

# Full Accuracy of Machine Learning for Differentiation between Optic Neuropathies and Pseudopapilledema

**Jin Mo Ahn**

Soongsil University

**Sangsoo Kim**

Soongsil University

**Kwang-Sung Ahn**

PDXen

**Sung-Hoon Cho**

PDXen

**Ungsoo Kim** (✉ [ungsookim@kimeye.com](mailto:ungsookim@kimeye.com))

Kim's Eye Hospital <https://orcid.org/0000-0003-2373-6240>

---

## Research article

**Keywords:** Machine Learning; Pseudopapilledema; Optic neuropathy; Optic disc swelling

**Posted Date:** May 8th, 2019

**DOI:** <https://doi.org/10.21203/rs.2.318/v2>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

---

**Version of Record:** A version of this preprint was published on August 9th, 2019. See the published version at <https://doi.org/10.1186/s12886-019-1184-0>.

# Abstract

**Background:** This study is to evaluate the accuracy of machine learning for differentiation between optic neuropathies and pseudopapilledema (PPE). **Methods :** Two hundred and ninety-five images of optic neuropathies, 295 images of PPE, and 779 control images were used. Pseudopapilledema was defined as follows: cases with elevated optic nerve head and blurred disc margin, with normal visual acuity (>0.8 Snellen visual acuity), visual field, color vision, and pupillary reflex. The optic neuropathy group included cases of ischemic optic neuropathy (177), optic neuritis (48), diabetic optic neuropathy (17), papilledema (22), and retinal disorders (31). We compared four machine learning classifiers (our model, GoogleNet Inception v3, 19-layer Very Deep Convolution Network from Visual Geometry group (VGG), and 50-layer Deep Residual Learning (ResNet)). Accuracy and area under receiver operating characteristic curve (AUROC) were analyzed **Results:** The accuracy of machine learning classifiers ranged from 95.89% to 98.63% (our model: 95.89%, Inception V3: 96.45%, ResNet: 98.63%, and VGG: 96.80%). A high AUROC score was noted in both ResNet and VGG (0.999). **Conclusions:** Machine learning techniques can be combined with fundus photography as an effective approach to distinguish between PPE and elevated optic disc associated with optic neuropathies.

## Background

Pseudopapilledema (PPE) is defined as an optic nerve with an elevated optic disc and blurred margins that is similar to papilledema or disc swelling associated with various optic neuropathies.<sup>1</sup> Although PPE is a benign condition, it should be differentiated from other optic neuropathies to reduce the need for unnecessary examination and to provide precise diagnosis, prognosis and therapeutic options to the patients. Recently, multi-modal imaging analysis including B-scan ultrasonography, fundus photography, autofluorescence, fluorescein angiography, and optical coherence tomography (OCT) have provided useful information for exact diagnosis of PPE.<sup>2-4</sup> However, the exact differentiation is still difficult.

Machine learning is the use of artificial computer intelligence to enable computers to learn automatically, without being programmed. In ophthalmology, machine learning has been used to analyze various disorders such diabetic retinopathy age-related macular degeneration, and glaucoma.<sup>5-7</sup> We investigated the accuracy and sensitivity of machine learning for differentiation between PPE, optic neuropathies and normals.

## Methods

### 1. Patients

Pseudopapilledema was defined as follows: cases with an elevated optic nerve head and blurred disc margins, with normal visual acuity (>0.8 Snellen visual acuity), visual field, color vision, and pupillary reflex. Only those patients who did not change their optic nerve head and visual function for more than one year were included in the present study. The optic neuropathies group includes 177 cases of ischemic optic neuropathy, 48 of optic neuritis, 17 of diabetic optic neuropathy, 22 of papilledema, and 31 of retinal disorders such as central retinal vein occlusion or posterior uveitis (Figure 1-a). Normal controls were enrolled from routine examination without any abnormal findings and visual problems.

## 2. Data Preparation

Fundus photographs of normal and glaucoma patients were collected from Kim's Eye Hospital. Fundus photography was performed using a non-mydratic auto fundus camera (AFC-330, Nidek, Japan). A total of 1,369 images were obtained, including 295 images of optic neuropathies, 295 of PPE, and 779 normal control images. The obtained images were scaled to a fixed width of 500 pixels while keeping the aspect ratio constant. To remove variations in lightning and brightness of images, the local average color was subtracted using Gaussian filtering.<sup>8</sup> Finally, pixels of each image was normalized to have 0 mean and 1 standard deviation. In order to produce fixed-size input necessary for machine learning models, photos were cropped at the region of optic nerve with a size of 240 240 pixels. Figure 1-b shows the schematic view of the image pre-processing step. The entire set of 1,369 images were split into an 876-image training dataset for training the model, a 274-image validation dataset for validation of the model while training, and a 219-image test dataset for evaluation of the final model. The validation dataset was generated by a random split of 20% of the entire dataset; the test dataset was generated by a random split of 20% of the remaining images after validation split (Table 1). Normal and PPE patients exhibited normal findings on red-free RNFL photography (Vx-10; Kowa Optimed, Inc., Tokyo, Japan), OCT (Cirrus HD-OCT, Carl Zeiss Meditec Inc., Dublin, CA and Spectralis, Heidelberg Engineering, Heidelberg, Germany), and visual field testing (Humphrey 740 visual field analyzer, Carl Zeiss Meditec Inc., Dublin, CA).

## 3. Convolutional Neural Network

### Data Augmentation

Since the images comprised a small dataset, we applied augmentation to each image to overcome overfitting. Each image was cropped at all four corners as well as in the middle, generating five images with a fixed size of 224 224 pixels. This cropping process was repeated again after flipping the image, thereby generating 10 images per photograph. Figure 1-c shows the schematic overview of data augmentation step. Data augmentation can help overcome overfitting by showing the computer an image from various views to aid in decision-making.<sup>9</sup>

### Training Model

We have constructed a convolutional neural network, using Google's Tensorflow deep learning framework as backend.<sup>10</sup> In order to produce best working model, an optimum set of working hyper-parameters are needed. These hyper-parameters include learning rate, activation function, patch size, filter size, number of fully connected layers, and number of hidden nodes in each fully connected layer. However, trying out all possible combinations of hyper-parameters is very time consuming and computationally expensive. Many methods have been proposed for hyper-parameter tuning such as grid search, random search,<sup>11</sup> genetic algorithm,<sup>12</sup> and Bayesian optimization.<sup>13</sup> We implemented Bayesian optimization for our hyper-parameter tuning process using python package Scikit-Optimize. Seven hyper-parameters were tuned using Bayesian optimization including number of convolution layers, number of convolution filters, number of convolution

patch size, number of fully connected layers, number of hidden nodes in each fully connected layer, activation function (rectifier linear unit, exponential linear units, hyperbolic tangent), and learning rate. Max pooling layers were fixed after every convolutional layer with patch size 2 2 and stride 2, and dropout layers with rate 0.5 were fixed after every fully connected layer. Mattern kernel was used for Bayesian optimization and expected improvement was used for acquisition function. The best hyper-parameters were selected after 100 rounds of updating the Gaussian process model. Figure 2 shows a schematic view of hyper-parameter tuning process. The training was conducted again with the selected hyper-parameters with Adam optimizer<sup>14</sup> and cross entropy as a loss function until the average loss of validation data for each epoch started to increase.

## Transfer learning

We conducted transfer learning,<sup>15</sup> which involved training our data with a predefined (trained) existing model using three well-known convolutional neural networks. These include GoogleNet Inception v3,<sup>16</sup> 19-layer Very Deep Convolution Network from Visual Geometry group (VGG) and 50-layer Deep Residual Learning also known as ResNet.<sup>17, 18</sup> These networks were trained using approximately 1.2 million images from ImageNet Large-Scale Visual Recognition Challenge. We have used ImageNet trained weights as an initial starting weight instead of random weights and modified the fully connected layers of the three networks to fit our classification needs. Bayesian optimization was used to tune the hyper-parameters. Four hyper-parameters were tuned including number of fully connected layers, number of hidden nodes, activation function, and learning rate. Dropout layers with rate 0.5 were fixed after every fully connected layer. Additional fine-tuning was conducted after hyper-parameter tuning using Adam optimizer and cross entropy as a loss function. Training was considered finished when the average loss of validation data for each epoch started to increase.

# 4. Evaluation

The model obtains an image as input and outputs the probability that the image represents a photograph of a normal subject, or one with PPE or papilledema. Since we used augmented data (10 images per photography), we generated 10 probabilities from a single image. By averaging these probability values, we obtained a single probability that the image is normal, or depicts PPE or papilledema (Fig 3-b). Using this strategy, we evaluated our model as well as GoogleNet Inception v3, VGG, and ResNet transferred model. Also, we have calculated micro-averaged sensitivity and specificity of each model and generated ROC (receiver operating characteristic) curve which indicates overall performance of how well the models classify images into three groups (Normal, PPE, papilledema).

## Results

Table 2 shows the summarized results of our model and transfer learning model. After hyper-parameter tuning, our model exhibited 3 convolution layers and 5 fully connected layers. The first convolution layer had a patch size of 2 2 with 27 filters, the second had a patch size of 1 2 12 with 40 filters, and the last

convolution layer had a patch size of 29 29 with 35 filters. Max pooling layer was applied after every convolutional layer with a patch size of 2 2 and a stride of 2. The fully connected layers consisted of 366, 177, 512, 159, and 133 hidden nodes, respectively. A dropout rate of 0.5 was used in fully connected layers. Rectifier linear unit was used as an activation function.

Figure 3-a shows the schematic architecture of our CNN model. The Inception V3 model exhibited 1 fully connected layer with 60 hidden nodes along with Rectifier linear unit as an activation function. The VGG model exhibited 3 fully connected layers with each layer having 512 hidden nodes and Exponential linear unit as an activation function. The ResNet model had 1 fully connected layer with 325 hidden nodes with hyperbolic tangent as an activation function. All the transferred models had dropout layer with dropout rate 0.5 after every fully connected layer. Further, all models used softmax layer as a classification layer. The best performing model based on test accuracy was the ResNet transfer learned model. The ROC curve for each model is depicted in Figure 4-a and the confusion matrix based on test data for each model is depicted in Figure 5. At the cost of 0.007 difference of AUROC, our model used the least number of parameters (11,636,096) among models (Figure 4-b). In addition, the validation loss graph showed that validation loss reached zero level at around epoch 16 (Figure 4-c).

## Discussion

This study suggests that machine learning techniques can be combined with fundus photography as an effective approach to distinguish between PPE and elevated optic disc related with optic neuropathies.

We have used 3 state-of-the-art convolutional neural networks, including GoogleNet Inception V3, VGG, and ResNet. In addition, we used pre-trained weights from ImageNet Large-Scale Visual Recognition Challenge as the initial parameter to train our model instead of random weights; this is a popular method since these initial parameters are already optimized for detecting natural images such as edges and curves,<sup>19</sup> thus solving the issue of overfitting when not many data are available. We have also trained our own model from scratch using Bayesian optimization as the hyper-parameter tuning process. As depicted in Table 2, transferred models outperformed our model based on test accuracy but, the difference, based on AUROC, was small. Between our model and the best performing ResNet transferred model, our model used far less parameters.. Although minimizing number of parameters is not a critical point in building a deep learning model, if one can get a similar result using far less parameters, one can reduce the computational time and also benefit when deploying the model on a web service, machines and etc.

Overfitting, which refers to models performing well on the trained data but not well on unseen data is a common issue when a small dataset is used to train the model.<sup>20</sup> Since our dataset consisted of only 1,369 images, there might have been a possibility of overfitting. However, we addressed this issue by incorporating regularization techniques such as adding dropout layers and data augmentation. Dropout randomly corrupts hidden nodes between layers which changes the detail of the model every training iteration.<sup>21</sup> Thus, this process leads to a more generalized model when a sufficient number of training iterations are given. Data augmentation allows the machine to learn an image from different views. This technique can also help

overcome the issue of small training dataset by generating many augmented images. The validation loss graph for our model indicates that our model reached minimum, which is an indication of an optimal model.<sup>9</sup>

The accuracy in this study was 95.89%. Chang MY. et al.<sup>22</sup> reported that differentiation using fluorescein angiography was 97% and it showed similar performance with this study. However, the differentiation method using machine learning is much safer and easier than FAG. In addition, other modalities such as B-scan ultrasonography, fundus photography and OCT revealed high misinterpretation rate.

Even though we have generated our own model and used well-known, state-of-the-art convolutional neural networks, getting an insight into how a machine classifies a fundus photo as normal or disease status can be a challenging task. Therefore, further study is needed into visualizing the convolution layers and filters to get an idea of how machines classify images. Further, larger scale datasets may help validate our findings.

## Conclusions

Machine learning techniques can be combined with fundus photography as an effective approach to distinguish between PPE and elevated optic disc associated with optic neuropathies.

## Declarations

Ethics approval and consent to participate

All procedures performed in studies involving human participants were in accordance with the ethical standards of the local ethics committee of Kim's Eye Hospital and with the 2013 Helsinki declaration. The need for consent was waived by IRB of Kim's Eye Hospital.

Consent for publication : not applicable

Availability of data and material : The datasets analyzed during the current study and trained models are available from the corresponding author on reasonable request.

Competing interests : None

Funding : None

Authors' contributions All authors have read and approved the manuscript,

Name

Location

Role

Contributions

JMA

Soongsil University

Author

Design and conceptualized study; analyzed the data; drafted the manuscript

SSK

Soongsil University

Author

Design and conceptualized study; analyzed the data, supervision of the manuscript

KSA

PDXen

Author

Design and conceptualized study; analyzed the data, supervision of the manuscript

SHC

PDXen

Author

Design and conceptualized study; analyzed the data, supervision of the manuscript

USK

Kim's Eye Hospital, Konyang University

Author

Design and conceptualized study; analyzed the data, supervision of the manuscript

Acknowledgements None

## List Of Abbreviations

OCT : optical coherence tomography

PPE : pseudopapilledema

VGG : Visual Geometry group

## References

1. Trick GL, Bhatt SS, Dahl D, Skarf B. Optic disc topography in pseudopapilledema: a comparison to pseudotumor cerebri. *J Neuroophthalmol* 2001;21:240-4.
2. Aghsaei Fard M, Okhravi S, Moghimi S, Subramanian PS. Optic Nerve Head and Macular Optical Coherence Tomography Measurements in Papilledema Compared With Pseudopapilledema. *J Neuroophthalmol* 2018.
3. Thompson AC, Bhatti MT, El-Dairi MA. Bruch's membrane opening on optical coherence tomography in pediatric papilledema and pseudopapilledema. *J AAPOS* 2018;22:38-43 e3.
4. Saenz R, Cheng H, Prager TC, Frishman LJ, Tang RA. Use of A-scan Ultrasound and Optical Coherence Tomography to Differentiate Papilledema From Pseudopapilledema. *Optom Vis Sci* 2017;94:1081-1089.
5. Gulshan V, Peng L, Coram M, et al. Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs. *JAMA* 2016;316:2402-2410.
6. Kim SJ, Cho KJ, Oh S. Development of machine learning models for diagnosis of glaucoma. *PLoS One* 2017;12:e0177726.
7. Rahimy E. Deep learning applications in ophthalmology. *Curr Opin Ophthalmol* 2018;29:254-260.
8. Ebner M. Color constancy based on local space average color. *Machine Vision and Applications* 2009;20:283-301.
9. Wong SC, Gatt A, Stamatescu V, McDonnell MD. Understanding data augmentation for classification: when to warp? *arXiv preprint arXiv:160908764* 2016.
10. Girija SS. *Tensorflow: Large-scale machine learning on heterogeneous distributed systems*. 2016.
11. Bergstra J, Bengio Y. Random search for hyper-parameter optimization. *Journal of Machine Learning Research* 2012;13:281-305.
12. Man K-F, Tang K-S, Kwong S. Genetic algorithms: concepts and applications [in engineering design]. *IEEE transactions on Industrial Electronics* 1996;43:519-534.
13. Snoek J, Larochelle H, Adams RP. Practical bayesian optimization of machine learning algorithms. *Advances in neural information processing systems*, 2012:2951-2959.
14. Kinga D, Adam JB. A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*, 2015.
15. George D, Shen H, Huerta E. Deep Transfer Learning: A new deep learning glitch classification method for advanced LIGO. *arXiv preprint arXiv:170607446* 2017.

16. Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z. Rethinking the inception architecture for computer vision. Proceedings of the IEEE conference on computer vision and pattern recognition, 2016:2818-2826.
17. Conneau A, Schwenk H, Barrault L, Lecun Y. Very deep convolutional networks for text classification. arXiv preprint arXiv:1606.01781 2016.
18. Wu S, Zhong S, Liu Y. Deep residual learning for image steganalysis. Multimedia tools and applications 2018;77:10437-10453.
19. Russakovsky O, Deng J, Su H, et al. Imagenet large scale visual recognition challenge. International Journal of Computer Vision 2015;115:211-252.
20. Dietterich T. Overfitting and undercomputing in machine learning. ACM computing surveys (CSUR) 1995;27:326-327.
21. Krizhevsky A, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 2012:1097-1105.
22. Chang MY, Velez FG, Demer JL, et al. Accuracy of Diagnostic Imaging Modalities for Classifying Pediatric Eyes as Papilledema Versus Pseudopapilledema. Ophthalmology 2017;124:1839-1848.

## Tables

Table 1. Sample number for Convolutional Neural Network.

	Normal	Pseudopapilledema	Papilledema	Total
Entire Data	779	295	295	1,369
Training Data	505	197	174	876
Validation Data	155	53	66	274
Test Data	119	45	55	219

Table 2. Evaluation of our model and transferred model

	Our Model		Inception V3		ResNet		VGG	
	Ensemble Accuracy	AUROC						
Training Data	100%	1.0	100%	1.0	100%	1.0	100%	1.0
Validation Data	96.35%	0.989	98.18%	0.993	98.18%	0.996	97.81%	0.996
Test Data	95.89%	0.992	96.35%	0.997	98.63%	0.999	96.80%	0.999

# Figures

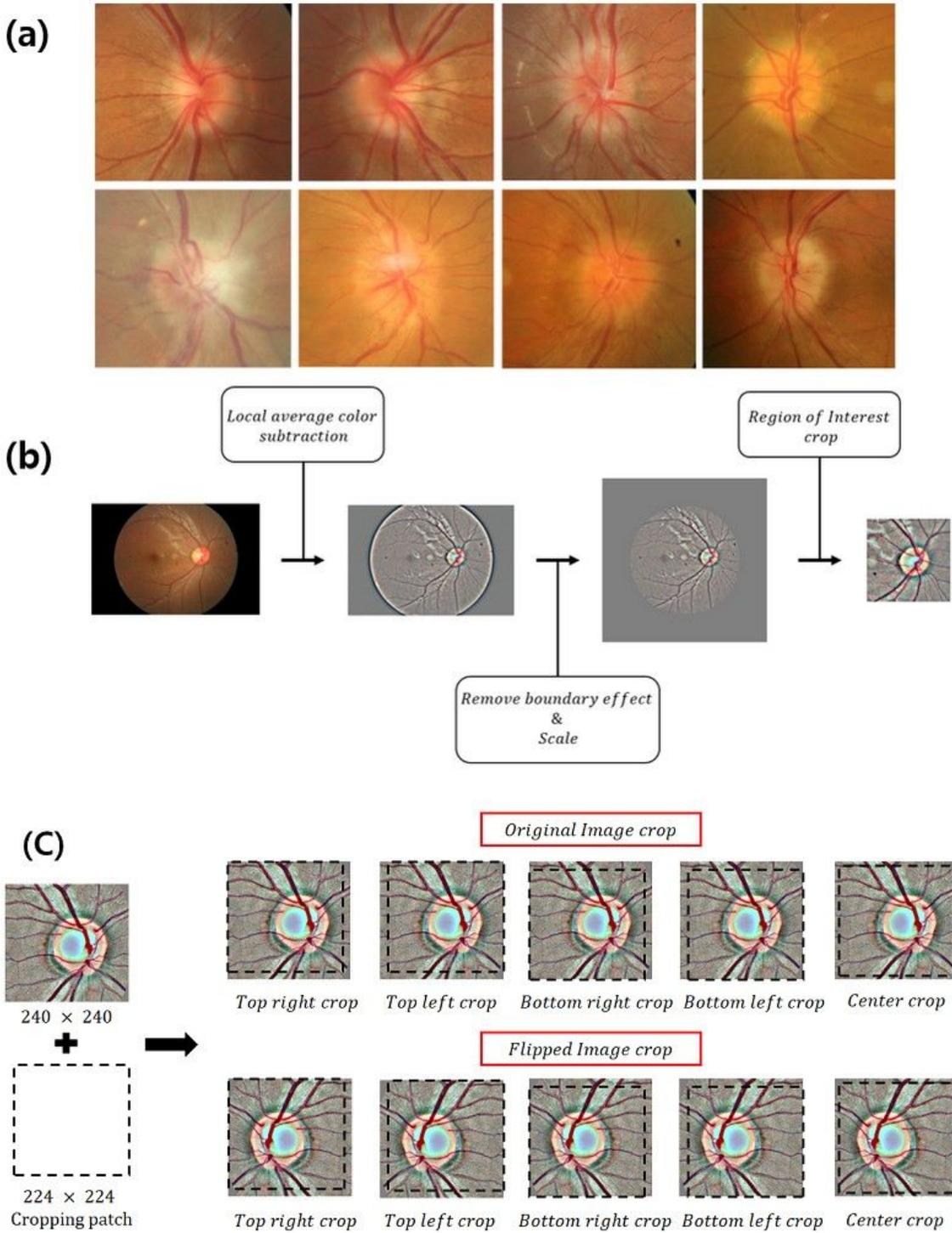


Figure 1

1-a. Optic disc findings in fundus photography. Various features from pseudopapilledema (upper low) and swollen disc from optic neuropathies (lower low) Figure 1-b. Schematic view of image pre-processing process. Figure 1-c. Schematic view of data augmentation process.

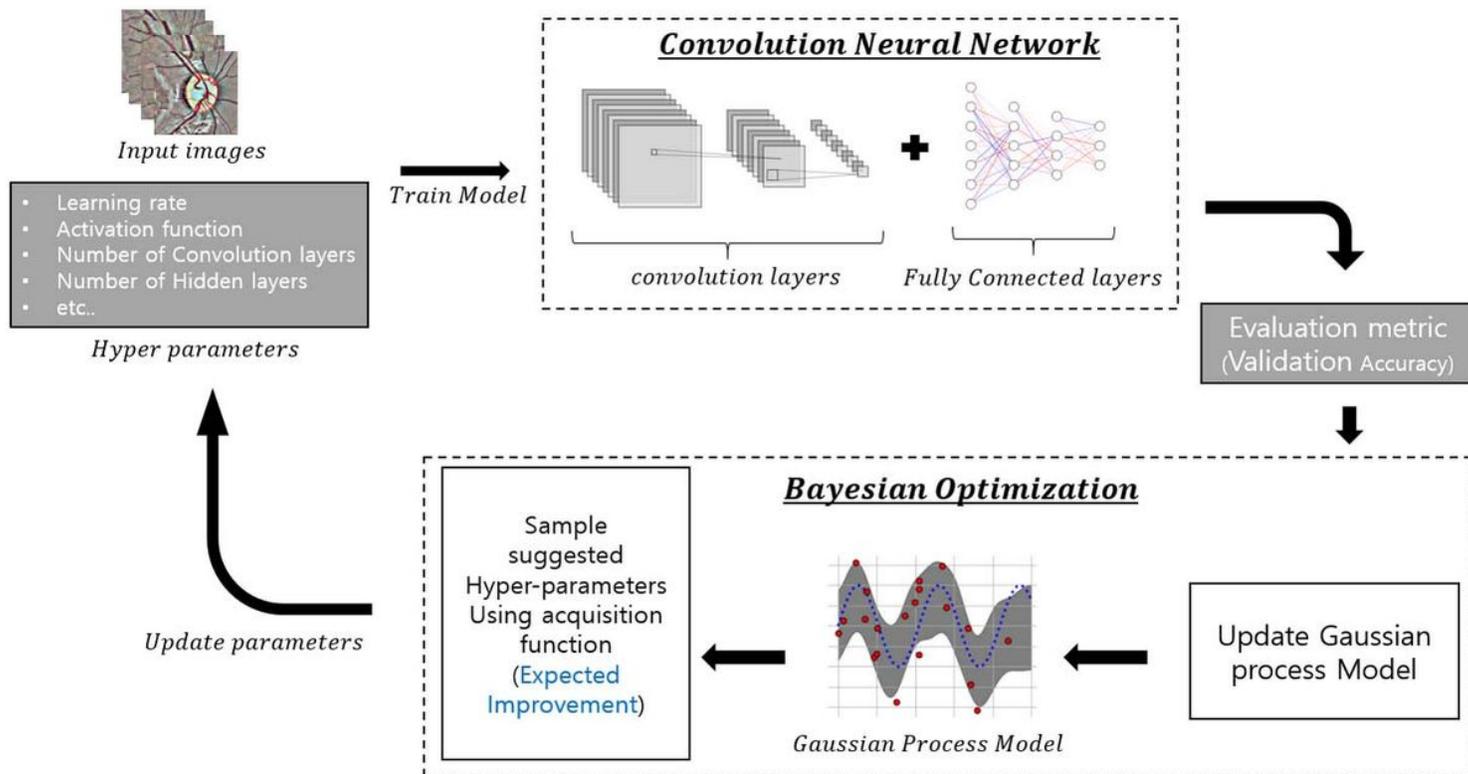


Figure 2

Schematic view of Bayesian optimization. Seven hyper-parameters were tuned using Bayesian optimization: learning rate, activation function, number of convolution layers, convolution patch size, filter size, number of fully connected layers, and number of hidden nodes in each fully connected layer.

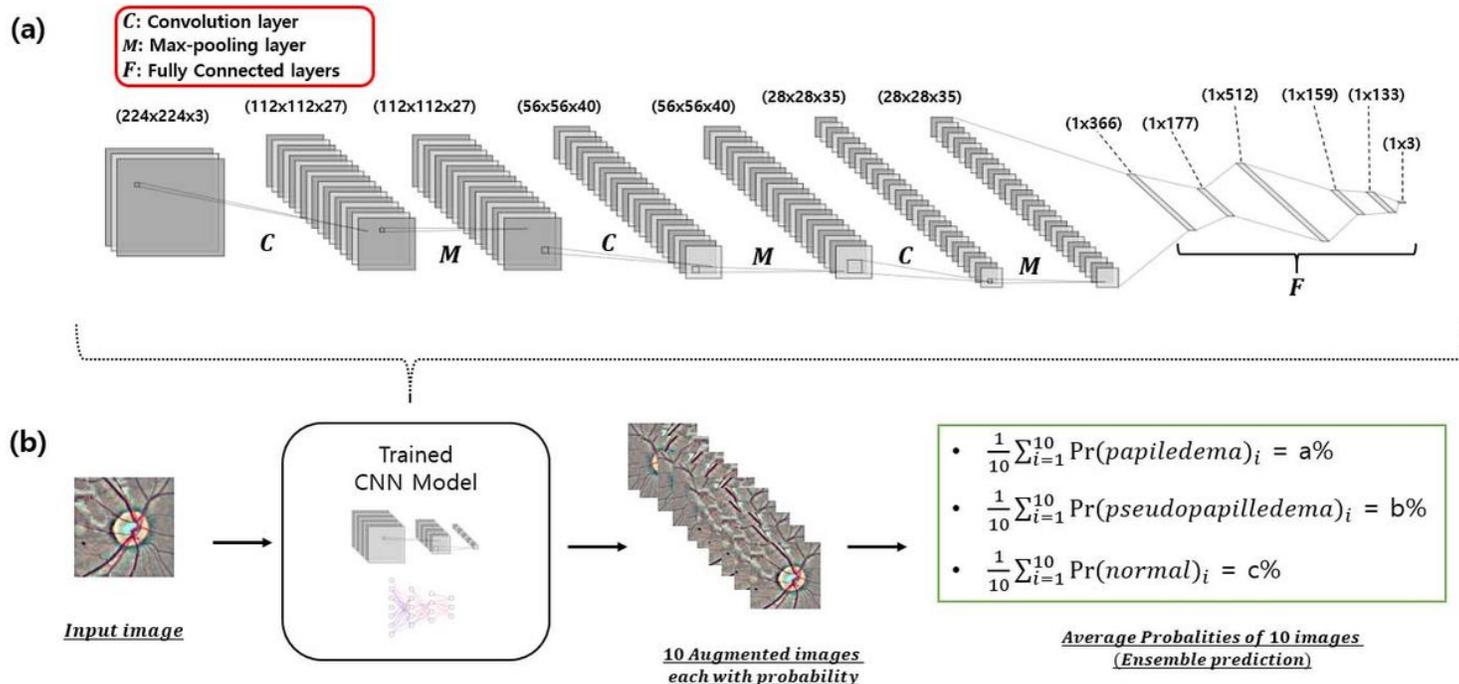
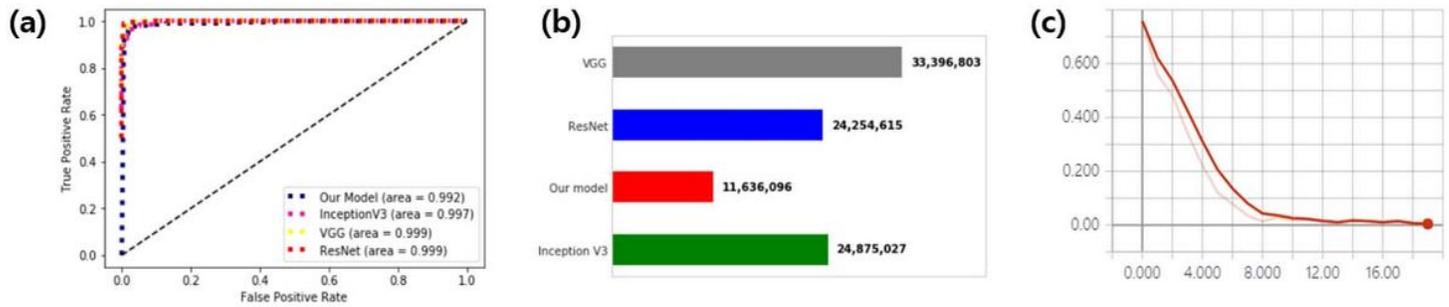


Figure 3

3-a. Schematic view of our model. It consists of 3 convolutional layer each with max pooling layer followed by 5 fully connected layers and a softmax layer. 3-b. Evaluation process. Ten augmented images were averaged to give a single probability for each class.



**Figure 4**

4-a. Receiver operating characteristic curve. 4-b. Number of parameter comparison between models. 4-c. Loss graph for our model. Y-axis indicates loss for validation data and X-axis indicates number of epoch.

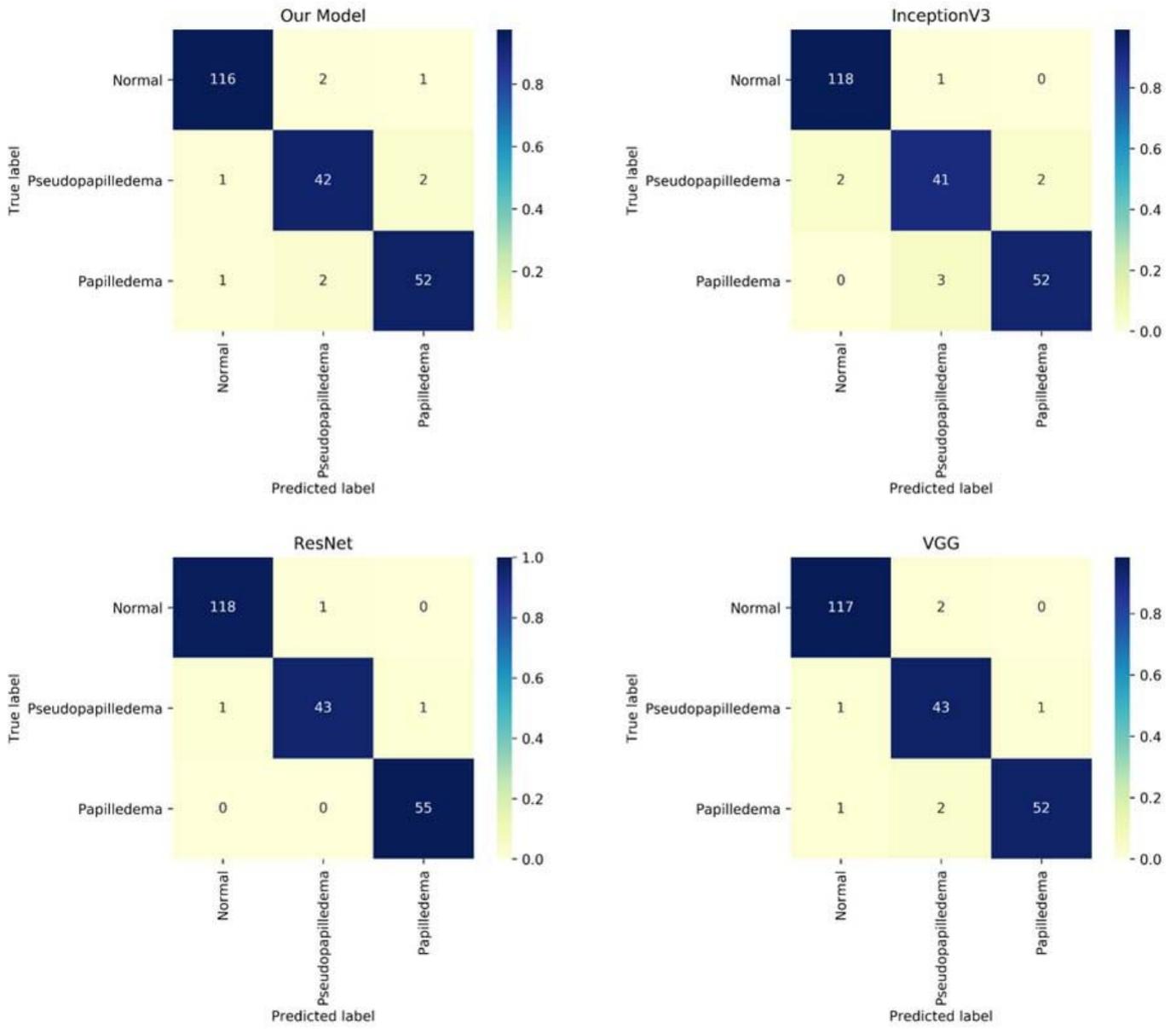


Figure 5

Confusion matrix.