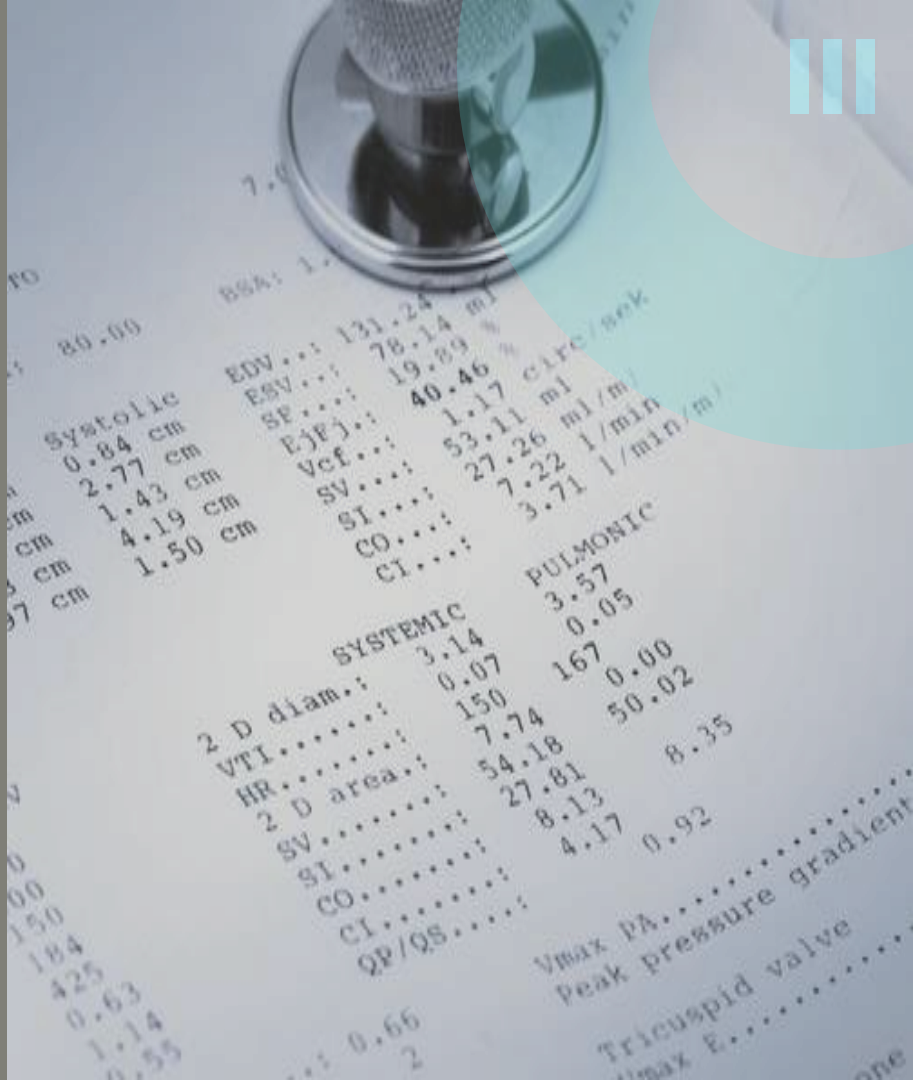
Implementation



05. Results

Results

The initial model is EfficientNet B7, where weights are initialized with the values pre-trained on ImageNet and no image preprocessing beyond resizing is applied. Fine-tuning this model on the 2019 data set by training the last classification layer achieves the average cross-entropy loss of 0.5703 and kappa score of 0.8546 on the validation folds. When applied to the test sample, the initial model achieves kappa of 0.8742, which is substantially worse compared to the benchmarks. During the work on the project, we iteratively test different refinements and track the model performance within cross-validation to make decisions about including certain modifications.



Results

The empirical experiments include:

- Tuning the parameters of image preprocessing functions. This includes cropping functions and smoothing functions, data augmentations and image size. The results of the empirical experiments suggest that cropping images and using augmentations is crucial to achieve a better predictive performance.
- Testing three different EfficientNet architectures: B0, B4, and B7. The B7 architecture demonstrates the best performance and selected as a final model architecture. At the same time, model B0 demonstrates the fastest training and inference, which might be useful when deploying a model for screening new retina photographs.
- Testing unfreezing of different layer combinations. We find that the best performance is achieved when all model weights are optimized during the pre-training stage and only the last classification layer weights are optimized during the fine-tuning stage. Unfreezing more layers during fine-tuning does not improve the performance as the sample size is too small.
- Selecting the optimizer and the learning rate. We try two optimization algorithms: Adam and stochastic gradient descent (SGD). According to the results, using Adam results in a better performance. In addition, we find that using a learning rate scheduler by iteratively decreasing the learning rate after every 5th epoch provide a slight performance improvement.

The final model is selected based on the out-of-fold Cohen's kappa score. The best model achieves a kappa score of 0.8742 on the validation folds.

06.

REFERENCES

10

- I. Jacob Cohen. A coefficient of agreement for nominal scales. Educational and psychological measurement, 20(1):37–46, 1960.
- II. Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. Jama, 316(22):2402–2410, 2016.
- III. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- IV. Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. arXiv preprint arXiv:1905.11946, 2019.
- V. RL Thomas, S Halim, S Gurudas, S Sivaprasad, and DR Owens. Idf diabetes atlas: A review of studies utilising retinal photography on the global prevalence of diabetes related retinopathy between 2015 and 2018. Diabetes research and clinical practice, page 107840, 2019.
- VI. Tien Yin Wong and Neil M Bressler. Artificial intelligence with deep learning technology looks into diabetic retinopathy screening. Jama, 316(22):2366–2367, 2016.
- VII. McHugh, Mary L. (2012). "Interrater reliability: The kappa statistic". Biochemia Medica. 22 (3): 276–282. doi:10.11613/bm.2012.031. PMC 3900052. PMID 23092060.