

Machine Learning Engineer Nanodegree

Capstone Project

Alejandro Trujillo

May 27th, 2018

I. Definition

Project Overview

This project focuses on leveraging the power of convolutional neural networks (CNNs) and transfer learning models to predict the diagnosis of Optical Coherence Tomography (OCT) scans.

Optical coherence tomography is a non-invasive imaging test. OCT uses light waves to take cross-section pictures of your retina.

These diagnosis provide treatment guidance for multiple diseases, including:

- Choroidal neovascularization (CNV)
- Drusen
- Diabetic Macular Edema (DME)

The goal of this project is to develop a model that can enhance and aid physicians in their diagnosis of CNV, Drusen and DME diseases.

Problem Statement

Problem

Approximately 30 million OCT scans are performed each year, and the analysis and interpretation of these images takes up a significant amount of time (Swanson and Fujimoto, 2017).

Well trained ophthalmologists can do an outstanding job at classifying any disease depicted in a OCT scan. However, long hours at work and eye fatigue can hinder their diagnosis. Even one misdiagnosed case can hinder the relationship with their patients and their reputation. Here I present an approach to aid ophthalmologists streamline their diagnosis and further enhance their relationship with their patients.

Intended Solution

The goal is to provide a machine learning model utilizing convolutional neural networks to predict the diagnosis of OCT scans with over 70% accuracy.

Strategy

The tasks involved in accomplishing this goal are the following:

1. Download and preprocess the OCT scans dataset
2. Create a CNN benchmark model with at least 1% accuracy to improve upon
3. Create a transfer learning model with at least 70% accuracy utilizing one of the following pretrained models:
 - a. VGG-19
 - b. ResNet-50
 - c. InceptionV3
 - d. Xception
4. Evaluate the models using their accuracy.

Metrics

The dataset is divided in training, validation and testing to prevent overfitting.

Each dataset has the following total number of images:

- Training has 5,082 images
- Validation has 938 images
- Test has 1,000 images

Additionally the following metrics are used to quantify the performance of the benchmark and optimized algorithms.

- **Training time:** Measures the time that the algorithm takes to learn the patterns within the training set.
 - Important to keep training time as small as possible without impacting accuracy to provide a quicker feedback loop and try other models or approaches.
- **Accuracy:** Indicates the degree to which the resulted output matches the correct label.
 - Important to keep as high as possible to provide to the user confidence that the model is producing a correct diagnosis.
- **Optimizer:** RMSprop divides the learning rate by an exponentially decaying average of squared gradients.

II. Analysis

Data Exploration

The original dataset is comprised of 84,495 total X-Ray images.

This large dataset had to be reduced to only 7,020 total images due to the limited resources available at the time. (8GB of RAM, Intel Core i7-4700HQ @ 2.40GHz and HDD Toshiba 1TB 5400 RPM)

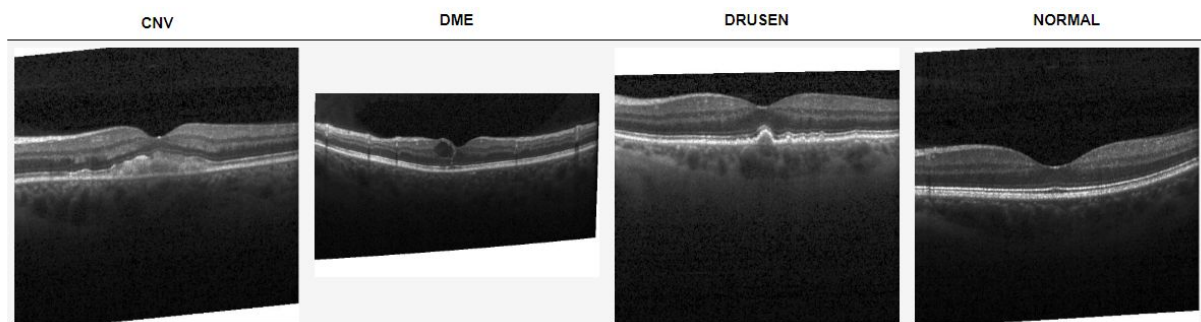
The reduced dataset of 7,020 total images was divided in the following sections

| | CNV | DME | DRUSEN | NORMAL | TOTAL |
|------------|-------|-----|--------|--------|--------------|
| Training | 1,636 | 939 | 1,229 | 1,278 | 5,082 |
| Validation | 252 | 217 | 210 | 259 | 938 |
| Testing | 250 | 250 | 250 | 250 | 1,000 |

Exploratory Visualization

As seen in the images below the different categories could display similar traits.

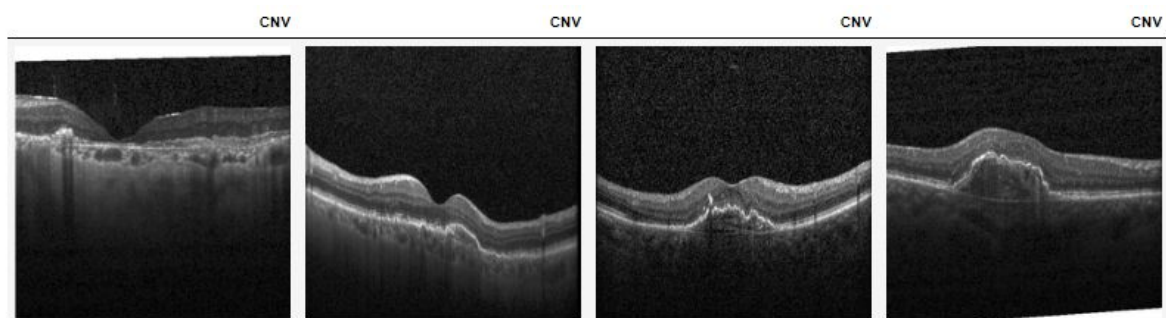
It is hard for someone without medical diagnosis background to identify outliers in the data, due to the specialized field of ophthalmology and radiology.



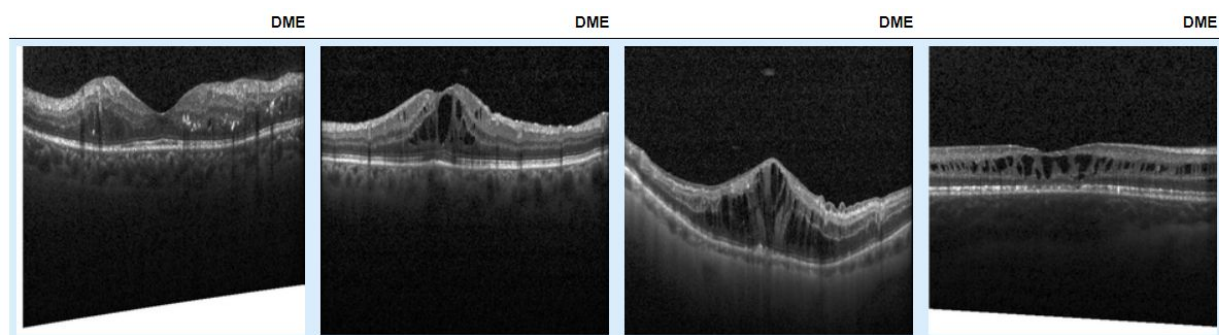
Nonetheless, when looking at many images of the different categories, one can see some distinct traits between the different categories.

Here is an example of the previous statement.

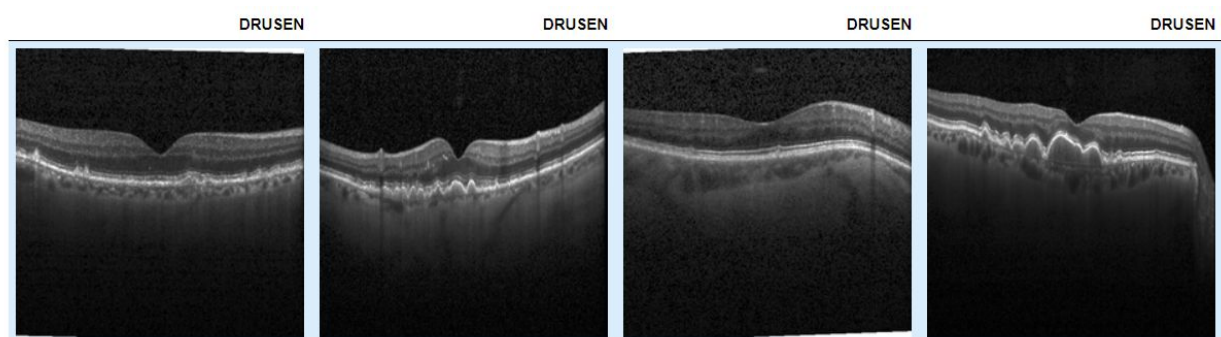
Distinctive CNV scans



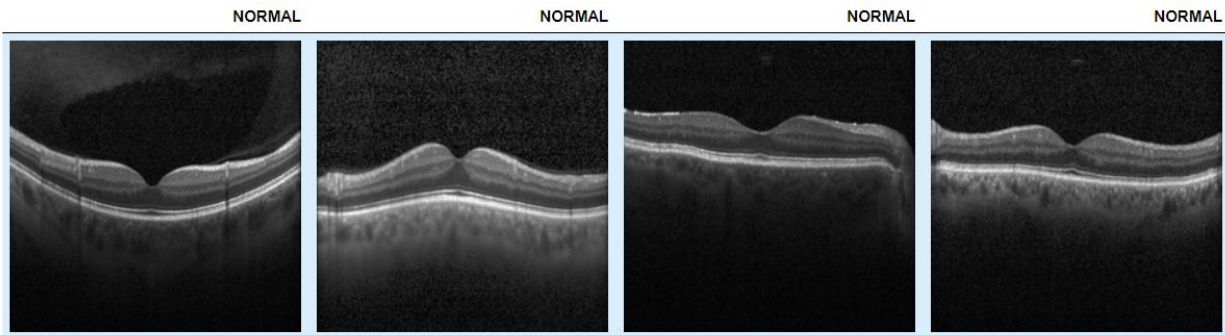
Distinctive DME scans



Distinctive DRUSEN scans



Distinctive NORMAL scans



A convolutional neural network should be able to analyse these images, identify the different patterns and produce an architecture robust enough to produce at least 70% accuracy.

Algorithms and Techniques

Processing the data.

After downloading the data. The initial step for data processing is to resize each image to 224x224 pixels on disk. This is with the purpose of optimizing computing resources and to accelerate training time. In a future step, the model will utilize this data dimension as its input shape to produce the CNN architecture.

Convolutional Neural Network (CNN)

I chose to use a CNNs architecture for my image classification project because of their outstanding image classification performance and frameworks available that facilitate the use of it, such as Keras and TensorFlow.

As stated by multiple sources and research papers on the internet:

“Convolutional neural networks (CNN) have recently shown outstanding image classification performance in the large scale visual recognition challenge (ILSVRC2012). The success of CNNs is attributed to their ability to learn rich mid-level image representations as opposed to hand-designed low-level features used in other image classification methods.” [1]

Additionally, an interesting fact to ponder is that a few years ago the following statement was made by Microsoft researchers.

“Based on a survey of machine vision literature and vision experts at Microsoft Research, we believe classification accuracy of better than 60% will be difficult without a significant advance in the state of the art.”^[2]

Even though that statement was made for CAPTCHA mechanisms a few years ago, it is interesting to see the improvements that have happened in the machine learning field.

Nowadays, the results that can be accomplished in the image classification reign with a limited amount of resources are impressive.

The CNNs architecture presented in this project is able to obtain at least 70% accuracy running on a 3 year old laptop with 8GB of RAM, Intel Core i7-4700HQ @ 2.40GHz and HDD Toshiba 1TB 5400 RPM .

Benchmark model

For the benchmark model, I built a basic CNN architecture from scratch. This model's initial goal was to achieve more than 1% accuracy.

The benchmark architecture will be as follows.

1. Convolutional 2D layer with kernel size 2 and ReLU activation function.
2. Max pooling 2D layer with pool size 2.
3. Convolutional 2D layer with kernel size 2 and ReLU activation function.
4. Max pooling 2D layer with pool size 2.
5. Convolutional 2D layer with kernel size 2 and ReLU activation function.
6. Max pooling 2D layer with pool size 2.
7. Dropout layer with 0.3 probability.
8. Flatten
9. Dense with 4 nodes and softmax activation function.

This benchmark model is used to compare the results of the optimized transfer learning model.

After its implementation, the benchmark model produced a test accuracy of **58.80%**

Parameters

The benchmark model is configured with the following parameters:

- Optimizer = rmsprop
- Loss = categorical_crossentropy
- Metrics = accuracy
- Epochs = 7

Time to train

The benchmark model took around 20 minutes to train on 5,082 images.

III. Methodology

Data Preprocessing

Minimal data preprocessing was required.

The only steps necessary to build the benchmark model and the transfer learning model are described below:

1. Resize all images to 224x224 on disk
2. Load the image from disk with a target size of 224x224 pixels
3. Finally convert the image into a 4D tensor with shape (1, 224, 224, 3).

Implementation

Transfer Learning

For the transfer learning model, I chose to utilize the pretrained Inception V3 model, with weights pre-trained on ImageNet.

In a deep learning architecture higher layers identify the most general features of the dataset, such as edges.

Deep layers identify more distinctive features, as the deep learning architecture grows

The premise for utilizing transfer learning is to reduce training time and improve accuracy. This should be possible by utilizing the pre-trained layers and focus on the deeper layers that identify most specific traits.

An image of the full InceptionV3 architecture can be found in the following link: [InceptionV3 Architecture](#)

Parameters

The transfer learning model is configured with the following parameters:

- Optimizer = rmsprop
- Optimizer = Stochastic Gradient Descent:
 - Learning rate = 0.0001
 - Momentum = 0.9
- Loss = categorical_crossentropy
- Metrics = accuracy
- Epochs = 7

Time to train

The time to train was roughly 20 minutes, slightly better than the benchmark model.

Accuracy

The transfer learning model produced a test accuracy of **81.90%**

Complications

- Creating an Amazon EC2 g3.8xlarge instance was the initial option to explore for creating the model. This would have enabled the model to be trained in much more data, and possibly output an accuracy of over 95%. However, a couple days were spent requesting Amazon to allow me to have such instance. Amazon's first response was to decline my request. Only after appealing their decision and after 4 days in total my request was granted but by then I had already created the model and achieved the expected results locally.

- It is unfeasible to utilize all the images of the original dataset given the computing resources available at the time.
 - The implementation phase utilized 100% of ram and disk causing the computer to become unresponsive many times and requiring restart.
 - The dataset had to be reduced to fit in the computer resources and to train the model.
 - Also, each image in the reduced dataset had to be resized on disk to the expected dimensions to optimize the disk read time.
- The first 200 layers of the InceptionV3 architecture had to be freezed, meaning utilize the existing weights, to achieve the above 70% accuracy. The top 110 layers of the architecture where trained on the OCT dataset. It is unclear why training the top 110 layers of this architecture was still faster than the benchmark model.
- Several days were spent preparing the computing resources, gathering data, researching articles about transfer learning, building prototypes. All this effort has enabled this project to be successful in building a model that can make predictions with over 70% accuracy.

Coding process

Getting familiar with Keras pretrained models can take some effort, research and trial and error. I decided to focus on their pretrained InceptionV3 model on the ImageNet dataset. I spent several hours in the “Utilize InceptionV3 pretrained model to predict OCT diagnosis” section of my InceptionV3 jupyter notebook trying different approaches until I was able to compute the expected accuracy.

Some of the approaches I researched were:

- Image Augmentation.
- Freezing and Unfreezing different InceptionV3 pretrained layers
- Increasing or reducing the dataset
- Increasing the number of epochs from 10 to 100.

A rough draft of my prototyping can be found in here: [Prototyping](#)

Refinement

The following parameters were refined and explored during implementation until the expected accuracy result was achieved.

- **Image Augmentation:** The following parameters were tried during prototyping but not implemented in the final solution because they didn't produce notable accuracy gains and the expected accuracy was achieved without requiring further image augmentation:
 - rotation_range=40,
 - width_shift_range=0.2,
 - height_shift_range=0.2,
 - rescale=1./255,
 - shear_range=0.2,
 - zoom_range=0.2,
 - horizontal_flip=True,
 - fill_mode='nearest'
 - Full list of parameters and definitions can be found here: <https://keras.io/preprocessing/image/>
- **Freezing and Unfreezing different InceptionV3 pretrained layers.**
 - Originally only 249 layers were frozen on the InceptionV3 model. This number had to be reduced to only 200 to allow the deeper layers of the model to be trained on the new dataset.
- **Increasing or reducing the dataset**
 - The total number of 84,495 images was reduced to only 5,082 images given the limited resources available.
 - Initially the number of images was reduced to only 1,000 images. The model was not able to identify specific features with such a reduced dataset. The model was only able to produce a top accuracy score of 53.1306% with this dataset.

- **Increasing the number of epochs from 10 to 100.**
 - Increasing the number of epochs didn't have much impact in the final accuracy score. Even with just 4 or 5 epochs the model was able to compute the necessary traits to produce an acceptable accuracy.

IV. Results

Model Evaluation and Validation

Final Model

The final model utilizes the pretrained InceptionV3 architecture on the ImageNet dataset.

Only the first 200 layers of the model are frozen with their initial weights. The deeper 110 layers are trained with the OCT dataset.

The full architecture of the InceptionV3 model can be found here: [InceptionV3 Architecture](#)

Architecture

The top layer of the pretrained architecture was replaced by the following layers, to allow the model to make predictions on the new dataset:

1. GlobalAveragePooling2D
2. Dense layer with 1024 nodes and ReLU activation function
3. Dense layer with 4 final nodes and Softmax activation function to make classifications.

Initially the InceptionV3 architecture was compiled using the following parameters:

- optimizer='rmsprop',
- loss='categorical_crossentropy'

After being compiled and 200 of its initial layers frozen the model went through a second compilation process producing better accuracy and less validation loss.

The second compilation process made use of the following parameters:

- Optimizer = Stochastic Gradient Descent with learning rate 0.0001 and momentum=0.9
- loss='categorical_crossentropy'

It is unclear why this second compilation process produced much better accuracy and less validation loss.

Final Results

The model only took about 20 minutes to train and it successfully produced a accuracy of over 70% consistently.

Expectations

The final model has been tested with various inputs and it has shown to generalize well to unseen data.

The model is robust enough to produce over 70% accuracy consistently as long as is trained with at least 5,082 training images.

The results from the model can be trusted for this exercise. However, consulting with an experienced opthamologist is necessary to take this model to a production environment.

Justification

Benchmark vs Transfer Learning Model

The following table shows the difference between the benchmark and the transfer learning model.

During trials. The benchmark model accuracy oscillates between 55% and 70% accuracy.

| | Total Training Images | Epochs | Time to Train | Accuracy |
|-------------------------|-----------------------|--------|---------------|----------------|
| Benchmark | 5,082 | 7 | 12 min 33sec | 58.80% |
| Transfer Learning Model | 5,082 | 5 | 22 min 46 sec | 81.90 % |

Final Solution

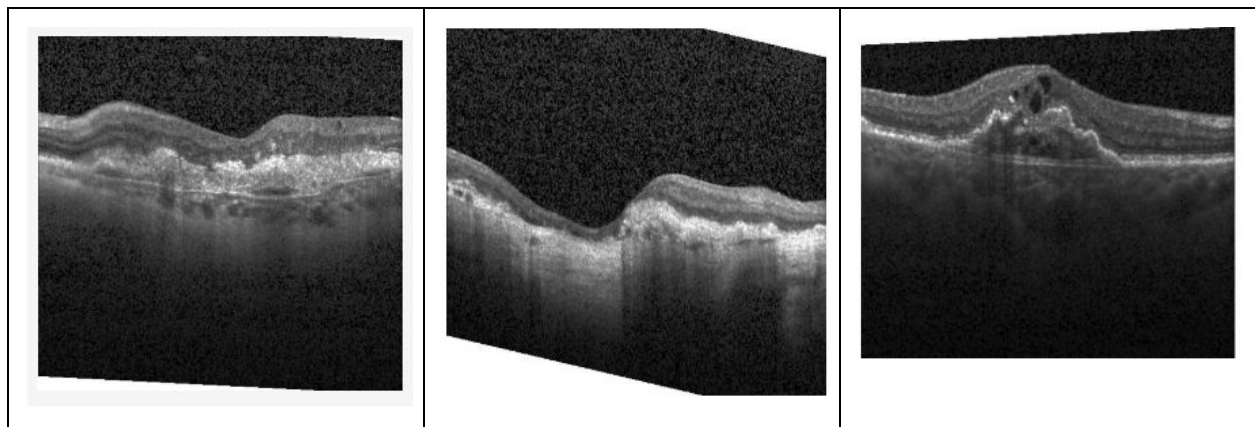
The final transfer learning solution achieves over 70% accuracy consistently. This final solution is robust enough to have solved the problem statement for this project.

V. Conclusion

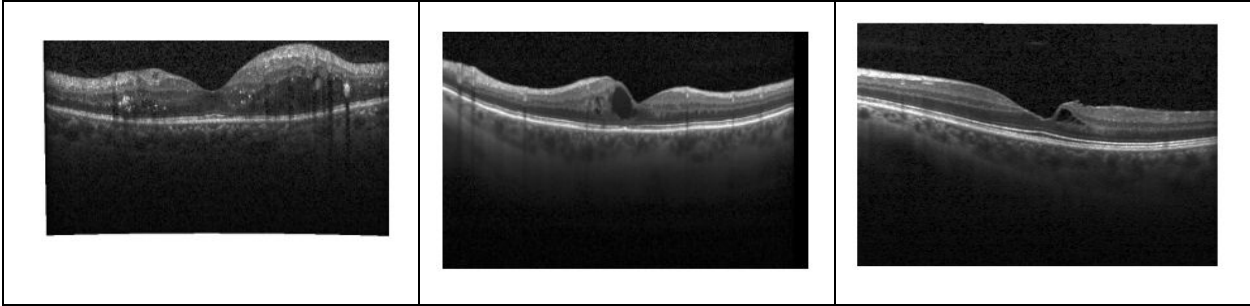
Free-Form Visualization

An important quality of the project is that it achieves a good level of accuracy. Here are a few examples where the project made the correct diagnosis prediction for a OCT scan.

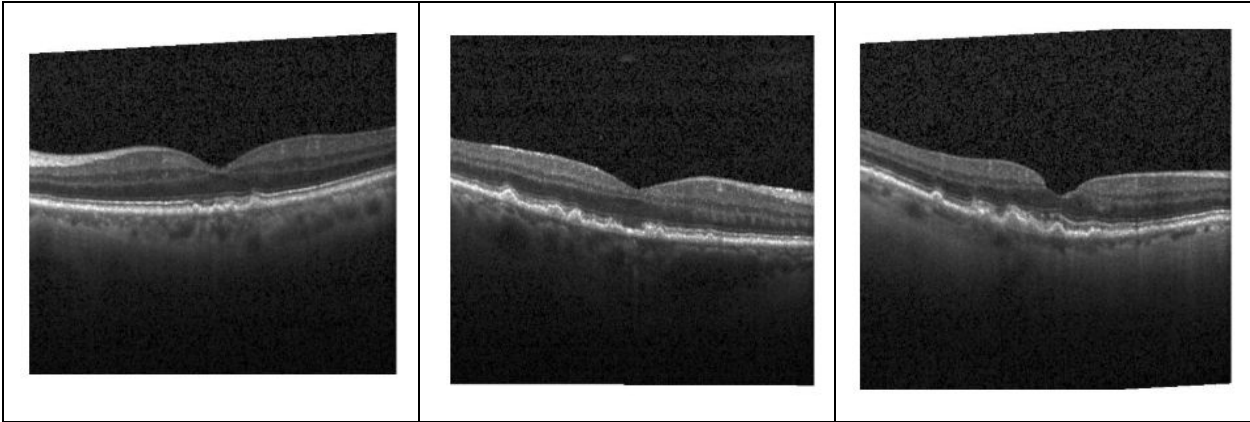
CNV



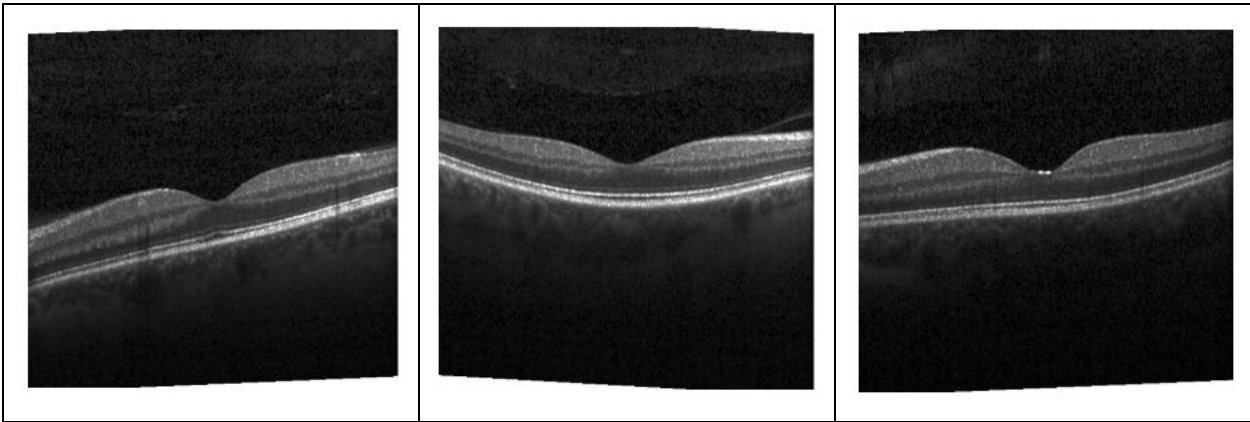
DME



DRUSEN



NORMAL



Reflection

Initially I went to kaggle to look for an idea for my capstone project. I knew I wanted to focus on the medical field and fortunately I found this idea of creating a classifier for OCT scans interesting.

The initial post can be found here: ([Kaggle OCT Scans](#))

This post states a good objective and it piqued my interest. So I decided to use it as my capstone project.

There were many difficulties in completing this project.

Preparing all the tools in my local computer required roughly 8hrs. The tools required where CUDA, TensorFlow and Keras. I had the opportunity to setup these tools at one of my Udacity Connect sessions.

It took me 8 hrs because some links to get the tools were broken or I installed the newer versions of CUDA which are not yet supported by TensorFlow.

The first difficulty started in getting the full dataset. The kaggle article doesn't have the data organized nor it has the full dataset. However, in it I found references that point to the source Mendeley.com. Fortunately, this source contains the full dataset already labeled.

Then I started processing the data. It was clear from the beginning that my local computer didn't had the necessary resources to process this large dataset.

I wanted to train the model using the full dataset but the resources available at the time didn't allow me to make use of the full dataset.

Therefore, I requested an EC2 instance capable of such computing power but my initial request was denied by Amazon. Only after appealing the decision was my request granted. But by then 4 days had gone by and I had already built a successful model locally.

I spent many days researching how to build a proper transfer learning model by reading many articles online, spending time in Udacity ML classroom and reading Keras documentation and examples. After all my research, I began prototyping and here is where a lot of my time was spent trying different configurations and approaches.

Training and testing the benchmark and transfer learning model took about 1 hr for both. So I had to make sure that I was implementing a reasonable CNN architecture and using the correct parameters to optimize my time.

I found very interesting that creating an acceptable classifier for OCT scans, and other X-Ray images, is very feasible nowadays. I would be very interested in seeing the positive impact that models like this can have in enhancing the medical field. Also I would assume physicians will be open-minded to adopt machine learning methods to streamline their diagnosis.

Also, these predictive models can help optimize a physician schedule by providing an accurate and quick diagnosis.

This final transfer learning model fits my expectations for the problem since it achieves above 70% accuracy and it only takes about 20 minutes to train.

Also, I was very impressed by achieving over 50% accuracy by the benchmark model. I initially thought that it would be difficult to get over 10% accuracy in such a basic CNN architecture and with only 5,082 training images.

I believe that to take this model to a production setting a professional opthamologist should be consulted to see what level of accuracy is necessary to produce production-ready diagnosis.

I believe there is a lot of potential for machine learning algorithms in the medical field and undoubtedly new business will be built in this field.

Improvement

Image augmentation should be explored to a greater extent to test if it helps achieve better accuracy.

Training the model on the full dataset should also be explored, as it should provide an accuracy improvement. Many cloud providers have instances specialized for machine learning, this option should be explored to possibly utilize the full dataset for training.

Additional classification capabilities for predicting other conditions, such as bone fractures and MRIs, could enhance the user experience as well.

This algorithm could be used to build a website. This website could help in exposing the the power of machine learning to ophthalmologists, and explore the enhancement of assisted diagnosis by a machine learning algorithm.

References:

- ^[1] Oquab, M., Bottou, L., Laptev, I. and Sivic, J. (2018). Learning and Transferring Mid-Level Image Representations using Convolutional Neural Networks. [online] Openaccess.thecvf.com. Available at: http://openaccess.thecvf.com/content_cvpr_2014/papers/Oquab_Learning_and_Transferring_2014_CVPR_paper.pdf [Accessed 28 May 2018].
- ^[2] Elson, J., Howell, J., Douceur, J. and Saul, J. (2018). Asirra: A CAPTCHA that Exploits Interest-Aligned Manual Image Categorization. [online] Microsoft.com. Available at: <https://www.microsoft.com/en-us/research/wp-content/uploads/2007/10/CCS2007.pdf> [Accessed 28 May 2018].
- Turbert, D. (2018). What Is Optical Coherence Tomography?. [online] American Academy of Ophthalmology. Available at: <https://www.aao.org/eye-health/treatments/what-is-optical-coherence-tomography> [Accessed 28 May 2018].
- En.wikipedia.org. (2018). Choroidal neovascularization. [online] Available at: https://en.wikipedia.org/wiki/Choroidal_neovascularization [Accessed 28 May 2018].
- Porter, D. (2018). What Are Drusen?. [online] American Academy of Ophthalmology. Available at: <https://www.aao.org/eye-health/diseases/what-are-drusen> [Accessed 28 May 2018].
- Bausch.com. (2018). Diabetic Macular Edema – Information about Diabetic Macular Edema or DME : Bausch + Lomb. [online] Available at: <http://www.bausch.com/your-eye-concerns/diseases-and-disorders/diabetic-macular-edema> [Accessed 28 May 2018].