

## PROJECT SPECIFICATION

# Pneumonia Detection from Chest X-Rays

## Exploratory Data Analysis

CRITERIA	MEETS SPECIFICATIONS
The student can create visualizations of the metadata that inform model training	<ul style="list-style-type: none"> <li>Students create distributions of diseases and comorbidities in their dataset</li> <li>Students create distributions of basic demographics of the patients who make up their datasets (such as age, gender, patient position, etc.)</li> <li>Students can use the above distributions to draw conclusions about how they will need to set up their model training</li> </ul>
The student can visualize relevant properties of pixel-level data	<ul style="list-style-type: none"> <li>Students use python's imshow to visualize medical images during EDA</li> <li>Students create distributions of intensity values of the pixel-level data within images and compare them both <b>within and across</b> diagnoses</li> <li>Students use both of these methods of inspecting images to draw meaningful conclusions about what their model will train on</li> </ul>

## Model Building & Training

CRITERIA	MEETS SPECIFICATIONS
The student creates an	<ul style="list-style-type: none"> <li>Students create a set of training data and a set of validation data that each have the appropriate</li> </ul>

CRITERIA	MEETS SPECIFICATIONS
appropriate train-test split of the data	proportions of positive and negative cases for their intended use (training and validation)
The student implements appropriate data augmentation to their training data	<ul style="list-style-type: none"> <li>• Student implements a class such as ImageDataGenerator from Keras to augment their training data only</li> <li>• Student <b>should not</b> augment testing/validation data</li> <li>• Student uses types of augmentation that are appropriate for medical imaging. There are no required types of augmentation</li> <li>• Students should normalize the imaging data so the model weights do not go to infinity.</li> </ul>
The student evaluates the performance of their model using the appropriate statistics	<ul style="list-style-type: none"> <li>• Student monitors the training progress of their model using log loss</li> <li>• Student changes training parameters to avoid overfitting and compares performances of different training paradigms.</li> <li>• Student trains enough epochs until the loss is "stable"</li> <li>• After training, student uses precision, recall, and F1 score to actually evaluate the utility of their model.</li> <li>• Find a threshold to classify if an image is pneumonia or not.</li> <li>• Students should show precision-recall curve and a curve of F1-score vs. threshold</li> </ul>
The student can integrate their model with real-world medical imaging data	<ul style="list-style-type: none"> <li>• Student can check DICOM header for <b>image position, image type and body part</b> on ALL .dcm files to check validity for their model using the pydicom python package.</li> <li>• Student can read imaging data in from a .dcm file, preprocess the image and feed it into their model using the pydicom python package.</li> </ul>

## FDA Description and Validation Plan

CRITERIA	MEETS SPECIFICATIONS
The student can describe the intended population and the clinical impact of their model	<ul style="list-style-type: none"><li>• Student should provide an intended use statement</li><li>• Student should point to data from their EDA to describe who their algorithm is indicated for and what the clinical setting is in which their algorithm would be used</li><li>• Student should describe limitations of their algorithm and how false positives or false negatives might affect a patient</li></ul>
The student can describe how their model was designed and trained	<ul style="list-style-type: none"><li>• Students provide a flowchart or architecture diagram of their model</li><li>• Students should describe the DICOM checks they use before sending an image through their algorithm</li><li>• Students should describe the preprocessing steps they use.</li><li>• Students should describe the architecture of the classifier</li><li>• Students should describe augmentation and its parameters used</li><li>• Students describe the parameters used for training</li><li>• Students should show the behavior of training and validating loss</li><li>• Students should describe the performance statistics and threshold used in final validation</li></ul>
The student can describe the dataset used to train the algorithm and how the ground truth	<ul style="list-style-type: none"><li>• Students should provide information for the training set</li><li>• Students should provide information for the validating set</li><li>• Students should describe how the ground truth of the NIH dataset is created, the benefit and limitations.</li></ul>

ground truth was created CRITERIA	MEETS SPECIFICATIONS
The student describes how they would create a FDA Validation set, ground truth, and what performance metric they would hold their algorithm to for FDA validation of their model.	<ul style="list-style-type: none"> <li>• Students should describe the <b>ideal</b> dataset that they would receive from a clinical partner for their FDA Validation Dataset</li> <li>• Students should describe how they would <b>ideally</b> create ground truth for this FDA Validation Dataset</li> <li>• Students should describe the performance metric and <b>the metric value</b> that they would hold their algorithm to, supported by literature</li> </ul>

---

### Suggestions to Make Your Project Stand Out!

- Create some of your own custom image augmentation (such as different image filtering techniques) rather than solely using those predefined by Keras' ImageDataGenerator.
  - Try creating two 'nested' models to specifically predict pneumonia. One that predicts pneumonia and/or infiltrates at the top level, and then a second model that specifically predicts pneumonia from the positive cases returned by the first model.
  - Have your model output a class activation map in addition to a single binary prediction of pneumonia. This map will help a clinician to understand what the model is detecting as probable pneumonia in each image.
-

