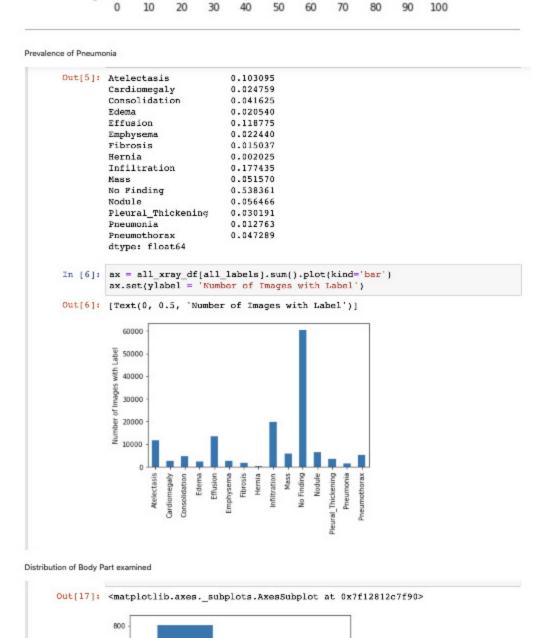
Search or jump to... Pull requests Issues Marketplace Explore □ k0f1 / Pneumonia-algorithm ♦ Code ① Issues ☼ Pull requests ⊙ Actions শ Projects ☐ Wiki ① Security ☑ Insights ⑥ Settings P Branch: master - Pneumonia-algorithm / FDA_Submission.md Go to file · · · k0f1 Update on FDA Latest commit 1eaaaaa 2 hours ago 🔞 History Al. 1 contributor 278 lines (143 sloc) | 18.8 KB Raw Blame 🖵 🖋 🗓 **FDA Submission** Your Name:: Kofi Ofuafor Name of your Device:: Pneumonia detection algorithm Algorithm Description 1. General Information Intended Use Statement: More than 1 million adults are hospitalized with pneumonia and around 50,000 die from the disease every year in the US alone. This algorithm is intended for use in assisting the clinician in acute care setting with a speedy detection of pneumonia within a chest x-ray. In otherwords it is to be used a a screening test for Pneumonia. The predicate device for this algorithm is a CADx device. Indications for Use: This algorithm is used as a Pneumonia classifier in a chest x-ray taken in acute medical or surgical admission units, accident and emergency and acute care wards. It is intended for use in people aged 20-70 years with no prior history of presence of edema, pleural effusion and pulmonary infiltrations Device Limitations: The device study was based on images from people of age 20-70 and was significantly less accurate in the presence of concurrent disease and therefore not recommended for use in the catgory of patients were there is prior presence of disease in a coloumn num_prior_positive in the demographic data. The algorithm was not traing in patients over 80 years old. Clinical Impact of Performance: · Performance statistics o Accuracy: 81% o Precision: 0.16666666666666666 o Sensitivity: 1 o Specificity: 0 For this algorithm, missing a detection of Pneumonia is not acceptable as human life and depends on it. False negative may slow down clinicians urgency to look at the patient xray for a physical examination. Therefore False Negative results have more impact than a False positive. Accordingly, the threshold for the algorithm has been weighed in favour of recall(sensitivity). When a test with high recall returns a negative result, you can be pretty confident that the result is truely negative. Recall does not take into account FP though, so you may still be labelling alot of negative cases as positive. So recall are good for screening tests. 2. Algorithm Design and Function Algorithm Fow chart The check_dicom function reads in a .dcm file, checks the important fields for our device, and returns a numpy array of just the imaging data. Dicom flow Chart Check dicom flow chart _ Architecture flow chart_ Architecture Diagram → Dense layers ← — The function, check_dicom(), is required to check the image type(MODALITY), Body Part Examined (Chest), and Image position (PA or AP) for each DICOM image and check if the input to the algorithm is valid to be predicted by our algorithm or not. DICOM Checking Steps: There are three steps: 1. Extract the desired attributes *Modality: DX The image type DX is a digital radiography *Body Part Examined: Must be Chest *Patient Position: AP or PA The image position could be AP or PA *Study Description: label This is the label for the chest x-ray. 2. Next the check_dicom function checks each image if the input to the algorithm is valid to be predicted upon by our device. 3. Finally, valid images are included as input and a numpy array of the image data is returned. Preprocessing Steps: The preprocess_image function used here takes in the image value returned from check_dicom, and the image mean, standard deviation as well as image size corresponding with that of VGG16 input image size of (1,224,224, 3). These inputs are used to standardize and resize the image which is returned as a processed image. Image augmentation techniques were also applied prior to training the model especially to training data. CNN Architecture: The pretrained model was loaded and the first 17 layers were rendered untrainable Here is the pretrained model Model: 'model_1' Layer (type) Param # Output Shape input_l (InputLayer) (None, 224, 224, 3) #2 block1_conv1 (Conv2D) (None, 224, 224, 64) 1792 block1_conv2 (Conv2D) (None, 224, 224, 64) 36928 Con.,Layer Block 1 blockl_pool (MaxPooling2D) (None, 112, 112, 64) block2_conv1 (Conv2D) (None, 112, 112, 128) 73856 blockZ_conv2 (Conv2D) (Mone, 112, 112, 128) Con. Layer Block2 block2_pool (MaxPooling2D) (None, 56, 56, 128) block3_conv1 (Conv2D) (None, 56, 56, 256) 295168 block3_conv2 (Conv2D) (Mone, 56, 56, 256) 590080 590080 #10 block3_conv3 (Conv2D) (None, 56, 56, 256) block3_pool (MexPooling2D) (None, 28, 28, 256) block4_conv1 (Conv2D) (Mone, 28, 28, 512) 1180160 block4_conv2 (Conv2D) 2359808 (None, 28, 28, 512) Con Layer Block 4 block4_conv3 (Conv2D) (None, 28, 28, 512) 2359808 block4_pool (MaxPooling2D) (None, 14, 14, 512) block5_conv1 (Conv2D) (None, 14, 14, 512) 2359808 block5_conv2 (Conv2D) (None, 14, 14, 512) 2359808 Con Layer Block 5 2359808 block5_conv3 (Conv2D) (None, 14, 14, 512) block5_pool (MaxPooling2D) (Mone, 7, 7, 512) Total params: 14,714,688 Trainable params: 14,714,688 Non-trainable params: 0 Seven fully connected layers were added to the pretrained model as classifiers. Of the 7 layers, three were Dropout layers to prevent input_1 False block1_conv1 False block1_conv1 False block1_conv2 False block1_pool False block2_conv1 False block2_conv2 False block2_conv2 False block3_conv1 False block3_conv2 False block3_conv3 False block3_pool False block4_conv1 False block4_conv2 False block4_conv3 False block4_pool False block5_conv1 False block5_conv2 False Model: 'sequential_1' Layer (type) Output Shape Param # model_1 (Model) (None, 7, 7, 512) 14714688 flatten_1 (Flatten) (None, 25088) (None, 25088) dropout_1 (Dropout) 0 dense_1 (Dense) (Mone, 1000) 25089000 Fully Connected classifier dropout_2 (Dropout) (None, 1000) dense_2 (Dense) (None, 500) 500500 dropout_3 (Dropout) (None, 500) 0 dense_3 (Dense) (None, 1) 501 Total params: 40,304,689 Trainable params: 27,949,009 Non-trainable params: 12,354,880 3. Algorithm Training · Types of augmentation used during training o horizontal_flip = True, o vertical_flip = False, o height_shift_range= 0.1, o width_shift_range=0.1, o rotation_range=20, o shear_range = 0.1, o zoom_range=0.1) Batch size o Training: 16 o Validation: 32 o Performance evaluation: 100 Optimizer learning rate o Learning rate: 5e-6. For 5 trainable layers. Decaying rate per epoch: 0.5 Layers of pre-existing architecture that were frozen 17 Layers of pretrained model were frozen and therefore not trainable block1_conv2 False block1_pool False block2_conv1 False block2_conv2 False block2_pool False block3_conv1 False block3_conv2 False block3_conv3 False block3_pool False block4_conv1 False block4_conv3 Felse block4_conv3 Felse block4_pool False block5_conv1 False block5_conv2 Felse Layers of pre-existing architecture that were fine-tuned * Fully connected layers were fine-tuned. Layers added to pre-existing architecture One one convulated was layers were fine tune and four layres were added for training Layer (type) Output Shape Param # model_1 (Nodel) (None, 7, 7, 512) 14714688 flatten_1 (Flatten) (None, 25088) dropout_1 (Dropout) (None, 25088) dense_1 (Dense) 25089000 (None, 1000) Fully Connected classifier dropout_2 (Dropout) (None, 1000) dense_2 (Dense) (None, 500) 500500 dropout_3 (Dropout) (None, 500) dense_3 (Dense) (None, 1) 501 Total params: 40,304,689 Trainable params: 27,949,809 Non-trainable params: 12,354,880 Algorithm training performance visualization Train and validation loss/accuracy versus epochs. Training Loss and Accuracy on Dataset 0.85 0.80 -5 0.70 0.65 0.55 - __ 0.50 -The model is converging and this can be observed from the decrease loss and validation loss with epochs. It reached 50% validation accuracy in just 1 epoch. The plot of how the model distinguih between the classes. In [109]: # fpr, tpr, threshold = plot_suc(valY, pred_Y)
plot_roc_curve(performances['ground_truth'], performances['probability']) 1.0 - Pneumonia (AUC:0.57) 0.8 False Positive Rate This is plot of two important parameters for the calculation of F1 Score. In [110]: # plot precision v recall
plot_precision_recall_curve(performances['ground_truth'], performances['precision_recall_curve(performances['ground_truth']) - Pneumonia (AP Score:0.24) 0.8 0.4 0.2 This is a plot of F1 Scores v Thresholds. In [62]: # Plot F1 Score plot_fls(valY, pred_Y) 0.20 0.10 0.05 0.00 0.35 Threshold 0.25 Highest f1 score is 0.34 and this corresponds to a threshold value of 0.38 Best_threshold = 38 The decision model for this threshold was based the calculation two thresholds, from F1 scored weighted by a precision value of 0.8 and recall value of 0.8 and the point of highest f1 score both converged on recall weighted performance It is important for the algorithm to pick all postive cases of pneumonia. Missing cases of pneumonia is not accepetable. Recall is more valuable under the circumstance where a negative test means you have a high chance of not having the disease. 4. Databases (For the below, include visualizations as they are useful and relevant) Description of Training Dataset: Data set was split 80:20 percent in favour of training data By using the len of positive cases of pneumonia to randomly select equal number of negative cases. Training data set Not Pneumonia Sizes of train and validation data after discarding excess data Training data set 1200 600 Not Pneumonia Description of Validation Dataset: Validation after split 80:20 with 20% of the entire data in the validation dataset By finding using random sample to add a negative data to the in the ration of 1:4 postive to negative cases Validation data set Validation set must contain 20:80 ratio (positive to negative) Validation data set 1200 600 -5. Ground Truth The NIH dataset radiologist reports were mined using Natural Language Processing (NLP) to create the disease labels from the associated radiology reports. It is expensive and laborous to hire a radiologist to label images often require several of them through a voting mechanism to obtain ground truth. NLP is fast inexpensive considerd to be >90% accurate. It can be associated with erroneous labelling. Is consideredd to be the silver standard level of ground truth, with gold standard being lung biopsy which is impractical. Llimitations Labels obtained by natural language are not a substitute for a radiologist acquired labels as they are not as reliable. It requires a large amount of data to obtain a decent level of accuracy. There is the ethical consideration as to whether whoc will be liable for any inaccuracy arising from disease labels obtained by NLP. These limitations can impact negatively on the model's ability to detect Pneumonia acurately. Inaddition, concurrent diseases with similar intensity values can mimic pneumonia. Often the diagnosis of pnueumonia is coupled with history physical signs and blood results all of which are not available to the model. 6. FDA Validation Plan Patient Population Description for FDA Validation Dataset: You will simply describe how a FDA Validation Plan would be conducted for your algorithm, rather than actually performing the assessment. Describe the following: The patient population that you would request imaging data from from your clinical partner. Make sure to include: Age ranges Sex Type of imaging modality Body part imaged Prevalence of disease of interest Any other diseases that should be included or excluded as comorbidities in the population Provide a short explanation of how you would obtain an optimal ground truth Provide a performance standard that you choose based on this Gender distribution 800 700 600 500 400 300 200 100 Σ Ags Distribution Ages group trained on 25000 20000 15000

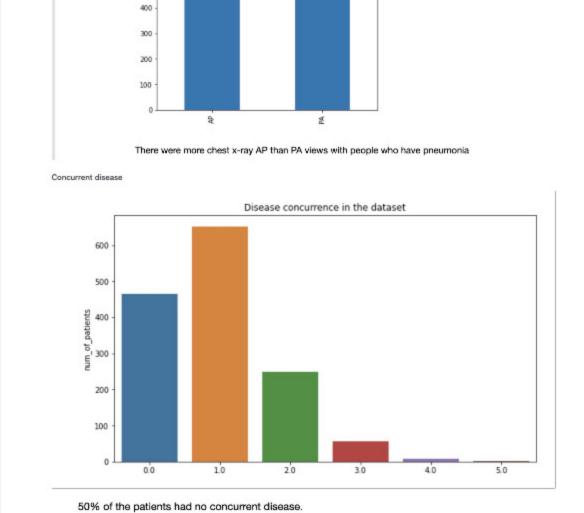
⊙ Unwatch 1 🖒 Star 0 😲 Fork 0



10000

5000

700 -600 -



Ground Truth Acquisiion methodology The NIH dataset radiologist reports were mined using Natural Language Processing (NLP) to

Silver standard The silver standard involves hiring several radiologists to each make their own diagnosis of an image. The final diagnosis is then determined by a voting system across all of the radiologists' labels for each image. Note, sometimes radiologists' experience levels

Often times, the gold standard(a test with the highest sensitivity and accuracy) is unattainable for an algorithm developer.

1. Biopsy-based labeling. Limitations: Impracticable, difficult and expensive to obtain.
2. NLP extraction. Limitations: may not be accurate.
3. Expert (radiologist) labeling. Limitations: expensive and requires a lot of time to come up with labeling protocols.
4. Labeling by another state-of-the-art algorithm. Limitations: may not be accurate.

Algorithm Performance Standard: Recall weighted threshold, was used to minimise the number of false negatives to the bearest minimum.

create the disease labels from the associated radiology reports.

are taken into account and votes are weighted by years of experience.

The other options to available to establish the ground truth to compare my algorithm are:

Best F1 scores of my model is 0.34 compares favorably with a radiologist average F1 score of 0.38 by Pranav Rajpurkar and Jeremy Irvin et al

Inc. Terms Privacy Security Status Help Contact Bithub Pricing API Training