

Project Proposal

Objective

The purpose of our project is to apply machine learning for the diagnosis of fourteen distinct thoracic diseases through the classification of chest X-ray images and interpretation of patients' pertinent medical information, such as their symptoms and clinical history. The objective is to train and develop a model that achieves a diagnosis accuracy greater than that of practicing radiologists.

Justification

Chest X-rays are the most frequently used medical diagnostic images, with over two billion X-rays taken annually. Additionally, there is a global shortage in radiologists capable of diagnosing these radiographs. This insufficiency is highlighted in countries such as Rwanda and Liberia, where populations of twelve and four million are served by eleven and two radiologists, sequentially. Consequently, there are not enough radiologists to properly analyze and diagnose the large numbers of chest X-rays.

Additionally, the diseases that are diagnosed through chest X-rays tend to be the most serious and deadliest diseases. These diseases include pneumonia, which hospitalized over a million Americans in the last year alone; tuberculosis, which currently affects $\frac{1}{4}$ of the world's population; and lung cancer, which is the deadliest cancer for men and second deadliest cancer for women.

These factors combined make a large, automated method for the diagnosis of these X-rays imperative. Not only would a computer be able to process and diagnose the X-rays at a faster rate than radiologists, but it would also be immune to human limitations, such as fatigue or inattentiveness, thus boosting its accuracy. A model with the capability of providing an accurate diagnosis would be able to be implemented in a plethora of environments. For example, in the aforementioned African countries of Rwanda and Liberia, this model could theoretically replace radiologists and assist in providing diagnoses; even in the United States, this model could serve to significantly expedite the process of diagnosis and lower healthcare costs.

Description

The first step in this project will be to replicate the DualNet model developed by researcher's at Philip's Research Institute in Cambridge, Massachusetts. DualNet functions by accepting two inputs: a frontal X-ray image and a lateral X-ray image. Subsequently, the inputs are processed through two separate convolutional neural networks, both of which use a modified DenseNet-121 architecture. The outputs of the two networks are concatenated together; subsequently, a final, fully connected layer provides a fused output with a vector of probabilities associated with each disease. Two DualNet models will be trained; one model for posteroanterior and lateral X-ray inputs and another for anteroposterior and lateral X-ray inputs. The models will be trained on the MIMIC-CXR dataset, which contains frontal and lateral X-rays with positive and negative labels for fourteen diseases. Post-replication, a new network will be developed to account for patient symptoms and history.

Limitations

There are several potential limitations to be encountered in this project. The first likely limitation is not receiving access to the code of the DualNet model. Not receiving their code would provide an axiomatic disadvantage in replicating their work; however, this limitation is not terminal as a new model could still be developed. Another possible limitation is the magnitude of the dataset that we will use to train the model; our dataset contains 371,920 chest X-ray images, making it the largest dataset of chest X-rays ever published. This massive number of X-rays will take a very long time to process and will likely need substantial computing/graphics power. However, this limitation is likely to be negated by the utilization of Amazon Web Services (AWS), a cloud-based computing platform. The last foreseeable limitation is the inability to fuse the replicated DualNet model with a new network that accounts for patient symptoms and clinical history. Since this aspect of the project has never before

been done before, the specific method that will be most conducive to a successful model remains to be seen. Experimentation with various network architectures will be needed to surpass this limitation.

Feasibility Study

Currently available resources

Mr. Writer, Dr. Gantz and Dr. Crowe are all AOS personnel with experience in machine learning. Mr. Writer is our primary mentor for the development and progression of this project as he can provide guidance and expertise on numerous aspects of our project. Additionally, we have already established contact with Dr. Alistair Johnson, a researcher at MIT, and Pranav Rajpurkar, a doctoral candidate at Stanford University, and received valuable information as to the procedure that we will follow in the development of our model.

I currently have a MacBook Pro laptop with a 3.3 GHz Intel Core i7 processor and an Intel Iris Graphics 550 graphics processor.

We have received access to the MIMIC-CXR dataset and the Indiana University chest X-ray collection. The former dataset will be used to train the replicated DualNet model, while the latter dataset will be used to train the added network that interprets patient symptoms and history.

I am familiar with the basic syntax of Python and have learned the Linear Algebra and Multivariate Calculus needed for machine learning.

Additional resources needed

We will need to establish contact with researchers at the Philip's Research Institute in an effort to receive more information about their model and potentially receive their code, as well. This is important as their model is an important and essential building block for our project. Specifically, we will need to learn more detail about their development and training process so that their model can be successfully replicated and subsequently improved upon.

A computer with faster processing speed and more powerful GPU may be necessary. An example of a laptop that meets these criteria is the Alienware Area-51m laptop, which starts at \$1949.99. Additionally, AWS will be used in place of the traditional external GPUs. AWS will be utilized during the development and training of the model to make the process more convenient and allow work to be done at home.

No additional supplies will be needed.

In depth knowledge of Python and its syntax will be needed, in addition to a familiarity with machine learning in general. To accomplish this, I will take online courses in Python and machine learning over the summer and will practice developing my own convolutional neural networks to gain familiarity. Additionally, I will complete a course in principal component analysis (PCA), a statistical procedure that is integral to neural networks.

The budget for this project is estimated to be around \$2000 if a more powerful laptop is necessary. If not, the budget is likely to be approximately \$100.

Risk assessment

There are no likely safety issues presented in this project. However, safety measures, such as keeping liquids away from laptops, will be taken.

Alternatives explored to achieve objective

An alternative that was explored to achieve our objective was to create separate networks that do not fuse together into a unified output. Essentially, each network would provide its own diagnosis independent of the other networks. We chose not to pursue this method for multiple reasons. Each network providing an independent output would lead to conflicting predictions, causing the user to be unable to effectively use the model for the purpose of a diagnosis. Additionally, real-life diagnoses made by doctors simultaneously account for frontal X-rays, lateral X-rays, and patient information simultaneously; thus, any model seeking to rival doctors in accuracy must also do the same. Thus, it was decided that the three networks (frontal X-ray, lateral X-ray, and patient information) would be distinct but would lead to a single output.