

Machine Learning Engineer Nanodegree

Eugene Medina

eu63n3m@outlook.com

July 28, 2020

Capstone Project Proposal

Convolutional Neural Network Model For Predicting COVID-19 Presence in Chest X-Rays

Domain Background

Machine learning usefulness in the medical field has gained traction over the years and the technology, although not yet fully mature, is simply making healthcare smarter. Machine learning, being a subset of artificial intelligence, may be familiar to many in use cases such as speech recognition, fraud detection, learning associations to improve online shopping experience, predictions of weather in a specific area, or the identification of exoplanets that are technically classified as habitable. However, machine learning has been proven to have life-impacting potential in healthcare – particularly in the field of medical diagnosis. My objective in this capstone project proposal is to present one of the many possible means of determining the presence of Covid-19 in patients from a chest X-ray. Although this is just a proof of concept and by no means considered as a diagnostic tool, I think that the basic foundation built is that machine learning in healthcare can have a significant impact on the global fight against Covid-19.

Similar projects have had successful attempts to use machine learning in the analysis of chest X-rays. One such project is published by CodeMD found in this [link](#) [1]. Another useful project is from CodeSpeedy found [here](#) [2]. Perhaps the nearest inspiration to this project is the work of Blake VanBerlo and Matt Ross published in Towards Data Science found in this [link](#) [3], which explains how we can discover features of Covid-19 on chest X-rays and predict the presence of it using machine learning.

Problem Statement

Chest X-ray (CXR) interpretations have been very manual in the past, with specialist physicians having to look at X-ray pictures of patients as they are presented. Doctors have

used the ABCDE approach to carry out structured interpretation of a chest X-ray: **A**irway, **B**reathing, **C**ardiac, **D**iaphragm, and **E**verything else. Although the procedure has been deemed accurate by the medical community, resources have recently been overwhelmed when the Covid-19 pandemic struck. Doctors need tools to accurately, and with equally important mandate, quickly diagnose patients who are suspected to have Covid-19 – the disease caused by SARS-CoV-2 (designated a new type of coronavirus) – to provide appropriate healthcare and treatment. Test kits have been in limited supply and new studies suggest that chest X-rays and chest CT scans can help diagnose the disease.

Despite the success and promising efforts in applying machine learning in the diagnosis of Covid-19 through CT scans, I believe the fact remains that Covid-19 is an infection of a global scale and is likely to be experienced in communities of varying sizes, including those that are considered remote. X-rays are more accessible in these communities as they are inexpensive and quick to perform. They are generally more accessible to healthcare providers in most regions. This is not intended to replace existing means of testing for Covid-19 such as molecular (polymerase chain reaction or PCR test) and serological (testing for antibodies in blood and tissues) tests. I am optimistic that the data that will be derived as a result of the prototype algorithm applied to the model may help provide insights to the global effort in the identification and treatment of Covid-19.

This project aims to leverage the use of a supervised machine learning technique, particularly classification, where the output will have a defined label.

Datasets and Inputs

The construction of a public open dataset for chest X-rays and CT scans which are positive or suspected of Covid-19 or other types of coronavirus and bacterial pneumonias such as MERS, SARS and ARDS, have since been established in both the medical and scientific communities around the world when the pandemic struck. Data are still being continually added through indirect collections from hospitals and physicians.

I plan to use the Mila Covid-19 image dataset found [here](#). Patients who have been confirmed with Covid-19 through polymerase chain reaction (PCR) test with pneumonia may present X-ray images with a pattern that is only moderately characteristic for the human eye [4].

Another useful dataset can also be found [here](#), and is intended to be used for research purposes only [5].

The Kaggle competition dataset [RSNA Pneumonia Detection Challenge](#) will also be used [6]. This contains several de-identified chest X-rays and includes a label indicating whether the image shows evidence of pneumonia.

Please note that due to the scarcity of publicly available chest X-rays of severe Covid-19 cases, we will typically have an unbalanced dataset. The dataset utilized will have more chest X-rays that are not indicative of the presence of Covid-19 compared to chest X-rays that are

indicative of the presence of Covid-19. When a class is underrepresented, in this case data for Covid-19-positive chest X-rays, accuracy can be misleadingly high. Therefore, class imbalancing methods need to be applied to mitigate the effects of one class severely outweighing others. We can later introduce to the user an option in the configuration file to weigh the underrepresented class more heavily in addition to the calculated weights from the datasets.

Solution Statement

Once the datasets and inputs are pre-processed and cleaned up, I am planning to implement the solution using an open-source machine learning model, covid-cxr. It successfully predicts the Covid-19 presence from chest X-rays. This model is a continuing effort of the Artificial Intelligence Research and Innovation Labs at the city of London, Ontario, Canada. Covid-cxr is a deep convolutional neural network and allows both binary and multi-class classification. According to the documentation, the binary classifier was trained on approximately 1000 Covid-19-negative, and 76 Covid-19-positive chest X-rays. Given the relatively small set of training data, the model was able to achieve a remarkable and encouraging metrics on the test set – an AUC (area under the curve) of 0.9633 and a sensitivity of 0.8750.

To make the predictions more believable and the data output more trustworthy, I also intend to apply a well-known explainability algorithm known as LIME (or Local Interpretable Model-agnostic Explanations, PDF document found [here](#) [7]).

Benchmark Model

The objective of this model is to prove that machine learning can be implemented to determine Covid-19 evidence from chest X-rays and should not be treated as a medical diagnostic tool. A considerable clinical trial and experimentation coupled with tons of data, collaboration with both the scientific and medical experts are necessary if this is to become a diagnostic tool.

A very important thing to note is that we need to have explainable AI that enables the model to explain its predictions. If we want clinicians to have confidence in this model, we must have some means of comparing the result and ensure that the model is not picking up meaningless correlations. The use of LIME (Local Interpretable Model-agnostic Explanations) will be leveraged to serve as a benchmark model as this is explainable, extensible and has been well-tested. By using LIME, we can tell if the model is learning unintended bias. We can then eliminate that bias through additional features moving forward.

Covid-cxr is a good reference model to be used as a benchmark for the proposed project. This existing model allows both binary and multi-class classification. It has also provided results that are remarkably accurate - the binary classifier model's performance on the test set having an AUC of 0.9633, a sensitivity (recall) of 0.875, and accuracy of 0.92. Conclusions and

overall performance of the proposed project will be drawn from the results and compared against the existing model.

Evaluation Metrics

Training metrics can be visualized in TensorBoard, or simply tabulated and interpreted on a graph once a training experiment is completed. Binary classification output will be displayed. Binary classification involves detection whether a chest X-ray shows evidence of Covid-19. The classifier is trained to assign X-ray images to either a non-Covid-19 class or a Covid-19 class. To compare the results of our model, results gathered from the neural network Covid-cxr AUC (area under the curve) and sensitivity will be utilized, as well as precision and accuracy. We can use these established metrics to evaluate the effectivity of our model.

Project Design

As a good practice for any machine learning project, data collected must be subjected to pre-processing and cleanup first. After the dataset has been pre-processed, a Pandas data frames will be created for the filenames and labels, then corresponding images of the dataset will be saved. A neural network model will then be trained, and its weights saved. LIME will be used to generate interpretable explanations for the model predictions on the dataset. Since model will not be deployed in production, or used as an official tool for medical diagnosis, an opportunity exists for further experimentation. We can train multiple models and save the best one to use. The introduction of hyperparameters will be an important part of the standard machine learning workflow – we can later introduce different combinations of hyperparameters on the model's performance metrics. Once a trained model is produced, we can obtain predictions and explanations for a list of images. Prediction results can then be saved in a .csv file containing image file names, predicted class, model output probabilities, and a file name of its corresponding explanation.

References

- [1] CodeMD <https://www.codemd.co.uk/xray-classifier/#:~:text=Machine%20Learning%20with%20Chest%20X-Rays%20Chest%20X-Rays%20are,highlight%20those%20with%20a%20high%20risk%20of%20pathology>
- [2] CodeSpeedy <https://www.codespeedy.com/detection-of-covid-19-from-chest-x-ray-images-using-machine-learning/>
- [3] Towards Data Science <https://towardsdatascience.com/investigation-of-explainable-predictions-of-covid-19-infection-from-chest-x-rays-with-machine-cb370f46af1d>
- [4] Joseph Paul Cohen. Postdoctoral Fellow, Mila, University of Montreal <https://github.com/ieee8023/covid-chestxray-dataset>

- [5] Core COVID-Net Team <https://github.com/agchung/Figure1-COVID-chestxray-dataset>
- [6] Radiological Society of North America, et. al. Kaggle <https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/data>
- [7] LIME pdf documentation <https://arxiv.org/pdf/1602.04938.pdf>